

1 **Title:** Diagnostic Evidence GAUGE of Single cells (DEGAS): A flexible deep-transfer learning
2 framework for prioritizing cells in relation to disease

3 **Authors:** Travis S. Johnson^{1,2,3}, Christina Y. Yu^{1,2}, Zhi Huang⁴, Siwen Xu⁵, Tongxin Wang⁶,
4 Chuanpeng Dong⁵, Wei Shao¹, Mohammad Abu Zaid¹, Xiaoqing Huang³, Yijie Wang⁶,
5 Christopher Bartlett⁷, Yan Zhang^{2,8}, Brian A. Walker⁹, Yunlong Liu^{5,10}, Kun Huang^{1,3,10,11*}, Jie
6 Zhang^{10*}

7

8 **Affiliations**

9 ¹Department of Medicine, Indiana University School of Medicine

10 ²Department of Biomedical Informatics, College of Medicine, The Ohio State University

11 ³Department of Biostatistics and Health Data Science, Indiana University School of
12 Medicine

13 ⁴School of Electrical and Computer Engineering, Purdue University

14 ⁵Center for Computational Biology and Bioinformatics, Indiana University School of
15 Medicine

16 ⁶Department of Computer Science, Indiana University

17 ⁷Nationwide Children's Hospital

18 ⁸The Ohio State University Comprehensive Cancer Center (OSUCCC - James)

19 ⁹Division of Hematology Oncology, Melvin and Bren Simon Comprehensive Cancer
20 Center, Indiana University

21 ¹⁰Department of Medical and Molecular Genetics, Indiana University School of Medicine

22 ¹¹Regenstrief Institute

23 * To whom correspondence should be addressed (jizhan@iu.edu or kunhuang@iu.edu)

24

25 **Abstract**

26 We propose *DEGAS* (Diagnostic Evidence GAUge of Single cells), a novel deep transfer
27 learning framework, to transfer disease information from patients to cells. We call such
28 transferrable information “impressions,” which allow individual cells to be associated with
29 disease attributes like diagnosis, prognosis, and response to therapy. Using simulated data and
30 ten diverse single cell and patient bulk tissue transcriptomic datasets from Glioblastoma
31 Multiforme (GBM), Alzheimer’s Disease (AD), and Multiple Myeloma (MM), we demonstrate the
32 feasibility, flexibility, and broad applications of the *DEGAS* framework. *DEGAS* analysis on
33 newly generated myeloma single cell transcriptomics led to the identification of *PHF19^{high}*
34 myeloma cells associated with progression.

35

36 **Keywords**

37 Prognostic models, Survival, Cox proportional hazards, Single cell RNA sequencing, scRNA-
38 seq, Machine Learning, Deep learning, Transfer learning, Multiple Myeloma, Alzheimer’s
39 Disease

40

41 **Background**

42 The emergence of single cell RNA sequencing (scRNA-seq) in 2009 has revolutionized the
43 medical research community with single cell level resolution, providing a much deeper
44 understanding of transcriptomic heterogeneity in tissues and diseases. Now that scRNA-seq is a
45 standard part of the biomedical research toolbox, increasing numbers of scRNA-seq studies have
46 been published [1, 2], and databases have quickly accumulated with scRNA-seq data, such as
47 Hemberg lab [3], scRNASeqDB [4], SCPortalen [5], Allen Institute Cell Types Database, and the
48 NCBI Gene Expression Omnibus (GEO) [5]. Many methods have been developed to analyze
49 scRNA-seq data, the most notable being *Seurat*, which includes ways to cluster and normalize
50 cell expression as well as perform integrative analysis with other data types (e.g., CITE-seq and
51 ATAC-seq) [6]. These methods are important for understanding many prognostic and diagnostic

52 disease attributes in scRNA-seq data. Here we use “disease attributes” as a broad term inclusive
53 of many types of information and labeling such as diagnostic information, disease subtypes,
54 disease status, prognostic information like survival, and responses to therapy. For *Seurat* and
55 similar methods, while cell types/clusters can be identified and associated with disease attributes
56 [7-10], individual cells are unable to be associated in the same manner. This may result in failing
57 to identify subsets of cells associated with disease attributes, especially if the disease-associated
58 cells cluster together with non-disease-associated cells.

59

60 Currently, disease associated cell types can be identified by transferring molecular heterogeneity
61 information from cells to patients using single cell expression deconvolution [11-13]. However,
62 this approach is limited as it focuses on the changes in relative abundance of subtypes of cells
63 instead of transcription changes of these cells. The resolution of the cell subtyping is constrained
64 by the clustering experiment. Therefore, novel machine learning methods that can transfer
65 information from patients to cells and identify latent links between them are sorely needed to
66 leverage the relative strengths of single cell and patient level data. For example, in cancer studies,
67 bulk transcriptomic data is ideal for studying inter-tumor heterogeneity and scRNA-seq is ideal for
68 studying intra-tumor heterogeneity. However, such integration faces numerous challenges since
69 different data modalities and different data sources can have different characteristics in terms of
70 quantity, quality, distribution and resolution [1]. For instance, it is common to find studies with a
71 large number of patient samples for bulk tissue RNA sequencing (RNA-seq), whereas studies
72 with scRNA-seq data usually contain a small number of patient samples. Most scRNA-seq
73 experiments generate a large number of cells per sample, making the scaling of such data to
74 multiple tissue samples computationally difficult [1]. On the other hand, a large patient sample
75 size is often required for statistical studies such as prediction of disease attributes [14]. If
76 traditional methods were used, the resulting scRNA-seq data could end up with cell numbers on
77 the scale of millions making such studies more difficult.

78

79 To address these challenges, previous studies have directly established associations of diseases
80 with cell types derived from scRNA-seq without using deconvolution. These methods mainly
81 utilize unsupervised methods and focused primarily on the number of differentially expressed
82 genes (DEGs) in a given cell type corresponding to DEGs related to some disease attribute [15,
83 16]. For example, Gawel *et al.* used enrichment of the cell cluster specific DEGs and multicellular
84 disease models (MCDMs) to visualize the cell types for prioritization [7]. *Muscat* identified DEGs
85 between treatment groups in scRNA-seq samples which were used to identify cell types related
86 to sample treatment [17]. Alternatively, k nearest neighbor (*kNN*) graphs have been used to
87 identify cell types that undergo transcriptional changes related to biological perturbations [18].
88 The cell type prioritization tool *Augur* did not primarily rely on DEGs, but still focused the biological
89 resolution to the cell type level [19]. They trained classifiers on each cell type with respect to the
90 disease state of the tissue from which those cells were sampled. The accuracy of the classifier in
91 each cell type was used to prioritize its relation to the disease state of interest [19]. These methods
92 rely on either prior knowledge to calculate enrichment of DEGs or require scRNA-seq data from
93 both disease and normal samples. Furthermore, all of these existing methods are reliant on
94 accurately defining the cell types within a scRNA-seq experiment. In summary, these methods
95 assign disease associations to the previously defined cell types and not to the individual cells.

96

97 To address such challenges as prioritizing individual cells in relation to disease with
98 considerations on sample size and computational cost, we established the combined deep
99 learning and transfer learning framework called *DEGAS* (Diagnostic Evidence GAuge of Single
100 cells) to integrate scRNA-seq and bulk tissue transcriptomic data with the goal to transfer clinical
101 information from patients to cells. The ability of *DEGAS* to assign patient-level disease attributes
102 to single cells, among other functions, provides a flexible and useful tool to prioritize cells, cell
103 types, patients, and patient subtypes in relation to disease attributes. In this paper, we focus on

104 the most relevant use case of associating disease attributes from patients to individual cells since
105 there is no current state-of-the art technique to perform this task.

106

107 We use transcriptomic data as an example where bulk expression is referred to as patients and
108 scRNA-seq is referred to as cells. The rationale behind the *DEGAS* framework is that scRNA-seq
109 data and patient-level transcriptomic data (e.g., RNA-seq with clinical information) share the same
110 feature space (i.e., common set of genes). In addition, a natural connection exists between the
111 two data types that can be leveraged to further identify the associations between patients and
112 cells. Viewing this association as a graph (**Fig. 1**), we can connect the disease attributes in
113 patients to individual cells, via a latent representation of the common feature space (selected
114 genes). This latent representation fitting two datasets can be learned using a transfer learning
115 technique called domain adaptation [20-23]. Domain adaptation applies linear or non-linear
116 transformations on the features for both datasets so that their distributions are similar after the
117 transformations. Our biological intuition is thus: the expression patterns of genes in cells and
118 tissues should carry a portion of the same biological patterns such as molecular pathways,
119 signaling cascades, and/or metabolic processes, making the information learned from this portion
120 of gene expression patterns transferable between patients and cells. Our hypothesis is that the
121 latent representation learned from these shared gene expression patterns will be simultaneously
122 predictive of patient disease attributes and cellular subtypes. Similar hypotheses are already
123 adopted to transfer information between different single cell experiments [6, 24-28] and to transfer
124 information from bulk transcriptomic cell type atlases to single cell experiments [29].

125

126 In our *DEGAS* framework, we incorporate patient-level disease attributes information with cell
127 type information from disparate datasets to perform cell prioritization on scRNA-seq data. These
128 disease associations in cells can be attributed to disease-related biological perturbations
129 identified in the patients. This novel deep transfer learning approach simultaneously trains a

130 model on single cell data and patient data along with their labels and learns a representation in
131 which the cells and patients occupy the same latent space. Multitask learning, also known as
132 parallel transfer learning, is precisely designed to achieve these two goals. Used extensively in
133 computer vision, multitask learning learns a low dimensional representation of the input data to
134 optimally address multiple tasks. Examples of such application in medical science include
135 predicting benign versus malignant tumor samples and subclassification in breast cancer
136 histology images [30, 31]. In this paper, we further extend this line of research to include datasets
137 with patient disease attributes that can be trained simultaneously so that the disease attributes
138 can be transferred (or cross-mapped) between single cells and patients. Specifically, our
139 framework enables knowledge learned from patients using deep learning models to be transferred
140 to single cells and vice versa. The major advantages of our transfer learning framework are that
141 the single-cell gene expression data and clinical bulk gene expression data can come from
142 different patient cohorts of the same disease without matched data while the disease associations
143 can still be directly assigned to individual cells. This flexibility not only presents an ingenious way
144 to integrate molecular omics data analysis in different levels, but also virtually merges them into
145 the same cohort, which makes studying a broad variety of heterogeneous diseases possible.

146

147 Various types of workflows can integrate the *DEGAS* framework, which can be tailored to user
148 preference and data availability. These workflows consist of preprocessing, formatting data,
149 training *DEGAS* models using the *DEGAS* framework, predicting disease associations in cells
150 using the *DEGAS* framework, and downstream analysis (**Fig. 1A**). The *DEGAS* framework in its
151 simplest form can be broken into three tasks during model training: 1) correctly labeling cells with
152 a cellular subtype using multitask learning; 2) correctly assigning clinical labels to patients using
153 multitask learning; and 3) generating a latent space in which patients and cells are comparable
154 using domain adaptation (**Fig. 1B**). To perform *DEGAS* analysis, first we select representative
155 gene features that are predictive of cell type, predictive of patient disease attributes, and present

156 at measurable levels in both scRNA-seq and bulk transcriptomic data. Secondly, we apply deep
157 learning models to learn the latent representation of the single-cell and patient-level transcriptomic
158 data, with the goal to simultaneously minimize cell type classification error, patient disease
159 attribute prediction error, and the differences between cells and patients in their latent
160 representation. Finally, the patient-level disease attributes such as survival and clinical subtypes
161 is predicted in the single cells using the patient label output layer and cell types are predicted in
162 patients using the cell type output layer (**Fig. 1C**). We call these transferrable label probabilities
163 “impressions” since information from gene expression of disparate data types and studies can be
164 extracted and the characteristics from one data type can be mapped to another. These
165 impressions of disease attributes in single cells can be wide ranging characteristics of the patient
166 samples but must be categorical or time to event. The most interesting of them that can be used
167 in *DEGAS* are disease status, disease subtype, survival, and response to therapy. Disease status,
168 subtype, and survival were used in our current experiments but there would also be much utility
169 in identifying cells associated with poor response to treatment as the data become available.
170 Furthermore, we emphasize the ability to make predictions of patient disease attributes in
171 individual cells since there is a lack of such method to perform this task to the best of our
172 knowledge. *DEGAS* is developed as a generalizable model generating deep transfer learning
173 framework that can be applied to any disease data as long as the data contain clinical information
174 for a cohort of patients or a separate clustering analysis result on sets of cells from single cell
175 level omic experiments of the same disease. Since there is not an inherent limitation to the use
176 of transcriptomic data, *DEGAS* can be potentially expanded to accommodate other modalities of
177 data with proper normalization steps.

178

179 To demonstrate the feasibility and effectiveness of the *DEGAS* framework, we first tested it on
180 simulated data and glioblastoma (GBM) transcriptomic data, which contain ground-truth labels of
181 cell types on single cell gene expression data and clinical labels for patient bulk tissue gene

182 expression data. Then we applied *DEGAS* to multiple Alzheimer's disease (AD) gene expression
183 datasets from Mount Sinai/JJ Peters VA Medical Center, Allen Institute, Grubmann *et al.* [32], and
184 Mathys *et al.* [15] in which certain cell type changes (microglia and neuron) are largely known [33-
185 39]. Finally, as an exploratory tool, we applied *DEGAS* to study multiple myeloma (MM)
186 transcriptomic data, where the disease associated subtypes of cells are largely unknown.

187
188 MM is a late stage of myeloma that stems from the proliferation of aberrant clonal plasma cells
189 in the bone marrow that secrete monoclonal immunoglobulins and is the second most common
190 blood cancer in the United States [40]. Patient level transcriptomic data for MM has been widely
191 available for some time and has been used to identify subtypes of MM with different prognoses
192 [41]. However, only recently has scRNA-seq become available for MM [9, 42, 43] and few
193 studies have identified the most high-risk subtypes of cells [9]. Here we combined our newly
194 generated late-stage myeloma scRNA-seq data from four local samples and bulk tissue data
195 from the Multiple Myeloma Research Foundation CoMMpass study, and then applied *DEGAS* to
196 infer clinical impressions for myeloma cell subtypes and successfully identified a *PHF19^{high}*
197 myeloma cell subgroup associated with a high-risk of progression.

198

199 **Methods**

200 ***Experimental design and datasets***

201 For a *DEGAS* cell prioritization experiment, one scRNA-seq dataset, one bulk expression dataset,
202 and patient sample labels (matched with the bulk data samples) are required as input. After
203 feature selection and scaling (see ***Feature selection and scaling***) of the raw input expression
204 data, there should be two expression matrices with rows corresponding to samples/cells and
205 matching columns corresponding to genes. The bulk patient sample labels should be one-hot
206 encoded in a matrix with rows corresponding to each sample and the columns corresponding to
207 each class of label. For survival sample labels the first column should be time and the second

208 column should be the even indicator (1 event and 0 censored). If cell labels are also available,
209 they should also be one-hot encoded with each row corresponding to a cell and each column
210 corresponding to a class of label. The *DEGAS* models can be trained and predicted on these
211 formatted data (**Fig. 1A**).

212

213 In this study we analyzed simulated data and data from three different diseases, GBM, AD, and
214 MM, to test the *DEGAS* framework and apply it for novel discoveries. The simulation, GBM, and
215 AD experiments were primarily used as validation datasets since the ground truth is known. The
216 simulated data were generated so that cell types are directly related to disease status in patients.
217 For GBM data, we used scRNA-seq data for five tumors from Patel *et al.* [44] and microarray data
218 for the GBM TCGA cohort [45] (**Table 1**). For AD data, we used human scRNA-seq from Allen
219 Institute for Brain Science (AIBS) Cell Types Database (<https://celltypes.brain-map.org/>) and AD
220 patient RNA-seq from the Mount Sinai/JJ Peters VA Medical Center Brain Bank (MSBB) study
221 [46] (**Table 1**).

222

223 **Table 1. Summary of the clinical features in each patient cohorts used in training.** * Final
224 age category is >90 years.

Glioblastoma Multiforme TCGA	
Feature	Details
Sex	74 Male, 37 Female
Age (years)	Range: 14-83, Mean: 56, Median: 58
Clinical GBM subtype	34 Classical, 33 Mesenchymal, 9 Neural, 35 Proneural
Alzheimer's Disease MSBB	
Feature	Details
Sex	90 Male, 131 Female

Age (years)	Range: 61-90+, Mean* > 82, Median = 84
AD diagnosis	135 AD, 86 Control
Multiple Myeloma MMRF	
Feature	Details
Sex	387 Male, 260 Female
Age (years)	Range: 27-93, Mean: 64, Median: 64
Relapse-free survival time (days)	Range: 13-1753, Mean: 665.4, Median: 629 200 patients progressed

225

226

227 We further expanded our inquiry into MM, which served as a discovery study. Since the plasma
228 cell subtypes are less understood in relation to MM clinical outcomes, we aimed to identify
229 subtypes of plasma cells associated with worse prognosis. We first utilized 647 CD138⁺-
230 enriched bone marrow patient samples from the Multiple Myeloma Research Foundation
231 CoMMpass study (MMRF). These data were generated as part of the Multiple Myeloma
232 Research Foundation Personalized Medicine Initiatives (<https://research.themmr.org>). The
233 dataset consisted of tumor tissue RNA-seq data and corresponding clinical information including
234 progression free survival (PFS) time and survival status. PFS was defined as the time taken for
235 a patient to relapse, progress, or die after treatment of the initial tumor. The demographic
236 information of the MMRF patients are shown in **Table 1**. The first scRNA-seq data used in this
237 study were generated by us using samples consisting of CD138⁺ plasma cells purified from
238 bone marrow from four myeloma patients including two MM patients.

239

240 There were six total samples collected from myeloma patients. Of these, four samples passed
241 initial quality control checks. Sample 1 and 6 were dropped due to sample degradation and data
242 quality issues. This in turn left with four usable samples, *i.e.*, samples 2, 3, 4, and 5 for our
243 study. The low number of patients was a good test case considering most scRNA-seq
244 experiments frequently have few patients. The single cells were sequenced using 10x
245 Genomics and Illumina NovaSeq6000 sequencer. *CellRanger 2.1.0*
246 (<http://support.10xgenomics.com/>) was utilized to process the raw sequence data. Briefly,
247 *CellRanger* used *bcl2fastq* (<https://support.illumina.com/>) to demultiplex raw base sequence
248 calls generated from the sequencer into sample-specific FASTQ files. The FASTQ files were
249 then aligned to the human reference genome GRCh38 with RNA-seq aligner *STAR*. The aligned
250 reads were traced back to individual cells and the gene expression level of individual genes
251 were quantified based on the number of UMIs (unique molecular indices) detected in each cell.
252 The filtered gene-cell barcode matrices generated by *CellRanger* were used for further analysis.
253 A second publicly available myeloma scRNA-seq dataset was used for validation, which
254 consisted of NHIP (normal control), MGUS (monoclonal gammopathy of undetermined
255 significance), SMM (smoldering multiple myeloma), and MM [42]. A second bulk tissue dataset
256 was used for validating the proportional hazards modeling. This dataset consisted of bulk
257 expression profiling by microarray of CD138+ plasma cells with overall survival (OS) information
258 for 559 MM patients [41]. The detailed information of the four datasets is shown in **Table 2**.

259

260 **Table 2. Overview of all datasets used in the analysis.** *The simulated patients were
261 generated from the splatter simulated cells by combining known proportions of cell types.
262 “None” is used to denote the lack of labels for the cells/samples in a given dataset. †Cells were
263 down-sampled from the total number of cells because some cell types were over-represented.

Study	Dataset	Sample size	Data type	Attribute
-------	---------	-------------	-----------	-----------

Simulation	Simulated cells*	5000 cells	scRNA-seq	Cell type
	Simulated patients*	600 patients	RNA-seq	Disease status
Glioblastoma	Patel <i>et al.</i> , 2014	532 cells (5 patients)	scRNA-seq (SMART-seq)	None
	TCGA GBM	111 patients	Microarray	GBM subtype
Alzheimer's disease	AIBS	47,396 cells (11 patients)	scRNA-seq (SMART-seq)	Brain cell types
	Grubman <i>et al.</i> , 2019	13,214 cells (12 patients)	snRNA-seq (10x Genomics)	AD and normal brain cell types
	Mathys <i>et al.</i> , 2019	5288 cells [†] (48 patients)	snRNA-seq (10x Genomics)	AD and normal brain cell types
	MSBB	682 samples (221 patients)	RNA-seq	AD diagnosis
Multiple myeloma	MMRF	647 patients	RNA-seq	PFS
	IUSM	22,968 cells (4 patients)	scRNA-seq (10x Genomics)	Subtype cluster (Subtype 1-5)
	Ledergor <i>et al.</i> , 2019	13,440 cells (35 patients)	scRNA-seq (MARS-seq)	Malignancy (NHIP, MGUS, SMM, MM)
	Zhan <i>et al.</i> , 2006	559 patients	Microarray	OS

264

265 ***Transfer learning using DEGAS***

266 Several types of labels including Cox proportional hazards, patient classification, and cell type
 267 classification, along with maximum mean discrepancy (MMD), a technique used to match
 268 distributions across different sets of data [22], were combined to create the multitask transfer
 269 learning framework *DEGAS*.

270

271 The first step was to find a set of gene expression features that were both informative of cell type
 272 and of patient disease attribute (*e.g.*, recurrence). The intersection of high variance genes found
 273 in the scRNA-seq and bulk expression data of patient samples are used for further analysis. The

274 definition of this gene set is up to the user but *Seurat-CCA*, *LASSO* selection, and even statistical
275 tests such as t-test and f-test can be used to define the gene set. Since these features are the
276 same between patients and single cells, the patients and cells share the same input layer. This
277 makes it possible to predict proportional hazard and cell type regardless of the input sample type
278 (patient or single cell data).

279

280 All experiments in this manuscript use a five-bootstrap aggregated three-layer DenseNet-based
281 implementation of *DEGAS*, but the simplest form of the *DEGAS* framework is a single layer
282 network. In our description of the overall architecture below (shown in **Fig. 1B,C**), we used a
283 single layer network for the purpose of simplicity. The following **Eq. 1** can nevertheless be
284 extrapolated to multiple layers and architectures, some of which we have already included in our
285 open-source software package. First, a hidden layer was used to transform the genes into a lower
286 dimension using a sigmoid activation function (**Eq. 1**). Where X represents an input expression
287 matrix, θ_{Hidden} represents the hidden layer weights, and b_{Hidden} represents the hidden layer bias.

$$288 \quad f_{Hidden}(X) = \text{sigmoid}(X^T \theta_{Hidden} + b_{Hidden}) \quad \text{Eq. 1}$$

289

290 Next, output layers were added for both the patient output and for the single cell output. For the
291 single cells, there could be classification output or no output. No output means there are no known
292 labels for the single cells to match. Similarly, patients could have Cox proportional hazard output,
293 classification output, or no output (implying no known labels for patients).

294

295 The Cox proportional hazards estimates consisted of a linear transformation to a single output
296 followed by a sigmoid activation function (**Eq. 2**):

$$297 \quad f_{Cox}(X) = \text{sigmoid}(f_{Hidden}(X)^T \theta_{Cox} + b_{Cox}), \quad \text{Eq. 2}$$

298 where the variable X represents an input expression matrix, θ_{Cox} represents the Cox proportional
 299 hazard layer weights [47], and b_{Cox} represents the Cox proportional hazard layer bias. The
 300 classification output consisted of a transformation to the same number of outputs as the number
 301 of labels, *i.e.*, patient subtypes, cellular subtypes, using a softmax activation function (**Eq. 3**).

$$302 \quad f_{class}(X) = \text{softmax}(f_{Hidden}(X)^T \theta_{class} + b_{class}), \quad \text{Eq. 3}$$

303 θ_{class} represents the classification layer weights and b_{class} represents the classification layer bias.

304

305 To train the *DEGAS* model, we need to compute three types of loss functions for the Cox
 306 proportional hazards output, classification output, and MMD [22] respectively. The Cox
 307 proportional hazards loss [47] was calculated only for the patient expression data (X_{Pat}) using the
 308 followup period (C), and event status (t) (**Eq. 4**). Similarly, the patient classification loss was only
 309 calculated for the patient data (X_{Pat}) using the patient labels (Y_{Pat}). Alternatively, the cellular
 310 classification loss was only calculated for the single cell expression data (X_{Cell}) and true subtype
 311 label (Y_{Cell}) (**Eq. 5**). The MMD loss was calculated between the patient expression data (X_{Pat}) and
 312 the single cell expression data (X_{Cell}) (**Eq. 6**), which is the key for mapping the distributions of the
 313 data representations between the single-cell and patient bulk tissue data.

$$314 \quad Loss_{Cox} = \sum_{C(i)=1} (f_{Cox}(X_{Pat})_i - \sum_{t_j \geq t_i} (\exp(f_{Cox}(X_{Pat})_j))) \quad \text{Eq. 4}$$

$$315 \quad Loss_{class} = \frac{1}{n} \sum_{i=1}^n (\sum (Y_{type,i} - f_{class}(X_{type})_i)) \text{ where } type \in \{Pat, Cell\} \quad \text{Eq. 5}$$

$$316 \quad Loss_{MMD} = MMD(X_{Cell}, X_{Pat}) \quad \text{Eq. 6}$$

317 Besides the three losses, we also add a L_2 -regularization loss term to constrain for the complexity
 318 of the model. The overall loss function was the weighted sum of the four types of loss using the
 319 hyper-parameters λ_0 (single cell loss function), λ_1 (patient loss function), λ_2 (MMD loss), and λ_3
 320 (regularization loss), so that the importance of each loss term and regularization term could be
 321 adjusted (**Eq. 7**):

$$322 \quad Loss_{ClassCox} = \lambda_0 Loss_{class} + \lambda_1 Loss_{Cox} + \lambda_2 Loss_{MMD} + \lambda_3 \|\theta\|_2^2. \quad \text{Eq. 7}$$

323

324 To address more diverse scenarios, we can also adapt **Eq. 7** for two classification outputs (**Eq.**
325 **8**), a single classification output without patient disease attribute (**Eq. 9**), a single classification
326 output without cell type label (**Eq. 10**), or a single Cox output without cell type label (**Eq. 11**):

$$327 \quad LOSS_{ClassClass} = \lambda_0 LOSS_{Class} + \lambda_1 LOSS_{Class} + \lambda_2 LOSS_{MMD} + \lambda_3 \|\theta\|_2^2, \quad \mathbf{Eq. 8}$$

$$328 \quad LOSS_{ClassBlank} = \lambda_0 LOSS_{Class} + \lambda_2 LOSS_{MMD} + \lambda_3 \|\theta\|_2^2, \quad \mathbf{Eq. 9}$$

$$329 \quad LOSS_{BlankClass} = \lambda_1 LOSS_{Class} + \lambda_2 LOSS_{MMD} + \lambda_3 \|\theta\|_2^2, \quad \mathbf{Eq. 10}$$

$$330 \quad LOSS_{BlankCox} = \lambda_1 LOSS_{Cox} + \lambda_2 LOSS_{MMD} + \lambda_3 \|\theta\|_2^2. \quad \mathbf{Eq. 11}$$

331

332 In summary, a common hidden layer was used to merge the single cells and patient data. Next,
333 an output layer was added to predict the proportional hazards or classes of the patient samples
334 [47]. The loss function for the proportional hazards prediction or patient classification was back-
335 propagated across both layers for each patient. The single cells also had an output layer
336 consisting of a softmax output to predict the cellular subtype of each cell. Error was back-
337 propagated across both layers from the label output for each cell. Finally, a model was learned
338 that can model both the single cells and the patients. To perform this task, we utilized the MMD
339 method [22] to reduce the differences between patients and cells in a low dimensional
340 representation. Both single cell and patient bulk tissue data were combined into a single group
341 such that the MMD loss was minimized between patient bulk tissue data and single cell data from
342 multiple patients. Because there are many different combinations of these outputs, *i.e.*, single cell
343 output followed by patient output, we implemented ClassCox, ClassClass, ClassBlank,
344 BlankClass, and BlankCox models based on equations (7)-(11) in the current version but intend
345 to provide more options in the future.

346

347 To keep the analyses consistent, we used the same network architecture and hyperparameters
348 throughout all of the experiments. Specifically, we used a three-layer DenseNet architecture
349 bootstrap aggregated five times such that **Eq. 1** would consist of a DenseNet instead of a single
350 layer feedforward network and five such models were trained. The same set of hyper-parameters
351 were used in all of the experiments in this study, except for the robustness to hyper-parameters
352 experiment, where they were intentionally altered to test the influences on the output results.
353 These are considered the default hyper-parameters in the *DEGAS* package but can be changed.
354 They are: training steps 2000, single cell batch size 200, patient batch size 50, hidden layer nodes
355 50, drop-out retention rate 50%, single cell loss weight (λ_0) 2, patient loss weight (λ_1) 3, MMD
356 loss weight (λ_2) 3, and L_2 -regularization weight (λ_3) 3.

357

358 ***Feature selection and scaling***

359 There are already multiple feature selection techniques available in a wide range of general
360 statistical packages and scRNA-seq packages. For this reason, *DEGAS* does not focus mainly
361 on feature selection, data cleaning, scRNA-seq clustering, but rather on transferring clinical traits
362 from patient to cells for the purpose of prioritizing those cells. For these reasons, a wide range of
363 feature selection techniques can be used before the *DEGAS* framework is applied.

364

365 Data from scRNA-seq experiments are generally very sparse. As a result, there are few genes
366 with viable expression for any given cell. Due to this, it is necessary to perform feature selection
367 to remove genes that are lowly expressed or have very low variance. When we select for high
368 variance and expressed genes in the bulk expression data, more genes are filtered out. After the
369 intersection of these two gene sets of expressed and high variance genes, we are left with less
370 than 1000 genes. It is worth noting that such number of gene features is comparable to Seurat
371 analysis, when usually hundreds to a couple of thousand highly variable genes are selected. The
372 feature selection steps were tailored to each dataset because the data sparsity and variance vary

373 greatly from one another, thus the tailored selection insured that enough genes with high enough
374 variability were available to train on. The feature selection steps are described individually in each
375 of the simulated, GBM, AD, and MM experiment sections.

376

377 For each experiment, the final feature scaling steps were consistent. The gene expression was
378 converted to sample-wise z-scores because it allows the genes to be more comparable between
379 samples and has been performed in multiple other studies [27, 48-50]. As the input to our deep
380 learning models, we scaled these z-scores to a range of [0,1]. This form of z-score scaling and
381 [0,1] scaling is commonly used in machine learning and deep learning to help model training [51-
382 53]. We follow this same convention for our deep learning models.

383

384 ***Disease association scores***

385 The final *DEGAS* output is either the output of a sigmoid or a softmax activation. For these
386 reasons, it can be useful to convert the [0,1] label output to an association score which can be
387 interpreted like a correlation coefficient. For these reasons, the output probability matrix from
388 *DEGAS* can be converted to a [-1,1] value using the *toCorrCoeff* function in the *DEGAS* package.
389 This function transforms the [0,1] output value matrix (P) with k labels to [-1,1] using **Eq. 12**.

$$390 \quad \text{disease association} = 2 \left(\frac{P - \frac{1}{k}}{2 - \frac{2}{k}} + \frac{1}{2} \right) - 1 \quad \text{Eq. 12}$$

391

392 ***Validating DEGAS using Simulated single cell data***

393 First we generated 5,000 single cells in four cell types where the cell type 4 had two subtypes
394 (cell type 4 disease and cell type 4 normal). Each of these five groups described above contains
395 1,000 cells. We split randomly these cells into 2 parts with 2,000 cells used for patient bulk tissue
396 data generation and 3,000 cells to use directly as single cell data. The 2,000 single cells used to
397 generate 600 patients across three different experiments (designated as simulation 1, 2, and 3)

398 where in simulation 1 the cell type 1 is associated with disease, in simulation 2 only the cell type
399 4 disease is associated with disease, and in simulation 3 the entire cell type 4 is associated with
400 disease. Each patient bulk tissue data was generated by randomly combining 400 single cells
401 using the proportions in **Table 3**.

402

403 **Table 3. Patient cellular makeup for simulation experiments.** The abbreviations are:

404 Simulation (sim), Normal (N), and Disease (D). The high-risk cell types are in bold.

	Cell type 1	Cell type 2	Cell type 3	Cell type 4N	Cell type 4D
Patients sim1D	50.0%	16.6%	16.6%	16.6%	00.0%
Patients sim1N	25.0%	25.0%	25.0%	25.0%	00.0%
Patients sim2D	25.0%	25.0%	25.0%	00.0%	25.0%
Patients sim2N	25.0%	25.0%	25.0%	25.0%	00.0%
Patients sim3D	16.6%	16.6%	16.6%	30.0%	20.0%
Patients sim3N	25.0%	25.0%	25.0%	25.0%	00.0%

405

406 We then performed 10-fold cross validation by training the *DEGAS* ClassClass models using cell
407 type and disease attribute. A total of 1000 gene features were used during training. We evaluated
408 the model capacity for mapping patient labels on patients and cell type labels on single cells using
409 PR-AUC and ROC-AUC. We then recapitulated the known cell type associations in each
410 simulation by overlaying disease association onto the simulated cells. As a comparison, we also
411 deconvoluted the patients using the 4 cell types using least squares. Deconvolution should be
412 able to correctly identify the cells of interest in simulation 1 and simulation 3. In contrast, cell type
413 prioritization using *Augur* [19] should be able to correctly identify the disease associated cell types
414 in simulation 2. In the simulation 1 *Augur* experiment, cell type 1, cell type 2, cell type 3, and cell
415 type 4 normal were randomly assigned to the disease or normal groups. In the simulation 2 *Augur*
416 experiment, cell type 1, cell type 2, and cell type 3 were randomly assigned to disease or normal
417 groups. The cell type 4 disease cells were all assigned to the disease group and the cell type 4
418 normal cells were all assigned to the normal group. In the simulation 3 *Augur* experiment, cell
419 type 1, cell type 2, and cell type 3 were randomly assigned to disease and normal groups. The

420 cell type 4 cells assigned to the disease group consisted of 60% cell type 4 normal and 40% cell
421 type 4 disease and cell type 4 cells assigned to the normal group consisted of 100% cell type 4
422 normal. These cell type proportions match those in the simulation 3 patients used by *DEGAS*.
423 The *Augur* output for each cell type is an ROC-AUC score that reflects how much a cell type
424 changes transcriptionally between disease and normal samples. To make the comparison fair
425 between our two methods, we use the output of our algorithm scaled from [0,1] where 0.5 implies
426 no association, 0 implies a negative association, and 1 implies a positive association. ROC-AUC
427 is on the same scale. In this way we compare the strength of signal between *Augur* and our
428 method to identify that cell type 4 has cell-intrinsic changes related to disease.

429

430 ***Validating DEGAS using GBM data***

431 The scRNA-seq data from the Patel *et al.* study [44] were downloaded from NCBI Gene
432 Expression Omnibus (GSE57872). The single cell expression values were previously normalized
433 to TPM containing 5,948 genes with $mean(\log_2(TMP)) > 4.5$ retained in the data table. The top 20%
434 variance genes were retained for training. These values were converted to z-scores then
435 standardized to a range of [0,1] for each sample. The TCGA GBM microarray expression data
436 was downloaded from *Firebrowse* (<http://firebrowse.org/>). Microarray data were used since it
437 contains more patient samples for training with GBM subtype information than RNA-seq data.
438 Likewise, the top 20% variance genes were retained for training and these expression values
439 were converted to z-scores then standardized to a range of [0,1] for each sample. The GBM
440 subtype labels for the TCGA patients were downloaded from Verhaak *et al.* [54]. The intersection
441 of genes between single cells and patients (199 genes) were used for the final model training.
442 Since subtype labels were only available for the GBM patient samples, we trained a BlankClass
443 *DEGAS* model (**Eq. 10**). This model minimizes the MMD loss between single cells and patients
444 while minimizing the classification loss only in GBM patients. We split the dataset into 10 groups
445 and performed 10-fold cross-validation by leaving out a single patient group during training. After

446 cross-validation, we converted the [0,1] *DEGAS* output to an association [-1,1] using the *DEGAS*
447 *toCorrCoeff* function. These association scores were overlaid on the GBM single cells and now
448 referred to as GBM subtype association scores because GBM subtype from patients is overlaid
449 on single cells. We plotted these association scores stratified by GBM subtype for each tumor
450 individually. We then compared the proportions of these cell types to the previously defined GBM
451 types from the original publication were marked with red boxes. We also visualized the GBM
452 subtypes association in single cells by calculating a low dimensional representation using *tSNE*
453 and overlaying the *kNN* smoothed GBM subtype associations. To make the scatter plots of cells
454 and patients more informative, *kNN* smoothing was used by averaging each point's GBM subtype
455 association value with its five nearest neighbors in *tSNE*. The model performance was shown with
456 the PR-AUC and ROC-AUC for each of the GBM subtype labels in the TCGA patients from cross-
457 validation.

458

459 In a second analysis on the GBM scRNA-seq and bulk expression data, using the same input
460 features, we overlaid risk derived from the overall survival in the TCGA GBM cohort onto the
461 individual cells from the Patel et al. study [44]. GBM has an extremely low 5-year survival rate
462 resulting only three patients being censored. We introduced more censoring in the data by
463 generating a uniformly distributed random vector of censoring times in the range 1 to 1063 days,
464 where 1063 days is the 90th percentile of survival times. If the censor time was lower than the
465 survival time, the patient was censored at that time instead of having an event at their true survival
466 time. We then trained 10 BlankCox *DEGAS* models based on the patient survival input during 10-
467 fold cross validation. The output from these *DEGAS* models were *kNN* smoothed based on the
468 *tSNE* coordinates using the *DEGAS knnSmooth* function and converted to death associations
469 using the *DEGAS toCorrCoeff* function. To highlight the differences in death association of cells,
470 these associations were centered to 0 using the *DEGAS centerFunc* function. We evaluated the

471 accuracy of the labels in patients using a rank-sum test based on the cox output in the GBM
472 patients.

473

474 ***Validating DEGAS and exploration using AD data***

475 For AD datasets, we were primarily interested in identifying known relationships between cell
476 types and AD diagnosis. For these reasons, we downloaded all of the adult Human scRNA-seq
477 data from the AIBS. Only inhibitory neurons, excitatory neurons, oligodendrocytes, astrocytes,
478 microglia, and oligodendrocyte progenitor cells (OPCs) were retained in the analysis due to the
479 extremely low sample sizes for the remaining cell types. The inhibitory and excitatory neuron
480 groups were merged into a single neuron group. These data were then \log_2 transformed,
481 converted to sample-wise z-scores, and then standardized to [0,1] by each sample. In the primary
482 analysis, only the top 50 up-regulated DEGs for each cell type (calculated by *Seurat*) were
483 retained in the single cell data (see **RESULTS**). In a distinct secondary analysis, features were
484 selected with >25% non-zero samples and top 20% variance genes (see **Supplementary**
485 **Materials**). The labels for the single cells consisted of the major cell types listed above. The AD
486 brain data was downloaded from Mount Sinai/JJ Peters VA Medical Center Brain Bank
487 (<https://www.synapse.org/#!/Synapse:syn3157743>). Each of the RNA-seq samples were either
488 from an AD patient's brain sample or a normal control brain sample. The binary disease attribute
489 of AD case or normal were used as the label for the model. Like in the previous experiment, the
490 RNA-seq values were \log_2 transformed, converted to sample-wise z-scores, and standardized to
491 [0,1] for each sample. The top 50% variance genes were retained for training to keep the feature
492 set larger. The intersection of the patient genes and single cell genes (Primary analysis: 169
493 genes, Secondary analysis: 456 genes) were using to train the final models. Using the cell type
494 classification for each AIBS single cell and the AD/normal classification for each MSBB patient
495 we were able to train a *DEGAS* ClassClass model (**Eq. 8**). The performance was evaluated using
496 10-fold cross-validation by leaving out each group during training once. As in the GBM

497 experiments, we converted the *DEGAS* output to an association using the *DEGAS toCorrCoeff*
498 function for each single cell so that each single cell now had an AD association. Correlation
499 analysis was performed on AD association scores for different cells with each cell type by taking
500 the median score and calculating the p-value by treating it as a correlation. In addition, single cells
501 were plotted overlaid with *kNN* smoothed AD association. Furthermore, to evaluate *DEGAS*
502 performance, PR-AUC and ROC-AUC were computed for the single cells during cross-validation
503 for each cell type in the single cell data. Similarly, AD diagnosis PR-AUC and ROC-AUC were
504 computed from the MSBB patient RNA-seq. For both the primary and secondary AIBS analysis,
505 DEGs were identified for the high AD association astrocytes and microglia based on the median
506 AD association then compared to their respective disease associated astrocyte (DAA) [55]
507 (**Supplementary File 1**), human Alzheimer's microglia (HAM) gene markers [56]
508 (**Supplementary File 2**), or disease associated microglia (DAM) gene markers [57]
509 (**Supplementary File 3**). A detailed description of these gene lists can be found in the
510 **Supplementary Materials DAA, HAM, and DAM markers** section.

511
512 To further highlight the cellular associations to AD, we also performed experiments using a
513 scRNA-seq dataset from Grubman *et al.* [32]. Since this dataset was sparser, genes were used
514 with >25% non-zero samples then the top 50% variance genes were selected from these. For the
515 MSBB data, the same initial feature selection was used (top 50% variance). The same
516 normalization and standardization procedure as the AIBS scRNA-seq and MSBB were used
517 again. The intersecting genes between Grubman *et al.* scRNA-seq constituted the final feature
518 set (61 genes). 10-fold cross validation was performed using a ClassClass model and the AD
519 associations were overlaid onto the Grubman *et al.* scRNA-seq in the same fashion as the
520 previous experiment. In addition, a targeted analysis on only the microglia cells was performed.
521 A single BlankClass model was trained using the same 61 features on the entire Grubman *et al.*
522 microglia scRNA-seq and MSBB RNA-seq. For both analyses, the AD associations were overlaid

523 onto the cells, AD associations were compared between cells from AD and normal patient
524 samples, and DEGs were identified for the high AD association astrocytes and microglia based
525 on the median AD association then compared to their respective DAA [55] (**Supplementary File**
526 **1**), HAM gene markers [56] (**Supplementary File 2**), or DAM gene markers [57] (**Supplementary**
527 **File 3**). For the targeted analysis on only microglia, correlation tests were performed between AD
528 associations and HAM gene markers [56] (**Supplementary File 2**). Also, DEGs were identified
529 for the high AD association microglia based on the median AD association then compared to the
530 HAM gene markers [56] (**Supplementary File 2**) and DAM gene markers [57] (**Supplementary**
531 **File 3**).

532
533 Lastly, *DEGAS* analysis was performed on the Mathys *et al.* scRNA-seq dataset [15]. In this
534 analysis, the same gene set as the AIBS Primary analysis, *i.e.*, all overlapping genes (157 genes)
535 were used as input features. 1000 cells or all cells if total number was less than 1000 were
536 sampled from each cell type since some cell types were over-represented. The same
537 normalization and standardization procedure was used as the previous analyses. 10-fold cross
538 validation was performed using these cells from Mathys *et al.* and the MSBB patient RNA-seq
539 data using cell type and patient AD status as outcomes respectively. These outcomes represent
540 a ClassClass *DEGAS* model. From the cross-validation results, the ROC-AUCs and PR-AUCs for
541 each cell type label and the patient AD status were calculated. AD associations were calculated
542 in the same fashion as all previous analyses. The Disease associations were then compared with
543 AD status of the scRNA-seq donors and across the cell types. DEGs were identified for the high
544 AD association astrocytes and microglia based on the median AD association then compared to
545 their respective DAA [55] (**Supplementary File 1**), HAM gene markers [56] (**Supplementary File**
546 **2**), or DAM gene markers [57] (**Supplementary File 3**).

547

548 ***Preprocessing of MM scRNA-seq***

549 The scRNA-seq data were first combined into a dataset using *Seurat-CCA* [28]. This initial dataset
550 integration allowed conserved subtypes of cells to be identified across datasets. All four patient
551 dataset counts were loaded into a *Seurat* object. *Seurat* normalized, scaled, removed poor quality
552 cells, and identified high variance genes. Using the union of high variance genes, multi-canonical
553 correlation analysis was run across all four datasets, the subspaces were aligned across patients,
554 the aligned single cells were plotted with *tSNE* [58], and clusters of cells were identified. The raw
555 expression values for the high variance genes identified by *Seurat* were \log_2 transformed,
556 converted to z-scores, and then scaled to [0,1].

557

558 Furthermore, each IUSM scRNA-seq patient was individually clustered using *Seurat* to check the
559 replicability of the clusters and were plotted with *UMAP* [59]. We used Rand, Fowlkes and
560 Mallows's index (FM), and Jaccard index (JI) to measure the cluster consistency between single
561 patient clustering experiments and the merged all-patient clustering results. The four single
562 patient clustering results, one for each IUSM scRNA-seq patient, were used as input into
563 *BERMUDA* [25] to visualize and evaluate the original *Seurat* clustering.

564

565 ***Preprocessing of MMRF patient data***

566 MMRF patients with bulk tissue RNA-seq and clinical data were used in MM analysis. We used
567 PFS as the disease attribute of interest. TPM values for the MMRF patient gene expression data
568 and the PFS data were used as the input for *DEGAS*, these values were \log_2 transformed,
569 converted to z-scores, and scaled to [0,1]. The union of the features (502 genes) identified by
570 *Seurat* in the single cell data and the features selected in the MMRF patient data were used as
571 the final feature set. The features retained in the MMRF data were identified by fitting an elastic-
572 net Cox model [60] to the TPM values based on the PFS.

573

574 ***Evaluate DEGAS performance on MM datasets***

575 PR-AUC and AUC were calculated for each of the output labels for the single cells and for patient
576 labels if a classification output was used for the patient data. Cox proportional hazard output was
577 used on patients, a log-rank test was calculated for each patient so that the hazard ratio and p-
578 value could be evaluated based on patient stratification by median proportional hazard.
579 Additionally, the same models were used to predict risk in the GSE2658 dataset which had
580 information on OS. The output for each GSE2658 sample averaged across all 10 *DEGAS* models
581 and stratified by median risk to show the robustness of the cox output across datasets.

582

583 ***Identification of CD138+ cell types associated with MM prognosis***

584 The single cells from MM patients can be assigned proportional hazards based on the MMRF Cox
585 output of the model. Each single cell in the validation set was assigned progression association
586 by feeding those samples through the Cox output layer. In this way, we can infer the association
587 with progression risk of specific cell types as well as the cell type enrichment contained in each
588 MMRF sample. Since the Cox output is a proportional hazard, we centered the outputs to zero
589 for each step of cross validation to produce a PFS association using the *DEGAS centerFunc*. We
590 plotted these relationships and conducted Student's t-tests on the subtype vs. PFS association in
591 IUSM single cells, PFS association vs. MM malignancy from Ledergor *et al.*, and subtype 2
592 enrichment vs. MM malignancy from Ledergor *et al* [42].

593

594 ***Analysis of differential gene expression in prognostic cell types***

595 T-tests were calculated cell subtype 1 vs all cell subtypes and cell subtype 2 vs. all cell subtypes
596 using the batch corrected gene expression values from *Seurat*. These values were stored in
597 **(Supplementary File 4 and Supplementary File 5)** respectively. For the marker set of *PHF19*,
598 *HELLS*, *EZH2*, *TYMS*, *ZWINT*, and *MKI67* we performed t-tests for each patient individually.

599

600 ***Evaluation of DEGAS robustness to hyper-parameters in GBM***

601 Using the GBM dataset, we evaluated the robustness of DEGAS model outputs to hyper-
602 parameters by repeating 10-fold cross-validation 100 times with randomly generated hyper-
603 parameters following a uniform distribution. The range of hyper-parameters used in training
604 consisted of training steps 1,000-3,000, single cell batch size 100-300, patient batch size 20-100,
605 hidden features 10-100, drop-out retention rate 0.1-0.9, Cell loss weight (λ_0) held at constant 2,
606 Patient loss weight (λ_1) 0.2-5, MMD loss weight (λ_2) 0.2-5, L_2 -regularization weight (λ_3) 0.2-5.

607

608 Using these outputs, we performed two tests. One was to evaluate the loss in performance based
609 on changing the hyper-parameters where performance was measured with ROC-AUC among the
610 TCGA GBM patients labeled by patient GBM subtype (Mesenchymal, Classical, Proneural,
611 Neural). In this test, we calculated the spearman correlation and plotted the scatter plot between
612 the AUC of each of the four GBM subtype labels and the hyper-parameters used.

613

614 Next, we evaluated whether or not the correct GBM subtype labels (Mesenchymal, Classical,
615 Proneural, Neural) could be recapitulated in the GBM scRNA-seq tumors that had known GBM
616 subtypes (MGH26: Proneural, MGH28: Mesenchymal, MGH29: Mesenchymal, MGH30:
617 Classical). To do this for each tumor (MGH26, MGH28, MGH29, MGH30), the rank of the
618 correct label was calculated by calculating the mean of each GBM subtype association across
619 all of the cells in that tumor. This resulted in each of the 100 random hyper-parameters having a
620 rank for each GBM subtype for each of the GBM scRNA-seq tumors (4 highest ranked, 1 lowest
621 ranked). Ideally all GBM scRNA-seq tumors would have a rank of 4 indicating the correct GBM
622 subtype was ranked the highest regardless of hyper-parameters. Similarly, we also calculated
623 the spearman correlation and plotted the scatter plot between correct label rank and the hyper-
624 parameters used.

625

626 ***Evaluation of domain adaptation for DEGAS disease association transfer***

627 We evaluate the necessity for domain adaptation to transfer disease associations to single cells
628 using 30 total experiments. These experiments evaluated disease associations in cells by
629 training with MMD loss vs. those without MMD loss for a variety of biases added between the
630 cells and patients. It is important to highlight the fact that without bias between different
631 datasets, in this case cells and patients, there is no need for domain adaptation. Practically in
632 real transcriptomic data, there will always be bias between datasets. For these reasons we
633 added bias for these 30 experiments. These experiments were conducted for every combination
634 of MMD loss (with and without MMD), simulation (three simulations), and cellular subtype (five
635 total subtypes since cell type 4 has two subtypes) totaling 30 combinations. The experiments
636 were conducted as follows. In each experiment, the counts of 300 cells from a given subtype
637 were aggregated together and multiplied by 1000 constituting a large systematic bias
638 associated with a single subtype. This bias vector was added to all of the patients in the given
639 simulation, both disease and normal. A single three-layer DenseNet *DEGAS* model with five-fold
640 bootstrap aggregation was trained on all the cells and all the patients then the disease
641 associations were predicted in the cells. We evaluated error by subtracting the expected
642 disease association from the predicted disease associations, e.g., cell type 1 in simulation 1
643 should be 1. We then compared the error rates between the *DEGAS* models with and without
644 MMD using a t-test.

645

646 ***Evaluation of regularization in DEGAS performance***

647 Regularization is an important method in machine learning to prevent model overfitting. Here we
648 utilized three such techniques to prevent overfitting, namely, L_2 -regularization, dropout, and
649 bootstrap aggregation. Since all of these techniques may work better or worse in different
650 scenarios, we perform a simple experiment where all of these regularization techniques are
651 removed and compared with the regularized results. We performed experiments using each of
652 the simulated datasets. To evaluate the robustness of our models we performed 10-fold cross

653 validation in each simulation. The simulated cells were split into 10 groups and the simulated
654 patients were split into 10 groups. For each fold of cross validation, our default *DEGAS* three-
655 layer DenseNet model with L_2 -regularization, dropout, and bootstrap aggregated 5 times was
656 trained then a three-layer DenseNet *DEGAS* model was trained on the same data without L_2 -
657 regularization, dropout, and bootstrap aggregation. Both models were then used to predict the
658 patient disease attributes in the holdout group of patients, the cell types in the hold out group of
659 cells, and the patient disease attributes in the cells. We compare the performances using ROC-
660 AUC and PR-AUC for patient disease status in patients and cell type in cells. Furthermore, we
661 evaluate the label transfer of patient labels to cells by calculating the error based on the
662 expected cell type association for each cell. We compare between the regularized and
663 unregularized error in cells with a t-test.

664

665 **Results**

666 ***DEGAS clinical impression framework***

667 In this study, we applied *DEGAS* to integrate and analyze scRNA-seq, bulk gene expression, and
668 clinical data (**Fig. 1**) from simulated data as well as three different diseases: GBM, AD, and MM.
669 The simulated, GBM, and AD datasets primarily served as validation to demonstrate the feasibility
670 and universality of the *DEGAS* transfer learning approach since the ground truth of the simulated
671 data was known, the correct GBM subtypes were known, and neuron loss with microglia gain in
672 AD brains were also known. We then further expand our study to MM data, which serves as the
673 discovery dataset, since the myeloma cell subtypes and high-risk factors related to MM are not
674 as well understood at the single-cell level. In the MM study, we applied *DEGAS* on patient data
675 from the Multiple Myeloma Research Foundation CoMMpass study (MMRF) and scRNA-seq data
676 that we generated from myeloma patients. Our aim was to identify the cell subtypes using the
677 impressions of progression risk on the single cells. We then applied the results to two separate
678 MM validation datasets, one of which contained plasma cells from normal bone marrow (NHIP),

679 two MM precursor conditions - monoclonal gammopathy of undetermined significance (MGUS)
680 and smoldering multiple myeloma (SMM), and MM. We tested if *DEGAS* assignment of
681 progression risk to cell subtypes were higher for more malignant conditions. An additional external
682 validation dataset of patient level expression data with OS was used to evaluate whether the
683 patient stratification learned by *DEGAS* was robust enough to be generalized to an external
684 survival dataset.

685

686 ***DEGAS correctly identifies high-risk cell types and subtypes in simulated data***

687 To evaluate *DEGAS* in a controlled context, 5,000 single cells were generated with *Splatter* [61]
688 (**Fig. 2A**) where 2,000 of the cells were held-out to generate simulated patients. Using this group
689 of held-out cells, 600 simulated patients were generated by aggregating sets of 400 simulated
690 cells (**Fig. 2B-D**). We conducted three simulation experiments, denoted Simulation 1, Simulation
691 2, and Simulation 3, where the single cells were aggregated in known proportions for each patient
692 so that we could generate a “disease” patient group with different cellular composition than the
693 “normal” patient group (**see Methods**). To highlight the utility of *DEGAS*, the experiments were:
694 Simulation 1: cell type 1 is enriched in disease patients (**Fig. 2B**); Simulation 2: one subtype of
695 cell type 4, *i.e.*, cell type 4 disease, is enriched in disease patients (**Fig. 2C**); and Simulation 3:
696 both subtypes of cell type 4 are enriched in disease patients (**Fig. 2D**).

697

698 Please note that the optimal number of clusters for the simulated single cells would be
699 determined to be four based on a standard scRNA-seq workflow (*i.e.*, *tSNE* followed by *K-*
700 *Medoids* where optimal cluster number is selected based on average silhouette width) (**Fig. 2E**).
701 This would cluster the cells into the four cell types while ignoring the two subtypes in the cell
702 type 4 (**Fig. 2F**). As a result, deconvolution algorithms will not be able to detect the subtype
703 level risk associations. Fortunately, cell type prioritization algorithms like *Augur* can detect these
704 changes within cell types due to disease. However, for situations that do not have a new cell

705 type or missing cell type in the disease (simulation 1), *Augur* cannot detect the association
706 between cell type 1 and disease since there is no disease associated cell type change (**Fig.**
707 **2G**). *Augur* can detect the disease-associated cell type 4 in simulation 2 (**Fig. 2H**). In simulation
708 3 where there is a mix of disease and normal subtypes for the cell type 4 in the disease group,
709 *Augur* again has difficulty in identifying the cell type 4 disease association (**Fig. 2I**). In contrast
710 to *Augur*, deconvolution can easily identify the correct cell type for Simulation 1 (**Fig. 2J**) and
711 Simulation 3 (**Fig. 2L**) but not for simulation 2 (**Fig. 2K**). In comparison, *DEGAS* not only
712 identified the correct cell type and subtypes in each experiment, it also correctly detected all of
713 the simulated disease associations (**Fig. 2G-L**). Additionally, *DEGAS* had high precision-recall
714 area under the curve (PR-AUC) predicting disease status of simulated patients (0.96-0.98)
715 (**Table S1**) and almost perfectly predicted the cell type of simulated cells (~1.0) during cross-
716 validation (**Table S2**). Since *DEGAS* directly assigns disease risk to cells, many of the problems
717 with cell type level analyses can be avoided and the correct groups of cells can be identified by
718 overlaying impressions of disease risk.

719

720 ***DEGAS correctly mapped single cells to corresponding GBM subtypes***

721 We first demonstrate *DEGAS* in a straightforward case to show the performance of our framework
722 using real data from GBM. We use single-cell data from Patel *et al.* [44], in which researchers
723 assigned four major GBM tumor subtypes (Proneural, Mesenchymal, Classical, and Neural) to
724 the scRNA-seq data obtained from five GBM tumors. Of the five tumor samples, four had been
725 labeled in the original publication with a single subtype based on the major proportion of cells
726 assigned to each GBM subtype. For GBM bulk tumor tissue expression data, we obtained
727 microarray data for 111 GBM patients from The Cancer Genome Atlas (TCGA), for which the
728 same labels of GBM subtypes were also provided. The OS was also available in a subset of 109
729 patients. As the simplest form of validation, we used these two datasets as input for the *DEGAS*
730 model to test if it could re-identify the same GBM subtypes for both single cells and the TCGA

731 GBM cohort simultaneously. Then we overlaid OS-derived death associations onto the cells to
732 visualize their association with OS. The resulting *DEGAS* models also proved to be accurate with
733 high PR-AUCs (0.79-0.97) when predicting each of the GBM subtypes in the TCGA patients
734 during 10-fold cross validation (**Table S3**). The OS BlankCox *DEGAS* models were able to stratify
735 the patients into high and low risk groups based on median patient risk (log-rank p-value < 0.05).
736 *DEGAS* correctly re-identified the same labels for all four tumors by overlaying GBM subtypes
737 associations on each single cell, as indicated by the groups of cell subtypes with the highest
738 association score determined by the median value (indicated with a red box) (**Fig. 3A-D**). For the
739 fifth tumor sample, MGH31, it was labeled as a combination of multiple GBM subtypes in the
740 original study, so we did not use it in our evaluation although *DEGAS* identified mesenchymal as
741 its most associated GBM subtype (**Fig. 3E**). Additionally, these relationships can be visualized by
742 plotting the single cells and overlaying the GBM subtype association or OS-derived death
743 association. It is clear that MGH28 and MGH29 have a high association with the mesenchymal
744 GBM subtype (**Fig. S1A**) and contain populations of cells with high death associations (**Fig. 3F**).

745

746 ***DEGAS identifies increased microglia, reduced neuron populations, DAAs and DAMs***

747 Aside from GBM, AD also has well documented characteristics that can be used as a test bed for
748 *DEGAS*. Specifically, there is a well-documented reduction in neurons [36-38], increase in
749 microglia [33-35, 39], and more recently, AD subtypes of astrocytes [55] and microglia [56, 57].
750 AD brain scRNA-seq data was obtained from the AIBS and bulk AD RNA-seq were retrieved from
751 MSBB [46]. During 10-fold cross-validation, *DEGAS* models for both primary and secondary AIBS
752 analyses achieved high AD diagnosis status PR-AUC (0.82 and 0.76) in MSBB patients (**Table**
753 **S4**) and high cell type prediction PR-AUCs (>0.99) for AIBS single cells (**Table S5**).

754

755 From the AIBS primary analysis *DEGAS* results, we confirmed that at the single cell level, the
756 AD associations were negative in neurons as previously described [62], which is shown by the

757 dark shade of neurons compared to other cell types (**Fig. 3G, Table 4**). In contrast, we
758 observed positive AD associations in microglia cells (**Fig. 3G, Table 4**). A strength of the
759 *DEGAS* framework is that it can detect intra-cell type differences in disease risk. Within the
760 astrocyte cell type, we identified an astrocyte subtype that had a positive association with AD
761 (**Fig. 3G**) that corresponded to the *Astro L1 FGFR3 FOS* subtype from the AIBS brain cell atlas
762 (*i.e.*, *FOS* is a DAA marker) [63] (**Fig. 3H**), had up-regulated DAA marker *GFAP* (**Fig. 3I**) [55],
763 and was enriched for DAA markers (OR = 30.93, Fisher's exact p-value < $2.2 \cdot 10^{-16}$, **Table S6**).
764 Furthermore, the high AD association microglia were enriched for DAM markers (OR = 17.07,
765 Fisher's exact p-value = $2.11 \cdot 10^{-10}$, **Table S7**). In the secondary AIBS analysis using high
766 variance genes, we again identified the strong negative AD association for neurons (**Fig. S2A**),
767 positive AD association in microglia (**Fig. S2A**), high AD association astrocytes enriched for
768 DAA markers (OR = 5.65, Fisher's exact p-value < $1.66 \cdot 10^{-8}$, **Fig. S2B,C, Table S8**), and high
769 AD association microglia enriched for DAM markers (OR = 14.34, Fisher's exact p-value <
770 $4.01 \cdot 10^{-11}$, **Table S9**). When we performed *DEGAS* analysis on a separate dataset from
771 Grubman et al. [32] with single cells from both AD and normal brains, we found that the major
772 cell types from AD brains were significantly more associated with AD than their counterparts in
773 normal brains as judged by median value (**Fig. 3J**).

774

775 **Table 4 Comparison of AD association scores in single cells between cell types as**
776 **visualized in Fig. 3G.** The *DEGAS* models were trained using neuron, oligodendrocyte,
777 astrocyte, OPC, and microglia cell types. The single cells were split into groups based on their
778 cell type and the mean AD associations of each cell type was evaluated as a correlation. The
779 neuron and microglia groups are bolded to highlight their much higher mean AD association. P-
780 values are calculated by treating the association score as a pearson correlation coefficient.

Cell type	Cell-type mean association	Number of cells	p-value
-----------	----------------------------	-----------------	---------

Neuron	-0.35	1329	$<2.2 \cdot 10^{-16}$
Oligodendrocyte	0.05	1795	$3.42 \cdot 10^{-2}$
Astrocyte	0.03	809	$3.94 \cdot 10^{-1}$
OPC	-0.12	738	$1.09 \cdot 10^{-3}$
Microglia	0.22	741	$1.42 \cdot 10^{-9}$

781

782 In the Grubman *et al.* scRNA-seq data, the astrocytes in AD brains were highly positively
783 associated with AD (AD association = 0.22, pearson correlation p-value = $7.89 \cdot 10^{-7}$) whereas
784 the astrocytes in normal brains were negatively associated with AD (AD association = -0.06,
785 pearson correlation p-value = $1.06 \cdot 10^{-2}$, **Fig. 3J**). Astrocytes from AD brains also expressed
786 *GFAP* at greater levels than astrocytes from normal brains (t-test p-value $< 2.20 \cdot 10^{-16}$) and high
787 AD association astrocytes were significantly enriched for DAA markers (OR = 21.90, Fishers
788 exact p-value = $2.21 \cdot 10^{-12}$, **Table S10**). Furthermore, the high AD association microglia were
789 moderately enriched for DAM markers (OR = 4.15, Fishers exact p-value = $4.11 \cdot 10^{-2}$, **Table**
790 **S11**). This provides evidence for DAA and DAM cells in the Grubman *et al.* dataset.

791

792 DAM and HAM marker enriched high AD association cells were independently identified in the
793 targeted analysis of the Grubman *et al.* microglia cells (**Fig. 3K**). AD associations were higher in
794 cells derived from AD patient samples than Normal patient samples (**Fig. 3L**, t-test p-value =
795 $6.66 \cdot 10^{-12}$), HAM up-regulated markers were more likely to be significantly positively correlated
796 to AD association than HAM down-regulated markers (**Fig. 3M**, t-test p-value = $2.63 \cdot 10^{-3}$). The
797 HAM marker APOE [56, 57] was positively correlated with AD association (**Table S12**,
798 PCC=0.18, p-value = $1.15 \cdot 10^{-4}$). High AD association microglia were significantly enriched for
799 HAM markers (OR = 21.47, Fishers exact p-value = $6.33 \cdot 10^{-4}$, **Table S13**) and DAM markers
800 (OR = 11.52, Fishers exact p-value = $1.50 \cdot 10^{-11}$, **Table S14**). It is important to note that there
801 was no overlap between the input feature set used to train the DEGAS model and HAM marker

802 genes that were identified, which shows DEGAS is a useful tool to identify disease associated
803 cells within a single cell type even without prior knowledge of marker genes.

804

805 After applying *DEGAS* to the Mathys *et al.* scRNA-seq dataset, the *DEGAS* models achieved
806 high AUCs for patient AD status (0.77), patient AD status PR-AUC (0.81), cell types (>0.98), as
807 well as cell type PR-AUCs (0.82-0.99) during cross validation (**Table S15-16**). The positive AD
808 association of microglia and negative AD association of neurons were recapitulated (**Fig. S3A**,
809 **Table S17**). Within the astrocyte cluster, there existed a subset of astrocytes with higher AD
810 association (**Fig. S3B**). High AD association astrocytes were significantly enriched for DAA
811 markers (OR = 14.75, Fishers exact p-value = $3.16 \cdot 10^{-15}$, **Table S18**). A closer comparison of
812 the scRNA-seq revealed that the top 10% AD association astrocytes, had 2.5 times higher
813 GFAP expression than the other astrocytes (t-test p-value = $6.36 \cdot 10^{-8}$, **Fig. S3C**). In fact, like
814 the Grubman *et al.* analysis, the AD association scores were higher in cells coming from AD
815 patients than normal patients for every cell type in the Mathys *et al.* analysis (**Table S17**).

816 Notably, we see increased AD association in AD derived astrocytes and microglia likely
817 representing DAAs and HAMs respectively (**Table S17**). Furthermore, high AD association
818 microglia were highly enriched for DAM markers (OR = 19.35, Fishers exact p-value < $2.2 \cdot 10^{-16}$,
819 **Table S19**) and high AD association in astrocytes correlated well with neuritic plaque count, a
820 marker for disease severity in AD patients (PCC = 0.22, p-value = $6.36 \cdot 10^{-12}$, **Table S17**).

821 Again, the Mathys *et al.* analysis provides another example to demonstrate that *DEGAS*
822 recapitulates the findings from the AIBS and Grubman *et al.* analyses and shows that *DEGAS*
823 models can capture cell type level as well as intra-cell type differences in disease association.

824

825 ***Identification of plasma cell subtypes in CD138+ scRNA-seq of MM***

826 In the MM study, unlike the previous two datasets, there were no predefined cell type labels, but
827 *DEGAS* was still capable of analyzing such data and give clinical perspective to the clusters of

828 cells in the MM scRNA-seq data. In order to cluster cells into groups, we first used Seurat [28], a
829 commonly used scRNA-seq data analysis tool, to merge and cluster all the CD138+ bone marrow
830 cells from four patients (two SMM and two MM) whose samples were collected at the IUSM. Using
831 *Seurat*, five major clusters of cells were identified (**Fig. 4A**). Cluster 1 consisted of the majority of
832 the cells in each sample and was most likely the main clone in each of the patients. Cluster 2 was
833 present in many of the patients and is described in detail after the *DEGAS* analysis. Cluster 3 and
834 5 were only present in patient 2 representing possible subclones in patient 2. Cluster 4 was shared
835 between multiple patients. These five clusters were used as the subtype labels in the *DEGAS*
836 framework. We verified these cell clusters by clustering cells from each patient individually with
837 *Seurat* and another scRNA-seq normalization tool *BERMUDA* (Batch Effect ReMoval Using Deep
838 Autoencoders) [25] for all four patients. We found that the individual clustering results closely
839 mirrored the *Seurat*-CCA clusters (**Fig. S4A-D, Table S20**) and that the subtype 2 was consistent
840 across all MM patients using *BERMUDA* (**Fig. S4E**). For bulk tissue data from MMRF, the clinical
841 outcomes of PFS for 647 patients were used as the patient-level input to *DEGAS* and overlaid
842 onto the CD138+ single cells from the four IUSM patients (**Fig. 4B**).

843

844 ***DEGAS patient stratification and cell type classification on MM***

845 A *DEGAS* model was trained on IUSM patient scRNA-seq data with subtype labels defined above
846 and MMRF patients with bulk tissue data and PFS information. The performance metrics were
847 calculated via 10-fold cross-validation. It is worth noting that for PR-AUC, random no skill
848 classifiers will achieve a performance equal to the percentage of the class of interest and in the
849 case of uncommon classes like subtype 4, the random classifier performance will be close to zero
850 (0.02). When predicting cellular subtype label in single cells, *DEGAS* was able to achieve a PR-
851 AUC between 0.44-0.98 for all of the five CD138+ cellular subtypes identified in the above scRNA-
852 seq data while the PR-AUC for subtype 2 reached 0.91 (**Table S21**). The receiver operating curve
853 AUCs (ROC-AUCs) were between 0.90-0.98 for these five subtypes (**Table S21**). Due to class

854 imbalance some of the subtypes did not perform as well as others based on PR-AUC but all of
855 the PR-AUCs were substantially greater than a purely random model. Aside from correctly
856 classifying the single cells, *DEGAS* was able to stratify the MMRF patients into high and low risk
857 groups based on median progression risk (log-rank p-value = $4.72 \cdot 10^{-10}$, **Fig. 4C**). We then
858 applied the trained model on an external patient transcriptomic dataset from Zhan *et al.* [41] for
859 validation. We demonstrated that the Cox proportional hazards portion for patient OS time of the
860 *DEGAS* model was robust across datasets, and the impressions extracted from the *DEGAS*
861 framework were capable of stratifying patients into low- and high-risk groups in the validation
862 dataset (log-rank p-value = $1.12 \cdot 10^{-3}$, **Fig. 4D**).

863

864 ***DEGAS identifies CD138+ cellular subtypes with high progression association***

865 The MM scRNA-seq data provided an example of an exploratory analysis with *DEGAS* which
866 can be used to generate hypotheses for future studies. The *DEGAS* model for the MM study
867 transfers clinical impressions to single cells (*i.e.*, single cells were directly assigned a
868 progression association score), as well as transfers cellular/molecular impressions to patients
869 (*i.e.*, patients are assigned subtype enrichment score). We found that the subtype 2 cells were
870 the most associated with prognosis (**Fig. 4B**) based on the *DEGAS* results. Specifically, the
871 subtype 2 cells were associated with a shorter time to progression (**Fig. 4E**, t-test p-value <
872 $2.2 \cdot 10^{-16}$). On an external validation scRNA-seq dataset from Ledergor *et al.* [42], the
873 progression association increased from NHIP (no disease) to SMM (**Fig. 4F**, t-test p-value =
874 $1.50 \cdot 10^{-2}$) and MM (**Fig. 4F**, t-test p-value = $1.70 \cdot 10^{-2}$), which is consistent with the order of
875 precursor conditions for MM (NHIP → MGUS → SMM → MM). In addition, the enrichment of the
876 subtype 2 cells increased from NHIP to near-MM stage SMM (**Fig. 4G**, t-test p-value = $3.10 \cdot 10^{-2}$)
877 and MM (**Fig. 4G**, t-test p-value = $3.40 \cdot 10^{-2}$).

878

879 ***MM prognostic subtypes have distinct gene signatures***

880 Differential gene expression analysis was performed between subtype 2 and all other subtypes
881 (**Supplementary File 5**), and we found that subtype 2 had significantly up-regulated *PHF19*
882 expression in all four of the patients (**Fig. 4H**). *PHF19* is a known marker for malignant disease in
883 MM [64]. Besides *PHF19*, its associated markers such as *HELLS*, *EZH2*, *TYMS*, *ZWINT*, and
884 *MKI67* were also significantly up-regulated in subtype 2. These results suggested the possible
885 existence of a more malignant CD138+/*PHF19*^{high} subpopulation of plasma cells represented by
886 the subtype 2 cluster. It is important to notice that the gene feature set that was used as input into
887 *DEGAS* only contained the *HELLS* gene, which further highlights the ability of *DEGAS* to predict
888 high-risk cellular subtypes that can be further studied.

889

890 ***DEGAS is robust to hyper-parameter choice***

891 To assess the robustness of *DEGAS*, we also analyzed how the hyper-parameter choices
892 influence its results using a set of 100 randomly generated hyper-parameters with 10-fold cross-
893 validation on each set of those 100 sets of hyper-parameters on the GBM datasets. The hyper-
894 parameters that we evaluated include: the number of training steps, batch size for single cells,
895 batch size for patients, number of hidden layer nodes, drop-out retention rate (the percentage of
896 nodes randomly retained at the hidden layer), patient loss weight, MMD loss weight, and L_2 -
897 regularization weight. The detailed information about the range of hyper-parameters that were
898 randomly sampled can be found in subsection titled *Evaluation of DEGAS robustness to hyper-*
899 *parameters* in the **Methods** section while the default parameters used for all previous experiments
900 can be found in the subsection titled *Transfer learning using DEGAS*. We discovered that among
901 the eight hyper-parameters, the majority of them did not significantly affect the ROC-AUC for
902 predicting GBM subtypes in TCGA GBM patients with the exception of three hyperparameters –
903 namely the drop-out retention rate, number of hidden layer nodes, and L_2 -regularization weight
904 with spearman correlation p-value < 0.1 (**Fig. S5, Table S22**). Similarly, the majority of hyper-
905 parameters did not significantly affect the correct assignment of subtype to GBM scRNA-seq

906 tumor, except for a few exceptions in training steps, patient loss weight, and MMD loss weight
907 with spearman correlation p-value < 0.1 (**Fig. S6, Table S22**). We therefore suggest users to keep
908 default settings for at least patient loss weight, MMD loss weight, and L_2 -regularization weight.
909 The percentage of GBM subtype labels ranking in the top two predicted labels improves from 74%
910 to 82% if the default parameters or greater values are used for patient loss weight, MMD loss
911 weight, and L_2 -regularization weight (**Fig. S7**).

912

913 ***Domain adaptation improves DEGAS disease association transfer***

914 Without any bias, MMD and no MMD performances were not different from one another. After
915 bias was added, MMD did improve the ability of *DEGAS* to transfer disease associations onto
916 cells (**Fig. S8**). MMD is important for our algorithm because the bias added to the patients
917 represents the types of systematic bias that are present between bulk and single cell
918 transcriptomic data. In the example of simulation 2 with cell type 2 bias added, it is clear that all
919 of the patients tended to cluster adjacent to the cell type 2 cluster (**Fig. S8A**). We defined high-
920 risk cells in this example as cells with a disease association >0.2 on a $[-1, 1]$ scale. Once the
921 *DEGAS* model had been trained and the disease associations overlaid onto the cells, the
922 *DEGAS* model trained without MMD predicted many cells in cell types other than cell type 1 as
923 being high-risk (**Fig. S8B**). In contrast, the *DEGAS* model trained with MMD only identified cell
924 type 1 cells opposed to other cell types as high-risk (**Fig. S8C**). Over all 30 experiments, we
925 found that the disease association error was lower in the *DEGAS* models with MMD than in the
926 *DEGAS* models without MMD (t-test p-value $< 2.2 \cdot 10^{-16}$, **Fig. S8D**). When the cells were
927 ordered by their error, there was no experiment where the *DEGAS* model without MMD
928 consistently outperformed the *DEGAS* model with MMD (Kolmogorov-Smirnov p-value $< 2.2 \cdot 10^{-16}$,
929 **Fig. S8E**).

930

931 ***Regularization improves the robustness of DEGAS models***

932 Regularization is an important part of the *DEGAS* model which prevents the data from being
933 overfit. Without regularization, *DEGAS* models perform worse during cross-validation (**Table.**
934 **S24-25, Fig. S9**). There is no case where an unregularized model performed better than a
935 regularized model in predicting patient labels during cross validation. Specifically, in simulation
936 3, the unregularized models performed 5% worse in PR-AUC when predicting patient labels in
937 patients (**Table. S24**). Similarly, the unregularized *DEGAS* models performed 9% worse in PR-
938 AUC when predicting cell type labels in cells (t-test p-value = $4.34 \cdot 10^{-3}$, **Table. S25**). The
939 regularization also improved the transfer of disease associations to the cells in 2/3 simulations
940 (**Fig. S9**).

941

942 **Discussion**

943 In this work, we developed the transfer learning framework *DEGAS* to integrate scRNA-seq and
944 patient-level transcriptomic data in order to infer the transferrable “impressions” between patient
945 characteristics in single cells and cellular characteristics in patients. Using transfer learning, we
946 trained a model with both scRNA-seq and patient bulk tissue gene expression data, then reduced
947 the differences between the distributions of the representations for the two data types in the final
948 hidden layer of our model via domain adaptation. This process allows information about patient
949 disease attributes as well as cell types to be transferred between the two data types. We focus
950 on the transfer of patient disease attributes to cells because there are far fewer available methods
951 addressing this task than deconvolution. We tested and validated the *DEGAS* framework on
952 datasets from one simulation and two diseases: GBM, which contained ground truth tumor
953 subtype labels, and AD, which contained ground truth cell type-disease associations.

954

955 These experiments on validation datasets demonstrate the necessity for *DEGAS* especially as it
956 relates to the current methods that rely on accurate clustering, cell type annotation, or case-
957 control scRNA-seq. For datasets that contain case and control scRNA-seq data, tools like *Augur*

958 are very effective to prioritize cell types. When no patient level transcriptomic data is available but
959 case-control scRNA-seq is available, tools like *Augur* should be used since *DEGAS* requires
960 patient level transcriptomic data. If patient level transcriptomic data and single cell transcriptomic
961 data are available and there is a necessity to overlay disease associations onto individual cells,
962 then only *DEGAS* can be used. Furthermore, if the scRNA-seq dataset does not contain case and
963 control samples then *DEGAS* needs to be used instead of *Augur* since *Augur* requires case and
964 control samples. The *DEGAS* framework in this sense can be used in a wide variety of study
965 designs as long as there is scRNA-seq and patient transcriptomic data.

966

967 Another challenging issue in scRNA-seq analysis is that it is difficult to determine the best
968 clustering options. In our simulation examples, we can determine that the correct number of
969 clusters based on average silhouette width would be four clusters. However, if the number of
970 clusters was increased in the clustering algorithm there would be a stronger correlation between
971 some clusters and disease. Therein lies the challenge – should the clustering results be optimized
972 to reflect the relative transcriptomic signals or should they be optimized to create the greatest
973 correlations with disease state? Furthermore, the different resolutions of clusters may capture
974 different correlations with disease. For these reasons, assigning disease associations directly to
975 cells alleviates some of these problems with cluster resolution decisions. Assigning disease
976 associations directly to cells not only solves the cluster resolution problem but also allows
977 simultaneous identification of cell-intrinsic and cell proportional changes.

978

979 The *DEGAS* algorithm can identify both cell-intrinsic changes and cell proportional changes as
980 demonstrated in the simulation examples and the AD study. In simulation 1, the disease is
981 associated with proportional changes in cell type 1. In simulation 2, the disease is associated
982 with a cell-intrinsic change of cell type 4. In simulation 3, there are both cell-intrinsic changes
983 and cell proportional changes in cell type 4. In the AD experiments, two of the single-cell

984 datasets did include data from both AD and normal brains [15, 32]. The cells that came from the
985 AD patients tended to have a higher association with AD, which indicates the detection of cell-
986 intrinsic changes. The importance of cell prioritization at the individual cell level is highlighted in
987 simulation and AD examples. Simulation 2 shows an example where cell level associations are
988 necessary due to clustering results that do not capture the disease associations. Specifically,
989 there are cases where cells will cluster together but have dissimilar associations to disease. If
990 the cells of cell type 4 are not evaluated individually, the association of the cell type 4 disease
991 subtype with disease could be lost. In the AD example, the astrocyte cell type is overall not
992 associated with AD. However, a subset of astrocytes expressing markers for DAAs were found
993 to have a positive disease association while still clustering with the astrocytes that were not
994 associated with disease. Similarly, microglia cells are broadly positively associated with AD but
995 the highest AD association microglia were enriched for DAM markers. When a targeted analysis
996 was performed on only microglia from AD and normal brains, highest AD association microglia
997 were enriched for both HAM and DAM markers. These examples show how DEGAS can identify
998 disease associated cells that cluster within a larger cell type.

999
1000 In short, the *DEGAS* analysis on AD data further validated our model by correctly identifying the
1001 decreased neuron and increased microglia proportions in AD patients. Aside from these known
1002 characteristics of AD pathology, we also identified a *GFAP*⁺ astrocyte subtype taken from normal
1003 human brain tissue that is associated with AD and is supported from AD mouse models [55]. We
1004 further validated this by finding that *GFAP* expression in Astrocytes was significantly increased in
1005 Astrocytes taken from AD patients and concluded that there may be an expansion of this Astrocyte
1006 subtype in AD. This is also a convincing example of the utility of *DEGAS* as it assigned disease
1007 association at the single cell level, allowing us to identify intra-cell type differences in disease risk
1008 that constitute disease-associated cells.

1009

1010 For the GBM single cell patient cohort, each GBM tumor, from which scRNA-seq data was
1011 generated, had a GBM subtype label [44]. The *DEGAS* results showed that the majority of cells
1012 in each tumor were labeled with the same GBM subtype as previously defined in Patel *et al.* [44].
1013 Specifically, *DEGAS* correctly mapped Proneural, Mesenchymal, Classical, and Neural GBM
1014 subtypes to single cells in four GBM tumor samples. This experiment also shows the broad
1015 applicability of the model since the single cells had no labels and the patient samples had
1016 multiclass labels. *DEGAS* is highly flexible and allows for different categories of output labels to
1017 be combined, which may include but are not limited to classification labels, Cox proportional
1018 hazard, and even no labels. This allows for a wide variety of applications to adopt the *DEGAS*
1019 framework so that impressions are not limited to only one type of disease attribute.

1020

1021 To explore disease with less understood cellular subtypes, we applied *DEGAS* to multiple MM
1022 datasets. The models were able to assign PFS metrics to individual cells and subtype populations
1023 of CD138+ cells identified by cell type clustering methods Seurat [28] and *BERMUDA* [25]. Among
1024 the identified subtypes of cells, subtype 2 was the most consistent between patients visualized
1025 by *BERMUDA* (**Fig. S4E**). Furthermore, we found that the subtype 2 cell population appeared to
1026 have a gradient of cells moving away from the main subtype 1 group, possibly associated with a
1027 certain degree of differentiation (**Fig. S4A-D**). We did experience a lower PR-AUC for subtype 4
1028 than the other subtypes used during model training. However, this subtype was extremely
1029 uncommon in the samples and as a result the random PR-AUC would be close to zero making
1030 the PR-AUC of 0.44 well above random. Considering that subtype 4 was not found to be highly
1031 associated with progression, the lower PR-AUC did not greatly affect our interpretation of the
1032 data, which mainly focused on subtype 1 and subtype 2. We believe that *DEGAS* could be
1033 improved for highly imbalanced data.

1034

1035 Upon further examination, we found evidence that the subtype 2 cells may represent a population
1036 of malignant plasma cells expressing high levels of *PHF19*. *PHF19* is known to play a role in
1037 hematopoietic stem cell state and differentiation [65-67] and is a marker for aggressive disease
1038 in MM [64]. Furthermore, knock down of *PHF19* has been shown to shift myeloma cells into a less
1039 proliferative state [64]. The subtype 2 cells express *SDC1* (also known as *CD138*) and showed
1040 significantly increased *PHF19* expression in comparison to the other subtypes. Since all of the
1041 IUSM MM cells in our study had already been FACS sorted for CD138+, it is possible we have
1042 identified a subpopulation of CD138+/*PHF19*^{high} cells in MM tumors. This could prove a useful
1043 finding since currently the association between *PHF19* and tumor aggressiveness is at the patient
1044 level whereas our results imply that only a fraction of malignant plasma cells in a MM tumor
1045 actually overexpress *PHF19*.

1046

1047 This subtype could be targeted using precision immunotherapies that are not restricted to a single
1048 patient since the CD138+/*PHF19*^{high} cells (*i.e.*, subtype 2) were found to be present in multiple
1049 (3/4) patients. Of the three patients with detectable levels of subtype 2 in the CD138+ fraction,
1050 two patients (patient 2 and patient 4) had relapsed MM at time of biopsy and the other patient
1051 (patient 5) was SMM at biopsy and later progressed to MM. The other patient (patient 3) had little
1052 to no detectable subtype 2 cells in the CD138+ fraction and was SMM at time of biopsy and has
1053 not progressed to MM. These signs again seem to indicate a common cellular phenotype
1054 associated with progression in MM.

1055

1056 Based on the validated results in a variety of disease data analyses, we find that *DEGAS* has
1057 broad applications in virtually all diseases with available patient-level and single cell level omic
1058 data. The tensorflow [68] machine learning code is integrated with a simple R package interface
1059 (<https://github.com/tsteelejohnson91/DEGAS>) which will facilitate researchers to manipulate
1060 scRNA-seq and bulk expression data on their own.

1061

1062 **Conclusion**

1063 *DEGAS* is a powerful transfer learning tool for integrating different levels of omic data and
1064 identifying the latent molecular relationships between populations of cells and disease
1065 attributes, which we refer to as impressions. We validated the *DEGAS* framework on simulated
1066 data, GBM and AD by showing *DEGAS* models were capable of accurately predicting patient
1067 characteristics at single-cell level. We then leveraged this transfer learning approach on MM
1068 data and identified a CD138+*IPHF19*^{high} subtype population in MM that was significantly
1069 associated with disease progression. This subtype contains unique RNA profiles and gene
1070 correlations that could be both leveraged as a prognostic biomarker and possibly targeted
1071 directly to reduce the risk of progression. We believe that *DEGAS* can be a powerful solution to
1072 overcome the challenge of integrating patient single-cell data with bulk tissue data so that
1073 researchers can identify populations of cells associated with an disease attribute of interest.
1074 Furthermore, *DEGAS* can accommodate flexible data types. This makes it a highly general
1075 framework that can be applied in multiple diseases and data types to identify cellular
1076 populations that are associated with prognosis or treatment response, or to identify specific
1077 patient groups with certain cell subtypes for personalized treatment.

1078

1079 **Acknowledgements**

1080 We thank the Center for Computational Biology and Bioinformatics for the computational
1081 resources and work space to complete the research. We also thank the MMRF for the data
1082 generated as part of the Multiple Myeloma Research Foundation Personalized Medicine
1083 Initiatives (<https://research.themmr.org> and www.themmr.org), the Allen Institute for Brain
1084 Science for the data generated as part of their cell types database, and Mount Sinai/JJ Peters
1085 VA Medical Center for the data generated as a part of their brain bank.

1086

1087 **Funding**

1088 National Institutes of Health NLM-NRSA Fellowship F31LM013056 to TSJ and The Ohio State
1089 University (Columbus, OH) and departmental start-up funding from the Indiana University
1090 School of Medicine (Indianapolis, IN) to KH and TSJ.

1091

1092 **Author contributions**

1093 TSJ, CYY, JZ, and KH conceived and designed the project. TSJ, CYY, SX performed the
1094 analyses. TSJ and ZH designed the software package. TSJ, CYY, XH, SX, CD, MA, YW, CB,
1095 YZ, YL, JZ, BW, and KH interpreted the results. ZH, TW, WS, YW, and CB provided technical
1096 guidance. TSJ, CYY, JZ, and KH wrote the manuscript. JZ and KH supervised the project.

1097

1098 **Competing interests**

1099 The authors declare that this research was conducted in the absence of any commercial or
1100 financial relationships that could be construed as a potential conflict of interest.

1101

1102 **Data and materials availability**

1103 The *DEGAS* R package is freely available on GitHub
1104 (<https://github.com/tsteelejohnson91/DEGAS>). A minimum reproducible example of the IUSM
1105 myeloma scRNA-seq data is also deposited on Github. Our complete myeloma scRNA-seq has
1106 been deposited on GEO (GSE161722). All other data are publicly available.

1107

1108 **Figure legends**

1109 **Fig. 1 A workflow diagram of the *DEGAS* framework. A)** The workflow for a typical experiment
1110 with *DEGAS*. Note that *DEGAS* is not meant to replace the abundant packages available to load,
1111 preprocess, select features, cluster, and visualize scRNA-seq data. It is rather meant to augment
1112 these packages to assign disease associations to cells. **B)** The scRNA-seq and patient

1113 expression data are preprocessed into expression matrices. Next, Bootstrap aggregated
1114 DenseNet *DEGAS* models are trained using both single cell and patient disease attributes using
1115 a multitask learning neural network that learns latent representation reducing the differences
1116 between patients and single cells at the final hidden layer using maximum mean discrepancy
1117 (MMD). **C)** The output layer of this model can be used to simultaneously infer disease attribute
1118 impressions in single cells and cellular composition impressions in patients.

1119

1120 **Fig. 2 Simulation study and baseline comparisons of *DEGAS* framework.** **A)** 5,000 simulated
1121 cells from *Splatter* with 4 cell types where one of the cell types has two subtypes. Cell type 4 is
1122 composed of two subtypes that are specific to either disease or normal patients. 2,000 of these
1123 cells were used to generate the 600 simulated patients in **B-D** and 3,000 were used as the cell
1124 input to our *DEGAS* models. **E)** Optimal cluster number (4 clusters) based on average silhouette
1125 width for the 3,000 cells not used to generate patients. **F)** The same 3,000 cells used as the
1126 cellular input colored by their cluster. **G)** *DEGAS* comparison to *Augur* in simulation 1. **H)** *DEGAS*
1127 comparison with *Augur* in simulation 2. **I)** *DEGAS* comparison with *Augur* in simulation 3. **J-L)**
1128 *DEGAS*-calculated disease association from each simulation overlaid onto 3,000 cells. The violin
1129 plot in the bottom left corner is deconvolution cell type proportion for cell type 1 in simulation 1
1130 patients (**J**), cell type 4 proportion in simulation 2 patients (**K**), and cell type 4 proportion in
1131 simulation 3 patients (**L**).

1132

1133 **Fig. 3 *DEGAS* validation in GBM and AD.** *DEGAS* output of the distribution of GBM subtypes
1134 in single cells from five GBM tumors. Four of the five tumors had known GBM subtype information
1135 from Patel *et al.* (MGH26: Proneural, MGH28: Mesenchymal, MGH29: Mesenchymal, and
1136 MGH30: Classical, indicated by red boxes) which were recapitulated by *DEGAS*. The subtype
1137 information for the tumors, MGH26, MGH28, MGH29, and MGH30 were derived from Patel *et al.*
1138 where MGH31 did not have a clearly defined GBM subtype. The association of cells assigned to

1139 each subtype were plotted for each tumor; **A)** MGH26, **B)** MGH28, **C)** MGH29, **D)** MGH30 and **E)**
1140 MGH31. Median values are marked by a diamond in each of the violin plots. **F)** The death
1141 association centered around 0 is overlaid on all of the single cells from the five tumors (indicated
1142 by color). **G)** *DEGAS* output of AD association for each single cell. The AD association score is
1143 indicated by the color and is overlaid onto AIBS single cells. This plot shows the negative AD
1144 association in neuron cells and positive AD association in Microglia. **H-I)** There also appeared to
1145 be a subpopulation of astrocytes with positive AD association. The astrocytes were plotted
1146 separately and colored by AIBS Astrocyte subtypes (**H)** and *GFAP* expression, a disease-
1147 associated astrocyte marker (**I**). **J)** Comparison of *DEGAS*-derived AD associations for single
1148 cells from AD and Normal control samples from Grubman *et al.* **K-M)** Targeted analysis of
1149 microglia from Grubman *et al.* including the AD associations overlaid onto microglia (**K**), AD
1150 association comparing AD status of patient sample from which the cells were sampled (**L**), and
1151 PCC between AD association with HAM marker genes comparing up- and down-regulated HAM
1152 marker genes (**M**). Significance values: n.s. (not significant), • (0.1), * (0.05), ** (0.01), *** (0.001).

1153
1154 **Fig. 4 Association between subtypes and progression risk in MM.** IUSM CD138+ scRNA-
1155 seq subtype clusters generated from *Seurat* colored by **A)** cluster, *i.e.*, subtype and **B)**
1156 progression association. **C)** Kaplan-Meier curves of PFS from cross-validation for the MMRF
1157 patients stratified by median proportional hazard. **D)** Kaplan-Meier curves of OS from Zhan *et al.*
1158 external dataset stratified by median proportional hazard. **E)** Progression association for IUSM
1159 CD138+ subtypes **F)** Progression association for NHIP, MGUS, SMM, and MM in the external
1160 dataset Ledergor *et al.* **G)** Subtype 2 enrichment for NHIP, MGUS, SMM, and MM in the
1161 external dataset Ledergor *et al.* NHIP: normal hip bone marrow, MGUS: monoclonal
1162 gammopathy of undetermined significance, SMM: smoldering multiple myeloma, MM: multiple
1163 myeloma. Significance values: • (0.1), * (0.05), ** (0.01), *** (0.001). All plots were generated

1164 using the default parameters for the *DEGAS* package described in the section of **Methods**:

1165 *Transfer learning using DEGAS*.

1166

1167 **References**

- 1168 1. Lahnemann, D., et al., *Eleven grand challenges in single-cell data science*. *Genome Biol*,
1169 2020. **21**(1): p. 31.
- 1170 2. Ma, A., et al., *Integrative Methods and Practical Challenges for Single-Cell Multi-omics*.
1171 *Trends Biotechnol*, 2020. **38**(9): p. 1007-1022.
- 1172 3. Kiselev, V.Y., A. Yiu, and M. Hemberg, *scmap: projection of single-cell RNA-seq data*
1173 *across data sets*. *Nat Methods*, 2018. **15**(5): p. 359-362.
- 1174 4. Cao, Y., et al., *scRNASeqDB: A Database for RNA-Seq Based Gene Expression Profiles in*
1175 *Human Single Cells*. *Genes (Basel)*, 2017. **8**(12).
- 1176 5. Abugessaisa, I., et al., *SCPortalen: human and mouse single-cell centric database*. *Nucleic*
1177 *Acids Res*, 2018. **46**(D1): p. D781-D787.
- 1178 6. Stuart, T., et al., *Comprehensive integration of single-cell data*. *Cell*, 2019. **177**(7): p.
1179 1888-1902. e21.
- 1180 7. Gawel, D.R., et al., *A validated single-cell-based strategy to identify diagnostic and*
1181 *therapeutic targets in complex diseases*. *Genome Med*, 2019. **11**(1): p. 47.
- 1182 8. Chen, S., et al., *Single-cell analysis reveals transcriptomic remodellings in distinct cell*
1183 *types that contribute to human prostate cancer progression*. *Nature Cell Biology*, 2021.
1184 **23**(1): p. 87-98.
- 1185 9. Jang, J.S., et al., *Molecular signatures of multiple myeloma progression through single*
1186 *cell RNA-Seq*. *Blood cancer journal*, 2019. **9**(1): p. 1-10.
- 1187 10. Maynard, A., et al., *Therapy-induced evolution of human lung cancer revealed by single-*
1188 *cell RNA sequencing*. *Cell*, 2020. **182**(5): p. 1232-1251. e22.
- 1189 11. Cobos, F.A., et al., *Benchmarking of cell type deconvolution pipelines for transcriptomics*
1190 *data*. *Nature communications*, 2020. **11**(1): p. 1-14.
- 1191 12. Johnson, T.S., et al., *Combinatorial analyses reveal cellular composition changes have*
1192 *different impacts on transcriptomic changes of cell type specific genes in Alzheimer's*
1193 *Disease*. *Sci Rep*, 2021. **11**(1): p. 353.
- 1194 13. Johnson, T.S., et al., *Spatial cell type composition in normal and Alzheimers human*
1195 *brains is revealed using integrated mouse and human single cell RNA sequencing*. *Sci*
1196 *Rep*, 2020. **10**(1): p. 18014.
- 1197 14. Jung, S.-H. and S.-C. Chow, *On sample size calculation for comparing survival curves*
1198 *under general hypothesis testing*. *Journal of biopharmaceutical statistics*, 2012. **22**(3): p.
1199 485-495.
- 1200 15. Mathys, H., et al., *Single-cell transcriptomic analysis of Alzheimer's disease*. *Nature*,
1201 2019. **570**(7761): p. 332-337.
- 1202 16. Rossi, M.A., et al., *Obesity remodels activity and transcriptional state of a lateral*
1203 *hypothalamic brake on feeding*. *Science*, 2019. **364**(6447): p. 1271-1274.

- 1204 17. Crowell, H.L., et al., *muscat detects subpopulation-specific state transitions from multi-*
1205 *sample multi-condition single-cell transcriptomics data*. Nat Commun, 2020. **11**(1): p.
1206 6077.
- 1207 18. Burkhardt, D.B., et al., *Quantifying the effect of experimental perturbations at single-cell*
1208 *resolution*. Nat Biotechnol, 2021. **39**(5): p. 619-629.
- 1209 19. Skinnider, M.A., et al., *Cell type prioritization in single-cell data*. bioRxiv, 2019: p.
1210 2019.12.20.884916.
- 1211 20. Kouw, W.M. and M. Loog, *A Review of Domain Adaptation without Target Labels*. IEEE
1212 Transactions on Pattern Analysis and Machine Intelligence, 2021. **43**(3): p. 766-785.
- 1213 21. Haroon, D.R., S. Szedmak, and J. Shawe-Taylor, *Canonical correlation analysis: An*
1214 *overview with application to learning methods*. Neural computation, 2004. **16**(12): p.
1215 2639-2664.
- 1216 22. Gretton, A., et al., *A kernel two-sample test*. Journal of Machine Learning Research,
1217 2012. **13**(Mar): p. 723-773.
- 1218 23. Andrew, G., et al. *Deep canonical correlation analysis*. in *International conference on*
1219 *machine learning*. 2013. PMLR.
- 1220 24. Zhang, F., Y. Wu, and W. Tian, *A novel approach to remove the batch effect of single-cell*
1221 *data*. Cell discovery, 2019. **5**(1): p. 1-4.
- 1222 25. Wang, T., et al., *BERMUDA: a novel deep transfer learning method for single-cell RNA*
1223 *sequencing batch correction reveals hidden high-resolution cellular subtypes*. Genome
1224 Biol, 2019. **20**(1): p. 165.
- 1225 26. Tran, H.T.N., et al., *A benchmark of batch-effect correction methods for single-cell RNA*
1226 *sequencing data*. Genome biology, 2020. **21**(1): p. 1-32.
- 1227 27. Johnson, T.S., et al., *LAMBDA: label ambiguous domain adaptation dataset integration*
1228 *reduces batch effects and improves subtype detection*. Bioinformatics, 2019. **35**(22): p.
1229 4696-4706.
- 1230 28. Butler, A., et al., *Integrating single-cell transcriptomic data across different conditions,*
1231 *technologies, and species*. Nat Biotechnol, 2018. **36**(5): p. 411-420.
- 1232 29. Aran, D., et al., *Reference-based analysis of lung single-cell sequencing reveals a*
1233 *transitional profibrotic macrophage*. Nature immunology, 2019. **20**(2): p. 163-172.
- 1234 30. Araujo, T., et al., *Classification of breast cancer histology images using Convolutional*
1235 *Neural Networks*. PLoS One, 2017. **12**(6): p. e0177544.
- 1236 31. Bardou, D., K. Zhang, and S.M. Ahmad, *Classification of Breast Cancer Based on*
1237 *Histology Images Using Convolutional Neural Networks*. IEEE Access, 2018. **6**: p. 24680-
1238 24693.
- 1239 32. Grubman, A., et al., *A single-cell atlas of entorhinal cortex from individuals with*
1240 *Alzheimer's disease reveals cell-type-specific gene expression regulation*. Nature
1241 neuroscience, 2019. **22**(12): p. 2087-2097.
- 1242 33. Holtman, I.R., et al., *Induction of a common microglia gene expression signature by*
1243 *aging and neurodegenerative conditions: a co-expression meta-analysis*. Acta
1244 Neuropathol Commun, 2015. **3**: p. 31.
- 1245 34. Hemonnot, A.L., et al., *Microglia in Alzheimer Disease: Well-Known Targets and New*
1246 *Opportunities*. Front Aging Neurosci, 2019. **11**: p. 233.

- 1247 35. Glass, C.K., et al., *Mechanisms underlying inflammation in neurodegeneration*. Cell,
1248 2010. **140**(6): p. 918-34.
- 1249 36. Donev, R., et al., *Neuronal death in Alzheimer's disease and therapeutic opportunities*. J
1250 Cell Mol Med, 2009. **13**(11-12): p. 4329-48.
- 1251 37. DeKosky, S.T. and S.W. Scheff, *Synapse loss in frontal cortex biopsies in Alzheimer's*
1252 *disease: correlation with cognitive severity*. Ann Neurol, 1990. **27**(5): p. 457-64.
- 1253 38. de Wilde, M.C., et al., *Meta-analysis of synaptic pathology in Alzheimer's disease reveals*
1254 *selective molecular vesicular machinery vulnerability*. Alzheimers Dement, 2016. **12**(6):
1255 p. 633-44.
- 1256 39. Akiyama, H., *Inflammatory response in Alzheimer's disease*. Tohoku J Exp Med, 1994.
1257 **174**(3): p. 295-303.
- 1258 40. Institute, N.C., *Cancer Statistics*, N.C. Institute, Editor. 2019: Cancer.gov.
- 1259 41. Zhan, F., et al., *The molecular classification of multiple myeloma*. Blood, 2006. **108**(6): p.
1260 2020-8.
- 1261 42. Ledergor, G., et al., *Single cell dissection of plasma cell heterogeneity in symptomatic*
1262 *and asymptomatic myeloma*. Nat Med, 2018. **24**(12): p. 1867-1876.
- 1263 43. Cohen, Y.C., et al., *Identification of resistance pathways and therapeutic targets in*
1264 *relapsed multiple myeloma patients through single-cell sequencing*. Nature medicine,
1265 2021: p. 1-13.
- 1266 44. Patel, A.P., et al., *Single-cell RNA-seq highlights intratumoral heterogeneity in primary*
1267 *glioblastoma*. Science, 2014. **344**(6190): p. 1396-401.
- 1268 45. Cancer Genome Atlas Research, N., *Comprehensive genomic characterization defines*
1269 *human glioblastoma genes and core pathways*. Nature, 2008. **455**(7216): p. 1061-8.
- 1270 46. Wang, M., et al., *The Mount Sinai cohort of large-scale genomic, transcriptomic and*
1271 *proteomic data in Alzheimer's disease*. Sci Data, 2018. **5**: p. 180185.
- 1272 47. Ching, T., X. Zhu, and L.X. Garmire, *Cox-nnet: An artificial neural network method for*
1273 *prognosis prediction of high-throughput omics data*. PLoS Comput Biol, 2018. **14**(4): p.
1274 e1006076.
- 1275 48. Couturier, C.P., et al., *Single-cell RNA-seq reveals that glioblastoma recapitulates a*
1276 *normal neurodevelopmental hierarchy*. Nat Commun, 2020. **11**(1): p. 3406.
- 1277 49. Guo, M., et al., *SINCERA: A Pipeline for Single-Cell RNA-Seq Profiling Analysis*. PLoS
1278 Comput Biol, 2015. **11**(11): p. e1004575.
- 1279 50. Iacono, G., R. Massoni-Badosa, and H. Heyn, *Single-cell transcriptomics unveils gene*
1280 *regulatory network plasticity*. Genome Biol, 2019. **20**(1): p. 110.
- 1281 51. Grus, J., *Data science from scratch: first principles with python*. 2019: O'Reilly Media.
- 1282 52. Ioffe, S. and C. Szegedy. *Batch normalization: Accelerating deep network training by*
1283 *reducing internal covariate shift*. in *International conference on machine learning*. 2015.
1284 PMLR.
- 1285 53. Juszczak, P., D. Tax, and R.P. Duin. *Feature scaling in support vector data description*. in
1286 *Proc. asc*. 2002. Citeseer.
- 1287 54. Verhaak, R.G., et al., *Integrated genomic analysis identifies clinically relevant subtypes of*
1288 *glioblastoma characterized by abnormalities in PDGFRA, IDH1, EGFR, and NF1*. Cancer
1289 Cell, 2010. **17**(1): p. 98-110.

- 1290 55. Habib, N., et al., *Disease-associated astrocytes in Alzheimer's disease and aging*. Nature
1291 Neuroscience, 2020. **23**(6): p. 701-706.
- 1292 56. Srinivasan, K., et al., *Alzheimer's Patient Microglia Exhibit Enhanced Aging and Unique*
1293 *Transcriptional Activation*. Cell Rep, 2020. **31**(13): p. 107843.
- 1294 57. Keren-Shaul, H., et al., *A Unique Microglia Type Associated with Restricting Development*
1295 *of Alzheimer's Disease*. Cell, 2017. **169**(7): p. 1276-1290 e17.
- 1296 58. Maaten, L.v.d. and G. Hinton, *Visualizing data using t-SNE*. Journal of machine learning
1297 research, 2008. **9**(Nov): p. 2579-2605.
- 1298 59. Becht, E., et al., *Dimensionality reduction for visualizing single-cell data using UMAP*. Nat
1299 Biotechnol, 2018.
- 1300 60. Friedman, J., T. Hastie, and R. Tibshirani, *Regularization Paths for Generalized Linear*
1301 *Models via Coordinate Descent*. J Stat Softw, 2010. **33**(1): p. 1-22.
- 1302 61. Zappia, L., B. Phipson, and A. Oshlack, *Splatter: simulation of single-cell RNA sequencing*
1303 *data*. Genome Biol, 2017. **18**(1): p. 174.
- 1304 62. Fu, H., et al., *Tau Pathology Induces Excitatory Neuron Loss, Grid Cell Dysfunction, and*
1305 *Spatial Memory Deficits Reminiscent of Early Alzheimer's Disease*. Neuron, 2017. **93**(3):
1306 p. 533-541 e5.
- 1307 63. Xu, J., et al., *Multimodal single-cell/nucleus RNA sequencing data analysis uncovers*
1308 *molecular networks between disease-associated microglia and astrocytes with*
1309 *implications for drug repurposing in Alzheimer's disease*. Genome research, 2021: p. gr.
1310 272484.120.
- 1311 64. Mason, M.J., et al., *Multiple Myeloma DREAM Challenge reveals epigenetic regulator*
1312 *PHF19 as marker of aggressive disease*. Leukemia, 2020. **34**(7): p. 1866-1874.
- 1313 65. Bagger, F.O., S. Kinalis, and N. Rapin, *BloodSpot: a database of healthy and malignant*
1314 *haematopoiesis updated with purified and single cell mRNA sequencing profiles*. Nucleic
1315 acids research, 2019. **47**(D1): p. D881-D885.
- 1316 66. Lara-Astiaso, D., et al., *Chromatin state dynamics during blood formation*. science, 2014.
1317 **345**(6199): p. 943-949.
- 1318 67. Vizán, P., et al., *The Polycomb-associated factor PHF19 controls hematopoietic stem cell*
1319 *state and differentiation*. Science advances, 2020. **6**(32): p. eabb2745.
- 1320 68. Abadi, M., et al. *Tensorflow: A system for large-scale machine learning*. in *12th {USENIX}*
1321 *Symposium on Operating Systems Design and Implementation ({OSDI} 16)*. 2016.
1322







