

**SINGLE-CELL TRANSCRIPTIONAL PROFILING OF
HUNTINGTON'S DISEASE IN HUMAN AND MOUSE
MODELS**

by

SERGIO SEBASTIAN PINEDA

B.S. Electrical Engineering
Texas A&M University, 2018

SUBMITTED TO THE DEPARTMENT OF ELECTRICAL
ENGINEERING AND COMPUTER SCIENCE IN PARTIAL
FULFILLMENT OF THE REQUIREMENTS FOR THE DEGREE OF
MASTER OF SCIENCE IN ELECTRICAL ENGINEERING AND
COMPUTER SCIENCE

AT THE
MASSACHUSETTS INSTITUTE OF TECHNOLOGY

SEPTEMBER 2020

© Massachusetts Institute of Technology 2020. All rights reserved.

Signature of Author

Department of Electrical Engineering and Computer Science

August 28, 2020

Certified by

Manolis Kellis

Professor of Electrical Engineering and Computer Science

Thesis Supervisor

Certified by

Myriam Heiman

Associate Professor of Neuroscience

Thesis Supervisor

Accepted by

Leslie A. Kolodziejcki

Professor of Electrical Engineering and Computer Science

Chair, Department Committee on Graduate Students

Single-Cell Transcriptional Profiling of Huntington's Disease in Human and Mouse Models

by

Sergio Sebastian Pineda

Submitted to the Department of Electrical Engineering and Computer Science,
on August 28, 2020, in partial fulfillment of the
requirements for the degree of
Master of Science in Electrical Engineering And Computer Science

Abstract

Uncovering the mechanisms that drive cell type-specific vulnerability differences in Huntington's disease (HD) is an imperative prerequisite to finding a reliable therapeutic target. The most affected regions and cell types likely possess vulnerability factors or lack protective factors possessed by other cell types that lead to their early dysfunction and enhanced loss, but such factors remain elusive. In order to characterize the cell type-specific responses induced by mutant huntingtin (mHTT), we applied single-nucleus RNA sequencing to profile gene expression changes in human HD and two commonly used mouse models of HD. In the process, we produced the first molecular atlas of the human neostriatum, identified new and previously unobserved cell types, and examined the molecular conservation of this brain region across species. We also developed and documented new computational methods and techniques for curating and analyzing single-cell data from post-mortem human brain samples.

Thesis Supervisor: Manolis Kellis

Title: Professor of Electrical Engineering and Computer Science

Thesis Supervisor: Myriam Heiman

Title: Associate Professor of Neuroscience

Contents

1	Introduction	9
2	Experimental Formulation	13
2.1	Experimental Cohorts	13
2.2	Novel Strategies for SnRNA-seq Data Analysis	14
3	Single-Cell Resolution Profiling of the Neostriatum	17
3.1	The Mouse Neostriatum	17
3.2	The Human Neostriatum	20
3.3	Molecular Conservation of the Neostriatum Within and Across Species	22
4	Analysis of Gene Expression Changes in HD and Mouse Models of HD	27
4.1	Methods for Differential Gene Expression Analysis	27
4.2	Differential Gene Expression Analysis of HD Mouse Models	31
4.3	Differential Gene Expression Analysis of Human HD	37
5	Discussion	43

List of Figures

3-1	ACTIONet plot of cell types identified in the neostriatum of 9 week-old R6/2 by snRNA-seq [n = 107,412 cells across seven R6/2 and eight isogenic control mice].	18
3-2	ACTIONet plot of cell types identified in the neostriatum of 6 month-old zQ175 by snRNA-seq [n = 60,790 cells across four zQ175 and four isogenic control mice].	19
3-3	ACTIONet plot of cell types identified in the caudate and putamen of human HD and control individuals by snRNA-seq [n = 128,797 cells across 25 tissue samples spanning seven HD and seven unaffected control individuals].	21
3-4	Average fraction of SPNs across grades of HD relative to control samples.	22
3-5	Similarity of neostriatal cell types by transcription-wide correlation of gene expression across all cell types identified in 6 month-old control mice.	23
3-6	Similarity of neostriatal cell types by transcription-wide correlation of gene expression across all cell types identified in human unaffected control samples.	24
3-7	Similarity of neostriatal cell types across species by transcription-wide correlation of ortholog gene expression across all cell types identified in human and mouse control samples.	26
4-1	Comparison of DEGs across methods in each experimental cohort. Spearman correlation by log-fold change.	30
4-2	Top 7 most statistically significant up- and down-regulated genes in SPNs and other highly abundant cell types of R6/2.	32

4-3	Top 7 most statistically significant up- and down-regulated genes in SPNs and other highly abundant cell types of zQ175.	33
4-4	Pearson’s correlation of DEGs by log-fold change across cell types in R6/2 [FDR-adjusted p -value < 0.001 and \log_2 fold-change > 0.1]	34
4-5	Pearson’s correlation of DEGs by log-fold change across cell types in zQ175 [FDR-adjusted p -value < 0.001 and \log_2 fold-change > 0.1]	35
4-6	Pearson’s correlation of DEGs by Z-score between R6/2 and zQ175.	36
4-7	Top 7 most statistically significant up- and down-regulated genes in SPNs and other highly abundant cell types in Human HD.	39
4-8	Differential expression of reactive astrocyte marker genes from Liddelow et al (2017) in astrocytes of Human HD and HD mouse models.	40
4-9	Pearson’s correlation of DEGs by log-fold change across cell types in Human HD [FDR-adjusted p -value < 0.001 and \log_2 fold-change > 0.1]	41
4-10	Comparison of DEGs between Human HD and HD mouse models by method [Pearson’s correlation of Z-scores].	42

Chapter 1

Introduction

Huntington's disease (HD) is a monogenic, autosomal-dominant neurodegenerative disease caused by the expansion of the CAG repeat tract in exon 1 of the huntingtin (*HTT*) gene. Whereas the CAG repeat in *HTT* normally encodes for approximately 20 consecutive glutamine residues in the huntingtin protein, HD-associated expansions in 1 in 10,000 individuals lead to 36+ CAG repeats in one of the *HTT* alleles. The resulting mutant isoform of the huntingtin protein [1] leads to progressive motor dysfunction, cognitive decline, brain atrophy, and ultimately premature death [2]. HD pathology exhibits a pattern of cell type- and region-specific vulnerability, starting with the selective loss of inhibitory neurons in the striatum.

Vonsattel et al. [3] devised a classification of pathological progression of HD along five grades (0-4), with grade 0 displaying no neuropathological abnormalities, grade 1 displaying 50% loss of striatal neurons, and grade 4 displaying 95% loss of striatal neurons. In particular, the majority of neuronal loss results from the death of a specific class of GABAergic neurons commonly known as spiny projection neurons (SPNs) or medium spiny neurons (MSNs) which can be broadly categorized into two sub-classes based on their projections and preferential expression of D1-type or D2-type dopamine receptors. Among SPNs, *DRD2*-expressing 'indirect pathway' SPNs (iSPNs) have been shown to be affected at earlier stages of HD than *DRD1*-expressing 'direct pathway' SPNs (dSPNs) [4–6]. In addition to this known differential vulnerability by cell type, SPNs in different subregions of the neostriatum display region-specific differential vulnerability in HD. Striatal neurodegeneration follows a

caudo-rostral gradient with the dorsal neostriatum, and specifically the caudate nucleus, suffering from greater neuronal loss than the rest of the neostriatum. Moreover, a subset of SPNs are compartmentalized into earlier-born and functionally distinct “patches” called striosomes [7] which have been hypothesized to show earlier vulnerability compared to those in the surrounding matrix [8].

Although it is still not fully understood how trinucleotide CAG repeat expansions in the mutant huntingtin gene (*mHTT*) lead to the death of striatal neurons, *mHTT* has been linked to both toxic gain-of-function as well as loss-of-function and effects on diverse cellular processes that include synaptic neurotransmission, neurotrophin signaling, mitochondrial function, DNA damage repair, axonal and vesicular transport, calcium signaling, and gene transcription [1]. In support of gene transcription dysregulation as having an important role in Huntington’s disease pathobiology, mutant Huntingtin protein is most toxic when localized in the nucleus [9, 10] and the degree of transcriptional dysregulation in HD correlates with brain area vulnerability in human HD samples [11].

High throughput gene expression assays have been widely applied as a means to study the molecular mechanisms underlying virtually every disease. Over the last two decades, numerous groups have applied such methods to study transcriptional changes in Huntington’s disease in both humans and mouse models by way of expression microarrays [11] or high-throughput RNA-sequencing [12]. While much about the disease has been learned through these studies, they have focused on gene expression changes in bulk tissue, which obscures cell-type specific mechanisms underlying selective vulnerability. More recent studies have attempted to investigate these mechanisms by employing cell type-specific profiling methods such as translating-ribosome affinity purification [12, 13], but these methods are lower throughput and require the use of transgenics and hence cannot be used to study disease mechanisms in humans.

Uncovering the genetic and molecular drivers of cell type-specific vulnerability differences in HD is an imperative prerequisite to finding a reliable therapeutic target or method of intervention. The most vulnerable regions and cell types likely possess vulnerability factors or lack protective factors possessed by other cell types that lead to their early dysfunction and enhanced loss, but such factors remain elusive. In an effort to broadly characterize striatal

cell types and their specific responses to mHTT, we used droplet-based single-nucleus RNA sequencing (snRNA-seq) [14, 15] to profile gene expression changes in human HD and mouse models of HD. In doing so, we produced high resolution molecular atlases of the human and mouse neostriata, identified novel and previously unobserved cell types, and assessed the physiological and molecular conservation of this brain region across species. We also developed and adapted new computational techniques for curating and analyzing single-cell data from post-mortem human brain samples. This study across mouse and human cell types revealed a previously unrecognized complexity of population-specific gene expression changes in response to mHTT, as well as a conservation of many of these gene expression changes in cell types that have not been previously classified as vulnerable in HD.

Chapter 2

Experimental Formulation

2.1 Experimental Cohorts

To gain cell type-specific resolution in gene expression profiling, we used single-nucleus RNA-sequencing to study two of the most widely used mouse models of HD:

- The R6/2 model which contains a knock-in of human *HTT* exon 1 with 160 CAG repeats in a B6/CBA mixed background [16]. The product of the m*HTT* knock-in in this model produces a truncated N-terminal fragment isoform of mHTT that localizes to the nucleus, similar to what is observed in human HD patients. This model was harvested at 9 weeks of age which represents an early/mid phenotypic time point.
- The zQ175 delta neo (zQ175DN, henceforth referred to as zQ175) model, which has a chimeric full-length huntingtin gene with the endogenous *Htt* exon 1 substituted by a human m*HTT* exon 1 knock-in containing approximately 190 CAG repeats in a C57BL/6J background [17–20]. This is a much less aggressive model compared to the R6/2 and has a life expectancy of 22-24 months, compared to 11-12 weeks for the R6/2. This model was harvested at 6 months of age, representing an early phenotypic time point and motivated by a bulk tissue study by Langfelder et al. [12] which reported this as the first robust symptomatic time point for transcriptional dysregulation in zQ175.

In both cohorts, litter mates that lacked the *mHtt* transgene were used as experimental controls.

Additionally, we collected and profiled 25 post-mortem caudate nucleus and putamen samples from 14 human HD grade 2-4 and unaffected control individuals in order to study gene expression changes directly relevant to the human disease and for comparison against the mouse models.

2.2 Novel Strategies for SnRNA-seq Data Analysis

Our experience with single-cell experiments across humans and mice and our reanalysis of studies by other groups has led to observation that human single-cell gene expression studies tend to have issues with reproducibility [13, 21–23]. We hypothesize that the variable nature of post-mortem human tissue samples, especially brain tissue, means that a number of non-disease-associated factors may affect gene expression signatures and exert tremendous influence on the results of downstream analyses.

Analysis of our own human cohorts spanning multiple diseases and brain regions has led to the identification of factors such as RNA integrity, autolysis time, homogenization media, dissection technique, and harvesting and banking protocols which can induce gene expression artifacts of sufficiently large magnitude such that they mask real biological signals. In tissue samples of acceptable quality, we would expect that several of these factors might have a strong effect on gene expression on a fraction of the cell population and that these changes would be visible and perhaps statistically significant at the single cell level, but masked at the bulk-level, resulting in more robust and reproducible results for bulk studies of identical cohorts. We have also found that non gene-expression associated artifacts such as cell type differential abundance are present in single-nucleus experiments and are affected by endogenous factors such as nuclear size and mass on top of the biases introduced by different sample preparation protocols. These observations challenge the widespread misconception that droplet-based single-nucleus assays are an unbiased method for cellular profiling. In consideration of the challenges and pitfalls that accompany our experimental approach, we have devised a number of computational and statistical tools and techniques to mitigate the

effects of such biases in our own studies.

Prior to cluster analysis we employed novel strategies for cell filtering and gene selection that were determined to be essential for accurate annotation (particularly for exceptionally rare cell populations) and reproducible downstream analysis. A typical snRNA-seq analysis workflow would start by pruning cells with low unique molecular identifier (UMI) counts, low library complexity, and high mitochondrial RNA (mtRNA) content. This is insufficient given that nuclei originating from smaller-sized cells will have lower RNA content, droplets containing multiple cells or nuclei will usually have a high library complexity, and high mtRNA abundance can be a disease phenotype. The last is critically important in our case since mtRNA release is a known HD phenotype and high mtRNA content is expected in the most vulnerable cell types [13]. This would also affect nuclear profiles since cytosolic mtRNA transcripts loaded onto endoplasmic reticulum-bound ribosomes will be pulled down with the rest of the nucleus during nuclear isolation. It must be stressed that although these conventional quality control steps are wholly insufficient, they may still capture an appreciable fraction of low quality cells if thresholds are chosen carefully and in the context of the experimental cohort.

Prior to cell type annotation, we took an aggressive cell filtering approach that resulted in a high dropout rate of cells in lower quality samples. Following the conventional workflow, we filtered low library complexity cells and those with UMI counts low enough to be indistinguishable from empty droplets, as well as cells with dramatically higher mtRNA percentages such that they were indistinguishable from noise, lacking clear expression of cell type-specific marker genes. Decomposition by our archetypal analysis for cell type identification (ACTION) algorithm [24] identified archetypes spanning multiple clusters, some of which clearly corresponded to doublets and others which showed a non-cell type- or disease-specific stress response signature, all of which were removed. Further manual doublet detection involved searching for high cell type-specific marker gene expression in clusters where they should not otherwise be present. We then used a network-based filtering approach to prune cells with low k -core number and low node centrality after successive iterations of network reconstruction until no such nodes were produced. Lastly, groups of cells that continued to show a strong batch effect following multiple nearest-neighbor (MNN) [25, 26] correction and batch

orthogonalization [unpublished] were removed. At each iteration we dropped genes that were no longer expressed in a certain percentage of our cell population.

Chapter 3

Single-Cell Resolution Profiling of the Neostriatum

3.1 The Mouse Neostriatum

Cluster analysis of the mouse data sets using our network-based cell type identification framework (ACTIONet) [27] recovered at least 19 distinct cell types and cellular sub-types [Fig. 3-1, 3-2]. Present in all samples across both HD models and their respective controls were major expected cell types including SPNs, astrocytes, oligodendrocytes, microglia, oligodendrocyte progenitor cells (OPCs) as well as numerous sub-types of GABAergic and cholinergic interneurons. Immediately distinguishable were two pairs of large sub-clusters of SPNs, the smaller ones of which corresponded to SPNs situated in the striosomal compartments and the larger ones to those in the surrounding matrix. Also recovered in expectedly smaller quantities were the rarer neurovascular populations (endothelial and mural). As consequence of the proximity of the striatum to the sub-ventricular zone we were also able to recover two populations of ependymal (ciliated and secretory) cells and a large population of cells expressing a number of neurogenesis-related genes, likely corresponding to neuronal and glial progenitors destined for the rostral migratory stream. It should be noted that this population was expectedly smaller in the 6 month-old zQ175 cohort compared to the younger R6/2 (9 weeks) group. Also present across all samples were two SPN-like populations: one recently identified independently by us and other groups (Neur.Adarb2) [13, 28, 29] and was

distinguished by high expression of *Adarb2* and *Foxp2*, and another novel sub-type of iSPN (SPN.D2.Other) with a unique molecular signature that was defined by high expression of *Alk*.

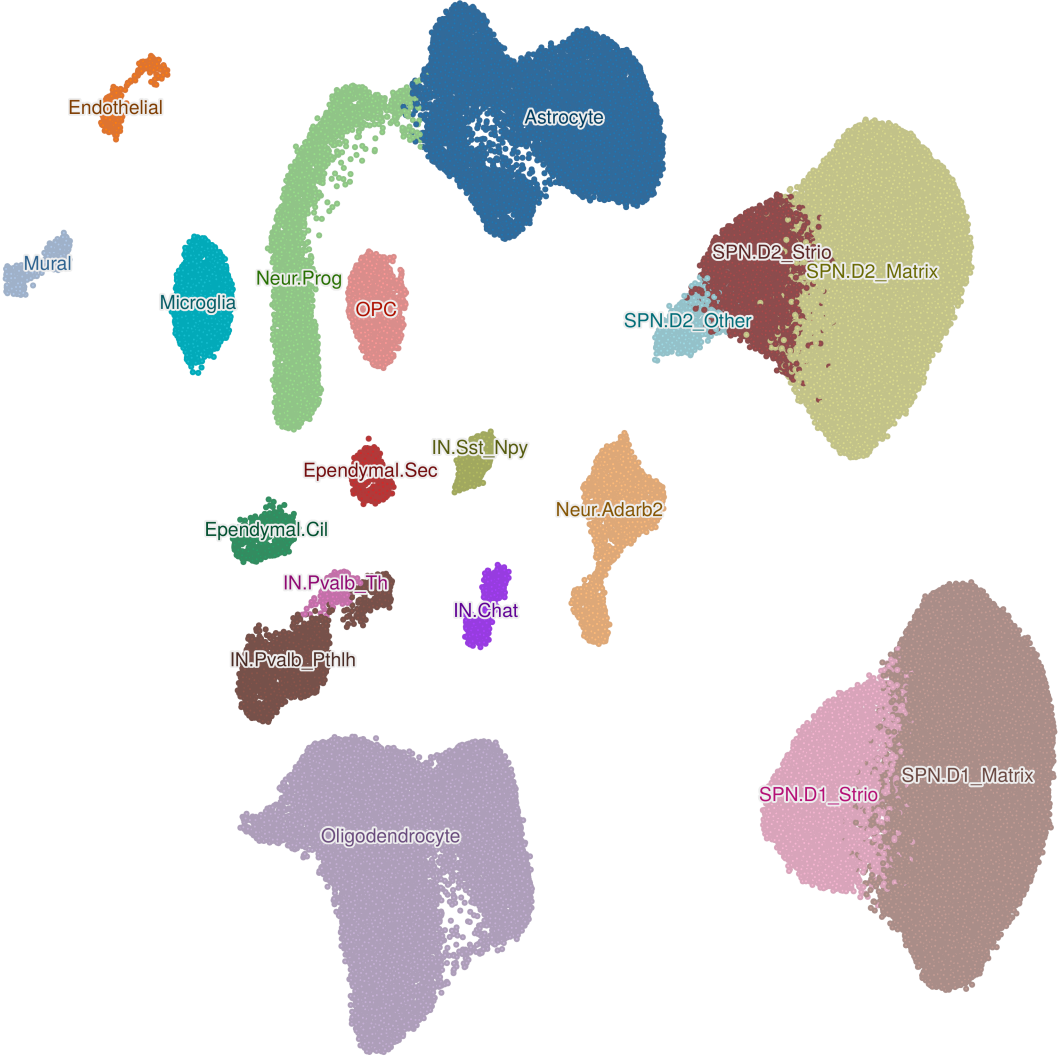


Figure 3-1: ACTIONet plot of cell types identified in the neostriatum of 9 week-old R6/2 by snRNA-seq [n = 107,412 cells across seven R6/2 and eight isogenic control mice].

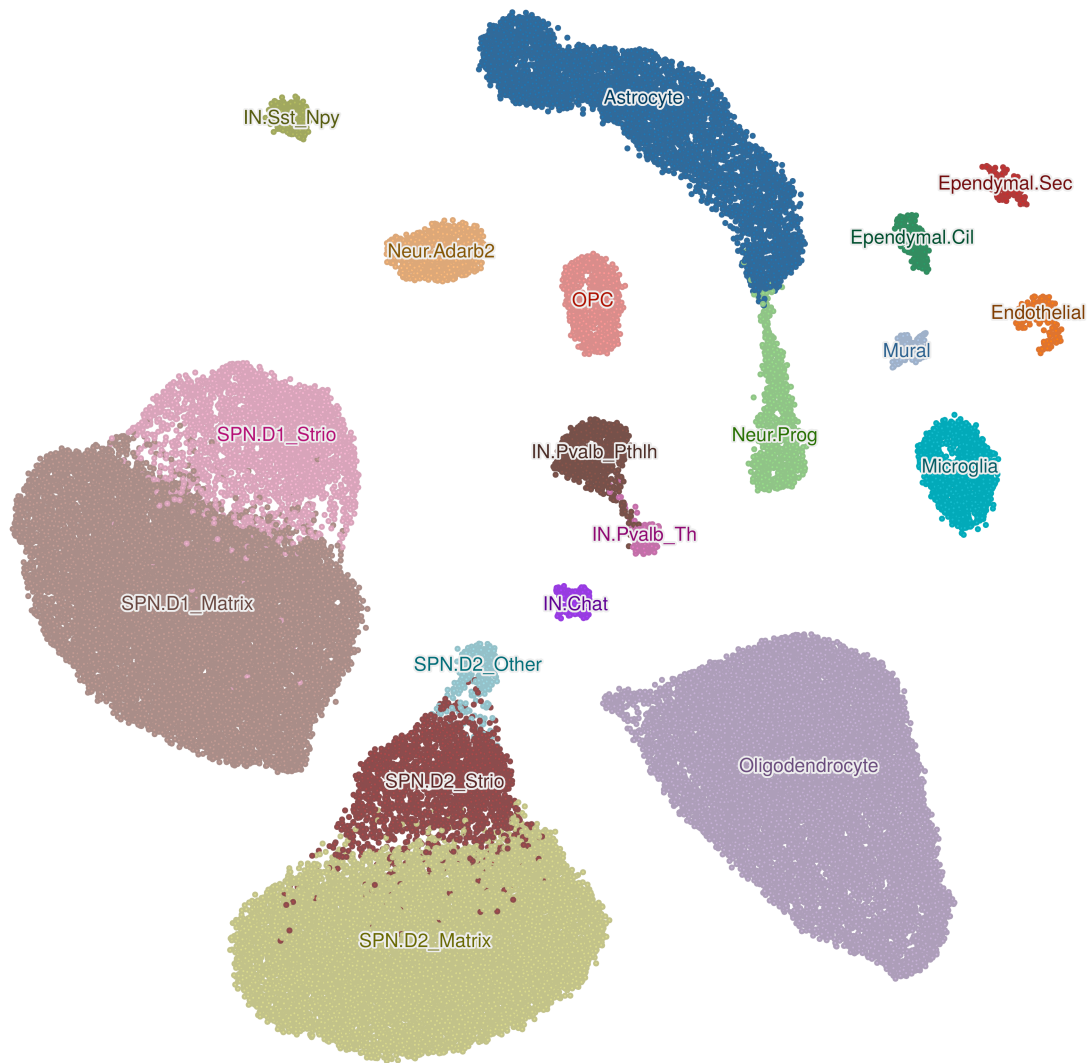


Figure 3-2: ACTIONet plot of cell types identified in the neostriatum of 6 month-old zQ175 by snRNA-seq [n = 60,790 cells across four zQ175 and four isogenic control mice].

3.2 The Human Neostriatum

A concomitant analysis of the human dorsal neostriatum revealed 29 unique populations [Fig. 3-3]. All cell types and sub-types identified in mice were recovered in human with the unnotable exception of neuronal progenitors given the absence of convincing evidence of adult neurogenesis in the human sub-ventricular area. The consistent clusters included the novel iSPN sub-type (SPN.D2.Other) which displays a remarkably similar orthologous transcriptional signature. The deeper sequencing coverage of the human samples allowed us to recover rare cellular sub-populations that we were not able to observe in the mouse data sets. One peculiar cluster which showed high expression of numerous overlapping interneuronal markers was further partitioned, revealing seven small clusters corresponding to a number of interneuronal sub-types including NPY^+ and NPY^- somatostatin-expressing GABAergic interneurons, and 5HT3aR interneurons (some of which co-expressed VIP , and others which did not). We also found interneuron populations that have not been previously reported to exist in the neostriatum such as basket and chandelier cells that possess transcriptional profiles remarkably similar to their cortical counterparts, as well as two groups of the recently discovered rosehip interneurons [30], again showing expression profiles nearly identical to the major rosehip sub-populations found in the cortex [unpublished data]. The high clustering resolution afforded to us by ACTIONet also allowed us to partition the small endothelial and mural clusters and identify tiny sub-populations of each corresponding to pericytes, arteriole smooth muscle cells, venous endothelial cells, and arterial endothelial cells, demonstrating clear evidence of the same neurovascular zonation previously reported in mouse [31]. Our analysis produced a cluster of $CD8^+$ T-cells, though it is unclear if these are similar to the memory T-cells recently reported to reside in the brain [32] or if they originate from the blood vessels that were homogenized along with the rest of the tissue. Additional populations recovered include fibroblasts and a large population of $CLSTN2^+/GABRR2^+$ neurons of ambiguous anatomical origin (hereafter referred to as $CLSTN2$ -neurons). Aside from the $CLSTN2$ -neurons, most of which originated from putamen samples, no other cell population showed a regional bias in their distribution. It must be stated that the difference in recovered cell types between our human and mouse data sets is likely a result of differences

in sequencing depth and that from this data there is no evidence to suggest that any of the populations observed only in the human data set is human-specific, as several of these interneuronal and neurovascular populations have been reported in mice by more targeted approaches [31, 33].

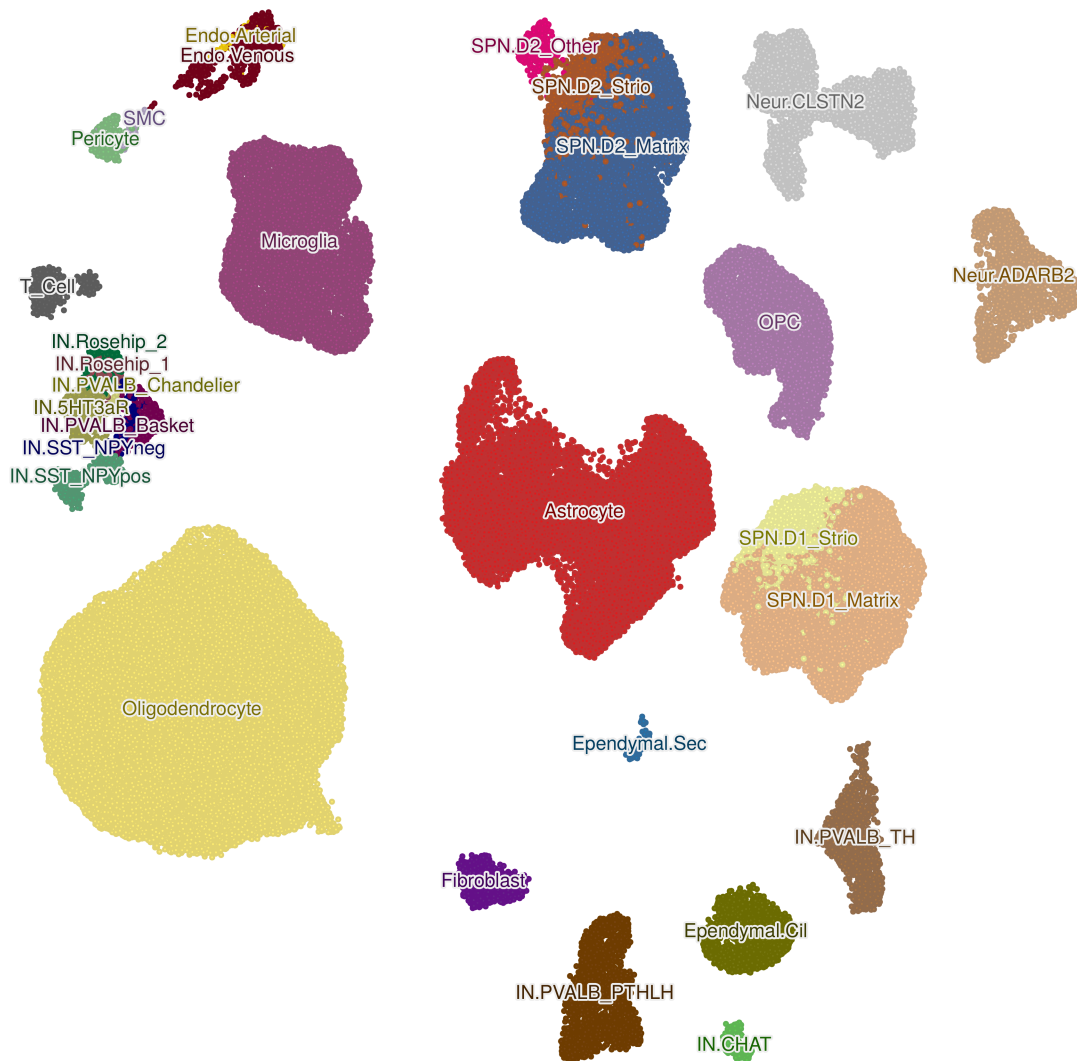


Figure 3-3: ACTIONet plot of cell types identified in the caudate and putamen of human HD and control individuals by snRNA-seq [n = 128,797 cells across 25 tissue samples spanning seven HD and seven unaffected control individuals].

As expected [34], dSPNs and iSPNs contributed approximately equally to the SPN population in control tissue, but the number of iSPNs decreased with human HD grade progression in concordance with reports of enhanced iSPN vulnerability in human HD that were based upon marker gene expression studies [5] [Fig. 3-4]. We make note of this observation in light of substantial evidence in support of the relative depletion of iSPNs vs. dSPNs and the dramatic differences observed across our samples, but stress the fact that the biased nature of droplet-based single-cell methods means that differential abundance is not an accurate or reliable means of assessing relative differences in population sizes.

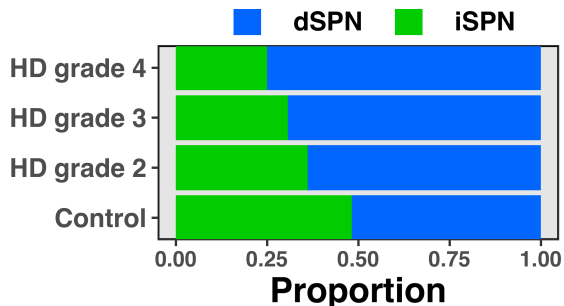


Figure 3-4: Average fraction of SPNs across grades of HD relative to control samples.

3.3 Molecular Conservation of the Neostriatum Within and Across Species

In order to validate our annotations and to better understand how the transcriptional profiles of each cell type compare to one another, we computed pairwise Pearson’s correlation coefficients of the mean expression profiles across all identified cell types within the human control samples and the wild type litter mates from the zQ175 cohort. Hierarchical clustering of the correlation matrices produced major groupings consistent across human and mouse corresponding to interneurons, SPNs and SPN-like cell types, glia, and neurovascular cell types [Fig. 3-5, 3-6]. This analysis revealed the interesting nature of *ADARB2*⁺ SPN-like cells. This cell type generally clusters alongside GABAergic interneurons, produc-

ing a nearly identical spatial arrangement under our algorithm in both human and mouse snRNA-seq studies. This cell type appears to express interneuron markers, namely *ADARB2* and *TACR1*, yet simultaneously exhibits high expression of SPN markers *DRD1*, *PPP1R1B*, *FOXP2*, and *BCL11B*. Hierarchical clustering in both human and mouse places this cell type directly between the SPNs and interneurons, but closer to the former. Remarkably, these neurons show the greatest degree of transcriptome-wide similarity to the novel iSPN group, which are themselves much more distinct from canonical SPNs than the matrix and striosome sub-types are among themselves.

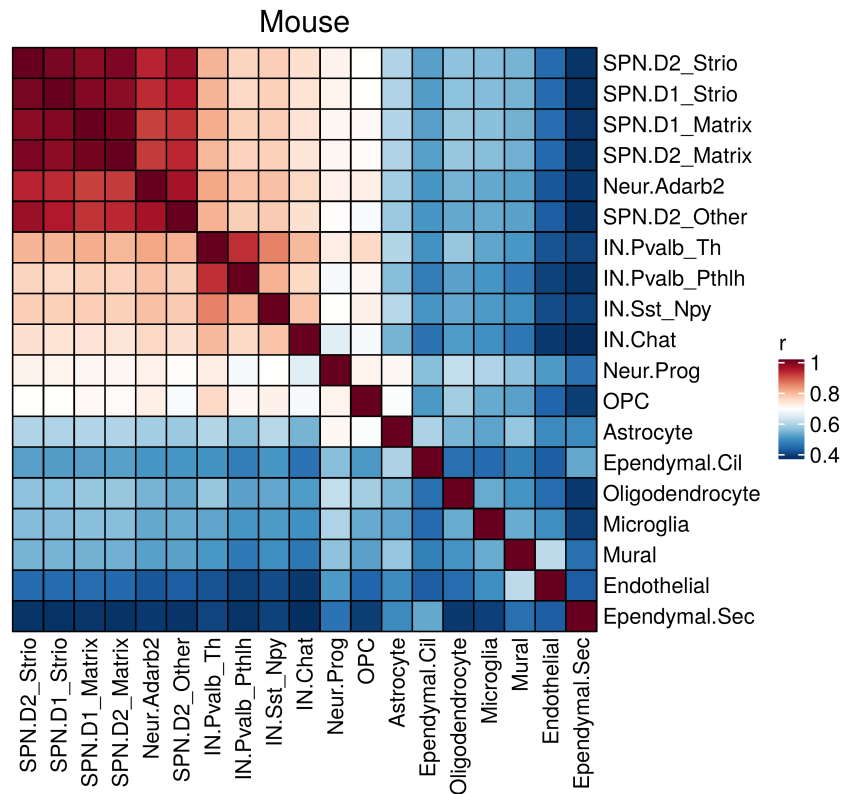


Figure 3-5: Similarity of neostriatal cell types by transcription-wide correlation of gene expression across all cell types identified in 6 month-old control mice.

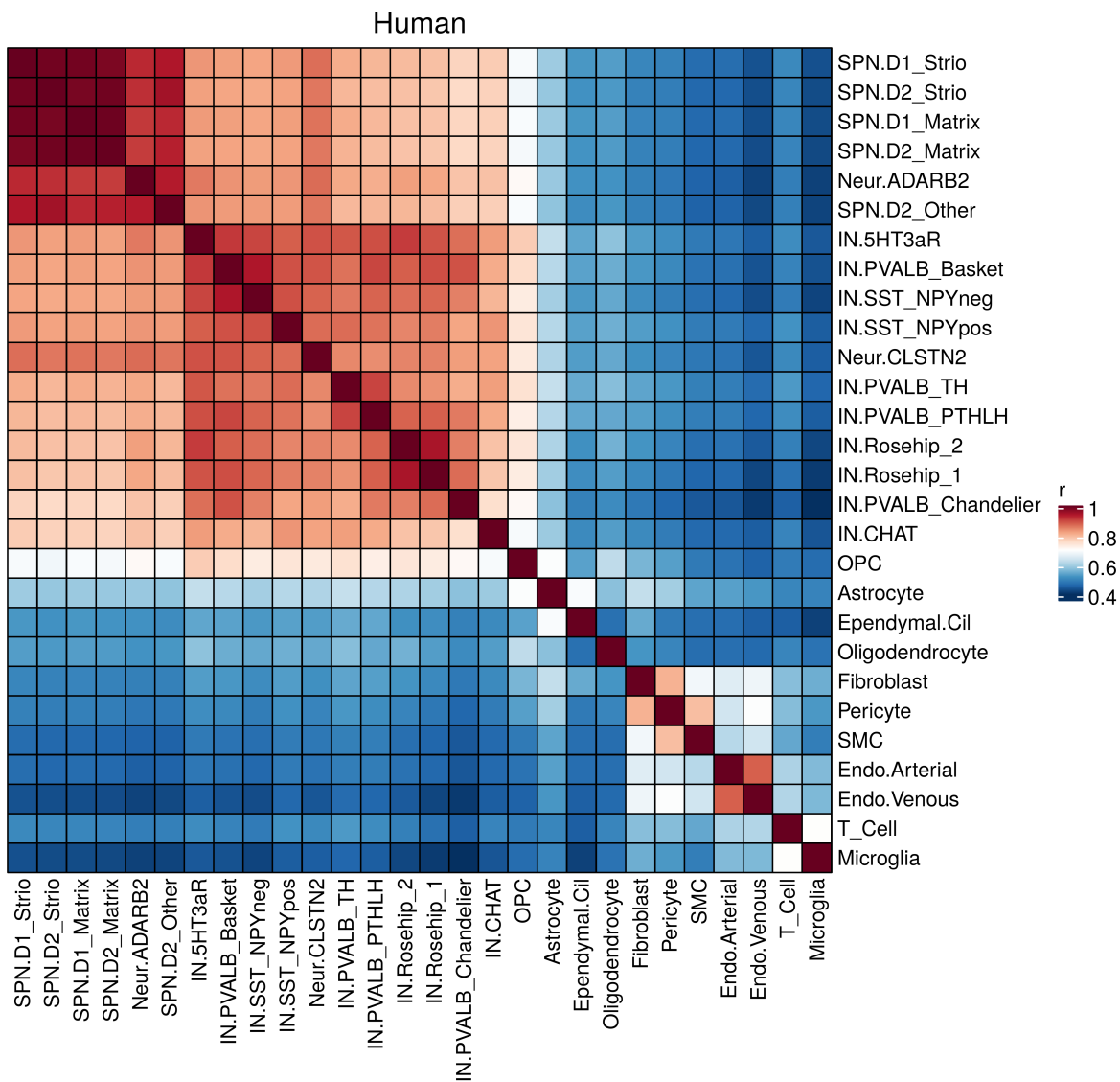


Figure 3-6: Similarity of neostriatal cell types by transcription-wide correlation of gene expression across all cell types identified in human unaffected control samples.

Next, we performed the same analysis, this time comparing cell type similarity across human and mouse [Fig. 3-7]. This analysis revealed a one-to-one mapping of most cell types demonstrating the expected conservation across species. Most populations in mouse mapping to multiple cell types in human were restricted to those that are likely present in mice but could not be readily resolved (e.g. endothelial, mural, and several *Pvalb*⁺ neuronal sub-types). Neural progenitors in mouse, *CD8*⁺ T-cells in human, and the potentially non-striatal *CLSTN2*-neurons, also in human, had no matching counterpart in the other species in this data set. The same hierarchy of cell type groupings seen within human and mouse data sets was recapitulated when comparing across species.

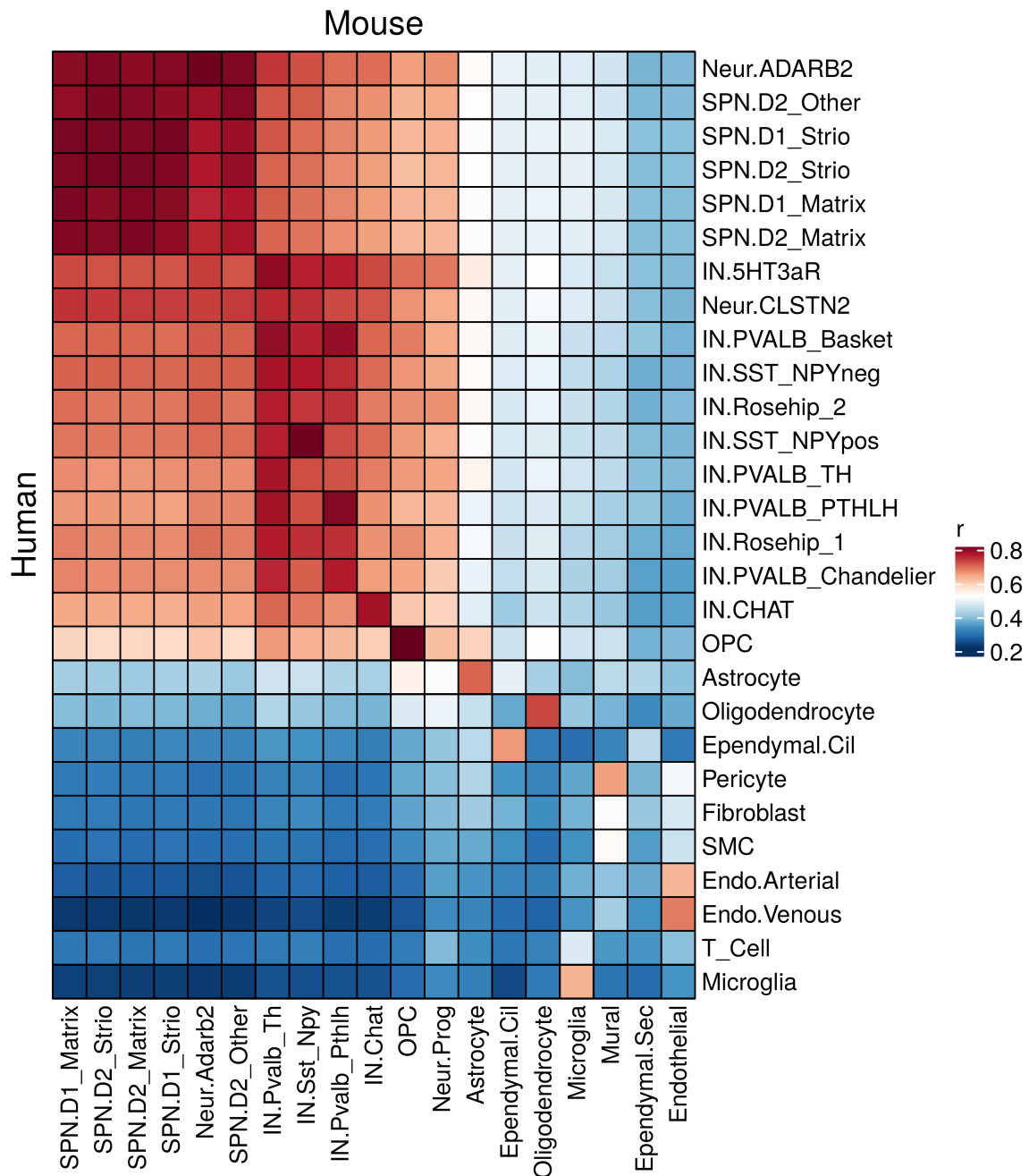


Figure 3-7: Similarity of neostriatal cell types across species by transcription-wide correlation of ortholog gene expression across all cell types identified in human and mouse control samples.

Chapter 4

Analysis of Gene Expression Changes in HD and Mouse Models of HD

4.1 Methods for Differential Gene Expression Analysis

Having identified all cellular populations of interest, we next sought to examine cell type-specific HD-associated changes in gene expression. Gene selection for differential gene expression (DGE) analysis was done by retaining only genes present in at least a certain percentage of cells in every sample within the same cell type. This criteria was acceptable since we did not intend to focus non-autosomal genes in this study. Cell types were excluded from DGE analysis if they were not present in sufficient quantities in at least three controls and three HD/HD model samples or if less than 1000 unique genes were retained in that population. We then compared three approaches for calling differentially-expressed genes (DEGs):

- The commonly used Wilcoxon rank-sum test.
- A multi-resolution pseudo-bulk approach implemented in DESeq2 [35].
- A multi-resolution, inverse variance-weighted pseudo-bulk approach implemented in limma [36].

The canonical implementations of each method have unique advantages and caveats. Pairwise tests such as the Wilcoxon rank-sum test and Student's t -test are common in

snRNA-seq DGE analysis [13, 21, 37] since they scale well with the number of observations and can take advantage of the large sample sizes produced by droplet-based sequencing protocols. The major disadvantages to these methods are that they cannot readily account for covariates, are not robust to outliers, and do not compute differential expression in a transcriptome wide context. Furthermore, each cell is considered a unique observation leading to heavily inflated p -values, and unless sub-sampled, batches with greater cell counts will dominate the results. Pseudo-bulk methods, on the other hand, are functionally equivalent to applying bulk RNA-seq approaches to a sorted cell population and have been shown in several instances to outperform single-cell methods when applied to single cell data [38, 39] since they permit complex experimental design, mask the effects of low quality cells within samples, and take advantage of proven normalization techniques. Unfortunately, current pseudo-bulk strategies still present some major drawbacks:

1. The sample size is reduced from the number of cells to the number of samples, effectively erasing the statistical power conferred by snRNA-seq.
2. If the data set is properly curated, intra-sample single-cell variability is itself an important source of information which is masked at the bulk level.
3. Samples composed of many high quality cells will be given the same weight as those composed of few low quality cells, eliminating yet another compelling feature of snRNA-seq.

To address these concerns we developed a novel pseudo-bulk DGE approach inspired by ensemble learning methods used in machine learning. In our multi-resolution pseudo-bulk method, cell type-specific bulk profiles within each sample are constructed by sorting gene-wise counts in decreasing order and discretizing the counts into n overlapping bins of increasing sizes. We compute the sum, expectation, and variance at each resolution and construct n gene-by-sample pseudo-bulk matrices for input to the appropriate pipelines (e.g. limma, DESeq2, edgeR). We call differentially-expressed genes at each resolution, taking the mean log-fold change across all resolutions and combine p -values which can then be adjusted accordingly. We supply the matrix of inverse, gene-wise, intra-sample single-cell

variances at each resolution as weights for model fitting by the limma algorithm. This enables an additional mode of outlier detection and correction by adjusting the weights and removing the contribution of genes with extreme median absolute deviations from the sample median variance, and those with non-positive variance. The result is a more robust and reproducible set of differentially-expressed genes with the statistical power expected of a single-cell experiment. The p -values do not decrease arbitrarily with the number of resolutions, but instead converge to a finite limit.

Samples with negligible quantities of low quality cells, with few or no covariates, and little inter-sample variance and differential cell abundance would see similar results using any of the above methods. Such cohorts are virtually non-existent in human studies, but often the norm for *in vitro* studies and those using isogenic mouse lines, and in fact we see exactly this robustness when comparing DGE results across methods for the R6/2 and zQ175 cohorts. All three methods produce consistent gene sets (mean cell type-wise $r > 0.86$) [Fig. 4-1 middle, bottom] in the mouse model studies, but vary more so for human (mean cell type-wise $r \in 0.69 - 0.8$) [Fig. 4-1 top]. These differences are present at single resolution in the pseudo-bulk trials, albeit with greatly reduced power and more variable log-fold changes. Such differences are expected given that each method takes in a different input and employs a different statistical model to identify DEGs in this highly noisy data set. Nevertheless, we observed decent overlap across the various methods within cells of the same type, and when comparing across the most statistically significant and most drastically altered (in terms of log-fold change) genes in human we saw many of the same genes come up across methods. We found our inverse-weighted, multi-resolution, linear-model approach to produce the most robust and reproducible results across this and other human data sets as well as in mouse models, and Wilcoxon predictably produced the least reproducible results. We used the set of DEGs generated by multi-resolution limma for downstream analysis.

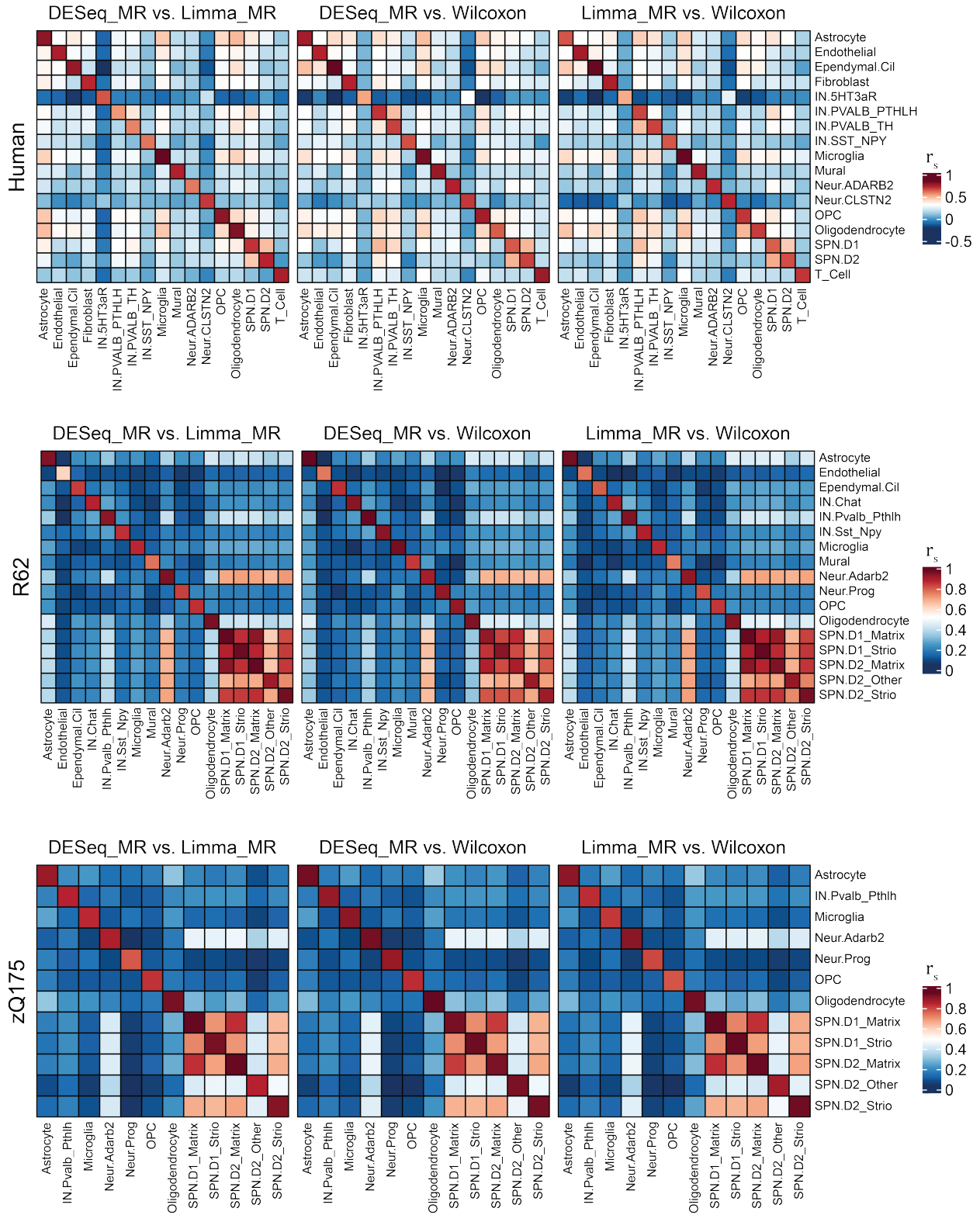


Figure 4-1: Comparison of DEGs across methods in each experimental cohort. Spearman correlation by log-fold change.

4.2 Differential Gene Expression Analysis of HD Mouse Models

We compared gene expression in HD model vs. control striatum for all sufficiently abundant cell types, which revealed differentially-expressed genes in dSPNs, iSPNs, *Adarb2*-neurons, astroglia, microglia, OPCs, and oligodendrocytes, as well as other striatal cell types in response to either form of mHTT at these HD model time points. Several SPN marker genes such as *Pde10a*, *Penk*, and *Bcl11b* were among the most significantly down-regulated genes in many of the striatal cell types suggesting that mHTT affects at least some common transcriptional regulator(s) that are present across all of these populations, including glial cells [Fig. 4-2, 4-3].

Comparison of all DEGs across cell types and across models revealed that iSPNs and dSPNs showed the most similar transcriptional responses for both down-regulated and up-regulated genes, with a $r > 0.9$ between all striosomal and matrix SPN sub-types [Fig. 4-4, 4-5]. The novel iSPN population and *Adarb2*-neurons showed the next closest similarity with $r > 0.8$. Interestingly, in some cases “other” iSPNs showed a more similar response to *Adarb2*-neurons than to the canonical SPN sub-types. Astrocytes, oligodendrocytes, and *Pvalb*⁺ interneurons showed the next closest, but lower, similarity. These striking results indicate that dSPNs and iSPNs display very similar responses to both an exon 1 fragment of mHTT and a full length knock-in *mHtt* with the same repeat size, and that there likely exists a component of this response that is shared across other non-SPN neuronal cell types. Comparisons of relative dysregulation across cell types reconfirmed that SPNs show the greatest degree of dysregulation with *Adarb2*-neurons showing the next highest, followed by *Pvalb*⁺ interneurons. From our mouse data, we did not observe a significant difference in magnitude of dysregulation between striosome and matrix SPNs.

The relative number of dysregulated genes in SPNs of the R6/2 cohort was twice that of the zQ175, which is coherent with this model representing a later phenotypic time point. Comparison across mouse models showed strong conservation of the mHTT-induced transcriptional response in SPNs, *Adarb2*⁺ neurons and oligodendrocytes, with canonical SPNs showing the greatest degree of overlap [Fig. 4-6]. Surprisingly DEGs in oligodendrocytes

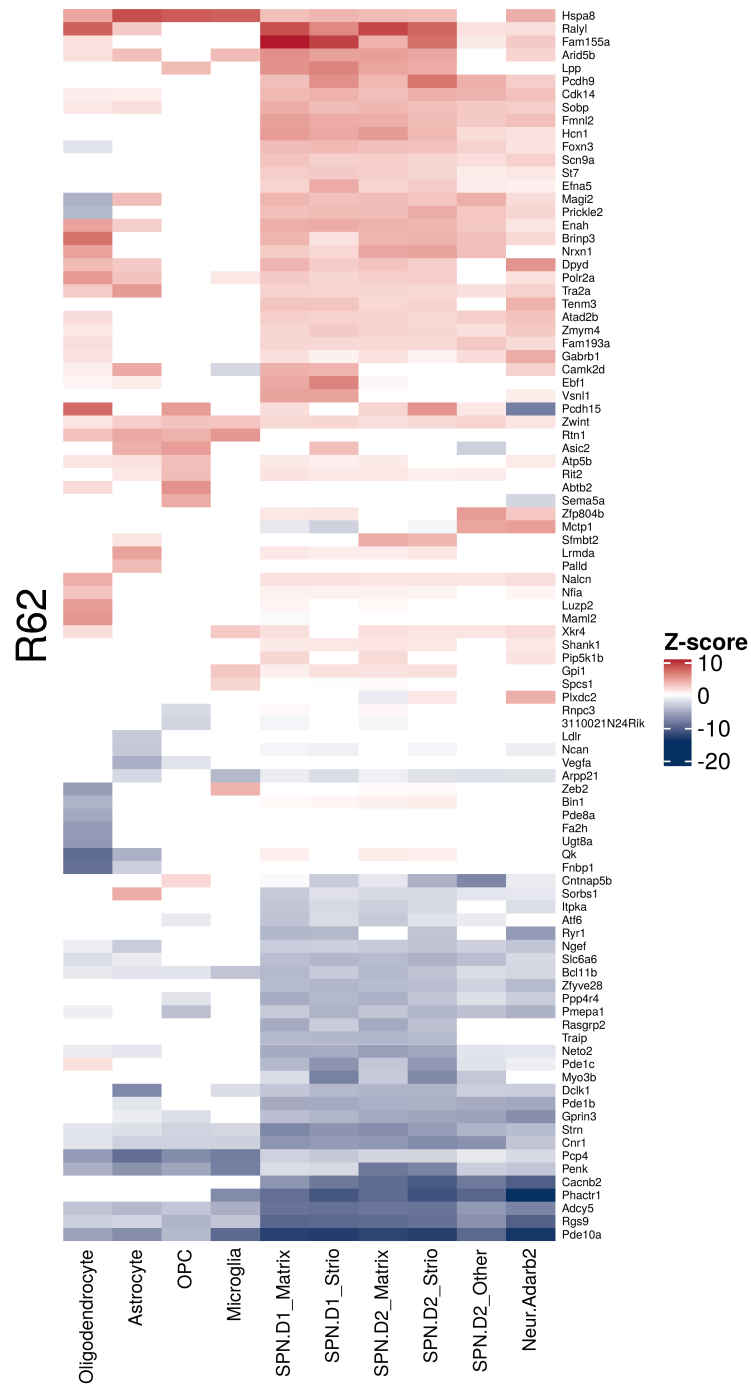


Figure 4-2: Top 7 most statistically significant up- and down-regulated genes in SPNs and other highly abundant cell types of R6/2.

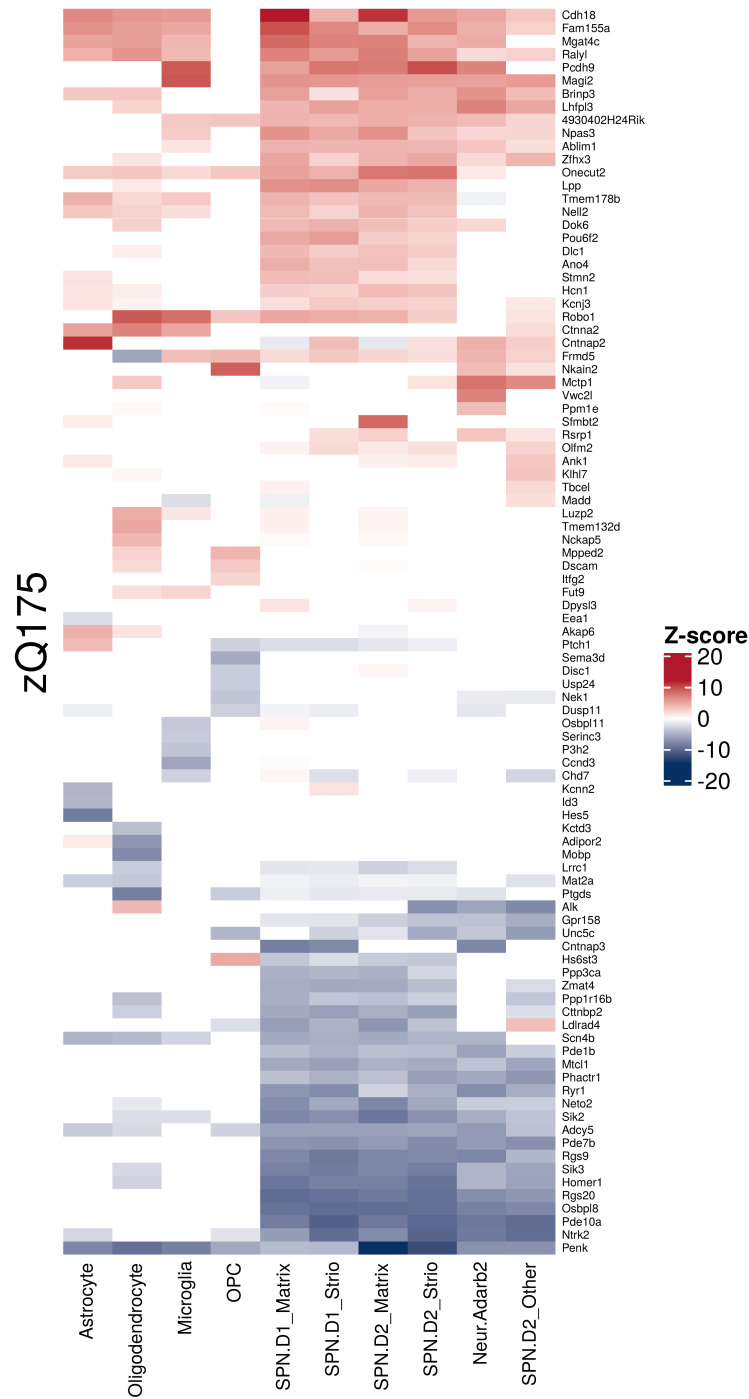


Figure 4-3: Top 7 most statistically significant up- and down-regulated genes in SPNs and other highly abundant cell types of zQ175.

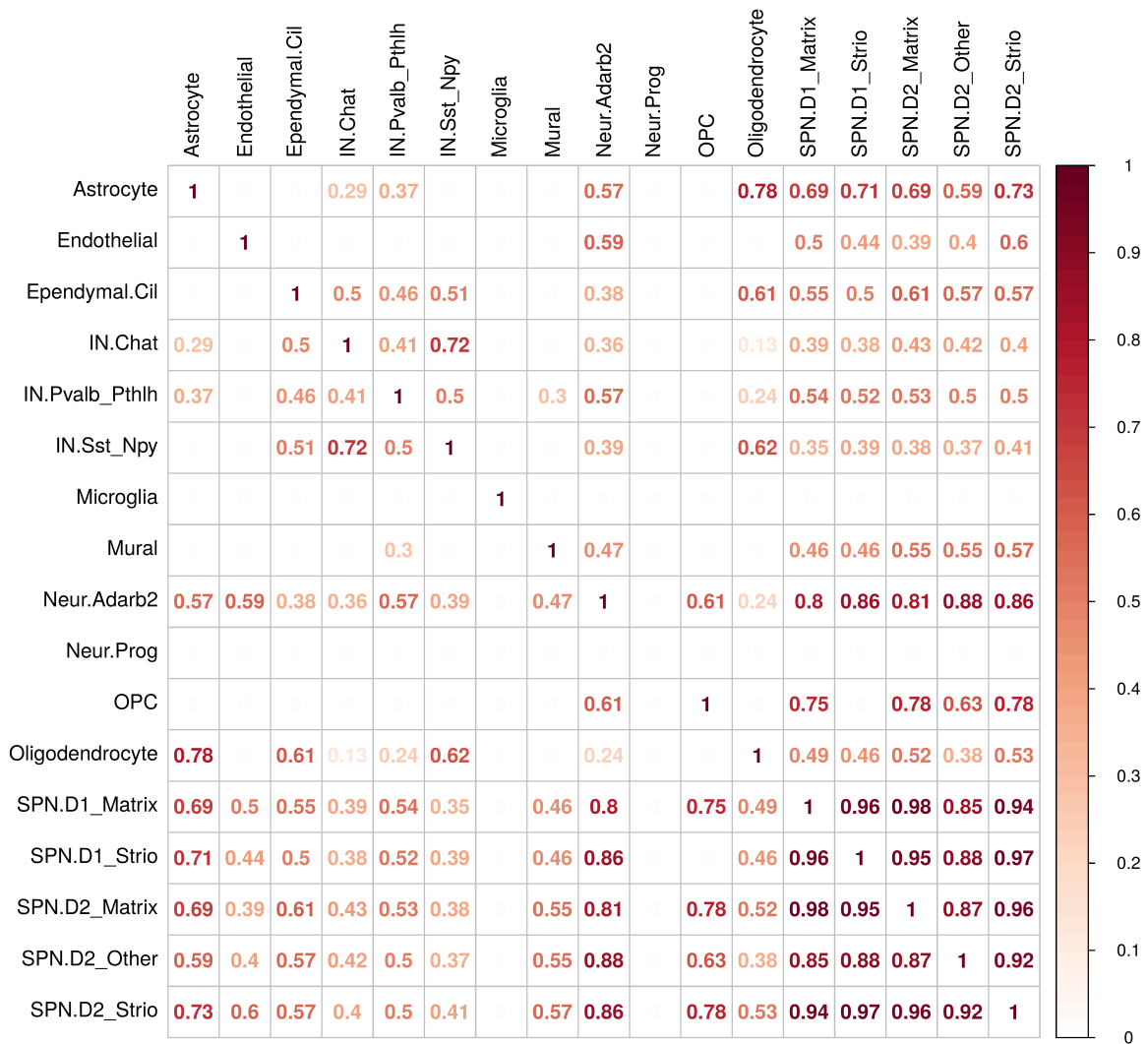


Figure 4-4: Pearson's correlation of DEGs by log-fold change across cell types in R6/2 [FDR-adjusted p -value < 0.001 and \log_2 fold-change > 0.1]

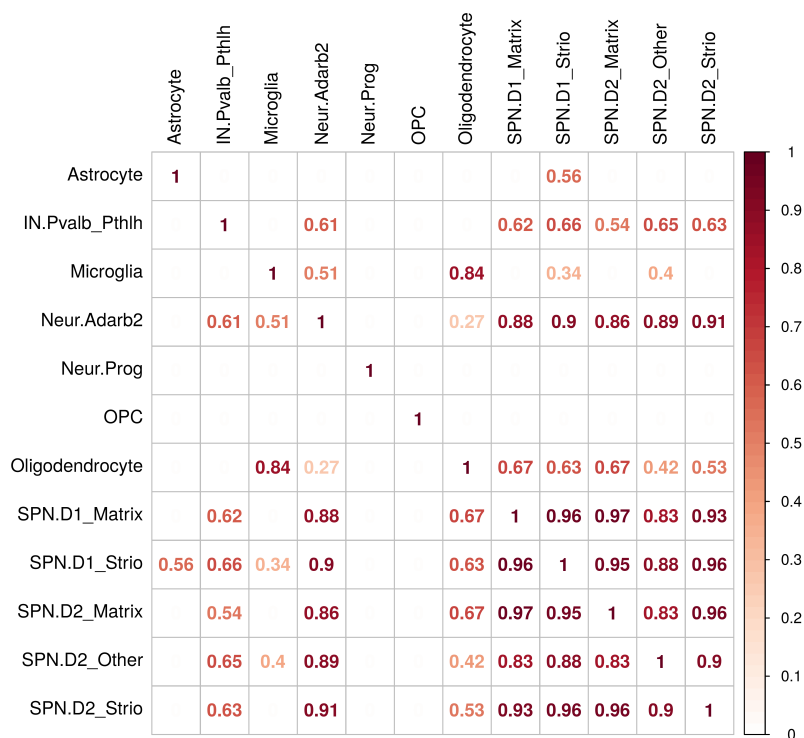


Figure 4-5: Pearson's correlation of DEGs by log-fold change across cell types in zQ175 [FDR-adjusted p -value < 0.001 and \log_2 fold-change > 0.1]

showed equal or stronger correlation across mouse models than “other” iSPNs and *Adarb2*-neurons. Other cell types showed little or no overlap in DEGs across models of HD. Despite *Pvalb*⁺/*Pthlh*⁺ interneurons showing a high degree of dysregulation in both models, the sets of perturbed genes seemed to be mostly unique to each model.

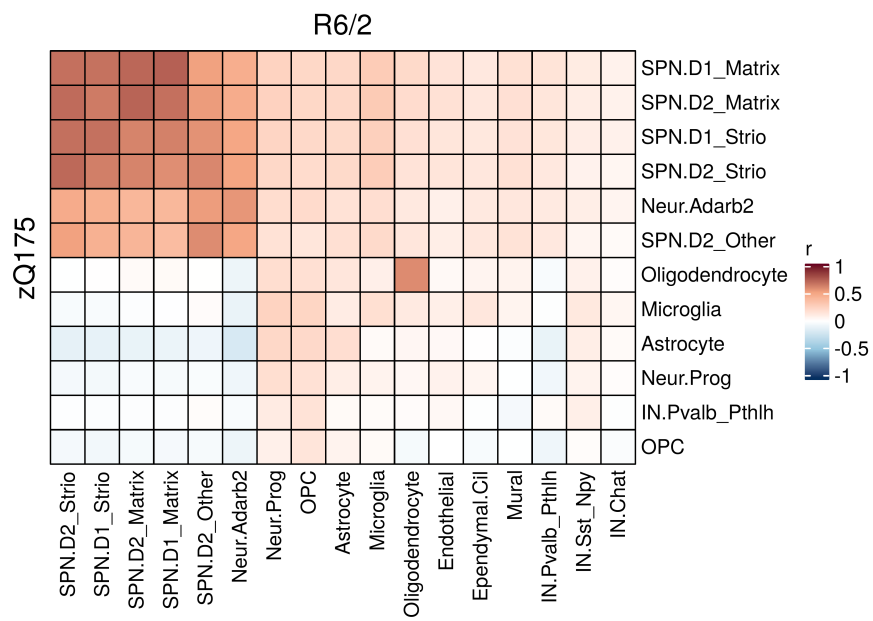


Figure 4-6: Pearson's correlation of DEGs by Z-score between R6/2 and zQ175.

4.3 Differential Gene Expression Analysis of Human HD

We then compared gene expression in the HD vs. control caudate and putamen for all sufficiently abundant cell types, which revealed DEGs for neuronal, glial, and neurovascular cell types as well as ciliated ependymal cells and T-cells. Immediately, we could see that there was a much higher proportion of dysregulated genes across all cell types compared to the mouse models. While this is partially expected given that our samples are representative of a late-stage, and in some cases end-stage, phenotypic time point, the post-mortem nature of these samples means that a sizable component of the signal likely stems from autolysis- and RNA-degradation-induced gene expression changes. Therefore, extensive care needs to be taken when drawing conclusions from these results. Prior to the application of our aggressive filtering strategy, we were able to reproduce an expression signature reported in other similar post-mortem human studies and observed in unpublished data from our other disease cohorts [13, 21, 22]. These genes include, but are not limited to, *RASGEF1B*, *INO80D*, *DHFR*, *LINGO1*, *SLC26A3*, and *HSP90AA1* as well as numerous ferroptosis pathway-associated genes. We have found this gene set to be a remarkably robust signature of low cell quality, but its prevalence among numerous studies may erroneously lead some to believe that it is a signature of a common underlying disease mechanism, when in reality it is a reproducible artifact. Our cell and gene selection strategies were able to detect and remove a considerable portion of this signal, resulting in a more reliable set of DEGs.

From the remaining genes, we again found a number of SPN marker genes to be down-regulated across numerous cell types including *PDE10A*, and the striosomal marker *TAC1* as well as *OTOF* and *ASIC2*, which are markers of *ADARB2*-neurons [Fig. 4-7]. A few dysregulated SPN markers, namely *PENK* and *ADORA2A*, showed statistically significant upregulation in contradiction of prior studies [11, 12]. Alternative DGE analysis using a multi-resolution implementation of DESeq2 and a simple Wilcoxon rank-sum test yielded inconsistent results for these two marker genes as well as for *GPR6*. In several instances these genes were flagged as outliers by outlier detection methods, leading us to believe that although statistically significant, the reported expression changes for these genes are unre-

liable for our cohort. Also among the most down-regulated genes in SPNs were contactin family members *CNTN3* and *CNTN5* as well as *PHACTR1*, of which the latter two were also among the most down-regulated in SPNs of mouse models. Junctophilin 4 (*JPH4*), which is co-expressed with and may function in a coordinated manner with *JPH3* (CT-G/CAG repeat expansions in which cause Huntington’s disease-like-2 disease) in neurons [40] was also found to be downregulated in human and R6/2, but not in zQ175. Unlike the mouse models which showed relatively symmetric dysregulation, the majority of DEGs in human samples trended downwards. Among the few up-regulated DEGs, the two showing the highest expression change in SPNs were humanin-like 12 (*MTRNR2L12*) and *DHFR*. Although we have identified *DHFR* as a potential artifact across multiple studies, it is also a well-established HD age-of-onset modifier [41] and remains among the top genes after applying our gene filtering measures, thus it is possible that upregulation of *DHFR* in SPNs is a true HD-associated response. We were not able to perform separate DEG analyses for the striosome and matrix SPN sub-populations due to a lack of sufficient striosome SPNs in our HD samples. This could be a consequence of the hypothesized enhanced vulnerability of these sub-populations leading to further depletion, but from this data we can make no claims about SPN sub-type-specific differences in pathological burden in humans.

A comparison of HD caudate and putamen dysregulated genes to genes linked to SNPs associated with earlier HD age-of-onset [41] revealed that several of these SNP-linked genes, including *FAN1*, *MSH3*, and *MLH3* were downregulated in one or both SPN types, but future studies will be needed to correlate these findings by human HD grade, and thus to human HD disease progression. Finally, an analysis of reactive astroglia markers [Fig. 4-8][42] revealed a molecular signature of pan-reactive, but not A1- or A2-reactive astroglia in HD. This signature was absent in mouse models, consistent with the lack of evidence of reactive astrogliosis in these models [43].

Comparison of DEGs across cell types again confirmed that SPNs exhibit a similar transcriptional response in HD ($r = 0.89$) [Fig. 4-9]. Surprisingly, the cell types that showed the most similar response to SPNs in human HD were not *ADARB2*-neurons, but *PVALB/TH* and *PVALB/PTHLH* interneurons which shared many DEGs with dSPNs ($r > 0.83$), but less so with iSPNs ($r > 0.67$). These results, which we also observed in both mouse models,

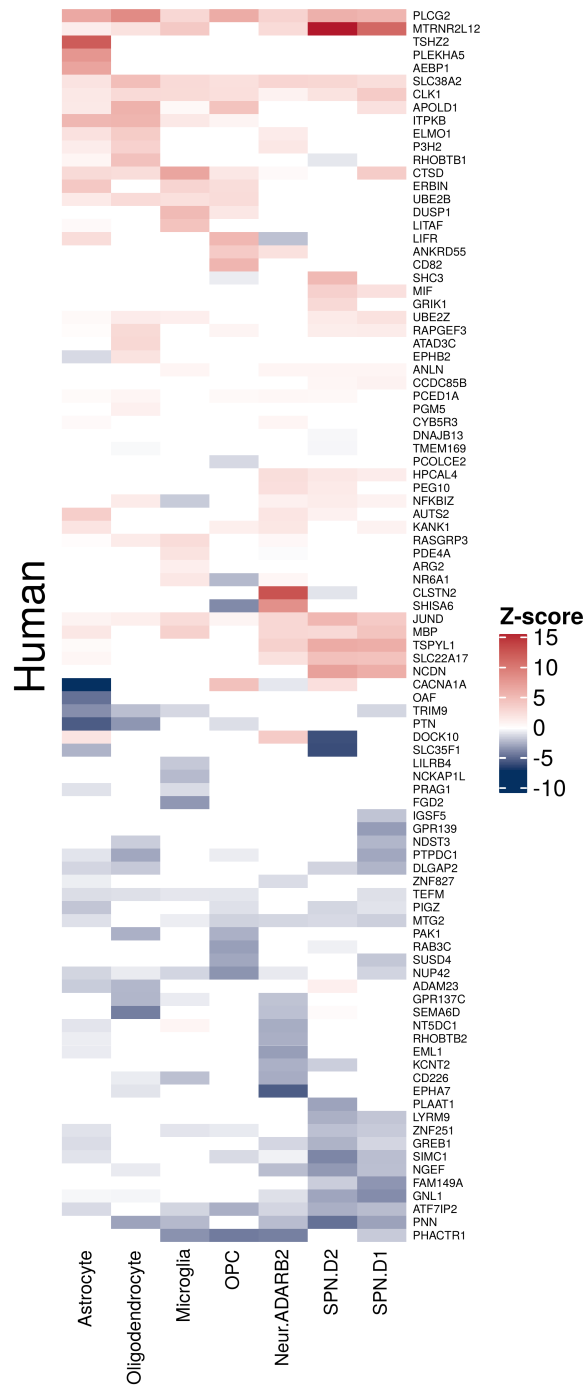


Figure 4-7: Top 7 most statistically significant up- and down-regulated genes in SPNs and other highly abundant cell types in Human HD.

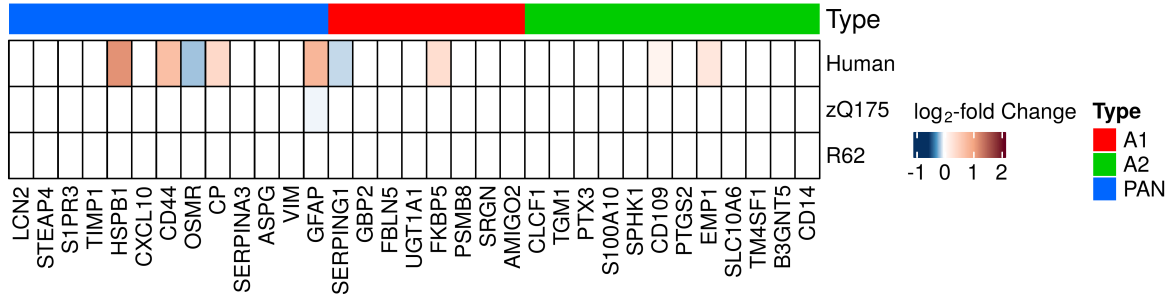


Figure 4-8: Differential expression of reactive astrocyte marker genes from Liddelov et al (2017) in astrocytes of Human HD and HD mouse models.

suggest that *PVALB*⁺ interneurons may be more affected in HD than previously believed. Several of the top dysregulated SPN genes were also observed to be similarly dysregulated in *ADARB2*-neurons, although transcriptome-wide changes between these cells and SPNs showed far less overlap compared to the mouse models.

Comparison on DEGs between human and mouse models showed no overlap when compared against either R6/2 or zQ175. Alternative DGE analysis using multi-resolution DESeq2 and Wilcoxon rank-sum test produced conflicting results [Fig. 4-10]. Cross-species comparison by pairwise Wilcoxon showed some overlap between human and mouse SPNs and *Adarb2*-neurons and to a lesser extent astrocytes and mural cells (the latter only vs. R6/2). The greatest similarity observed was between mouse iSPNs and human *ADARB2*-neurons ($r < 0.33$) when compared against R6/2. Comparison by DESeq2 showed weak correlation ($r < 0.25$) only between mouse SPNs and human *ADARB2*-neurons and some overlap between endothelial cells. These results are not entirely unexpected due to human sample variability and the fact that the human cohort represents a late- and end-stage time point compared to an early-mid time point in the mouse models. There was also a large difference in sequencing depth between human and mouse samples, potentially introducing a disruptive batch effect. Therefore, we can make no concrete claims about the face validity of these models. Further studies using well-curated grade 0 and grade 1 samples will be necessary to accurately assess how well these models can recapitulate cell type-specific gene expression changes.

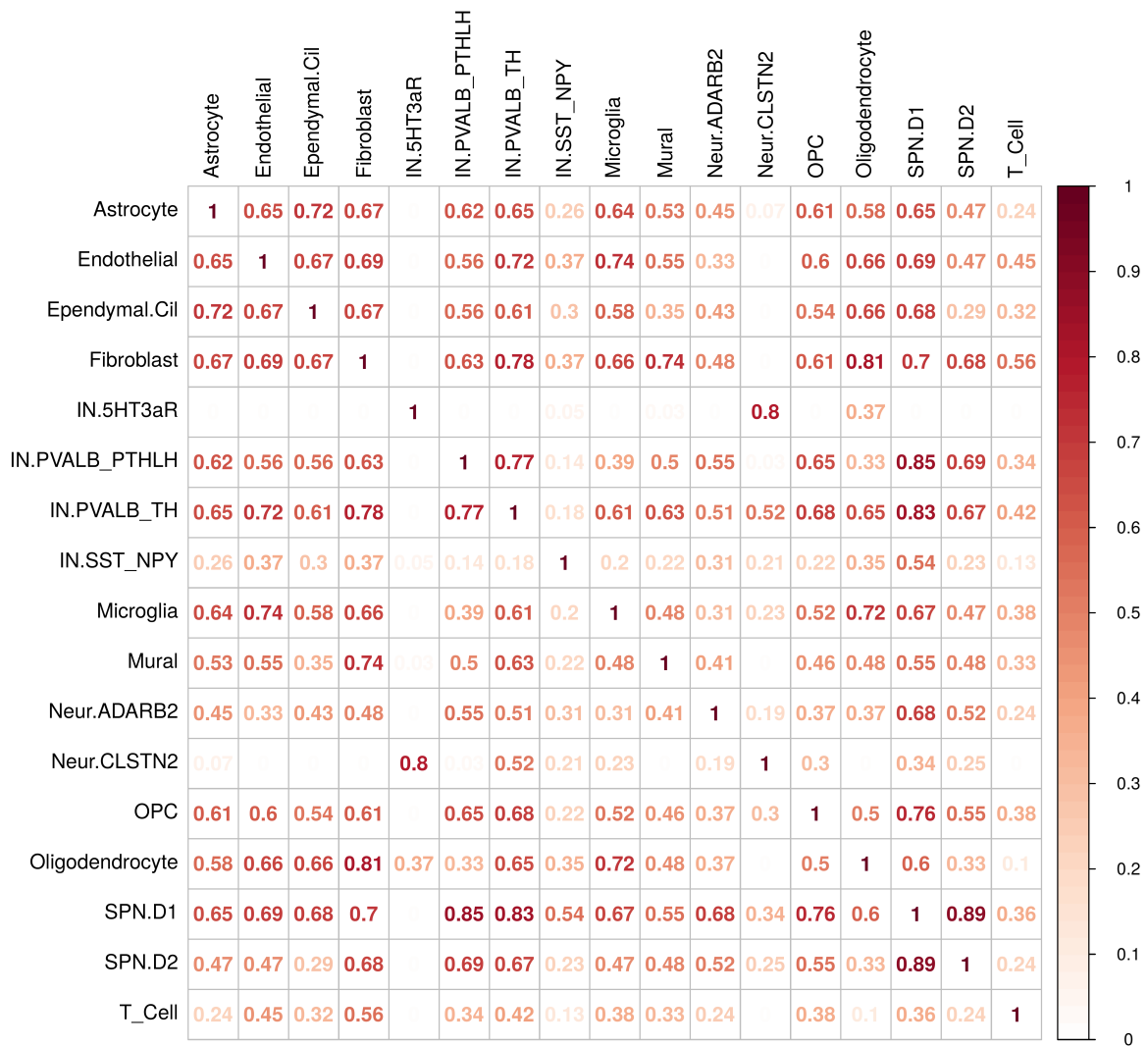


Figure 4-9: Pearson's correlation of DEGs by log-fold change across cell types in Human HD [FDR-adjusted p -value < 0.001 and \log_2 fold-change > 0.1]

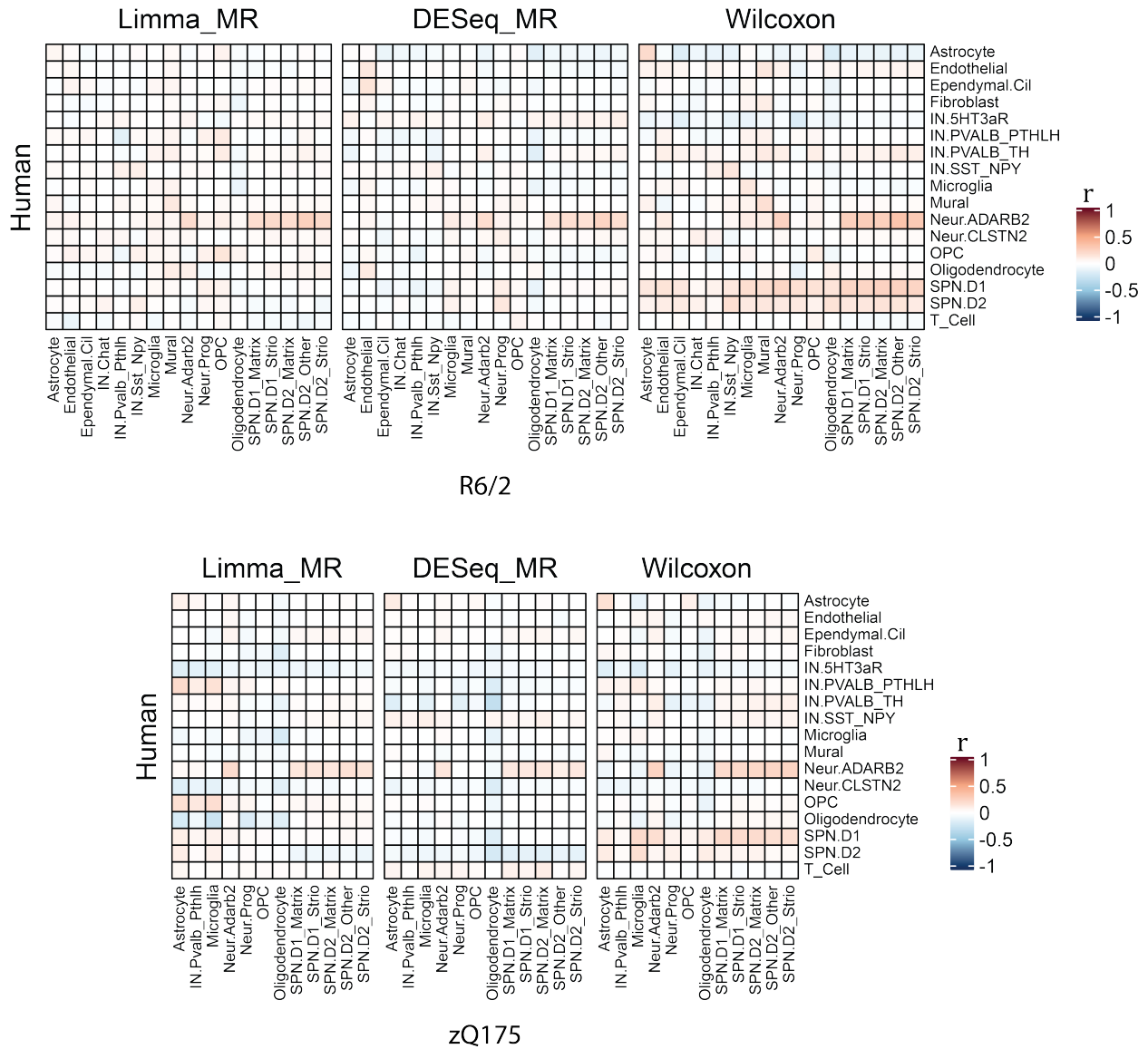


Figure 4-10: Comparison of DEGs between Human HD and HD mouse models by method [Pearson's correlation of Z-scores].

Chapter 5

Discussion

We report the first single-cell resolution molecular atlas of the human neostriatum, and the highest resolution profile on the mouse neostriatum to date. Our analysis required the development and application of novel strategies for cell clustering, data curation, and differential expression analysis which we document here. In particular, we describe a new ensemble learning-inspired pseudo-bulk method for differential gene expression analysis of single-cell data that produces robust and reproducible DEGs with high statistical power and requires only a few dozen high quality cells. Using these new strategies, we were able to construct and analyze gene expression profiles of matrix and striosomal SPNs as well as a novel iSPN sub-type found in both human and mouse. We were also able to identify and transcriptionally characterize multiple previously unobserved interneuron populations and confirm the existence of neurovascular zonation in the human striatum. We found strong evidence of inter-species conservation of cell types and transcriptional profiles.

We also report the first single-cell dissection of Huntington’s disease in human and two commonly used mouse models. Transcriptional changes and dysregulation have long been associated with HD [9–11], but cell type-specific studies of gene expression across striatal cell types in HD models and human HD have been difficult to perform due to the technical challenges posed by the diversity of intermingled striatal cell types. By studying human HD and HD mouse models with snRNA-seq we revealed both non-cell type-specific and cell type-specific responses that are induced by mHTT. We observed significant overlap in expression changes for SPNs within both human and mouse and an impressive degree

of similarity among these and *ADARB2*-neurons, although the latter also demonstrates a severe and unique response of its own. Our studies suggest that our newly discovered iSPNs, the recently identified *ADARB2*-neurons, and perhaps even certain populations of parvalbumin-expressing GABAergic interneurons could be vulnerable in HD. Comparison across mouse models revealed a concordant response to mHTT in SPNs, *Adarb2*-neurons, and oligodendrocytes, but not for any other recovered cell population.

Analysis of the astroglial population revealed significant upregulation of pan-reactive markers, but not A1 or A2, in human HD astroglia but not in mouse models. Moreover, our human HD studies suggest an enhanced loss of striatal striosome markers and a grade-dependent loss of iSPNs vs. dSPNs among the surviving SPNs, and we note that several genes near age-of-onset modifier SNPs identified by the GeM-HD Consortium were altered in expression in dSPNs and iSPNs of human HD tissue. However, much larger-scale snRNA-seq human HD studies will be needed to correlate these findings by human HD grade, and thus to human HD disease progression, as well as to more fully characterize gene expression changes present in the rarer cell types studied here.

In conclusion, our gene expression studies across mouse models of HD and human HD not only create a first catalog of the cell type-specific molecular changes induced by various forms of mHTT, but also yielded high-resolution molecular atlases that will serve as an invaluable resource to the striatal biology community by providing a reference profile for the diverse cellular populations found in the striata of humans and mice.

Bibliography

1. Zuccato, C., Valenza, M. & Cattaneo, E. Molecular mechanisms and potential therapeutic targets in Huntington's disease. *Physiological Reviews* (2010).
2. Martin, J. B. & Gusella, J. F. Huntington's disease: Pathogenesis and management. *New England Journal of Medicine* (1986).
3. VONSATTEL, J.-P. *et al.* Neuropathological Classification of Huntington's Disease. *Journal of Neuropathology and Experimental Neurology* (1985).
4. Reiner, A. *et al.* Differential loss of striatal projection neurons in Huntington disease. *Proceedings of the National Academy of Sciences of the United States of America* (1988).
5. Albin, R. L. *et al.* Preferential loss of striato-external pallidal projection neurons in presymptomatic Huntington's disease. *Annals of Neurology* (1992).
6. Richfield, E. K., Maguire-Zeiss, K. A., Vonkeman, H. E. & Voorn, P. Preferential loss of preproenkephalin versus preprotachykinin neurons from the striatum of Huntington's disease patients. *Annals of Neurology* (1995).
7. Graybiel, A. M. & Ragsdale, C. W. Histochemically distinct compartments in the striatum of human, monkey, and cat demonstrated by acetylthiocholinesterase staining. *Proceedings of the National Academy of Sciences of the United States of America* (1978).
8. Hedreen, J. C. & Folstein, S. E. Early loss of neostriatal striosome neurons in huntingtons disease. *Journal of Neuropathology and Experimental Neurology* (1995).
9. Benn, C. L. *et al.* Contribution of nuclear and extranuclear polyQ to neurological phenotypes in mouse models of Huntington's disease. *Human Molecular Genetics* (2005).

10. Saudou, F., Finkbeiner, S., Devys, D. & Greenberg, M. E. Huntingtin acts in the nucleus to induce apoptosis but death does not correlate with the formation of intranuclear inclusions. *Cell* (1998).
11. Hodges, A. *et al.* Regional and cellular gene expression changes in human Huntington's disease brain. *Human Molecular Genetics* (2006).
12. Langfelder, P. *et al.* Integrated genomics and proteomics define huntingtin CAG length-dependent networks in mice. *Nature Neuroscience* (2016).
13. Lee, H. *et al.* Cell Type-Specific Transcriptomics Reveals that Mutant Huntingtin Leads to Mitochondrial RNA Release and Neuronal Innate Immune Activation. *Neuron* (2020).
14. Klein, A. M. *et al.* Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell* (2015).
15. Macosko, E. Z. *et al.* Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. *Cell* (2015).
16. Mangiarini, L. *et al.* Exon I of the HD gene with an expanded CAG repeat is sufficient to cause a progressive neurological phenotype in transgenic mice. *Cell* (1996).
17. Heikkinen, T. *et al.* Characterization of Neurophysiological and Behavioral Changes, MRI Brain Volumetry and 1H MRS in zQ175 Knock-In Mouse Model of Huntington's Disease. *PLoS ONE* (2012).
18. Menalled, L. B., Sison, J. D., Dragatsis, I., Zeitlin, S. & Chesselet, M. F. Time course of early motor and neuropathological anomalies in a knock-in mouse model of Huntington's disease with 140 CAG repeats. *Journal of Comparative Neurology* (2003).
19. Menalled, L. B. *et al.* Comprehensive Behavioral and Molecular Characterization of a New Knock-In Mouse Model of Huntington's Disease: ZQ175. *PLoS ONE* (2012).
20. Franich, N. R. *et al.* Phenotype onset in Huntington's disease knock-in mice is correlated with the incomplete splicing of the mutant huntingtin gene. *Journal of Neuroscience Research* (2019).

21. Mathys, H. *et al.* Single-cell transcriptomic analysis of Alzheimer’s disease. *Nature* (2019).
22. Schirmer, L. *et al.* Neuronal vulnerability and multilineage diversity in multiple sclerosis. *Nature* (2019).
23. Velmeshev, D. *et al.* Single-cell genomics identifies cell type-specific molecular changes in autism. *Science* (2019).
24. Mohammadi, S., Ravindra, V., Gleich, D. F. & Grama, A. A geometric approach to characterize the functional identity of single cells. *Nature Communications* (2018).
25. Haghverdi, L., Lun, A. T. L., Morgan, M. D. & Marioni, J. C. Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nature Biotechnology* (2018).
26. Hie, B., Bryson, B. & Berger, B. Efficient integration of heterogeneous single-cell transcriptomes using Scanorama. *Nature Biotechnology* (2019).
27. Mohammadi, S., Davila-Velderrain, J. & Kellis, M. A multiresolution framework to characterize single-cell state landscapes. *Nature Communications* (2020).
28. Märtn, A. *et al.* A Spatiomolecular Map of the Striatum. *Cell Reports* (2019).
29. Saunders, A. *et al.* Molecular Diversity and Specializations among the Cells of the Adult Mouse Brain. *Cell* (2018).
30. Boldog, E. *et al.* Transcriptomic and morphophysiological evidence for a specialized human cortical GABAergic cell type. *Nature Neuroscience* (2018).
31. Vanlandewijck, M. *et al.* A molecular atlas of cell types and zonation in the brain vasculature. *Nature* (2018).
32. Smolders, J. *et al.* Tissue-resident memory T cells populate the human brain. *Nature Communications* (2018).
33. Muñoz-Manchado, A. B. *et al.* Diversity of Interneurons in the Dorsal Striatum Revealed by Single-Cell RNA Sequencing and PatchSeq. *Cell Reports* (2018).

34. Gerfen, C. The Neostriatal Mosaic: Multiple Levels Of Compartmental Organization In The Basal Ganglia. *Annual Review of Neuroscience* (1992).
35. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biology* (2014).
36. Ritchie, M. E. *et al.* Limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Research* (2015).
37. Wertz, M. H. *et al.* Interleukin-6 deficiency exacerbates Huntington’s disease model phenotypes. *Molecular Neurodegeneration* (2020).
38. Soneson, C. & Robinson, M. D. Bias, robustness and scalability in single-cell differential expression analysis. *Nature Methods* (2018).
39. Crowell, H. L. *et al.* On the discovery of subpopulation-specific state transitions from multi-sample multi-condition single-cell RNA sequencing data. *bioRxiv* (2020).
40. Nishi, M., Sakagami, H., Komazaki, S., Kondo, H. & Takeshima, H. Coexpression of junctophilin type 3 and type 4 in brain. *Molecular Brain Research* (2003).
41. Lee, J.-M. M. *et al.* CAG Repeat Not Polyglutamine Length Determines Timing of Huntington’s Disease Onset. *Cell* (2019).
42. Liddelow, S. A. *et al.* Neurotoxic reactive astrocytes are induced by activated microglia. *Nature* (2017).
43. Diaz-Castro, B., Gangwani, M. R., Yu, X., Coppola, G. & Khakh, B. S. Astrocyte molecular signatures in Huntington’s disease. *Science Translational Medicine* (2019).