

# Gene-regulatory circuitry of disease risk and progression

by

Carles Boix

B.S., Molecular Biology with certificates in Quantitative and Computational Biology  
Princeton University (2015)

Submitted to the Program of Computational and Systems Biology  
in partial fulfillment of the requirements for the degree of

Doctor of Philosophy in Computational and Systems Biology

at the

Massachusetts Institute of Technology

May 2022

© Massachusetts Institute of Technology 2022. All rights reserved.

Author.....  
Carles Boix  
Computational and Systems Biology Graduate Program  
May 10<sup>th</sup>, 2022

Certified by.....  
Manolis Kellis  
Professor of Computer Science, Massachusetts Institute of Technology  
Thesis Supervisor

Accepted by.....  
Christopher Burge  
Professor of Biology, Massachusetts Institute of Technology  
Director, Computational and Systems Biology Graduate Program



# Gene-regulatory circuitry of disease risk and progression

by  
Carles Boix

Submitted to the Program of Computational and Systems Biology  
on May 18th, 2022 in partial fulfillment of the requirements for the degree of

Doctor of Philosophy in Computational and Systems Biology

## Abstract

Complex diseases act heterogeneously through a remarkable diversity of cellular and functional outcomes across the human body, primarily through epigenetic organization and gene regulation. Genetics is a powerful tool to shed light on genes involved in disease, but we need maps of gene regulation and function at the tissue and cell type resolution to better understand disease mechanisms. Increasingly high resolution measurements of cellular epigenomes and transcriptomes allow us to observe this cellular heterogeneity at scale. To systematically model these, we require scalable statistical tools that can interpret and model gene regulation, its machinery, and complex transcriptional states. In this thesis, I build references of gene regulation and function in health and disease to interpret disease-linked genomic loci and develop methods to learn context-specific representations in order to understand how a single genome yields robust and diverse transcriptional outcomes through modularity of biological functions, how this heterogeneity is maintained, and how it breaks down over time and in disease.

In my first project, I build an integrative annotated reference of the human epigenome, systematically integrating multiple annotation projects covering hundreds of human cell lines, tissues, and states. I use this reference to map and dissect non-coding disease loci, specifically map multiple pleiotropic disease loci in coronary artery disease and dissect a locus showing tissue-specific gene involvement. In my second project, I model Alzheimer's disease (AD) progression across affected brain regions using single-cell transcriptomics, identify specifically vulnerable neuronal populations by brain region along the disease trajectory, and uncover pathways and neuronal circuits that may mediate AD vulnerability. To analyze large-scale transcriptomic references, I develop a fast and scalable method for calling high-resolution gene expression modules from single-cell data, use it to map the complex and modular glial changes underlying AD, and highlight metabolic and immune switches in cognitive impairment. In my third project, I investigate somatic mosaicism as a source of cellular dysfunction and to uncover missing genetic determinants and mechanisms in AD. To do so, I develop methods to map mosaic burden in individual cells jointly with expression and find increased cell type-specific somatic mosaic burden in dementia which I map to specific pathways and genes implicated in AD. Finally, I model neuronal trajectories through neurodegeneration in human brains and mouse models of AD, developing methods for building disease-driven pseudotime trajectories and mapping transcriptomic changes along paths to neuronal senescence driven by DNA damage accumulation.

Thesis Supervisor: Manolis Kellis, PhD

Title: Professor of Computer Science, MIT; Associate Member, Broad Institute of MIT & Harvard;  
Principal Investigator, Computer Science and Artificial Intelligence Lab of MIT

## Acknowledgements

This work would not be possible without the support and collaboration of many amazing people at and outside MIT. In particular, I would like to thank Manolis Kellis, under whose mentorship I've been incredibly lucky to work. Your unrelenting support, persistence, and kindness has buoyed me throughout my time at MIT. You create a remarkable, open, and highly collaborative environment and I am very thankful that I was able to work with you in your lab while at MIT.

My work on Alzheimer's disease would not have been possible without the expertise and mentorship of Li-Huei Tsai and members of the Tsai group - in particular Hansruedi Mathys, Leyla Akay, Gwyneth Welch, and Vishnu Dileep. Similarly, I have been lucky to work with great collaborators through ENCODE, including Jacob Schreiber, Chuck Epstein, Seth Strattan, Idan Gabdank, Anshul Kundaje, Bill Noble, and Brad Bernstein. Thank you for taking me in and for sharing your expertise and excitement in this research.

I have to thank Kellis lab graduate students past and present, including Max Wolf, Alvin Shi, Khoi Nguyen, Abhishek Sarkar, Kunal Bhutani, Benjamin James, Eloi Schmauch, Djuna von Maydell, and Neil Band, for being great friends in and outside of the lab and for your openness and frankness in discussing research, sharing your expertise, and challenging each other. I hope I was for you at least a small part of what you have been for me.

I would not have learned or enjoyed this work as much as I did without the mentorship of many others in the Kellis lab, including Wouter Meuleman, Yongjin Park, Maria Kousi, Patty Purcell, Shahin Mohammadi, Lei Hou, Na Sun, Suvi Kuosmanen, Leandro Agudelo, and Bo Zhao. Thank you for sharing your joy for biology and machine learning with me. You have been amazingly thoughtful, reliable, and kind and have inspired and supported me as both mentors and collaborators.

Wouter, I'm glad we overlapped in the lab, if only for a short time - you taught me how to think creatively and thoughtfully at the intersection of multiple fields. It's always good to see you - in addition to your acerbic humor, you have always been refreshingly direct and given me great advice. I count myself lucky that you have continued to be a mentor and a friend to me.

Yongjin, thank you for being a guiding hand for me in machine learning and statistics, introducing me to such a wide range of perspectives on and ways for modeling and thinking about biological data. You were always excited to discuss new ideas for methods and research directions and gave up a lot of your time to teach and mentor me.

Maria, you caught me at my lowest point in this journey and rekindled my love for biology. You taught me to think clearly and write concisely about a problem, when to be persistent, and when to be creative. Thank you for your clarity and patience through the rollercoaster of our work together.

Singing has been a constant source of community, joy, and creativity for me at MIT. Thank you to all of the amazing friends that I've had the chance to sing with over the years, including Thao Nguyen, Srinivasan Raghuraman, Allison Hamilos, and so many others, who are forthright, hilarious, and all amazingly talented. None of this is possible at MIT without the fantastic and wacky Bill Cutter, who fosters a great tight-knit choral family. Thank you Bill - you have always brightened my day and expanded my boundaries in many unexpected ways - I certainly did not come to MIT to act and sing opera, but I am so thankful I had a chance to do so here. I also want to thank Kerry Deal, who taught me far outside of my comfort zone and tempered my analytical side, showing me that you can't perform if you overthink.

Most importantly, I would like to thank my parents and my brothers for their love and support. I am incredibly lucky to have them. I would not be anywhere near where I am today without my mother's efforts, I can only emulate her inner strength, drive, patience, and generosity towards others.

This thesis is dedicated to her and to my father, for grounding me, for giving me perspective, and for teaching me to live in good humor and appreciate the absurdity and beauty of life.



# Table of Contents

<b>Abstract</b>	<b>3</b>
<b>Acknowledgements</b>	<b>4</b>
<b>Table of Contents</b>	<b>5</b>
<b>0 Introduction</b>	<b>8</b>
0.1 Mapping the epigenomic basis of cellular heterogeneity	8
0.2 Focusing on brain functional and transcriptional heterogeneity in health and disease	9
0.3 Alzheimer’s disease pathology and progression at the cellular and molecular level	10
0.4 Genetic risk factors of complex disease - from common to somatic variants	12
0.5 Thesis outline	14
<b>1 Regulatory genomic circuitry of human disease loci by integrative epigenomics</b>	<b>15</b>
1.1 Introduction	16
1.2 Generation of EpiMap, an epigenome atlas for interpreting gene-regulation and disease dissection	17
1.3 Characterization of the human epigenomic landscape across EpiMap	20
1.4 Learning enhancer modules, targets, regulators from the epigenomic atlas	27
1.5 Methodology for epigenomic datasets and processing	34
1.6 Epigenomic imputation methodology	35
1.7 Imputation as a tool for QC	37
1.8 Epigenomic annotations methodology	37
1.9 Learning latent representations of chromatin regulators	40
1.10 ENCODE Imputation Challenge summary and overview	51
1.11 Results of the ENCODE Imputation Challenge	55
1.12 Discussion	64
1.13 Supplementary Figures	65
<b>2. Dissecting disease-associated common variation and GWAS loci using a functional epigenomic atlas</b>	<b>81</b>
2.1 Introduction	82
2.2 Interpreting GWAS loci	82
2.3 GWAS-tissue co-enrichments	86
2.4 Dissecting multifactorial traits	88
2.5 Dissection of a multi-gene pleiotropic coronary artery disease-associated locus	93
2.6 Methodology for GWAS enrichment analysis	97
2.7 Discussion	98
2.8 Supplementary Figures	99
<b>3. A multi-region single-cell-resolution atlas of the aging human brain</b>	<b>107</b>
3.1 Introduction	108

3.2 A 1.3M nucleus transcriptional atlas of Alzheimer’s disease progression across six brain regions	109
3.3 Cohort selection and single-nucleus RNA-seq	112
3.4 snRNA-seq processing, QC, and annotation	113
3.5 Excitatory neurons are structured into highly brain region-specific subtypes	114
3.6 A unique thalamus-specific inhibitory neuron subtype	115
3.7 Glial subtype, regional, and functional diversity is annotated by gene expression modules	117
3.8 Methodology for gene expression and regulon modules	120
3.9 Discussion	123
3.10 Supplementary Figures	125
<b>4. Single-cell multi-region dissection of Alzheimer’s disease progression maps AD-specific regional selective neuronal vulnerability</b>	<b>182</b>
4.1 Introduction	183
4.2 Glial compositional and expression changes in AD	183
4.3 Expression modules separate AD-associated changes across glial cell types	184
4.4 Identification of selectively vulnerable excitatory neuron subtypes	193
4.5 Cell type composition analyses	199
4.6 Differential gene expression	200
4.7 Discussion	201
<b>5. Single-cell mosaicism analysis reveals cell-type-specific somatic mutational burden in Alzheimer’s Dementia</b>	<b>203</b>
5.1 Introduction	204
5.2 Cohort assembly and variant detection	205
5.3 Cohort and single-nucleus SMART-seq2 methods	207
5.4 Enrichment for somatic changes in AD patients	208
5.5 Methodology for mutation calling	208
5.6 Mutational Burden calculations	209
5.7 Cell-type identification through transcriptomic analysis	210
5.8 Methodology for transcriptional profiling and clustering	211
5.9 Cell-type specific enrichment of mosaic changes in the aging brain and in AD	212
5.10 Increased mutational burden correlates with expression changes at the single-cell level	213
5.11 AD-specific mosaic mutational burden concentrates in AD-related genes and pathways	214
5.12 Gene and pathway enrichment methodology	216
5.13 Discussion	217
5.14 Supplementary Figures	220
<b>6. Temporal modeling of Alzheimer’s disease progression</b>	<b>238</b>

6.1 Introduction	239
6.2 Temporal progression of neurons in the CK-p25 mouse model of AD towards senescence	239
6.3 DNA damage in the CK-p25 mouse model	248
6.4 Disease pseudotime by regressing out core axes of non-disease variability	249
<b>7. Conclusion</b>	<b>255</b>
7.1 Building gene regulatory models by moving from snapshot references to perturbations	256
7.2 Modeling variant function by resolving disease-risk across the allele frequency spectrum	256
7.3 Systematically learning biological components using representation learning	257
<b>References</b>	<b>259</b>

## **0 Introduction**

The human body is built from an incredibly heterogeneous complement of cells, all of whose functions and forms are derived from just an estimated ~20,000 protein-coding genes encoded in just 6.4 billion pairs of nucleotides present in tens of trillions of cells in each individual<sup>1</sup>.

Cellular heterogeneity, which is robustly wired and maintained across our lifespans, arises from incremental evolutionary processes of selection that reuse, modify, and rewire existing functions to improve and build new capacities. Crucially, cellular functions must efficiently use cell resources as well as be robust to genetic perturbations and stressors in the maintenance of cell identities. Due to these evolutionary pressures shaping cellular and functional heterogeneity, this balancing act of robustness, function, and efficiency through a gradual and greedy search is built in a complex and modular manner. Understanding biology through a modular lens offers hope for systematically understanding the massive diversity of biological molecular and cellular processes and states. Recent advances in biological methods enable us to quantify cells in human organs at the level of independent cells and to measure how different cell types modify the genome epigenetically to create functional cellular diversity.

My thesis work asks how a single genome yields a robust and diverse cellular heterogeneity through modularity of biological functions and how these attributes are maintained, perturbed, and break down over time. To do so, I build large-scale biological references of human cells and tissues in health, aging, analyze how they are perturbed in disease, and develop algorithms and statistical tools for their analysis. In doing so, I show that biology can inform our algorithmic and statistical approaches to understanding cellular heterogeneity and its maintenance. Just as modularity and robustness are core principles of biological organization, we develop algorithms that leverage and learn this modularity to understand cellular heterogeneity, its maintenance, and its dysregulation in disease.

### **0.1 Mapping the epigenomic basis of cellular heterogeneity**

A core way in which the human cells robustly maintain and specify this cellular diversity is through epigenetic organization of their DNA, modifying nucleosome residues and recruiting DNA binding and chromatin regulating proteins that can selectively activate and inactivate gene transcription. Mapping out the targeting, maintenance, and downstream effects on genes of these factors and modifications is key to understanding how genes are regulated to yield desired cell state outcomes either in cellular specification, creating highly diverse and long-lasting neurons, balancing cellular turnover in organs such as the digestive tract lining, or engineering finely calibrated structures such as the inner ear; or in response to stimuli, ranging from defense responses to viruses to cell-cell communication under localized stress to metabolic adaptation in nutrient starvation.

Understanding how transcriptional gene regulation works in the human genome is particularly challenging due to its size and complexity. About 98-99% of the genome is not directly protein coding - genes are spaced by large stretches of noncoding DNA that can be up to hundreds of kilobases in size and contain elements with roles in gene regulation and chromosome conformation<sup>2</sup>. Tissue-specific expression of genes is often coordinated by many gene-distal regulatory elements (enhancers) acting at long ranges from the gene in question. In addition, enhancer and promoter activation is modulated by combinations of DNA-binding and associated factors, both with and without specific DNA target recognition. These include over 1,500 factors, principally DNA-binding transcription factors (TFs)<sup>3,4</sup> and chromatin regulators and remodelers (CRs), all of which work combinatorially to specify a multitude of possible cellular and functional outcomes.

Advances in epigenomic mapping using techniques such as ChIP-seq, ATAC-seq, and DNase-seq, which profile the accessibility of DNA, modification of histones, or presence of

specific proteins along the genome, have enabled researchers to begin to map out histone modifications and the genome-wide localizations of DNA-associated factors. Large-scale experimental mapping<sup>2,5-7</sup> and integration of histone modification marks and DNA accessibility have helped annotate diverse classes of gene-regulatory annotations, including distal-acting/tissue-specific enhancers and proximal-acting/mostly-constitutive promoters<sup>8,9</sup>. These systematic, genome-wide profiles of gene regulation machinery across cell types reveal significant biological modularity that has been used to annotate the genome and enhance our understanding of gene regulatory circuitry. For example, the systematic evaluation of combinations of histone modifications has led to widely-used annotations of the human non-coding genome's functional roles and their heterogeneity across cell states<sup>8,9</sup>.

While epigenomic atlases have proven invaluable for systematically modeling and understanding the roles of the non-coding genome in gene regulation and disease, these maps also have considerable limitations. In particular, they miss many disease-relevant tissues, have variable quality, are prone to experimental noise and methodological variation between protocols, labs, antibody lots, reagents, batches, computational processing pipelines, software versions, integration pipelines. Moreover, consortia building these maps focus on a set of common marks across samples and often exclude samples missing some marks or marks missing in some samples, thus reducing biological space coverage. In the end, these efforts often only profile a few marks in many samples, or many marks in few samples due to cost limitations as the number of desirable experiments (including ChIP-seq of many DNA-associated factors, and conformation and accessibility measurements) and interesting tissues and cell states (from tissue-level to cell subtype-level and under a number of perturbations) grows.

While the continuing development of these maps should aid us to identify consistent and biologically relevant patterns of gene regulation, at the current pace, previous approaches to interpret data of increasingly higher dimension and size will not suffice. More recent efforts to build single-cell transcriptomic<sup>10</sup> and epigenomic<sup>11</sup> atlases of the human body have only exacerbated this trend, and such efforts are not properly compounding decades of data and research knowledge. Due to improvements in technology and continuing profiling for specific applications, much of this data is being left behind and sits unused even when it may still contain un-explored biology or could be used to boost newer cellular maps. In addition, modeling approaches that were previously reasonable to analyze a set of histone modification ChIP-seq tracks or bulk-level transcriptomic data will miss higher-level biological patterns extractable from high-dimensional data covering many samples, assay types, timepoints, and specific perturbations. Therefore, this growth in experimental data must be matched by algorithms that can handle the breadth and dimensionality of genome-wide epigenomic data and that learn and leverage biological modularity to understand human gene regulation and its perturbations in disease and health. Our work in this thesis pairs building, boosting, and interpreting these reference maps of cellular heterogeneity side-by-side with the development of resources, tools, and methods for their analysis. Our work shows how novel insights can be derived using simple linear algebra as well as recent advances in deep learning and Bayesian modeling to better understand the core axes of variation and structure of these large epigenomic and transcriptomic datasets.

## **0.2 Focusing on brain functional and transcriptional heterogeneity in health and disease**

Epigenetically organized cellular heterogeneity leads to thousands of diverse cell types, which form the building blocks of highly complex organs and systems in the human body. Among these, the human brain performs some of the most complex higher-order functions in our body while striking a remarkable balance between efficiency of function and robustness to

perturbation. The high diversity of functions and plasticity in the brain, which has undergone extensive evolutionary changes from our common ancestors, must often be maintained in individual cells across our lifespan. While recent reports suggest that hippocampal neurons have limited potential for adult neurogenesis<sup>12,13</sup>, the vast majority of individual neurons in the brain are post-mitotic and are thought to last through adult life<sup>14</sup>. In addition to robustly maintaining their identity and function over very long timescales, cells in the brain must operate efficiently, as the brain is highly energy-dependent, consuming an estimated 20% of our caloric intake<sup>15,16</sup>.

On either side of the age range, the brain is host to both a range of early-onset developmental neuropsychiatric disorders, believed to arise in part from abnormal specification and balance of cellular identity<sup>17,18</sup>, and a number of aging-associated neurodegenerative diseases<sup>19</sup>, which each affect specific regions of the brain, disrupt higher-order cognitive functions, and show rapid progression after their late onset. Certain neurodegenerative diseases have long been linked to specific neuronal subpopulations in the brain<sup>20,21</sup>, and recent research has proposed that whole neuronal circuits play roles in neurodegenerative disease vulnerability in specific regions as well as its establishment and progression through brain regions<sup>22</sup>. In addition, while genetic determinants of these diseases have been proposed through family, trio, and GWAS studies, among others<sup>23–26</sup>, it is still unclear how these diseases are triggered in age-dependent manners and rapidly progress across brain regions.

These unique requirements of cellular longevity, efficiency, robustness, and maintenance, coupled with the large number of diseases and disorders resulting from its dysregulation make the brain fascinating and crucially important to study through the lens of cellular heterogeneity and its maintenance. Until recently, molecular characterization of the brain to understand its transcriptomic and epigenomic heterogeneity was limited to very low-throughput techniques or bulk sequencing assays aggregating millions of cells at a time. However, recent technological advances can now measure the transcriptomic and epigenomic states of hundreds of thousands of individual cells at a time (single cell or sc techniques). These novel approaches enable the systematic quantification of the cellular composition and states of the brain and its changes at the level of independent cells in disease. These techniques, in conjunction with established and novel computational methodologies and approaches, promise to help map out cellular heterogeneity across brain regions and are particularly exciting for understanding its maintenance in neurodegeneration and its dysregulation in neurodevelopmental disease.

### **0.3 Alzheimer's disease pathology and progression at the cellular and molecular level**

Our work in this thesis has focused on Alzheimer's disease, a prevalent neurodegenerative disorder characterized by a stereotypical progression through brain regions that is not well understood at the level of molecular mechanism. Alzheimer's disease (AD) and other related neurodegenerative disorders are marked by pathological protein aggregates showing a stereotypical progression through human brain regions in neurodegeneration<sup>21,27,28</sup>.

Understanding the molecular underpinnings of the regional specificity, consistency, and cross-region progression of these disorders has enormous therapeutic potential for targeting both early stages and cross-region progression in these disorders. However, this is hindered in part by a lack of comprehensive molecular characterization and systematic dissection of these brain regions in disease.

The pathological hallmarks of AD include extracellular amyloid plaque deposits and intracellular neurofibrillary tangles (NFT) in neurons. AD is staged according to regional progression of NFTs, which first appear in the entorhinal cortex (Braak stages I-II), then the hippocampus and thalamus (III-IV), and finally progress to the neocortex (V-VI), a sequence that is most often synchronous with cognitive decline from mild memory impairment into severe dementia<sup>21,29–31</sup>.

Studies of AD genetic risk have identified multiple strong-effect variants that often lead to early disease onset, including the e4 allele of APOE and variants in APP, PSEN1, and PSEN2, which directly affect amyloid plaque build up. With the exception of APOE-e4, the majority of cases are sporadic, possibly coming from unresolved environmental factors or driven by the combined effect of many weak-effect variants across genes in a polygenic manner. These studies have highlighted the role of genetic variants linked to microglia (brain resident immune cells) in the AD brain, but AD-implicated genes are expressed across all brain cells. Amyloid precursor protein (APP), for example, is highly expressed in oligodendrocytes, which have primary roles in neuronal myelination and metabolic support. APOE itself is expressed in both microglia and synapse and scaffolding-associated astrocytes, and a number of genes are expressed across many cell types, such as CLU or PLCG2. A detailed transcriptomic map of these cell types in AD-affected regions that can put these genetic risks and other AD-implicated genes into a cellular and pathway context is therefore essential to systematically understand the disease from before its onset to memory loss and to late-stage dementia.

Cross-tissue single-cell atlases of human development and disease have recently proven very insightful for understanding aging<sup>10</sup> and diseases such as COVID-19<sup>32</sup>. Although many of the AD-relevant regions have been studied at the molecular level<sup>33–37</sup>, a comprehensive multi-region single-cell resolution atlas incorporating information from dozens of individuals is missing. This cross-region atlas is essential to ask detailed questions about the consistency of AD's progression across the brain. In particular, it remains to be determined how cellular diversity varies between regions and how it is affected by the progression of Alzheimer's disease. Furthermore, the transcriptional modules controlling the biology of the diverse cell types in each region, the common and brain region specific alterations of gene expression in AD, and the transcriptional signatures distinguishing vulnerable and resilient cell types remain to be determined. While many previous studies have focused on neuronal changes in early-affected regions and have identified single vulnerable excitatory neuron subtypes in Alzheimer's disease<sup>38–45</sup>, region-specific glial diversity and changes in AD are less well understood. In addition, no study has comprehensively tracked subtype vulnerability across the AD-involved brain regions<sup>21,29–31</sup> or between subtypes at different stages of AD. In addition, these studies have not been able to systematically ask what factors may mediate vulnerability at the molecular level, or whether vulnerable subtypes show transcriptomic changes from relatively more resilient subtypes or at different timepoints.

These analyses highlight how diseases perturb and dysregulate cellular heterogeneity, and enable researchers to identify affected cells and pathways at different stages of the disease. However, in Alzheimer's disease in particular, we are also interested in better mapping the trajectories of single neurons through disease in response to increasing pathological burden and glial cell dysregulation. In AD and other diseases, we are interested in understanding at what resolution these changes happen and in which order, and whether there are and how cells arrive at multiple disease-related cell states.

Single-cell maps of development often model cellular transitions through developmental trajectories by ordering cells along a "pseudotime" axis that correspond with time points and changing cell states. While these methods can make the strong assumption that the major axes of variation will coincide with developmental trajectory, this is often not the case with transcriptomes of cells from disease. In AD in particular, we must either disentangle or separately model the highly heterogeneous neuronal subtype-driven effects from disease-driven changes. For modeling disease, there is a pressing need for algorithms that can disentangle axes of variation in these large transcriptomic maps. Technique that can properly jointly model cell subtype and disease could also be expanded to model additional covariates, such as sex, age, or genetic background, and their development is therefore highly relevant to a number of diseases. In AD, developing these models and maps of neuronal trajectories can allow us to

track states before neuronal degeneration, prioritize targets and timepoints best suited to manipulation for arresting or treating neurodegeneration, and identify resilient states and trajectories towards these states that could be manipulated to alleviate and reverse cognitive decline.

#### **0.4 Genetic risk factors of complex disease - from common to somatic variants**

Uncovering the genetic basis of disease risk informs our understanding of how disease perturbs and dysregulates cellular states, highlighting key risk genes and pathways (such as APOE in AD), quantifying their effects on disease, pinpointing cell types of action, and predicting individuals' disease risk. Studies of genetic disease risk factors and variants, paired with detailed epigenomic and transcriptomic maps, have proven invaluable for elucidating the molecular basis of complex traits. Recently, genome-wide association studies (GWAS) have been extremely successful in discovering more than 100,000 genomic loci containing common single-nucleotide polymorphisms (SNPs) associated with complex traits and disease-related phenotypes, providing a very important starting point for the systematic dissection of the molecular mechanism of human disease<sup>46,47</sup>.

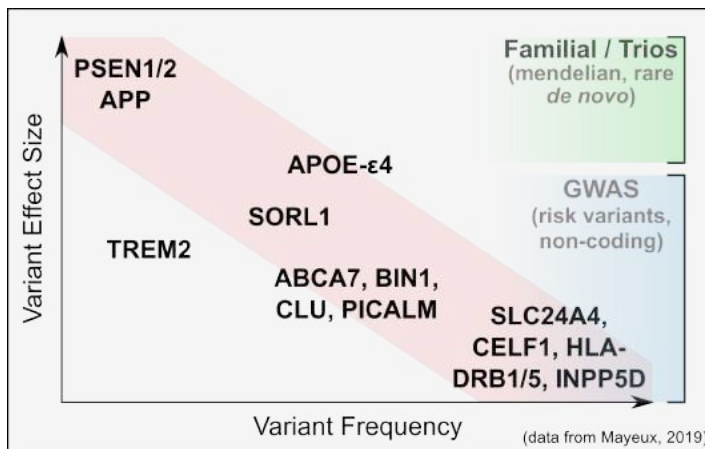
Transcriptomic maps of disease-implicated tissues have enabled researchers to pinpoint the pathways and start to uncover mechanisms of action for bona-fide, well supported GWAS-prioritized genes<sup>37,48–50</sup>. However, the vast majority of these genetic associations remain devoid of any mechanistic hypothesis underlying their molecular and cellular functions, as more than 90% lie outside protein-coding exons and likely play non-coding roles in gene-regulatory regions whose circuitry remains unresolved<sup>51,52</sup>. Maps of the non-coding genome across cell states reveal preferential localization (enrichment) of trait-associated genetic variants in tissue-specific gene-regulatory elements<sup>5,52–58</sup> and fine-map likely causal genetic variants in enriched annotations<sup>56,59–61</sup>. For example, GWAS risk loci from Alzheimer's disease are highly enriched in the enhancers of microglia, the primary immune cells of the brain<sup>5,62</sup>. Together with dissection of specific loci, these epigenomic analyses have highlighted the fundamental roles of microglia in plaque clearance, inflammation, and cell communication in AD<sup>48,49,63</sup>.

Despite their success and ubiquity, GWAS analyses, downstream functional fine-mapping, and systematic locus dissection only shed a partial light on risk factors and molecular mechanisms of complex diseases. In particular, these analyses are only well powered to detect disease associations for common alleles in the human population, and likely bias discovery towards loci that aren't highly selected against and have smaller effects in disease. Previous work mapping out pedigrees of affected families and more recent work using deep exome sequencing and trios (mother-father-child) to detect novel and strong-effect inherited variation has proven more useful for identifying high effect variants, which are often rarer in the general population<sup>23,64,65</sup>.

However, these genetic risk analyses may often miss environmental factors and are not powered to properly integrate or detect the potential effects of many weak-effect alleles. In addition, identified genetic risk factors may all act at early points of disease progression, while disease mechanisms late in disease (e.g. neurodegeneration after microglial inflammation in AD) remain under-explored. Finally, late-onset neurodegenerative disorders, including Alzheimer's Dementia (AD), Parkinson's disease, Huntington's disease, and Amyotrophic Lateral Sclerosis are primarily sporadic. These disorders may arise through the combined effect of multiple environmental risk factors and many weak-effect common alleles<sup>66,67</sup>, though heritable forms have also been described, stemming from a small number of strong-effect rare alleles<sup>68–71</sup>. Hindering the search for effective therapeutic interventions, our current knowledge of genetic contributors is still greatly incomplete, in part due to their complex etiology, heterogeneous level of affectedness, diverse symptoms, and the interplay between genetic and environmental factors<sup>72,73</sup>.



Several genetic contributors to AD in particular have been elucidated, including: the  $\epsilon 4$  haplotype allele of the *ApoE* gene locus<sup>74</sup>; at least 25 loci from genome-wide association studies (GWAS)<sup>66,67</sup>, including *BIN1*<sup>75</sup>, *PICALM*<sup>76</sup>, *CLU*<sup>76,77</sup>, *CR1*<sup>77</sup>; a small number of genes discovered from family pedigrees, involving *APP*, *PSEN1*, *PSEN2*<sup>68-71</sup>; and an increasingly-large set of coding and non-coding loci from ongoing exome-sequencing and whole-genome sequencing (WGS) studies<sup>78,79</sup>. With the exception of the *ApoE* locus, genetic results in AD follow the commonly-observed inverse relationship between allele frequency and effect size, in the continuum between high-frequency weak-effect variants (from GWAS) and low-frequency strong-effect variants (from family studies), due to the impact of selection, maintaining strong-effect variants at low frequencies (**Fig. 0-1**)<sup>80</sup>. Despite these major advances however, the vast majority of AD cases lack any detectable genetic basis.



**Figure 0-1.** Illustration of relationship between variant frequency against variant effect size in detecting genetic risk factors for Alzheimer's disease and dementia<sup>80</sup>.

Beyond inherited genetic variants, genetic contributors to disease can include somatic mutations incurred either during early embryonic development<sup>81</sup> (thus affecting a large fraction of cells and tissues), in later stages of embryogenesis (thus confined to a specific tissue or cell type), or during childhood, adulthood, and late-life (thus affecting only a small number of cells). These mutations can stem from environmental exposures, disease-associated DNA damage, disruption of DNA repair processes, or through the natural process of error-prone mitotic cell divisions. The phenomenon describing such somatic mutations is referred to as mosaicism, given that the affected individual is a mosaic of cells of different genetic makeup, including both risk and protective alleles of disease loci. Unlike inherited mutations, late-stage somatic mutations are not subject to the strong selective pressures acting on germline variation, and thus can result in strong gene-perturbing effects, which can have a disproportionate impact on cellular phenotype. Increased somatic mutational burden in individual cells can increase the chance of single-cell dysregulation. Increasing numbers of dysregulated cells can subsequently manifest as tissue- and organ-level disruptions, providing an additional path for disease onset and progression, beyond inherited genetic components and environmental contributors.

Though significant progress has been made in studying the relevance of mosaicism in disease, with prominent examples from cancer and clonal hematopoietic malignancies<sup>82-84</sup>, congenital heart disorders<sup>85</sup>, and skin disorders<sup>86</sup>, the inaccessibility of brain tissue has hindered the study of mosaicism in neurodevelopmental and neurodegenerative disease<sup>87,88</sup>. Bulk DNA sequencing has been used to recognize high-frequency somatic mutations in autism spectrum disorders<sup>89-91</sup> and in hemimegalencephaly<sup>92</sup>, from their partial representation in WGS reads. For DNA-damage-repair disorders, a single-cell DNA sequencing study showed increased mutations in nine cases vs. controls<sup>93</sup>, profiling three to six neurons from each affected

individual. Thus, it remains unclear whether somatic mosaicism plays a role in neurodegeneration more broadly (outside DNA-damage-repair disorders), whether increased mutations are also found in non-neuronal cells (which are increasingly recognized to play important roles in AD<sup>37</sup>), and how DNA perturbations relate to gene expression changes at the single-cell level, and ultimately to phenotypic effects in brain function and cognition.

## **0.5 Thesis outline**

This thesis aims to understand how a single genome yields robust and diverse outcomes through modularity of biological functions, how this heterogeneity is maintained, and how it breaks down over time or in disease. Biological modularity that arises from evolutionary processes - balancing robustness and efficiency in novel functions - informs our algorithmic approaches, as we develop algorithms that leverage and learn this modularity in context-aware manners to understand cellular heterogeneity, its maintenance, and its dysregulation in disease.

In our work, we computationally integrate and annotate epigenomic assays to create the largest-to-date reference of epigenomic cellular heterogeneity for use in understanding gene regulation and mapping disease loci (**Chapter 1: EpiMap**). We develop and apply methods for learning gene regulatory modularity and use them to enhance and organize biological data to understand how the epigenomic landscape is shared and maintained, build improved references, and identify context-specific components in the maintenance of chromatin modifications (**Chapter 1: EpiMap, ENCODE Imputation challenge, chromatin regulators**). We develop methods to prioritize diseases to specific tissues and cell types of action, highlighting the dysregulation of cellular heterogeneity and gene regulation in disease, and show how disease-associated loci can act in various pleiotropic manners (**Chapter 2: EpiMap, coronary artery disease**).

We study cellular heterogeneity in the brain, one of the most functionally diverse and recently evolved organs. We carry out a census of cell types at the single-cell level across brain regions, develop new methods to learn modular patterns of gene usage at high resolution and in huge datasets to annotate functions and region-specific differences, and identify functions that are shared across or swap cell types of action in regions developed at different evolutionary timepoints (**Chapter 3: Cross brain-region atlas in the aging human brain**). We use these references to track how normal functions are disrupted in Alzheimer's disease at a cell type-specific level and to systematically evaluate circuits and pathways mediating neuronal vulnerability in neurodegeneration (**Chapter 4: AD progression across brain regions**).

We measure the breakdown of this cellular heterogeneity in conjunction with DNA damage, developing tools to quantify the accumulation of both single nucleotide and structural variants in mouse models of AD (**Chapter 5: CK-p25 DNA damage**) and in individuals with late-stage dementia (**Chapter 5: Somatic mosaicism in AD**). We study how individual cells respond to disease, studying senescence and the breakdown of neuronal cellular identity after DNA damage in mouse models of AD (**Chapter 6: CK-p25 neuronal senescence**). We develop methods that separate out disease-relevant variation from that linked to cellular identity and other biological covariates in order to place cells on a disease pseudotime trajectory and understand how neurons and other cells in the brain respond to increasing AD pathological burden (**Chapter 6: Disease pseudotime trajectories**).

# 1 Regulatory genomic circuitry of human disease loci by integrative epigenomics

Parts of this chapter have been presented/published as:

## **Regulatory genomic circuitry of human disease loci by integrative epigenomics**

*(Nature, 2021)*

Carles A. Boix<sup>1-3</sup>, Benjamin T. James<sup>1,2</sup>, Yongjin P. Park<sup>1,2,4</sup>, Wouter Meuleman<sup>5</sup>, Manolis Kellis<sup>1,2\*</sup>

1. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA, USA. 2. Broad Institute of MIT and Harvard, Cambridge, MA, USA. 3. Computational and Systems Biology Program, Massachusetts Institute of Technology, Cambridge, MA, USA. 4. Department of Pathology and Laboratory Medicine, Statistics, University of British Columbia, Vancouver, BC, Canada. 5. Altius Institute for Biomedical Sciences, United States. Correspondence to: e-mail: [manoli@mit.edu](mailto:manoli@mit.edu)

## **Learning hierarchically composable components of chromatin regulation and remodeling**

*(presented as a poster at the ENCODE meeting in 2018)*

Carles A. Boix, Wouter Meuleman, Charles Epstein, Bradley E. Bernstein, Manolis Kellis

## **The ENCODE Imputation Challenge: A critical assessment of methods for cross-cell type imputation of epigenomic profiles**

*(in preparation)*

Jacob Schreiber\*, Carles A. Boix\*, Jin wook Lee, Hongyang Li, Yuanfang Guan, Chun-Chieh Chang, Jen-Chien Chang, Alex Hawkins-Hooker, Bernhard Schölkopf, Gabriele Schweikert, Mateo Rojas Carulla Arif Canakoglu Francesco Guzzo, Luca Nanni, Marco Masseroli, Mark James Carman, Pietro Pinoli, Chenyang Hong, Kevin Y. Yip, Jeffrey P. Spence, Sanjit Singh Batra, Yun S. Song, Shaun Mahoney, Zheng Zhang, Wuwei Tan, Yang Shen, Yuanfei Sun, Minyi Shi, Jess Adrian, Richard Sandstrom, Nina Farrell, Jessica Halow, Kristen Lee, Lixia Jiang, Xinqiong Yang, J. Seth Stratten, Michael Snyder, Manolis Kellis, William Stafford Noble, Anshul Kundaje, ENCODE Imputation Challenge Participants

## 1.1 Introduction

Genome-wide association studies (GWAS) have been extremely successful in discovering more than 100,000 genomic loci containing common single-nucleotide polymorphisms (SNPs) associated with complex traits and disease-related phenotypes, providing a very important starting point for the systematic dissection of the molecular mechanism of human disease<sup>46,47</sup>. However, the vast majority of these genetic associations remain devoid of any mechanistic hypothesis underlying their molecular and cellular functions, as more than 90% lie outside protein-coding exons and likely play non-coding roles in gene-regulatory regions whose circuitry remains unresolved<sup>51,52</sup>.

Large-scale experimental mapping<sup>2,5-7</sup> and integration of histone modification marks and DNA accessibility helped annotate diverse classes of gene-regulatory annotations, including distal-acting/tissue-specific enhancers and proximal-acting/mostly-constitutive promoters<sup>8,9</sup>. These maps help elucidate the molecular basis of complex traits by revealing preferential localization (enrichment) of trait-associated genetic variants in tissue-specific gene-regulatory elements<sup>5,52-58</sup> and fine-mapping likely causal genetic variants in enriched annotations<sup>56,59-61</sup>.

However, these maps also have limitations: they miss many disease-tissues, have variable quality, are prone to experimental noise and methodological variation between protocols, labs, antibody lots, reagents, batches, computational processing pipelines, software versions, integration pipelines. Moreover, consortia that require common marks across samples often exclude samples missing some marks or marks missing in some samples, thus reducing biological space coverage, and often only profile few marks in many samples, or many marks in few samples due to cost limitations.

Here, we overcome many of these limitations and present a new human epigenome reference, EpiMap (for Epigenome Integration across Multiple Annotation Projects). We infer multi-mark chromatin state annotations<sup>8</sup>, and a high-resolution enhancer annotation that combines DNA accessibility and multiple enhancer states. We group enhancers into modules of common activity patterns, and infer candidate upstream regulators and downstream gene functions for each module based on motif and gene ontology enrichments. We also infer enhancer target genes using a machine learning approach (**Fig. 1-1a**). We integrate this high-resolution gene-regulatory circuitry with genetic associations, revealing traits with epigenomic enrichments, and predict causal variants and tissue-specific target genes. We distinguish unifactorial, multifactorial, and polyfactorial traits based on the diversity of their enriched tissues, and partition polyfactorial trait loci by their enriched tissues, revealing distinct biological processes and disease comorbidity patterns. We also distinguish monotropic vs. pleiotropic loci, and find that top-scoring loci frequently have multiple predicted driver variants, converging through multi-enhancer common-target, multi-gene common-tissue, or multi-gene multi-tissue pleiotropy patterns. Our results demonstrate the utility of dense, rich, multi-dimensional, high-resolution epigenomic and gene-circuitry annotations for gene-regulatory, complex-trait-dissection, and disease-locus-mechanism studies, resulting in unprecedented scale, scope, and coverage of biological space and disease complexity.

In addition to building a new atlas of the human epigenome through integration and imputation of epigenomic profiles across hundreds of biosamples, we designed and ran a public competition to challenge teams around the world to build next generation genomic imputation algorithms (the ENCODE Imputation Challenge). The challenge consisted of two major stages, included 50+ new experiments for evaluating method performance, and we collected and assessed over 20 cross-cell type imputation methods from almost 200 participants in the challenge. We describe regimes of under and over-performance of specific methods, challenges in evaluation of epigenomic imputation methods, and a new quantile normalization scheme for

epigenomic tracks.

Finally, we developed a general factorization method for learning biologically-relevant and genomic context-specific components in large datasets of epigenomic profile tracks, which we applied to DNA-associated factors to understand how chromatin regulators and transcription factors pattern across the epigenome in five human cell lines. Using this method, we identified genomic context-specific combinatorial interactions between chromatin regulators and between regulators and histone marks and transcription factors reported by ENCODE.

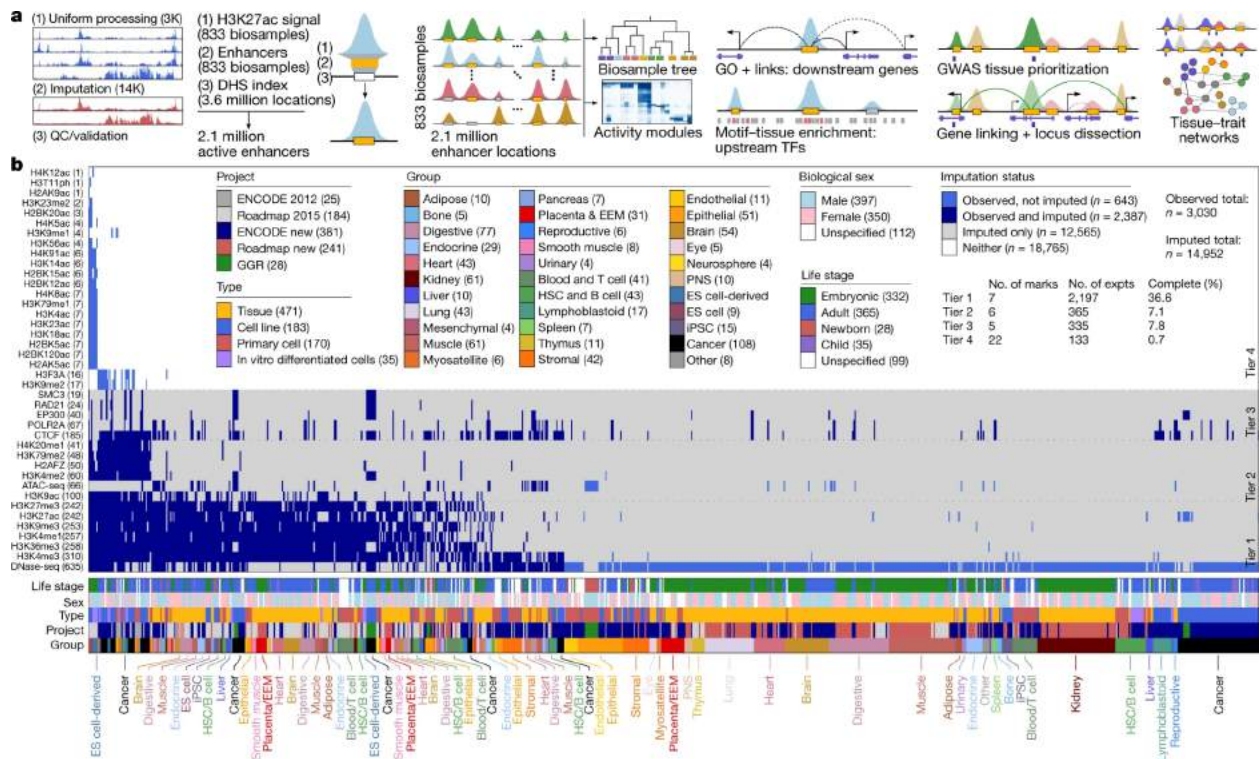
## **1.2 Generation of EpiMap, an epigenome atlas for interpreting gene-regulation and disease dissection**

We uniformly reprocessed 3,030 observed<sup>2,5-7</sup> tracks across 859 biosamples (406 ENCODE<sup>2</sup>, 425 Roadmap Epigenomics<sup>5</sup>, 28 GGR<sup>7</sup>) spanning 18 epigenomic assays, and computationally imputed<sup>94-96</sup> 14,952 tracks (**Fig. 1-1b, Supplementary Fig. 1-S1, Supplementary Table S1**), available for download and interactive visualization<sup>97</sup> at <http://compbio.mit.edu/epimap/>.

Our imputed tracks matched held-out observed tracks, both visually across randomly-selected regions (**Fig. 1-2a,b**), and quantitatively with >85% peak recovery and >75% average genome-wide correlation for punctate marks (59% of tracks) genome-wide (**Supplementary Fig. 1-S2**). Imputation was robust even with few supporting datasets, and performed best when target datasets showed >50% average correlation to their 10 nearest datasets, which held for 98% of single-assay samples (**Supplementary Fig. 1-S3**).

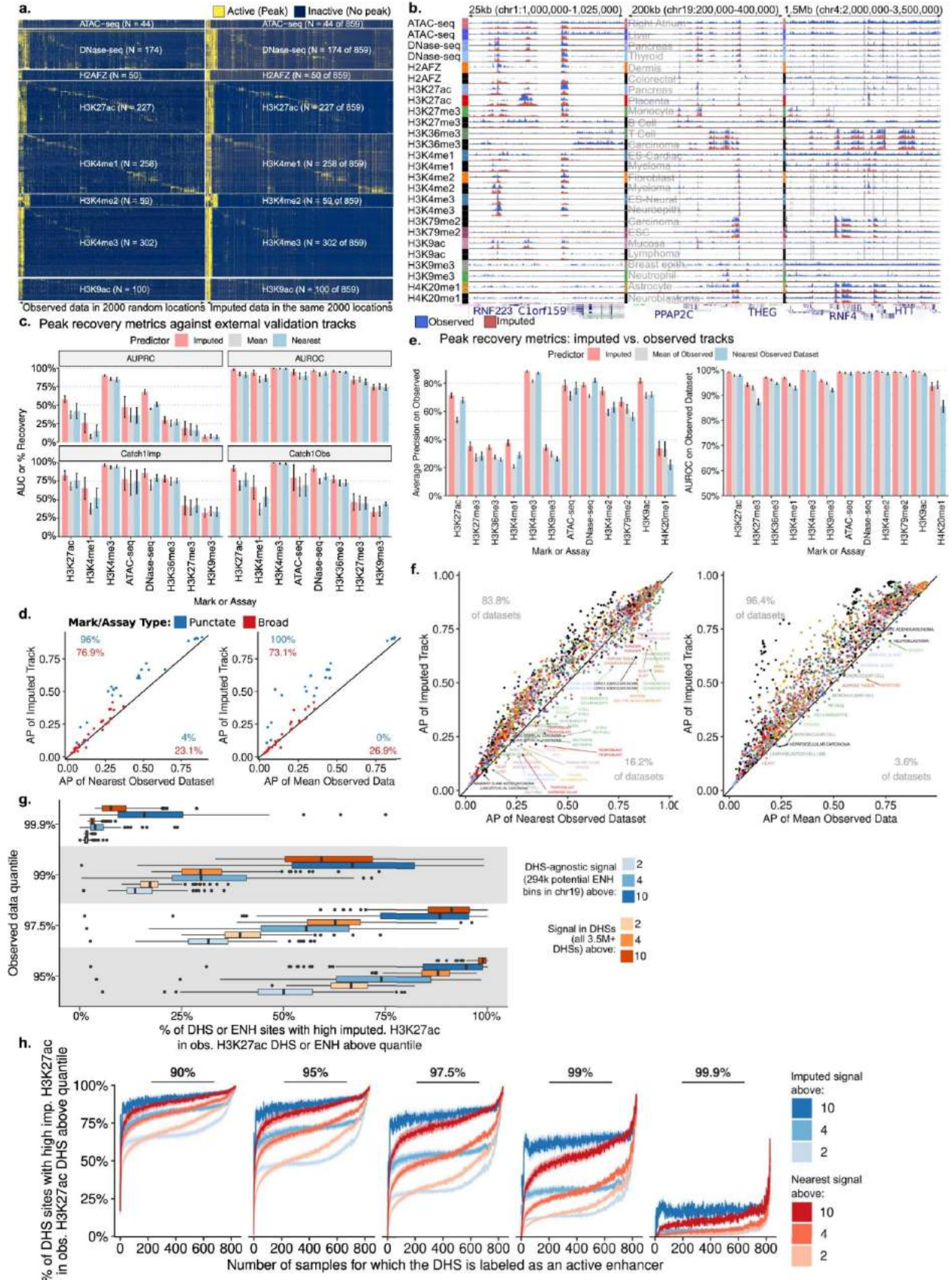
Imputed data also matched independent post-data-freeze experiments, outperforming “average-signal” and “nearest-track” benchmarks (in practice only knowable after generating the target track) for 96% of punctate marks and 77% of broad marks both genome-wide and specifically focusing on rare events (**Fig. 1-2c-h**).

Disagreement between imputed and observed tracks helped flag 138 potentially-problematic datasets, which independently also showed markedly-lower QC scores (**Supplementary Fig. 1-S2a-c**), and revealed potential sample or antibody swaps (**Supplementary Fig. 1-S4-5**) some of which were independently flagged by the data producers. Subtracting imputed-track signal from observed-track signal revealed 13 experiments with potential antibody cross-reactivity or secondary specificity (**Supplementary Fig. 1-S6-S8**). We removed from subsequent analyses the 138 flagged datasets and 442 tracks based solely on ATAC-seq or low-quality DNase-seq data, resulting in 2,850 observed and 14,510 imputed marks across 833 biosamples used in the remainder of this work.



**Figure 1-1. EpiMap resource overview. a.** We create a compendium of over 17k epigenomic tracks across 18 marks by uniform processing and imputation and use these to call chromatin states for 833 biosamples and active enhancer states over 2.1M DHSs. We learn unsupervised representations of the enhancer activities that enable us to call enhancer-downstream genes, upstream regulators, and prioritize, dissect, and compare hundreds of GWAS traits and thousands of loci. **b.** Data matrix across 859 samples (columns) and 40 assays (rows), ordered by number of experiments (parentheses) and colored by metadata.





**Figure 1-2. Imputation validation.** **a.** Heatmap of paired observed and imputed signal intensity across all punctate Tier 1 and Tier 2 assays across 2000 highest-max-signal bins among 5000 randomly-selected 25bp bins. Samples (rows) and bins (columns) are clustered and diagonalized using maximum imputed signal intensity, with broadly-active regions shown first. **b.** Paired observed (blue) and imputed (red) tracks for all Tier 1 and Tier 2 assays in three regions at different resolutions for randomly-selected samples. Each row shows a single track across three different resolutions. Full tracks at <https://epigenome.wustl.edu/epimap>. **c.** Genome-wide imputation performance metrics for predicting 51 external validation tracks across 8 assays in 14 biosamples (average precision, AUROC predicting top 1% of observed data and peak recovery of top 1% Imputed or Observed with top 5% Observed or Imputed, respectively) in chr1, shown for either the appropriate imputed track, the best-matching of the other observed tracks, or the observed signal average. **d.** Scatter comparison of average precision (AP) of imputed data with either nearest observed track or signal average in punctate (blue) and broad (red) marks. **e.** Genome-wide imputation performance metrics (AP, AUROC) for predicting observed tracks (evaluated on all observed tracks with an imputed prediction) in chr19, shown for either the appropriate imputed track, the best-matching of the other observed tracks, or the observed signal average. **f.** Scatter comparison of average precision (AP) of imputed data with either nearest observed track or signal average across all datasets, colored by sample group. Cases where the nearest sample or the mean heavily outperformed the imputation are labeled (points with over 25%, for nearest, or 10%, for mean, greater average precision than the imputed track). **g.** Sample-specific percentage of the 2M DHSs with imputed H3K27ac above a certain cutoff that are also in the top 10%, 5%, 2.5%, 1%, and 0.1% of 3.6M DHSs by matched observed datasets. **h.** Sample-specific percentage of the 2M DHSs with imputed (blue) or nearest observed (red) H3K27ac above a certain cutoff that are also in the top 10%, 5%, 2.5%, 1%, and 0.1% of 3.6M DHSs by matched observed datasets, partitioning the DHSs by the number of samples in which each DHS is called as an active enhancer.

### 1.3 Characterization of the human epigenomic landscape across EpiMap

The resulting compendium of 833 high-quality reference epigenomes, grouped into 33 tissue categories, represents a major increase in biological space coverage, with 75% of biosamples (624 of 833) corresponding to new biological specimens. Observed and imputed data co-clustered, but imputed datasets better captured the continuity between and within different sample categories, and more clearly revealed sample-type groups, likely driven both by their cleaner signal tracks and their sheer number. Moreover, distances between imputed datasets were less affected by technical covariates (**Supplementary Fig. 1-S9**) and were more consistent with sample groupings (**Supplementary Fig. 1-S10**).

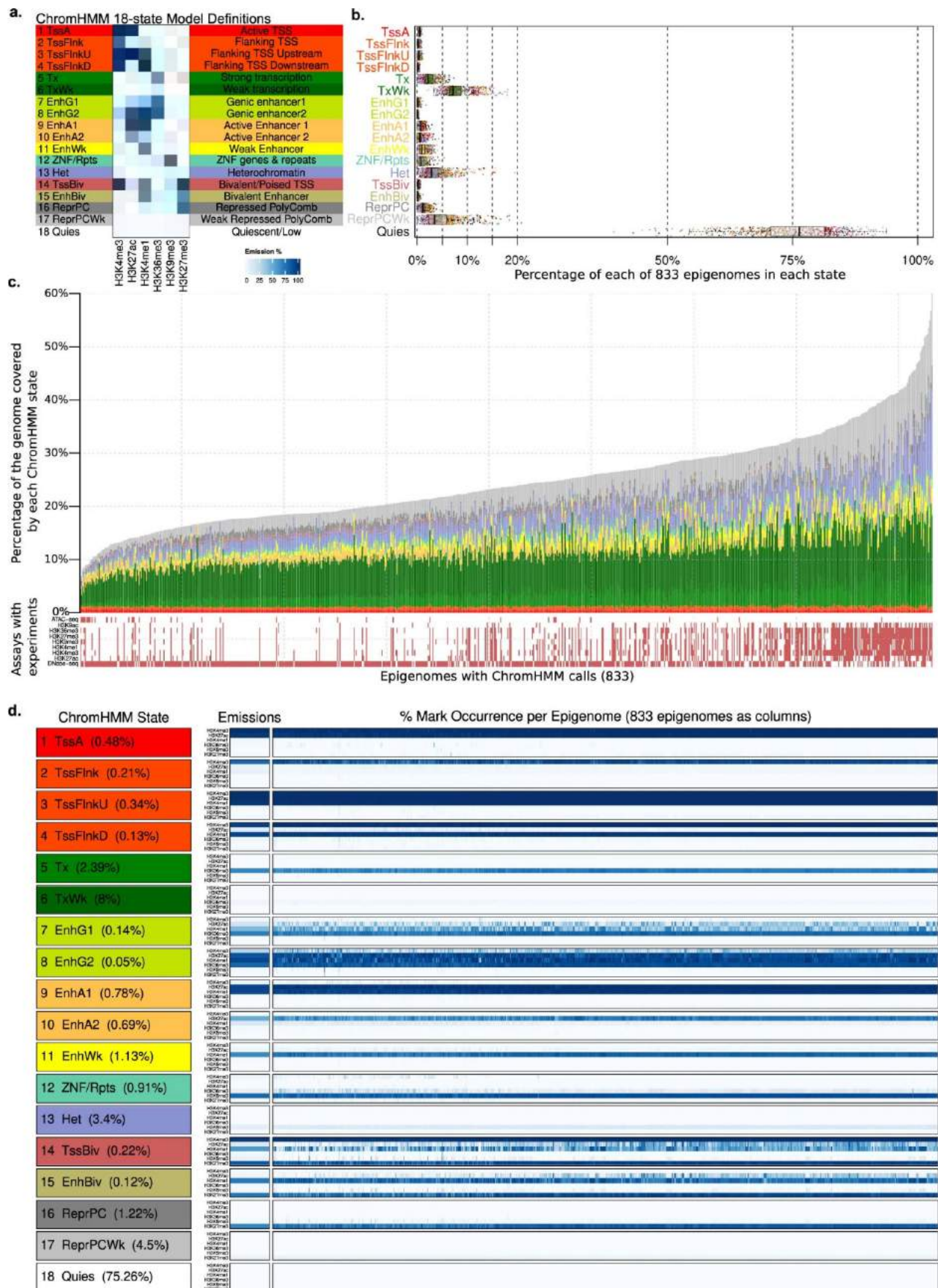
Hierarchical and two-dimensional-embedding clusterings of multiple marks using both genome-wide and relevant-region-specific correlation patterns (**Fig. 1-3a-c**) grouped biosamples firstly by lifestage (adult vs. embryonic) and sample type (complex tissues vs. primary cells vs. cell lines), and secondly by distinct groups of brain, blood, immune, stem-cell, epithelial, stromal, and endothelial biosamples within them. Active marks (H3K27ac, H3K4me1, H3K9ac) primarily grouped biosamples by differentiation lineage (blood, immune, spleen, thymus, epithelial, stromal, endothelial) and tissue type (lung, kidney, heart, muscle, brain), while repressive marks (H3K27me3, H3K9me3) instead captured lifestage (pluripotent, iPSC-derived, embryonic) (**Fig. 1-3b, Supplementary Fig. 1-S11**), consistent with previous studies<sup>98,99</sup>. Donor sex was not a primary factor in sample grouping.

We annotated genome-wide locations of 18 chromatin states<sup>2,5,8</sup> in all 833 biosamples using combinations of histone modifications, including multiple types of enhancer, promoter, transcribed, bivalent, and repressed regions (**Fig. 1-4a**), using a mixture of observed and scaled



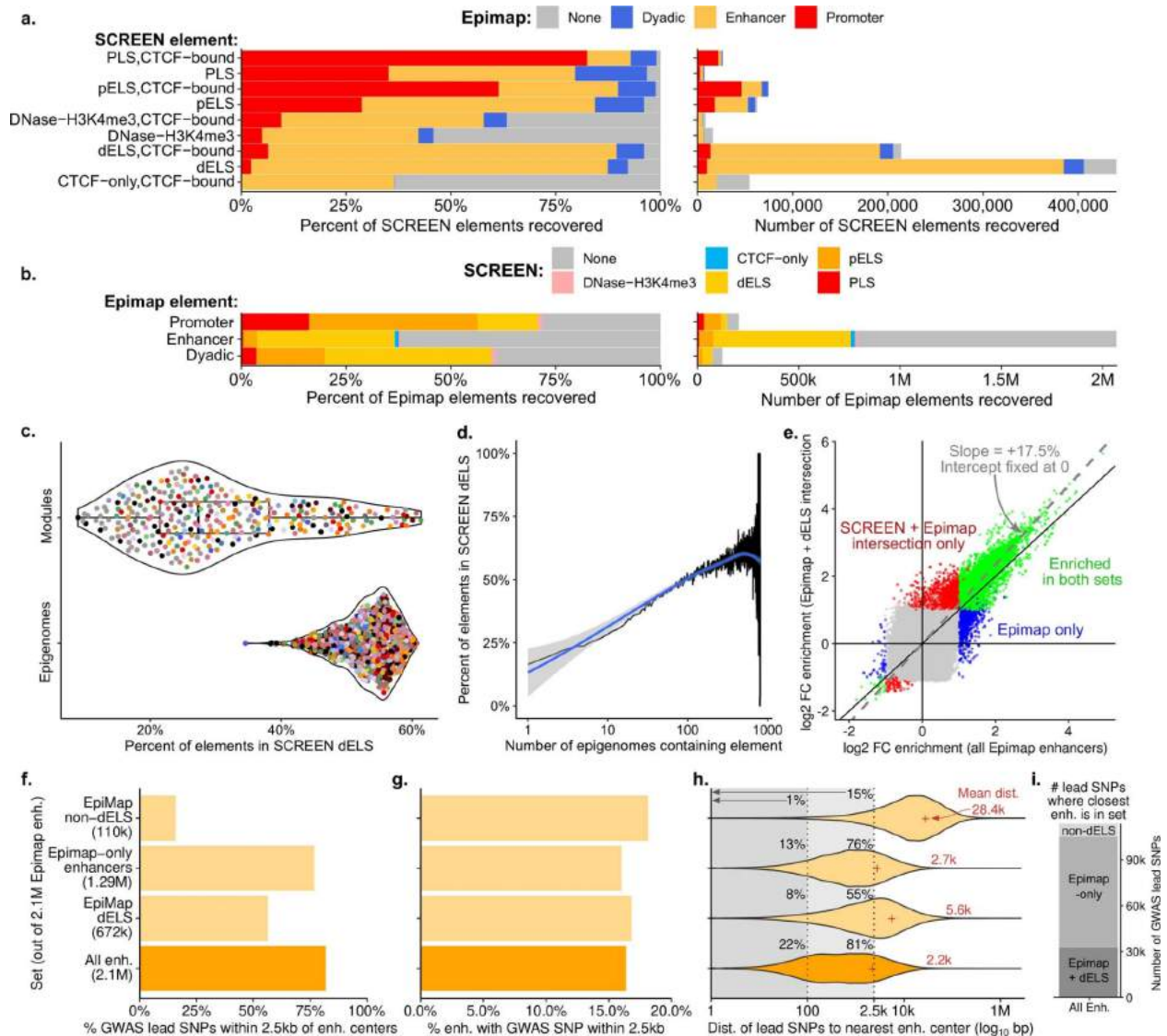


matrices shown. Clustering conducted on the fused matrix (left panel, constructed as in main figure). Observed data availability matrix (grey is available, white is unavailable) is shown for the top nine marks and accessibility assays by number of observed datasets. **b.** Two-dimensional embeddings of Tier 1 and 2 marks colored by tissue group, using Spearman correlation within matched chromatin states. Arrows point from the center of mass of all biosamples to that of the specified group. **c.** Hierarchical clustering of 833 biosamples based on enhancer activity distances (**Supplementary Fig. 1-S12**). Subtrees enriched for specific sample types are highlighted and labeled (colors). Samples are labeled by metadata in the outer ring (**Supplementary Table S1**).

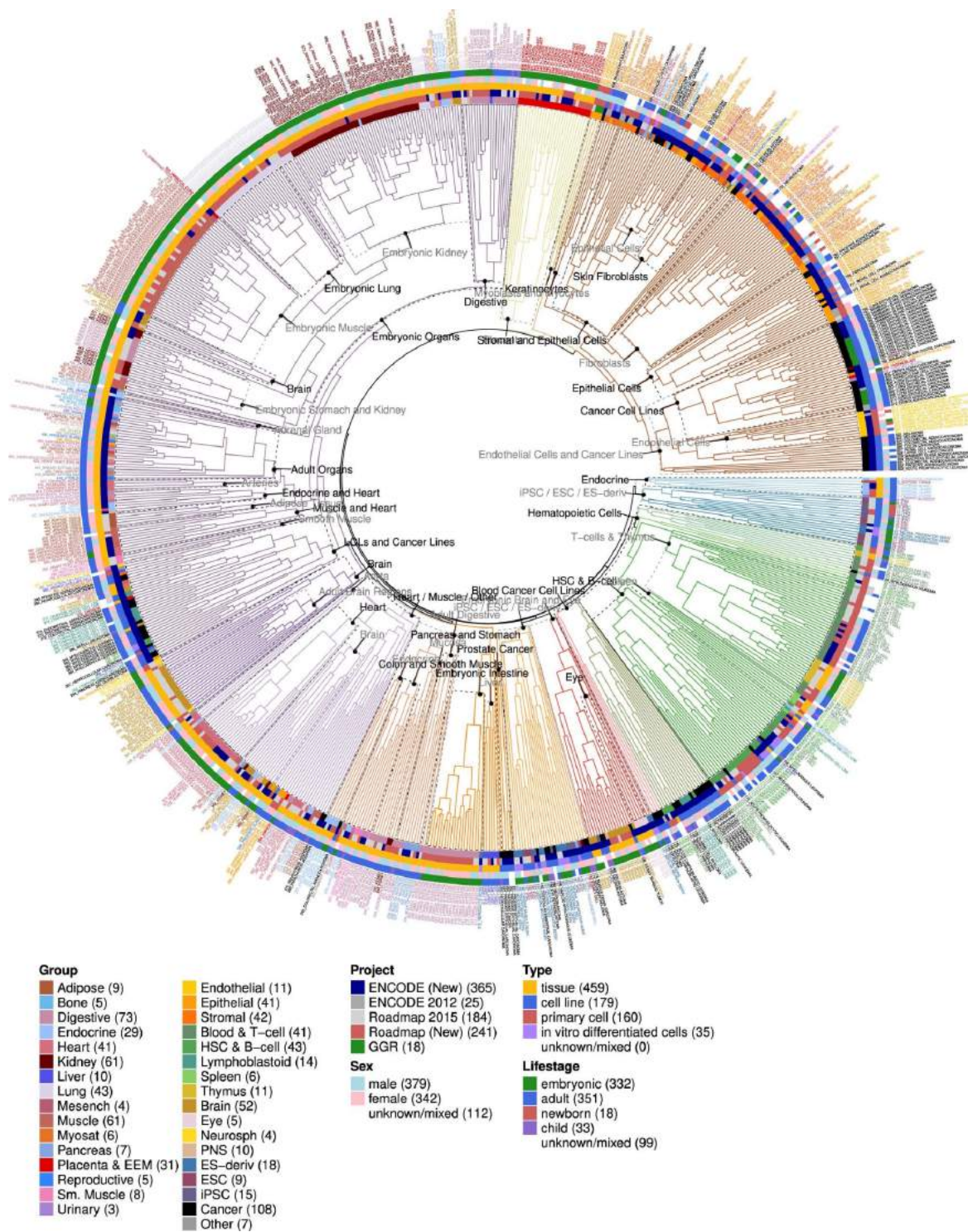


**Figure 1-4. Chromatin states.** **a.** Epigenomic state mnemonics for ChromHMM 18-state model (left) with emissions matrix (center) and state definitions (right). The 18-state model was trained on Roadmap data for the Roadmap 2015 paper. **b.** Distributions of per-state genome coverage (box-plots) across 833 biosamples (points, colored by tissue group) according to the ChromHMM 18-state model annotations. **c.** Genome coverage for each 18-state-model ChromHMM state across 833 biosamples after QC. Lower panel shows availability of 9 top marks, ordered by the number of observed datasets. Biosamples are ordered by percent of the genome not annotated as quiescent. **d.** Comparison of per state (left panel) model emissions (middle panel) against mark occurrence in state calls (right panel) across 833 biosamples (columns in right panel). Observed occurrence matched the emissions closely, with three exceptions, corresponding to bivalent chromatin states and transcribed enhancers (cumulatively covering 0.48% of the genome on average), which showed discrepancies for 12.1% of biosamples on average, likely stemming from their low frequency in the genome, and the frequent co-occurrence of H3K27ac and H3K4me1.





**Figure 1-5. Comparison with SCREEN.** **a.** Recovery of each category of SCREEN elements by each category of EpiMap elements, percentage (left) and number (right). **b.** Recovery of each category of EpiMap elements by each category of SCREEN elements, percentage (left) and number (right). **c.** Percentage of EpiMap enhancers in dELS in each of 300 modules and 833 epigenomes. **d.** Percentage of EpiMap enhancers in dELS by number of epigenomes containing element (blue line represents loess smooth). **e.** Comparison of motif-module log<sub>2</sub> fold change enrichments for all enhancers and for the intersection of enhancers and dELS. **f,g,h,i.** Comparison of enhancer sets within EpiMap enhancers (intersections with dELS, non-dELS, unique, and all enhancers), showing percentage of 113k pruned GWAS catalog lead SNPs within 2.5kb of enhancer centers (**f**), percent of enhancers with a GWAS SNP within 2.5kb of their center (**g**), the distribution of distances of GWAS SNPs to their nearest enhancer within each set (**h**), and the number of SNPs for which the nearest enhancer fell into each of the constituent sets (**i**).



**Figure 1-6. Enhancer-sharing sample tree.** Tree constructed from complete-linkage hierarchical clustering of Jaccard similarity matrix of enhancer activity across biosamples. 50 subtrees covering the full tree are annotated with their major subgroups. Tree is cut and colored to create 20 clusters for purposes of visualization. Leaves are labeled with metadata and reduced sample names and colored according to their tissue group. Metadata (heatmap track) from inside-out: tissue group, project, sample type, and donor sex and lifestage.

#### 1.4 Learning enhancer modules, targets, regulators from the epigenomic atlas

For each high-resolution active-enhancer region, we defined H3K27ac-based local activity levels across 833 biosamples, and used them to group enhancers into 300 enhancer modules (**Fig. 1-7a,b, 2-8a-c, Supplementary Fig. 1-S14-15**), including 290 tissue-specific modules (1.8M enhancers, 88%, activity in 2% of biosamples on average) and 10 broadly-active modules (251k enhancers, 12%, activity across 77% of sample categories on average).

Enhancer modules showed highly-significant, high-resolution, tissue-specific gene-ontology enrichments for neighboring genes (**Fig. 1-8, 2-9 and Supplementary Fig. 1-S16**), including: ion channels (brain modules); camera-type eye development (eye); neural precursor cell proliferation (neurosphere); endothelial proliferation, hemidesmosomes, and digit morphogenesis (endothelial, stromal, and epithelial); and organ development and morphogenesis (embryonic).

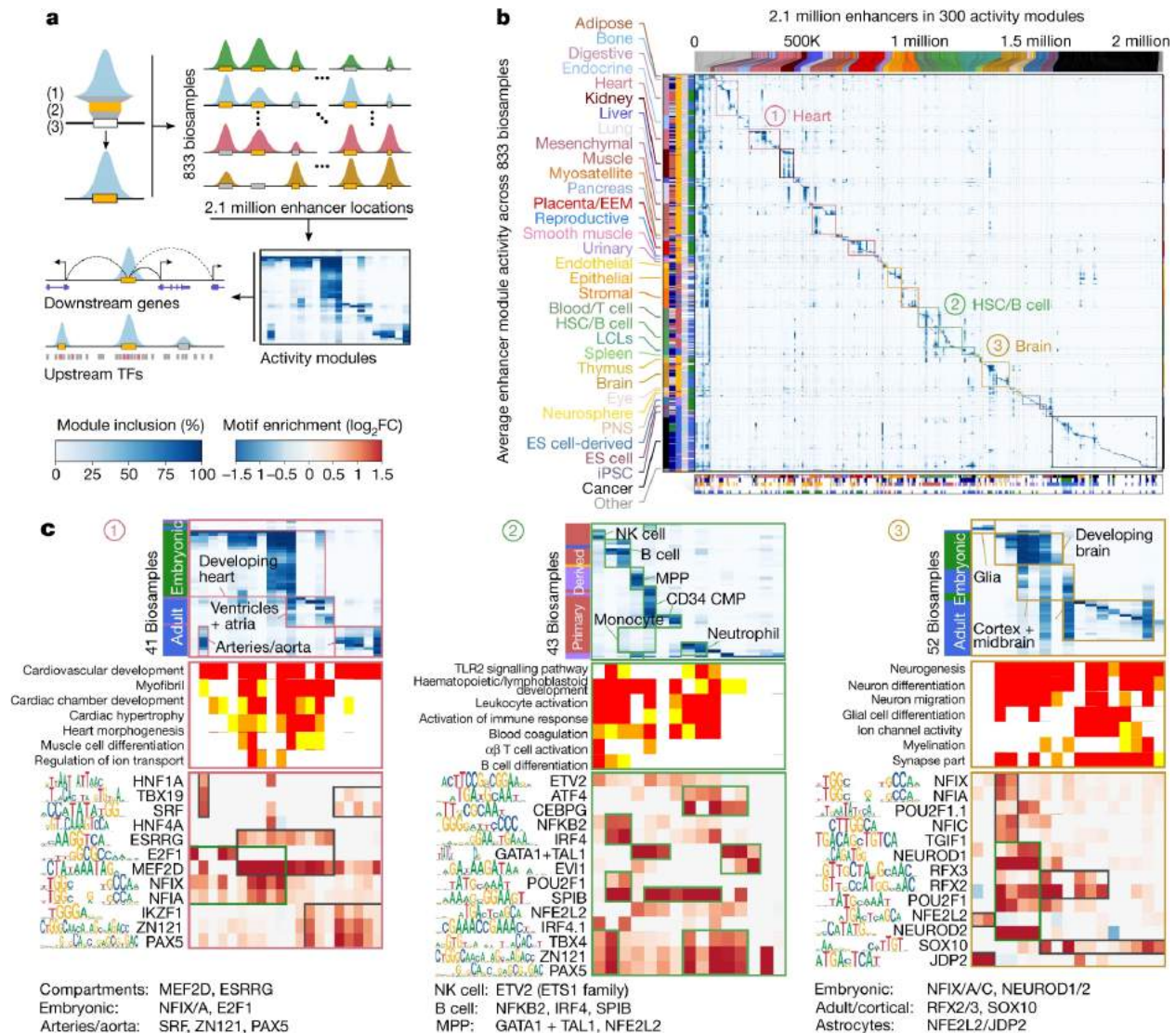
We predicted 3.3M tissue-specific enhancer-gene links by combining epigenomic-transcriptional correlation and genomic proximity, each gene linked to 13 enhancers, and each enhancer to 1.5 genes, at a median distance of 42k bp. Links were ~6-fold more specific than enhancers, and sample-specific links spanned larger distances than constitutive links (**Fig. 1-10a-f**). Our links outperformed previous linking approaches, using both gene-set enrichment metrics and curated gold-standards (**Fig. 1-10g,h**) and greatly expanded the biosamples with predicted links (from 127 to 833).

We predicted upstream regulators for 273 modules (91%), implicating 1175 motifs grouped into 160 archetypes<sup>102</sup> (**Fig. 1-8, 2-9, and Supplementary Fig. 1-S17**), including 152 tissue-specific motif archetypes (enriched in 6 modules on average) and 8 broadly-enriched (enriched in 53 modules on average). Specific motifs include: *GATA* and *SPI1* (blood/immune)<sup>103</sup>; *NEUROD2* and *RFX4* (brain/PNS)<sup>104,105</sup>; *KLF4* for digestive tissues<sup>106</sup>; and *TEAD3* for placenta, myosatellite, and epithelial cells<sup>107</sup>.

Broadly-enriched motifs revealed highly connected, combinatorially acting master regulators, including: *HNF1A* in liver, kidney, and pancreas (with *NR5A2*)<sup>108</sup>; *AP-1/JDP2* in immune, bone, and cancer<sup>109</sup>; and *TEAD3*, paired alternately with *MYF6* (myosatellite), *TFAP2A* (placenta), and *AP-1* (stromal) (**Fig. 1-9d,e**).

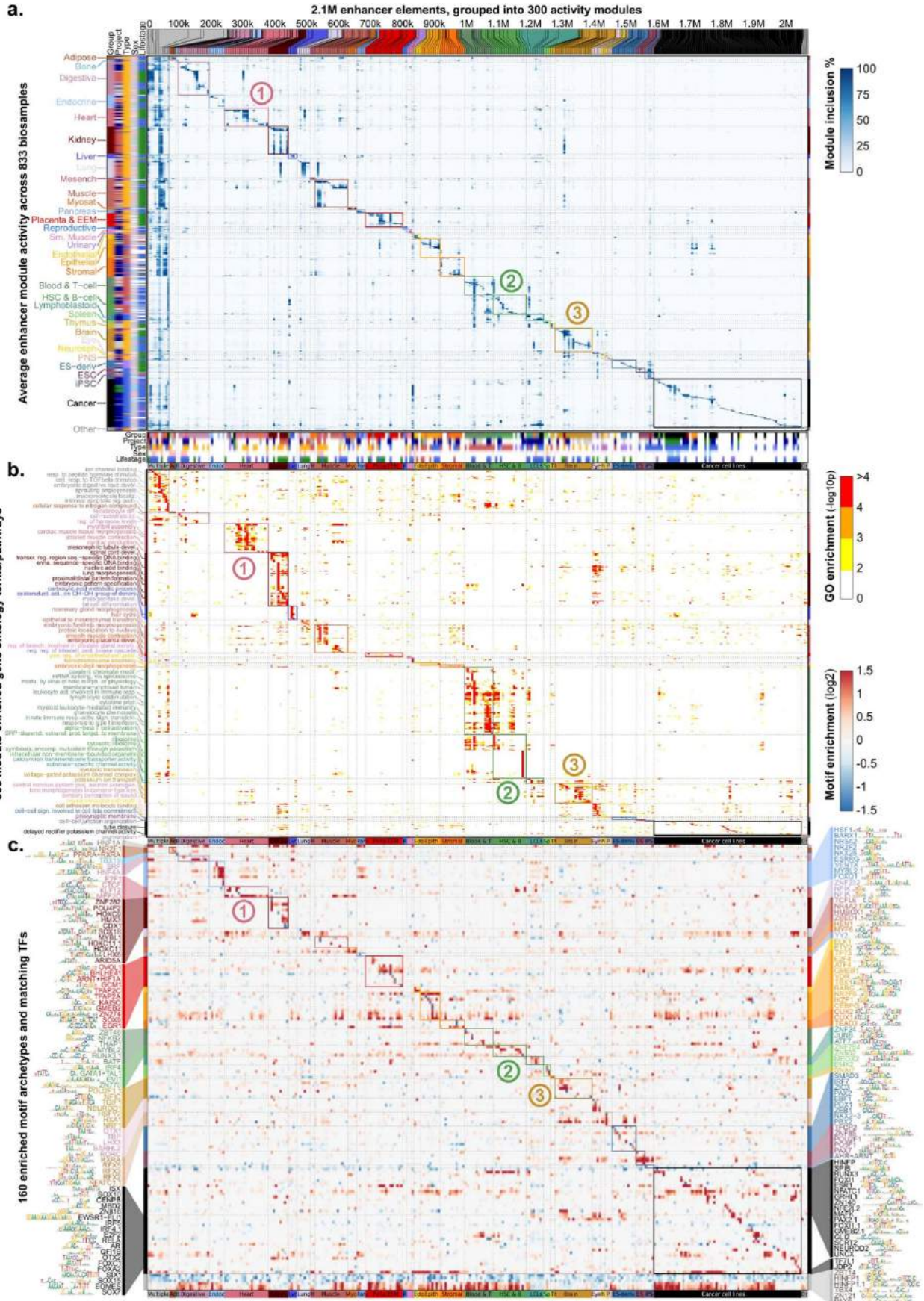
Motif enrichments often partitioned tissue categories into developmental-stage-specific and tissue-specific subgroups (**Fig. 1-7c, 2-9a-c**), including: heart into embryonic heart (*NFIX*, *E2F1*), aorta/arteries (*SRF*, *PAX5*), vs. heart chambers (*MEF2D*, *ESRRG*); brain into embryonic (*NFIX*, *NEUROD2*), adult brain (*RFX2*, *SOX10*), vs. astrocytes (*NFE2L2*, *JDP2*); and hematopoietic cells into NK cells (*ETV2*), B cells (*NFKB2*, *SPIB*), vs. multipotent-progenitors (*GATA1*, *NFE2L2*).





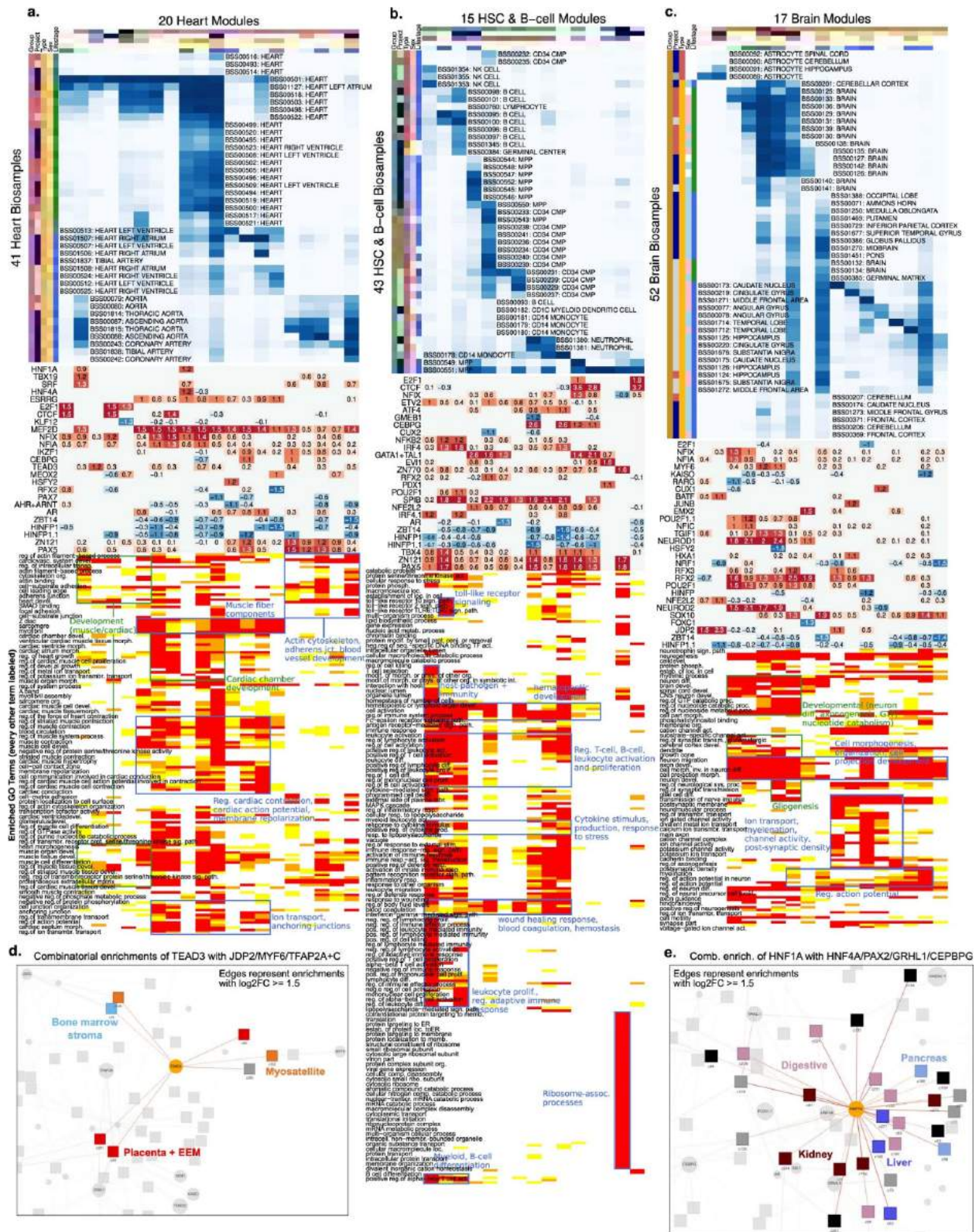
**Figure 1-7. Enhancer module circuitry.** **a.** Overview of gene regulatory module clustering. Full module breakdown in **Fig. 1-8** and online at <https://compbio.mit.edu/epimap>. **b.** Clustering of 2.1M enhancer elements (top) into 300 modules (columns) using enhancer activity levels (heatmap) across 833 samples (rows), quantified by H3K27ac levels within accessible enhancer chromatin states. Bottom panel shows enrichment of each module for each metadata annotation, highlighting 34 groups of modules (separated by dotted lines): 33 sample-type-specific (colored boxes) and 1 multiply-enriched (left-most). **c.** Subsets of enhancer module centers (top panels) and motifs (bottom panels) for heart, brain, and hematopoietic cell samples (top, rows), selected GO terms (middle, rows), and selected motifs (bottom, rows) in modules (columns) with maximal enrichment in each of the three sample categories. Full subsets in **Fig. 1-9**.





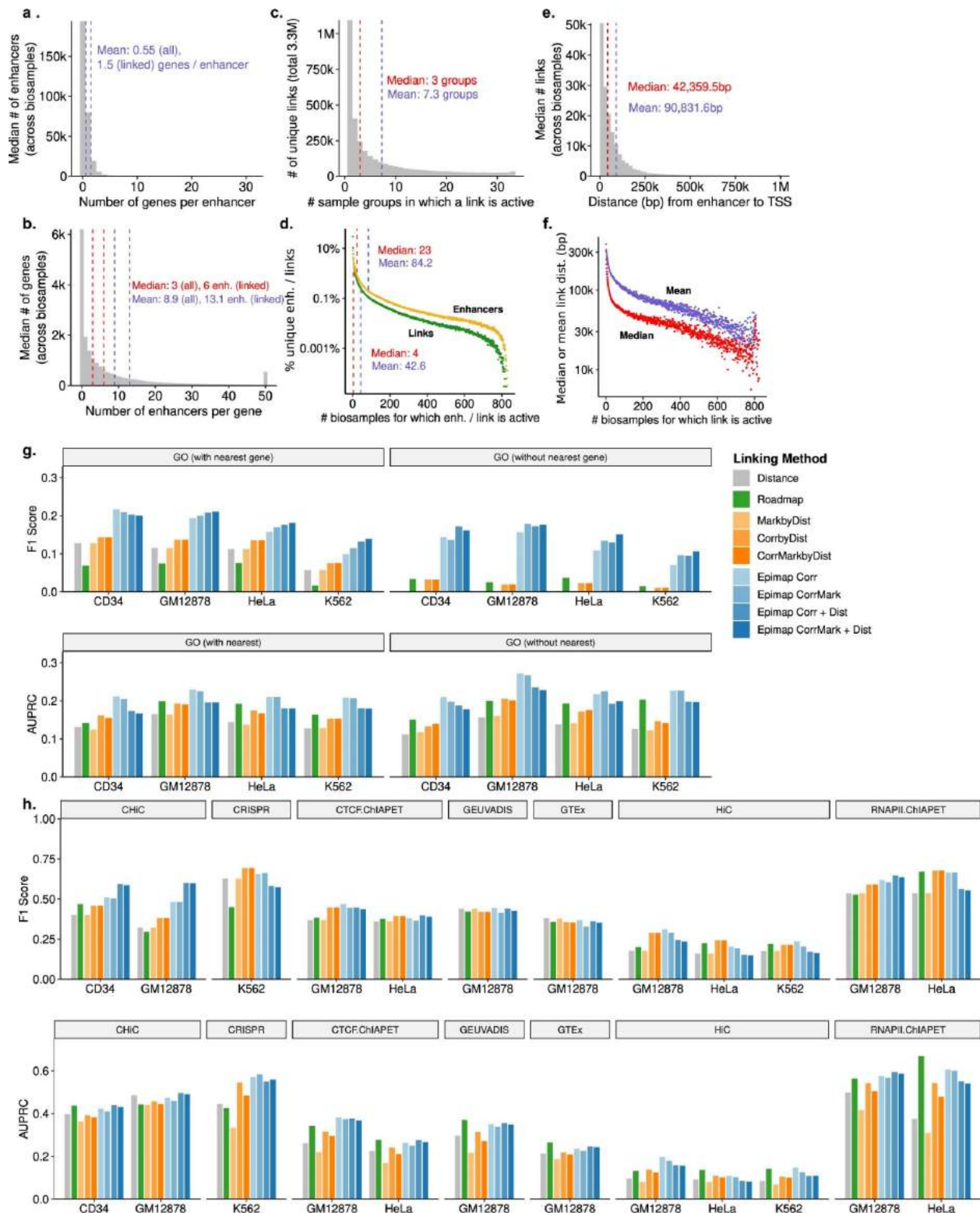
**1-8. Expanded enhancer module circuitry.** **a.** Clustering of 2.1M enhancer elements (top) into 300 modules (columns) using enhancer activity levels (heatmap) across 833 samples (rows), quantified by H3K27ac levels within accessible enhancer chromatin states. Bottom panel shows enrichment of each module for each metadata annotation, highlighting 34 groups of modules (separated by dotted lines): 33 sample-type-specific (colored boxes) and 1 multiply-enriched (left-most). **b.** Gene ontology<sup>110,111</sup> (GO) enrichments (heatmap) for each module (columns) across 865 terms (rows) with  $P < e^{-4}$ . GO terms colored by maximal enrichment group. Only 36 representative terms are shown, chosen by a bag-of-words approach within each tissue group. **c.** Motif enrichment (heatmap) for each module (columns) across 160 motif clusters (rows) with enrichment  $\log_2 FC > 1$ . Motifs colored by module of maximal enrichment.





**Figure 1-9. Expanded module enrichments and motif networks.** a-c. Enhancer module circuitry for a. heart, b. brain, and c. hematopoietic cells expanding Figure 1-7d subsets. From top to bottom, we show module centers for all samples in each group against all modules whose

maximal inclusion lies in the group, motifs with over two-fold enrichment (text:  $\log_2FC$ ), and GO enrichments for each module with tissue-category-specific enrichments. **d,e.** Snapshots of the motif-module network highlighting TEAD3 and HNF1A. Edges represent enrichments with  $\log_2FC \geq 1.5$ .



**Figure 1-10. Linking statistics and validation. a-f.** Enhancer-gene linking statistics. **a.** Barchart of median number of genes per enhancer per bin across biosamples. **b.** Barchart of median number of enhancers per gene per bin across biosamples. **c.** Total number of sample groups in which a unique link is active, out of 3.3M unique links. **d.** Percent of unique elements

(enhancers or links) for which the element is active in a given number of biosamples. **e.** Median distance between the enhancer and TSS of a link per distance bin across biosamples. **f.** Mean or median link distance for all enhancers active in a given number of biosamples. All means represented by blue dashed lines and text and medians by red dashed lines and text. **g-h.** Comparison (prediction F1 score and AUPRC) of gene-enhancer link predictions (blue) with distance, activity by distance, correlation by distance, correlation and activity by distance, and Roadmap Epigenomic links on **g.** functional gene set-based links and **h.** physical, genetic, and perturbation-based links across four cell lines.

### 1.5 Methodology for epigenomic datasets and processing

**Primary data sources and metadata information:** We analyzed 3,030 datasets, including 2329 epigenomic ChIP-seq datasets, 635 DNase-seq datasets, and 66 ATAC-seq datasets from ENCODE at <https://www.encodeproject.org/>, released as of Sept. 24th 2018. These marks include: Tier 1 assays: DNase-seq, H3K4me1, H3K4me3, H3K27ac, H3K36me3, H3K9me3, and H3K27me3; Tier 2 assays: ATAC-seq, H3K9ac, H3K4me2, H2AFZ, H3K79me2, and H4K20me1; Tier 3 assays: POLR2A, EP300, CTCF, SMC3, RAD21; and Tier 4 histone marks: 16 non-imputed histone acetylation marks, 4 methylation marks (H3K9me2, H3K79me1, H3K9me1, H3K23me2), H3F3A, and H3T11ph. We assigned unique sample IDs to each unique combination of: extended biosample summary, donor, sex, age, and lifestage, wherever each attribute was available. We removed samples with genetic perturbations, and kept only samples with appropriately matched ChIP-seq controls. We provide a metadata matrix including the mapping between ENCODE accessions and our unique sample IDs (**Supplementary Table S1**, also at [compbio.mit.edu/epimap](http://compbio.mit.edu/epimap)). We mapped the 111 Roadmap biosamples and 16 ENCODE 2012 biosamples to any of our biosamples with overlapping dataset accessions if the accessions were used in the flagship Roadmap epigenomics analysis. This mapping assigned 25 samples to ENCODE 2012 and 184 samples to Roadmap 2015, some of which were merged multi-donor samples in Roadmap, out of the final 833 samples that passed QC. These were merged into 16 and 111 tissue types respectively in the Roadmap 2015 publication.

**Uniform data processing:** We downloaded one alignment file per replicate, prioritizing filtered alignments aligned with BWA in hg19 whenever possible. We uniformly processed ChIP-seq and DNase-seq datasets according to the processing pipelines established by the Roadmap Epigenomics Consortium<sup>5</sup>. Briefly, we filtered out improperly paired and non-uniquely mapped reads, truncated reads to 36 base pairs, filtered out a blacklist of low complexity and artifact regions (ENCODE accession ENCSR636HFF), and filtered reads against a mappability track of uniquely mappable regions for 36bp reads<sup>112</sup>. Truncating read-lengths inevitably missed some repetitive regions that the longer reads could have helped resolve, but helped to avoid potential biases from alignment differences, as over two-thirds of the datasets had read-lengths of 36 or lower (**Supplementary Fig. 1-S9**). We converted bam files to tagAlign, used liftOver<sup>113</sup> to map GRCh38 alignments to hg19, and pooled all experiments within each ID and assay combination. We subsampled pooled ChIP-seq data sets to a maximum of 30 million reads and DNase-seq and ATAC-seq data sets to a maximum of 50 million reads. We used SPP to estimate fragment length. In cases with extremely low fragment length in ATAC-seq and DNase-seq datasets we used the average fragment length (73) from the average of the rest of the tracks. We generated  $-\log_{10}$  p-value signal tracks against matched whole cell extract (WCE) for both ChIP-seq and accessibility data sets using the MACS2<sup>114</sup> and the SPP<sup>115</sup> peak caller and cross correlation analysis to identify the proper fragment length as in the Roadmap analysis.



## 1.6 Epigenomic imputation methodology

**Imputation:** We carried out epigenomic imputation on 859 unique biosamples using ChromImpute<sup>94</sup> for a total of 10,778 imputed datasets over thirteen Tier 1 and Tier 2 assays using predictors trained on all 35 epigenomic assays across 859 samples. We additionally imputed 4,345 datasets for the five DNA associated factors, using only the 35 epigenomic assays as features to train predictors with ChromImpute. We provide all imputed and processed observed tracks along with tracksets for the 833 QCed samples at <https://epigenome.wustl.edu/epimap><sup>97</sup>.

**Quality control:** For imputation quality control (QC) and validation, we compare observed tracks to imputed tracks when both were available (i.e. when at least two original observed datasets were available for that biosample). We calculated all imputation QC metrics from the original ChromImpute publication<sup>94</sup>, including genome-wide correlation, % imputed and observed peak recovery, and AUC for all pairs of imputed and observed tracks. In addition to quantitative metrics, we visually inspected epigenomic predictions as part of our quality control. We show (**Fig. 1-2b**) three dense and varied regions of different resolutions (25kb, 200kb, 1.5Mb) for each of two randomly chosen samples containing both observed and imputed tracks for each assay. We calculated epigenomic profile quality metrics NSC, RSC, and read depth for all datasets and compared these to the imputation QC metrics (tables in **Supplementary Table S1**). We flagged low-quality tracks by detecting the elbow in the ranked correlation metrics, which we calculated as the point where the change in correlation exceeded 5% of the correlation. Validation on external datasets was carried out on 51 experimental tracks across 8 marks and assays from ENCODE after our data freeze, similarly subsampled to 30M (marks) and 50M reads (accessibility), which we re-mapped from hg38 to hg19 and evaluated on fully re-mapped 200bp bins (90.1%) in chr1 (**Fig. 1-2c-d**). For data homogeneity analysis, we restricted the data to only biosamples in each mark with both observed and imputed data (**Supplementary Fig. 1-S10**).

**Sample and antibody swap detection:** To systematically identify both potential sample or antibody (Ab) swaps and poor quality experiments, we computed the correlation of each observed experiment against all 10,734 imputed tracks for histone marks and assays (all imputed tracks before removing samples by QC). We then calculated the average correlation among the top 10 most similar tracks to each observed track. We flagged potential Ab swaps by comparing the average correlation against samples of the putative mark against those computed for other marks. We fit an OLS model to each mark comparison, flagged datasets with residuals greater than 3 standard deviations of the average correlation, and visually confirmed 7 Ab swaps (6 low-quality tracks). Similarly, we flagged potential sample swaps by comparing the correlation between imputed and observed tracks against the average correlation in the top 10 tracks in the same mark. We fit an OLS model and flagged datasets with residuals greater than 3 standard deviations of the residuals distribution. We report 19 potentially swapped samples, of which 5 were also flagged as low-quality tracks (**Supplementary Fig. 1-S8**).

**Secondary reactivities:** In addition to genome-wide QC of imputed tracks, we also focused on the specific differences between observed and imputed tracks. For each observed mark, we generated a genome-wide 'delta' track, computed as the difference in signal intensity between observed and imputed data, re-scaling imputed tracks to match signal intensity properties of the observed tracks, as observed tracks showed a general bias for higher intensity. Some of these 'delta' tracks showed surprisingly high correlations with 'primary' tracks of non-putative marks, indicating potential secondary antibody reactivities. In order to flag these reactivities, we compared the average correlation of each of the delta tracks to the top 10 closest imputed

tracks for each mark. As with antibody swaps, we fit an OLS model in each mark combination to flag outliers. We flagged 19 tracks and report 13 after visual inspection as potential secondary reactivities or single replicate swaps (e.g. in the case of DNase-seq) (**Supplementary Fig. 1-S7 and S8**). We noted that some cases showed clear difference tracks that don't match available antibodies, suggesting that the secondary reactivity is not a common mark in our compendium.

**Biological space coverage:** To evaluate the similarity of imputed and observed tracks across samples, we calculated the pairwise genomic correlations between all pairs of imputed and observed signal tracks. We hierarchically clustered each individual mark's imputed or observed correlation matrix using Ward's method. We averaged all imputed matrices for the six main marks (H3K27ac, H3K4me1, H3K4me3, H3K36me3, H3K27me3, H3K9me3) to create a fused correlation matrix, which we similarly clustered. We plotted the hierarchically clustered tree for the fused matrix alongside the metadata information for each biosample using the circlize R package<sup>116</sup>.

Additionally, we calculated mark-specific spearman correlations restricted to relevant features within all observed and imputed tracks per mark. We mapped each of 13 marks to its top state by emission probability in the ChromHMM 25 state model and any other states with emission probability over 80%. For ATAC-seq, we use the same region list as DNase-seq. For each mark, we averaged and reduced each 25bp signal track to any 200bp regions that were labeled as one of the states associated with the mark in any of the 127 imputed Roadmap biosamples under the 25-state model<sup>5,94</sup>. We calculated the spearman correlation between sets of these region-restricted mark signal tracks and generated similarity matrices across all datasets for a mark. Using these spearman correlation matrices on all observed and imputed signal tracks, we computed UMAP dimensionality reductions for each mark and assay using with the uwot R package<sup>117</sup> with the default parameters except for `n_neighbors=250`, `min_dist=0.25`, and `repulsion_strength=0.25`.

**Additional QC:** We first constructed genome-wide browser visualization tracks and manually inspected hundreds of genomic regions, and we found very strong agreement between imputed datasets and high-quality observed datasets (**Fig. 1b, Supplementary Table S1**). Across large regions, imputed tracks accurately captured the density, location, intensity, and cell-type-specificity of both active and repressed marks, both punctate and broad. At the highest resolution, they captured individual enhancer elements, precise boundaries in histone modification peaks, and finer-grained features of epigenomic assays, such as dips in H3K27ac peaks indicative of nucleosome displacement. We also generated summary visualizations that capture thousands of randomly-selected areas of the genome for each of 8 punctate epigenomic assays across all 859 samples, and once more found strong and consistent signal concordance, both across broadly-active and tissue-specific areas (**Fig. 1c**).

We also quantitatively evaluated the similarity between imputed tracks and observed tracks using genome-wide correlation and recovery of peaks (**Supplementary Fig. 1-S1,S2**). For the 1212 observed datasets on 8 punctate marks (59% of tracks), we found 75% average genome-wide correlation and 85% peak recovery, indicating high-quality reconstructions, even though imputed datasets were generated without any knowledge of observed data for the imputed mark in the target cell type. For the 839 observed datasets on 5 broad marks (41% of tracks), we found 53% genome-wide correlation and 62% peak recovery, as expected for broad marks that do not have clearly-defined peaks, and whose intensity may vary within broad marked domains without specific functional implications.



## 1.7 Imputation as a tool for QC

We next used the agreement between imputed and observed tracks as a quality control (QC) metric to detect and report potentially problematic datasets. While most datasets showed comparable levels of agreement, 138 datasets (6.7% of 2,051) showed sharply lower scores, indicative of possible experimental errors or lower quality data. Indeed, these experiments also showed markedly lower ENCODE QC scores, including consistently-lower read counts and worse strand cross-correlation metrics (**Supplementary Fig. 1-S2,S3**).

For 11 of the flagged datasets (7.9% of 138), the low imputed-observed agreement stemmed from sample swaps or antibody swaps, a type of problem that is not captured by traditional QC metrics. We flagged swaps by comparing the observed tracks against the average correlation in the top 10 most similar imputed tracks per mark, and included an additional 15 samples not in the original 138 as potential swaps (see Methods, **Supplementary Fig. 1-S3,S4**). Several of the datasets were also flagged by the data producers as swapped, and the remaining ones have now been shared with the data producers and the broader community as part of **Supplementary Table S1**.

We also used the difference between observed and imputed datasets to recognize 15 experiments with potential antibody cross-reactivity or secondary specificity. For example, after subtracting the imputed H3K36me3 signal from the observed H3K36me3 signal in a T-helper 17 cell (id: BSS01695), we find that the residual signal track shows 55% correlation with H3K4me1, suggesting antibody cross-reactivity, secondary reactivity due to other reasons, or antibody swaps for one experiment of multiple (**Supplementary Fig. 1-S6-8**). The most common secondary reactivities observed involved dimethyl groups as either the primary or secondary reactivity (in 10 datasets), in particular from H3K4me1 to H3K4me2 for 2 datasets.

Lastly, we used QC of observed and imputed tracks to remove problematic datasets from our subsequent analyses. First, we removed the 138 flagged observed datasets, and only used imputed data for these samples. Second, we removed all imputed data for 22 samples (374 tracks) that were based solely on ATAC-seq data, as imputed marks in these samples showed overall lower correlation with other samples, possibly due to sparse overlap between ATAC-seq and other marks. Third, we removed all imputed and observed data for 4 samples that were based solely on low quality DNase-seq (68 tracks). This resulted in a total of 2,850 observed and 14,510 imputed marks across 833 samples used in the remainder of this work.

## 1.8 Epigenomic annotations methodology

**Chromatin state annotations:** We computed epigenomic annotations on 3,533 imputed and 1,465 observed datasets for six marks on 833 samples using ChromHMM with the fixed 18-state model from Roadmap<sup>5</sup> with the same mnemonics and colors. We use observed data wherever possible, except in cases with no observed data or where observed data was removed in QC. The table of signal tracks used to calculate annotations is available as **Supplementary Table S2**. Observed data was binarized from signal tracks with a  $-\log_{10}$  p-value signal cutoff of 2. In order to binarize imputed data and facilitate the comparison with observed data, we established mark-specific binarization cutoffs. We first separately calculated the overall probability distributions of all imputed and observed tracks for each mark. Then for each mark, we set the imputed binarization cutoff value to the value of the quantile matching the quantile in observed data for the  $-\log_{10}$  p-value > 2 cutoff. We used liftOver<sup>13</sup> to map all 833 (after QC) ChromHMM annotations to GRCh38 and provide these alongside hg19 annotations and binarized imputed and observed datasets and as tracksets at <https://epigenome.wustl.edu/epimap>.

**Defining active enhancers:** We define active enhancers as the intersection of DHS (DNase I Hypersensitive Sites) regions with enhancer annotations and high H3K27ac signal (average

signal  $> 2$  in the region containing the DHS  $\pm 100$ bp). We define DHS regions from an index list of 3,568,912 DHS consensus locations determined from 733 DNase-seq experiments. We intersect these regions with the 833 imputed enhancer annotations (states 7,8,9,10,11, and 15 in the 18-state model). This results in 2,842,995 regions with at least one enhancer annotation in any biosample. Finally, we intersect this matrix with the H3K27ac signal in the  $\pm 100$ bp region encompassing each DHS from the same tissue-specific imputed and observed datasets used to calculate the ChromHMM annotations. This procedure results in 2,356,914 active enhancer regions. We created an equivalent promoter element region using the promoter annotations (states 1,2,3,4, and 14 in the 18-state model). We noticed that a number of regions share both enhancer and promoter annotations. As a conservative cutoff, we assign all regions to either enhancers or promoters if 75% or more of its active occurrences are labeled as that type of element (**Supplementary Fig. 1-S13**). This final thresholding procedure yields 2,069,090 enhancers, 204,104 promoters, and 122,358 dyadic elements (neither specifically promoter or enhancer). Matrices and enhancer locations are available through <http://compbio.mit.edu/epimap>.

For all images using tissue group order, including ChromHMM tracks and modules heatmaps, groups are ordered alphabetically within six major groups: Tissue/Organs (Adipose, Bone, Digestive, Endocrine, Heart, Kidney, Liver, Lung, Mesench, Muscle, Myosat, Pancreas, Placenta & EEM, Reproductive, Sm. Muscle, and Urinary), Other Primary Cells (Endothelial, Epithelial, and Stromal), Blood + Immune (Blood & T-cell, HSC & B-cell, Lymphoblastoid, Spleen, and Thymus), Nervous system (Brain, Eye, Neurosph, and PNS), Stem (ES-deriv, ESC, iPSC), and Other (Cancer, Other).

**Defining enhancer modules:** In order to define enhancer modules, we clustered the binary enhancer matrix defined by intersecting enhancer annotations with DHS regions and with average centered and flanking ( $\pm 100$  bp) H3K27ac signal above a  $-\log_{10}$  pval of 2 using the k-centroids algorithm with the Jaccard distance with the number of clusters set to  $k=300$ . The average module contains 6,897 enhancers, and the largest module (enumerating constitutive elements) contains 93,554 enhancer regions. In all heatmap plots of module centers (and associated enrichment figures), we diagonalize the matrix by ordering each column in the heatmap (module centers) by the biosample which contributes the maximal signal. All columns which have signal over 25% in more than 50% of rows are shown first. We use this diagonalization procedure for all diagonalized heatmaps. We colored each module by the tissue group which contains its maximal signal. Modules highlight sample groupings and organize according to cell type and tissue. Major groups are ordered alphabetically within six major groups and samples are ordered within groups according to Ward method's clustering of the Jaccard distance of the modules centers matrix. We performed enrichment on the module centers against the metadata of included samples (signal over 25%) by the hypergeometric test, and show enrichments with  $-\log_{10}p > 2$  (**Figure 2b**).

**Gene ontology enrichment:** We performed gene ontology (GO) enrichments on each enhancer module using GREAT v3.0.0 for the Biological Process, Cellular Component, and Molecular Function ontologies<sup>111</sup>. We analyzed and visualized the results in the same manner as in the Roadmap core paper<sup>5</sup>. We only considered enrichments of 2 or greater with a multiple testing-corrected p-value less than 0.01. For Figure 4C, we reduced the enrichments by modules matrix to terms with a maximal  $-\log_{10}p$ -value  $> 4$  that were enriched in less than 10% of modules. The full enrichment matrix is shown as **Supplementary Figure 1-S16**. As in the case of the diagonalized module centers, we labeled each term according to the module containing its maximal signal. We used a bag of words approach (as described in Roadmap<sup>5</sup>) to pick 72 representative terms out of 865 total terms for **Fig. 1-8b** such that each tissue group has at least

one term and the rest are representatively allocated across groups.

**Motif enrichment:** We performed motif enrichment analysis across enhancer modules as described in the Roadmap paper<sup>5,118</sup>. Briefly, we measured enrichment of 1,690 motifs consisting of the JASPAR 2018<sup>119</sup> core non-redundant vertebrate motifs, the HOCOMOCO v11<sup>120</sup> human motif set, and the SELEX motifs by Jolma et al (2013)<sup>121</sup>. We computed enrichments for each of the 1,690 motifs relative to a joint DHS and intergenic background, additionally controlled by 100 shuffled motifs for each motif. We report the motif with the highest enrichment in any module for each of 286 previously identified motif archetypes<sup>102</sup>. We only report motifs with a max log<sub>2</sub> FC of at least 1, resulting in 160 motif archetypes (corresponding to 1,175 total motifs), which we show with their PWM logos against all 300 modules as **Fig. 1-8c**.

**Enhancer-gene linking:** We predicted enhancer-gene links for each biosample using the Pearson correlation between gene expression and nearby enhancer's (within 1Mb) histone mark activity for six marks (H3K27ac, H3K4me1, H3K4me2, H3K4me3, and H3K9ac). We pre-computed correlations between all genes and nearby enhancers across the 304 biosamples with paired expression data. A negative set of correlations for each enhancer was computed using random genes in a different chromosome. We predicted links for each biosample and ChromHMM enhancer state separately (states E7, E8, E9, E10, E11, and E15). Predictions were made by training an XGBoost classifier on the positive set of all valid links against their paired negative links, using pre-computed correlations and distance to TSS as features, and keeping all links with probability above 5/7<sup>122</sup>.

We validated enhancer-gene links using curated gold-standard data<sup>123</sup> in CD34, GM12878, HeLa, and K562 cells (**Fig. 1-10**). We compared four sets of correlation-based predictions (alone or with H3K27ac and H3K4me1 activity, and with and without distance-based re-scaling) against distance alone, enhancer-gene links from Roadmap, and H3K27ac correlation and/or activity times distance (calculated using EpiMap tracks and enhancers in compared epigenomes)<sup>124</sup>. For methods without a threshold value, such as distance alone, only the nearest/highest score gene is utilized for each as a cutoff value for F1. In addition, we created a GO-based gold-standard set of links from GO terms that were enriched within enhancer clusters by GREAT<sup>111</sup>. For each GO term per cluster, we added enhancer-gene links for enhancers within 1Mb of at least two genes in the GO term. Negative link sets were constructed by taking physical and eQTL negative link sets that were also non-GO enriched.

**Data availability:** We provide all imputed and processed observed tracks along with ChromHMM annotations and tracksets for the 859 imputed and final 833 QCed samples at <https://epigenome.wustl.edu/epimap><sup>97</sup>. All other processed and intermediate datasets, including metadata (**Supplementary Tables S1 and S2**), flagged samples, annotations, DHS locations, enhancer and promoter definitions, enhancer and promoter matrix, modules, and matched RNA-seq data can be found at <http://compbio.mit.edu/epimap>. We also provide an interactive data and analysis browser through the website including biosample and track exploration, creation of custom track hubs, modules and motifs enrichments, and per-GWAS dissections for each of the GWAS and their lead SNPs (**Supplementary Fig. 1-27**)<sup>125</sup>.

ChromImpute can be found at <http://www.biolchem.ucla.edu/labs/ernst/ChromImpute/>. Analysis was performed with R (3.5 and 3.6) and Python 3.7. Analysis code is available from <http://compbio.mit.edu/epimap/>.

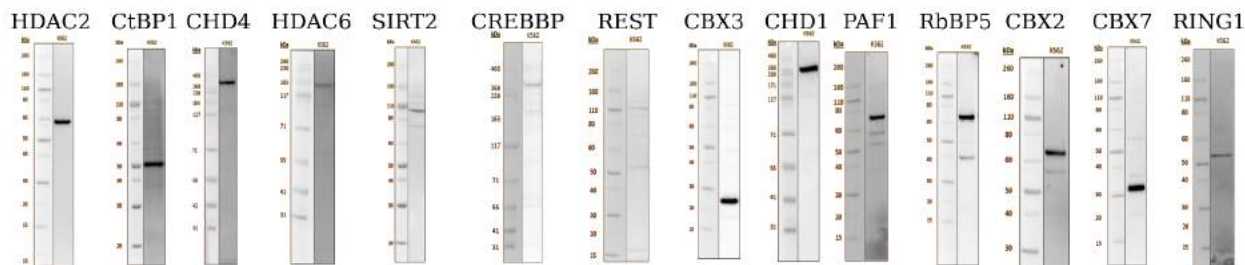
## 1.9 Learning latent representations of chromatin regulators

**Summary:** Cellular differentiation depends on selective activation and inactivation of gene transcription. Active and inactive gene states are associated with distinct patterns of chromatin modification at key residues in the histone amino-terminus, including acetylation and methylation, among others. ChIP-seq has enabled the genome-wide characterization of chromatin states associated with these modifications. We have employed ChIP-seq to map the genome wide localization of the enzymatic machinery that accomplishes histone modification, so called chromatin regulators. We have used, in five cell types, a panel of ~50 antibodies to chromatin regulatory enzymes, their cofactors, and components of chromatin remodeling complexes. Combinatorial patterns of histone modification and factor localization suggest new insights into both the action of these regulators and the components of regulatory complexes.

We develop a method to learn sparse latent components of genomic profiles across the genome using a variational auto-encoder framework. Using this method, we identify genomic context-specific combinatorial interactions between chromatin regulators and between regulators and histone marks and transcription factors reported by ENCODE.

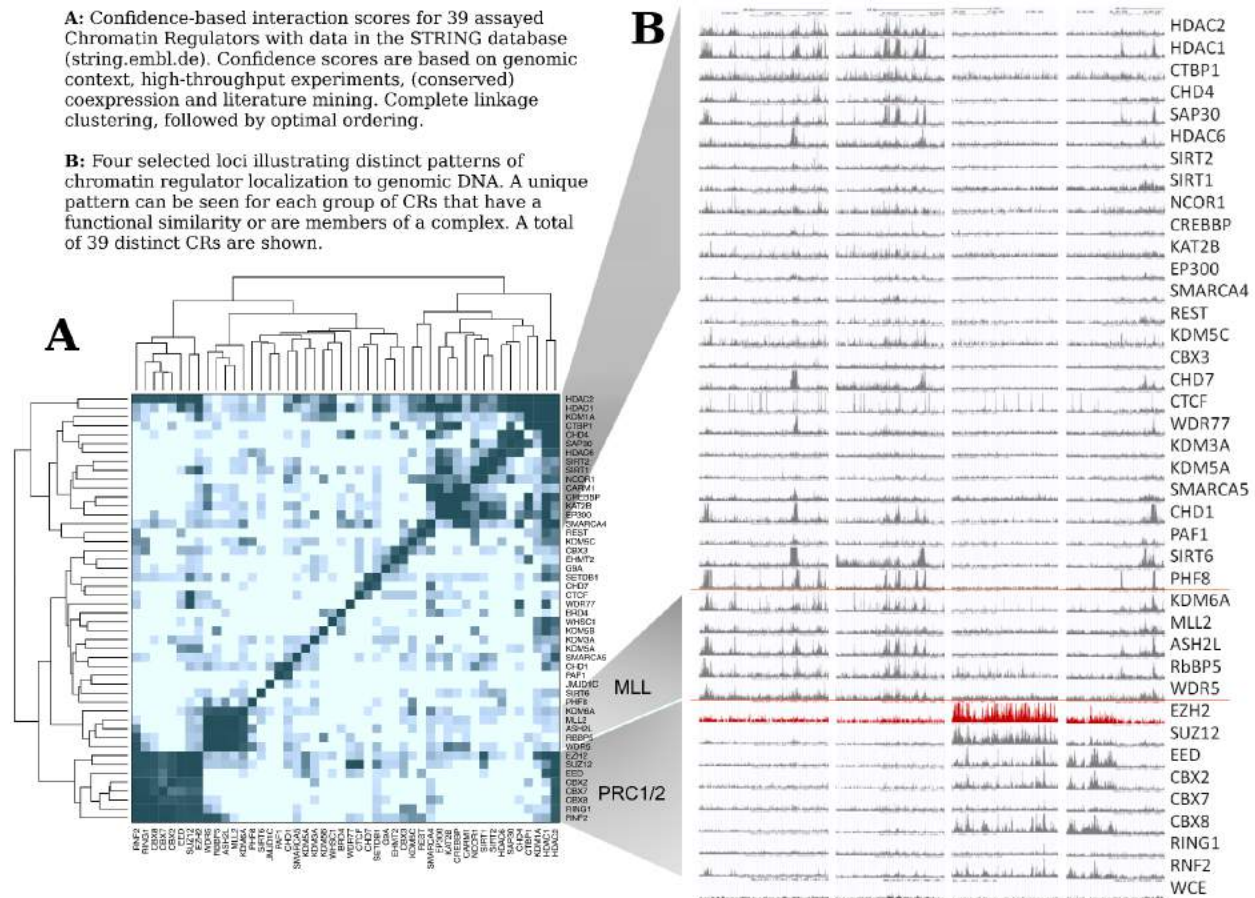
We report both high-resolution chromatin state models that incorporate measurements of both chromatin modifications and the occupancy by reported genomic components and their hierarchical interactions and a dissection of the modular, context-dependent activities of chromatin regulators. These complexes are frequently mutated in cancer, and our study should provide a resource for better understanding how perturbations of chromatin regulators can contribute to the dysregulation of cellular differentiation that accompanies neoplasia.

**Data and Data Processing:** ChIP-seq was performed in a panel of five cell types (A549, HepG2, H1-hESC, K562, and GM12878) for 50 chromatin regulators. Antibodies were screened and validated and ChIP-seq is performed in the cell lines of interest (**Fig. 1-11**). Validation was done by western blot and comparison of ChIP-seq tracks to known co-factors.



**Figure 1-11. Selected westerns from antibody validation prior to chromatin regulator ChIP-seq.**

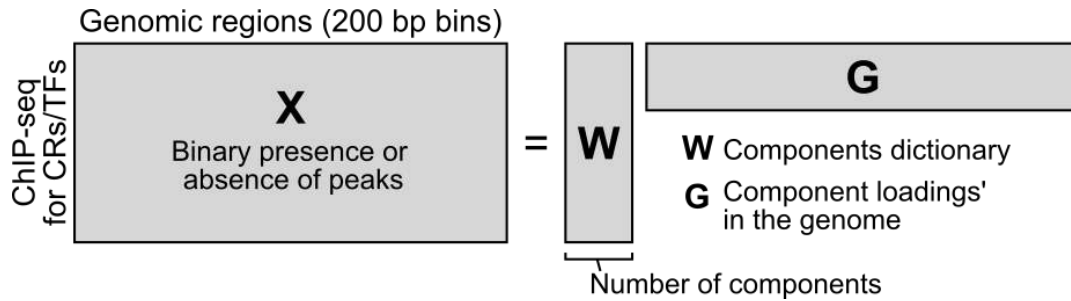
Reads were mapped to hg19 using BWA, trimmed, and filtered for duplicates. Unmappable sites were removed and the quality of mapping was evaluated for each replicate. We performed replicate concordance analysis, followed by the removal of samples with low concordance. Replicates were pooled, subsampled to uniform depth (30M reads), and strand cross-correlation analysis was performed to evaluate ChIP-seq peak size and data quality. Pooled data was both binarized by peak calling and turned into a continuous track by comparison to whole cell extract control using MACS. Binary data consisted of peak calls segmented into 200bp bins across the genome (which aligned with ChromHMM bins), while continuous data (fold change over control) was segmented into 25bp bins. In addition to chromatin regulators, ENCODE transcription factor and histone mark ChIP-seq experiments were processed analogously to facilitate the inter-experiment comparisons, for a total of 630 datasets across 5 cell types (**Fig. 1-12**).



**Figure 1-12. Comparison of PPI interactions with ChIP-seq tracks. A.** Interaction scores for 39 assayed chromatin regulators with data in the STRING database (string.embl.de)<sup>126</sup>. Confidence scores are based on genomic context, high-throughput experiments, conserved coexpression, and literature mining. Matrix is ordered by complete linkage clustering followed by optimal ordering. MLL, PRC1/2, and deacetylase complexes are highlighted. **B.** Four selected loci illustrating distinct patterns of chromatin regulator localization to genomic DNA. A unique pattern can be seen for each group of CRs that have a functional similarity or are members of a complex. A total of 39 distinct CRs are shown.

**Non-negative matrix factorization model:** We developed a non-negative matrix factorization (NMF) approach in order to learn components of chromatin regulation from ChIP-seq tracks of CRs and TFs. The factorization approach enabled us to learn components (or factors) which could combine in context-specific and hierarchical ways across the epigenome. The non-negative constraints yielded interpretable components which were directly understandable from the component itself.





**Figure 1-13. Illustration of NMF model approach.** We factorize a ChIP-seq data matrix ( $X$ ) of binarized peak calls for a set of profiled chromatin regulators (CRs) and transcription factors (TFs) across the genome in 200 bp bins. We learn a sparse dictionary ( $W$ ) in which each column (component) selects a constrained set of chromatin regulators and transcription factors. We also learn where these components appear in the genome with a jointly learned sparse component loading matrix ( $G$ ).

More concretely, given a matrix  $X$  (experiments by genome bins), composed of binary peak calls (in  $\{0, 1\}$ ) across the genome in 200 bp bins, we developed a method to learn a dictionary of components ( $W$ ) and their loadings across the genome ( $G$ ) by non-negative factorization (Fig. 1-13). We minimized the binary cross-entropy loss on the reconstruction  $z$  with sparsity and orthogonality inducing penalties, as expressed in the following loss function:

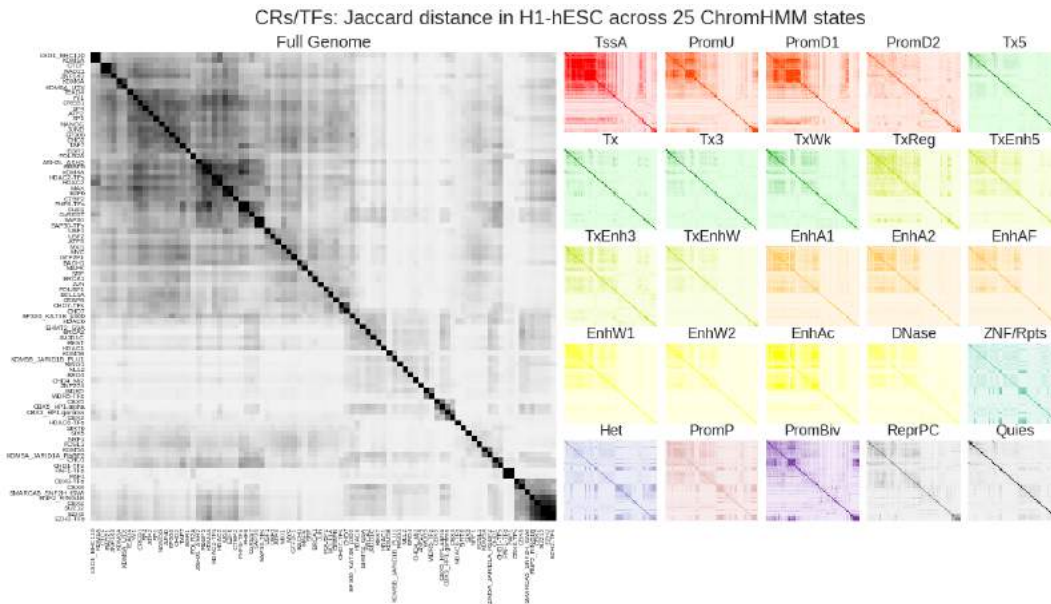
$$\arg \min_{W, G, \beta > 0} |X - \sigma(WG)| + \lambda_W \|W\|_1 + \lambda_G \sum_i \frac{\|G_i\|}{|X_i|} + \lambda_{ortho} \|GG^T - I_K\|_1 + |H - \sigma(\beta G)| + \lambda_\beta \|\beta\|_1$$

In particular, we penalized loading multiple complexes proportionally to the number of factors at a site and promoted orthogonality of loadings  $G$  and uniqueness of factors  $W$ . Additionally, we jointly trained to predict histone marks ( $H$ ) using the component loadings. Inference was performed using stochastic gradient variational bayes (SGVB) in an variational autoencoder (VAE) framework<sup>127,128</sup>, which can perform efficient mini-batch training, easily training the full epigenome, and accommodate missing data, different data modalities, and mixed objectives (e.g. reconstruct of CRs/TFs and jointly predict prediction of histone marks). Model parameter selection was to be done by performing hyperparameter optimization to minimize BCE reconstruction error while maximizing the consistency of the learned dictionary, as measured by concordance of components learned by 1000 bootstraps of the algorithm with different initial conditions. Here I present an overview of the results, insights, and core directions of the project in the hopes that it can be useful for future research.

**Inferring context-specific components of chromatin regulator binding:** Chromatin regulators and remodelers (CRs) are DNA factors mainly tasked with modifying histones and remodeling nucleosomes in order to reshape or maintain the epigenomic landscape of cells undergoing differentiation, activated by specific signaling pathways, or just maintaining their cellular identity<sup>129</sup>. Due to their varied roles in controlling gene regulation through, for example, epigenomic activation<sup>130-133</sup> or repression<sup>134-138</sup>, these factors are highly implicated in cancers and in developmental disorders<sup>90,139-141</sup>. Their roles in regulating epigenomic features make understanding where and how they affect the epigenome crucial to understanding and treating a wide range of disorders.

Importantly, CRs rarely have a DNA binding affinity, unlike transcription factors (TFs), making understanding how these are localized to the genome even more complex than in the case of

TFs. Much fruitful research on these has focused on identifying protein complexes containing these CRs in order to pinpoint units with specific epigenome localization and function<sup>136,138,142,143</sup>. However, most studies focus on one set of factors and cannot systematically evaluate interactions between many factors jointly. More importantly, the protein, chromatin, and histone interactions that mediate localization of CRs can often be context dependent, with regulators switching partners depending on TF availability, cell type, or epigenomic context (**Fig. 1-14**).



**Figure 1-14. Context-specific colocalization patterns by epigenomic state.** Jaccard distance for chromatin regulators/remodelers (CRs) and transcription factors (TFs) binarized ChIP-seq peaks in H1-hESC cells across the full genome (left) or within each of 25 ChromHMM states (right).

Our systematic panel of over 50 CRs in five cell types, jointly analyzed with ENCODE TFs and histone modification ChIP-seq data (for 630 datasets total) offers a view into these complex and context-dependent associations of chromatin regulators (**Fig. 1-12**)<sup>129</sup>. In addition, we developed a non-negative matrix factorization model (see Methods, above) to learn context-specific combinatorial patterns of histone modification and factor localization over the epigenome that suggested new insights into both the action of these regulators and the components of regulatory complexes.

We learned factorizations and dictionaries for each cell type (A549, GM12878, HepG2, H1-hESC, K562) separately (**Fig. 1-15a-e**). We identified on the order of 50-100 components which included very few singletons or redundant components due to our regularization and orthogonality penalty. These components separated well established chromatin regulator complexes in similar locations, such as the CoREST complex (CoREST, HDAC1/2, KDM1A/LSD1)<sup>143</sup> and the heterochromatin-establishing complex (CBX3/HP1-gamma and CBX5/HP1-alpha) (**Fig. 1-15b-c**)<sup>144</sup>. Our factorizations also separated the PRC2 complex (EZH2, EED, and SUZ12)<sup>138</sup> from PRC1 (CBX2, CBX8, and RNF2/RING1B)<sup>136</sup>, both of which are primarily localized to and maintain repressed (H3K27me3-marked) regions in the epigenome (**Fig. 1-15d**).

We identified combinatorial associations forming higher order components within the learned dictionaries, including a joint CoREST and HP1 component with LSD1/ZNF143 and CBX3/5 that was localized to a mixture of heterochromatin and active regions (**Fig. 1-15c**). The same

combination, with the addition of HDAC2, was localized almost exclusively to heterochromatin sites, possibly reflecting specific turnover, deactivation, and silencing of previously accessible and active regions, a key step in cell fate specification (**Fig. 1-15c**). The CoREST factors showed particularly high combinatorial patterning, also associating with reprogramming factors POU5F1 and NANOG as well as with YY1, MAX, and CREB1 (**Fig. 1-15b**). Both of these components were identified at promoter sites and possibly mark different patterns of promoter activation and gene regulation. The most striking example of this higher order organization was shown by EP300, a gene-activating acetyltransferase, which appeared in particularly varied contexts, both in enhancer-related components with different sets of TFs and in promoters with the ASH2L/RBBP5 regulator complex (**Fig. 1-15e**)<sup>130,133</sup>.

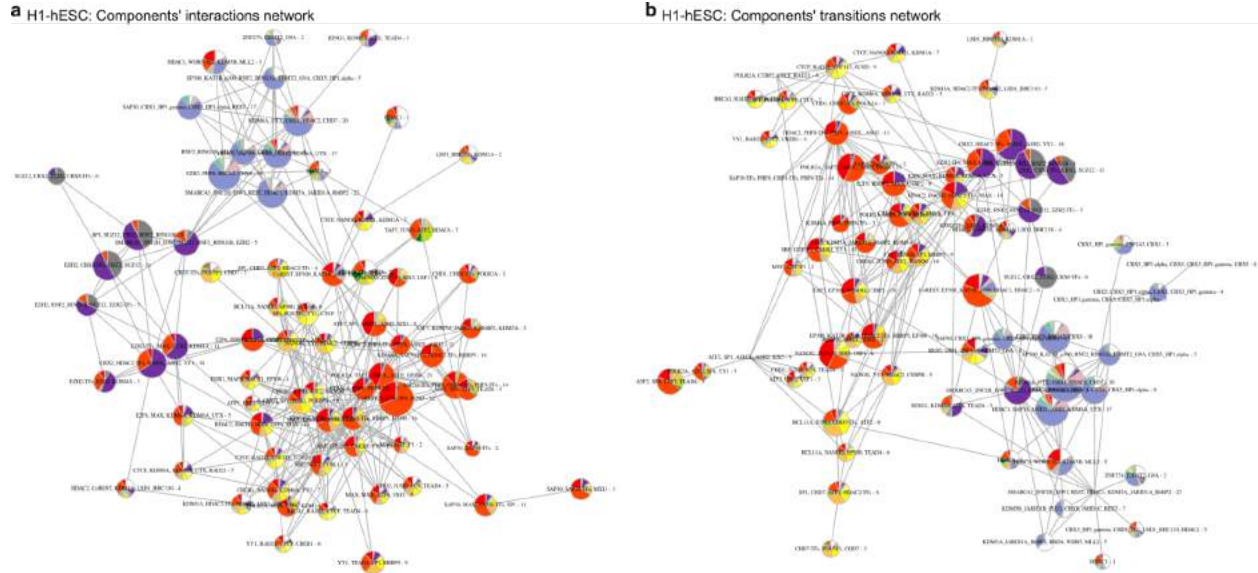
In the case of H1-hESC cells, which due to their pluripotent state contain a strikingly large number of bivalent promoters with both repressive and activating histone marks, we learned additional, context-specific components for these complexes that also contained promoter or enhancer-specific factors. For example, we identified PRC1 factors in two components with promoter-associated factors such as MAX, CREB1, and YY1, suggesting a colocalization of the two sets of factors in embryonic cell lines, despite their disparate functions (**Fig. 1-15d**). Strikingly, we identified a component that was enriched in active and bivalent promoters and included repressive EZH2, which methylates H3K27me3, side by side with KDM6A and KDM4A, which demethylate H3K27me3 and silencing mark H3K9me3, respectively<sup>145,146</sup>. This component also included ASH2L/ASH2 and RBBP5, which are involved in methylation of H3K4 and epigenetic activation of enhancers and promoters. This unexpected finding highlighted the interplay between maintenance of repressive and activating marks and demonstrates that the bivalency of embryonic promoters extends beyond just the specific histone marks at these sites and involves multiple chromatin regulating factors that work at cross purposes but are jointly positioned at these locations (**Fig. 1-15d**).





**Figure 1-15. Dictionary learning uncovers partners and shows state-specific modes of binding for regulators.** **a**, Dictionary of chromatin regulation components learned in H1-hESC on chromatin regulators and associated transcription factors. Chromatin state percentages (**a**) and raw values (**b-e**) for components and regulators as annotated by chromHMM are shown on the bottom and right margins, respectively. Chromatin regulators are annotated in red. **b-e**, Components associated with specific factors within the dictionary. **b**, LSD1 (KDM1A/BHC110) an H3K4 and H3K9 demethylase, works in concert with HDAC1/2 and UTX (KDM6A) but is also recruited to heterochromatic regions with HP1 (CBX1/3/5). **c**, CBX3 occurs with canonical HP1 components CBX5 (and CBX1 in other cell types), but also with ZNF143. **d**, EZH2, involved in PRC2 and methylation of H3K27me3, is found in putatively bivalent promoters (red) and active promoters and repressive regions. **e**, EP300, a gene-activating acetyltransferase, appears in enhancer-related components with different sets of TFs and in promoters with the ASH2L/RBBP5 regulator complex.

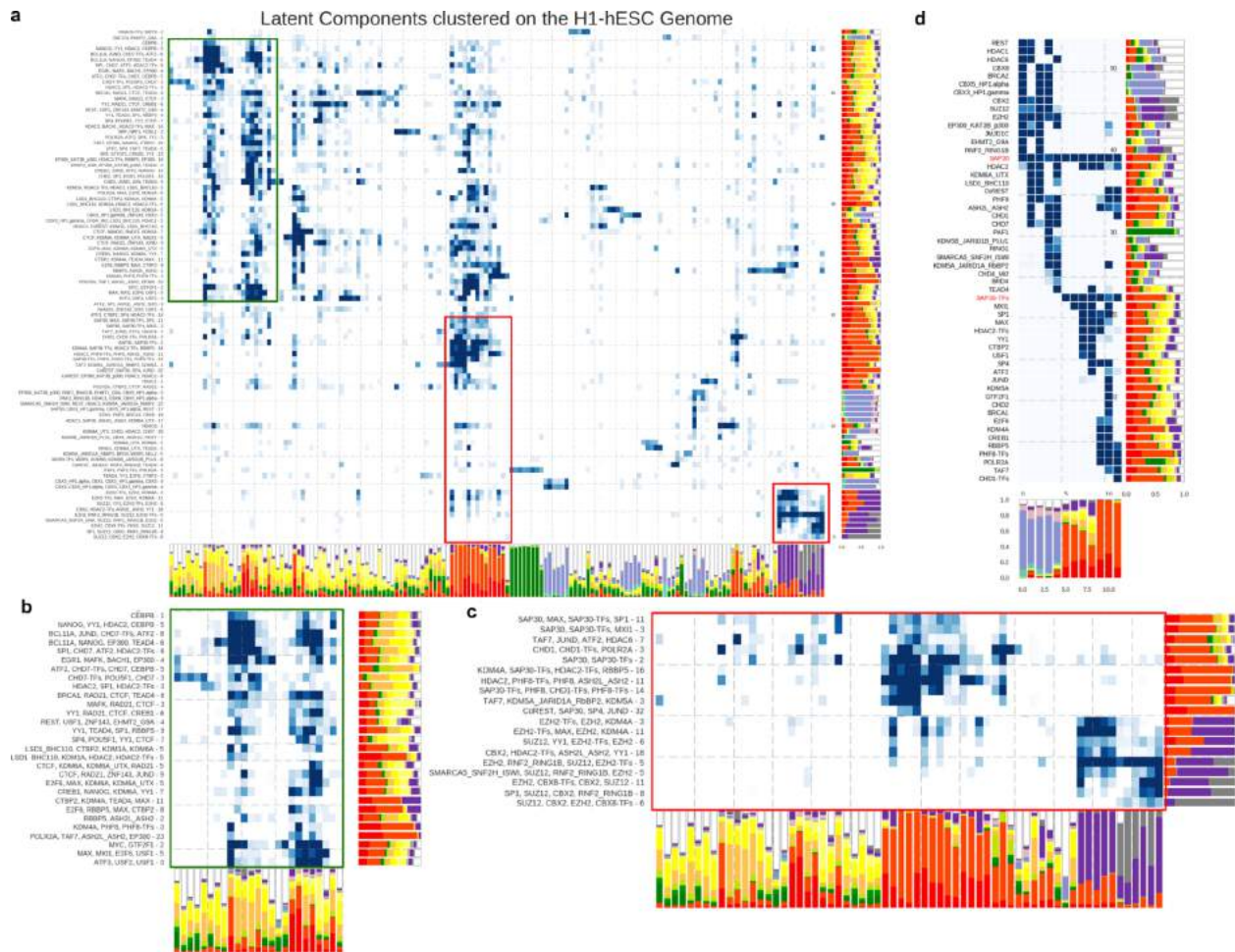
To better understand these component interactions and the order in which they occur in the genome, we built a pair of networks that represent how often components are in the same genomic bin (transitions network, **Fig. 1-16a**) and how often they fall into sequential genomic bins (transitions network, **Fig. 1-16b**). Components that shared similar states and the same factors were almost always connected in the network, and networks revealed similar higher order interactions to the ones captured by the dictionary (e.g. EZH2 + MAX in bivalent promoter states). In the interactions network, SAP30 was involved in complexes across the network, whereas repressive and heterochromatin factors were restricted to specific states and parts of the network, suggesting that SAP30 and its associated complex either act orthogonally to chromatin state or are found in transitions between chromatin states, therefore explaining their promiscuous appearance (**Fig. 1-16a**)<sup>147</sup>. In the transition network, we identified multiple CTCF/RAD21-including components connecting promoter states and enhancer states<sup>148</sup>. Hub nodes in the network marked core transition states, including multiple states at the boundary between promoters to bivalent promoters and between bivalent promoters and repressed states, which may represent a transition between bivalent or active promoters and broad repression around a promoter (**Fig. 1-16b**). A CoREST component including histone acetylase EP300 and deacetylases HDAC1/2 was also highly connected to repressed, silenced, and active areas of the genome, and together with its unorthodox composition could represent a transition state between acetylated and deacetylated regions of the genome present in order to maintain sharp boundaries on activating marks (**Fig. 1-16b**).



**Figure 1-16. H1-hESC component networks.** **a**, Network of component interactions. Interactions are calculated by the Jaccard similarity of two components, e.g. number of shared active locations over total active locations between the two. Each node is a component, represented by a pie chart of ChromHMM states showing the saliency of regions assigned to each of the components. **b**, Network of component transitions, where transitions are counted as the number of times that one component follows another in a subsequent 200bp bin along the genome. Nodes represent components, as in (a).

To further uncover higher-order interactions between components and understand how components load into the epigenome as a whole, we clustered the learned loadings of the learned components (G matrix) over the full H1-hESC epigenome matrix using k-centroids clustering (**Fig. 1-17a**). Overall, the epigenome clustered primarily by its chromatin state, as annotated independently by ChromHMM, with the majority of clusters showing highly specific state enrichment for promoter, enhancer, transcribed, heterochromatic, repressive, and bivalent promoter states. Components also predicted histone marks as well as or better than the compendium of CRs and TFs alone (**Fig. 1-18a-c**), suggesting that little information content is lost relative to the original data when viewed through the lens of histone marks and their epigenetic contexts.





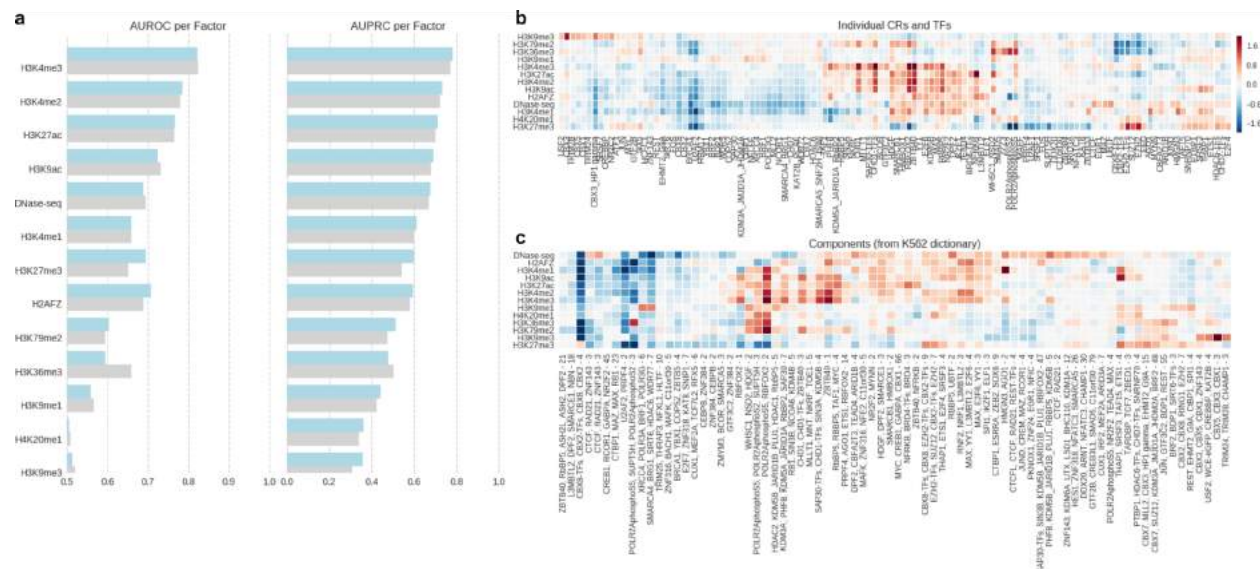
**Figure 1-16. Dictionary components clustered over the genome reveal higher order interactions between components and building blocks of chromatin states. a**, Factors against clusters shown with ChromHMM percentages of each margin. Clusters were calculated using k-centroids. **b-c**, Two insets from panel (a) are highlighted. **b**, Promoters (red), which include de-repressive demethylases UTX and LSD1 versus enhancers (orange/yellow), which contain CHD7 and TFs such as NANOG and POU5F1. **c**, Repressive and bivalent elements are marked by either PRC1 and PRC2 or EZH2 with KDM6A (H3K27me3 demethylase), and SAP30 characterizes a set of inactive promoters. **d**, SAP30 components show high diversity of chromatin state roles, ranging from heterochromatin to, separately, weakly and strongly active promoters.

Across all classes of chromatin states, promoters and enhancers loaded the widest range of different CR and TF factors and were annotated by the largest number of learned components. On the other hand, transcribed states were annotated solely by the learned PAF1 / POLR2A component, which is a known complex crucial for transcription elongation. We identified multiple separate groups of heterochromatic states - those solely annotated by canonical heterochromatin establishing CBX3/CBX5 (HP1), a set of states containing both CBX3/5 in higher order interactions with other factors (such as HDAC1, ZNF143, LSD1, and RBP5), and states marked mostly by EHMT2/G9A and ZNF274, which were highly enriched for heterochromatin and ZNF-enriched states (**Fig. 1-17a**). Interestingly, we identified a small set of genomic clusters that were highly enriched for genic enhancers (light green) and were marked

by a component containing chromatin helicase CHD1 and POL2RA, suggesting a role for CHD1 in maintaining enhancers within genes through active transcription (**Fig. 1-17a**)<sup>142,149</sup>.

While promoters and enhancers showed much more complex modes of chromatin regulator and transcription factor localization, we were able to use the learned components to pick out separate modes of chromatin regulation within these regions. For example, promoters (red) included de-repressive demethylases UTX and LSD1 whereas enhancers (orange/yellow) contained CHD7 and TFs such as NANOG and POU5F1 (**Fig. 1-17a-c**). Enhancers also commonly contained the cohesin complex, including RAD21, CTCF, and SMC3 (not profiled in H1-hESC), in a variety of components with other factors including YY1 and CREB1 (**Fig. 1-17b**). Promoters were further marked by the RBBP5/ASH2L complex as well as the CoREST complex and specific deacetylases (HDAC1/2) (**Fig. 1-17c**).

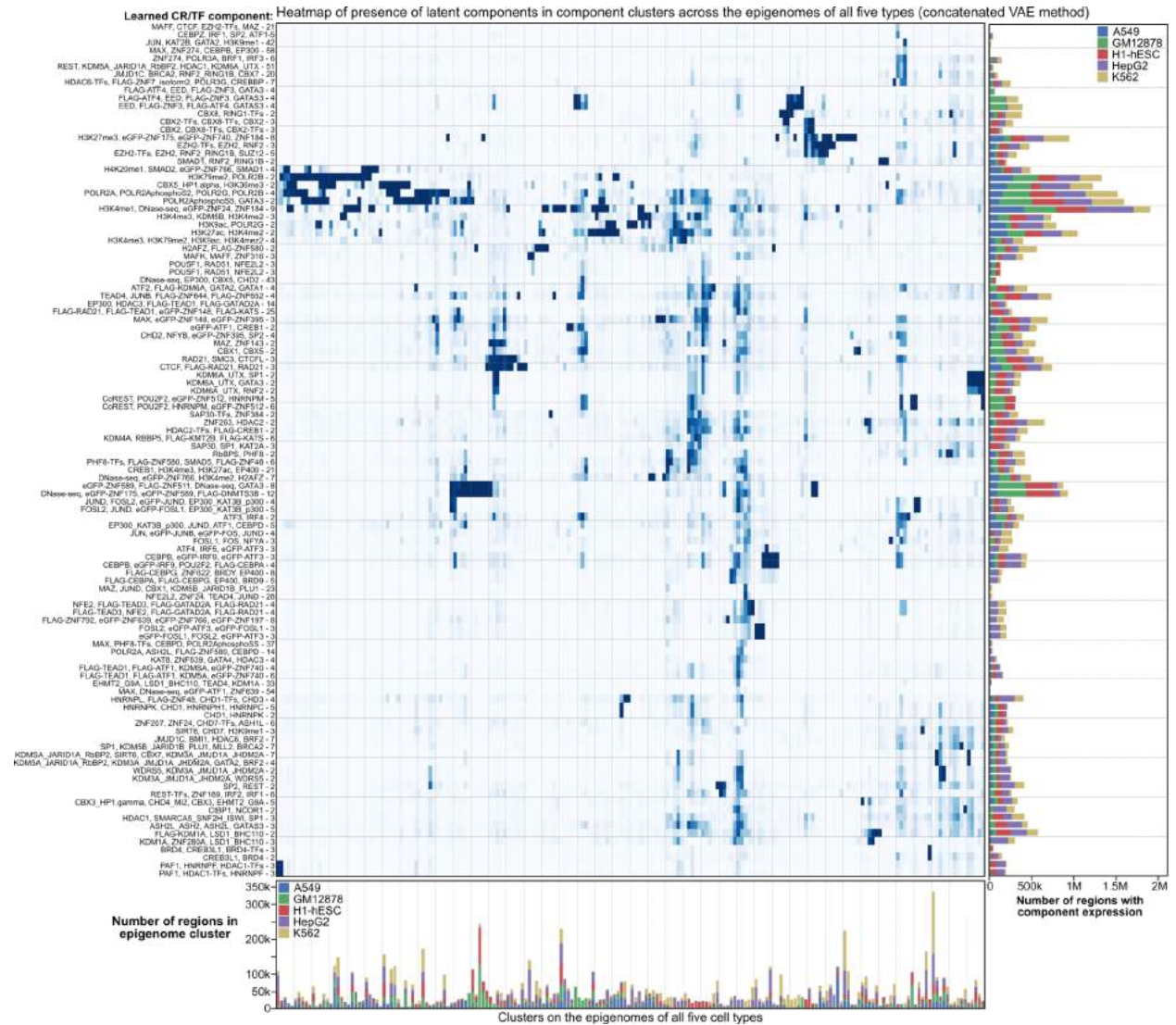
Additionally, while repressive and bivalent elements are marked by either PRC1 and PRC2 or EZH2, often with KDM6A (H3K27me3 demethylase, as noted above), SAP30 occupancy characterized a set of inactive promoters (**Fig. 1-17d**). Components with SAP30, a core histone deacetylase complex member, appeared in a highly diverse set of chromatin state roles, ranging from heterochromatin in conjunction with REST, deacetylases and HP1 factors to weakly (YY1, CTBP2, USF1) or strongly (PHF8, POLR2A, TAF7, CHD1) active promoters (**Fig. 1-17d**).



**Figure 1-18. Components predict histone marks as well as or better than the compendium of CRs and TFs.** Predicting of 13 histone marks was done using logistic regression and either individual factors or components for the K562 epigenome. **a**, auROCs and auPRCs for CRs/TFs (grey) and components (lightblue). **b-c**, Coefficients for regression for (b) CRs/TFs and (c) components.

In addition to cell type-specific models, we also learned joint components of chromatin regulators and transcription factors over all five cell types. While the CR and TF panels were not evaluated evenly across cell types and there were missing ChIP-seq experiments in each dataset, we concatenated the epigenomes for each of the five cell types and were able to leverage the ability of the neural network variational approach to train with missing data. For this model, we also added histone marks as an additional modality for factors. We found broad agreement in the majority of components learned in this joint manner, identifying many of the previously mentioned complexes (e.g. cohesin, CoREST, PRC1, and PRC2) (**Fig. 1-19**).

Specific transcription factors clustered by missingness and dataset availability in this analysis including ATF3, FOSL1, and GATA2D, which were all only measured in HepG2 and K562. The analysis revealed a set of epigenomic locations with specific occupancy patterns for the different cell types and showed consistency in both immediate and higher-order partnerships across the cell types. Future work along these directions should more explicitly model missingness in the data to enable us or others to make strong conclusions about cell type-specific chromatin regulation patterns.



**Figure 1-19. Learning dictionaries over multiple cell types despite missing data.** Five epigenomes clustered together on a common basis of jointly learned components. All cell line epigenomes are used to train a common dictionary and separate, cell line-specific loadings. We mask missing data, balance the training objective by data available, and share a common component-associated coefficient matrix for predicting histone marks across cell types. Cell lines are A549, GM12878, H1-hESC, HepG2, K562 - colors on margins are in the same order, where A549 is blue.

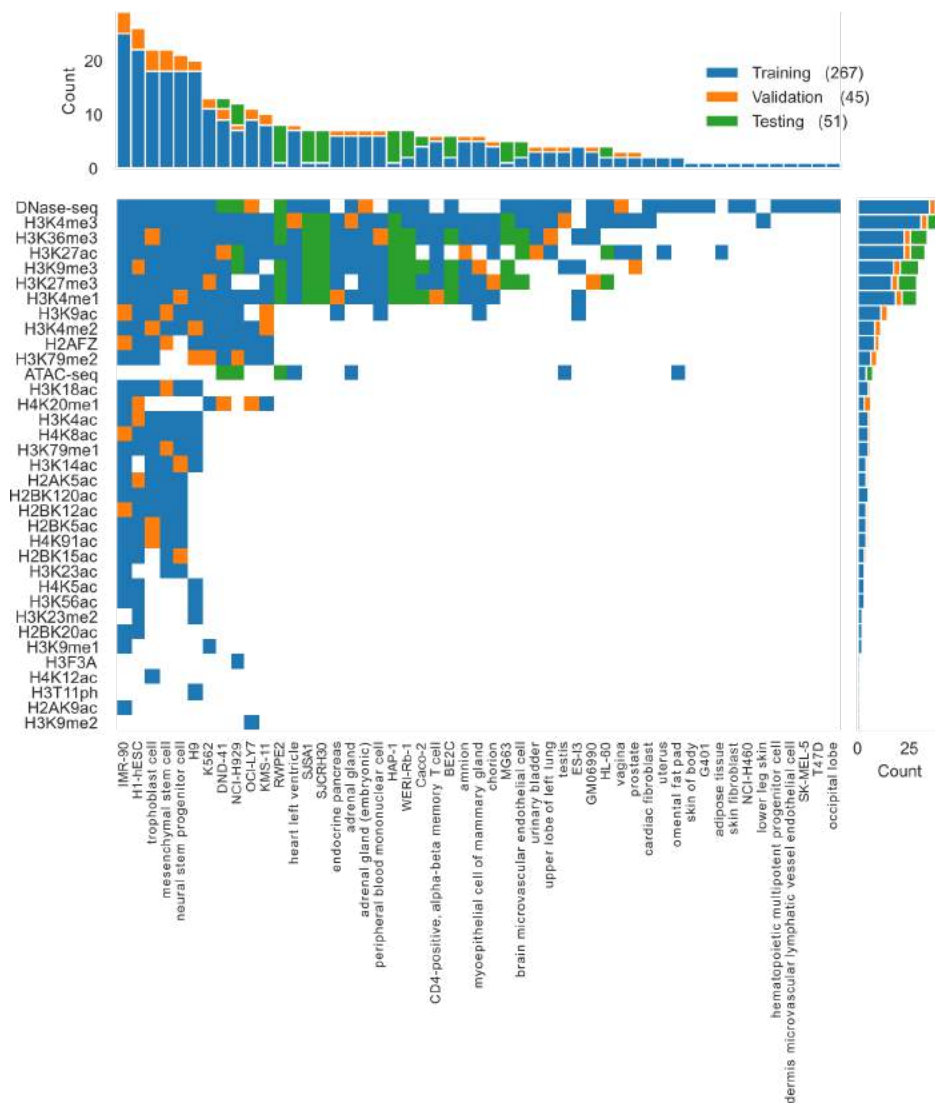


### 1.10 ENCODE Imputation Challenge summary and overview

**Summary:** Although functional genomics experiments are invaluable for understanding cellular mechanisms, comprehensively performing all such experiments, even across a fixed set of sample and assay types, is infeasible in practice. A promising alternative to performing experiments exhaustively is to, instead, perform a core set of experiments and subsequently use machine learning methods to impute the remaining experiments. However, questions remain as to the quality of the imputations, the best approaches for performing imputations, and even what performance measures meaningfully evaluate performance of such models. In this work, we address these questions by comprehensively analyzing imputations from 23 imputation models submitted to the ENCODE Imputation Challenge. We find that measuring the quality of imputations is significantly more challenging than reported in the literature, and is confounded by three factors: major distributional shifts that arise because of differences in data collection and processing over time, the amount of available data per cell type, and redundancy among performance measures. Our systematic analyses suggest several steps that are necessary, but also simple, for fairly evaluating the performance of such models, as well as promising directions for more robust research in this area.

**Challenge overview:** We hosted the ENCODE Imputation Challenge (<https://www.synapse.org/encodeimpute>), a public challenge for imputing epigenomic profiles, which began on February 20th, 2019, and concluded on August 14th, 2019. The challenge evaluated how well predictive models could impute held-out epigenomics experiments using other functional genomics experiments and nucleotide sequence as input (see challenge site for more details). Overall, we acquired 267 data sets from the ENCODE Portal to use as the training set, 45 data sets from the ENCODE Portal to use as a validation set, and performed 51 new experiments as a prospective evaluation to use as a test set (**Fig. 1-20**).

The challenge was divided into two stages. In the first stage, participants were provided with the training and validation data sets as well as a real-time public leaderboard of performance on the held-out validation set. Team BrokenNodes and Hongyang\_Li\_and\_Yuanfang\_Guan tied for first place at the conclusion of the first stage. In the second stage, the teams were allowed to re-organize, and participants were encouraged to refine their models using lessons learned from the first stage. The winners of the second stage, and of the entire challenge, were the top three teams based on performance on the held-out prospective test set, to which the teams did not have access at any point in the challenge.



**Figure 1-20. The ENCODE Imputation Challenge data matrix.** The matrix shows the experiments used in the challenge, colored based on whether they were in the training set (blue), the validation set (orange), or the blind test set (green). White squares indicate that an experiment has not yet been performed. The marginal bar plots show the number of experiments in each assay and cell type.

The challenge was well attended with 196 people signing up on Synapse. Eight teams submitted results for the first round. After teams merged before the second round, 23 imputation models were submitted. Of these models, only one (UIOWA Michaelson) did not submit the full set of required imputations. Although our method for calculating team rankings as a part of the challenge accounted for missing imputations, our subsequent analyses excluded this model.

**Performance Measures:** Prior to the start of the challenge, we specified nine different performance measures to be used in the challenge. These performance measures included (1) the genome-wide mean-squared-error (MSE), (2) the genome-wide Pearson correlation, (3) the genome-wide Spearman correlation, (4) the MSE calculated in promoter regions (defined as +/-2kb from the start of GENCODEv38 annotated genes<sup>1</sup>), (5) the MSE calculated in gene bodies

from GENCODEv38 annotated genes, (6) the MSE calculated in enhancer regions as defined by FANTOM5 annotated permissive enhancers<sup>150</sup>, (7) the MSE weighted at each position by the variance of the experimental signal for that assay across the training set, (8) the MSE at the top 1% of genomic positions ranked by experimental signal, and (9) the MSE at the top 1% of genomic positions ranked by predicted signal. We note that 8 and 9 make a calculation similar to recall and precision, respectively.

We used a multi-stage process, originally developed for the ENCODE Transcription Factor Binding Challenge, to aggregate these performance measures into a single score to determine the challenge winners. First, ten equally-sized bootstraps were drawn from the pool of all genomic positions, and each of the nine performance measures was calculated for each team on each of the bootstraps for each experiment. For each bootstrap-experiment pair, the scores were converted to rankings across teams for each performance measure, and these rankings were then averaged across performance measures. This resulted in a score for each team in each bootstrap-experiment pair. This score was then converted back into a ranking over teams for each bootstrap-experiment pair. Next, these rankings were aggregated across experiments

by calculating  $\frac{1}{|E|} \sum_{e \in E} \min(0.5, r_e)$  where  $E$  is the set of all experiments,  $e$  is an individual experiment, and  $r_e$  is a team's ranking on experiment  $e$ . Finally, a rank was calculated across teams for each bootstrap, and the 90th percentile score, i.e. the second-best bootstrap rank, was used to determine the winners. This procedure is implemented at [https://github.com/ENCODE-DCC/imputation\\_challenge](https://github.com/ENCODE-DCC/imputation_challenge).

**Baseline Methods:** The methods submitted by the participants were compared to two baseline methods. The first baseline was the average activity, which is a straw-man imputation approach that simply predicts the average training set signal at each position in the genome across all cell types for a given assay type<sup>96,151</sup>. Consequently, this approach cannot make cell type-specific predictions. However, it represents the simple rule that regions of the genome that always exhibit peaks in signal and that regions of the genome that never exhibit peaks will continue to do so in other cell types. The second baseline was the Avocado model, using the same model architecture and training procedure described by Schreiber et al.<sup>152</sup>. Importantly, this model was not tuned for this data set---it was applied as-is using the default settings and hyperparameters.

Although we had initially expected that ChromImpute would serve as a baseline in this challenge, for logistical reasons ChromImpute was not applied to the challenge data until well after the challenge concluded. Because the participants did not have access to these predictions, as they did with the other two baselines, we did not include ChromImpute in the original rankings or analysis. However, we have included a ranking of methods that includes ChromImpute in a re-analysis of the challenge participants using the six measures used to evaluate the original ChromImpute method<sup>94</sup>, as a reference. These measures emphasize the relative shape of signals, and included Pearson correlation, three measures quantifying percentage overlap between positions exhibiting high signal, and AUC measures for predicting peaks in observed signal from imputed signal values and vice versa. In order to obtain team ranks on these measures, we first ranked each team's prediction for each test track on each measure separately. We then averaged ranks across metrics and re-assigned integer ranks in each track for each team. Each team's final rank was then computed from the average of their predictions' track ranks for the 51 test tracks.

**Quantile Normalization:** We developed a three-step quantile normalization method for normalizing signals across genomics experiments. Importantly, the normalization is done separately for signal in peak and background regions (as defined by MACSv2 peak calls<sup>114</sup>) to

account for peaks spanning differing proportions of the genome across cell types. In the first step, quantiles are derived separately from each training set experiment. That is, if there are  $N$  training set experiments,  $M_p$  peak quantile bins, and  $M_b$  background quantile bins, one would extract  $Q_p \in R^{N, M_p}$  and  $Q_b \in R^{N, M_b}$ . Quantiles are extracted by ranking all signal values for an experiment (in peaks or outside of peaks, respectively), binning those ranks into either  $M_p$  or  $M_b$  equally sized bins, and assigning to each bin the average signal value from positions within the bin. In the second step, an average is taken across experiments for each quantile bin to construct reference quantiles  $R_p \in R^{M_p}$  and  $R_b \in R^{M_b}$ . Finally,  $R_p$  and  $R_b$  are applied to the test set tracks, with  $R_p$  being applied only within signal peaks and  $R_b$  being applied only within background regions. Because peak regions are more complex and span a larger dynamic range than the background, we set  $M_p$  to be 1000 and  $M_b$  to be 50. Given that this procedure is designed to combat distributional shift, we note that it should be applied to test set experiments before evaluation.

**Data Processing:** We processed the DNase and ATAC-seq experiments using a uniform pipeline<sup>153</sup>. First, FASTQ files containing read sequences and quality scores for the training and validation sets experiments were downloaded from the ENCODE Portal, and FASTQs for the test set experiments were acquired from our own experiments. For ATAC-seq experiments (but not DNase-seq), we first trimmed adapters and then mapped reads to the hg38 reference human genome using the Bowtie2 aligner<sup>154</sup>. After mapping, reads were filtered to remove unmapped reads and mates, non-primary alignments, reads failing platform/vendor quality checks, and PCR/optimal duplicates (-F 1804). Reads mapping reliably to more than one location (MAPQ < 30), i.e., multi-mapping reads, were removed. Duplicate reads were then marked with Picard MarkDuplicates and removed<sup>155</sup>. For single-end DNase data sets, a single read was chosen from a set of duplicate reads, whereas for paired-end data sets, reads were chosen if either end was unique. Although this is the standard approach for processing these types of data, this step had unintended consequences for the challenge, as described in the results below. Reads from biological and technical replicates were merged. We normalized the sequencing depth across data sets by subsampling them to a maximum of 50 million reads (after excluding reads mapping to mitochondria). Although there are several ways to represent the signal from sequencing experiments, e.g., read-counts and fold-change, we chose to use the statistical significance of the fold-change to be consistent with previous imputation literature. We used the MACSv2 peak caller to compute the fold-enrichment and statistical significance. MACSv2 was applied to smoothed counts (150 bp smoothing window) of read-starts (5' ends of reads) at each position in the genome relative to the expected number of reads from a local Poisson-simulated background distribution. We provided the genome-wide signal tracks containing the statistical significance of enrichment (i.e., the  $-\log_{10}$  p-values) at each base pair in the genome. The processing pipeline is open-source and available at <https://github.com/ENCODE-DCC/atac-seq-pipeline>.

Next, we processed the histone ChIP-seq experiments using the ENCODE processing pipeline. For each experiment we downloaded FASTQ files from the ENCODE Portal for at least two replicate experiments and a control experiment. All reads were mapped to the hg38 reference human genome using the BWA aligner<sup>156</sup>. After mapping, the process was similar to the ATAC-seq/DNase-seq pipeline. Reads were filtered to remove unmapped reads and mates, non-primary alignments, reads failing platform/vendor quality checks, and PCR/optimal duplicates (-F 1804). Multi-mapping reads (MAPQ < 30) were also removed. Duplicates were identified using Picard MarkDuplicates and subsequently removed, with the same issues as mentioned for DNase data sets. Reads from the biological and technical replicates were then

merged. We normalized the sequencing depth across data sets by subsampling each to a maximum of 50 million reads. We used the MACSv2 peak caller to calculate fold-enrichment and statistical significance of counts of extended ChIP-seq reads (reads were extended in the 5' to 3' direction based on the predominant fragment length), relative to the number of extended reads from the control experiment. The statistical significance of the enrichment was computed using a local Poisson null distribution whose mean parameter is estimated from the control experiment. For the purposes of this challenge, we provided the genome-wide signal tracks containing the statistical significance of enrichment (i.e., the  $-\log_{10}$  p-values) at each base pair in the genome. The processing pipeline is open-source and available at <https://github.com/ENCODE-DCC/chip-seq-pipeline2>.

### 1.11 Results of the ENCODE Imputation Challenge

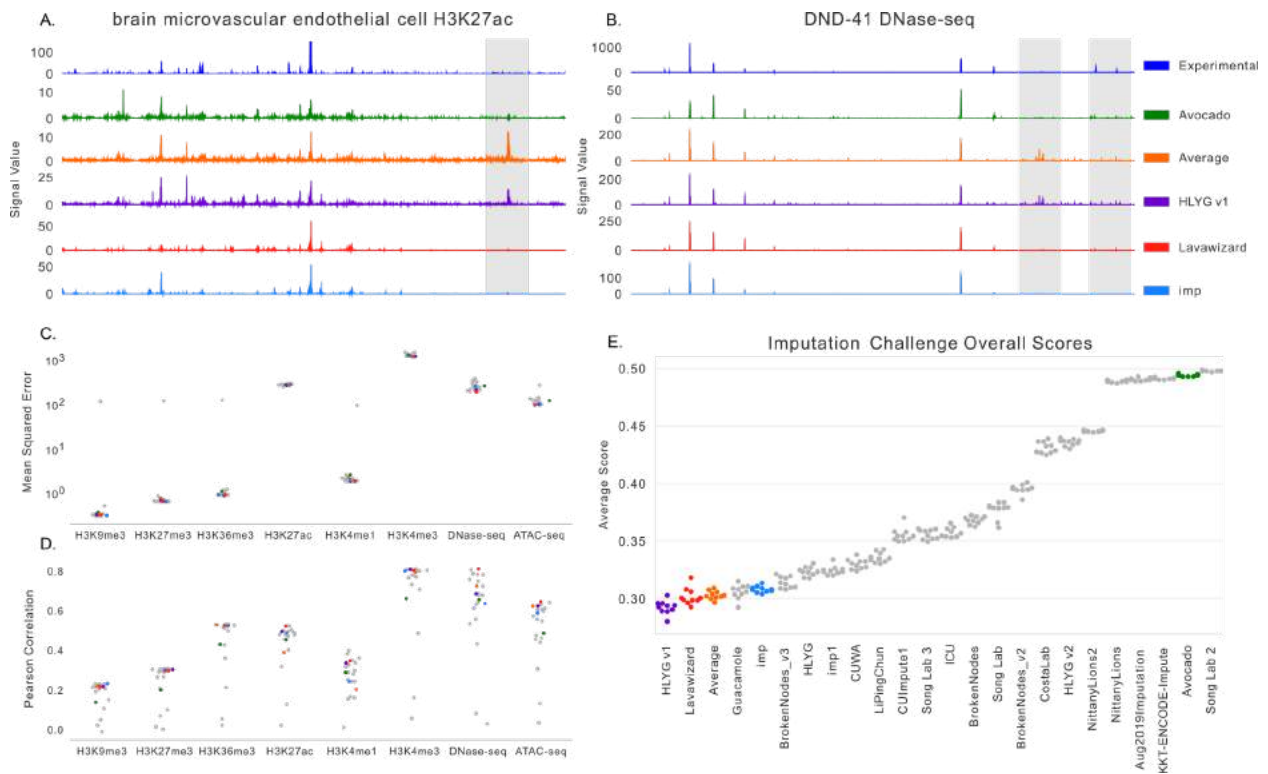
Participants submitted 23 models to the second stage of the challenge. Each group was allowed to submit up to three models to encourage inclusion of unorthodox solutions with at least one submission. As a result, the models encompassed a diverse range of strategies (see **Table 1-1**). The models differed primarily along three axes. The first axis was the signal preprocessing, with almost every method further preprocessing the data from the given  $-\log_{10}$  signal p-values. The second axis was the data sources used to construct input features. Most methods followed previously published methods by only using assay measurements as inputs (denoted “functional” in **Table 1-1**). However, five of the methods used nucleotide sequence as input, eight methods used the average activity baseline, and three used Avocado’s imputations. The third axis was the manner in which the underlying tensor structure of the data was modeled. Some methods explicitly modeled the data as a tensor (e.g., imp and Lavawizard), whereas other methods only implicitly modeled the structure through rule-based approaches or similarity methods (e.g., the HLYG and KNN-based approaches).

Name	Model	Norm	Inputs			
			Sequence	Functional	Average	Avocado
Aug2019Impute						
BrokenNodes/v2	KNN	arcsinh		✓		
BrokenNodes v3	KNN	arcsinh		✓		✓
CostaLab v2						
CUImpute1/CUWA/ICU	ensemble	arcsinh		✓	✓	✓
Guacamole/Lavawizard	DTF	arcsinh		✓	✓	
HLYG/v1/v2	GBT	quantile	✓	✓	✓	
imp/imp1	DTF+AE	Cauchy		✓		
KKT-ENCODE	CNN	arcsinh	✓			
LiPingChun		arcsinh			✓	
NittanyLions	KNN			✓		
NittanyLions2	KNN	quantile		✓		
SongLab	CNN	log1p		✓		
SongLab2	HMM			✓		
SongLab3	CNN	log1p		✓	✓	✓
UIOWA	CNN	quantile	✓	✓		

**Table 1-1. Methodologies of imputation methods.** The table lists the modeling strategies and input features used by each of the models, as reported by the teams. The models include k-nearest neighbors (KNN), deep tensor factorization (DTF), autoencoders (AE), convolutional neural networks (CNN), hidden Markov models (HMM), and gradient-boosted decision trees (GBT). The authors of Aug2019Impute and CostaLab v2 did not describe their methods



An initial inspection of the imputations revealed that most methods captured the general shape of the signal well. Examples drawn from H3K27ac in brain microvascular endothelial cells and DNase-seq in DND-41 cells (**Fig. 1-21A/B**) suggest two sources of error: the misprediction of a small number of peaks relative to the total number of true peaks, and the misprediction of the precise signal value within correctly predicted peaks. Focusing on the misprediction of peaks, we noted that some methods made similar mistakes as the average activity baseline, whereas others made similar mistakes as the Avocado baseline (gray highlights in **Fig. 1-21A/B**). Unsurprisingly, methods that used Avocado's imputations as input had the highest genome-wide correlation with Avocado's predictions (it is worth noting that CUImpute1 only used Avocado's imputations for some, but not all, assays). In contrast, methods that explicitly used the average activity did not always exhibit higher correlation with it than other methods. This finding suggests that, because the average activity can be directly derived from the training set, many types of models are able to implicitly learn it even when not explicitly trained on it.



**Figure 1-21. Results from the ENCODE Imputation Challenge.** (A) The H3K27ac signal for brain microvascular endothelial cells that is observed (in blue), from baseline methods, and from the winning three teams in the challenge. (B) The same as (A) except showing DNase-seq signal in DND-41 cells. (C) The average MSE for each method across test set tracks and bootstraps but partitioned by assay type. (D) The same as (C) except for Pearson correlation. (E) The overall score, calculated as described in the performance measures section of the methods, across all test set tracks and performance measures shown for each bootstrap for each team. The baseline methods and winners are colored.

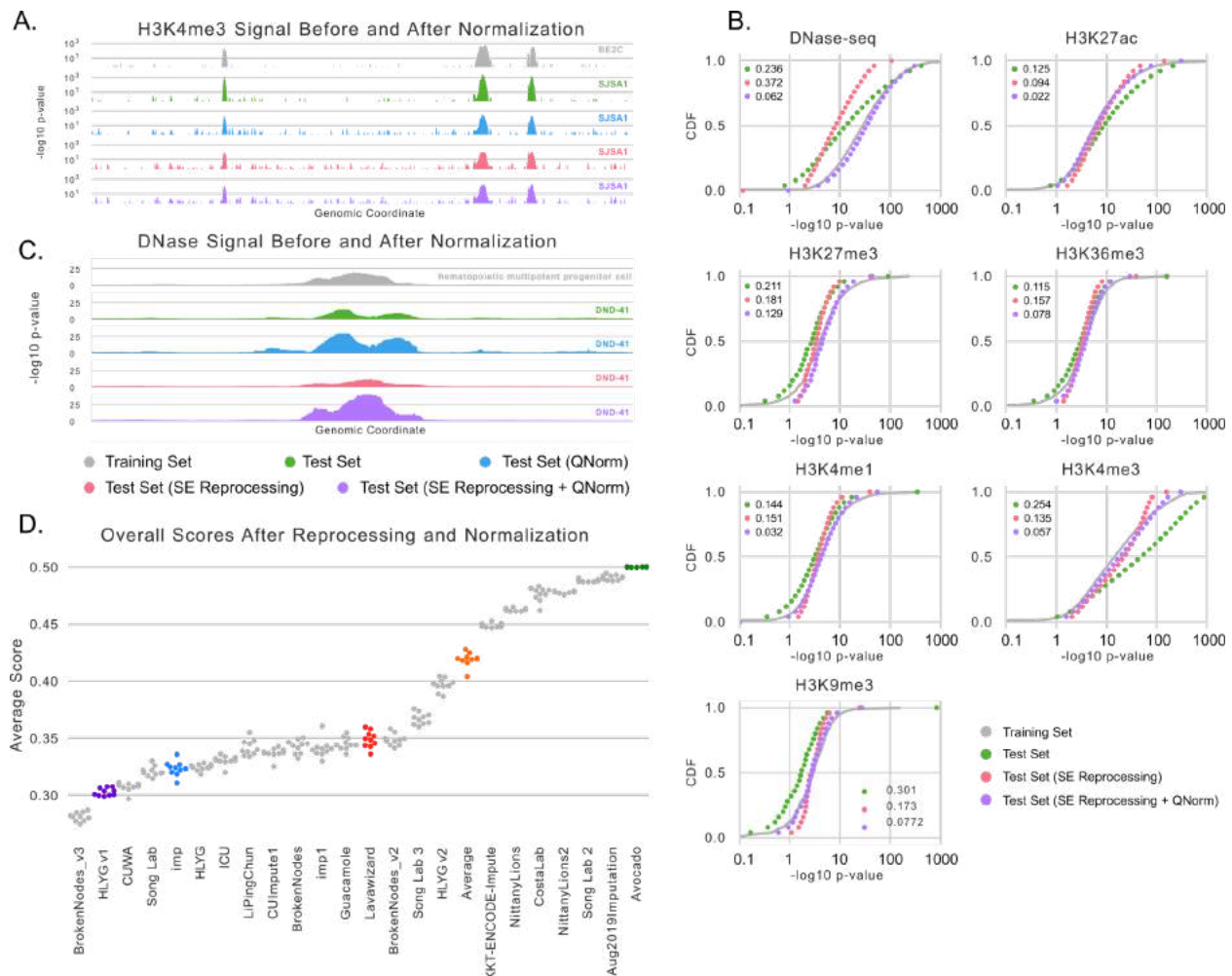
Next, we comprehensively evaluated the methods using a battery of performance measures that were specified at the beginning of the challenge (see Performance Measures, above). We found that performance on these measures depended heavily on the imputed assay (**Fig. 1-21C/D**).

For instance, most models exhibited four orders of magnitude higher MSE on H3K4me3 than on H3K9me3. However, several assays that exhibited the highest MSE also exhibited the highest Pearson correlation, indicating that the scale of MSE across assays is likely more related to the dynamic range of the assay rather than the accuracy of the imputations. Unsurprisingly, a projection of all imputed and experimental tracks clustered predominately by assay type (Silhouette Score = 0.4601), as opposed to by cell type (SS = -0.4028) or imputation method (SS = -0.3133). Accordingly, we used a rank-based transform to account for differences in dynamic range when calculating global performance measures across experiments to ensure that assays with large dynamic ranges didn't dominate the evaluation. After calculating the global performance of each method, we found that there was a gradient of methods that performed increasingly well, and a set of methods that performed relatively poorly (**Fig. 1-21E**). The best performing methods, and hence the winners of the challenge, were Hongyang Li and Yuangfang Guan v1 (abbreviated as "HLYGv1") in first place, Lavawizard and Guacamole (two similar methods from the same team) tied for second place, and imp in third place.

Given the diverse modeling strategies of the winning teams, our primary take-away from these results is that there does not appear to be a single key insight that led to good overall performance on the measures used in the challenge. HLYGv1 used nucleotide sequence as input, but so did KKT-ENCODE and UIOWA Michaelson; all three models submitted by Hongyang Li and Yuanfang Guan used gradient boosted trees (GBTs), yet their models exhibited both good and poor performance. However, these results do suggest certain models to be wary of: convolutional neural networks and k-nearest neighbor models underperformed deep tensor factorization (DTF) and GBT models. This is likely because the similarities used by KNN models are a less sophisticated version of the representations learned by tensor factorization approaches, and that the specific structure presented in the data is not well modeled by simple applications of convolutions.

However, when we compared model performance to the baseline methods, we made two important observations. First, almost every team outperformed the Avocado baseline, as one might expect because the participants had access to the Avocado model and predictions during the development process, and because the default settings were used for Avocado despite them being tuned for significantly larger amounts of training data. Second, the average activity baseline performs extremely well, coming in third in our ranking and first place in five of the nine performance measures used. Both of these observations are a reversal from the first round in the challenge, where Avocado outperformed all the participants but almost all the participants outperformed the average activity baseline. This reversal in performance between the two baselines is partially because the evaluation setting changed from overrepresenting well characterized cell types to focusing on poorly characterized ones and, as we will see later, partially due to the performance measures used for the challenge.

**Accounting for distributional shift:** A visual inspection of the test set experiments revealed significant distributional differences in peak signal values between the training and test sets for some assays (**Fig. 1-22A**). Most obviously, the signal values within H3K4me3 peaks from test set experiments were generally much higher than the signal values within peaks from training and validation set experiments (**Fig. 1-22B**). Although one would expect a locus to exhibit different signal in different cell types because of real biology, one would also expect that the distribution of signal values within peaks across entire experiments would be similar for experiments of the same assay. Because distributional shifts have major ramifications for the scale-based performance measures used in the challenge, we next investigated the source of these distributional differences.



**Figure 1-22. Distributional shift and quantile normalization.** (A) Experimental signal measuring H3K4me3 in BE2C cells from an unnormalized training set experiment (gray), an unnormalized test set experiment in SJS1A1 cells (green), the test set signal after quantile normalization (blue), the test set signal after single-end reprocessing (red), and the test set signal after single-end reprocessing and quantile normalization (purple). (B) Distributions of signal values within peaks in chr16/17 for each reprocessed assay across the unnormalized training set (gray), the unnormalized test set (green), the single-end reprocessed test set (red), and the single-end reprocessed and quantile-normalized test set (purple). The KS statistics between the training set distribution and the test set distributions are shown in the legends and the CDFs are summarized using 25 dots for visualization purposes. (C) An example locus that exhibits a DNase peak in both the training and test sets. (D) A re-scoring of the challenge participants against single-end reprocessed and quantile-normalized test set signal.

After considering several potential covariates that could explain this distribution shift, including multiple measures of experimental quality, we found that the primary driver was a subtle difference in how the test set experiments were processed. By design, the test set experiments were performed during the challenge to ensure a truly prospective evaluation. However, experimental methods have changed in the many years since the training data were collected. Most notably, collecting paired-end data is now the standard approach for ENCODE data sets because the procedure yields higher quality data and is now cheap enough for broad usage;

however, almost all of the training set experiments predate this switch and involve single-end data. The processing of single-end and paired-end data is largely similar, but a crucial difference occurs in the deduplication step. Specifically, deduplication of single-end reads using PICARD allows the mapping of only one read start to each position on the genome on each strand<sup>155</sup>. In contrast, deduplication of paired-end data can result in more than one read-start per position on each strand because read-pairs are only removed if the read start of both ends are duplicates. Consequently, the number of reads mapping within peaks from paired-end data can be significantly higher than what one would get using single-end data. Importantly, the shift is not simply caused by paired-end data being higher quality, as we first explored, but rather differences in the deduplication step.

We confirmed that differences in processing, rather than differences in data quality, explained the distributional shift by reprocessing the paired-end data sets (except for ATAC-seq which requires paired-end data) as single-end data. Specifically, for each paired-end experiment in the test set we concatenated the fastq files of reads from both ends and ran the same single-end processing pipeline that was run on the other single-end experiments in the challenge. We found that the reprocessed data had distributions of peak signal values significantly closer to the training set, as measured by the Kolmogorov-Smirnov (KS) statistic, for four of the histone modification assays including H3K4me3 (**Fig. 1-22B**). The remaining two histone modification assays already resembled the training set before reprocessing. However, we found that the distribution of DNase-seq peak signal values had a larger KS-statistic after reprocessing than before. This is likely because 21 of the 38 training set experiments contained paired-end data, which would shift the distribution of signal values in the training set up. Although the most principled next step would be to reprocess all of the experiments used in the challenge and subsequently re-training and re-evaluating each submission, this analysis was not possible because we did not require that participants submit code that could retrain their models on new data sets. Given no perfect solution, we chose to continue with the single-end reprocessed test set tracks for our subsequent analyses.

We found that reprocessing the histone modification data significantly reduced the distributional shift but did not perfectly correct it. The remaining differences are likely related to small changes in experimental protocol over time, such as improvements in sequencing technology, antibodies used, and read lengths measured. A general-purpose correction for the remaining differences is to explicitly quantile normalize the data such that the signal values in the testing experiments exhibit the same signal distribution as those in the training experiments. Quantile normalization is powerful because it is a non-linear method, in contrast to min-max or z-score scaling, and has been extensively applied to genomics data sets, including those measuring bulk gene expression<sup>157</sup>, single-cell RNA-seq<sup>158</sup>, and ChIP-seq data when combined with a spike-in reference<sup>159</sup>. We account for differing proportions of the genome exhibiting peaks across cell types by separately quantile normalizing the signal within peaks and the signal in background regions. Finally, because the distribution of signal is significantly different across assays, we apply this quantile normalization to each assay separately. After normalization, we confirmed that the distribution of within-peak signal values was almost identical to the distribution of within-peak training signal values across all assays (**Fig. 1-22B**), even for the DNase-seq experiments.

In theory, one could apply quantile normalization to the original paired-end test set data and, by definition, produce signal values with the same distribution without the need for reprocessing. However, when looking at a representative DNase peak, we found that the reprocessed data was not a simple monotonic transform of the original data (**Fig. 1-22C**). Specifically, the paired-end data exhibited a peak shape unlike that observed in the single-end data, and simply quantile normalizing the signal does not fix the differences in shape. More comprehensively, when considering a 10 Mbp region of chr1 on each of the 48 reprocessed experiments, we

clearly observed that paired-end data is not a monotonic transformation of single-end data. Although the assays associated with activity, such as H3K4me3 and DNase-seq, exhibit Spearman correlations up to 0.938 between the paired-end and single-end processed signals, repressive marks exhibit Spearman correlations as low as 0.037 and the average Spearman correlation across all tracks was only 0.453. Further, even though some assays exhibit high correlation, this value is inflated by the large number of low-signal values and, indeed, the largest variability comes at loci with high signal values.

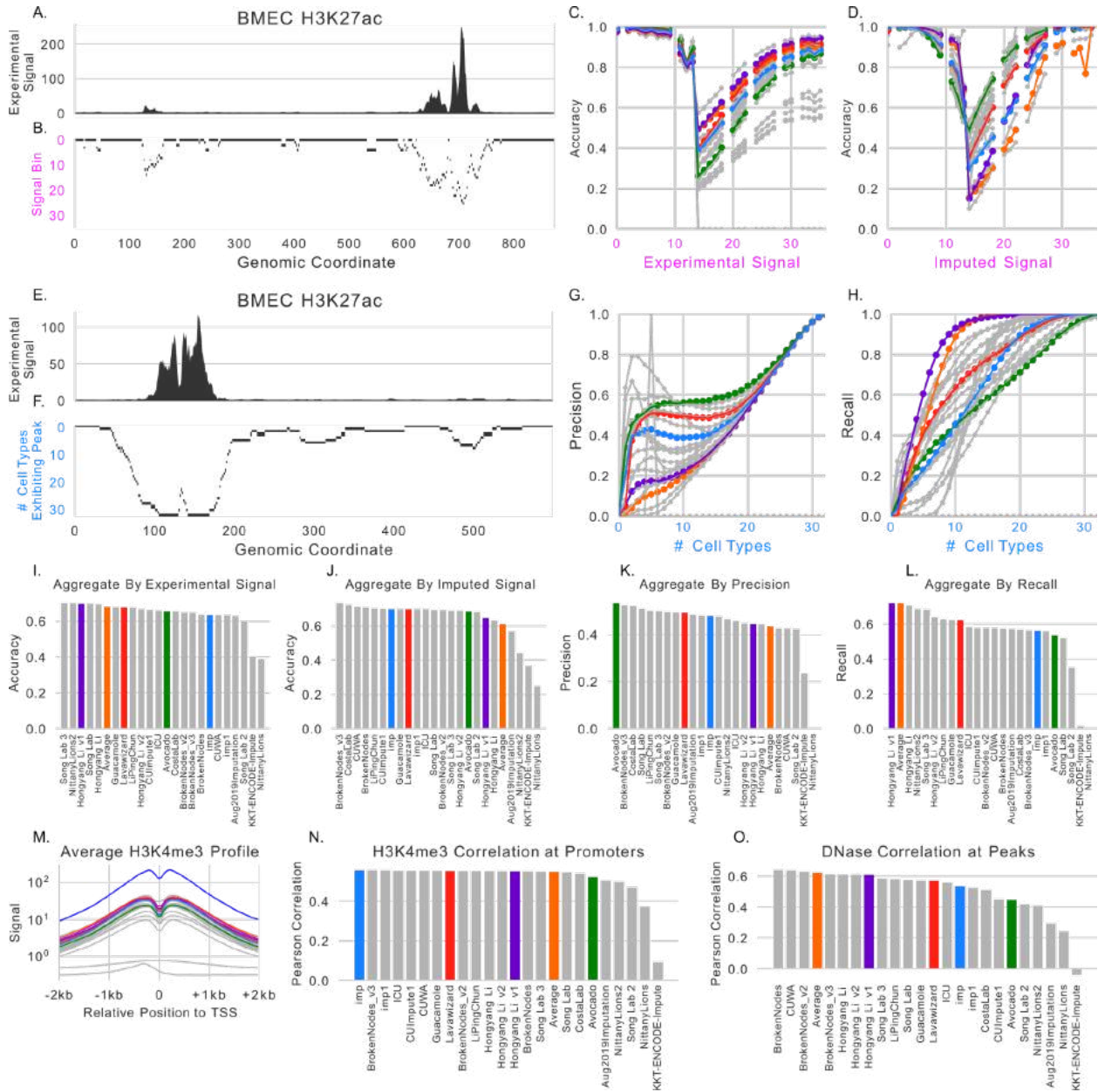
Moving forward with our method of reprocessing the test data using single-end settings and then quantile normalizing to correct the remaining differences, we next re-scored the originally submitted imputations (**Fig. 1-22D**). We observed that the number of methods outperforming the average activity baseline increased from two to 16 and that BrokenNodes v3 rose from sixth place to first place in the rankings. Although HLYGv1 remained within the top three, the other two winners descended in the rankings. This might be explained by HLYGv1 using quantile normalization, albeit a slightly different version than the one we used, during training. Interestingly, many of the methods performed similarly to each other, reinforcing the idea found in the original challenge that there is not necessarily one way to do imputation. Indeed, the best performing model is a simple KNN-based approach using arcsinh-transformed data and the second best performing model uses gradient-boosting trees on quantile transformed data. Critically, we note that it would not be fair to use these rankings to declare challenge winners because we did not give the teams an opportunity to retrain or tune their methods on the transformed data. Rather, our take-away is that the distributional shift is partially responsible for the good performance of the average activity baseline but does not fully explain it.

**Designing more informative performance measures:** Although the measures used in the challenge were devised to rank methods independently for each experiment based on their genome-wide (or across large portions of the genome) performance, this property meant that they ultimately exhibited a high degree of redundancy with each other. Essentially, by uniformly weighting all positions along the genome, methods with low genome-wide MSE were likely to have low MSE within promoters, gene bodies, or the top 1% of signal as well. Exacerbating this issue, MSE-based measures were disproportionately confounded by the large distributional shift described in the previous section in comparison to the shape-based measures. Illustrating this, we found that most of the residual---sometimes over 99% in H3K4me3 assays---came at correctly-predicted peaks. Realizing this weakness, we next designed three new types of performance measures that, respectively, reweighted genomic bins based on signal strength, considered multiple experiments simultaneously, and focused on shape within active areas. All evaluations in this section are done against the reprocessed, quantile-normalized test set signal.

A strategy for measuring performance in a complementary way to uniformly-weighted genome-wide performance is to explicitly calculate the performance with respect to the magnitude of the observed, or imputed, signal (**Fig. 1-23A/B**). Rather than being limited by considering only the top 1% bin of signal, such as by using the `mse1obs` or `mse1imp` measures, considering all signal bins provides a finer-grained view of model performance. As an example, if the imputations exhibit high accuracy when the imputed signal is high, then one may be confident that predicted peaks are correct when using imputations for which there is no corresponding experimental data; in contrast, if the imputations exhibit low accuracy when the imputed signal is high, then one might be more skeptical of peak calls but more trusting of regions not called as peaks, e.g. facultative peaks that are not active in the studied cell types. Although any measure can be partitioned by signal magnitude, we focus on accuracy between binarized imputed signal and peak calls for the experimental signal. Accuracy was excluded from the original set of performance measures because the sparsity of peaks can make it difficult to interpret genome-wide; in this setting, we anticipate accuracy to be more valuable in the signal bins where one might reasonably find a peak. Importantly, we did not use rank-based



classification measures (e.g., AUROC or AUPR) here, because once the signal is partitioned by strength, applying a rank-based measure to each bin is less meaningful than when applied genome-wide.



**Figure 1-23. Additional performance measures.** (A) Experimentally observed signal for H3K27ac in brain microvascular endothelial cells. (B) An example of partitioning the track from A into logarithmically-spaced bins (the rows). (C) The accuracy between binarized imputations and MACS2 peak calls for each signal bin when using the experimental signal to define the bins. (D) The same as C except using the imputed signal to define the bins. (E) The same as A but a different locus. (F) The same as B except calculating bins using the number of cell types that each locus exhibits a peak in. (G) The precision of the binarized imputed signal against MACS2 peak calls when evaluated separately for each bin. (H) The same as H except the recall instead of the precision. (I) The average area under the curves, calculated as shown in C, across all test

set tracks for each participant. (J) The average area under the curves calculated as shown in D across all test set tracks for each participant. (K) The precision score calculated in the same manner as I/J. (L) The same as K, except the recall score. (M) The average H3K4me3 profile of experimental (blue), quantile-normalized (magenta), and imputed signals at strand-corrected promoters. (N) The average Pearson correlation between imputed and quantile-normalized signal across all promoters and H3K4me3 test set tracks. (O) The average Pearson correlation between imputed and quantile-normalized signal across all observed DNase peaks and DNase test set tracks.

When we partitioned genomic loci based on experimental signal, we found that model performance aggregated across all tracks generally falls into three regimes: (1) when the imputed or experimental signal is low, the accuracy is high, (2) when the imputed or experimental signal is between 1 and 10 the accuracy severely drops, and (3) when the imputed or experimental signal is high, the accuracy returns to being high (**Fig. 1-23C**). Although the second regime includes ambiguous peak calls, it also includes the most difficult to call peaks (thus, the relatively low accuracy) and should be emphasized by performance measures. When focusing on H3K27ac signal in brain microvascular endothelial cells we can also see that there is a flipping of rankings between the first and second regime; imp1 and LavaWizard both outperform the average activity and HLYGv1 when the experimental signal is low, but perform significantly lower when the experimental signal is higher.

Interestingly, the ranking of methods when partitioning genomic loci using imputed signal seems to be a reversal from using experimental signal. However, the ranking of methods is significantly different when partitioning genomic loci using the imputed signal instead of the experimental signal (**Fig. 1-23D**). HLYGv1 and the average activity are among the top performers when partitioning by experimental signal but are among the worst performers when partitioning by imputed signal. An explanation for this flip is that these approaches measure notions of precision and recall, respectively, which have a known trade-off. Because the average activity is essentially a union of peaks across cell types in the training set, it will have a high recall but a low precision. Methods, such as HLYGv1, that rely too heavily on the average activity will exhibit the same tradeoff.

A straightforward way to condense these curves into a single value for a performance measure is to take the average value across the curve. This value is essentially a re-weighting of genome-wide accuracy that uniformly values each bin of signal values rather than each locus, and so will downweight the more common low signal value loci and upweight the less common higher signal values ones. Notably, the winners of the ENCODE Imputation Challenge did not perform the best across all test set experiments when partitioning by either experimental signal or by imputed signal (**Fig. 1-23E/F**). Indeed, the top two performers when partitioning by experimental signal (Song Lab 3 and NittanyLions2) came in 12th and 19th respectively in the original evaluation.

A primary source of error for imputation models comes from loci that exhibit functional activity in some, but not all, cell types<sup>151</sup>. Evaluating whether the imputations can distinguish between cell types that do and do not exhibit signal at a given locus is crucial for ensuring that the imputations are cell type-specific. However, because traditional genome-wide performance measures treat each experiment independently, they cannot explicitly evaluate this property. To better understand how well these methods can identify what cell types loci are active in, for each assay we partitioned genomic positions by the number of experiments that exhibit a peak for that assay and then evaluated each partition separately. For example, if a locus exhibited a DNase-seq peak in 3 out of 5 cell types, that locus would be grouped for evaluation with other loci that also exhibited DNase-seq peaks in 3 out of 5 cell types (**Fig. 1-23G**).

We observe trends that are reminiscent of partitioning loci by signal strength. As the number of

cell types that exhibit peaks increases, so too does the precision and recall of the methods (**Fig. 1-23H/I**). This indicates that, generally, imputation methods are better at predicting peaks at loci that are peaks in many cell types than they are at predicting cell type-specific activity.

Interestingly, we noted that several methods had peaks in precision when the number of cell types the peak was expressed in is low. Given that performance was extremely variable in this regime, we think that focusing on this measure in future studies will be useful when comparing models. Consistent with the role that the average activity plays as essentially the union of peaks across cell types, we see that it has a low aggregate precision score across all test set tracks but has the second highest aggregate recall score (**Fig. 1-23J/K**). Put another way, the average activity is very good at identifying peaks that are common across many cell types but very poor at identifying the cell types that cell type-specific peaks occur in. Somewhat surprisingly, the Avocado baseline had the highest aggregate precision score, but the challenge winners that most resemble it (Lavawizard and imp) did not exhibit the most similar performance.

The performance measures that have been proposed so far predominantly involve genome-wide calculations, even if they involve re-weighting loci contributions. An alternate form of performance measure is to focus on specific forms of biochemical activity at loci that are known to be relevant. The MSEProm, MSEEnh, and MSEGene measures attempt to quantify this by focusing on promoters, enhancers, and gene bodies respectively, but measure the performance of all assays at these loci. Rather, we next investigate two more performance measures that follow the reasoning of Ernst et al. that only specific assays should be measured at these loci<sup>94</sup>.

The first measure evaluates the shape of H3K4me3 signal at promoter regions. This histone modification is known to be enriched at promoter elements and is indicative of active transcription. Further, after correcting for the strand of the promoter, the mark exhibits a distinctive bimodal pattern (**Fig. 1-23L**). We reasoned that focusing on the ability to recapture this shape would provide an orthogonal evaluation to the other performance measures proposed so far. We calculated the average Pearson correlation between the imputed signal and the quantile-normalized experimental signal across all genes for all test set tracks measuring H3K4me3. The average correlations showed that many methods outperformed the average activity baseline. However, consistent with previous results, none of the challenge winners were in the top five according to this measure.

The second measure evaluates the shape of DNase signal at observed DNase peaks. We anticipated that recapturing the shape of DNase signal would be more challenging because DNase does not exhibit a pattern that is as consistent as H3K4me3 at promoter regions. Further, the subtle patterns encoded in DNase signal can be useful for deciphering the precise regulatory role that the underlying nucleotide sequence is playing. Consistent with predicting DNase signal being a more challenging task, we found that methods exhibited a wider range of performances than they did with H3K4me3 prediction. We also found that only two methods outperformed the average activity baseline. This might initially be counterintuitive, because chromatin accessibility is fairly cell type-specific. However, because this evaluation is limited to observed DNase peaks, methods are not being penalized for incorrectly predicting that non-peak regions are exhibiting peaks. This observation indicates that loci may retain the shape of their peaks across cell types (when binned at 25 bp resolution).

## 1.12 Discussion

In this work, we presented a comprehensive map of the human epigenome, EpiMap, encompassing ~15k epigenomic tracks across 833 distinct biological samples that greatly expand the coverage of both embryonic and adult tissues and cells. We combined observed and imputed datasets across 18 epigenomic marks, and combined them to jointly annotate and distinguish diverse classes of chromatin states, including enhancer, promoter, transcribed, repressed, and quiescent regions. We extensively validated the high quality of our annotations and found that they outperformed stringent benchmarks, using both held-out and external experimental datasets for validation.

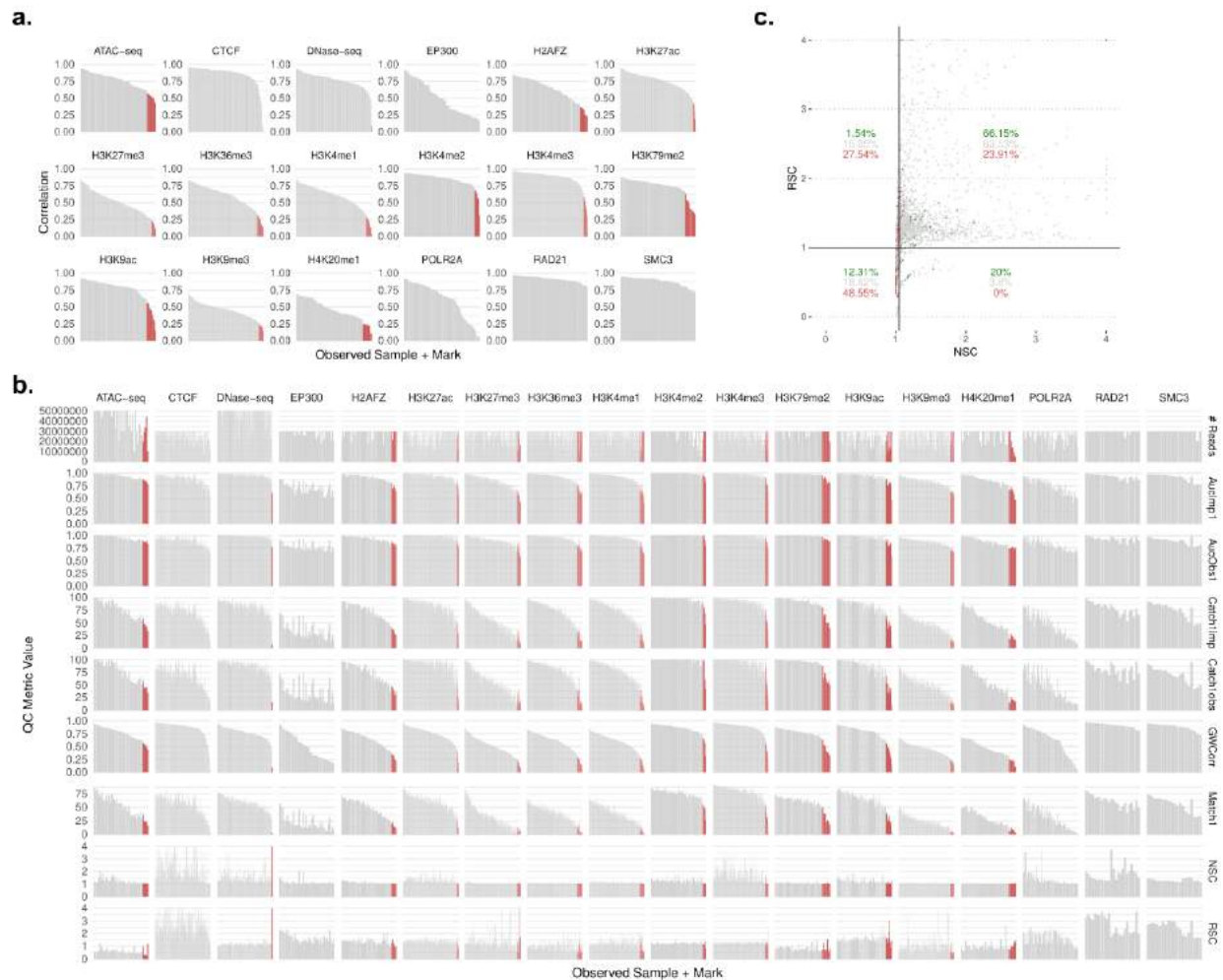
We used this resource to assemble a comprehensive view of human genome circuitry across primary tissues, cells, and cell lines, annotating: 2.1M high-resolution gene-regulatory regions; their activity patterns across 833 biosamples; their enriched regulatory motifs, motif combinations, and putative upstream regulators responsible for their co-regulation; their enriched gene functions and biological pathways that they likely control; and their tissue-specific target genes. Our high-resolution enhancer annotations provide a highly-concentrated view of the gene-regulatory landscape, yielding many gene-regulatory insights but covering only 0.8% of the genome in each sample and only 13% total across all samples. Our linking revealed the high number of enhancers controlling each gene and the high tissue-specificity of long-range enhancer-gene links. Our upstream regulator analysis revealed a highly combinatorial and hierarchical view of gene regulation, with a small number of master regulators (eg. *RFX1-5*, *GRHL1*, *HNF1A*, *AP-1*) interacting with diverse partners in different tissues to define tissue-specific gene-regulatory programs.

Our collection has several limitations: tissue samples are not at single-cell resolution; we do not consider donor genotype or phenotype; imputation may result in increased homogeneity and miss rare sample-specific events; and we still miss many tissues, environmental and stimulation conditions, and developmental stages. In addition, our work relies on highly-quality imputation of epigenomic tracks to fill in missing assays in under-characterized conditions and samples.

To spur the development of new methods for the imputation and modeling of epigenomic profiles across the genome, we ran the ENCODE imputation challenge, where teams across the world competed to innovatively develop new imputation methods for genomic tracks. In the course of the challenge, we identified specific and pervasive issues for evaluating performance and in comparing new, deeper sequenced and paired-end datasets to historical single-end sequencing datasets and we proposed a normalization scheme to alleviate these issues.

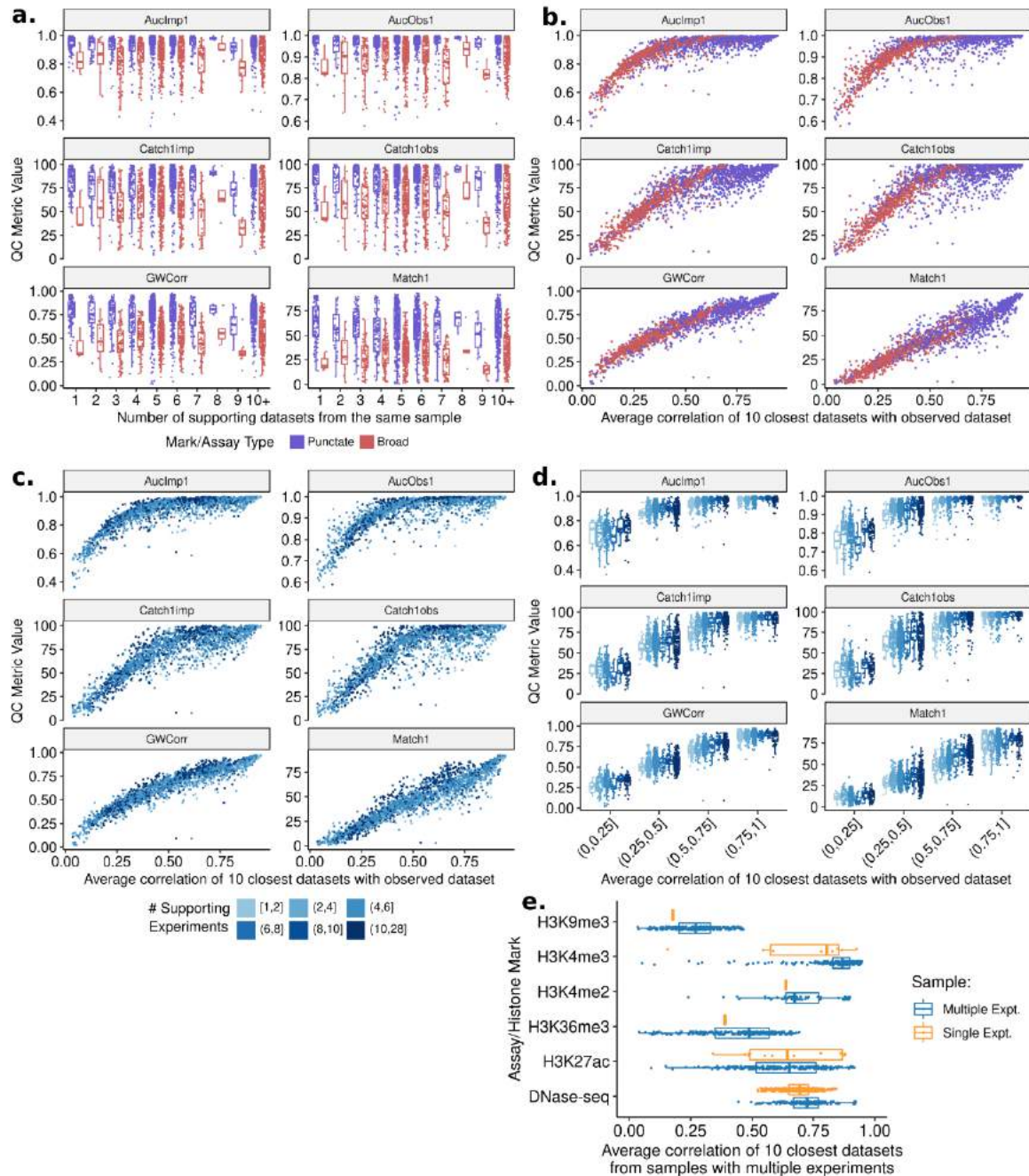
In our work on the factorization of chromatin regulators, we also illustrated how representation learning methods can be used to systematically annotate and interpret large-scale genomic datasets. We built a scalable factorization framework with variational bayes using VAEs that we used to interpret how chromatin regulators bind in complexes in context-dependent manners across the genome. Our modeling work is readily extendable and applicable in the context of any number of high-dimensional genomic datasets to identify structured and genome-localized components of epigenomic machinery. Together with the EpiMap resource, our work enables many future studies: hierarchical and multi-resolution tree-based analyses of gene regulation and GWAS; machine-learning-based gene circuitry and combinatorial regulatory motif analyses<sup>160,161</sup>.

### 1.13 Supplementary Figures

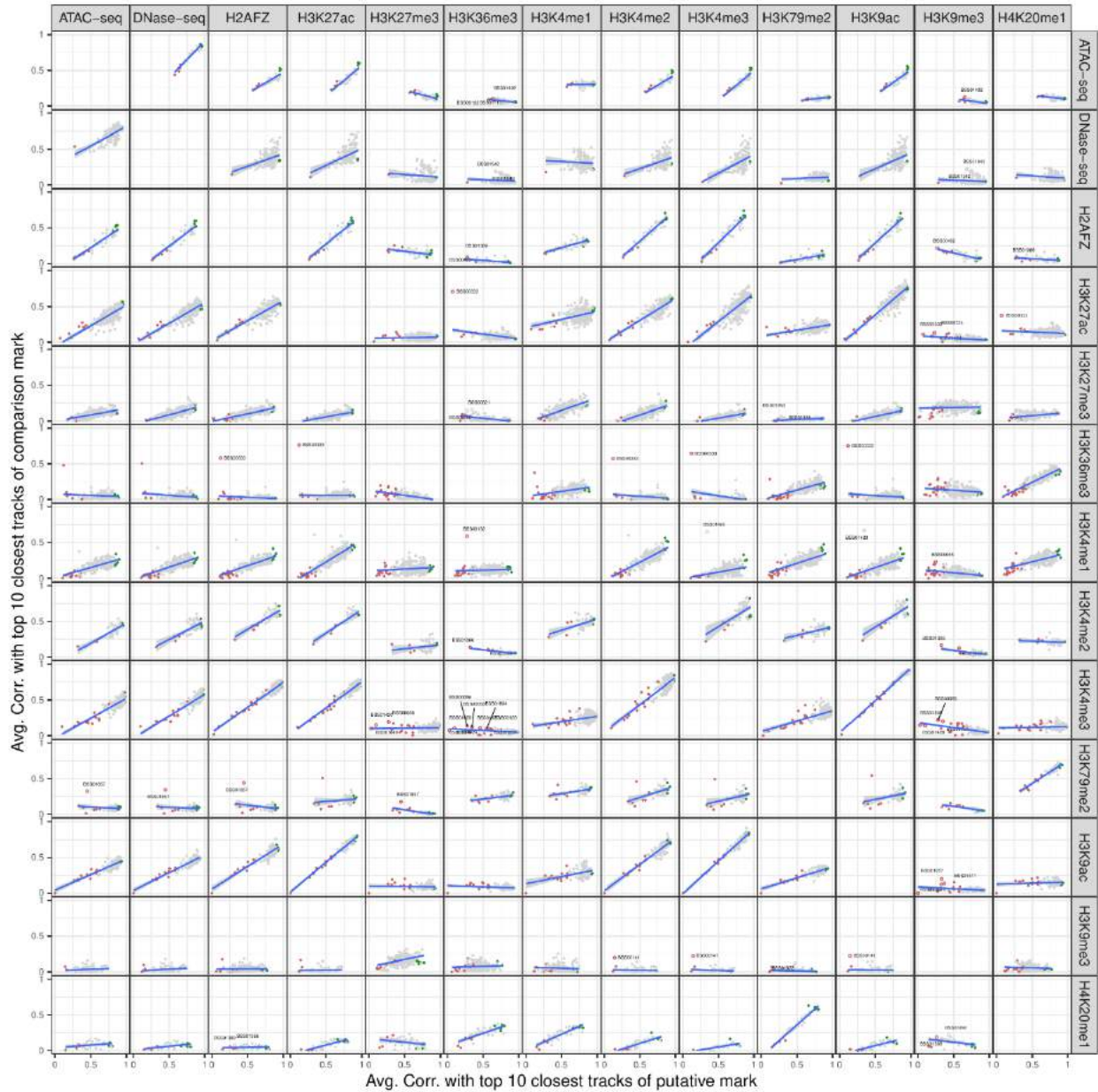


**Figure 1-S2: a.** Track agreement between observed and imputed in the Tier 1 marks, where both are available, by genome wide correlation. Top (green) and bottom (red) tracks are labeled according to automated elbow discovery in the ranked list. **b.** All imputation metrics (from original ChromImpute paper: correlation, AUC, and peak recovery) and QC metrics (reads, NSC, RSC) across all 18 imputed Tier 1, 2, and 3 assays, ordered according to correlation between imputed and observed within each assay. Bottom tracks for Tier 1 and 2 assays are labeled in red, as above. **c.** Imputation QC metrics reflect external ChIP-seq quality metrics (NSC and RSC). Normalized Strand Cross-correlation coefficient (NSC) against Relative Strand Cross-correlation coefficient (RSC) for observed tracks. Top (green) and bottom (red) tracks are labeled according to automated elbow discovery in the imputed vs. observed correlation ranked list. Poor agreement tracks are strongly clustered in the lower left quadrant of tracks failing QC. Critical values subdivide the plot for NSC (1.05) and RSC (1).

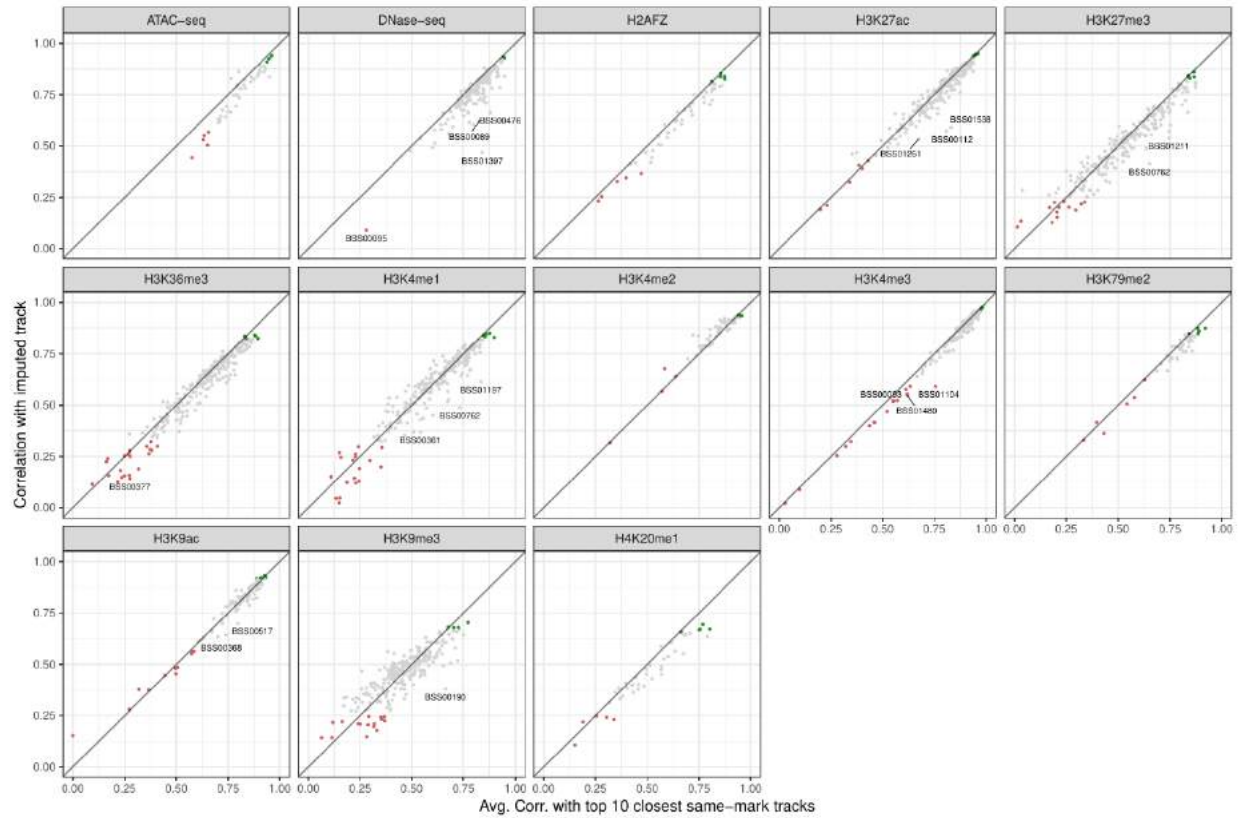




**Figure 1-S3: Imputation quality by level of sample-internal and sample-external support.** **a.** Imputation QC metrics in punctate (blue) and broad (red) marks by number of other supporting datasets from the same sample (sample-internal support). **b.** Imputation QC metrics for punctate and broad marks against the average correlation of the 10 closest datasets with the observed dataset. (sample-external support). **c-d.** Metrics against continuous (**c**) and quartiles of (**d**) sample-external support, colored by level of sample-internal support, showing that imputation metrics are dominated by the level of sample-external support. **e.** Level of sample-external support from samples with multiple experiments for samples with only one experiment (orange) or samples with multiple experiments (blue), by available mark or assay. Almost all samples with only DNase-seq have strong cross-sample support from other samples with multiple experiments (average correlation of nearest samples above 50%).

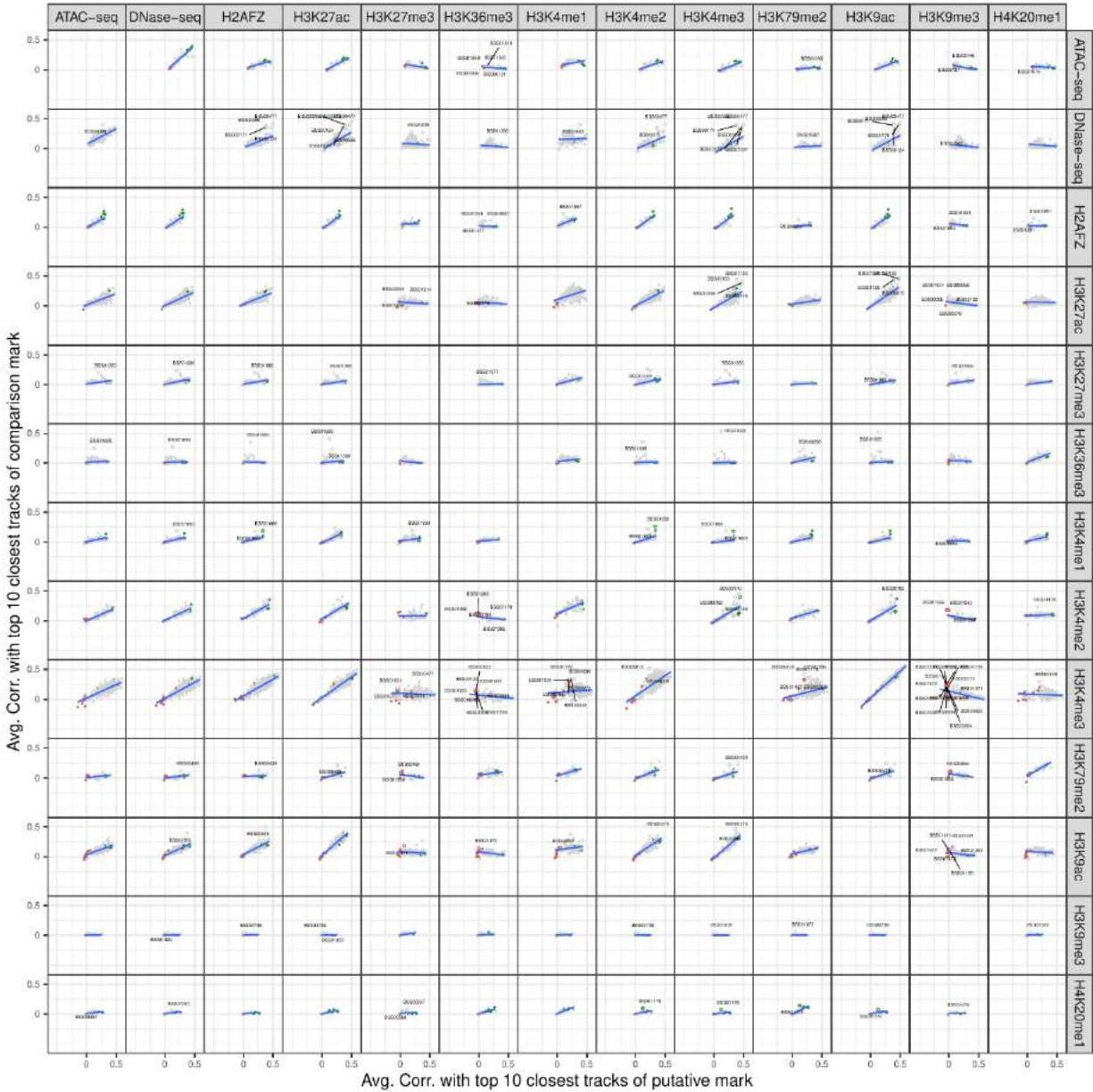


**Figure 1-S4: Imputed-observed agreement can be used to systematically flag antibody swaps.** Each panel shows the average correlation with the closest 10 tracks within the putative mark (each row) for each observed dataset against all other Tier 1 and 2 marks and assays. We used the overall mark-mark trend to flag outliers for visual inspection.

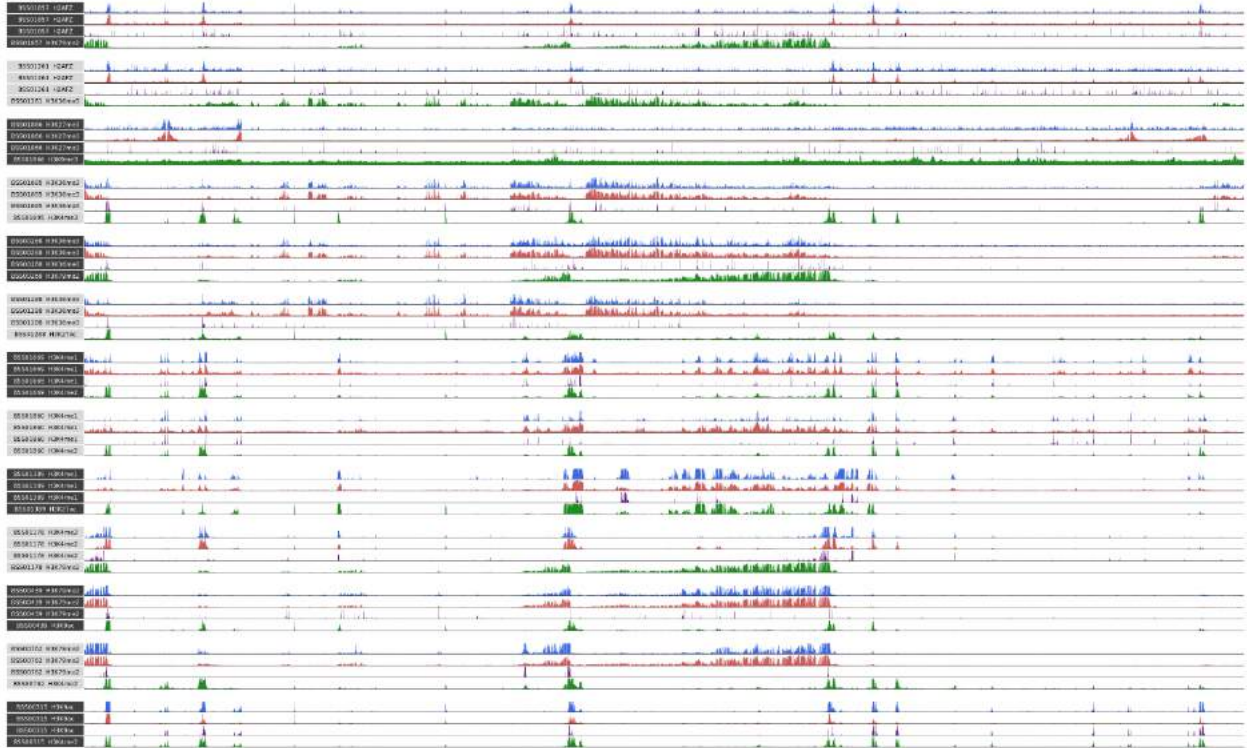


**Figure 1-S5: Imputed-observed agreement can be used to flag sample swaps.** We compare the imputed-observed correlation to the average imputed-observed correlation within the top 10 closest samples for the putative mark and flagged outliers for visual inspection (after removing antibody swaps).





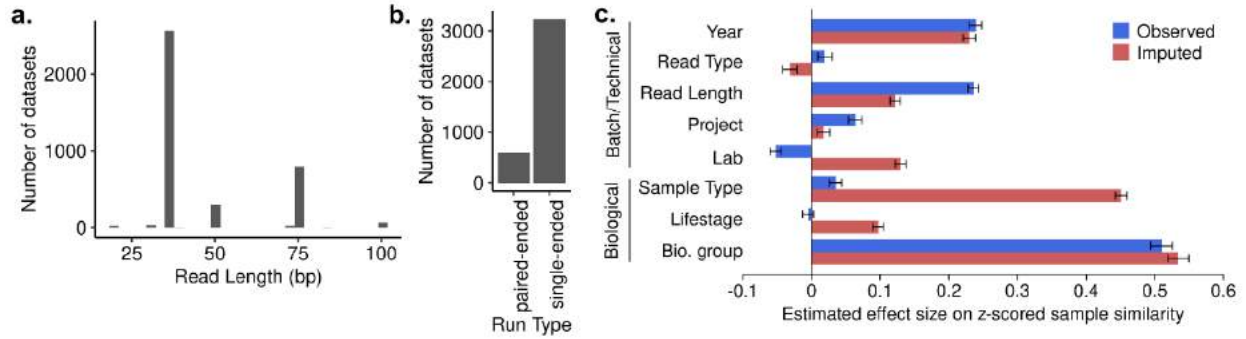
**Figure 1-S6: Imputed-observed agreement can be used to systematically flag secondary reactivities of antibodies.** Each panel shows the average correlation of the difference between a pair of observed and imputed datasets with closest 10 tracks within the putative mark (each row) for each dataset against all other Tier 1 and 2 marks and assays. As with antibody swaps, we used the overall mark-mark trend to flag outliers for visual inspection.



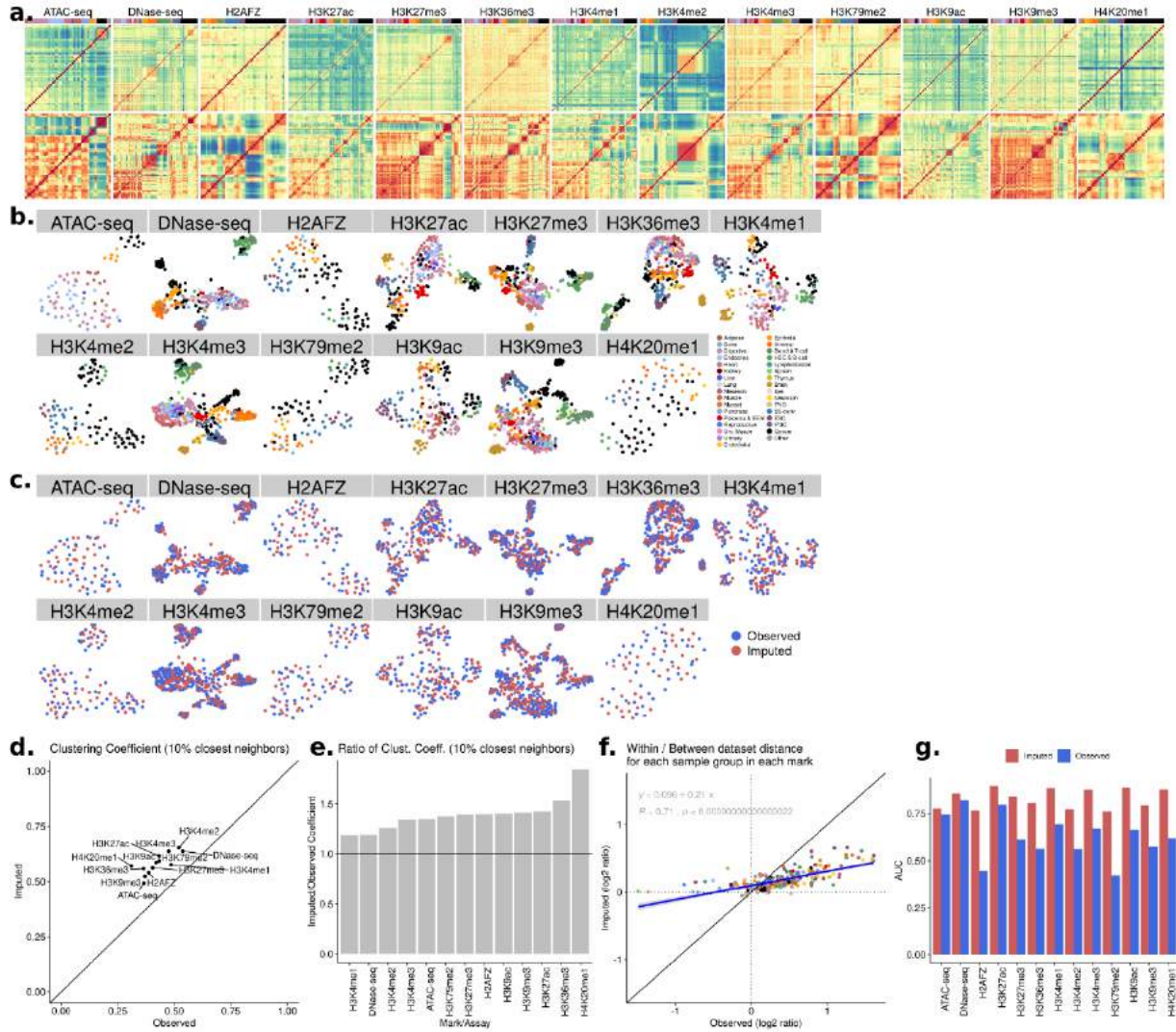
**Figure 1-S7:** We can use imputation to identify antibodies with secondary reactivities. Track sets for 13 flagged samples highlight a disagreement from secondary reactivity (or sample swapping of one experiment out of multiple replicates), where the difference between observed and imputed best correlate with an external histone mark or assay. Each track set shows the Difference (purple) from Observed (blue) to Imputed (red) and the best Match (green) to the difference (a mark in the same sample) in chromosome 1 from 1.5Mb to 2Mb (chosen to show a range of diverse elements and marks).





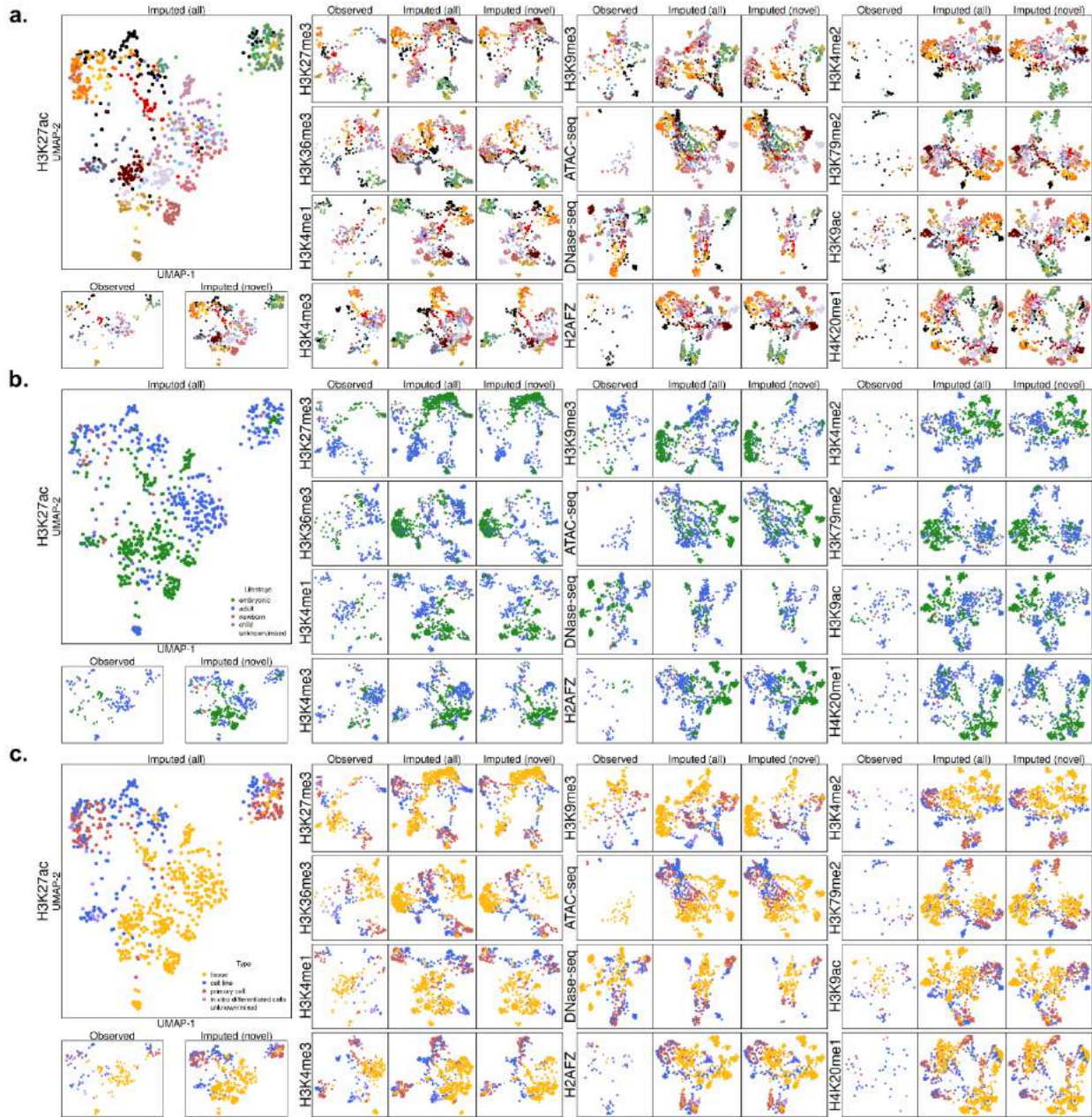


**Figure 1-S9: a.** Read length distribution in observed datasets before processing. Over two-thirds of observed datasets have read-length of 36bp or lower. **b.** Distribution of sequencing run type in observed datasets. **c.** Estimated effect sizes of technical and biological covariates on the genome-wide correlation between two imputed or two observed datasets in the same mark. We regressed indicator variables for each covariate against per-mark z-scored correlations of all pairwise combinations of either observed or imputed datasets in each mark. Imputed datasets show lower effect sizes from all technical covariates except for lab or origin and higher effect sizes for all biological covariates than observed datasets. Error bars represent two standard errors.

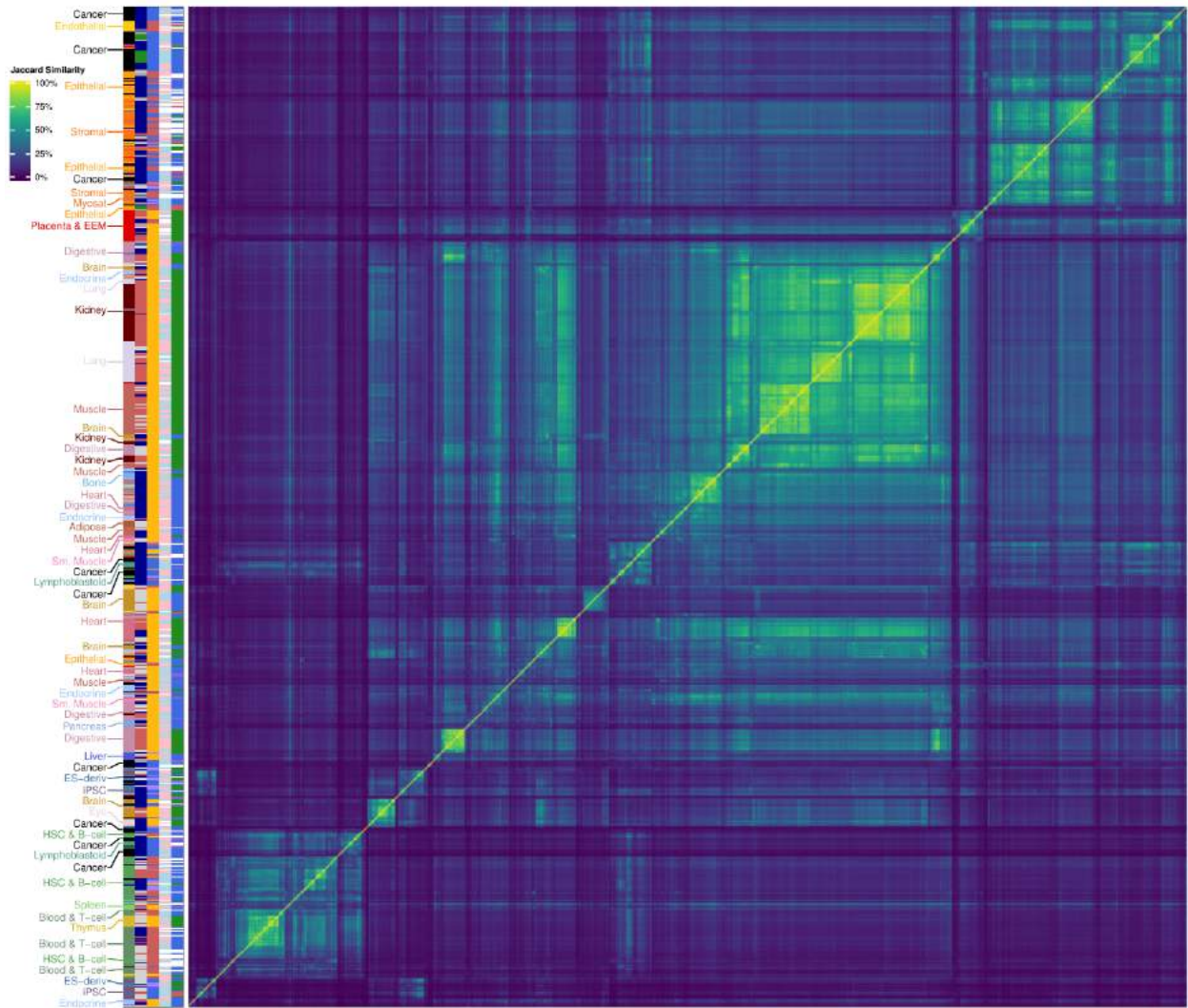


**Figure 1-S10: Data clustering and homogeneity in samples with both observed and imputed data.** **a.** Per-mark heatmap of Spearman correlation for observed (top) and imputed (bottom) data in samples with both observed and imputed data in the mark or assay. **b,c.** Joint imputed and observed per-mark dimensionality reduction (UMAP) from Spearman correlation between all imputed and observed tracks in the subset of samples with both datatypes, colored by **b.** sample group and **c.** datatype. **d.** Per-mark observed versus imputed clustering coefficients. Clustering coefficients were computed on per-mark networks for either observed or imputed datasets constructed by connecting each dataset with its top decile of closest neighbors. **e.** Ratio of imputed to observed clustering coefficients. **f.** Ratio of within-group to between-group distance in observed versus imputed data for each sample group (colors) and each mark. Imputed and observed data show concordant ratios (Pearson  $R=0.72$ , Spearman  $\rho=0.71$ ) **g.** AUC for predicting whether a pair of marks is in the same sample group from their pairwise distance. Imputed data outperforms observed data across all marks and assays and achieves AUC > 0.75 for most marks.



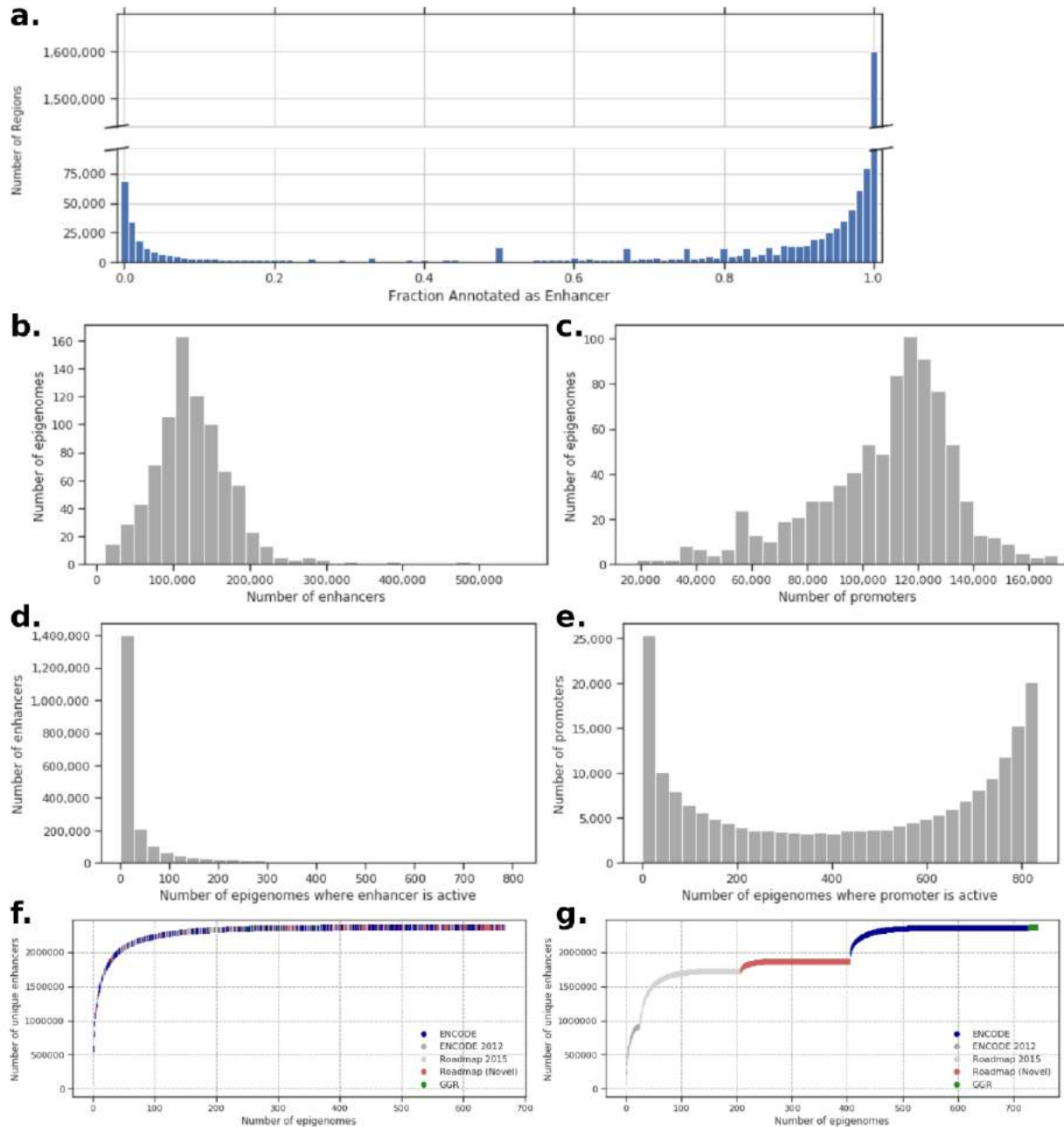


**Figure 1-S11: a-c.** Joint UMAP embeddings of observed and imputed data within each Tier 1 and 2 mark/assay. Separately, observed, imputed samples, and all novel imputed tracks are plotted for each mark/assay and colored according to **a.** tissue group, **b.** biological lifestage, or **c.** sample type. Imputed UMAP highlights differences in cell types: H3K27ac clusters hematopoietic cells and tissues (green in panel **a**) closely, reflecting lineage, whereas H3K27me3 clusters iPSC, ESCs, and derived cells (purple/blue in panel **a**), reflecting differentiation stage. UMAP embedding was calculated from spearman correlation of tracks within regions marked by mark or assay relevant states within the Roadmap compendium.



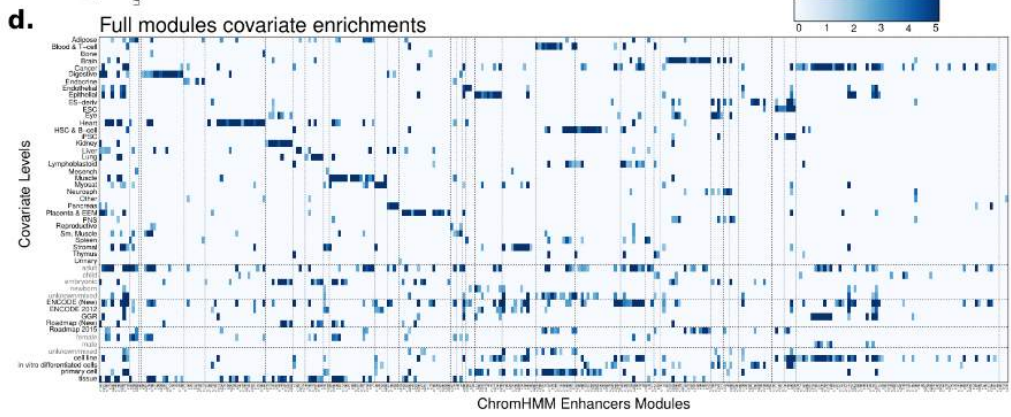
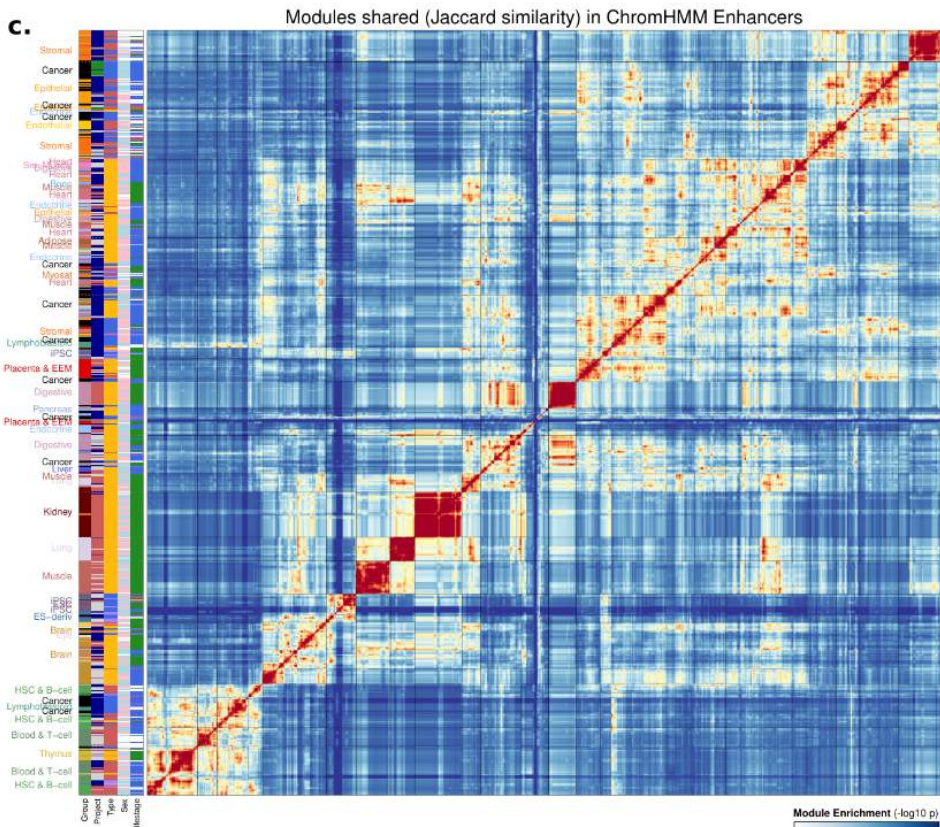
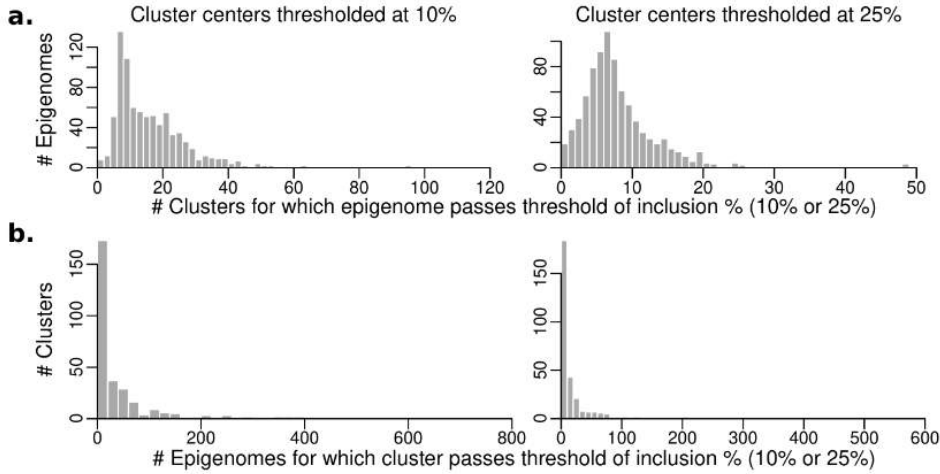
**Figure 1-S12:** Jaccard similarity matrix (heatmap) across 833 epigenomes (metadata on left) from binarized enhancer activity matrix (2.1M enhancers by 833 epigenomes). Similarity matrix clustered by complete-linkage clustering. Consecutive blocks of at least six samples from the same group are labeled on the left.



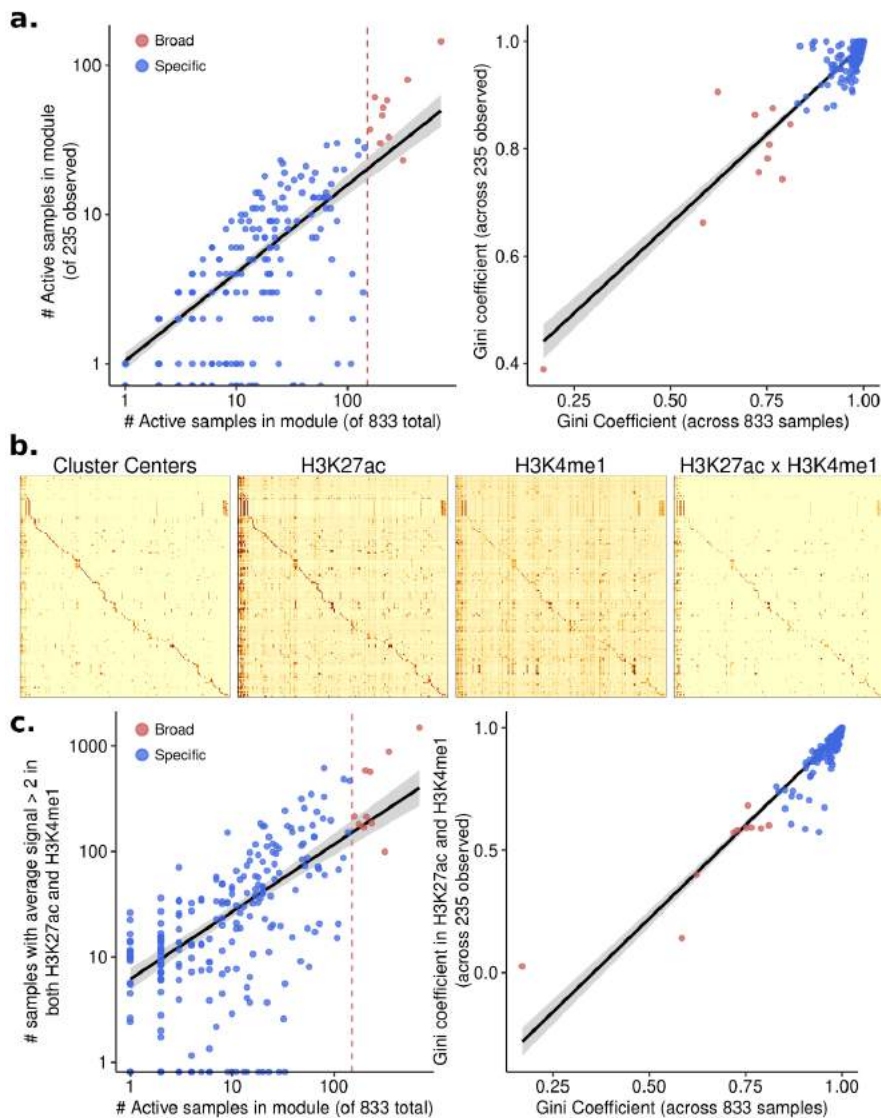


**Figure 1-S13: a.** Number of DHS sites annotated as enhancers instead of promoters across 833 epigenomes. Most regions are either labeled enhancer (at least 75% of occurrences are enhancers) or promoter (at least 75% of occurrences are promoters) across all of their active occurrences. Using these cutoffs, we labeled 2,069,090 enhancers, 204,104 promoters, and 122,358 dyadic elements (neither specifically promoter or enhancer). **b,c.** Histogram of number of active enhancers (**b**) and promoters (**c**) per sample (row-margins of binary activity matrix). **d,e.** Histogram of the number of samples for which each enhancer (**d**) or promoter (**e**) is an active enhancer (column-margins of binary activity matrix) **f.** Rarefaction curve for enhancer recovery for all enhancer and dyadic elements (2.3M total) across 833 samples (points, colored by project of origin). Curve was created by iteratively adding the sample contributing the most novel active elements until all 2.3M elements were accounted for. **g.** Rarefaction curve for enhancer recovery, shown in order of project completion, for all enhancer and dyadic elements

(2.3M total) across 833 samples (points, colored by project of origin). Curve was created by iteratively adding the sample contributing the most novel active elements until all 2.3M elements were accounted for. Samples were considered in order of project publication/completion, only taking samples from the next project when the current project did not contribute any more enhancers.



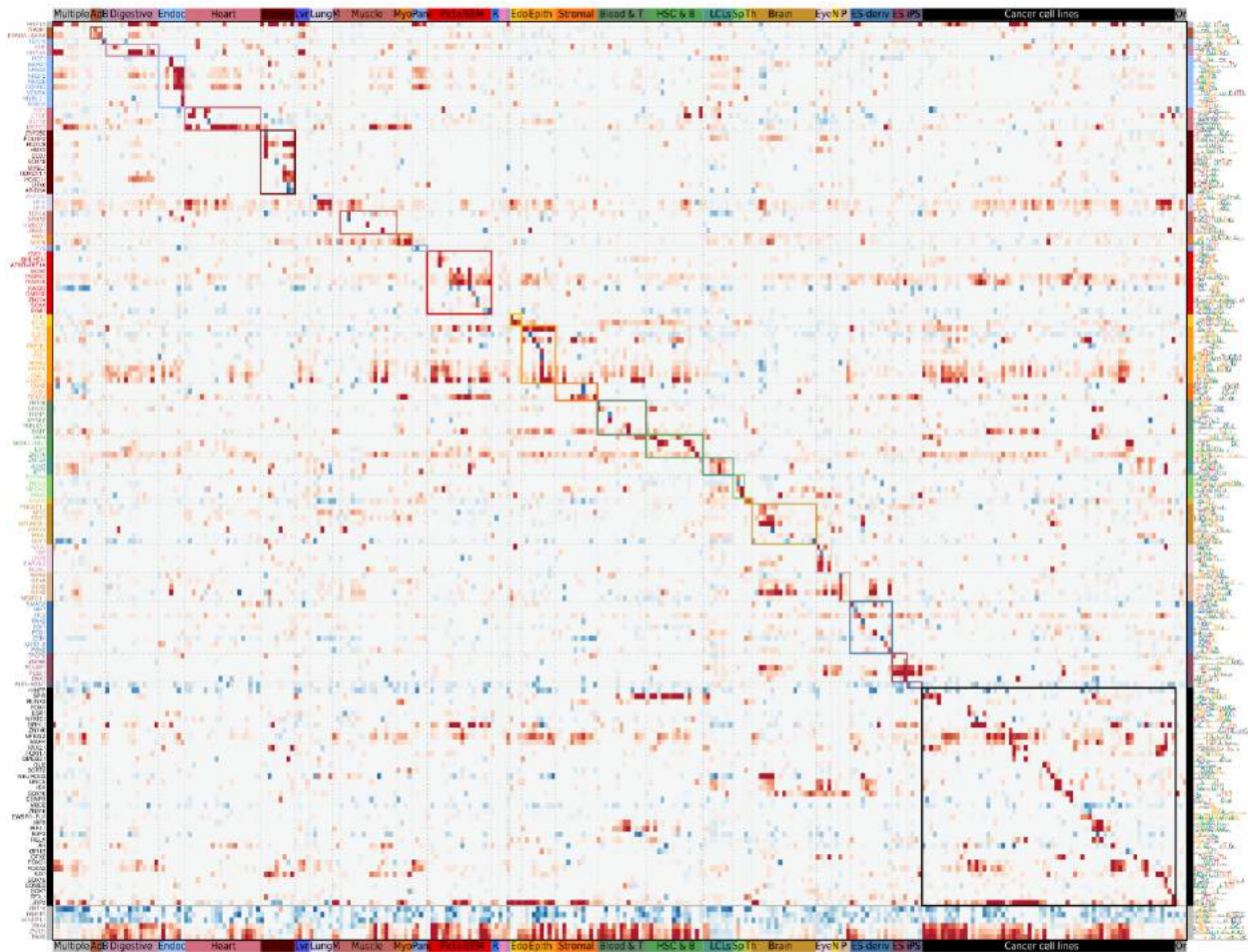
**Figure 1-S14: a,b.** Module inclusion distributions according to two different module inclusion cutoffs of 10% and 25% (columns). **a.** Number of modules (N=300) for each of 833 epigenomes with module inclusion of at least 10% or 25% (top). **b.** Number of epigenomes (N=833) for each of 300 modules at the same inclusion cutoffs. **c.** Sample similarity (heatmap) across 833 epigenomes by their number of shared modules. (jaccard similarity, intersection over union of modules of each pair of epigenomes). Similarity matrix clustered by Ward's method. Consecutive blocks of at least six samples from the same group are labeled on the left. **d.** Full module enrichments across main metadata facets of group, project, type, sex, and lifestage, for which the reduced version is shown as a panel in **Figure 2a**. Significance assessed according to hypergeometric test on the module centers matrix with module inclusion of at least 25%.



**Figure 1-S15: Observed data validation of broad and specific modules. a.** Validation with 235 observed H3K27ac datasets. Number of samples (left) and Gini coefficient (right) per module in the full data versus the observed-only H3K27ac data. Broad modules are defined as modules with over 150 active samples in the full data (dashed line) and labeled in red. Despite the sample composition differences in observed vs. full data, the fraction of active samples is



quite consistent between the datasets (left,  $R^2=0.802$ , right  $R^2=0.82$ ). For the 10 broad modules (red) from the full data, all were active in at least 9.7% of observed samples and 24% on average, seven were the most broadly active in terms of active samples in observed, and nine were in the bottom ten modules by Gini coefficient (least unequal, or most broadly expressed). **b.** Side-by-side heatmaps of module centers, average observed H3K27ac, average observed H3K4me1, and the product of H3K27ac and H3K4me1 in 205 datasets with observed data for both marks. **c.** Number of samples (left) and Gini coefficient (right) as in (a), for the product of observed H3K27ac and H3K4me1 in 205 datasets with observed data for both marks.



**Figure 1-S17:** Full motif enrichments for 160 motif families with enrichments with  $\log_2FC \geq 1$  on modules (**Figure 2c**).



## 2. Dissecting disease-associated common variation and GWAS loci using a functional epigenomic atlas

Parts of this chapter have been presented/published as:

### Regulatory genomic circuitry of human disease loci by integrative epigenomics

(*Nature*, 2021)

Carles A. Boix<sup>1-3</sup>, Benjamin T. James<sup>1,2</sup>, Yongjin P. Park<sup>1,2,4</sup>, Wouter Meuleman<sup>5</sup>, Manolis Kellis<sup>1,2\*</sup>

1. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA, USA. 2. Broad Institute of MIT and Harvard, Cambridge, MA, USA. 3. Computational and Systems Biology Program, Massachusetts Institute of Technology, Cambridge, MA, USA. 4. Department of Pathology and Laboratory Medicine, Statistics, University of British Columbia, Vancouver, BC, Canada. 5. Altius Institute for Biomedical Sciences, United States. Correspondence to: e-mail: [manoli@mit.edu](mailto:manoli@mit.edu)

### Single-cell dissection of live human hearts in ischemic heart disease and heart failure reveals cell-type-specific driver genes and pathways

(*in review*)

Suvi Linna-Kuosmanen<sup>1,2,3,\*</sup>, Eloi Schmauch<sup>1,2,3,\*</sup>, Kyriakitsa Galani<sup>1,2</sup>, Carles A. Boix<sup>1,2</sup>, Lei Hou<sup>1,2</sup>, Tiit Örd<sup>3</sup>, Anu Toropainen<sup>3</sup>, Lindsey K. Stolze<sup>4</sup>, Elamaran Meibalan<sup>2,10</sup>, Julio C. Mantero<sup>1,2</sup>, Ashley Renfro<sup>1,2</sup>, Johannes Ojanen<sup>3</sup>, Leandro Z. Agudelo<sup>1,2</sup>, Maija Hollmen<sup>5</sup>, Juhon Jalkanen<sup>6</sup>, Jarmo Gunn<sup>6,7</sup>, Pasi Tavi<sup>3</sup>, Casey E. Romanoski<sup>4</sup>, Calum A. MacRae<sup>2,8,9,10</sup>, Minna U. Kaikkonen<sup>3</sup>, Guillermo Garcia-Cardena<sup>2,9,11</sup>, Tuomas Kiviniemi<sup>6,7,8</sup>, and Manolis Kellis<sup>1,2\*</sup>

<sup>1</sup> Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA, 02139, USA. <sup>2</sup> Broad Institute of MIT and Harvard, Cambridge, MA, 02142, USA. <sup>3</sup> A. I. Virtanen Institute for Molecular Sciences, University of Eastern Finland, 70211 Kuopio, Finland. <sup>4</sup> Department of Cellular and Molecular Medicine, The College of Medicine, The University of Arizona; Tucson, AZ 85721, USA. <sup>5</sup> Medicity Research Laboratories, University of Turku, 20500 Turku, Finland. <sup>6</sup> Heart Center, Turku University Hospital, 20521 Turku, Finland. <sup>7</sup> Department of Medicine, University of Turku, 20500 Turku, Finland. <sup>8</sup> Cardiovascular Medicine and Network Medicine Division, Brigham and Women's Hospital, Harvard Medical School, Boston, MA 02115, USA. <sup>9</sup> Harvard Stem Cell Institute, Cambridge, MA 02138, USA. <sup>10</sup> Department of Medicine, Harvard Medical School, Boston, MA 02115, USA. <sup>11</sup> Center for Excellence in Vascular Biology, Department of Pathology, Brigham and Women's Hospital and Harvard Medical School, Boston, MA 02115, USA. \*These authors contributed equally: Suvi Linna-Kuosmanen, Eloi Schmauch. ✉ Corresponding authors, email: [suvil@mit.edu](mailto:suvil@mit.edu), [manoli@mit.edu](mailto:manoli@mit.edu)

## 2.1 Introduction

Annotating the molecular basis of human disease remains an unsolved challenge, as 93% of disease loci are non-coding, and gene-regulatory annotations highly incomplete<sup>46,47</sup>. In Chapter 1, we presented EpiMap, a compendium of 10,000 epigenomic maps across 800 samples, which we use to define chromatin states, high-resolution enhancers, enhancer modules, upstream regulators, and downstream target genes.

Here, we use this resource to annotate 30,000 genetic loci associated with 540 traits<sup>5,52,53</sup>, predicting trait-relevant tissues, putative causal nucleotide variants in enriched-tissue enhancers, and candidate tissue-specific target genes for each. We partition multifactorial traits into tissue-specific contributing factors with distinct functional enrichments and disease-comorbidity patterns, and reveal both single-factor monotropic and multi-factor pleiotropic loci. Top-scoring loci frequently have multiple predicted driver variants, converging through multiple common-target enhancers, multiple common-tissue genes, or multiple genes/tissues, indicating extensive pleiotropy. We detail the dissection of coronary artery disease (CAD) GWAS loci and provide additional examples of disease-implicated genes linked through GWAS locus prioritization to schizophrenia (CACNA1A) and breast cancer (NTN4).

We further applied our EpiMap resource in conjunction with single-cell transcriptomic data from live human hearts to dissect a GWAS risk locus implicated in coronary artery disease and ischemic heart failure. The GWAS locus showed pleiotropic activity, implicating two genes (JCAD and SVIL) in different tissues (artery endothelial cells and smooth muscle cells) through multiple prioritized variants and specific transcription factors. Overall, our results and in particular our case studies in specific disease demonstrate the importance of dense, rich, high-resolution epigenomic annotations for complex trait dissection.

## 2.2 Interpreting GWAS loci

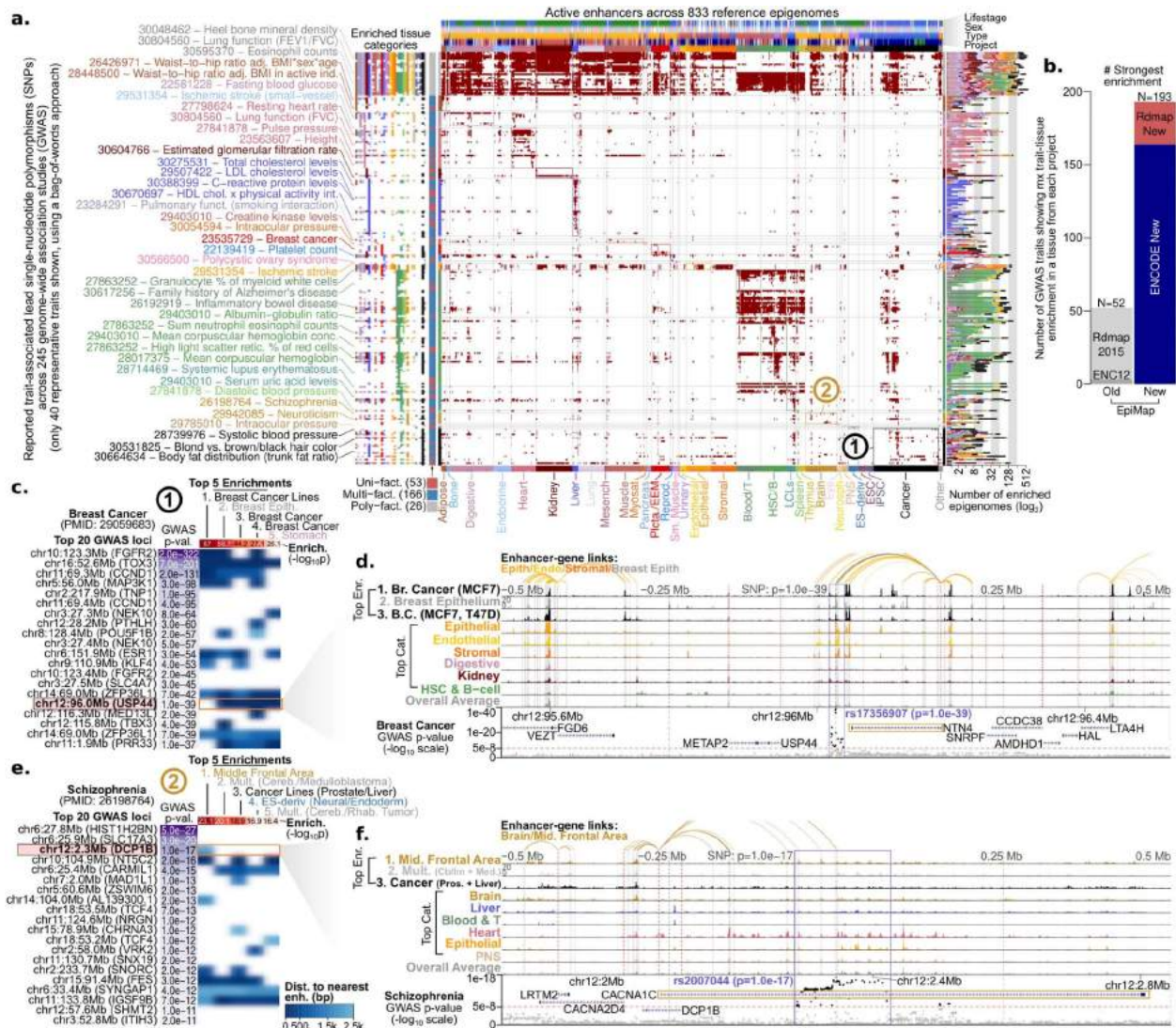
We used our 2.1M enhancer annotations and their tissue specificity to interpret genetic variants associated with complex traits<sup>5,52,53</sup>. We compiled a compendium of 803 well-powered genome-wide association studies<sup>162</sup> (GWAS) with  $\geq 10$  significant loci and  $\geq 10,000$  cases (15% of 5454 GWAS publications) that capture  $\geq 70,000$  GWAS loci (63% of NHGRI-EBI catalog<sup>47</sup>).

We found 17,658 significant enhancer-level trait-tissue enrichments, enabling fine-mapping of candidate driver SNPs in  $\geq 27,000$  loci (39%) from 245 traits in tissue-enriched enhancers (adjusted p-value  $< 1\%$ ) (**Fig. 2-1a, Supplementary Fig. 2-S1-2**). New biosamples captured the strongest GWAS enrichment in 79% of cases (193 of 245) and the only significant enrichment in 24% of cases (N=58), and our annotations captured 2.5-fold more GWAS studies than DNase alone (245 vs. 97) (**Fig. 2-1b, 2-2a-d**). Using SNP-level tests, we found 15,866 trait-tissue enrichments from 256 traits at 2% FDR. The SNP-level analysis represents a  $\sim 4.5$ -fold increase from the 54 traits enriched in H3K27ac and 58 traits enriched in H3K4me1 reported by the Roadmap Epigenomics project, which used a 2% FDR and a more stringent SNP pruning methodology on an older GWAS catalog.

To capture common enrichments of similar biosamples, we also calculated enhancer-level trait enrichments for enhancer modules, resulting in  $\sim 3$ -fold more enriched traits (717 vs. 245) but 38% fewer SNPs in enriched annotations (SNP-level  $\sim 2.5$ -fold reduced at 104 vs. 256) (**Supplementary Fig. 2-S3**). Instead, using the hierarchical enhancer-sharing tree (**Fig. 2-3c**) to reveal the appropriate tissue resolution of GWAS enrichment in subtree-vs-parent enhancer-GWAS comparisons, revealed 540 nominally enriched traits and 32,532 SNPs in enriched annotations (**Fig. 2-2e, Supplementary Fig. 2-S1**).

Our epigenomic enrichments and enhancer-gene links yielded new biological insights on

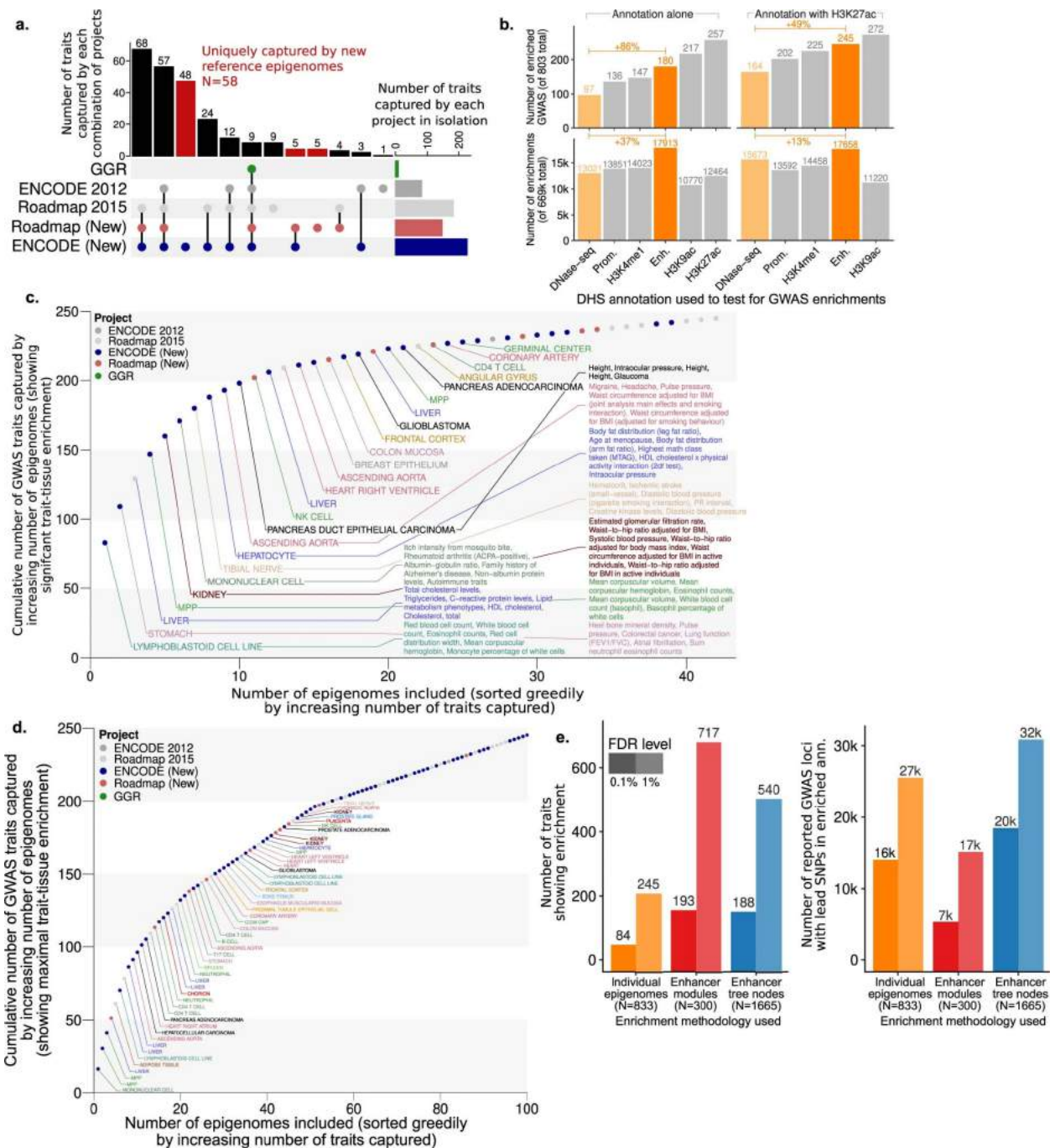
disease loci, with many compelling examples. For breast cancer GWAS<sup>163</sup>, enriched in epithelial/cancer biosamples (Fig. 2-1c), the highly-localized rs17356907 genetic signal ( $p=10^{-39}$ , rank #16) localizes precisely in a narrow epithelial/cancer enhancer nearest to *USP44* but linked instead to tumorigenesis-/angiogenesis-implicated *NTN4* (Fig. 2-1d). For schizophrenia GWAS<sup>164</sup>, maximally-enriched in mid-frontal cortex (Fig. 2-1e), the diffuse rs2007044 genetic signal ( $p=10^{-17}$ , rank #3) overlaps a broad set of enhancers nearest to the *DCP1B* promoter, all of which linked to neuropsychiatric-disorder-implicated calcium-channel *CACNA1C*, suggesting multiple causal variants contributing jointly<sup>165</sup> (Fig. 2-1f). We provide an interactive website for exploring >30,000 additional loci across >500 traits.



**Figure 2-1. GWAS tissue-prioritization.** **a.** Trait-tissue enrichment (center, heatmap) between reported lead single-nucleotide polymorphisms (SNPs) from 245 genome-wide association studies (rows) and accessible active enhancers across 833 biosamples (columns) (adjusted nominal  $p$ -value  $< 1e-3$ ). Enriched tissue groups (left) and number of enriched biosamples (right) shown for each trait. Only 40 representative traits labeled, using a bag-of-words approach (full list of traits in **Supplementary Fig. 2-S2**). Traits colored by sample with maximal trait-tissue enrichment. **b.** Contribution of each project to the maximum GWAS trait-tissue enrichment for the 245 traits with significant enrichments. **c.** Enhancer overlaps with top 20 lead SNPs for

breast cancer for top 5 enrichments on the enhancer tree. **d.** 1Mb locus centered on breast cancer lead SNP rs17356907 (chr12:96027759,  $p=1.0e-39$ ) showing H3K27ac signal (middle panel) in enhancer DHSs for the top three enriched subtrees in the enhancer tree, six selected tissue categories, and overall average signal. TSSs indicated by red dashed lines. Top panel shows enhancer gene-links for SNP-proximal enhancers for the top enrichments and across the locus for epithelial, endothelial, and stromal cell biosamples. Linked enhancers are highlighted in grey. Bottom panel shows gene models and breast cancer GWAS summary statistics, with SNPs below  $p=5e-8$  in grey. **e.** Enhancer overlaps with top 20 lead SNPs for schizophrenia for top 5 enrichments on the enhancer tree. **f.** 1Mb locus centered on schizophrenia lead SNP rs2007044 (chr12:2344960,  $p=1.0e-17$ ), with links (top), H3K27ac signal (middle), and summary statistics as in panel **d**, with links for brain biosamples.





**Figure 2-2. GWAS prioritization statistics.** **a.** Number of traits (y-axis) with significant GWAS trait-tissue enrichments for each combination (column) of projects (rows). **b.** Comparison of GWAS enrichments found (top) and number of significant trait-tissue pairs SNPs in significantly-enriched annotations (bottom) using different annotations within DHSs either without (left) or with H3K27ac signal (right). DNase-seq signal alone enriches for far fewer GWAS than enhancer states alone or with H3K27ac. **c.** Increase in the cumulative number of GWAS traits (y-axis) with significant trait-tissue enrichments with increasing numbers of biosamples (x-axis), ordered to maximize the number of novel trait annotations captured with each new biosample. Top 25 samples labeled and colored by tissue group, with top 6 GWAS



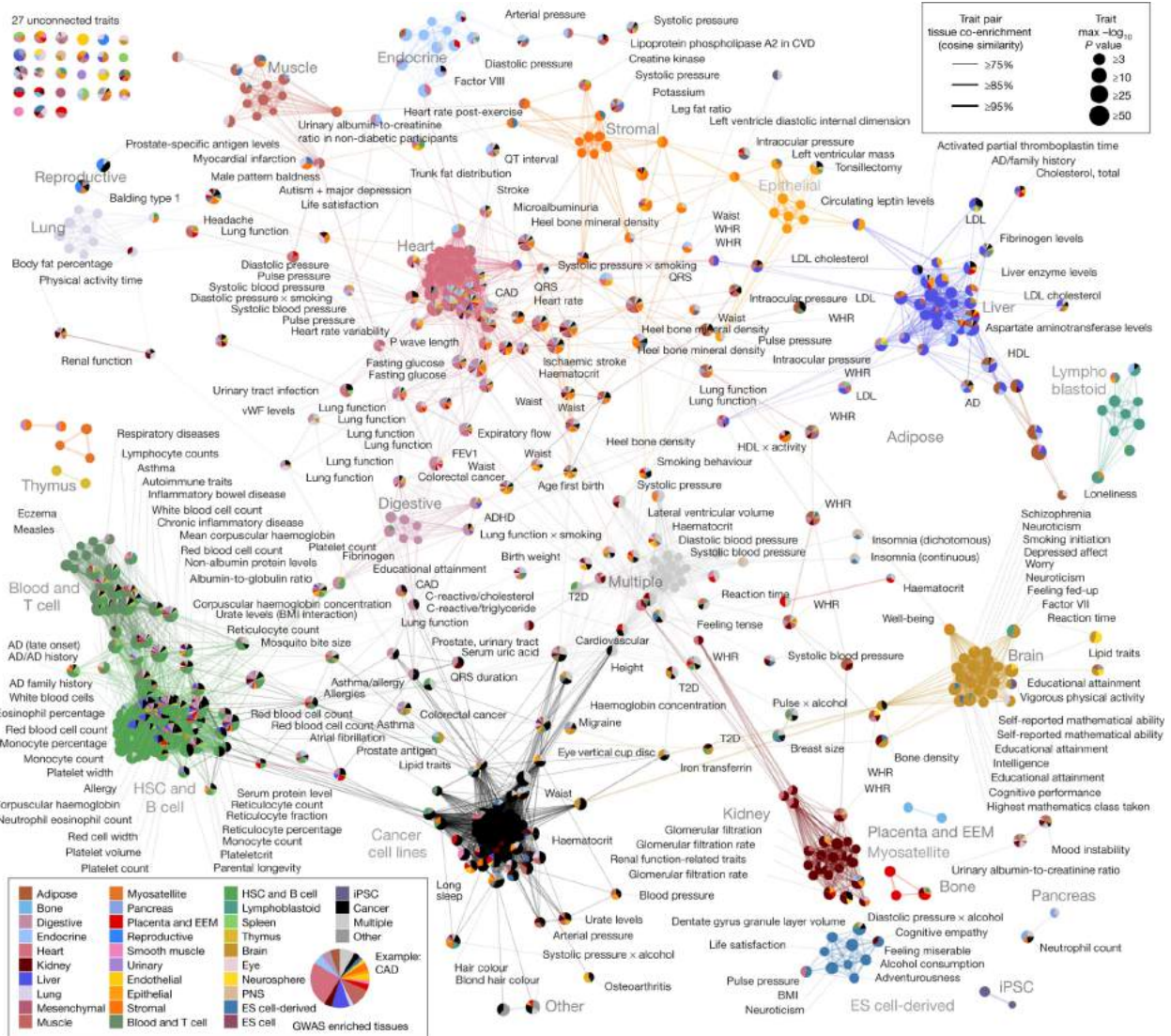
traits shown for the first 8 samples. Points colored by project. All 245 traits are captured after inclusion of 42 samples. **d.** Increase in the cumulative number of GWAS traits (y-axis) with maximal trait-tissue enrichments with increasing numbers of biosamples (x-axis). All 245 traits are captured after inclusion of 100 samples. **e.** Comparison of GWAS enrichments found (y-axis, left) and number of lead SNPs in significantly-enriched annotations (y-axis, right) using different methodologies (x-axis) for two adjusted nominal p-value cutoffs (shades).

### 2.3 GWAS-tissue co-enrichments

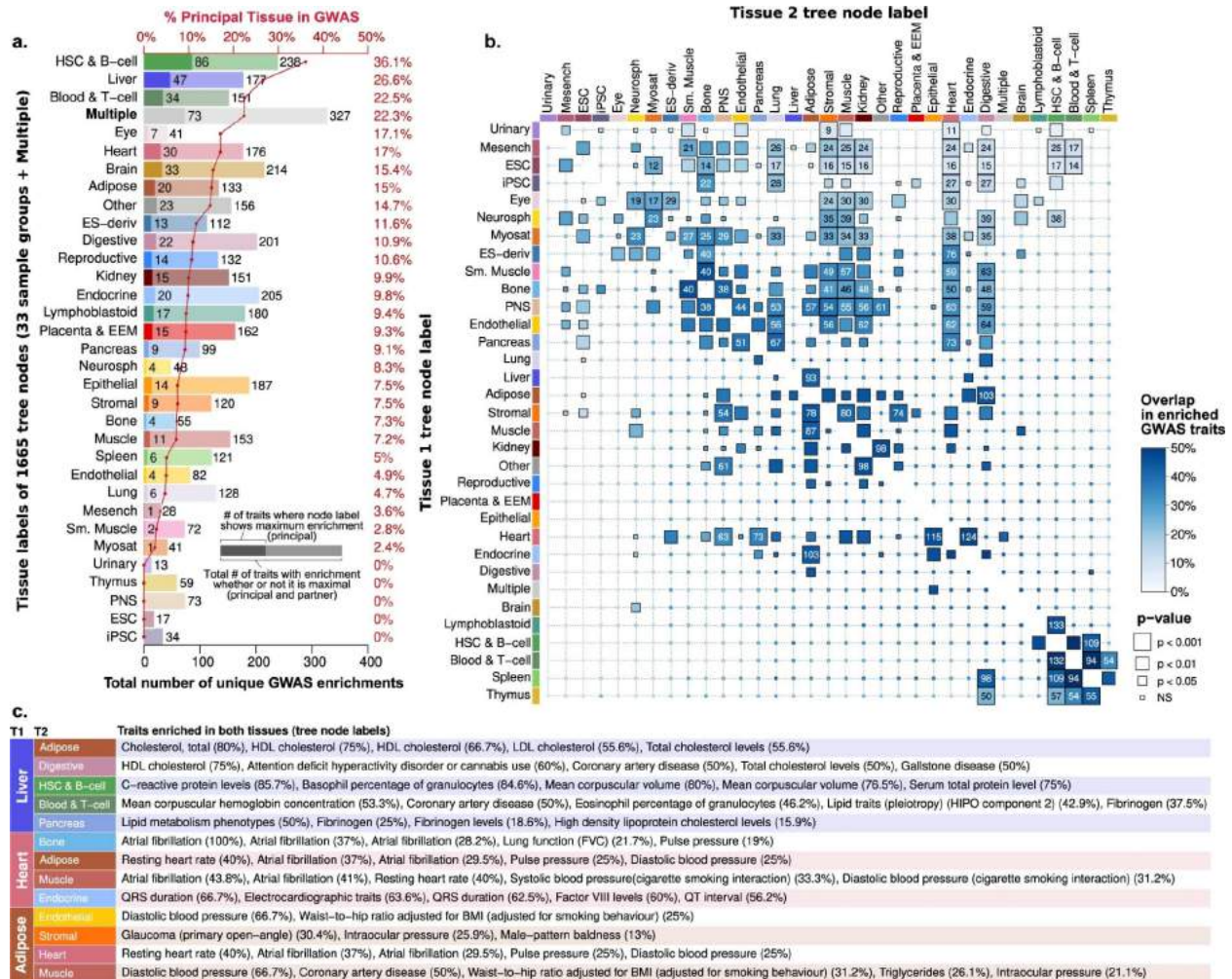
We then studied trait-tissue trait-trait and tissue-tissue epigenome-GWAS co-enrichment patterns to gain insights into their complex interactions. First, we used the number of distinct tissue categories enriched in each trait (**Fig. 2-1a**, left; **Supplementary Data S1**) to distinguish: 56 ‘unifactorial’ traits (22%) with most enriched nodes in only one tissue group (eg. QT interval in heart, educational attainment in brain, hypothyroidism in immune cells); 192 ‘multifactorial’ traits (79%) enriched in 5 tissue categories on average (eg. Alzheimer’s disease in immune and brain<sup>63</sup>; waist-to-hip ratio<sup>166</sup> in adipose, muscle, kidney, and digestive tissues); of which 26 ‘polyfactorial’ traits (11%) enriched in 14 tissue categories on average (including coronary artery disease (CAD)<sup>167</sup> in 19 tissue groups, including liver, heart, adipose, muscle, and endocrine samples).

Second, we used trait co-enrichment patterns in the same tissues to cluster GWAS traits with similar properties. The resulting network (**Fig. 2-3**) showed a small number of densely-connected communities of primarily unifactorial traits (eg. cognitive/psychiatric traits in brain/neurons; heart-beat intervals in heart; cholesterol in liver; filtration in kidney; immune traits in T-cells; and blood cell counts in hematopoietic cells), with multifactorial connectors between them (eg. CAD between heart/endocrine/liver; HDL/triglycerides between liver/adipose; lung-function between lung/heart/digestive; blood-pressure between heart and endocrine/endothelial/liver; cell-count between liver and digestive), and polyfactorial traits between these communities (eg. waist-hip-ratio, heel-bone mineral density) (Supplementary Fig. 2-S4-22). Many biologically-meaningful similarities in this epigenomic-co-enrichment-based network are missed by a genetic-overlap-based network (934 edges;  $\geq 5\%$  common loci; 10kb resolution), which only captures 5% of epigenomic-co-enrichment edges (283 of 5547) (Supplementary Fig. 2-S6-25).

Third, we used co-enrichment properties of pairs of tissues in the same traits to distinguish “principal” tissues (e.g. immune cells, liver, heart, brain, adipose) that showed consistently-higher enrichments vs. “partner” tissues (e.g. digestive, lung, muscle, epithelial) for the same GWAS traits, suggesting they play driver vs. auxiliary roles (**Fig. 2-4a**). Specific principal-partner tissue pairs co-occurred more frequently than expected (**Fig. 2-4b**), and revealed biologically-meaningful traits where they likely co-act (**Fig. 2-4c**), including: liver with adipose (for cholesterol traits), with digestive (for gallstone) and with blood cells (for serum protein levels); and adipose with endothelial (for waist-to-hip ratio), with heart (for atrial fibrillation), and with muscle (for blood pressure).



**Figure 2-3. Trait-trait network.** Network across 538 traits (by per-node adjusted nominal p-value) by similarity of epigenetic enrichments (cosine sim.  $\geq 0.75$ ), laid out using the Fruchterman-Reingold algorithm. Traits (nodes) are colored by contributing groups (pie chart by fraction of  $-\log_{10} p$ , size by maximal  $-\log_{10} p$ ) and interactions (edges) by the group with maximal dot product of enrichments between two traits. Redundant node names indicate different GWAS studies (full names for non-singleton nodes in **Supplementary Fig. 2-S5**).



**Figure 2-4. Tissue-tissue GWAS relationships.** Principal and partner tissue enrichments. **a.** For each tree node label (rows), the number of GWAS traits (black x-axis, bottom) showing maximum enrichment in that tree node (dark bars, principal tissue) or any enrichment in that tree node (light bars, partner tissue), and the percentage of tissue-enriched traits for which the tissue shows the maximal enrichment (red x-axis, top) across 540 traits. **b.** Overlap in enriched GWAS traits between pairs of tissues with maximal enrichment in the trait (principal tissue, rows) and lower enrichment in the same trait (partner tissue, columns), using tree node labels. **c.** Top traits in significant interactions for selected tissue pairs (Liver, Endocrine, Muscle, Heart, Adipose, PNS). For each pair of co-enriched tissue groups we reported the top 5 GWAS by their percent of significant enrichments coming from either group.

## 2.4 Dissecting multifactorial traits

We next used our epigenomic annotations to partition multifactorial traits into tissue-specific components, by studying functional and disease enrichments for distinct subsets of enhancer-overlapping SNPs in each enriched overlapping tissue (**Fig. 2-5a-d, 2-6, Supplementary Fig. 2-S9**).

For example, the 339 CAD-associated SNPs lying in enriched-tissue enhancers partitioned into: 195 heart-enhancer SNPs enriched in artery, cardiac, and vessel morphogenesis; 171

endocrine-enhancer SNPs in lipid homeostasis; 169 liver-enhancer SNPs in cholesterol/lipid metabolism and transport; 122 adipose-enhancer SNPs in axon guidance and focal adhesion, consistent with adipose tissue innervation processes; and 112 ES-derived-muscle-overlapping SNPs, enriched in septum morphogenesis, chamber and aorta development.

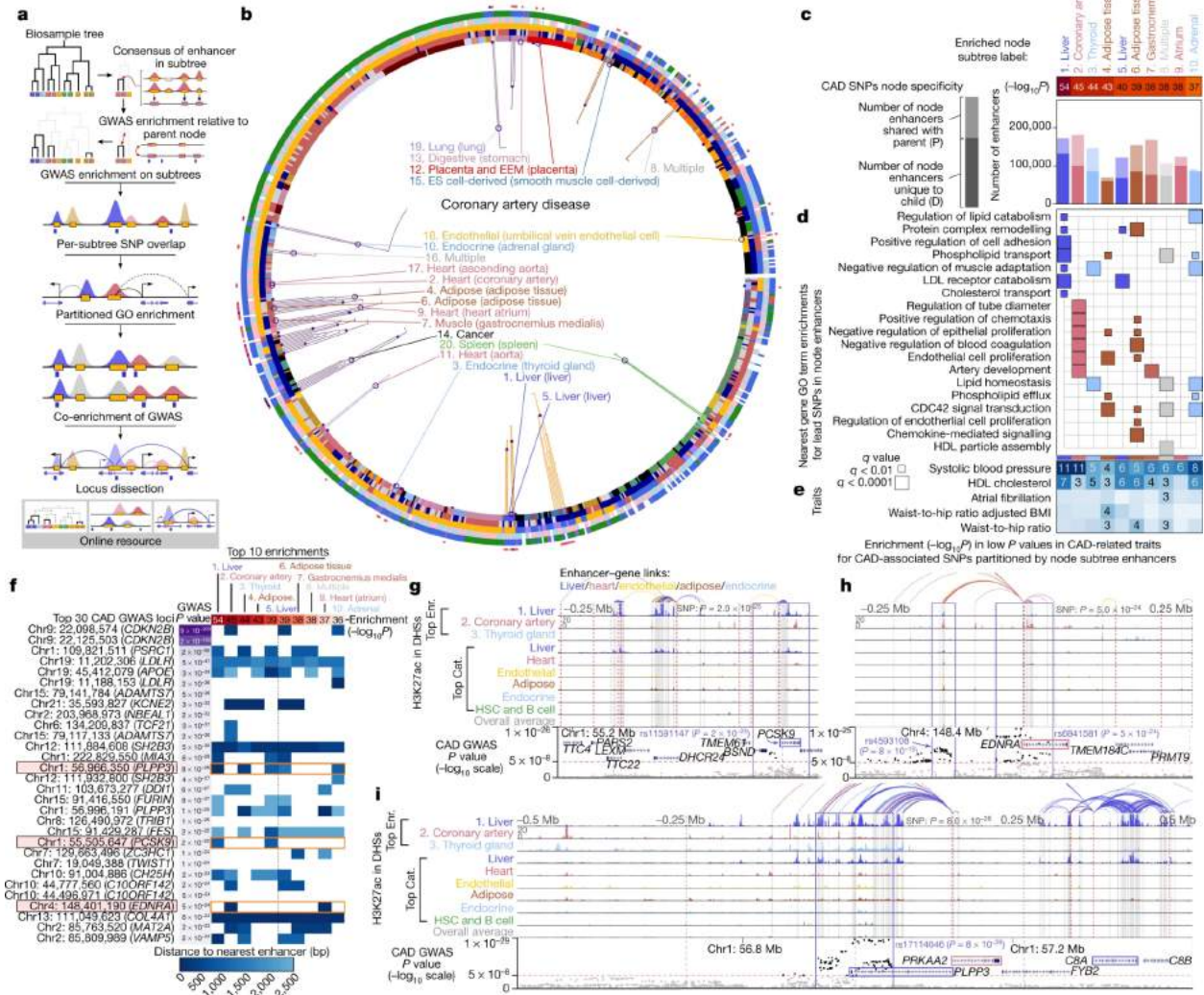
These partitions also showed distinct co-associations (**Fig. 2-5e**). For example: heart/muscle/endothelial-enhancer CAD SNPs co-associated with high blood pressure and atrial fibrillation; liver/endocrine-enhancer CAD SNPs with systolic blood pressure; adipose-enhancer CAD SNPs with waist-to-hip ratio; and liver/adipose/endocrine CAD SNPs in HDL cholesterol.

Individual multifactorial-trait loci included both single-tissue and multi-tissue loci (**Fig. 2-5f**). Some CAD loci overlapped only heart enhancers (eg. *EDNRA*, *TCF21*, *ADAMTS7*), some only liver (eg. *PCSK9*), some lacked any enhancer overlaps (possibly acting at non-enhancer levels of regulation, or in uncaptured tissues/conditions), and many overlapped enhancers active in multiple tissues (eg. *LDLR*, *APOE*, *SH2B3*, *COL4A1*), suggesting multiple mechanisms of action even at the single-locus level.

For example, the liver-only CAD-associated locus near LDL-cholesterol-regulator *PCSK9*<sup>168</sup> (rs11591147,  $p=2e-25$ , rank #21) showed a strong liver-specific signal, and has liver-specific enhancer-gene links to *PCSK9* (**Fig. 2-5g**). Heart-only CAD-associated 250-kb locus near *EDNRA* contains two separate associations, TSS-centered rs6841581 ( $p=5e-24$ , rank #27) and enhancer-centered rs4583018 ( $p=8e-15$ , rank #66), both in strong coronary-artery enhancers, both linked to *EDNRA* through strong artery links, putatively reflecting multiple functional variants<sup>165</sup> converging on the same target gene in the same tissue (**Fig. 2-5h**).

Even seemingly single-tissue loci sometimes showed second-tissue signals: the 1Mb rs17114046 locus ( $p=8e-28$ , rank #14, **Fig. 2-5i**), showed primarily liver activity with multiple SNP-overlapping enhancers linked to liver-expressed *PLPP3*, whose liver-specific deletion increases atherosclerosis<sup>169</sup>; however our liver-specific links also implicated liver-produced complement factor *C8A*<sup>170</sup>, and our heart/muscle-specific links implicated *PRKAA2*, an AMP-kinase subunit involved cardiac metabolism<sup>171</sup>, which are both biologically-relevant, highlighting that even individual loci may be pleiotropic, a property repeatedly found for many top-scoring loci.

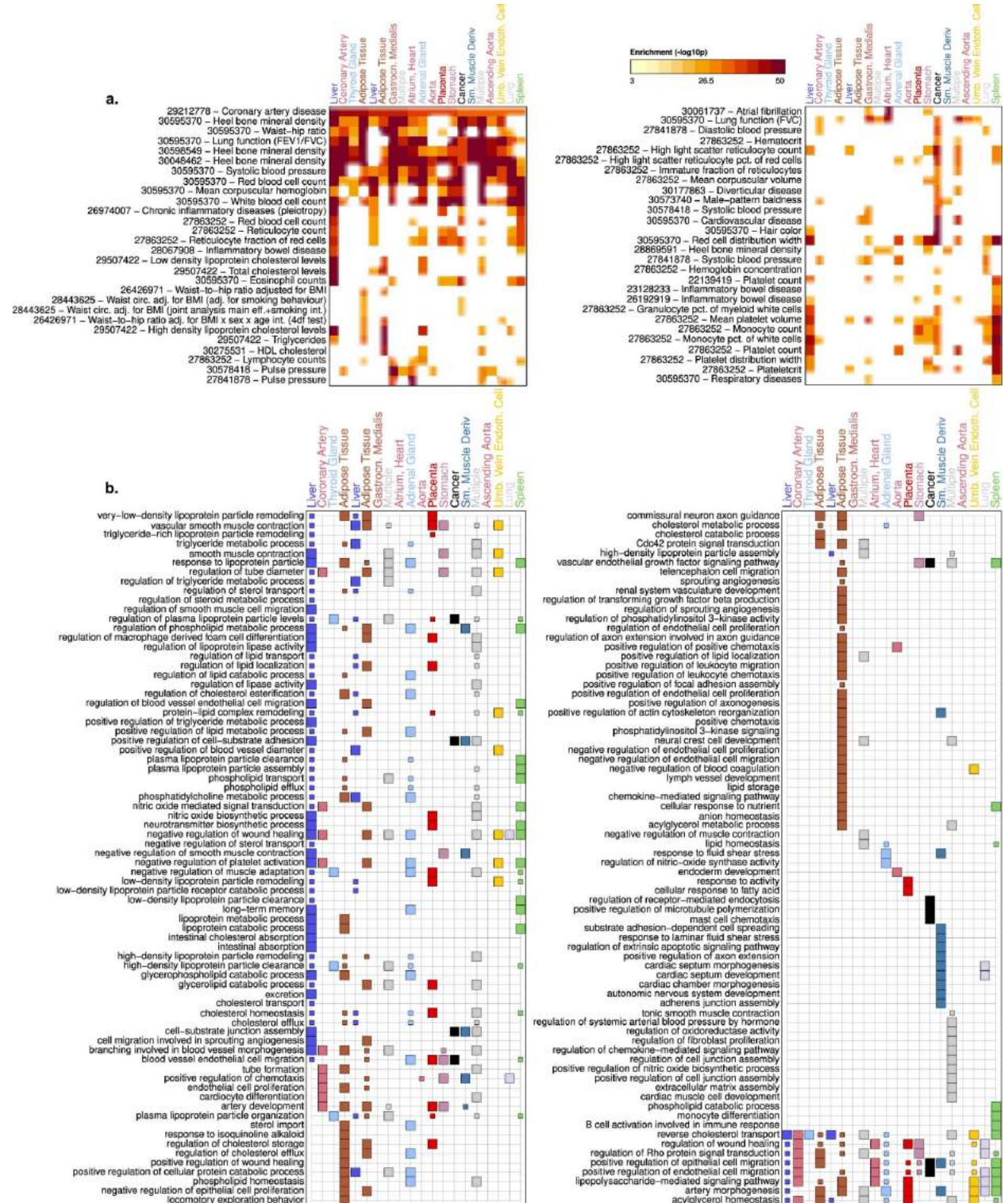




**Figure 2-5. Partitioning of polyfactorial traits.** **a.** Workflow for dissection of GWAS epigenomic enrichment using the biosample tree (**Figure 2-1d**). Additional trait enrichments, SNP assignments, links, and their corresponding loci are available through <https://compbio.mit.edu/epimap>. **b.** Epigenomic enrichments for Coronary Artery Disease (CAD, PubMedID 29212778) on enhancer-sharing tree. Nodes passing adjusted nominal p-value  $< 1e-3$  are labeled by rank, category, and components, and subtrees are shown (large circles = top 20 nodes by  $-\log_{10}p$ ). Leaves annotated by metadata and number of enriched parent nodes (outer, red=1, black=2). **c.** Top 10 enriched nodes for CAD with nominal p-values (heatmap) and shared enhancer set sizes (barplot) with number at the subtree (full bar) and number of differential enhancers between the node and its parent (tested set, dark bar). **d.** GO enrichments of node enhancers with lead SNPs (nearest expressed genes), colored by each node's tissue group and diagonalized (over-representation test). **e.** Enrichment for significant loci in overlap of CAD loci with loci from five related traits, within enriched enhancers in each node (heatmap,  $-\log_{10}p$  of one-tailed Mann-Whitney test against each trait's loci in enhancer annotations). **f.** Enhancer overlaps with CAD GWAS' top 30 lead SNPs for top 10 enrichments on the enhancer tree. **g, h, i.** Loci centered on CAD lead SNPs with links (top), H3K27ac signal (middle), and summary statistics, for lead SNPs rs11591147 (chr1:55,505,647,  $p=2e-25$ ) (**g**),



rs6841581 (chr4:148,401,190,  $p=5e-24$ ) (h), and rs17114046 (chr1:56,966,350,  $p=8e-28$ ) (i). Loci show (top) enhancer-gene links for SNP-proximal enhancers for the top enrichments and across the locus for labeled categories (linked enhancers in grey); (middle) H3K27ac signal in enhancers for top three enriched subtrees, six selected tissue categories, and average; (bottom) genes (TSS - red lines) and CAD GWAS summary statistics, with SNPs below  $p=5e-8$  in grey.



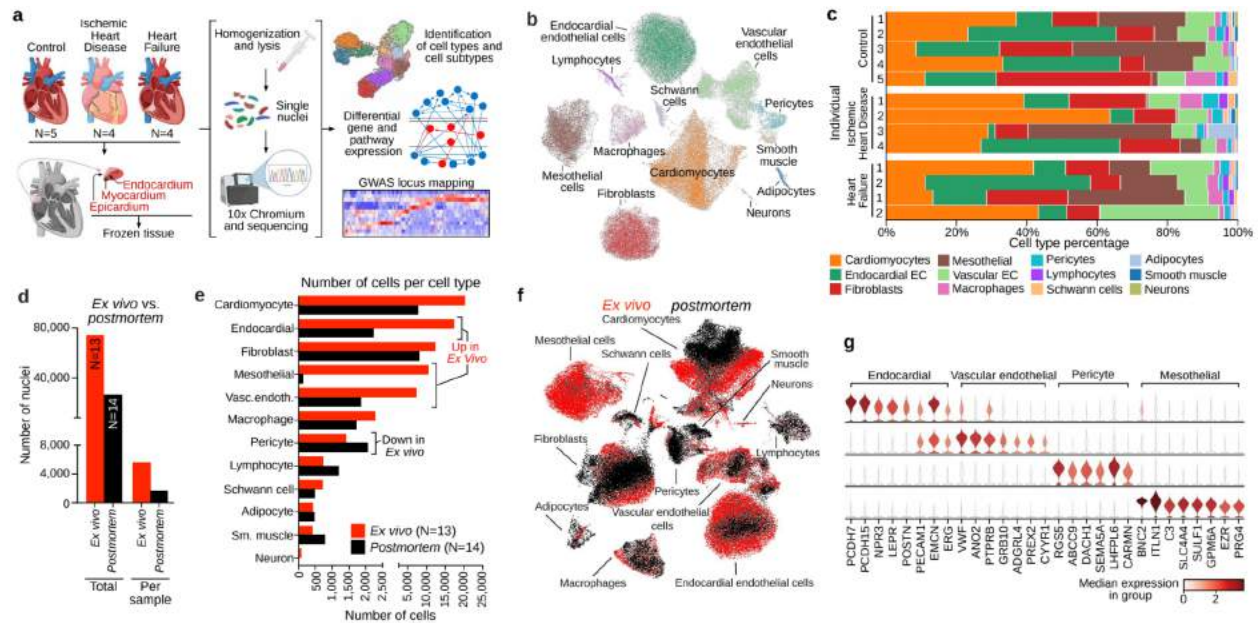
**Figure 2-6. Extended CAD dissection.** **a.** Shared enrichments with CAD for 56 of 803 traits sharing at least two enrichments with the top 20 enriched nodes of CAD. Matrix is diagonalized according to maximal enrichment (nominal p-value, only enrichments passing adjusted nominal p-value < 1e-3 are shown). **b.** Extended GO terms for CAD lead SNPs in enriched nodes. All go terms with at least  $-\log_{10}q > 2$  enriched in less than 25% of nodes.

## 2.5 Dissection of a multi-gene pleiotropic coronary artery disease-associated locus

**Summary:** Ischemic heart disease, the single most common cause of death worldwide, has an annual death rate of over 9 million people. Although genome-wide association studies promote deeper understanding of its causal mechanisms, resolution at the cellular and molecular level requires identification of the cell-type-specific circuits enabling dissection of driver variants, genes, and signaling pathways in normal and diseased tissues. Here, we provide the first single-cell dissection of the cell types and disease-associated gene expression changes in the living human heart, using cardiac biopsies from open-heart surgery of control, ischemic heart disease, and ischemic and non-ischemic heart failure patients. We identify 12 major cell types and define markers for each of them, providing the first extensive reference set for the live heart. Major cell types include cardiovascular cells (cardiomyocytes, endothelial cells, fibroblasts), rarer cell types (B lymphocytes, neurons, Schwann cells), and rich populations of previously understudied layer-specific epicardial and endocardial cells. We reveal substantial differences in disease-associated gene expression at the cell subtype level, with arterial pericytes playing central role in ischemic heart disease and heart failure pathogenesis thus, demonstrating the importance of high-resolution cellular subtype mapping in gaining mechanistic insight into human cardiovascular disease.

**Cell-level transcriptomic profiling of live human heart:** We obtained live (*ex vivo*) cardiac tissue biopsies from 13 donors collected during open-heart surgery from the right atrial appendage of 5 control, 4 ischemic heart disease, 2 ischemic heart failure and 2 non-ischemic heart failure patients (i.e., pressure-overload heart failure caused by aortic stenosis), with each sample containing all layers of the heart wall, namely the epicardium, myocardium and endocardium (**Fig. 2-7a**). We collected nuclei successfully from all samples and profiled them using snRNA-seq by 10X Genomics, resulting in 76,087 nuclei grouped in 12 main cell types (**Fig. 2-7b**) annotated using available marker genes<sup>172-175</sup> and gene ontology enrichment of biological processes and represented across all samples (**Fig. 2-7c**).

In previous studies<sup>172</sup>, endocardial endothelial cells and epicardial mesothelial cells were barely recovered, comprising only 6.8% and 0.5% of all cells, respectively, whereas here, endocardial and epicardial cells were among the most abundant cell types comprising 23% and 14% of the cells, respectively. The distribution of other cell types corresponded to expected ratios in the cardiac tissue, given the collection of all three layers of the heart wall (**Fig. 2-7a**). We also found that our live tissue datasets show better cell enrichment, significantly stronger expression signals, and better separation of all cell types and cell subtypes than tissue collected after brain or cardiovascular death (**Fig. 2-7d-e**), consistent with better preservation of tissue properties. The strongest differences between *postmortem* and living tissue were found in endocardial and vascular endothelial cells, which form the inner layers of the heart and blood vessels, respectively, and epicardial mesothelial cells, which form the outer layer of the heart (**Fig. 2-7e-f**), suggesting that these cell populations are most degraded in *postmortem* tissue.



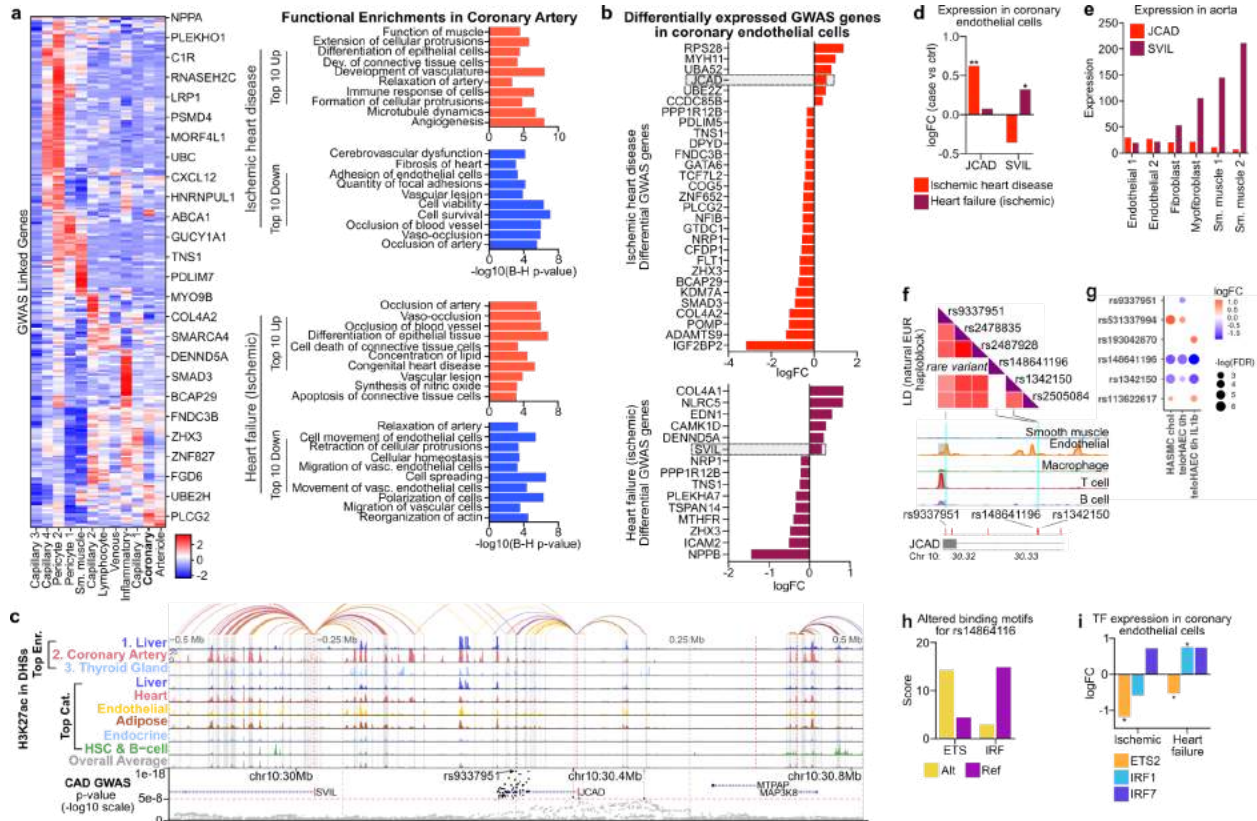
**Figure 2-7. Identification of marker genes and comparison of *ex vivo* and *postmortem* tissues.** **a.** *Ex vivo* samples obtained from the right atrial appendage from 13 individuals and 3 groups: control ( $N = 5$ ), ischemic heart disease ( $N = 4$ ), and heart failure, containing both Ischemic heart failure ( $N = 2$ ) and non-ischemic heart failure ( $N = 2$ ) samples. snRNA-seq was processed using the 10X Chromium chemistry (v3). Image was created using BioRender. **b.** Uniform manifold approximation and projection (UMAP) embedding of the 76,087 nuclei distinguishes populations of 12 main cell types. **c.** Cell-type distribution for each sample. **d.** Comparison of collected nuclei in *ex vivo* and *postmortem*<sup>172</sup> right atrial tissue by total number of samples and per sample. **e.** Number of cells per cell type for the *ex vivo* and *postmortem* samples. **f.** UMAP of the *postmortem* (black) and *ex vivo* (red) cell distributions. **g.** Cell-type-specific marker genes for endocardial endothelial cells, vascular endothelial cells, pericytes, and epicardial mesothelial cells, in addition to endothelial marker genes shared between endocardial and vascular endothelial cells.

Our living tissue samples revealed several novel marker genes in these cell types that can aid verification of affected pathways and improve directed targeting and isolation of these cells (Fig. 2-7g). For example, for adult human endocardial endothelial cells, we found several new marker genes that are not expressed in vascular endothelial cells, including protocadherins (*PCDH7*, *PCDH15*), and *CGNL1*, involved in cell-cell recognition, and tight junction assembly and maintenance (Fig. 2-7g).

**GWAS target prioritization and locus dissection:** To identify enhancers that mediate effects of disease-associated genetic variants to genes in ischemic heart disease, we linked enhancers to their target genes and mapped coronary artery disease associated single nucleotide polymorphisms (SNPs) to the enhancers in each main cell type and in the vascular subtypes. Generally, these GWAS-linked genes were enriched in common pathways that mediated disease-relevant changes within the specific context of each cell type's function. For example, in ischemic heart disease, the coronary endothelial GWAS-linked genes enriched for cellular rearrangement, movement and differentiation, angiogenesis, and immune response, whereas in ischemic heart failure, they enriched for increased vascular occlusion, differentiation, cell death, lipid concentration, and nitric oxide synthesis (Fig. 2-8a). Pericyte 2 and capillary 4 populations



were enriched for most GWAS-linked genes in the vascular subtypes and cardiomyocytes and pericardial mesothelial cells in the main cell types.



**Figure 2-8. GWAS gene mapping and pathway enrichment.** **a.** Expression of genes linked to coronary artery disease GWAS loci in vascular cells. IPA's 'Diseases and Biological Functions' enrichments are shown for coronary artery endothelial cells in ischemic heart disease and heart failure groups. 'Top up' shows top enrichments with positive Z-score (function increases based on differential gene expression) and 'Top down' top enrichments with negative Z-score (function decreases). All enrichments are statistically significant based on corrected p-values. **b.** Statistically significantly ( $p$ -value  $< 0.05$ ) differentially expressed GWAS-linked genes in coronary artery endothelial cells for ischemic heart disease and ischemic heart failure. **c.** EpiMap linking of the *JCAD/SVIL* locus<sup>176</sup>. **d.** *JCAD* and *SVIL* expression in coronary endothelial cells in ischemic heart disease and ischemic heart failure patients compared to the control group. **e.** *JCAD* and *SVIL* expression in aortic cell (sub)types based on scRNA-seq<sup>177</sup>. **f.** Illustration of LD SNPs of the natural European haploblock combined with the scATAC-seq data from human aorta generated using LDlink<sup>178</sup>. **g.** Allele-specific enhancer activity measured with STARR-seq in telHAECs under basal and inflammatory conditions (6h IL1b) and HASMCs subjected to cholesterol loading for 24h. SNPs demonstrating significant changes ( $FDR < 0.1$ ) in enhancer activity are shown **h.** Transcription factor binding motifs altered by rs148641196. Position weight matrix scores for reference and alternate are shown. **i.** Changes in transcription factor (TF) *ETS2*, *IRF1*, and *IRF7* gene expression in coronary endothelial cells in ischemic heart disease and ischemic heart failure.

Among the top differentially expressed GWAS genes ( $p < 0.05$ ) in coronary endothelial cells





## 2.6 Methodology for GWAS enrichment analysis

We pruned the NHGRI-EBI GWAS catalog<sup>162</sup> (downloaded from <https://www.ebi.ac.uk/gwas/docs/file-downloads> on May 3rd, 2019) using a greedy approach: within each trait + PMID combination, we ranked associations by their significance (p-value), and added SNPs iteratively if they were not within 5kb of previously added SNPs (vs. 1Mb SNP pruning used in Roadmap Epigenomics). We also removed all associations in the HLA locus (for hg19: chr6:29,691,116-33,054,976). This reduced the catalog from 121k to 113k associations. Finally, we reduced the catalog to 803 unique GWAS (from 5454 GWAS) with an initial sample size of at least 10,000 cases or individuals (wherever cases and controls were not annotated) and 10 GWAS SNPs after pruning. In enrichment analyses, we considered enhancers as intersecting with GWAS SNPs if the SNP was within 2.5kb of the enhancer midpoint.

### **Flat GWAS-epigenome enrichments and module-based GWAS-epigenome enrichments:**

We performed the hyper-geometric test to evaluate GWAS enhancer-level enrichments on flat epigenomes and on modules. For these flat enrichments, we compare each number of SNP-enhancer intersections for each enhancer set (flat epigenome or module) to the full set of intersections in all M enhancers. As above, we correct for multiple testing for each GWAS and enhancer set combination by computing and correcting with null association p-values for flat epigenomes and modules using the null catalogs generated for the tree enrichment. In addition, we performed SNP-level hyper-geometric enrichment tests at the level of captured SNPs: for each GWAS and enhancer set combination, we compared the number of GWAS SNPs captured by the set vs. the total number of GWAS catalog SNPs captured by the set against the total number of GWAS SNPs vs. the total number of GWAS catalog SNPs. We corrected SNP-level enrichment p-values with the `p.adjust` function in R with the BH method. Rarefaction curves were calculated on the flat epigenome enrichments by iteratively adding the sample that was either significantly enriched or the maximal enrichment for the most remaining GWAS until all GWAS were accounted for (**Fig. 2-2c-d**).

**Tree-based GWAS-epigenome enrichments:** We constructed a tree by hierarchically clustering the Jaccard similarity of the binary enhancer-by-epigenomes matrix using complete-linkage clustering. Then, for each node in the tree, we calculated its consensus epigenomic set, defined as the set of all enhancers present in all leaves of the subtree, such that each node's set was a superset of that of its parent. For each GWAS, we asked whether the novel consensus enhancers at a node were significantly enriched for lead SNPs by comparing the enrichment between each node and its parent as measured by the likelihood ratio test between two logistic regressions.

Briefly, for each GWAS catalog unique trait and PubMedID, we find all intersections of its pruned SNPs with  $M=2,069,090$  enhancers. Then  $\mathbf{Y}$  is an indicator vector of size  $M$  which shows the intersected enhancers. We find all consensus enhancers (intersection of epigenomes in the subtree) in the node of interest (vector  $\mathbf{X}_N$ ) and in its parent ( $\mathbf{X}_P$ ). All vectors are  $1 \times M$ . We calculate  $\mathbf{X}_D = \mathbf{X}_N - \mathbf{X}_P$  (specific enhancers), which is also in  $\{0, 1\}^{(1 \times M)}$  as each node contains a superset of its parent's enhancers. We then calculate the following two logistic regressions: M1:  $\mathbf{Y} \sim \mathbf{X}_P + 1$ ; M2:  $\mathbf{Y} \sim \mathbf{X}_P + \mathbf{X}_D + 1$ . We calculated the log-likelihood difference and applied the likelihood ratio test to test whether adding the specific enhancers (M2) is significantly different from the parent model (M1). To correct for p-value inflation on a per GWAS and node basis, we generated 1000 null GWAS for each lead SNP set size by shuffling the trait associations across GWAS locations, giving 243,000 null GWAS in total. We used these catalogs to compute the null association p-values for each permuted GWAS and used the 0.1% and 1% top quantiles as adjusted nominal p-value cutoffs.

For the CAD example, GO terms<sup>181</sup> were calculated using the nearest gene of each enhancer hit

by a lead SNP. We pruned genes to expressed genes by calculating average RNA-seq profiles for each tissue group and excluded genes that had  $\log_2$  FPKM of less than 2 in the average RNA-seq of each sample's group. 341 of 833 samples have matched RNA-seq, which we list in addition to releasing the processed data at <http://compbio.mit.edu/epimap>. We kept only the GO terms that were significant in 25% or less of nodes and report the top 2 GO terms per node in **Fig. 2-5d**.

For locus dissections (in NTN4, CACNA1C, EDNRA, and PLPP3), we found the nearest active enhancer to each lead SNP in each node (within 2.5kb), plot H3K27ac signal in the 2.1M enhancers only, and both **(1)** directly-mapped links originating at one of the enhancers near a lead SNP in the top three enriched epigenomes or **(2)** any links in the locus present in at least half of the samples in one of the selected tissue groups.

**Tissue similarity:** We assigned each internal node in the tree to a unique tissue if over 50% of the subtree's leaves came from that tissue and as "Multiple" if the subtree was not majority one tissue. We assigned tissue labels to 641 of 832 (77%) internal nodes where the majority of leaves corresponded to a single group. Using these assignments, we created a tissue by GWAS matrix by adding the  $-\log_{10}$  p-values for each tissue node set from all of the GWAS enrichments on the tree. We binarized this matrix and computed the Jaccard similarity across tissues to calculate a tissue similarity matrix. To assess significance of tissue overlap, we compared each overlap value against the overlaps from 10,000 permuted enrichments. We collapsed each permuted matrix into a tissue by GWAS matrix to compute the overlaps under the null. We performed the permutations for each tissue against other tissues by shuffling the enrichment p-values on the node by GWAS matrix. Specifically, we (a) binarized the enrichment matrix, (b) fixed the column of the group of interest and (c) permuted the remainder of the matrix keeping its row and column marginals the same and then (d) calculated the cosine distance between the permuted and the original matrix of enrichments.

**Cross-GWAS Network:** To evaluate the cross-GWAS similarity, we normalized the tissue by GWAS matrix for each GWAS to obtain the proportion of significance attributed to each tissue for each GWAS (**Supplementary Fig. 2-S4**). We reduced the matrix to 538 significant GWAS with at least 20,000 cases (or individuals when no cases were specified) with an adjusted nominal p-value of less than 1% at the per-node and per-GWAS size level. We created a GWAS-GWAS network using the cosine distance matrix as an adjacency matrix, keeping 5,547 links with a cosine distance of 0.25 or less. We used the Fruchterman-Reingold algorithm to lay out the graph<sup>182</sup>. We use the tissue by GWAS matrix to color links according to the maximum tissue in the product between each pair of nodes and to color nodes according to the maximal tissue for each node (**Supplementary Fig. 2-S5**).

In order to compare the epigenetic network to trait genetic similarity, we binned SNPs in the GWAS catalog into 10kb windows starting from the beginning of each chromosome. We counted the number of intersecting bins between two traits and kept any trait pairs with Jaccard similarity of at least 1%. To compare this to the epigenetic network, we plotted only links in the epigenetic network that coincide with any SNP-sharing GWAS pairs. Additionally, we plotted the heatmaps of the tree enrichments distance matrix and the genetic similarity matrix side by side, first organized by hierarchically clustering the enrichments matrix and then by clustering the genetic similarity matrix (**Supplementary Fig. 2-S6-7**).

## 2.7 Discussion

In this work, we presented a comprehensive map of the human epigenome, EpiMap, encompassing ~15k epigenomic tracks across 833 distinct biological samples that greatly expand the coverage of both embryonic and adult tissues and cells. Our work provides

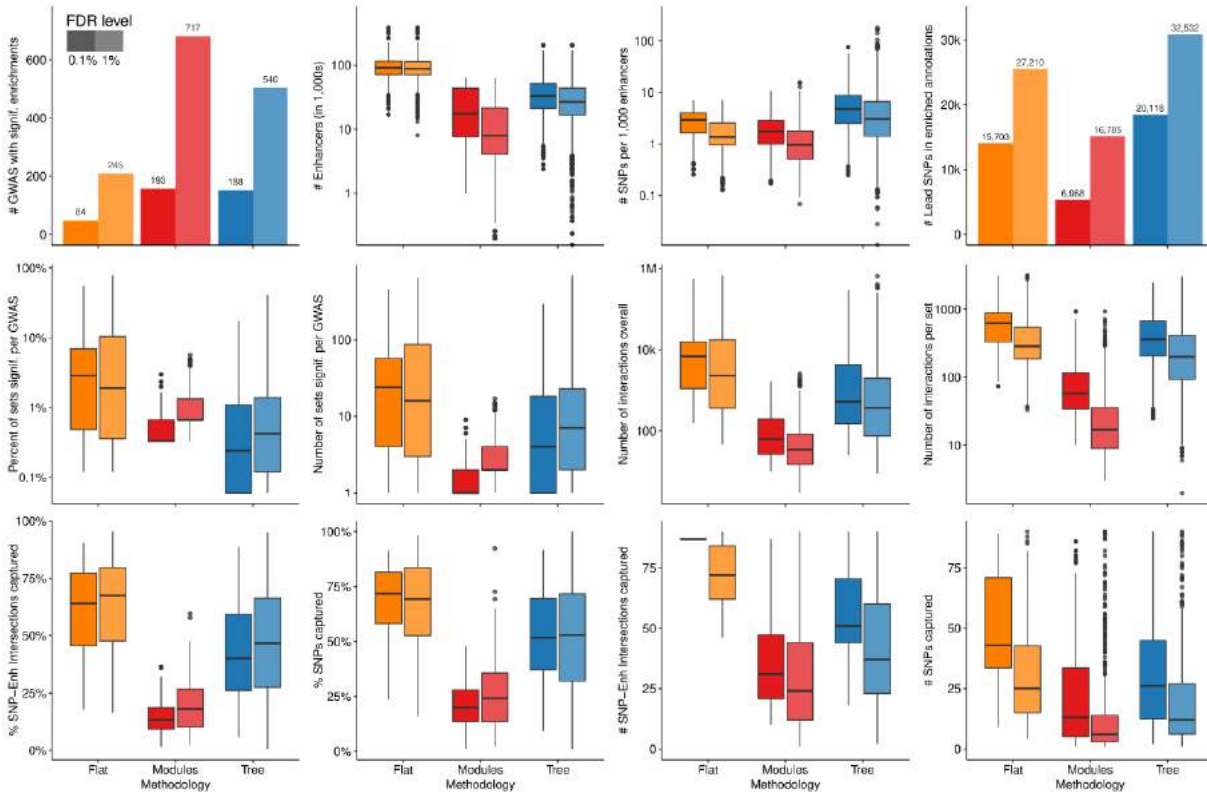
high-resolution molecular dissections of complex traits and human disease circuitry. We found significant epigenomic enrichments for 540 GWAS-studied traits implicating 30,000 SNPs in tissue-enriched enhancers, used trait- and tissue- co-enrichment patterns to annotate tissue partnerships and trait pleiotropy, and to partition disease SNPs into tissue-specific functional components. For individual GWAS loci, our work provides mechanistic insights at unprecedented scale. We highlight specific examples of GWAS dissections fitting different models of GWAS mechanistic basis at increasing levels of complexity, from the typically sought single-enhancer-to-gene, to multiple-enhancer convergence on a single target<sup>165</sup>, to multi-gene through multi-tissue pleiotropy in a single locus.

Beyond the specific examples highlighted in our figures, we also provide a rich interactive supplementary website (**Supplementary Fig. 2-27**) for our study (at [compbio.mit.edu/epimap](http://compbio.mit.edu/epimap)), enabling detailed interactive exploration of: functional and motif enrichments of 300 enhancer modules; motif-tissue networks and enrichments; GWAS enrichments for 540 traits against our biosample tree; GWAS enriched-tissue enhancer-SNP overlaps and target gene predictions; and 30,000 disease-locus visualizations with putative driver SNPs, enhancers, tissue, and tissue-specific target genes. These enable detailed hypothesis generation for experimental follow-up in countless studies of gene regulation and disease, which we illustrated in an example locus dissection (JCAD/SVIL) in coronary artery disease.

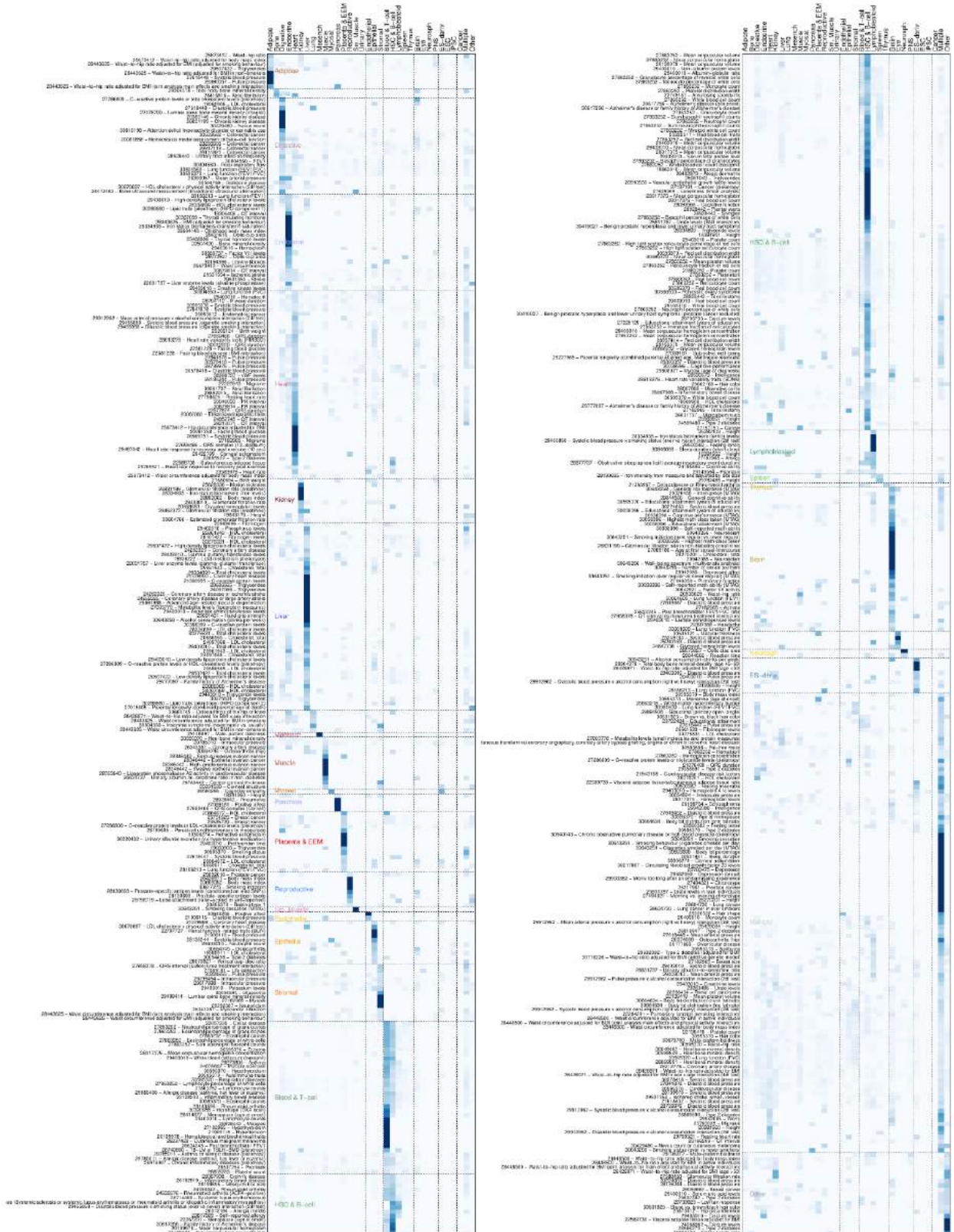
Our work enables many future studies: hierarchical and multi-resolution tree-based analyses of gene regulation and GWAS; machine-learning-based gene circuitry and combinatorial regulatory motif analyses<sup>160,161</sup>; more sophisticated network analyses of our tissue-trait, trait-trait, and tissue-tissue relationships; and guiding the experiment prioritization, methodological development, and validation experiments, which can continue to further our understanding of gene regulation and human disease circuitry.



## 2.8 Supplementary Figures

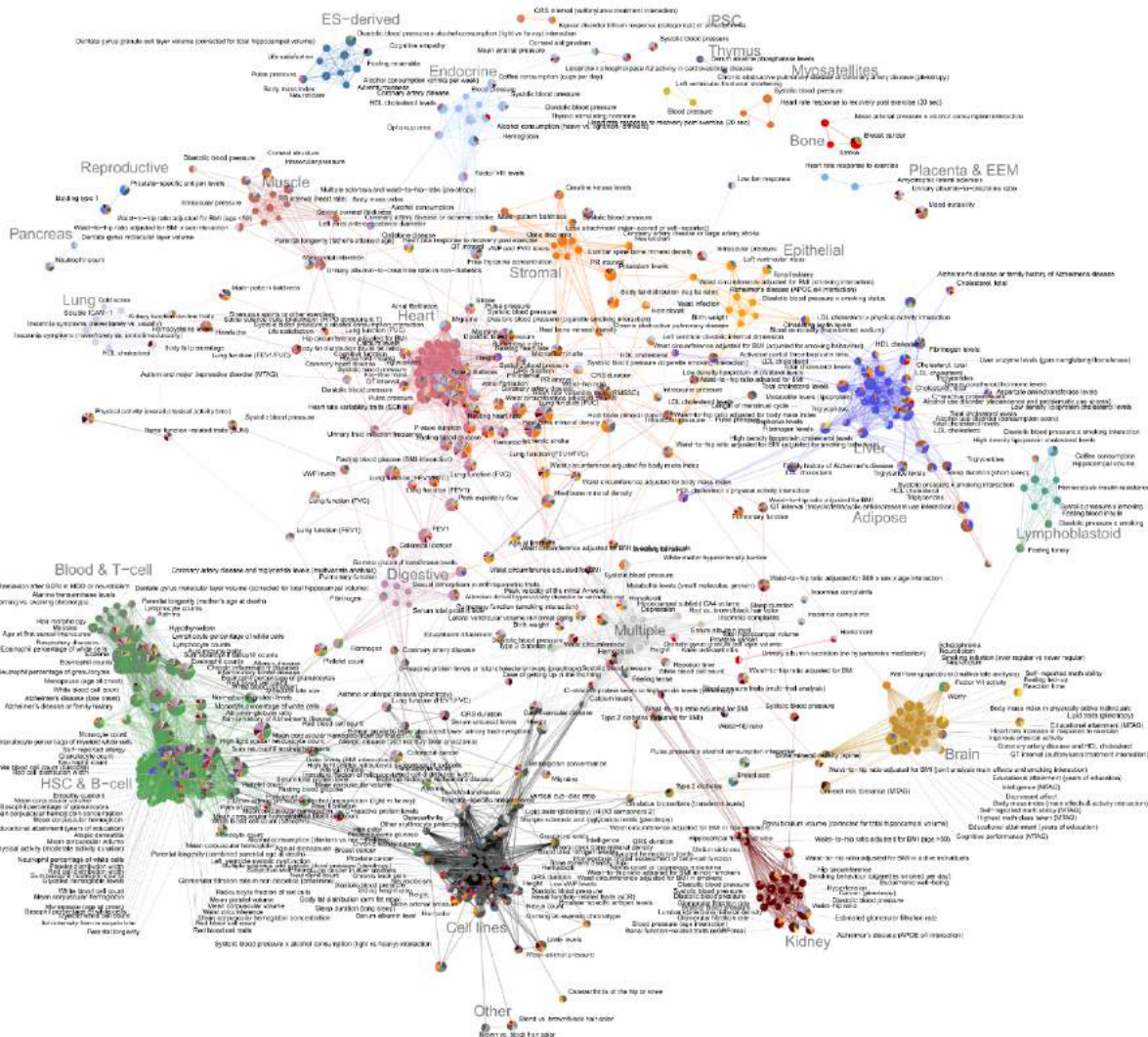


**Figure 2-S1:** Extended method validation figures comparing epigenomic GWAS enrichments using different methodologies (x-axis) for two adjusted nominal p-value cutoffs (shades). First and last figures from the top row are shown as **Fig. 2-2e**.



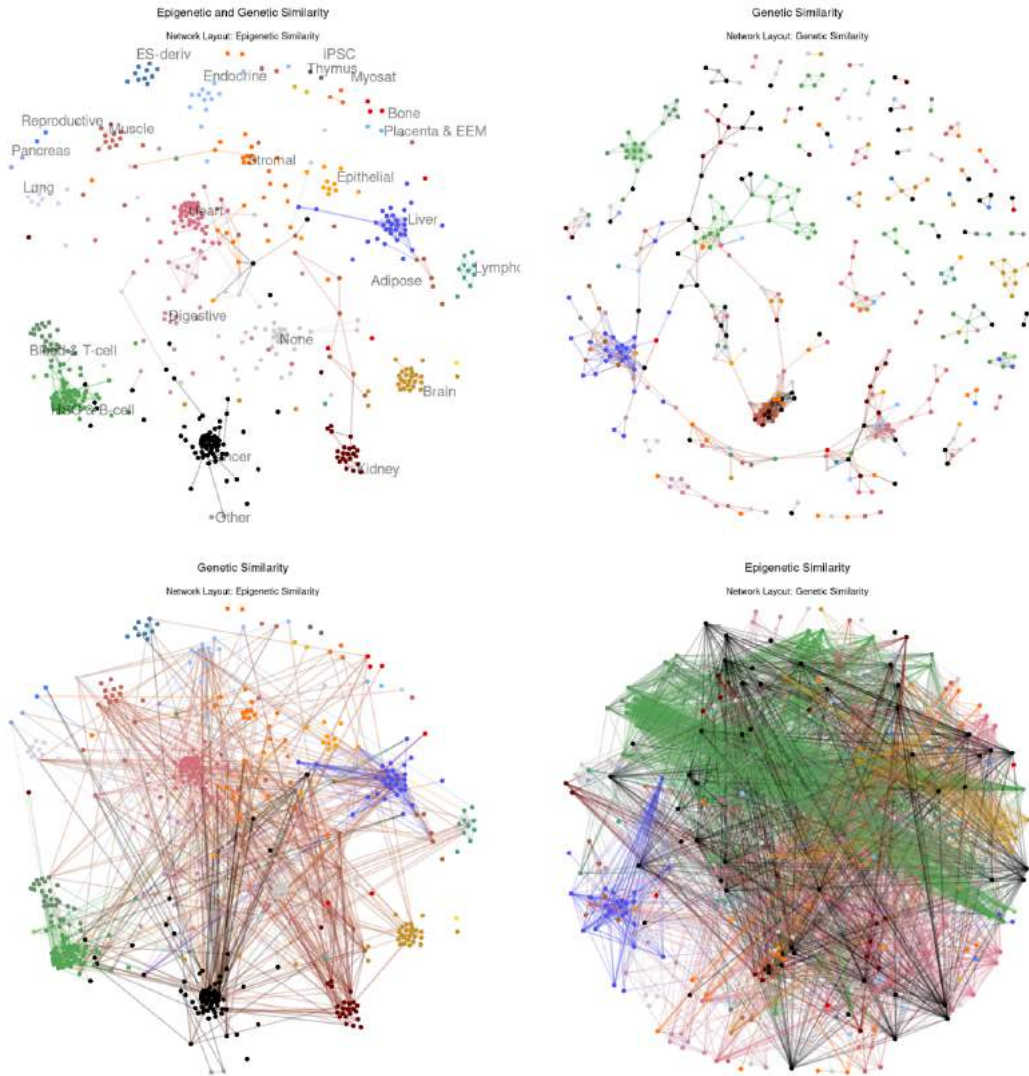
**Figure 2-S4:** Overall tissue-level prioritization of 540 traits at adjusted nominal p-value < 1%. Heatmap represents trait (rows) vs. tissue matrix and is split for visibility. Traits are diagonalized

according to their top tissue enrichment, and values are the trait-normalized and tissue aggregated  $-\log_{10}p$ -values.



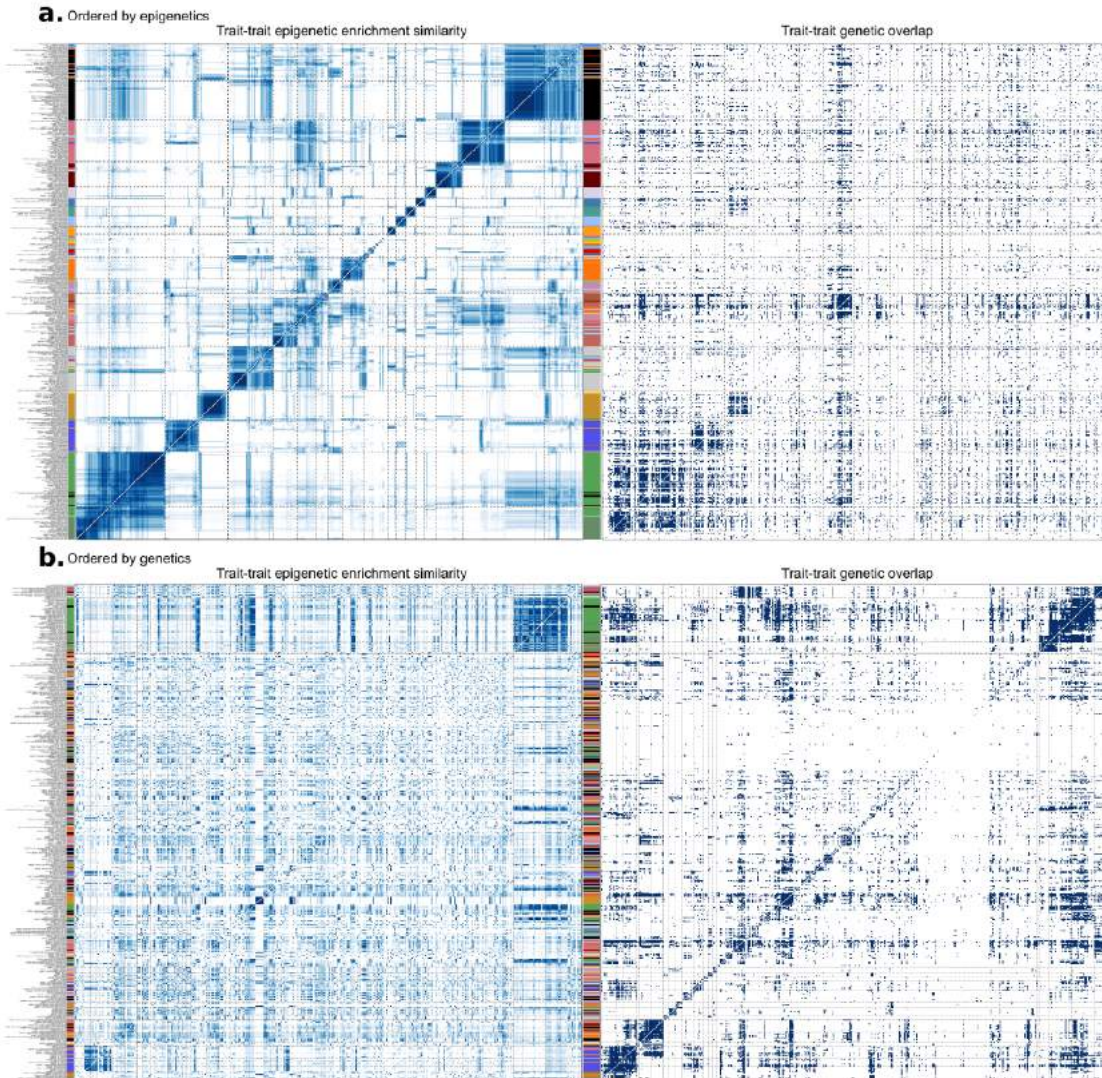
**Figure 2-S5:** Trait-trait network (as Figure 5) across 511 traits by similarity of epigenetic enrichments (cosine sim.  $\geq 0.75$ ), laid out using the Fruchterman-Reingold algorithm. Traits (nodes) are colored by contributing groups (pie chart by fraction of nominal  $-\log_{10}p$ , size by maximal  $-\log_{10}p$ ) and interactions (edges) by the group with maximal dot product of enrichments between two traits. All 511 traits labeled.





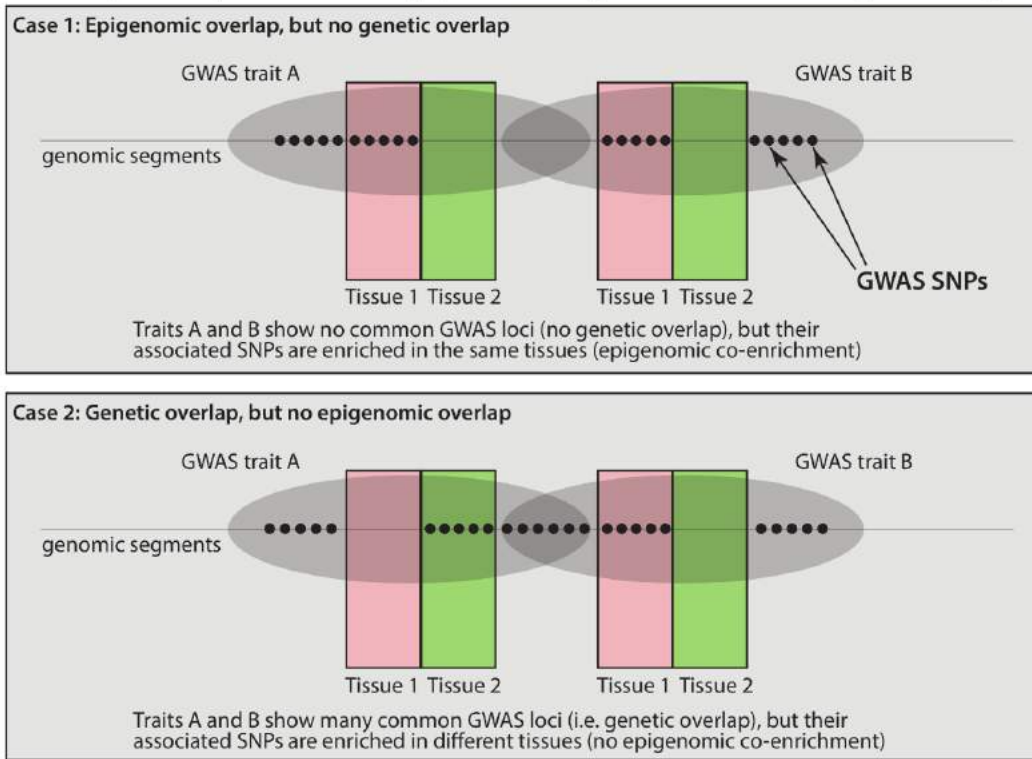
**Figure 2-S6:** Comparison networks with genetics and epigenetics. (top left) Edges with both high epigenetic similarity and any genetic overlap on epigenetic similarity layout. (bottom left) All trait pairs with genetic overlap (>5% jaccard similarity of lead SNPs overlapping when binned into 10k bp bins starting from the start of each chromosome, network using epigenetic layout. (top right) Network of trait pairs with any genetic overlap laid out by genetics. (bottom right) All trait pairs with high epigenetic similarity on the genetic layout network





**Figure 2-S7:** **a.** Trait-trait similarity matrix by cosine similarity of epigenetic enrichments (left) and by 0/1 genetic similarity (defined as >5% jaccard similarity of lead SNPs overlapping when binned into 10k bp bins starting from the start of each chromosome). Matrices are ordered by hierarchical clustering according to Ward's method on the epigenetic matrix. **b.** Epigenetic similarity (left) and genetic similarity (right) matrices, as above. Matrices ordered by hierarchical clustering according to Ward's method on the genetic similarity matrix.

Two simple example cases where epigenomic and genetic trait-trait overlaps can disagree



**Figure 2-S8:** Two example cases where the epigenetic and the genetic trait-trait links may not agree.



### 3. A multi-region single-cell-resolution atlas of the aging human brain

Parts of this chapter have been presented/published as:

#### **Single-cell multi-region dissection of Alzheimer's disease progression reveals selectively vulnerable neurons**

*(in preparation)*

Hansruedi Mathys<sup>1-3\*</sup>, Carles Boix<sup>4-6\*</sup>, Leyla Akay<sup>1,2</sup>, Ayesha Ng<sup>1,2</sup>, Xueqiao Jiang<sup>1,2</sup>, Kyriakista Galani<sup>4,5</sup>, Julio Mantero<sup>4,5</sup>, Neil Band<sup>4,5,7</sup>, David A. Bennett<sup>8</sup>, Li-Huei Tsai<sup>1,2,5+</sup>, Manolis Kellis<sup>4,5+</sup>

<sup>1</sup>Picower Institute for Learning and Memory, MIT <sup>2</sup>Department of Brain and Cognitive Sciences, MIT <sup>3</sup>University of Pittsburgh  
<sup>4</sup>Computer Science and Artificial Intelligence Laboratory, MIT <sup>5</sup>Broad Institute of MIT and Harvard <sup>6</sup>Computational and Systems  
Biology Program <sup>7</sup>University of Oxford <sup>8</sup>Rush Alzheimer's Disease Center. \*Equal contributors. +Correspondence to:  
[manoli@mit.edu](mailto:manoli@mit.edu) and [lhtsai@mit.edu](mailto:lhtsai@mit.edu)



### 3.1 Introduction

Alzheimer's disease (AD) is a highly prevalent neurodegenerative disorder that displays a consistent progression through brain regions and is not well understood at the level of molecular mechanism. To address this, we profiled the transcriptome of six different brain regions from a cohort of 48 individuals spanning AD progression, capturing over 1.3 million cells. Our multi-region single-cell molecular atlas of the aged human brain showcases brain cellular diversity and reveals that the composition of major cell types strongly varies across anatomical regions with a predominance of neurons in neocortical regions and glial cells in the entorhinal cortex, the hippocampus, and the anterior thalamus. We identify highly region-specific subtypes of excitatory and inhibitory neurons and glia, including an inhibitory interneuron population that is unique to the thalamus and distinct from the canonical parvalbumin, somatostatin, and VIP subclasses of inhibitory neurons. We discover thalamic astrocyte subtypes that are not observed in the neocortex and are characterized by the expression of a unique combination of genes encoding neurotransmitter transporters, ion channels, and signaling receptors. We develop a novel computational method to identify gene expression modules at the atlas scale and across multiple resolutions, revealing regional glial programs.

Alzheimer's disease (AD) and other related neurodegenerative disorders are marked by pathological protein aggregates showing a stereotypical progression through human brain regions in neurodegeneration<sup>21,27,28</sup>. Understanding the molecular underpinnings of the regional specificity, consistency, and cross-region progression of these disorders has enormous therapeutic potential for targeting both early stages and cross-region progression in these disorders. However, this is hindered in part by a lack of comprehensive molecular characterization and systematic dissection of these brain regions in disease.

The pathological hallmarks of AD include extracellular amyloid plaque deposits and intracellular neurofibrillary tangles (NFT) in neurons. AD is staged according to regional progression of NFTs, which first appear in the entorhinal cortex (Braak stages I-II), then the hippocampus and thalamus (III-IV), and finally progress to the neocortex (V-VI), a sequence that is most often synchronous with cognitive decline from mild memory impairment into severe dementia<sup>21,29-31</sup>.

Cross-tissue single-cell atlases of human development and disease have recently proven very insightful for understanding aging<sup>10</sup> and diseases such as COVID-19<sup>32</sup>. Although many of the AD-relevant regions have been studied at the molecular level<sup>33-37</sup>, a comprehensive multi-region single-cell resolution atlas incorporating information from dozens of individuals is missing. This cross-region atlas is essential to ask detailed questions about the consistency of AD's progression across the brain. In particular, it remains to be determined how cellular diversity varies between regions and how it is affected by the progression of Alzheimer's disease. Furthermore, the transcriptional modules controlling the biology of the diverse cell types in each region, the common and brain region specific alterations of gene expression in AD, and the transcriptional signatures distinguishing vulnerable and resilient cell types remain to be determined.

While many previous studies have focused on neuronal changes in early-affected regions and have identified single vulnerable excitatory neuron subtypes in Alzheimer's disease<sup>38-45</sup>, region-specific glial diversity and changes in AD are less well understood. In addition, no study has comprehensively tracked subtype vulnerability across the AD-involved brain regions<sup>21,29-31</sup> or between subtypes at different stages of AD. In addition, these studies have not been able to systematically ask what factors may mediate vulnerability at the molecular level, or whether vulnerable subtypes show transcriptomic changes from relatively more resilient subtypes or at different timepoints.

Here, we present a systematic molecular cell atlas of the human brain spanning six distinct anatomical regions affected at different stages of Alzheimer's disease (AD) progression to understand the healthy brain and its disease-associated changes<sup>183</sup>. We profile the transcriptomes of over 1.3 million individual nuclei from the entorhinal cortex, hippocampus, anterior thalamus, and the neocortical angular gyrus, mid temporal cortex, and prefrontal cortex across 283 post-mortem brain samples from 48 individuals. We annotate highly region-specific neuronal and glial subtype diversity across regions and discover a novel interneuron subtype and region-specific glial subtypes. We annotate identity programs and regulons for each subtype and annotate regional functional differences and interpret AD changes across cell types using a novel highly-scalable multi-resolution method to detect gene expression modules in atlas scale data.

We systematically identify cognitive decline-associated vulnerable excitatory neuronal subtypes depleted in a region-specific manner by AD stage, which are co-depleted in a known hippocampus-entorhinal cortex circuit and share marker genes which may mediate their vulnerability. We map out transcriptional changes in glia and neurons, identifying an early-AD glial metabolic shift as well as shared neuronal expression changes, a subset of which may be involved in neuronal resilience to neurodegeneration. Overall, our work maps the rich neuronal and glial subtype and functional diversity across AD-implicated brain regions including novel subtypes, identifies pathway-level alterations and dysregulation in a region-specific manner in AD via a new and general methodology, and systematically annotates region-specific AD-vulnerable subtypes and mechanisms of resilience and vulnerability to neurodegeneration.

### **3.2 A 1.3M nucleus transcriptional atlas of Alzheimer's disease progression across six brain regions**

To characterize the molecular progression of Alzheimer's disease (AD) across brain regions, we performed single-nucleus RNA-sequencing (snRNA-seq) on nuclei isolated from 283 post-mortem brain samples across six brain regions from a cohort of 48 participants in the Religious Order Study (ROS) or the Rush Memory and Aging Project (MAP) (**Fig. 3-1a**). We selected individuals across a range of AD progression, from non-AD and early-AD (N=20, Braak Stages I-II) to middle (N=14, Stages III-IV) and late-AD (N=14, Stages V-VI). After quality control, the combined dataset contained 1,353,075 single nuclei from six brain regions including: the early-AD-affected entorhinal cortex (EC, 221k cells, Stages I-II); the hippocampus (HC, 221k) and anterior thalamus (TH, 208k), which are affected in Stages III-IV; and later-affected neocortical angular gyrus (AG, 220k), midtemporal cortex (MT, 227k), and prefrontal cortex (PFC, 255k, Stages V-VI) (**Supplementary Fig. 3-S1a**).

We annotated 76 high-resolution cell types, grouped into 14 major cell type classes, which were largely conserved across individuals (**Fig. 3-1b**). These include 32 excitatory (Exc, 436k cells, 32.2% of total) and 23 inhibitory neuronal subtypes (Inh, 160k, 11.8%), oligodendrocytes (409k cells, 30.2%), oligodendrocyte precursor cells (OPCs, 81k, 6.0%), astrocytes (127k, 9.4%), immune cells, including microglia, CNS-associated macrophages (CAMs), and T cells (107k, 7.9%), and vascular cells (33k, 2.45%) including endothelial (End), ependymal (Epd), choroid plexus ependymal (CPEC), smooth muscle cells (SMCs), fibroblast (Fib), and pericyte (Per) subpopulations (**Fig. 3-1b,c**). We identified clusters, marker genes, cell type identity programs<sup>184</sup>, and transcriptional regulons and compared our atlas with previously published data across species (**Supplementary Figs. 3-S1-6, Supplementary Tables 1-3**)<sup>185-187</sup>.

To gain insight into the cellular architecture of the human brain we investigated differences in the major cell type composition across the six regions analyzed. The fraction of neuronal cells increased significantly from the anterior thalamus (14.4% neurons) to the three-layer allocortical hippocampus (32.2%), the entorhinal peri-allocortex (36.6%), and the six-layered neocortical

regions (AG, MT, and PFC, 58.9%), the most recent part of the cerebral cortex to evolve (**Fig. 3-1c-e, Supplementary Fig. 3-S7**). Both inhibitory and excitatory neurons were significantly more abundant in the neocortex and least abundant in the anterior thalamus while glial cells (Ast, Oli, OPC, and Mic) tended to be less abundant in neocortical samples across individuals (**Fig. 3-1d-f**). The differences in the composition of major cell types between regions were reproducibly observed across study participants, irrespective of the disease status of the cohort analyzed (**Supplementary Fig. 3-S8a-d**), suggesting that the variability in the major cell type composition between regions is a fundamental characteristic of the human brain rather than a consequence of AD pathology. The observed dominance of glial cells over neurons in the thalamus versus the cortex is in agreement with previous studies reporting that the glia to neuron ratio varies between brain regions<sup>188</sup> and is especially high in the thalamus<sup>189</sup>. These results also aligned with previous snRNA-seq studies in mouse brain<sup>190,191</sup>, which showed lower neuronal abundance and higher oligodendrocyte fraction in the thalamus relative to the cortex (**Supplementary Fig. 3-S8**).





**Figure 3-1. Single-nucleus RNA-seq interrogation of six distinct regions of the aged human brain.** **a.** Cohort and snRNA-seq profiling summary, covering 283 samples across 6 brain regions from 48 individuals along a spectrum of Alzheimer's disease progression from ROSMAP. **b.** Joint UMAP of 1.3M cells across 14 major cell types colored and labeled by 76 high-resolution subtypes. **c-d.** Joint UMAP, colored by major cell type (**c**) and region of origin (**d**). **e.** Barplot of regional composition of major cell types. **f.** Relative enrichment of major cell types across regions by quasi-binomial regression (FDR-corrected p-values, \*\*\* $p < 0.001$ , \*\* $< 0.01$ , \* $< 0.05$ ). **g-h.** Global breakdown, region composition, enrichment, and number of cells for excitatory (**g**) and inhibitory (**h**) neuronal subtypes. **i-j.** Neuronal subtype UMAPs for excitatory (**i**) and inhibitory (**j**) neurons. **k.** Glutamatergic vs. GABAergic scores for all neuron subtypes. **l.** Gene expression heatmap for top four markers per inhibitory subclass, averaged at the sample by subclass level (columns). **m.** Functional enrichments for marker genes of each inhibitory neuron subclass (adj. p-value  $< 0.01$ ). Top unique enrichments shown (1-2 subclasses, term with  $< 500$  genes). **n.** RNAscope validation of FOXP2 and MEIS2 as markers of the unique thalamus subtype, with quantification (left) performed with wilcoxon test and representative images on the right.

### 3.3 Cohort selection and single-nucleus RNA-seq

**Cohort selection from ROSMAP:** We selected a total of 48 individuals from the Religious Orders Study and Rush Memory and Aging Project (ROSMAP) that are both ongoing longitudinal clinical-pathologic cohort studies of aging and Alzheimer's disease. The studies include clinical data collected annually, detailed post-mortem pathological evaluations, and extensive genetic, epigenomic, transcriptomic, proteomic, and metabolomic bulk-tissue profiling (DA et al., 2018). For the purpose of this study, individuals were selected based on the modified NIA-Reagan diagnosis of Alzheimer's disease and the Braak stage score (Braak stage 0,1,2,  $n=20$ ; Braak stage 3,4,  $n=14$ ; Braak stage 5,6,  $n=14$ ), with 26 individuals having a positive pathologic diagnosis of AD and 22 individuals having a negative pathologic diagnosis of AD. Details of the clinical and pathological data collection methods have been previously reported <sup>29-31,183,192</sup>. Individuals were balanced between sexes (male:female ratios 13:13 in AD, 11:11 in NoAD), matched for age (medians 86.6 years (AD) and 86.0 years (NoAD)), and postmortem interval (medians 5.9 hours (AD) and 6.3 hours (NoAD)). Informed consent was obtained from each subject, and the Religious Orders Study and Rush Memory and Aging Project were approved by an Institutional Review Board (IRB) of Rush University Medical Center.

**Isolation of nuclei from frozen post-mortem brain tissue:** The protocol for the isolation of nuclei from frozen post-mortem brain tissue was adapted from a previous study<sup>37</sup>. All procedures were carried out on ice or at 4 °C. In brief, post-mortem brain tissue was homogenized in 700  $\mu$ l homogenization buffer (320 mM sucrose, 5 mM  $\text{CaCl}_2$ , 3 mM  $\text{Mg}(\text{CH}_3\text{COO})_2$ , 10 mM Tris HCl pH 7.8, 0.1 mM EDTA pH 8.0, 0.1% IGEPAL CA-630, 1 mM  $\beta$ -mercaptoethanol, and 0.4 U  $\mu$ l<sup>-1</sup> recombinant RNase inhibitor (Clontech)) using a Wheaton Dounce tissue grinder (15 strokes with the loose pestle). Then the homogenized tissue was filtered through a 40- $\mu$ m cell strainer, mixed with an equal volume of working solution (50% OptiPrep density gradient medium (Sigma-Aldrich), 5 mM  $\text{CaCl}_2$ , 3 mM  $\text{Mg}(\text{CH}_3\text{COO})_2$ , 10 mM Tris HCl pH 7.8, 0.1 mM EDTA pH 8.0, and 1 mM  $\beta$ -mercaptoethanol) and loaded on top of an OptiPrep density gradient (750  $\mu$ l 30% OptiPrep solution (30% OptiPrep density gradient medium, 134 mM sucrose, 5 mM  $\text{CaCl}_2$ , 3 mM  $\text{Mg}(\text{CH}_3\text{COO})_2$ , 10 mM Tris HCl pH 7.8, 0.1 mM EDTA pH 8.0, 1 mM  $\beta$ -mercaptoethanol, 0.04% IGEPAL CA-630, and 0.17 U  $\mu$ l<sup>-1</sup> recombinant RNase inhibitor) on top of 300  $\mu$ l 40% OptiPrep solution (40% OptiPrep density gradient medium, 96 mM sucrose, 5 mM  $\text{CaCl}_2$ , 3 mM  $\text{Mg}(\text{CH}_3\text{COO})_2$ , 10 mM Tris HCl pH 7.8, 0.1 mM EDTA pH 8.0, 1 mM  $\beta$ -mercaptoethanol, 0.03% IGEPAL CA-630, and 0.12 U  $\mu$ l<sup>-1</sup> recombinant

RNase inhibitor). The nuclei were separated by centrifugation (5 min, 10,000 g, 4 °C). A total of 100 µl of nuclei was collected from the 30%/40% interphase and washed with 1 ml of PBS containing 0.04% BSA. The nuclei were centrifuged at 300g for 3 min (4 °C) and washed with 1 ml of PBS containing 0.04% BSA. Then the nuclei were centrifuged at 300g for 3 min (4 °C) and re-suspended in 100 µl PBS containing 0.04% BSA. The nuclei were counted and diluted to a concentration of 1,000 nuclei per microliter in PBS containing 0.04% BSA.

**Droplet-based snRNA-seq:** For droplet-based snRNA-seq, libraries were prepared using the Chromium Single Cell 3' Reagent Kits v3 according to the manufacturer's protocol (10x Genomics). The generated snRNA-seq libraries were sequenced using NextSeq 500/550 High Output v2 kits (150 cycles) or NovaSeq 6000 S2 Reagent Kits.

### 3.4 snRNA-seq processing, QC, and annotation

**snRNA-seq data preprocessing:** Gene counts were obtained by aligning reads to the GRCh38 genome using Cell Ranger software (v.3.0.2) (10x Genomics)<sup>193</sup>. To account for unspliced nuclear transcripts, reads mapping to pre-mRNA were counted. After quantification of pre-mRNA using the Cell Ranger count pipeline, the Cell Ranger aggr pipeline was used to aggregate all libraries (without equalizing the read depth between groups) to generate a gene-count matrix. The Cell Ranger 3.0 default parameters were used to call cell barcodes. We used SCANPY<sup>194</sup> to process and cluster the expression profiles and infer cell identities. We kept only protein coding genes and filtered out cells with over 20% mitochondrial or 5% ribosomal RNA, leaving 1.47M cells over 48 individuals and 283 samples across all regions. We further filtered the dataset to the top 5000 most variable genes and used them to calculate the low dimensional embedding of the cells (UMAP) (default parameters, using 50 principal components and 15 nearest neighbors) and clusters using with the Leiden clustering algorithm at a high resolution (15), giving 337 preliminary clusters<sup>195</sup>. We separately called doublets using DoubletFinder and flagged and removed clusters with strong doublet profiles and clusters showing strong individual-specific batch effects, leaving a final dataset of 1.35M cells<sup>196</sup>.

**Cell type annotations:** For the UMAP visualization of individual major cell type classes (excitatory neurons, inhibitory neurons, astrocytes, oligodendrocytes, OPCs, immune cells), the SCTransform-based integration workflow of Seurat was used to align data from individual samples, using default settings<sup>197,198</sup>. We selected the set of relevant principal components based on Elbow plots. We annotated cell types using previously published marker genes and single-cell RNA-sequencing data<sup>34,37,185,199,200</sup> (M et al., 2018). To annotate cell types based on previously published single-cell RNA-sequencing data (Allen Institute's cell types database; <https://portal.brain-map.org/atlas-and-data/rnaseq/human-multiple-cortical-areas-smart-seq>), we used three separate approaches. First, Spearman rank correlation coefficients between the average expression profiles of neuronal subpopulations previously defined by the Allen Brain Institute<sup>185</sup> and the neuronal subtypes identified in this study were computed using the `cor` function in R. Second, to project annotations of neuronal subpopulations previously defined by the Allen Brain Institute onto the neuronal cells analyzed in this study, we followed the integration and label transfer workflow of Seurat<sup>198</sup>. Third, we determined cell type marker genes based on data published by the Allen Brain Institute<sup>185</sup> using the `FindAllMarkers` function from Seurat (Wilcoxon Rank Sum test with Bonferroni correction for multiple testing; adjusted  $P < 0.05$ ) and computed module scores for each cell type marker gene set across all neuronal cells analyzed in this study using the `AddModuleScore` function of Seurat. To further annotate cell types, we determined marker genes using the `FindAllMarkers` function from Seurat (Wilcoxon Rank Sum test with Bonferroni correction for multiple testing; adjusted  $P < 0.05$ ). We tested only genes that were detected in a minimum of 25% of the cells within the cell type (`min.pct = 0.25`) and that showed, on average, at least a 0.25-fold difference (log-scale) between the cells of the

cell type and all remaining cells (logfc.threshold = 0.25). Marker genes of the high-resolution cell types or states were determined separately for each major cell type class.

**Cell cycle scores and global properties of gene expression:** G2/M and S phase cell cycle scores were determined using the function CellCycleScoring in Seurat. Histograms showing the distribution of the G2/M- and the S phase module scores in each major cell class were generated using Prism 9 software. The statistical analyses comparing the number of genes detected per cell and the number of unique transcripts (UMIs) detected per cell between cell types was performed using Prism 9 software.

**Integration of external data sources:** Single-cell transcriptomic data from the human lateral geniculate nucleus (dLGN)<sup>201</sup> was obtained from the Allen Brain Institute (<https://portal.brain-map.org/atlas-and-data/rnaseq/comparative-lgn>). Single-cell transcriptomic data from multiple cortical areas and the hippocampal formation of the mouse brain<sup>202</sup> were obtained from the Allen Brain Institute (<https://portal.brain-map.org/atlas-and-data/rnaseq/mouse-whole-cortex-and-hippocampus-10x>). Single-cell transcriptomic data across nine regions in the adult mouse brain<sup>190</sup> were obtained from the McCarroll and Macosko Labs (<http://dropviz.org/>). Single-cell transcriptomic data from the mouse nervous system<sup>191</sup> were obtained from the Linnarsson lab (<http://mousebrain.org/downloads.html>). The external data sets and the human multi-region data presented in this study were integrated using the reciprocal PCA (RPCA) pipeline with default parameters in Seurat ([https://satijalab.org/seurat/articles/integration\\_rpca.html](https://satijalab.org/seurat/articles/integration_rpca.html)). The integration of single-cell data was performed separately for astrocytes, excitatory neurons, and inhibitory neurons. For the integration of GABAergic neurons, the single-cell transcriptomic data from multiple cortical areas and the hippocampal formation of the mouse brain<sup>202</sup> was downsampled to 50,000 GABAergic neurons. For the integration of glutamatergic neurons, the human multi-region data set was downsampled to 5000 neurons per high resolution cell type. The mouse cortical data<sup>202</sup> was downsampled to 50,000 glutamatergic neurons. The frontal cortex, posterior cortex, hippocampus, and thalamus data of the DropViz dataset were combined and downsampled to 50,000 neurons. Downsampling of the data was performed using the Seurat function “subset” with default parameters. The comparison of microarray data from different human brain regions was performed using the Differential Search tool of the Allen Brain Atlas data portal (<https://human.brain-map.org/microarray/search>). The thalamus was selected as the target structure and compared to the cerebral cortex as the contrast structure. The differential search results including the fold change values and p-values of the top 2000 genes were downloaded from the data portal.

### **3.5 Excitatory neurons are structured into highly brain region-specific subtypes**

We first studied the diversity of excitatory neuron subtypes, the majority of which were highly region-specific to either the hippocampus (7 subtypes), entorhinal cortex (9), or the anterior thalamus (3) (**Fig. 3-1g,i, Supplementary Fig. 3-S9-13**). The remaining excitatory subtypes were predominantly shared across the neocortical regions (AG, MT, and PFC) although they were sometimes found in the hippocampus and entorhinal cortex, especially for subtypes in deeper layers such as L6 CT (corticothalamic) and L6b (**Fig. 3-1g, Supplementary Fig. 3-S12**). Within non-neocortical regions, hippocampal subtypes included neurons from the highly structured CA1 and CA2/CA3 subfields and dentate gyrus (DG), and the more entorhinal cortex-proximal subiculum and para/pre-subiculum areas, which were sometimes captured in EC samples (**Fig. 3-1g,i, Supplementary Fig. 3-S12**)<sup>34</sup>. Entorhinal cortex-specific subtypes clustered separately from neocortical subtypes for the same layers, shared markers in common, including RELN, TOX3, and GPC5, and contained subtypes from both the lateral (L2 RELN GPC5) and medial (L2 TOX3 POSTN) entorhinal cortex (**Fig. 3-1g,i, Supplementary Fig.**

**3-S11**)<sup>203</sup>. Finally, excitatory neurons in the anterior thalamus were predominantly composed (74%) of a single thalamus-specific subtype (NXPH1 RNF220) while other thalamic subtypes were shared with the hippocampus and entorhinal cortex (**Fig. 3-1g, i**). Overall, we identified regional specificity in neuronal subtypes, markers, and composition that was highly consistent across individuals and in cross-species comparisons (**Supplementary Fig. 3-S13**)<sup>190,191,202</sup>.

In the neocortex, we captured transcriptionally distinct near- (L5/L6 NP) and far-projecting (L5 ET: extra-telencephalic) excitatory neurons, the latter of which contained a much higher number of expressed genes and transcripts per cell than most other excitatory subtypes (**Supplementary Fig. 3-S9**), possibly related to the pronounced differences in dendritic morphology and axonal projection between L5 ET and IT neurons<sup>204</sup>. Genes distinguishing the extra-telencephalic projection neurons from near-projecting neurons were involved in neuron projection development and included adhesion-related FAT3 and KIRREL3, neurite-outgrowth genes NRG1 and NTNG1, and axonal navigation-mediators ROBO2 and SLIT2 (**Supplementary Fig. 3-S10a, Supplementary Table 4**). Scoring each neuronal cell based on these genes yielded low scores for inhibitory subtypes, which are largely locally projecting<sup>205</sup>, and the second highest scores for subiculum neurons, which sometimes form long-range projections to the mammillary body (**Supplementary Fig. 3-S10b,c**)<sup>206</sup>. Interestingly, we predicted POU3F1 as a transcriptional regulator of both L5 ET and subiculum programs, suggesting that it may play a role in the formation of long-range projections (**Supplementary Fig. 3-S10d,e**).

To uncover signatures controlling excitatory subtype functions, we identified 18 expression programs, most of which were subtype-specific and overlapped with marker genes, with the exception of subtype-shared program 18 which included immediate early genes NPAS4, FOS, FOSB, EGR1, ARC, and JUNB, possibly encoding a response to neuronal activity (**Supplementary Fig. 3-S12e,f, Supplementary Table 2**). The major thalamic subtype (NXPH1 RNF220) specifically used programs 14 and 15, which included many regulators of intra- and intercellular signaling, such as genes involved in calcium binding (CALB1/2, CASQ2, FSTL4, PLCB4), axon guidance (NTNG1, GDNF, NOTCH2, SEMA3F), and in neurotransmitter reception (GRIN2D, CHRNA2, ACHE, DRD2, OPRK1, OPRM1, HTR7, HRH3) (**Supplementary Fig. 3-S12e,f**). Using SCENIC, we predicted thalamic transcriptional regulators LHX9, SOX2, SHOX2, and TCF7L2, each whose expression and regulon was significantly higher in the thalamus than in other subtypes across individuals (**Supplementary Fig. 3-S10e-i**)<sup>187</sup>. We validated our results by independently analysing microarray data from different human brain regions<sup>33</sup>. Supporting our single-cell analysis, we found that 15 out of the 20 aforementioned marker genes of NXPH1 RNF220 neurons were expressed at a significantly higher level in the thalamus compared to the cerebral cortex (**Supplementary Fig. 3-S13**). Furthermore, integration of our data with previously published single cell data from different regions of the mouse brain<sup>190,191,202</sup> revealed that the identified separation between thalamic and cortical excitatory neurons in the human brain is evolutionarily conserved (**Supplementary Fig. 3-S13**). Overall, NXPH1 RNF220 neurons express a unique combination of genes involved in intra- and intercellular signaling that distinguish them from previously described cortical excitatory neurons and may be important for proper input and output processing of the anterior thalamus in the context of reciprocal connectivities with various regions such as frontal cortical areas, hippocampal subregions, and hypothalamic nuclei<sup>207</sup>.

### **3.6 A unique thalamus-specific inhibitory neuron subtype**

We next examined whether inhibitory neurons show the same extent of diversity across brain regions as excitatory neurons, annotating 23 inhibitory subtypes across regions with subtype-distinguishing marker genes (**Fig. 3-1h, Supplementary Fig. 3-S13-19**). The large majority of inhibitory subtypes were observed in all five cortical regions, in contrast to excitatory



neurons, indicating a high level of cortical conservation (**Fig. 3-1h, Supplementary Fig. 3-S15a-i**). Despite their presence across the cortical regions, some inhibitory subtypes showed region-skewed composition, with many subtypes enriched in the neocortex, including PVALB HTR4 and CUX2 MSR1. In contrast, SST subtype L6 SST NPY was predominantly found in the EC and HC and EC-enriched PVALB subtype GPC5 RIT2 was marked by GPC5 as with multiple EC excitatory subtypes, suggesting that there are significant differences in inhibitory neuron composition between the neocortex and evolutionarily older structures of the cortex (**Fig. 3-1h, Supplementary Fig. 3-S15i-o**). Additionally, in the HC, EC, and MT, medial (MGE) ganglionic eminence-derived neurons (SST and PVALB) were significantly more abundant than caudal (CGE) ganglionic eminence-derived GABAergic neurons (VIP and LAMP5), but these two major clades were not significantly different in the PFC (**Supplementary Fig. 3-S16l**).

While the majority of inhibitory subtypes were present across the five cortical regions, the same subtypes were nearly absent in the anterior thalamus, which instead contained a unique, thalamus-specific GABAergic neuron type marked by MEIS2 and FOXP2 (**Fig. 3-1h,l**). Genes associated with the thalamic inhibitory neuron program (Inh GEP-22) included genes involved in neurite outgrowth such as the semaphorins SEMA3C and SEMA3E, DISC1, and SPON1 and receptors for serotonin (HTR2A), choline (CHRM2), acetylcholine (CHRNA3), and glutamate (GRM3) (**Fig. 3-1l-m, Supplementary Fig. 3-S15-16, Supplementary table 2**). The program also comprised many transcription factors, including MEIS2, FOXP2, LEF1, and OTX2, which overlapped with SCENIC in predicting LEF1 and FOXP2 as regulators for the thalamic population (**Supplementary Fig. 3-S16g**), supporting the subtype's uniqueness in transcriptional determination, receptors, signaling, and adhesion relative to other inhibitory subtypes. To validate the localization and specificity of markers of the thalamic inhibitory neuron subtype we performed *in situ* hybridization for both FOXP2 and MEIS2 with GAD2 on anterior thalamus and prefrontal cortex *post-mortem* brain samples from four individuals from our cohort, and found significant thalamus-specific co-localization of both marker genes with GAD2 (**Fig. 3-1n**).

To determine whether thalamic differences are conserved evolutionarily, we integrated our data with previously published single cell data from different regions of the mouse brain and the human lateral geniculate nucleus (dLGN)<sup>190,191,201,202</sup>. Human and mouse thalamic inhibitory neurons clustered together and separately from neocortical inhibitory neurons, and multiple subtype markers were recapitulated both in dLGN interneurons and in mouse thalamic inhibitory neurons, suggesting that thalamic-cortical inhibitory neuron differences share similar transcriptional programs with mice and in at least one other thalamic region (**Supplementary Fig. 3-S18**).

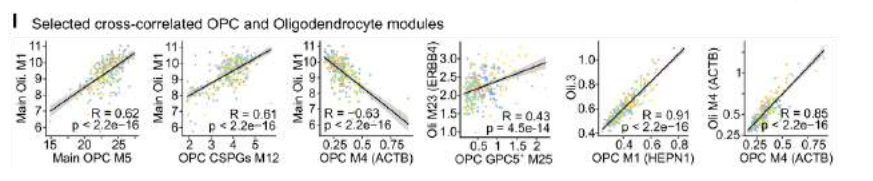
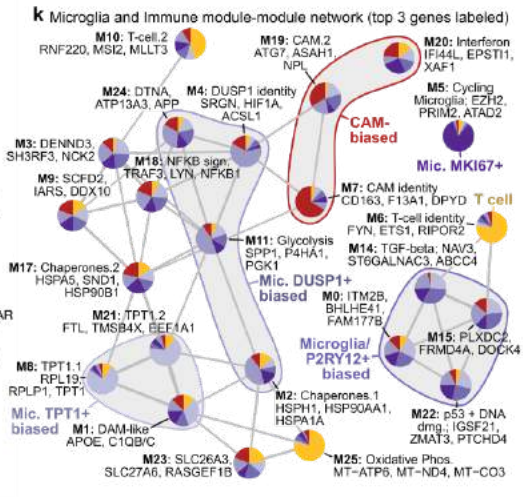
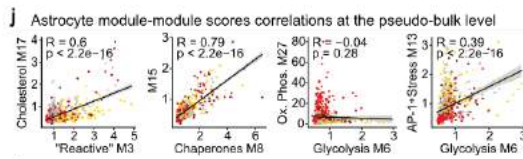
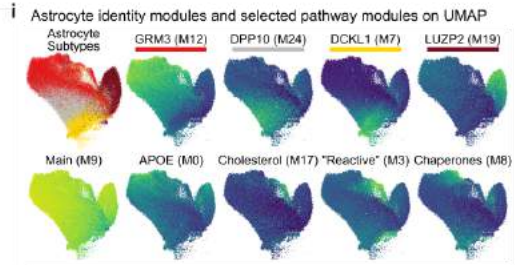
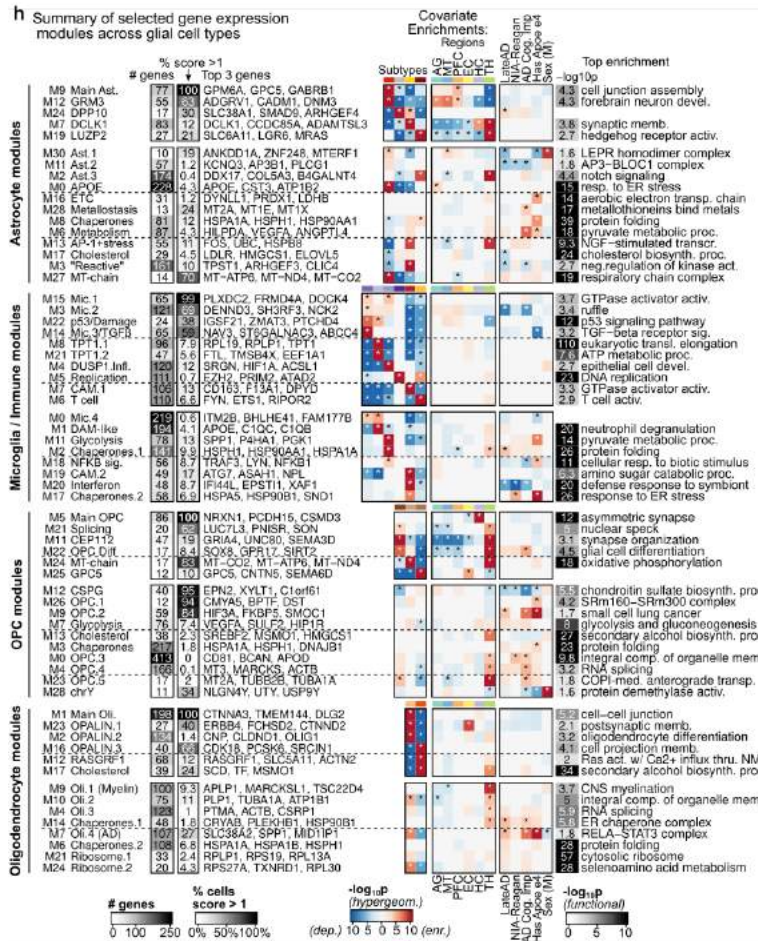
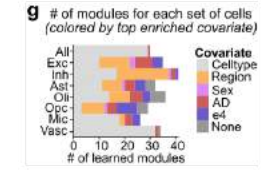
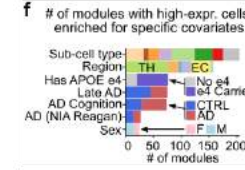
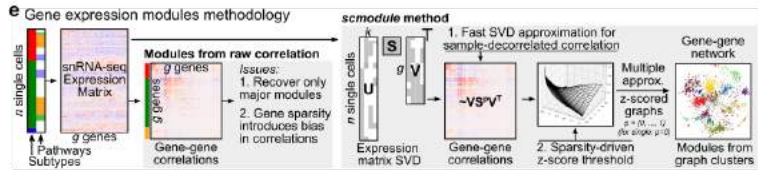
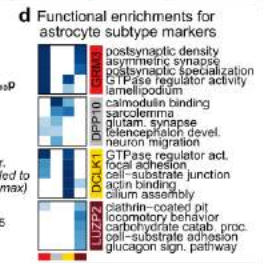
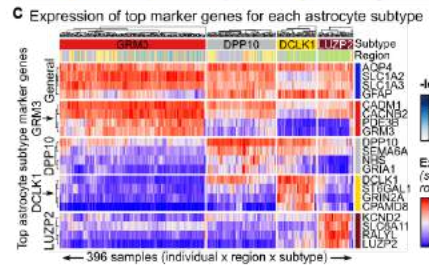
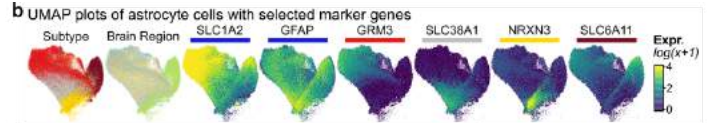
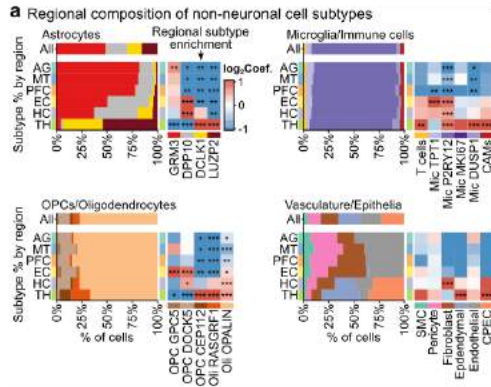
As thalamic MEIS2 neurons surprisingly expressed several typically glutamatergic neuron genes, such as CAMK2D, RORB, and SNAP25, we determined glutamatergic- and GABAergic module scores for every neuronal cell to further explore the chimeric nature of this subtype (**Supplementary Fig. 3-S15-16, Supplementary Table 4**). As expected, in the cortex excitatory neuron subtypes exhibited high glutamatergic scores and low GABAergic scores and inhibitory subtypes showed the inverse pattern (**Fig. 3-1k**). GABAergic and glutamatergic scores were significantly negatively correlated both across all subtypes (**Fig. 3-1k**) and within excitatory and inhibitory neurons, suggesting a common regulatory mechanism controlling the activity of the two gene expression programs across neuronal subtypes (**Supplementary Fig. 3-S15h-j**). Interestingly, both thalamic MEIS2 inhibitory and NXPH1 excitatory neurons fell between the cortical excitatory and inhibitory clusters, both in module space and principal component space (**Fig. 3-1k, Supplementary Fig. 3-S15a-c**), suggesting that they are less polarized than their cortical counterparts with regards to the expression of cortical glutamatergic versus GABAergic programs.

### 3.7 Glial subtype, regional, and functional diversity is annotated by gene expression modules

Given the striking compositional differences for neuronal subtypes across brain regions, we asked whether glial cells exhibited similar differences across regions. We identified multiple transcriptionally distinct subsets for each major non-neuronal cell type and determined their characteristic marker genes (**Fig. 3-2a, Supplementary Fig. 3-S19-24**). Among glial cell types, astrocytes had the highest regional heterogeneity, containing both highly neocortex- (GRM3 and DPP10) and thalamus-enriched subtypes (LUZP2 and DCLK1) (**Fig. 3-2a-b, Supplementary Fig. 3-S19**). Almost all astrocytes expressed AQP4, SLC1A2, and SLC1A3, but cortical astrocytes were enriched for glutamate processing and transport, with the GRM3 subtype expressing GRID1, GRIA2, and GLUL and the DPP10 subtype expressing GRIA1. Hippocampus and thalamus-enriched DCLK1 astrocytes had lower glutamate transporter activity, were enriched for focal adhesion, and uniquely expressed both glutamate receptor GRIN2A and cell-cell adhesion gene NRXN3 (as opposed to neocortical NRXN1 and CADM1) (**Fig. 3-2a-d, Supplementary Fig. 3-S19**). Thalamic LUZP2 astrocytes expressed GABA uptake genes SLC6A1 and SLC6A11 at much higher levels than other subtypes, even though the proportion of inhibitory neurons was not markedly higher in the thalamus. However, the unique thalamic MEIS2/FOXP2 interneurons shared multiple markers with thalamus-absent neocortex GRM3 astrocytes, including GRM3, MEIS2, and VAV3, also top markers for mouse thalamus-cortical interneuron differences, suggesting that astrocytes in evolutionarily newer regions share some functions with inhibitory neurons in older regions (**Supplementary Fig. 3-S22**)<sup>190,191,202</sup>.

To systematically identify gene expression programs driven by factors such as subtype identity, brain region, age of individual, and AD status at multiple resolutions, we developed a method to identify gene expression modules using gene-gene correlations in atlas-scale snRNA-seq datasets (**Fig. 3-2e, Supplementary Fig. 3-S25**). Compositional effects dominate gene-gene correlations in mixed-cell and subtype datasets, and naive methods to call modules from this signal only recover the most common and high-expression gene expression modules. To account for these effects, instead of computing correlations directly from the expression matrix, we efficiently approximate a sample-decorrelated gene-gene correlation matrix, in a similar manner to zero-phase component analysis<sup>208</sup>. To correct for inflated estimates in pairs of high-sparsity genes, we use bivariate splines to apply a consistent z-score cutoff across tested gene-gene pairs based on their sparsity. Finally, we build a graph using the z-score thresholded matrix as an adjacency matrix and use leiden clustering to identify gene modules<sup>195</sup>.

We used our method to identify modules both across all cells in the atlas and for each major cell type separately (each with an average of 33 modules with  $\geq 10$  genes) and recovered modules expressed to varying degrees, ranging from an adhesion module found in all astrocytes to a cell cycle module in just 0.7% of microglia (**Supplementary Fig. 3-S25-36, Supplementary Table 5**). Learned modules captured diverse aspects of cell identity, enriching for subtype identity (205 modules), brain region (156, with 77 thalamus and 34 entorhinal cortex-specific), AD status (73), and covariates such as APOE genotype (78) and sex (24) (**Figures 2f-g, Supplementary Table 5**).



**Figure 3-2. Glial diversity + modules.** **a.** Global breakdown, region composition, enrichment, and number of cells for non-neuronal subtypes across brain regions. **b.** UMAP plots for astrocyte cells by astrocyte subtype, region of origin, and expression of selected marker genes. **c.** Gene expression heatmap for top four markers per astrocyte subclass and of selected broad astrocyte markers (blue), averaged at the sample by subclass level (columns). **d.** Functional enrichments for marker genes of each astrocyte subclass (adj. p-value < 0.01). Top unique enrichments shown (1-2 subclasses, terms with < 500 genes). **e.** Gene expression modules framework. **f.** Number of modules enriched for each covariate across all module sets (hypergeometric test, p-value < 0.001). Barplots are colored by the covariate level for which the modules are enriched (or by the modules set of origin for cell subtype). **g.** Number of learned modules in each set of modules, colored by their top enriched covariate. **h.** Overview of selected gene expression modules across major glial cell types, showing the module name, number of genes, percent expression, top module genes, enrichments by subtype and region (hypergeometric test), and the top functional enrichment for each module. **i.** Localization of astrocyte identity and pathway modules on the astrocyte UMAP. **j.** Scatterplots and correlation of scores for selected pairs of astrocyte modules. Each dot represents the module expression scores for a subtype in a specific sample and is colored by the astrocyte subtype. **k.** Microglial and immune modules network from correlation of module pairs at the subtype by sample level (edges shown where FDR-adjusted p-value < 0.05). Nodes are colored by module's relative expression in each of the microglial and immune subtypes and groups highlight sets of subtype-biased modules. **l.** Scatterplots and correlation of scores for selected pairs of OPC and oligodendrocyte modules. Each dot represents scores at the sample-level and is colored by region of origin.

We identified 32 modules in astrocytes, many specifically enriched for brain region and subtype (**Fig. 3-2f, Supplementary Fig. 3-S26, 3-S30**). These included a broadly-active astrocyte program (M9) with GPM6A and GPC5, enriched for cell junction assembly, and subtype and region-specific identity programs such as LUZP2 and thalamus-associated M19 (SLC6A11, LGR6, MRAS), which enriched for hedgehog signaling, M12 (GRM3: forebrain neuron development), and M7 (DCLK1: synaptic membrane) (**Fig. 3-2g**). Other modules spanned a diverse set of functions, including metallostasis, RNA splicing, glycolysis, oxidative phosphorylation, and cholesterol biosynthesis (**Fig. 3-2f,g**). Module-module correlations across samples revealed co-expressed programs, such as reactive astrocytes (M3, marked by TPST1, CLIC4, and EMP1) with cholesterol biosynthesis module M17 ( $r=0.60$ ), and glycolysis (M6) with AP-1 module M13 ( $r=.39$ , including FOS/JUN and ubiquitin), a pair potentially co-expressed in astrocytes under metabolic stress (**Fig. 3-2g,h, Supplementary Fig. 3-S35a**).

In contrast, microglia and other immune cells, which included T cells (2.7% of immune cells) and macrophages (CAMs, 3.9%) in addition to four microglia subtypes, showed little regional specificity (**Fig. 3-2a, Supplementary Fig. 3-S21a-g, 3-S22-24, 3-S26, 3-S31**). Immune modules contained identity programs for T cells (M6: FYN and ETS1), macrophages (M7: CD163 and F13A1), and DNA replication (M5) for cycling MKI67+ microglia (0.3%). We often identified multiple highly subtype-specific programs for microglia, including proinflammatory (M4: SRGN, ACSL1, and hypoxia factor HIF1A), glycolysis (M11), and unfolded protein response (M2) modules in DUSP1+ microglia (4.1%). Similarly, TPT1+ microglia (5.4%) were marked by modules enriched for translation (M8), ATP metabolism (M21), and a disease associated microglia, or DAM-like, signature (M1: APOE, C1Q, and MHC-II) (**Fig. 3-2i**)<sup>49</sup>. In addition, we identified multiple pathway-driven modules present across immune cells, including NF-KB signaling (M18), interferon activation (M20), and modules enriched for p53 signaling and DNA damage response (M22), and TGF-beta signaling (M14) across microglia (**Fig. 3-2i**,



### Supplementary Fig. 3-S31).

For OPCs and oligodendrocytes, we identified thalamus-enriched subtypes (OPC CEP112 and Oli RASGRF1) with minor transcriptomic subtype differences to neocortex-enriched subtypes (OPC DOCK5/GPC5 and Oli. OPALIN) (**Fig. 3-2g, Supplementary Fig. 3-S20, 3-S22-24, 3-S26, 3-S29-30, 3-S14a-f**). Unlike in immune cells, oligodendrocyte-lineage modules showed high regional specificity, including thalamus-enriched modules for oligodendrocyte myelination (M9) and OPC differentiation (M22) (**Fig. 3-2f**). Interestingly, two OPC modules, thalamus and CEP112-enriched M11 and entorhinal cortex and GPC5-enriched M25, were marked by synapse-associated genes such as neural adhesion-related SEMA3D, SEMA6D, and CNTN5, and glutamate receptor GRIA4, suggesting a role for OPC sensing and response to neuronal activity in specific brain regions (**Fig. 3-2f**). Modules unique to these cell types also included OPC chondroitin sulfate biosynthesis M12 (marked by XYLT1); oligodendrocyte selenoamino acid metabolism M24; and OPC glial differentiation-enriched M22, which was marked by oligodendroglial remodeling factors SOX8 and SIRT2 as well as GPR17, possibly guiding OPCs to remodel into oligodendrocytes when myelin sheath stress is detected<sup>209,210</sup> (**Fig. 3-2f**).

Higher-order module-module correlations across the cell types included modules associated with specific sample attributes (e.g. AD, region) but also co-expressed modules with overlapping gene sets encoding similar functions. We hierarchically clustered modules across the cell types and identified broad shared programs for cholesterol biosynthesis (C10), chaperones (C5), ribosomes (C1, C2), ER protein processing (C7), oxidative phosphorylation (C18), synapse interaction (C16), and glycolysis and response to hypoxia (C20) (**Supplementary Fig. 3-S36a-b**). The oligodendrocyte-lineage cell types showed particularly cross-correlated module activity, including the main oligodendrocyte module (M1) with the OPC chondroitin sulfate module (M12,  $r=0.61$ ) and the main OPC module (M5,  $r=0.62$ ) with which it shared neuron-interacting genes including NCAM1/2 and DLG2 with M5 (**Fig. 3-2i, Supplementary Fig. 3-S35b-c, 3-S36c**). Interestingly, OPC M4 and oligodendrocyte M4 ( $r=0.83$ ) were highly correlated, shared 40+ genes (including actins ACTB/ACTG1 and translation factors EIF4A1/EIF5A1/EEF2), and enriched for mRNA binding, nuclear speckles, and stress granules, but were strongly anti-correlated with the main oligodendrocyte module (OPC-M4 to Oli-M1,  $r=-0.63$ ), which enriched for axonogenesis, cell-cell junctions, and nodes of ranvier, suggesting that while oligodendrocyte M1 represents homeostatic oligodendrocytes and neuronal adhesion, the speck-associated M4 modules may mark a different, potentially dysfunctional state of oligodendroglial cells (**Fig. 3-2i, Supplementary Fig. 3-S32-33, 3-S35b-c, 3-S36c**)<sup>211</sup>.

### 3.8 Methodology for gene expression and regulon modules

**Gene expression modules using ZCA (scmodule):** We would like to find gene expression modules by calculating gene-gene correlations in single-cell data and using these to detect communities of similarly expressed genes. However, in single-cell data, which often contains an unbalanced composition of cell types, modules computed using this approach will be dominated by the most common cell type markers and pathways. In addition, correlation values will often be inflated for pairs of sparsely expressed genes. We developed a method which accounts for these pitfalls to call multi-resolution gene expression programs from single-cell data using an SVD-based approximation of zero-phase component analysis (ZCA) and gene sparsity-dependent thresholding. Our method is available as the scmodule python library from <http://github.com/KellisLab/scmodule/>.

**ZCA for single-cell data:** The preprocessing transformations alternately called decorrelation, whitening, or sphering transform a matrix  $\mathbf{X}$  with a matrix  $\mathbf{W}$  such that the covariance of  $\mathbf{XW}$  is the identity matrix<sup>212</sup>. In particular, zero-phase component analysis (ZCA) is the transformation

which maximizes the similarity of the transformed data to the original, which is achieved by setting  $\mathbf{W} = \mathbf{C}^{-1/2}$ , where  $\mathbf{C}$  is the covariance of  $\mathbf{X}$ .

In single-cell data, given a count matrix  $\mathbf{X}$  with  $n$  cells (rows) and  $g$  genes (columns), we would like to perform ZCA decorrelation on the samples as a preprocessing step for calling modules. Computing and storing the  $n$  by  $n$  sample-wise covariance  $\mathbf{C}_n = \mathbf{X}\mathbf{X}^T / g$  is prohibitively expensive for modern datasets (with  $n > 1e6$ ), even without centering  $\mathbf{X}$ . Instead, we can analytically approximate the covariance with the SVD of  $\mathbf{X}^T = \mathbf{U}_n \mathbf{S}_n \mathbf{V}_n^T$  as  $\mathbf{C}_n \approx (\mathbf{U}_n \mathbf{S}_n \mathbf{V}_n^T)^T (\mathbf{U}_n \mathbf{S}_n \mathbf{V}_n^T) / g = \mathbf{V}_n \mathbf{S}_n^2 \mathbf{V}_n^T / g$  and therefore  $\mathbf{C}_n^{-1/2} = g^{1/2} \mathbf{V}_n \mathbf{S}_n^{-1} \mathbf{V}_n^T$ . The ZCA transformation can then be computed as  $\mathbf{X}_{ZCA} = \mathbf{C}_n^{-1/2} \mathbf{X} = g^{1/2} \mathbf{V}_n \mathbf{S}_n^{-1} \mathbf{V}_n^T \mathbf{X}$  before calculating the covariance of  $\mathbf{X}_{ZCA}$  for downstream analysis. While this approximation is already tractable for small single-cell datasets, we may not be able to compute the matrix multiplications or centered SVD for larger datasets. Here, we can use the SVD of  $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ , which is commonly calculated in single-cell analyses, to approximate  $\mathbf{C}_n^{-1/2}$  as  $g^{1/2} \mathbf{U}\mathbf{S}^{-1} \mathbf{U}^T$  and  $\mathbf{X}_{ZCA} = g^{1/2} \mathbf{U}\mathbf{S}^{-1} \mathbf{U}^T \mathbf{X}$ . From this, the non-centered covariance of  $\mathbf{X}_{ZCA}$  is  $\mathbf{C}_{ZCA} = \mathbf{X}_{ZCA}^T \mathbf{X}_{ZCA} / n = g * (\mathbf{U}\mathbf{S}^{-1} \mathbf{U}^T \mathbf{X})^T (\mathbf{U}\mathbf{S}^{-1} \mathbf{U}^T \mathbf{X}) / n$ . By substituting the SVD in for  $\mathbf{X}$ , this reduces to  $\mathbf{C}_{ZCA} = g * \mathbf{V}\mathbf{V}^T / n$ , which is a simple and efficient approximation for very large single-cell data. Because this approximation commonly drops out the largest identity program in the data due to the decorrelation approach, we allow computing  $\mathbf{C}_{ZCA} = g * \mathbf{V}\mathbf{S}^P \mathbf{V}^T / n$ , for any  $P$ , in order to tune the relative involvement of the larger eigenvalue components of the SVD.

**Thresholding correlation and computing modules:** To control for inflated correlation estimates in highly sparse genes, we bin the correlations for each pair of genes according to their sparsities ( $\log_{10}$  fraction of cells expressing the gene). We calculate the mean and standard deviation for each 2D bin and smooth the estimates by fitting bivariate splines to the binned statistics, weighted by the log number of examples in the bin. We use the smoothed estimates to z-score the correlation matrix, which we threshold with a single z-score cutoff to create an adjacency matrix for a gene-gene graph. Graphs are laid out using the Fruchterman-Reingold algorithm and we remove connected components with fewer than 4 genes<sup>182</sup>. We then use the leidenalg package and the Leiden algorithm with an RBCConfiguration vertex partition to cluster the graph into gene modules<sup>195</sup>. To robustly estimate modules for each set of cells in our analyses, we first performed a grid search for the optimal number of SVD components for cells of that type. We then computed the z-scored matrices for each of 10 bootstraps, selecting 90% of batches for each bootstrap and using only genes expressed in over 5% of cells in the full dataset for the cell type. We thresholded the average of the bootstrapped z-score estimates with  $z = 4.5$  to build a graph. In order to balance the contributions of modules across the compositional spectrum, we calculated and thresholded separate graphs for eigenvalue powers  $P = 0, 0.25, 0.5, 0.75, \text{ and } 1$  and combined them using multi-graph leiden clustering to call modules with leiden resolution=3. Although we sometimes identify smaller modules, here we only report modules with at least 10 genes.

**Module enrichments and network:** Module enrichments for cell subtypes and brain regions were performed using the hypergeometric test by calculating whether cells with a module score above one standard deviation from the mean score were significantly enriched in a specific subtype or region. For plotting scores against other modules, averaged modules scores to either the subtype by sample level (within the same major cell type) or at the sample level alone (across cell types) and calculated Pearson correlations and p-values using R's cor.test function. To create the module-module network across microglial and immune modules, we calculated module-module score Pearson correlations using the logged module scores at the aggregated subtype by sample level, using a one-sided test with an adjusted p-value  $< 0.01$  as a cutoff<sup>ff213</sup>.

**Gene expression programs using consensus non-negative matrix factorization (cNMF):**

Gene expression programs underlying both cell-type identity and cellular activities were determined following the consensus non-negative matrix factorization (cNMF) analysis pipeline established by Kotliar et al. using default parameters<sup>184</sup>. The number of components (K) to use for cNMF was determined based on a diagnostic plot showing the stability of the solution and the Frobenius reconstruction error as a function of K. To reduce runtime and working memory requirements the data were downsampled using the Seurat function “subset” with default parameters. The data were downsampled to 200 cells per major cell type. For the cNMF analysis at the level of high resolution cell types, the analysis was performed separately on excitatory neurons, inhibitory neurons, and astrocytes. For these analyses the data were downsampled to 2000 cells per astrocyte subtype and 1000 cells per excitatory- and inhibitory neuron subtype. Statistical significance of the overlap between the top 200 genes of a gene expression program and cell type marker genes was computed using Fisher’s exact tests.

**SCENIC analysis and computation of regulon module scores:** The gene regulatory network analysis was performed using pySCENIC with default parameters<sup>186</sup>. To reduce runtime and working memory requirements the data were downsampled to 2000 cells per major cell type. For the SCENIC analysis at the level of high resolution cell types, the analysis was performed separately on excitatory neurons, inhibitory neurons, and astrocytes and the data were downsampled to 1000 cells per high resolution cell type. To identify the top cell-type specific regulons, we calculated regulon specificity scores as described by Suo et al. and ranked the regulons based on their regulon specificity score<sup>214</sup>. Finally, we calculated the “activity” of each regulon in each cell using the AddModuleScore function of Seurat. The calculation of regulon module scores for major cell types was performed on a random sample of 50% of the cells (676,537 cells). For the analysis at the level of high resolution cell types, the regulon module scores were determined based on all the cells of a major cell type class. For the statistical analysis of differences in the activity of regulons between cell types, the average regulon module score per individual and major cell type or high resolution cell type was computed, respectively. The statistical analyses comparing the regulon module score activity was performed using Prism 9 software.

**Analysis of GABAergic- and glutamatergic module scores:** GABAergic- and glutamatergic module scores across all neuronal cell types were determined based on a set of GABAergic- and glutamatergic neuron marker genes, respectively, using the AddModuleScore function of Seurat. The sets of GABAergic- and glutamatergic neuron marker genes were determined based on the Human Multiple Cortical Areas SMART-seq data set from the Allen Brain Institute (<https://portal.brain-map.org/atlas-and-data/rnaseq/human-multiple-cortical-areas-smart-seq>). We identified marker genes using the FindAllMarkers function from Seurat (Wilcoxon Rank Sum test with Bonferroni correction for multiple testing; adjusted  $P < 0.05$ ). We tested only genes that were detected in a minimum of 25% of the cells within the cell type (min.pct = 0.25) and that showed, on average, at least a 0.25-fold difference (log-scale) between the cells of the cell class of interest and all remaining cells (logfc.threshold = 0.25). To quantify the intermediate character of thalamic Exc and Inh neurons, we first computed the average GABAergic- and glutamatergic module score values for each neuronal cell type and for each individual. Then we used the resulting data to determine the first principal component (PC1) scores (the coordinates of the individual observations on the first principal component axis) using the princomp function in R. The ridgeline plot showing the distribution of PC1 score for each neuronal cell type was generated using the ggplot2 package in R. To determine the association between the average glutamatergic- and the GABAergic module score across neuronal cell types we performed a simple linear regression analysis using Prism 9 software.

**Analysis of extratelencephalic projection neuron module scores:** Marker genes

significantly upregulated in extratelencephalic projection neurons (Exc L5 ET) compared to near-projecting excitatory neurons in layers 5 and 6 (Exc L5/6 NP) were determined using the FindAllMarkers function from Seurat (Wilcoxon Rank Sum test with Bonferroni correction for multiple testing; adjusted  $P < 0.05$ ). We tested only genes that were detected in a minimum of 25% of the cells within the cell type (min.pct = 0.25) and that showed, on average, at least a 0.25-fold difference (log-scale) between Exc L5 ET cells and Exc L5/6 NP cells (logfc.threshold = 0.25). The Exc L5 ET module score was computed based on the identified marker genes using the AddModuleScore function of Seurat. To determine the Exc L5 ET module scores across excitatory and inhibitory neurons, cells were downsampled to 2000 cells per high resolution cell type.

### 3.9 Discussion

In this work, we presented a comprehensive atlas of the aged human brain across brain regions impacted along Alzheimer's disease (AD) progression, comprising over 1.3 million single nuclei in six brain regions from each of 48 individuals along a range of AD progression. Our analysis of multiple brain regions from a large cohort of individuals enabled us to annotate cellular regional diversity, identify a novel neuronal subtype and region-specific glial subtypes, develop a method to detect functional expression modules in atlas-scale snRNA, track regional transcriptomic changes across early-to-late AD, pinpoint vulnerable region-specific excitatory neuronal cell types depleted along AD progression, and propose molecular and connectome pathways mediating excitatory neuron vulnerability and resilience to AD, gaining novel insight into fundamental aspects of the human brain and how it is affected in AD.

We systematically annotated neuronal and glial subtypes across regions with markers, expression programs, and predicted transcriptional regulators. We found significant compositional heterogeneity across regions, including a glial-dominated cellular architecture in the thalamus<sup>188,189</sup>, an increased number of medial (SST and PVALB) over caudal ganglionic eminence GABAergic neurons in the hippocampus, entorhinal cortex, and midtemporal cortex, and highly region-specific excitatory neuron subtypes, in line with previous reports in mice<sup>190,191,202,215</sup> and in some independent regions in humans<sup>34,203</sup>, which our study extends to the human thalamus and across cortical regions. We discovered a type of thalamic GABAergic neuron (MEIS2/FOXP2) that is molecularly distinct from the canonical subclasses of inhibitory neurons and is highly region-specific, unlike other GABAergic neurons<sup>216</sup>. The developmental origin of these neurons remains to be established, but they may arise from the embryonic midbrain<sup>217</sup>, as in the case of recently reported mouse thalamic interneurons which share its expression of transcription factor SOX14 (Sox14). In contrast to other glial cell types that are shared across regions, we identify region specific astrocyte subtypes in the human brain, consistent with a recent study on the mouse cortex<sup>218</sup>. We also observe a spectrum of astrocyte cell states in the human cortex but find the most striking differences between neocortical- and thalamic astrocytes, which include a battery of genes distinguishing cortical astrocytes that center around their ability to sense, transport, and metabolize glutamate.

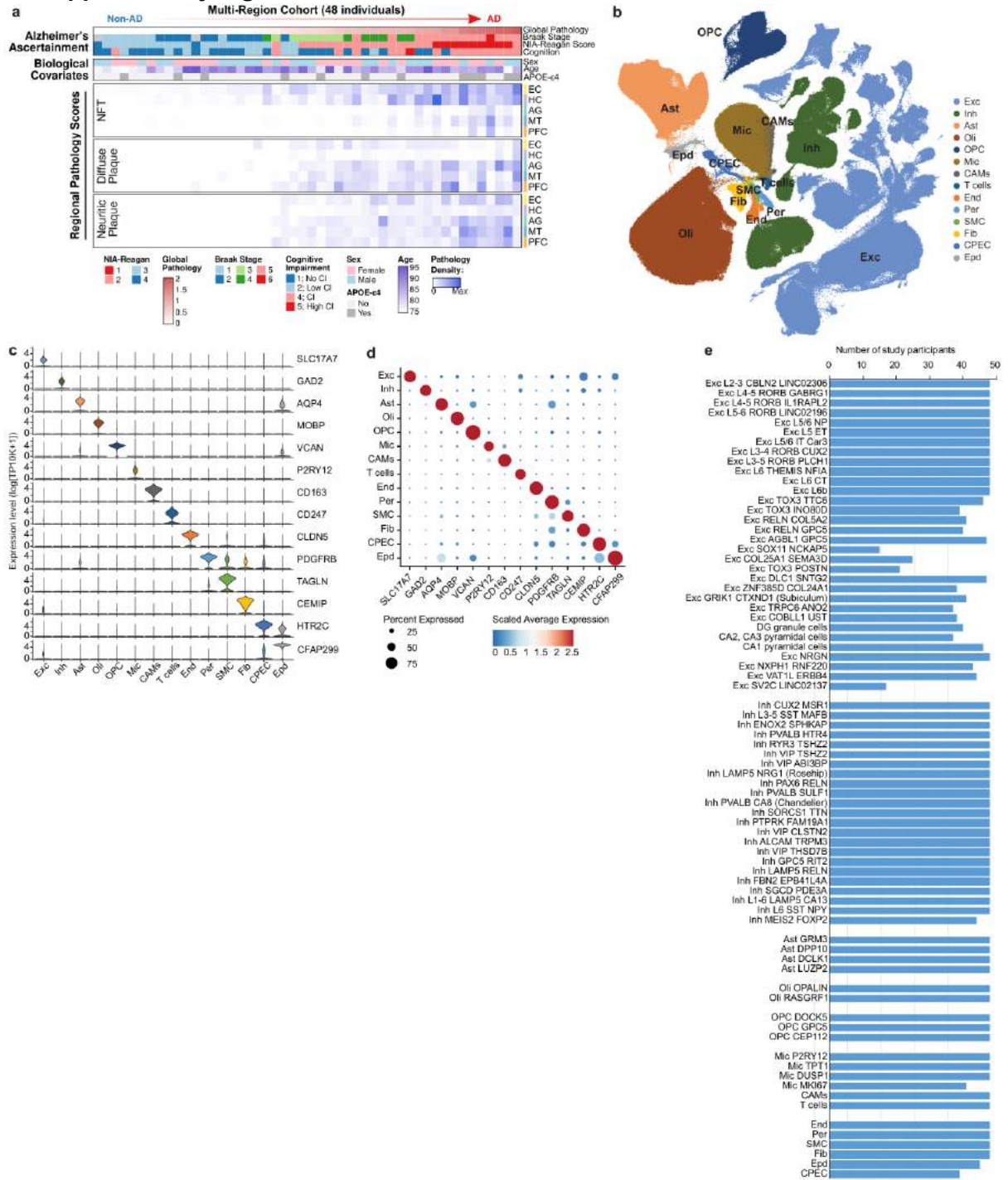
We developed a fast and scalable method to identify gene expression modules at multiple resolutions in atlas-scale snRNA-seq data that we used to annotate cellular functions across our atlas. Our analysis uncovered a huge diversity of cellular functions expressed at different levels of prevalence which highly co-varied within- and across-cell types and regions, including subtype specific proinflammatory and DAM-like microglial responses<sup>49</sup>, co-expression of astrocyte cholesterol synthesis with reactivity<sup>219</sup> and glycolysis with stress response factors, and many pairs of oligodendrocyte lineage modules, including a pair of possibly subtype-determining entorhinal cortex-enriched modules in progenitors and oligodendrocytes. In addition to examples outlined in the main text, we have created an interactive supplementary website for our study ([http://compbio.mit.edu/ad\\_multiregion](http://compbio.mit.edu/ad_multiregion)), allowing detailed exploration of annotations,



expression, functional modules, and AD transcriptomic changes to enable experimental hypothesis generation for targeted follow-up in the study of AD and brain architecture.

Our work provides a cross-cutting view of brain regional heterogeneity with particular applications in Alzheimer's disease across regions and disease timepoints, enabling many future studies: of neuronal development, specialization, and specification; application of our modules framework across scRNA-seq analyses and methodological development for other single-cell modalities; targeted experimental follow-up for molecular changes in Alzheimer's in each cell type; comprehensive joint investigation of Alzheimer's-affected circuitry and of determinants of neuronal vulnerability; and investigation of mechanisms of neuronal vulnerability and resilience, with applications in potential therapeutics, furthering our understanding of brain development, architecture, and changes in Alzheimer's disease progression.

### 3.10 Supplementary Figures



Extended Data Figure 1

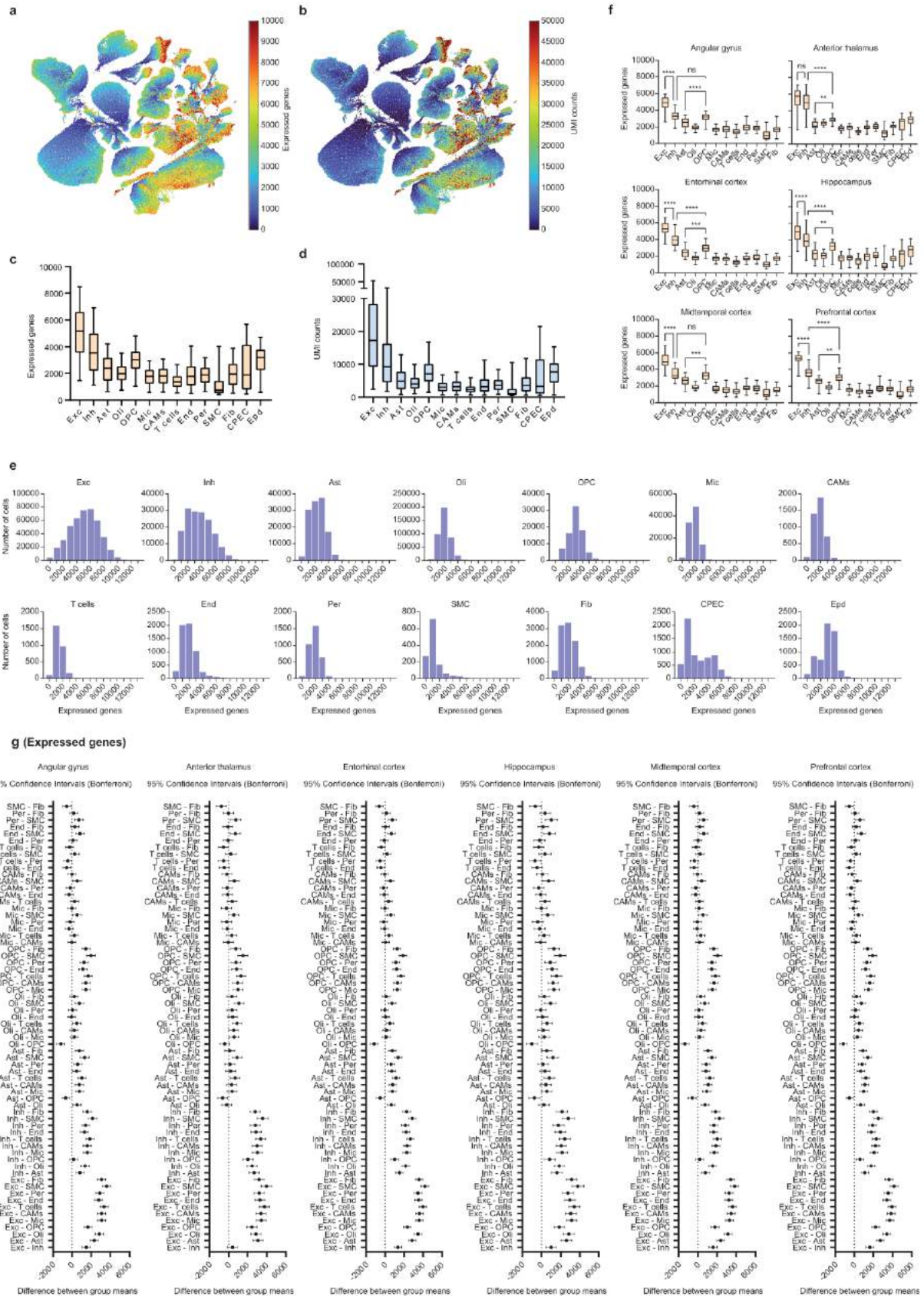
**Supplementary Figure 3-S1. Overview of the study cohort and major cell type annotations.** **a**, Cohort metadata overview. A total of 283 post-mortem brain tissue samples from 24 male and 24 female study participants were analysed across Alzheimer's disease

progression (AD). Top two panels show metadata at the individual level and bottom three panels show region-specific pathology measurements of neurofibrillary tangle burden (nft), neuritic plaque burden (plaq\_n), and diffuse plaque burden (plaq\_d). Individuals (columns) are ordered according to their global AD pathology burden. **b**, Uniform manifold approximation and projection (UMAP) plot of all annotated cells coloured by major cell type. **c**, Violin plot showing the expression of known cell type marker genes in each major cell type. **d**, Dot plot showing the scaled average expression of known cell type marker genes in each major cell type. **e**, Representation of individuals across cell types. The bar plot shows the number of study participants contributing at least one cell to the cell type indicated. Most cell type clusters contain cells from the majority of study participants.





**Supplementary Figure 3-S2. Cortical excitatory and inhibitory neuron subtypes align with previously identified neuronal subpopulations.** **a**, Spearman correlation between the average expression profiles of cortical excitatory neuron subpopulations previously defined by the Allen Brain Institute (rows) and the neocortical excitatory neuron subtypes identified in this study (columns). **b**, Projection of cortical excitatory neuron subpopulation annotations previously defined by the Allen Brain Institute (rows) onto the neocortical excitatory neuron subtypes identified in this study (columns). The heat map shows the proportion of cells (rows) that were labeled with the cell type alias label (columns) indicated. **c**, Glutamatergic neuron subclass module scores across cortical excitatory neuron cell types identified in this study (columns). The module scores were computed on the basis of the expression of glutamatergic neuron subclass marker gene sets. **d**, Module scores of intratelencephalic projection neuron subpopulations previously defined by the Allen Brain Institute (rows) across the corresponding cortical excitatory neuron cell types identified in this study (columns). **e**, Spearman correlation between the average expression profiles of inhibitory neuron subpopulations previously defined by the Allen Brain Institute (rows) and the inhibitory neuron subtypes identified in this study (columns). **f**, Projection of inhibitory neuron subpopulation annotations previously defined by the Allen Brain Institute (rows) onto the inhibitory neuron subtypes identified in this study (columns). The heat map shows the proportion of cells (rows) that were labeled with the cell type alias label (columns) indicated. **g**, GABAergic neuron subclass module scores across inhibitory neuron cell types identified in this study (columns). The module scores were computed on the basis of the expression of GABAergic neuron subclass marker gene sets. **h-k**, Module scores of GABAergic neuron cell types previously defined by the Allen Brain Institute (rows) across the corresponding inhibitory neuron cell types identified in this study (columns) ((**h**) LAMP5 subclass of GABAergic neurons; (**i**) SST subclass of GABAergic neurons; (**j**) VIP subclass of GABAergic neurons; (**k**) PVALB subclass of GABAergic neurons).



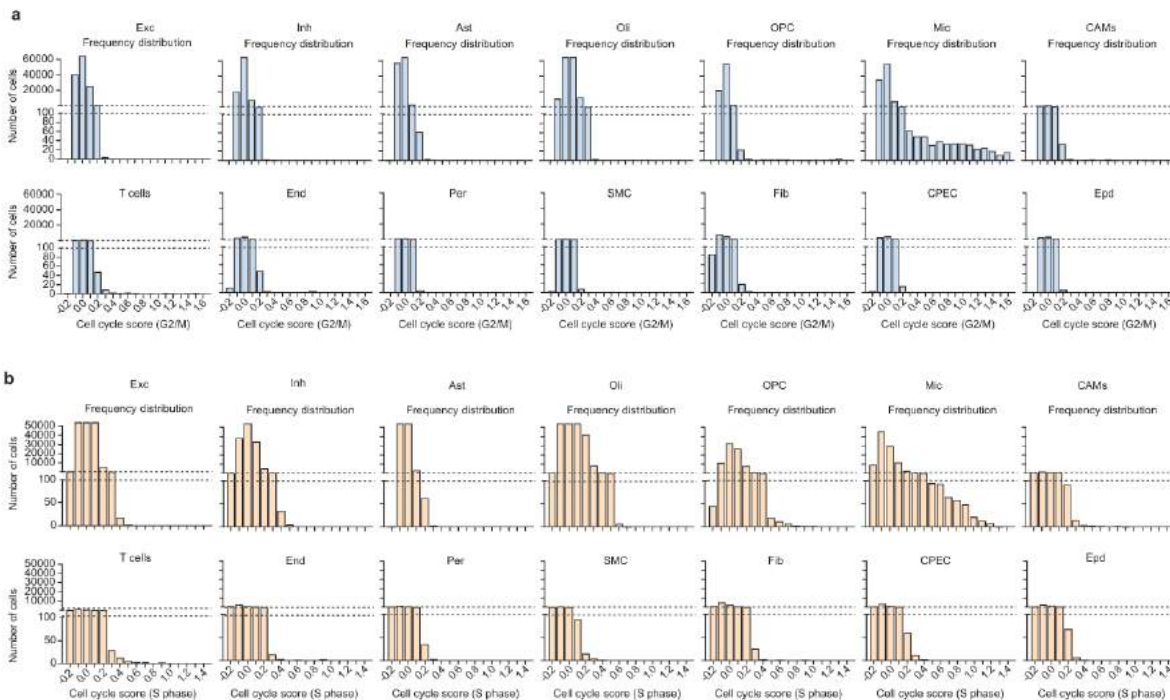
Extended Data Figure 3

**Supplementary Figure 3-S3. Global properties of gene expression: number of expressed genes and transcripts per cell.** **a-b**, UMAP plots of all annotated cells showing the distribution of **(a)** the number of genes detected per cell and **(b)** the number of unique transcripts (UMIs) detected per cell. **c-d**, Box plots of **(c)** the number of genes and **(d)** the number of unique transcripts (UMIs) detected per cell across all major cell types. Within each box, horizontal lines denote median values; boxes extend from the 25th to the 75th percentile of each group's distribution of values; whiskers extend from the 5th to the 95th percentile. **e**, Histograms showing the distribution of the number of genes detected per cell for each major cell type. **f-g**, Mean number of genes detected per cell per individual and major cell type across the six brain regions analyzed. **(f)** Box plots of the mean number of genes detected per cell per individual. Within each box, horizontal lines denote median values; boxes extend from the 25th to the 75th percentile of each group's distribution of values; whiskers extend from the 5th to the 95th percentile. \*\*\*\*P < 0.0001, \*\*\*P < 0.001, \*\*P < 0.01, \*P < 0.05; ns, P > 0.05 (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). **(g)** Difference between group means with 95% confidence intervals (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test).

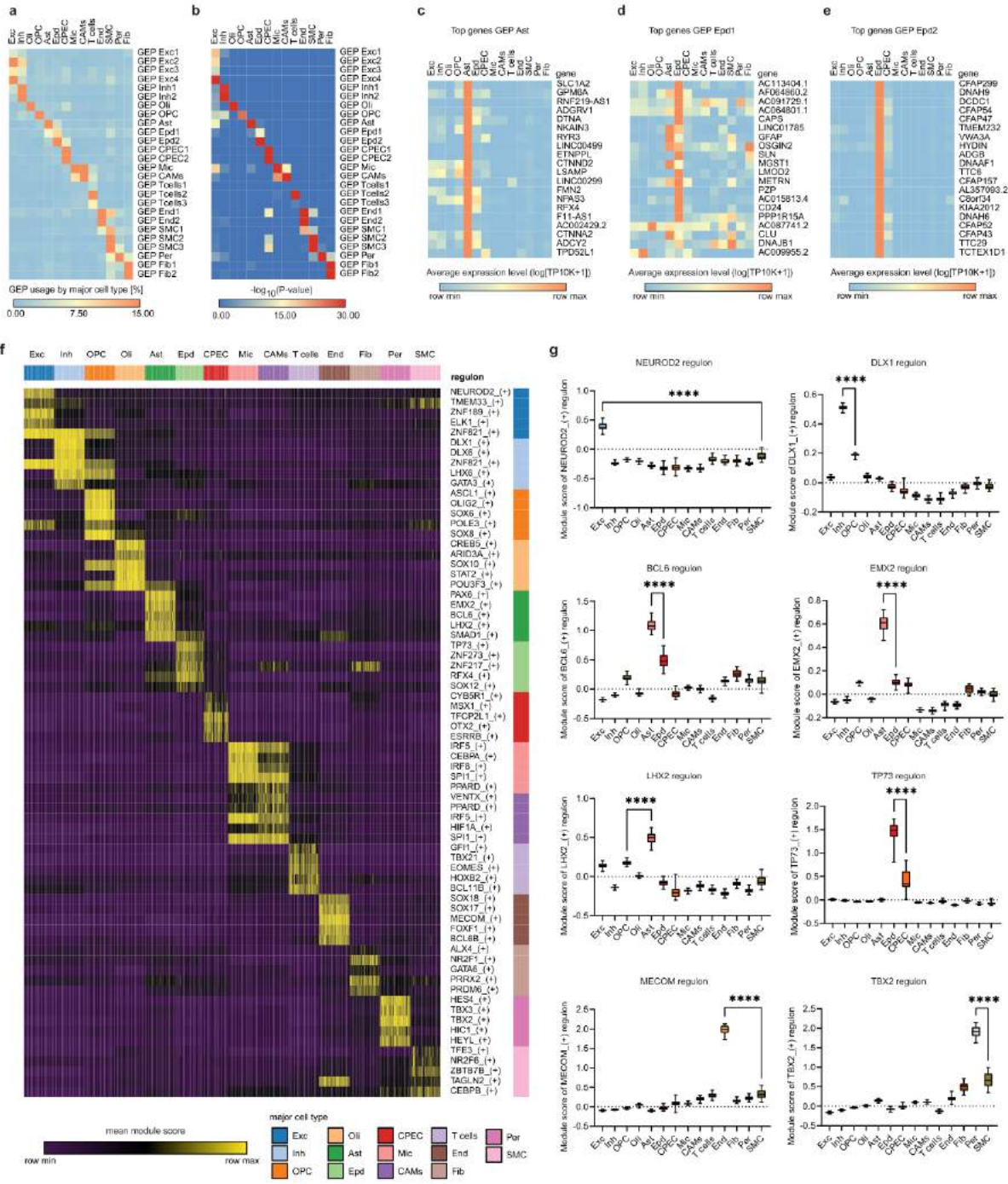




**Supplementary Figure 3-S4. Global properties of gene expression: number of transcripts per cell. a-b,** Mean number of unique transcripts (UMIs) detected per cell per individual and major cell type across the six brain regions analyzed. **(a)** Box plots of the mean number of unique transcripts detected per cell per individual. Within each box, horizontal lines denote median values; boxes extend from the 25th to the 75th percentile of each group's distribution of values; whiskers extend from the 5th to the 95th percentile. \*\*\*\*P < 0.0001, \*\*\*P < 0.001, \*\*P < 0.01, \*P < 0.05; ns, P > 0.05 (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). **(b)** Difference between group means with 95% confidence intervals (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). **c-e,** Mean number of unique transcripts (UMIs) detected per cell and major cell type in three previously published single cell studies of the mouse brain. Box plots show the mean number of unique transcripts detected per cell. Within each box, horizontal lines denote median values; boxes extend from the 25th to the 75th percentile of each group's distribution of values; whiskers extend from the 5th to the 95th percentile. \*\*\*\*P < 0.0001, \*\*\*P < 0.001, \*\*P < 0.01, \*P < 0.05; ns, P > 0.05 (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). The heat maps show the statistical significance (signed  $-\log_{10}(\text{Bonferroni-corrected P values})$ ) of the difference in the mean number of transcripts per cell between major cell types. The signed P value score is the  $-\log_{10}(\text{Bonferroni-corrected P value})$  with the sign of the difference (in the mean number of transcripts per cell) between the cell type indicated in the row minus the cell type indicated in the column.

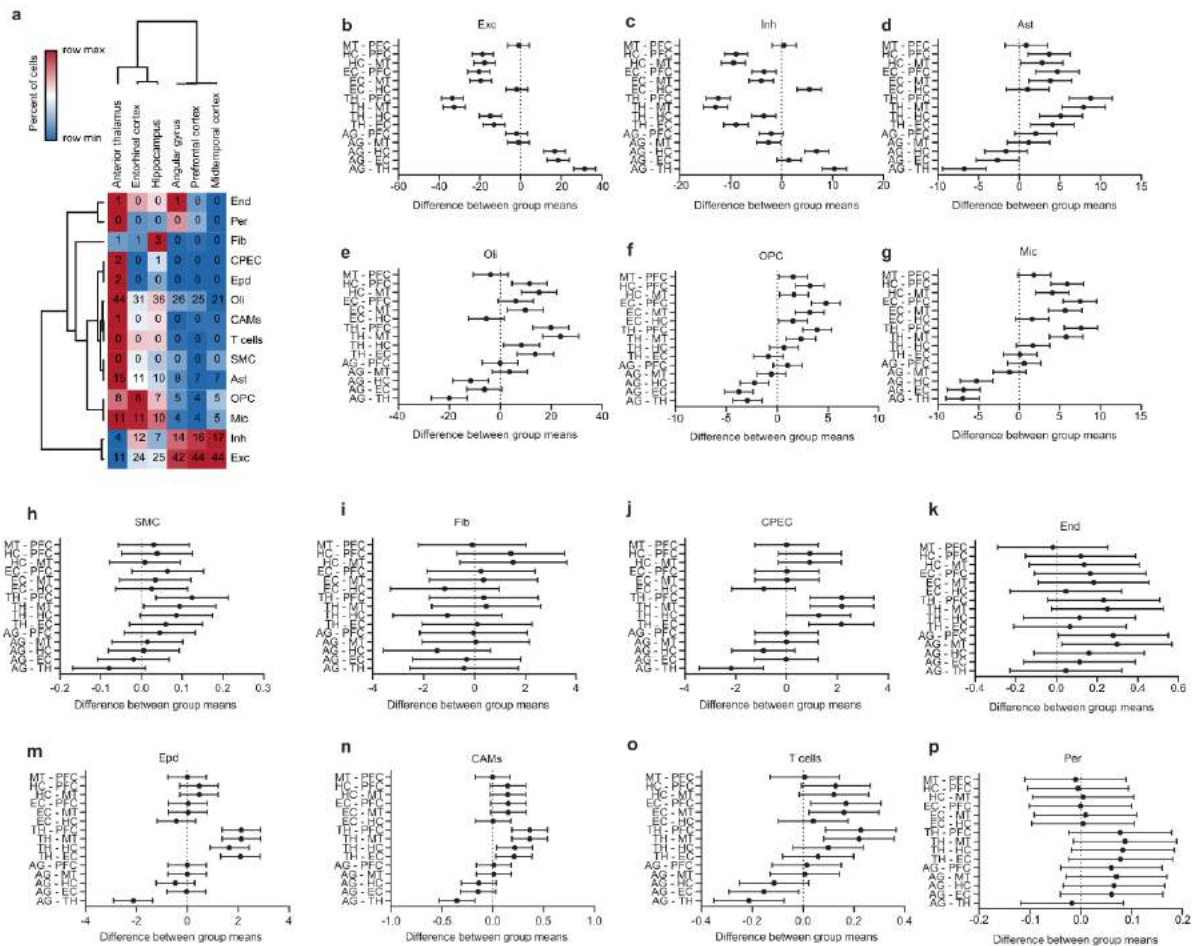


**Supplementary Figure 3-S5. Cell cycle scores. a-b,** Histograms showing the distribution of **(a)** the G2/M- and **(b)** the S phase module scores. The module scores were computed based on the expression of gene sets associated with the S phase and the G2/M phase of the cell cycle.



**Supplementary Figure 3-S6. Gene expression programs and key regulators underlying major cell type identity in the human brain.** **a**, Heat map showing percent usage of all gene expression programs (GEPs) (rows) in all major cell types (columns). **b**, Overlap (one-sided Fisher's exact test) between the top 200 genes associated with the gene expression programs indicated (rows) and the major cell type marker genes (columns). The P values have been adjusted for multiple hypothesis testing;  $-\log_{10}(\text{Bonferroni-corrected P values})$  are shown. **c-e**, Relative expression level of the top 20 genes associated with the gene expression programs (**c**) GEP Ast (preferentially used by astrocytes), (**d**) GEP Epd 1 (used by astrocytes and ependymal

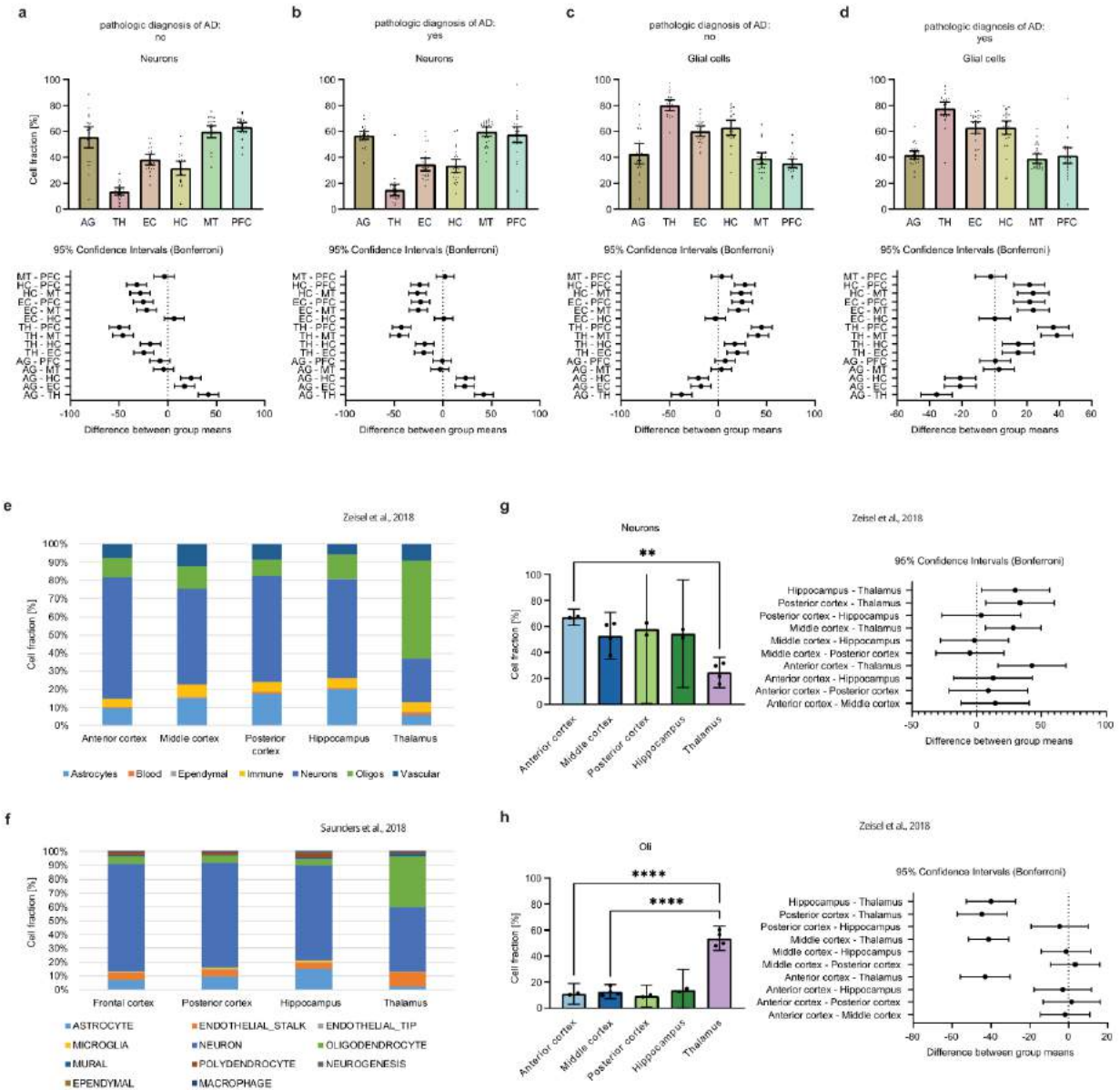
cells), and **(e)** GEP *Epd 2* (preferentially used by ependymal cells) across all major cell types. **f**, Identification of major cell type-specific transcription factor regulons. The heat map shows the module score of the top 5 transcription factor regulons (rows) for each major cell type across all individuals and major cell types (columns). **g**, Activity of selected transcription factor regulons across the major cell types. The box plots show the distribution of the average transcription factor regulon module score per individual and major cell type. Within each box, horizontal lines denote median values; boxes extend from the 25th to the 75th percentile of each group's distribution of values; whiskers extend from the 5th to the 95th percentile. \*\*\*\* $P < 0.0001$ , \*\*\* $P < 0.001$ , \*\* $P < 0.01$ , \* $P < 0.05$ ; ns,  $P > 0.05$  (ordinary one-way ANOVA corrected for multiple comparisons using Šidák's multiple comparisons test).



Extended Data Figure 7

**Supplementary Figure 3-S7. Relative abundance of major cell types across regions. a**, Relative abundance of major cell types across regions. The heat map shows the proportion of each major cell type relative to all the cells isolated from a brain region. **b-p**, Relative abundance of selected major cell types across brain regions (AG, HC, MT, PFC:  $n = 48$ ; TH:  $n = 45$ ; EC:  $n = 46$ ). Difference between brain region group means with 95% confidence intervals (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test).

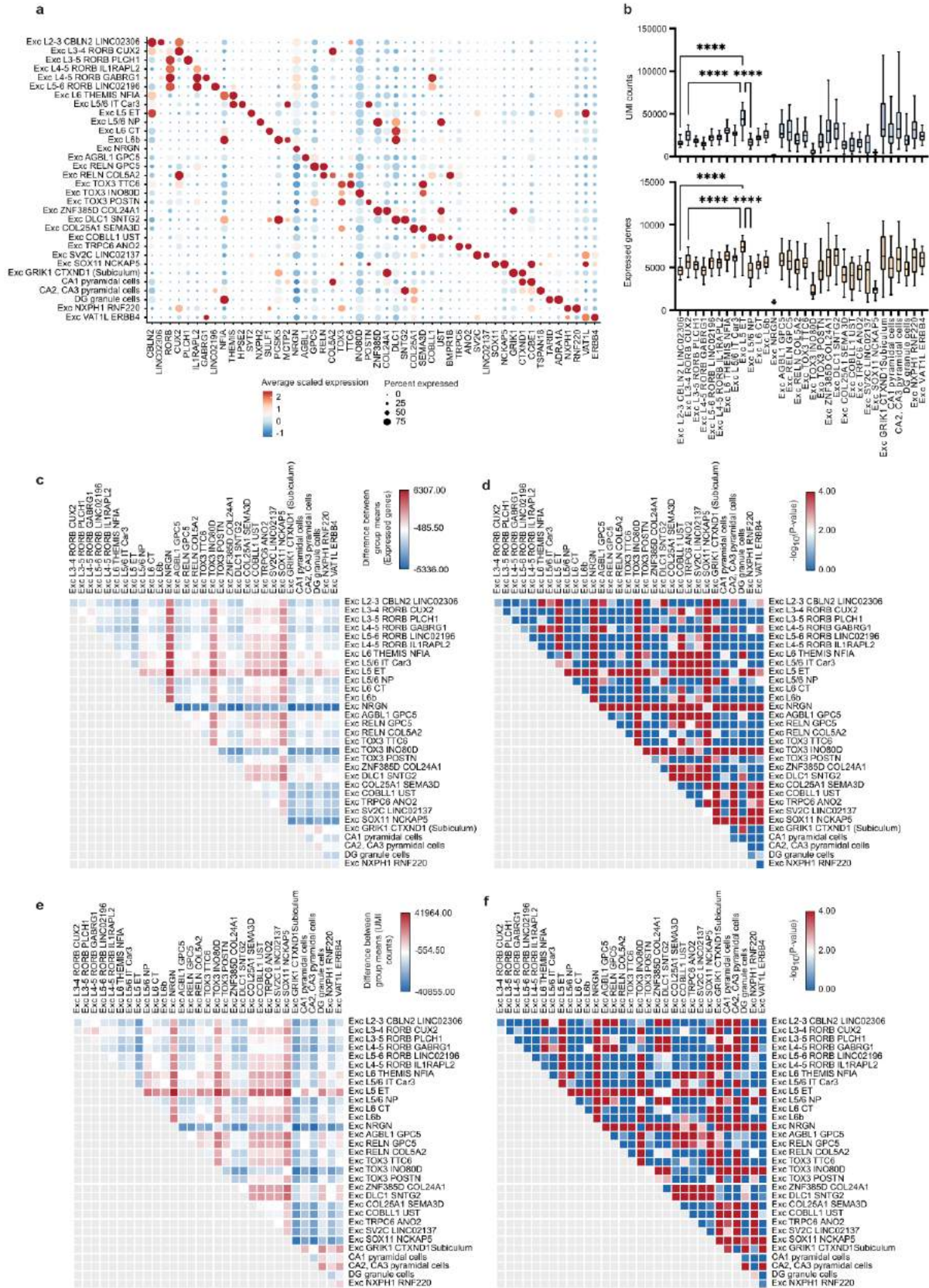




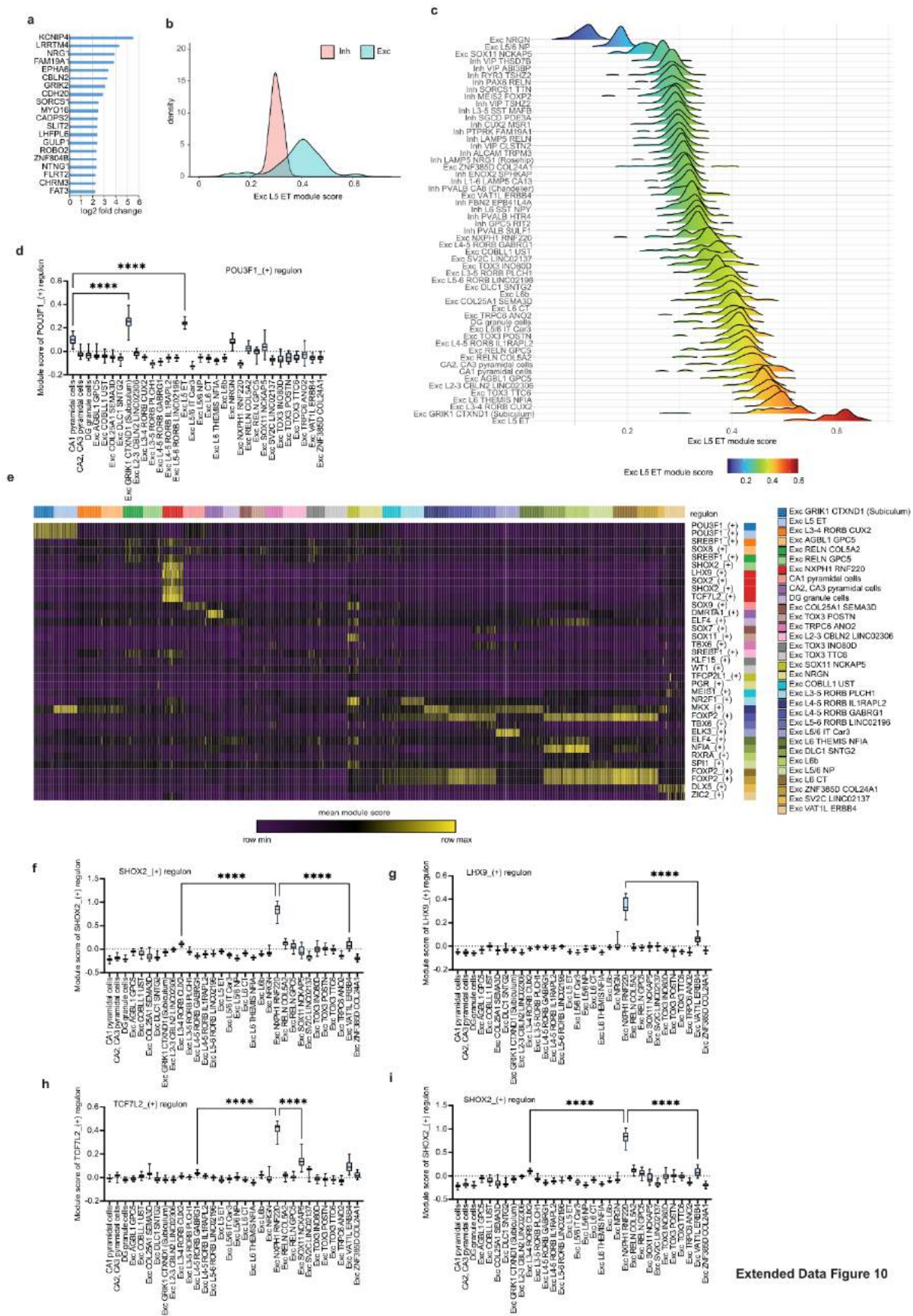
**Supplementary Figure 3-S8. Relative abundance of neurons and glial cells across brain regions.** **a-d**, Relative abundance of neurons and glial cells across brain regions in individuals with (**b,d**) and without (**a,c**) a pathologic diagnosis of Alzheimer’s disease. The bar plots show the mean fraction of neurons (**a,b**) and glial cells (**c,d**) per individual and brain region ((**a,c**) AG, EC, HC, MT, PFC: n = 22; TH: n = 21; (**b,d**) AG, HC, MT, PFC: n = 26; TH: n = 24; EC: n = 24). The fraction of neurons and glial cells was computed relative to all the cells isolated from a brain region of an individual. Data are expressed as mean with 95% confidence intervals and individual data points are shown. Lower panels: difference between group means with 95% confidence intervals (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni’s multiple comparisons test). **e-h**, Relative abundance of neurons and glial cells across regions of the mouse brain. The bar plots show the proportion of the cell types indicated relative to all the cells isolated from a brain region (based on the studies by (**e**) Zeisel et al. 2018 and (**f**) Saunders et al. 2018). (**g,h**) Relative abundance of (**g**) neurons and (**h**) oligodendrocytes



across regions of the mouse brain. The bar plots show the mean fraction of cells of the cell types indicated (Anterior cortex, Posterior cortex, Hippocampus:  $n = 2$ ; Middle cortex and thalamus:  $n = 4$ ). Data are expressed as mean with 95% confidence intervals and individual data points are shown. \*\*\*\* $P < 0.0001$ , \*\*\* $P < 0.001$ , \*\* $P < 0.01$ , \* $P < 0.05$ ; ns,  $P > 0.05$  (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). Right panels: difference between group means with 95% confidence intervals (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test).



**Supplementary Figure 3-S9. Marker genes, number of expressed genes and transcripts per cell across excitatory neuron subtypes.** **a**, Marker genes of excitatory neuron subtypes. The dot plot shows the average scaled expression of the marker genes indicated across all excitatory neuron subtypes. **b**, Mean number of genes (lower panel) and transcripts (upper panel) detected per cell and individual. Box plots show the distribution of the mean number of genes or transcripts detected per cell and individual. Within each box, horizontal lines denote median values; boxes extend from the 25th to the 75th percentile of each group's distribution of values; whiskers extend from the 5th to the 95th percentile. \*\*\*\*P < 0.0001, \*\*\*P < 0.001, \*\*P < 0.01, \*P < 0.05; ns, P > 0.05 (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). **c,d**, Difference in the number of expressed genes per cell between excitatory neuron subtypes. The heat map in panel (**c**) shows the difference between group means shown in the lower panel of (b) and the heat map in panel (**d**) shows the  $-\log_{10}$ (Bonferroni-corrected P values) (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). **e,f**, Difference in the number of transcripts per cell between excitatory neuron subtypes. The heat map in panel (**e**) shows the difference between group means shown in the upper panel of (**b**) and the heat map in panel (**f**) shows the  $-\log_{10}$ (Bonferroni-corrected P values) (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test).

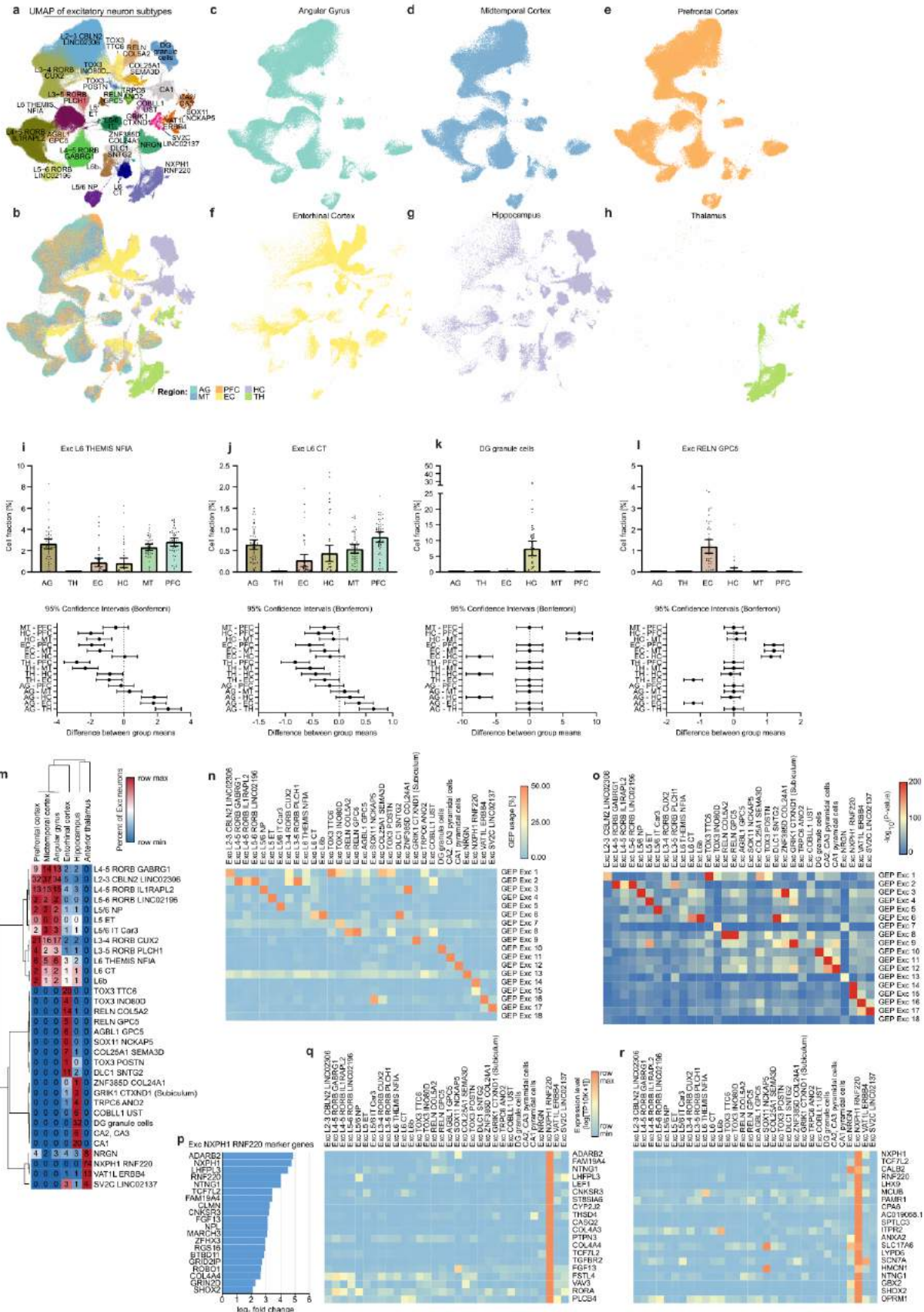


Extended Data Figure 10



**Supplementary Figure 3-S10. Gene expression programs and key regulators underlying excitatory neuron diversity.** **a**, Top 20 genes upregulated in extratelencephalic projection neurons (Exc L5 ET) compared to near-projecting excitatory neurons in layers 5 and 6 (Exc L5/6 NP). **b**, Density plot showing the distribution of the Exc L5 ET module score for excitatory and inhibitory neurons. The Exc L5 ET module score was computed on the basis of the expression of genes significantly upregulated in extratelencephalic projection neurons (Exc L5 ET) compared to near-projecting excitatory neurons in layers 5 and 6 (Exc L5/6 NP). **c**, Ridgeline plot showing the distribution of the Exc L5 ET module score for each neuronal cell type. **d**, Activity of the POU3F1 regulon across the excitatory neuron subtypes. The box plots show the distribution of the average regulon module score per individual and cell type. Within each box, horizontal lines denote median values; boxes extend from the 25th to the 75th percentile of each group's distribution of values; whiskers extend from the 5th to the 95th percentile. \*\*\*\*P < 0.0001, \*\*\*P < 0.001, \*\*P < 0.01, \*P < 0.05; ns, P > 0.05 (ordinary one-way ANOVA corrected for multiple comparisons using Šídák's multiple comparisons test). **e**, Identification of excitatory neuron subtype-specific transcription factor regulons. The heat map shows the mean module score of the top transcription factor regulons (top 5 regulons for Exc NXPH1 RNF220 neurons) (rows) across all individuals and excitatory neuron subtypes (columns). **f-i**, Activity of selected transcription factor regulons across the excitatory neuron subtypes. The box plots show the distribution of the average transcription factor regulon module score per individual and excitatory neuron subtype. Within each box, horizontal lines denote median values; boxes extend from the 25th to the 75th percentile of each group's distribution of values; whiskers extend from the 5th to the 95th percentile. \*\*\*\*P < 0.0001, \*\*\*P < 0.001, \*\*P < 0.01, \*P < 0.05; ns, P > 0.05 (ordinary one-way ANOVA corrected for multiple comparisons using Šídák's multiple comparisons test).

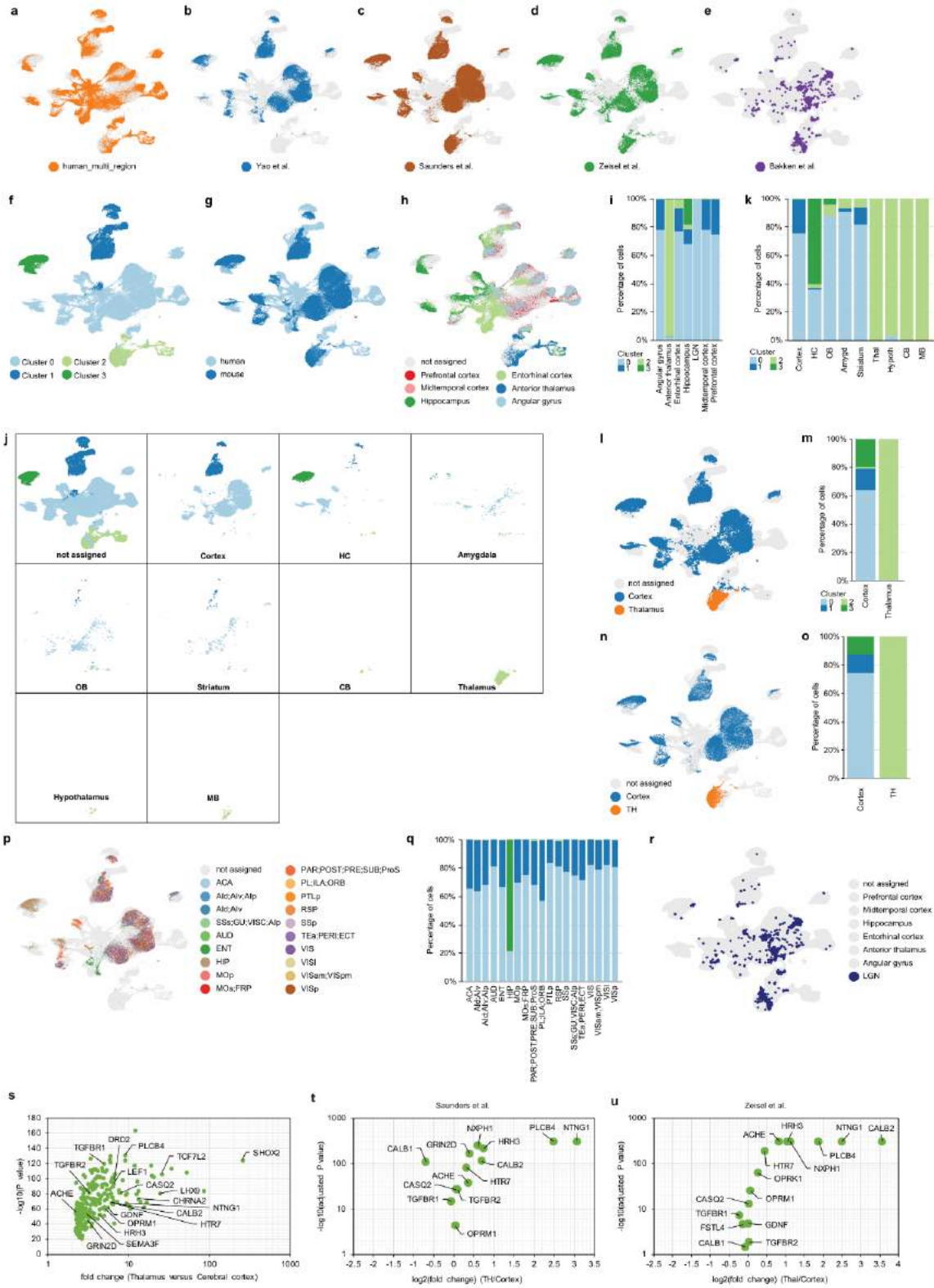




Extended Data Figure 12

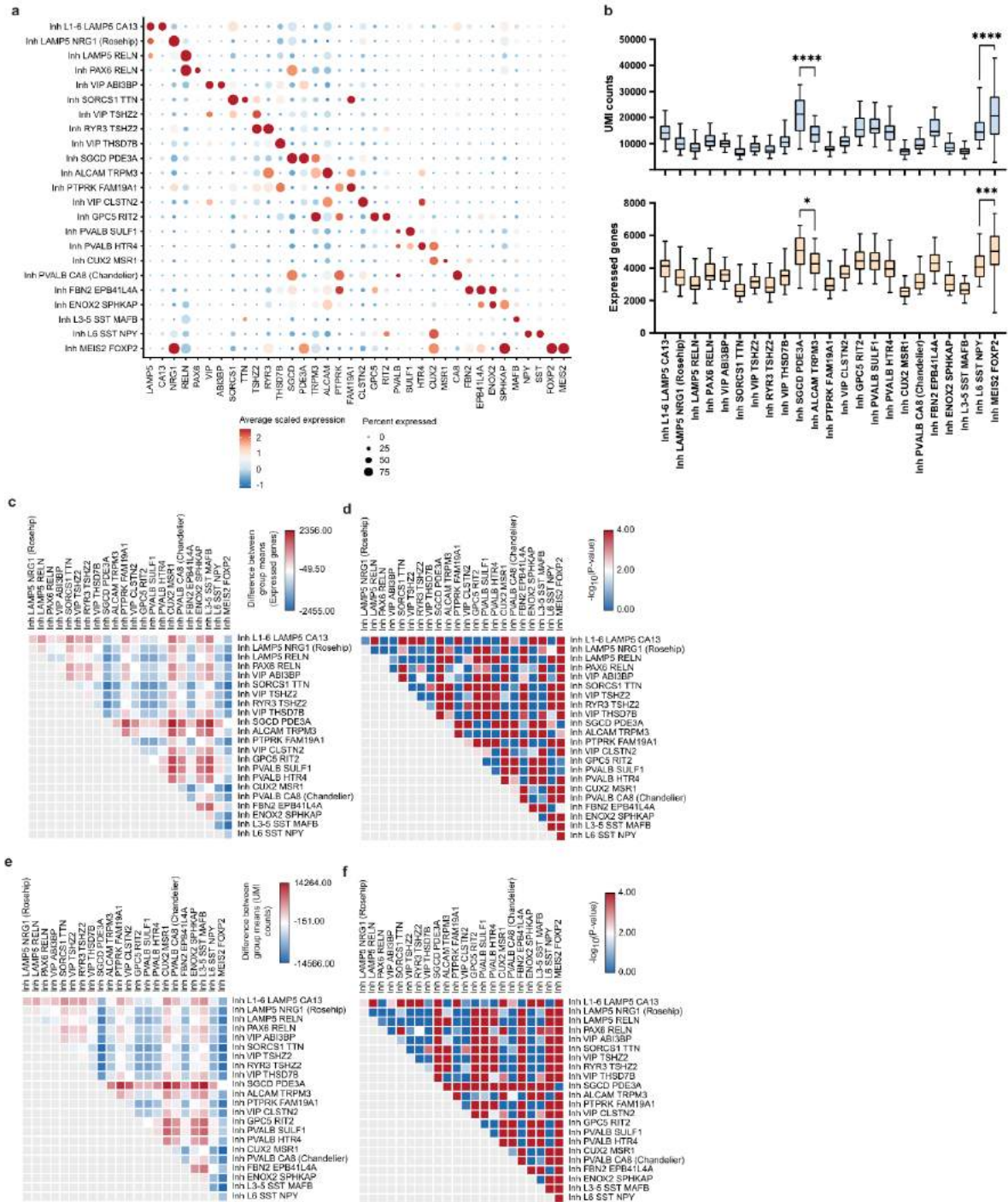
**Supplementary Figure 3-S12. Excitatory neuron composition across regions and gene expression programs.** **a**, UMAP visualization of excitatory neurons coloured by cell type. Cells isolated from all six regions are shown. **b**, UMAP visualization of excitatory neurons coloured by the brain region the cells were isolated from. **c-h**, UMAP visualization of excitatory neurons separated by the brain regions indicated. **i-l**, Relative abundance of selected excitatory neuron subtypes across brain regions. The bar plots show the mean fraction of cells of the subtype indicated per individual and brain region (AG, HC, MT, PFC: n = 48; TH: n = 45; EC: n = 46). The fraction of cells of a subtype was computed relative to all the cells isolated from a brain region of an individual. Data are expressed as mean with 95% confidence intervals and individual data points are shown. Lower panels: difference between group means with 95% confidence intervals (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). **m**, Relative abundance of excitatory neuron subtypes across brain regions. The heat map shows the proportion of each excitatory neuron subtype relative to all excitatory neurons isolated from a brain region. **n**, Heat map showing percent usage of all excitatory neuron gene expression programs (GEPs) (rows) in all excitatory neuron subtypes (columns). **o**, Overlap (one-sided Fisher's exact test) between the top 200 genes associated with the gene expression programs indicated (rows) and the excitatory neuron marker genes (columns). The P values have been adjusted for multiple hypothesis testing;  $-\log_{10}(\text{Bonferroni-corrected P values})$  are shown. **p**, Top 20 marker genes of Exc NXPH1 RNF220 neurons. The bar plot shows the  $\log_2$  fold change of the expression level of the genes indicated in Exc NXPH1 RNF220 neurons compared to all other excitatory neuron subtypes. **q,r**, relative expression level of the top 20 genes associated with the gene expression programs GEP Exc 14 (q) and 15 (r) (preferentially used by Exc NXPH1 RNF220 neurons) across all excitatory neuron subtypes.





Extended Data Figure 13

**Supplementary Figure 3-S13. Aligning human and mouse glutamatergic neurons. a-h,** UMAP visualization of glutamatergic neurons based on integrated analysis of five transcriptomic (scRNA-seq and snRNA-seq) datasets. Cells and nuclei are coloured by dataset (**a-e**), the clusters assigned based on the integrated analysis (**f**), species (**g**), and human brain region (**h**). **i**, Cluster composition. Bar plot showing the proportion of cells (y axis) from each human brain region (x axis) assigned to each of the four clusters identified based on the integrated analysis. **j**, UMAP visualization of glutamatergic neurons from the Zeisel et al. dataset separated by the brain regions indicated. Cells are coloured by the clusters assigned based on the integrated analysis. **k**, Cluster composition of the glutamatergic neurons from the Zeisel et al. dataset. Bar plot showing the proportion of cells (y axis) from each brain region (x axis) assigned to each of the four clusters identified based on the integrated analysis. **l-q**, Cluster composition (**m,o,q**) and UMAP visualization of glutamatergic neurons (**l,n,p**) from the Zeisel et al. dataset. (**l,m**), the Saunders et al. dataset (**n,o**), and the Yao et al. dataset (**p,q**). Cells are coloured by the brain regions indicated. Bar plots show the proportion of cells (y axis) from each brain region (x axis) assigned to each of the four clusters identified based on the integrated analysis. **r**, UMAP visualization of glutamatergic neurons from the Bakken et al. dataset. Nuclei are coloured by the brain region indicated. **s**, Volcano plot showing the fold change in the expression of selected genes between the thalamus and the cerebral cortex based on the analysis of microarray data from different human brain regions (Hawrylycz et al., 2012) using the Differential Search tool of the Allen Brain Atlas data portal. **t,u**, Volcano plots showing the fold change in the expression of selected genes in glutamatergic neurons between the thalamus and the cortex based on the Saunders et al. dataset (**t**) and the Zeisel et al. dataset (**u**).

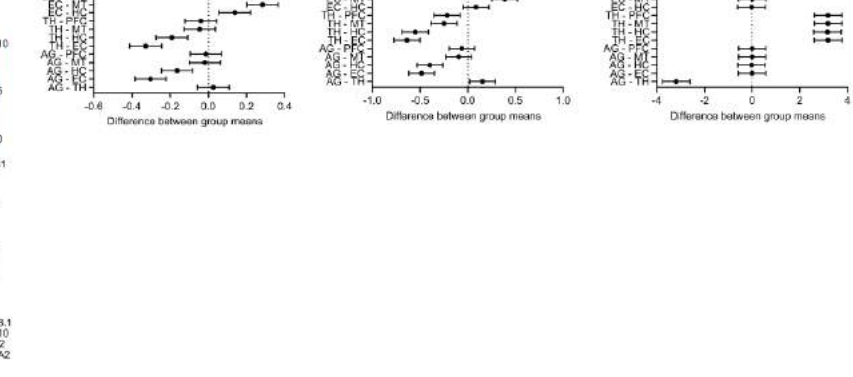
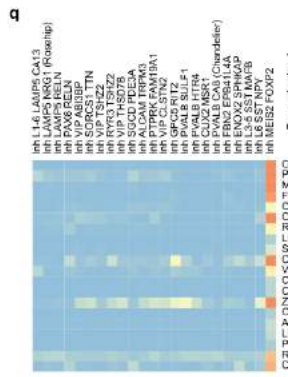
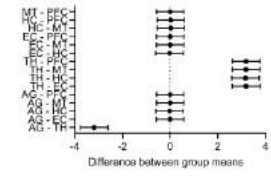
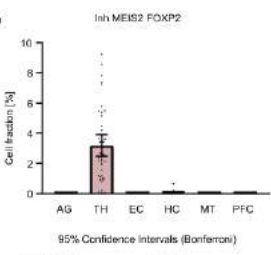
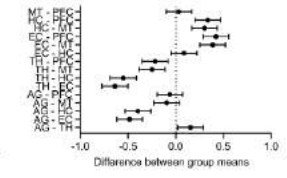
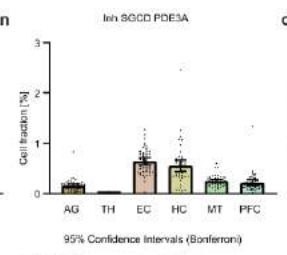
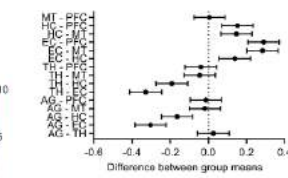
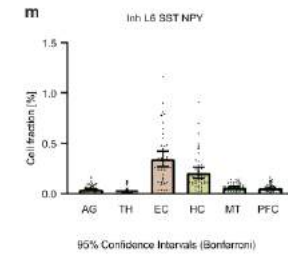
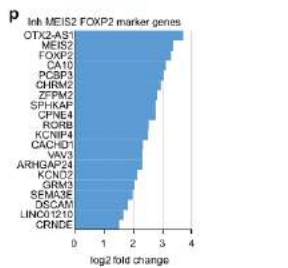
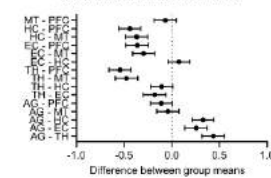
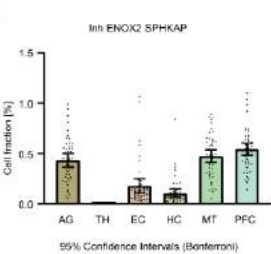
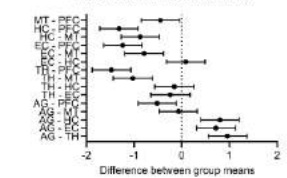
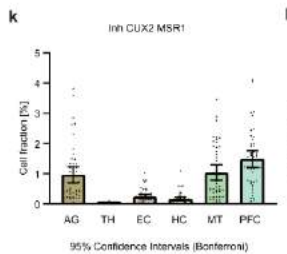
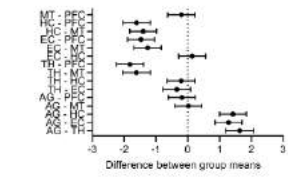
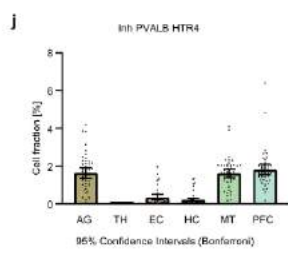
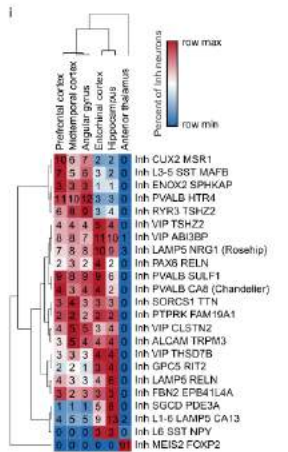
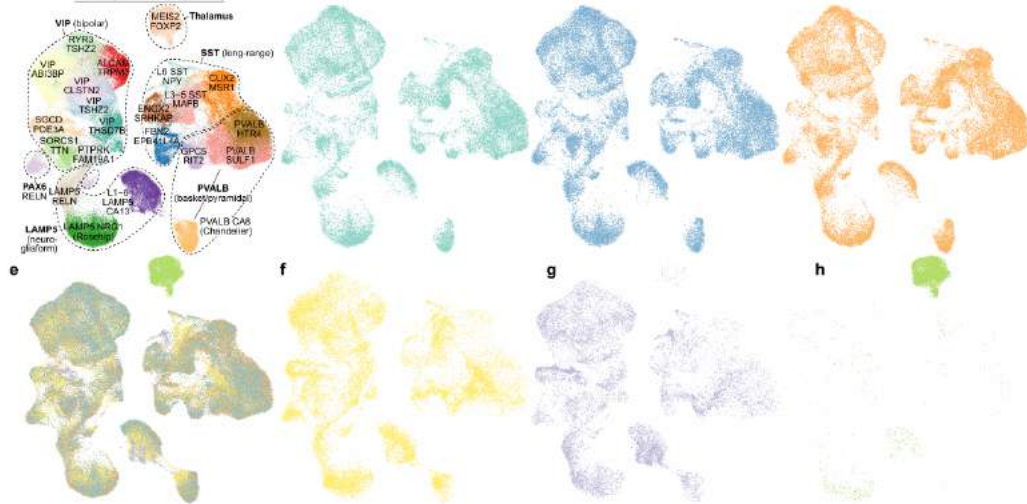


**Supplementary Figure 3-S14. Marker genes, number of expressed genes and transcripts per cell across inhibitory neuron subtypes. a**, Marker genes of inhibitory neuron subtypes. The dot plot shows the average scaled expression of the marker genes indicated across all inhibitory neuron subtypes. **b**, Mean number of genes (lower panel) and transcripts (upper panel) detected per cell and individual. Box plots show the distribution of the mean number of genes or transcripts detected per cell and individual. Within each box, horizontal lines denote median values; boxes extend from the 25th to the 75th percentile of each group's distribution of

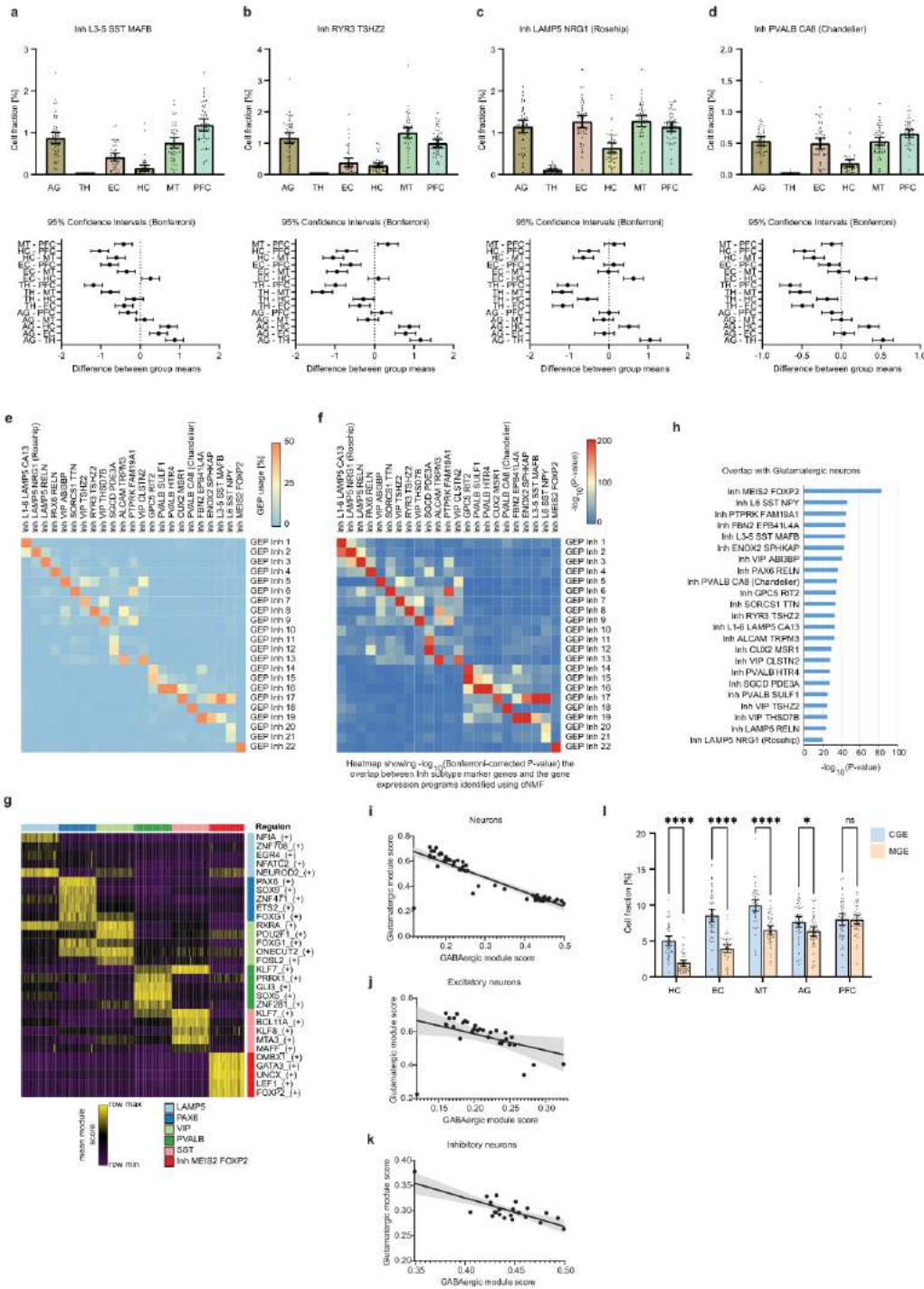
values; whiskers extend from the 5th to the 95th percentile. \*\*\*\*P < 0.0001, \*\*\*P < 0.001, \*\*P < 0.01, \*P < 0.05; ns, P > 0.05 (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). **c,d**, Difference in the number of expressed genes per cell between inhibitory neuron subtypes. The heat map in panel (**c**) shows the difference between group means shown in the lower panel of (**b**) and the heat map in panel (**d**) shows the  $-\log_{10}$ (Bonferroni-corrected P values) (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). **e,f**, Difference in the number of transcripts per cell between inhibitory neuron subtypes. The heat map in panel (**e**) shows the difference between group means shown in the upper panel of (**b**) and the heat map in panel (**f**) shows the  $-\log_{10}$ (Bonferroni-corrected P values) (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test).



**a** UMAP of inhibitory neuron subtypes, by subclass **b**



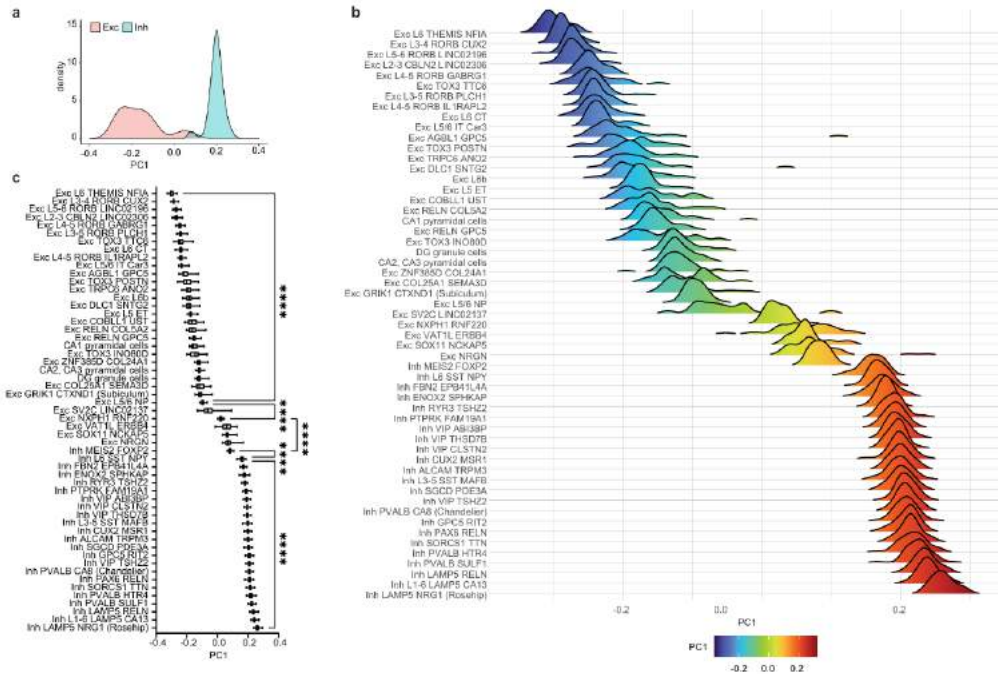
**Supplementary Figure 3-S15. Inhibitory neuron diversity across different regions of the human brain.** **a**, UMAP visualization of inhibitory neurons coloured by cell type. Cells isolated from all six regions are shown. **b**, UMAP visualization of inhibitory neurons coloured by the brain region the cells were isolated from. **c-h**, UMAP visualization of inhibitory neurons separated by the brain regions indicated. **i**, Relative abundance of inhibitory neuron subtypes across brain regions. The heat map shows the proportion of each inhibitory neuron subtype relative to all inhibitory neurons isolated from a brain region. **j-o**, Relative abundance of selected inhibitory neuron subtypes across brain regions. The bar plots show the mean fraction of cells of the cell types indicated per individual and brain region (AG, HC, MT, PFC: n = 48; TH: n = 45; EC: n = 46). The fraction of cells of a cell type was computed relative to all the cells isolated from a brain region of an individual. Data are expressed as mean with 95% confidence intervals and individual data points are shown. Lower panels: difference between group means with 95% confidence intervals (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). **p**, Top 20 marker genes of Inh MEIS2 FOXP2 neurons. The bar plot shows the  $\log_2$  fold change of the expression level of the genes indicated in Inh MEIS2 FOXP2 neurons compared to all other inhibitory neuron subtypes. **q**, Expression level of the top 20 genes associated with the gene expression program GEP Inh 22 (preferentially used by Inh MEIS2 FOXP2 neurons) across all inhibitory neuron subtypes. **r**, GABAergic- and glutamatergic module scores across all neuronal cell types. The module scores were computed on the basis of the expression of GABAergic- (x-axis) and glutamatergic (y-axis) marker gene sets.



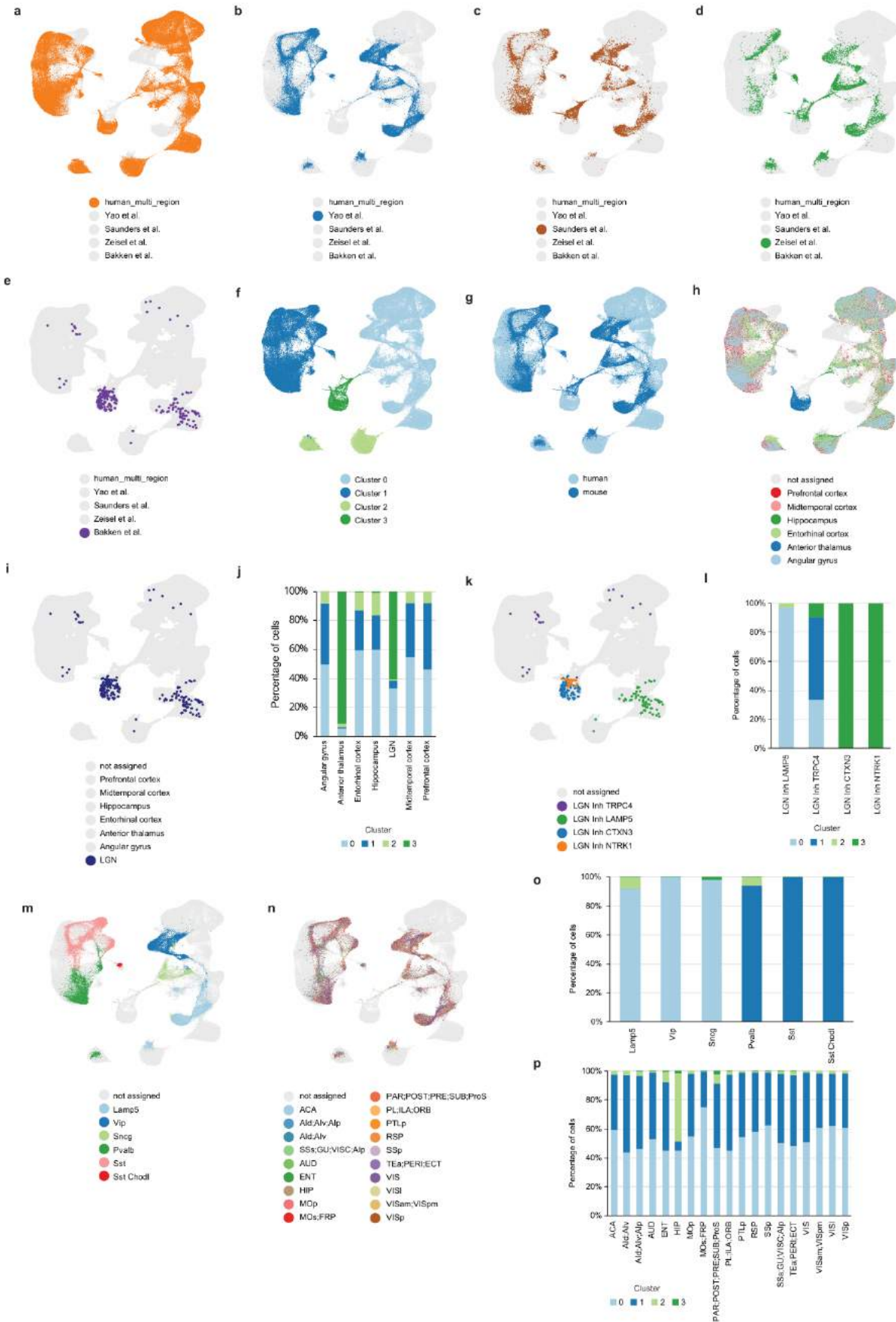
**Supplementary Figure 3-S16. Gene expression programs and key regulators underlying inhibitory neuron diversity.** a-d, Relative abundance of selected inhibitory neuron subtypes across brain regions. The bar plots show the mean fraction of cells of the subtype indicated per individual and brain region (AG, HC, MT, PFC: n = 48; TH: n = 45; EC: n = 46). The fraction of cells of a subtype was computed relative to all the cells isolated from a brain region of an individual. Data are expressed as mean with 95% confidence intervals and individual data points are shown. Lower panels: difference between group means with 95% confidence

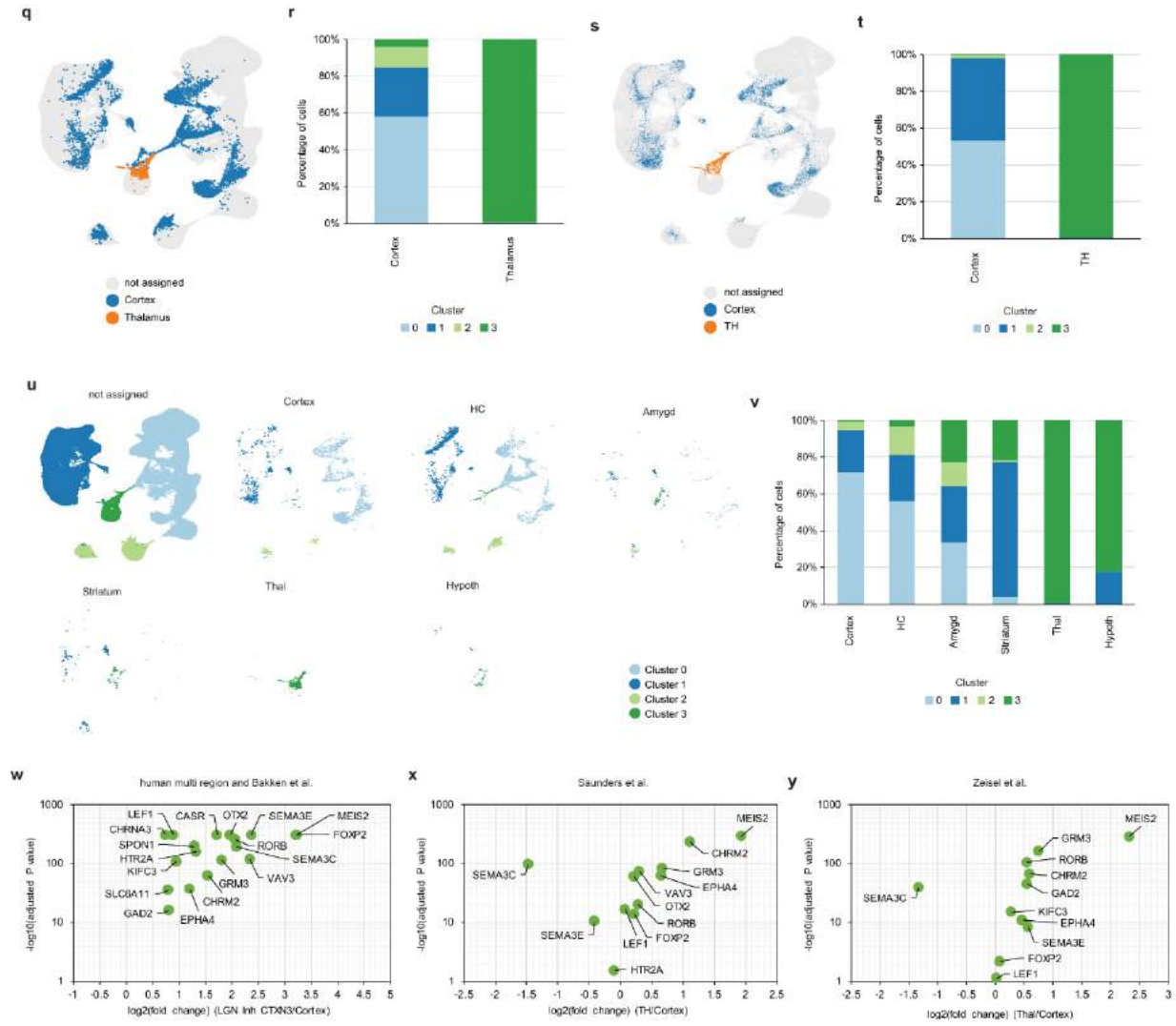
intervals (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). **e**, Heat map showing percent usage of all inhibitory neuron gene expression programs (GEPs) (rows) in all inhibitory neuron subtypes (columns). **f**, Overlap (one-sided Fisher's exact test) between the top 200 genes associated with the gene expression programs indicated (rows) and the inhibitory neuron marker genes (columns). The P values have been adjusted for multiple hypothesis testing;  $-\log_{10}(\text{Bonferroni-corrected P values})$  are shown. **g**, Identification of inhibitory neuron subclass-specific transcription factor regulons. The heat map shows the module score of the top 5 transcription factor regulons (rows) for each subclass across all individuals and subclasses of inhibitory neurons (columns). **h**, Overlap (one-sided Fisher's exact test) between inhibitory neuron subtype marker genes and marker genes of glutamatergic neurons. The P values have been adjusted for multiple hypothesis testing;  $-\log_{10}(\text{Bonferroni-corrected P values})$  are shown. **i-k**, Linear regression analysis showing a negative association between the glutamatergic- and the GABAergic module score across **(i)** all neuronal cell types, as well as across **(j)** excitatory neuron subtypes and **(k)** inhibitory neuron subtypes. **l**, Relative abundance of inhibitory neurons originating from the medial (MGE) ganglionic eminences (SST and PVALB) and the caudal (CGE) ganglionic eminence (VIP, PAX6, and LAMP5) across brain regions. The bar plots show the mean fraction of cells per individual and brain region (AG, HC, MT, PFC:  $n = 48$ ; TH:  $n = 45$ ; EC:  $n = 46$ ). The fraction of cells was computed relative to all the cells isolated from a brain region of an individual. Data are expressed as mean with 95% confidence intervals and individual data points are shown. \*\*\*\* $P < 0.0001$ , \*\*\* $P < 0.001$ , \*\* $P < 0.01$ , \* $P < 0.05$ ; ns,  $P > 0.05$  (two-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test).





**Supplementary Figure 3-S17. Polarization of glutamatergic and GABAergic neuronal subtypes.** **a**, The first principal component that was determined based on the average glutamatergic- and GABAergic module scores of all neuronal cell types separates Exc and Inh neuron subtypes. The density plot shows the distribution of the PC1 scores (the coordinates of the individual observations on the first principal component axis) of excitatory and inhibitory neuron subtypes. **b**, Ridgeline plot showing the distribution of the PC1 scores described in **(a)** for each neuronal cell type. **c**, Statistical analysis of the distribution of PC1 scores described in **(a)**. The box plot shows the distribution of the mean PC1 score per cell type and study participant. Within each box, horizontal lines denote median values; boxes extend from the 25th to the 75th percentile of each group's distribution of values; whiskers extend from the 5th to the 95th percentile. \*\*\*\* $P < 0.0001$ , \*\*\* $P < 0.001$ , \*\* $P < 0.01$ , \* $P < 0.05$ ; ns,  $P > 0.05$  (ordinary one-way ANOVA corrected for multiple comparisons using Šidák's multiple comparisons test).

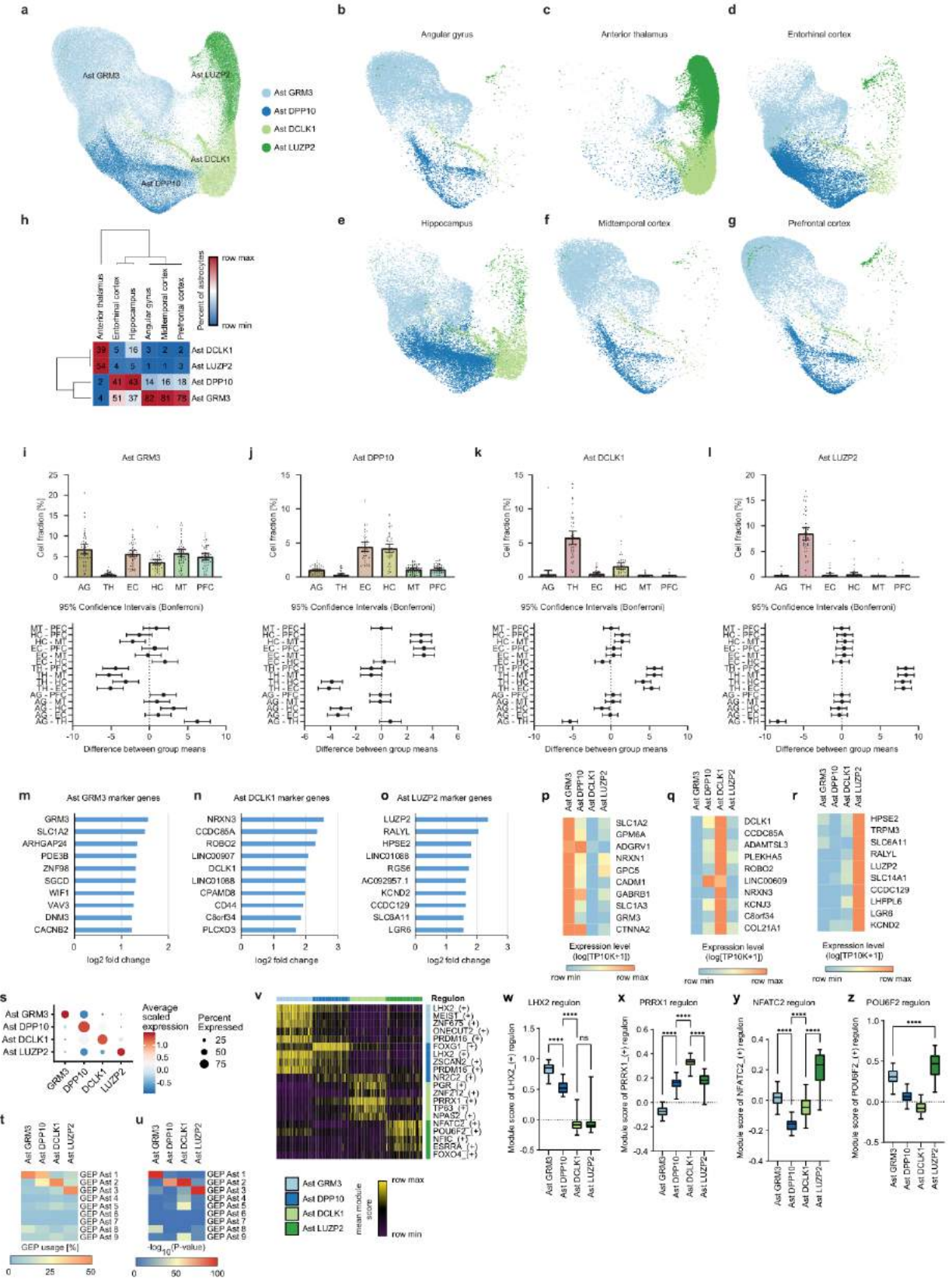




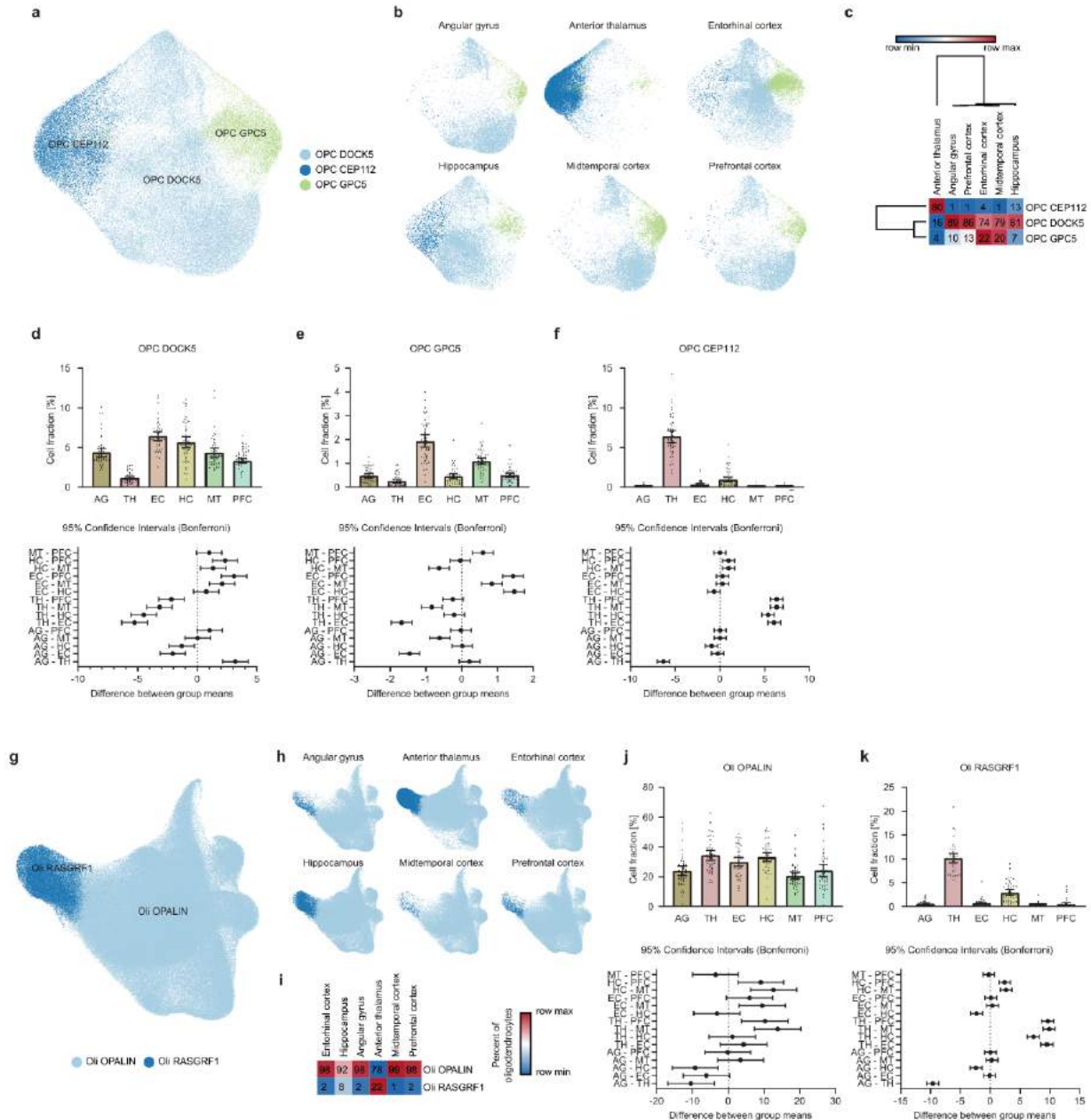
**Supplementary Figure 3-S18. Aligning human and mouse GABAergic neurons.** **a-i**, UMAP visualization of GABAergic neurons based on integrated analysis of five transcriptomic (scRNA-seq and snRNA-seq) datasets. Cells and nuclei are coloured by dataset (**a-e**), the clusters assigned based on the integrated analysis (**f**), species (**g**), and human brain region (**h,i**). **j**, Cluster composition. Bar plot showing the proportion of cells (y axis) from each human brain region (x axis) assigned to each of the four clusters identified based on the integrated analysis. **k**, UMAP visualization of GABAergic neurons from the Bakken et al. dataset coloured by the GABAergic neuron subtypes defined by Bakken et al. **l**, Cluster composition. Bar plot showing the proportion of cells (y axis) of each GABAergic neuron subtype defined by Bakken et al. assigned to each of the four clusters identified based on the integrated analysis. **m-t**, Cluster composition (**o,p,r,t**) and UMAP visualization of GABAergic neurons (**m,n,q,s**) from the Yao et al. dataset. (**m-p**), the Zeisel et al. dataset (**q,r**), and the Saunders et al. dataset (**s,t**). Cells are coloured by the brain regions or GABAergic neuron subclass indicated. Bar plots show the proportion of cells (y axis) from each brain region or GABAergic neuron subclass (x axis) assigned to each of the four clusters identified based on the integrated analysis. **u**, UMAP visualization of GABAergic neurons from the Zeisel et al. dataset separated by the brain regions indicated. Cells are coloured by the clusters assigned based on the integrated analysis. **v**,

Cluster composition of the GABAergic neurons from the Zeisel et al. dataset. Bar plot showing the proportion of cells (y axis) from each brain region (x axis) assigned to each of the four clusters identified based on the integrated analysis. **w-y**, Volcano plots of differential expression analysis of GABAergic neurons comparing the human lateral geniculate nucleus (dLGN) versus the human cortex (**w**), and comparing the mouse thalamus versus the mouse cortex based on the Saunders et al. dataset (**x**) and the Zeisel et al. dataset (**y**).





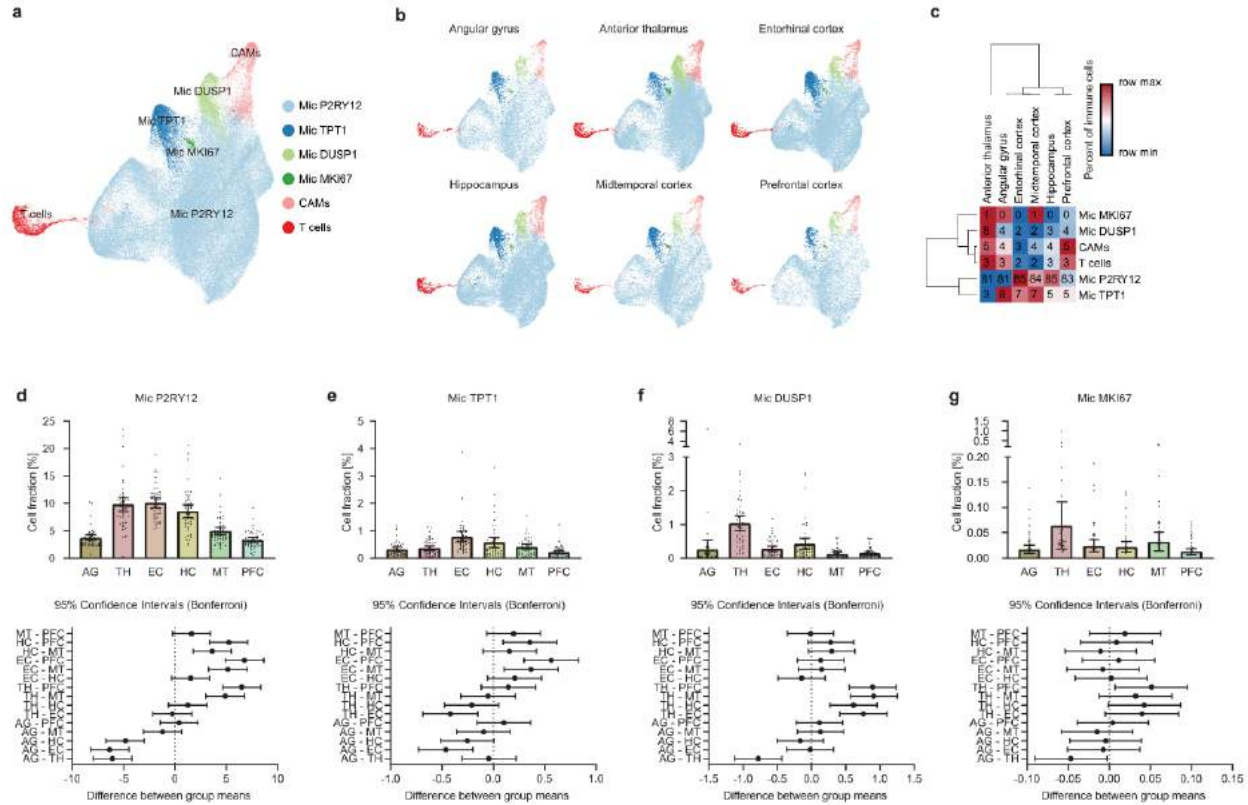
**Supplementary Figure 3-S19. Astrocyte diversity across different regions of the human brain.** **a**, UMAP visualization of astrocytes coloured by subtype. Cells isolated from all six regions are shown. **b-g**, UMAP visualization of astrocytes separated by the brain regions indicated. **h**, Relative abundance of astrocyte subtypes across brain regions. The heat map shows the proportion of each astrocyte subtype relative to all astrocytes isolated from a brain region. **i-l**, Relative abundance of astrocyte subtypes across brain regions. The bar plots show the mean fraction of cells of the subtype indicated per individual and brain region (AG, HC, MT, PFC: n = 48; TH: n = 45; EC: n = 46). The fraction of cells of an astrocyte subtype was computed relative to all the cells isolated from a brain region of an individual. Data are expressed as mean with 95% confidence intervals and individual data points are shown. Lower panels: difference between group means with 95% confidence intervals (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). **m-o**, Top 10 marker genes of the astrocyte subtypes Ast GRM3 (**m**), Ast DCLK1 (**n**), and Ast LUZP2 (**o**). The bar plots show the  $\log_2$  fold change in expression between the astrocyte subtype indicated and all other astrocyte subtypes. **p-r**, Relative expression level of the top 10 genes associated with the gene expression programs GEP Ast 1 (preferentially used by the astrocyte subtype Ast GRM3) (**p**), GEP Ast 2 (preferentially used by the astrocyte subtype Ast DCLK1) (**q**), and GEP Ast 3 (preferentially used by the astrocyte subtype Ast LUZP2) (**r**) across all astrocyte subtypes. **s**, Marker genes of astrocyte subtypes. The dot plot shows the average scaled expression of the marker genes indicated across all astrocyte subtypes. **t**, Heat map showing percent usage of all astrocyte gene expression programs (GEPs) (rows) in all astrocyte subtypes (columns). **u**, Overlap (one-sided Fisher's exact test) between the top 200 genes associated with the gene expression programs indicated (rows) and the astrocyte marker genes (columns). The P values have been adjusted for multiple hypothesis testing;  $-\log_{10}(\text{Bonferroni-corrected P values})$  are shown. **v**, Identification of astrocyte subtype-specific transcription factor regulons. The heat map shows the mean module score of the top 5 transcription factor regulons (rows) across all individuals and astrocyte subtypes (columns). **w-z**, Activity of selected transcription factor regulons across the astrocyte subtypes. The box plots show the distribution of the average transcription factor regulon module score per individual and astrocyte subtype. Within each box, horizontal lines denote median values; boxes extend from the 25th to the 75th percentile of each group's distribution of values; whiskers extend from the 5th to the 95th percentile. \*\*\*\*P < 0.0001, \*\*\*P < 0.001, \*\*P < 0.01, \*P < 0.05; ns, P > 0.05 (ordinary one-way ANOVA corrected for multiple comparisons using Šidák's multiple comparisons test).



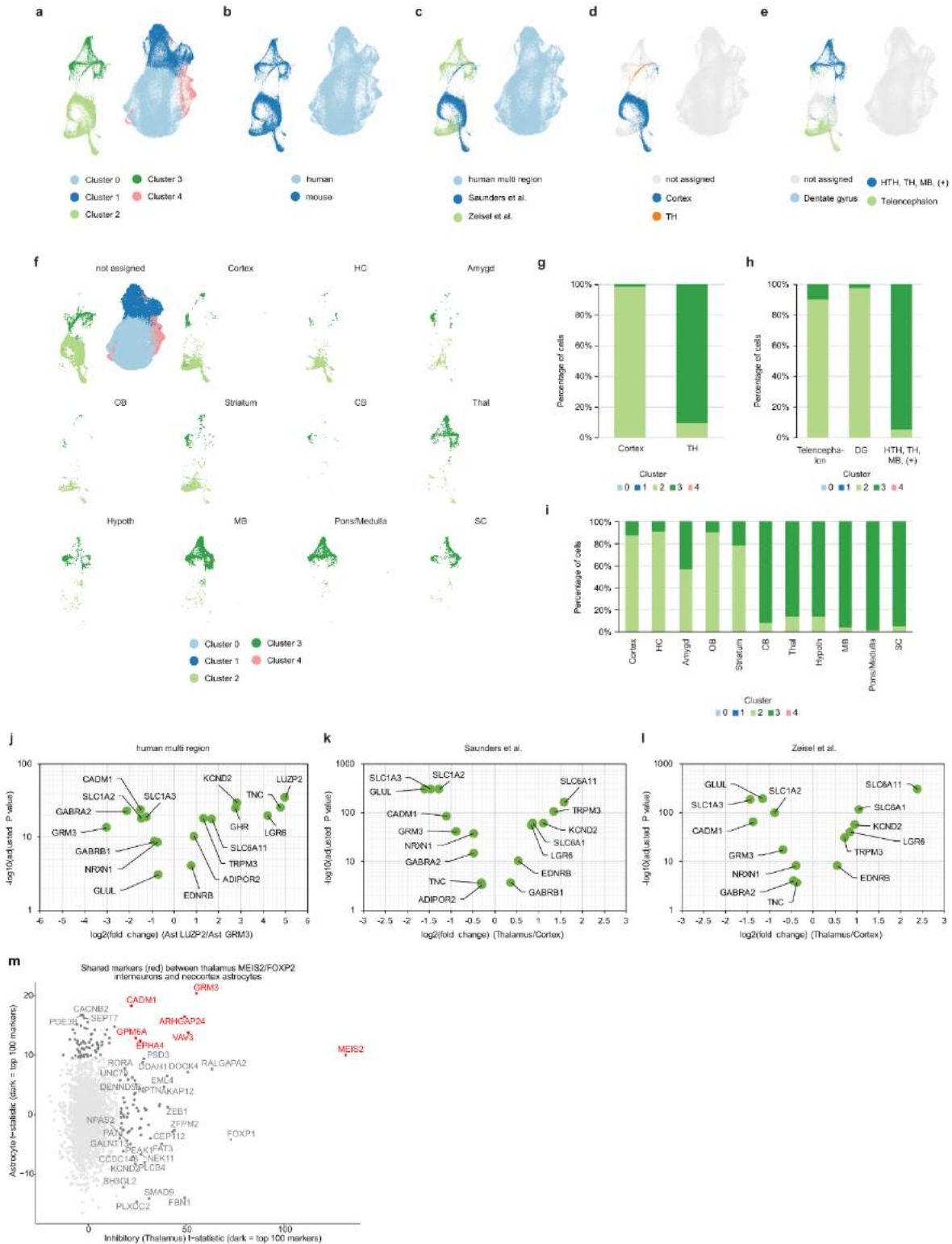
**Supplementary Figure 3-S20. Aligning human and mouse astrocytes.** **a-f**, UMAP visualization of astrocytes based on integrated analysis of three transcriptomic (scRNA-seq and snRNA-seq) datasets. Cells and nuclei are coloured by the clusters assigned based on the integrated analysis (**a**), species (**b**), and dataset (**c**). **d,e**, UMAP visualization of astrocytes from the Saunders et al. dataset (**d**) and the Zeisel et al. dataset (**e**). Cells are coloured by the brain region indicated. **f**, UMAP visualization of astrocytes from the Zeisel et al. dataset separated by the brain regions indicated. Cells are coloured by the clusters assigned based on the integrated analysis. **g-i**, Cluster composition. Bar plot showing the proportion of astrocytes (y axis) from each brain region (x axis) assigned to each of the five clusters identified based on the integrated analysis. The panels show the results for the Saunders et al. dataset (**g**) and the Zeisel et al. dataset (**h,i**). **j-l**, Volcano plots of differential expression analyses comparing the human astrocyte subtypes Ast LUZP2 and Ast GRM3 (**j**), and comparing thalamic versus cortical

astrocytes isolated from the mouse brain based on the Saunders et al. dataset (**k**) and the Zeisel et al. dataset (**l**). **m**, Shared markers between neocortex astrocytes and thalamus interneurons (Inh MEIS2 FOXP2). Scatter plot of t-test for each subtype against other cells of the same class (vs. astrocytes and vs. inhibitory neurons), with top 100 genes colored in gray and overlap colored in red.



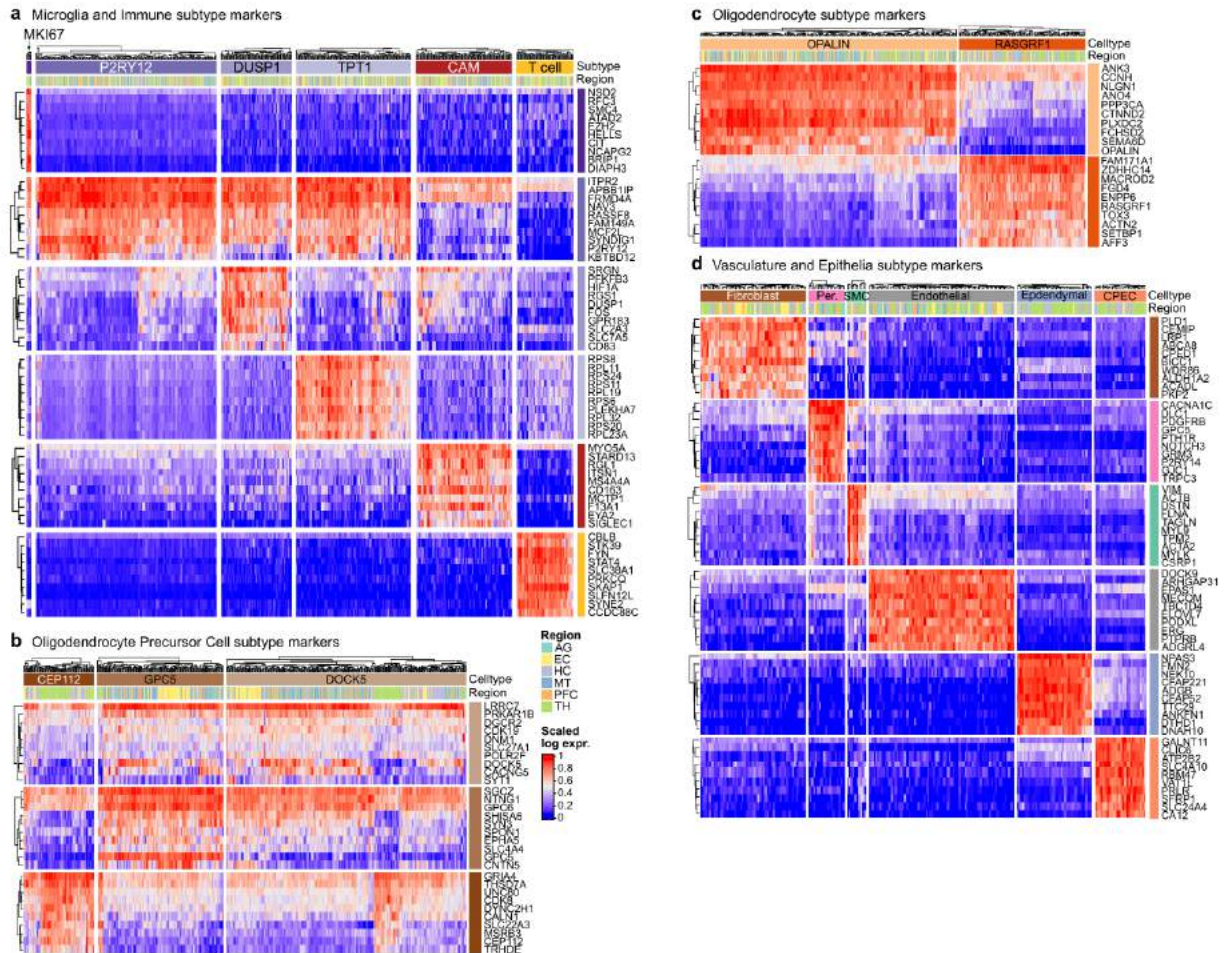


**Supplementary Figure 3-S21. Oligodendrocyte lineage diversity across different regions of the human brain.** **a**, UMAP visualization of OPCs coloured by subtype. Cells isolated from all six regions are shown. **b**, UMAP visualization of OPCs separated by the brain regions indicated. **c**, Relative abundance of OPC subtypes across brain regions. The heat map shows the proportion of each OPC subtype relative to all OPCs isolated from a brain region. **d-f**, Relative abundance of OPC subtypes across brain regions. The bar plots show the mean fraction of cells of the cell types indicated per individual and brain region (AG, HC, MT, PFC:  $n = 48$ ; TH:  $n = 45$ ; EC:  $n = 46$ ). The fraction of cells of a cell type was computed relative to all the cells isolated from a brain region of an individual. Data are expressed as mean with 95% confidence intervals and individual data points are shown. Lower panels: difference between group means with 95% confidence intervals (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test). **g**, UMAP visualization of oligodendrocytes coloured by subtype. Cells isolated from all six regions are shown. **h**, UMAP visualization of oligodendrocytes separated by the brain regions indicated. **i**, Relative abundance of oligodendrocyte subtypes across brain regions. The heat map shows the proportion of each oligodendrocyte subtype relative to all oligodendrocytes isolated from a brain region. **j,k**, Relative abundance of oligodendrocyte subtypes across brain regions. The bar plots show the mean fraction of cells of the cell types indicated per individual and brain region (AG, HC, MT, PFC:  $n = 48$ ; TH:  $n = 45$ ; EC:  $n = 46$ ). The fraction of cells of a cell type was computed relative to all the cells isolated from a brain region of an individual. Data are expressed as mean with 95% confidence intervals and individual data points are shown. Lower panels: difference between group means with 95% confidence intervals (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test).

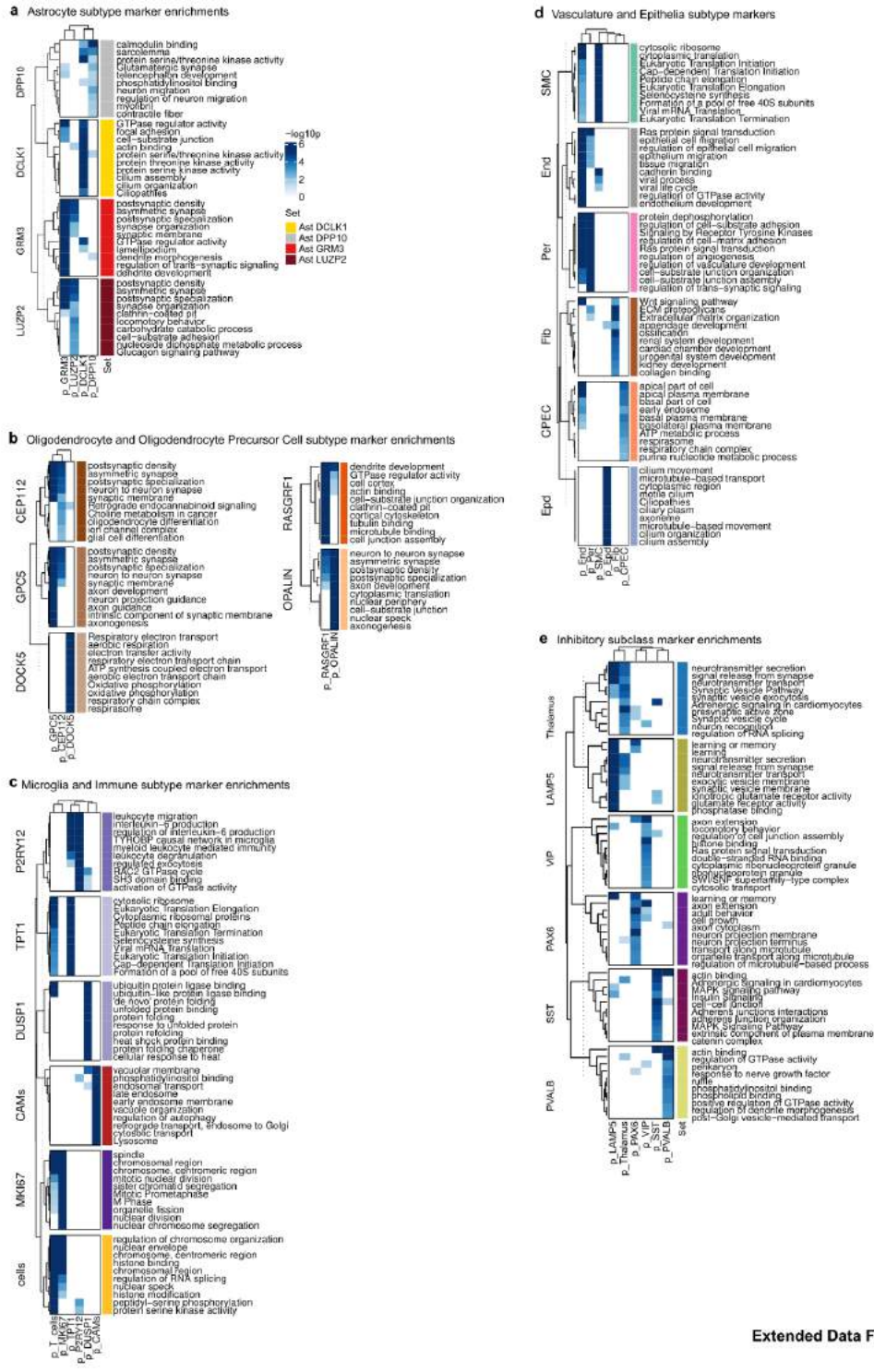


**Supplementary Figure 3-S22: Microglia and immune cell diversity across brain regions. a**, UMAP visualization of microglia and immune cells coloured by cell type. Cells isolated from all six regions are shown. **b**, UMAP visualization of microglia and immune cells separated by the

brain regions indicated. **c**, Relative abundance of microglia and immune cells across brain regions. The heat map shows the proportion of each subtype relative to all microglia and immune cells isolated from a brain region. **d-g**, Relative abundance of microglia and immune cells across brain regions. The bar plots show the mean fraction of cells of the subtypes indicated per individual and brain region (AG, HC, MT, PFC: n = 48; TH: n = 45; EC: n = 46). The fraction of cells of a subtype was computed relative to all the cells isolated from a brain region of an individual. Data are expressed as mean with 95% confidence intervals and individual data points are shown. Lower panels: difference between group means with 95% confidence intervals (ordinary one-way ANOVA corrected for multiple comparisons using Bonferroni's multiple comparisons test).

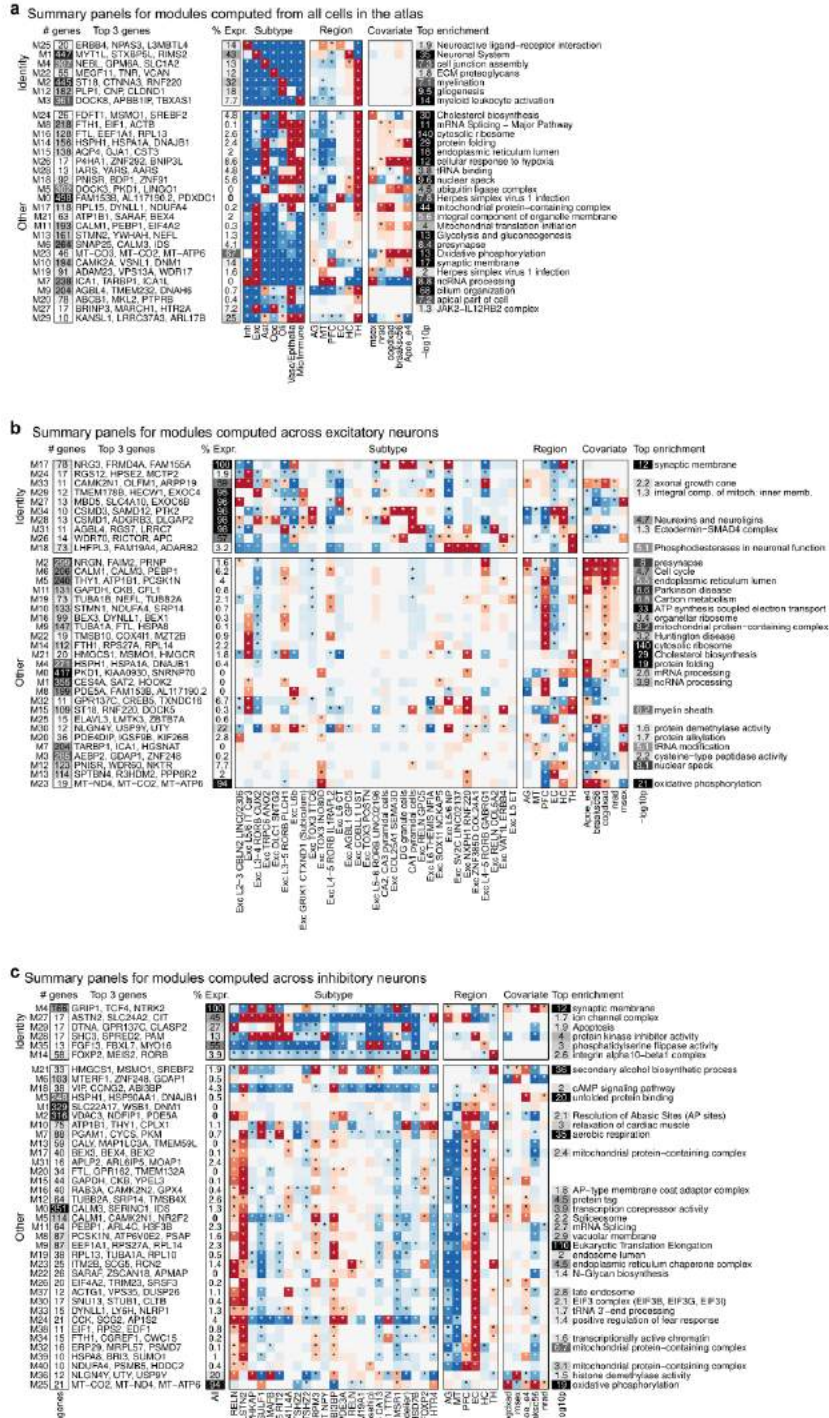






Extended Data Figure 25

**Supplementary Figure 3-S24. Marker gene enrichments.** a-e, Unique enrichments for marker genes for astrocytes (a), oligodendrocytes and OPCs (b), microglia and immune cells (c), vasculature and epithelia (d), and inhibitory subclasses (e). Top 10 terms that are unique to two subclasses or fewer and that have fewer than 500 genes in the term are shown. Adjusted p-values are shown as  $-\log_{10}p$ -values capped at 5.

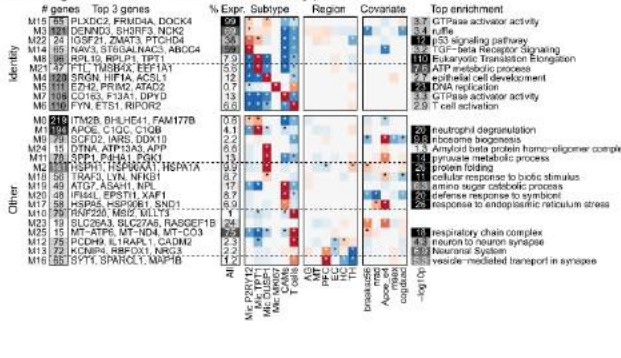


**Supplementary Figure 3-S25. Module summary panels across modules computed from all cells and from neuronal cell types. a-c, Overview of gene expression modules with at least 10 genes each across major non-neuronal cell types, showing the module name, number of genes, percent expression, top module genes, enrichments by subtype and region, and the top**

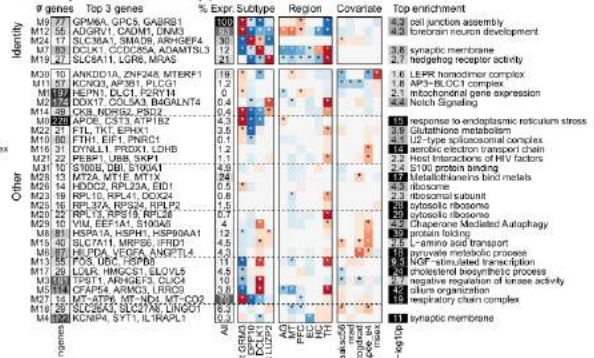


functional enrichment for each module. Percent expression is the percent of cells whose average expression (log1p, normalized) of the module is above 1. Covariate enrichments are performed by hypergeometric test, comparing the intersection of cells with z-scored module expression of at least 1 vs. with  $z < 1$  against a particular level of a covariate of interest (e.g. cells from the entorhinal cortex region or cells of a specific subtype). Panels summarize modules found in all cells in the atlas (a), excitatory neurons (b), and inhibitory neurons (c). OPCs (d), and vasculature and epithelia (e). All modules except vasculature and epithelia modules are split into identity vs. other, where identity modules are highly enriched in a single subtype and have an average expression greater than 1 (log1p, normalized) for over 50% of the subtype's cells.

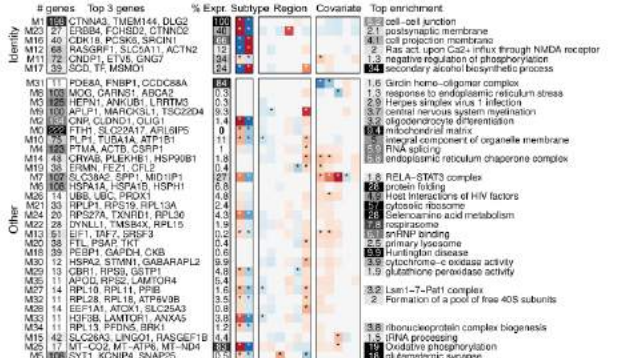
**a** Microglial and immune cell modules summary panels



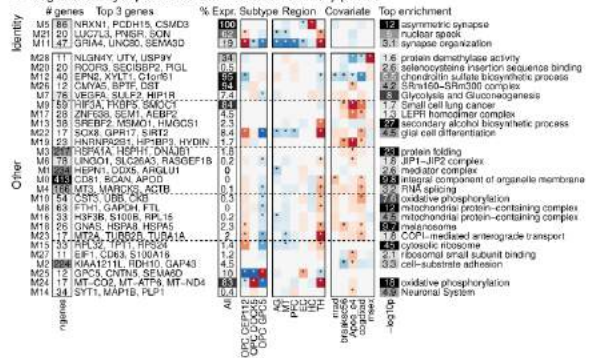
**b** Astrocyte modules summary panels



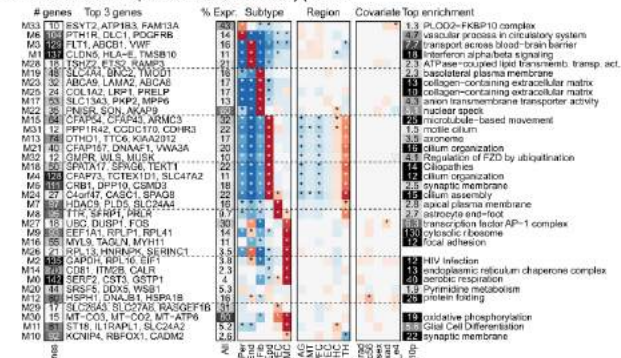
**c** Oligodendrocyte modules summary panels



**d** Oligodendrocyte precursor cell modules summary panels

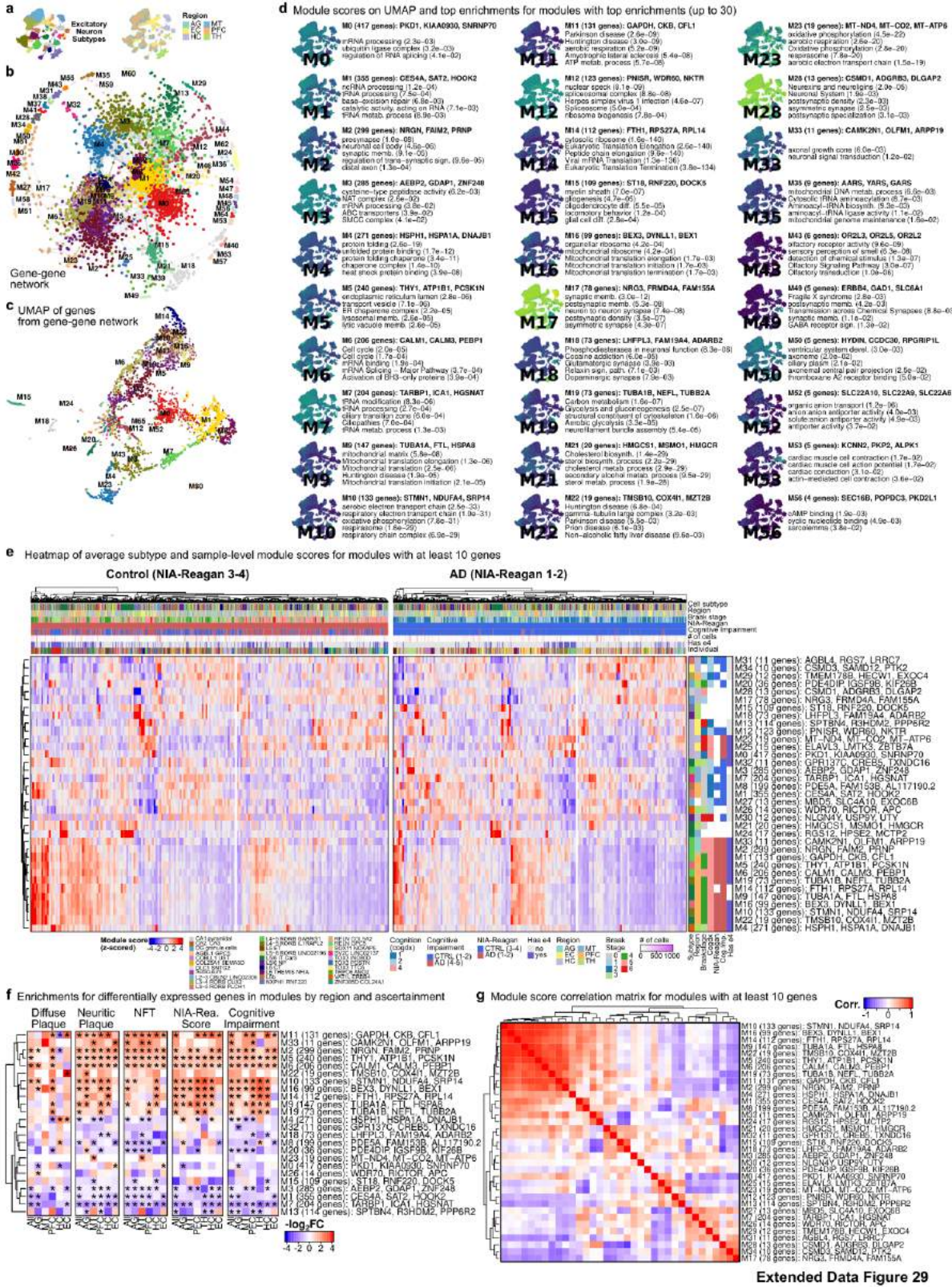


**e** Vasculature and epithelia modules summary panels



**Supplementary Figure 3-S26. Module summary panels across modules computed from non-neuronal cell types. a-e**, Overview of gene expression modules with at least 10 genes each across major non-neuronal cell types, showing the module name, number of genes, percent expression, top module genes, enrichments by subtype and region (hypergeometric test), and the top functional enrichment for each module. Percent expression is the percent of cells whose average expression ( $\log_{10} p$ , normalized) of the module is above 1. Covariate enrichments are performed by hypergeometric test, comparing the intersection of cells with z-scored module expression of at least 1 vs. with  $z < 1$  against a particular level of a covariate of interest (e.g. cells from the entorhinal cortex region or cells of a specific subtype). Panels summarize modules found in microglia and immune cells (**a**), astrocytes (**b**), oligodendrocytes (**c**), OPCs (**d**), and vasculature and epithelia (**e**). All modules except vasculature and epithelia modules are split into identity vs. other, where identity modules are highly enriched in a single subtype and have an average expression greater than 1 ( $\log_{10} p$ , normalized) for over 50% of the subtype's cells.

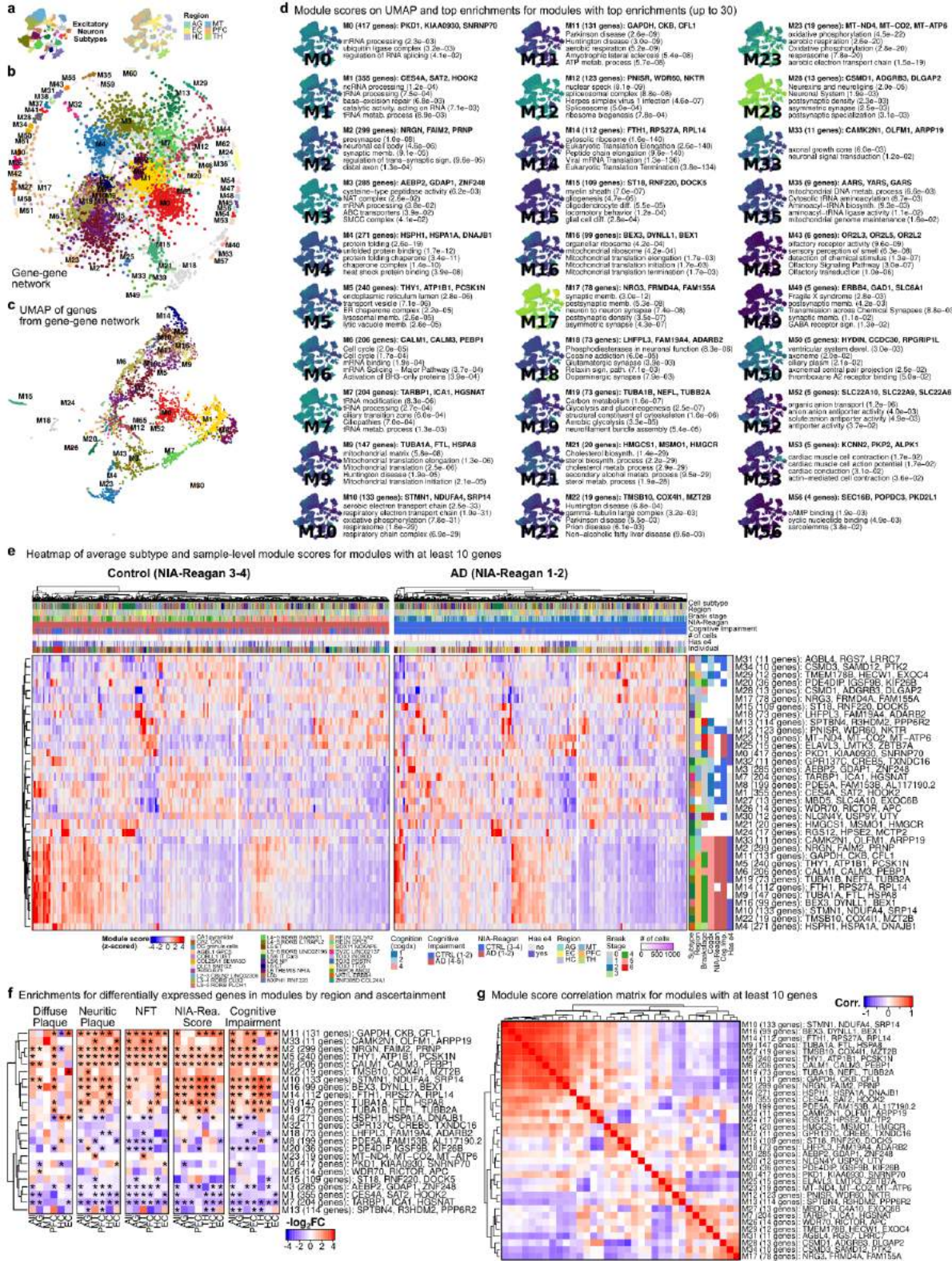




**Supplementary Figure 3-S27. Broad modules computed using all cells in the atlas. a, UMAP plots of cells used to detect the broad modules across all cells, colored by major cell type**

and by region of origin. **b**, Gene-gene network from modules method, where each point is a gene, edges connect pairs of genes with above-threshold correlation, and genes are colored by their leiden cluster (module) in the graph. **c**, UMAP of genes from estimated gene-gene correlation from sample-decorrelated matrix. Each point is a gene, colored by the leiden clusters (modules) as in **(b)**. **d**, Modules scores on the cell UMAP from **(a)** and the names and p-values of the top gene set enrichments (term size < 500) for the 30 modules with the top enrichments (or top 30 in order of size). Module scores calculated as the average (normalized, log<sub>10</sub>p) expression of all of the genes in the module. Modules are sorted by number of genes and labeled with the top three genes that are most correlated to the overall module score. **e**, Heatmap of average subtype and sample-level module scores for modules with at least 10 genes. Module scores are z-scored for the heatmap. Top panel annotates columns, giving the properties of each subtype + sample column. Side panel shows hypergeometric enrichments for top z-scoring samples ( $z > 1$ ) in each metadata column. **f**, Heatmap of hypergeometric enrichments of up (red) or down (blue) AD DEGs in modules for DEGs computed under different conditions and brain regions. Only modules with significant enrichments are shown and rows are clustered hierarchically by Euclidean distance. **g**, Correlation matrix heatmap for module scores for modules with at least 10 genes, with rows and columns hierarchically clustered by Euclidean distance.





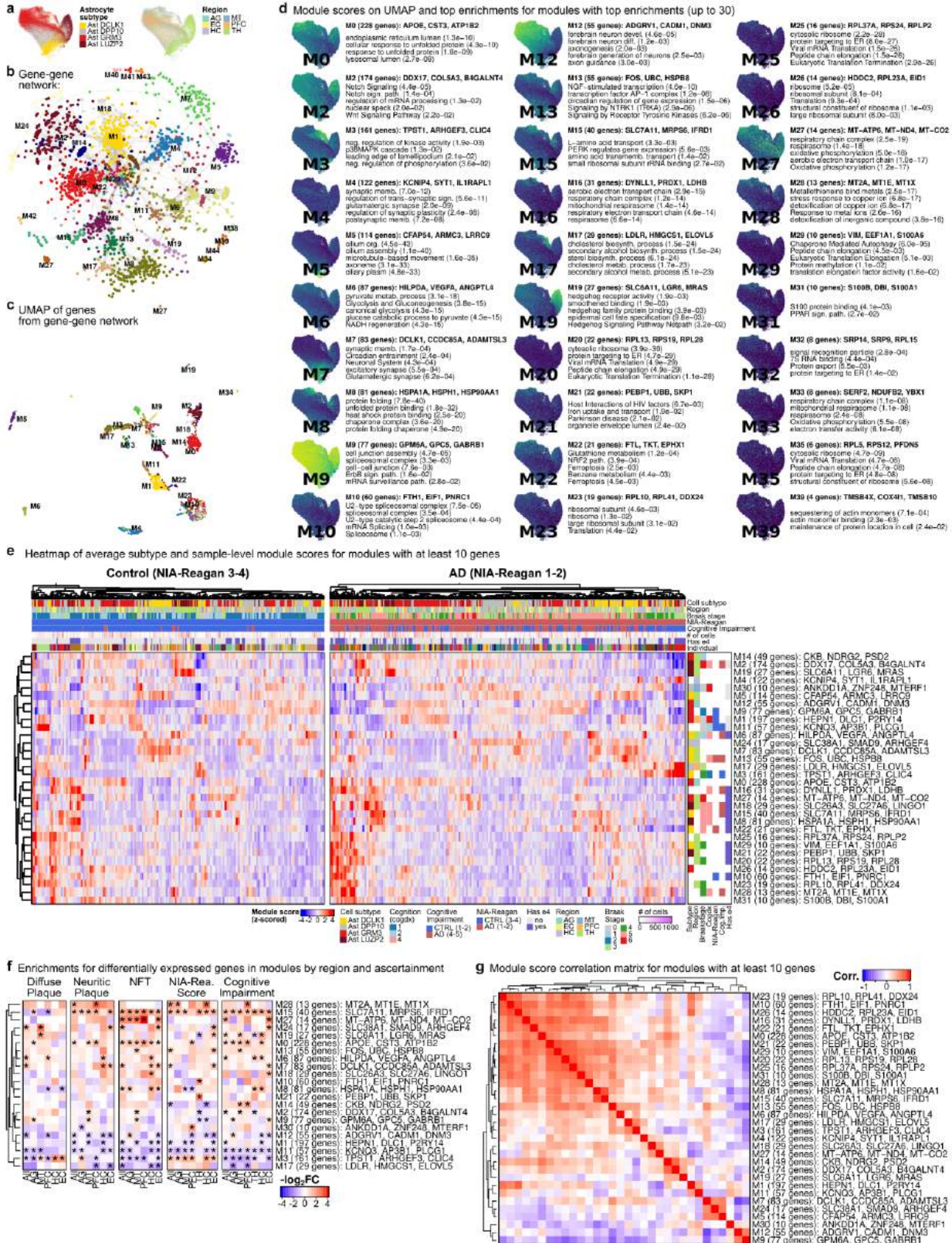
**Supplementary Figure 3-S28. Excitatory modules. a**, UMAP plots of cells used to detect modules in excitatory neurons alone, colored by cell subtype and by region of origin. **b**, Gene-gene network from modules method, where each point is a gene, edges connect pairs of

genes with above-threshold correlation, and genes are colored by their leiden cluster (module) in the graph. **c**, UMAP of genes from estimated gene-gene correlation from sample-decorrelated matrix. Each point is a gene, colored by the leiden clusters (modules) as in **(b)**. **d**, Modules scores on the cell UMAP from **(a)** and the names and p-values of the top gene set enrichments (term size < 500) for the 30 modules with the top enrichments (or top 30 in order of size). Module scores calculated as the average (normalized, log1p) expression of all of the genes in the module. Modules are sorted by number of genes and labeled with the top three genes that are most correlated to the overall module score. **e**, Heatmap of average subtype and sample-level module scores for modules with at least 10 genes. Module scores are z-scored for the heatmap. Top panel annotates columns, giving the properties of each subtype + sample column. Side panel shows hypergeometric enrichments for top z-scoring samples ( $z > 1$ ) in each metadata column. **f**, Heatmap of hypergeometric enrichments of up (red) or down (blue) AD DEGs in modules for DEGs computed under different conditions and brain regions. Only modules with significant enrichments are shown and rows are clustered hierarchically by Euclidean distance. **g**, Correlation matrix heatmap for module scores for modules with at least 10 genes, with rows and columns hierarchically clustered by Euclidean distance.



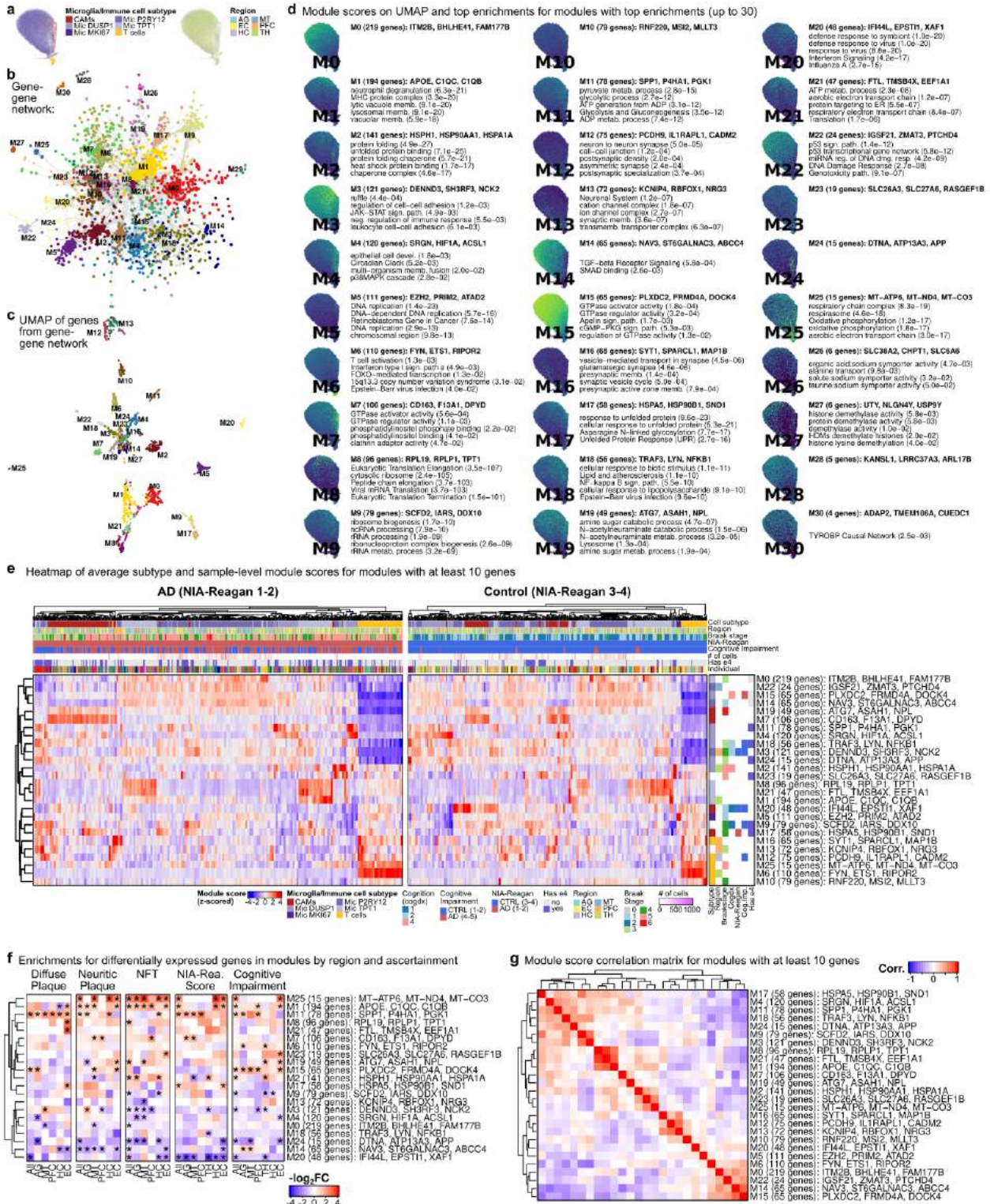






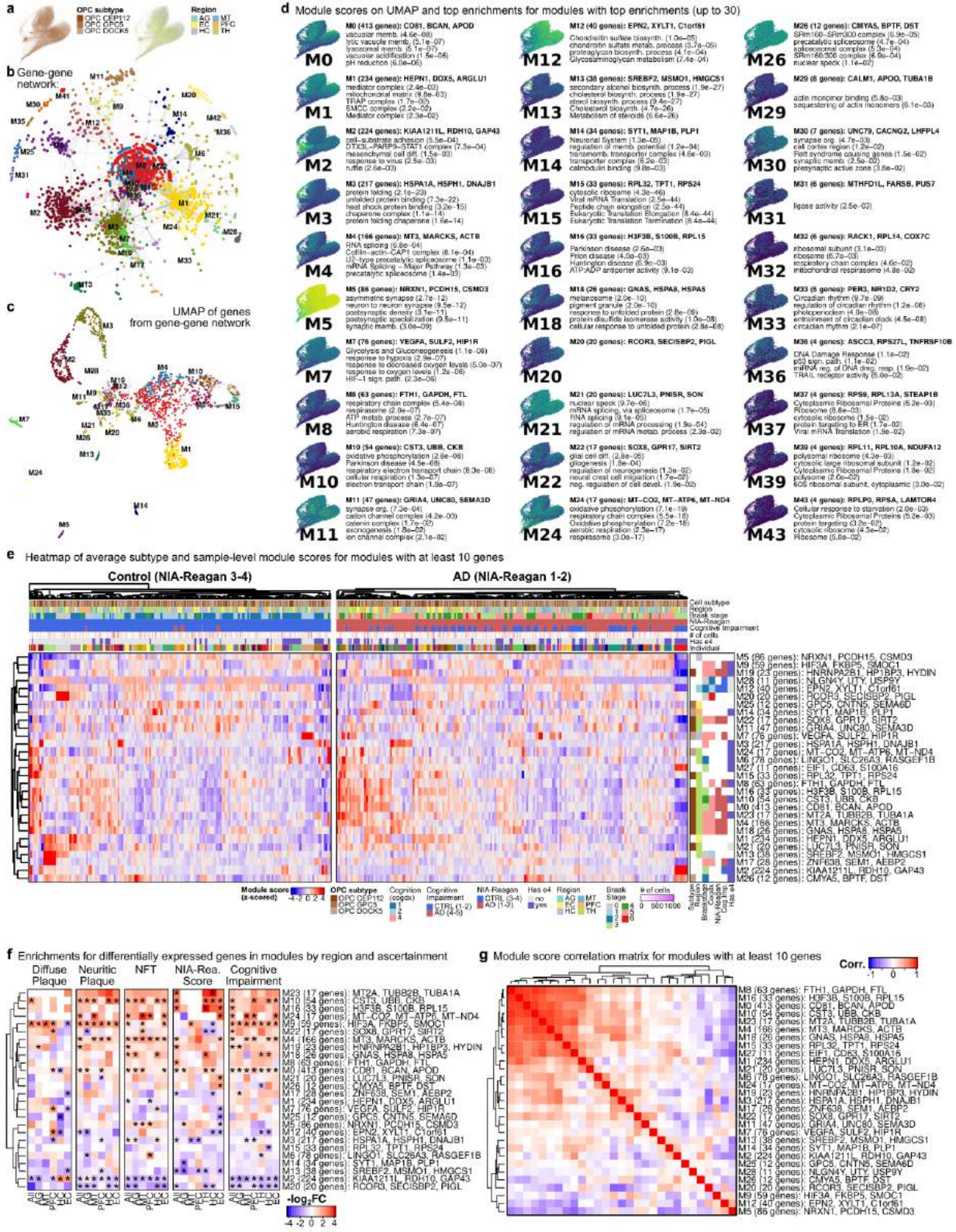
Supplementary Figure 3-S30. Astrocyte modules. Breakdown as for Ext. Fig. 3-S29.





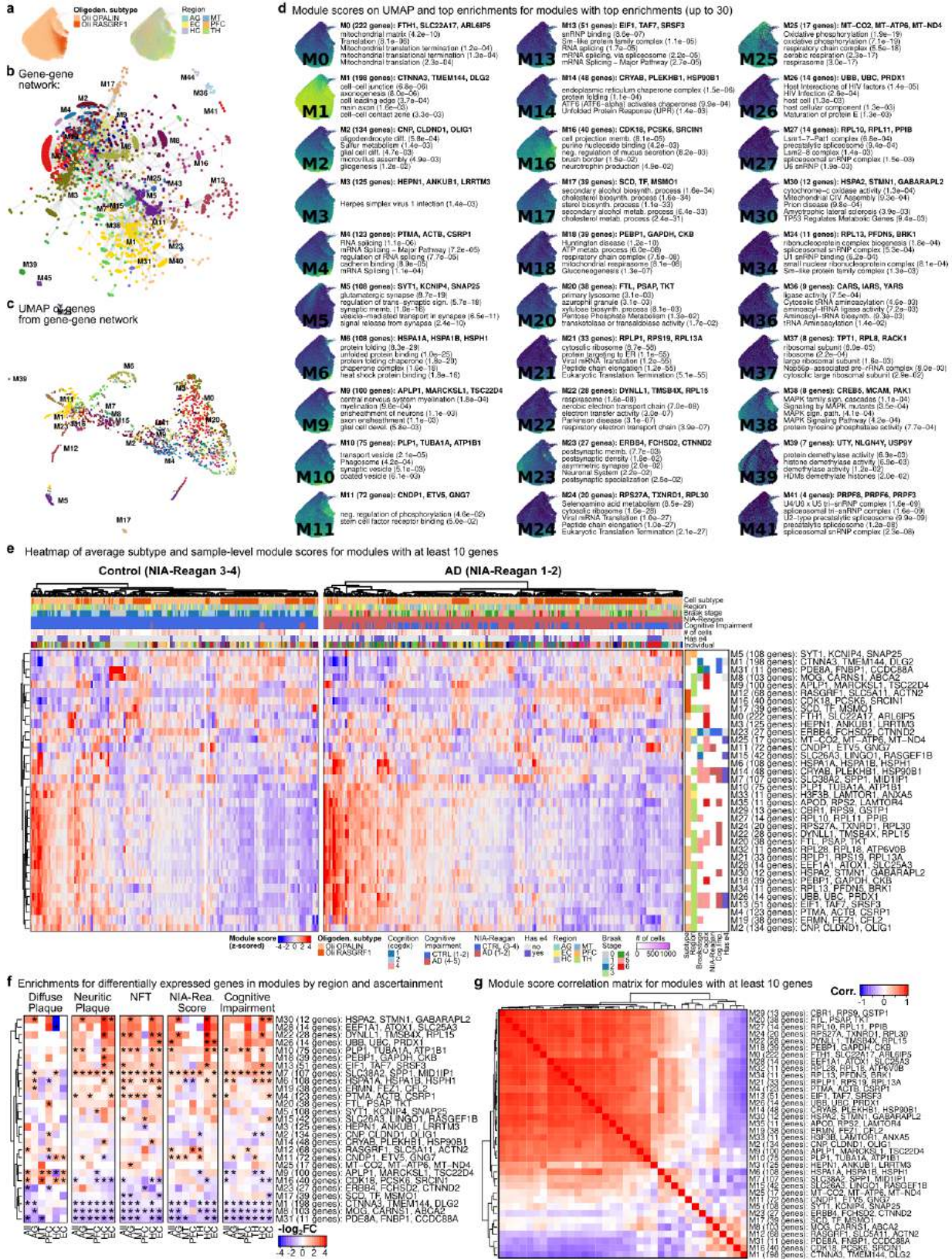
Supplementary Figure 3-S31. Microglia/Immune modules. Breakdown as for Ext. Fig. 3-S29.





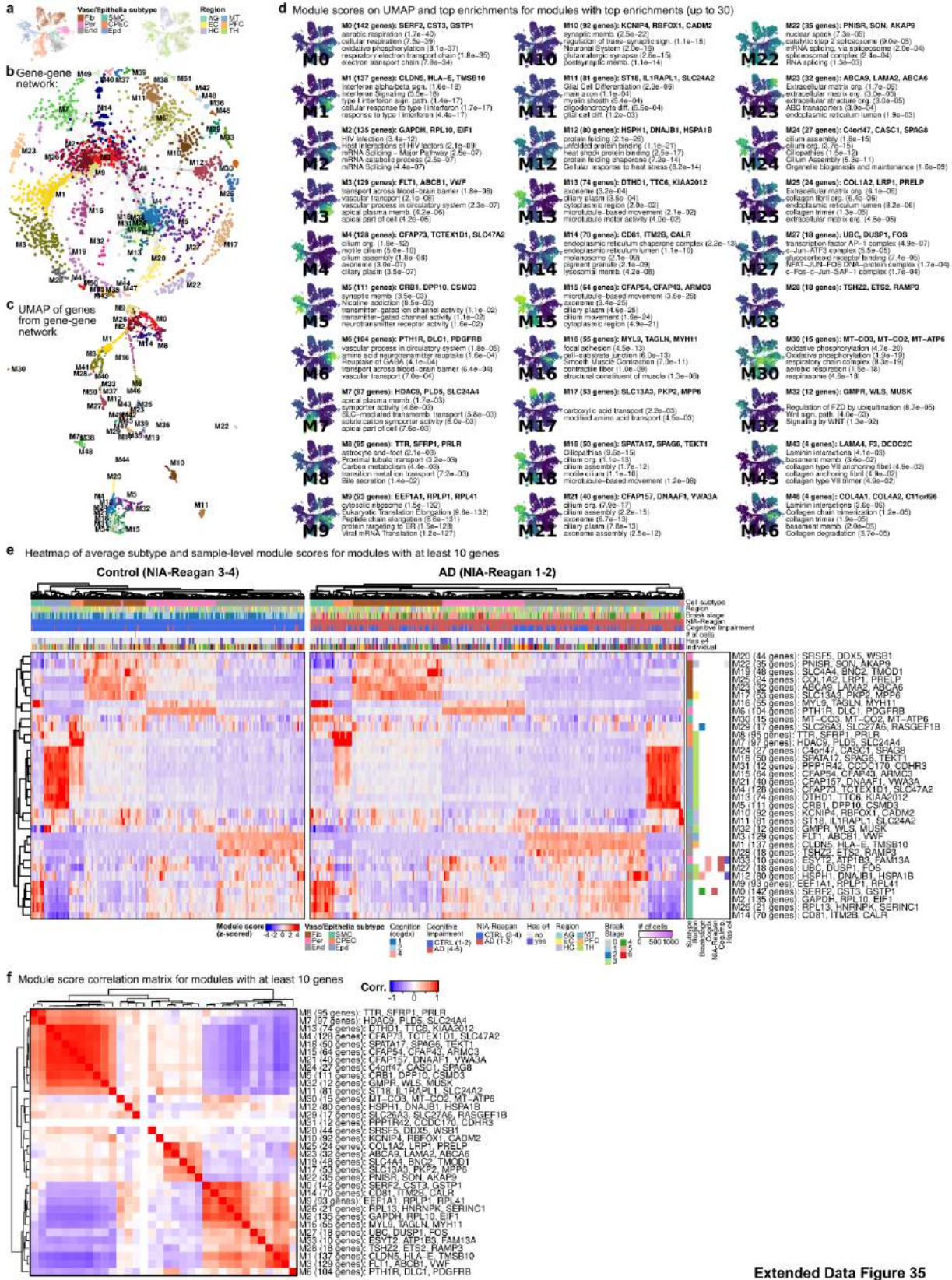
Supplementary Figure 3-S32. OPC modules. Breakdown as for Ext. Fig. 3-S29.





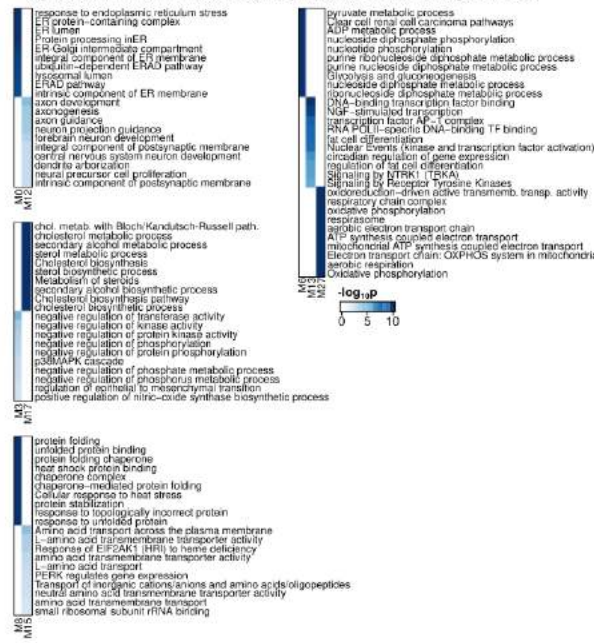
Supplementary Figure 3-S33. Oligodendrocyte modules. Breakdown as for Ext. Fig. 3-S29.



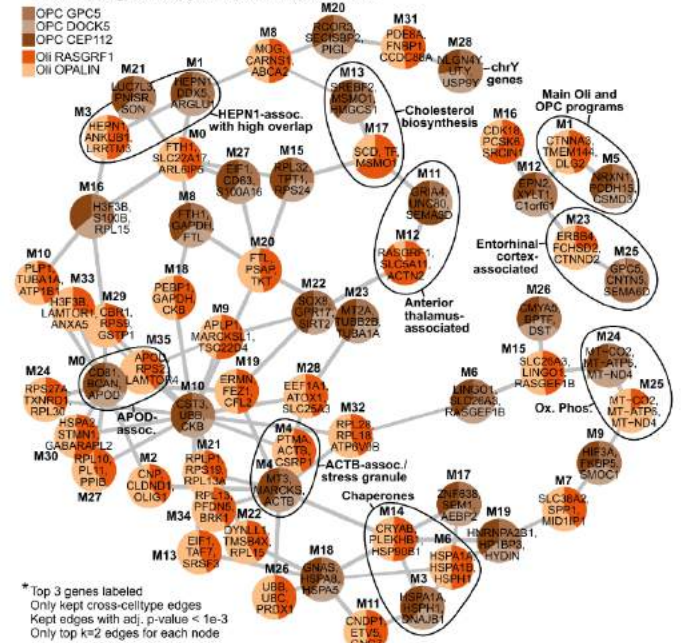




**a** Top functional enrichments for selected pairs of correlated astrocyte modules

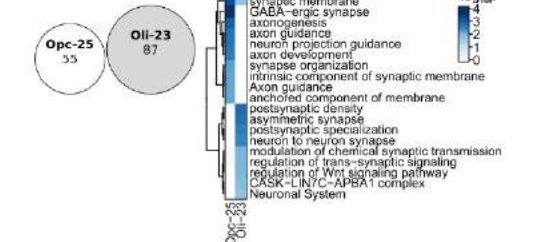


**b** OPC and Oligodendrocyte module-module network\*

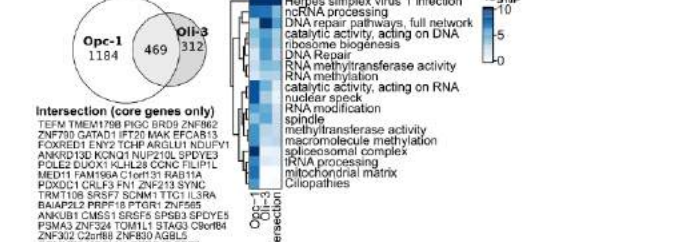


**c** Enrichments and gene set overlaps for selected correlated OPC/Oli. modules

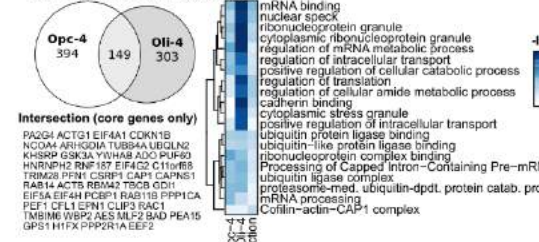
Entorhinal cortex-associated modules



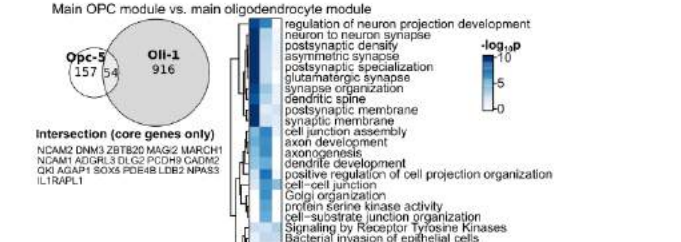
HEPN1-associated modules for OPCs/Oligodendrocytes



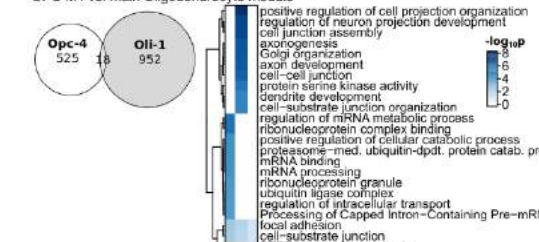
Actin / nuclear speck / stress granule modules



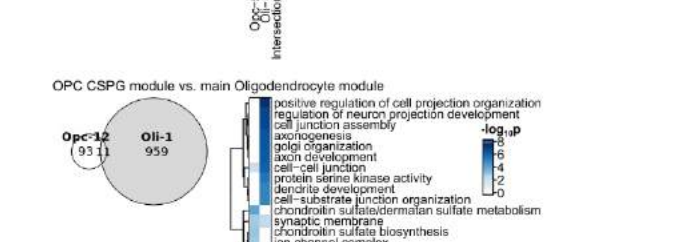
Main OPC module vs. main oligodendrocyte module



OPC M4 vs. main Oligodendrocyte module



OPC CSPG module vs. main Oligodendrocyte module



**Supplementary Figure 3-S35. Extended enrichments for specific gene expression modules.** **a**, Functional enrichments for selected astrocyte modules from **Fig. 2j**, showing top 10 functional enrichments for each pair of compared correlated modules (and for M6, M13, M27 together). Only terms with fewer than 500 genes shown. **b**, Module-module network for OPC (brown) and oligodendrocyte (orange) modules from correlation of module pairs at the sample level, keeping only edges between modules detected in different cell types (and where correlation has FDR-adjusted p-value < 0.01). Edges were further pruned to keep the top two edges for each node and the undirected graph was simplified to merge repeated edges. Nodes are colored by module's relative expression in each of the OPC or oligodendrocyte subtypes and groups highlight pairs of modules with similar function or localization. **c**, Enrichments and gene set overlaps for selected correlated OPC and oligodendrocyte modules from (**b**). Gene set overlaps for the full set of genes assigned to each of the modules shown as a scaled venn diagram (euler diagram), with intersecting core genes below. Functional enrichments are shown for each module's genes and their intersection (up to 10 terms per set).





for modules with at least 10 genes from all sets of modules (263 modules in total). Heatmaps are ordered by the hierarchical clustering of the correlation matrix and cuts represent 20 clusters cut from the hierarchical clustering dendrogram. Left and right side bars label rows by their modules set of origin (major cell type colors and gray for all cells). The most commonly shared genes in selected clusters of modules are shown on the right of the gene set overlap heatmap. **b**, Functional enrichments for each cluster of modules for the shared genes (>2 modules) in each cluster. Up to 10 enrichments shown, ordered by p-value, labeled by their source and only keeping terms with fewer than 500 genes.

## **4. Single-cell multi-region dissection of Alzheimer's disease progression maps AD-specific regional selective neuronal vulnerability**

Parts of this chapter have been presented/published as:

### **Single-cell multi-region dissection of Alzheimer's disease progression reveals selectively vulnerable neurons**

*(in preparation)*

Hansruedi Mathys<sup>1-3\*</sup>, Carles Boix<sup>4-6\*</sup>, Leyla Akay<sup>1,2</sup>, Ayesha Ng<sup>1,2</sup>, Xueqiao Jiang<sup>1,2</sup>, Kyriakista Galani<sup>4,5</sup>, Julio Mantero<sup>4,5</sup>, Neil Band<sup>4,5,7</sup>, David A. Bennett<sup>8</sup>, Li-Huei Tsai<sup>1,2,5+</sup>, Manolis Kellis<sup>4,5+</sup>

<sup>1</sup>Picower Institute for Learning and Memory, MIT <sup>2</sup>Department of Brain and Cognitive Sciences, MIT <sup>3</sup>University of Pittsburgh <sup>4</sup>Computer Science and Artificial Intelligence Laboratory, MIT <sup>5</sup>Broad Institute of MIT and Harvard <sup>6</sup>Computational and Systems Biology Program <sup>7</sup>University of Oxford <sup>8</sup>Rush Alzheimer's Disease Center. \*Equal contributors. +Correspondence to: [manoli@mit.edu](mailto:manoli@mit.edu) and [lhtsai@mit.edu](mailto:lhtsai@mit.edu)

## 4.1 Introduction

Alzheimer's disease (AD) is a highly prevalent neurodegenerative disorder, whose consistent progression through brain regions is not well understood at the level of molecular mechanism. To address this, we profiled the transcriptome of six different brain regions from a cohort of 48 individuals spanning AD progression, capturing over 1.3 million cells. Our multi-region single-cell molecular atlas of the aged human brain showcases brain cellular diversity and reveals that the composition of major cell types strongly varies across anatomical regions with a predominance of neurons in neocortical regions and glial cells in the entorhinal cortex, the hippocampus, and the anterior thalamus. We identify highly region-specific subtypes of excitatory and inhibitory neurons and glia, including an inhibitory interneuron population that is unique to the thalamus and distinct from the canonical parvalbumin, somatostatin, and VIP subclasses of inhibitory neurons. We discover thalamic astrocyte subtypes that are not observed in the neocortex and are characterized by the expression of a unique combination of genes encoding neurotransmitter transporters, ion channels, and signaling receptors. We develop a novel computational method to identify gene expression modules at the atlas scale and across multiple resolutions, revealing regional glial programs and metabolic changes in AD. Finally, we identify region-specific depleted excitatory subtypes, gene expression programs that may mediate their vulnerability, and transcriptomic changes that mediate excitatory neuron resilience and degeneration in AD. Our work charts a spatial map of AD, providing insight into cellular response and resilience to disease progression, pinpoints the weakest elements of the brain's neuronal circuitry that are depleted in patients with AD, and provides a valuable reference for future studies of the human brain.

## 4.2 Glial compositional and expression changes in AD

We investigated how cellular composition is affected throughout AD progression (**Fig. 4-1a**, binomial regression, versus individual-level AD status and region-level pathology measurements)<sup>220</sup>. At the level of major cell type, we did not observe significant compositional changes, but saw slight decreases in the number of both excitatory (OR=0.94, individuals stratified by NIA-Reagan score<sup>221</sup>), inhibitory neurons (OR=0.93), and OPCs (OR=0.85) as well as an increase in the number of oligodendrocytes (OR=1.14, adj. p-value=0.01) and vascular cells (OR=1.24), mostly driven by changes in the EC, HC, and PFC regions, especially in later stages of AD (**Fig. 4-1a**). In astrocytes and microglia/immune cells, which did not change at the major cell type level, we detected a modest enrichment of macrophages, DUSP1 microglia, and of HC and TH-enriched DCLK1 astrocytes (OR=1.39, adj. p-value=0.04) in AD<sup>222</sup>, and a depletion of cycling and of TPT1 microglia (OR=0.72, adj. p-value=0.03), especially in EC and HC. Vascular subtypes were all moderately enriched in AD, with a significant increase in the fraction TH-enriched ependymal cells (**Fig. 4-1b, Fig. 4-2a**)<sup>223</sup>.

To analyze broad changes in glial expression and function in Alzheimer's disease, we computed differentially expressed genes (DEGs) for each major cell type across regions, against both individual-level AD status (early-to-mid/late AD, early/mid-to-late AD, and cognitive impairment) and region-level pathology measurements (neuritic plaque, diffuse plaque, and neurofibrillary tangles: NFT) using a negative binomial linear mixed model framework adjusting for both known covariates and potential unknown batch effects (**Fig. 4-1c, Fig. 4-2**)<sup>224,225</sup>. Cognitive impairment showed the highest number of DEGs across cell types, followed by NFT density, and diffuse plaque showed relatively few changes (**Fig. 4-1c**). At the region level, excitatory and inhibitory neuronal DEGs clustered by region, with relatively high within-region overlap across conditions, suggesting that diverse neuronal subtypes are also differently affected by AD. Glial DEGs clustered by condition, with high overlap for DEGs computed relative to NFT and neuritic plaque density and cognitive impairment, suggesting that transcriptomic responses to these pathologies are closest to end-stage cognitive phenotypes (**Fig. 4-2b**).



A large fraction of DEGs were also shared across cell types (in 3+ cell types: 972 genes in NFT, ~35% per cell type), and showed up-regulation of response to unfolded proteins (CLU, HSP90AB1, CANX), RNA splicing regulation (SRSF3/6, FXR1), and ubiquitin binding (UBA1, UBB, UBC), and cohesin loading (RAD21, SMC3, STAG2, NIPBL) (**Fig. 4-1d**). Top cell type-specific DEGs were broadly consistent across regions (**Fig. 4-1e**), and included up-regulation of phagosome components and antigen presentation in microglia, including well-established DEGs such as APOE, HLA-DRA, and CD74 as well as TFs PPARG (microglia polarization) and FOXP1; of glial cell differentiation in OPCs (OLIG1/2 and MARCKS); and of intracellular transport and cadherin binding in oligodendrocytes (ANLN, CDH20, SVIP); and of synaptic signaling and transport vesicles in excitatory neurons (NCDN, NEUROD2, NPTX1) (**Fig. 4-1d,e, Fig. 4-2**).

We asked where and when GWAS-linked genes play a role in Alzheimer's pathology, intersecting DEGs for regional pathology measurements with 149 known familial AD and GWAS genes<sup>50,77</sup>. We found that 49% of genes were DE for least one cell type (**Fig. 4-1f**) and multiple genes showed region-specific expression, including lipid transporter ABCA7 (thalamus-enriched)<sup>226</sup>, zinc-finger protein ZNF655 (entorhinal cortex)<sup>227</sup>, and complement receptor CR1 (neocortex)<sup>77</sup>. GWAS genes were top expressed (75 genes) and differential (30) in microglia and 25 genes were differentially expressed in at least three cell types, including up-regulated CLU, PLCG2, and SORT1 and down-regulated DENND6A (**Fig. 4-3**). Among all cell types, astrocytes and microglia responded the most to diffuse plaques and showed the largest differential changes in regions with high neuritic plaque density, while neurons and oligodendrocyte-lineage cells showed more changes relative to NFT density.

#### **4.3 Expression modules separate AD-associated changes across glial cell types**

We next used the modules as a framework to interpret the large number of DEGs by grouping them into functional and identity programs (**Fig. 4-1g**). Across glia, we found a number of commonly AD-perturbed pathways, including chaperone-containing modules, which were both highly correlated across cell types and highly enriched for up-regulated DEGs in samples with high NFT density. Similarly, we identified independent glycolysis (astrocytes, microglia, and OPCs) and oxidative phosphorylation (all cell types) modules, each strongly cross-cell type correlated and highly upregulated in AD across ascertainment variables (**Fig. 4-1g, Supplementary Fig. 3-S36, 4-S2**).

Astrocyte-specific responses to AD included up-regulation of plaque-responding/interacting genes in M2 (ITM2B, ITM2C, CLU), of stress-associated transcription factors (FOS, JUND) and ubiquitin-related genes (UBC, USPs) in M9, and of zinc and copper metallostasis genes (MT1E/G/M/X and MT2A) in M5. Normal astrocyte functions, including extracellular matrix, adhesion, and neurogenesis-related genes were depleted in AD (M1, M15) while the DCLK1 and DPP10 identity modules (M7, M24) were enriched for up-regulated AD DEGs, consistent with DCLK1 compositional increase in AD (**Fig. 4-1g, Supplementary Fig. 3-S27f**).



**Figure 4-1. Gene module alterations.** **a.** Compositional changes for major cell types in AD by quasi-binomial regression with FDR-correction.  $\log_2$  OR shown both for each AD variable across regions (left) and for each region in late-AD (right). **b.** Compositional changes for glia in AD, as in (a). **c.** Number of DEG enrichments for each cell type for both region-level pathology measurements and individual-level AD ascertainment. **d.** Top 5 pathways for DE genes in each major cell type and for DEGs shared across cell types ( $\geq 3$  cell types) versus regional NFT levels. **e.** Regional effect sizes for the top eight up/down DEGs in glia versus individual-level cognitive impairment status. **f.** Expression level by region/subtype and effect sizes of top DE GWAS genes across regional pathology measurements. 29 unique genes are shown (top 15 per pathology measurement). **g.** DEG enrichments for gene expression modules in glia. **h.** Scatterplots and correlation of scores for selected pairs of microglia/immune AD-associated modules. Each dot represents scores at the subtype and sample-level and is colored by region of origin. **i.** Effect sizes of DEGs in the glial glycolysis modules, in each of the three glial subtypes with a glycolysis module and for each AD ascertainment variable. **j.** Correlation of glial glycolysis and oxidative phosphorylation module scores at the sample level. **k.** Heatmap of individual-level average module scores for five glial metabolism-associated modules, with columns clustered by Euclidean distance. Bottom panel shows APOE e4 carriers and cognitive status for each individual. **l.** Scatterplots and correlation of sample-level scores for two pairs of oligodendrocyte/OPC modules associated with glial metabolism. **m.** Boxplots of module scores at the individual-level split by APOE e4 carrier status for the glial metabolism-associated modules.

In immune cells, we detected a consistent AD up-regulation of genes in DAM-like M1, including APOE, CD81, C1Q, CD74, TREM2, GRN, MHC-II genes, and lysosomal genes CST3, LAPTM5, and multiple cathepsins, and of genes within the shared microglia/macrophage module M15, including cell cycle driver CCND3 and GWAS genes SORL1, PICALM, and MEF2C. Microglia also up-regulated oxidative phosphorylation (M25) and chaperones (M2) in response to NFT, glycolysis (M11), and down-regulated interferon expression (M20, including IFI44L, XAF1) (**Fig. 4-1g, Supplementary Fig. 3-S28f**). These AD DEG-enriched changes did not co-occur at the sample-level, as modules showed little correlation with each other across regions and individuals, or at the cell level, as M1 (DAM) was TPT1+ enriched and M11 (glycolysis) and M2 (chaperones) were DUSP1+ enriched (**Fig. 4-5f, 3h**). Notably, while increases in DUSP1 microglia were tied to M2 and M11 DEGs, subtype composition did not drive differential expression, as expression of M1 genes increased while its subtype (TPT1+) was depleted in AD (**Fig. 4-1a**). Regional measurements of AD pathology separated modules, as M15 (shared program), M11 (glycolysis), and M20 (interferon) responded to high burden of diffuse plaque whereas M1 (DAM), M2 (chaperones), and M25 (ox. phos.) instead responded the most to NFT burden (**Fig. 4-1g**).

Glycolysis modules in astrocytes and OPCs were also highly co-expressed with the microglial module, and each contained a core set of genes including PDK1, PFKL, PFKP, PGK1, TP11, LDHA, DDIT4, VEGFA, and BNIP3L, which showed the highest AD up-regulation in astrocytes and microglia and in regions with diffuse plaque (**Fig. 4-1i-j**). In addition, all non-neuronal cell types, regardless of the presence of metabolism-associated modules, up-regulated core glycolysis driver GAPDH and mitophagy-regulating BNIP3L (and HK2 in microglia<sup>228</sup>) in response to NFTs and in cognitively impaired individuals, suggesting that the metabolic shift in AD may be tied to mitophagy (**Fig. 4-1i, Supplementary Fig. 3-S36, 4-S3**)<sup>229</sup>. Notably, astrocytes uniquely up-regulated DDIT4 and VEGFA, side-by-side with genes that suppress fatty-acid metabolism (ANGPTL4) and promote lipid droplet storage of fatty-acids (HILPDA),

while microglia up-regulated multiple glycogen-related genes (GBE1, UGP2, PYGL), highlighting the complexity of the glial metabolic response to AD (**Fig. 4-1i-j**).

While we did not observe a significant increase of OPC glycolysis in AD, OPC M9 (marked by FKBP5, HIF3A) was AD DEG-enriched across conditions, in particular with respect to diffuse plaque (**Fig. 4-1g,k-l**). This module was strongly correlated with genes enriched for energy-related signaling (insulin, AMPK, adipocytokine, FoxO) which included FOXO1, INSR, and AKT3, and up-regulated EGFR, PIK3R1, and ADIPOR2 in AD (**Fig. 4-4c**). The module was most correlated with oligodendrocyte M7 ( $r=0.63$ ), which shared its enrichment for adipocytokine signaling (ADIPOR2, ACACB) and up-regulated many genes in AD including glutamine (SLC38A2) and zinc (SLC39A11) importers, apoptosis/mitophagy-inhibitor BCL2, and lipogenesis-promoting MID1IP1 (**Fig. 4-1g, 3l, Fig. 4-4**). Interestingly, all of the metabolic modules were similarly expressed across regions in each individual and many had significantly higher expression in e4 carriers, suggesting that the e4 allele may accentuate the glial metabolic shift in AD (**Fig. 4-5f, 3m**). We have provided an interactive website for exploring these modules and the atlas as a whole at [http://compbio.mit.edu/ad\\_multiregion](http://compbio.mit.edu/ad_multiregion).

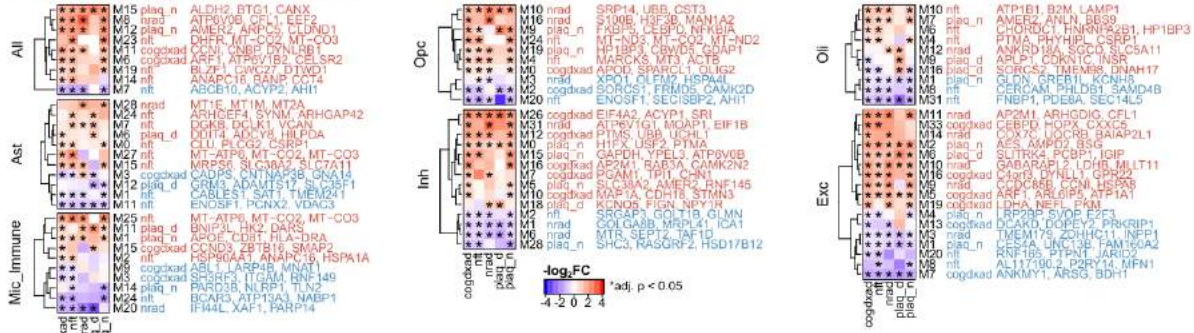




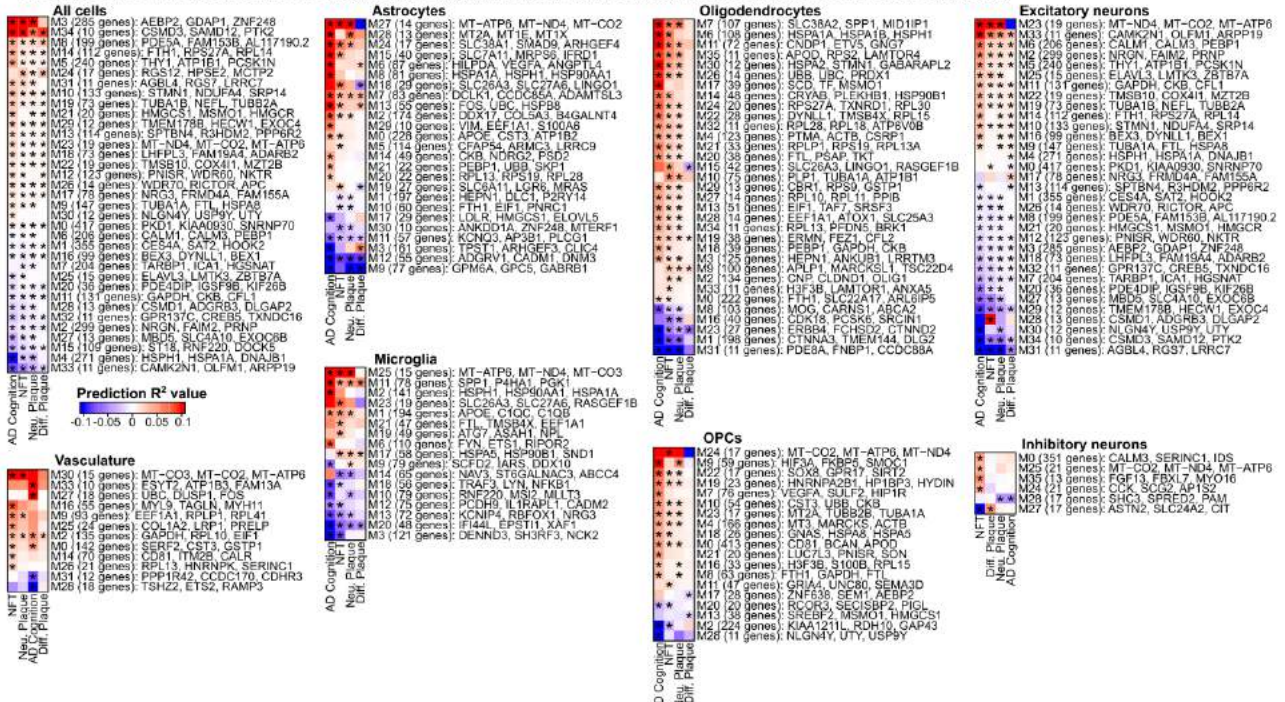
**Figure 4-2. Compositional and expression changes in AD in major cell types.** **a**, Compositional differences in glial in AD by NIA-Reagan score 1-2, late Braak stage (5-6), impaired cognition (4-5), or presence of plaques or neurofibrillary tangles. Right panel shows the per-region compositional changes in late AD. Gray boxes indicate interactions that are not computed, with fewer than 500 cells across individuals. **b**, Overlap of AD DEGs in each major cell type for each combination of region and condition (AD variable). DEG overlap computed by Jaccard distance and rows/columns hierarchically ordered by Euclidean distance. **c**, Top functional enrichments for unique DEGs and DEGs shared across cell types ( $\geq 3$  cell types) across major cell types by condition (up to the top 5 terms with fewer than 500 genes shown per gene set).



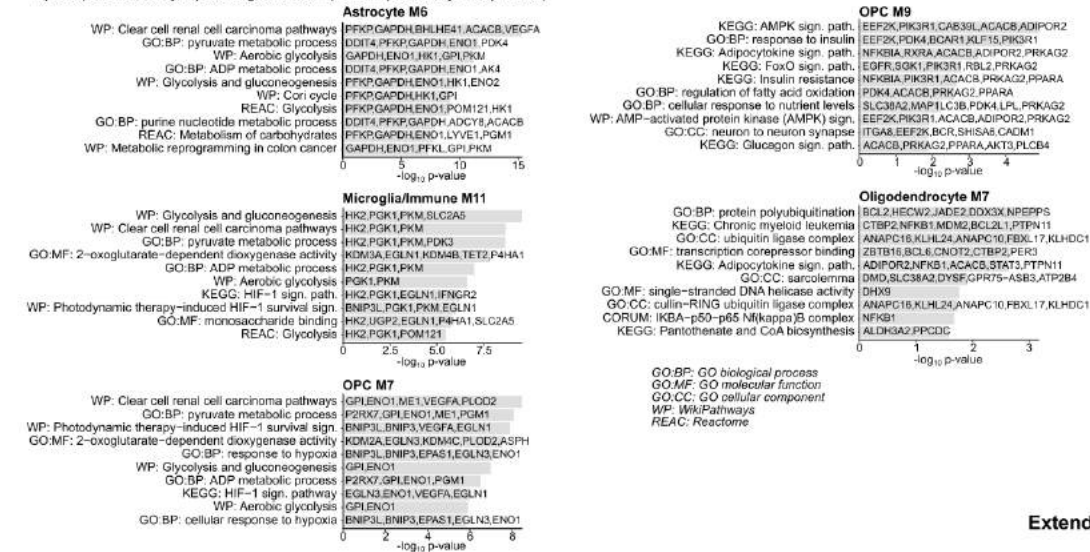
**a AD DEG enrichments in gene expression modules**



**b R-squared values for predictions of AD status (cognitive level or regional pathology variables) from gene expression module scores at the sample level**



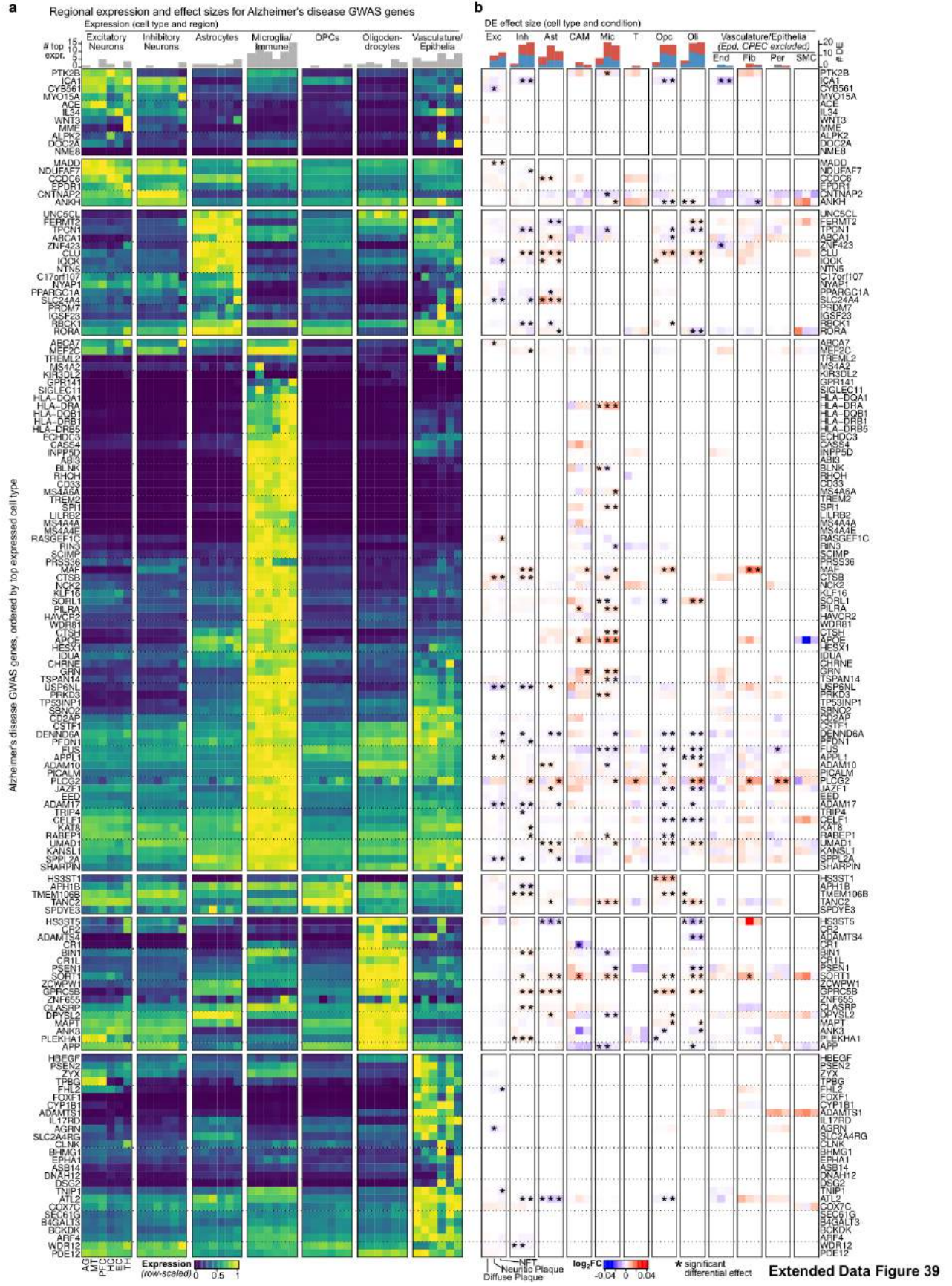
**c Top 10 functional enrichments for all genes associated with each glial metabolism-associated module**  
Up to top 5 differentially expressed genes shown per term (ordered by DEG p-values)



Extended Data Figure 40

**Figure 4-3. Module scores versus Alzheimer's disease status.** **a**, Heatmap of hypergeometric enrichments of up (red) or down (blue) AD DEGs in modules for DEGs in all sets of modules across all regions, by AD condition. Only modules with at least two significant enrichments are shown and rows are clustered hierarchically by Euclidean distance. **b**, Heatmap of r-squared values for predictions of AD status (cognitive level or regional pathology variables) from gene expression module scores at the sample level. Only modules with significant predictions are shown. rows are clustered hierarchically by Euclidean distance. **c**, Top 10 functional enrichments for all genes associated with each glial metabolism-associated module. Only terms with fewer than 500 genes shown. Terms were pruned by removing terms that had Jaccard similarity > 0.75 with another more significantly enriched term.





Extended Data Figure 39

**Figure 4-4. Alzheimer's disease GWAS-linked genes in the multi-region atlas. a,** Expression level by region/subtype and effect sizes of 150 Alzheimer's disease candidate risk genes from Alzheimer's disease GWAS risk loci. **b,** Differential effect sizes and significance for each candidate risk gene in each minor cell type across regional pathology measurements. Ependymal cells and CPEC cells were excluded as the thalamus does not have regional pathology measurements.

#### 4.4 Identification of selectively vulnerable excitatory neuron subtypes

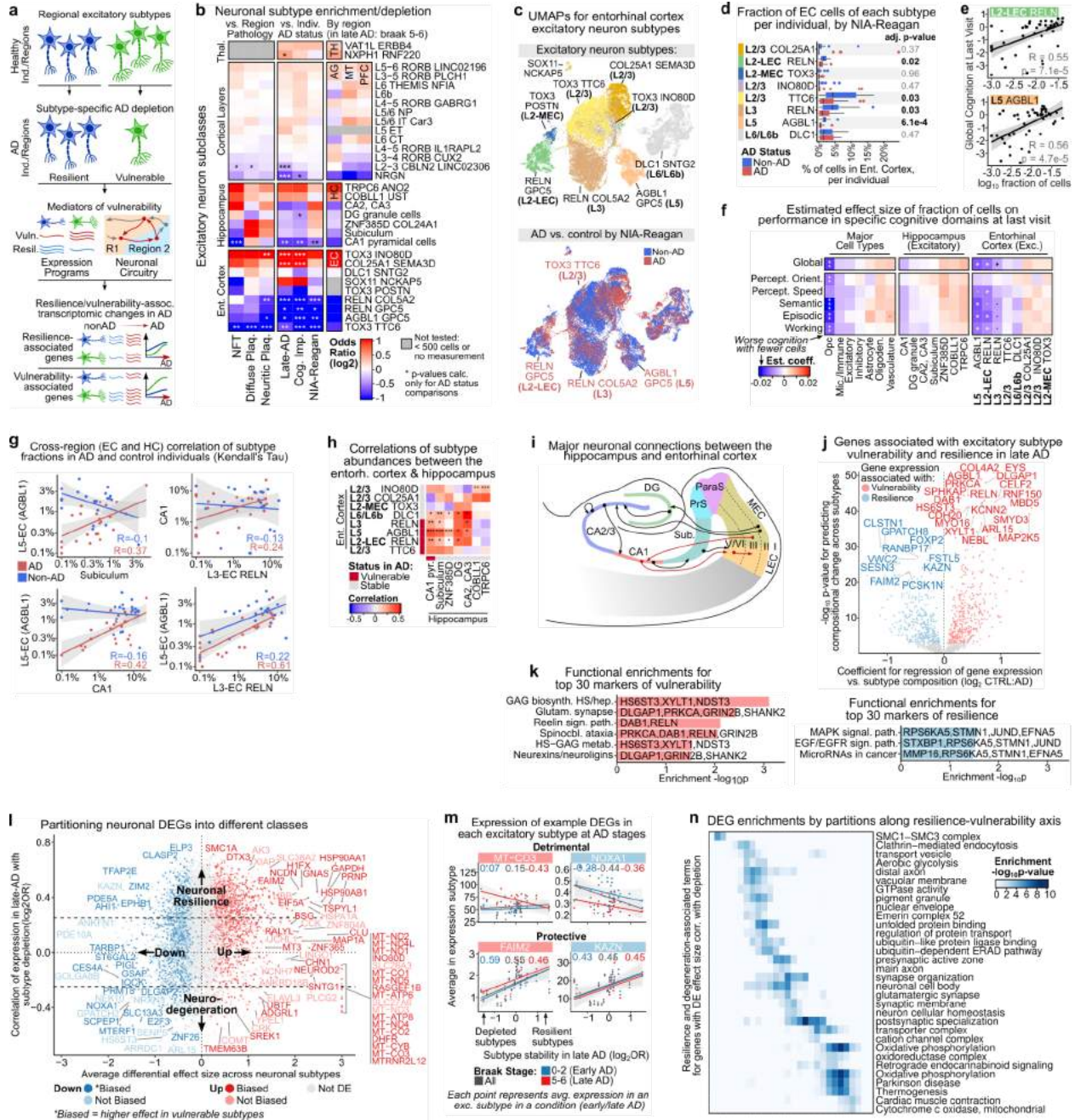
Finally, we asked how the regional progression of AD affected the diverse, region-specific neuronal subtypes identified in our atlas (**Fig. 4-5a**). Excitatory subtype composition was significantly affected in AD, with one hippocampus-specific and four entorhinal cortex-specific subtypes significantly depleted across AD individuals and with L2-3 neocortical subtypes significantly depleted in late AD, consistent with AD's regional progression (**Fig. 4-5b**). Inhibitory neurons were depleted less and later on than excitatory neurons, with multiple subtypes, including two SST subtypes, significantly depleted in late AD and in samples from regions with high NFT density (**Fig. 4-6a,b**). The five depleted hippocampus and entorhinal cortex excitatory subtypes (OR=0.38-0.66 by NIA-Reagan score), included L2 RELN+ lateral EC, L3 RELN+, L5, and L2/3 (TOX3 TTC6) neurons in the entorhinal cortex, as well as CA1 pyramidal neurons in the hippocampus (**Fig. 4-5c,d, Supplementary Fig. 3-S11**)<sup>40</sup>. Individuals with lower percentages of these vulnerable excitatory neuron subtypes performed significantly worse on cognitive tasks, with strongest observed impacts on episodic and global cognitive function for GPC5+ neurons (L2 RELN+ and L5 AGL1+) (**Fig. 4-5e,f**)<sup>29</sup>. Notably, while overall excitatory fraction was not predictive of cognition, depletion of OPCs across regions, and in particular in non-neocortex regions, was significantly associated with impaired cognition (**Fig. 4-1a, 4f, Fig. 4-6c**).

Given that these subtypes lie in highly interconnected regions, we next asked whether projection-connected neuronal subtypes between regions were coordinately depleted or expressed shared programs that might mediate subtype vulnerability. We found that the abundances of vulnerable connected neuronal subtypes were correlated between the hippocampus and entorhinal cortex in AD but not control individuals, with the strongest effects seen for entorhinal L5 against L5-projecting subiculum (Kendall's tau=0.37 in AD; -0.1 in control individuals) and CA1 (0.42/-0.16) hippocampus neurons and for CA1 against L2-LEC (0.26/-0.07) and L3 RELN+ (0.24/-0.13) entorhinal cortex neurons, both of which project in part to CA1 (**Fig. 4-5g-i, Supplementary Fig. 3-S20d,e**)<sup>20,39,230,231</sup>. Separately, we found 391 genes whose expression in control individuals was predictive of subtype vulnerability (log<sub>2</sub>OR in late-AD) (**Fig. 4-1a, 4j**), including reelin signaling pathway genes RELN and DAB1, which transduce signal through APOE receptors VLDLR and ApoER2<sup>232</sup>, kinase-associated MAP2K5, PRKCA, and SPHKAP, and heparan sulfate proteoglycan (HSPG) biosynthesis genes HS6ST3, XYLT1, and NDST3, a pathway which has been implicated in Alzheimer's pathogenesis through the identification of protective APOE variants (**Fig. 4-5j,k, Fig. 4-6**)<sup>24,233,234</sup>. We also found 346 resilience-associated genes enriched for MAPK and EGF signaling pathways, including CLSTN1, which mediates APP transport<sup>235</sup>, and PCSK1N, which has been proposed to impede amyloid aggregation, both of which were among 44 resilience-associated genes up-regulated in AD (**Fig. 4-5j,k, Fig. 4-6**)<sup>236</sup>. Interestingly, while RELN expression, typically restricted to inhibitory neurons, was highly specific to two EC excitatory subtypes, its downstream partner DAB1 was present across subtypes and its expression was highly correlated with that of VLDLR (r=0.56) and FYN (0.52), which phosphorylates DAB1 in reelin signaling (**Fig. 4-6g**). Combined expression of these last three genes was highest in L2-LEC, L2-MEC, and CA1 pyramidal cells, and in L5/6 NP and L2-3 neurons in the neocortex, suggesting that while depleted neocortical

L2-3 and hippocampus CA1 neurons do not express RELN, they lie downstream of its signaling (**Fig. 4-6h**). Overall, these AD-specific co-depletions and marker genes for subtype vulnerability suggest that the circuit formed by the hippocampal CA1 subfield and entorhinal cortex L2/3 and L5 could be studied jointly as an AD-vulnerable connected component in the context of specific expression programs such as HSPG and reelin signaling.

We then computed DEGs for each subtype to ask both what transcriptomic changes the vulnerable subtypes incur at different stages of AD and whether these changes are consistent in excitatory neurons across regions (see Methods) (**Fig. 4-7a-b**)<sup>224,225</sup>. Neuronal DEGs showed high consistency across vulnerable subtypes, up-regulating pathways including ER protein processing, oxidative phosphorylation, and iron uptake as well as AD risk-associated CLU, PLCG2, BIN1, APP, and CTSB across conditions (**Fig. 4-7b**). To understand when differential expression changes occur in neurodegeneration, we partitioned DEGs shared across all subtypes (DE in  $\geq 25\%$  of subtypes) into depletion-associated (higher expression in vulnerable subtypes in late-AD) versus resilience-associated genes (higher in resilient subtypes) (**Fig. 4-5l,m**). Resilience-associated DEGs were enriched for a diverse set of functions, including ubiquitin ligase binding (UBB, UBC, UBA52), heat shock family chaperones (HSP90AA1, HSP90AB1), ER protein processing (CALR, CANX, B4GAT1), and mediators of neuronal death (FAIM2, XIAP, CDK5R1, SNCB, PRNP), whereas depletion-associated DEGs were only highly enriched for mitochondrial oxidative phosphorylation but included CRK and NEUROD2, both associated with reelin signaling (**Fig. 4-5m, Fig. 4-7**)<sup>38,39,41</sup>. Finally, the DE effect sizes of resilience-associated genes were often higher in vulnerable cell types, possibly due to selective loss of non-expressing cells or an up-regulation of these genes in remaining neurons (**Fig. 4-7**). Interestingly, these genes with biased DE effects showed additional enrichments for aerobic glycolysis (including PGK1, LDHB, and SLC2A3) and clathrin mediated endocytosis (including AP2M1/AP2S1, OCRL, and COPS8), suggesting two pathways that may enable vulnerable neurons to resist neurodegeneration (**Fig. 4-5n**).

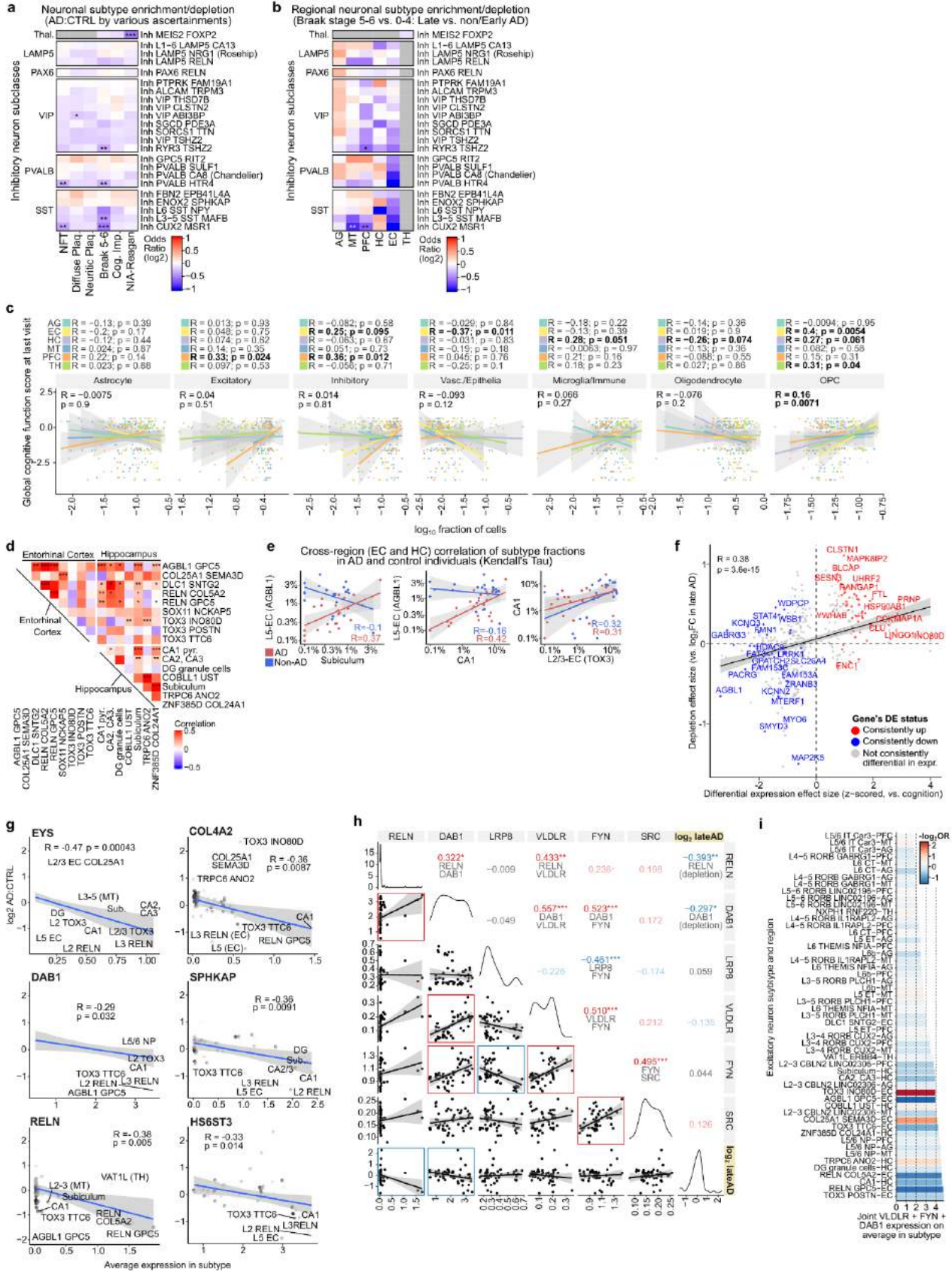




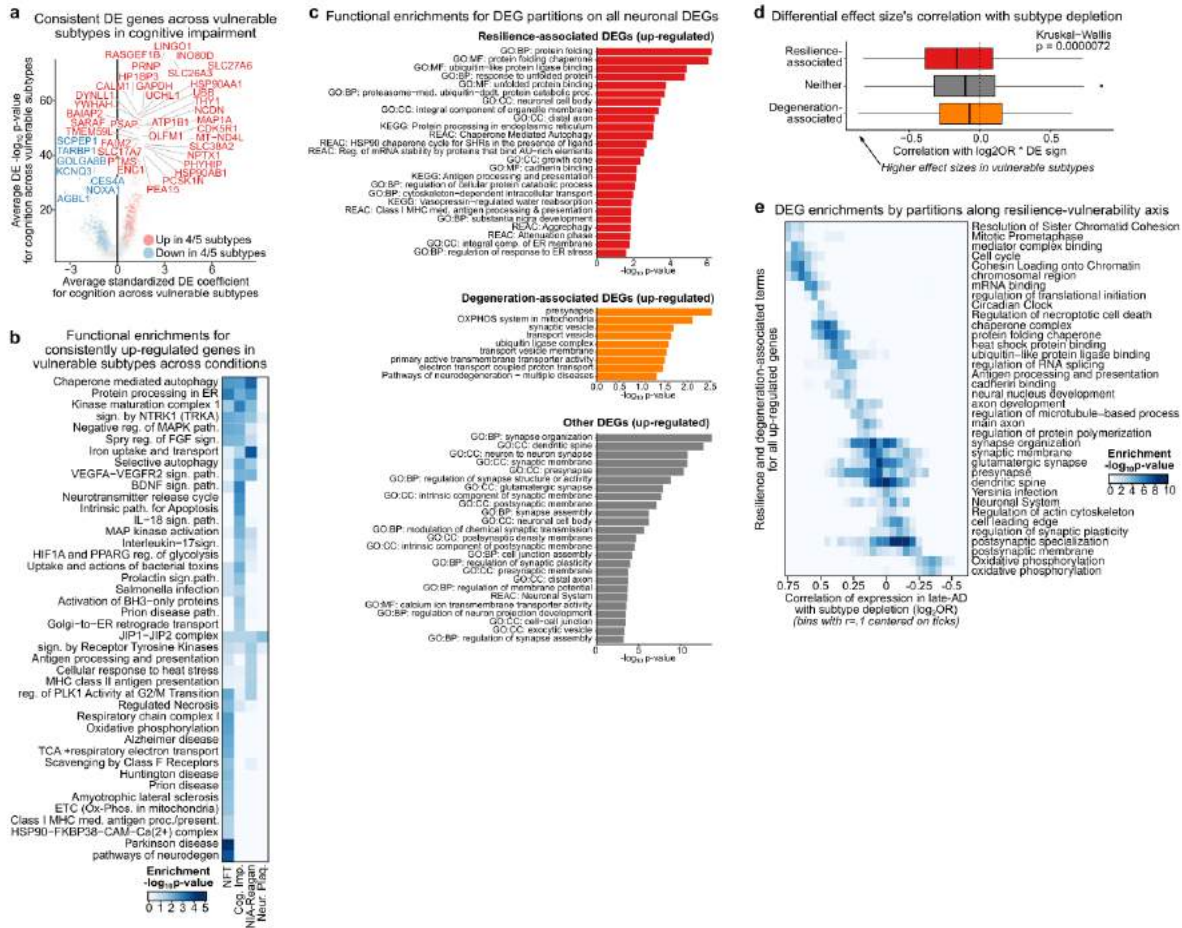
**Figure 4-5. Neuronal Vulnerability in Alzheimer's Disease.** **a.** Overview of neuronal vulnerability analyses. We identify AD-vulnerable neuronal subtypes and putative mediators of vulnerability and tie transcriptomic changes to neuronal resilience and vulnerability to AD. **b.** Compositional differences in excitatory neuron subtypes in AD by NIA-Reagan score 1-2, late Braak stage (5-6), impaired cognition (4-5), or presence of plaques or neurofibrillary tangles. **c.** Entorhinal cortex excitatory subtype UMAP by neuronal subtype and by AD (NIA-Reagan 1-2 vs. 3-4). **d.** Boxplots of fraction of neurons in each individual, split by AD status, FDR-corrected p-values from one-sided wilcoxon test. **e.** Estimated effect size of cell fraction ( $\log_{10}$ ) on scores for performance in various cognitive domains at last visit and combined scores from all domains (global). Linear regression FDR-corrected p-values (\*\*p-adj. < 0.01, \* < 0.05, dot is < 0.1). **f.** Correlation between global cognition at last visit and cell fraction for two AD-vulnerable



entorhinal cortex subtypes. **g.** Scatter plot and correlations (Kendall's tau) of subtype fraction between four pairs of vulnerable neuronal subtypes between the hippocampus and entorhinal cortex. **h.** Correlation between subtype fraction between the hippocampus and entorhinal cortex in the same individuals (\*p-adj. <0.05). **i.** Hippocampus and entorhinal cortex diagram highlighting putatively vulnerable pathways. **j.** Genes associated with subtype vulnerability or resilience across all brain regions. Linear regression between normalized sample+subtype-level gene expression and  $\log_2$  odds ratio for late-AD (Braak stages 5-6 vs. 1-4) with FDR-corrected p-values. **k.** Functional enrichments and intersected genes for top 30 markers of vulnerability and resilience (terms with < 500 genes). **l.** Scatterplot of each tested gene's average DE effect size versus its expression's correlation with subtype depletion in late-AD. **m.** Examples of resilience-associated and depletion-associated genes (from l). Each dot represents the gene's average expression in an excitatory subtype in either early-AD, late-AD, or across all individuals. Linear regression lines and CI are shown for each set of points. **n.** Functional enrichments for each identified class of neuronal DEGs (terms < 500 genes).



**Figure 4-6. Neuronal compositional changes and vulnerability in Alzheimer's disease.** **a**, Compositional differences in inhibitory neuron subtypes in AD by NIA-Reagan score 1-2, late Braak stage (5-6), impaired cognition (4-5), or presence of plaques or neurofibrillary tangles. Gray boxes indicate interactions that are not computed due to MEIS2 FOXP2 specificity to the thalamus, where we do not have measured regional scores. **b**, Compositional differences in inhibitory neuron subtypes in late AD (Braak Stage 5-6 vs. 1-4) in each region. Gray boxes indicate interactions that are not computed due to subtype regional specificity. **c**, Scatterplots of major cell type  $\log_{10}$  fraction in each of the six regions (colors) against global cognitive function performance at last visit. Lines represent best linear regression fits with 95% CIs, overall correlation and p-value are given in black in the scatter plot and by-region correlations and p-values are given above the plots, with nominally significant effects in bold. **d**, Full correlation matrix between subtype fraction between the hippocampus and entorhinal cortex in the same individuals, as described in the methods (\*p-adj. <0.05). **e**, Example subtype fraction interactions conditioned on AD status (NIA-Reagan). **f**, Scatterplot of genes' average differential expression effect size across neuronal subtypes versus their effect size in predicting subtype compositional changes (association with vulnerability/resilience). Consistently up/down DEGs are colored in red/blue. **g**, Example genes predictive of subtype vulnerability. Scatterplots show average expression in the subtype across individuals against the effect size of the depletion or enrichment in AD as measured by the  $\log_2$  odds-ratio for late-AD, as in the Methods. **h**, Pairwise scatter plots for gene expression at the subtype level for genes in the reelin signaling pathway, showing strong co-expression between RELN, DAB1, VLDLR, and FYN. **i**, Barplot of joint average expression of VLDLR, FYN, and DAB1 in each neuronal subtype, colored by its compositional change ( $\log_2$ OR) in lateAD.



**Figure 4-7. Neuronal DEGs and DEG partitions.** **a**, Consistent DEGs in AD for AD-vulnerable neuronal subtypes. **b**, Functional enrichments for consistently up-regulated genes in vulnerable subtypes across conditions. **c**, Functional enrichments for neuronal partitions (top 250 genes by effect size in each set). **d**, Boxplots showing the correlation of genes' differential effect sizes across subtypes with subtype depletion. **e**, Functional enrichments on genes binned across the resilience-vulnerability spectrum as defined by lateAD expression correlation with subtype depletion.

#### 4.5 Cell type composition analyses

**Analysis of cell type composition differences between brain regions:** For comparing the relative abundance of major cell types across brain regions, the fraction of cells of a major cell type class was computed relative to all the cells isolated from a region. For the statistical analysis of cell type composition differences between brain regions, we also computed the relative abundance of major cell type classes separately for each study participant. To this end, the fraction of cells of a major cell type class was computed relative to all the cells isolated from a brain region of an individual. At the level of high resolution cell types or cell states two distinct measures of relative abundance were computed. First, the relative abundance of each subtype of a major cell class was computed as the proportion of a subtype relative to all cells of the corresponding major cell class isolated from a brain region. Second, for the statistical analysis of differences between brain regions, the fraction of cells of a subtype was computed relative to all the cells isolated from a brain region of an individual. The statistical analyses comparing the



relative abundance of major cell types and subtypes between brain regions was performed using Prism 9 software.

**Analysis of cell type composition alterations in AD:** For each high resolution cell type and brain region the overrepresentation of cells isolated from individuals annotated with the corresponding trait of interest was computed using the hypergeometric cumulative distribution function in MATLAB (MathWorks). More specifically, the extreme upper tail probabilities were computed using the hypergeometric cumulative distribution function `hygecdf(x, M, K, N, 'upper')`. We calculated compositional differences in AD vs. control individuals by modeling the number of cells of a certain cell type or subtype in a specific sample (individual by region) relative to the total number of cells using a quasi-binomial regression model. We modeled AD status by binary ascertainment variables (cogdx 4-5, NIA-Reagan score 1-2, Braak stage 5-6 vs others, as well as any detected presence of nft, neuritic plaque, or diffuse plaque in the region) while adjusting for brain region and sex. We used the `emmeans` package in R to assess significance of the regression contrasts and used `p.adjust` with the `fdr` method to adjust p-values. We modeled the effects of fraction of cells on cognitive performance in multiple domains with gaussian linear regression of cognitive performance on last visit versus  $\log_{10}$  fraction of cells in the subtype or major cell type jointly with covariates for age, sex, Apoe-e4, and post-mortem interval, with FDR p-value correction (`p.adjust` in R).

**Cross-region fraction comparisons and markers of neuronal vulnerability:** We compared the fractional abundances of pairs of neuronal subtypes between two regions using Kendall's tau in AD individuals only (NIA-Reagan score). Significance was assessed using beta regression (R library `betareg`) controlling for sex, APOE genotype, post-mortem interval, and age of death, and we adjusted p-values using `p.adjust` in R with the `fdr` method. We identified markers associated with excitatory neuron subtype vulnerability by performing linear regression to predict the  $\log_2$  odds ratio of each subtype's depletion in late AD based on its expression at the subtype-aggregate level ( $\log(X+1)$ , averaged normalized expression in each subtype by region by individual batch), controlling for age, sex, and post-mortem interval and adjusting p-values with `p.adjust` (`fdr`).

#### 4.6 Differential gene expression

**Differentially expressed genes in Alzheimer's disease:** We performed differential expression analyses with three separate methods: MAST, Nebula, and Wilcoxon testing<sup>224,225</sup>. For all methods, we subset the tested genes to only genes present in over 20% of cells. For MAST and Nebula, we calculated and included in the regression the top 10 components of unwanted variation calculated using RUV on the pseudo-bulk level data (individuals by regions). For these methods, we also included as covariates the individual's sex, age of death, and post-mortem interval, each cell's counts-per-gene and number of captured genes, and where applicable, the high-resolution cell subtypes and the brain region. For Nebula we used a poisson mixed-model on the count data with an offset of the  $\log_{10}$  total counts per cell. For MAST and wilcoxon, we normalized each cell to a total library size of 10,000 counts. We ran wilcoxon tests on both the cell and individual levels. We adjusted p-values for multiple testing in all cases by using the `p.adjust` function in R with the `fdr` method. For our final set of differential genes in each analysis, we took all genes that were significant (adjusted p-value > 0.05) and concordant in both the MAST and Nebula results. We also provide the results for wilcoxon tests, but did not use these to determine concordant results as they do not correct for many covariates. We computed differential expression against five AD ascertainment variables: continuous measurements of nft, `plaq_n`, and `plaq_d` measured in each region except the thalamus (excluded from these analyses) and binary cognitive impairment (cogdx control=1,2,3 vs. AD=4,5) and NIA-Reagan score classifications (control=3,4 vs. AD=1,2). We provide differential expression results for

each of the 14 major cell types (with T cells, CNS macrophages, and each vascular subtype separately) for all regions jointly and for each region separately. We also provide results for each of the excitatory subtypes either in its most prevalent region for EC, HC, or TH subclasses, or across the neocortex for neocortical subtypes (**Supplementary Table 7**).

**DEG module enrichments:** To assess the enrichment of up/down/non-DE genes in each module, we first assigned each tested DE gene to its closest module by correlation to the module's average expression profile. We then performed a hypergeometric enrichment test for the number of genes in a category (up/down/not DE) assigned to the module, against the total number of tested genes assigned to the module, the total number in the category and the total number tested, and corrected p-values using  $p.adjust$  (BH).

**Neuronal DEG partitions:** In order to partition neuronal DEGs into resilience and vulnerability or degeneration-associated subclasses, we calculated each genes' average expressions and DE effect sizes at the subtype-level and compared these to the relative depletion of the subtypes. For each gene that was DE in at least 25% of all neuronal subtypes, we calculated the correlation of its average subtype expression in late AD with each subtype's compositional stability ( $\log_2$  OR in late AD) across excitatory subtypes, separating protective or resilience-associated genes (correlation  $> 0.2$ ) from detrimental or vulnerability-associated genes (correlation  $< -0.2$ ). We calculated functional enrichments on neuronal DEG partitions using the top 250 genes ordered by effect size in each category. We further separated DEGs with higher effect sizes in vulnerable subtypes from those with similar effect sizes across all neuronal subtypes by calculating the correlation of their differential effect sizes in each subtype with that subtype's depletion ( $\log_2$  OR in late AD). To perform enrichments along the continuum of resilient-vulnerable genes, we kept only genes with biased effect sizes (effect size correlation  $< -0.2$ ) and binned them along the axis of expression correlation (window size 0.2 for at 0.05 intervals) and performed functional enrichments for all bins jointly.

**DEG and module pathway enrichments:** We performed DEG enrichments for each differential expression run using the `gprofiler2` package in R, with multi-query for up and down-regulated genes, as unordered queries, a p-value cutoff of 0.05, and using GO, REAC, WP, KEGG, and CORUM as annotation sources, and kept enriched terms with fewer than 500 genes. Module enrichments were performed in the same manner, using the core genes identified for each module.

#### 4.7 Discussion

In this work, we presented a comprehensive atlas of the aged human brain across brain regions impacted along Alzheimer's disease (AD) progression, comprising over 1.3 million single nuclei in six brain regions from each of 48 individuals along a range of AD progression. Our analysis of multiple brain regions from a large cohort of individuals enabled us to annotate cellular regional diversity, identify a novel neuronal subtype and region-specific glial subtypes, develop a method to detect functional expression modules in atlas-scale snRNA, track regional transcriptomic changes across early-to-late AD, pinpoint vulnerable region-specific excitatory neuronal cell types depleted along AD progression, and propose molecular and connectome pathways mediating excitatory neuron vulnerability and resilience to AD, gaining novel insight into fundamental aspects of the human brain and how it is affected in AD.

In AD, we identified significant transcriptomic changes across cell types, and interpreted them through the lens of our functional modules, identifying metabolic shifts in glia, up-regulation of specific plaque- and tangle-mediating pathways, and coordinated changes in oligodendrocyte lineage cells in conjunction with a depletion of OPCs in AD. AD risk genes were highly perturbed in AD, with the largest number of changes in microglia in line with their enrichment for

GWAS signal, but a significant proportion were perturbed across multiple cell types, including amyloid-interacting CLU<sup>237</sup>, and some showed region-specific expression, including thalamus-enriched lipid transporter ABCA7<sup>226</sup>.

We tracked specific glial and neuronal compositional alterations in the course of AD and identified multiple AD-vulnerable excitatory neuron subtypes in the early-affected entorhinal cortex and hippocampus<sup>21</sup>, and in the neocortex in late AD only, that were significantly depleted in brains from AD patients and whose depletion predicted cognitive decline, including documented entorhinal cortex layer II (L2) and hippocampal CA1 subfield neurons<sup>41-45</sup>. We find that selectively vulnerable neurons in the entorhinal cortex are composed of multiple transcriptionally distinct excitatory subtypes, including two RELN-positive subtypes<sup>38,39</sup>, RORB-positive L5 (AGBL1/GPC5)<sup>40</sup>, and TOX3 and TTC6-expressing L2/3 neurons. We provide evidence that neurons predicted to be part of the same circuit connecting the hippocampus and the entorhinal cortex are co-depleted and express shared markers of vulnerability, including genes involved in reelin-signaling pathway and in heparan sulfate proteoglycan biosynthesis. Vulnerable subtypes are altered similarly to other excitatory neurons in the AD brain but differences highlight pathways that may enable neuron resilience in AD including aerobic glycolysis and clathrin mediated endocytosis.

Our work provides a cross-cutting view of brain regional heterogeneity with particular applications in Alzheimer's disease across regions and disease timepoints, enabling many future studies: of neuronal development, specialization, and specification; application of our modules framework across scRNA-seq analyses and methodological development for other single-cell modalities; targeted experimental follow-up for molecular changes in Alzheimer's in each cell type; comprehensive joint investigation of Alzheimer's-affected circuitry and of determinants of neuronal vulnerability; and investigation of mechanisms of neuronal vulnerability and resilience, with applications in potential therapeutics, furthering our understanding of brain development, architecture, and changes in Alzheimer's disease progression.

## 5. Single-cell mosaicism analysis reveals cell-type-specific somatic mutational burden in Alzheimer's Dementia

Parts of this chapter have been presented/published as:

### Single-cell mosaicism analysis reveals cell-type-specific somatic mutational burden in Alzheimer's Dementia

*(in review)*

Maria Kousi<sup>1,2,\*</sup>, Carles Boix<sup>1-3,\*</sup>, Yongjin P. Park<sup>1,2</sup>, Hansruedi Mathys<sup>2,4,5</sup>, Samuel Sledzieski<sup>1,2</sup>, Zhuyu Peng<sup>4,5</sup>, David A. Bennett<sup>6</sup>, Li-Huei Tsai<sup>2,4,5</sup>, Manolis Kellis<sup>1,2,6,+</sup>

<sup>1</sup> Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA, USA. <sup>2</sup> Broad Institute of MIT and Harvard, Cambridge, MA, USA. <sup>3</sup> Computational and Systems Biology Program, Massachusetts Institute of Technology, Cambridge, MA, USA. <sup>4</sup> Picower Institute for Learning and Memory, Cambridge, MA, USA. <sup>5</sup> Department of Brain and Cognitive Sciences, Cambridge, MA, USA. <sup>6</sup> Rush Alzheimer's Disease Center, Rush University Medical Center, Chicago, IL, USA. \* Equal contributors. + Correspondence to: manoli@mit.edu



## 5.1 Introduction

Despite significant advances in identifying genetic drivers of neurodegenerative disorders, the majority of affected individuals lack molecular genetic diagnosis, with somatic mutations proposed as one potential contributor to increased risk. Here, we report the first cell-type-specific map of somatic mosaicism in Alzheimer's Dementia (AD), using 4,014 cells from prefrontal cortex samples of 19 AD and 17 non-AD individuals. We integrate full-transcript single-nucleus RNA-seq (SMART-Seq) with matched individual-level whole-genome sequencing to jointly infer mutational events and the cell-type in which they occurred. AD individuals show increased mutational burden, localized in excitatory neurons, oligodendrocytes, astrocytes and disease-associated "senescent" cells. High-mutational-burden cells showed mutational enrichment and similar single-cell expression profiles in AD cases versus non-AD individuals, indicating cellular-level genotype-to-phenotype correlation. Somatic mutations are specifically enriched for known AD genes, and implicate biologically meaningful cell-type specific processes, including: neuronal energy regulation, endocytic trafficking (*NEFM*), lipid metabolism (*CNP*, *CRYAB*), proteostasis (*USP34*), cytoskeleton, and microtubule dynamics (*MACF1*).

Late-onset neurodegenerative disorders, including Alzheimer's Dementia (AD), Parkinson's disease, Huntington's disease, and Amyotrophic Lateral Sclerosis are primarily sporadic, arising through the combined effect of multiple environmental risk factors and many weak-effect common alleles<sup>66,67</sup>, though heritable forms have also been described, stemming from a small number of strong-effect rare alleles<sup>68-71</sup>. Hindering the search for effective therapeutic interventions, our current knowledge of genetic contributors is still greatly incomplete, in part due to their complex etiology, heterogeneous level of affectedness, diverse symptoms, and the interplay between genetic and environmental factors<sup>72,73</sup>.

Several genetic contributors to AD have been elucidated, including: the *E4* haplotype of the *ApoE* gene locus<sup>74</sup>; at least 25 loci from genome-wide association studies (GWAS)<sup>66,67</sup>, including *BIN1*<sup>75</sup>, *PICALM*<sup>76</sup>, *CLU*<sup>76,77</sup>, *CR1*<sup>77</sup>; a small number of genes discovered from family pedigrees, involving *APP*, *PSEN1*, *PSEN2*<sup>68-71</sup>; and an increasingly-large set of coding and non-coding loci from ongoing exome-sequencing and whole-genome sequencing (WGS) studies<sup>78,79</sup>. With the exception of the ApoE locus, genetic results in AD follow the commonly-observed inverse relationship between allele frequency and effect size, in the continuum between high-frequency weak-effect variants (from GWAS) and low-frequency strong-effect variants (from family studies), due to the impact of selection, maintaining strong-effect variants at low frequencies. Despite these major advances however, the vast majority of AD cases lack any detectable genetic basis.

Beyond inherited genetic variants, genetic contributors to disease can include somatic mutations incurred either during early embryonic development<sup>81</sup> (thus affecting a large fraction of cells and tissues), in later stages of embryogenesis (thus confined to a specific tissue or cell type), or during childhood, adulthood, and late-life (thus affecting only a small number of cells). These mutations can stem from environmental exposures, disease-associated DNA damage, disruption of DNA repair processes, or through the natural process of error-prone mitotic cell divisions. The phenomenon describing such somatic mutations is referred to as mosaicism, given that the affected individual is a mosaic of cells of different genetic makeup, including both risk and protective alleles of disease loci. Unlike inherited mutations, late-stage somatic mutations are not subject to the strong selective pressures acting on germline variation, and thus can result in strong gene-perturbing effects, which can have a disproportionate impact on cellular phenotype. Increased somatic mutational burden in individual cells can increase the chance of single-cell dysregulation. Increasing numbers of dysregulated cells can subsequently manifest as tissue- and organ-level disruptions, providing an additional path for disease onset

and progression, beyond inherited genetic components and environmental contributors.

Though significant progress has been made in studying the relevance of mosaicism in disease, with prominent examples from cancer and clonal hematopoietic malignancies<sup>82–84</sup>, congenital heart disorders<sup>85</sup>, and skin disorders<sup>86</sup>, the inaccessibility of brain tissue has hindered the study of mosaicism in neurodevelopmental and neurodegenerative disease<sup>87,88</sup>. Bulk DNA sequencing has been used to recognize high-frequency somatic mutations in autism spectrum disorders<sup>89–91</sup> and in hemimegalencephaly<sup>92</sup>, from their partial representation in WGS reads. For DNA-damage-repair disorders, a single-cell DNA sequencing study showed increased mutations in nine cases vs. controls<sup>93</sup>, profiling three to six neurons from each affected individual. Thus, it remains unclear whether somatic mosaicism plays a role in neurodegeneration more broadly (outside DNA-damage-repair disorders), whether increased mutations are also found in non-neuronal cells (which are increasingly recognized to play important roles in AD<sup>37</sup>), and how DNA perturbations relate to gene expression changes at the single-cell level, and ultimately to phenotypic effects in brain function and cognition.

Here, we report the first single-cell joint map of somatic mutation and gene expression patterns in AD, using 4,014 cells from prefrontal cortex samples of 36 human *post mortem* brains with matched WGS data, including 19 AD and 17 non-AD individuals based on cognitive clinical diagnosis. We profile full-transcript single-nucleus RNA sequencing with SMART-Seq, and develop a new method to jointly infer somatic mutations and single-cell gene expression levels, enabling us to map mosaic mutations onto specific cell types and cellular states in which they occur. We find that AD individuals show higher overall mutational burden, higher damaging-variant burden for missense and stop-codon mutations, and higher cell-type-specific burden for excitatory neurons, oligodendrocytes, astrocytes, and a newly-described class of “senescent” cells. Somatic mutations are enriched for known AD genes, and for specific novel and biologically-meaningful genes and pathways, including: neuronal energy regulation, endocytic trafficking, lipid metabolism, proteostasis, cytoskeleton, and microtubule dynamics. Overall, our study uncovers the role of cell-type specific mutational burden in AD, reveals candidate new AD target genes and pathways, and provides a general framework for understanding somatic variation and its phenotypic associations at the single-cell level in human disease.

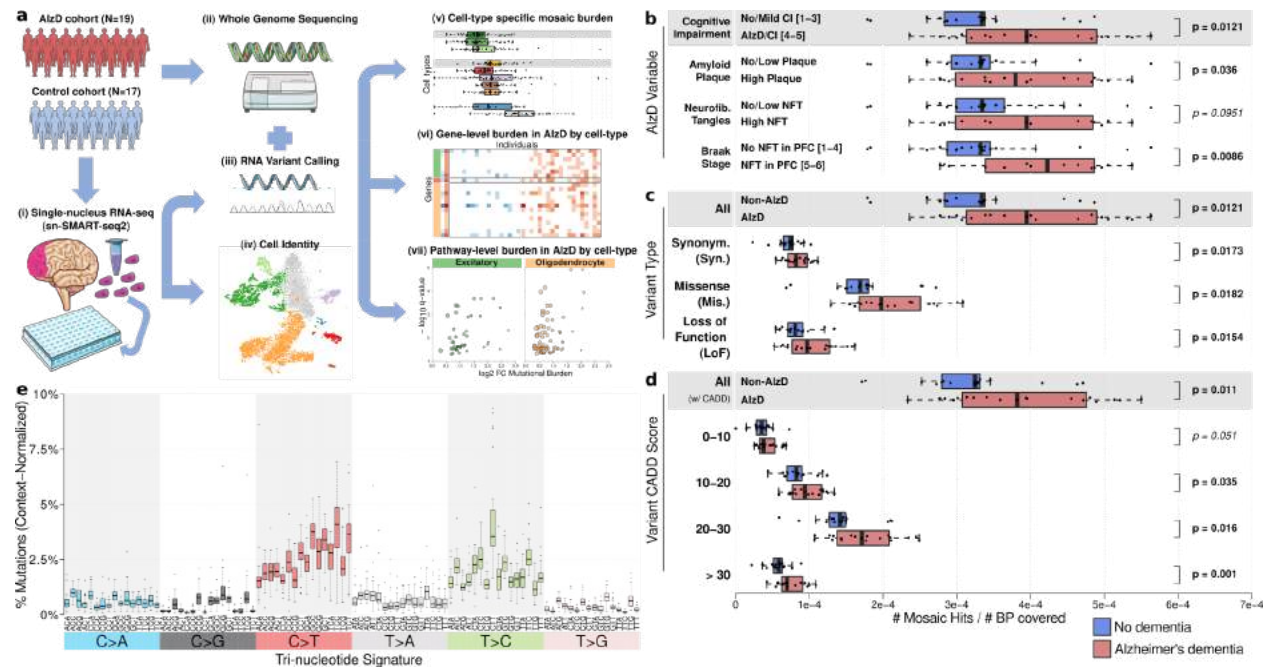
## 5.2 Cohort assembly and variant detection

We obtained *post-mortem* human prefrontal cortex brain samples from 47 Religious Orders Study and Rush Memory and Aging Project (ROSMAP)<sup>238</sup> individuals (**Fig. 5-1a**), classified into 22 AD and 25 non-AD individuals, of which 36 also have whole-genome-sequencing (30x coverage)<sup>239</sup>.

We used well-based single-nucleus full-transcript profiles (Smart-Seq2 v2.0<sup>240</sup>) across all 47 individuals, and obtained 6,180 cells (5,421 after QC). We used known marker gene expression to annotate 2,121 oligodendrocytes (39.1%), 1,170 excitatory neurons (21.6%), 255 microglia (4.7%), 221 inhibitory neurons (4.1%), 220 astrocytes (4.1%), 94 OPCs (1.7%), and 40 endothelial cells (0.7%). We also identified a cluster of 1,300 cells (24%) that we refer to here as “senescent” with reduced transcription, loss of cell-type-identity markers, and expression of a subset of senescence-associated markers (AXL, SENP7), and positioned between neurons and oligodendrocytes in the two-dimensional t-SNE cell embedding. However, these did not express all classical senescence markers (e.g. CDKN1A/2A), suggesting a possibly distinct neuronal and oligodendrocyte pre-apoptotic state<sup>241</sup>. We only analyzed somatic mutations for the subset of 36 individuals with whole-genome sequence (encompassing 4,040 cells: 935 excitatory neurons, 167 inhibitory neurons, 159 astrocytes, 188 microglia, 74 oligodendrocyte progenitor cells (OPC), 1382 oligodendrocytes, 32 endothelial cells, 1103 senescent cells; see Methods).

For the subset of 36 individuals with whole-genome sequencing data (4,014 cells), we inferred somatic mutations in each of their cells (GATK HaplotypeCaller, using germline-depth-adjusted thresholds, **Supplementary Fig. 5-S1**), merged all mutations, excluded known AD SNPs and germline variants, and kept only exonic mutations present in >20% of reads for that cell (**see Methods**), resulting in an initial set of 66,664 mutational calls across all cells and all individuals. After excluding blacklisted regions/genes, RNA editing sites, mutations shared across >1 individuals, ExAC variants with MAF>1%, hyper-mutated cells, and merging calls within 1-3bp of another, resulting in 55,447 mutational calls (see Methods).

These mutations show high deleteriousness, indicative of relaxed selection, and consistent with their somatic occurrence in only one tissue and a subset of cells. The majority (51.3%) were amino-acid changing (missense), 24.1% loss-of-function (nonsense/stop, frameshifting, splice-disrupting), 21.5% amino-acid preserving (synonymous), and 3.1% other (inter/intragenic) (**Supplementary Fig. 5-S2a**), with consistent frequencies across individuals and sequencing depths (**Supplementary Fig. 5-S2b,c**). They also showed high mutation deleteriousness (CADD)<sup>242</sup> scores (average=22, top 0.63% of genome), and their nonsynonymous to synonymous rate<sup>243</sup> deviates significantly from the baseline neutral selection ( $dN/dS=2.5$ ) (**Supplementary Fig. 5-S2e**).



**Figure 5-1. Somatic mutation burden in Alzheimer's Dementia (AD).** **a. Overview.** Across 36 post-mortem prefrontal cortex samples from 19 AD and 17 non-AD individuals (i), we profile full-transcript single-nucleus RNA-seq (SMART-seq2) (ii) and whole-genome sequencing (iii) to jointly infer mosaic mutations and cell identity (iv) across 4,014 cells, which we use to infer cell-type-specific mosaic burden (v), gene-level burden (vi) and pathway-level burden (vii). **b-d.** Increased mosaic mutational burden (x-axis) between AD (red) and non-AD (blue) individuals for different ascertainment variables (**b**), types of mutation (**c**), and predicted variant effect (**d**, CADD score). **e.** Normalized frequency (y-axis) of somatic mutations for different single-nucleotide variant tri-nucleotide contexts (columns) across individuals (points) highlights aging-related signatures (C-to-T followed by T-to-C).

### 5.3 Cohort and single-nucleus SMART-seq2 methods

**Cohort and metadata.** We assembled a cohort of 47 individuals, 22 with Alzheimer's dementia and 25 age and sex matched persons with no-dementia from the ROSMAP cohorts<sup>37,183</sup>. The cases had a clinical diagnosis of Alzheimer's dementia, mild cognitive impairment, or no cognitive impairment based on a detailed clinical evaluation as previously described<sup>30,244</sup>. For mosaic mutational burden analysis we additionally split the cohort by overall histopathology reported levels of  $\beta$ -amyloid load, PHFtau tangle density, and Braak stage, where stages 1-4 indicate the absence of NFTs in prefrontal cortex (PFC) and stages 5-6 indicate the presence of NFTs in the PFC (**Supplementary Fig. 5-S3A-G** and **Supplementary Table S5**)<sup>19,31</sup>.

#### **Isolation of nuclei from frozen postmortem brain tissue for SMART-Seq2 plate-based single-nucleus RNA-sequencing.**

We adapted the protocol for the isolation of nuclei from frozen postmortem brain tissue from Swiech et al.<sup>245</sup>. All procedures were carried out on ice or at 4 °C. Briefly, postmortem brain tissue was homogenized in 2 ml Homogenization Buffer (320 mM Sucrose, 5 mM CaCl<sub>2</sub>, 3 mM Mg(Ac)<sub>2</sub>, 10 mM Tris-HCl pH 7.8, 0.1 mM EDTA pH 8.0, 0.1% IGEPAL CA-630, 1 mM  $\beta$ -mercaptoethanol, 0.4 U/microliter Recombinant RNase Inhibitor (Clontech)) using a Wheaton Dounce Tissue Grinder (10 strokes with the loose pestle). 3 ml of Homogenization Buffer was added (final volume 5 ml) and the homogenized tissue was incubated on ice for 5 minutes. Then the homogenized tissue was filtered through a 40  $\mu$ m cell strainer, mixed with an equal volume of Working Solution (83% OptiPrep™ Density Gradient Medium (Sigma-Aldrich), 5 mM CaCl<sub>2</sub>, 3 mM Mg(Ac)<sub>2</sub>, 10 mM Tris HCl pH 7.8, 0.1 mM EDTA pH 8.0, 1 mM  $\beta$ -mercaptoethanol) and loaded on top of an OptiPrep density gradient (10 ml 29% OptiPrep solution (29% OptiPrep™ Density Gradient Medium, 134 mM Sucrose, 5 mM CaCl<sub>2</sub>, 3 mM Mg(Ac)<sub>2</sub>, 10 mM Tris-HCl pH 7.8, 0.1 mM EDTA pH 8.0, 1 mM  $\beta$ -mercaptoethanol, 0.04 % IGEPAL CA-630, 0.17 U/microliter Recombinant RNase Inhibitor) on top of 5 ml 35 % OptiPrep solution (35% OptiPrep™ Density Gradient Medium, 96 mM Sucrose, 5 mM CaCl<sub>2</sub>, 3 mM Mg(Ac)<sub>2</sub>, 10 mM Tris-HCl pH 7.8, 0.1 mM EDTA pH 8.0, 1 mM  $\beta$ -mercaptoethanol, 0.03 % IGEPAL CA-630, 0.12 U/microliter Recombinant RNase Inhibitor)). The nuclei were separated by ultracentrifugation using an SW32 rotor (20 minutes, 9000 rpm, 4 °C). 3 ml of nuclei were collected from the 29%/35% interphase and washed with 15 ml ice cold PBS containing 0.5% BSA and 2 mM EDTA. The nuclei were centrifuged at 100 g for 5 minutes (4 °C) and resuspended in 1 ml PBS containing 0.5% BSA and 2 mM EDTA. Then the nuclei were stained by adding two drops of NucBlue Live ReadyProbes Reagent (ThermoFisher Scientific, catalog number: R37605) and passed through a 40  $\mu$ m cell strainer (Falcon Cell Strainers, Sterile, Corning, product 352340). Single Hoechst-positive nuclei were sorted into 96-well plates (Eppendorf twin.tec PCR Plate 96, catalogue number 951020401) containing 5 microliters of Buffer TCL (QIAGEN, catalogue number 1031576) per well containing 1% beta-mercaptoethanol (Sigma-Aldrich, catalogue number M6250), snap frozen on dry ice, and then stored at -80 °C before whole transcriptome amplification, library preparation and sequencing.

**Reverse transcription, whole-transcriptome amplification, library construction, sequencing.** Single-nucleus RNA sequencing libraries were generated based on the SMART-Seq2 protocol<sup>246</sup> with the following modifications. RNA from single cells was first purified with Agencourt RNAClean XP beads (Beckman Coulter, catalogue number A63987) before oligo-dT primed reverse transcription with Maxima H Minus Reverse Transcriptase (Thermo Fisher, catalogue number EP0753), which was followed by 21 cycle PCR amplification using KAPA HiFi HotStart ReadyMix (KAPA Biosystems, catalogue number KK2602) with subsequent Agencourt AMPure XP bead (Beckman Coulter, catalogue number A63881) purification. Libraries were tagged using the Nextera XT DNA Library Preparation Kit (Illumina, catalogue number FC-131-1096) and the Nextera XT Index Kit v2 Sets A,B,C, and D according to the



manufacturer's instructions with minor modifications. Specifically, reactions were run at one fourth the recommended volume, the tagmentation step was extended to 10 minutes, and the extension time during the PCR step was increased from 30 seconds to 60 seconds. Libraries from 192 cells with unique barcodes were combined and sequenced on the Illumina HiSeq 2000 platform at the MIT BioMicro Center.

**Read processing.** We aligned, uniformly-processed, and called mutations on each of 6,180 single nuclei from 47 individuals (average of 131.5 per individual) following GATK best practices pipeline for RNA-seq<sup>247</sup>. We first aligned 40-bp single-end reads from each sequenced single nucleus separately with the STAR aligner<sup>248</sup> against the b37 genome with decoy contigs using a two pass alignment (options: `--outFilterMultimapNmax 20 --alignSJoverhangMin 8 --alignSJDBoverhangMin 1 --alignIntronMin 20 --alignIntronMax 1000000`). We then used Picard tools<sup>155</sup> to revert and merge the alignment with unaligned reads and marked duplicates on the merged bam file. We identified and removed alignments on decoy contigs, sorted and fixed NM, MD, and UQ tags with Picard tools, filtered duplicates, unmapped, and non-primary alignment reads, and split reads by Ns in their CIGAR string. Finally, to generate a bam file for transcriptomics and mutation calling we ran two passes of GATK<sup>247</sup> BaseRecalibrator for each cell against the b37 genome with decoy contigs, using dbSNP<sup>249</sup> v138 and the Mills 1000G indels<sup>250</sup>.

#### 5.4 Enrichment for somatic changes in AD patients

Summing across all cell types, we found that somatic changes are enriched in AD vs. non-AD brains, when ascertained using cognitive impairment (Cogdx 4-5 vs. 1-3,  $p=0.0121$ ; **Fig. 5-1b**),  $\beta$ -amyloid load ( $p=0.036$ ), and Braak stage<sup>251</sup> which measures neurofibrillary tangle (NFT) abundance across cortical areas (stage 5-6 marking NFT diffusion across neocortex vs. stage 1-4 marking NFTs restricted to entorhinal cortex and hippocampus,  $p=0.0086$ ) (**Fig. 5-1b**). However, overall PHFtau NFT density<sup>252</sup> did not show a significant difference in mutational burden ( $p=0.0951$ ), highlighting the importance of clinical diagnosis and precise localization measures (see **Supplementary Fig. 5-S3**).

We found higher somatic burden in AD vs. non-AD across all types of mutations, including synonymous (+18.1%,  $p=0.0173$ ), missense (+22.6%,  $p=0.0182$ ), and LoF (+25.7%,  $p=0.0154$ ) (**Fig. 5-1c**), with the enrichment in AD being greater with stronger deleteriousness scores, from +16% for CADD 10-20 ( $p=0.035$ ) to +22.3% for CADD 20-30 ( $p=0.016$ ), to +32% for CADD>30 ( $p=0.001$ ) (**Fig. 5-1d**). Across both AD and non-AD, changes were enriched for mutational signatures of the aging brain (C-to-T and T-to-C transitions), suggesting that burden was driven by differences in mutational abundance, rather than distinct mutational processes between AD and non-AD (**Fig. 5-1e, Supplementary Fig. 5-S4a-c**)<sup>93,253</sup>

#### 5.5 Methodology for mutation calling

**WGS.** For 37 of the 47 individuals, we obtained germline variants from ROSMAP (at <https://www.synapse.org/#!Synapse:syn10901595>)<sup>183,239</sup>. Variants were profiled using whole genome sequencing (WGS) at 30x coverage with 150bp paired-end reads, aligned on the GRCh37 genome, and processed using the GATK best practices pipeline, including marking duplicates, indel realignment, and base quality recalibration, and called using the HaplotypeCaller and GenotypeGVCFs<sup>247</sup>. To create a germline list to filter against our calls, we conservatively retained all variants with a depth of at least 3 and a non-reference allele fraction of at least 20% using vcfutils. Of the 37 WGS datasets, 17 were performed on the prefrontal cortex, 11 on another brain region, and 9 on blood.

**snRNA-seq.** In order to call potential mosaic hits, we ran GATK's HaplotypeCaller on each cell independently, merged all vcfs using GenomicsDBImport on each chromosome, ran

GenotypeGVCFs on the merged cells in each individual using dbSNP v138 as a database of known SNPs, and concatenated vcfs with GatherVcfs<sup>247</sup>. Using GATK's VariantFiltration, we filtered out any variants corresponding to the individual's germline calls or in low-complexity regions<sup>254</sup>. We compared known germline variants to potential mosaic variants in the sn-RNA-seq data to establish additional variant filters and kept variants with QD > 15 at DP = 3, QD > 18 at DP= 4, QD > 20 at DP=4, and with FS < 10, SOR > .5, MQ < 61, or DP < 3 (**Supplementary Fig. 5-S1A-C**). We kept any variants with over 20% allelic balance that also overlapped exons of protein coding genes in the GENCODE annotation, resulting in 66,664 variants across all cells (keeping 24.3% of 273,508 preliminary calls). We further filtered variants by flagging: **(a)** sites in blacklisted regions (chrY, chrM, decoy contigs, and scaffolds), or **(b)** in blacklisted genes (HLA locus, IgG genes, and pseudogenes); **(c)** RNA editing sites from RADAR v2, DARNED, or REDportal<sup>255-257</sup>; **(d)** sites that appeared in more than 1 individual and potentially represented alignment artifacts; **(e)** sites in cells filtered out in the final transcriptome analysis; **(f)** calls with only one alternate supporting read; **(g)** 326 calls matching variants in ExAC with MAF > 1%<sup>258</sup> (**Supplementary Fig. 5-S1D**); **(h)** sites from 31 hyper-mutated flagged cells; and **(i)** merging calls found within 1-3bp of each other in the same samples into single indels (reduced 3,105 to 1,459 calls). This resulted in 56,608 (84.9% from 66,664 and 20.7% of 273k) final, filtered variants in 37 individuals with matched WGS, with mutations in 3,044 of 5,821 kept cells) for downstream analysis (**Supplementary Fig. 5-S1A-E**). For all mosaic burden analyses, we kept 55,447 variants (97.9%) from the 36 individuals with WGS and cognitive impairment scores (cogdx variable) between 1 and 5.

**Variant annotation.** We annotated variants with SnpEff v4.3<sup>259</sup> using the GENCODE v28lif37 b37 genome annotation<sup>260</sup>. We grouped variants into synonymous (12,166, 21.5% of 56,608 variants from the 37 individuals), missense (29,071, 51.3%), loss of function (LoF, 13,623, 24.1%), and other variants (1,748, 3.1%). LoF variants included frameshift, start or stop variants, indels, and splice variants (**Supplementary Fig. 5-S2A-C**). We used ExAC to annotate the MAF of 3,636 rare matching variants (6.42% of filtered calls), with 3,283 (90.3%) at less than 0.01% and 1,880 (51.7%) at less than 0.001%<sup>258</sup> (**Supplementary Fig. 5-S1D**). Finally, we annotated 53,789 (95%) variants with CADD scores<sup>261</sup>, of which 5,861 (10.9%) have scores up to 10, 13,484 (25.1%) between 10 and 20, 24,007 (44.6%) between 20 and 30, and 10,437 (19.4%) above 30 (**Supplementary Fig. 5-S2D**).

**Mutational signatures.** We computed the observed trinucleotide signatures and normalized signatures within each sample by enumerating the possible DNA tri-nucleotide contexts in genomic locations covered by snRNA-seq at a callable depth for each cell. We estimate and report the average % for each type of single-nucleotide variant (SNV) for all 36 individuals as **Figure 1E** and split by case/control status according to cognitive diagnosis in **Supplementary Fig. 5-S5A, S5B** (clipped at 10%). Mutational signature decomposition was performed relative to the COSMIC<sup>253</sup> mutational signatures using deconstructSigs<sup>262</sup> in R (**Supplementary Fig. 5-S5C**).

## 5.6 Mutational Burden calculations

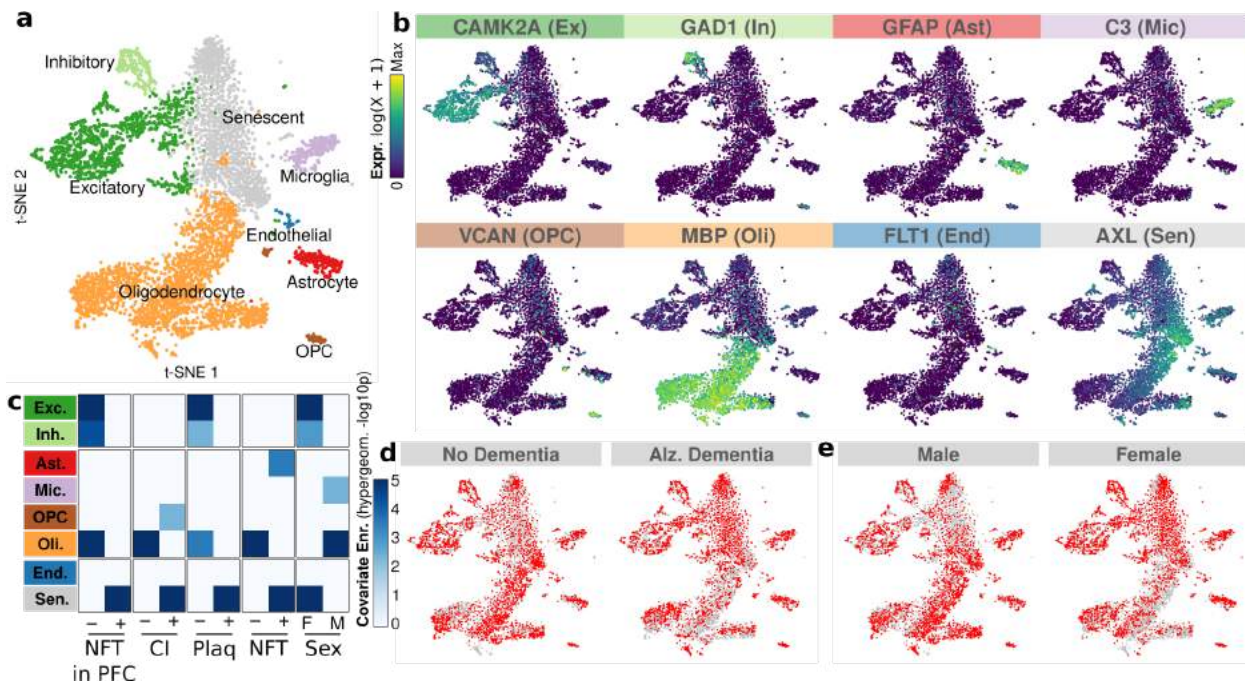
**Individual-level burden.** We performed all mosaic mutational burden analyses on the 36 of the 37 individuals with WGS and mosaic variant calls. For burden analyses, we removed one individual of the 37 who had a cognitive score of 6 (other dementia). The 36 individuals contained 19 cases and 17 controls by cognitive diagnosis (scores of 1-3 vs. 4-5). Amyloid level split the burden cohort 20 to 16 (amyloid > 1), tangles split 15 to 21 (tangles > 5), and the presence of tangles in PFC split 14 to 22 (braaksc 1-4 vs. 5-6). For each individual, we estimated the average mosaic mutational burden as the number of mosaic calls divided by the number bases covered at a callable depth in the snRNA-seq data. We compared the distribution

mosaic burden in cases to controls using a one-sided t-test. For variant type and CADD divisions, we computed burden estimates per individual per variant class and compared cases and controls, as defined by cognitive diagnosis, within synonymous, missense, and LoF variants, as well as in each CADD score bin.

**Cell-level burden.** For the cell-level analysis, we estimated the mutational burden for each individual in each cell type and in all neurons (label “Neuronal”: excitatory and inhibitory neurons) and all glia (label “Glial”: aggregating astrocytes, microglia, OPCs, and oligodendrocytes). We used one-sided t-tests to determine significance both between major cell groups (neuronal to glial and to senescent and glial to senescent), and within cell-types between cases and controls as defined by cognitive diagnosis. Additionally, we evaluated the effect of sex (male/female) on mosaic mutational burden both over all cell types in aggregate and at the cell-level by two-sided t-tests (**Supplementary Fig. 5-S8**). To evaluate the effect of age-of-death in a cell-type-dependent manner, we fit a linear regression to the per-individual burden averages within either neuronal ( $p=0.608$ ) or glial cells ( $p=0.127$ ). We also evaluated the interaction of age and cell-type by comparing the additive model to a model with the interaction term (comparison by anova,  $p=0.55$ ).

### 5.7 Cell-type identification through transcriptomic analysis

We next analyzed these mutations in the context of their cell-type-identity, leveraging the unique nature of our dataset. We annotated the cell type of our 4,040 cells across 36 whole-genome-sequenced individuals with mutational burden information, in the context of all 5,421 cells from all 47 individuals (**Fig. 5-2a, Supplementary Fig. 5-S5a**). We used known marker genes to annotate seven well-established cell types, including: excitatory neurons (*CAMK2A*, *NRGN*, *SLC17A7*), inhibitory neurons (*GAD1*, *GAD2*), astrocytes (*AQP4*, *GFAP*), microglia (*C3*, *CD74*, *CSF1R*), oligodendrocytes (*MBP*, *MOBP*, *PLP1*), OPCs (*PDGFRA*, *VCAN*), and endothelial cells (*FLT1*, *CLDN5*) (**Fig. 5-2b, Supplementary Fig. 5-S5b-d**), with each cluster showing substantial contributions by multiple individuals (**Supplementary Fig. 5-S6, 5-S7**). We also annotated a transcriptionally-distinct cluster of senescent cells (**Fig. 5-2a**), bearing resemblance to neurons and oligodendrocytes, but lacking clean cell-type identity (**Supplementary Fig. 5-S5c-e, 5-S8**). Senescent cells were enriched in AD individuals (across all ascertainment parameters, **Fig. 5-2c**), and also in female individuals (**Fig. 5-2c,e**). Oligodendrocytes were depleted in cases and in female individuals (**Fig. 5-2c-d**), likely reflecting higher oligodendrocyte vulnerability in AD, and consistent with our recently-reported AD-associated and sex-specific down-regulation of myelination pathways<sup>37</sup>.



**Figure 5-2. Transcriptomic landscape of Alzheimer's dementia.** **a.** Gene expression landscape (t-SNE) and cell type annotation (colors) across 5,421 cells from all 47 individuals. **b.** Representative cell-type-specific marker gene expression (color scale). **c.** Enrichment (heatmap) for each cell type (rows) across phenotypic classes (columns) and sex. **d,e.** Expression landscape (t-SNE) of cells by phenotype (d) and sex (e).

## 5.8 Methodology for transcriptional profiling and clustering

**Cell identities.** For each of 6,180 cells, we used HTSeq-count<sup>263</sup> on its filtered and recalibrated bam file to compute the cell's transcriptomic coverage over each gene's exons in GENCODE gene annotation (v28 lifted to b37). We used SCANPY<sup>194</sup> to process and cluster the expression profiles and infer cell identities. We kept only 19,765 protein coding genes detected in at least 3 cells and filtered out 37 cells with less than 100 expressed genes. We identified and filtered out 322 additional cells showing very strong individual-specific batch effects, leaving 5,821 cells over 47 individuals. We used the filtered dataset to calculate the low dimensional embedding of the cells (t-Stochastic Neighbor Embedding: t-SNE) (default parameters, perplexity=50), built a nearest neighbors graph (n=10), and clustered it with the Louvain clustering (resolution=2), giving 24 preliminary clusters. We then manually assigned clusters based on the following 2-3 major marker genes per class: Neuronal: *GRIN1*, *SNAP25*, *SYT1*; Excitatory neurons: *CAMK2A*, *NRGN*, *SLC17A7*; Inhibitory neurons: *GAD1*, *GAD2*; Astrocytes: *AQP4*, *GFAP*; Microglia: *C3*, *CD74*, *CSF1R*; Oligodendrocytes: *MBP*, *MOBP*, *PLP1*; Oligodendrocyte progenitor cells (OPCs): *PDGFRA*, *VCAN*; Endothelial: *FLT1* and *CLDN5*. We merged clusters sharing marker genes to obtain 9 final clusters, defining two broad neuronal subtypes (1,170 excitatory and 221 inhibitory cells), four glial clusters (220 astrocytes, 255 microglia, 2121 oligodendrocytes, and 94 OPCs), 40 endothelial cells, 400 cells with strong individual-specific batch effect and cancer signatures (*DNMT3A*, *COL6A3*), and 1,300 senescent cells, marked by *CARD8*, *FAM126A*, *IRX2*, *ALK*, *SENP7*, and *GMFB* and lower overall transcription (**Supplementary Figure 5-S8 and Supplementary Table S6**). In burden analyses, where we used the 36 individual cohort, we kept 4,440 cells corresponding to these individuals (Ex: 935, Inh: 167, Ast: 159, Mic: 188, OPC: 74, Oli: 1,382, End: 32, Batch: 400, Sen: 1,103). As a final



step of QC, we removed the 400 cells showing batch effect. In our previous study<sup>37</sup>, we did not observe a senescent cell population, likely due to the protocol's selection bias for barcoded cells with higher numbers of UMIs.

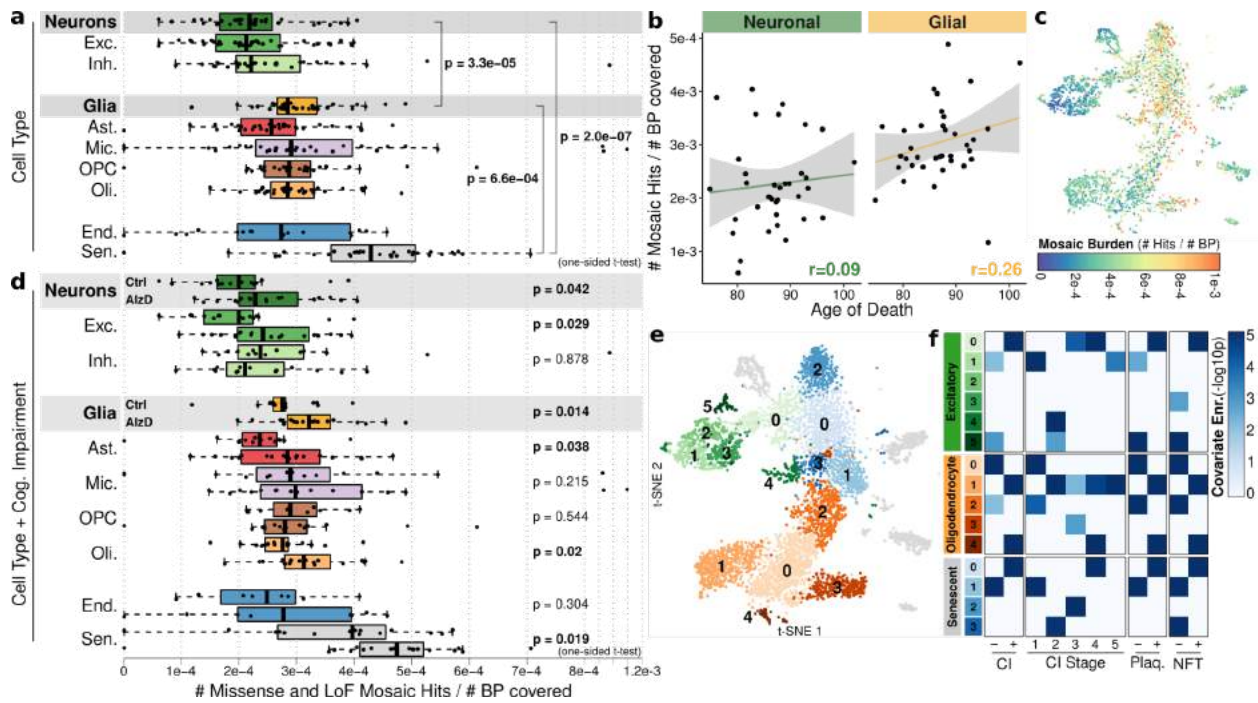
**Sub-clustering.** We further sub-clustered the cell types with the largest number of cells ( $N \geq 1,000$ , excitatory neurons, oligodendrocytes, and senescent clusters) by Louvain clustering (resolution=0.3). To further define the clusters and sub-clusters, we performed marker-gene discovery on the clusters against all cells and on the sub-clusters against their overall cluster with scanpy, removing unannotated and MT-like genes. We show the top 10 marker genes per cluster in **Fig. 5-2F** and per subcluster in **Supplementary Fig. 5-S10A-C**.

**Covariate enrichment.** We calculated the enrichment of each level of each of the major covariates (sex, cognitive impairment level, cognitive impairment score, NFT level, amyloid plaque level, braak stage, and presence/absence of NFTs in the pre-frontal cortex) in each cluster and sub-cluster by the hypergeometric test, and show enrichments with  $-\log_{10}p > 2$  (**Fig. 5-2E and 5-3F**).

### **5.9 Cell-type specific enrichment of mosaic changes in the aging brain and in AD**

We next used these cell type annotations to evaluate the cell-type-specific mutational burden of each cell type, first focusing on overall burden irrespective of disease status. We found that glial cells show a 34.6% higher mutational burden than neuronal cells ( $p=3.3e-5$ ), consistent with their continued proliferation, in contrast to post-mitotic neurons (**Fig. 5-3a**). Mosaic burden for both neuronal and glial cells correlated with age ( $r=0.09$  and  $r=0.26$ ) (**Fig. 5-3b**), consistent with the enrichment we found for aging-associated mutational signatures (**Fig. 5-1e**). Senescent cells harbored a 1.81-fold higher mutational load than neurons ( $p=2e-7$ ) and 1.34-fold higher than glia ( $p=6.6e-4$ ), suggesting that mutational load may be in part contributing to their senescent status. Male and female subjects did not show significant differences in mutational burden (**Supplementary Fig. 5-S9**).

We next studied cell-type-specific differences in mutational burden between AD and non-AD individuals. Those with AD showed significantly-enriched mutational burden in neurons (24% increase,  $p=0.042$ ), glia (18.5%,  $p=0.014$ ), and senescent cells (34%,  $p=0.019$ ) (**Fig. 5-3c-d**), with the neuronal signal driven primarily by excitatory neurons (30%,  $p=0.029$ ), and the glial signal by astrocytes (24%,  $p=0.038$ ) and oligodendrocytes (17.5%,  $p=0.02$ ). These results indicate that our observed overall increase in mutational burden of AD individuals does not stem solely from a difference in cell type composition, but is instead reflected in each of several cell types individually. Not surprisingly, these results indicate that neurons (which have been the primary focus of previous studies of brain mosaicism<sup>88,93</sup>) are not the sole target of somatic mutations in AD, and that glial cells can show equally-strong mutational burden and disease enrichment. Our results also show that somatic variants implicate both neuronal and glial cell types, in contrast to weak-effect common variants from GWAS at the other extreme of the allele-frequency spectrum, which enrich exclusively in microglial enhancers<sup>62,264</sup>.



**Figure 5-3. Mutational burden distribution.** **a.** Mutational burden distribution (x-axis) across individuals (points) for each cell type (rows). **b.** Mosaic burden (y-axis) by age (x-axis) with linear model fit (line) and 95% confidence interval (shade). **c.** Mosaic burden (shade) across cells with mutations (points). **d.** Mosaic burden differences between Non-AD (top) vs. AD (bottom) individuals by cell type (color). **e.** Sub-clustering of most-abundant cell types. **f.** Phenotypic enrichment for cell subclusters (hypergeometric test vs. main cluster).

### 5.10 Increased mutational burden correlates with expression changes at the single-cell level

For all major cell types, greater mutational burden at the single-cell level correlated with altered single-cell expression patterns (**Fig. 5-3e-f**, **Supplementary Fig. 5-S11**), indicating a relationship between genetic perturbations (somatic mutations) and phenotypic outcome (gene expression) at the cellular level. Unbiased expression-based sub-clustering (**Supplementary Fig. 5-S10**) of excitatory neurons, oligodendrocytes, and senescent cells without any guidance from mutational information (**Fig. 5-3f**), resulted in cell states with significant differences in burden (**Supplementary Fig. 5-S11a**).

Moreover, greater mutational burden correlated with both greater enrichment in AD vs. non-AD individuals (**Supplementary Fig. 5-S11b**), and increased gene-expression proximity to senescent cells (**Fig. 5-3e**), indicative of a gradient of mutational burden across cell-state space (**Fig. 5-3c**). For example, excitatory neuron subcluster Exc0 vs. Exc1 showed significantly-higher mutational burden ( $P=3.6e-9$ ), significantly-higher AD vs. non-AD enrichment ( $P=2.98e-7$ ), and highest proximity in gene-expression space to senescent cells (**Fig. 5-3e**). Similarly, Sen0 vs. Sen1 showed higher burden ( $P=0.026$ ), higher AD vs. non-AD enrichment ( $p=1.13e-7$ ), and gene expression patterns more distant from other cell types. For oligodendrocytes, Oli1 showed higher burden and higher AD vs. non-AD enrichment, though its expression patterns were shifted away from the senescent cluster, possibly indicative of disease response as we previously reported<sup>37</sup>.

Beyond these cluster-level correlations between greater burden and gene expression changes,

we found a cell-level correlation between the two, with a gradient of mutational burden along the expression-based two-dimensional manifold t-SNE clustering of cells (**Fig. 5-3c**), with mutational burden highest for senescent cells and progressing to lower burden in the periphery of the t-SNE plot. Even within each subtype, the previously-noted greater burden in AD individuals continued to hold (**Supplementary Fig. 5-S11c**), but was only significant for the largest cell states due to the decrease in power.

### 5.11 AD-specific mosaic mutational burden concentrates in AD-related genes and pathways

In addition to showing higher overall mutational burden, AD individuals showed a strong concentration of somatic mutations in specific genes and biological processes, which was not seen for normotypic individuals, providing insights on how they contribute to AD pathology in individuals carrying these mutations. In total, 21 genes showed gene-level increased mutational burden in AD (**Fig. 5-4a, Supplementary Figure 5-S12; Supplementary Table S1**). In excitatory neurons, seven genes showed greater mutational burden in AD, with roles in: protein degradation (*USP34*, *MYCBP2*), DNA damage (*SETX*)<sup>265</sup>, axon growth and transport (*NEFM*)<sup>266</sup>, endoplasmic reticulum tubule turnover (*CCPG1*)<sup>267</sup>, and synaptic function (*CYFIP2*)<sup>268</sup>, *PCDH9*<sup>269,270</sup>). In oligodendrocytes, 13 genes with greater mutational burden in AD include: myelin-associated genes (phosphodiesterase *CNP*)<sup>271</sup>, chaperone *CRYAB*)<sup>272</sup>; dendrite development (*DOCK9*)<sup>273</sup>; cytoskeleton/scaffold components (axonal growth gene *FEZ1*)<sup>274</sup>; membrane-cytoskeleton linker *ANK2*)<sup>275</sup>; microtubules-actin linker *MACF1*)<sup>276,277</sup>, which also showed increased burden in astrocytes); and protein aggregate regulation ( $A\beta$ -downregulated *ATP5B*)<sup>278</sup>; endosomal protein-aggregate mediator *ZFYVE16*; and glial-specific tauopathy inclusion mediator *CRYAB*)<sup>279,280</sup>).

We next studied the individual cells where these mutations were found in the context of the tSNE gene expression landscape, pinpointing both the cell type and the cell state of each mutation-harboring cell for each gene, and the clinical status (AD vs. non-AD) of the individuals where the mutation was detected (**Fig. 5-4b, Supplementary Fig. 5-S12-S14, Supplementary Table S2**). Indeed, we find that strong-effect missense and frame-shift mutations from these genes are concentrated in cells from relevant cell types, and that these cells come from AD individuals, consistent with the AD-specific and cell-type-specific enrichment of the mutational burden for these genes. Moreover, we find that the cells showing highest expression and highest coverage for these genes indeed match the cell types with AD-specific mutational enrichment, indicating that these enrichments are not an artifact of high-burden cells or senescent cells, which do not contribute any of the observed mutations. Even within excitatory neurons and oligodendrocytes, mutations in these genes are found in cells from sub-clusters corresponding to senescent-distal lower-burden cells, rather than the senescent-proximal higher-burden sub-clusters. These properties provide a strong affirmation of our discovery strategy, and indicate that these gene-enriched mutational enrichments are less likely to represent solely a consequence of dysregulation leading to increased mutation in these pathways.

We next searched for pathways that show greater mutational burden in AD individuals, even if individual genes within them did not rise to significance alone. First, we tested a set of 12 known AD genes and found significant mutational burden in AD, both overall ( $q=0.024$ ) and in a cell-type-specific manner with burden concentrating in oligodendrocytes (23 mutations vs. 14,  $q=0.024$ ) and astrocytes (8 vs. 2,  $q=0.025$ ) (**Fig. 5-4c**), indicating that inherited-variant-containing genes are also recurrent targets of somatic mutations. Second, we tested 658 Reactome pathways with sufficient coverage for greater AD vs. non-AD mutational burden, and found 49 significantly-enriched pathways for excitatory-neuron mutations, and 80

significantly-enriched pathways for oligodendrocyte mutations (vs. only 19 and 4 pathways in controls) (**Fig. 5-4d, Supplementary Figure 5-S15 and Supplementary Table S3**).

For neuronal mutations, enriched pathways with greater AD burden included: neuron energy regulation (glycosylation, sialic acid metabolism); protein turnover and degradation (neddylation); DNA damage repair (excision repair); stress response (hypoxia response, syndecans, EPHA-mediated growth cone collapse); transport (mitochondrial import, protein localization, vesicle transport), and diverse signaling pathways. Oligodendrocyte mutations showed significantly greater AD enrichment in pathways involving: lipid metabolism (fatty acid metabolism); signaling (Notch1/3); endocytic vesicular transport (Golgi-to-ER retrograde, anterograde, endosomal sorting, clathrin-mediated, COPI); cell cycle (kinetochores); DNA damage (p53 phosphorylation); and proteostasis (HSP90 chaperones).

We also evaluated overall mutational burden in the aged human brain, across both AD and non-AD individuals, to search for genes with significantly-greater and significantly-lower burden. Greater-burden genes were enriched in: microtubule maintenance (*ATM, USP33, LYST, MARK1, CEP350, PCM1, BICD1, CTNNA1*) and synaptic cleft (*NLGN1, LGI1*) for neurons; DNA binding and gene expression regulation (*SETX, ZNF573, DENND4B, JMJD1C*) for glial mutations; and genomic stability, cytoskeleton dynamics and cell survival across cell types (**Fig. 5-4e, Supplementary Table S4, Supplementary Fig. 3-S16-S17**).

Lower-burden genes were enriched in: axon targeting (*STXBP1, UCHL1*) and RNA binding (*COLGB1, HSP90AB1, PEBP1*) in neurons; cell-identity genes of myelination (*ABCA2, CNTN2, MBP, PLP1, MAG, MYO5A, CLU*) and microtubule-associated processes (*ATRX, CNTN2, HSPA1B, TUBA1A, GAPDH, TUBB4A*) in glia; and housekeeping and cell-identity marker genes across cell types. Mutational depletion might reflect a combination of decreased mutation rate across genomic regions, differential targeting of repair machinery proteins, or cellular-level selective pressure that deplete mutations leading to cell death.

Glia and neurons showed mutational enrichment and depletion in different genes, and their enriched-mutation pathway were largely distinct. However, they both showed mutational depletion in tau-regulating proteins (*CLU, GPI, HSP90AB1, STXBP1* in neurons,  $p_{\text{adj}}=0.042$ ; *CLU, HSP90AA1, ACTB, BIN1* in glia,  $p_{\text{adj}}=0.019$ ), with *CLU*, a GWAS<sup>76</sup> and neuropathology<sup>281</sup> AD-associated secreted chaperone that inhibits protein aggregation, showing mutational depletion in both.





individual cell, annotated each covered gene in each cell with its coverage and number of hits as one case. We kept only protein-coding mutations annotated by SnpEff as missense or LoF variants. We reduced our dataset to 41,755 missense and LoF variants from 4,014 cells from the 36 individuals with paired WGS that were not filtered out in mosaic calling or transcriptome QC. Within each major cell type (and jointly for all neurons, all glial cells, and all cells), we created a stratified dataset by stratifying by sex (male/female), age (17 individuals up to 86.9 and 19 greater than 86.9), and by binning each cell by gene's (sample's) coverage. We tested all genes with at least one mutation and 10,000 bp of coverage in cases and controls. For each tested gene, we sampled 10,000 null genes from the stratified dataset, matching the gene's coverage in each cell, and compared the gene's burden to that of null examples to obtain a significance p-value. We used the *qvalue*<sup>282</sup> R package to perform FDR correction. Per-cell-type pathway enrichment was performed using *gprofiler*<sup>283</sup> for an over-representation test for genes with *q*-value < 0.05 in each cell type (**Supplementary Table S3**).

**Gene-level mosaic burden in disease by cell type.** Second, to discover genes that were significantly over-burdened in disease, we compared the relative rates of mosaic events in AD individuals versus controls as defined by clinical ascertainment of cognitive impairment. Within each cell type (and jointly for all neurons, all glial cells, and all cells), we tested all genes with at least 10,000 bp of coverage both in AD cases and controls and with at least one mutation in cases. We performed a two-sided exact binomial test for the number of mutations in the number of covered bp in cases against the estimated rate of mutation in controls, adding a pseudocount of one mutation to both cases and controls. FDR correction was performed using the *qvalue*<sup>282</sup> R package. We also performed a permutation test as a secondary test for over-burdened genes in AD individuals. We performed the stratification as in the case of the disease-agnostic analysis, further stratifying the dataset by the estimated burden in each gene, agnostic to disease status (binned into 0-10<sup>-5</sup>, 10<sup>-5</sup>-10<sup>-4</sup>, 10<sup>-4</sup>-10<sup>-3</sup>, and 10<sup>-3</sup> or greater). We tested all genes flagged as significantly over- or under-burdened by the binomial test. For each tested gene, we sampled 10,000 matched null examples from the stratified dataset and compared the overall ratio of hits in AD individuals to controls to the ratios from null examples in order to obtain a significance p-value (**Supplementary Table S4**). We evaluated the enrichment of mutations for bona-fide AD genes (*ABCA7*, *APOE*, *APP*, *BIN1*, *CLU*, *CR1*, *PICALM*, *PLD3*, *PSEN1*, *PSEN2*, *SORL1*, and *TREM2*) in aggregate in each cell type and overall using a binomial test and FDR correction as above.

**Pathway-level burden in disease by cell type.** We evaluated pathway-level differences in mosaic mutational burden in cases versus controls on the Reactome<sup>284</sup> pathway maps obtained from the *reactome.db* R package<sup>285</sup>. In each cell type and in the aggregated cell types, we aggregated mosaic mutations and coverage from the gene-level to the pathway-level and kept all pathways with least 10,000 bp of coverage both in AD cases and controls, at least one mutation, between 25 and 250 genes, and where at least 90% of the pathway genes are expressed in the tested cell type. To compare each pathway's rate of cases to controls, we performed a binomial test followed by FDR calculation for multiple testing correction as in the gene-level analysis (**Supplementary Table S5**).

### 5.13 Discussion

Our study constitutes the first report of cell-type-specific mutational burden in affected-vs-unaffected individuals in AD, by leveraging single-cell full-transcript RNA sequencing technology (SMART-seq2) and matched whole-genome-sequencing to jointly infer somatic coding variants and their corresponding cell type and cell state, thus enabling us to understand the functional impact of mutational changes on genes and pathways.

This approach allowed us for the first time to partition somatic burden into the cell types and

states where it occurs, and to survey the mosaic burden across multiple neuronal and glial cell types (while previous studies focused exclusively on neurons<sup>93,286</sup>). This revealed that overall mutational burden is higher in glial cells than in neurons, consistent with continued glial cell divisions in contrast to post-mitotic neurons. It also revealed that overall mutational burden, irrespective of disease status, concentrates in distinct genes and pathways between glial cells and neuronal cells, while other genes and pathways are depleted for mutations. Along with neuronal mutations that directly impact neuronal function, glial mutations can play important roles in brain changes, through indirect effects on neuronal functions via glial cells, and through direct effects on glial functions.

Our paper is also the first to report disease-enriched somatic mutational burden for a neurodegenerative disorder, where accumulation of somatic mutation burden across a lifetime of exposures is particularly relevant. We found that strong-effect somatic mutations in AD cluster in specific cell types, with excitatory neurons, astrocytes and oligodendrocytes emerging as primary drivers of AD-associated somatic mutations, highlighting the combined roles of neuronal cells, glial cells, and their interactions in AD<sup>287,288</sup>.

This greater mutational burden in AD also concentrated in specific genes and pathways, which differed between cell types. This included AD-associated genes and pathways from GWAS, including cholesterol metabolism, endosomal trafficking, and cytoskeletal systems. This also included AD-associated genes and pathways from rare variant and whole-genome-sequencing studies, including in protein aggregation, degradation, neurofilament tangles, and tau-associated proteins<sup>289–291</sup>. Our somatic burden analysis also captured novel pathways that were not previously genetically implicated, including oligodendrocyte myelination, post-translational modification, microtubule dynamics, and neuronal energy regulation through glycosylation, and stress response. These provide a new set of actionable targets for therapeutic development, which may have direct strong effects on neuronal function, but which previously lacked a genetic basis, as they lie outside the previously-studied range of the allele frequency spectrum from common to rare alleles.

In addition to elucidating the roles of neuronal and glial cell types in AD, our work highlights the importance of a new “senescent” cell population in the aging brain, which may represent a distinct cell type associated with aging, and was strongly enriched in AD. Even though these senescent cells are positioned in the two-dimensional tSNE embedding between neurons and oligodendrocytes, they show a distinct set of marker genes than either cell type, and instead might represent a collection of cells from different cell types that have lost their identity. As a group, senescent cells showed the largest mutational burden, both overall, and specifically in AD, which may have contributed to them losing their differentiated identity and entering a low-transcription and possibly pre-apoptotic state.

Our approach also allowed us to capture not only cell type identity, but also cell state, by relating the single-cell gene expression profile of each cell to the specific set of mutations it carries, and thus connect genotype to phenotype at single-cell resolution. Indeed, we saw a gradient of somatic mutational load along single-cell gene-expression space, indicating a strong correlation between the two, both between cell sub-clusters, and within each sub-cluster, indicating a continuum of mutational burden vs. expression dysregulation. For example, neuronal cells closest in gene expression to the senescent cells also showed increased mutational burden, while neuronal cells furthest in gene expression showed the lowest mutational burden. Similarly, oligodendrocyte gene expression sub-clusters correlated with mutational burden and distance from senescent cells. This concordance between accumulation of mutations and gene expression profiles at the cellular level may indicate that mutational burden accumulation results in phenotypic alterations; conversely, dysregulation at the single-cell level might also lead to

increased mutations via DNA damage, decreased mismatch repair, and increased double-stranded breaks, which were previously-implicated in AD<sup>292</sup>.

It is challenging to distinguish whether mutational burden drives pathology, results from cellular dysfunction, or both. In the first case, high mutational burden can lead to cellular dysregulation, gene expression changes, and cellular pathology. Conversely, cellular dysfunction caused by other mechanisms (eg. protein aggregation, impaired protein turnover, DNA repair disruption) can lead to increased DNA damage and somatic mutation accumulation. Both processes may be acting simultaneously as well, with increased dysregulation leading to increased mutation accumulation and vice versa.

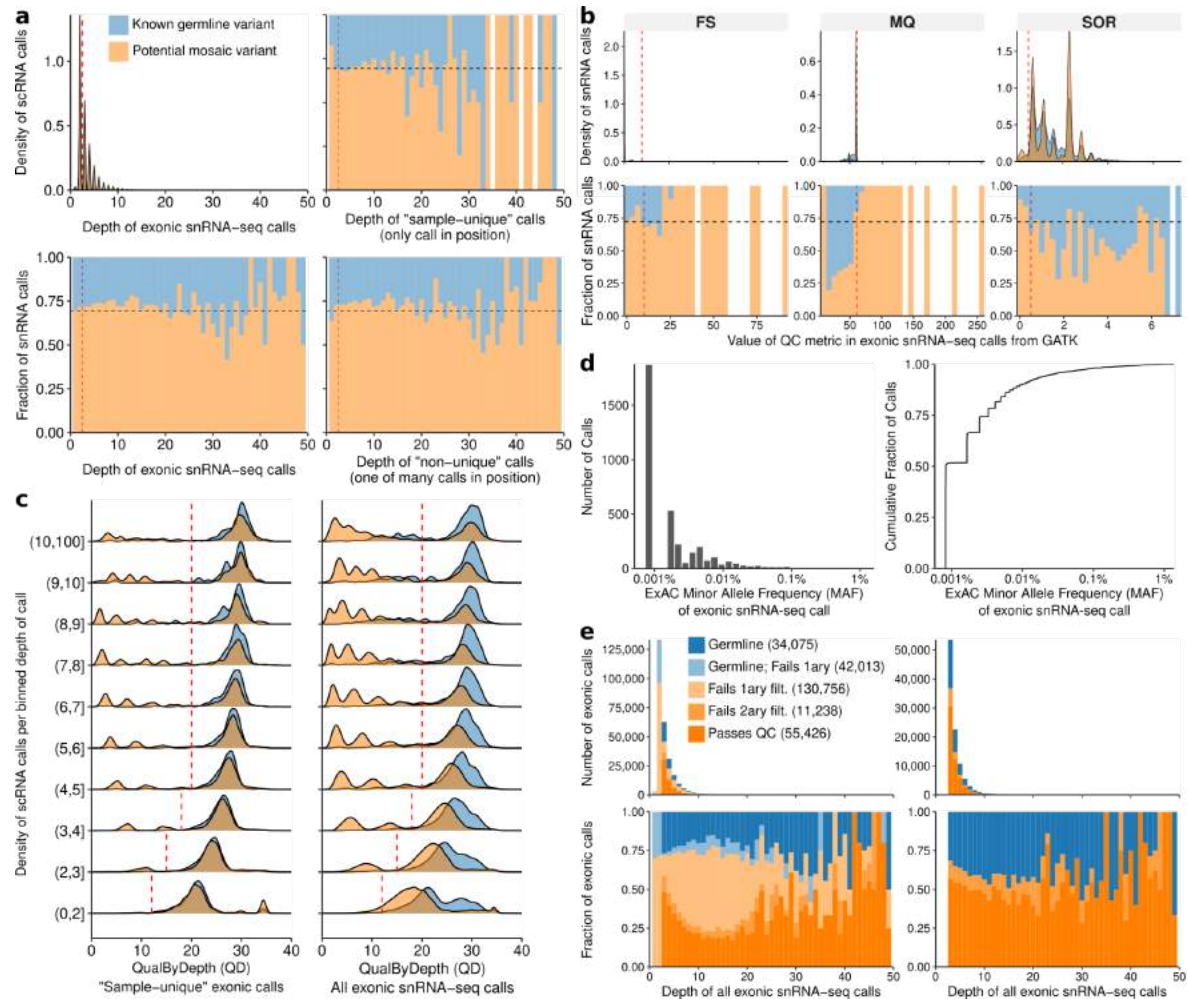
It is important to note that our approach only captures a fraction of the mosaic landscape that exists in aging and AD. First, we only annotate mosaicism in expressed transcripts, which has the advantage of focusing our sequencing efforts on strong-effect exonic mutations and on expressed genes, but does not capture non-coding regions, low-abundant transcripts, or alleles that lack expression. Second, RNA-based mutational inference can be biased due to transcription-associated mutations, transcription-associated repair, and the decreased stability of RNA relative to DNA, and thus is not well-suited for estimating an overall rate of somatic mutation; as such, we instead focus on relative differences here. Third, we do not capture the mosaic burden and cellular pathology incurred by cells that have already apoptosed, and even for senescent cells that exhibit low levels of transcription, our RNA-based approach likely misses many of their mutations in unexpressed segments of their genome.

Going forward, these limitations can be remedied in part by the use of technologies that simultaneously profile DNA and RNA from the same cell, and thus can decouple expression from genotype, and more directly ascertain how mosaic mutations relate with cellular state and gene expression levels. In addition, spatial transcriptomic approaches can enable direct comparison of cellular phenotypes in the context of surrounding pathology and their somatic mutations of neighboring cells.

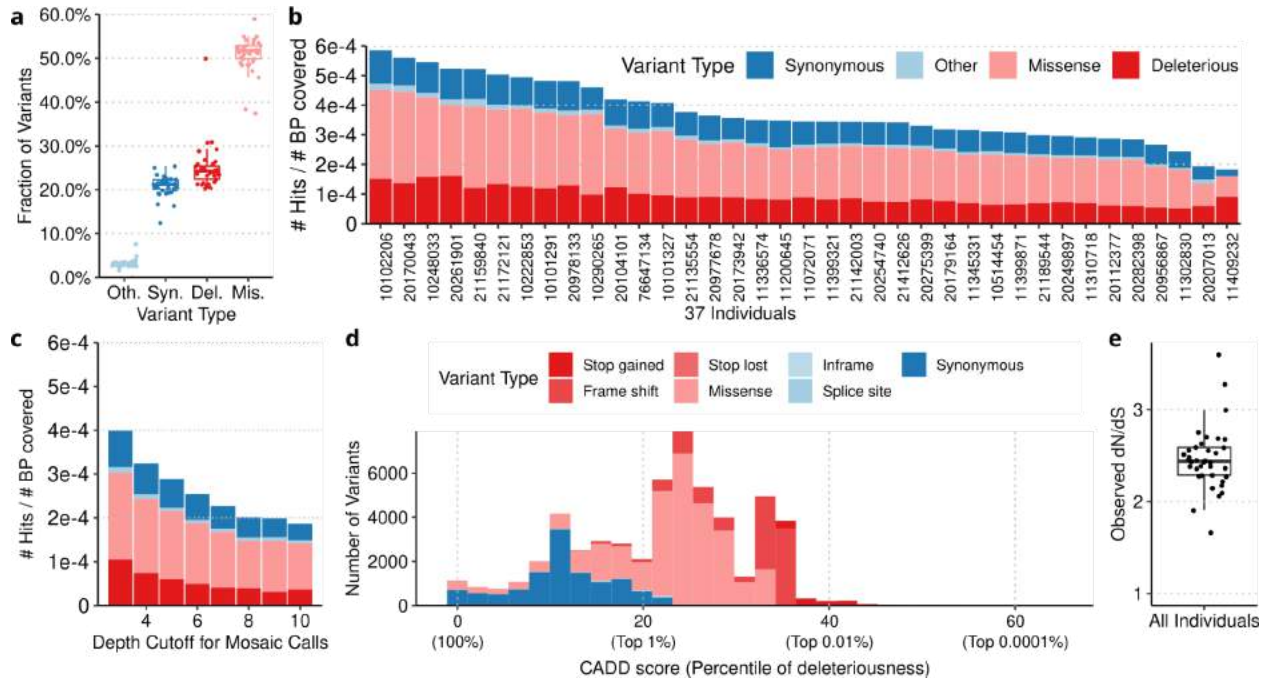
Our study provides a foundational framework for analyzing such datasets, and initial insights on the relationship between somatic mutational burden, cellular state, and dysregulation of specific genes and pathways contributing to AD pathology. Our initial findings can be reinforced and refined as new technological advances make joint profiling of DNA and RNA from large numbers of cells commonplace in the study of the aging human brain in health and disease.



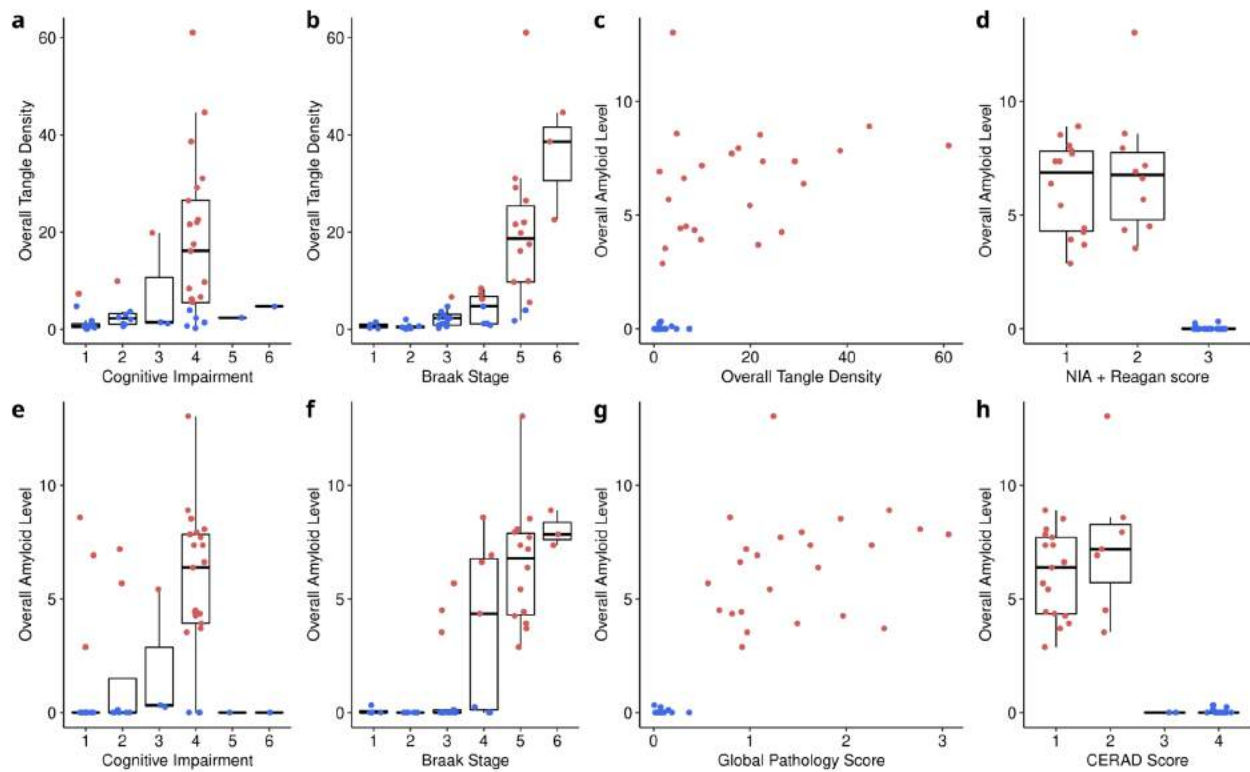
## 5.14 Supplementary Figures



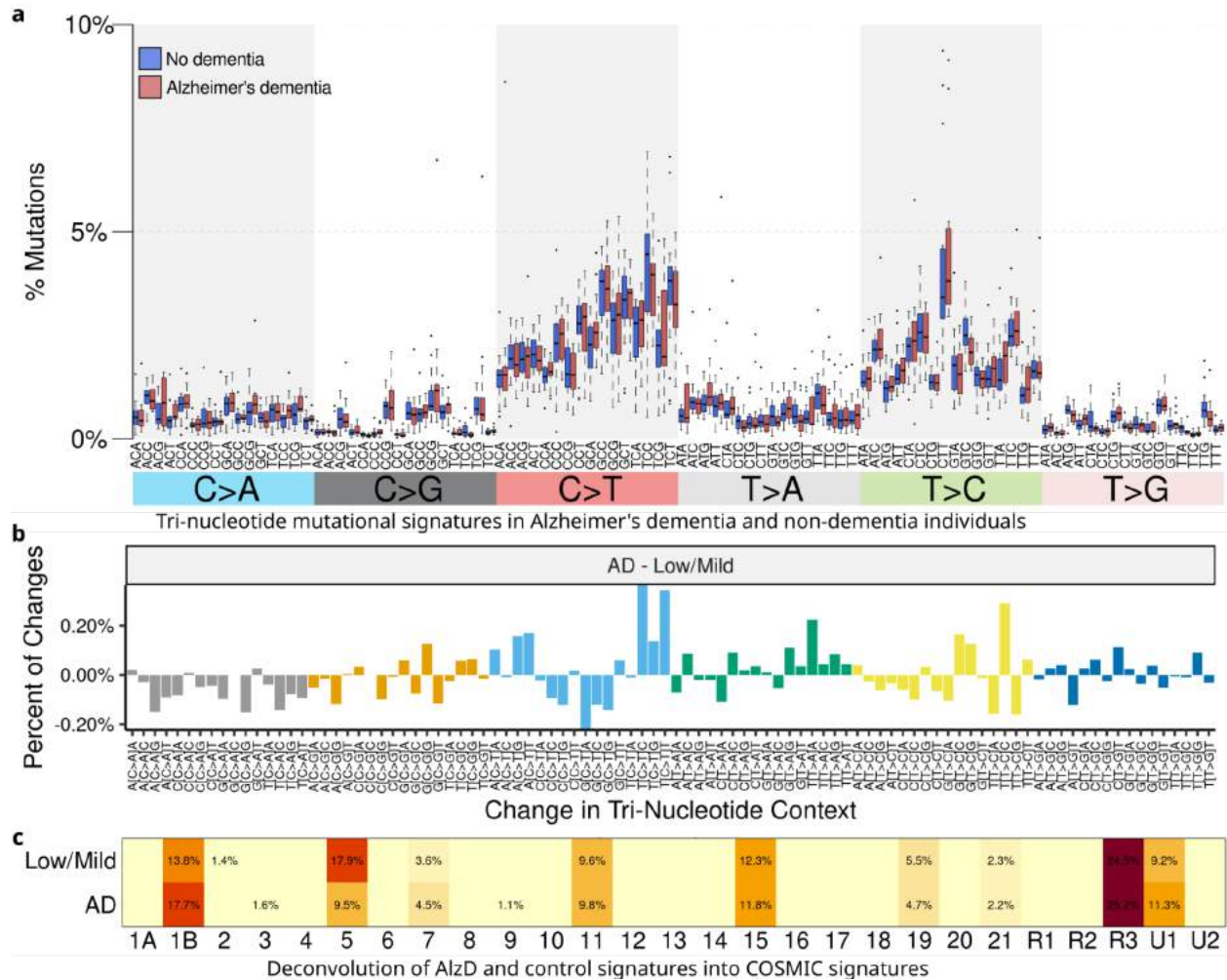
**Supplementary Figure 5-S1. Variant filtering.** **a.** Distribution of calls from snRNA-seq by depth of call (DP), split into known germline variants by intersection with WGS or potential mosaic variants. Top-left: density of calls by depth; bottom-left: stacked histogram showing relative fraction of germline and potential mosaic calls per depth; top-right and bottom-right: breakdown of calls into sample-unique calls, where the call is supported across multiple samples (cells) in the individual versus non-unique calls. Horizontal line labels average fraction of potential mosaic calls relative to all calls. **b.** Comparison of FS (FisherStrand), MQ (RMSMappingQuality), and SOR (StrandOddsRatio) variant filtration metrics between known germline calls and potential mosaic calls in order to establish filtration cutoffs (vertical lines). **c.** Comparison of quality by depth (QD) metric at 10 depth cutoffs between validated calls corresponding to germline versus potential mosaic calls. As QD is evaluated on all samples in one individual and many germline variants are found across multiple cells, cutoffs are established on sample-unique calls, for which the germline distribution is comparable to that of potential mosaic events. **d.** Minor Allele Frequency (MAF) distribution of calls found in ExAC with MAF < 1% (left) histogram of number of calls by MAF (right) cumulative distribution of all calls by MAF, where over 50% are at less than 0.001% MAF. **e.** Breakdown of calls by overall number (top) and relative fraction (bottom) including all calls (left) or just those passing the 1st round of filtering (right). 55,426 variants pass QC (in addition to 1,459 indel variants from merging calls, not shown).



**Supplementary Figure 5-S2. Mutational burden by variant type.** **a.** Distribution of variant type breakdown at the individual level, where each point is one individual, for missense (average: 51.3%), synonymous (21.5%), LoF (24.1%), and other (3.1%) mosaic variants. **b.** Overall mutational rate (number of mosaic hits over the number of bases covered at a callable depth) by individual and variant type. **c.** Overall mutational rate by depth of mosaic call and variant type. **d.** Distribution of CADD score against variant type for mosaic calls. **e.** Distribution of dN/dS rates across individuals.

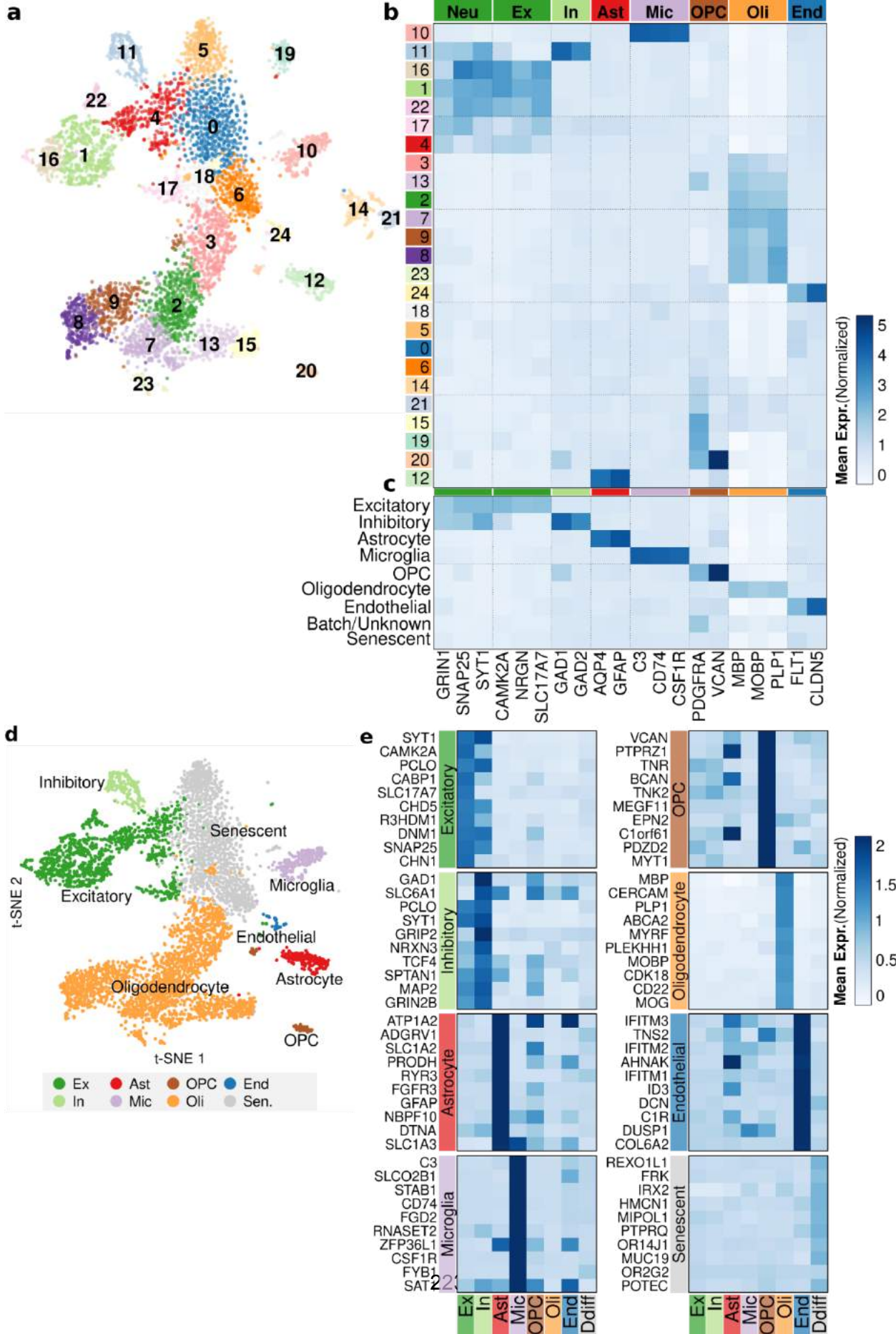


**Supplementary Figure 5-S3. Cohort Ascertainment Measures.** Comparison of Alzheimer's dementia and disease ascertainment measures over the full 47 individual cohort used for transcriptional analysis. **a-b.** Comparisons of cognitive impairment (CI) stage and Braak stage against overall tangles levels (colored by high tangles level in red, low in blue) **c-h.** Comparisons against overall amyloid level, colored by amyloid level (high in red, low in blue) of overall tangle density, NIA/Reagan score, CI stage, Braak Stage, global pathology score, and CERAD score.

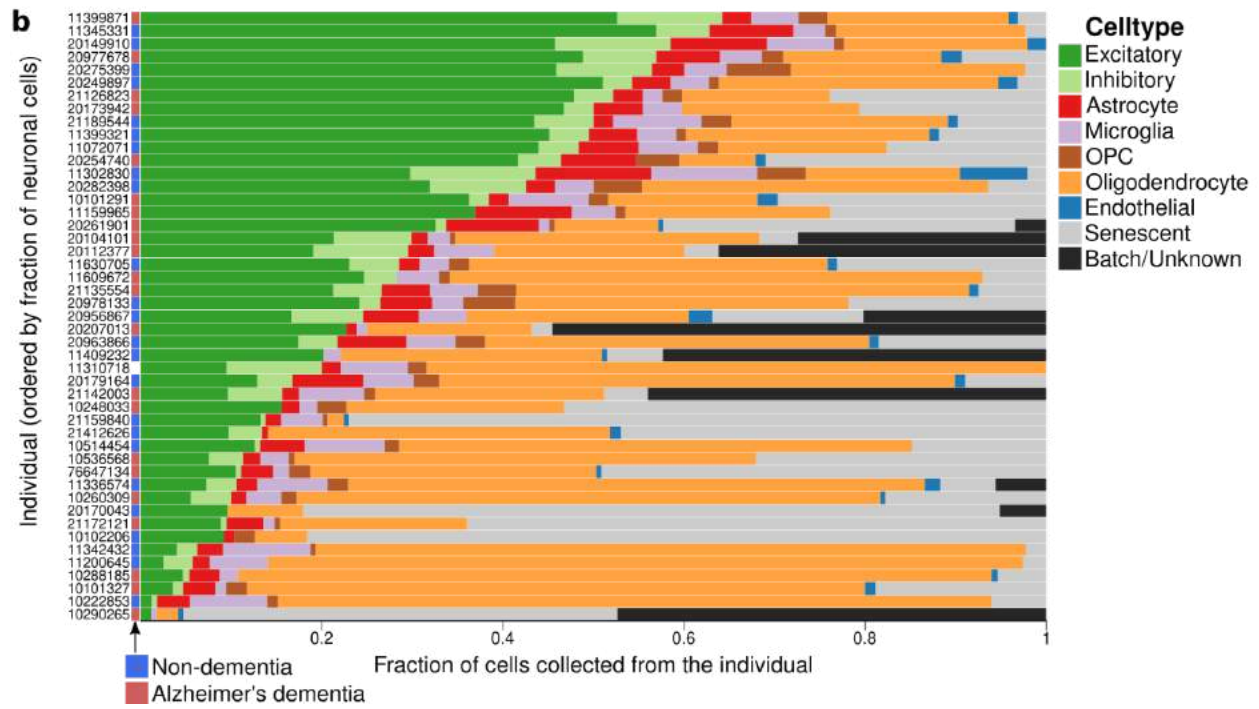
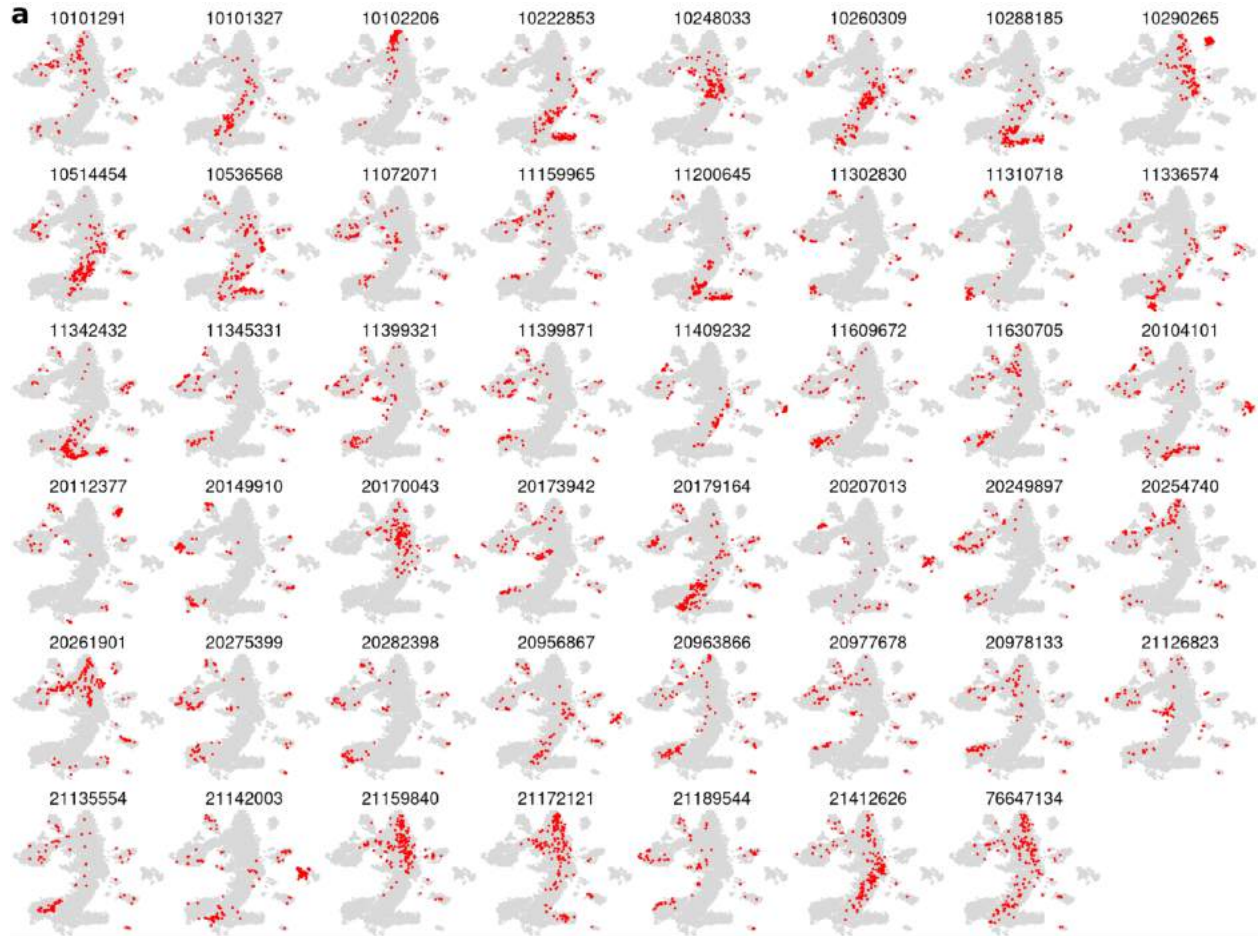


**Supplementary Figure 5-S4. Mutational signatures in disease.** **a.** Distribution of mutational signatures by cognitive impairment status (non-dementia in blue and AD in red), where each point is an individual. No significant differences were found for any signature after p-value adjustment for multiple testing. **b.** Difference in the overall mutational signatures between Alzheimer's dementia and non-dementia individuals. **c.** Mutational signature deconvolution of the overall dementia and non-dementia signatures (rows) according to the COSMIC v1 signatures (columns).

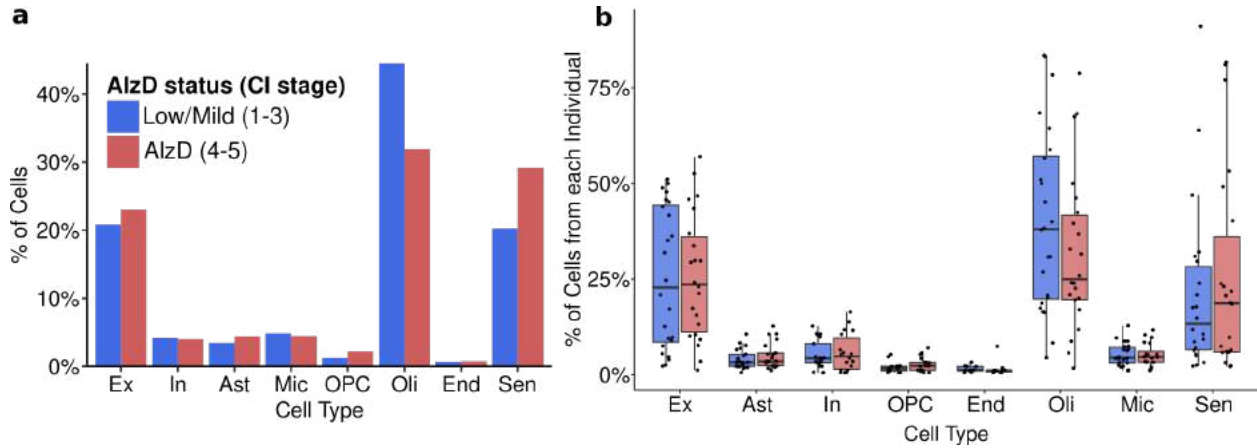




**Supplementary Figure 5-S5. Cell type identification.** **a.** Preliminary louvain clusters for the t-SNE of 5,821 cells. **b.** Comparison of known brain cell type markers to louvain clusters. Sub-clusters 14 and 19-21 are removed in further analysis. **c.** Comparison of marker genes to cell types merged from louvain clusters. **d.** t-SNE plot for 5,421 cells from 47 individuals across 8 merged cell clusters (1,170 excitatory and 221 inhibitory cells), four glial clusters (220 astrocytes, 255 microglia, 2121 oligodendrocytes, and 94 OPCs), 40 endothelial cells, 400 cells with strong individual-specific batch effect and cancer signatures (DNMT3A, COL6A3), and 1,300 senescent cells. **e.** Learned marker genes for each cell type. The top 10 markers are shown for each cell type against their average expression from each cluster.

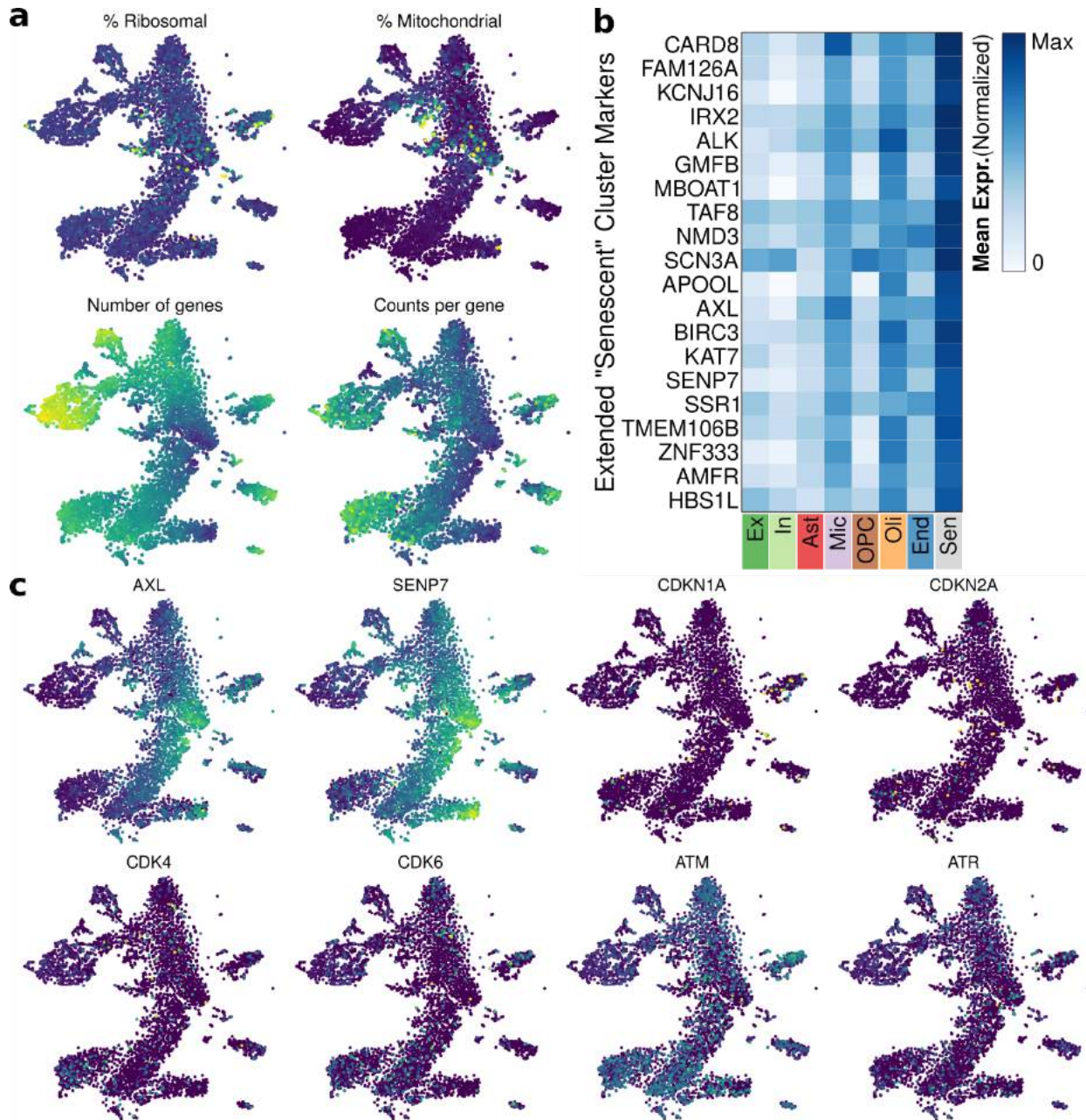


**Supplementary Figure 5-S6. Cell type breakdown by individual.** **a.** Locations of each of the 47 individuals' cells on the t-SNE of 5,821 cells before removal of batch cells. **b.** Cell fractions per individual for 8 major cell types and one batch. Individuals are labeled as Alzheimer's dementia or non-dementia as ascertained by their cognitive impairment status.

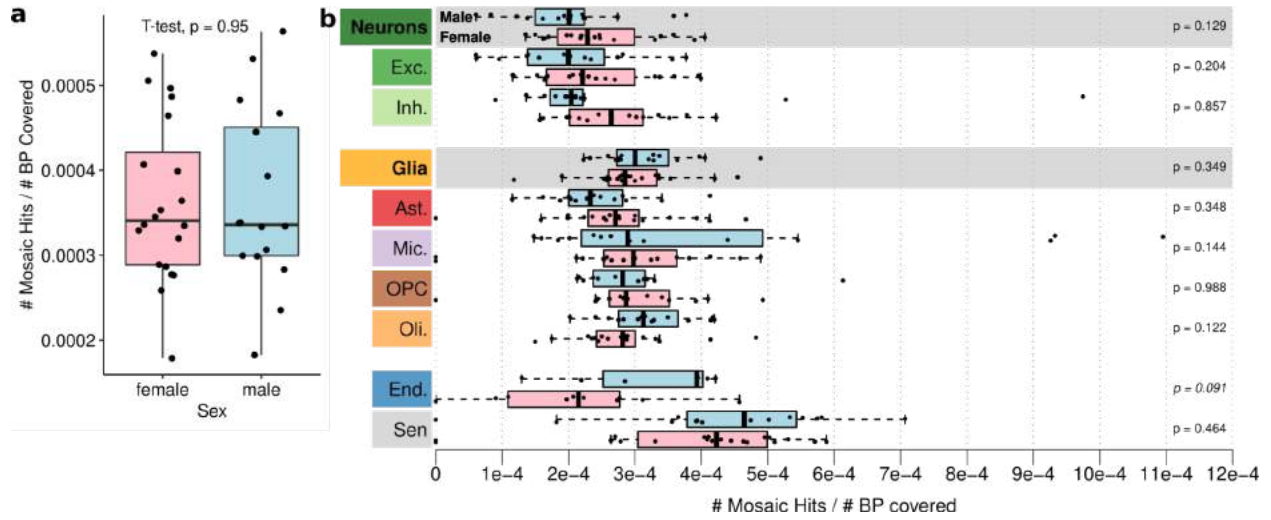


**Supplementary Figure 5-S7. Cell fractions by disease status.** **a.** Overall percentage of cells from all AD individuals (22, cogdx in 4-5, red) and control individuals (24, cogdx in 1-3, blue) coming from each cell type. One of the 47 individuals has cogdx level 6 and is not included. **b.** Per-individual fraction of each cell type.

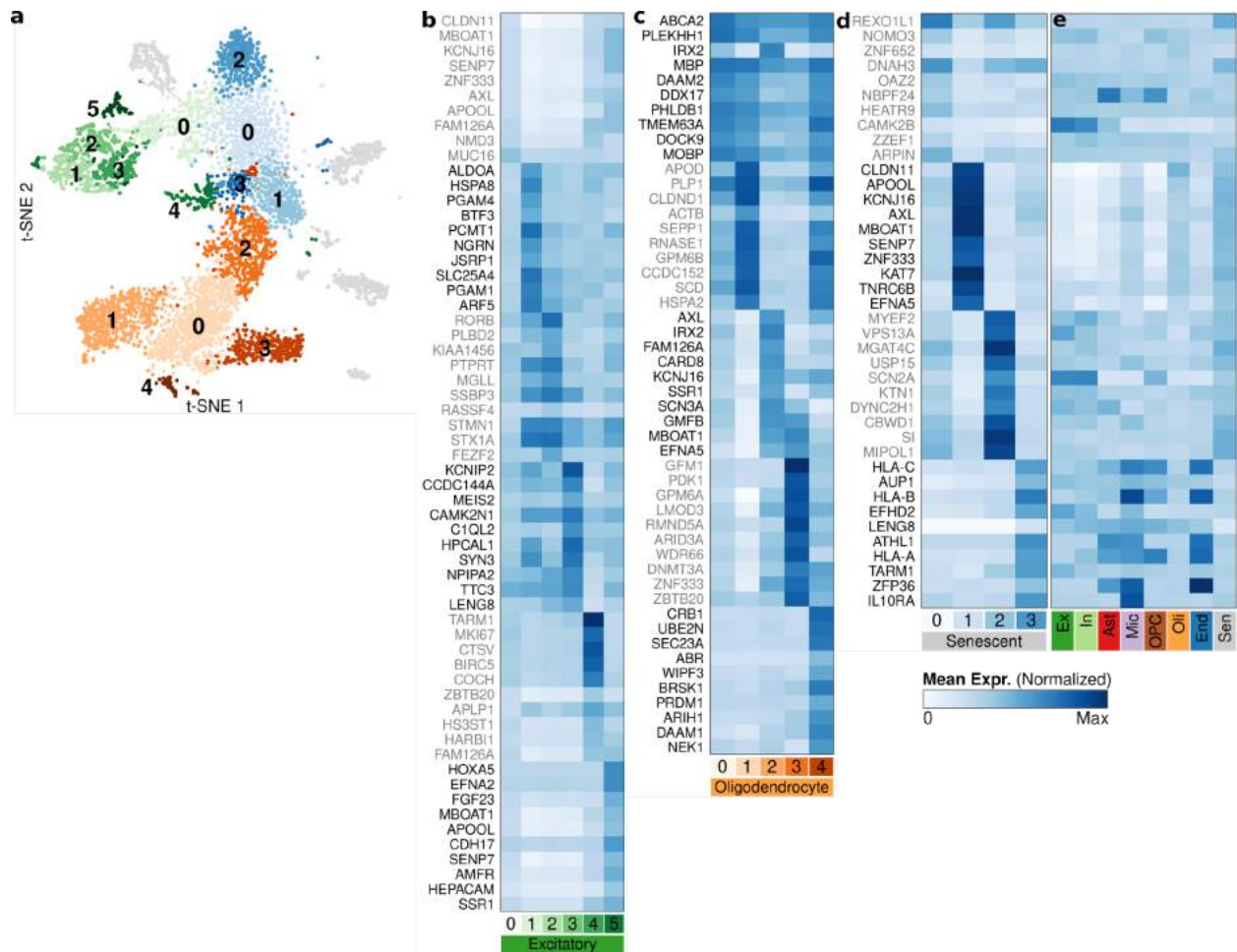




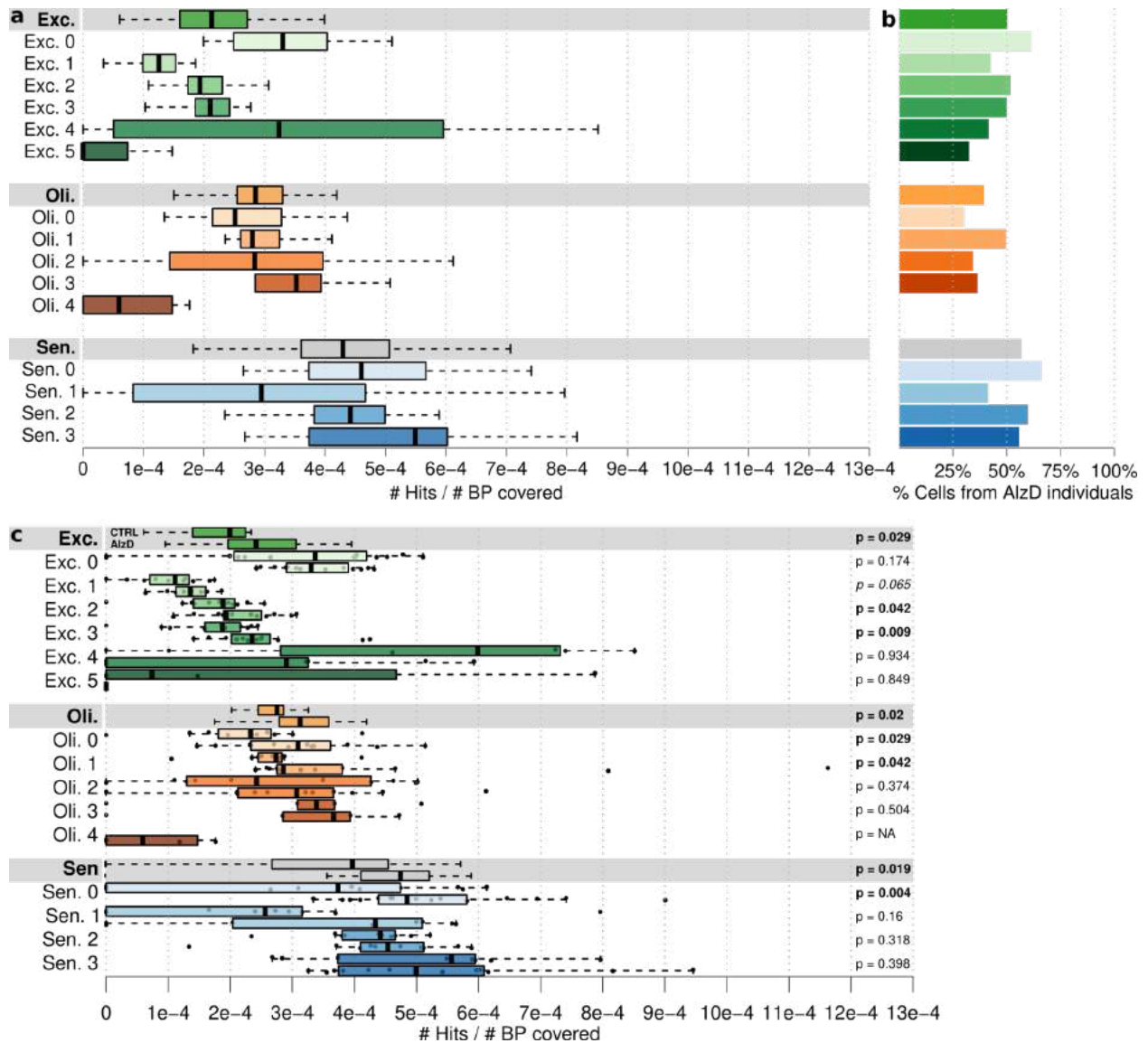
**Supplementary Figure 5-S8. Extended markers for t-SNE and "senescent" cluster. a.** Percentage of reads from ribosomal or mitochondrial genes, number of genes, and number of reads per gene in each cell. "Senescent" cluster cells show overall lower expression and less counts per gene. **b.** Top 20 "senescent" cluster markers where at least half of the cells show greater than average expression of the gene. Marker genes were computed for the cluster against all cells by a wilcoxon test. **c.** Classical senescence markers, including CDKN1A and CDKN2A, which regulate CDK4 and CDK6, show low signal across the cells in the t-SNE (dark = low expression).



**Supplementary Figure 5-S9. Mosaic mutational burden by sex.** **a.** Mosaic mutational burden (number of missense and loss-of-function mosaic hits over number of callable bases) by sex across all cells in 36 individuals. P-value is reported for two-sided t-test. **b.** Per-cell-type mosaic mutational burden by sex (female in pink, male in blue), for each cell type and neuronal and glial cells overall. P-values for two-sided t-tests are reported.

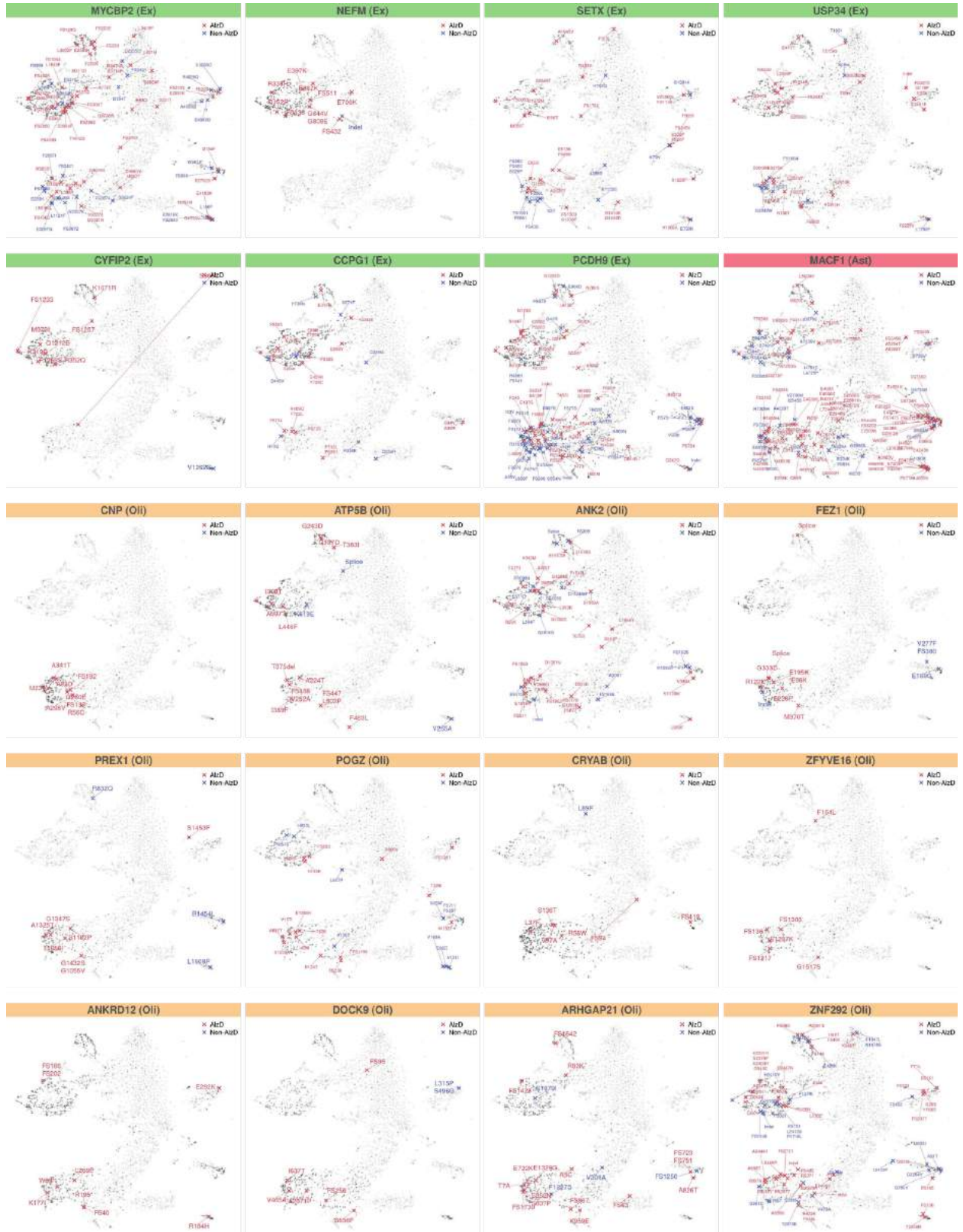


**Supplementary Figure 5-S10. Major cell type subclusters and marker genes. a.** Subclusters for excitatory neurons, oligodendrocytes, and senescent cells. **b-d.** Average expression of the top 10 markers for each excitatory, oligodendrocyte, or senescent subcluster. Marker genes were computed for each subcluster against all cells in the main cluster by a wilcoxon test. **e.** Average expression of each of senescent subclusters' markers in the major cell types.



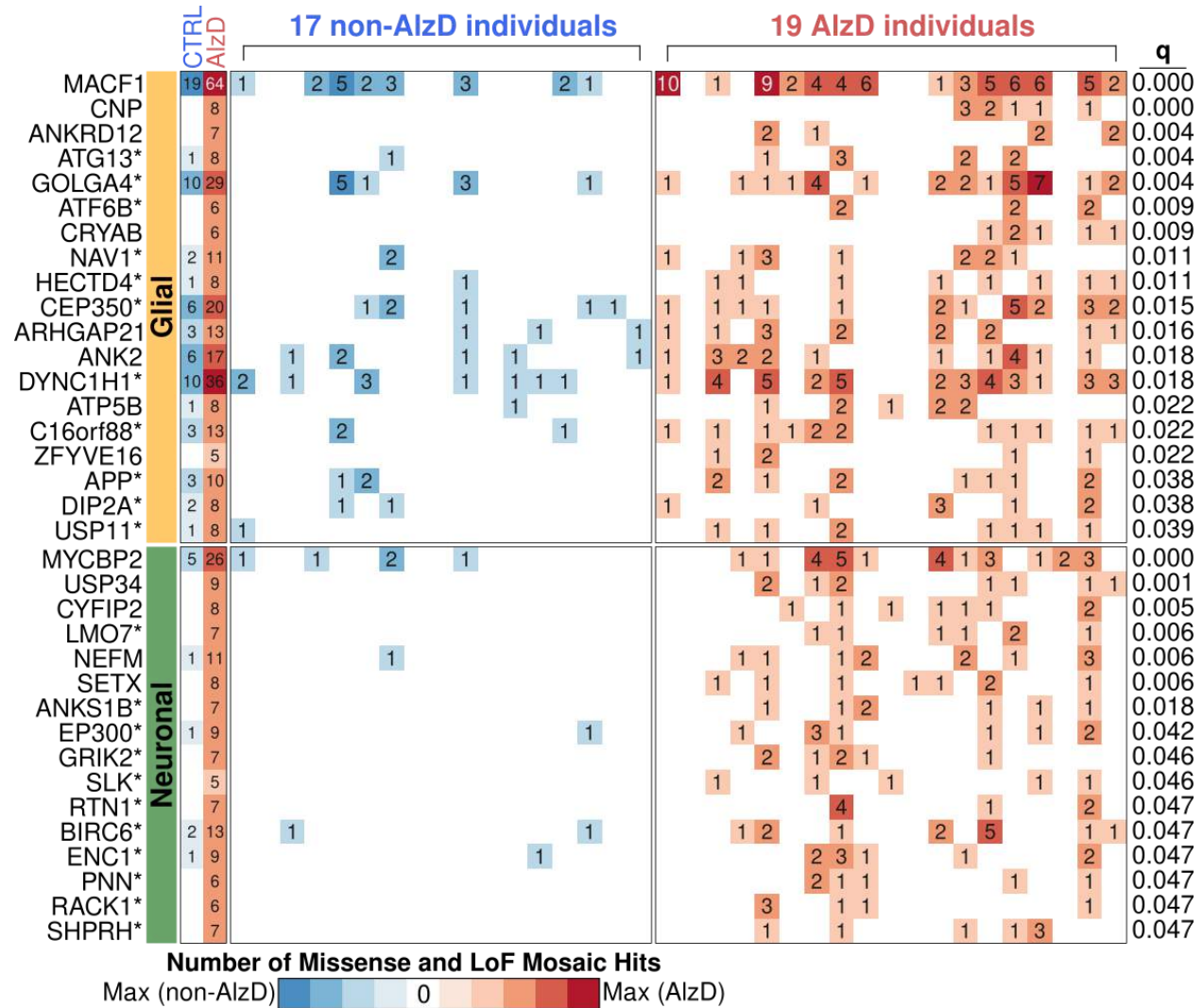
**Supplementary Figure 5-S11. Mosaic mutational burden in subclusters in disease. a.** Mosaic mutational burden (number of missense and loss-of-function mosaic hits over number of bases at a callable depth) for all 36 individuals in each of excitatory neurons, oligodendrocytes, and senescent cells and their subclusters. **b.** Percentage of cells in each set from Alzheimer's dementia individuals. **c.** Mosaic mutational burden by subcluster and Alzheimer's dementia status for each of excitatory neurons, oligodendrocytes, and senescent cells and their subclusters.



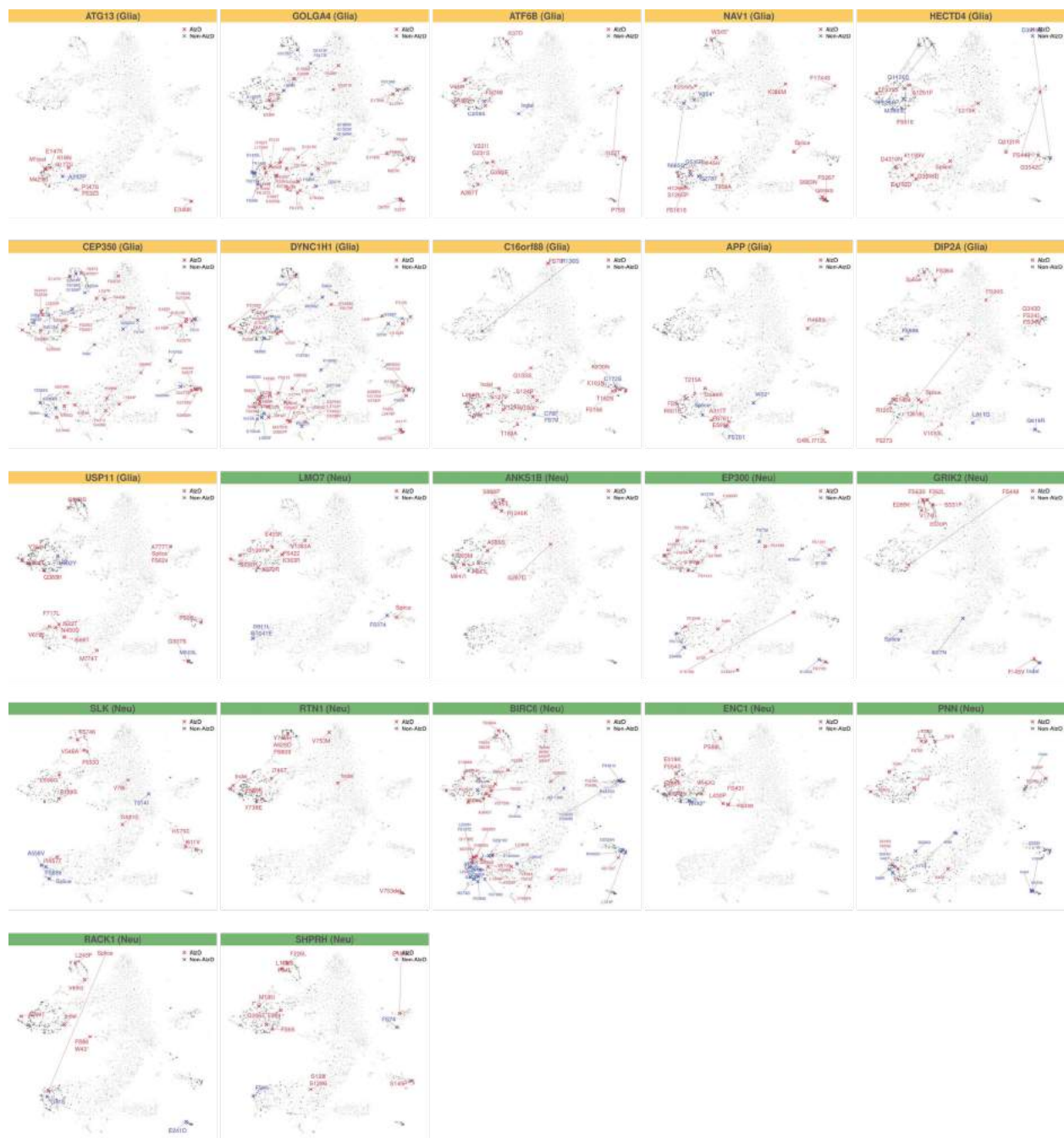


**Supplementary Figure 5-S12. Genes with significantly higher burden in AD.** Mosaic hits in each over-burdened gene in excitatory neurons, astrocytes, and oligodendrocytes are labeled

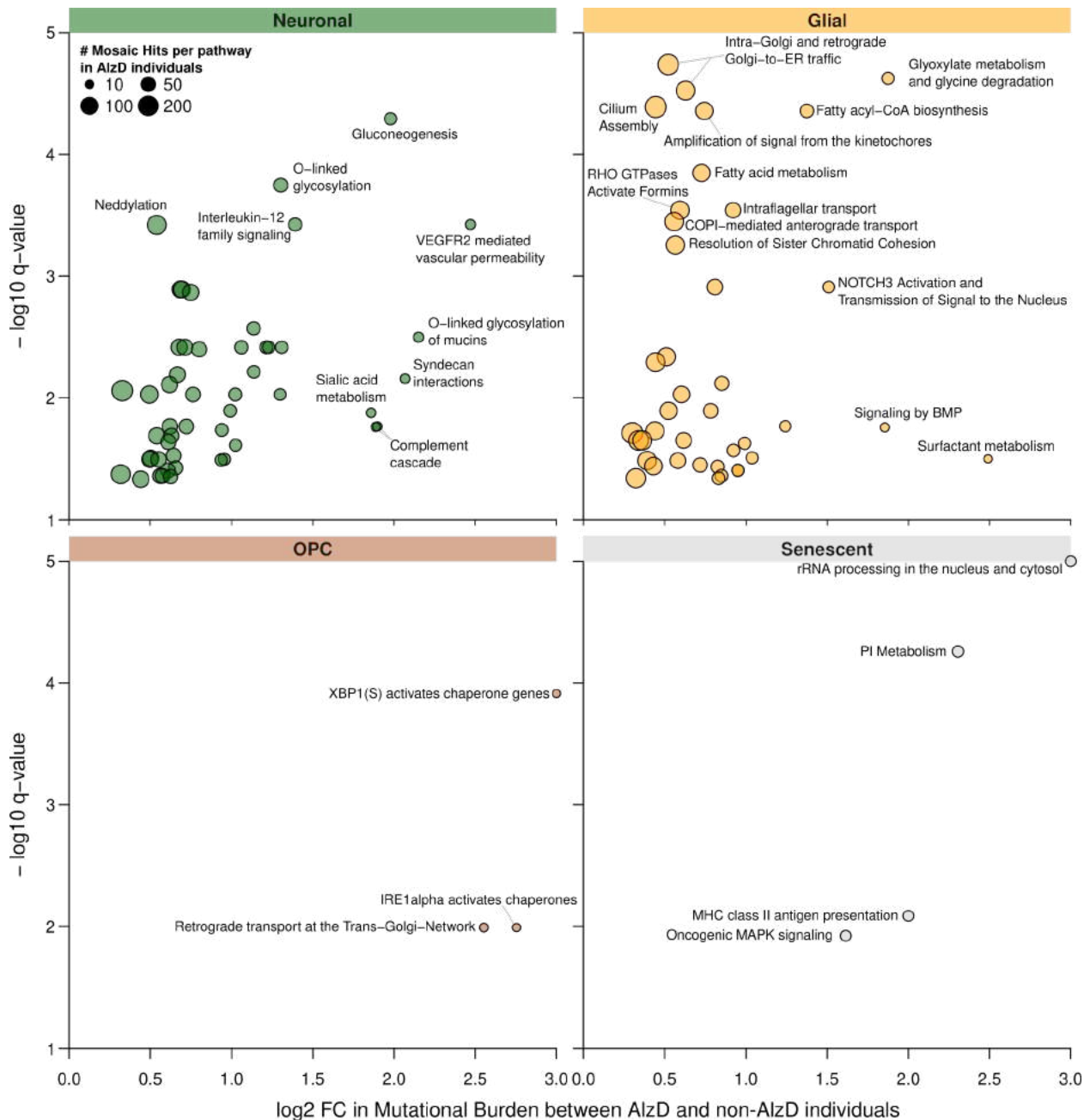
with their change and colored by their occurrence in cells from individuals with Alzheimer's dementia or no dementia. The t-SNE is colored according to coverage of the gene, with darker colors indicating higher coverage. All missense and LoF mutations used in burden calculations are shown, even in cell types where the gene is not enriched for burden.



**Supplementary Figure 5-S13. Neuronal and Glial genes enriched for AD burden.** Genes showing increased mosaic mutational burden (missense and loss-of-function hits, binomial test and FDR correction,  $q < 0.05$  shown) in Alzheimer's dementia versus non-dementia individuals in aggregated glial cells and neuronal cells. Heatmaps show number of mutations (left) overall in AD and non-AD individuals, (center) in each non-AD individual, and (right) in each AD individual in the cohort. Stars indicate genes not highlighted in a single-cell type analysis.

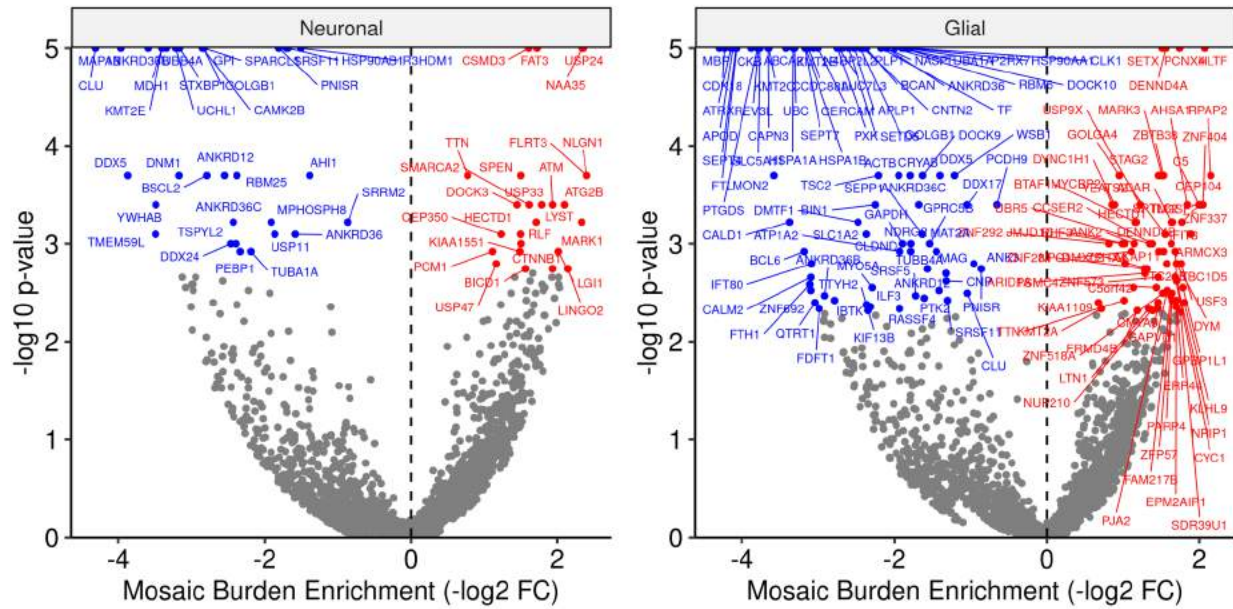


**Supplementary Figure 5-S14. Genes with significantly higher burden in AD in neuronal and glial cells.** Mosaic hits in each over-burdened gene in neuronal and glial cells are labeled with their change and colored by their occurrence in cells from individuals with Alzheimer’s dementia or no dementia. The t-SNE is colored according to coverage of the gene, with darker colors indicating higher coverage. All missense and LoF mutations used in burden calculations are shown, even in cell types where the gene is not enriched for burden.

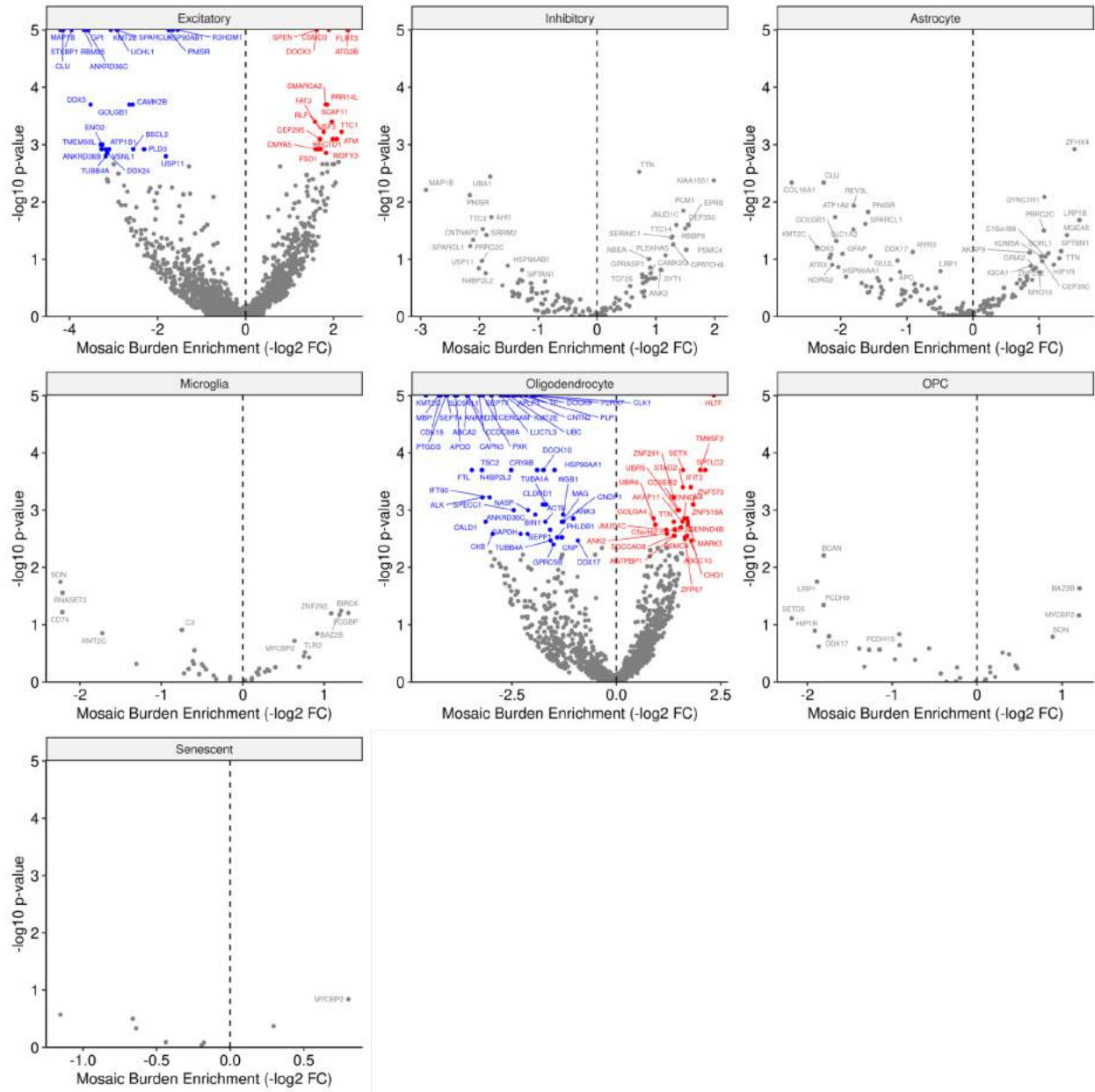


**Supplementary Figure 5-S15. Pathway-level enrichments in disease.** Pathways with increased mosaic mutational burden in AD versus non-AD individuals in either neuronal cells, glial cells, OPCs, or senescent cells (missense and loss-of-function hits, log<sub>2</sub> fold-change in burden against -log<sub>10</sub>q-value, size corresponds to number of hits in the pathway). Pathways with q < 0.001 or log<sub>2</sub>fc > 1.5 are labeled. Only pathways withstanding FDR correction with q < 0.05 for a binomial test against burden in non-AD individuals are shown (see Methods).





**Supplementary Figure 5-S16.** Genes significantly depleted (blue) or enriched (red) for mosaic mutational burden in a disease-agnostic manner in neuronal cells or glial cells (missense and loss-of-function hits, log<sub>2</sub> fold-change in burden against -log<sub>10</sub>p-value; genes with q < 0.05 are labeled).



**Supplementary Figure 5-S17.** Genes significantly depleted (blue) or enriched (red) for mosaic mutational burden in a disease-agnostic manner in each of the captured cell types (missense and loss-of-function hits,  $\log_2$  fold-change in burden against  $-\log_{10} p\text{-value}$ ; genes with  $q < 0.05$  are labeled in red/blue; genes with nominal  $p < 0.1$  are labeled in gray).

## 6. Temporal modeling of Alzheimer's disease progression

Parts of this chapter have been presented/published as:

### **Neurons burdened by DNA double strand breaks incite microglia activation through antiviral-like signaling in neurodegeneration**

*(in review)*

Gwyneth Welch<sup>1,2</sup>, Carles Boix<sup>3</sup>, Eloi Schmauch<sup>3,4,5</sup>, Jose Davila-Velderrain<sup>3</sup>, Matheus B. Victor<sup>1,2</sup>, Vishnu Dileep<sup>1,2</sup>, Lorenzo Bozzelli<sup>1,2</sup>, Qiao Su<sup>6</sup>, Jemmie Cheng<sup>1,2</sup>, Audrey Lee<sup>1,2</sup>, Noelle Leary<sup>1,2</sup>, Andreas Pfenning<sup>6</sup>, Manolis Kellis<sup>3,4</sup>, and Li-Huei Tsai<sup>1,2,4\*</sup>

1. Picower Institute for Learning and Memory, Massachusetts Institute of Technology, Cambridge, MA, USA. 2. Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA, USA. 3. Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA, USA. 4. Broad Institute of Harvard and MIT, Cambridge, MA, USA. 5. A. I. Virtanen Institute for Molecular Sciences, University of Eastern Finland, Kuopio, Finland. 6. Carnegie Mellon University Departments of Computational Biology and Biology and Neuroscience Institute, Pittsburgh, PA, USA.

\*Correspondence to: [lhtsai@mit.edu](mailto:lhtsai@mit.edu)

### **Single-cell multi-region dissection of Alzheimer's disease progression reveals selectively vulnerable neurons in Alzheimer's disease**

*(in preparation)*

Hansruedi Mathys<sup>1-3\*</sup>, Carles Boix<sup>4-6\*</sup>, Leyla Akay<sup>1,2</sup>, Ayesha Ng<sup>1,2</sup>, Xueqiao Jiang<sup>1,2</sup>, Kyriakista Galani<sup>4,5</sup>, Julio Mantero<sup>4,5</sup>, Neil Band<sup>4,5,7</sup>, David A. Bennett<sup>8</sup>, Li-Huei Tsai<sup>1,2,5+</sup>, Manolis Kellis<sup>4,5+</sup>

<sup>1</sup>Picower Institute for Learning and Memory, MIT <sup>2</sup>Department of Brain and Cognitive Sciences, MIT <sup>3</sup>University of Pittsburgh <sup>4</sup>Computer Science and Artificial Intelligence Laboratory, MIT <sup>5</sup>Broad Institute of MIT and Harvard <sup>6</sup>Computational and Systems Biology Program <sup>7</sup>University of Oxford <sup>8</sup>Rush Alzheimer's Disease Center. \*Equal contributors. \*Correspondence to:

[manoli@mit.edu](mailto:manoli@mit.edu) and [lhtsai@mit.edu](mailto:lhtsai@mit.edu)

### **Disentangling axes of single-cell variation to build disease-driven pseudotime trajectories**

*(in preparation)*

Carles Boix<sup>1-3\*</sup>, Neil Band<sup>1,2,4\*</sup>, ..., David A. Bennett<sup>5</sup>, Li-Huei Tsai<sup>2,6,7</sup>, Manolis Kellis<sup>1,2+</sup>

<sup>1</sup>Computer Science and Artificial Intelligence Laboratory, MIT <sup>2</sup>Broad Institute of MIT and Harvard <sup>3</sup>Computational and Systems Biology Program <sup>4</sup>University of Oxford <sup>5</sup>Rush Alzheimer's Disease Center <sup>6</sup>Picower Institute for Learning and Memory, MIT

<sup>7</sup>Department of Brain and Cognitive Sciences, MIT. \*Equal contributors. \*Correspondence to: [manoli@mit.edu](mailto:manoli@mit.edu)

## 6.1 Introduction

In our cross-brain region atlas of AD (discussed in Chapters 3 and 4), we identified specific characteristics of AD progression across regions at the level of excitatory neuron subtypes and brain regions and layers. Given the consistent brain region progression of Alzheimer's disease and the massive excitatory neuron subtype diversity across these regions identified in Chapter 3, we searched for and identified specific excitatory subtypes affected and depleted in Alzheimer's disease (Chapter 4). These showed a specific progression of depletion, with subtypes in multiple layers and structures of the entorhinal cortex and hippocampus affected first, followed by degeneration of specific neocortical layers (L2 - PFC/MT). These changes were strongly linked to increased cognitive impairment and pathology burden in AD. Using these subtypes, we identified consistent transcriptomic changes in AD, and were able to separate changes into resilience-associated and degeneration-associated transcriptomic changes and genes.

These last results hinted at a diversity of neuronal states within AD that are orthogonal to excitatory neuron identities. Given these results, we are interested in better mapping potential trajectories of neuronal vulnerability at the single-cell level, understanding at what resolution these changes happen and in which order, and whether there are multiple AD-related states and subclasses of neurons in AD or in other neurodegenerative disease. Understanding neuronal trajectories of depletion can inform us as to how specific neuronal subtypes are vulnerable to AD - whether they are further along AD progressions or whether they are less tolerant to being in late AD-related stages. In addition, we would like to ask whether cells from specific individuals are biased to coming from certain parts of these trajectories and whether these attributes are linked to individuals' cognitive decline and regional AD progression.

Developing these models and maps of neuronal trajectories can allow us to track states before neuronal degeneration, prioritize targets and timepoints best suited to manipulation for arresting or treating neurodegeneration, and identify resilient states and trajectories towards these states that could be manipulated to alleviate and reverse cognitive decline.

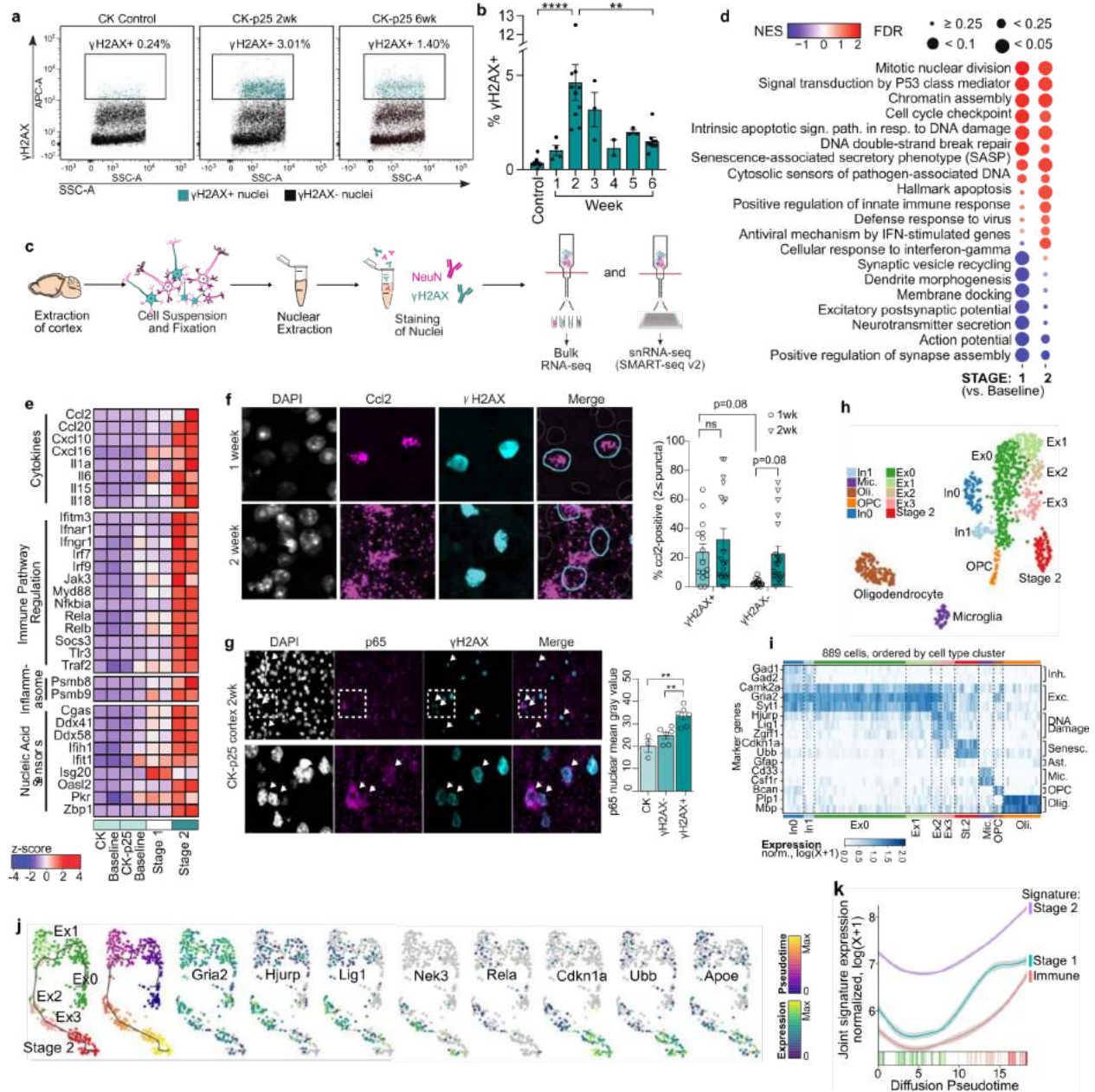
## 6.2 Temporal progression of neurons in the CK-p25 mouse model of AD towards senescence

**Summary:** DNA double strand breaks (DSBs) are linked to aging, neurodegeneration, and senescence. However, the role played by neurons burdened with DSBs in disease-associated neuroinflammation is not well understood. Here, we isolate neurons harboring DSBs from the CK-p25 mouse model of neurodegeneration through fluorescence-activated nuclei sorting (FANS), and characterize their transcriptomes using single-nucleus, bulk, and spatial sequencing techniques. We find that neurons harboring DSBs enter a late-stage DNA damage response marked by the activation of senescent and antiviral-like immune pathways. We map out a temporal trajectory for neurons entering senescent states in this AD mouse model that can inform our models of human neurodegeneration.

**DSB-bearing neurons in the CK-p25 mouse model of neurodegeneration:** We utilized the CK-p25 mouse model of inducible neurodegeneration to understand how DSB-bearing neurons contribute to disease development. In these mice, the *CamkII* promoter drives the expression of the neurotoxic protein fragment p25 through a doxycycline-off system<sup>293</sup>. P25 is the calpain-cleaved product of p35, an activator of cyclin dependent kinase 5 (Cdk5). Previously, we determined that the first pathologies observed in these mice are increased DSBs in neurons<sup>294</sup> and activation of microglia<sup>295</sup>. These pathologies occur 1-2 weeks after the onset of p25 expression when mice are taken off doxycycline (dox). Intracellular amyloid-beta accumulation occurs as early as 2-3 weeks after induction<sup>296</sup>. Neuronal loss, learning deficits, and tau hyperphosphorylation are also observed in the following 4-12 weeks<sup>293,296,297</sup>. These pathological



events are inducible and occur in a highly concerted and predictable manner<sup>62</sup>.



**Figure 6-1. Neurons marked by DNA DSBs activate inflammatory signaling at early stages of neurodegeneration.**

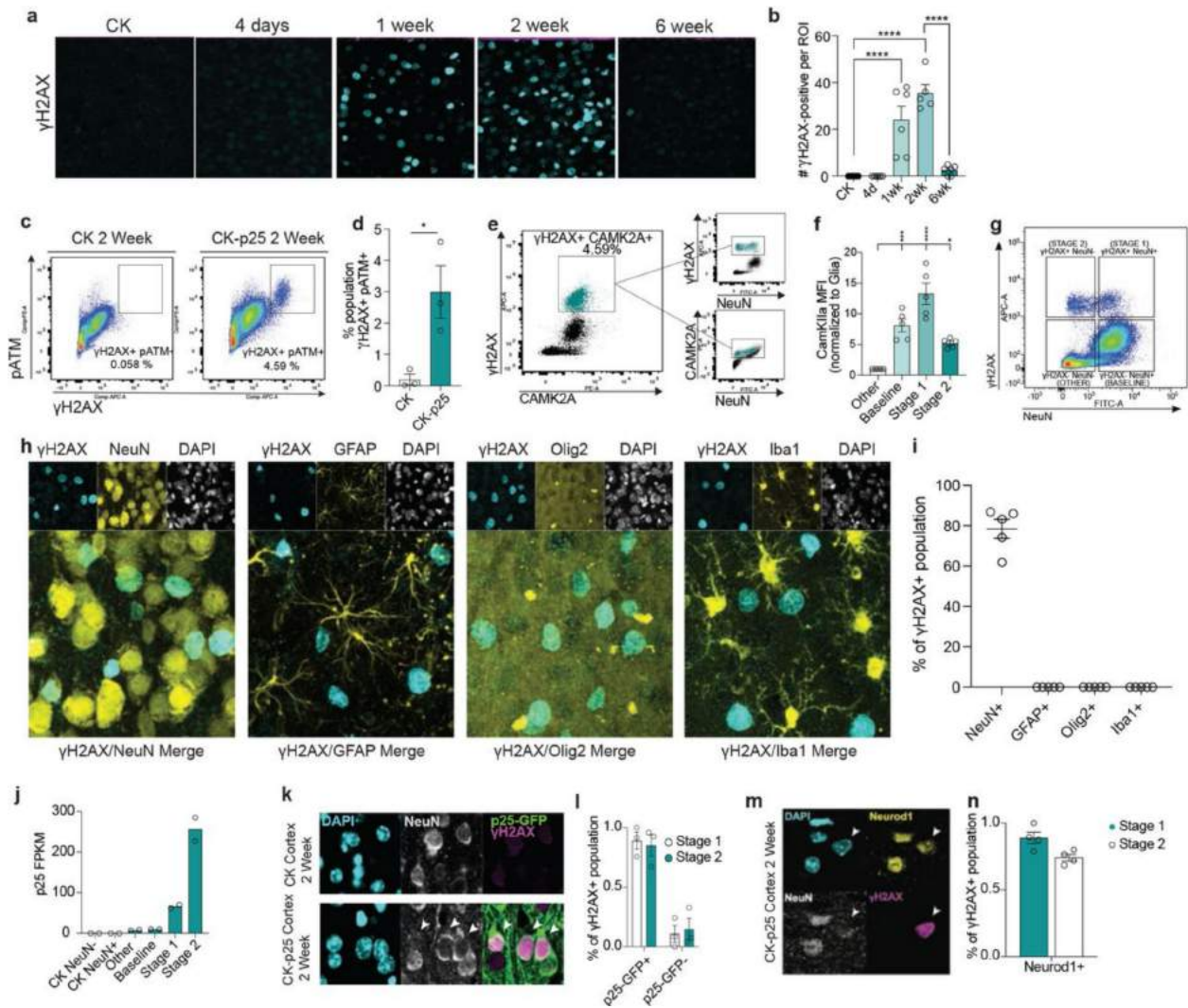
**a.** Flow cytometry dot plots of  $\gamma$ H2AX+ nuclei from CK and CK-p25 cortex.  $\gamma$ H2AX+ are highlighted in turquoise, and percent total population is indicated above the gating box. **b.** CK-p25 mice were taken off dox for 1 through 6 weeks, then cortical nuclei were extracted and stained for  $\gamma$ H2AX. Each datapoint represents percent  $\gamma$ H2AX+ nuclei from one mouse cortex. **c.** RNA seq workflow for the gated populations. Bulk RNA-seq: Cortices were extracted from CK and CK-p25 mice (n=2 per genotype for bulk RNA-seq, 2-week time point only). Tissue was homogenized, then nuclei were extracted and stained for NeuN and  $\gamma$ H2AX. 50,000 nuclei were sorted from each gated population. snRNA-seq: Cortices were extracted from CK and CK-p25 mice (n=3 per genotype $\times$ time point for SMART-seq, 1-week and 2-week timepoints). Tissue was homogenized, then nuclei were extracted and stained for NeuN and

$\gamma$ H2AX. 48 nuclei were sorted for each gated CK-p25 population per mouse, and 32 nuclei were sorted for each gated CK population per mouse. **d.** Differential gene ontology terms from gene set enrichment analysis (GSEA) of bulk RNAseq data. Left column depicts GSEA results from Stage 1 vs. CK-p25 Baseline contrast. Right column depicts GSEA results from Stage 2 vs. CK-p25 Baseline contrast. Color indicates normalized enrichment score (NES). Size indicates false discovery rate (FDR). **e.** Heatmap of differentially expressed genes belonging to inflammatory gene sets from bulk RNA-seq data. Each column represents one mouse. Gene sets are organized by biological function. **f.** (left): Representative images of Ccl2 RNAscope combined with  $\gamma$ H2AX immunofluorescent staining. Imaging was performed on 1wk and 2wk CK-p25 cortex. (right): Quantification of the number of  $\gamma$ H2AX+ and  $\gamma$ H2AX- cells with  $2 \leq$  Ccl2 puncta. Each datapoint represents average %Ccl2-positive cells in one image from one mouse. 4-3 images were taken per mouse. (Ccl2 1wk n=4, Ccl2 2wk n=5) **g.** Representative image of p65 immunostaining in CK-p25 cortex at two-week time point. White arrowheads indicate  $\gamma$ H2AX+ nuclei with nuclear p65 immunoreactivity. Quantification of mean p65 intensity for  $\gamma$ H2AX+ and  $\gamma$ H2AX- nuclei. Each data point represents the average p65 nuclear mean gray value of 20-60 nuclei from one mouse. **h.** UMAP of gated populations from CK and CK-p25 cortex at 1-week and 2-week timepoints. Colors indicate cell type annotation. **i.** Marker gene expression for each cell type cluster. Columns represent 985 individual cells ordered by cell type cluster. **j.** Trajectory analysis of Ex0,1,2,3, and Stage 2 neurons. Cells from each cluster (indicated by color annotation) are ordered across pseudotime (indicated by color gradient). **k.** Trajectory analysis of Ex0,1,2,3, and Stage 2 neurons. Smoothed gene signature expression across pseudotime. Stage 1=Significantly upregulated genes from Stage 1 vs. Baseline contrast. Stage 2= Significantly upregulated genes from Stage 2 vs. Baseline contrast. Immune= Genes belonging to immune gene ontologies that were significantly upregulated in Stage 2 neurons. padj<0.05, Log2 fold-change  $\geq$ 1.0. Error bars represent standard error of mean (S.E.M.); \*\*\*\*P<0.0001, \*\*\*P<0.001, \*\*P<0.01, \*P<0.05; One-way ANOVA with Tukey's test for multiple comparisons (**b,f**). Two-way ANOVA followed by Sidak's test for multiple comparisons (**g**). Data are pooled from 4 independent experiments (**b**). Data are representative of two independent experiments (**f,g**).

The phosphorylation of the histone variant H2A.X ( $\gamma$ H2AX) by ATM kinase occurs rapidly after DSB detection. This post-translational histone modification is essential for efficient DSB recognition and repair, and is a robust DSB biomarker<sup>298</sup>. Using flow cytometry, we identified a distinct population of  $\gamma$ H2AX-positive nuclei at the two-week time point in the CK-p25 cortex, but not in the CK control cortex (**Fig. 6-1a**). CK mice express the *CamkII*-driven tetracycline transactivator (tTA), but not p25. A timeline analysis was also performed to determine when  $\gamma$ H2AX-positive nuclei begin to accumulate in the CK-p25 cortex.  $\gamma$ H2AX-positive nuclei were detectable as early as one week after induction ( $1.038 \pm 0.2627\%$  population), peaked at two weeks ( $4.614 \pm 0.9416\%$  population), and gradually decreased thereafter (**Fig. 6-1b**). The significant reduction in  $\gamma$ H2AX-positive nuclei at six weeks corresponds with previous observations of neuronal loss in this model, suggesting that  $\gamma$ H2AX-positive cells degenerate by six weeks<sup>293,296,297</sup>. We were able to observe similar  $\gamma$ H2AX-positive population dynamics using immunofluorescence microscopy (**Fig. 6-2a,b**). We did not detect  $\gamma$ H2AX-positive nuclei after only four days off dox, suggesting that one week is the earliest time point at which this population appears. To validate that DNA damage response pathways were active in these nuclei, we confirmed co-immunoreactivity for phosphorylated ATM kinase (pATM) (**Fig. 6-2c,d**).

We found that all  $\gamma$ H2AX-positive nuclei stained positively for CamkIIa, a marker of excitatory neurons (**Fig. 6-2e,f**). Interestingly, at the two week time point, only  $57.37 \pm 8.822\%$  of

$\gamma$ H2AX-positive nuclei also stained positive for NeuN, a general marker for forebrain neurons (Fig. 6-2e,g). Degradation of neuron cell type identity is observed within aging and Alzheimer's disease patients<sup>299,300</sup>, and loss of NeuN expression has been proposed to be an indicator of declining neuronal health<sup>301</sup>. Indeed, reduced NeuN expression or reduced NeuN immunoreactivity is an established feature of neuronal damage<sup>302–305</sup>. Furthermore,  $\gamma$ H2AX-positive nuclei did not overlap with markers of microglia (Iba1), astrocytes (GFAP), or oligodendrocytes and oligodendrocyte precursor cells (Olig2), indicating that they were not of glial origin (Fig. 6-2h,i). Thus, we classified nuclei based on their NeuN and  $\gamma$ H2AX immunoreactivity. This resulted in four distinct populations:  $\gamma$ H2AX-negative NeuN-positive (“Baseline” neurons),  $\gamma$ H2AX-positive NeuN-positive (Stage 1 neurons),  $\gamma$ H2AX-positive NeuN-negative (Stage 2 neurons), and  $\gamma$ H2AX-negative NeuN-negative (putative non-neuronal cells, which we refer to as “Other”) (Fig. 6-2g). We used fluorescence activated nuclei sorting (FANS) followed by bulk RNA sequencing to transcriptionally profile each population (Fig. 6-1c). All nuclei were collected for sequencing at the two-week time point, which was when we observed the peak density of  $\gamma$ H2AX-positive nuclei.



**Figure 6-2.** a. Representative images of  $\gamma$ H2AX immunostaining in CK and CK-p25 cortex over a 6-week timeline analysis. b. Average number of  $\gamma$ H2AX+ nuclei quantified per image at each timepoint. Each data point represents one mouse. c. Flow cytometry dot plot of  $\gamma$ H2AX+ pATM+

nuclei in CK and CK-p25 cortex. **d.** Quantification of  $\gamma$ H2AX+ pATM+ nuclei in CK and CK-p25 cortex. Each data point represents percent  $\gamma$ H2AX+ pATM+ nuclei for one mouse. **e.** Flow cytometry dot plot of  $\gamma$ H2AX, Camk2a, and NeuN immunoreactivity in nuclei from CK-p25 cortex.  $\gamma$ H2AX+ Camk2a+ nuclei are gated in the left graph. This gated population is highlighted in turquoise in the two graphs to the right. **f.** Camk2a median fluorescent intensity (MFI) is quantified for Glia, Baseline, Stage 1, and Stage 2 populations. MFI values are normalized to the Glia population. **g.** Representative dot plot of  $\gamma$ H2AX and NeuN immunoreactivity in 2-week induced CK-p25 cortex. **h.** Representative images of cell type and  $\gamma$ H2AX immunostaining in the 2wk-induced CK-p25 cortex. Celltype markers from left to right: neurons (NeuN), astrocytes (GFAP), oligodendrocytes and oligodendrocyte precursor cells (Olig2), microglia (Iba1). **i.** Quantification of percent of  $\gamma$ H2AX-positive nuclei overlapping with each cell type stain. Each dot represents one mouse. 50-100  $\gamma$ H2AX-positive nuclei were analyzed for each mouse and cell type marker. **j.** Expression p25 in CK and CK-p25 gated populations. One datapoint represents fragments per kilobase per million (FPKM) from a gated population from one mouse. **k.** Representative images of GFP,  $\gamma$ H2AX, and NeuN immunostaining from CK and CK-p25 cortex at the 2wk time point. Arrowheads indicate  $\gamma$ H2AX+ nuclei that express p25-GFP regardless of NeuN immunoreactivity. **l.** Quantification of GFP expression across Stage 1 and Stage 2 populations. Each data point represents percent  $\gamma$ H2AX+ population from Stage 1 or Stage 2 from one CK-p25 mouse. **m.** Representative images of Neurod1,  $\gamma$ H2AX, and NeuN from CK-p25 cortex at two-week timepoint. White arrowhead indicates a  $\gamma$ H2AX+ nucleus with Neurod1 immunoreactivity but not NeuN immunoreactivity. **n.** Quantification of Neurod1 expression across Stage 1 and Stage 2 populations. Each data point represents percent  $\gamma$ H2AX+ population from Stage 1 or Stage 2 from one CK-p25 mouse. Error bars represent standard error of mean (S.E.M.); \*\*\*\*P<0.0001, \*\*\*P<0.001, \*P<0.05, n.s. not significant; One-way ANOVA with Tukey's test for multiple comparisons (b,f). Student's T-test (d). Data are averages of 4 images per mouse (b). Data are representative of at least 2 independent experiments (d). Data are pooled from 2 independent experiments (f).

Both Stage 1 and Stage 2 populations expressed more p25 transgene compared to other populations (**Fig. 6-2j**). CK-p25 mice exclusively express p25 in excitatory forebrain neurons, as it relies on the expression of tTA driven by the CamkII promoter. Because the p25 transgene is fused to GFP, we performed staining to validate this finding. We found that  $85.74 \pm 0.055$  % of  $\gamma$ H2AX-positive cells, regardless of NeuN immunoreactivity, were also GFP-positive (**Fig. 6-2k,l**). To verify Stage 2 neuronal identity independent of transgene expression, we also stained for Neurod1, a neuronal transcription factor. We found that  $81.96 \pm 0.046$  % of  $\gamma$ H2AX-positive cells stained positively for Neurod1, regardless of NeuN immunoreactivity (**Fig. 6-2m,n**). Together, our neuron marker and transgene expression analysis indicated that both Stage 1 and Stage 2 populations were likely excitatory neurons engaged in a DSB response.

**Differential gene expression in DSB-bearing neurons:** We performed two differential expression analyses to characterize gene expression changes in these populations: Stage 1 vs. CK-p25 Baseline, and Stage 2 vs. CK-p25 Baseline. Large transcriptional changes were observed in both contrasts. 3,031 upregulated and 717 downregulated transcripts were identified in Stage 1, and 5,055 upregulated and 3,792 downregulated transcripts were identified in Stage 2 ( $\log_2$  fold change  $\geq 1$ , adjusted p-value < 0.05). Gene set enrichment analysis (GSEA)<sup>306</sup> revealed that both Stage 1 and Stage 2 neurons displayed a significant enrichment for genes implicated in DSB repair, apoptotic signaling, and cell cycle, and a significant reduction in synaptic function processes (**Fig. 6-1d**). These pathways were previously found deregulated in both the CK-p25 mouse and AD human brain tissue<sup>62,294,295,307,308</sup>.



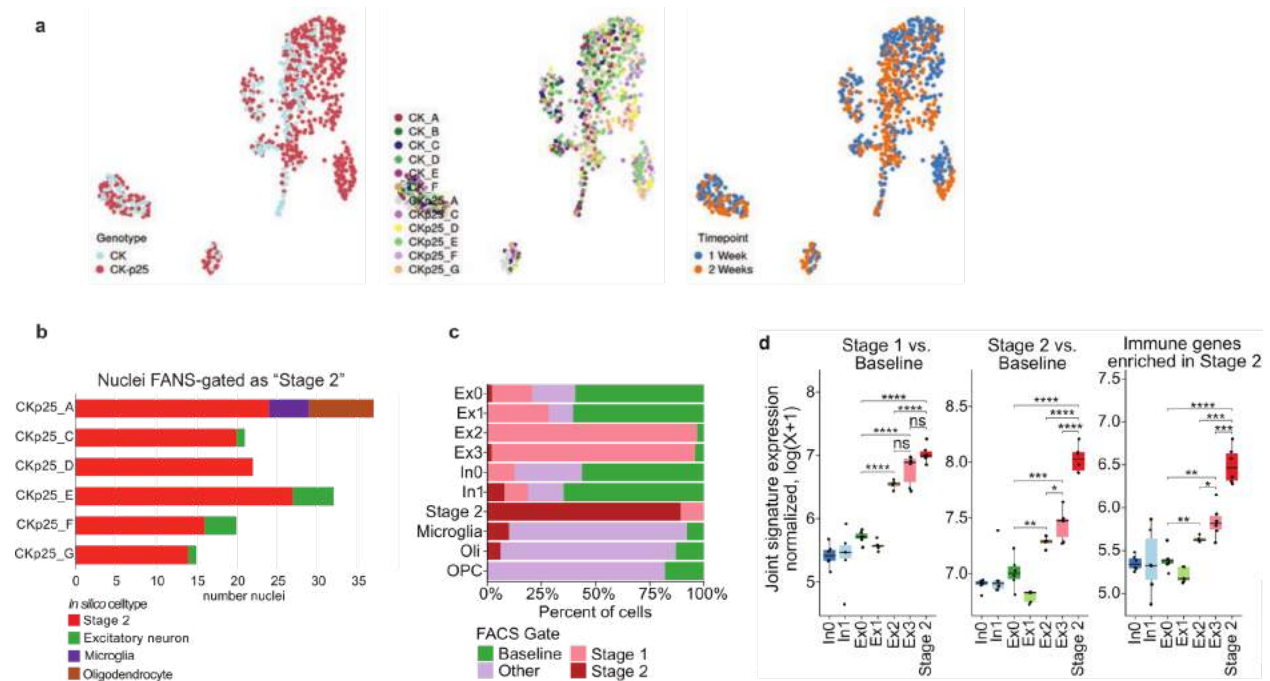
Remarkably, we also found that a number of innate immune pathways were enriched in Stage 1 and Stage 2 neurons. This included gene ontology terms ‘Senescence-Associated Secretory Phenotype (SASP)’, ‘cytosolic sensors of pathogen-associated DNA’, and ‘positive regulation of innate immune response’ (**Fig. 6-1d**). Upon closer inspection, we observed an enrichment of genes linked to DSB-mediated activation of SASP signaling, particularly in Stage 2 neurons. This included nucleic acid sensors such as *Cgas* and *Zbp1*, NFkB subcomponents *Rela* and *Relb*, SASP factors *Il6*, *Il15*, *Ccl2*, *Cxcl10*, and interferon stimulated genes *Isg15* and *Ifitm3* (**Fig. 6-1e**). Notably, a number of these genes are expressed in neurons following viral infection<sup>309,310</sup>. These genes are also classic hallmarks of senescence, suggesting that DSB accumulation elicits senescent and anti-viral-like responses in neurons.

Because the CK-p25 model is characterized by the development of type-I interferon-reactive microglia specifically by two weeks induction<sup>295</sup>, we wanted to determine if DSB-bearing neurons express cytokines before microglia. To do this, we performed RNAscope multiplexed fluorescent in-situ hybridization in the CK-p25 cortex after one and two weeks of induction. We focused our profiling on *Cxcl10* and *Ccl2* because these pro-inflammatory chemotactic molecules were highly expressed in Stage 2 neurons and are known to be secreted by neurons upon viral infection<sup>309,311–314</sup>. We also performed RNAscope for the excitatory neuron marker *Camk2a* to again confirm if  $\gamma$ H2AX-positive nuclei were of neuronal origin. In agreement with our RNAseq and flow cytometry data,  $98.81 \pm 1.19\%$  of  $\gamma$ H2AX-positive nuclei were also *Camk2a*-positive. This high positivity rate reveals that  $\gamma$ H2AX-positive nuclei are nearly exclusively neurons. RNAscope analysis also revealed that DSB-bearing neurons were the only cells to express *Cxcl10* and *Ccl2* at the one-week time point. In contrast, *Cxcl10* and *Ccl2* gene expression was highly enriched in other cell types at the two-week time point (**Fig. 6-1f**). This indicated that expression of these chemokines in  $\gamma$ H2AX-positive neurons precedes their expression in glial cells, and that cytokine secretion from DSB-bearing neurons may be an early mechanism of glial cell recruitment and activation in the CK-p25 brain.

To identify the master regulators of DSB-associated neuronal immune signaling, we performed transcription factor enrichment analysis using Enrichr<sup>315</sup>. A subset of immune pathway genes was extracted from the significantly upregulated genes in Stage 2 neurons for analysis. The overlap between this immune gene module and transcription factor target genes was then calculated. Multiple subunits of the NFkB complex were consistently enriched across Enrichr transcription factor libraries, including *Rela*, *Nfkb1*, and *Relb* (Supplementary Table 1). Notably, the NFkB transcription factor plays a well-established role in SASP activation and the DSB response<sup>316–318</sup>. We chose to focus on *Rela*, also known as p65, because it is a core member of the canonical NFkB complex, and it was most frequently enriched in the Enrichr analysis. We stained for p65 in two-week CK and CK-p25 cortices. NFkB is normally sequestered in the cytosol, but translocates to the nucleus to form an active complex upon cellular insult. Nuclear p65 intensity was significantly higher in  $\gamma$ H2AX-positive neurons compared to other  $\gamma$ H2AX-negative cells (**Fig. 1g**), supporting evidence of increased NFkB transcriptional activity.

**Single cell RNA-sequencing of DSB-bearing neurons:** Notably, inflammatory gene expression was lower in Stage 1 neurons compared to Stage 2 neurons (**Fig. 6-1e**), and more Stage 1 neurons were identified at the one-week time point compared to Stage 2 neurons. This suggested that the Stage 2 population may develop after Stage 1, which aligned with our initial prediction when defining these distinct cellular states based on NeuN-reactivity. To better understand the relationship between Stage 1 and Stage 2 neurons, we performed single nucleus RNA-sequencing on each FANS-gated population at both one and two-week timepoints (**Fig. 6-1e**). A total of 1,357 single nucleus libraries were prepared using SMARTseq2 chemistry. Following quality control measures (see methods), 889 libraries remained for downstream

analysis (**Fig. 6-1h, 6-3a**). Cells were classified into major cell type clusters based on marker gene expression, resulting in the identification of 521 excitatory neurons, 71 inhibitory neurons, 131 oligodendrocytes, 33 oligodendrocyte precursor cells (OPCs), and 50 microglia. We did not detect any astrocytes.

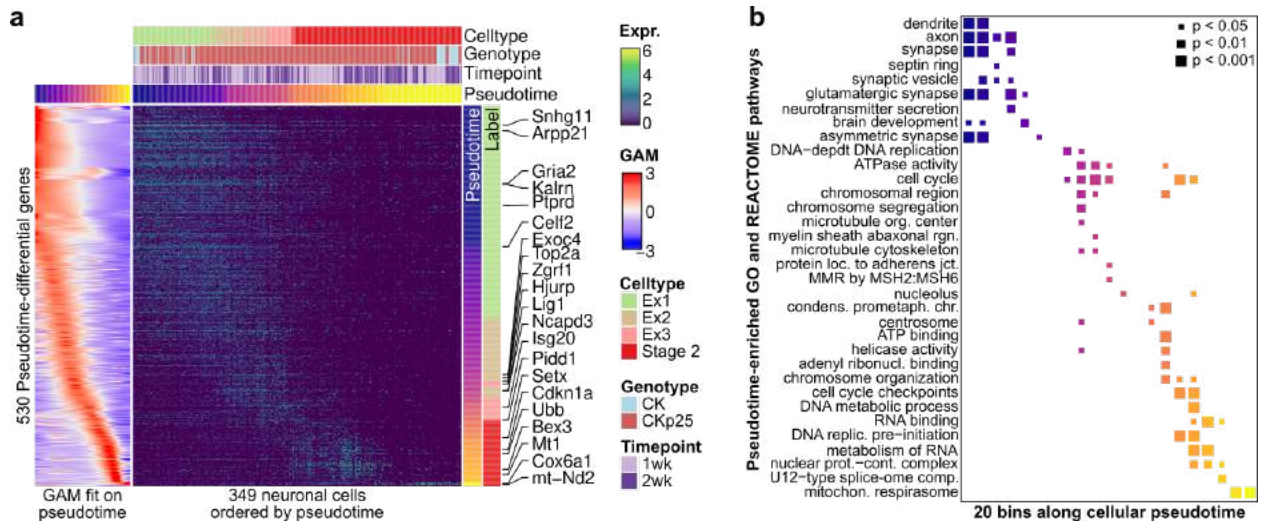


**Figure 6-3.** **a.** UMAPs labeled by genotype, mouse id, and time point. **b.** Distribution of Stage 2-gated nuclei in in silico cell type clusters, stratified by biological replicate. **c.** Percent FANS label distribution across scRNA cell type clusters. Bar colors refer to the FANS gate label. **d.** Bulk RNA-seq gene signature enrichment in neuron cell type clusters. Stage 1 signature enrichment (left), Stage 2 signature enrichment (middle), and immune genes from Stage 2 enrichment (right). Each data point indicates average gene expression across all cells from one mouse. Only CK-p25 mice were used for this analysis. \*\*\*\* $P < 0.0001$ , \*\*\* $P < 0.001$ , \*\* $P < 0.01$ , \* $P < 0.05$ , n.s. not significant; One-way ANOVA with Tukey's test for multiple comparisons (**d**).

Remarkably, the majority of Stage 2-gated nuclei formed their own cell type cluster (**Fig. 6-1h**). While we did identify 13 Stage 2-gated nuclei in microglia and oligodendrocyte clusters, all of the nuclei in question came from one CK-p25 mouse (**Fig. 6-3b,c**). None of the other five CK-p25 mice had Stage 2-gated nuclei fall into glial clusters. Combined with our previous cell type immunostaining analysis (**Fig. 6-2h,i**), we concluded that the immune signature identified in the Stage 2 population was not likely to be driven by contamination from microglia or oligodendrocytes. Compared to other neuronal clusters, the *in silico* Stage 2 cluster was significantly enriched for the FANS-gated Stage 1 and Stage 2 gene signatures (**Fig. 6-3d**). Interestingly, the Stage 2 cluster expressed only moderate levels of the excitatory neuron cluster markers *Camk2a*, *Gria1*, and *Syt2*, and lacked marker genes of other canonical cell types such as inhibitory neurons (*Gad1*, *Gad2*), astrocytes (*GFAP*), microglia (*Cd33*, *Csf1r*), oligodendrocytes (*Plp1*, *Mbp*), and OPCs (*Bcan*). Instead, the Stage 2 cluster expressed marker

genes indicative of senescence, including *Cdkn1a*, and *Ubb* (**Fig. 6-1i**). These distinctive cell type markers further suggested to us that the Stage 2 cluster represented a population of DSB-bearing neurons engaged in a senescence-like inflammatory response.

To further examine Stage 1 and Stage 2 cell type heterogeneity, we sub-clustered all neuronal cells. This resulted in the identification of four excitatory neuron subclusters (Ex0-3), and two inhibitory neuron subclusters (In0, In1) (**Fig. 6-1h**). Stage 1 cells were enriched in subclusters Ex2 and Ex3 compared to Ex0 and Ex1 (**Fig. 6-3c**). Subclusters Ex2 and Ex3 were also marked by DNA repair-associated genes, such as *Hjurp*, *Lig1*, and *Zgrf1*, further suggesting that Ex2 and Ex3 cells were engaged in a DSB repair response (**Fig. 6-1i**).



**Figure 6-4. a.** Heatmap of smoothed generalized additive model fits and expression of differentially expressed genes across cells (columns) as ordered by diffusion pseudotime. Genes (rows) are ordered by their maximal point along the pseudotime and annotated by their maximal cell type for pseudotime and cell label. **b.** Pseudotime enriched GO and REACTOME pathways for 20 bins of genes along cellular pseudotime.

Interestingly, these DNA repair-associated genes were not expressed in the Stage 2 cluster. Combined with the enrichment for the classical senescence marker *Cdkn1a*, this suggested that Stage 2 cells may be engaged in a later-stage response to DSBs. This motivated us to perform trajectory analysis. The single cell analysis package Monocle3<sup>319</sup> was used to order cells from subclusters Ex0,1,2,3 and Stage 2 along a pseudotime, which was generated through a learned trajectory of subcluster gene expression differences (**Fig. 6-1j**, **Fig 6-4a,b**). Both CK and CK-p25 cells were used for this analysis. Next, we identified genes that changed as a function of pseudotime (**Fig 6-4a**). We observed sequential enrichment of Stage 1 and Stage 2 gene signatures along pseudotime (**Fig. 6-1k**). The expression of Stage 2 signature genes specifically involved in immune response (Immune) were also enriched at the end of pseudotime (**Fig. 6-1k**). We found that genes associated with neuronal identity such as *Gria2* and *Kalrn* were most highly expressed at the beginning of the trajectory (**Fig. 6-1j**). Meanwhile, genes associated with DNA repair like *Hjurp*, *Lig1*, and *Nek3* peaked along the middle, and *Cdkn1a* expression peaked at the end of the trajectory. We also observed late-stage expression of *Apoe*, supporting previous observations that *Apoe* expression increases in neurons following injury<sup>320</sup>.

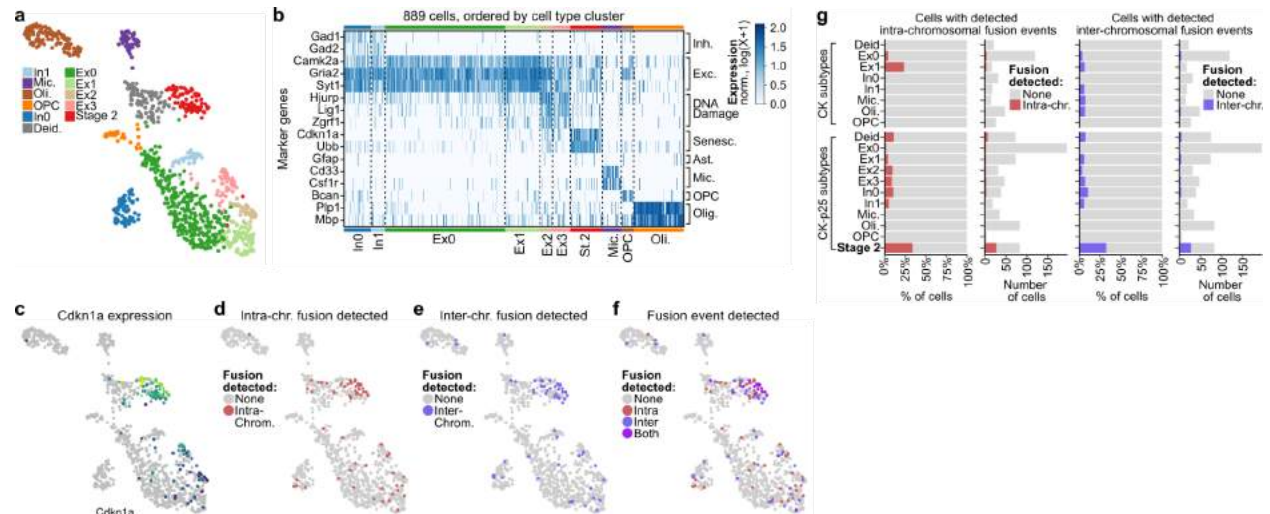




### 6.3 DNA damage in the CK-p25 mouse model

In section 6.2, we modeled the neuronal neurodegeneration and senescence trajectories for the CK-p25 mouse model of Alzheimer's disease. The mouse model, which contains a vector that expresses AD-implicated p25 (CDK5R1) under a neuronal promoter, leads to the accumulation of double stranded breaks (DSBs) in p25-expressing neurons followed by senescence and AD-like cognitive and neurodegenerative phenotypes. Modeling neuronal trajectories through this process in conjunction with DNA damage provides a window into how DNA damage in AD may lead to neuronal senescence and death<sup>296</sup>.

We identified distinct populations of  $\gamma$ H2AX-positive nuclei (DSB-containing) at the two-week time point after p25 induction in the CK-p25 cortex, but not in the CK control cortex. We divided these cells into  $\gamma$ H2AX+ and NeuN+ (Stage 1) and  $\gamma$ H2AX+ but NeuN- (Stage 2) cells according to the presence of the DSB marker ( $\gamma$ H2AX) and the neuronal identity marker (NeuN), where Stage 2 cells represented late-stage damaged neurons (see section 6.2). To understand the transcriptional responses of these cells after p25 induction, we performed single-cell full transcript sequencing (Smart-seq2) on FACS sorted cells from these populations as well as non- $\gamma$ H2AX neurons and glia, obtaining 1,008 cells from 5 CK and 5 CK-p25 mice after QC (Fig. 6-5a). We identified three separate Stage 1 populations (Ex1, Ex2, and Ex3), whose marker genes were enriched for DNA-repair and chromatin-associated factors, including DNA ligase Lig1 and Holliday junction recognition protein Hjurp (Fig. 6-5b). In contrast, Stage 2 cells contained traditionally senescence-associated markers such as Cdkn1a and Ubb (Fig. 6-5b-c).



**Figure 6-5. AD mouse model shows increased gene fusions in Stage 2 cells.** **a**, UMAP of 1,008 Smart-seq2 cells from 5 CK and 5 CK-p25 mice, colored by cell subtype (Ex: Excitatory neurons, In: Inhibitory neurons, Mic: Microglia, Oli: Oligodendrocytes, OPC: Oligodendrocyte precursor cells, Deid: Deidentified cells). **b**, Expression of subtype marker genes and genes related to senescence and DNA damage for 889 cells with consistent cellular identity. **c**, Expression of Cdkn1a on the cell UMAP. **d-f**, UMAP of cells labeled by whether a fusion was detected in the single-nucleus transcriptomic data. Colors label either an intra-chromosomal event (red, **d**), an inter-chromosomal fusion event (blue, **e**), or both detected the same cell (purple, **f**). **g**, Barplots showing the percentage and number of cells of each subtype in either CK or CK-p25 mice with intra (left) or inter-chromosomal fusions (right).

In order to align these transcriptional events to the establishment and repair of DSBs, we performed gene fusion detection to search for sites with putative incorrectly resolved DSB repairs. We realigned all reads in each cell independently using the STAR aligner and the mm10 CTAT toolkit and prioritized *de novo* chimeric splice junctions for potential gene fusion events using STAR-Fusion with default parameters<sup>248,321,322</sup>. We filtered out potential fusion events that contained promiscuously aligning snoRNA Snhg11 and separated inter-chromosomal versus intra-chromosomal called gene fusion events. Stage 2 neurons, out of all CK-p25 neurons, were significantly enriched for both intra ( $p=4.6e-6$ , Fisher's exact test) and inter-chromosomal events ( $p=1e-6$ ) and very significantly enriched for the joint set of fusions ( $p=2.5e-11$ ) (**Fig. 6-5d-g**). In addition, cells often contained both inter and intra-chromosomal events ( $p=8e-8$ , Fisher's) (**Fig. 6-5f**).

However, we did not detect a significant enrichment of potential fusions in Stage 2 clusters (Ex1, Ex2, and Ex3), and we did not see a significant increase of fusions by week from p25 expression induction (1 or 2 weeks). Given that Stage 1 cells were also marked by  $\gamma$ H2AX and showed high expression of DNA repair genes, it is likely that these cells are at an intermediate stage with many active break sites and few resolved breaks. These results suggest that the Stage 2 cells are uniquely burdened by incorrectly resolved DNA breaks, detectable as gene fusion events from single-cell full transcriptome sequencing.

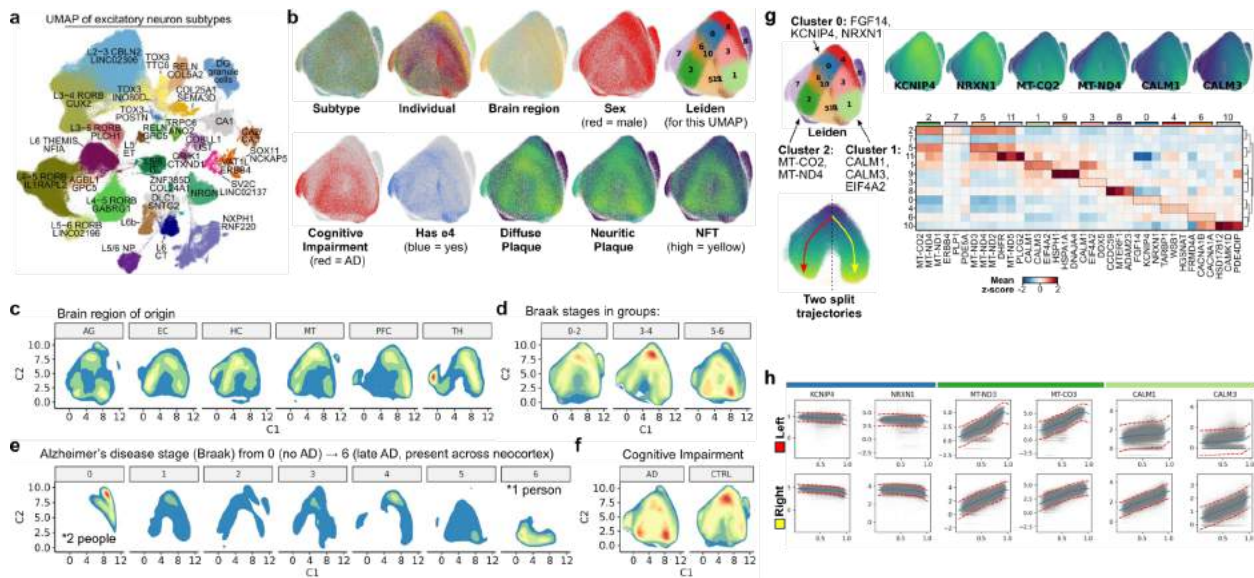
This work allows us to jointly model DNA damage and its transcriptional response in neurons on the path towards senescence. However, it is limited in two important ways that open up future research directions. First, we can only detect a fraction of the incorrectly resolved junctions - only if they rearrange expressed genes and are an incorrect resolution. Therefore, it is possible that the Stage 1 or 2 cells contain many more breaks than are resolved, or possibly that earlier DSBs can be correctly resolved while later DSBs overwhelm the cell and are incorrectly resolved. Second, it may be the case that some neurons are able to resolve DSBs and return to a baseline neuronal state instead of entering a senescent state. While we have unprecedented resolution in describing these dynamics given our single-cell, multi-stage and trajectory approach, it is hard to distinguish and thus model neurons that may be returning to a baseline state (Ex0). Therefore, future work to shed light on neurodegeneration should quantify the dynamics and transitions between neuronal states on the path to DNA damage-induced neuronal senescence.

#### **6.4 Disease pseudotime by regressing out core axes of non-disease variability**

Traditional pseudotime models for single-cell data have been applied to developmental datasets and certain disease-associated datasets to great success. The prototypical algorithm follows a set number of steps, which include (1) computing a cell-cell representation such as a UMAP or a k-NN cell-cell graph, (2) finding clusters of cells, (3) linking cell clusters, (4) draw a minimal-spanning path between the cell clusters, (5) reducing the path to a trajectory, (6) assigning each cell to the closest point in the path to give it a "pseudotime score". Examples of algorithms following these steps include Monocle<sup>319</sup>, Slingshot<sup>323</sup>, and PAGA<sup>324</sup>. The power of these methods lies in their ability to order cells across and within clusters, but they rely on the cell-cell graph being dominated by the attribute of interest, whether it be developmental trajectory or disease. These methods will fail in cases when batch-associated variation drives cell-cell clustering. Therefore, these algorithms are not directly applicable to the majority of disease datasets from neuronal data as neuronal cells cluster primarily by neuronal subtype, even if they share similar disease-associated changes (**Figure 6-6a**).

Here, we present a simple and intuitive approach to calculate unified pseudotime trajectories across subtype clusters within an atlas-scale single-cell dataset by performing batch correction on unwanted axes of cellular variability to promote disease-related variation in the dataset. To

remove subtype-associated transcriptional signatures, we processed single-cell RNA-seq data as in standard analyses (see Chapters 3-5), but applied ComBat<sup>325</sup> batch correction for cell subtypes after normalizing cells. In addition, we performed batch-balanced k-NN to obtain a balanced nearest neighbors graph across subtype batches after calculating the PCA of the ComBat-corrected matrix. This procedure yielded a cell-cell representation (**Figure 6-6b**) that did not cluster by subtype but did show individual-specific and AD-relevant structure, as cognitively impaired individuals and e4 carriers clustered towards the bottom of the UMAP. Despite never explicitly looking to cluster by disease, we also saw a top to bottom gradient for AD variables when plotting the density of cells (as contour maps) for individual-level cognitive impairment and Braak stage but not for sample-level brain region of origin (**Figure 6-6c-f**). In particular, the contour plots for Braak stage groups and cognitive impairment showed clear enrichment for non-AD individuals in leiden cluster 0 and progression towards two separate AD-associated subclusters (1 and 2) in late AD individuals.



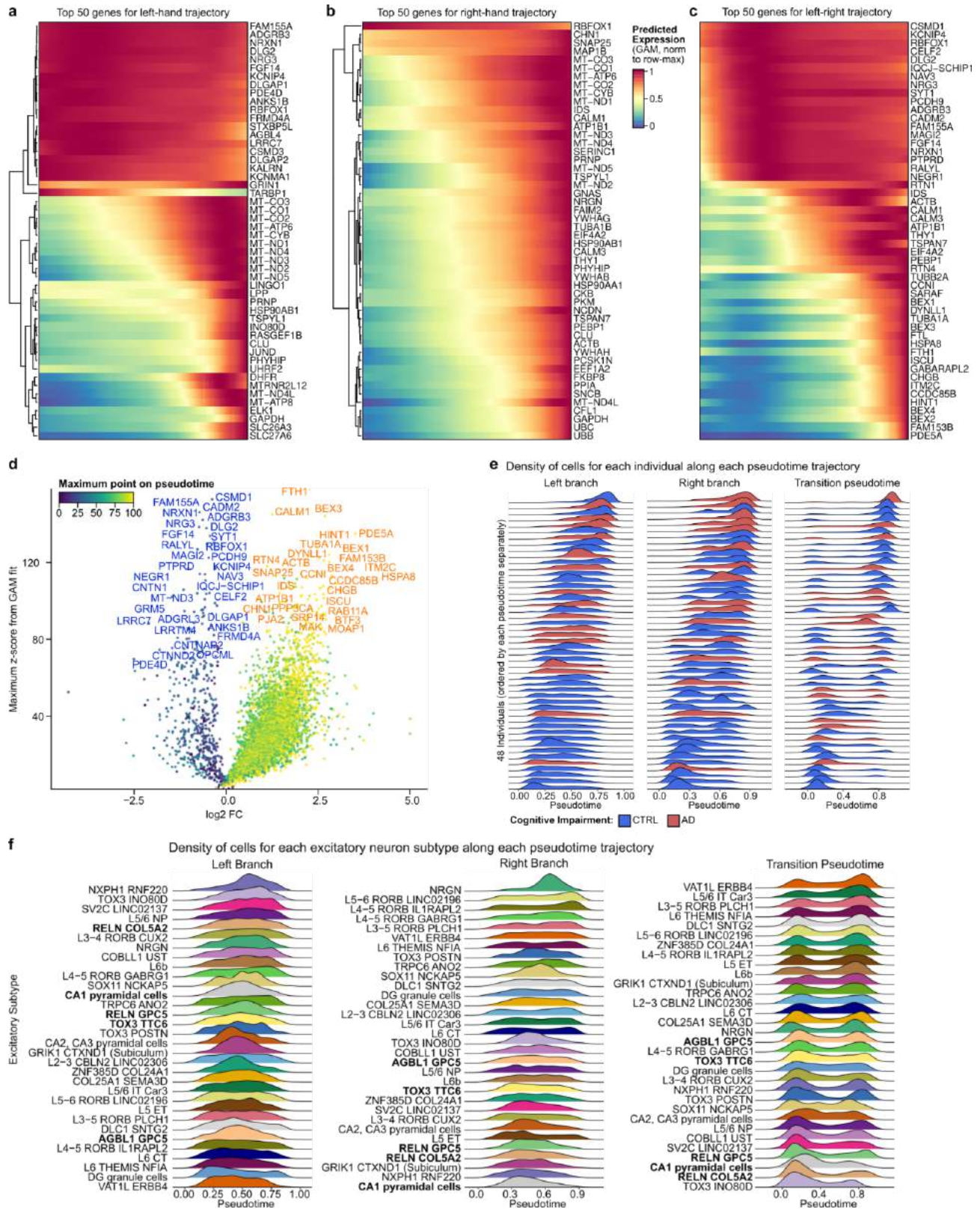
**Figure 6-6. A neuronal subtype-agnostic cell-cell graph and UMAP.** **a**, UMAP of excitatory neuron subtypes from the multi-region dataset, with ~435k cells across six brain regions and 48 individuals. **b**, Subtype-agnostic UMAP, with ~435k neurons. Each panel represents a separate metadata covariate. **c-f**, Contour density maps for cells within each metadata category for brain region of origin (**c**), braak stages (**d-e**), and cognitive impairment stages (**f**). **g**, Leiden clusters and split trajectory on the subtype-agnostic UMAP (left) and UMAPs and heatmap for leiden cluster marker genes on the subtype-agnostic UMAP (right). **h**, Scatterplots of marker gene expression along left-hand side (top) and right-hand side (bottom) trajectories. Expression is batch-corrected (log1p) and lines represent GAM fit with 95% confidence interval.

**Subtype-agnostic UMAP contains multiple AD-related states:** We clustered the subtype-agnostic UMAP and determined marker genes for each leiden cluster (t-tests, **Figure 6-6g**). The non-AD subcluster (0) was marked by core neuronal function genes including ion channel KCNIP4 and neural adhesion NRXN1, and the two AD-associated subclusters (1 and 2) represented separate states with different marker genes. The right-hand state (2) was marked by mitochondrial genes including MT-CO2 and MT-ND4 whereas the left-hand state was marked by calcium-signaling associated CALM1 and CALM3 and translation factor EIF4A2 (**Figure 6-6g**). We plotted the two selected marker genes for each major subcluster (0-2) along cells

from either side of a branching AD pseudotime ( $0 \rightarrow 2$  or  $0 \rightarrow 1$ ) derived from the subtype-agnostic UMAP (**Figure 6-6g-h**). Certain genes, such as CALM1/3, showed a trajectory of increasing expression in only one side of the split AD trajectory, suggesting that the two trajectories should be further analyzed to understand differences in the global neuronal trajectories in AD (**Figure 6-6h**).

**Building pseudotime trajectories towards AD:** Given our observations of two AD-related states and one non-AD state, we built two parallel pseudotime trajectories, between  $0 \rightarrow 2$  (right) and  $0 \rightarrow 1$  (left-hand side). To do this, we first removed all clusters biased towards one individual over others (>20% cells) and removed the intermediate cluster (5), which we analyzed separately below. We then fit generalized additive models (GAMs) of splines to each gene in each of the two trajectories, and computed p-values at their locations of maximal differential expression along the trajectory in order to rank genes' changes along each trajectory (**Figure 6-7a-b**). The vast majority of significantly differential genes along the trajectories showed end-to-end expression changes, with few intermediate or midpoint peaks in gene expression trajectories. The left-hand side trajectory in particular was marked by few gene expression changes and increased mitochondrial gene expression as well as up-regulation of fatty acid and chloride transporters SLC27A6 and SLC26A3 and down-regulation of HIV-linked RNA binding protein TARBP1 (**Figure 6-7a**). The right-hand side trajectory showed broader gene up-regulation, including translation factors (EEF1A2 and EIF4A2), chaperones (HSP90AB1, HSP90AA1, CLU), ubiquitin (UBB, UBC), and other aggregate-associated proteins (PRNP, SNCB) (**Figure 6-7b**). Given increased up-regulation of genes in this state, it may represent a strongly up-regulated state in response to Alzheimer's progression.



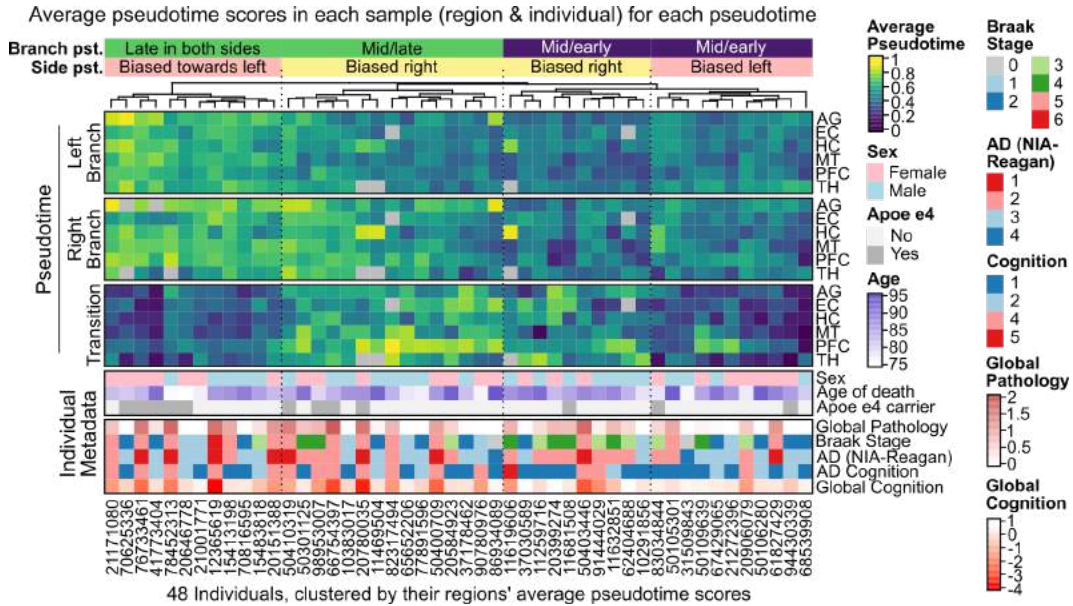


**Figure 6-7. Genes and cell scores for three pseudotime trajectories.** a, Heatmap of the GAM fits along the left-hand side ( $0 > 2$ ) pseudotime, showing the top 50 genes (rows, by GAM

p-value) along pseudotime, with the rows normalized to row-maximum values. **b**, Heatmap of top 50 GAM fits for the right-hand side ( $0 > 1$ ) pseudotime, as in **(a)**. **c**, Heatmap of top 50 GAM fits for the horizontal ( $2 > 1$ ) pseudotime across AD states, as in **(a)**. **d**, Volcano plot of  $\log_2FC$  vs. maximum z-score from GAM fit for gene in the horizontal pseudotime **(c)**. Top genes are annotated and colored by their point of maximum expression on the pseudotime (labels are blue - 2/right and orange - 1/left). **e**, Density of cells for each individual along each pseudotime trajectory, where rows are individuals, are ordered independently in each panel, and are colored by the individual's cognitive impairment status. **f**, Density of cells for each excitatory neuron subtype along each pseudotime trajectory. Rows are colored by subtype colors and ordered by average pseudotime scores of points in the subtype.

**State-change genes from AD state traversal:** Just as each trajectory between non-AD and AD states (1 and 2) described possible the changes that neurons may undergo in disease progression, the trajectory from clusters 1 to 2 or its inverse revealed genes that may track with state-change within AD. We computed the GAM fits for the cells between states 1 and 2, as ordered by the left-right axis, and visualized the changes of the top 50 genes (**Figure 6-7c**) across a binned "horizontal" pseudotime. Overall, most genes tend to show higher expression in the right state (1) and genes with higher expression in state 2 (left-hand side) still have high expression in state 1, unlike genes specific to state 1, whose expression is restricted to this state. Within the top 50 genes (by GAM fit p-value), we find enrichments for iron homeostasis (including ICSU and FTH1), NOGO signaling (RTN1, RTN4), calcium signaling, (CALM1, CALM3), brain expressed family members BEX1-4, and stress-granule associated ACTB and translation factors (**Figure 6-7c**). These genes are not necessarily top markers for either trajectory, suggesting that they specifically mark transitions between the two AD-related subclusters. Interestingly, phosphodiesterase PDE5A was the top differential gene along the horizontal trajectory (towards state 1) and family member PDE4D is differential in the opposite direction (**Figure 6-7d**). These phosphodiesterases, which downregulate the concentration of intracellular cyclic nucleotides (PDE5A: cGMP and PDE4D: cAMP), are currently top targets for Alzheimer's treatment, with multiple PDE5 inhibitors currently in randomized control trials for Alzheimer's disease<sup>326</sup>.

**Individuals and subtypes on the pseudotime trajectories:** We next asked whether cells from AD individuals consistently scored high on AD-related pseudotime trajectories. Plotting the density of cell scores from each individual separately, we found that both branches ordered individuals by AD status by their cells' locations on the trajectories (**Figure 6-7e**). The right-hand side trajectory best tracked with individual-level cognitive impairment. The horizontal trajectory did not match the cognitive status of the individuals but did show high heterogeneity by individual, where some individuals were highly left- or right-hand side biased. We scored the neurons in each region in each individual to study AD progression across regions within individuals (**Figure 6-8**). We found that individuals showed quite consistent average pseudotime values across regions, either indicating that neurons are similarly affected across regions, or that there are strong individual-level batch effects (or true biological effects) that affect the inference of these pseudotime trajectories.



**Figure 6-8.** Heatmap of average pseudotime score across excitatory neurons in each individual in each region for the three pseudotime trajectories. Bottom panels show individual-level annotations. Columns are clustered by euclidean distance and complete linkage.

When we repeated this analysis segmenting by neuronal subtypes, vulnerable subtypes were enriched for cells in earlier stages of right-side trajectory. The left branch scores did not track subtype vulnerability. Interestingly, vulnerable subtypes were particularly depleted for high scores on the early right branch and were enriched for being on the left-hand side of the horizontal pseudotime trajectory. Jointly these observations suggested that vulnerable subtypes either do not fully reach the end stage of the right branch or die/disappear upon reaching these transcriptional states. Therefore, our observations of subtype vulnerability in AD, our split pseudotime trajectories, and the two branching signatures developed in this section could be used to explore a potential mechanism for neurodegeneration in Alzheimer's whereby cells muster a specific response to AD, transition between two states of neuronal response to AD, and may enter a degenerative/dying state from either a specific state (or both) (**Figure 6-7e**).

## 7. Conclusion

In this thesis, I presented a series of analyses of cellular heterogeneity, its regulation, and its maintenance in health, aging, and disease. We built a map of the human epigenome to understand gene regulation and map disease loci and worked to develop methods to interpret, model, and impute epigenomic profiles of the human genome (**Chapter 1**). We used our methods and references to map and dissect non-coding disease loci, identifying multiple pleiotropic disease loci in coronary artery disease and dissecting an example showing tissue-specific gene involvement (**Chapter 2**). We built a transcriptomic reference of multiple human brain regions in age and in Alzheimer's disease at the level of individual cells, identified novel neuronal subtypes and evolutionarily-shared programs, and developed a framework for fast and scalable methods to interpret large scale single-cell atlases (**Chapter 3**). We investigated Alzheimer's disease progression across brain regions, mapped AD genetic risk genes and changes in functional modules across cell types and regions, and identified specific neuronal subtypes that are particularly susceptible to neurodegeneration (**Chapter 4**). We proposed somatic mosaicism as a new source of cell-level genetic risk for neurodegeneration, developed methods to map mosaic burden in individual cells jointly with expression, and found increased cell type-specific somatic mosaic burden in Alzheimer's disease (**Chapter 5**). Finally, we modeled cellular trajectories through neurodegeneration in human brains and mouse models of AD in conjunction with DNA damage accumulation and developed tools for building disease-driven pseudotime trajectories (**Chapter 6**).

Broadly, my work gives insights into how a single genome yields robust and diverse outcomes through modularity of biological functions, how this heterogeneity is maintained, and how it breaks down over time and in disease. These are very wide research directions, and it is to be expected that despite decades of research we are still far from comprehensive and clear mechanistic understandings of gene regulation, the genetic basis of disease, and neurodegeneration. The insights gained from my thesis work have relied crucially on recent cutting-edge advances in sequencing, microfluidics, and machine learning that are enabling researchers to systematically map out and characterize the human body in health and disease at increasingly high resolution.

These tools massively expand our view of biological processes but alone are not sufficient for developing mechanistic models in biology or reaching a deeper understanding of specific biological phenomena. In addition, these tools are biased by what we can measure - the advent of sequencing-driven approaches for measuring biological phenomena has been readily applied to and scaled up transcriptomic and epigenomic (and more recently proteomic) measurements. Given these low-hanging fruits in enumeration and characterization, in comparison, gene regulation and genetics-adjacent fields are currently relatively underexplored. These include the study of biology at multiple resolution, from protein structure, modification, and interaction to metabolism at the molecular, cellular, and systems level; and the study of biology at multiple time-scales: models of protein and cellular component kinetics at short temporal resolutions, signaling cascades and cell-cell communication at medium time-scales, and environmental exposures at long temporal resolutions. In addition, regardless of the domain of experimental methodological scientific advances, we do not yet have computational tools that can scale up knowledge synthesis and hypothesis generation in biology.

However, beyond enumeration of biological processes, single-cell and genomics approaches, in particular when coupled with high-throughput perturbations and advances in natural language processing and representation learning, will enable novel directions that can tackle large questions in gene regulation, genetics, and disease research. As we move from building maps of cellular heterogeneity to models of gene regulation and interrogating the impact of genetic variation on these models using perturbation experiments at scale and across time, we will need better models of both gene regulation and of the effects of genetic variants on these networks.



This will require increasingly sophisticated machine learning models for disentangling effects of perturbations on cells over time, theoretical and applied models of the effects of human genetic variation, algorithms for robust data synthesis and interpretation at scale, and - in the longer term - for hypothesis generation for targeted downstream analyses and systematic exploration of biological variation. Given exciting recent developments in high-throughput experimentation, statistical genetics, and machine learning, I map out below how these might converge in future research on gene regulation, disease genetics, and neurodegeneration.

### **7.1 Building gene regulatory models by moving from snapshot references to perturbations**

Genomic profiling at the single-cell level has enabled the systematic characterization of cellular biological processes, providing snapshots of cellular and genomic states at unprecedented scale. However, static maps of cellular processes, absent stimulation, temporal trajectory, or broader context, are not a sufficient or particularly efficient manner to fully model these processes. From an engineering standpoint, fully fledged models of human gene regulation should be able to model biological variation at many scales, from cellular identity and development to responses under stimulation or stress to slight changes in the presence of specific genetic variation. From a biological standpoint, these approaches only scratch the surface of describing mechanisms and are not flexible tools for downstream experimentation and hypothesis-driven research.

Currently, single-cell maps for disease are under-explored, as tools to properly map out consistent variation corresponding to different biological and technical covariates are not available. Even though these references often contain cells at many temporal points within the disease and with varying responses to their surrounding microenvironment, we lack the computational tools to systematically identify their molecular profiles as independent signals. In Alzheimer's disease in particular, poor and low-resolution phenotyping that is not at the level of individual cells means that much of the heterogeneity is lost. Approaches like the pseudotime trajectory method presented in section 6.4, which seeks to separate out disease-relevant variation from that driven by subtypes or technical effects, will be critical to extract more than surface level analyses in these datasets. In particular, recent developments in representation learning focusing on disentanglement<sup>327</sup>, contrastive learning, adversarial frameworks, and compositional representations<sup>328</sup> should be extended and applied to these datasets in order to both separate and comprehensively enumerate robust and consistent sources of variation in the data.

Excitingly, recent work has started to incorporate both existing and de novo genetic perturbations into single-cell transcriptomic and genomic maps - for example through perturb-seq and census-seq<sup>329-332</sup> - allowing researchers to measure the landscape around otherwise homogeneous snapshots of cell state. Randomly editing cells in cultured human cell lines to knock genes and regulatory elements out (KO), down (KD), or to partially alter their sequences alone or in combinations is already enabling quantitative measurements of how alterations to specific genes percolate through expression networks in context-specific manners to affect biological pathways and their dynamics<sup>329,333</sup>. Past contemporary algorithmic approaches building mostly undirected graphs, the direct perturbations and size of these analyses will require and benefit from computational methods, such as graph neural networks, that can learn from, build, and model large graph structures jointly with perturbations and other covariates.

### **7.2 Modeling variant function by resolving disease-risk across the allele frequency spectrum**

Genetic perturbations will help build gene regulatory networks and map out biological pathways in context-specific manners, but KO/KD experiments are blunt-edged approaches and more sophisticated models of variant function and effect size will be needed to properly understand the effects of human variation on cellular and systems-level function. Parallel development of assays containing genetic perturbations corresponding to existing human variation will aid in quantifying variant impact and add nuance to gene effects in networks beyond the current capabilities of KO/KD experiments. These approaches already include mapping expression in response to perturbations in cell lines with a varied genetic background<sup>332</sup> and single-cell level eQTL studies in specific tissues<sup>334</sup>. In modeling disease, these approaches will benefit from single-cell references in large cohorts, which should shed light on both how strong-effect risk variants affect cellular heterogeneity and how many weak-effect variants in high-PRS (polygenic risk score) individuals combine to elevate disease risk overall.

In order to properly model variant effects on gene regulation, there is a pressing need for computational and theoretical models that can resolve and reconcile human variation with experimental genetic perturbations and that can extrapolate to unseen variation. Experimental measurements will have to be combined with computational predictions of variant effect size that draw from sequence, protein structure, and gene regulatory networks as well as models of variant function and effect size for existing human variation across effect sizes and allele frequencies. As discussed in the main text, these methods should incorporate the range of variant effect size, from gene knock-outs to non-coding, and model the diversity of possible effects, as different variants within the same gene may affect different protein domains, interacting partners, yield different conformations, and operate in different functions.

In particular, statistical genetics models of variant function and disease pleiotropy will have to disentangle these factors, separating evolutionary pressures and genetic background from variant effect sizes and gene impact in cell types and pathways acting in disease to quantify variant impact in disease. Current statistical genetics approaches are not yet able to handle and disentangle these multiple lines of evidence but recent developments in protein structure modeling<sup>335</sup> and predicting epigenomic changes from sequence<sup>160,336</sup> show promise for modeling variant effects in a standalone manner. However, current black-box approaches for modeling variant effects still perform poorly at predicting the effects of underrepresented sequences not in the human genome and proteome. These models will have to build on deep learning advances - such as for robust and adversary proof models - to faithfully predict outside of currently mapped human variation and build robust and interpretable models of sequence binding, protein structure, protein interactions, and downstream regulatory effects.

### **7.3 Systematically learning biological components using representation learning**

While the continuing development of tools for measuring and quantifying molecular processes in the cell and genome should aid us to identify consistent and biologically relevant patterns of gene regulation, at the current pace, previous approaches to interpret data of increasingly higher dimension and size will not suffice. More recent efforts to build single-cell transcriptomic<sup>10</sup> and epigenomic<sup>11</sup> atlases of the human body have only exacerbated this trend, and such efforts are not properly compounding decades of data and research knowledge.

Due to improvements in technology and continuing profiling for specific applications, much of this data is being left behind and sits unused even when it may still contain un-explored biology or could just be used to boost newer cellular maps. As other technologies for measuring different modalities of cellular systems at multiple resolutions and time scales come of age, this growth in experimental data will have to be matched by algorithms that can handle its breadth and dimensionality. In particular, even if we can successfully generate transcriptomic and genomic maps of homeostatic and perturbed states at scale and can build context and cell

type-specific models of gene regulation using genetic perturbations, we are far from algorithmic approaches that can take in and synthesize this data properly.

To match and push the pace of experimental advances, we are interested in developing biological machine learning tools along three lines - disentangling sources of variation<sup>327</sup>, learning interpretable representations<sup>328</sup>, and performing knowledge synthesis for downstream hypothesis generation. Disentangling methods, with a rich history going from PCA<sup>337</sup> to more recent information theoretic variational autoencoder approaches<sup>327</sup>, are essential for biological data analysis as methods for systematically uncoupling sources of biological and technical variation (also see chapters 3.7 and 6.4). Separating biological and technical effects in computational biology has been an area of constant research<sup>338-340</sup>, but as datasets increase in size and number of modalities, there is a need for more general and powerful methods to separate many sources of variation in a single analysis. In parallel and potentially concurrently, these must be paired with machine learning approaches to learn comprehensive and interpretable representations for these sources of variation (as in chapter 1.9). Finally, we identify a potential avenue for future automated hypothesis and knowledge generation within gene regulation and genomics. If used jointly with methods for disentangling variation and representation learning, synthesis methods leveraging natural language processing advances<sup>341</sup> and transfer learning<sup>342</sup> could in the future map biological data representations to existing knowledge and identify constrained but unmodeled areas of variation, and aid researchers in prioritizing and targeting subsequent experiments.

While fully automated hypothesis generation and synthesis, which would require high-level artificial intelligence, are still far out of reach for biological knowledge generation, recent advances in machine learning present exciting prospects for researchers to perform computer-aided systematic exploration and synthesis of biology at scale.

## References

1. Frankish, A. *et al.* GENCODE 2021. *Nucleic Acids Res.* **49**, D916–D923 (2021).
2. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
3. Lambert, S. A. *et al.* The Human Transcription Factors. *Cell* **175**, 598–599 (2018).
4. Chen, T. & Dent, S. Y. R. Chromatin modifiers and remodellers: regulators of cellular differentiation. *Nat. Rev. Genet.* **15**, 93–106 (2014).
5. Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–330 (2015).
6. Stunnenberg, H. G., International Human Epigenome Consortium & Hirst, M. The International Human Epigenome Consortium: A Blueprint for Scientific Collaboration and Discovery. *Cell* **167**, 1897 (2016).
7. Genomics of Gene Regulation. *Genome.gov*  
<https://www.genome.gov/Funded-Programs-Projects/Genomics-of-Gene-Regulation>.
8. Ernst, J. & Kellis, M. Discovery and characterization of chromatin states for systematic annotation of the human genome. *Nat. Biotechnol.* **28**, 817–825 (2010).
9. Hoffman, M. M. *et al.* Unsupervised pattern discovery in human chromatin structure through genomic segmentation. *Nat. Methods* **9**, 473–476 (2012).
10. Tabula Muris Consortium. A single-cell transcriptomic atlas characterizes ageing tissues in the mouse. *Nature* **583**, 590–595 (2020).
11. Zhang, K. *et al.* A single-cell atlas of chromatin accessibility in the human genome. *Cell* **184**, 5985–6001.e19 (2021).
12. Zhao, C., Deng, W. & Gage, F. H. Mechanisms and functional implications of adult neurogenesis. *Cell* **132**, 645–660 (2008).
13. Moreno-Jiménez, E. P. *et al.* Adult hippocampal neurogenesis is abundant in neurologically healthy subjects and drops sharply in patients with Alzheimer’s disease. *Nat. Med.* **25**, 554–560 (2019).
14. Magrassi, L., Leto, K. & Rossi, F. Lifespan of neurons is uncoupled from organismal lifespan. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 4374–4379 (2013).
15. Clarke, D. D. & Sokoloff, L. Basic neurochemistry: molecular, cellular and medical aspects. *Lippincott-Raven, Philadelphia* (1999).
16. Raichle, M. E. & Gusnard, D. A. Appraising the brain’s energy budget. *Proc. Natl. Acad. Sci. U. S. A.* **99**, 10237–10239 (2002).
17. Lim, J. S. *et al.* Brain somatic mutations in MTOR cause focal cortical dysplasia type II leading to intractable epilepsy. *Nat. Med.* **21**, 395–400 (2015).
18. D’Gama, A. M. *et al.* Mammalian target of rapamycin pathway mutations cause hemimegalencephaly and focal cortical dysplasia. *Ann. Neurol.* **77**, 720–725 (2015).
19. Schneider, J. A., Arvanitakis, Z., Bang, W. & Bennett, D. A. Mixed brain pathologies account for most dementia cases in community-dwelling older persons. *Neurology* **69**, 2197–2204 (2007).
20. Braak, H. & Braak, E. The human entorhinal cortex: normal morphology and lamina-specific pathology in various diseases. *Neurosci. Res.* **15**, 6–31 (1992).
21. Braak, H. & Braak, E. Staging of Alzheimer’s disease-related neurofibrillary changes. *Neurobiol. Aging* **16**, 271–8; discussion 278–84 (1995).
22. McGregor, M. M. & Nelson, A. B. Circuit Mechanisms of Parkinson’s Disease. *Neuron* **101**, 1042–1056 (2019).
23. Sherrington, R. *et al.* Alzheimer’s disease associated with mutations in presenilin 2 is rare and variably penetrant. *Hum. Mol. Genet.* **5**, 985–988 (1996).
24. Arboleda-Velasquez, J. F. *et al.* Resistance to autosomal dominant Alzheimer’s disease in an APOE3 Christchurch homozygote: a case report. *Nat. Med.* **25**, 1680–1683 (2019).



25. Cichon, S. *et al.* Genomewide association studies: history, rationale, and prospects for psychiatric disorders. *Am. J. Psychiatry* **166**, 540–556 (2009).
26. Homsy, J. *et al.* De novo mutations in congenital heart disease with neurodevelopmental and other congenital anomalies. *Science* **350**, 1262–1266 (2015).
27. Sonnen, J. A. *et al.* Pathological correlates of dementia in a longitudinal, population-based sample of aging. *Ann. Neurol.* **62**, 406–413 (2007).
28. Jack, C. R., Jr *et al.* Hypothetical model of dynamic biomarkers of the Alzheimer's pathological cascade. *Lancet Neurol.* **9**, 119–128 (2010).
29. Wilson, R. S. *et al.* Temporal course and pathologic basis of unawareness of memory loss in dementia. *Neurology* **85**, 984–991 (2015).
30. Bennett, D. A. *et al.* Natural history of mild cognitive impairment in older persons. *Neurology* **59**, 198–205 (2002).
31. Bennett, D. A., Schneider, J. A., Wilson, R. S., Bienias, J. L. & Arnold, S. E. Neurofibrillary tangles mediate the association of amyloid load with clinical Alzheimer disease and level of cognitive function. *Arch. Neurol.* **61**, 378–384 (2004).
32. Delorey, T. M. *et al.* COVID-19 tissue atlases reveal SARS-CoV-2 pathology and cellular targets. *Nature* **595**, 107–113 (2021).
33. Hawrylycz, M. J. *et al.* An anatomically comprehensive atlas of the adult human brain transcriptome. *Nature* **489**, 391–399 (2012).
34. Habib, N. *et al.* Massively parallel single-nucleus RNA-seq with DroNc-seq. *Nat. Methods* **14**, 955–958 (2017).
35. Lake, B. B. *et al.* Integrative single-cell analysis of transcriptional and epigenetic states in the human adult brain. *Nat. Biotechnol.* **36**, 70–80 (2018).
36. Grubman, A. *et al.* A single-cell atlas of entorhinal cortex from individuals with Alzheimer's disease reveals cell-type-specific gene expression regulation. *Nat. Neurosci.* **22**, 2087–2097 (2019).
37. Mathys, H. *et al.* Single-cell transcriptomic analysis of Alzheimer's disease. *Nature* **570**, 332–337 (2019).
38. Chin, J. *et al.* Reelin depletion in the entorhinal cortex of human amyloid precursor protein transgenic mice and humans with Alzheimer's disease. *J. Neurosci.* **27**, 2727–2733 (2007).
39. Kibro-Flatmoen, A. *et al.* Re-emphasizing early Alzheimer's disease pathology starting in select entorhinal neurons, with a special focus on mitophagy. *Ageing Res. Rev.* **67**, 101307 (2021).
40. Leng, K. *et al.* Molecular characterization of selectively vulnerable neurons in Alzheimer's disease. *Nat. Neurosci.* **24**, 276–287 (2021).
41. Stranahan, A. M. & Mattson, M. P. Selective vulnerability of neurons in layer II of the entorhinal cortex during aging and Alzheimer's disease. *Neural Plast.* **2010**, 108190 (2010).
42. Kordower, J. H. *et al.* Loss and atrophy of layer II entorhinal cortex neurons in elderly people with mild cognitive impairment. *Ann. Neurol.* **49**, 202–213 (2001).
43. Van Hoesen, G. W., Hyman, B. T. & Damasio, A. R. Entorhinal cortex pathology in Alzheimer's disease. *Hippocampus* **1**, 1–8 (1991).
44. Gómez-Isla, T. *et al.* Profound loss of layer II entorhinal cortex neurons occurs in very mild Alzheimer's disease. *J. Neurosci.* **16**, 4491–4500 (1996).
45. Price, J. L. *et al.* Neuron number in the entorhinal cortex and CA1 in preclinical Alzheimer disease. *Arch. Neurol.* **58**, 1395–1402 (2001).
46. Visscher, P. M. *et al.* 10 Years of GWAS Discovery: Biology, Function, and Translation. *Am. J. Hum. Genet.* **101**, 5–22 (2017).
47. Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* **47**, D1005–D1012 (2019).
48. Ulland, T. K. *et al.* TREM2 Maintains Microglial Metabolic Fitness in Alzheimer's Disease.

- Cell* **170**, 649–663.e13 (2017).
49. Keren-Shaul, H. *et al.* A Unique Microglia Type Associated with Restricting Development of Alzheimer's Disease. *Cell* vol. 169 1276–1290.e17 (2017).
  50. Wightman, D. P. *et al.* A genome-wide association study with 1,126,563 individuals identifies new risk loci for Alzheimer's disease. *Nat. Genet.* **53**, 1276–1282 (2021).
  51. Gallagher, M. D. & Chen-Plotkin, A. S. The Post-GWAS Era: From Association to Function. *Am. J. Hum. Genet.* **102**, 717–730 (2018).
  52. Ward, L. D. & Kellis, M. Interpreting noncoding genetic variation in complex traits and human disease. *Nat. Biotechnol.* **30**, 1095–1106 (2012).
  53. Ernst, J. *et al.* Mapping and analysis of chromatin state dynamics in nine human cell types. *Nature* **473**, 43–49 (2011).
  54. Maurano, M. T. *et al.* Systematic localization of common disease-associated variation in regulatory DNA. *Science* **337**, 1190–1195 (2012).
  55. Cantor, R. M., Lange, K. & Sinsheimer, J. S. Prioritizing GWAS results: A review of statistical methods and recommendations for their application. *Am. J. Hum. Genet.* **86**, 6–22 (2010).
  56. Farh, K. K.-H. *et al.* Genetic and epigenetic fine mapping of causal autoimmune disease variants. *Nature* **518**, 337–343 (2015).
  57. Dimas, A. S., Deutsch, S. & Stranger, B. E. Common regulatory variation impacts gene expression in a cell type-dependent manner. (2009).
  58. Ward, L. D. & Kellis, M. HaploReg: a resource for exploring chromatin states, conservation, and regulatory motif alterations within sets of genetically linked variants. *Nucleic Acids Res.* **40**, D930–4 (2012).
  59. Pickrell, J. K. Joint analysis of functional genomic data and genome-wide association studies of 18 human traits. *Am. J. Hum. Genet.* **94**, 559–573 (2014).
  60. Gusev, A. *et al.* Partitioning heritability of regulatory and cell-type-specific variants across 11 common diseases. *Am. J. Hum. Genet.* **95**, 535–552 (2014).
  61. Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nat. Genet.* **47**, 1228–1235 (2015).
  62. Gjoneska, E. *et al.* Conserved epigenomic signals in mice and humans reveal immune basis of Alzheimer's disease. *Nature* **518**, 365–369 (2015).
  63. Henstridge, C. M., Hyman, B. T. & Spires-Jones, T. L. Beyond the neuron-cellular interactions early in Alzheimer disease pathogenesis. *Nat. Rev. Neurosci.* **20**, 94–108 (2019).
  64. Purcell, S. M. *et al.* A polygenic burden of rare disruptive mutations in schizophrenia. *Nature* **506**, 185–190 (2014).
  65. De Rubeis, S. *et al.* Synaptic, transcriptional and chromatin genes disrupted in autism. *Nature* **515**, 209–215 (2014).
  66. Lambert, J. C. *et al.* Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer's disease. *Nat. Genet.* **45**, 1452–1458 (2013).
  67. Kunkle, B. W. *et al.* Author Correction: Genetic meta-analysis of diagnosed Alzheimer's disease identifies new risk loci and implicates A $\beta$ , tau, immunity and lipid processing. *Nat. Genet.* **51**, 1423–1424 (2019).
  68. Kelleher, R. J., 3rd & Shen, J. Presenilin-1 mutations and Alzheimer's disease. *Proceedings of the National Academy of Sciences of the United States of America* vol. 114 629–631 (2017).
  69. Cai, Y., An, S. S. A. & Kim, S. Mutations in presenilin 2 and its implications in Alzheimer's disease and other dementia-associated disorders. *Clin. Interv. Aging* **10**, 1163–1172 (2015).
  70. Chartier-Harlin, M. C. *et al.* Early-onset Alzheimer's disease caused by mutations at codon 717 of the beta-amyloid precursor protein gene. *Nature* **353**, 844–846 (1991).
  71. Goate, A. *et al.* Segregation of a missense mutation in the amyloid precursor protein gene

- with familial Alzheimer's disease. *Nature* **349**, 704–706 (1991).
72. Cummings, J. L. Cognitive and behavioral heterogeneity in Alzheimer's disease: seeking the neurobiological basis. *Neurobiol. Aging* **21**, 845–861 (2000).
  73. Bettens, K., Sleegers, K. & Van Broeckhoven, C. Genetic insights in Alzheimer's disease. *Lancet Neurol.* **12**, 92–104 (2013).
  74. Saunders, A. M. *et al.* Association of apolipoprotein E allele epsilon 4 with late-onset familial and sporadic Alzheimer's disease. *Neurology* **43**, 1467–1472 (1993).
  75. Seshadri, S. *et al.* Genome-wide analysis of genetic loci associated with Alzheimer disease. *JAMA* **303**, 1832–1840 (2010).
  76. Harold, D. *et al.* Genome-wide association study identifies variants at CLU and PICALM associated with Alzheimer's disease. *Nat. Genet.* **41**, 1088–1093 (2009).
  77. Lambert, J.-C. *et al.* Genome-wide association study identifies variants at CLU and CR1 associated with Alzheimer's disease. *Nat. Genet.* **41**, 1094–1099 (2009).
  78. Bis, J. C. *et al.* Whole exome sequencing study identifies novel rare and common Alzheimer's-Associated variants involved in immune response and transcriptional regulation. *Mol. Psychiatry* (2018) doi:10.1038/s41380-018-0112-7.
  79. Raghavan, N. S. *et al.* Whole-exome sequencing in 20,197 persons for rare variants in Alzheimer's disease. *Ann Clin Transl Neurol* **5**, 832–842 (2018).
  80. Mayeux, R. Putative risk factors for Alzheimer's disease. in *Alzheimer's Disease* 39–50 (CRC Press, 2019).
  81. Ju, Y. S. *et al.* Somatic mutations reveal asymmetric cellular dynamics in the early human embryo. *Nature* **543**, 714–718 (2017).
  82. Loh, P.-R. *et al.* Insights into clonal haematopoiesis from 8,342 mosaic chromosomal alterations. *Nature* **559**, 350–355 (2018).
  83. Genovese, G., Jaiswal, S., Ebert, B. L. & McCarroll, S. A. Clonal hematopoiesis and blood-cancer risk. *The New England journal of medicine* vol. 372 1071–1072 (2015).
  84. Genovese, G. *et al.* Clonal hematopoiesis and blood-cancer risk inferred from blood DNA sequence. *N. Engl. J. Med.* **371**, 2477–2487 (2014).
  85. Manheimer, K. B. *et al.* Robust identification of mosaic variants in congenital heart disease. *Hum. Genet.* **137**, 183–193 (2018).
  86. Groesser, L. *et al.* Postzygotic HRAS and KRAS mutations cause nevus sebaceous and Schimmelpenning syndrome. *Nat. Genet.* **44**, 783–787 (2012).
  87. Poduri, A., Evrony, G. D., Cai, X. & Walsh, C. A. Somatic Mutation, Genomic Variation, and Neurological Disease. *Science* vol. 341 1237758 (2013).
  88. McConnell, M. J. *et al.* Intersection of diverse neuronal genomes and neuropsychiatric disease: The Brain Somatic Mosaicism Network. *Science* **356**, (2017).
  89. D'Gama, A. M. *et al.* Targeted DNA Sequencing from Autism Spectrum Disorder Brains Implicates Multiple Genetic Mechanisms. *Neuron* **88**, 910–917 (2015).
  90. Iossifov, I. *et al.* The contribution of de novo coding mutations to autism spectrum disorder. *Nature* **515**, 216–221 (2014).
  91. Neale, B. M. *et al.* Patterns and rates of exonic de novo mutations in autism spectrum disorders. *Nature* **485**, 242–245 (2012).
  92. Mirzaa, G. M. *et al.* Association of MTOR Mutations With Developmental Brain Disorders, Including Megalencephaly, Focal Cortical Dysplasia, and Pigmentary Mosaicism. *JAMA Neurol.* **73**, 836–845 (2016).
  93. Lodato, M. A. *et al.* Aging and neurodegeneration are associated with increased mutations in single human neurons. doi:10.1101/221960.
  94. Ernst, J. & Kellis, M. Large-scale imputation of epigenomic datasets for systematic annotation of diverse human tissues. *Nat. Biotechnol.* **33**, 364–376 (2015).
  95. Durham, T. J., Libbrecht, M. W., Howbert, J. J., Bilmes, J. & Noble, W. S. PREDICTD PaRallel Epigenomics Data Imputation with Cloud-based Tensor Decomposition. *Nat.*

- Commun.* **9**, 1402 (2018).
96. Schreiber, J., Durham, T., Bilmes, J. & Noble, W. S. Avocado: a multi-scale deep tensor factorization method learns a latent representation of the human epigenome. *Genome Biol.* **21**, 81 (2020).
  97. Li, D., Hsu, S., Purushotham, D., Sears, R. L. & Wang, T. WashU Epigenome Browser update 2019. *Nucleic Acids Research* vol. 47 W158–W165 (2019).
  98. Calo, E. & Wysocka, J. Modification of enhancer chromatin: what, how, and why? *Mol. Cell* **49**, 825–837 (2013).
  99. Becker, J. S., Nicetto, D. & Zaret, K. S. H3K9me3-Dependent Heterochromatin: Barrier to Cell Fate Changes. *Trends Genet.* **32**, 29–41 (2016).
  100. Meuleman, W. *et al.* Index and biological spectrum of human DNase I hypersensitive sites. *Nature* **584**, 244–251 (2020).
  101. ENCODE Project Consortium *et al.* Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature* **583**, 699–710 (2020).
  102. Vierstra, J. *et al.* Global reference mapping of human transcription factor footprints. *Nature* **583**, 729–736 (2020).
  103. Laslo, P. *et al.* Multilineage transcriptional priming and determination of alternate hematopoietic cell fates. *Cell* **126**, 755–766 (2006).
  104. Blackshear, P. J. *et al.* Graded phenotypic response to partial and complete deficiency of a brain-specific transcript variant of the winged helix transcription factor RFX4. *Development* **130**, 4539–4552 (2003).
  105. Olson, J. M. *et al.* NeuroD2 is necessary for development and survival of central nervous system neurons. *Dev. Biol.* **234**, 174–187 (2001).
  106. Katz, J. P. *et al.* The zinc-finger transcription factor Klf4 is required for terminal differentiation of goblet cells in the colon. *Development* **129**, 2619–2628 (2002).
  107. Jacquemin, P., Martial, J. A. & Davidson, I. Human TEF-5 is preferentially expressed in placenta and binds to multiple functional elements of the human chorionic somatomammotropin-B gene enhancer. *J. Biol. Chem.* **272**, 12928–12937 (1997).
  108. Tanaka, T. *et al.* Dysregulated expression of P1 and P2 promoter-driven hepatocyte nuclear factor-4 $\alpha$  in the pathogenesis of human cancer. *The Journal of Pathology: A Journal of the Pathological Society of Great Britain and Ireland* **208**, 662–672 (2006).
  109. Wagner, E. F. & Eferl, R. Fos/AP-1 proteins in bone and the immune system. *Immunol. Rev.* **208**, 126–140 (2005).
  110. Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat. Genet.* **25**, 25–29 (2000).
  111. McLean, C. Y. *et al.* GREAT improves functional interpretation of cis-regulatory regions. *Nat. Biotechnol.* **28**, 495–501 (2010).
  112. Karimzadeh, M., Ernst, C., Kundaje, A. & Hoffman, M. M. Umap and Bismap: quantifying genome and methylome mappability. *Nucleic Acids Res.* **46**, e120 (2018).
  113. Rosenbloom, K. R. *et al.* The UCSC Genome Browser database: 2015 update. *Nucleic Acids Res.* **43**, D670–81 (2015).
  114. Feng, J., Liu, T., Qin, B., Zhang, Y. & Liu, X. S. Identifying ChIP-seq enrichment using MACS. *Nat. Protoc.* **7**, 1728–1740 (2012).
  115. Kharchenko, P. V., Tolstorukov, M. Y. & Park, P. J. Design and analysis of ChIP-seq experiments for DNA-binding proteins. *Nat. Biotechnol.* **26**, 1351–1359 (2008).
  116. Gu, Z., Gu, L., Eils, R., Schlesner, M. & Brors, B. circlize Implements and enhances circular visualization in R. *Bioinformatics* **30**, 2811–2812 (2014).
  117. Melville, J. uwot: The Uniform Manifold Approximation and Projection (UMAP) Method for Dimensionality Reduction. (2019).
  118. Kheradpour, P. & Kellis, M. Systematic discovery and characterization of regulatory motifs in ENCODE TF binding experiments. *Nucleic Acids Res.* **42**, 2976–2987 (2014).



119. Khan, A. *et al.* JASPAR 2018: update of the open-access database of transcription factor binding profiles and its web framework. *Nucleic Acids Res.* **46**, D1284 (2018).
120. Kulakovskiy, I. V. *et al.* HOCOMOCO: towards a complete collection of transcription factor binding models for human and mouse via large-scale ChIP-Seq analysis. *Nucleic Acids Res.* **46**, D252–D259 (2018).
121. Jolma, A. *et al.* DNA-binding specificities of human transcription factors. *Cell* **152**, 327–339 (2013).
122. Liu, Y., Sarkar, A., Kheradpour, P., Ernst, J. & Kellis, M. Evidence of reduced recombination rate in human regulatory domains. *Genome Biology* vol. 18 (2017).
123. Moore, J. E., Pratt, H. E., Purcaro, M. J. & Weng, Z. A curated benchmark of enhancer-gene interactions for evaluating enhancer-target gene prediction methods. *Genome Biol.* **21**, 17 (2020).
124. Fulco, C. P. *et al.* Activity-by-contact model of enhancer–promoter regulation from thousands of CRISPR perturbations. *Nature Genetics* vol. 51 1664–1669 (2019).
125. Chang, W. *et al.* Shiny: web application framework for R. *R package version 1*, (2017).
126. Franceschini, A. *et al.* STRING v9.1: protein-protein interaction networks, with increased coverage and integration. *Nucleic Acids Res.* **41**, D808–D815 (2013).
127. Kingma, D. P. & Welling, M. Auto-Encoding Variational Bayes. *arXiv [stat.ML]* (2013).
128. Lee, D. D. & Seung, H. S. Learning the parts of objects by non-negative matrix factorization. *Nature* **401**, 788–791 (1999).
129. Ram, O. *et al.* Combinatorial patterning of chromatin regulators uncovered by genome-wide location analysis in human cells. *Cell* **147**, 1628–1639 (2011).
130. Cao, F. *et al.* An Ash2L/RbBP5 heterodimer stimulates the MLL1 methyltransferase activity through coordinated substrate interactions with the MLL1 SET domain. *PLoS One* **5**, 1–11 (2010).
131. Wu, L. *et al.* ASH2L regulates ubiquitylation signaling to MLL: Trans-Regulation of H3 K4 methylation in higher eukaryotes. *Mol. Cell* **49**, 1108–1120 (2013).
132. Shi, Y. *et al.* Histone demethylation mediated by the nuclear amine oxidase homolog LSD1. *Cell* **119**, 941–953 (2004).
133. Hennig, A. K., Peng, G. H. & Chen, S. Transcription Coactivators p300 and CBP Are Necessary for Photoreceptor-Specific Chromatin Organization and Gene Expression. *PLoS One* **8**, (2013).
134. Minc, E., Courvalin, J.-C. & Buendia, B. HP1 $\gamma$  associates with euchromatin and heterochromatin in mammalian nuclei and chromosomes. *Cytogenet. Genome Res.* **90**, 279–284 (2000).
135. Creppe, C., Palau, A., Malinverni, R., Valero, V. & Buschbeck, M. A Cbx8-Containing Polycomb Complex Facilitates the Transition to Gene Activation during ES Cell Differentiation. *PLoS Genet.* **10**, (2014).
136. Gil, J. & O’Loughlen, A. PRC1 complex diversity: Where is it taking us? *Trends Cell Biol.* **24**, 632–641 (2014).
137. Schoeftner, S. *et al.* Recruitment of PRC1 function at the initiation of X inactivation independent of PRC2 and silencing. *EMBO J.* **25**, 3110–3122 (2006).
138. Margueron, R. & Reinberg, D. The Polycomb complex PRC2 and its mark in life. *Nature* **469**, 343–349 (2011).
139. Feng, W. *et al.* The chromatin remodeler CHD7 regulates adult neurogenesis via activation of sox2 transcription factors. *Cell Stem Cell* **13**, 62–72 (2013).
140. Engelen, E. *et al.* Sox2 cooperates with Chd7 to regulate genes that are mutated in human syndromes. *Nat. Genet.* **43**, 607–611 (2011).
141. Zaidi, S. *et al.* De novo mutations in histone-modifying genes in congenital heart disease. *Nature* **498**, 220–223 (2013).
142. Tai, H. H. *et al.* CHD1 associates with NCoR and histone deacetylase as well as with RNA

- splicing proteins. *Biochem. Biophys. Res. Commun.* **308**, 170–176 (2003).
143. You, A., Tong, J. K., Grozinger, C. M. & Schreiber, S. L. CoREST is an integral component of the CoREST- human histone deacetylase complex. *Proceedings of the National Academy of Sciences* vol. 98 1454–1458 (2001).
  144. Maison, C. & Almouzni, G. HP1 and the dynamics of heterochromatin maintenance. *Nat. Rev. Mol. Cell Biol.* **5**, 296–304 (2004).
  145. Guerra-Calderas, L., González-Barrios, R., Herrera, L. A., de León, D. C. & Soto-Reyes, E. The role of the histone demethylase KDM4A in cancer. *Cancer Genetics* vol. 208 215–224 (2015).
  146. Gažová, I., Lengeling, A. & Summers, K. M. Lysine demethylases KDM6A and UTY: The X and Y of histone demethylation. *Molecular Genetics and Metabolism* vol. 127 31–44 (2019).
  147. Zhang, Y. *et al.* SAP30, a Novel Protein Conserved between Human and Yeast, Is a Component of a Histone Deacetylase Complex. *Molecular Cell* vol. 1 1021–1031 (1998).
  148. Hansen, A. S., Pustova, I., Cattoglio, C., Tjian, R. & Darzacq, X. CTCF and cohesin regulate chromatin loop stability with distinct dynamics. *Elife* **6**, (2017).
  149. Smolle, M. *et al.* Chromatin remodelers Isw1 and Chd1 maintain chromatin structure during transcription by preventing histone exchange. *Nat. Struct. Mol. Biol.* **19**, 884–892 (2012).
  150. Andersson, R. *et al.* An atlas of active enhancers across human cell types and tissues. *Nature* **507**, 455–461 (2014).
  151. Schreiber, J., Singh, R., Bilmes, J. & Noble, W. S. A pitfall for machine learning methods aiming to predict across cell types. *Genome Biol.* **21**, 282 (2020).
  152. Schreiber, J., Bilmes, J. & Noble, W. S. Completing the ENCODE3 compendium yields accurate imputations across a variety of assays and human biosamples. *Genome Biol.* **21**, 82 (2020).
  153. Lee, J. *et al.* ENCODE ATAC-seq pipeline. (2020).
  154. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
  155. Wysocker, A., Tibbetts, K. & Fennell, T. Picard tools version 1.90. <http://picard.sourceforge.net> (Accessed 14 December 2016) **107**, 308 (2013).
  156. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
  157. Zhao, Y., Wong, L. & Goh, W. W. B. How to do quantile normalization correctly for gene expression data analyses. *Sci. Rep.* **10**, 15534 (2020).
  158. Townes, F. W. & Irizarry, R. A. Quantile normalization of single-cell RNA-seq read counts without unique molecular identifiers. *Genome Biol.* **21**, 160 (2020).
  159. Bonhoure, N. *et al.* Quantifying ChIP-seq data: a spiking method providing an internal reference for sample-to-sample normalization. *Genome Res.* **24**, 1157–1168 (2014).
  160. Lee, D. *et al.* A method to predict the impact of regulatory variants from DNA sequence. *Nat. Genet.* **47**, 955 (2015).
  161. Moyerbrailean, G. A. *et al.* Which Genetics Variants in DNase-Seq Footprints Are More Likely to Alter Binding? *PLoS Genet.* **12**, e1005875 (2016).
  162. MacArthur, J. *et al.* The new NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog). *Nucleic Acids Res.* **45**, D896–D901 (2017).
  163. Michailidou, K. *et al.* Association analysis identifies 65 new breast cancer risk loci. *Nature* **551**, 92–94 (2017).
  164. Goes, F. S. *et al.* Genome-wide association study of schizophrenia in Ashkenazi Jews. *Am. J. Med. Genet. B Neuropsychiatr. Genet.* **168**, 649–659 (2015).
  165. Lupien, M., Markowitz, S. & Scacheri, P. C. Combinatorial effects of multiple enhancer variants in linkage disequilibrium dictate levels of gene expression to confer susceptibility to common traits. *Genome* (2014).
  166. Winkler, T. W. *et al.* The Influence of Age and Sex on Genetic Associations with Adult Body

- Size and Shape: A Large-Scale Genome-Wide Interaction Study. *PLoS Genet.* **11**, e1005378 (2015).
167. van der Harst, P. & Verweij, N. Identification of 64 Novel Genetic Loci Provides an Expanded View on the Genetic Architecture of Coronary Artery Disease. *Circ. Res.* **122**, 433–443 (2018).
  168. Abifadel, M. *et al.* Mutations in PCSK9 cause autosomal dominant hypercholesterolemia. *Nat. Genet.* **34**, 154–156 (2003).
  169. Busnelli, M. *et al.* Liver-specific deletion of the P1pp3 gene alters plasma lipid composition and worsens atherosclerosis in apoE<sup>-/-</sup> mice. *Sci. Rep.* **7**, 44503 (2017).
  170. Lappegård, K. T. *et al.* A vital role for complement in heart disease. *Mol. Immunol.* **61**, 126–134 (2014).
  171. Arad Michael, Seidman Christine E. & Seidman J.G. AMP-Activated Protein Kinase in the Heart. *Circ. Res.* **100**, 474–488 (2007).
  172. Litviňuková, M. *et al.* Cells of the adult human heart. *Nature* **588**, 466–472 (2020).
  173. Suryawanshi, H. *et al.* Cell atlas of the foetal human heart and implications for autoimmune-mediated congenital heart block. *Cardiovasc. Res.* **116**, 1446–1457 (2020).
  174. Tucker, N. R. *et al.* Transcriptional and Cellular Diversity of the Human Heart. *Circulation* **142**, 466–482 (2020).
  175. Wang, L. *et al.* Single-cell reconstruction of the adult human heart during heart failure and recovery reveals the cellular landscape underlying cardiac function. *Nat. Cell Biol.* **22**, 108–119 (2020).
  176. Boix, C. A., James, B. T., Park, Y. P., Meuleman, W. & Kellis, M. Regulatory genomic circuitry of human disease loci by integrative epigenomics. *Nature* vol. 590 300–307 (2021).
  177. Wirka, R. C. *et al.* Atheroprotective roles of smooth muscle cell phenotypic modulation and the TCF21 disease gene as revealed by single-cell analysis. *Nat. Med.* **25**, 1280–1289 (2019).
  178. Machiela, M. J. & Chanock, S. J. LDlink: a web-based application for exploring population-specific haplotype structure and linking correlated alleles of possible functional variants. *Bioinformatics* **31**, 3555–3557 (2015).
  179. Douglas, G. *et al.* A key role for the novel coronary artery disease gene JCAD in atherosclerosis via shear stress mechanotransduction. *Cardiovasc. Res.* **116**, 1863–1874 (2020).
  180. Zhao, C. *et al.* Supervillin promotes tumor angiogenesis in liver cancer. *Oncol. Rep.* **44**, 674–684 (2020).
  181. Yu, G., Wang, L.-G., Han, Y. & He, Q.-Y. clusterProfiler: an R package for comparing biological themes among gene clusters. *OMICS* **16**, 284–287 (2012).
  182. Csardi, G., Nepusz, T. & Others. The igraph software package for complex network research. *InterJournal, Complex Systems* **1695**, 1–9 (2006).
  183. Bennett, D. A. *et al.* Religious Orders Study and Rush Memory and Aging Project. *J. Alzheimers. Dis.* **64**, S161–S189 (2018).
  184. Kotliar, D. *et al.* Identifying gene expression programs of cell-type identity and cellular activity with single-cell RNA-Seq. *Elife* **8**, (2019).
  185. Hodge, R. D. *et al.* Conserved cell types with divergent features in human versus mouse cortex. *Nature* **573**, 61–68 (2019).
  186. Van de Sande, B. *et al.* A scalable SCENIC workflow for single-cell gene regulatory network analysis. *Nat. Protoc.* **15**, 2247–2276 (2020).
  187. Aibar, S. *et al.* SCENIC: single-cell regulatory network inference and clustering. *Nat. Methods* **14**, 1083–1086 (2017).
  188. Azevedo, F. A. C. *et al.* Equal numbers of neuronal and nonneuronal cells make the human brain an isometrically scaled-up primate brain. *J. Comp. Neurol.* **513**, 532–541 (2009).
  189. Pakkenberg, B. & Gundersen, H. J. Total number of neurons and glial cells in human brain

- nuclei estimated by the disector and the fractionator. *J. Microsc.* **150**, 1–20 (1988).
190. Saunders, A. *et al.* Molecular Diversity and Specializations among the Cells of the Adult Mouse Brain. *Cell* **174**, 1015–1030.e16 (2018).
  191. Zeisel, A. *et al.* Molecular Architecture of the Mouse Nervous System. *Cell* **174**, 999–1014.e22 (2018).
  192. Bennett, D. A., Schneider, J. A., Tang, Y., Arnold, S. E. & Wilson, R. S. The effect of social networks on the relation between Alzheimer’s disease pathology and level of cognitive function in old people: a longitudinal cohort study. *Lancet Neurol.* **5**, 406–412 (2006).
  193. Zheng, G. X. Y. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat. Commun.* **8**, 14049 (2017).
  194. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol.* **19**, 15 (2018).
  195. Traag, V. A., Waltman, L. & van Eck, N. J. From Louvain to Leiden: guaranteeing well-connected communities. *Sci. Rep.* **9**, 5233 (2019).
  196. McGinnis, C. S., Murrow, L. M. & Gartner, Z. J. DoubletFinder: Doublet Detection in Single-Cell RNA Sequencing Data Using Artificial Nearest Neighbors. *Cell Syst* **8**, 329–337.e4 (2019).
  197. Hafemeister, C. & Satija, R. Analyzing scRNA-seq data with the sctransform and offset models. (2020).
  198. Seurat, R. package (v2. 3.4 and v3. 1.3) Butler *et al.* (2018). *Stuart et al* (2019).
  199. Kierdorf, K., Masuda, T., Jordão, M. J. C. & Prinz, M. Macrophages at CNS interfaces: ontogeny and function in health and disease. *Nat. Rev. Neurosci.* **20**, 547–562 (2019).
  200. Velmeshev, D. *et al.* Single-cell genomics identifies cell type-specific molecular changes in autism. *Science* **364**, 685–689 (2019).
  201. Bakken, T. E. *et al.* Single-cell and single-nucleus RNA-seq uncovers shared and distinct axes of variation in dorsal LGN neurons in mice, non-human primates, and humans. *Elife* **10**, (2021).
  202. Yao, Z. *et al.* A taxonomy of transcriptomic cell types across the isocortex and hippocampal formation. *Cell* **184**, 3222–3241.e26 (2021).
  203. Franjic, D. *et al.* Molecular Diversity Among Adult Human Hippocampal and Entorhinal Cells. *bioRxiv* 2019.12.31.889139 (2020) doi:10.1101/2019.12.31.889139.
  204. Kalmbach, B. E. *et al.* Signature morpho-electric, transcriptomic, and dendritic properties of human layer 5 neocortical pyramidal neurons. *Neuron* **109**, 2914–2927.e5 (2021).
  205. Kepecs, A. & Fishell, G. Interneuron cell types are fit to function. *Nature* **505**, 318–326 (2014).
  206. Allen, G. V. & Hopkins, D. A. Mamillary body in the rat: topography and synaptology of projections from the subicular complex, prefrontal cortex, and midbrain tegmentum. *J. Comp. Neurol.* **286**, 311–336 (1989).
  207. Jankowski, M. M. *et al.* The anterior thalamus provides a subcortical circuit supporting memory and spatial navigation. *Frontiers in Systems Neuroscience* vol. 7 (2013).
  208. Wainberg, M. *et al.* A genome-wide atlas of co-essential modules assigns function to uncharacterized genes. *Nat. Genet.* **53**, 638–649 (2021).
  209. Dziedzic, A., Miller, E., Saluk-Bijak, J. & Bijak, M. The GPR17 Receptor-A Promising Goal for Therapy and a Potential Marker of the Neurodegenerative Process in Multiple Sclerosis. *Int. J. Mol. Sci.* **21**, (2020).
  210. Fang, N. *et al.* Sirt2 epigenetically down-regulates PDGFR $\alpha$  expression and promotes CG4 cell differentiation. *Cell Cycle* **18**, 1095–1109 (2019).
  211. Lester, E. *et al.* Tau aggregates are RNA-protein assemblies that mislocalize multiple nuclear speckle components. *Neuron* **109**, 1675–1691.e9 (2021).
  212. Kessy, A., Lewin, A. & Strimmer, K. Optimal Whitening and Decorrelation. *Am. Stat.* **72**, 309–314 (2018).



213. Student. Probable Error of a Correlation Coefficient. *Biometrika* **6**, 302–310 (1908).
214. Suo, S. *et al.* Revealing the Critical Regulators of Cell Identity in the Mouse Cell Atlas. *Cell Rep.* **25**, 1436–1445.e3 (2018).
215. Tasic, B. *et al.* Shared and distinct transcriptomic cell types across neocortical areas. *Nature* **563**, 72–78 (2018).
216. Allaway, K. C. *et al.* Genetic and epigenetic coordination of cortical interneuron development. *Nature* **597**, 693–697 (2021).
217. Jager, P. *et al.* Dual midbrain and forebrain origins of thalamic inhibitory interneurons. *Elife* **10**, (2021).
218. Batiuk, M. Y. *et al.* Identification of region-specific astrocyte subtypes at single cell resolution. *Nature Communications* vol. 11 (2020).
219. Liddelow, S. A. *et al.* Neurotoxic reactive astrocytes are induced by activated microglia. *Nature* **541**, 481–487 (2017).
220. Karran, E. & De Strooper, B. The amyloid cascade hypothesis: are we poised for success or failure? *Journal of Neurochemistry* vol. 139 237–252 (2016).
221. Consensus recommendations for the postmortem diagnosis of Alzheimer’s disease. The National Institute on Aging, and Reagan Institute Working Group on Diagnostic Criteria for the Neuropathological Assessment of Alzheimer’s Disease. *Neurobiol. Aging* **18**, S1–2 (1997).
222. Vázquez-Vélez, G. E. *et al.* Doublecortin-like Kinase 1 Regulates  $\alpha$ -Synuclein Levels and Toxicity. *J. Neurosci.* **40**, 459–477 (2020).
223. Yang, A. C. *et al.* A human brain vascular atlas reveals diverse mediators of Alzheimer’s risk. *Nature* (2022) doi:10.1038/s41586-021-04369-3.
224. Finak, G. *et al.* MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. *Genome Biol.* **16**, 278 (2015).
225. He, L. *et al.* NEBULA is a fast negative binomial mixed model for differential or co-expression analysis of large-scale multi-subject single-cell data. *Commun Biol* **4**, 629 (2021).
226. Aikawa, T., Holm, M.-L. & Kanekiyo, T. ABCA7 and Pathogenic Pathways of Alzheimer’s Disease. *Brain Sci* **8**, (2018).
227. Bis, J. *et al.* Correction: Whole exome sequencing study identifies novel rare and common Alzheimer’s-Associated variants involved in immune response and transcriptional regulation (Molecular Psychiatry,(2018), 10.1038/s41380-018-0112-7). *Mol. Psychiatry* (2019).
228. McCoy, M. K., Kaganovich, A., Rudenko, I. N., Ding, J. & Cookson, M. R. Hexokinase activity is required for recruitment of parkin to depolarized mitochondria. *Hum. Mol. Genet.* **23**, 145–156 (2014).
229. Zhang, J. & Ney, P. A. Role of BNIP3 and NIX in cell death, autophagy, and mitophagy. *Cell Death Differ.* **16**, 939–946 (2009).
230. Witter, M. P. & Moser, E. I. Spatial representation and the architecture of the entorhinal cortex. *Trends Neurosci.* **29**, 671–678 (2006).
231. Amunts, K. *et al.* Cytoarchitectonic mapping of the human amygdala, hippocampal region and entorhinal cortex: intersubject variability and probability maps. *Anat. Embryol.* **210**, 343–352 (2005).
232. Lussier, A. L., Weeber, E. J. & Rebeck, G. W. Reelin Proteolysis Affects Signaling Related to Normal Synapse Function and Neurodegeneration. *Front. Cell. Neurosci.* **10**, 75 (2016).
233. Liu, C.-C. *et al.* APOE3-Jacksonville (V236E) variant reduces self-aggregation and risk of dementia. *Sci. Transl. Med.* **13**, eabc9375 (2021).
234. Stoopschinski, B. E. *et al.* Specific glycosaminoglycan chain length and sulfation patterns are required for cell uptake of tau versus  $\alpha$ -synuclein and  $\beta$ -amyloid aggregates. *J. Biol. Chem.* **293**, 10826–10840 (2018).

235. Vagnoni, A. *et al.* Calsyntenin-1 mediates axonal transport of the amyloid precursor protein and regulates A $\beta$  production. *Hum. Mol. Genet.* **21**, 2845–2854 (2012).
236. Hoshino, A. *et al.* A novel function for proSAAS as an amyloid anti-aggregant in Alzheimer's disease. *J. Neurochem.* **128**, 419–430 (2014).
237. Wojtas, A. M. *et al.* Loss of clusterin shifts amyloid deposition to the cerebrovasculature via disruption of perivascular drainage pathways. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E6962–E6971 (2017).
238. Bennett, D. A. *et al.* Religious Orders Study and Rush Memory and Aging Project. *Journal of Alzheimer's Disease* vol. 64 S161–S189 (2018).
239. De Jager, P. L. *et al.* A multi-omic atlas of the human frontal cortex for aging and Alzheimer's disease research. *Sci Data* **5**, 180142 (2018).
240. Enge, M. *et al.* Single-Cell Analysis of Human Pancreas Reveals Transcriptional Signatures of Aging and Somatic Mutation Patterns. *Cell* **171**, 321–330.e14 (2017).
241. Kritsilis, M. *et al.* Ageing, Cellular Senescence and Neurodegenerative Disease. *Int. J. Mol. Sci.* **19**, (2018).
242. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet.* **46**, 310–315 (2014).
243. Savova, V. *et al.* Genes with monoallelic expression contribute disproportionately to genetic diversity in humans. *Nat. Genet.* **48**, 231–237 (2016).
244. Bennett, D. A. *et al.* Decision rules guiding the clinical diagnosis of Alzheimer's disease in two community-based cohort studies compared to standard practice in a clinic-based cohort study. *Neuroepidemiology* **27**, 169–176 (2006).
245. Swiech, L. *et al.* In vivo interrogation of gene function in the mammalian brain using CRISPR-Cas9. *Nat. Biotechnol.* **33**, 102–106 (2015).
246. Picelli, S. *et al.* Full-length RNA-seq from single cells using Smart-seq2. *Nat. Protoc.* **9**, 171–181 (2014).
247. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
248. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
249. Sherry, S. T. *et al.* dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res.* **29**, 308–311 (2001).
250. Mills, R. E. *et al.* An initial map of insertion and deletion (INDEL) variation in the human genome. *Genome Res.* **16**, 1182–1190 (2006).
251. Braak, H. & Braak, E. Neuropathological staging of Alzheimer-related changes. *Acta Neuropathol.* **82**, 239–259 (1991).
252. Bukar Maina, M., Al-Hilaly, Y. K. & Serpell, L. C. Nuclear Tau and Its Potential Role in Alzheimer's Disease. *Biomolecules* **6**, 9 (2016).
253. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**, 415–421 (2013).
254. Li, H. Toward better understanding of artifacts in variant calling from high-coverage samples. *Bioinformatics* **30**, 2843–2851 (2014).
255. Ramaswami, G. & Li, J. B. RADAR: a rigorously annotated database of A-to-I RNA editing. *Nucleic Acids Res.* **42**, D109–13 (2014).
256. Kiran, A. & Baranov, P. V. DARNED: a DAtabase of RNa EDiting in humans. *Bioinformatics* **26**, 1772–1776 (2010).
257. Lo Giudice, C., Tangaro, M. A., Pesole, G. & Picardi, E. Investigating RNA editing in deep transcriptome datasets with REDIttools and REDIportal. *Nat. Protoc.* **15**, 1098–1131 (2020).
258. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *BioRxiv* 30338 (2016).
259. Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff. *Fly* vol. 6 80–92 (2012).

260. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res.* **47**, D766–D773 (2019).
261. Rentzsch, P., Witten, D., Cooper, G. M., Shendure, J. & Kircher, M. CADD: predicting the deleteriousness of variants throughout the human genome. *Nucleic Acids Res.* **47**, D886–D894 (2019).
262. Rosenthal, R., McGranahan, N., Herrero, J., Taylor, B. S. & Swanton, C. DeconstructSigs: delineating mutational processes in single tumors distinguishes DNA repair deficiencies and patterns of carcinoma evolution. *Genome Biol.* **17**, 31 (2016).
263. Anders, S., Pyl, P. T. & Huber, W. HTSeq—a Python framework to work with high-throughput sequencing data. *Bioinformatics* **31**, 166–169 (2015).
264. Nott, A. *et al.* Brain cell type-specific enhancer-promoter interactome maps and disease-risk association. *Science* **366**, 1134–1139 (2019).
265. Suraweera, A. *et al.* Senataxin, defective in ataxia oculomotor apraxia type 2, is involved in the defense against oxidative DNA damage. *J. Cell Biol.* **177**, 969–979 (2007).
266. Campos-Melo, D., Hawley, Z. C. E. & Strong, M. J. Dysregulation of human NEFM and NEFH mRNA stability by ALS-linked miRNAs. *Mol. Brain* **11**, 43 (2018).
267. Yao, R.-Q., Ren, C., Xia, Z.-F. & Yao, Y.-M. Organelle-specific autophagy in inflammatory diseases: a potential therapeutic target underlying the quality control of multiple organelles. *Autophagy* 1–17 (2020) doi:10.1080/15548627.2020.1725377.
268. Tiwari, S. S. *et al.* Alzheimer-related decrease in CYFIP2 links amyloid production to tau hyperphosphorylation and memory loss. *Brain* **139**, 2751–2765 (2016).
269. Yan, J. *et al.* Hippocampal transcriptome-guided genetic analysis of correlated episodic memory phenotypes in Alzheimer’s disease. *Front. Genet.* **6**, 117 (2015).
270. Bruining, H. *et al.* Genetic Mapping in Mice Reveals the Involvement of Pcdh9 in Long-Term Social and Object Recognition and Sensorimotor Development. *Biol. Psychiatry* **78**, 485–495 (2015).
271. Silva, P. N. *et al.* CNP and DPYSL2 mRNA expression and promoter methylation levels in brain of Alzheimer’s disease patients. *J. Alzheimers. Dis.* **33**, 349–355 (2013).
272. Ousman, S. S. *et al.* Protective and therapeutic role for alphaB-crystallin in autoimmune demyelination. *Nature* **448**, 474–479 (2007).
273. Kuramoto, K., Negishi, M. & Katoh, H. Regulation of dendrite growth by the Cdc42 activator Zizimin1/Dock9 in hippocampal neurons. *J. Neurosci. Res.* **87**, 1794–1805 (2009).
274. Chen, X. *et al.* Novel schizophrenia risk factor pathways regulate FEZ1 to advance oligodendroglia development. *Transl. Psychiatry* **7**, 1293 (2017).
275. Goellner, B. & Aberle, H. The synaptic cytoskeleton in development and disease. *Dev. Neurobiol.* **72**, 111–125 (2012).
276. Kakinuma, T., Ichikawa, H., Tsukada, Y., Nakamura, T. & Toh, B.-H. Interaction between p230 and MACF1 is associated with transport of a glycosyl phosphatidyl inositol-anchored protein from the Golgi to the cell periphery. *Exp. Cell Res.* **298**, 388–398 (2004).
277. Zaoui, K., Benseddik, K., Daou, P., Salaün, D. & Badache, A. ErbB2 receptor controls microtubule capture by recruiting ACF7 to the plasma membrane of migrating cells. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 18517–18522 (2010).
278. Cha, M.-Y. *et al.* Mitochondrial ATP synthase activity is impaired by suppressed O-GlcNAcylation in Alzheimer’s disease. *Hum. Mol. Genet.* **24**, 6492–6504 (2015).
279. Seidel, K. *et al.* The HSPB8-BAG3 chaperone complex is upregulated in astrocytes in the human brain affected by protein aggregation diseases. *Neuropathol. Appl. Neurobiol.* **38**, 39–53 (2012).
280. Dabir, D. V., Trojanowski, J. Q., Richter-Landsberg, C., Lee, V. M.-Y. & Forman, M. S. Expression of the small heat-shock protein alphaB-crystallin in tauopathies with glial pathology. *Am. J. Pathol.* **164**, 155–166 (2004).
281. Bertrand, P., Poirier, J., Oda, T., Finch, C. E. & Pasinetti, G. M. Association of

- apolipoprotein E genotype with brain levels of apolipoprotein E and apolipoprotein J (clusterin) in Alzheimer disease. *Brain Res. Mol. Brain Res.* **33**, 174–178 (1995).
282. Bass, A. & Storey, J. *qvalue R package*. (Github).
283. Raudvere, U. *et al.* g:Profiler: a web server for functional enrichment analysis and conversions of gene lists (2019 update). *Nucleic Acids Res.* **47**, W191–W198 (2019).
284. Jassal, B. *et al.* The reactome pathway knowledgebase. *Nucleic Acids Res.* **48**, D498–D503 (2020).
285. Ligtenberg, W. reactome. db: A set of annotation maps for reactome. *R package version 1*, (2018).
286. Lodato, M. A. *et al.* Somatic mutation in single human neurons tracks developmental and transcriptional history. *Science* **350**, 94–98 (2015).
287. Oksanen, M. *et al.* PSEN1 Mutant iPSC-Derived Model Reveals Severe Astrocyte Pathology in Alzheimer's Disease. *Stem Cell Reports* **9**, 1885–1897 (2017).
288. Nasrabady, S. E., Rizvi, B., Goldman, J. E. & Brickman, A. M. White matter changes in Alzheimer's disease: a focus on myelin and oligodendrocytes. *Acta Neuropathologica Communications* vol. 6 (2018).
289. Matsuyama, S. S. & Jarvik, L. F. Hypothesis: microtubules, a key to Alzheimer disease. *Proc. Natl. Acad. Sci. U. S. A.* **86**, 8152–8156 (1989).
290. Weingarten, M. D., Lockwood, A. H., Hwo, S. Y. & Kirschner, M. W. A protein factor essential for microtubule assembly. *Proc. Natl. Acad. Sci. U. S. A.* **72**, 1858–1862 (1975).
291. Grundke-Iqbal, I. *et al.* Abnormal phosphorylation of the microtubule-associated protein tau (tau) in Alzheimer cytoskeletal pathology. *Proc. Natl. Acad. Sci. U. S. A.* **83**, 4913–4917 (1986).
292. Madabhushi, R. *et al.* Activity-Induced DNA Breaks Govern the Expression of Neuronal Early-Response Genes. *Cell* vol. 161 1592–1605 (2015).
293. Cruz, J. C., Tseng, H.-C., Goldman, J. A., Shih, H. & Tsai, L.-H. Aberrant Cdk5 activation by p25 triggers pathological events leading to neurodegeneration and neurofibrillary tangles. *Neuron* **40**, 471–483 (2003).
294. Kim, D. *et al.* Deregulation of HDAC1 by p25/Cdk5 in neurotoxicity. *Neuron* **60**, 803–817 (2008).
295. Mathys, H. *et al.* Temporal Tracking of Microglia Activation in Neurodegeneration at Single-Cell Resolution. *Cell Reports* vol. 21 366–380 (2017).
296. Cruz, J. C. *et al.* p25/Cyclin-Dependent Kinase 5 Induces Production and Intraneuronal Accumulation of Amyloid beta In Vivo. *Journal of Neuroscience* vol. 26 10536–10541 (2006).
297. Fischer, A., Sananbenesi, F., Pang, P. T., Lu, B. & Tsai, L.-H. Opposing roles of transient and prolonged expression of p25 in synaptic plasticity and hippocampus-dependent memory. *Neuron* **48**, 825–838 (2005).
298. Mah, L.-J., El-Osta, A. & Karagiannis, T. C.  $\gamma$ H2AX: a sensitive molecular marker of DNA damage and repair. *Leukemia* **24**, 679–686 (2010).
299. Caldwell, A. B. *et al.* Dedifferentiation and neuronal repression define familial Alzheimer's disease. *Science Advances* vol. 6 (2020).
300. Mertens, J. *et al.* Age-dependent instability of mature neuronal fate in induced neurons from Alzheimer's patients. *Cell Stem Cell* vol. 28 1533–1548.e6 (2021).
301. Yousef, A. *et al.* Neuron loss and degeneration in the progression of TDP-43 in frontotemporal lobar degeneration. *Acta Neuropathol Commun* **5**, 68 (2017).
302. Unal-Cevik, I., Kiliç, M., Gürsoy-Ozdemir, Y., Gurer, G. & Dalkara, T. Loss of NeuN immunoreactivity after cerebral ischemia does not indicate neuronal cell loss: a cautionary note. *Brain Res.* **1015**, 169–174 (2004).
303. McPhail, L. T., McBride, C. B., McGraw, J., Steeves, J. D. & Tetzlaff, W. Axotomy abolishes NeuN expression in facial but not rubrospinal neurons. *Exp. Neurol.* **185**, 182–190 (2004).



304. Wu, K.-L. *et al.* Loss of neuronal protein expression in mouse hippocampus after irradiation. *J. Neuropathol. Exp. Neurol.* **69**, 272–280 (2010).
305. Collombet, J.-M. *et al.* Early reduction of NeuN antigenicity induced by soman poisoning in mice can be used to predict delayed neuronal degeneration in the hippocampus. *Neurosci. Lett.* **398**, 337–342 (2006).
306. Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 15545–15550 (2005).
307. Yang, Y., Geldmacher, D. S. & Herrup, K. DNA replication precedes neuronal cell death in Alzheimer's disease. *J. Neurosci.* **21**, 2661–2668 (2001).
308. Busser, J., Geldmacher, D. S. & Herrup, K. Ectopic cell cycle proteins predict the sites of neuronal cell death in Alzheimer's disease brain. *J. Neurosci.* **18**, 2801–2807 (1998).
309. Klein, R. S. *et al.* Neuronal CXCL10 directs CD8+ T-cell recruitment and control of West Nile virus encephalitis. *J. Virol.* **79**, 11457–11466 (2005).
310. Préhaud, C., Mégret, F., Lafage, M. & Lafon, M. Virus infection switches TLR-3-positive human neurons to become strong producers of beta interferon. *J. Virol.* **79**, 12893–12904 (2005).
311. Sokolova, A. *et al.* Monocyte chemoattractant protein-1 plays a dominant role in the chronic inflammation observed in Alzheimer's disease. *Brain Pathol.* **19**, 392–398 (2009).
312. El Khoury, J. *et al.* Ccr2 deficiency impairs microglial accumulation and accelerates progression of Alzheimer-like disease. *Nat. Med.* **13**, 432–438 (2007).
313. Chai, Q., She, R., Huang, Y. & Fu, Z. F. Expression of neuronal CXCL10 induced by rabies virus infection initiates infiltration of inflammatory cells, production of chemokines and cytokines, and enhancement of blood-brain barrier permeability. *J. Virol.* **89**, 870–876 (2015).
314. Conductier, G., Blondeau, N., Guyon, A., Nahon, J.-L. & Rovère, C. The role of monocyte chemoattractant protein MCP1/CCL2 in neuroinflammatory diseases. *J. Neuroimmunol.* **224**, 93–100 (2010).
315. Kuleshov, M. V. *et al.* Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res.* **44**, W90–7 (2016).
316. Janssens, S. & Tschopp, J. Signals from within: the DNA-damage-induced NF- $\kappa$ B response. *Cell Death Differ.* **13**, 773–784 (2006).
317. McCool, K. W. & Miyamoto, S. DNA damage-dependent NF- $\kappa$ B activation: NEMO turns nuclear signaling inside out. *Immunol. Rev.* **246**, 311–326 (2012).
318. Chien, Y. *et al.* Control of the senescence-associated secretory phenotype by NF- $\kappa$ B promotes senescence and enhances chemosensitivity. *Genes & Development* vol. 25 2125–2136 (2011).
319. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nature Biotechnology* vol. 32 381–386 (2014).
320. Wang, C. *et al.* Gain of toxic apolipoprotein E4 effects in human iPSC-derived neurons is ameliorated by a small-molecule structure corrector. *Nature Medicine* vol. 24 647–657 (2018).
321. Haas, B. *et al.* STAR-Fusion: Fast and Accurate Fusion Transcript Detection from RNA-Seq. *BioRxiv* (2017) doi:10.1101/120295.
322. Liu, F. *et al.* Systematic comparative analysis of single-nucleotide variant detection methods from single-cell RNA sequencing data. *Genome Biol.* **20**, 242 (2019).
323. Street, K. *et al.* Slingshot: cell lineage and pseudotime inference for single-cell transcriptomics. *BMC Genomics* vol. 19 (2018).
324. Wolf, F. A. *et al.* PAGA: graph abstraction reconciles clustering with trajectory inference through a topology preserving map of single cells. *Genome Biology* vol. 20 (2019).
325. Johnson, W. E., Evan Johnson, W., Li, C. & Rabinovic, A. Adjusting batch effects in

- microarray expression data using empirical Bayes methods. *Biostatistics* vol. 8 118–127 (2007).
326. García-Osta, A., Cuadrado-Tejedor, M., García-Barroso, C., Oyarzábal, J. & Franco, R. Phosphodiesterases as therapeutic targets for Alzheimer's disease. *ACS Chem. Neurosci.* **3**, 832–844 (2012).
  327. Burgess, C. P. *et al.* Understanding disentangling in  $\beta$ -VAE. *arXiv [stat.ML]* (2018).
  328. Lotfollahi, M. *et al.* Compositional perturbation autoencoder for single-cell response modeling. *bioRxiv* 2021.04.14.439903 (2021) doi:10.1101/2021.04.14.439903.
  329. Replogle, J. M. *et al.* Mapping information-rich genotype-phenotype landscapes with genome-scale Perturb-seq. *bioRxiv* 2021.12.16.473013 (2021) doi:10.1101/2021.12.16.473013.
  330. Dixit, A. *et al.* Perturb-Seq: Dissecting Molecular Circuits with Scalable Single-Cell RNA Profiling of Pooled Genetic Screens. *Cell* **167**, 1853–1866.e17 (2016).
  331. Adamson, B. *et al.* A Multiplexed Single-Cell CRISPR Screening Platform Enables Systematic Dissection of the Unfolded Protein Response. *Cell* **167**, 1867–1882.e21 (2016).
  332. Mitchell, J. M. *et al.* Mapping genetic effects on cellular phenotypes with 'cell villages'. *bioRxiv* 2020.06.29.174383 (2020) doi:10.1101/2020.06.29.174383.
  333. Ursu, O. *et al.* Massively parallel phenotyping of variant impact in cancer with Perturb-seq reveals a shift in the spectrum of cell states induced by somatic mutations. *bioRxiv* 2020.11.16.383307 (2020) doi:10.1101/2020.11.16.383307.
  334. Yazar, S. *et al.* Single-cell eQTL mapping identifies cell type-specific genetic control of autoimmune disease. *Science* **376**, eabf3041 (2022).
  335. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).
  336. Avsec, Ž. *et al.* Effective gene expression prediction from sequence by integrating long-range interactions. *Nat. Methods* **18**, 1196–1203 (2021).
  337. Pearson, K. LIII. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* **2**, 559–572 (1901).
  338. Leek, J. T. svaseq: removing batch effects and other unwanted noise from sequencing data. *Nucleic Acids Res.* gku864 (2014).
  339. McCarthy, D. J., Campbell, K. R., Lun, A. T. L. & Wills, Q. F. Scater: pre-processing, quality control, normalization and visualization of single-cell RNA-seq data in R. *Bioinformatics* **33**, 1179–1186 (2017).
  340. Tibshirani, R. Variance stabilization and the bootstrap. *Biometrika* **75**, 433–444 (1988).
  341. Vaswani, A. *et al.* Attention is All you Need. *Adv. Neural Inf. Process. Syst.* **30**, (2017).
  342. Pan, S. J. & Yang, Q. A Survey on Transfer Learning. *IEEE Trans. Knowl. Data Eng.* **22**, 1345–1359 (2010).