

<sup>1</sup> Multigranular single-cell analysis reveals inflammatory interactions driving  
<sup>2</sup> macular degeneration

<sup>3</sup> Manik Kuchroo<sup>1\*</sup>, Marcello DiStasio<sup>2\*</sup>, Eric Song<sup>3\*</sup>, Eda Calapkulu<sup>3</sup>, Le Zhang<sup>4</sup>, Maryam Ige<sup>5</sup>, Amar H.  
<sup>4</sup> Sheth<sup>5</sup>, Madhvi Menon<sup>6</sup>, Alexander Tong<sup>7</sup>, Abhinav Godavarthi<sup>8</sup>, Yu Xing<sup>3</sup>, Scott Gigante<sup>9</sup>, Holly  
<sup>5</sup> Steach<sup>10</sup>, Jessie Huang<sup>7</sup>, Guillaume Huguet<sup>11,16</sup>, Janhavi Narain<sup>12</sup>, Kisung You<sup>17</sup>, George Mourgos<sup>2,3</sup>,  
<sup>6</sup> Rahul M. Dhodapkar<sup>5</sup>, Matthew J. Hirn<sup>13,14</sup>, Bastian Rieck<sup>15</sup>, Guy Wolf<sup>11,16</sup>, Smita Krishnaswamy<sup>7,17§</sup>,  
<sup>7</sup> Brian P. Hafler<sup>2,3,18§</sup>

<sup>1</sup> Department of Neuroscience, Yale University, New Haven, CT, <sup>2</sup> Department of Pathology, Yale University, New Haven, CT,  
<sup>3</sup> Department of Ophthalmology and Visual Science, Yale University, New Haven, CT, <sup>4</sup> Department of Neurology, Yale University,  
New Haven, CT, <sup>5</sup> Yale School of Medicine, New Haven, CT, <sup>6</sup> Division of Infection, Immunity and Respiratory Medicine, University  
of Manchester, Manchester, United Kingdom, <sup>7</sup> Department of Computer Science, Yale University, New Haven, CT, <sup>8</sup> Department of  
Applied Math, Yale University, New Haven, CT, <sup>9</sup> Computational Biology, Bioinformatics Program, Yale University, New Haven,  
CT, <sup>10</sup> Department of Immunobiology, Yale University School of Medicine, New Haven, CT, <sup>11</sup> Mila – Quebec AI institute, Montréal,  
Quebec, Canada, <sup>12</sup> Department of Computer Science, Rutgers University, New Brunswick, NJ, <sup>13</sup> Department of Computational  
Mathematics, Science and Engineering; Michigan State University, East Lansing, MI, <sup>14</sup> Department of Mathematics, Michigan State  
University, East Lansing, MI, <sup>15</sup> Department of Biosystems Science and Engineering, ETH Zurich, Switzerland, <sup>16</sup> Department of  
Mathematics and Statistics, Université de Montréal, Montréal, Quebec, Canada, <sup>17</sup> Department of Genetics, Yale University, New  
Haven, CT, <sup>18</sup> Broad Institute of MIT and Harvard, Cambridge, MA.

<sup>8</sup>  
<sup>9</sup> \* These authors contributed equally. § These authors jointly supervised the work. Correspondence to Brian Hafler, 300 George Street,  
<sup>10</sup> New Haven, CT, 06511. Email: [brian.hafler@yale.edu](mailto:brian.hafler@yale.edu). Smita Krishnaswamy, 333 Cedar St,  
<sup>11</sup> New Haven, CT 06520. E-mail: [smita.krishnaswamy@yale.edu](mailto:smita.krishnaswamy@yale.edu).

<sup>12</sup> **Abstract**

<sup>13</sup> Discovery of effective therapies for neurodegeneration represents a major obstacle, in part due to challenges in  
<sup>14</sup> disease monitoring and drug access to the central nervous system (CNS). Age-related macular degeneration  
<sup>15</sup> (AMD) is a neurodegenerative disease of the CNS, affecting the retina, where anti-VEGF therapy is the only  
<sup>16</sup> effective intervention. Due to commonalities in disease progression and pathophysiology, AMD may represent a  
<sup>17</sup> uniquely accessible model to investigate therapies for other neurodegenerative diseases, leading us to examine  
<sup>18</sup> whether pathways of disease progression are shared across neurodegenerative conditions. We used single-nucleus  
<sup>19</sup> RNA sequencing to profile lesions from the macula of 17 retinas with varying degrees of AMD pathology as  
<sup>20</sup> well as controls. To identify populations of cells implicated in disease, we created a robust computational  
<sup>21</sup> pipeline, which first identified subsets of cells enriched in particular phases of disease progression and then  
<sup>22</sup> created signatures of disease pathogenesis in a cell type specific manner. By applying this pipeline to single-cell  
<sup>23</sup> profiling of human AMD lesions, we identified activated microglial and astrocyte populations that are enriched  
<sup>24</sup> in the early phase of disease. Examining single-cell data generated from Alzheimer's disease and progressive  
<sup>25</sup> multiple sclerosis, we found a similar activation profile in transcriptional populations enriched in the early phases  
<sup>26</sup> of these neurodegenerative diseases. In late-stage AMD, our pipeline identified a microglia-to-astrocyte signaling  
<sup>27</sup> axis mediated by IL-1 $\beta$  which drives VEGFA expression and angiogenesis characteristic of disease pathogenesis.  
<sup>28</sup> We validated this mechanism using *in vitro* and *in vivo* assays, identifying a possible new therapeutic target for  
<sup>29</sup> AMD and possibly other neurodegenerative conditions. Thus, due to the shared glial states across diseases, the  
<sup>30</sup> retina provides a novel human system for investigating therapeutic approaches in multiple neurodegenerative  
<sup>31</sup> diseases.

<sup>32</sup> **Introduction**

<sup>33</sup> AMD is a neurodegenerative disease of the retina that affects 196 million individuals worldwide and has a  
<sup>34</sup> significant impact on patient's quality of life [1]. Similar to other neurodegenerative diseases of the central  
<sup>35</sup> nervous system (CNS), such as Alzheimer's disease (AD) and progressive multiple sclerosis (MS), AMD can  
<sup>36</sup> be categorized into stages. Initially, in the early, 'dry' stage of AMD, extracellular amyloid-beta containing  
<sup>37</sup> deposits known as drusen accumulate in the retina, leading to the activation of glia [2]. In advanced, 'neovascular'  
<sup>38</sup> AMD, angiogenesis and fibrosis driven by vascular endothelial growth factor (VEGF) cause photoreceptor  
<sup>39</sup> and vision loss [3]. In MS and AD, glial dysregulation is associated with neuronal damage and progressive

40 neurologic impairment [4, 5]. This raises the question of whether pathogenic glia activation states are shared  
41 across neurodegeneration, and whether the human retina can be used as a model for interventions targeting glial  
42 for similar neurodegenerative diseases.

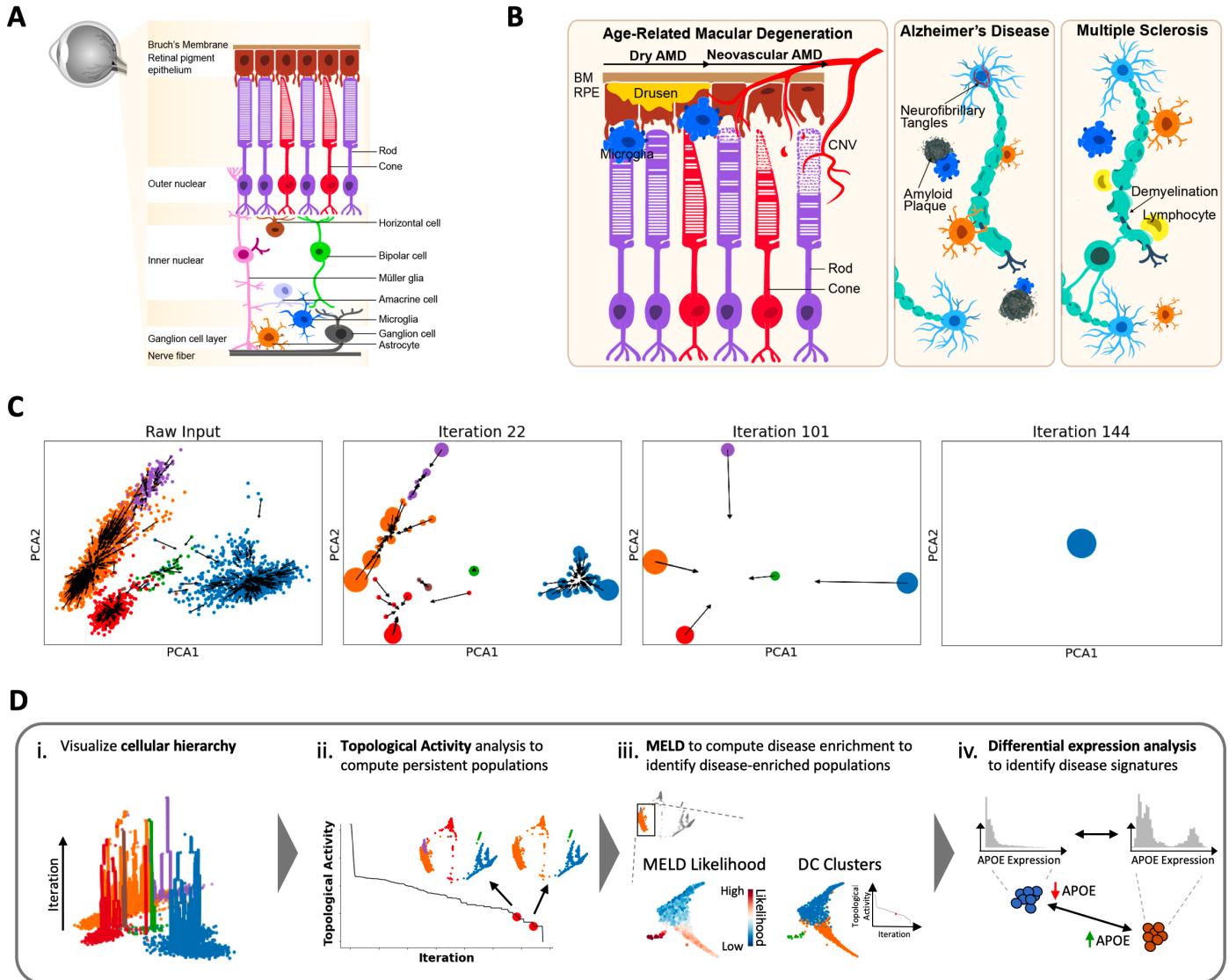
43 While single-cell transcriptomics has given insight into the cellular perturbations in AD and MS [4–7], a  
44 single-cell transcriptomic analysis of AMD has not been performed. To identify cell types and states enriched  
45 across stages of AMD, we performed massively parallel microfluidics-based single nucleus RNA-sequencing  
46 (snRNA-seq) to create the first single-cell transcriptomic dataset of AMD pathology, comprising 71,063 cells  
47 across multiple stages of disease. In such large data sets, identifying cellular populations that drive disease  
48 and could be targeted for therapeutic benefit remains a challenge with current approaches. This often occurs  
49 because pathogenic populations may be a small subset of a recognized compartment of the tissue. Thus, it  
50 can be challenging to identify such populations among the noise and complexity present in single-cell data. To  
51 address this, we developed a topologically-inspired machine learning suite of tools called Cellular Analysis with  
52 Topology and Condensation Homology (CATCH). At the center of this framework is a pathogenic-population  
53 discovery pipeline whose key component is a method called diffusion condensation [8]. Diffusion condensation  
54 identifies groups of similar cells across scales systematically to discover subpopulations of interest within a data  
55 diffusion framework. In this approach, cells are iteratively pulled towards the weighted average of their neighbors  
56 in high dimensional gene space, slowly eliminating variation. When cells come close to each other, diffusion  
57 condensation merges them together, creating a new cluster. When combined with a single cell differential  
58 abundance method MELD [9], diffusion condensation can identify distinct subpopulations associated with disease  
59 progression. This represents an improvement over clustering tools that partition the data based on metrics of  
60 cluster interconnectedness. Since this approach identifies specific disease-enriched populations, condition-specific  
61 signatures can be built and compared across neurodegenerative conditions, helping build a common understanding  
62 of shared disease mechanisms.

63 Using the CATCH pipeline, we identified two populations of activated glia, one microglial subset and one  
64 astrocyte subset, enriched in the early phase of dry AMD. These subsets were characterized by signatures of  
65 phagocytosis, lipid metabolism and lysosomal functions. Surprisingly, by reapplying our pipeline to AD [4] and  
66 MS [5] single-cell data sets, we identified the same signatures in the early phases of multiple neurodegenerative  
67 diseases, indicating a common mechanism for glial activation in the early phase of neurodegeneration. The  
68 microglia and astrocyte expression signatures were validated in human retinal and brain tissue. In late  
69 stage, neovascular AMD, CATCH identified an inflammasome expression signature in microglia as well as a

70 proangiogenic signature in astrocytes. Through computational receptor-ligand interaction analysis, we identified a  
71 key signaling axis between microglia-derived IL-1 $\beta$  and pro-angiogenic astrocytes, the driver of neovascularization  
72 and photoreceptor loss in advanced disease in AMD [3]. Through a combination of human induced pluripotent  
73 stem cell (iPSC)-derived astrocyte stimulation assays, *in vivo* mouse experiments, and analysis of postmortem  
74 human AMD retinal samples, we validated this pro-angiogenic microglial-astrocyte axis mediated by IL-1 $\beta$  in  
75 late-stage neovascular AMD. As inflammasome and glial IL-1 $\beta$  signaling are important in AD and MS [10–12],  
76 these pathways represent glial molecular signatures shared between neurodegenerative conditions that affect the  
77 retina and the brain. This study offers both a framework for identifying disease-affected cellular populations and  
78 disease signatures from complex single cell data as well as key insights into the shared drivers of neurodegeneration.

79 **Results**

80 **CATCH efficiently identifies, characterizes and compares disease-enriched populations in**  
 81 **complex single-cell transcriptomic data**



**Fig. 1: Overview of neurodegenerative disease processes and the topological diffusion condensation approach.** (A) Sketch of retina cross section showing layers and major cell types. (B) Illustration of the role of innate immune cells in neurodegenerative disease pathogenesis. In the dry stage of AMD, there is accumulation of extracellular drusen debris between Bruch's membrane (BM) and the retinal pigment epithelium (RPE), leading to activation of glia. In the neovascular late-stage of AMD, VEGF-mediated choroidal neovascularization (CNV) develops, which can lead to vision loss through rod and cone photoreceptor cell death. Accumulation of extracellular plaques and intracellular neurofibrillary tangles in Alzheimer's disease and myelin damage in progressive multiple sclerosis are both accompanied by microglia (blue) and astrocyte (orange) activation. (C) Visual description of cellular condensation process undertaken by diffusion condensation across four granularities. Points are moved to and merged with their nearest neighbors as determined by a weighted random walk over the data graph. Over many successive iterations, cells collapse, denoting cluster identity at various iterations. (D) The coarse graining process described in (C) creates hundreds of granularities of clusters which can be analyzed in meaningful ways: i) we can visualize the hierarchy of clusters computed by diffusion condensation, to identify the merging behavior across granularities; ii) we can identify meaningful, persistent partitions of the data by performing topological activity analysis; iii) in conjunction with MELD [9], we can scan across these meaningful granularities to identify resolutions that optimally split disease-enriched populations of cells from healthy populations of cells and finally; iv) we can compute differentially enriched genes between populations of interest.

82 As parts of the central nervous system (CNS), the retina contains many different functional layers and distinct  
 83 strata that are occupied by a highly diverse set of cell types and states (Fig. 1A). Furthermore, as a component  
 84 of the CNS, the retina shares features with the brain at the level of cell biology and degenerative pathology (Fig.

85 1B). Similar to AMD, MS and AD have defined disease phases, each with an early or acute active, and a late  
86 or chronic inactive disease stage [13–15]. To identify pathogenic cellular states enriched in AMD, and relate  
87 them to states found in AD and MS, we performed massively parallel microfluidics-based snRNA-seq to profile  
88 lesions from the macula of 11 retinas with varying degrees of AMD pathology and 6 control samples, creating a  
89 single-cell view of AMD pathology. We then applied a pipeline, CATCH, to parse this dataset into meaningful  
90 groupings of cell-types and states to identify pathogenic mechanisms of disease which may be shared across  
91 neurodegenerative conditions. We used snRNA-seq for our analysis, which has been shown to perform well for  
92 sensitivity and cell-type classification as compared to scRNA-seq [16]. snRNA-seq has the added advantages  
93 that it minimizes gene expression changes resulting from tissue dissociation as well as minimizes challenges in  
94 dissociation for tissues such as the retina and brain

95 Cells can exist in various transcriptional states, which naturally fall into a hierarchy or organization. Within  
96 this hierarchy, cells of a more similar functional niche, for instance microglia and astrocytes, are more closely  
97 related to one another than cells of a more disparate niche, for instance microglia and endothelial cells. Learning  
98 this hierarchy from data is important to the development of a systematic understanding of biological function  
99 and can provide insight into mechanisms of disease pathogenesis. As cell types may be differentially affected by  
100 disease, the simultaneous identification and characterization of abundant classes of cells at coarse granularity as  
101 well as rare cell types or states at fine granularity provides a comprehensive framework for defining, modeling, and  
102 understanding specific cellular pathways in disease. While biological data has structure at many different levels of  
103 granularity, most clustering methods offer one or just a few levels of granularity. These few levels of granularity  
104 can create inaccurate identifications of disease-associated cellular states. To address this, we developed CATCH,  
105 a framework that combines the principles of data manifold geometry with computational topology to create  
106 a better understanding of cellular states across granularities. While the core component of CATCH, diffusion  
107 condensation [8], and its mathematical properties [17] have been established and used to identify multigranular  
108 structure in biomedical datasets [18], it has not been applied to single cell transcriptomic data. Here, we adapted  
109 and built a pipeline around diffusion condensation to systematically sweep through all possible granularities of  
110 the cellular hierarchy to identify pathogenic populations, and infer mechanisms of neurodegeneration.

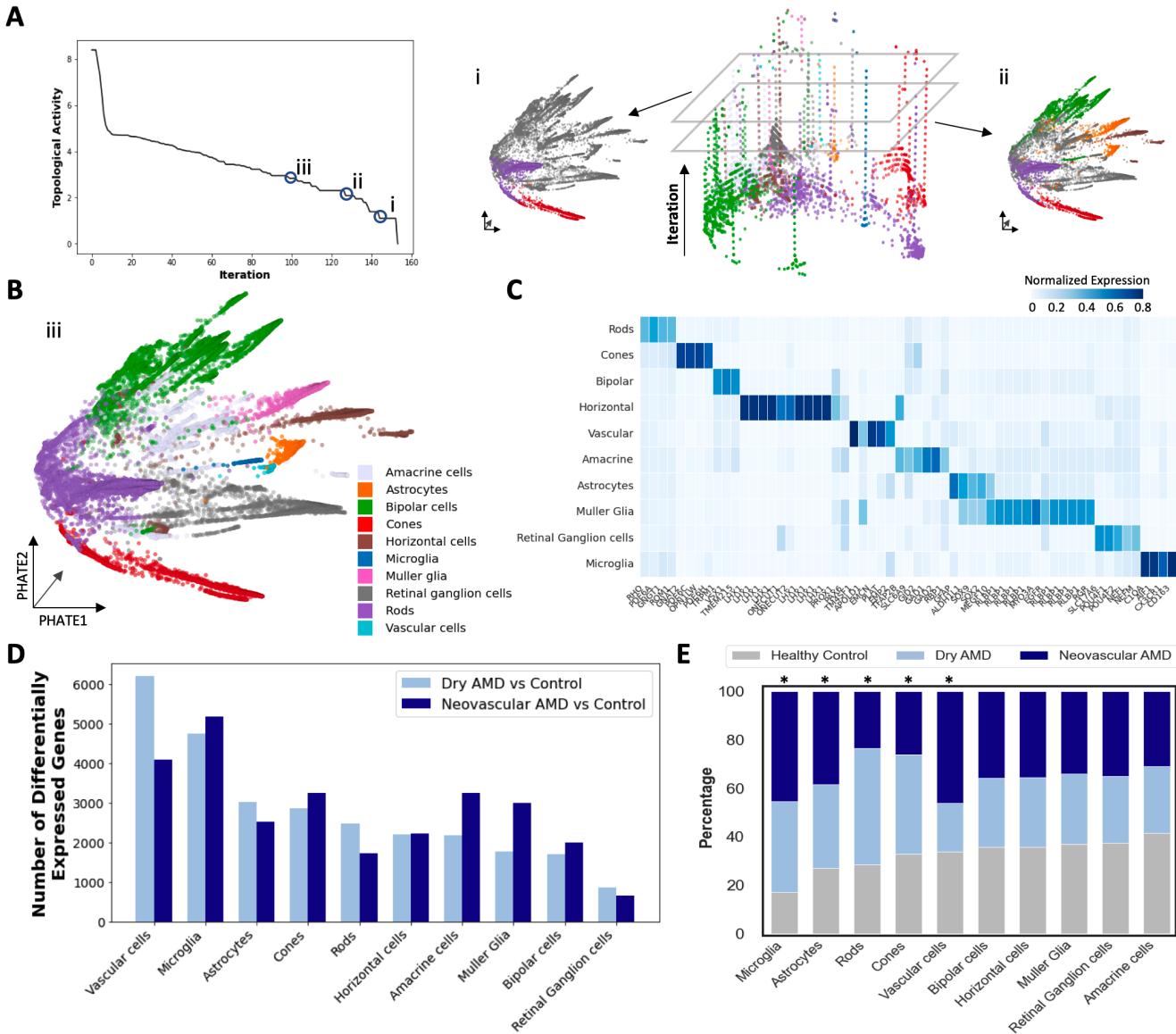
111 To learn the cellular hierarchy from complex single cell transcriptomic data, we adapted diffusion condensation  
112 to efficiently move cells towards their most similar neighbors in terms of their transcriptomic profile across  
113 successive iterations. When cells collapse into one another, diffusion condensation merges them together, thereby  
114 clustering them at a specific level of granularity (Fig. 1C). By slowly condensing and then merging similar cells,

115 diffusion condensation effectively learns how cells relate to one another over hundreds of levels of granularity.  
116 Since diffusion condensation does not force cells to merge at any given iteration, as done by other hierarchical  
117 clustering approaches, the length of time a cell, or cluster of merged cells, remains persistent denotes not only  
118 their transcriptomic interrelatedness but also their uniqueness from other cells. Cells that take only a few  
119 iterations to merge are very similar to one another, while cells that take a significant number of iterations to  
120 merge are more different in their overall transcriptomic profile. This approach is fundamentally separate from  
121 popular community detection clustering methods based metrics such as modularity and silhouette score, which  
122 optimize cluster labels based on network interconnectedness. Diffusion condensation is a coarse graining approach  
123 which slowly merges similar populations together across scales. This feature of the algorithm allows us to perform  
124 downstream analysis and identify populations enriched in disease states.

125 The CATCH framework utilizes the persistence characteristic of diffusion condensation to learn and analyze  
126 the cellular hierarchy to identify pathogenic transcriptomic states and to create robust signatures of disease from  
127 single cell data. The cellular hierarchy is visualized to identify the hierarchical and persistence structure of the  
128 data ([Fig. 1D-i](#)). Meaningful granularities of the cellular hierarchy are identified through topological activity  
129 analysis, an analysis that identifies highly persistent and stable granularities for downstream characterization  
130 ([Fig. 1D-ii](#)). With this analysis, we identify clusters that isolate cells found disproportionately in pathogenic or  
131 healthy samples using the single cell enrichment analysis method MELD [9] ([Fig. 1D-iii](#)). Finally, we create rich  
132 signatures of disease by identifying differentially expressed genes in pathogenic populations of cells using a fast  
133 modification to Earth Mover’s Distance that leverages the cellular hierarchy ([Fig. 1D-iv](#)).

134 For additional details on each component of the CATCH pipeline, including the adaptions to diffusion  
135 condensation, visualization of the cellular hierarchy, topological activity analysis and our implementation of  
136 differential expression analysis, see methods.

137 To validate the computational analysis, we perform ablation studies on each component of the CATCH  
138 pipeline ([Extended Data Fig. 2](#)). To benchmark our approach’s ability to find known cell populations, we  
139 perform rigorous comparisons against established clustering methods using 40 synthetic and real world datasets  
140 ([Extended Data Fig. 1](#) and [Extended Data Fig. 3](#)). Finally, we show the ability of this pipeline to identify rare  
141 cell types ([Extended Data Fig. 5](#)) and signatures of disease populations in real single cell data ([Extended Data](#)  
142 [Fig. 9](#)). For an overview of computational analysis and comparisons, see methods.



**Fig. 2: Single-nucleus RNA-seq profiling of the macula from human individuals with varying stages of AMD pathology.** (A) (left) Topological activity analysis of human retina single cell data across all condensation iterations. By computing gradients on topological activity (see Methods), we identify three granularities at which persistent partitions of the data occur (represented by resolutions i, ii and iii), and select them for downstream analysis. (right) Condensation process of AMD single cell data visualized across iterations (from bottom to top) with the most coarse-grained granularity clusters visualized on PHATE embedding: resolution i. represents the most coarse-grained clusters and resolution ii. represents the second most coarse-grained clusters. (B) Populations identified at the finest granularity identified by topological activity analysis (resolution iii.) were visualized and all populations were assigned a cell type based on which cell type gene signature they displayed the highest expression of. (C) Cell type specific genes visualized along with average normalized expression of known cell type specific marker genes. All major retinal cell types were identified by CATCH process described in A-B. (D) Differentially expressed genes identified by Wasserstein Earth Mover's Distance (EMD) between cells from early-stage dry and late-stage neovascular AMD lesions and cells from control retinas on a cell type specific basis. Number of significantly differentially expressed genes between control and AMD cells reported in a cell type and stage specific manner (FDR corrected p-value < .1). Cell types sorted by most differential genes between dry AMD and control comparison. Vascular cells, microglia and astrocytes have the most differentially expressed genes in dry AMD compared to control samples. (E) Bar chart indicates the contribution of cell types in each cluster from control, dry AMD and neovascular AMD samples. Microglia and astrocytes are the most statistically significantly enriched cell types in AMD, while rods and cones are the most depleted cell types in neovascular AMD. Vascular cells are the most enriched cell type in the neovascular AMD condition. All statistics were computed using two-sided multinomial test with multiple comparisons correction ( $p < 0.01$ ).

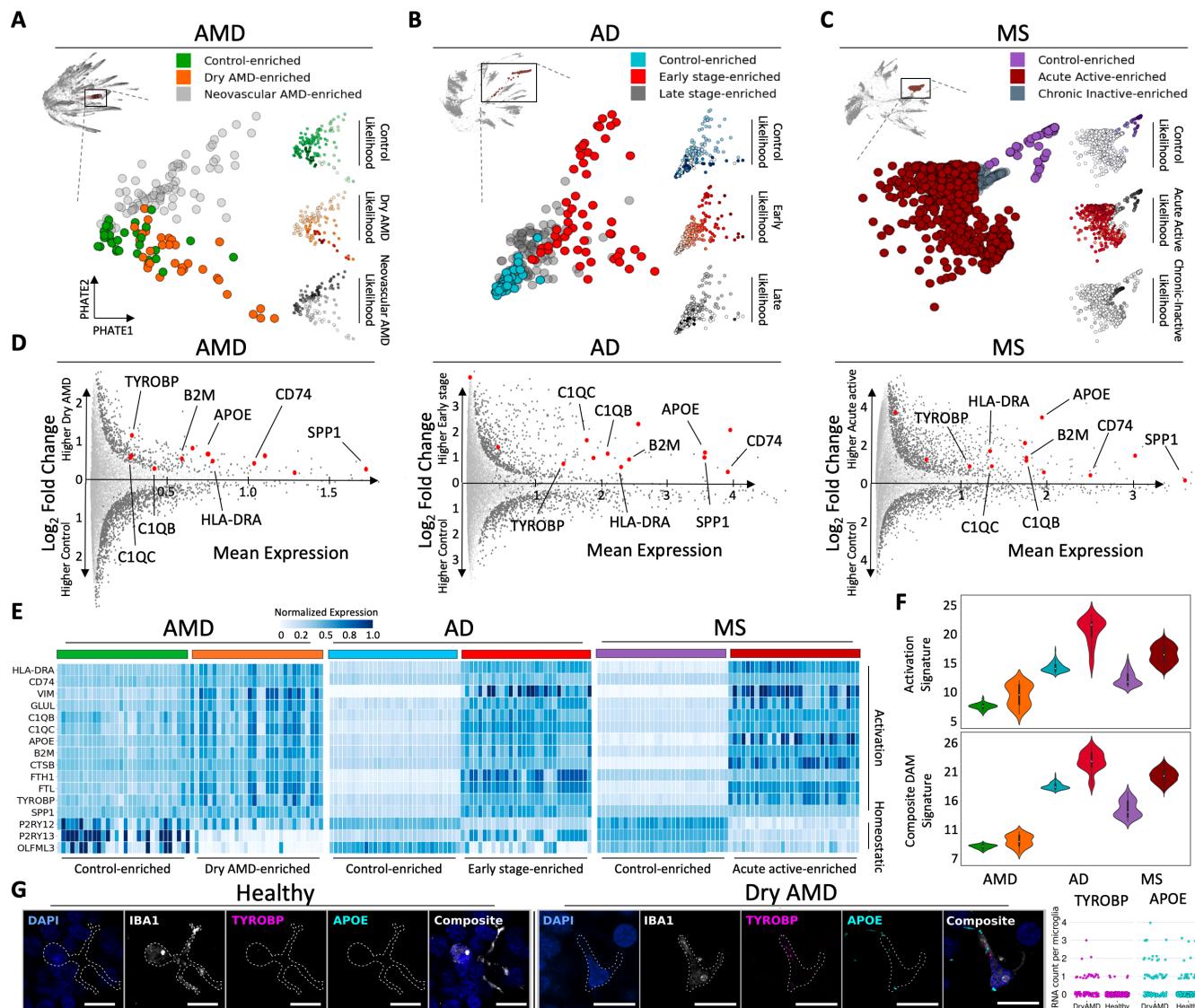
144 We applied CATCH to the AMD snRNA-seq dataset to identify the major cell types present in the control  
 145 and AMD samples. We performed topological activity analysis and identified three granularities of the cellular  
 146 hierarchy for downstream analysis (granularities with low activity and high persistence). We visualized the

snRNA-seq dataset using PHATE and the CATCH-defined clusters at the coarsest two identified granularities (Fig. 2A). When visualizing the third granularity, we observed a number of clusters, which we categorized as cell types based on the expression of previously established cell type specific marker genes [19] (Extended Data Fig. 4A) (see Methods). Using this approach, we identified neuronal cell types, including retinal ganglion cells, horizontal cells, bipolar cells, rod photoreceptors, cone photoreceptors, and amacrine cells, as well as rare non-neuronal cell types, including microglia, astrocytes, Müller glia, and vascular cells (Fig. 2B-C). To determine if these populations could be found with established approaches, we applied Louvain [20] clustering to the AMD single cell data. Louvain revealed 22 populations at coarse granularity, and 40 populations at fine granularity (Extended Data Fig. 5A-B). Across both resolutions, however, rare innate immune cell types such as microglia, astrocytes and Müller glia, were not identified with the Louvain method, with markers specific for these cell types not localizing to any one cluster. Finally, to demonstrate the ability of CATCH to identify meaningful populations of cells across granularities, we further explored subtypes of bipolar cells, a diverse set of interneurons that transmits signals from rod and cone photoreceptors to retinal ganglion cells [21–23]. By analyzing a coarse granularity of the bipolar cells, we identified the first two major subtypes, ON-center and OFF-center (Extended Data Fig. 4B). By analyzing a finer granularity, we identified all 12 major subtypes of cells based on the expression of cell subtype specific marker genes (Extended Data Fig. 4C-E).

To identify cell types implicated in AMD pathogenesis in an unbiased manner, we applied condensation-based differential expression analysis to the CATCH-identified cell types. By comparing the cells that originated from retinas with either dry or neovascular AMD to the cells from control retinas, we identified differentially expressed genes using Earth Mover's Distance within each cell type (set FDR corrected p-value < .1 across all comparisons) [24]. By analyzing the number of differentially expressed genes across all cell types, we found that vascular cells, microglia, and astrocytes had the greatest number of differentially expressed genes across stages of AMD compared to control samples (Fig. 2D). Furthermore, we performed abundance analysis to identify if certain cell types were significantly more enriched in either dry or neovascular AMD. This analysis revealed a statistically significant increase in the proportion of microglia and astrocyte nuclei from donors with both dry and neovascular AMD compared to control samples (two-sided multinomial test, p-value < .01) (Fig. 2E). Furthermore, there was a statistically significant enrichment of vascular cells in neovascular AMD, highlighting the importance of vascular cells in the development of pathological angiogenesis present at that stage of disease (two-sided multinomial test, p-value < .01). There was a relative decrease in abundance of both rod and cone photoreceptors in advanced neovascular AMD, consistent with the known loss of photoreceptors in the advanced

177 stage of disease (two-sided multinomial test, p-value < .01)([Fig. 2E](#)). These findings suggest that non-neuronal  
178 cell types including microglia, astrocytes, and vascular cells are important cell types in AMD pathogenesis, with  
179 not only the most transcriptional alterations but also changes in abundance during AMD progression.

180 Microglial activation signature identified in dry AMD is shared across the early phase of  
 181 multiple neurodegenerative diseases



**Fig. 3:** Fine grain analysis of microglia reveals a shared activation signature enriched in the early phase of three different neurodegenerative diseases. (A) 141 microglia identified by diffusion condensation at coarse granularity (upper left) can be further subdivided into three clusters at fine granularity, each enriched for cells from a different disease-state. Disease state enrichment was calculated using MELD (right) for each condition: Control (top), dry AMD (middle) and neovascular AMD (bottom), with higher MELD likelihood scores shown with darker colors. A resolution of the condensation homology which optimally isolated MELD likelihood scores from each condition was identified using topological activity analysis. Microglia are revisualized using PHATE. (B) As in panel A, three subsets of 288 microglia are found in AD with diffusion condensation and topological activity analysis, each enriched for cells from a different stage of pathology as computed by MELD (right). Cells are revisualized with PHATE. (C) As in panel A, three subsets of 1263 microglia are found in MS with diffusion condensation and topological activity analysis, each enriched for cells from a different stage of disease as computed by MELD (right). Cells are revisualized with PHATE. (D) Differential expression analysis between control-enriched and early or acute active disease-enriched microglia across neurodegenerative diseases reveals a shared activation pattern in early disease (increased expression of *TYROBP*, *B2M*, *APOE*, *CD74*, *SPP1*, *HLA-DRA*, *C1QB*, *C1QC*). Significance values (dark grey) were assigned to genes based on FDR corrected p-value < .1. (E) Heatmap demonstrating differences in expression of the neurodegenerative shared activation pattern and a homeostatic signature between control-enriched and early or acute active disease-enriched microglia across neurodegenerative diseases. Color conventions are as in panels A-C. Rows correspond to genes and columns represent individual cells. We have plotted 40 cells from each dataset selected through random sampling to reveal the difference between control-like and early disease-like cellular states. (F, upper) Composite microglial activation signature for the neurodegenerative shared activation pattern in control-enriched and early or acute active disease-enriched microglia across neurodegenerative diseases (y-axis - gene expression of signature). (F, lower) Disease associated microglia (DAM) signature (from [7]) for control-enriched and early or acute active disease-enriched microglia across neurodegenerative diseases. Color conventions are as in panels A-C. (y-axis - gene expression of signature) (G) Micrographs of combined in-situ RNA hybridization and IBA1 immunofluorescence demonstrating elevated expression of key components of the neurodegenerative shared activation pattern (*TYROBP* and *APOE*) in IBA1-positive cells, a marker of microglia, from retinas with dry AMD (right group) compared to control retinas (left group). All scale bars = 10µm. The average number of puncta identified per IBA1-positive cell for *TYROBP* was  $0.28 \pm 0.05$  in dry AMD (n=191) vs.  $0.02 \pm 0.01$  for control (n=464; p < 1e-10; Chi-square test for 0 vs. >0). The average number of puncta identified per IBA1-positive cell for *APOE* was  $0.57 \pm 0.09$  in dry AMD vs.  $0.14 \pm 0.03$  for control (p < 1e-08; Chi-square test for 0 vs. >0).

182 While microglia activation states and their dynamics have been identified in mouse models of AD [7] and related  
183 expression states found in humans [25], it is not well understood to what extent these states and dynamics are  
184 shared across human neurodegenerative diseases. The study of microglia in the CNS has been difficult due  
185 to their rarity, requiring focused enrichment strategies [7, 25]. With the ability of CATCH to sweep across all  
186 hierarchies of clusters, we can identify subpopulations of rare cell types at fine granularity to perform a rigorous  
187 and in-depth analysis of cellular states. To identify microglial subpopulations enriched in specific phases of AMD  
188 and build transcriptomic signatures of disease, we identified CATCH granularities that isolated high MELD  
189 likelihood scores computed for control, dry, and neovascular AMD conditions. We computed MELD likelihood  
190 scores for each condition on all microglia in AMD (Fig. 3A). Next, we identified a granularity highlighted  
191 by topological activity analysis that partitioned regions of high disease likelihood from regions of low disease  
192 likelihood (see Methods). With this approach, we identified three clusters, each enriched for a different condition:  
193 a cluster enriched for cells from control samples, a cluster enriched for cells from early, dry AMD samples, and a  
194 cluster enriched for cells from late-stage, neovascular AMD samples (Fig. 3A).

195 To identify signatures of AMD present in microglia during the early stage of dry disease pathogenesis, a phase  
196 in which microglia have been previously implicated [2], we performed differential expression analysis between  
197 control-enriched and the dry AMD-enriched clusters. Analyzing the top most differentially expressed genes  
198 (FDR corrected p-value < .1) between these subpopulations, a clear activation signature appeared in the early,  
199 dry AMD-enriched cluster, including *APOE*, *TYROBP*, and *SPP1* (Fig. 3D), genes known to play a role in  
200 neurodegeneration [7]. The association of TYROBP and APOE were validated on sections of human retinal  
201 macula by simultaneous immunofluorescence for IBA1, a microglia-associated gene, and in situ hybridization for  
202 TYROBP and APOE. On sections of human retinal macula, IBA1-positive cells from patients with dry AMD  
203 showed enrichment relative to controls for gene transcripts from *TYROBP* and *APOE*, indicating polarization of  
204 a subset of microglia towards the neurodegenerative microglial phenotype in early disease (Fig. 3G). Increased  
205 expression of *TYROBP* and *APOE* in microglia was also identified using in situ hybridization on lesions from  
206 human brain tissue with early stage AD and early progressive MS compared with controls (Extended Data  
207 Fig. 7C).

208 Due to the similarity between this activation state and a previously defined disease-associated microglial state  
209 described in mice [7, 26], we performed a comprehensive analysis of microglial states in two other neurodegenerative  
210 diseases, AD and progressive MS. Applying the CATCH approach to snRNA-seq data from AD [4] and MS [5],  
211 we identified all major cell types based on the expression of cell type specific marker genes (Extended Data

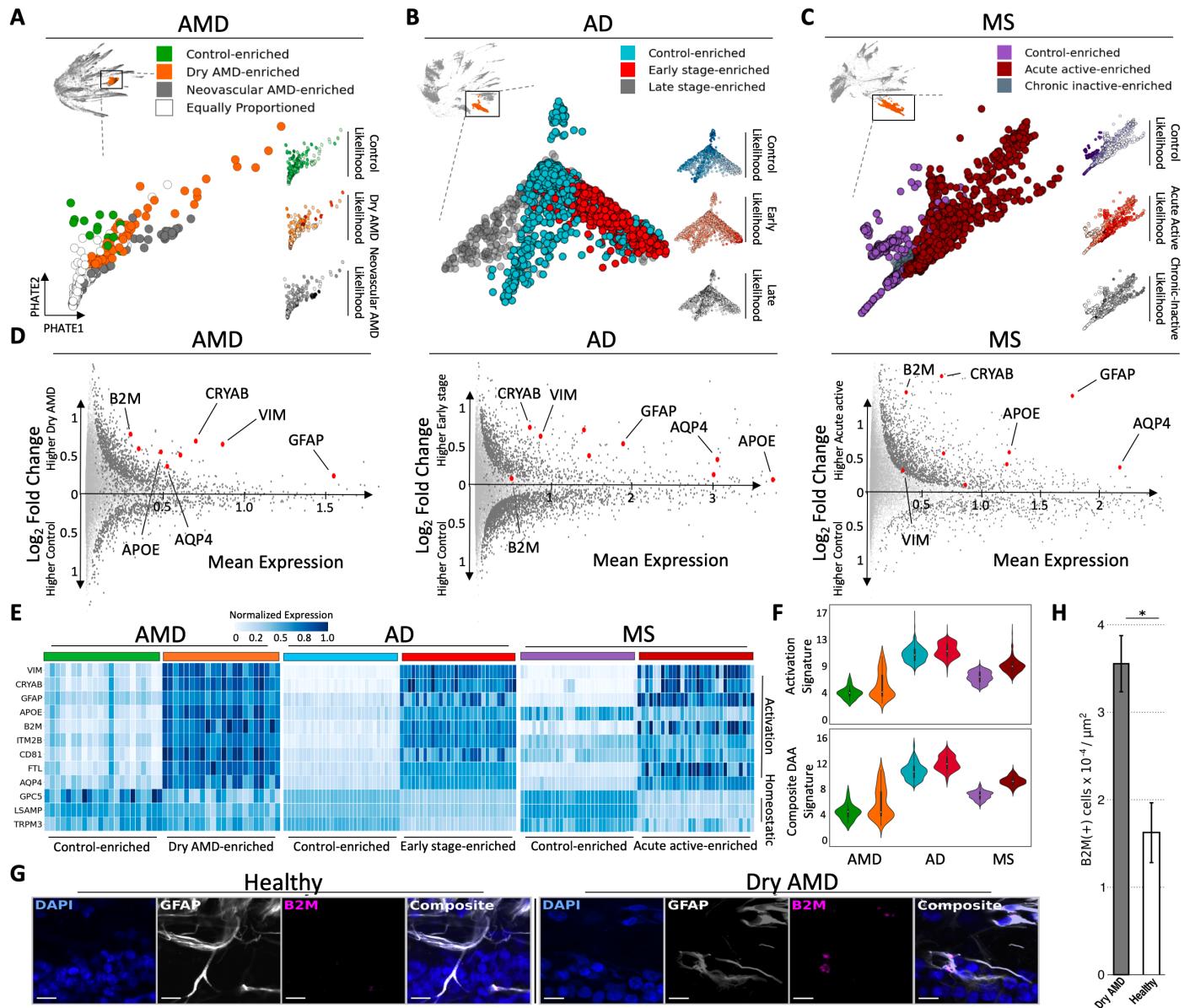
212 Fig.6A-D). As in AMD, enrichment analysis revealed that microglia were significantly enriched in AD and MS  
213 when compared to control brain tissue (Extended Data Fig.6E-F). Similar to our analysis of AMD identifying  
214 disease-phase specific transcriptomic states, we applied MELD and topological activity analysis to microglia in  
215 the AD and MS datasets and identified three clusters of microglia in each disease: a cluster enriched for cells  
216 from control brain tissue; a cluster enriched for cells from early-stage AD tissue or acute active MS lesions; and  
217 a cluster enriched for cells from late stage AD tissue or chronic inactive MS lesions (Fig. 3B,C). Differential  
218 expression analysis between the control-enriched and the early disease-enriched clusters yielded a common shared  
219 activation profile in all three diseases when analyzing the top differentially expressed genes (Fig. 3D, middle and  
220 right panels) (FDR corrected p-value < .1)).

221 To understand the early disease enriched microglial populations, we visualized the microglial activation  
222 signature (*CD74*, *SPP1*, *VIM*, *FTL*, *B2M*) (*APOE*, *TYROBP*, *CTSB*) (*C1QB* and *C1QC*) as well a homeostatic  
223 signature (*P2RY12*, *P2RY13*, and *OLFML3*) on control-enriched and early disease-enriched clusters from  
224 neurodegenerative diseases (Fig. 3E). A clear divergence is seen between the expression pattern of the homeostatic  
225 signature in control-enriched populations and early disease-enriched populations across conditions. With higher  
226 expression of activation genes and lower expression of homeostatic genes, the early activated population of  
227 microglia display a divergent polarization state. We built a composite microglial activation signature and mapped  
228 it onto the clusters along with a previously described disease-associated microglia signature found in an AD  
229 mouse model [7]. The early stage of neurodegenerative disease-enriched clusters displayed higher expression of  
230 both signatures compared with the control-enriched clusters (Fig. 3F with expression values ranging from 5 to  
231 25 for our activation signature and 7 to 26 for DAM signarure).

232 This shared neurodegenerative microglial phenotype across AMD, MS, and AD involves upregulation of  
233 multiple genes implicated in studies of neurodegenerative disease risk. These include *APOE*, a key regulator  
234 of the transition between homeostatic and neurotoxic states in microglia [27] strongly implicated in risk for  
235 AD [28,29] and AMD [30]; *TYROBP* which encodes the TREM2 adaptor protein DAP12, mutations of which are  
236 implicated in a frontal lobe syndrome with AD-like pathology [31] and expression of which is upregulated in white  
237 matter microglia in MS lesions; *SPP1* (osteopontin), implicated in microglial activation in brains affected by  
238 MS [32] and AD [33]; and *CTSB*, encoding the major protease in lysosomes cathepsin-B, which is upregulated in  
239 microglia responding to β-amylloid plaques in AD [33]. Initiation of the pathologic accumulation of extracellular  
240 material occurs by different means in these three neurodegenerative diseases. However, the finding that microglial  
241 phagocytic, lipid metabolism, and lysosomal activation pathways are upregulated in the early or acute active

<sup>242</sup> stage of all three diseases suggests a convergent role for dysregulation in microglia directed towards clearance of  
<sup>243</sup> extracellular deposits of debris.

244 **Astrocyte activation signature identified in dry AMD is shared across the early phase of**  
 245 **multiple neurodegenerative diseases**



**Fig. 4: Fine grain analysis of astrocytes reveals a shared activation signature enriched in the early phase of neurodegenerative diseases.** (A) 474 astrocytes identified by diffusion condensation at coarse granularity (upper left) can be further subdivided into three clusters at fine granularity, each enriched for cells from a different stage of neurodegenerative disease. Disease state enrichment was calculated using MELD (right) for each condition: Control (top), dry AMD (middle) and neovascular AMD (bottom), with higher MELD likelihood shown with darker colors. A resolution of the condensation homology which optimally isolated MELD likelihood scores from each condition was identified using topological activity analysis. Astrocytes are visualized using PHATE. (B) As in panel A, three subsets of 2,361 astrocytes are found in AD with diffusion condensation and topological activity analysis, each enriched for cells from a different stage of disease as computed by MELD (right). Astrocytes are visualized with PHATE. (C) As in panel A, three subsets of 5,469 astrocytes are found in MS with diffusion condensation and topological activity analysis, each enriched for cells from a different stage of MS as computed by MELD (right). Astrocytes are visualized with PHATE. (D) Differential expression analysis between control-enriched and early stage of neurodegenerative disease-enriched clusters across neurodegenerative diseases reveals a shared activation pattern in the early stage of disease. This signature includes *B2M*, *CRYAB*, *VIM*, *GFAP*, *AQP4*, *APOE*, *ITM2B*, *CD81*, *FTL*. Significance values (dark grey) were assigned to genes with FDR corrected p-value < .1. (E) Heatmap demonstrating differences in astrocyte expression of the neurodegenerative shared activation pattern and a homeostatic signature between control-enriched and early or acute active disease-enriched astrocytes across neurodegenerative diseases. Color conventions are as in panels A-C. Rows correspond to genes and columns represent individual cells. We have plotted 40 cells from each data-set selected through random sampling to reveal the difference between control-like and early disease-like cellular states. (F) Composite astrocyte activation signature (top) and disease-associated astrocyte signature (DAA) for the neurodegenerative shared activation pattern in control-enriched cluster and early disease-enriched cluster across neurodegenerative diseases. Color conventions are as in panels A-C. (G) Micrographs of combined in-situ RNA hybridization and GFAP immunofluorescence showing more abundant *B2M* expression in astrocyte-rich retinal layers from dry AMD retina when compared to control. All scale bars = 10 $\mu\text{m}$ . (H) Bar plot showing density of *B2M* transcripts in retina samples affected by dry AMD and control.

246 While astrocyte transcriptomic states and dynamics have been established in mouse models of AD, astrocyte  
247 profiles have not been profiled in human AMD lesions at a single-cell resolution [6]. As our initial analysis  
248 implicated astrocytes in disease pathogenesis (Fig. 2D-E), we performed similar cross-disease analysis within the  
249 astrocyte populations using the CATCH method. Using MELD and topological activity analysis, we identified  
250 four clusters of astrocytes at fine granularity within the diffusion condensation hierarchy: a cluster enriched for  
251 cells from control samples, a cluster enriched for cells from patients with early, dry AMD, a cluster enriched for  
252 cells from patients with late-stage neovascular AMD and a cluster with equal numbers of cells from all three  
253 conditions (Fig. 4A). When comparing the transcriptomic profiles of cells within the dry AMD-enriched and  
254 control-enriched astrocyte populations, key activation and degeneration associated genes, such as *GFAP*, *VIM*,  
255 and *B2M* were upregulated (Fig. 4D).

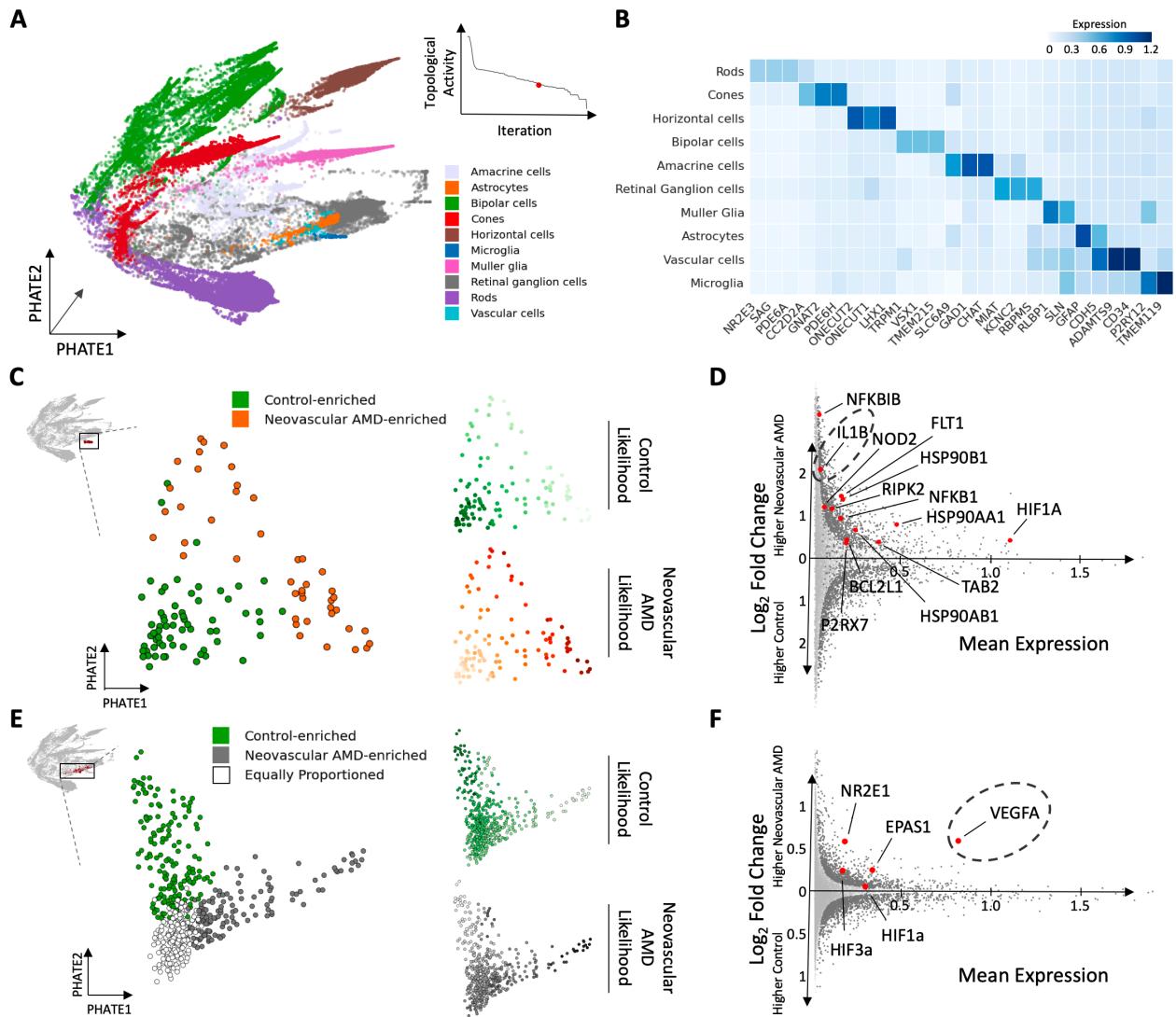
256 Using MELD and topological activity analysis, we identified clusters that isolated stage-specific populations  
257 within MS and AD astrocytes. In both diseases, we identified three clusters: a cluster enriched for cells from  
258 control brain tissue, a cluster enriched for cells from early-stage AD tissue or acute active MS lesions, and a  
259 cluster enriched for cells from late-stage AD tissue or chronic inactive MS lesions (Fig. 4B, C). By comparing  
260 the control-enriched and early disease-enriched clusters within each dataset using condensed transport, we  
261 identified a shared gene signature enriched in the early stage neurodegenerative disease subcluster across all  
262 three diseases (Fig. 4E). The integrated gene signature included markers of activated astrocytes, including  
263 *VIM*, *GFAP*, *CRYAB*, and *CD81* [34, 35], major histocompatibility complex (MHC) class I (*B2M*) [36, 37], iron  
264 metabolism (*FTH1* and *FTL*), a water channel component implicated in debris clearance (*AQP4*) [38], along  
265 with lysosomal activation and lipid and amyloid phagocytosis (*CTSB*, *APOE*). Of interest, many upregulated  
266 genes were shared between the microglial and astrocyte early stage activation signatures, suggesting common  
267 glial stress pathways become activated in neurodegeneration.

268 Similar to microglia, we mapped homeostatic (*GPC5*, *LSAMP*, *TRPM3*) and composite activation signatures  
269 (*B2M*, *CRYAB*, *VIM*, *GFAP*, *AQP4*, *APOE*, *ITM2B*, *CD81*, *FTL*) to early disease-enriched and control-enriched  
270 astrocyte clusters across neurodegenerative diseases. Similar to microglia, the composite activation signature and  
271 homeostatic signatures were divergently expressed by early enriched clusters (Fig. 4E,F upper with expression  
272 values ranging from 0 to 17). Using a recently published disease-associated astrocyte signature established in an  
273 AD mouse model [6], we built a composite activation signature and mapped that onto the early-disease and  
274 control-enriched clusters across conditions. The early disease-enriched clusters displayed higher expression of the  
275 disease-associated astrocyte (DAA) gene signature in addition to the composite activation signature (Fig. 4F,

276 lower with expression values ranging from 0 to 16).

277 To validate the astrocyte signature in tissue, we performed simultaneous GFAP immunofluorescence and  
278 RNA in situ hybridization for *B2M*, a component of MHC-I and member of the shared gene signature on sections  
279 of the human macula. The retinal layers occupied by GFAP-positive astrocytes (inner plexiform layer to inner  
280 limiting membrane) contained a higher density of *B2M* transcripts in retinas affected by dry AMD relative to  
281 control retina (p-value < 1e-03, two-sided Student's t-test) ([Fig. 4G, H](#)).

282 Microglia display inflammasome activation signature and astrocytes display pro-angiogenic  
 283 signature in late-stage neovascular AMD



**Fig. 5: Cell-type specific changes in gene expression during AMD disease progression.** (A) PHATE visualization of 46,783 nuclei isolated from neovascular AMD and control retinas [39]. CATCH analysis identified a resolution of the condensation homology, which isolated cell types. As in Figure 3, each cellular cluster was assigned a cell type identity based on which gene signature it expressed at the highest level. (B) CATCH identified cell types, as shown by the average normalized expression of known cell type specific marker genes. (C) Disease state enrichment was calculated using MELD (right) for each condition: Control (top), and neovascular AMD (bottom), with higher MELD likelihood shown with darker colors. A resolution of the condensation homology which optimally isolated MELD likelihood scores from each condition was identified using topological activity analysis. Microglia are revisualized using PHATE. Two subsets of microglial cells, one enriched for microglia from retinas with neovascular AMD and another from control retinas. (D) Differential expression analysis between control-enriched and neovascular disease-enriched microglial clusters revealed a different activation pattern in late disease. Significance values (dark grey) were assigned to genes with FDR corrected p-value < .1. This signature includes *NFKBIB*, *IL1B*, *NOD2*, *FLT1*, *HSP90B1*, *RIPK2*, *NFKB1*, *HSP90AA1*, *HIF1A*, *BCL2L1*, *P2RX7*, *TAB2*, *HSP90AB1*. (E) Disease state enrichment was calculated using MELD (right) for each condition: Control (top) and neovascular AMD (bottom) with higher MELD likelihoods shown with darker colors. A resolution of the condensation homology which optimally isolated MELD likelihood scores from each condition was identified using topological activity analysis. Astrocytes are revisualized using PHATE. CATCH identified three subsets of astrocyte cells, one enriched for astrocytes from neovascular retinas, another from control retinas and a third equally split between conditions. (F) Differential expression analysis between the control-enriched and neovascular disease-enriched astrocyte clusters reveals a different activation pattern in late-stage neovascular disease. Significance values (dark grey) were assigned to genes with FDR corrected p-value < .1. This signature includes *NR2E1*, *EPAS1*, *VEGFA*, *HIF1a*, *HIF3a*.

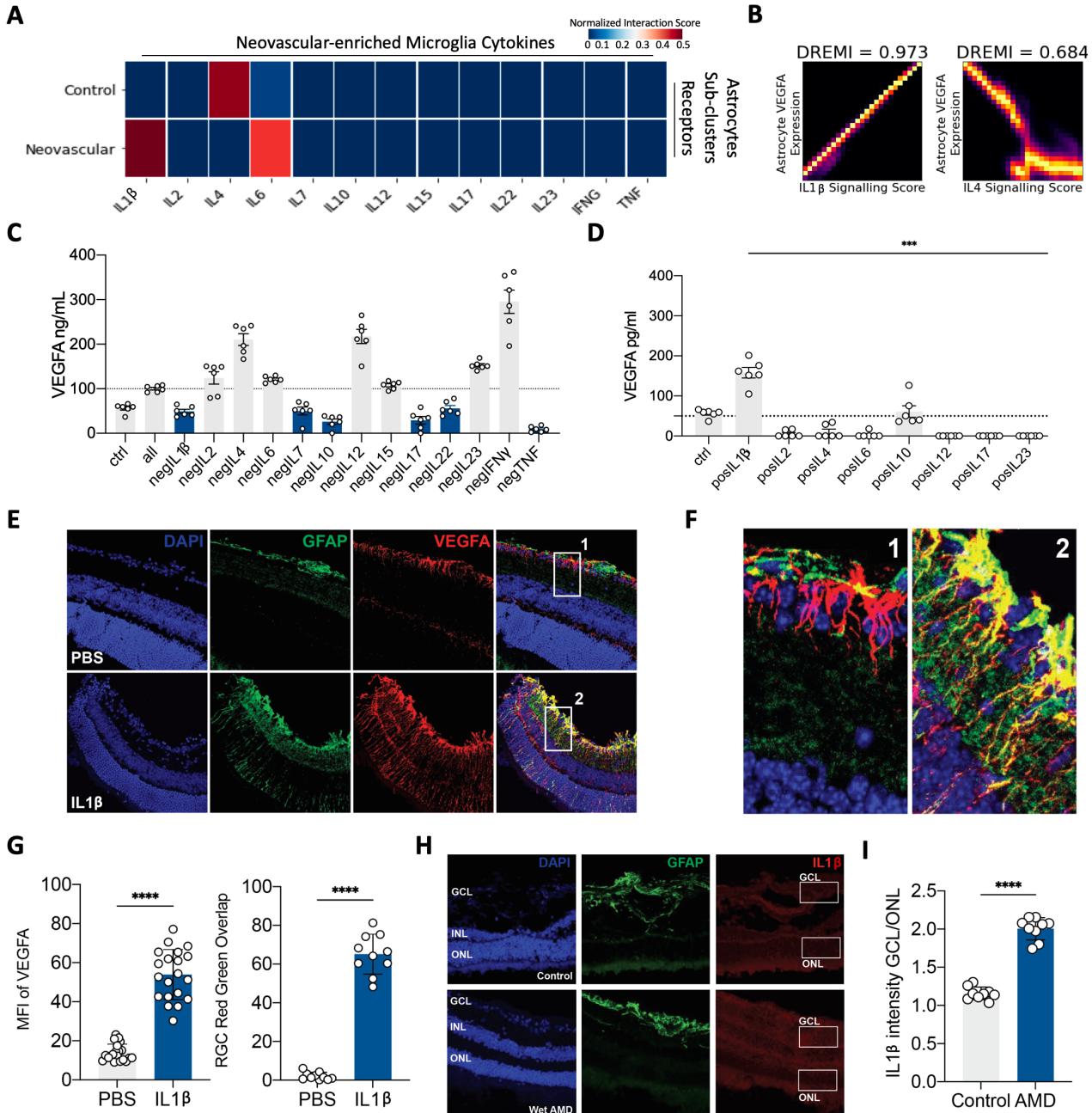
284 While glial activation signatures are shared during the early phase of multiple neurodegenerative disease, it is of  
 285 interest to understand if they persist or evolve in the late stage of neurodegenerative diseases. To understand  
 286 these glial activation dynamics across stages of AMD, AD and MS, we performed differential expression analysis

287 between the early stage of neurodegenerative disease-enriched clusters and the late stage of neurodegenerative  
288 disease-enriched clusters of astrocytes and microglia. Across both comparisons, molecular signatures present  
289 in the early stage of AMD, MS, and AD are not detected in microglia and astrocytes during the late stage  
290 of neurodegeneration ([Extended Data Fig.8A-B](#)), indicating transcriptional changes in glia during disease  
291 progression.

292 To examine the transcriptional changes in glia during progression from early dry to late-stage neovascular AMD  
293 pathology, we performed snRNASeq on three additional retinas from human donor retinas with neovascular AMD,  
294 and applied the CATCH analysis to 46,783 nuclei when combined with the previously sequenced samples. We  
295 identified a granularity of the CATCH hierarchy with low topological activity and assigned cell type labels based  
296 on the expression of cell type-specific gene signatures ([Fig. 5A,B](#)). Following the fine grained CATCH analysis,  
297 we identified two clusters of microglia: one cluster enriched for cells from control retinas and one cluster enriched  
298 for cells from late-stage, neovascular AMD retinas ([Fig. 5C](#)). To identify cell-type specific transcriptional changes  
299 in the subpopulation of microglia enriched in late stage neovascular AMD pathology, we performed condensation-  
300 based differential expression analysis between the control-enriched and the neovascular AMD-enriched clusters.  
301 Analyzing the top differentially expressed genes between these subpopulations (FDR corrected p-value < .1)  
302 revealed an inflammasome-related signature including *IL1B*, *NOD2*, and *NFKB1*. The pro-IL-1 $\beta$  protein requires  
303 both cleavage and release via inflammasome-mediated caspase activation and pyroptosis for bioactivity [[40](#)].  
304 Here, activation of inflammasome sensors and oligomerization into proteolytically active complexes may occur in  
305 response to a significant and lasting drop in oxygen tension or chronic lipid exposure [[40, 41](#)], both known to  
306 drive inflammasome activation via NLRP3 (NOD-, LRR- and pyrin domain-containing 3) ([Fig. 5D](#)). In late stage  
307 AD and MS alternative cellular stress-associated pathways were upregulated including transcriptional regulators  
308 of the ER stress response (*XBP1*) and their target genes involved in protein folding and transport (*HSPA1A*,  
309 *HSPA1B*, *HSP90AA1*) and glycosylation (*ST6GAL1* and *ST6GALNAC3*), as well as regulators of autophagy  
310 and proteostasis (*ATG7*, *MARCH1*, *USP53*). These signatures highlight a shared cellular stress induction.

311 Using the fine grained CATCH workflow, we identified two astrocyte subpopulations: one cluster enriched  
312 for cells from control retinal samples and one cluster enriched for cells from late-stage, neovascular AMD  
313 retinal samples ([Fig. 5E](#)). To identify signatures of AMD present in astrocytes during the late stage of disease  
314 pathogenesis, we performed condensation-based differential expression analysis between control-enriched and the  
315 neovascular AMD-enriched clusters. Analyzing the top differentially expressed genes (FDR corrected p-value<.1)  
316 between these subpopulations revealed elevation of *VEGFA*, *NR2E1*, and *HIF1A* expression ([Fig. 5F](#)), all of

317 which are regulators of cellular responses to low oxygen tension [42–44]. While VEGFA is known to be an  
318 important mediator of the abnormal blood vessel growth that characterizes late-stage neovascular AMD and is  
319 the target of current therapies for the treatment of disease [30, 45, 46], our data demonstrate in humans a specific  
320 subpopulation of retinal astrocytes that are a source of this signal.



**Fig. 6: Identifying cytokine regulators of astrocyte VEGFA secretion.** **(A)** Interaction analysis between diffusion condensation identified subtypes of astrocytes and neovascular-enriched microglia (detailed in Figure 5) computed with CellPhoneDB [47]. Interactions between cytokines produced from neovascular-enriched microglia were computed against cytokine-receptors on astrocyte subtypes. Interactions between specific cytokine-receptor pairs were added to produce a single cytokine interaction value for control and neovascular astrocyte subtypes. **(B)** DREM1 association analysis between astrocyte VEGFA expression, IL-1 $\beta$  signaling score, and IL-4 signaling score. Signaling scores for IL-1 $\beta$  and IL-4 were computing by adding receptor expression of IL-1 $\beta$  and IL-4 respectively neovascular-enriched astrocytes from Figure 5. **(C)** Conducted negative screen in human iPSC-derived astrocytes 24 hours after stimulation, subtracting one cytokine (e.g. 'negIL2') from the combinatorial pool to test its necessity in generating a VEGFA-producing astrocyte compared to vehicle control (ctrl). All represents stimulation with a mixture of cytokines (IL-1 $\beta$ , IL-2, IL-4, IL-6, IL-7, IL-10, IL-12, IL-15, IL-17, IL-22, IL-23, IFN $\gamma$ , TNF). VEGFA protein is measured using enzyme-linked immunosorbent assay (ELISA). **(D)** Conducted single cytokine positive screen in human iPSC-derived astrocytes to test the sufficiency of each cytokine to stimulate astrocyte VEGFA production. VEGFA protein levels are measured using ELISA 24 hours after stimulation with each cytokine compared to vehicle control (ctrl). **(E)** IL-1 $\beta$  or PBS was injected intravitreally into a mouse eye. Retinas were collected 72 hours later for immunofluorescent imaging. **(F)** Zoomed in images of regions indicated in E. **(G)** Quantification of mean fluorescence intensity (MFI) of VEGFA after injection of IL-1 $\beta$  or PBS in the mouse eyes after 72 hours (left) and quantification of amount of VEGFA and GFAP overlap in the ganglion cell layer of the mouse retina after injection of IL-1 $\beta$  or PBS (right). **(H)** Immunofluorescence imaging of human post-mortem control and neovascular AMD retinas. **(I)** Quantification of IL-1 $\beta$  intensity in the ganglion cell layer (GCL) over the outer nuclear layer (ONL) of the retina from F.

322 As microglia are known to influence astrocyte functional states through the secretion of soluble factors, we wanted  
323 to determine if microglia-derived cytokines could drive VEGFA expression from retinal astrocytes [48–50]. Since  
324 CATCH was able to isolate astrocyte and microglial states, we utilized CellPhoneDB interaction analysis [47] to  
325 create a putative list of possible microglia-derived cytokines that may interact with astrocytes to drive VEGFA  
326 expression (Fig. 6A). From this analysis, the neovascular-enriched microglia cluster interacted most significantly  
327 with astrocytes through IL-1 $\beta$  and IL-6, while in controls, microglia-astrocyte interaction was primarily mediated  
328 by IL-4. Furthermore, IL-1 $\beta$  interacted most significantly with the neovascular-enriched astrocyte subpopulation.  
329 Using conditional-Density Resampled Estimate of Mutual Information (DREMI), a method to identify non-linear  
330 associations in data [51], we find that IL-1 $\beta$  signaling on astrocytes was most significantly associated with  
331 astrocyte production of VEGFA. Meanwhile IL-4 signaling was most significantly associated with a decrease in  
332 astrocyte VEGFA production (Fig. 6B). We then set out to validate the cytokine regulators of astrocyte VEGFA  
333 production in an unbiased manner.

334 Cytokines are a part of a complex network of proteins that can produce additive, synergistic, or antagonistic  
335 effects. To demonstrate this relationship, we used two screening methods. We first used a combinatorial screening  
336 approach utilizing all cytokines identified in our snRNAseq dataset, removing one at a time to test its necessity  
337 in creating a VEGFA expressing astrocyte. Screening with human iPSC-derived astrocytes demonstrated that  
338 IL-1 $\beta$ , IL-10 and IL-17 are positive regulators of VEGFA production in these cells as their subtraction causes  
339 decreased VEGFA compared to human iPSC-derived astrocytes stimulated with all cytokines (Fig. 6C). We  
340 then tested the sufficiency of some of these cytokines being able to regulate VEGFA production by completing  
341 a single protein stimulation and noted that, interestingly, only IL-1 $\beta$  caused astrocyte VEGFA secretion (Fig.  
342 6D). Across both analyses, IL-1 $\beta$  positively regulated induction of VEGFA from astrocytes both *in vitro* (Fig.  
343 6C-D) and *in silico* (Fig. 6B). Our analysis of VEGFA regulation validated the computational prediction of IL-4  
344 being a negative regulator of VEGF-A production (Fig. 6B-C), showing the utility of our approach in identifying  
345 signaling interactions between cellular subsets identified with CATCH.

346 With identification of cytokine mediators of astrocyte VEGFA production, we validated our findings *in vivo*  
347 by injecting IL-1 $\beta$  intravitreally in a mouse. This resulted in upregulation of VEGFA (Fig. 6E, F). Not only  
348 was there an increase in the amount of VEGFA (Fig. 6G, right), there was an increase of overlapping signals of  
349 GFAP and VEGFA, indicative of astrocyte VEGFA activation and secretion (Fig. 6G, left), along with VEGFA  
350 expression extending from ganglion cell layer localization down to other layers of the retina. Altogether, this  
351 demonstrated the sufficiency of cytokines such as IL-1 $\beta$  to induce VEGFA secretion in astrocytes *in vitro* and *in*

352 *vivo*. Cytokines such as IL-1 $\beta$  are increased in the vitreous of patients with neovascular AMD [52], but source  
353 and the role of these cytokines in angiogenesis has not been explored. We undertook immunohistochemical  
354 staining for IL-1 $\beta$  in retinal samples from the macula of patients with AMD and healthy controls, observing that  
355 there was an increased amount of IL-1 $\beta$  intensity in the inner retinal layers, where astrocytes reside (Fig. 6I).  
356 Furthermore, upregulation of VEGFA was seen in these areas (Fig. 6G), indicating that the phenomenon we  
357 observe *in vitro* and in mice likely occurs in human neovascular AMD as well (Fig. 6G-I).

## 358 Discussion

359 Here, we used snRNA-seq to generate the first single-cell transcriptomic atlas of AMD during pathological  
360 progression, as well as develop a machine learning pipeline that allows for meaningful comparisons between cell  
361 types and states across diseases and phases. To generate rich signatures for cross-disease comparison among rare  
362 cellular subpopulations, we applied a topology-inspired suite of machine learning tools for single-cell analysis,  
363 ‘CATCH’, a tool that identifies cellular subpopulations enriched in a specific condition. This pipeline identified  
364 cell states enriched in disease, characterized pathogenic expression signatures, and predicted cellular interactions  
365 between pathogenic populations, uncovering potential therapeutic targets.

366 Using CATCH, we identified and characterized specific subpopulations of microglia and astrocytes enriched  
367 in the early stage of dry AMD displaying activation signatures related to phagocytosis, lipid metabolism, and  
368 lysosomal function. Surprisingly, we found similar populations of microglia and astrocytes in analyses of previously  
369 published AD and MS single-cell data. While initial inciting events likely differ between neurodegenerative  
370 conditions, lipid-rich extracellular plaques play a prominent role in each condition. It is likely that glial cells  
371 coordinate clearance of extracellular debris and, in turn, become activated. While the initial phagocytic clearance  
372 may be beneficial, glial activation has been shown to play a role in degeneration in AMD, AD, and MS. In later  
373 stages of disease, this shared activation landscape evolves. In advanced neovascular AMD, our analysis identified  
374 a microglia inflammasome-related signature that drives pro-angiogenic astrocyte polarization and pathologic  
375 neovascularization. Microglial inflammasome activation and subsequent IL-1 $\beta$  release could be mediated by  
376 a variety of signaling sensors. The NLRP3 sensor may be activated in response to a variety of stress signals,  
377 including extended lipid exposure or prolonged hypoxia, and has been previously implicated as a microglial  
378 driver of neurodegenerative immunopathology, making it a likely candidate [53]. Microglia are highly mobile cells  
379 and responsive to a wide variety of stimuli. While lineage tracing that definitively differentiates mononuclear  
380 phagocyte origin into circulating macrophages, tissue resident macrophages, and microglia remains challenging,

381 it is believed that the mononuclear phagocytes found at the apical side of the RPE in the vicinity of drusen,  
382 which induce activation of the inflammasome, come from all three populations [54]. Furthermore, emerging data  
383 suggests that the inflammasome and IL-1 $\beta$  have critical roles in promoting degeneration in MS and AD [10–12].  
384 Thus our results implicate this immune sensor in AMD as well.

385 This set of analyses has clear implications for potential therapeutics for AMD and other neurodegenerative  
386 diseases. Currently, anti-VEGF therapy is the primary intervention approved to treat AMD and is only effective  
387 in the most advanced stage of disease. Our unbiased analysis not only identified the cell-type specificity of  
388 VEGFA expression but also identified pathogenic signaling interactions that promote AMD disease progression.  
389 Currently, therapies that inhibit IL-1 $\beta$  are available and used in clinical practice for the treatment of other  
390 diseases. Inhibiting microglia-derived IL-1 $\beta$  in neovascular AMD could provide therapeutic benefit, preventing  
391 further neovascularization in advanced patients, or even preventing neovascularization before it begins in  
392 patients with earlier stages of disease. Since these mechanisms are shared across MS and AD, it is plausible  
393 that these interventions could provide benefit to patients suffering from other neurodegenerative conditions as  
394 well. Identifying promising therapeutic candidates to test in neurodegenerative disease clinical trials remains  
395 vitally important, and our data suggest that approaches targeting glia may be broadly applicable to multiple  
396 neurodegenerative diseases.

## 397 Acknowledgements

398 We would like to thank the retina donors and their families for their contribution to this work. Without their  
399 sacrifice, our study would not have been possible.

400 This work was supported by NEI K08-EY026652, the Thome Memorial Foundation, and the Doris Duke  
401 Charitable Foundation (to B.P.H.), by NIAID training grant 1F30-AI157270 (to M.K.) and by NIAID 5U19-  
402 AI089992-08 (to S.K.) and by NIGMS 1RO1-1355929 (to S.K. and G.W.). We thank the Advancing Sight  
403 Network and the Lions Gift of Sight Eye Bank for timely retrieval of donor eyes.

## 404 Author Contributions

405 Conception: M.K., M.M., S.K. B.P.H. ; Design of Work: M.K., M.D., E.C., S.K. B.P.H. ; Acquisition of Data:  
406 M.D., E.C., M.I., L.Z., M.M., Y.X., B.P.H. ; Analysis of Data: M.K., M.D., E.C., A.H.S., R.M.D., B.P.H. ;  
407 Interpretation of Data: M.K., M.D., E.C., A.S., B.R.; G.W.; S.K.; B.P.H ; Creation of New Software: M.K.,  
408 S.G., J.H. ; Writing - Drafting: M.K., M.D., E.C., B.R., B.P.H., S.K. ;

409 **Declaration of Interests**

410 Dr. Krishnaswamy is on the scientific advisory board of KovaDx and AI Therapeutics. Dr. Hafler receives  
411 research funding from Nayan Therapeutics.

412

413 **Methods**

414 **CATCH Analysis Details**

415 The CATCH framework constitutes a group of topologically inspired machine learning tools to identify, characterize  
416 and compare condition-enriched populations of cells across the cellular hierarchy. This framework is centered  
417 around the diffusion condensation process, which learns the structure of data across granularities. Beyond making  
418 significant adaptions to diffusion condensation, we have introduced new tools to help analyze the rich amount  
419 of multigranular information produced by diffusion condensation: cellular hierarchy visualization, topological  
420 activity analysis, automated cluster characterization and differential expression analysis.

421 In the following sections, we provide a thorough description of each aspect of CATCH. This includes detailed  
422 descriptions of the diffusion condensation process as well as its relationship with MELD, Wasserstein earth  
423 mover's distance and topological activity analysis. We complete this section with a rigorous set of comparisons  
424 to benchmark our method.

425 **Background in manifold learning and diffusion filters**

426 Many of the core concepts in diffusion condensation and its adaptions presented here are based on advances in  
427 manifold theory and graph filters. Typically,  $n$ -dimensional data  $X = \{x_1, \dots, x_N\}$  can be modeled as originating  
428 from a  $d$ -dimensional manifold  $\mathcal{M}^d$  collected via a nonlinear function  $x_i = f(z_i)$ . This is because data collection  
429 strategies (such as single cell RNA-sequencing) create high dimensional observations even when the intrinsic  
430 dimensionality is relatively low. Algorithms that use this manifold assumption [55–58] leverage the intrinsic, low  
431 dimensional geometry of the manifold to explore relationships in data. Diffusion maps [56] presented a framework  
432 that captures intrinsic manifold geometry using random walks that aggregate local relationships between data  
433 points to reveal nonlinear geometries. These local relationships, known as affinities, are constructed using a  
434 Gaussian kernel function:

$$\mathbf{K}(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{\varepsilon}\right) \quad (1)$$

435 where  $\mathbf{K}$  is an  $N \times N$  Gram matrix and bandwidth parameter  $\varepsilon$  which controls locality. A diffusion operator  
436 is defined as the row normalization of the  $N \times N$  Gram matrix  $\mathbf{K}$ :

$$\mathbf{P} = \mathbf{D}^{-1}\mathbf{K} \quad (2)$$

437 where  $\mathbf{D}(x_i, x_i) = \sum_j \mathbf{K}(x_i, x_j)$ . The diffusion operator matrix  $\mathbf{P}$  represents single-step transition probabilities  
438 for a Markovian random walk or diffusion process. Furthermore, as shown in [56], powers of this diffusion  
439 operator  $\mathbf{P}$  (represented as  $\mathbf{P}^t$  where  $t > 0$ ) represent a  $t$ -step random walk.

440 Recent works in data diffusion [24, 59–61] have shown that this framework proposed by [56] can be used as  
441 a low pass filter when the operator  $\mathbf{P}$  is directly applied to data features, effectively moving data points close  
442 to their diffusion neighbors on the manifold. This low pass filtering process effectively removes high frequency  
443 variation, or noise, and maintains only the principle low dimensional geometry of the data manifold.

#### 444 Overview of diffusion condensation and its limitations

445 Diffusion condensation is a dynamic process the builds upon previously established concepts in diffusion filters,  
446 diffusion geometry and topological data analysis. The algorithm slowly and iteratively moves points together in  
447 a manner that reveals the topology of the underlying geometry. The diffusion condensation approach involves  
448 two steps that are iteratively repeated until all points converge:

- 449 1. Compute a time inhomogeneous Markov diffusion operator from the data;  
450 2. Apply this operator to the data as a low-pass diffusion filter, moving points towards local centers of gravity.

451 As established in prior work [8, 17, 24], the application of the operator  $\mathbf{P}$  to a vector  $\mathbf{v}$  averages the values  
452 of  $\mathbf{v}$  over small neighborhoods in the data. When applied directly to a coordinate function, this application  
453 condenses points towards local centers of gravity as determined by bandwidth parameter  $\varepsilon$ , creating a filtered  
454 set of coordinates. In this process, if  $X(0) = X$  is the original dataset with diffusion operator  $\mathbf{P}_0 = \mathbf{P}$ , then  
455  $X(1) = \bar{X} = X(0) * \mathbf{P}_0$ . While previous applications of diffusion filters simply apply one iteration of this diffusion  
456 filtering process to data, we can iterate this process to further reduce the variability in the data by computing  
457 the Markov matrix  $\mathbf{P}_1$  using the coordinate-filtered  $X(1)$ . A new filtered coordinate representation  $X(2)$  is  
458 obtained by applying  $\mathbf{P}_1$  to the coordinate functions of  $X(1)$ . Initial applications of the diffusion operator  $\mathbf{P}$   
459 to  $X$  dampens high frequency variations in the coordinate function, efficiently moving similar points close to  
460 one another. Later applications dampen low frequency variation, moving similar groups of points towards one

461 another. A more complete explanation of diffusion condensation and its mathematical properties can be found  
462 in [8] and [17].

463 In its original form, the diffusion condensation process cannot be applied to scRNAseq data. While useful for  
464 general data analysis tasks, this process has limitations:

- 465 1. the approach does not work in the non-linear space of the single cell transcriptomic manifold;  
466 2. does not scale to even thousands of data points;  
467 3. does not identify granularities of the topology which meaningfully partition the cellular state space and  
468 4. does not identify pathogenic populations implicated in disease processes.

469 In this work, we address each of these limitations and further extend the framework to efficiently perform key  
470 single cell analysis tasks such as cluster characterization and differential expression analysis.

471 To address these concerns, we have made the following significant adaptions for application to single cell data:

- 472 1. Dynamically learn the geometry of the single cell manifold with each diffusion filter using  $t$ -step random  
473 walks optimized with spectral entropy;
- 474 2. Visualize learned hierarchy via embedding the condensation tree;
- 475 3. Use topological activity to identify meaningful granularities for downstream analysis;
- 476 4. Implement diffusion operator landmarking, weighted random walks and data merging to efficiently scale to  
477 thousands of cells;
- 478 5. Implement diffusion condensation with alpha decay kernel for automated cluster characterization and  
479 efficient computation of differentially expressed genes.

480 Manifold-intrinsic diffusion condensation learns cellular hierarchy from single-cell transcriptomic  
481 data

---

**Algorithm 1** Manifold-intrinsic Diffusion Condensation

---

**Input:** Cell-by-PC data matrix  $\mathbf{X}$ , initial kernel bandwidth parameter  $\varepsilon_0$  and merge threshold  $\zeta$   
**Output:** cluster labels by iteration

- 1:  $\mathbf{X}_0 \leftarrow \mathbf{X}, i \leftarrow 0$
- 2: **while** number of points in  $\mathbf{X}_i > 1$  **do**
- 3:   Merge data points  $a,b$  if  $\|\mathbf{X}_i(a) - \mathbf{X}_i(b)\|_2 < \zeta$ , where  $\mathbf{X}_i(a)$  is the  $a$ -th row of  $\mathbf{X}_i$
- 4:   Update the cluster assignment for each original data point based on merging
- 5:    $\mathbf{D}_i \leftarrow$  compute pairwise distance matrix from  $\mathbf{X}_i$
- 6:    $\mathbf{K}_i \leftarrow$  alpha-decay kernel affinity( $\mathbf{D}_i, \varepsilon_i$ )
- 7:    $\mathbf{P}_i \leftarrow$  row normalize  $\mathbf{K}_i$  to get a Markov transition matrix (*diffusion operator*)
- 8:    $\mathbf{t}_i \leftarrow$  spectral entropy of  $\mathbf{P}_i$
- 9:    $\mathbf{X}_{i+1} \leftarrow \mathbf{P}_i^{\mathbf{t}_i} \mathbf{X}_i$
- 10:    $\varepsilon_{i+1} \leftarrow \text{update}(\varepsilon_i)$
- 11:    $i = i + 1$
- 12: **end while**

---

482 Our implementation of diffusion condensation algorithm takes a cell-by-principal component matrix  $\mathbf{X}$   
483 (typically first 50 components) and computes a diffusion operator  $\mathbf{P}$ , representing the probability distribution  
484 of transitioning from one cell to another in a single step using a  $\alpha$ -decay kernel function with fixed bandwidth  
485  $\varepsilon$  (Alg. 1: Steps 5-7). While other manifold learning techniques abstract the data to a point where derived  
486 manifold-intrinsic features have an unclear relationship with gene expression, our approach learns the manifold  
487 while working in principal components, which have a clear relationship with genes. By using the principal  
488 components as the substrate for condensation, we can easily characterize clusters and perform differential  
489 expression analysis in gene expression space in downstream analysis.

490 Another key improvement we make in the condensation algorithm is raising  $\mathbf{P}$  to the power of  $t$  (rather  
491 than 1 as in [8]), simulating a  $t$ -step random walk over the data. This approach adaptively denoises and  
492 refines these transition probabilities across iterations such that transitions occur on the non-linear single cell  
493 manifold [24,39,56]. This  $t$ -step diffusion operator  $\mathbf{P}^t$  are applied to the input data, acting as a manifold-intrinsic  
494 diffusion filter, effectively replacing the coordinates of a point with the weighted average of its  $t$ -step diffusion  
495 neighbors. We track the values of  $t$  computed across iterations and perform an ablation study to show the  
496 necessity of adaptively tuning  $t$  in each iteration of the manifold-intrinsic diffusion condensation ([Extended Data](#)  
497 [Fig. 2A-B](#)). See Alg. 1 for pseudocode of this algorithm. When the distance between two cells falls below a  
498 distance threshold  $\zeta$ , cells are merged together, denoting them as belonging to the same cluster going forward  
499 (Alg. 1: Steps 3,4). It is important to note that in the original work, [8] did not merge points. This process is

500 then repeated iteratively until all cells have collapsed to a single cluster. This merging step, implemented in our  
501 manifold-intrinsic diffusion condensation approach, allows for the fast computation of the cellular hierarchy during  
502 coarse graining. When applying this manifold-intrinsic diffusion condensation process to single-cell transcriptomic  
503 data, we can see cells condense to cluster centroids across iterations, efficiently and rigorously learning the  
504 hierarchy of single-cells (Fig. 1C). Finally, through scalable implementation tricks, such as diffusion operator  
505 landmarking [62] and weighted random walks, we have allowed diffusion condensation to scale to thousands of  
506 single cells (Extended Data Fig.2F). Additional details on the selection of  $t$  as well as scalable implementation  
507 tricks can be found below.

508 **Learning manifold geometry dynamically with spectral entropy and  $t$ -step diffusion filters**

509 While the initial implementation of diffusion condensation was created to understand multigranular structure of  
510 linear data, single cells occupy a highly non-linear space requiring manifold learning strategies [24, 39, 56]. In  
511 single cell data, technical noise, such as drop out and variation, creates measurement artifacts. When building  
512 diffusion probabilities on this sort of noisy data, high transition probabilities can be calculated between unrelated  
513 cells inappropriately. Thus, directly working with  $\mathbf{P}$ , fails to acknowledge non-linearities and technical artifacts  
514 present within single cell data. Previous work in data diffusion has shown that raising the diffusion operator  $\mathbf{P}$   
515 to the power of  $t$  refines these transition probabilities, increasing the chance of transitioning to more related  
516 cells [24, 39, 56]. This powering step allows learning of the relevant non-linear geometry of the data manifold,  
517 allowing us to ignore spurious neighbors found in the ambient measurement space of cells and instead finding  
518 diffusion neighbors that lie on the single cell manifold.

519 As single cell datasets can often suffer from different types and scales of noise, previous approaches have found  
520 that the correct number of  $t$ -steps to take must be computed adaptively in a data dependent manner [24, 63].  
521 Previously proposed strategies to select  $t$  however, are often slow, as they require trial-and-error approach  
522 which rely upon the structure of the underlying dataset. In diffusion condensation, however, the structure of  
523 the underlying dataset continuously shifts between granularities due to the repeated application of diffusion  
524 filters, making the repeated computation of  $t$  necessary and through these techniques computationally unwieldy.  
525 Therefore, we propose to select  $t$  adaptively at each condensation iteration by using a *spectral entropy* based  
526 approach. Previously, it has been shown that powering the diffusion operator  $\mathbf{P}$  differentially effects the  
527 eigenvectors of the powered matrix. While the noisy, high frequency eigenvectors rapidly reduce to zero, the  
528 more informative, low frequency eigenvectors diminish much less rapidly [24]. We reason that there is a value of

529  $t$  which optimally reduces the noisy information from the high frequency eigenvectors while maintaining the  
530 maximum information from the low frequency, informative eigenvectors. To identify this point, we compute the  
531 spectral entropy of the diffusion probabilities  $\mathbf{P}$  when powered to different levels of  $t$ .

532 Spectral entropy is defined as the Shannon entropy of normalized eigenvalues, i.e.,

$$S(\mathbf{P}, t) = - \sum_i \psi_i^t \log(\psi_i^t), \quad (3)$$

533 As there is a degree of information loss with each increasing value of  $t$ , we try to identify the point at which  
534 this information loss curve stabilizes. While powering to low values of  $t$  rapidly decreases spectral entropy as large  
535 amount of noise diminish, powering to higher values of  $t$  only slowly reduces entropy due to the slower removal  
536 of information from informative, low frequency eigenvectors. Taking the point at which this stabilization occurs  
537 as done in [39], optimally allows us to adaptively select a value of  $t$  at each diffusion condensation iteration,  
538 allowing us to produce a diffusion filter which has learned the single cell manifold.

539 In fact, deriving  $t$  adaptively in a data driven manner is critical to learning the multigranular cluster structure  
540 of data. In order to illustrate this point, we generated synthetic single cell data using Splatter [64]. As can  
541 be seen, across differing amounts of variational and drop out noise, optimally selecting  $t$  via spectral entropy  
542 produces a better set of cluster labels than when setting  $t$  in a fixed, user-determined manner (Extended Data  
543 Fig. 3B). In fact, we can see that setting  $t$  to 1 does not learn the data manifold or the cluster structure of even  
544 fairly noiseless single cell data, revealing the need for selecting a high level of  $t$  in an adaptable, data-driven  
545 manner. Finally, we see that over successive condensation steps, the complexity of the data decreases and thus  
546 requires lower levels of  $t$  to learn (Extended Data Fig. 3A).

547 **Improving scalability with weighted random walks, landmarked diffusion operators and merged  
548 data points**

549 Repeated computation of a diffusion operator from high dimensional single cell data, powering of this diffusion  
550 operator to identify the optimal value of  $t$  followed by diffusion filter application via matrix multiplication is  
551 computationally expensive. Repeating these computations, potentially hundreds of times, as done by diffusion  
552 condensation is unwieldy. In fact, this approach, in its most basic implementation, scales very poorly to high  
553 dimensional single cell data with tens of thousands of features and potentially hundreds of thousands of cells. To  
554 improve computational efficiency, we perform the following steps:

- 555 1. Merge points together that fall below a preset distance threshold  $\zeta$  to create a cluster and weighting random  
556 walks to maintain effect of data density;

557     2. Compute compressed diffusion operator through landmarking [62] to efficiently compute spectral entropy  
558     as done in [39].

559     Collectively, these advances drastically improve the computational speed of diffusion condensation ([Extended](#)  
560     [Data Fig. 2F](#)). In practice, a complete cellular hierarchy of a 13,000 cell dataset can be analyzed within 6 minutes  
561     in a Google Colaboratory notebook (a service which provides 4-core 2-GHz CPU and 20 GB of RAM for free).

562     **Visualizing and analyzing condensation tree with topological activity analysis to identify meaningful**  
563     **granularities for downstream analysis**

564     Topological data analysis (TDA) is a powerful framework that learns and analyzes data across granularities. In  
565     TDA, one identifies related data points by identifying all pairs whose distance falls below a distance threshold  
566      $\delta$  in a distance matrix  $\mathbf{D}$ . Any pair of points that falls below this threshold is deemed to be part of the same  
567     *connected component* or cluster. As  $\delta$  increases, more cell pairs will be connected, quickly creating fewer connected  
568     components, or fewer larger clusters, at coarser granularities. In topological data analysis, *persistent homology*  
569     is a principled approach to track the connected components that are created and destroyed across a range  
570     of granularities. While diffusion condensation learns the multigranular structure of data through a cascade  
571     of non-linear diffusion filtration approach instead of an increasing distance threshold, these approaches are  
572     intuitively related.

573     We can study this diffusion condensation process either in a holistic manner, evaluating all granularities  
574     simultaneously, or in a detailed manner, by evaluating meaningful granularities independently. At a high level,  
575     the cellular hierarchy can be studied by visualizing the cellular hierarchy, containing all merges across all  
576     granularities. As manifold-intrinsic diffusion condensation operates in PCA dimensions, we practically implement  
577     this visualization by stacking the first two axes of  $\mathbf{X}_i \rightarrow \mathbf{X}_{i+1} \cdots \mathbf{X}_I$ , creating a hierarchical tree that summarizes  
578     the cluster structure of the data across granularities ([Fig. 1D-i](#)).

579     For more detailed analysis, we can cut this hierarchical tree at meaningful levels to identify granularities of  
580     clusters that optimally partition cells into meaningful clusters based on the data geometry. Using *persistent*  
581     *homology*, we define a *topological activity analysis*, a technique to analyze the creation and destruction of clusters  
582     across consecutive iterations ( $\mathbf{X}_i \rightarrow \mathbf{X}_{i+1}$ ) of the manifold-intrinsic diffusion condensation process. Topological  
583     activity analysis is a variation of the total persistence summary statistic often used to characterize topological  
584     activity in classical topological data analysis [65]. In this analysis framework, we summarize the merging of points  
585     during the condensation process and assign each cluster a topological ‘prominence’ value known as *persistence*.

586 Highly persistent components are taken to represent groups of cells that are similar in their transcriptional profile  
587 and distinct from other cells. These clusters, and their associated persistence values, are best represented using  
588 a ‘persistence barcode.’ This is a visualization [66] consisting of horizontal bars of different lengths; each bar  
589 corresponds to one topological feature—a subgroup of cells in our case—while the length of each bar depicts  
590 the persistence of that feature, directly indicating to what extent the feature is prominent. Assuming that the  
591 persistence barcode consists of a set of bars with end coordinates  $\mathcal{B} := \{b_1, \dots, b_k\}$ , we calculate an activity  
592 curve  $\mathbf{A}: \mathbb{R} \rightarrow \mathbb{N}$  defined by  $\mathbf{A}(i) := |\{b \in \mathcal{B} \mid b \leq i\}|$ , i.e., the number of topological features (cell clusters) that  
593 are active and independent at a given iteration  $i$ . This activity curve, first proposed by [67] and implemented  
594 by [68], allows us to identify iterations of rapid condensation as well as iterations of relative inactivity through the  
595 gradient of  $\mathbf{A}$ . Specifically, we are interested in contiguous segments in the preimage of  $\partial\mathbf{A}/\partial i = 0$ , which we refer  
596 to as  $i$ -segments. The length of an  $i$ -segment is the number of iterations for which there is no change in topological  
597 activity. Thus, the number of iterations for which  $\partial\mathbf{A}/\partial i = 0$  provides a principled way of selecting meaningful  
598 condensation granularities computed by the diffusion condensation process. Inspired by the nomenclature of  
599 persistent homology, we refer to the length of a  $i$ -segment of no topological activity as its persistence, meaning  
600 that we are looking for the most persistent of such topological activity segments.

## 601 Identification of disease-enriched populations in conjunction with MELD

602 While analysis of the cellular hierarchy will identify populations of related cells in an unbiased and multigranular  
603 manner, it does not use condition of origin information to identify cellular populations that are enriched in  
604 disease conditions of interest. While we can integrate cells from different disease conditions in our analysis, cells  
605 of a certain pathogenic transcriptomic state may be over represented in a submanifold of a given cell type. By  
606 comparing the cells of a particular type directly to each other based on condition of origin, we dilute out this  
607 enrichment information and lose important signal. In fact, identifying these pathogenic states and comparing  
608 them directly with clustering and differential expression tools has been shown to be a more powerful method to  
609 identifying condition-enriched cell states and expression signatures [9, 69]. We explore this point later in this  
610 section.

611 To take condition-specific information into consideration, we use MELD to identify cellular populations  
612 that are enriched or depleted in different disease phases [9]. MELD is a manifold-geometry based method of  
613 computing a likelihood score for each cell, indicating whether it is more likely to be seen in the normal or diseased  
614 sample. Finding a clustering method that separates these condition-enriched groups is a difficult problem that

615 needs to be performed to identify discrete cellular populations which can be thoroughly described. To rigorously  
 616 identify cell populations with strong disease-specific enrichment signals, we combine this cell-level MELD score  
 617 with information from our topological activity analysis to identify resolutions that produce stable clusters. Then  
 618 within this stable clustering, we identify populations that are enriched in differing disease conditions.

## 619 Automated cluster characterization via manifold-intrinsic diffusion condensation

620 While identification of pathogenic cellular states is critical, biologists are more interested in what defines these  
 621 populations. Most manifold learning methods visualize or cluster populations of interest, requiring further  
 622 expensive computation to characterize cell populations and discover differentially expressed genes. As our  
 623 approach continuously condenses the transcriptomic profiles of single cells to local cluster centroids in manifold  
 624 space, at any iteration, the transcriptomic states of the condensed data can be extracted at no additional  
 625 computational cost. To enhance this convergence to centroids we implement our diffusion condensation process  
 626 with an  $\alpha$ -decay kernel ([Extended Data Fig. 2C](#)). This kernel more strongly thresholds the conversion of distances  
 627 to affinities, closely resembling the box kernel, which accurately computes cluster centroids over the course of  
 628 main point merges. When diffusion condensation merges two cells together at a particular iteration, the newly  
 629 formed point lies close to the centroid of the original two cells in transcriptomic space. Under specific conditions,  
 630 the new point is exactly the cluster centroid as delineated in the Proposition below. First, we define the  $\alpha$ -decay  
 631 kernel as:

$$\mathbf{K}_\alpha(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^\alpha}{\varepsilon^\alpha}\right), \quad i, j = 1, \dots, N. \quad (4)$$

632 The standard Gaussian kernel function as shown in equation [1](#) has an  $\alpha$  of 2. The default  $\alpha$ -decay kernel  
 633 meanwhile uses a much higher value (default in our implementation is 40), which converts close distances  
 634 into affinities much more stringently ([Extended Data Fig. 2C](#)). As  $\alpha$  increases to infinity, this kernel function  
 635 converges almost completely to the box kernel. With this kernel, we are ready to state a set of conditions under  
 636 which the diffusion condensation process can be easily characterized.

637 **Proposition 1.** Assume there exists a unique global minimum non-zero distance  $\delta_i$  between points  $x_a, x_b$  at  
 638 each iteration  $i$ , with the next pair of points at distance at least  $\delta_i + \tau_i$  with  $0 < \tau_i$ . Note that  $x_a, x_b$  could have  
 639 multiplicity greater than 1, representing clusters of size  $> 1$ . Then set the bandwidth to  $\epsilon_i := \delta_i + \tau_i/2$  at each  
 640 iteration of the condensation process. For a large enough  $\alpha$ , the diffusion condensation process will maintain two  
 641 invariants for the first  $N - 1$  steps:

- 642     1. The number of points will be  $N - i$ ;  
643     2. Unique points will be located at the centroid of their cluster.

644     *Proof.* It is easy to verify (1) and (2) hold for step zero. For all  $i < N$  and for sufficiently large  $\alpha$ ,  $\mathbf{K}_\alpha(x_k, x_j)$   
645     becomes arbitrarily close to 1 for  $(k, j) \in \{(a, a), (a, b), (b, a), (b, b)\}$  and 0 otherwise. Exactly one merge occurs  
646     at each timestep between points at  $x_a$  and  $x_b$ . Given  $\mathbf{P}_i$  as described above, they merge to the point  $\frac{|x_a|x_a + |x_b|x_b}{|x_a| + |x_b|}$ ,  
647     i.e. the cluster centroid. By induction (1) and (2) hold for all  $i < N$ .  $\square$

648     In this setting, the condensation process always converges in exactly  $N - 1$  steps. In practice, we aim  
649     for much shorter convergence times as there are many fewer than  $N - 1$  interesting levels of clustering. For  
650     50,498 cells, we find a set of parameters that allow for convergence in 150 steps. For this reason we use a larger  
651     bandwidth  $\epsilon_i$  which leads to much faster convergence and gives cluster centers at each level that are close to but  
652     not exactly the cluster centroids of the points they represent. Another factor is the setting of the  $\alpha$  parameter.  
653     Since, manifold-intrinsic diffusion condensation operates in PC dimensions, the complete gene expression profile  
654     of cluster centroid  $x_{ab}$  can easily extracted by inverting the PC dimensions. We show that this point is not only  
655     mathematically true but also empirically true in practice ([Extended Data Fig.3C](#)).

## 656     Differential expression analysis via approximation of gene Wasserstein distance

657     Beyond cluster characterization, differential expression analysis is a critical method to identify signatures of  
658     pathogenic populations. Earth Mover's Distance (EMD), also known as 'optimal transport', typically manifested  
659     in 1D-Wasserstein distance, is a popular and established method to extract differentially expressed genes between  
660     clusters [24, [70–72](#)]. EMD, however, is computationally expensive, as it computes an optimal mapping between  
661     points, running in  $\tilde{O}(n^3)$  time. Previously, tree-based implementations like FlowTree [73] and QuadTree [74] have  
662     been able to closely approximate ground truth Wasserstein distance while significantly improving runtime by  
663     constraining the transport of points through the branches of a hierarchical tree [75]. Since diffusion condensation  
664     too produces a tree embedding of the data, we utilize tree based transport for differential expression.

665     EMD, or 1-D Wasserstein distance, is a measure of distance between two distributions. For a given ground  
666     distance, the Wasserstein distance between distributions can be thought of as the minimal total distance needed  
667     to move one distribution to the other. Let  $\mu, \nu$  be two distributions on a measurable space  $\Omega$  with metric  
668      $d(\cdot, \cdot)$ , and  $\Pi(\mu, \nu)$  be the set of joint distributions  $\pi$  on the space  $\Omega \times \Omega$ , such that for any subset  $\omega \subset \Omega$ ,  
669      $\pi(\omega \times \Omega) = \mu(\omega)$  and  $\pi(\Omega \times \omega) = \nu(\omega)$ . The 1-Wasserstein distance  $W_d$  also known as the earth mover's distance

670 (EMD) is defined as:

$$W_d(\mu, \nu) := \inf_{\pi \in \Pi(\mu, \nu)} \int_{\Omega \times \Omega} d(x, y) \pi(dx, dy). \quad (5)$$

When  $\mu, \nu$  are discrete distributions over points in  $\mathbb{R}^d$ , of size  $m, n$  respectively, this can be equivalently expressed in matrix notation as:

$$\begin{aligned} W_d(\mu, \nu) &:= \min_{\Pi \geq 0} \sum_{i=1}^m \sum_{j=1}^n \Pi_{ij} d(x_i, x_j) \\ \text{subject to: } &\sum_{i=1}^m \Pi_{ij} = \nu_j, \quad \forall j \in \{1, \dots, n\} \\ &\sum_{j=1}^n \Pi_{ij} = \mu_j, \quad \forall i \in \{1, \dots, m\} \end{aligned} \quad (6)$$

671 For general ground distances this is computable using the Hungarian algorithm in  $\tilde{O}(n^3)$  time. Intuitively, the  
 672 difficulty in computing the optimal transport is finding the map  $\Pi$  which optimizes the cost within the constraints.  
 673 However, for a tree metric, this optimal map is easy to compute in closed form because there is only a single  
 674 path (through the tree) between pairs of points. This single path between pairs of points results in a reduced  
 675 computational complexity of  $\tilde{O}(n)$ . This is best understood using the Kantorovich-Rubinstein dual form of the  
 676 Wasserstein distance:

$$W_d(\mu, \nu) = \sup_{f: \|f\|_L \leq 1} \int_{\Omega} f(x) d\mu - \int_{\Omega} f(x) d\nu \quad (7)$$

677 where the witness function  $f : \Omega \rightarrow \mathbb{R}$  and  $\|\cdot\|_L$  denotes the Lipschitz norm. This dual form holds under a few  
 678 minor conditions which hold for the spaces considered here. For more information see [76].

679 Given some rooted tree  $T$  with strictly non-negative edge lengths, we define the natural tree metric  $d_T(x, y)$   
 680 as the length of the unique path between nodes  $x, y$ . We denote the mass of a distribution on a subtree  $T_r$  rooted  
 681 at node  $r$  as  $\mu(T_r) = \sum_{x \in T_r} \mu(x)$ . For each node  $v \in T$  we denote its associated parent edge as  $e_v$  with weight  
 682  $w_v$ . In this setting, it is easy to construct the optimal witness function in eq. 7. Without loss of generality,  
 683 one starts at the root  $r$  and builds  $f$  such that  $f(r) = 0$  and for each edge  $e(u, v)$  where  $u$  is a parent of  $v$ ,  
 684  $f(v) = f(u) + w_e \cdot \text{sign}(\mu(T_v) - \nu(T_v))$ . Given this construction, it is easy to see that the Wasserstein distance  
 685 with tree ground distance has the following closed form:

$$W_{d_T}(\mu, \nu) = \sum_{v \in T} w_v |\mu(T_v) - \nu(T_v)|. \quad (8)$$

686 The question then comes to: what are useful tree metrics? An ideal tree metric that has low distortion  
 687 of Euclidean space and is scalable to high dimensions. QuadTree [74] is a tree metric algorithm designed to

688 approximate the optimal transport distance between discrete measures with Euclidean ground distance by  
 689 recursively partitioning space into hypercubes, but does not scale well with dimension. Specifically, assume,  
 690 without loss of generality, that the data lies in the  $[0, 1]^d$  hypercube, then at each level  $h \in [0, H)$  divide the space  
 691 into  $2^{dh}$  hypercubes with side length  $2^{-h}$ . This forms an H-level tree with each node representing a hypercube.  
 692 If the center of the hypercube is randomly shifted, then the QuadTree distance  $W_{d_{QT}}$  has distortion at most  
 693  $O(d \log 1/\tau)$  where  $\tau$  is the minimum distance between datapoints, i.e.

$$c \cdot (d \log \tau) W_{d_{QT}}(\mu, \nu) \leq W_{\|\cdot\|_2}(\mu, \nu) \leq C \cdot (d \log \tau) W_{d_{QT}}(\mu, \nu) \quad (9)$$

694 for some constants  $c, C$  in expectation [74].

695 However, QuadTree distance scales poorly as it is computed in  $O(Nd \cdot \log(d1/\tau))$ . In the high dimensional  
 696 setting, such as snRNAseq data, the poor scaling with respect to  $d$  both computationally and in the approximation  
 697 is undesirable. In this setting [75] suggests sampling trees using furthest point clustering [77]. Furthermore, [73]  
 698 implements FlowTree, a small modification to QuadTree that makes tree Wasserstein distances significantly more  
 699 accurate with the addition of small additional computational cost.

700 Drawing from both FlowTree and QuadTree, CATCH implements a new formulation of EMD over the  
 701 diffusion condensation tree. For two diffusion condensation clusters  $a, b$  located at  $C_a, C_b$  respectively we define  
 702 the *condensation-based Wasserstein approximation distance* between them as:

$$W_{CT}(a, b, T) = \|C_a - C_b\|_2 + \sum_{e(u,v) \in T_a} w_e \cdot a(T_v) + \sum_{e(u,v) \in T_b} w_e \cdot b(T_v) \quad (10)$$

703 where  $w_e := 2^{-h} \|C_v - C_u\|_2$  for edge  $e(u, v)$  at depth  $h$  and  $a(x), b(x)$  are defined as indicator functions of their  
 704 respective clusters.

705 This leads to the following proposition stating that no matter how close we are to the settings in Proposition 1,  
 706  $W_{CT}$  still represents a valid tree Wasserstein distance between clusters.

707 **Proposition 2.** The condensation-based Wasserstein distance approximation distance  $W_{CT}$ , for any diffusion  
 708 condensation tree  $T$ , defines a valid Wasserstein distance over a tree ground distance for any two clusters in that  
 709 tree.

710 *Proof.* We show this by constructing the associated tree metric  $d_{CT}$  on an arbitrary condensation tree  $T_{CT}$  and  
 711 conclude by showing that  $W_{d_{T_{CT}}}$  is equivalent to  $W_{CT}$ . Begin by rooting the tree at a node representing  $C_a$   
 712 with two children, the root of  $T_a$  named  $r_a$  and  $C_b$ . The edge  $e(C_a, r_a)$  has weight 0 and the edge  $(C_a, C_b)$  has

713 weight  $\|C_a - C_b\|_2$ . The node  $C_b$  will have a single child node the root of  $T_a$  named  $r_b$ , and is connected by an  
 714 edge of length zero. All other nodes will be defined as in  $T_a$  and  $T_b$  with associated edge weights.

It is easy to verify that the path measure over  $T_{CT}$  construction represents a valid distance  $d_{CT}$ . Finally, we verify that the Wasserstein distance with a ground distance of  $d_{CT}$  is equivalent to  $W_{CT}$  as defined in eq. 10. Indeed, because we added a skip connection in the tree to directly connect nodes  $a, b$  with an edge of length  $\|C_a - C_b\|_2$  and since  $a(T_v)$  for  $v \in T_b$  is always zero and vice versa, we have

$$\begin{aligned} W_{d_{CT}}(a, b) &= \sum_{e(u,v) \in T_{CT}} w_e |a(T_v) - b(T_v)| \\ &= w_{e(C_a, C_b)} |a(T_{C_b}) - b(T_{C_b})| + \sum_{e(u,v) \in T_a} w_e |a(T_v) - b(T_v)| + \sum_{e(u,v) \in T_b} w_e |a(T_v) - b(T_v)| \\ &= \|C_a - C_b\|_2 |0 - 1| + \sum_{e(u,v) \in T_a} w_e |a(T_v) - 0| + \sum_{e(u,v) \in T_b} w_e |0 - b(T_v)| \\ &= \|C_a - C_b\|_2 + \sum_{e(u,v) \in T_a} w_e \cdot a(T_v) + \sum_{e(u,v) \in T_b} w_e \cdot b(T_v) \\ &= W_{CT}(a, b, T). \end{aligned}$$

715

□

716 Note that  $W_{CT}$  does not calculate the Wasserstein distance over the same tree for each set of clusters, and  
 717 as shown in [73] this often improves the accuracy as compared. In addition, it is useful conceptually but not  
 718 essential that the cluster centers  $C_a, C_b$  are near the cluster centroids. In Proposition 1 we delineated the setting  
 719 where this holds exactly, but these parameters are impractical for our efficient computation requiring  $n - 1$   
 720 diffusion steps. Instead, we are satisfied with centers that are close to the centroids but are efficiently computable  
 721 in many fewer diffusion steps. Our formulation is similar to the standard Wasserstein distance with tree ground  
 722 distance as in eq. 8, but simplified and optimized for the case of comparing clusters which are elements of the  
 723 tree metric. We make two changes. First, we add a skip connection in the tree to directly connect nodes  $a, b$   
 724 with an edge of length  $\|C_a - C_b\|_2$  as in [73] which is empirically more faithful in their experiments and ours.  
 725 Next, we note that  $a(T_v)$  for  $v \in T_b$  is always zero and vice versa, thus simplifying the second and third terms.  
 726 These two optimizations give us an algorithm that is efficient in high dimensions and is effective empirically  
 727 ([Extended Data Fig. 1E](#) and ([Extended Data Fig. 2D](#)) across granularities ([Extended Data Fig. 2E](#)).

728 Using this intuition, CATCH is able to rapidly perform differential expression analysis by approximating  
 729 the Wasserstein metric on a per gene basis along the hierarchies generated by manifold-intrinsic diffusion  
 730 condensation. Leveraging our approach's ability to summarize transcriptomic landscapes with the  $\alpha$ -decay

731 kernel, we use multiple granularities of the cellular hierarchy to accurately approximate ground truth Wasserstein  
732 distance between genes and identify cluster-specific expression signatures [75] (Fig. 1D-iv). We show that this is  
733 empirically true with our comparisons ([Extended Data Fig. 2D](#) and [Extended Data Fig. 1F](#)).

734 Inspired by previous statistically sound methods of identifying differentially expressed genes, we implement a  
735 resampling-based approach to identify true differentially expressed genes [70, 78]. In this approach, we estimate  
736 false discovery rate (FDR), which is the expected proportion of rejected null hypotheses falsely for each gene's  
737 test statistic at a given significance level [70, 78]. To calculate FDRs from our Wasserstein values, we generate  
738 a null distribution by permuting the cluster labels (in practice 1000 times) and compute Wasserstein distance  
739 between the permuted classes each time. Using the median of permuted Wasserstein distances for each gene, we  
740 create a null distribution from which we can compute p-values per gene. The attained p-values are corrected  
741 using the Benjamini-Hochberg procedure [79].

#### 742 Comparison to other clustering algorithms on synthetic and real single cell data

743 We wanted to benchmark our CATCH approach against existing clustering strategies applied to single cell  
744 data. Using a combination of 40 synthetic single cell datasets as well as real single cell and flow cytometry  
745 data, we compared the clustering performance of our adapted implementation of diffusion condensation against  
746 Louvain and Leiden, multigranular clustering techniques often applied to single cell data, as well as Seurat's  
747 Shared Nearest Neighbors clustering algorithm and FlowSOM, state-of-art methods for clustering single-cell  
748 transcriptomic and flow cytometry data respectively.

749 Splatter is a simulator of realistic single cell data where ground truth cluster labels are known [64]. Using  
750 these ground truth labels, we generated increasingly noisy single cell datasets with two different types of biological  
751 noise: variation and drop out ([Extended Data Fig. 1A](#)). With each of these datasets, we follow the CATCH  
752 framework: first we compute and visualize the condensation homology ([Extended Data Fig. 1B](#)) before performing  
753 topological activity analysis to identify the top four most persistent granularities ([Extended Data Fig. 1C](#)) and  
754 then finally computing adjusted rand index, a common measure for determining clustering accuracy against a  
755 set of ground truth cluster labels ([Extended Data Fig. 1D](#)), keeping the highest score from our comparisons.  
756 Intriguingly, the most persistent population (iv), nearly always had the highest adjusted rand index score. Using  
757 this comparison approach we compared diffusion condensation to Louvain, Leiden, and Seurat's Shared Nearest  
758 Neighbors clustering algorithms across 40 synthetic single cell datasets. For Louvain and leiden of the comparison  
759 approach, four different resolutions of clusters were computed and compared, keeping only the comparison which

760 produced the highest adjusted rand index. Across both increasing levels of drop out and increasing amounts of  
761 variation, CATCH performed better than Louvain, Leiden, and Seurat's Shared Nearest Neighbors clustering  
762 algorithms across 10 different simulations. As noise increased to 0.7 and 0.9 drop out and 0.3 and 0.4 variation,  
763 CATCH outperformed other approaches in a statistically significant manner (two-sided t-test between CATCH  
764 and each of the other clustering approaches at each iteration, p-value<.01) ([Extended Data Fig. 1E](#)).

765 Next, we compared CATCH against Louvain and Leiden clustering approaches on real single cell data where  
766 multigranular clusters had been identified by an biological expert [80, 81]. First, we analyzed real single cell  
767 transcriptomic data generated from a developing zebrafish with known cell type cluster ground truths [80].  
768 We organized these cluster labels into multigranular cluster labels by first aggregating 18 cell types found in  
769 four tissue types before aggregating them into three germ layers. In this manner, we produced ground truth  
770 cluster labels across granularities. We then compared the top four most persistent CATCH granularities against  
771 multigranular clusters computed using Louvain and Leiden, again tuning the resolution parameter to produce  
772 ten different cluster labels. At all granularities of ground truth cluster labels, CATCH out performed Louvain  
773 and Leiden despite more granularities being computed for the comparison approaches ([Extended Data Fig. 3B](#)).

774 Finally, as flow cytometry gating analysis has long been held as the gold standard for cell type identification  
775 and comparison, we compared CATCH to other clustering approaches on flow cytometry data. Using 1.3 million  
776 cells generated from 30 patients, we compared the performance of CATCH to louvain, leiden and the flow  
777 cytometry clustering gold-standard FlowSOM [81]. Across all 30 comparisons, CATCH significantly outperformed  
778 other comparisons in a statistically significant way (two-sided t-test between CATCH and each of the other  
779 clustering approaches, p-value<.01) ([Extended Data Fig.3A](#))

780 **Automated cluster characterization and Earth Mover's Distance between genes in synthetic and**  
781 **real single cell data**

782 While manifold-intrinsic diffusion condensation implemented with an  $\alpha$ -decay kernel can theoretically approximate  
783 ground truth cluster characterizations and compute differentially expressed genes, we wanted to demonstrate  
784 this reasoning in synthetic and real single cell data. To empirically show that our condensation-based approach  
785 approximates EMD between two clusters, we compute EMD values between genes using Wasserstein optimal  
786 transport as well as our approximate approach on synthetic and real data using Gaussian and  $\alpha$ -decay kernel  
787 implementations of diffusion condensation. Using single cell data generated from splatter, we compute diffusion  
788 condensation and identified the granularity with the highest topological persistence using topological activity  
789 analysis. We then computed ground truth and approximate differential expression values by comparing every

790 cluster at this granularity with every other cluster. In our analysis, a total of 12,130,200 and 4,535,640 gene  
791 comparisons were computed using Gaussian and  $\alpha$ -decay approaches respectively. Comparing both Gaussian and  
792  $\alpha$ -decay approximate Wasserstein distances against ground truth per gene Wasserstein values, we can see the  
793 value in our  $\alpha$ -decay approach ([Extended Data Fig.2D](#)) as it approximates ground truth Wasserstein distance  
794 with a correlation coefficient of .979. Furthermore, our approach computed all 4,535,640 gene comparison in  
795 63 seconds while ground truth values were computed in 43,125 seconds, equating to a 684 fold increase in  
796 computational speed.

797 We repeated our comparison in real single cell data, again comparing both approaches to ground truth  
798 Wasserstein EMD values, this time across 10 granularities identified by topological activity analysis. As previously  
799 performed, at each granularity, all clusters were compared to all other clusters using each approach. Across  
800 all comparisons, a total of 10,166,640 and 2,541,660 comparisons were computed for the Gaussian and  $\alpha$ -decay  
801 implementations respectively. Again we see that  $\alpha$ -decay is critical to accurately capturing ground truth EMD  
802 values, with our  $\alpha$ -decay approach correlating highly with ground truth EMD while Gaussian approach was less  
803 correlated ([Extended Data Fig.1F](#)). Furthermore, we again see an increase in computational speed with our  
804 condensation based approach. In our weighted implementation, we are able to compute all 2,541,660 comparisons  
805 in 32 seconds, while ground truth EMD values were computed in 27,517 seconds, equating to a similar 860  
806 fold increase in computational speed. Next, we show that this correlation between ground truth EMD and  
807 condensation-based Wasserstein distance approximation is not a feature of cluster granularity as defined by  
808 number of cluster ([Extended Data Fig.3D](#)). Finally, we also use  $\alpha$ -decay and Gaussian implementations to  
809 compute and compare cluster characterizations to ground truth in real single cell data. Using the same set of  
810 clusters and granularities as previously computed, we see that  $\alpha$ -decay kernel again more accurately characterizes  
811 clusters than a Gaussian kernel ([Extended Data Fig.3C](#)).

812 **CATCH identifies ground truth differentially expressed genes from noisy single cell data**

813 Previously, disease signatures within a cell type have been determined by comparing cells' gene expression  
814 profiles based on their condition of origin. For instance, microglia would be separated into two groups based on  
815 condition of origin, either disease or healthy, which would then be compared. We believe that CATCH improves  
816 on this framework by first identifying disease-enriched states and then identifying differentially expressed genes  
817 between these states. This is because our procedure accounts for significant noise that can appear in single cell  
818 data to more purely identify cell states enriched in particular disease settings. In fact, previous studies have

819 validated that this approach identifies biological processes better than previous 'condition-of-origin' comparison  
820 approaches [9].

821 In order to highlight this point experimentally, we generated increasingly noisy synthetic single cell data  
822 with two clusters using splatter as described previously (Extended Data Fig.??A). By computing differentially  
823 expressed genes using Earth Mover's Distance using either cluster label or CATCH's pipeline, we tried to  
824 determine if we could recover ground truth differential expression values between clusters and recover ground  
825 truth differentially expressed genes as computed by Wasserstein Earth Mover's Distance. Unsurprisingly, both  
826 strategies (CATCH identified clusters based on MELD enrichment and cluster condition of origin) perfectly  
827 correlated with ground truth Earth Mover's Distance values and identified all differentially expressed genes with  
828 no false positives. Real single cell data, however, has high levels of noise with different cellular states often  
829 blending into one another on the manifold. When we increased noise and these clusters began to blend together,  
830 however, CATCH's pipeline not only better correlated with ground truth Earth Mover's Distance values, but  
831 also recovered more ground truth differentially expressed genes (Extended Data Fig.??B).

832 To illustrate this point in real single cell data, we performed differential expression analysis between microglia  
833 based on their condition of origin across all three neurodegenerative disease datasets. We reason that if our  
834 approach is more sensitive to identify differentially expressed genes, a less sensitive approach would not find  
835 as strong of a shared signature. After setting significance cutoffs based on our per gene false discovery rates,  
836 we identified significantly enriched genes in the early or acute active phase of each disease ([Extended Data Fig.](#)  
837 [9a](#)). However, across all comparisons, we identified significantly fewer differentially expressed genes in this cell  
838 type analysis (135, 68 and 416) than with our pipeline (618, 795 and 1551 for AMD, AD and MS respectively),  
839 indicating that the identification of pathogenic cellular subtypes with CATCH before comparison increases our  
840 ability to detect differentially expressed genes. In cross-disease comparisons among early stage neurodegenerative  
841 microglia, only 17 common genes were found, significantly less than the 168 common genes found with our  
842 pipeline. Of the common genes, only half of the activation signature was found (*APOE*, *B2M*, *FTH1*, *FTL*,  
843 *SPP1*).

844 Similar to our coarse grained microglial comparison, we compared the strength of our approach in astrocytes.  
845 After setting significance cutoffs based on our per gene q-values, we identified significantly fewer enriched genes  
846 (221, 271, and 886) than we found with our analysis (1444, 680, and 2278 genes for AMD, AD and MS respectively)  
847 ([Extended Data Fig. 9b](#)). In our cell type level analysis, only 28 common genes were found, significantly less  
848 than the 630 common genes found with our pipeline. Of the common genes, only half of the activation signature

849 was found (*AQP4*, *CD81*, *CRYAB*, *GFAP*).

850 Collectively, these comparisons reveal the sensitivity of this discovery pipeline for finding gene signatures and  
851 biologically meaningful relationships in noisy single cell gene expression data.

852 **Data and Software Availability:**

853 Single nucleus RNA sequence data will be deposited in the publicly available GEO data repository, and will be  
854 available from the authors upon request. The CATCH package, as implemented in python, is available for download  
855 with a guided tutorial on the Krishnaswamy Lab Github page: <https://github.com/KrishnaswamyLab/CATCH>.

856 **Other Computational Methods Details**

857 **Single-nucleus AMD RNA sequencing and pre-processing**

858 snRNA-seq data from macular samples, were processed according to the following steps. Sample demultiplexing  
859 and read alignment to the NCBI reference pre-mRNA GRCh38 was completed to map reads to both unspliced  
860 pre-mRNA and mature mRNA transcripts using CellRanger version 3.1.0. Gene and cell matrices from retinas  
861 with dry AMD (n=3), neovascular AMD (n=8), and healthy controls (n=6) were then combined into a single  
862 file. We prefiltered using parameters in scprep (v1.0.3, <https://github.com/KrishnaswamyLab/scprep>). Cells  
863 that contained at least 1400 unique transcripts were kept for further analysis to generate a cell by gene matrix  
864 containing 71,063 cells. Normalization was performed using default parameters with L1 normalization, adjusting  
865 total library size of each cell to 1000. Any cell with greater than 200 normalized counts of mitochondrial mRNA  
866 was removed. Batch correction was performed using Harmony (<https://github.com/immunogenomics/harmony>)  
867 to align batch effects introduced by sequencing batch, postmortem interval, sample acquisition location and 10X  
868 sequencing chemistry [82]. Raw data files for human snRNA-seq data will be available for download through  
869 GEO under an accession number to be assigned with no restrictions on data availability.

870 **Single-nucleus AD and MS RNA sequencing pre-processing**

871 snRNA-seq data for AD and MS was acquired from published sources [4,5]. Cells that contained at least 1000  
872 unique transcripts were kept for further analysis to generate a cell by gene matrix for each disease. Normalization  
873 was performed using scprep default parameters with L1 normalization, adjusting total library size of each cell to  
874 1000. Any cell with greater than 200 normalized counts of mitochondrial mRNA was removed. Batch correction  
875 was performed on MS data using Harmony (<https://github.com/immunogenomics/harmony>) to align batch  
876 effects introduced by sequencing batch, capture batch and sex.

877 **Cell type identification with CATCH**

878 All cell types were identified by performing topological activity analysis on the diffusion condensation calculated  
879 condensation homology. In order to identify cell types, we identified a resolutions with no topological activity  
880 which partitioned the cellular state space well and assigned each cluster to a cell type based on cell type specific  
881 marker genes.

882 **Interaction Analysis:**

883 Cell-cell ligand-receptor analysis was conducted on pre-processed snRNA expression data using the CellPhoneDB  
884 python package (<https://github.com/Teichlab/cellphonedb>, v2.1.4) [47]. Before conducting analysis, the package  
885 database of 834 curated ligand-receptor combinations and multi-unit protein complexes was supplemented with  
886 2557 ligand-receptor interactions found in the celltalker database (<https://github.com/arc85/celltalker>) [83].  
887 The in-built database-generate function was utilized to update the existing database. Our comprehensive  
888 user-generated database was invoked in each run of the CellPhoneDB statistical-analysis command function.

889 CellPhoneDB interaction maps were computed on differing inputs. First, disease phase enriched microglia and  
890 astrocytes with subcluster identity were run to identify signaling interactions between astrocyte and microglial  
891 activation states (Fig. 6B). The number of permutations was set to 2,000 and p-value threshold was set to 0.01.

892 **Biological Methods Details**

893 **Ethics Statement**

894 This study was approved by the Yale Institutional Review Board. We complied with all relevant ethical regulations  
895 for work with human participants, and all human tissue samples were obtained with informed consent prior to  
896 tissue collection.

897 **Human tissues**

898 Postmortem eyes for the Chromium Single Cell 3' assay ( $n = 17$ ) and medical records containing AMD disease  
899 stage were obtained from Advancing Sight Network (Alabama), Lions Gift of Sight Eye Bank (Minnesota), or  
900 the Yale Department of Pathology with a maximum post-mortem interval of 13 hours. Globes were examined for  
901 retinal disease by an ophthalmologist (B.P.H.) prior to dissection and dissociation of the samples. Retina for  
902 snRNA-seq was obtained from the unrelated human post-mortem donors that included normal, intermediate  
903 dry on AREDS2, and neovascular AMD stages (Supplementary Data Table S1). For each sample we profiled  
904 the macula, which is the region of the retina responsible for central vision and affected most severely by AMD

905 pathology. We identified three intermediate AMD samples from patients taking the AREDS2 eye vitamin and  
906 mineral supplement with drusen, a pathologic sign associated with the intermediate dry stage of the disease.  
907 Eight postmortem AMD samples had neovascularization in the advanced stage of the disease. Normal donors had  
908 no history of retinal disease. Additional clinical data for the subjects is given in Supplementary Data Table S2.

909 **Retinal dissection and solation of nuclei from frozen retinal tissue**

910 Globes were placed in RNAlater (ThermoFisher) and transported on ice. Trehpene punches (6 mm diameter)  
911 were used to isolate samples from the macula in the central retina, located away from the optic disc and major  
912 arterioles. For each punch of tissue, the retina was mechanically separated from the underlying retinal pigment  
913 epithelium-choroid, snap-frozen on dry ice and stored at -80°C. Nuclei were isolated and purified using the  
914 Nuclei EZ Prep Nuclei Isolation Kit (Sigma), following the manufacturer's protocol, with some modification. All  
915 procedures were carried out on ice or at 4°C. Briefly, frozen retinal tissue was subjected to dounce homogenization  
916 (25 times with pestle A followed by 25 times with tight pestle B) using the KIMBLE Dounce Tissue Grinder Set  
917 (Sigma) in 2ml EZ Lysis buffer. The sample was transferred to a 15ml tube with an additional 2ml EZ lysis  
918 buffer and incubated on ice for 5min. Following incubation, the sample was centrifuged at 500xg, 5min at 4°C.  
919 Supernatants were discarded, and the isolated nuclei were resuspended in 4ml EZ lysis buffer, incubated for  
920 5min on ice and centrifuged at 500xg for 5min at 4°C. Next, the nuclei were washed with 4ml ice-cold Nuclei  
921 Suspension Buffer (1X PBS containing 0.01% BSA and 0.1% RNase inhibitor), resuspended in 1ml Nuclei EZ  
922 Storage buffer and passed through a 40µM nylon cell strainer. The nuclei suspensions were counted with trypan  
923 blue prior to loading on the microfluidics platform.

924 **Droplet-based microfluids snRNA-seq**

925 Isolated nuclei from each macular sample was processed through microfluidics-based single nuclear RNA-seq.  
926 Single-cell libraries were prepared using the Chromium 3' v2 and v3 platforms (10X Genomics) following the  
927 manufacturer's protocol. Briefly, single nuclei were partitioned into Gel beads in Emulsion in the 10X Chromium  
928 Controller instrument followed by lysis and barcoded reverse transcription of RNA, amplification, shearing and 5'  
929 adapter and sample index attachment. On average, 7000 nuclei were loaded on each channel that resulted in the  
930 recovery of 4000 nuclei. Libraries were sequenced on the Illumina NextSeq 500 platform. After quality control  
931 preprocessing, snRNA-seq profiles were used in subsequent analyses. This dataset was corrected for batch effects  
932 across samples using the Harmony algorithm [82].

933 **In situ RNA hybridization and immunofluorescence**

934 To validate the gene expression differences, in situ hybridization was performed using RNAscope Multiplex  
935 Fluorescent V2 Assay (Advanced Cell Diagnostics, Hayward, CA, USA). Macula dissected from whole human  
936 globes were fixed in 4% paraformaldehyde (PFA) at 4°C overnight. Tissues were sequentially dehydrated with  
937 15% sucrose, then 30% sucrose before embedding in OCT, and frozen on dry ice. OCT molds were sectioned at 10  
938 µm thickness. RNA in situ hybridization was performed according to the manufacturer's protocol. Briefly, fixed  
939 frozen sections were baked at 60°C for 1 hr prior to incubation in 4% PFA for 10 mins and protease digestion  
940 pretreatment. Target probes were hybridized to an HRP-based temperature sensitive signal amplification system,  
941 followed by color development. Housekeeping genes *POLR2A*, *PPIB*, and *UBC* were used as internal-control  
942 mRNA ([Extended Data Fig. 7](#)); if probes for these mRNAs were not visualized, the sample was regarded  
943 as not available for gene expression study. The probes used include *APOE*, *TYROBP*, *B2M*, *VEGFA*, and  
944 *HIF1A* (Advanced Cell Diagnostics, Hayward, CA, USA). The slides were counterstained with DAPI during  
945 immunofluorescence protocol (see below). Positive staining was determined by fluorescent punctate dots in the  
946 appropriate channels in the nucleus and/or cytoplasm. Following RNA in situ hybridization protocol, fixed  
947 frozen sections were blocked with animal serum and incubated overnight at 4°C with primary antibodies (see  
948 antibody segment below). Secondary antibody incubation was for 1 hr at room temperature and cell nuclei were  
949 counterstained with DAPI. Images were captured immediately using a confocal microscope (Zeiss LSM800, Jena,  
950 Germany). The following antibodies against human antigens were used: GFAP (1:500, MA5-12023, Invitrogen)  
951 and Iba1 (1:500, 019-19741, Fujifilm). Antibodies were visualized with Alexa Fluor 488 (1:200, A-11001 /  
952 A-21208, Invitrogen).

953 **Mice**

954 Four- to eight-week-old mixed sex C57BL/6 mice were purchased from the National Cancer Institute and  
955 subsequently bred and housed at Yale University. All procedures used in this study (sex-matched, age-matched)  
956 complied with federal guidelines and the institutional policies of the Yale School of Medicine Animal Care and  
957 Use Committee.

958 **Cells**

959 iPSC-derived astrocyte cells were purchased from Brainxell.com (Brainxell, Madison Wisconsin). Cells were  
960 cultured according to provider's guidelines using 1:1 DMEM/F12 and Neurobasal medium with N2 supplement

961 (1x), Glutamax (0.5mM), Astrocyte supplement (1x), Fetal bovine serum (1%).

## 962 Cell culture

963 IPSC-derived astrocyte cells were cultured to a fully differentiated state before cytokine stimulation. Cytokines,  
964 (IL-1 $\beta$ , IL2, IL4, IL6, IL7, IL10, IL12, IL15, IL17, IL22, IL23, IFNG, TNF) were all purchased from PeproTech.com  
965 (Peprotech, Cranbury, NJ). For single cytokine stimulation, cells were stimulated with each cytokine at a  
966 concentration of 100ng/mL for 24 hours. For combinatorial cytokine stimulation, cocktail of all cytokines minus  
967 cytokine of interest was made with each cytokine concentration at 50ng/mL. Cells were stimulated for 24 hours  
968 before media was collected. Collected media was centrifuged at 1000xg to remove any cells and debris before  
969 performing an ELISA.

## 970 Enzyme-linked immunosorbent assay

971 Enzyme-linked immunosorbent assay (ELISA) was performed using a mouse VEGF-A ELISA Kit (Cusabio LLC)  
972 following the manufacturer's instructions.

## 973 Intravitreal injection

974 Mice were anaesthetized using a mixture of ketamine (50mg/kg) and xylazine (5mg/kg), injected intraperitoneally.  
975 Mice eyes were sterilized using betadine. A small hole was made at the lateral aspect of the limbus was made  
976 using a 33 gauge insulin syringe. Using a blunt end Hamilton syringe, 1  $\mu$ l of PBS or IL-1 $\beta$  (100ng) was injected  
977 at a 45 degree angle at the limbus intravitreally. Once the infusion was finished, syringe was left in place for a  
978 minute before removal of the syringe. Injection site was washed with sterile PBS and puralube vet ointment was  
979 applied to the eyes. Mice were monitored until full recovery.

## 980 Mice tissue processing and microscopy

981 Retinas were dissected, fixed in 2% PFA for one hour and immediately processed in a blocking solution (10%  
982 normal donkey serum, 1% bovine serum albumin, 0.3% PBS-Triton X-100) for overnight incubation at 4°C. For  
983 retina sections, primary antibodies were incubated overnight at 4°C, then washed five times at room temperature  
984 in PBS and 0.5% Triton X-100, before incubation with a fluoro-conjugated secondary antibody diluted in PBS  
985 and 0.5% Triton X-100 for 2 hours in room temperature. Sections were washed five times at room temperature,  
986 stained with DAPI and mounted before imaging. Confocal images were taken on a Leica SP8 microscope.  
987 Quantitative analysis was performed using either FIJI or ImageJ image-processing software (NIH or Bethesda)  
988 or Imaris 8 software (Oxford Instruments).

989 **References**

- 990 [1] Wong, W. L. *et al.* Global prevalence of age-related macular degeneration and disease burden projection for  
991 2020 and 2040: a systematic review and meta-analysis. *The Lancet Global Health* **2**, e106–e116 (2014).
- 992 [2] Mitchell, P., Liew, G., Gopinath, B. & Wong, T. Y. Age-related macular degeneration. *Lancet* **392**,  
993 1147–1159 (2018).
- 994 [3] Bird, A. C. *et al.* An international classification and grading system for age-related maculopathy and  
995 age-related macular degeneration. The International ARM Epidemiological Study Group. *Surv Ophthalmol*  
996 **39**, 367–374 (1995).
- 997 [4] Mathys, H. *et al.* Single-cell transcriptomic analysis of alzheimer's disease. *Nature* **570**, 332–337 (2019).
- 998 [5] Schirmer, L. *et al.* Neuronal vulnerability and multilineage diversity in multiple sclerosis. *Nature* **573**, 75–82  
999 (2019).
- 1000 [6] Habib, N. *et al.* Disease-associated astrocytes in Alzheimer's disease and aging. *Nat Neurosci* **23**, 701–706  
1001 (2020).
- 1002 [7] Keren-Shaul, H. *et al.* A Unique Microglia Type Associated with Restricting Development of Alzheimer's  
1003 Disease. *Cell* **169**, 1276–1290 (2017).
- 1004 [8] Brugnone, N. *et al.* Coarse graining of data via inhomogeneous diffusion condensation. In *2019 IEEE*  
1005 *International Conference on Big Data (Big Data)*, 2624–2633 (IEEE, 2019).
- 1006 [9] Burkhardt, D. B. *et al.* Quantifying the effect of experimental perturbations in single-cell rna-sequencing  
1007 data using graph signal processing. *bioRxiv* (2020).
- 1008 [10] Lemière, S. NLRP3 inflammasome activity as biomarker for primary progressive multiple sclerosis. *Nature*  
1009 *Reviews Neurology* **16**, 350–350 (2020). URL <https://doi.org/10.1038/s41582-020-0366-y>.
- 1010 [11] Zhang, Y., Dong, Z. & Song, W. NLRP3 inflammasome as a novel therapeutic target for alzheimer's disease.  
1011 *Signal Transduction and Targeted Therapy* **5** (2020). URL <https://doi.org/10.1038/s41392-020-0145-7>.
- 1012 [12] White, C. S., Lawrence, C. B., Brough, D. & Rivers-Auty, J. Inflammasomes as therapeutic targets for  
1013 alzheimer's disease. *Brain Pathology* **27**, 223–234 (2017). URL <https://doi.org/10.1111/bpa.12478>.

- 1014 [13] Faissner, S., Plemel, J. R., Gold, R. & Yong, V. W. Progressive multiple sclerosis: from pathophysiology to  
1015 therapeutic strategies. *Nat Rev Drug Discov* **18**, 905–922 (2019).
- 1016 [14] Huang, W.-J., Chen, W.-W. & Zhang, X. Multiple sclerosis: Pathology, diagnosis and treatments. *Exp.*  
1017 *Ther. Med.* **13**, 3163–3166 (2017).
- 1018 [15] Braak, H. & Braak, E. Neuropathological staging of Alzheimer-related changes. *Acta Neuropathol* **82**,  
1019 239–259 (1991).
- 1020 [16] Ding, J. *et al.* Systematic comparison of single-cell and single-nucleus RNA-sequencing methods. *Nat.*  
1021 *Biotechnol.* **38**, 737–746 (2020).
- 1022 [17] Huguet, G. *et al.* Time-inhomogeneous diffusion geometry and topology. In "" ( "", New Haven, Connecticut,  
1023 2022).
- 1024 [18] Moyle, M. W. *et al.* Structural and developmental principles of neuropil assembly in *c. elegans*. *Nature* **591**,  
1025 99–104 (2021). URL <https://doi.org/10.1038/s41586-020-03169-5>.
- 1026 [19] Menon, M. *et al.* Single-cell transcriptomic atlas of the human retina identifies cell types associated with  
1027 age-related macular degeneration. *Nat. Commun.* **10**, 4902 (2019).
- 1028 [20] Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities in large  
1029 networks. *Journal of Statistical Mechanics: Theory and Experiment* **2008**, P10008 (2008).
- 1030 [21] Shekhar, K. *et al.* Comprehensive classification of retinal bipolar neurons by (single-cell) transcriptomics.  
1031 *Cell* **166**, 1308–1323.e30 (2016).
- 1032 [22] Peng, Y.-R. *et al.* Molecular classification and comparative taxonomies of foveal and peripheral cells in  
1033 primate retina. *Cell* **176**, 1222–1237.e22 (2019).
- 1034 [23] Yan, W. *et al.* Cell atlas of the human fovea and peripheral retina. *Sci. Rep.* **10**, 9802 (2020).
- 1035 [24] van Dijk, D. *et al.* Recovering gene interactions from single-cell data using data diffusion. *Cell* **174**, 716 –  
1036 729.e27 (2018).
- 1037 [25] Srinivasan, K. *et al.* Alzheimer's patient microglia exhibit enhanced aging and unique transcriptional  
1038 activation. *Cell Reports* **31**, 107843 (2020). URL <https://doi.org/10.1016/j.celrep.2020.107843>.

- 1039 [26] Friedman, B. A. *et al.* Diverse Brain Myeloid Expression Profiles Reveal Distinct Microglial Activation  
1040 States and Aspects of Alzheimer's Disease Not Evident in Mouse Models. *Cell Rep* **22**, 832–847 (2018).
- 1041 [27] Krasemann, S. *et al.* The TREM2-APOE Pathway Drives the Transcriptional Phenotype of Dysfunctional  
1042 Microglia in Neurodegenerative Diseases. *Immunity* **47**, 566–581 (2017).
- 1043 [28] Corder, E. H. *et al.* Gene dose of apolipoprotein E type 4 allele and the risk of Alzheimer's disease in late  
1044 onset families. *Science* **261**, 921–923 (1993).
- 1045 [29] Lambert, J. C. *et al.* Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer's  
1046 disease. *Nat Genet* **45**, 1452–1458 (2013).
- 1047 [30] Fritzsche, L. G. *et al.* A large genome-wide association study of age-related macular degeneration highlights  
1048 contributions of rare and common variants. *Nat Genet* **48**, 134–143 (2016).
- 1049 [31] Satoh, J. I., Kino, Y., Yanaizu, M. & Saito, Y. Alzheimer's disease pathology in Nasu-Hakola disease brains.  
1050 *Intractable Rare Dis Res* **7**, 32–36 (2018).
- 1051 [32] van der Poel, M. *et al.* Transcriptional profiling of human microglia reveals grey-white matter heterogeneity  
1052 and multiple sclerosis-associated changes. *Nature Communications* **10**, 1139 (2019).
- 1053 [33] Sala Frigerio, C. *et al.* The Major Risk Factors for Alzheimer's Disease: Age, Sex, and Genes Modulate the  
1054 Microglia Response to A $\beta$  Plaques. *Cell Rep* **27**, 1293–1306 (2019).
- 1055 [34] Giovannoni, F. & Quintana, F. J. The Role of Astrocytes in CNS Inflammation. *Trends Immunol* **41**,  
1056 805–819 (2020).
- 1057 [35] Zamanian, J. L. *et al.* Genomic analysis of reactive astrogliosis. *J Neurosci* **32**, 6391–6410 (2012).
- 1058 [36] Bombeiro, A. L., Hell, R. C., Simões, G. F., Castro, M. V. & Oliveira, A. L. Importance of major  
1059 histocompatibility complex of class I (MHC-I) expression for astroglial reactivity and stability of neural  
1060 circuits in vitro. *Neurosci Lett* **647**, 97–103 (2017).
- 1061 [37] Ransohoff, R. M. & Estes, M. L. Astrocyte expression of major histocompatibility complex gene products  
1062 in multiple sclerosis brain tissue obtained by stereotactic biopsy. *Arch Neurol* **48**, 1244–1246 (1991).
- 1063 [38] Xie, L. *et al.* Sleep drives metabolite clearance from the adult brain. *Science* **342**, 373–377 (2013).

- 1064 [39] Moon, K. R. *et al.* Visualizing structure and transitions in high-dimensional biological data. *Nature*  
1065 *Biotechnology* **37**, 1482–1492 (2019).
- 1066 [40] Latz, E., Xiao, T. S. & Stutz, A. Activation and regulation of the inflammasomes. *Nature Reviews*  
1067 *Immunology* **13**, 397–411 (2013). URL <https://doi.org/10.1038/nri3452>.
- 1068 [41] Cantuti-Castelvetri, L. *et al.* Defective cholesterol clearance limits remyelination in the aged central nervous  
1069 system. *Science* **359**, 684–688 (2018).
- 1070 [42] Shweiki, D., Itin, A., Soffer, D. & Keshet, E. Vascular endothelial growth factor induced by hypoxia may  
1071 mediate hypoxia-initiated angiogenesis. *Nature* **359**, 843–845 (1992).
- 1072 [43] Zeng, Z. J. *et al.* TLX controls angiogenesis through interaction with the von Hippel-Lindau protein. *Biol*  
1073 *Open* **1**, 527–535 (2012).
- 1074 [44] Wang, G. L., Jiang, B. H., Rue, E. A. & Semenza, G. L. Hypoxia-inducible factor 1 is a basic-helix-  
1075 loop-helix-PAS heterodimer regulated by cellular O<sub>2</sub> tension. *Proc Natl Acad Sci U S A* **92**, 5510–5514  
1076 (1995).
- 1077 [45] Kliffen, M., Sharma, H. S., Mooy, C. M., Kerkvliet, S. & de Jong, P. T. Increased expression of angiogenic  
1078 growth factors in age-related maculopathy. *Br J Ophthalmol* **81**, 154–162 (1997).
- 1079 [46] Wong, T. Y., Liew, G. & Mitchell, P. Clinical update: new treatments for age-related macular degeneration.  
1080 *Lancet* **370**, 204–206 (2007).
- 1081 [47] Efremova, M., Vento-Tormo, M., Teichmann, S. A. & Vento-Tormo, R. CellPhoneDB: inferring cell–cell  
1082 communication from combined expression of multi-subunit ligand–receptor complexes. *Nature Protocols* **15**,  
1083 1484–1506 (2020).
- 1084 [48] Escartin, C. *et al.* Reactive astrocyte nomenclature, definitions, and future directions. *Nat Neurosci* **24**,  
1085 312–325 (2021).
- 1086 [49] Guttenplan, K. A. *et al.* Neurotoxic Reactive Astrocytes Drive Neuronal Death after Retinal Injury. *Cell*  
1087 *Rep* **31**, 107776 (2020).
- 1088 [50] Liddelow, S. A. *et al.* Neurotoxic reactive astrocytes are induced by activated microglia. *Nature* **541**,  
1089 481–487 (2017).

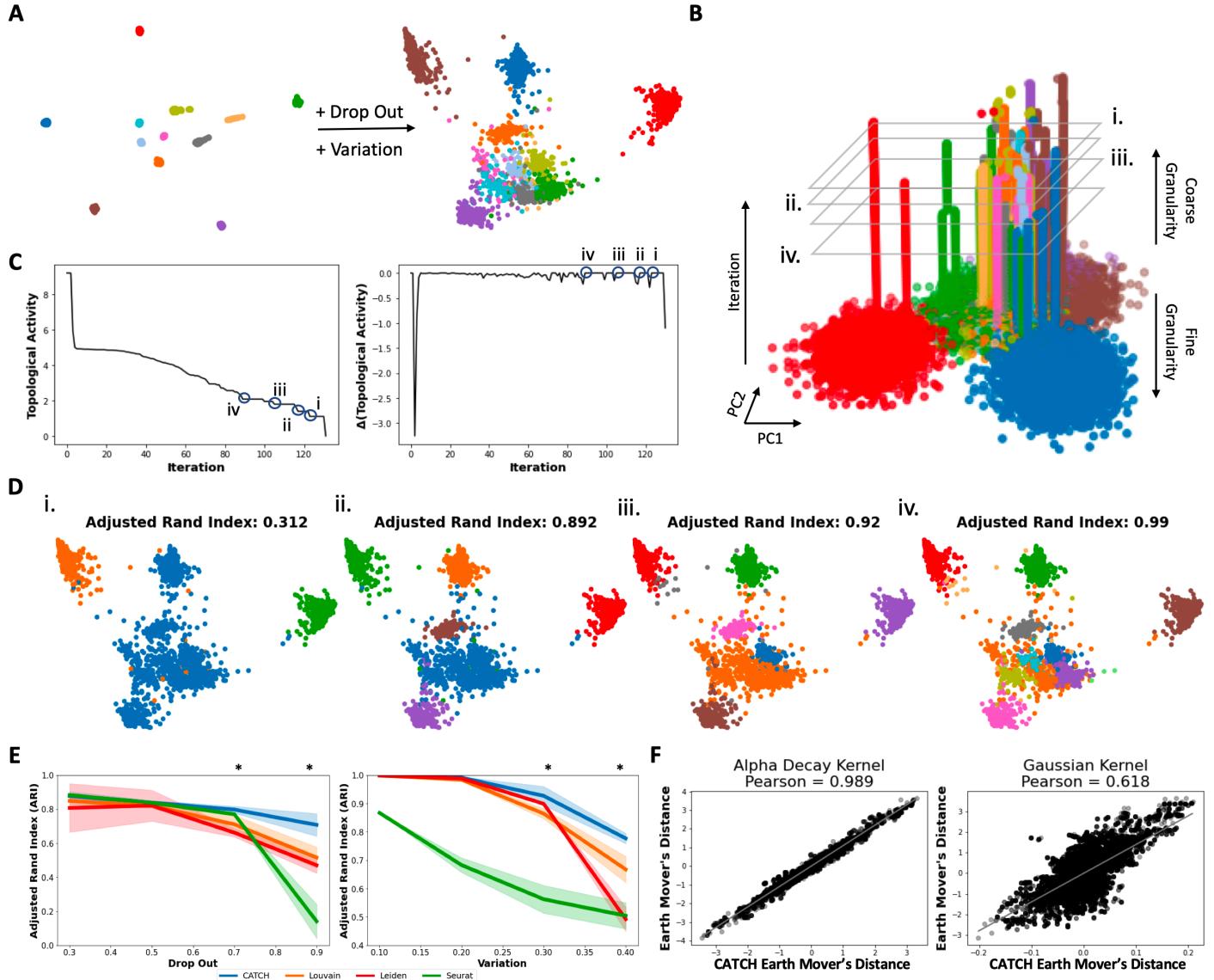
- 1090 [51] Krishnaswamy, S. *et al.* Conditional density-based analysis of t cell signaling in single-cell data. *Science*  
1091 **346**, 1250689–1250689 (2014).
- 1092 [52] Zhao, M. *et al.* Interleukin-1 $\beta$  Level Is Increased in Vitreous of Patients with Neovascular Age-Related  
1093 Macular Degeneration (nAMD) and Polypoidal Choroidal Vasculopathy (PCV). *PLoS One* **10**, e0125150  
1094 (2015).
- 1095 [53] Heneka, M. T., McManus, R. M. & Latz, E. Inflammasome signalling in brain function and neurode-  
1096 generative disease. *Nature Reviews Neuroscience* **19**, 610–621 (2018). URL <https://doi.org/10.1038/s41583-018-0055-7>.
- 1097 [54] Guillonneau, X. *et al.* On phagocytes and macular degeneration. *Prog Retin Eye Res* **61**, 98–128 (2017).
- 1098 [55] Moon, K. R. *et al.* Manifold learning-based methods for analyzing single-cell rna-sequencing data. *Current  
1099 Opinion in Systems Biology* **7**, 36–46 (2018).
- 1100 [56] Coifman, R. R. & Lafon, S. Diffusion maps. *Applied and computational harmonic analysis* **21**, 5–30 (2006).
- 1101 [57] Van Der Maaten, L., Postma, E. & Van den Herik, J. Dimensionality reduction: a comparative. *J Mach  
1102 Learn Res* **10**, 66–71 (2009).
- 1103 [58] Izenman, A. J. Introduction to manifold learning. *Wiley Interdisciplinary Reviews: Computational Statistics*  
1104 **4**, 439–446 (2012).
- 1105 [59] Lindenbaum, O., Stanley, J., Wolf, G. & Krishnaswamy, S. Geometry based data generation. In *Advances  
1106 in Neural Information Processing Systems*, 1400–1411 (MIT Press, 2018).
- 1107 [60] Gama, F., Ribeiro, A. & Bruna, J. Diffusion scattering transforms on graphs. In *International Conference  
1108 on Learning Representations* (ICLR, 2019). ArXiv:1806.08829.
- 1109 [61] Gao, F., Wolf, G. & Hirn, M. Geometric scattering for graph data analysis (2019). To appear in the  
1110 *Proceedings of the 36th International Conference on Machine Learning*, arXiv:1810.03068.
- 1111 [62] Gigante, S. *et al.* Compressed diffusion. In *2019 13th International conference on Sampling Theory and  
1112 Applications (SampTA)* (IEEE, 2019). URL <https://doi.org/10.1109/sampta45681.2019.9030994>.

- 1114 [63] Batson, J., Royer, L. & Webber, J. Molecular cross-validation for single-cell rna-seq. *bioRxiv* (2019). URL  
1115 <https://www.biorxiv.org/content/early/2019/09/30/786269>. <https://www.biorxiv.org/content/early/2019/09/30/786269.full.pdf>.
- 1117 [64] Zappia, L., Phipson, B. & Oshlack, A. Splatter: simulation of single-cell RNA sequencing data. *Genome*  
1118 *Biology* **18** (2017). URL <https://doi.org/10.1186/s13059-017-1305-0>.
- 1119 [65] Chen, C. & Edelsbrunner, H. Diffusion runs low on persistence fast. In *Proceedings of the IEEE International*  
1120 *Conference on Computer Vision* (ICCV), 423–430 (Curran Associates, Inc., Red Hook, NY, USA, 2011).
- 1121 [66] Ghrist, R. Barcodes: The persistent topology of data. *Bulletin of the American Mathematical Society* **45**,  
1122 61–75 (2008).
- 1123 [67] Rieck, B., Sadlo, F. & Leitte, H. Topological machine learning with persistence indicator functions. In Carr,  
1124 H., Fujishiro, I., Sadlo, F. & Takahashi, S. (eds.) *Topological Methods in Data Analysis and Visualization V*,  
1125 87–101 (Springer, Cham, Switzerland, 2020). [1907.13496](https://arxiv.org/abs/1907.13496).
- 1126 [68] O’Bray, L., Rieck, B. & Borgwardt, K. Filtration curves for graph representation. In *Proceedings of the*  
1127 *27th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, 1267–1275  
1128 (Association for Computing Machinery, New York, NY, USA, 2021).
- 1129 [69] Dann, E., Henderson, N. C., Teichmann, S. A., Morgan, M. D. & Marioni, J. C. Differential abundance  
1130 testing on single-cell data using k-nearest neighbor graphs. *Nature Biotechnology* (2021). URL <https://doi.org/10.1038/s41587-021-01033-z>.
- 1132 [70] Nabavi, S., Schmolze, D., Maitituoheti, M., Malladi, S. & Beck, A. H. EMDomics: a robust and powerful  
1133 method for the identification of genes differentially expressed between heterogeneous classes. *Bioinformatics*  
1134 **32**, 533–541 (2015). URL <https://doi.org/10.1093/bioinformatics/btv634>.
- 1135 [71] Wang, T. & Nabavi, S. Differential gene expression analysis in single-cell rna sequencing data. In *2017*  
1136 *IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 202–207 (2017).
- 1137 [72] Orlova, D. Y. *et al.* Earth Mover’s Distance (EMD): A True Metric for Comparing Biomarker Expression  
1138 Levels in Cell Populations. *PLOS ONE* **11**, e0151859 (2016).
- 1139 [73] Backurs, A., Dong, Y., Indyk, P., Razenshteyn, I. & Wagner, T. Scalable nearest neighbor search for optimal  
1140 transport (2020). [1910.04126](https://arxiv.org/abs/1910.04126).

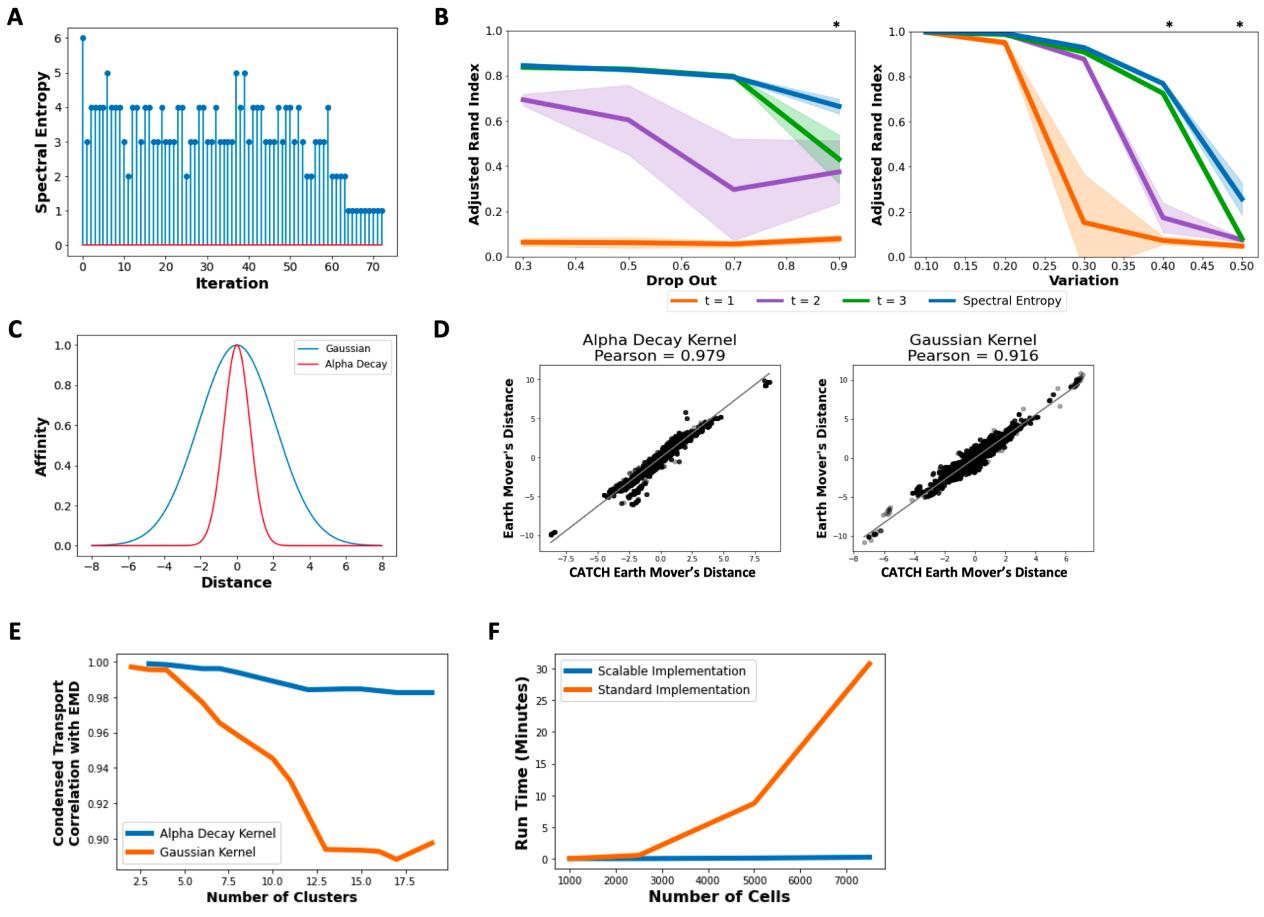
- 1141 [74] Indyk, P. & Thaper, N. Fast image retrieval via embeddings. In *3rd International Workshop on Statistical*  
1142 *and Computational Theories of Vision* (IEEE Computer Society Press, 2003).
- 1143 [75] Le, T., Yamada, M., Fukumizu, K. & Cuturi, M. Tree-sliced variants of wasserstein distances. In *Advances*  
1144 *in neural information processing systems*, 12304–12315 (Neural Information Processing Systems Foundation,  
1145 2019).
- 1146 [76] Peyré, G. & Cuturi, M. Computational Optimal Transport. *arXiv:1803.00567 [stat]* (2019). [1803.00567](https://arxiv.org/abs/1803.00567).
- 1147 [77] Gonzalez, T. F. Clustering to minimize the maximum intercluster distance. *Theoretical Computer Science*  
1148 **38**, 293–306 (1985). URL [https://doi.org/10.1016/0304-3975\(85\)90224-5](https://doi.org/10.1016/0304-3975(85)90224-5).
- 1149 [78] Storey, J. D. A direct approach to false discovery rates. *Journal of the Royal Statistical Society: Series B*  
1150 *(Statistical Methodology)* **64**, 479–498 (2002). URL <https://doi.org/10.1111/1467-9868.00346>.
- 1151 [79] Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: A practical and powerful approach to  
1152 multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)* **57**, 289–300 (1995).  
1153 URL <http://www.jstor.org/stable/2346101>.
- 1154 [80] Wagner, D. E. *et al.* Single-cell mapping of gene expression landscapes and lineage in the zebrafish embryo.  
1155 *Science* **360**, 981–987 (2018). <https://science.sciencemag.org/content/360/6392/981.full.pdf>.
- 1156 [81] Aghaeepour, N. *et al.* Critical assessment of automated flow cytometry data analysis techniques. *Nature*  
1157 *methods* **10**, 228–238 (2013).
- 1158 [82] Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with harmony. *Nature Methods*  
1159 **16**, 1289–1296 (2019).
- 1160 [83] Ramilowski, J. A. *et al.* A draft network of ligand–receptor-mediated multicellular signalling in human.  
1161 *Nature Communications* **6** (2015).

1162 **Extended data**

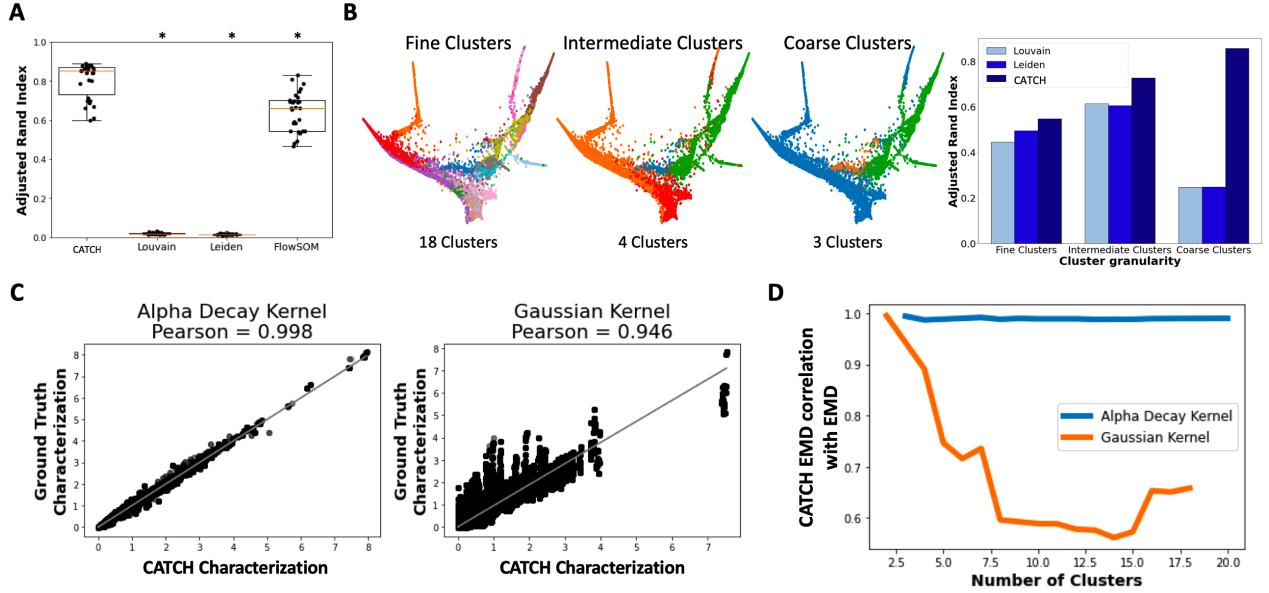
1163 **Extended Data Figures**



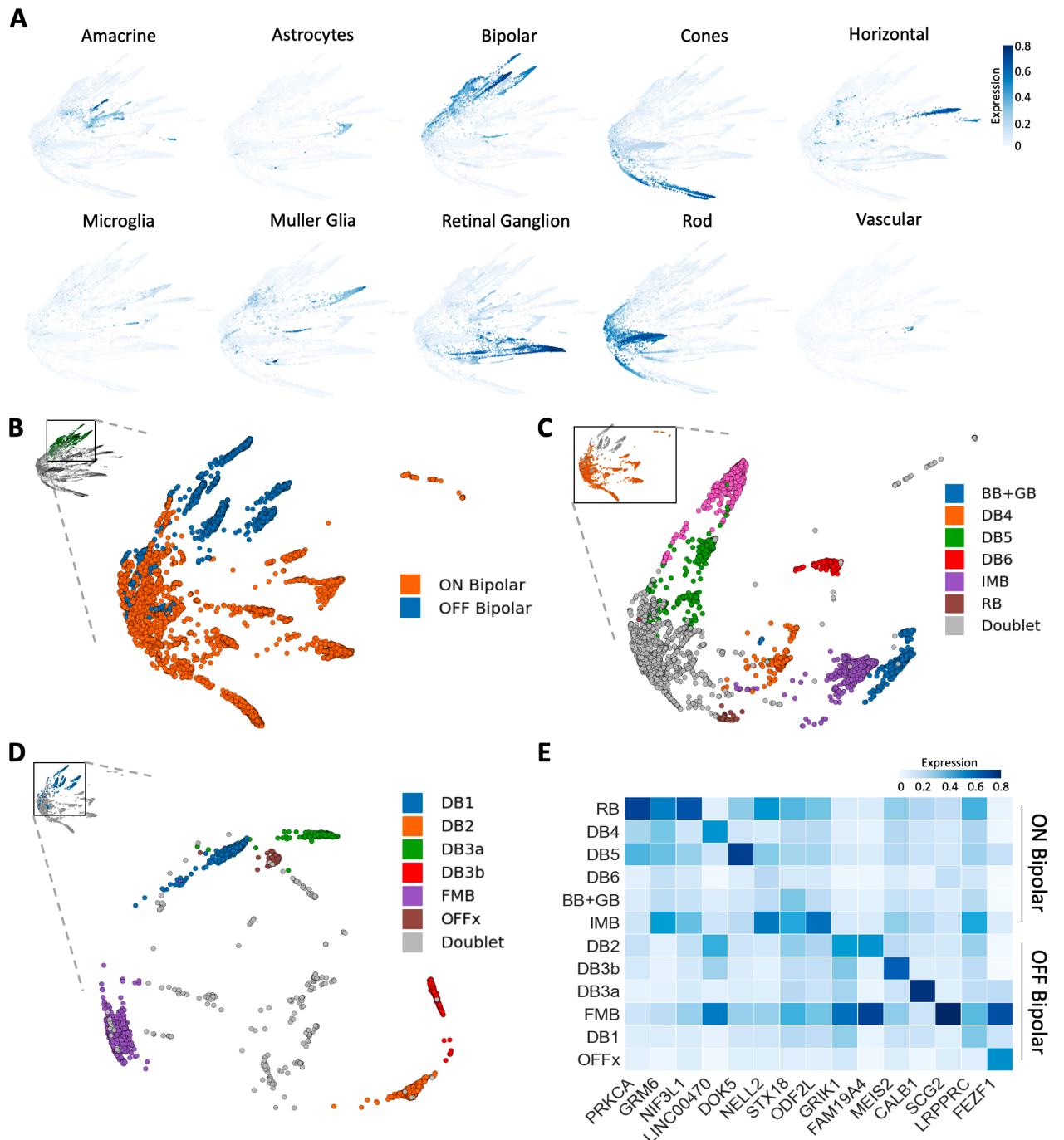
**Extended Data Fig. 1: Diffusion condensation identifies and characterizes populations of related cells across granularities.** (A) PHATE visualization of 10,000 cells generated from splatter [64]. Noiseless data is first generated on which ground truth clusters are computed (left). Two types of biological noise, variation and drop out, are simulated and the dataset is visualized with ground truth cluster labels highlighted (right). (B) condensation homology visualization of noisy splatter single cell data simulation. Four granularities are highlighted (represented as i.-iv.), illustrating 4 resolutions identified in (C) as meaningful. (C) Topological activity is first computed (left) before gradient analysis is performed on the topological activity curve (right). Resolutions identified by this analysis as being meaningful are highlighted (represented as i.-iv.). (D) Meaningful resolutions identified in (C) are represented with Adjusted Rand Index score when compared to ground truth cluster labels shown. Visualizations are arranged from the coarsest granularity of clusters (left) to the finest granularity (right). (E) Forty different splatter synthetic single cell datasets are simulated with either increasing amounts of drop out (left panel) or biological variation (right panel). Cluster labels are computed using a range of multiresolution clustering techniques (CATCH, Louvain, Leiden, Seurat). The top four most optimal resolutions from each algorithm are compared to ground truth cluster labels computed on noiseless data, with the highest Adjusted Rand Index score saved. This is repeated 10 times using different random seeds for each algorithm. Shading represents 2 standard deviations around the mean of each algorithm's performance. At the highest levels of dropout and variational noise, CATCH performs significantly better than Louvain, Leiden and Seurat ( $p < 0.05$ , two-sided Student's t-test, with multiple comparisons testing). (F) Condensation's approximation Earth Mover's Distance with an  $\alpha$ -decay kernel shows superior fidelity with ground truth Earth Mover's Distance on 4,360 single cell PBMCs measured on the 10X platform. The diffusion condensation algorithm was implemented with  $\alpha$ -decay (left panel) and Gaussian (right panel) kernels. In each implementation, topological activity analysis identified persistent clustering granularities. For this comparison, granularities with 20 or less clusters were analyzed by comparing all clusters to each other. This resulted in 10,166,640 and 2,541,660 comparisons for the  $\alpha$ -decay kernel and Gaussian implementations respectively.



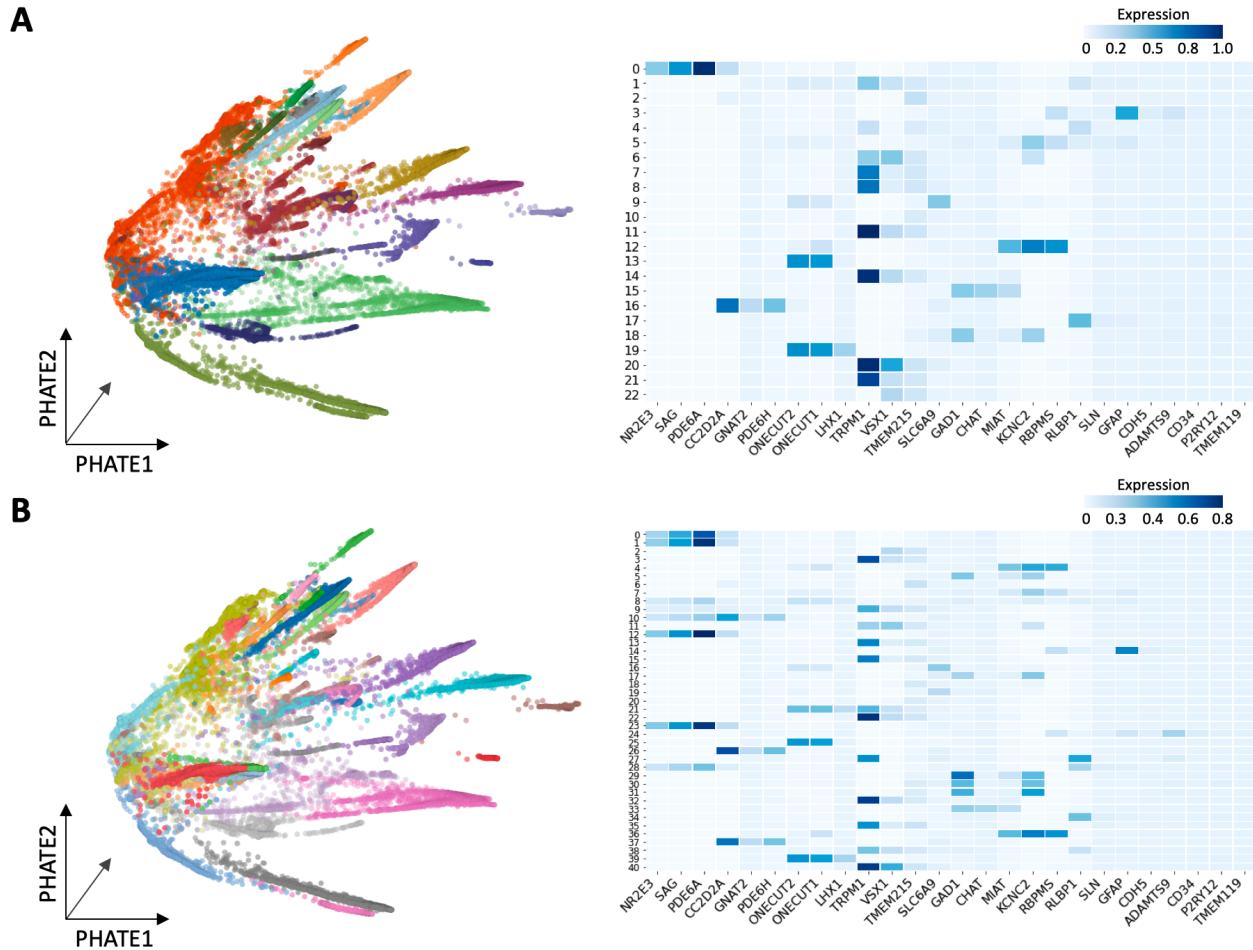
**Extended Data Fig. 2: Key advancements in CATCH clustering approach.** (A) Ideal number of t-steps calculated by spectral entropy per iteration when running diffusion condensation on 4,360 single cell PBMCs measured on the 10X platform. (B) Comparison of different CATCH implementations on synthetic splatter data with increasing levels of noise. Comparing implementations with fixed numbers of t steps set for every iteration ( $t=1,2,3$ ) against the final diffusion condensation approach which uses spectral entropy to tune  $t$  at every iteration. Adaptively tune  $t$  outperforms other implementations significantly as noise levels increase across noise types (two-sided Student's t-test,  $p < .05$ ). (C) Visualization of difference between Gaussian and  $\alpha$ -decay kernels. (D) Comparison of differentially expressed genes using CATCH condensed transport implemented with Gaussian kernel and  $\alpha$ -decay kernel against ground truth 1D-Wasserstein EMD distance. In each comparison, CATCH was run on 10,000 cells generated from splatter as done in part B. Topological activity analysis was used to compute salient granularities for downstream analysis. In each comparison, all salient granularities with less than 20 clusters was used. At each granularity, we computed differentially expressed genes between every combination of clusters. Across all comparisons, 10,249,140 and 4,535,640 comparisons for the Gaussian and  $\alpha$ -decay kernel implementations were computed respectively. (E) Comparing correlation between ground truth EMD values with condensed transport values across granularities using the same multi-cluster and multi-granular comparison strategy described in (D). Visualizing reported correlation values as a feature of cluster granularity (denoted by number of cluster on x-axis) for both Gaussian and  $\alpha$ -decay kernels. (F) Run time comparison between scalable and standard implementations of CATCH.



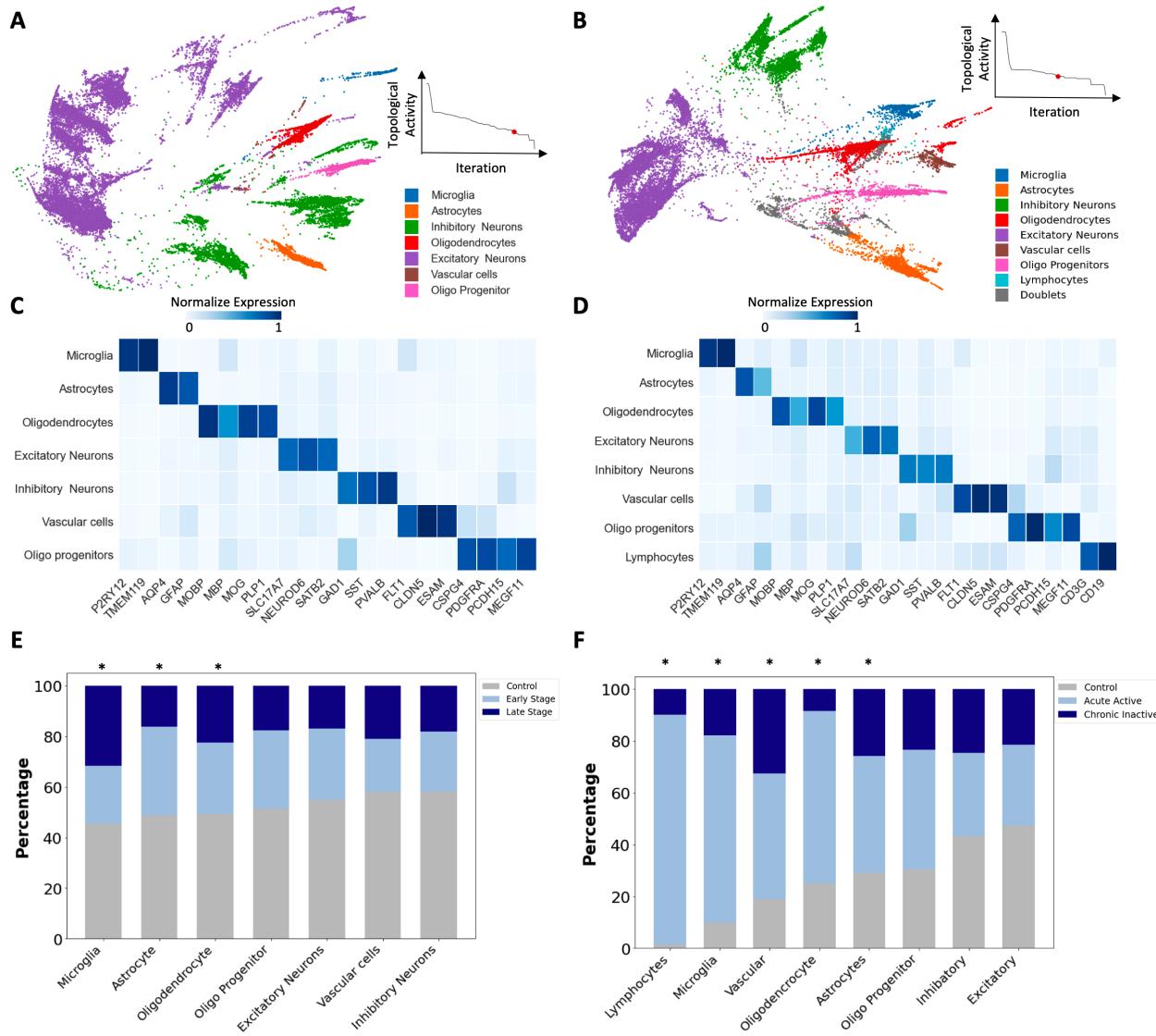
**Extended Data Fig. 3:** Comparison of diffusion condensation against differing implementations and other clustering approaches on synthetic and real single cell data. **(A)** Comparison of CATCH, Louvain, Leiden and FlowSOM on flow cytometry data with cluster labels have been identified through conventional gating analysis. Comparison was repeated on 1.3 million cells generated from 30 patients in the FlowCAP dataset (center line, median; box limits, upper and lower quartiles; whiskers, 1.5x interquartile range; points, all data). **(B)** Comparison of CATCH against multigranular clustering approaches, Louvain and Leiden, on real single cell data across clustering granularity. **(C)** Diffusion condensation cluster characterization implemented with  $\alpha$ -decay and Gaussian kernels compared to ground truth cluster characterizations across granularities on 4,360 PBMCs measured with 10X. **(D)** Comparing correlation between ground truth EMD values with condensed transport values across granularities using the same multi-cluster and multi-granular comparison strategy described in S1F on 10X single-cell data. Reported correlation values as a feature of cluster granularity (denoted by number of cluster on x-axis) for both Gaussian and  $\alpha$ -decay kernels, representing 12,061,332 and 3,373,476 comparisons respectively.



**Extended Data Fig. 4: CATCH identifies known subtypes of bipolar cells across multiple levels of granularity.** (A) Cell type specific signatures based on composite normalized expression of cell type specific marker genes visualized on PHATE. (B) CATCH identifies ON and OFF bipolar cell subsets at a granularity identified with topological activity analysis. (C) Finer grained analysis of ON bipolar cells reveals known subsets. (D) Finer grained analysis of OFF bipolar cells reveals known subsets. (E) CATCH reliably identifies established cell types, as shown by average normalized expression of known bipolar subset-specific marker genes.

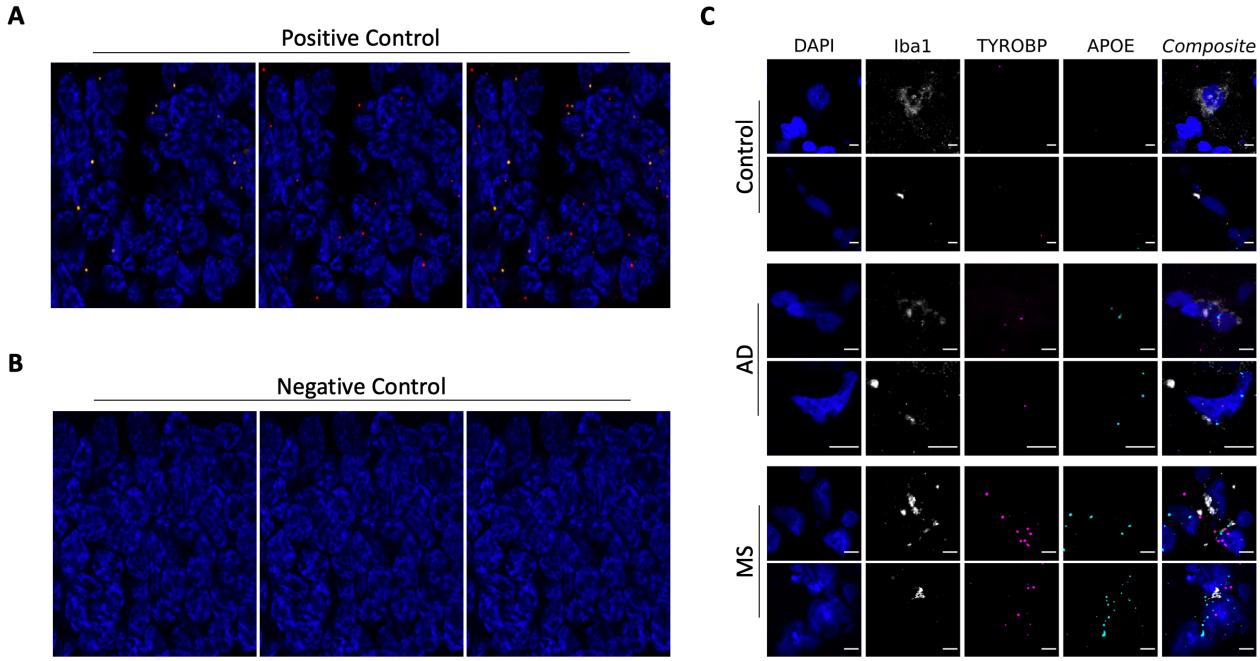


**Extended Data Fig. 5: Louvain does not identify rare glial populations across granularities.** (A) Visualization of 22 coarse grain clusters identified by Louvain. The identified populations are not able to identify all known cell types, as shown by the average normalized expression of known cell type specific marker genes. (B) Visualization of 40 fine grain clusters identified by Louvain. The identified populations at this granularity also do not resolve all known cell types, as shown by the average normalized expression of known cell type specific marker genes.

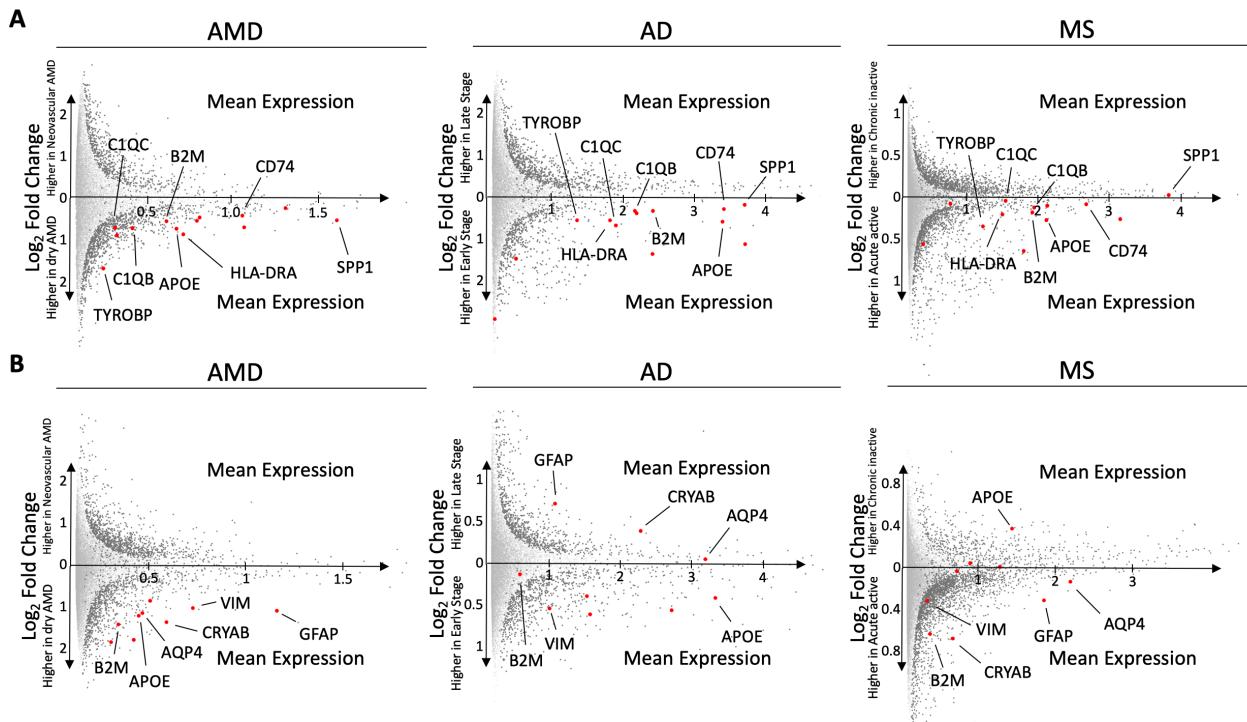


**Extended Data Fig. 6:**

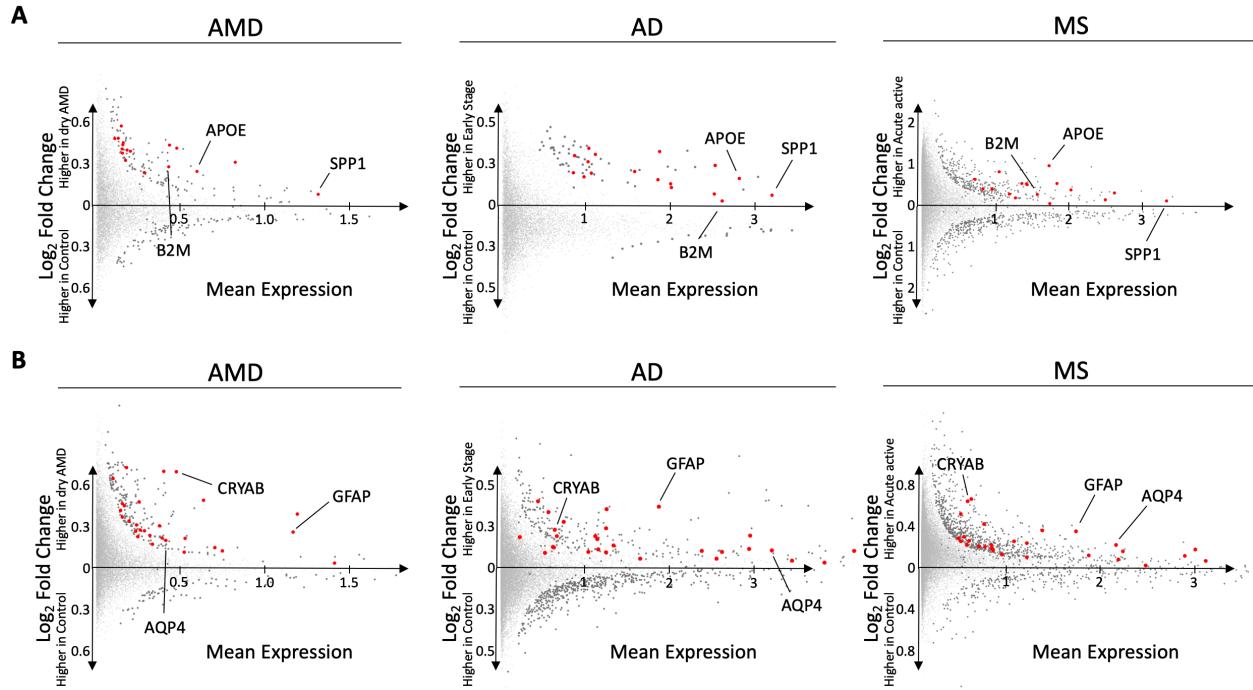
**CATCH analysis of AD and MS snRNAseq data reveals enrichment and activation of microglia and astrocytes in disease.** (A) 43,650 cells pooled from 48 AD patients and healthy donors. Samples were taken from disease free brain tissue and diseased brain tissue at early and late pathological stages. All major cell types were identified by CATCH via persistence analysis and visualized with PHATE [39]. Ideal CATCH granularity identified via topological activity analysis (right). (B) 46,796 cells pooled from 21 progressive MS patients and healthy donors. Samples were taken from disease free brain tissue and diseased brain tissue at acute and chronic stages of inflammation. All major cell types were identified by CATCH via persistence analysis and visualized with PHATE [39]. Ideal CATCH granularity identified via topological activity analysis (right). (C) CATCH reliably identifies cell types in AD brain tissue, as shown by average normalized expression of known cell type-specific marker genes. (D) CATCH reliably identifies cell types in MS brain tissue, as shown by average normalized expression of known cell type-specific marker genes. (E) Microglia and astrocytes are the most enriched cell types in AD using cross condition abundance analysis. (F) Microglia and astrocytes are significantly enriched in progressive MS using cross condition abundance analysis. In E,F: \* =  $p < 0.01$ , two-sided multinomial test with multitest correction.



**Extended Data Fig. 7:** Control probes for fluorescence *in situ* hybridization and validation of early activation microglial signature in MS and AD. Representative images of fluorescence *in situ* hybridization for (A) positive control probe (POL2RA labeled in red and UBC labeled in yellow) and (B) negative control probe (DapB labeled in yellow and red). (C) Representative images of *in situ* RNA hybridization of *APOE* (labeled in turquoise), *TYROBP* (labeled in pink) with simultaneous immunofluorescence of microglial marker IBA1 (white). Elevated expression of *APOE* and *TYROBP* is seen in IBA1-positive cells in microglia from brain tissue with early AD and early progressive MS, compared to healthy controls. Each row represents a sample from a different case. All scale bars = 10 $\mu$ m.



**Extended Data Fig. 8:** Shared activation signatures is diminished in advanced disease and replaced with disease-specific stress signature in microglia. (A) Comparing advanced or chronic inactive disease-enriched microglial cluster to early or acute active-enriched microglial cluster reveals a significant reduction in the microglial activation signature in later stages of neurodegeneration. (B) Comparing advanced or chronic inactive disease-enriched astrocyte cluster to early or acute active-enriched cluster reveals a significant reduction in the astrocyte activation signature in later stages of neurodegeneration.



**Extended Data Fig. 9: Cell type level differential expression analysis at coarse granularity across neurodegenerative diseases.** (A) Performing EMD-based differential expression analysis between microglia which originate from dry AMD patients and control subjects identified a gene signature enriched in the early stage of dry AMD. By performing similar differential expression analysis between microglia from brain samples from patients with early AD, acute active MS, and controls, we identified a shared activation signature of 17 genes. This common signature includes *APOE*, *SPP1* and *B2M* while all other genes are highlighted in red. FDR corrected p-value<.1 in all comparisons. (B) Performing EMD-based differential expression analysis between astrocytes which originate from dry AMD patients and control subjects identified a gene signature enriched in the early stage of dry AMD. By performing similar differential expression analysis between astrocytes from control and early AD samples and control and acute inflammation MS samples, we identify a shared activation signature of 28 genes. This common signature includes *GFAP*, *AQP4*, and *CRYAB* while all other genes are highlighted in red. FDR corrected p-value<.1 in all comparisons.

## 1164 Extended Data Tables

| Retina | Sex | Age | Postmortem Interval (Hrs) | Left/Right Eye | Condition            | Assay     |
|--------|-----|-----|---------------------------|----------------|----------------------|-----------|
| 1      | F   | 90  | 2                         | Right          | Control              | snRNA-seq |
| 2      | F   | 81  | 4                         | Left           | Control              | snRNA-seq |
| 3      | M   | 65  | 8                         | Left           | Control              | snRNA-seq |
| 4      | M   | 78  | 4                         | Left           | Control              | snRNA-seq |
| 5      | F   | 72  | 4                         | Left           | Control              | snRNA-seq |
| 6      | M   | 85  | 3                         | Right          | Control              | snRNA-seq |
| 7      | M   | 86  | 9                         | Left           | Intermediate Dry AMD | snRNA-seq |
| 8      | M   | 72  | 3                         | Left           | Intermediate Dry AMD | snRNA-seq |
| 9      | F   | 82  | 4                         | Right          | Intermediate Dry AMD | snRNA-seq |
| 10     | F   | 67  | 2                         | Left           | Neovascular AMD      | snRNA-seq |
| 11     | F   | 79  | 10                        | Right          | Neovascular AMD      | snRNA-seq |
| 12     | F   | 100 | 9                         | Left           | Neovascular AMD      | snRNA-seq |
| 13     | F   | 93  | 8                         | Right          | Neovascular AMD      | snRNA-seq |
| 14     | F   | 94  | 13                        | Left           | Neovascular AMD      | snRNA-seq |
| 15     | F   | 92  | <1                        | Right          | Neovascular AMD      | snRNA-seq |
| 16     | F   | 76  | 9                         | Left           | Neovascular AMD      | snRNA-seq |
| 17     | F   | 74  | 5                         | Right          | Neovascular AMD      | snRNA-seq |

Extended Data Table 1: Human retinal specimen details.

| Subject | Sex | Age | Smoking History | Hypertension | Diabetes | Treatment anti-VEGF | Treatment AREDS vitamins |
|---------|-----|-----|-----------------|--------------|----------|---------------------|--------------------------|
| 1       | F   | 90  | No              | No           | No       | No                  | No                       |
| 2       | F   | 81  | Unk.            | Unk.         | Unk.     | Unk.                | Unk.                     |
| 3       | M   | 65  | No              | No           | No       | No                  | No                       |
| 4       | M   | 78  | No              | Yes          | No       | No                  | No.                      |
| 5       | F   | 72  | No              | No           | No       | No                  | No                       |
| 6       | M   | 85  | No              | No           | No       | No                  | No                       |
| 7       | M   | 86  | Former smoker   | No           | No       | No                  | No                       |
| 8       | M   | 72  | No              | No           | No.      | No                  | Yes                      |
| 9       | F   | 82  | No              | No           | No       | No                  | Yes                      |
| 10      | F   | 67  | Yes             | Yes          | No       | Yes                 | No                       |
| 11      | F   | 79  | Unk.            | Unk.         | Unk.     | Unk.                | Unk.                     |
| 12      | F   | 100 | Unk.            | Unk.         | Unk.     | Unk.                | Unk.                     |
| 13      | F   | 93  | Unk.            | Unk.         | Unk.     | Unk.                | Unk.                     |
| 14      | F   | 94  | Unk.            | Unk.         | Unk.     | Unk.                | Unk.                     |
| 15      | F   | 92  | Unk.            | Unk.         | Unk.     | Unk.                | Unk.                     |
| 16      | F   | 76  | No              | Yes          | No       | Yes                 | Yes                      |
| 17      | F   | 74  | No              | No           | No       | Yes                 | Yes                      |

**Extended Data Table 2:** Human subject clinical information. ‘Unk.’ = Data not present in available records.

| Retina | Sex | Age | Postmortem Interval (Hrs) | Cells post QC | UMI/Cell | Genes/Cell |
|--------|-----|-----|---------------------------|---------------|----------|------------|
| 1      | F   | 90  | 2                         | 3069          | 6704     | 2545       |
| 2      | F   | 81  | 4                         | 3019          | 9107     | 2865       |
| 3      | M   | 65  | 8                         | 3934          | 10819    | 3005       |
| 4      | M   | 78  | 4                         | 1182          | 5255     | 2270       |
| 5      | F   | 72  | 4                         | 4085          | 5923     | 2281       |
| 6      | M   | 85  | 3                         | 1321          | 3052     | 1670       |
| 7      | M   | 86  | 9                         | 7848          | 5652     | 2360       |
| 8      | M   | 72  | 3                         | 6311          | 9587     | 3298       |
| 9      | F   | 82  | 4                         | 4814          | 6087     | 2415       |
| 10     | F   | 67  | 2                         | 1533          | 9609     | 3152       |
| 11     | F   | 79  | 10                        | 3465          | 4989     | 2084       |
| 12     | F   | 100 | 9                         | 5349          | 8565     | 2887       |
| 13     | F   | 93  | 8                         | 4568          | 8708     | 3239       |
| 14     | F   | 94  | 13                        | 281           | 2187     | 1378       |
| 15     | F   | 92  | <1                        | 6696          | 5426     | 2348       |
| 16     | F   | 76  | 9                         | 7476          | 6408     | 2658       |
| 17     | F   | 74  | 5                         | 6112          | 6321     | 2234       |

**Extended Data Table 3:** Human retinal specimen details.