

Single-cell differential splicing of Alzheimer’s Disease in 1.9 million cells across 416 individuals

by

Christian Hwa

B.S. Computer Science and Engineering
Massachusetts Institute of Technology, 2021

Submitted to the Department of Electrical Engineering and Computer
Science

in partial fulfillment of the requirements for the degree of

Master of Engineering in Electrical Engineering and Computer Science

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

September 2022

© Massachusetts Institute of Technology 2022. All rights reserved.

Author
Department of Electrical Engineering and Computer Science
August 5, 2022

Certified by.....
Manolis Kellis
Head, MIT Computational Biology Group
Principal Investigator, Computer Science and Artificial Intelligence Lab
Thesis Supervisor

Accepted by
Katrina LaCurts
Chair, Master of Engineering Thesis Committee

Single-cell differential splicing of Alzheimer’s Disease in 1.9 million cells across 416 individuals

by

Christian Hwa

Submitted to the Department of Electrical Engineering and Computer Science
on August 5, 2022, in partial fulfillment of the
requirements for the degree of
Master of Engineering in Electrical Engineering and Computer Science

Abstract

Alzheimer’s Disease (AD) is a polygenic disease with variable phenotypic response among patients. While single-cell differential expression data can provide important information regarding the disease’s pathology, it only captures steady-state cell behavior. By contrast, RNA velocity analysis allows inference of both direction and magnitude of regulatory changes occurring over disease progression. However, despite the success of existing RNA velocity analyses, the pace of disease progression remains uncharacterized, as previous work focused primarily on macro-scale trajectory diagrams and failed to capture trajectories of individual genes relative to global pathological changes.

Here, we use scRNA-seq profiling of 1.9 million cells from dorsolateral prefrontal cortex of 416 post-mortem human samples to model temporal regulatory changes of AD. We use RNA velocity to infer spliced and unspliced counts for each gene and a beta-binomial distribution to model the unspliced/spliced ratio, and we identify differentially-spliced genes across six brain cell types and differentially-spliced genes between AD and control cohorts. Of 16,948 expressed genes, 1158 genes show significant differential splicing (DS) effects with $FDR < 0.05\%$. We find that cell-type specific and globally-altered genes are significantly enriched in AD-associated processes, including dysregulation of the TNF pathway in astrocytes, dysregulation of PD-L1 expression across all cell-types, cellular transport, and lipid dyshomeostasis.

Overall, our study provides a novel method to analyze and compare the RNA velocity between AD and non-AD cohorts, and identifies a new set of differentially-spliced AD genes, when previous analyses focused primarily on gene expression changes.

Thesis Supervisor: Manolis Kellis

Title: Head, MIT Computational Biology Group

Principal Investigator, Computer Science and Artificial Intelligence Lab

Acknowledgments

I would like to thank my mom and dad for their unwavering love, support, and patience. I would also like to thank my sister for always being there; and Chewie and Zeus for their support. Everything I've done would not be possible without them.

I would also like to thank Dr. Liang He for being my mentor throughout the year. I am very grateful to have been able to work under you; this thesis would not be possible without your insights and guidance.

I would like to thank Dr. Yongjin Park for his guidance as well. This thesis would not have been finished without your help.

Special mention to Dr. Li-lun Ho for his efforts in data generation and experimental work. I would also like to thank all members of the Kellis lab for their openness and advice, especially Patty who made this all possible.

Finally, I would like to thank my advisor, Prof. Manolis Kellis, for his guidance and support. The work the group is doing is life changing and I am honored to have been given this opportunity.

Contents

1	Introduction	17
1.1	Neuropathology	18
1.2	Amyloid-beta accumulation	19
1.3	Single-cell RNA-seq	20
1.4	RNA Splicing and Velocity	21
1.5	Related Works	22
2	Aims	25
2.1	Identifying a set of differential splicing genes	25
2.2	Cell-type specific analysis and meta-analysis	26
2.3	Pathway enrichment analysis	26
2.4	Differential kinetics analysis	26
3	Dataset and Processing	27
3.1	scRNA ROSMAP Dataset	27
3.2	scRNA-seq processing	28
3.3	Dataset quality control and cell type differentiation	29
4	Beta-binomial regression reveals differentially spliced genes	33
4.1	Splicing count inference	33
4.2	Splicing analysis approach	34
4.2.1	Beta binomial distribution	34
4.2.2	Meta-analysis aggregation	34

4.3	Global splicing characteristics do not differ between AD and control cohorts	36
4.4	Beta-binomial regression reveals cell-type specific differentially spliced genes	38
4.4.1	Enrichment of SRRM2, CRY1, and ASPDH1 in Astrocytes . .	39
4.4.2	TRIM9, SRRM2, and COL11A1 differential splicing in excitatory neurons	41
4.4.3	Differential splicing of Z3CHAV1, XRCC4, and LPCAT2 in microglia	43
4.5	Meta-analysis reveals differential splicing of SRRM2, GOLGA4, and VPS13A	44
4.5.1	SRRM2 splicing dysregulation in AD	45
4.5.2	Differential splicing of GOLGA4 highlights the role of lipid transport in AD pathogenesis	48
4.5.3	Differential splicing of VPS13A reveals role of mitochondria in AD pathogenesis	49
5	Random priming effects in isoform analysis	51
6	Pathway enrichment analysis identifies AD-related enriched pathways from differentially spliced genes	53
6.1	Methodological Approach	53
6.2	Analysis reveals enrichment of spliceosome, homologous recombination, and reactive oxygen species in AD astrocytes	54
6.3	Pathway analysis reveals enrichment of TNF and PD-L1 in AD microglia populations	55
6.4	PD-L1, endocytotic, and insulin resistance pathways are enriched in excitatory neurons	57
6.5	Meta-analysis reveals enrichment of PD-L1, NF-kappa B, and adherens junction and close association with Parkinsons and prion pathways .	58

7	RNA velocity analysis reveals sub-population AD-enrichment	61
7.1	RNA velocity inference	61
7.2	AD cohort express less cellular sub-population diversity in astrocytes	62
7.3	Differential splicing genes show sub-population specific enrichment in AD cohorts	63
7.4	Astrocyte cluster specific analysis of potential drive genes	65
7.5	AD microglia cells exhibit less cluster diversity	66
7.6	Cluster localization differential between velocity and expression in OXR1 and PNISR	67
7.7	Microglia cluster specific AD driver genes	68
8	Conclusion	71
8.1	Summary	71
8.2	Future Work	72

List of Figures

3-1	Breakdown of the scRNA ROSMAP dataset as well as summary statistics regarding surrounding the dataset	28
3-2	Breakdown of the steps involved for removing doublets [5]	30
3-3	UMAP of all the cell types in the scRNA ROSMAP dataset, B: individual breakdown of the number of cells in each cell type	31
4-1	Global spliced/unspliced counts for excitatory, microglia, oligodendrocyte, and inhibitory cell types. Points that are colored red indicate control, while points that are colored blue indicate AD individuals . .	36
4-2	Box plot comparing percentage of spliced reads between AD and non-AD cohorts across 5 cell types	37
4-3	Differentially spliced genes between AD and control groups in astrocytes. Volcano plot of differentially spliced genes in astrocytes. Adjusted $p < 0.0001$	39
4-4	Hypothesized model of sleep disturbance and onset of AD	41
4-5	Differentially spliced genes between AD and control groups in excitatory neurons. Volcano plot of differentially spliced genes in astrocytes. Adjusted $p < 0.0001$	42
4-6	Differentially spliced genes between AD and control groups in microglia. Volcano plot of differentially spliced genes in astrocytes. Adjusted $p < 0.0005$	44
4-7	Raw spliced/unspliced counts for SRRM2 across astrocytes, excitatory, inhibitory, and microglia cell types	47

4-8	Raw spliced/unspliced counts for GOLGA4 across astrocytes, excitatory, inhibitory, and microglia cell types	48
4-9	Raw spliced/unspliced counts for VPS13A across astrocytes, excitatory, inhibitory, and microglia cell types	49
5-1	Scatterplot showing the relationship between gene body length and percentage of spliced reads. Highpercshortlen are genes where at least a 95% of the gene body is composed of exons	52
6-1	Top 10 enriched pathways in astrocytes given the differentially spliced astrocyte gene set (FDR adjusted p-values), significance and enrichment thresholds were both equal to 0.05	54
6-2	Top 10 enriched pathways in microglia given differentially spliced microglia gene set, significance and enrichment thresholds were both equal to 0.05. FDR adjusted p-values	56
6-3	The top 10 enriched pathways in excitatory neurons given the differentially spliced excitatory gene set as input, significance and enrichment thresholds were both equal to 0.05 with FDR adjusted p-values . . .	57
6-4	Top 10 enriched pathways across all cell types with FDR adjusted p-values, significance and enrichment thresholds were both equal to 0.05	59
7-1	Velocity vectors overlaid onto the cell's UMAP space for astrocytes. Input gene set was the top 2000 differentially spliced genes in astrocytes. A: vector field for non-AD cohorts. B: vector field for AD individuals	63
7-2	Gene specific analysis for SRRM2, GOLGA4, XRCC4, OXR1, and PNISR in astrocytes. Each gene has a spliced/unspliced plot, velocity vector field overlaid in UMAP space, and expression patterns mapped in UMAP space. A: plots for the control cohort. B: plots for AD individuals	65
7-3	Cluster specific driver genes for astrocytes in AD individuals. Computed by finding the partial gene likelihoods for each cluster.	66

7-4	Velocity vectors overlaid onto the cell's UMAP space for microglia. Input gene set was the top 2000 differentially spliced genes in microglia. A: vector field for control cohorts. B: vector field for AD individuals .	67
7-5	Gene specific analysis for SRRM2, GOLGA4, XRCC4, OXR1, and PNISR in microglia. Each gene has a spliced/unspliced plot, velocity vector field overlaid in UMAP space, and expression patterns mapped in UMAP space. A: plots for the control cohort. B: plots for AD individuals	68
7-6	Cluster specific driver genes for microglia in AD individuals. Computed by finding the partial gene likelihoods for each cluster.	69

List of Tables

4.1	Table representing the top 16 globally differentially spliced genes . . .	45
-----	---	----

Chapter 1

Introduction

Alzheimer's disease (AD) is a devastating neurodegenerative disease with progressive memory loss, motor disorders, and behavioral changes. These alterations of the patient's physiology are cellularly hallmarked by the accumulation of senile amyloid-beta plaques and neurofibrillary tangles. AD accounts for nearly 60-70% of cases of dementia, with typical life expectancy after diagnosis being three to nine years depending on the rate of progression. AD manifests itself into two broad categories: early-onset AD (EOAD) and late-onset AD (LOAD). Typically, patients diagnosed with AD ages 65 and younger are categorized as EOAD, while patients older than 65 years are categorized as LOAD. EOAD only represents a small fraction of AD cases, typically 5-6% of people with AD develop symptoms before the age of 65 [1].

As with most neurodegenerative diseases, AD progresses through multiple stages: first symptoms, early-stage, middle-stage, and late-stage. First symptoms are often mistaken for other causes such as stress. They manifest similar symptoms of mild cognitive impairment such as short-term memory loss, trouble learning new information, and acute impairment with abstract thinking. Early-stage symptoms include increased impairments towards learning and memory, as well as agnosia and apraxia. Patients also exhibit decreased vocabulary and impairment towards oral and written communication. In the middle-stage, patients are unable to perform simple daily tasks. Communication becomes difficult with impaired memory and vocabulary. Typically, it is at this stage where patients are unable to recognize family

members. Finally, in late-stage AD, patients are completely dependent on others to live. Patients typically experience complete loss of speech and are bedridden with a permanent caregiver. Death is generally not associated with the disease, but rather other factors such as pneumonia or infection. Furthermore, there have been reported cases of patients that exhibit early-stage symptoms before the typical first symptoms, highlighting the difficulty in diagnosing AD early.

Unfortunately, modern medicine does not have a cure or even a way to slow down the neurodegeneration caused by AD. Preventative measures against AD are still nascent in terms of results. While studies have attempted to identify links between AD and diet or lifestyle, there is difficulty in controlling for variables such as disease progression, comorbidities, and a lack of a standardized AD-specific biomarker.

1.1 Neuropathology

AD is characterized neuropathologically by neuronal and synaptic loss in the cerebral cortex and other regions. This degradation is easily observed using MRI and PET devices. There is a significant difference in sizes of specific brain regions between AD and non-AD samples. Furthermore, there is an observed accumulation of neurofibrillary tangles (NFTs) and amyloid-beta plaques in many brain regions that ultimately lead to the pathogenic cascade of AD [36]. NFTs are aggregates of hyperphosphorylated tau proteins. While they are a primary biomarker for AD, we still do not fully understand the role it plays in AD progression. Amyloid-beta plaques are mostly composed of amyloid-beta peptides, generally 36-43 amino acids in length. Amyloid-beta peptides are the byproduct of cleaving the amyloid precursor protein (APP). These peptides can aggregate and form variable oligomers. High concentrations of oligomers can induce misfolding in adjacent oligomers, instigating a misfolding cascade [36]. Interestingly the propagation of NFTs and amyloid-beta plaques are diametrically opposed: NFTs generally start at in subcortical regions before progressing to cortical regions, while amyloid-beta plaques generally begin accumulating in the cortical regions before moving to subcortical regions.

1.2 Amyloid-beta accumulation

While we do not have a definitive understanding of AD pathogenesis, one of the leading theories that researchers have supported for the past couple decades has been the amyloid cascade hypothesis: the deposition of amyloid-beta peptides in the brain triggers a pathogenic cascade that eventually leads to AD onset. The discovery of the APP mutation that led to the discovery of the amyloid-beta plaque cycle served as the foundation for the hypothesis. The hypothesis was further supported when studies discovered that autosomal dominant mutations in PSEN1 and PSEN2 led to forming catalytic active sites for gamma secretase [7]. Furthermore, one of the most known risk factors for AD is apolipoprotein E (APOE), which has three different alleles APOE2, APOE3, APOE4. Possession of each of the three alleles lead to drastically different genetic preconditions: APOE2 is thought to be protective, APOE3 is neutral, and APOE4 can decrease the age of onset for AD. While PDAPP transgenic mice studies have shown the APOE directly affects pathways related to the clearance of amyloid-beta, we still do not have a complete understanding on how the different APOE proteins interact with AD pathogenesis [19].

The amyloid cascade hypothesis appears to do a sufficient job in explaining AD progression, however, the hypothesis does not capture the interaction between amyloid-beta peptides and tau proteins. Specifically, certain neurodegenerative diseases share the same tau pathology as AD without the accumulation of amyloid-beta plaques, suggesting that tau proteins play a role in cellular loss as well. Studies in transgenic mice have also shown that a reduction in level of tau can improve behavioral issues mediated by amyloid-beta plaques. However, applying those studies to human AD remains a difficulty. Regardless, understanding the pathways that lead to amyloid-beta peptide deposition and the cellular responses remain important goals in our goal towards understanding AD.

1.3 Single-cell RNA-seq

When studying complex diseases, researchers need to be able to create a molecular view that allows them to quantify cellular responses. Single-cell next generation sequencing technology such solutions, and has allowed for unprecedented understanding into the cellular changes that occur. Single cell RNA enables scientists to circumvent the scale of studying the entire proteome, all proteins expressed in a single cell, by studying the protein encoding, mRNA transcripts, that ultimately lead to the assembly of said proteins [3].

Initially, RNA-sequencing (RNA-seq) was the genomic main assay that allowed for direct quantification and detection of messenger RNA (mRNA) molecules in a sample. Unfortunately, for effective RNA-seq, millions of cells need to be provided and pooled before further analysis [15]. This pooling step involves averaging across the cellular population, thus losing the ability for scientists to understand the cellular responses at the single cell level.

Single-cell RNA-seq (scRNA-seq) directly addresses that problem. scRNA-seq allows for studies to compare transcriptomes of different individuals and populations with single cell resolution. Thus, in the context of studying AD, this allows us to compare the behaviors of single cells between AD and non-AD individuals, hopefully leading to the identification of key regulatory systems that are disrupted. Furthermore, scRNA-seq also allows for the identification of rare cell subpopulations, potentially leading to the discovery of new intermediary cell states. As mentioned above, AD propagates through affected neuronal populations, raising the question: during pathogenesis, are there any temporary states that cells enter before death? With scRNA-seq, we can now hopefully answer such questions.

Although there are a myriad of different scRNA-seq techniques, they all follow the same basic principles. First, individual cells from the tissue sample are isolated using different techniques such as FACS, microfluidics, and dielectrophoretic digital sorting. Next, the cells are lysed while preserving mRNA transcripts. To capture specific mRNA transcripts, poly-T tail primers are normally used. This enables the

assay to capture polyadenylated RNA transcripts instead of ribosomal RNA. Next, reverse transcriptase is used to convert the poly-T primed mRNA transcript to cDNA [3]. In this step, it is common for assays to also include a unique molecular identifier (UMI) to each transcript, enabling single RNA transcript resolution as well as tracking metadata such as cellular origin. cDNA is then amplified using PCR, before being pooled and sequenced by next generation sequencing (NGS). Downstream analysis can then be performed by using different software packages for processing such as Seurat or CellRanger.

1.4 RNA Splicing and Velocity

While using scRNA-seq to study gene co-expression patterns can yield interesting results, gene expression counts merely represent a single snapshot in time. We are not able to track the series of state and associated regulatory changes that occur after observation. Instead, we can use RNA velocity, an estimation of the time derivative of the gene expression state, to predict the future cell states. RNA velocity can be directly estimated by classifying mRNA transcripts in scRNA-seq protocols as either spliced or unspliced, where the relative abundance of each type can give insight into the regulatory direction of the cell. The unspliced mRNA count represents an indicator of the spliced mRNA content in the cell's future states [23]. This basic relationship allows us to identify patterns of expression amongst different pathways where a decrease in unspliced mRNA represents a downregulation, whereas an increase in unspliced mRNA represents an upregulation. This can be done without an a priori information regarding cellular relationships. As a result, the cellular dynamics and future direction of cellular progression can be reconstructed by combining snapshot data with a causal model.

The abundance of spliced and unspliced mRNA counts can be modeled by a system of ordinary differential equations. Given transcription rate $\alpha(t)$, the splicing rate β , and degradation rate γ , the instantaneous average for unspliced counts, μ , and spliced, s , can be derived by the following [11]:

$$\frac{d\mu_u(t)}{dt} = \alpha(t) - \beta\mu_u(t)$$

$$\frac{d\mu_s(t)}{dt} = \beta\mu_u(t) - \gamma\mu_s(t)$$

Equation 1 abstracts the splicing process to a simple equation that intuitively makes sense. Unspliced counts are modeled as the transcription rate subtracted by the splicing rate, while the spliced counts are modeled by the splicing rate subtracted by the degradation rate [11]. Ultimately, this system of equation allows us to estimate degradation to splicing ratio, γ/β . Finally, with this ratio, we can estimate the velocity of given cell by the following [11]:

$$v_i := \frac{\Delta s_i}{\Delta t} = (\beta u_i - \gamma s_i) \propto u_i - \frac{\gamma}{\beta} s_i$$

Where s_i and u_i are cell specific counts, Δt is the time differential, and Δs_i is the change in spliced counts over the Δt increment.

Leveraging the ability to differentiate between unspliced and spliced mRNA transcripts at the single cell level, we can begin to dig deeper into the association between AD and altered splicing. Alternative splicing is an important posttranscriptional regulatory mechanism, allowing the human body to create a diverse set of complex protein structures. Studies have identified that in patients with AD, there is an unusual aggregation of pieces of spliceosomal complex nuclear RNA, suggesting that the main splicing machinery may be damaged in AD. As a result, there is great interest in better understanding how splicing is disrupted in AD, with the hopes of unearthing new information surrounding the disease’s pathogenesis.

1.5 Related Works

To investigate AD pathogenesis utilizing RNA velocity, I will be using a beta-binomial regression as the model to identify differential velocity genes. To our knowledge, the beta-binomial distribution has never been applied on single cell data in such a manner.

However, it has been applied for differential methylation analysis in DNA. Historically, when studying DNA methylation, one can use Fisher's exact test, algorithm-based smoothing, or even hidden Markov models; but they all lack the ability to account for variability between samples. The beta-binomial distribution accounts for variability by ensuring that the probability of success is not a fixed value, instead, it is randomly drawn from a beta distribution [8].

More formally, the beta-binomial can be obtained from $Binom(p, n)$ by assuming $p \sim Beta(\alpha, \beta)$. We can then construct the following probability mass function

$$Pr(M = m|n, \alpha, \beta) = \binom{n}{m} \frac{B(m + \alpha, n - m + \beta)}{B(\alpha, \beta)}$$

Where B is the beta function. If we let:

$$\pi = \frac{\alpha}{(\alpha + \beta)}, \gamma = \frac{1}{(\alpha + \beta + 1)}$$

then

$$E(M) = n\pi$$

$$Var(M) = n\pi(1 - \pi)(1 + (n - 1)\gamma)$$

Where π is the analog of the probability of success variable while γ represents the dispersion parameter [8].

To fit our scRNA-seq data to a beta-binomial distribution, I will model a spliced count as a success and an unspliced count as a failure. By comparing this spliced unspliced ratio for each gene in each cell type for every individual, we can hopefully identify a new set of genes that are differentially spliced.

Chapter 2

Aims

The main aim of this thesis is to be able to better understand the pathogenesis of AD through a differential velocity perspective. To our understanding, this type of analysis has never been done before. Achieving this main aim can be broken down into four steps:

1. Identifying a set of differential splicing genes
2. Cell-type specific analysis and meta-analysis
3. Pathway enrichment analysis
4. Differential kinetics analysis

2.1 Identifying a set of differential splicing genes

To identify a set of differential splicing genes, we implemented and tested a beta-binomial distribution model. Drawing inspiration from identifying DNA-methylation sites as mentioned in the related works section, we treat each cell as a sample, such that we can compare the ratio between spliced and unspliced counts. We use a beta-binomial distribution, as opposed to a naïve binomial distribution due to the ability for the beta-binomial distribution to account for noise.

Currently, single-cell protocols require amplification so that small amounts of mRNA can be detected. This step is typically performed using PCR to amplify

cDNA molecules that result from the reverse transcription of the poly-adenylated RNA molecules. These steps, in combination, can result in significant technical noise. As a result, it is important to having a model that accounts for this.

2.2 Cell-type specific analysis and meta-analysis

After validating the model from above, we apply the regression at the cell type specific layer. Enabling us to rank all 16,948 genes by their effect size. Meta-analysis can then be performed by aggregating all cell type specific rankings, enabling us to rank all genes across all cell types.

2.3 Pathway enrichment analysis

Using pathfindR, we performed pathway enrichment analysis from the sorted gene list. We performed this pathway enrichment analysis at the cell-type specific level as well as the aggregated meta-cell, thus enabling us to not only identify cell-type specific disrupted pathways, but also identify pathways that are disrupted across all cell-types in AD patients. We also use this step as a simplified validation process in that we ensure the top pathways disrupted are known, or upstream of, AD pathways.

2.4 Differential kinetics analysis

Here, we perform classical differential kinetics analysis to further understand how our differential splicing genes differ in AD cases compared to control. Using scVelo, we aim to find a dynamical model with the hopes of learning the full transcriptional dynamics of the splicing kinetics. We also aim to use different clustering algorithms to identify subclusters within specific cell-types, with the goal of identifying unique subclusters in AD cases. Furthermore, in addition to unique clustering, we leveraged the directionality of a velocity vector to identify differences in transcriptional direction in specific cell-types.

Chapter 3

Dataset and Processing

In the following section, we will introduce the single-cell Religious Order Study/Memory and Aging Project (ROSMAP) dataset we used for this study. The ROS dataset is a clinical-pathological study of memory, motor, and functional problems in religious Catholic individuals, both males and females, aged 65 years and older in the United States. On the other hand, the MAP is an epidemiological clinical-pathological study of dementia and other chronic diseases associated with aging in older persons. In ROSMAP, all participants agree to yearly blood draws as well as detailed clinical evaluation.

3.1 scRNA ROSMAP Dataset

In our single-cell ROSMAP dataset, we sequenced 427 individuals from ROSMAP, of which, we sampled 209 individuals from ROS and 218 individuals from MAP, composed of 215 Males and 212 Females. Of the 427 individuals, 238 (55.7%) were clinically diagnosed with AD while 189 (44.3%) were not diagnosed with AD. All samples were acquired post-mortem in the prefrontal cortex (PFC) (Figure 3-1). Note that after processing, our final dataset cohort was composed of 416 individuals instead of 427, shrinking the total number of cells from 2.3 million to 1.9 million.

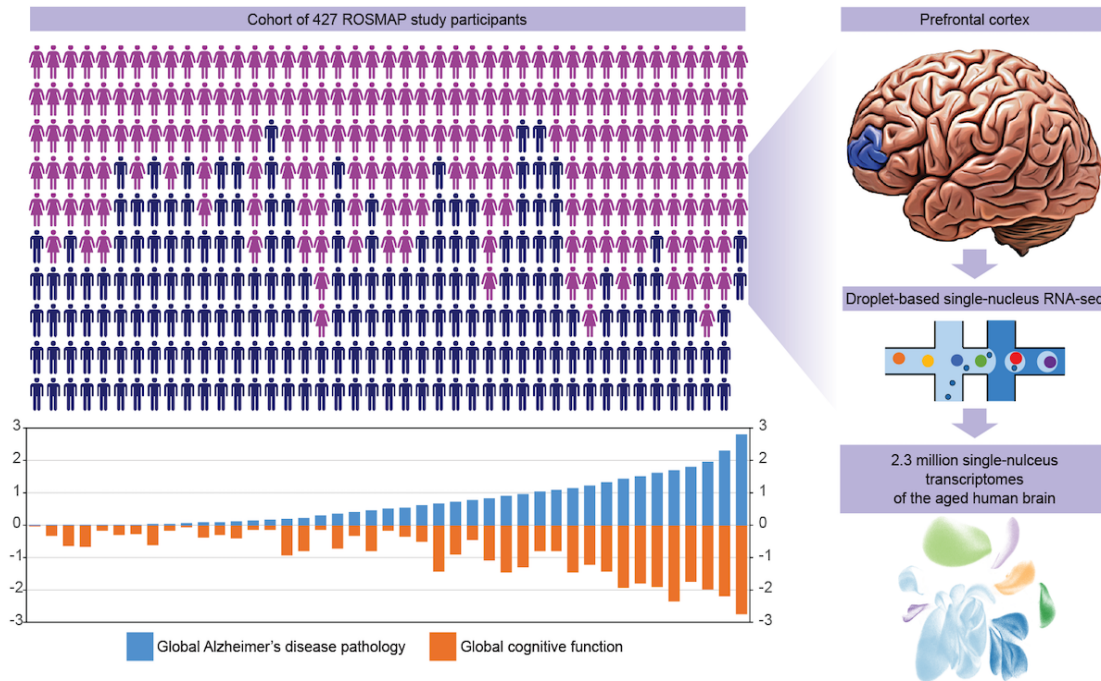


Figure 3-1: Breakdown of the scRNA ROSMAP dataset as well as summary statistics regarding surrounding the dataset

3.2 scRNA-seq processing

Droplet-based snRNA-seq libraries were prepared using a mix of V2 and V3 Chromium Single Cell 3' Reagent Kit and sequenced on an Illumina NovaSeq6000 at the MIT BioMicro Center. The raw sequence reads were aligned to pre-mRNA annotated Homo sapiens reference genome GRCh38 and counts were estimated using STAR-SOLO. We processed each individual BAM file with the following configuration: number of loci the reads map to: =1 for unique mappers, >1 for multimappers; multiple alignment index = 1, number of mismatches per paired alignment; local alignment scores where +1/-1 for matches/mismatches respectively, and penalties for indels and gaps; sequences and quality score of cell barcodes and UMIs for demultiplexing; sequence of entire barcode; quality of entire barcode; and the edit distance to the reference (number of mismatched + inserted + deleted bases). While we used STARSOLO as the processing software, we used CellRanger's V2 and V3 barcode whitelist in our analysis. Note that empty droplets were filtered from the results using CellRanger's

default EmptyDrops parameters. The output from this preprocessing pipeline provides us with the reads that match the transcript, splicing junctions, all counts from reads overlapping any given gene’s exon and introns.

On average, each sample from the dataset took on average six hours to process. Due to the high number of samples and the low processing throughput for each sample, split the cohort into 14 groups, thus enabling us to parallelize the process. Note that due to the high volume of genes and cells, additional optimizations were made in the processing loop to ensure that there was sufficient memory to process.

3.3 Dataset quality control and cell type differentiation

Given the scRNA ROSMAP data containing 2.3million cells, we first need to cluster the cells into their respective cell types. However, prior to clustering, we perform a quality control step to remove extreme outliers. Firstly, we removed all cells that expressed less than 500 genes with expressions in less than 50 cells. Next, cells with high concentrations of mitochondrial genes were also removed (maximum of 10%). We performed this step based on existing literature highlighting the upregulation of mitochondrial genes in low-quality (broken) cells, as a result, it is important to ensure that these cells are removed so they do not skew the data.

The next quality control step involves identifying and removing doublets. In the introduction, we outlined the general process behind scRNA-seq and the ability for the technology to study transcriptional heterogeneity at the single cell resolution. However, a common error with scRNA-seq involves two or more cells being considered a single cell; these technical artifacts are called doublets. Doublets have a higher chance of occurring when scRNA-seq assays profile a large number of cells, such that when they are placed in wells, groups of cells are tagged with the same unique barcode. To identify and remove these doublets from our dataset, we perform the following four steps:

1. Generate artificial doublets
2. Merge real data with artificial data
3. Perform PCA and finding proportion of artificial nearest neighbor (pANN)
4. Rank and order pANN values

First, using the software package DoubletFinder, we generated a set of artificial doublets based on our ROSMAP data. To generate artificial droplets, we combined gene expression profiles of two randomly selected droplets. Next, we append the dataframe containing the original ROSMAP data with the artificially created droplet dataframe. With this merged real-artificial dataframe, we performed principal component analysis, thus projecting the data into a lower dimensional space that still captures the relationship between the two. We use k-nearest neighbors clustering for every real cell in the real-artificial dataset in PC space. Once this is completed, we can then compute the proportion of artificial doublets in each cell's clustered nearest neighbors. Next, we ranked and ordered the cells based on their pANN values in an ascending manner, thereby allowing us to set a threshold pANN such that any pANN values higher than the threshold represent doublets, and values below the threshold represent valid cells. Finally, we can filter out all cells above the pANN threshold, providing us with a doublet free single cell dataset (Figure 3-2).

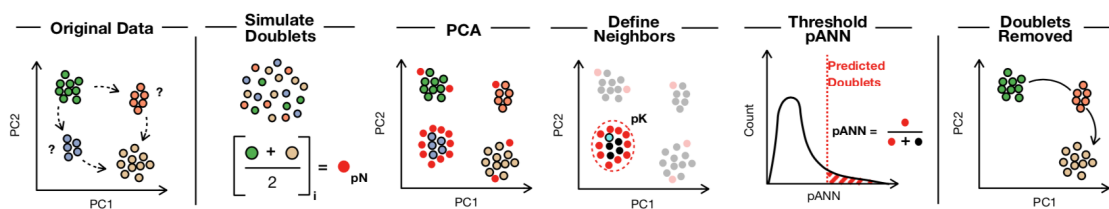


Figure 3-2: Breakdown of the steps involved for removing doublets [5]

To identify cell-type membership for each cell, we leverage known cell-type marker genes from Mathys et al [13]. They identified three marker genes for each major cell type: Microglia (mic), Astrocytes (ast), Excitatory (exc), Inhibitory (inh), Oligodendrocytes (oli), Oligodendrocytes progenitor cells (OPC), and Vasculature (vas). We

first used Seurat’s default clustering algorithm that is based on embedding the data into a KNN graph based on the Euclidean distance in PCA space. We then projected the data onto a UMAP using the same PCA inputs from the KNN step before. We then colored the UMAPs based on the expression profiles of all marker genes from each cell-type (Figure 3-3A). Based on the expression profile from each marker gene, we identified the cell-type membership for each cluster. For example, if Microglia cells have known marker genes A, B, and C, and the expression profile for cluster 7 are significantly increased for genes A, B, C, we assume that cluster 7 represents microglia cells in our dataset. Finally, we perform a secondary doublet quality control after the clustering step. In a perfect scenario, our clusters perfectly reflect the expected expression profiles for each marker gene in each cell type. However, there are cases where a cluster shows differential expression of marker genes from multiple cell types. We deal with this anomaly by assuming these cells are doublets and filter them out accordingly.

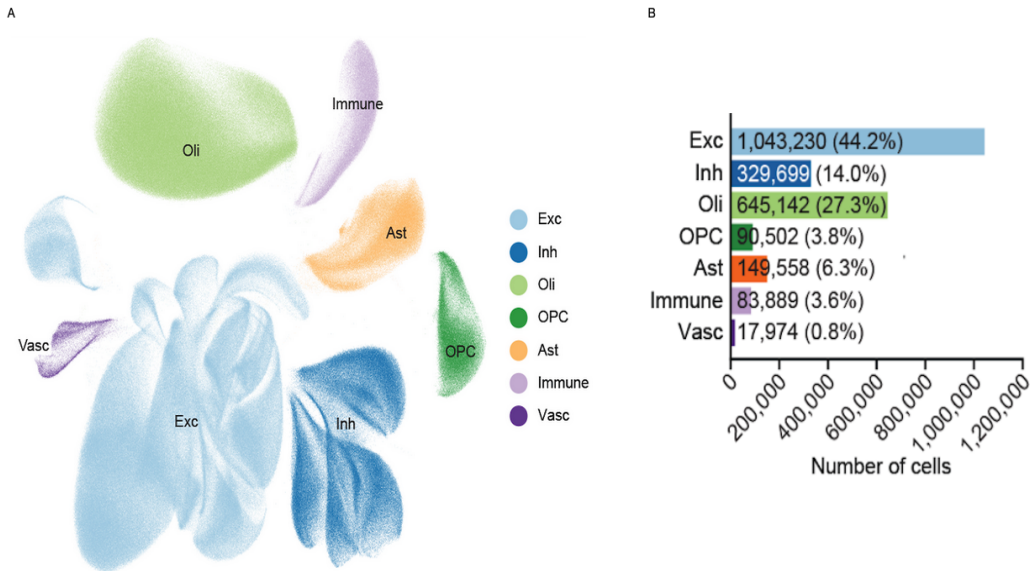


Figure 3-3: UMAP of all the cell types in the scRNA ROSMAP dataset, B: individual breakdown of the number of cells in each cell type

Chapter 4

Beta-binomial regression reveals differentially spliced genes

Once we have preprocessed the dataset, we can begin our beta-binomial based regression analysis with the goal of 1) confirming the validity of the beta-binomial regression for splicing analysis 2) confirming that there are no global splicing differences between AD and control groups 3) identifying a novel set of differentially spliced genes 4) understand to what extent these genes affect specific cell types or cause a global change.

4.1 Splicing count inference

To begin, we first need inferred spliced and unspliced counts for each cell and gene. This was accomplished by using the STARSOLO and Seurat packages. We used STARSOLO's default configuration to estimate the spliced and unspliced counts for each gene and cell. Next, we performed a pseudo-bulk aggregation of the data, in hopes of creating a matrix that is n-genes by 2. For each cell-type, we iterated over all individuals and extracted the relevant columns (cells), such that we have an n-genes by n-cell matrix that represents all individual cells across every individual for a given specific cell type and the counts for each gene. We then performed an aggregation step by summing across the rows, thereby collapsing all the cells into one summed

column. We are left with a column matrix where the rows represent genes, and the column represents either the spliced or unspliced summed total. Note that we filter out genes that have both their mean spliced and unspliced counts less than one.

Many software packages are available to estimate the spliced and unspliced counts, namely STARSOLO and Velocity. We chose STARSOLO because the results between the two packages were the same; furthermore, the STARSOLO supports fastq file formats which significantly increases processing speed compared to Velocity. More importantly, Velocity only provides support for aligned BAM files, whereas STARSOLO can use variable human reference genomes such as hg38 or hg19, ultimately providing more flexibility with respect to processing.

4.2 Splicing analysis approach

4.2.1 Beta binomial distribution

We implemented our beta-binomial regression in R, using AOD's implementation as a validation model. Counts were aggregated across each individual, creating an $n \times m$ where n = number of individuals, and m = number of genes matrix that was used as an input into the model. A $p \times 1$ where p = test parameters were also used as inputs into the model. We applied the model to multiple phenotypes such as: BraaK score, NFT score, AD, sex, age, and amyloid-beta score. We found that regression using the AD parameter showed the strongest statistical correlation.

4.2.2 Meta-analysis aggregation

While cell-type specific variations provide invaluable information, we also hope to identify genes that are differentially spliced across all cell-types. We achieved this by performing Fisher's combined probability test to fuse the p-value from the cell-type specific analyses. Specifically, given a list of p-values, we can combine them into one test statistic as follows:

$$\chi_{2k}^2 \sim -2\sum_{i=1}^k \log(p_i)$$

Where p_i represents the i th p-value in our list. As a result, when the p-values approach zero, it will result in a larger combined test statistic, suggesting that the null hypothesis may not be true for every i . Conversely, when the null hypothesis is true the combined statistic follows a chi-squared distribution with $2k$ degrees of freedom. With this information, we can easily calculate the p-value for the combined statistic.

The logic above hinges on the assumption that the combined test statistic follows a chi-squared distribution. This can be proven by the fact that under the null hypothesis for each i , the p-value follows that of a uniform distribution from 0 to 1. As a result, taking the negative log of a uniform distribution follows an exponential distribution as follows:

If we let $Y = \log(X)$

$$P(Y < y) = P(-\log(X) < y) = P(X > \exp(-y))$$

Given X is uniformly distributed between $[0, 1]$

$$P(Y < y) = P(X > \exp(-y)) = F(1 - \exp(-y))$$

Next, we can construct the chi-squared distribution with two degrees of freedom by scaling a value that follows the exponential distribution by a factor of two. Leaving us with k independent chi-squared values with two degrees of freedom such that the sum follows a chi-squared distribution with $2k$ degrees of freedom.

4.3 Global splicing characteristics do not differ between AD and control cohorts

To identify our new differential splicing gene set, we first must investigate the splicing characteristics of the entire dataset, ensuring that there is no global difference between AD and non-AD cohorts' spliced/unspliced ratio.

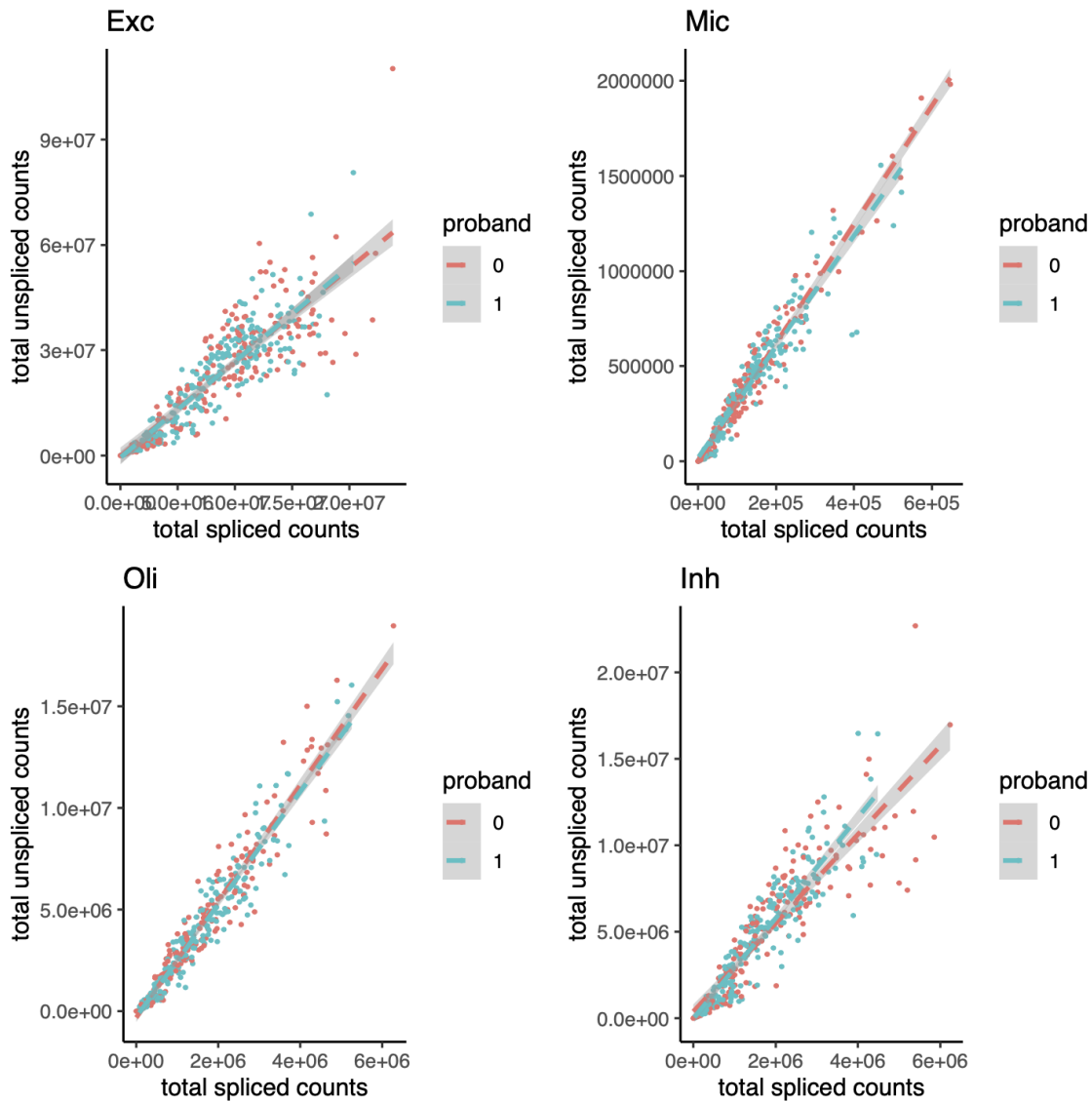


Figure 4-1: Global spliced/unspliced counts for excitatory, microglia, oligodendrocyte, and inhibitory cell types. Points that are colored red indicate control, while points that are colored blue indicate AD individuals

Initial analysis into the global splicing characteristics did not show any significant

difference between splicing counts in excitatory, microglia, inhibitory, and oligodendrocyte cell types (Figure 4-1). Interestingly, there appears to be a slight difference between the two cohorts for inhibitory neurons. The AD cohort appears to have decreased splicing when compared to the control group. While this wasn't expected, recent studies have shown inhibitory neurons to be one of the earliest cell types whose neuronal functions are dysregulated. Specifically, it is hypothesized that amyloid-beta plaque and tau protein accumulation have been linked to inhibitory neuronal dysfunction and overall activity of neuronal circuits [32].

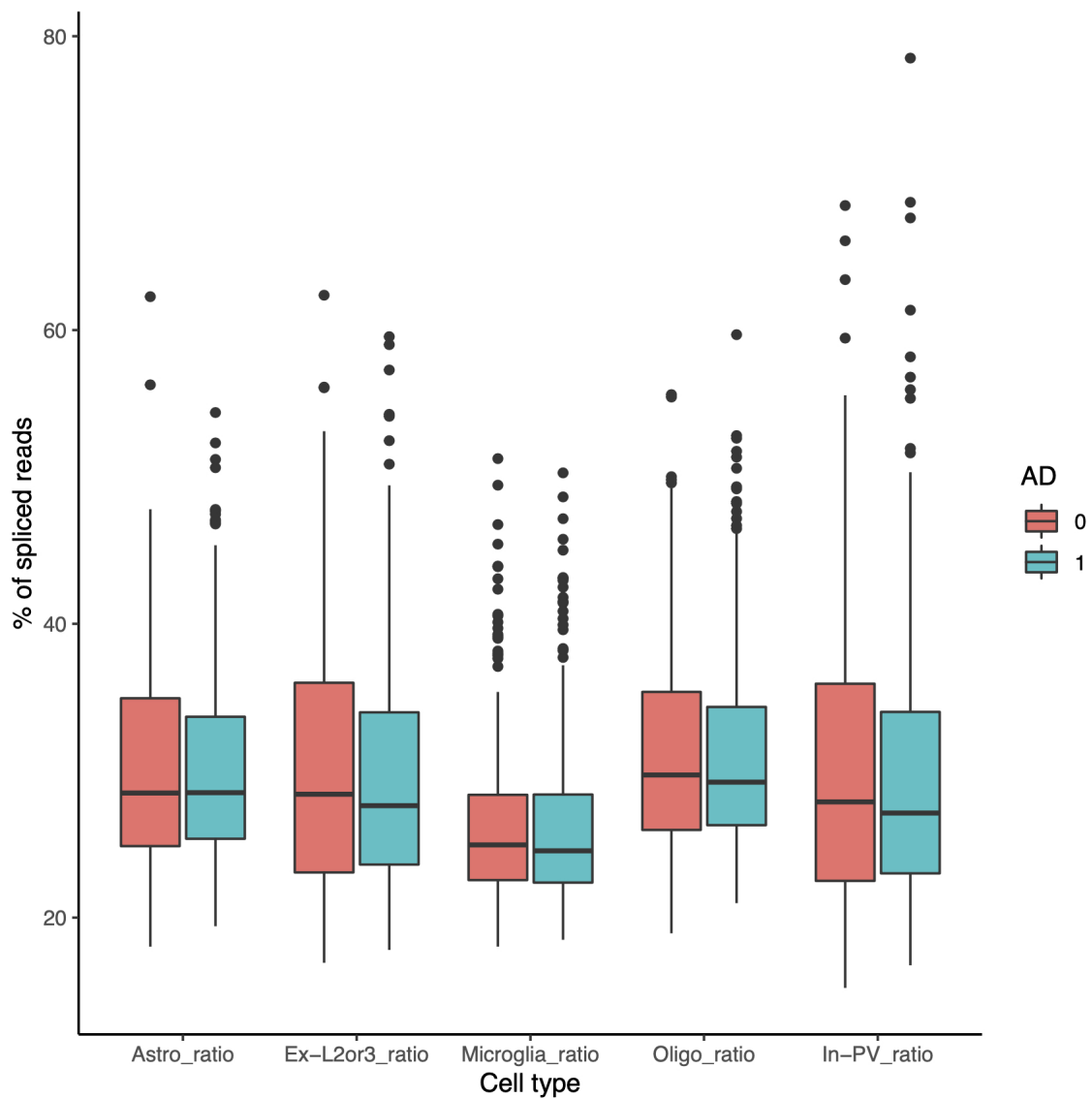


Figure 4-2: Box plot comparing percentage of spliced reads between AD and non-AD cohorts across 5 cell types

Next, we also wanted to model the global splicing characteristics with respect to the percentage of spliced reads. We do this since our final goal is to be able to have one indicator variable that describes the difference between AD and control cohorts. From this perspective, it is evident that there is no difference between AD and non-AD percent of spliced reads (Figure 4-2); however, we note that there are a significant number of outliers especially in the microglia and oligodendrocyte populations.

4.4 Beta-binomial regression reveals cell-type specific differentially spliced genes

Once we have confirmed that there is minimal difference between AD and non-AD global splicing ratios, we hope to identify specific genes that do not reflect this pattern. As mentioned above, due to the nature of scRNA-seq assays, there is significant experimental noise that needs to be dealt with. While we could have modeled the counts as a naïve binomial model, the beta-binomial overdispersion parameter allows us to consider the experimental noise as well as the variation in read coverage. As a result, we hypothesized that pairing spliced/unsliced counts from the same individual with a beta-binomial regression could provide more statistical power and hopefully uncover novel AD characteristics.

4.4.1 Enrichment of SRRM2, CRY1, and ASPDH1 in Astrocytes

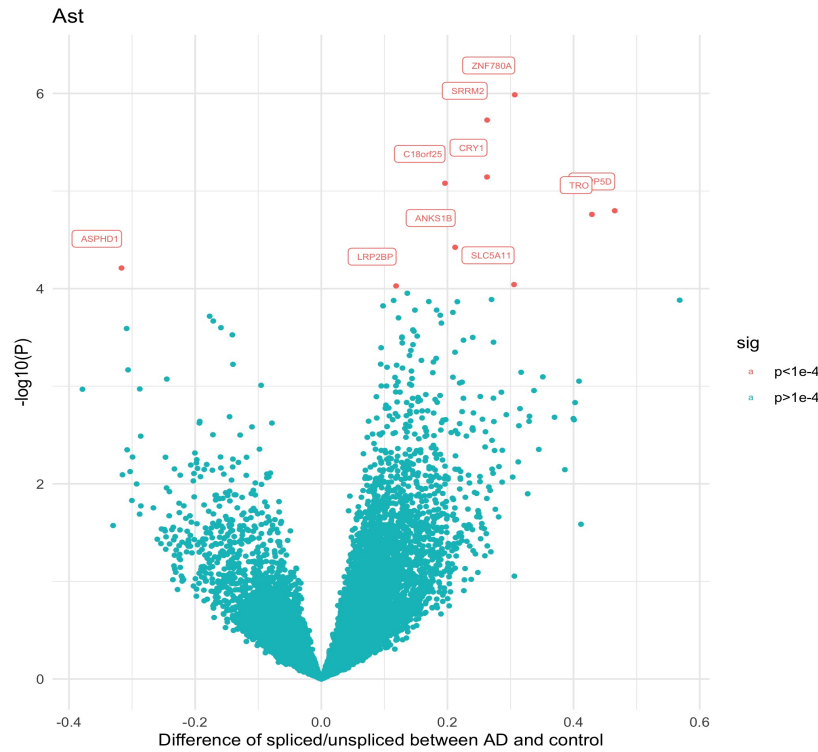


Figure 4-3: Differentially spliced genes between AD and control groups in astrocytes. Volcano plot of differentially spliced genes in astrocytes. Adjusted $p < 0.0001$

Among the top 10 most differentially spliced genes in astrocytes ($p=0.0001$), we identified SRRM2, CRY1, ASPDH1, and ZNF80A as genes of interest (Figure 4-3). SRRM2 is a known splicing factor that encodes the serine/arginine repetitive matrix protein 2. Studies have shown that there is consistent dysregulation of SRRM2 in cases of Parkinson's disease. The specific mechanisms and greater implication of SRRM2 dysregulation will be further elaborated on below.

CRY1 encodes a flavin adenine dinucleotide-binding protein that plays a key role in the circadian oscillator complex, which regulates the circadian rhythm. While we cannot pinpoint the exact role CRY1 plays in AD pathogenesis, we are able to identify CRY1 as significantly dysregulated in AD cases. This is consistent with our understanding of AD where disturbances in sleep cycles and circadian rhythms have been widely reported. Circadian clock disturbances clearly lead to dysfunction of

the hypothalamic suprachiasmatic nucleus (SCN) dysfunction (Figure 4-4) [9]. SCN dysfunction ultimately leads to sleep disturbance, which in turn results in increase of amyloid-beta plaque; finally resulting in harmful aggregations of amyloid plaques and tau proteins.

Unfortunately, we are not completely aware of the exact role ASPDH1 plays in cellular functions; however, it has been predicted that ASPDH1 promotes dioxygenase activity as well as peptidyl-amino acid modification. In Patel et. al's transcriptomic analysis of both early and late-stage AD, they identified ASPDH1 as one of 15 genes that were consistently differentially expressed across all AD samples in the frontal cortex [12]. Their results found strong association between ASPDH1 and AD pathologies such as plaque and NFT accumulation as well as nerve cell death. Unfortunately, it remains unclear as to whether these genes are causative or a response to the pathology. Nonetheless, it appears that ASPDH1 plays some role in AD pathogenesis.

Finally, ZNF780A is predicted to be involved in enabling DNA-binding transcription factory activity, as well as the regulation of RNA polymerase II (pol II). Robinson et. al demonstrated that pol II promotes binding of splicing machinery by comparing pol II with a viral RNA polymerase, T7. Their study showed that while both generated similar quantities of pre-mRNA, spliced mRNAs were created significantly more in the pol II system [27]. Indicating that while pol II may not be directly part of the splicing mechanisms in a cell, it still increased the efficiency of splicing machinery to find suitable binding sites on mRNA complexes. Their findings also agrees with our results where we observe consistently higher spliced counts in AD cohorts compared to control. While our study is not able to pinpoint the exact mechanism by which these genes cause AD onset, from the functional aspect of these genes and given our understanding of AD, they appear to play an important role in the disease's progression. Furthermore, these results also validate our novel application of the beta-binomial model in a splice/unsplice count context, as the significantly spliced genes are related to known AD mechanisms and pathways.

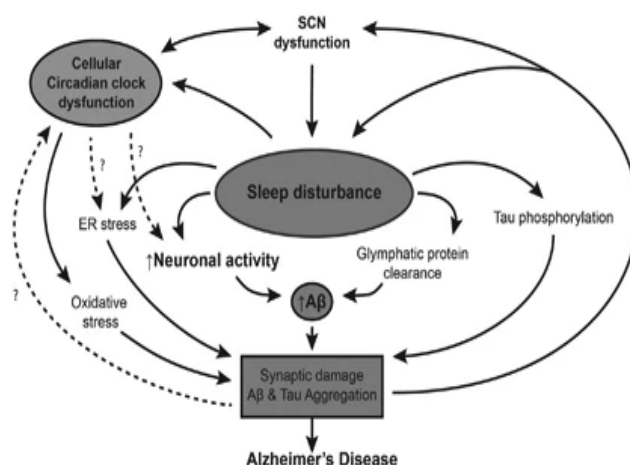


Figure 4-4: Hypothesized model of sleep disturbance and onset of AD

4.4.2 TRIM9, SRRM2, and COL11A1 differential splicing in excitatory neurons

Within the excitatory neuron cohort, we identified three genes of interest that are differentially spliced: TRIM 9, SRRM2, and COL11A1 (Figure 4-5). TRIM9 is a member of the tripartite motif (TRIM) family; while its exact function remains unknown, the TRIM family of proteins are involved in a broad range of processes, where alterations to the TRIM proteins result in pathological conditions associated with genetic diseases, viral infections, and cancers. Tanji et al. experimentally demonstrated that in healthy brains, there was a higher concentration of TRIM9 proteins in the cerebral cortex and hippocampus, regions of the brain associated with higher order functions. However, in individuals with neurodegenerative dementia, they found that TRIM9 immunoreactivity was significantly decreased in the same brain regions [21]. On the other hand, NFTs and senile plaques in AD did not demonstrate any TRIM9 immunoreactivity, suggesting that TRIM9 plays an important role in the formation of neurodegenerative pathologies. Our results agree with these findings: TRIM9 has a negative spliced/unsliced ratio, indicating the downregulation of the gene (Figure 4-5). Unfortunately, past studies were not able to pinpoint the exact role TRIM9 plays in the mechanisms that result in plaque accumulation; nonetheless, it is encouraging

that our results confer with the current understanding of AD.

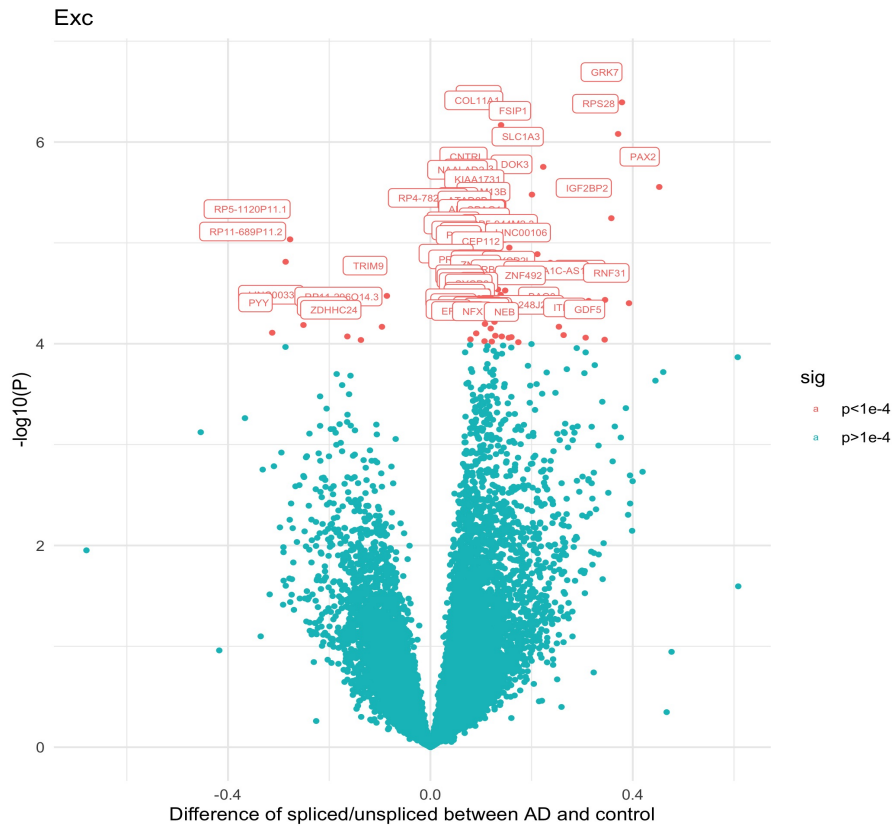


Figure 4-5: Differentially spliced genes between AD and control groups in excitatory neurons. Volcano plot of differentially spliced genes in astrocytes. Adjusted $p < 0.0001$

Once again, we notice that SRRM2 is differentially spliced in excitatory cell types; Section 4.5.1 will go into more detail regarding the suspected role of SRRM2 in AD pathogenesis. Lastly, COL11A1 is a known AD associated gene that encodes alpha chains for type XI collagen; COL11A1 also plays an important role in the P13K-Akt signaling pathway that is a known promoter of metabolism, proliferation, cell survival and growth. P13K-Akt enhances cell survival by inhibiting proapoptotic proteins and processes as well as inhibiting the expression of genes whose function mediates cell apoptosis [35] This is in line with our current understanding of AD where massive neuronal cell death is a hallmark pathology for neurodegenerative diseases in addition to apoptotic cell death. Consequently, we hypothesize that differential splicing of COL11A1 results in P13K-Akt signaling pathway dysregulation, ultimately inhibiting cell growth and survival.

4.4.3 Differential splicing of Z3CHAV1, XRCC4, and LPCAT2 in microglia

From microglia populations, we identified three differentially spliced AD genes: Z3CHAV1, XRCC4, and LPCAT2 (Figure 4-6). Z3CHAV1 encodes for a zinc finger antiviral protein that is a known RNA-binding protein that regulates the stability and translation of mRNAs. These proteins and their related pathways are key components of the cellular response of stress [31]. Relative levels of these RNA-binding proteins (RBP) are directly correlated to the accumulation of amyloid-beta and tau proteins in AD brains. Guo et. al's RBP accumulation study directly identified Z3CHAV1 as one of four genes that demonstrated significant increase in AD cases [26]. Additionally, protein network correlation analysis for RBP complexes identified modules related to the spliceosome showed stage specific enrichment in AD cases. This result agrees with our findings of splicing genes, such as SRRM2, as being differentially spliced.

XRCC4 encodes for proteins that interact with DNA-dependent protein kinases in the repair of double-strand breaks through non-homologous end joining (NHEJ). While age is a confounding factor for decreased NHEJ activity, Kanungo et. al demonstrated through age-matched brain comparison that there is a significant reduction in NHEJ activity in AD cases [18]. Unfortunately, it is not known whether there is a direct link between NHEJ levels and AD; however, understanding XRCC4 differential splicing patterns could yield new information linking DNA damage with AD pathologies.

In addition to amyloid-beta and tau protein accumulation, AD is also characterized by altered lipid metabolism in the brain. Brain lipids play a critical component in the development of cell membranes for neurons, astrocytes, and oligodendrocytes. Fatty acids bind themselves to lysophosphatidylcholine (LPC) before crossing the blood brain barrier (BBB) where the brain can synthesize the fatty acids; fatty acids bound to LPCs are found to be more efficient in crossing the BBB compared to other methods [29]. LPCAT2 plays an important role in this process by remodeling the fatty acid composition of membranes that enables the distribution of the bound fatty

acids to LPCs in the brain. As a result, differential splicing of LPCAT2 could result in the dysregulation of LPCs, and as a result, lead to the pathologies associated with AD.

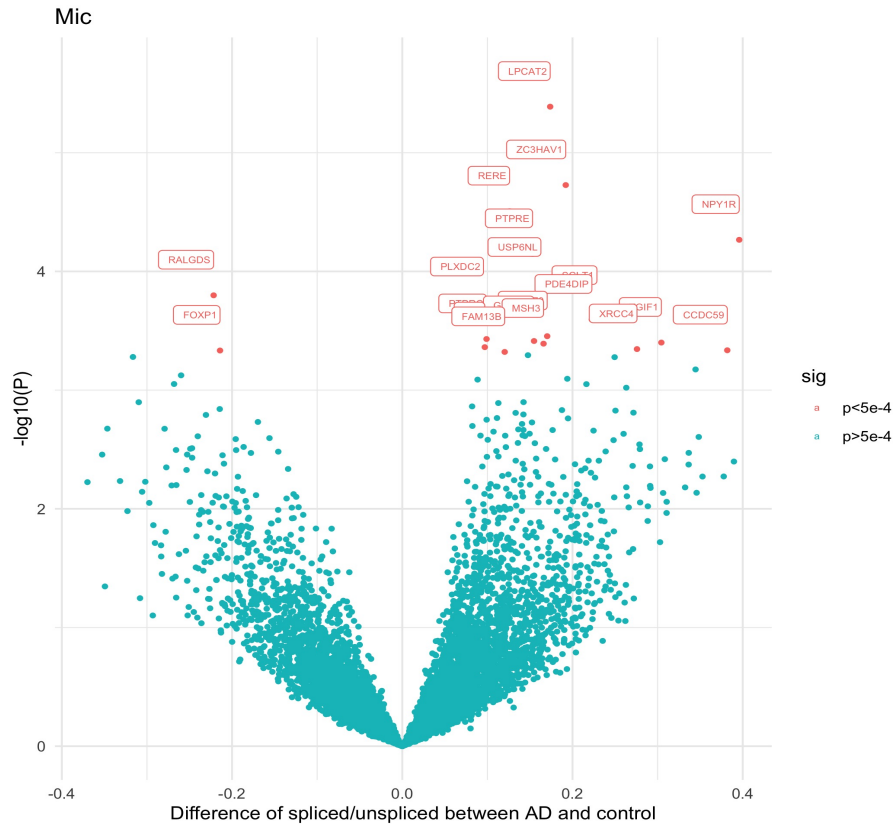


Figure 4-6: Differentially spliced genes between AD and control groups in microglia. Volcano plot of differentially spliced genes in astrocytes. Adjusted $p < 0.0005$

4.5 Meta-analysis reveals differential splicing of SRRM2, GOLGA4, and VPS13A

While analyzing cell type specific differential splicing genes offer insight into potential causes of known AD pathologies, we also wanted to gain a better understanding of genes that were consistently differentially spliced across all cell types (Table 4.1). Of the top genes that were differentially expressed, they could be categorized into three main distinct functional groups: splicing factor genes, known AD genes,

Gene	P-Value
SRRM2	1.61132184793516E-22
CACNA1A	1.46910163228038E-18
COL11A1	2.11318970577534E-14
CEP78	3.02967142294861E-15
VPS13A	6.09861448134243E-13
KIAA1731	2.46020750967306E-12
MACF1	2.5734753191161E-12
GINM1	8.46793502954224E-12
ANKRD36C	2.13382522046777E-11
PCF11	2.19692279203946E-11
NAALAD2	8.0352663719384E-12
GOLGA4	1.0032152746282E-10
PNPT1	2.12239605848419E-10
RPAP2	2.40377078321169E-10
CENPC	2.40479263251022E-10
N4BP2	3.00306098249708E-10

Table 4.1: Table representing the top 16 globally differentially spliced genes

and mitochondrial genes. Most of the top genes we identified were splicing factors. Interestingly, APOE was not in the top 100 differentially spliced genes. Perhaps indicating that transcriptional regulatory changes of APOE does not necessarily impact the overall ratio between spliced/unsliced counts, but rather lead to other biochemical cellular changes whose difference cannot be directly measured through a velocity perspective.

The following section will focus on what we identified as three key genes from the meta-analysis rankings. Namely SRRM2, GOLGA4, and VPS13C. All three of these genes were ranked within the top 15 differentially spliced genes across all cell types. Furthermore, each of these genes belong to the three identified major functional categories: splicing factors, known AD genes, and mitochondrial genes respectively.

4.5.1 SRRM2 splicing dysregulation in AD

SRRM2 is a protein that enables C2H2 zinc finger domain binding activity which is involved in mRNA splicing. Splicing is the process by which introns are skipped and exons are ligated, thereby creating proteins through a translation mechanism.

However, the human body introduces another layer of genomic complexity through alternative splicing. Alternative splicing is a change in the splicing mechanism where exons of pre-mRNA are ligated in different orders to form large transcripts with different protein sequences [28]. Studies have focused on studying abnormal alternative splicing events in neurodegenerative diseases such as Parkinson's disease, specifically highlighting six genes, of which SRRM2 is included, that are involved in alternative splicing in Parkinson's disease. Shehadeh et al. analyzed the blood microarray of patients with Parkinson's disease and found that SRRM2 was the only gene that was differentially upregulated in all their analyses. Their results indicated a downregulation of the long SRRM2 isoform and an upregulation of the short isoform [22]. Thus, understanding the effects of splicing and the resulting changes to the isoforms could be critical in understanding the underlying mechanisms of disease. However, due to limitations, which will be expanded in Chapter 5, we are not able to properly investigate the effects of alternative splicing and their related isoforms. Regardless, identifying SRRM2 differential splicing in AD cases, along with the fact that there has been extensive research into the role SRRM2 plays in Parkinson's disease progression, highlights the potential importance of understanding splicing characteristics in AD [20].

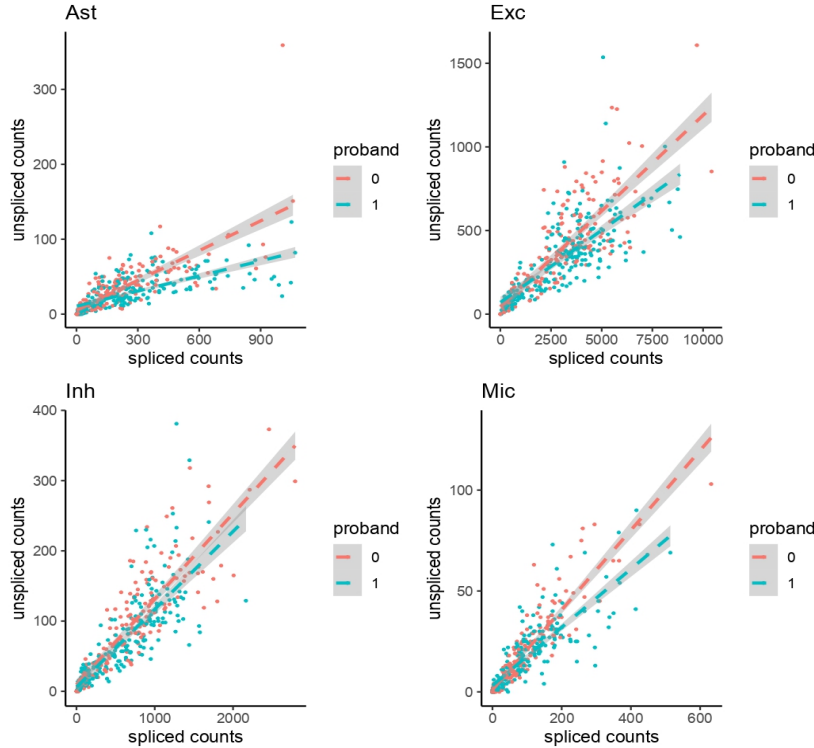


Figure 4-7: Raw spliced/unspliced counts for SRRM2 across astrocytes, excitatory, inhibitory, and microglia cell types

There is a significantly larger difference between spliced/unspliced ratio in AD patients in astrocytes, excitatory, and microglia types compared to inhibitory cells (Figure 4-7). These results confer with known splicing dysregulation in neurodegenerative diseases. Astrocytes play a pivotal role in brain homeostasis where the presence of amyloid-beta have been shown to disrupt neurotransmitter uptake and alter signaling. Interestingly, compared to the other three cell types, inhibitory neurons showed the smallest difference between the spliced/unspliced ratio of AD and control individuals [14]. While we are not exactly sure why this is the case, little is known about the pathology in inhibitory cell types of AD patients, whereas extensive research has been done with understanding excitatory neuron loss and downstream contributions of amyloid-beta accumulation and neuronal dysfunction. Regardless, our results indicate that, while not as severe, inhibitory neurons still exhibit slight differences in splicing ratios- suggesting that understanding alternative splicing in inhibitory neurons could yield novel information regarding AD pathogenesis.

4.5.2 Differential splicing of GOLGA4 highlights the role of lipid transport in AD pathogenesis

The Golgi apparatus participates in the glycosylation and more importantly the transport of proteins and lipids in the secretory pathway. GOLGA4 encodes for one family of proteins found in the Golgi, known as golgins. While the exact role of these GOLGA4 golgins are unknown, it has been hypothesized that they play a role in the pathway that allows for lipids and proteins to pass the Golgi membrane. Additionally, studies have also identified elevated levels of GOLGA4 in patients with AD compared to control in cerebrospinal fluid. Specifically, Kork et al. generated more than 200 hybridomas with the goal of finding an antibody that can differentiate between AD and non-AD cohorts. The study found one antigen that was able to differentiate the two cohorts, GOLGA4 [10]. While the exact pathways and mechanisms that connect GOLGA4 and AD pathogenesis remain unknown, the results support the idea that GOLGA4 serves as a marker gene in AD.

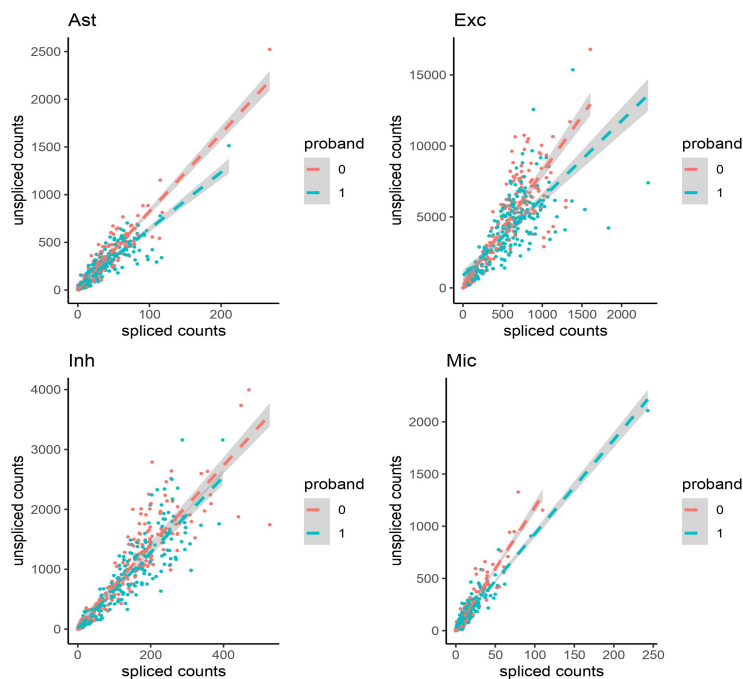


Figure 4-8: Raw spliced/unspliced counts for GOLGA4 across astrocytes, excitatory, inhibitory, and microglia cell types

Much like SRRM2, GOLGA4 differential splicing was significantly higher in as-

trocytes, excitatory, and microglia cell types compared to inhibitory neurons (Figure 4-8). In fact, there doesn't seem to be a significant difference between splicing ratios in inhibitory cell types. While we do not know of the exact reason for this phenomenon, it could be linked to the fact that cerebrospinal fluids increase excitability of pyramidal neurons without affecting inhibitory neurons [2]. As a result, possibly explaining the difference between the splicing patterns of GOLGA4 in excitatory and inhibitory cell types.

4.5.3 Differential splicing of VPS13A reveals role of mitochondria in AD pathogenesis

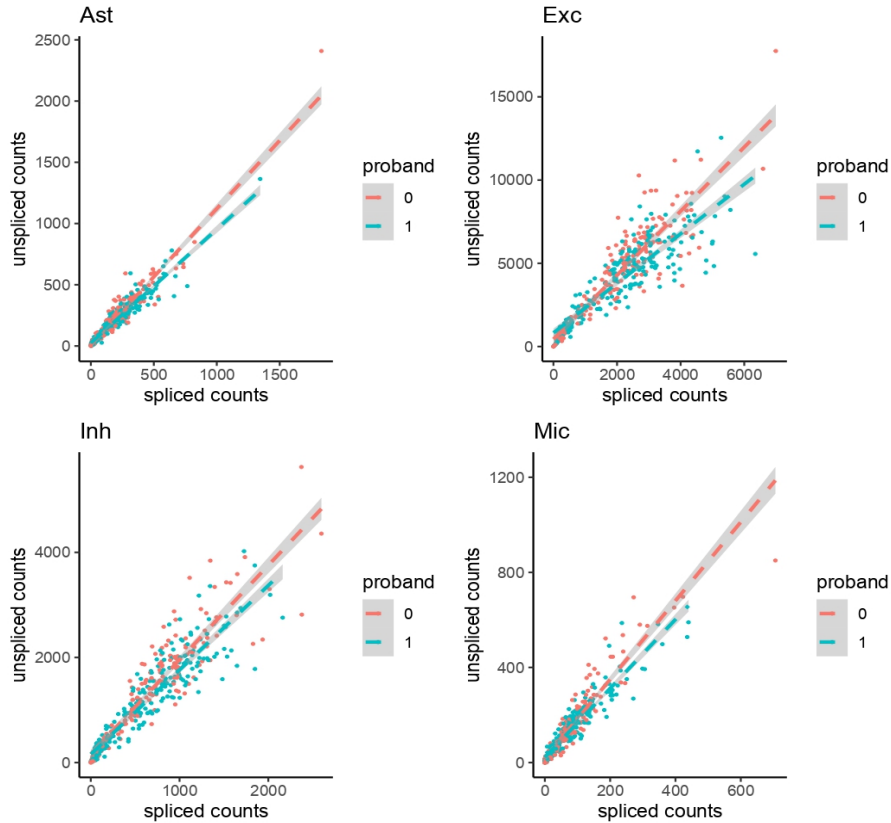


Figure 4-9: Raw spliced/unspliced counts for VPS13A across astrocytes, excitatory, inhibitory, and microglia cell types

VPS13A is a mitochondrial gene that encodes for a protein called chorein. In addition to the accumulation of amyloid plaques and tau proteins in AD pathology, evidence from studies imply that damage to the mitochondria may play a funda-

mental role in the pathogenesis of AD. The mitochondria is the main energy source providing ATP that enables cellular functions to occur. Enabling neuronal activity as well as regulatory mechanisms in the nervous system. Specifically, damage to the mitochondrial respiratory chain or defense system can hinder the ability for the organelle to defend against reactive oxygen species accumulation [6]. As a result, it is not surprising to find dysregulation of mitochondrial functions associated with neurodegenerative diseases such as AD.

Interestingly, compared to the results from SRRM2 and GOLGA4, the difference in ratios between inhibitory cell types and the other cell types are much closer (Figure 4-9). However, the AD cohort still exhibits consistently higher ratios of spliced counts. A potential cause for similarity between the differences could be attributed to the fact that the mitochondria are conserved organelles, and as a result, their effects are not necessarily localized to one cell type, but rather all cell types.

Chapter 5

Random priming effects in isoform analysis

Using IGV, we managed to analyze the distribution of the reads over the entire genome across all individuals aggregated on their AD phenotype. We noticed that AD patients consistently had fewer reads in the intronic regions when compared to control. As such, we thought intronic retention was occurring only in AD cases. However, analysis comparing gene length with respect to the percentage of spliced reads revealed an unexpected technical artifact: there is a strong inverse correlation between gene length and the percentage of spliced reads (Figure 5-1). As the gene body length increases, the percentage of spliced reads decreases. Furthermore, we also found that the spliced reads are biased towards genes where exons compose 95% of the total gene body length.

This inverse relationship between gene body length and percent of spliced reads indicates that intronic retention may not be the main explanation behind our observation, but rather a result of random priming. scRNA-seq assays depend on probes that attach to the poly-A tails of the gene for amplification and quantification to occur. However, in many intronic regions of genes, there are long poly-A sequences that occur naturally. Given that the probe is not able to differentiate between poly-A tails and poly-A sequences, we hypothesize that as gene length increases, the probability of a poly-A sequence occurring increases, thus increasing the chances for probes to at-

tach themselves these intronic poly-A sequences instead of poly-A tails. STARSOLO then classifies these randomly primed sequences as unspliced because they do not line up with known spliced gene sequences, thus artificially inflating the unspliced counts and deflating the spliced counts. Because of random priming, it is difficult for us to fully explore the transcript as well as exon expression. We suspect as improvements in single cell processing software and probe design occur, exploring these avenues could yield new information regarding disease pathogenesis.

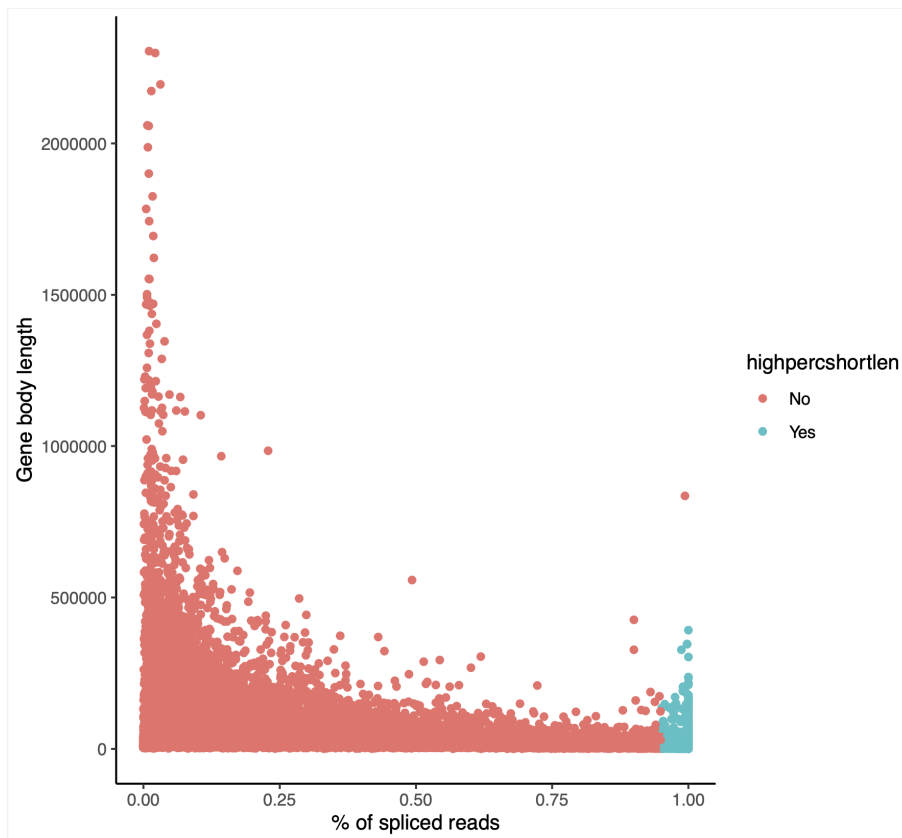


Figure 5-1: Scatterplot showing the relationship between gene body length and percentage of spliced reads. Highpercshortlen are genes where at least a 95% of the gene body is composed of exons

Chapter 6

Pathway enrichment analysis identifies AD-related enriched pathways from differentially spliced genes

6.1 Methodological Approach

We wanted to better understand what pathways are enriched given a set of differentially spliced AD genes. To do so, we performed analysis at the cell type level as well as the meta-cell type level in hopes of understanding which pathways are enriched globally, and which are enriched for only certain cell types.

We performed this analysis using the pathFindR package. P-value adjustments were computed with FDR, and final significance threshold of $p=0.05$ and an enrichment threshold 0.05. Identified pathways are based on the KEGG pathway database.

6.2 Analysis reveals enrichment of spliceosome, homologous recombination, and reactive oxygen species in AD astrocytes

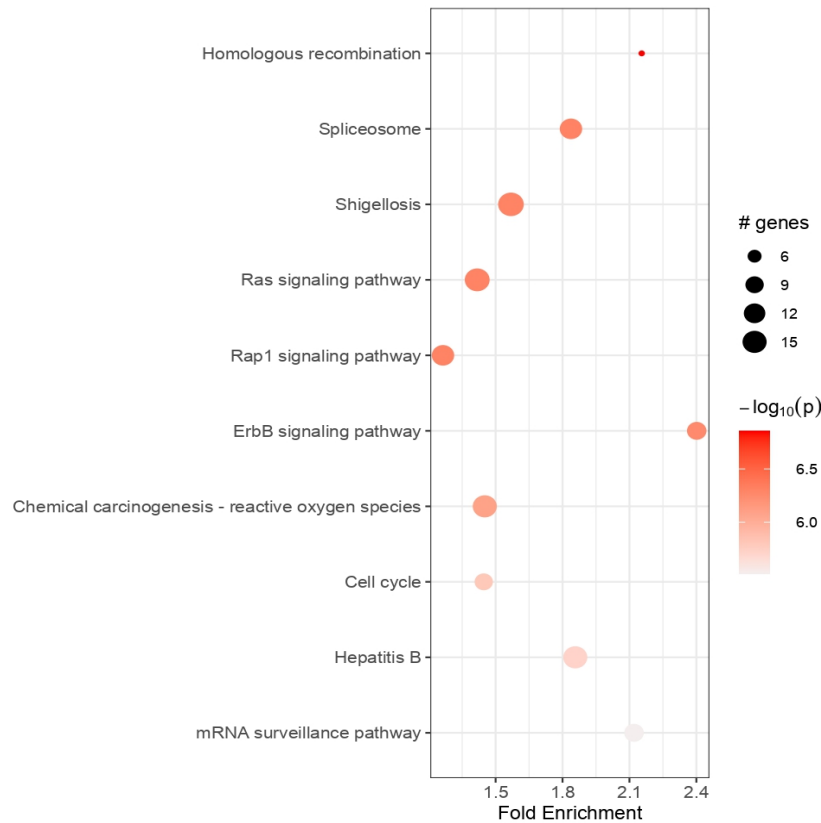


Figure 6-1: Top 10 enriched pathways in astrocytes given the differentially spliced astrocyte gene set (FDR adjusted p-values), significance and enrichment thresholds were both equal to 0.05

Pathway enrichment analysis of astrocytes identified three pathways of interest that are related with known AD pathways: spliceosome, homologous recombination, chemical carcinogenesis – reactive oxygen species (Figure 6-1). Given SRRM2 being differentially spliced across all cell types, it is not surprising to find the spliceosome pathway as being enriched. Unfortunately, this analysis does not shed light on the exact mechanisms by which the spliceosome pathway is dysregulated in AD. However, the results do support our hypothesis that differential splicing of splicing factors could play an important role in AD pathogenesis.

Homologous recombination is an important pathway that participates in the repair of DNA damage. Accumulation of damaged DNA is a known aging factor, and consequently a reliable biomarker for AD pathogenesis [34]. Furthermore, the enrichment of the reactive oxygen species is also expected in the presence of DNA damage. Our current understanding of DNA damage in AD suggests that accumulation of reactive oxygen species is the major source of such damage in AD brains. Lin et al. identified elevated levels of oxidative DNA lesions in pre-clinical AD brains, specifically pointing to increased levels of 8-hydroxy-2'-deoxyguanosine as the cause for DNA lesions from reactive oxygen species [34]. Ultimately, our pathway enrichment analysis supports existing evidence regarding AD pathogenesis through the dysregulation of homologous recombination pathway, which results in the inhibition of the brain's ability to repair DNA damage. Inhibited DNA repair mechanisms lead to the accumulation of reactive oxygen species in the brain, where it has been experimentally verified as a known AD pathology and a cause of amyloid-beta accumulation.

6.3 Pathway analysis reveals enrichment of TNF and PD-L1 in AD microglia populations

We identified enrichment of the TNF signaling and PD-L1 expression pathways in microglia as particularly interesting and related to AD (Figure 6-2). The tumor necrosis factor (TNF) pathway plays an important role in cell proliferation, apoptosis, and modulation of immune responses. TNF is a multifunction proinflammatory cytokine that helps in the regulation of lipid metabolism amongst other things. Due to the importance of TNF in cellular immune response as well as apoptosis, it has historically been of much interest in cancer research. However, recent studies have begun to link elevated TNF levels with decline in cognitive function in AD patients. Studies have exposed human microglial cells to amyloid-beta plaques that resulted in the release of high levels of TNF; furthermore, recent in vitro studies caused overexpression of TNF in neurons in the human hippocampus resulted in microglial and neuronal

death, suggesting TNF driven inflammation may have deleterious effects [4]. Finally, TNF also upregulates the expression of APP, the precursor molecule to amyloid-beta, and as a result, also causes the creation of large amounts of amyloid-beta plaques.

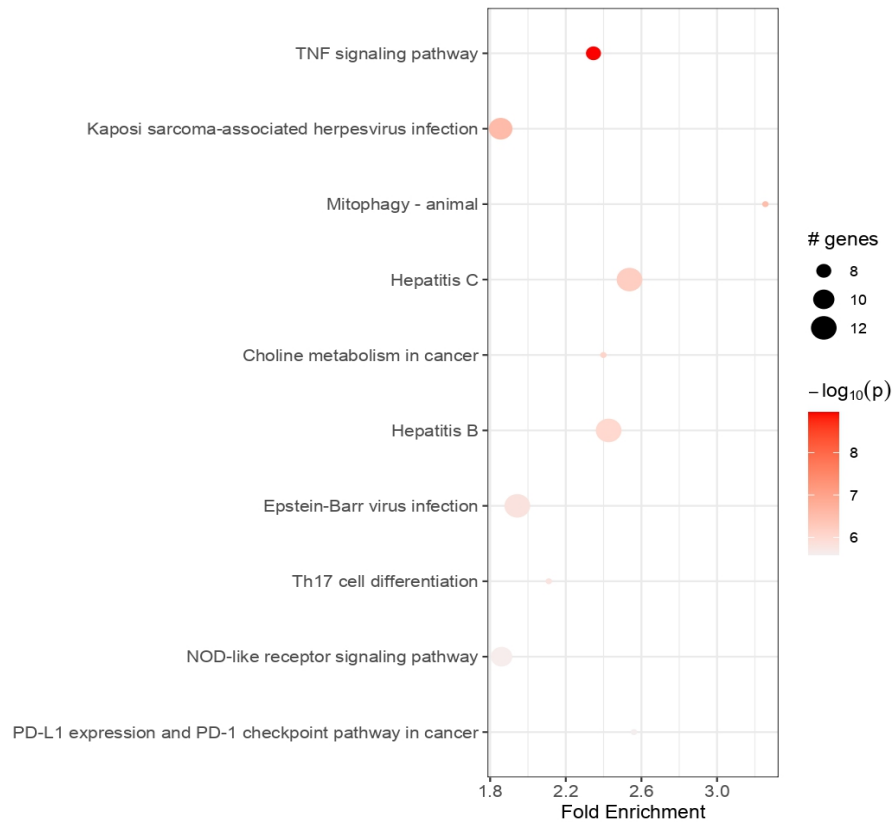


Figure 6-2: Top 10 enriched pathways in microglia given differentially spliced microglia gene set, significance and enrichment thresholds were both equal to 0.05. FDR adjusted p-values

PD-L1 is a key immune checkpoint that maintains immune response in the cell. Rosenzweig et al. experimentally confirmed that by blocking the PD-L1 pathway in animal AD and tauopathy models, they were able to cause an immune response that resulted in the modification of disease progression [25]. They provided AD mouse models with injections of anti-PD-L1 antibodies that resulted in the reversal of memory loss, decrease in the rate of neuronal loss, and retained learning ability. Because of the ability for PD-L1 targeted blocking to modify AD pathology, there has been recent interest in anti-PD-L1 targeted drug therapies. It is also encouraging that our novel single cell differential splicing method identified the PD-L1 pathway

as enriched in AD, as previous differential expression studies were not able to identify this pathway.

6.4 PD-L1, endocytotic, and insulin resistance pathways are enriched in excitatory neurons

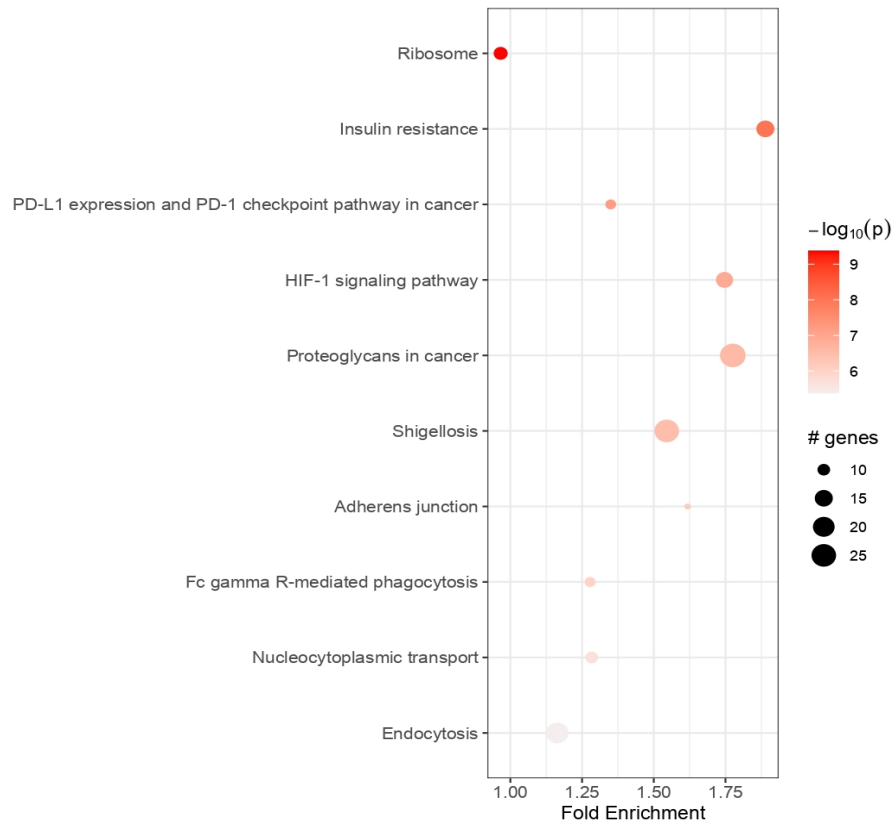


Figure 6-3: The top 10 enriched pathways in excitatory neurons given the differentially spliced excitatory gene set as input, significance and enrichment thresholds were both equal to 0.05 with FDR adjusted p-values

Pathway enrichment analysis revealed multiple pathways that were either directly related to AD or hypothesized to have causative relationship with AD pathogenesis: PD-L1, endocytosis, and insulin resistance (Figure 6-3). The PD-L1 pathway is one of the most enriched pathways across all the cell types, further supporting the results from Rosenzweig et al. where they demonstrated decreased AD pathologies in the presence of anti-PD-L1 antibodies as well as our results from section 6.3.

Enrichment of the endocytosis pathway also indicates the efficacy of our approach of understanding AD through a differential splicing perspective. Endocytosis is the mechanism by which the cell brings nutrients, lipids, and ligands from the cell membrane to the interior of the cell body. The endocytosis pathway is closely linked to amyloid-beta and tau pathologies in that they are the mechanism by which precursor amyloid-beta molecules are brought from the cell surface to the endosome where they are cleaved to produce amyloid-beta [17]. Consequently, we hypothesize that dysregulation of the endocytotic pathway paired with upregulation of APP production could play a major role in AD pathogenesis.

Finally, insulin resistance related pathways are widely known to be dysregulated in AD. Specifically, studies have identified insulin signaling as inhibited due to upregulation of insulin resistance in AD, which ultimately leads to the formation of NFTs in the brain. Interestingly, both AD and type II diabetes have strongly linked on molecular and clinical levels, suggesting the importance to understanding metabolic factors in AD pathogenesis [33]. While we currently do not have a clear understanding of the linking between the two, we can believe that insulin's role in regulating cellular inflammation could be important. Regardless, being able to identify multiple orthogonal pathways that all confer with our current understanding of AD pathogenesis proves promising for our differential splicing analysis.

6.5 Meta-analysis reveals enrichment of PD-L1, NF-kappa B, and adherens junction and close association with Parkinsons and prion pathways

Results from meta-analysis represents a set of globally enriched pathways; we identified five that were of particular interest: PD-L1 expression, adherens junction, NF-kappa B, Parkinson's, and prion pathways (Figure 6-4). In addition to seeing PD-L1 enriched in microglia and excitatory cell types, the PD-L1 pathway is also the top pathway enriched across genes that are differentially spliced in all cell types. This

further reinforces the hypothesis of PD-L1 dysregulation as a main factor behind AD immune response, but also provides hope for anti-PD-L1 therapies against AD.

The adherens junction is integral in maintaining the blood brain barrier. As mentioned in section X, impaired lipid metabolism is a known AD pathology; the blood brain barrier remains the main way by which lipids can enter the brain. Recent studies have identified elevated levels of occludin, a tight junction protein, as a potential marker for AD [24]. While the exact mechanism by which tight junction dysregulation causes AD pathologies is not known, all evidence suggests alterations to intercellular junctions in the blood brain barrier, which causes changes to the blood brain barrier permeability, have a potentially causative relationship with AD pathogenesis.

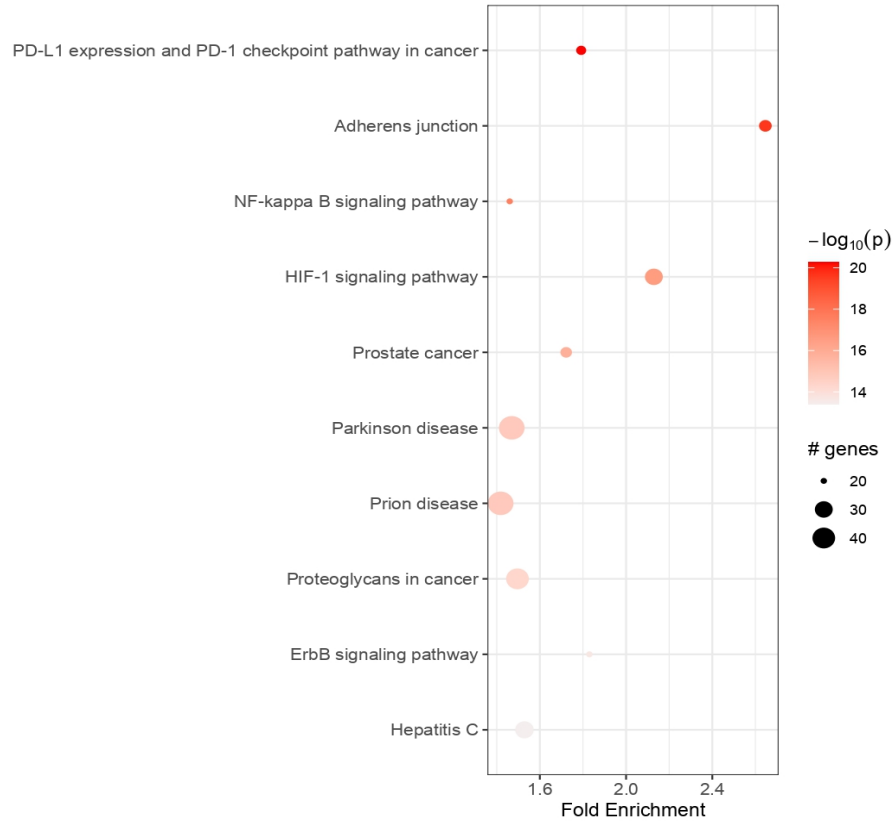


Figure 6-4: Top 10 enriched pathways across all cell types with FDR adjusted p-values, significance and enrichment thresholds were both equal to 0.05

Recent studies have identified NF-kappa B dysregulation in AD, specifically over-expression of NF-kappa B was consistent in AD cases. Furthermore, many of our cell-type specific pathways lay directly upstream of NF-kappa B expression. For ex-

ample, reactive oxygen species and TNFs are known inducers of NF-kappa B activity. Conversely, they found that suppression of NF-kappa B reduces apoptosis as well as amyloid-beta plaque deposits [30]. Additionally, disruptions in sleep cycles (a known AD pathology further elaborated on in section 4.4.1) lead to over expression NF-kappa B pathways. These results suggest that further research in understanding how the upstream changes affect NF-kappa B pathways may yield new information regarding AD pathogenesis, as well as understanding how NF-kappa B suppression based drug therapies could be a viable treatment approach.

Finally, the appearance of Parkinson's disease related pathways and prion disease related pathways further confirms the statistical power of our differential splicing analysis. While AD and Parkinson's are two different neurodegenerative diseases, they share a lot of pathological and epidemiological cues [16] As a result, it is not surprising to find Parkinson's disease related pathways as being significantly enriched in our global gene set. Furthermore, studies have also shown AD to exhibit prion disease-like behavior. It has been shown that cell surface prion glycoproteins are able to transduce neurotoxic signals from amyloid-beta oligomers, ultimately leading to neuronal death and cognitive impairment [16].

Chapter 7

RNA velocity analysis reveals sub-population AD-enrichment

In this section, we will be sharing our results from RNA Velocity analysis based on our cell type and meta-analysis differential splicing gene set. Given we identified multiple enriched AD-related pathways in the previous section, we wanted to better understand how these pathways are affected. Specifically, we wanted to know whether differential splicing is due to abrupt changes in gene expression, as well as whether genes are differentially spliced or expressed in specific cellular sub-populations in AD. To achieve this, we performed classical RNA velocity analysis combined with cell clustering.

7.1 RNA velocity inference

We used scVelo to implement RNA dynamical model. We first start by filtering out cells with low read counts in the gene-cell input matrix and selecting the top 2000 most expressed genes as well as the top 2000 differentially spliced genes. We do this to gain a better understanding of how our differentially spliced genes set affect molecular dynamics, and hopefully contrasting it with the results from the top 2000 most expressed genes to uncover potentially novel information regarding AD pathogenesis. Next, counts were normalized by the number of reads as well as log-transforming the

counts. We then performed PCA where the top 30 principal components were used to find k-nearest neighbors which were pooled to smooth the counts. We used the scVelo functions *velocity* and *recoverDynamics* to infer velocities. Next, by projecting the velocities to a lower dimensional space such as UMAP as well as combining it with the Louvain clustering algorithm, we can easily visualize the velocity vectors in UMAP space. Furthermore, we can subset specific genes to identify sub-clusters where the gene of interest has significantly higher velocities, as well as sub-clusters where the gene is differentially expressed.

Since the focus of the study is to further our understanding of AD pathogenesis through an RNA velocity perspective, we hopefully identify genes that are likely to show different kinetic regimes across multiple subpopulations. We expect different cell states to show different splicing kinetics as distinct cell lineages are normally governed by variations in the gene regulatory network. To achieve this, we leverage the dynamical model from above to perform a likelihood ratio test for different kinetics, thus enabling us to cluster cell types into their unique kinetic regimes where each regime can be fitted separately. We will use scVelo's differential kinetic test method, grouping by clusters, to identify specific genes whose kinetic behavior cannot be explained by the overall dynamics.

7.2 AD cohort express less cellular sub-population diversity in astrocytes

Mapping the top 2000 differentially spliced genes in astrocytes in UMAP space with the Louvain clustering algorithm revealed that the AD cohort exhibited less cluster diversity (Figure 7-1). Specifically, the AD cohort is composed of eight louvain clusters (Figure 7-1B) while the control cohort is composed of 11 louvain sub-clusters (Figure 7-1A). This finding suggests that for the same set of genes, AD clearly causes splicing dysregulation compared to control that results in the loss of normal RNA trajectory diversity. Furthermore, the control cohort also exhibits more diversity in the main

cluster (Figure 7-1A: clusters 0, 1, 2, 3, 4, 5, 7, 8, 9, 11), compared to the AD main cluster (Figure 7-1B: clusters 0, 1, 2, 3, 4, 6). We also note that there is a large difference between the UMAP representations of the two cohorts (Figure 17). This could possibly be explained by the fact that the input gene set we used was the top differentially spliced genes between cohorts, and not the quantitatively optimal gene set based on connectivity and min-count thresholds.

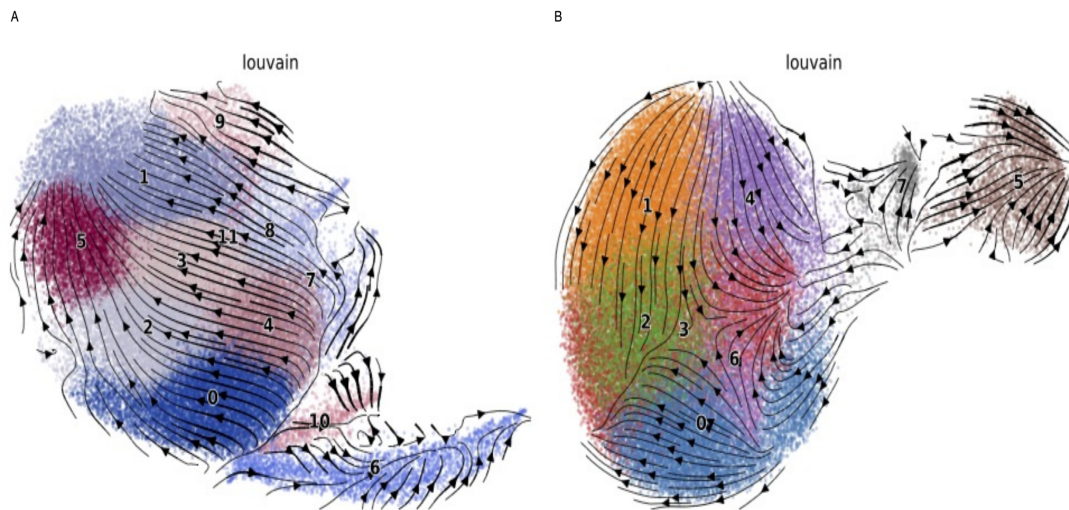


Figure 7-1: Velocity vectors overlaid onto the cell's UMAP space for astrocytes. Input gene set was the top 2000 differentially spliced genes in astrocytes. A: vector field for non-AD cohorts. B: vector field for AD individuals

7.3 Differential splicing genes show sub-population specific enrichment in AD cohorts

We selected five different genes in our gene specific analysis: two different splicing factors, SRRM2 and PNISR; known AD gene GOLGA4; DNA repair related gene, XRCC4; and OXR1, which plays a major role in regulating oxidative stress in the cell. With GOLGA4, we can see that both the velocity and expression patterns are distinctly different between the two cohorts (Figure 7-2A/B). Looking at the velocity fields, in the AD cohort, GOLGA4 is more localized. We can see that cluster 7 in the

AD cohort exhibits the highest expression value compared to the control cohort that shows clusters 0 and 4 (Figure 7-1) as sharing the highest expression. Interestingly, in the control cohort, the velocity vectors also display the highest values in both those clusters as well, whereas in the AD samples, the velocity appears to be localized in cluster 4 (Figure 7-1). Future work can be focused on better understanding the relationship between velocity trajectories and downstream expression patterns in cell type subclusters in AD. This pattern is also observed in XRCC4 where we see higher expression patterns in cluster 7 whereas the velocity vectors appear to be maximal in cluster 7 as well as clusters 4 and 0. Interestingly, in the case of SRRM2, the control cohort appears to have the velocity vectors relatively more localized- there is a larger difference between the minimal cluster (cluster 1) compared to the maximal (cluster 0). Whereas in the case of the AD cohort, the maximal velocity clusters are still localized around clusters 7 and 4, the relative difference between the clusters is significantly less. We hypothesize that this could be due to the observed global dysregulation of splicing factors in AD. We also observe this pattern with PNISR, a splicing factor gene, where in AD cohorts the relative difference in expression values in the main cluster is significantly smaller than the difference observed in the control cohort.

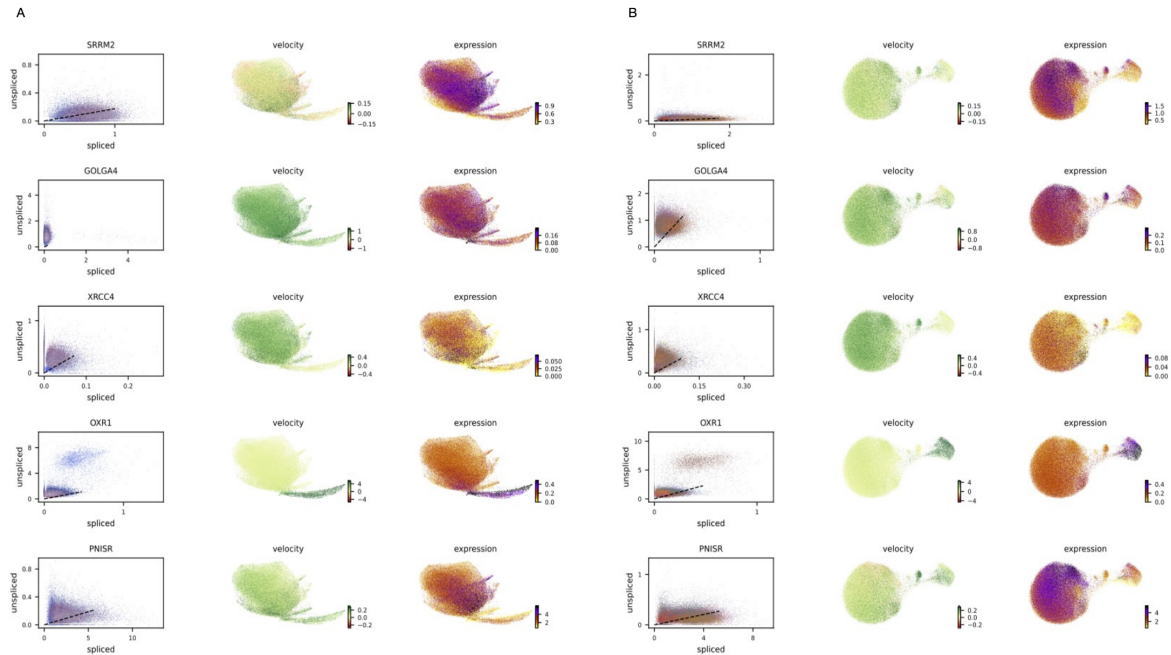


Figure 7-2: Gene specific analysis for SRRM2, GOLGA4, XRCC4, OXR1, and PNISR in astrocytes. Each gene has a spliced/unspliced plot, velocity vector field overlaid in UMAP space, and expression patterns mapped in UMAP space. A: plots for the control cohort. B: plots for AD individuals

7.4 Astrocyte cluster specific analysis of potential drive genes

Given the variation of both velocity flows and expression patterns in specific sub-clusters, we wanted to understand what the main driver genes are for each identified cluster (Figure 7-3). Interestingly, clusters 0, 1, 2, 3, 4, 6 share many of the top driver genes such as SLC1A4, AHCYL1, BCL6, and SYTL4. These clusters also correspond to the main cluster in UMAP space, suggesting these genes are important in characterizing AD specific astrocyte trajectories. On the other hand, clusters 5 and 7, which are the two clusters that also display a lot of localized activities in our genes of interest, do not share many of the top genes compared to the other clusters. We believe that better understanding how these driver genes influence AD progression with respect to splicing dysregulation could yield novel information regarding AD

pathogenesis.

	0	1	2	3	4	5	6	7
0	SLC1A3	SLC1A3	SLC1A3	GRAMD1C	SLC1A3	NPY1R	SLC1A3	RASGRP3
1	AHCYL1	AHCYL1	AHCYL1	BCL6	AHCYL1	SYT16	GRAMD1C	SHROOM4
2	BCL6	GRAMD1C	GRAMD1C	SLC1A3	GRAMD1C	PARM1	BCL6	MOG
3	GRAMD1C	BCL6	BCL6	AHCYL1	BCL6	GABRG2	COL27A1	FRMD4B
4	PMP2	SYTL4	SYTL4	PMP2	PLCE1	SCN3A	PMP2	TMCC3
5	FTL	PLCE1	PMP2	PLSCR1	ADAMTSL3	PPARGC1B	SYTL4	RCAN2
6	ITPKB-AS1	COL27A1	COL27A1	YBX3	TJP2	SLC38A11	CIT	HSD11B1
7	SYTL4	CLIP1	CIT	WARS	COL27A1	SNTG1	PLCE1	CDK6
8	COL27A1	PMP2	PLCE1	CLIP1	SYTL4	DNAH14	NR4A1	SLC5A11
9	CIT	CIT	NR4A1	ACTN1	SMAD9	FGF12	COL12A1	ENPP2

Figure 7-3: Cluster specific driver genes for astrocytes in AD individuals. Computed by finding the partial gene likelihoods for each cluster.

7.5 AD microglia cells exhibit less cluster diversity

In our microglia analysis, we also observe the same pattern seen in astrocytes: the control cohort has more diversity of clusters in UMAP space: 13 as opposed to 8 (Figure 7-4A/B). Interestingly, we noticed that cells in the AD cohort exhibit a unique feature in that there is seemingly a unique group of cells that follow a different trajectory (light purple in Figure 7-4B); however, the magnitude does not seem to be great enough for the Louvain algorithm to consider it as an independent cluster. Given that this phenomenon exists for both microglia and astrocyte cell types, we feel as though there is much work to be done in better understanding the decreased diversity of cluster trajectories between AD and non-AD cohorts.

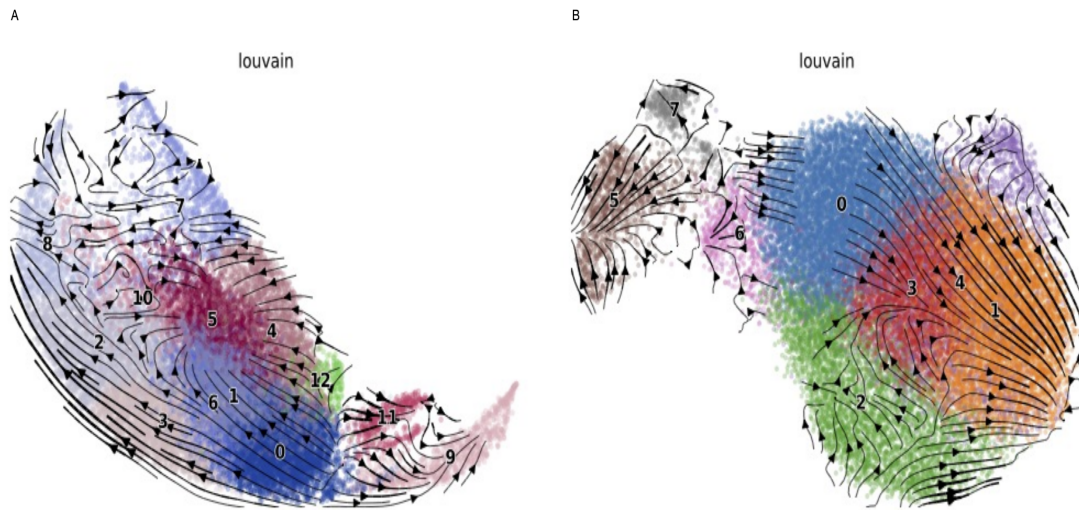


Figure 7-4: Velocity vectors overlaid onto the cell's UMAP space for microglia. Input gene set was the top 2000 differentially spliced genes in microglia. A: vector field for control cohorts. B: vector field for AD individuals

7.6 Cluster localization differential between velocity and expression in OXR1 and PNISR

Once again, we identified the same mismatch between maximal velocity clusters and maximal expression clusters in microglia AD cohorts (Figure 7-5). Specifically for OXR1 and PNISR, we notice that the maximal expression clusters are 5 and 0, while maximal velocity clusters are 5 and 6 (Figure 7-4). In the case for OXR1, while both velocity and expression clusters are both in cluster 5, velocity distribution in cluster 5 clearly is more localized within the cluster compared to the expression distribution (Figure 7-5). This suggests that there may be a unique sub-cluster within cluster 5 that arises because of AD pathology. PNISR also displays a similar pattern. We can see that the velocity trajectories are maximal in clusters 7 and 6, while expression values are maximized in cluster 0 (Figure 7-5B). Whereas in the control cohort both velocity and expression values are maximized in cluster 0 (Figure 7-5A).

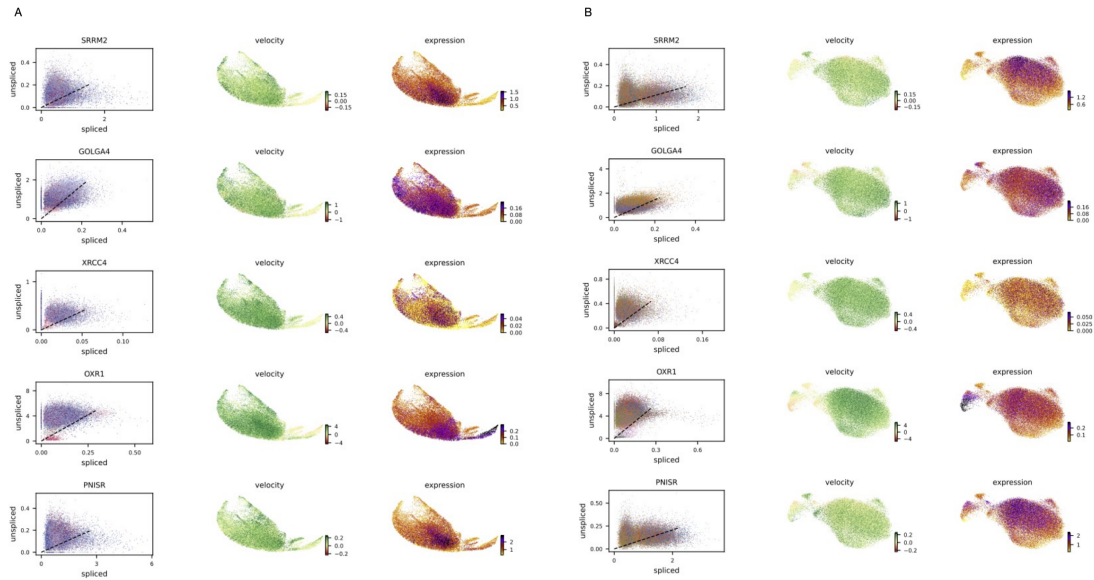


Figure 7-5: Gene specific analysis for SRRM2, GOLGA4, XRCC4, OXR1, and PNISR in microglia. Each gene has a spliced/unsplliced plot, velocity vector field overlaid in UMAP space, and expression patterns mapped in UMAP space. A: plots for the control cohort. B: plots for AD individuals

7.7 Microglia cluster specific AD driver genes

Analysis of driver genes for each identified cluster in AD microglia cohorts also revealed the same pattern as seen in astrocytes. Clusters 0, 1, 2, 3, 4, 6, which correspond to the main cluster in UMAP space share a lot of the top genes: TGFBR1, JAK3, and IL4R (Figure 7-6). Whereas cluster 5 and 7, which correspond to the standalone clusters in UMAP space, contain a significantly different set of driver genes. Interestingly, SHROOM4 was identified as the top gene in cluster 7 for both microglia and astrocyte populations. While we are not entirely sure why this is the case, SHROOM4 is a protein that is involved with cytoskeletal architecture, and as a result, dysregulation of this gene could result in fundamental changes to the cell's organization, thereby altering the trajectories of RNA transcripts.

	0	1	2	3	4	5	6	7
0	TGFBR1	JAK3	TGFBR1	TGFBR1	TGFBR1	SNTG1	TGFBR1	SHROOM4
1	JAK3	IL4R	JAK3	JAK3	IL4R	ELAVL3	JAK3	C10orf90
2	IL4R	TGFBR1	IL4R	IL4R	JAK3	SYT14	IL4R	MOG
3	CACNA1A	CACNA1A	BCL6	HIF1A	CACNA1A	AGBL4	GPR155	FGF1
4	HIF1A	BCL6	HIF1A	CACNA1A	BCL6	KCTD8	HIF1A	RFTN2
5	GPR155	SPP1	GPR155	GPR155	HIF1A	DZIP1	SYNE1	KCTD8
6	XAF1	HIF1A	PTPRC	PDK4	P4HA1	MGAT4C	TLE1	PPFIBP1
7	PTPRC	PTPRC	XAF1	PTPRC	HSP90AA1	DNAH6	GALNT3	AF131217.1
8	PDK4	NAMPT	PDK4	XAF1	BAG3	PTPRS	YES1	ITGB8
9	SMC4	P4HA1	SMC4	TXNIP	ASPHD1	RP11-761E20.1	MDFIC	ANO4

Figure 7-6: Cluster specific driver genes for microglia in AD individuals. Computed by finding the partial gene likelihoods for each cluster.

Chapter 8

Conclusion

8.1 Summary

In this study, we presented a novel method that leverages inferred spliced/unspliced count from RNA velocity along with a beta-binomial model in hopes of revealing differential splicing in AD progression. We first began by validating the statistical power of the model before applying it to the splicing counts of the ROSMAP cohort. We then identified cell-type specific differentially spliced genes as well as identifying a meta level set of genes that are differentially spliced across all cell types. We found that our differential splicing gene set confers with our current understanding of AD pathogenesis, as well as, revealing dysregulation of splicing factors that have been well studied in other neurodegenerative diseases but not AD. Our method was also validated by our results where known AD genes, such as GOLGA4, exhibited differential splicing characteristics. However, unfortunately due to random priming effects, we were not able to gain a full transcriptomic understanding of the AD related changes that occur.

Next, given our novel set of differential splicing genes, we wanted to gain a better understanding of the pathways most affected. Our pathway enrichment analysis revealed that many of the top pathways were either directly related to AD, such as reactive oxygen species and insulin resistance, or were splicing related pathways. These pathways also represent potential targets for future research in AD drug therapies.

Finally, to better understand the effects such differential splicing in AD, we performed RNA velocity analysis combined with cell clustering. We identified stronger cluster localization of velocity trajectories and expression values in the AD cohort and identified differential expression in specific cell sub-populations in AD cases. Overall, our findings point to differential splicing as an AD hallmark, but more importantly, showcases the efficacy of our novel differential splicing analysis in helping understand disease etiology.

8.2 Future Work

In this study, we presented a novel analysis method based on RNA velocity inferred splicing counts to identify differentially spliced genes in AD. While our results have conferred with current understanding of the disease, we believe there is still much work to be done. We believe that the next step will be experimental validation of our results. Specifically, we propose immunohistochemistry studies focusing on SRRM2 and other key genes we identified to better understand their function in vitro, spatial transcriptomic studies to observe if structural changes are a cause for difference in RNA velocity trajectories, and also RNAscope studies to gain a deeper understanding into the expression of our differentially spliced genes. We also believe collaboration with other groups centered around the computational validation and improvement of the model will yield major improvements.

Finally, we believe, on a higher level, more focus should be placed on better understanding random priming. Especially, as the cost of scRNA-seq studies go down, understanding the technical limitations and artifacts of the technology will not just save time, but also provide avenues of improvement. We hope that as the technology matures, we can gain a fuller understanding of the changes to the transcriptomic landscape as a result of AD.

Bibliography

- [1] Tomonori Aikawa. Abca7 and pathogenic pathways of alzheimer's disease. *Brain Sciences*, 2018.
- [2] Jonny Daborg Ilse Riebe Pontus Wasling Henrik Zetterberg Eric Hanse Andreas Bjorefeldt, Ulf Andreasson. Human cerebrospinal fluid increases the excitability of pyramidal neurons in the in vitro brain slice. *The Journals of Physiology*, 593, January 2015.
- [3] Sarah A Teichmann Tapio Lonnberg Ashrasful Haque, Jessica Engel. A practical guide to single-cell rna-sequencing for biomedical research and clinical applications. *Genome Medicine*, 2017.
- [4] Marwan N Sabbagh Boris Decourt, Debomoy K Lahiri. Targeting tumor necrosis factor alpha for alzheimer's disease. *Current Alzheimers Research*, 14(4), January 2018.
- [5] Zev J Gartner Christopher S McGinnis, Lyndsay M Murrow. Doubletfinder: Doublet detection in single-cell rna sequencing data using artificial nearest neighbors. *Cell Systems*, 2019.
- [6] Xueping Chen Danshen Zhang Chunyan Guo, Li Sun. Oxidative stress, mitochondrial damage and neurodegenerative diseases. *Neural Regeneration Research*, 8(21), July 2013.
- [7] Katarzyna Marta Zoltowska Thomas Enzlein Natalie S Ryan Antoinette O'Connor Maria Szaruga Elizabeth Hill Rik Vandenberghe Nick C Fox Lucia Chavez-Guitierrez Dieter Petit, Sara Gutierrez Fernandez. A profiles generated by alzheimer's disease causing psen1 variants determine the pathogenicity of the mutation and predict age at disease onset. *Molecular Psychiatry*, 2022.
- [8] Andrew D Smith Egor Dolzhenko. Using beta-binomial regression for high-precision differential methylation analysis in multifactor whole-genome bisulfite sequencing experiments. *BMC Bioinformatics*, 2014.
- [9] David M Holtzman Erik S Musiek, David D Xiong. Sleep, circadian rhythms, and the pathogenesis of alzheimer disease. *Experimental and Molecular Medicine*, 47(148), March 2015.

- [10] Anand Goswami Joachim Weis Gary Brook Alfred Yamoah Jaspen Anink Eleonora Aronica Stefan Fritz Carmen Huck Carola Schipke Oliver Peters Martin Tepel-Heidi Noels Vera Jankowski Felix Kork, Joachim Jankowski. Golgin a4 in csf and granulovacuolar degenerations of patients with alzheimer disease. *Neurology*, 91(19), November 2018.
- [11] Tara Chari Lior Pachter Gennady Gorin, Meichen Fang. Rna velocity unraveled. *bioRxiv*, 2022.
- [12] Charles Curtis Sang Hyuck Lee Claire Troakes Richard J.B Dobson Stephen J Newhouse Hamel Patel, Angela K Hodges. Transcriptomic analysis of probable asymptomatic and symptomatic alzheimer brains. *Brain, Behavior, and Immunity*, May 2019.
- [13] Zhuyu Peng Fan Gao Shahin Mohammadi Jennie Z Young Madhvi Menon Liang He Fatema Abdurrob Xueqiao Jiang Anthony J Martorell Richard M Ransohoff Brian P Hafler-David A Bennett Manolis Kellis Li-Huei Tsai Hansruedi Mathys, Jose Davila-Velderrain. Single-cell transcriptomic analysis of alzheimer’s disease. *Nature*, 2019.
- [14] Caitlin Davies Lewis Wilkins Desclan King Jamie Rose Makis Tzioras Jane H Tulloch Colin Smith Tara L. Spire-Jones Hatice Kurucu, Marti Colom-Cadena. Inhibitory synapse loss and accumulation of amyloid beta in inhibitory presynaptic terminals in alzheimer’s disease. *European Journal of Neurology*, July 2021.
- [15] Lee Ji Hyun. Bang Duhee Hwang, Byungjin. Single-cell rna sequencing technologies and bioinformatics pipelines. *Experimental and Molecular Medicine*, 2014.
- [16] Bigqian Liu Jiayi Zhou. Alzheimer’s disease and prion protein. *Intractable Rare Disease Research*, 2(2), May 2013.
- [17] Jiyeon Lee Floy R Stewart Deborah K Verges Luz M Silverio Guojun Bu Steven Mennerick David M Holtzman John R Cirrito, Jae-Eun Kang. Endocytosis is required for synaptic activity-dependent release of amyloid- in vivo. *Neuron*, 58(1), April 2008.
- [18] Jyotshna Kanungo. Dna-dependent protein kinase and dna repair: relevance to alzheimer’s disease. *Alzheimer’s Research Therapy*, 5(13), April 2013.
- [19] Su Wu Suizhen Lin Deanna Koger Cynthia DeLong Jeffrey C Hansen Patrick M Sullivan Steven M Paul Kelly R Bales, Feng Liu. Human apoe isoform-dependent effects on brain -amyloid levels in pdapp transgenic mice. *Journal of Neuroscience*, 2009.
- [20] Linda Koch. Altered splicing in alzheimer transcriptomes. *Nature Reviews Genetics*, 2018.

- [21] Fumiaki Mori Akiyoshi Kakita Hitoshi Takahashi Koichi Wakabayashi Kunikazu Tanji, Tetsu Kamitani. Trim9, a novel brain-specific e3 ubiquitin ligase, is repressed in the brain of parkinson's disease and dementia with lewy bodies. *Neurobiology of Disease*, 38(2), May 2010.
- [22] Liyong Wang Alexandra Guevara Carlos Singer Jeffery Vance Spyridon Papatropoulos Lina A Shehadeh, Kristine Yu. Srrm2, a potential blood biomarker revealing high alternative splicing in parkinson's disease. *PLoS One*, 5(2), February 2010.
- [23] Gioele La Manno. Rna velocity of single cells. *Nature*, 2018.
- [24] Bengt Winblad Ovidiu Alexandru Bajenaru Nenad Bogdanovic Mihaela Oana Romanitan, Bogdan O Popescu. Occludin is overexpressed in alzheimer's disease and vascular dementia. *Journal of Cellular Molecular Medicine*, 11(3), May 2007.
- [25] Afroditi Tsitsou-Kampeli Hadas Keren-Shaul Hila Ben-Yehuda Pierre WEill-Raynal Liora Cahalon Alex Kertser Kuti Baruch Ido Amit Assaf Weiner Michal Schwartz Neta ROsenzweig, Rar Dvir-Szternfeld. Pd-1/pd-l1 checkpoint blockade harnesses monocyte-derived macrophages to combat cognitive impairment in a tauopathy mouse model. *Nature Communications*, 10(465), January 2019.
- [26] Maotian Zhou-Sean R Kundinger-Marla Gearing James J Lah Allan I Levey Joshua M Shulman Nicholas T Seyfried Qi Guo, Eric B Dammer. Targeted quantification of detergent-insoluble rna-binding proteins in human brain reveals stage and disease specific co-aggregation in alzheimer's disease. *Frontiers in Molecular Neuroscience*, March 2021.
- [27] Richard Robinson. Rna polymerase promotes splicing, prevents degradation. *PLoS Biology*, 4(6), June 2006.
- [28] Shyh-Jer Huang Hung-Jen Chen Pin-Ru Chen Ya-Hsien Lin Yu-Chen Ho Wen-Lin Chang Chang-Hai Tsai Woei-Cherng Shyu Shinn-Zong Lin Ru-Huei Fu, Shih-Ping Liu. Aberrant alternative splicing events in parkinson's disease. *Cell Transplantation*, 2013.
- [29] Richard D Semba. Perspective: The potential role of circulating lysophosphatidylcholine in neuroprotection against alzheimer disease. *Advances in Nutrition*, 11(4), July 2020.
- [30] Ilias Kounatidis Simon Vann Jones. Nuclear factor-kappa b and alzheimer disease, unifying genetic and environmental risk factors from cell to humans. *Frontiers of Immunology*, 8, December 2017.
- [31] Paul Cheng Tanya Todorova, Florian Bock. Parp13 and rna regulation in immunity and cancer. *Trends in Molecular Medicine*, 21(6), June 2015.

- [32] Kostecka Z Hromadka T Vargova G, Vogels T. Inhibitory interneurons in alzheimer's disease. *Bratislava Medical Journal*, 2018.
- [33] Archana Juvekar Viplav Kshirsagar, Chetan Thingore. Insulin resistance: a connecting link between alzheimer's disease and metabolic disorder. *Metabolic Brain Disease*, 36, 2020.
- [34] Yan Gu Mathilda Jing Chow Jingyi Peng Kuncheng Zhao-Damu Tang Xiaozeng Lin, Anil Kasoor. Contributions of dna damage to alzheimer's disease. *International Journal of Molecular Science*, 21(5), March 2020.
- [35] Ying Hu Lei Zhang Ju Wang Yan Shi Hu, Juncai Xin. Analyzing the genes related to alzheimer's disease via a network and pathway-based approach. *Alzheimer's Research Therapy*, 9(29), April 2017.
- [36] Jennifer D McBride Sneha Narasimhan Hyesung Kim Lakshmi Changolkar Bin Zhang-Ronald J Gathagan-Cuiyong Yue Christopher Dengler Anna Stieber-Magdalena Nitla Douglas A Coulter Ted Abel Kurt R Brunden John Q Trojanowski Virginia M-Y Lee Zhuohao He, Jing L Guo. Amyloid- plaques enhance alzheimer's brain tau-seeded pathologies by facilitating neuritic plaque tau aggregation. *Nature Medicine*, 2017.