

1 AAnet resolves a continuum of spatially-localized cell states to unveil  
2 tumor complexity

3 Aarthi Venkat<sup>1#</sup> Scott E. Youlten<sup>2,#</sup> Beatriz P. San Juan<sup>3,4,#</sup> Carley Purcell<sup>3,4</sup> Matthew Amodio<sup>5,6</sup>  
4 Daniel B. Burkhardt<sup>2,7</sup> Andrew Benz<sup>7,8</sup> Jeff Holst<sup>9</sup> Cerys McCool<sup>3,4</sup> Annelie Mollbrink<sup>10</sup> Joakim  
5 Lundeberg<sup>10</sup> David van Dijk<sup>1,5,11</sup> Leonard D. Goldstein<sup>3,4</sup> Sarah Kummerfeld<sup>3,4,§</sup> Smita  
6 Krishnaswamy<sup>1,2,5,12,§</sup> Christine L. Chaffer<sup>3,4,§</sup>

7 #Co-first authors

8 §Co-senior authors

9 Correspondence: smita.krishnaswamy@yale.edu; c.chaffer@garvan.org.au

10 <sup>1</sup>*Computational Biology and Bioinformatics Program, Yale University, New Haven, CT, USA,* <sup>2</sup>*Department of Genetics, Yale School*  
11 *of Medicine, New Haven, CT, USA,* <sup>3</sup>*Kinghorn Cancer Centre, Garvan Institute of Medical Research, Darlinghurst 2010, Australia,*  
12 <sup>4</sup>*St Vincent's Clinical School, University of New South Wales Medicine, University of New South Wales, Darlinghurst 2010,*  
13 *Australia* <sup>5</sup>*Department of Computer Science, Yale University, New Haven, CT, USA,* <sup>6</sup>*Broad Institute at MIT and Harvard,*  
14 *Cambridge, MA, USA* <sup>7</sup>*Cellarity, Somerville, MA, USA,* <sup>8</sup>*Department of Mathematics, Yale University, New Haven, CT, USA,*  
15 <sup>9</sup>*University of New South Wales, Kensington, New South Wales, Australia,* <sup>10</sup>*SciLifeLab, KTH Royal Institute of Technology,*  
16 *Sweden* <sup>11</sup>*Department of Internal Medicine, Yale School of Medicine, New Haven, CT, USA,* <sup>12</sup>*Applied Mathematics Program, Yale*  
17 *University, New Haven, CT, USA*

## 18 Summary

19 Identifying functionally important cell states and structure within a heterogeneous tumor remains a significant  
20 biological and computational challenge. Moreover, current clustering or trajectory-based computational models  
21 are ill-equipped to address the notion that cancer cells reside along a phenotypic continuum. To address this, we  
22 present Archetypal Analysis network (AAnet), a neural network that learns key archetypal cell states within  
23 a phenotypic continuum of cell states in single-cell data. Applied to single-cell RNA sequencing data from  
24 pre-clinical models and a cohort of 34 clinical breast cancers, AAnet identifies archetypes that resolve distinct  
25 biological cell states and processes, including cell proliferation, hypoxia, metabolism and immune interactions.  
26 Notably, archetypes identified in primary tumors are recapitulated in matched liver, lung and lymph node  
27 metastases, demonstrating that a significant component of intratumoral heterogeneity is driven by cell intrinsic  
28 properties. Using spatial transcriptomics as orthogonal validation, AAnet-derived archetypes show discrete spatial  
29 organization within tumors, supporting their distinct archetypal biology. We further reveal that ligand:receptor  
30 cross-talk between cancer and adjacent stromal cells contributes to intra-archetypal biological mimicry. Finally,  
31 we use AAnet archetype identifiers to validate GLUT3 as a critical mediator of a hypoxic cell archetype harboring  
32 a cancer stem cell population, which we validate in human triple-negative breast cancer specimens. AAnet is a  
33 powerful tool to reveal functional cell states within complex samples from multimodal single-cell data.

## 34 Introduction

35 Cancer cells can dynamically change their functional state to facilitate survival [8, 15, 31]. This creates a  
36 phenotypic continuum of cell states within a tumor that can be viewed as a cell state landscape. In this model,  
37 dynamic gene expression changes enable movement about the landscape. Defining the breadth of cell states and  
38 their structural organization within a tumor is currently a significant biological and computational challenge, yet  
39 is likely to reveal critical opportunities to perturb cancer progression.

40 The two predominant approaches for characterizing cell state heterogeneity from single-cell transcriptomic  
41 data are clustering and trajectory inference [24, 50]. Clustering partitions the cellular state space into discrete cell  
42 types, and trajectory inference identifies continuous paths that define a pattern of cellular dynamics. However,  
43 when there are no clear delineations between cellular states, nor clear trajectories or lineage structure on the data  
44 manifold, neither approach suffices to map the cellular state space. Thus, to define biological similarities and  
45 differences within and between these tumors, there is a need for a method that can dissect cellular heterogeneity  
46 at single-cell resolution while maintaining continuous variation along the cell state continuum. For this purpose,  
47 we turn to archetypal analysis, a factor analysis technique first introduced by Cutler and Breiman [11]. Archetypal  
48 analysis first extracts factors that represent the "archetypes", or extreme states, of a dataset. Then, all datapoints  
49 can be described as a convex combination of archetypes. In other words, archetypal analysis models the data  
50 as a simplex, where the extreme points are corners and other points are on the faces or are internal to the  
51 simplex (Figure 1). However, viewing this technique geometrically makes it clear that most datasets would not  
52 be accurately described by such a simplex in the ambient data space.

53 This motivated the development of our approach — Archetypal Analysis network (AAnet) — which performs  
54 archetypal analysis by *learning* a simplicial representation of the data. AAnet is implemented as an autoencoder,  
55 i.e. a neural network that learns meaningful representations of the given data, but is regularized to perform  
56 archetypal analysis. Instead of fitting a simplex on the data space, the AAnet encoder learns the optimal  
57 transformation from the data space to a latent space bound by a simplex, and the decoder learns the transformation  
58 back to the data space. This means that AAnet learns archetypes and a representation of each datapoint as  
59 convex combination of archetypes through nonlinear dimensionality reduction. This adds flexibility to archetypal  
60 detection while preserving data geometry.

61 Triple-negative breast cancer (TNBC) is a particularly heterogeneous subtype of breast cancer as it is an  
62 amalgamation of all "other" breast cancers that cannot be classified as ER+, PR+, or HER2+. With a lack of  
63 specific markers to characterize TNBC, there are consequently a paucity of effective targeted therapies to treat  
64 it. Here, we use AAnet to deconvolute TNBC heterogeneity into biologically interpretable archetypes. Using  
65 novel single-cell RNA sequencing (scRNAseq) and spatial transcriptomics datasets modeling tumor formation  
66 and metastasis *in vivo*, AAnet identifies five archetypes in primary tumors, demonstrating they are reproducibly

defined by distinct cancer hallmarks. These include a proliferative archetype associated with cell cycle progression; an oxidative archetype associated with oxidative phosphorylation, ROS production, and adipogenesis; a hypoxic archetype enriched for enzymes associated with oxygen-independent glycolysis; a cell damage/death archetype that captures variation introduced by technical factors; and an immune-stimulatory archetype with enriched expression of HLA genes and cytokines. We validate these archetypes by showing they are recapitulated in metastases, are spatially organized, and colocalize with distinct microenvironmental cells types and metabolic niches. Moreover, in a cohort of 34 human TNBC samples, AAnet reveals significant archetypal heterogeneity between patients. Notably, we identify a subset of patients defined by the hypoxic archetype favoring residence in a cancer stem cell niche, and we validate GLUT3 as a critical regulator of that archetypal cell state. These findings highlight the powerful ability of AAnet to define biological function and organization within cancer, potentially aiding in classifying patients according to biological similarities within and across patient samples. Furthermore, AAnet defines core transcriptional programs driving distinct archetypes, thereby suggesting potential therapeutic opportunities to target them.

## 80 Results

### 81 Defining *in vivo* cellular heterogeneity in scRNAseq data of triple-negative breast cancer

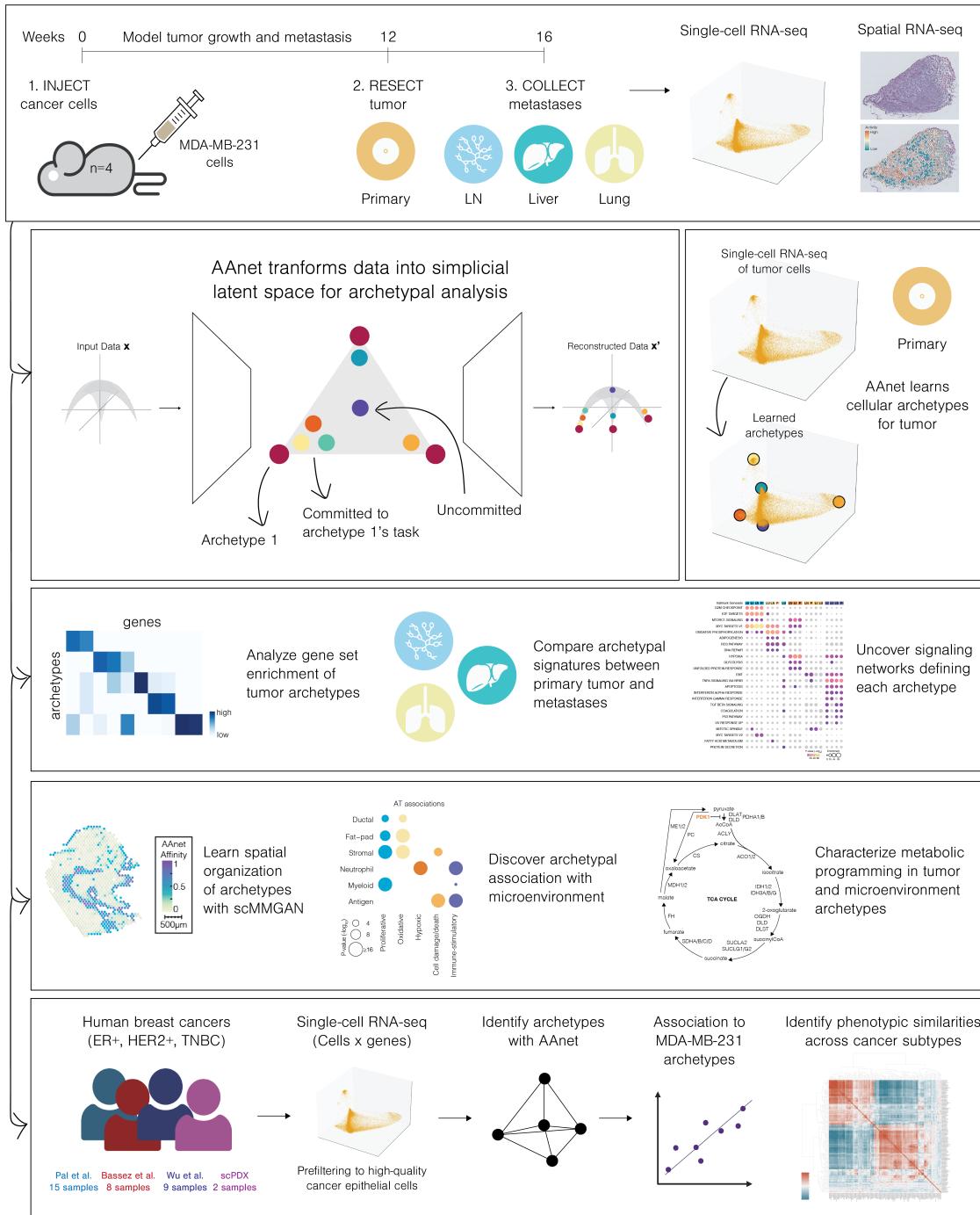
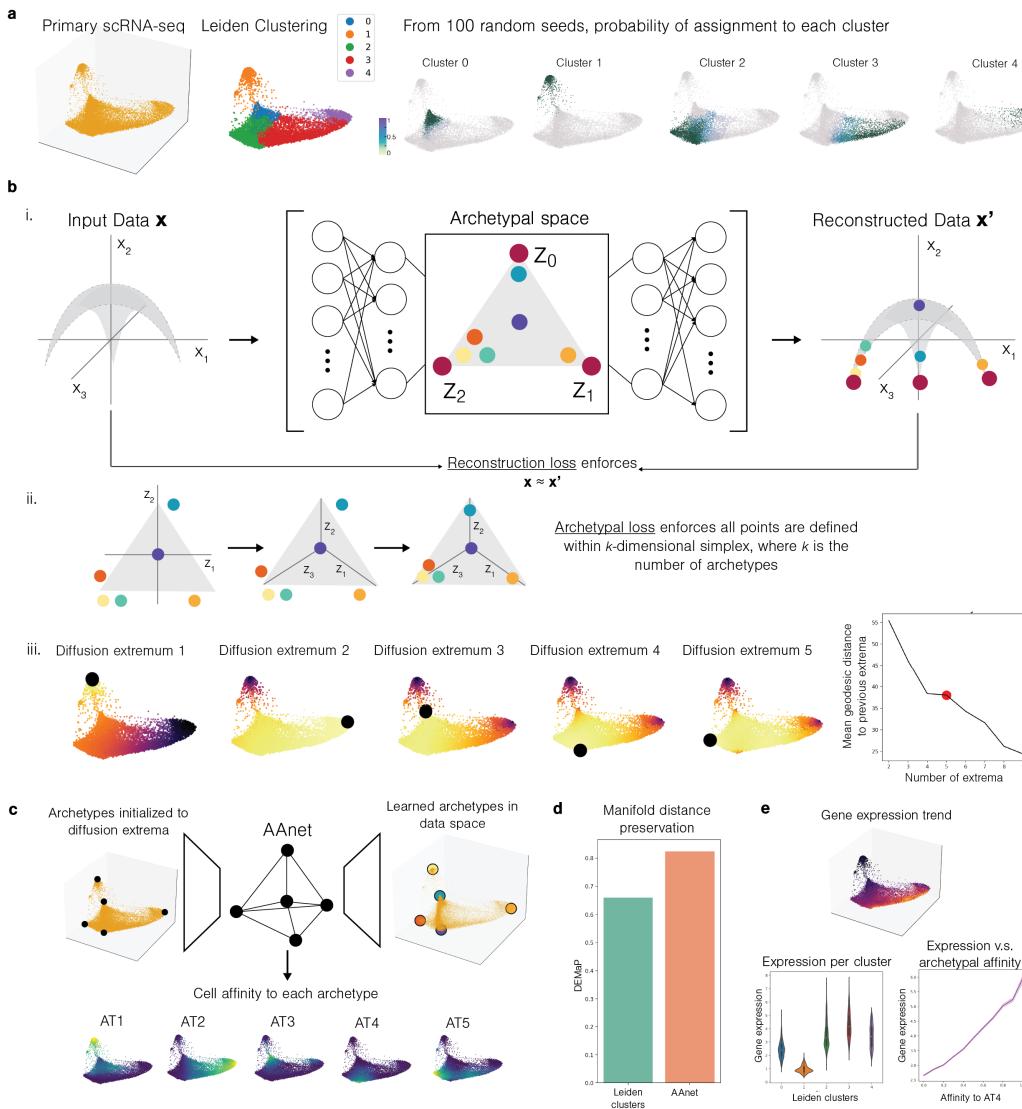


Figure 1. Overview of experimental design, AAnet, and downstream analyses.

## 82 Overview of AAnet

83 **Figure 2**



**Figure 2. Overview of AAnet architecture.** (a) Primary tissue visualization is a continuum of cells. Clustering the data with standard clustering tools 100 times (with no parameters changed) results in shifting boundaries between adjacent clusters. (b) i. AAnet learns transformation of data into a latent space shaped as a simplex. ii. Archetypal loss enforces that data lies in a latent space shaped like an  $k$ -dimensional simplex. iii. Diffusion extrema loss infers the extrema from the data geometry. The diffusion extrema also inform the number of archetypes for downstream analysis. (c) AAnet is initialized with diffusion extrema and learns affinity to each archetype and decoding of archetypes the data space. (d) Manifold distance preservation score (DEMaP) [32] of cluster representation versus AAnet representation. (e) AAnet latent space can be used to characterize continuous gene trends not easily characterizable with clustering.

84 We, and others, have shown that cancer cells reside along a phenotypic continuum [8]. The ability to identify  
 85 critical cell states along the continuum will garner insight into the molecular programs enabling cell state  
 86 adaptation, thereby facilitating therapeutic strategies to prevent it. For characterizing cell state heterogeneity  
 87 from scRNAseq data, clustering and trajectory inference are considered a standard part of single-cell workflows  
 88 and best practices [30]. However, we show that when the data lies on a continuum without latent cluster structure  
 89 (e.g., discrete cell types) or latent lineage structure (e.g., development axis), these approaches lack concordance  
 90 across methods and are limited in their ability to meaningfully characterize the cellular state space (Figure 2a,

91 Supplementary Figure 1a-c).

92 Archetypal analysis provides a framework to identify "archetypes", or extreme states, in a dataset and  
93 characterize the datapoints as a convex combination of archetypes [2, 11, 21, 38]. Currently, the biggest challenge  
94 of archetypal analysis is identifying the correct and relevant archetypes. Linear archetypal analysis methods,  
95 such as Principal Convex Hull Analysis (PCHA), fit a linear simplex onto the data in the ambient space [11]  
96 but fail to correctly identify the extrema when the extrema in the ambient space do not conform to the data  
97 geometry. These methods thus prove inflexible to more complex datasets, such as scRNASeq data.

98 To this end, we developed AAnet, a neural network for nonlinear archetypal analysis (Methods). AAnet  
99 learns a low-dimensional latent representation of the data as a regular simplex (Figure 2b i). This is achieved  
100 by regularizing the encoding layer of the neural network to encode points as convex or barycentric coordinates  
101 based on the archetypal points (Figure 2b ii). The autoencoder-style (i.e. encoder-decoder) architecture and  
102 the archetypal regularization together ensure that the model learns an accurate transformation to a simplex  
103 representation of the data, and it can decode data from the simplicial space to generate new data. In order to add  
104 further robustness to noise and accuracy to the model, we developed an approach (Methods) to identify extreme  
105 points based on the underlying data geometry, termed *diffusion extrema* (Figure 2b iii). We then use geodesic  
106 distances between diffusion extrema to choose the number of archetypes  $k$ , and we initialize the archetypes to  
107 the first  $k$  diffusion extrema at the beginning of training.

108 Once the model is trained, the simplicial latent representation can be used for exploration of the dataset.  
109 The vertices in the latent space, encoded by standard basis vectors, can be decoded to the archetypes in the gene  
110 space, enabling characterization and comparison of their expression profiles. Furthermore, the archetypal space  
111 coordinates provide an interpretable measure of each cell's affinity to each archetype (Figure 2c, Supplementary  
112 Figure 1g). Importantly, these archetypal affinities retain more information about cellular relationships than  
113 clustering while maintaining the interpretability of cell types. We show that the simplicial latent representation  
114 better preserves geodesic distances [32] than the cluster representation (Figure 2d), suggesting it could prove  
115 useful in tasks that depend on cluster annotations (e.g. [17, 27]). Finally, we can also represent signals, including  
116 gene expression, with respect to archetypal affinities, which allows characterization of continuous signals not  
117 possible with cluster-based enrichment analysis (Figure 2e).

## 118 Comparison of AAnet to other approaches on simulated data

119 To compare AAnet with existing approaches for characterizing cell state heterogeneity and archetypal analysis,  
120 we generated a nonlinearly-transformed tetrahedron, or a simplex with four vertices that are "ground truth"  
121 archetypes. We also simulated a nonlinear signal based on the true archetypal affinity to vertex four (Supple-  
122 mentary Figure 1a). Clustering (Supplementary Figure 1b) and trajectory inference (Supplementary Figure  
123 1c) approaches show disagreement in cluster and pseudotime assignments respectively and fail to capture the  
124 underlying relationship between the simulated signal and vertex four. Additionally, existing archetypal analysis  
125 methods cannot correctly infer the vertices of the tetrahedron and show worsening performance as we increase  
126 the nonlinearity of the tetrahedron transformation, suggesting that the linearity of these approaches is their  
127 major limitation (Supplementary Figure 1d).

128 By contrast, AAnet is able to infer the true vertices with nearly perfect performance at all levels of nonlinearity  
129 ( $\gamma_{extrema} = 5$ , Supplementary Figure 1e). Without the diffusion extrema loss ( $\gamma_{extrema} = 0$ ), AAnet shows  
130 better average performance than existing methods, though lacks robustness at very high degrees of nonlinearity  
131 (Supplementary Figure 1e). The approach to identify the number of archetypes based on diffusion extrema  
132 correctly identifies four vertices (Supplementary Figure 1f). Finally, AAnet captures interpretable archetypal  
133 affinities, and plotting the simulated signal against the inferred vertex 4 archetypal affinity shows AAnet is able  
134 to recapitulate the sinusoidal relationship (Supplementary Figure 1g).

## 135 Validation of AAnet on an antigen-specific CD8+ T cell dataset

136 To validate our method on real biological data, we leveraged a published single-cell dataset of tumor-specific  
137 CD8+ T cells [10]. (Supplementary Figure 2a).

138 Using a mouse model of lung adenocarcinoma designed to express neoantigens, Connolly et al. identified a  
139 reservoir of antigen-specific stem-like T cells in the tumor-draining lymph node (dLN), which then migrated to

140 the site of the tumor to terminally differentiate. Tumor-specific CD8+ T cells from early (8 weeks p.i.) and late  
141 (17 weeks p.i.) dLNs (top left), early and late tumors (top right), early dLNs and early tumors (bottom left),  
142 and late dLNs and late tumors (bottom right) (Supplementary Figure 2a) were co-embedded with PHATE [32].  
143 "Flags" (blue, white, green) were hand-annotated based on a continuum of expression of key immune-related  
144 marker genes in distinct regions of the single-cell embeddings. Given the utility of these flags to characterize  
145 the heterogeneity of the state space without demarcating boundaries between points, we refer to these flags as  
146 archetypes.

147 The white archetype was characterized by a naive CD8+ T cell signature based on high expression of *Sell*,  
148 *Lef1*, and *Ccr7*. The green archetype was characterized by a stem-like signature, defined by high expression of  
149 *Tcf7*, *Xcl1*, *Slamf6*. The blue archetype was characterized by an exhausted signature, with high expression of  
150 *Pdcld1*, *Havcr2*, and *Cd101*.

151 Here, we ran AAnet separately on each of the paired embeddings to determine if it could recapitulate these  
152 archetypes (Supplementary Figure 2b). First, we compared expression of each marker gene to characterize the  
153 archetypes. This resulted in annotation of one archetype in each of the paired embeddings as a naive archetype,  
154 one archetype in three out of four paired embeddings as stem-like, and one archetype in three out of four paired  
155 embeddings as exhausted (Supplementary Figure 2c). Three archetypes (LN AT2, Week 8 AT2, and Lung AT3)  
156 did not strongly express any of the markers. These correspond to the hand-annotated exhausted archetype in the  
157 lymph node, an uncharacterized part of the manifold, and the hand-annotated stem-like archetype in the lung,  
158 respectively. The authors note in the text that, while annotated, there was no prominent population of exhausted  
159 T cells in the lymph node and no prominent population of stem-like T cells in the lung. The uncharacterized  
160 archetype does not correspond to these three cell state extremes, possibly corresponding to an uncovered cell  
161 state.

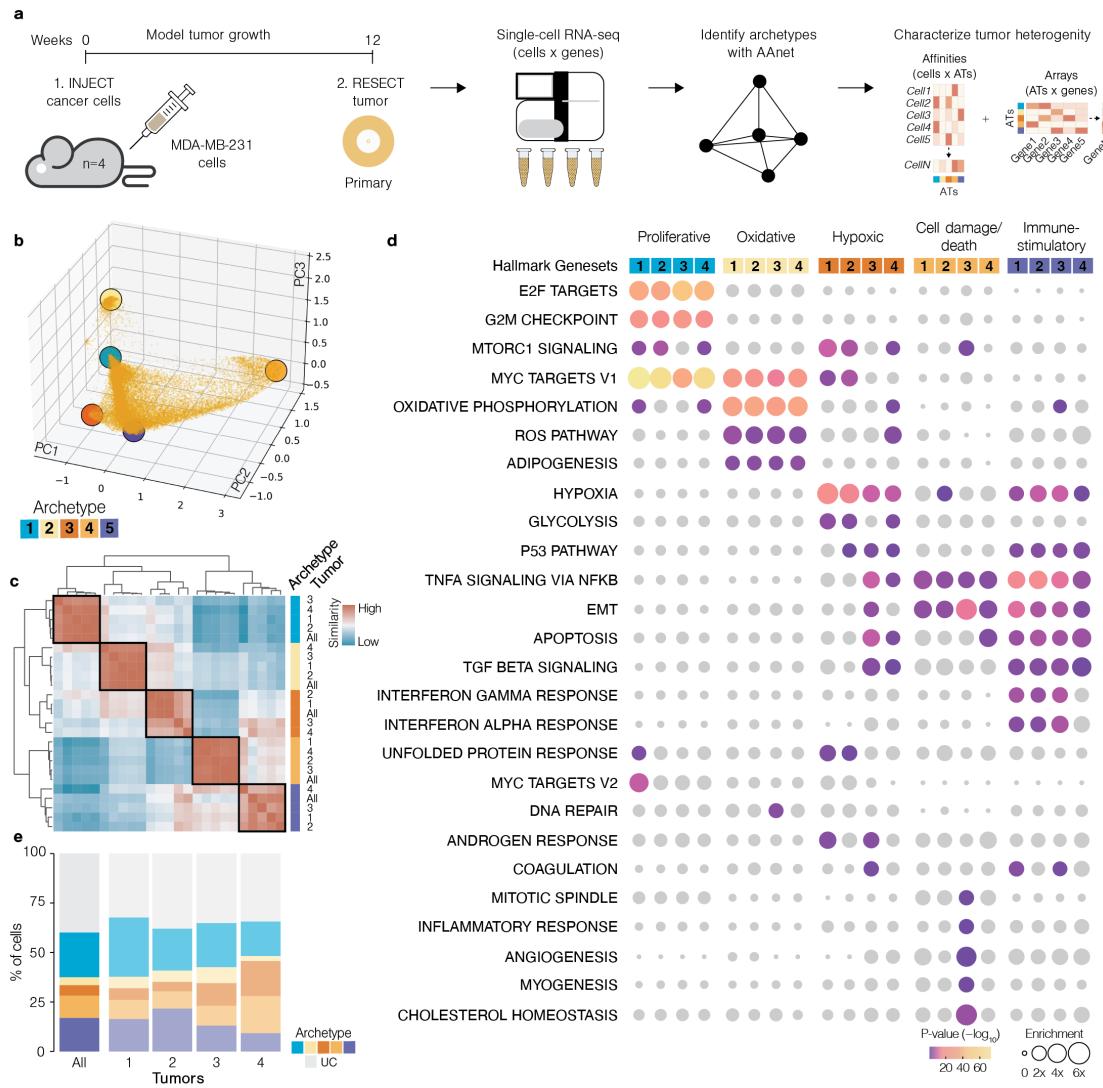
162 With the archetypes expressing the markers of interest, we computed pairwise cosine similarity across  
163 all measured genes (not only the markers of interest, as in the original work). This showed clear clustering  
164 corresponding to the naive, exhausted, and stem-like archetypes (Supplementary Figure 2d). This bolsters the  
165 finding in the paper by suggesting that these cell states share broader transcriptomic similarity not limited to  
166 the nine markers known to be associated with CD8+ T cell states.

167 Finally, to highlight the utility of the latent space learned by AAnet, we plotted the marker expression versus  
168 the latent space score for each archetype in the week 17 embedding (Supplementary Figure 2e). In all cases, the  
169 corresponding marker genes are upregulated, and the other marker genes are downregulated. Furthermore, we  
170 see non-linear dynamics of gene patterns corresponding to distance in the latent space, adding an additional  
171 layer of information through which to interpret the results.

172 Together, this analysis corroborates the use of AAnet for characterizing datasets with continuous and nonlinear  
173 structure with respect to archetypes. Furthermore, it validates the ability to compare archetypes across datasets  
174 to identify unified signatures of response.

175 **AAnet identifies five major archetypal expression states in primary tumor of triple-negative  
176 breast cancer model**

177 **Figure 3**



**Figure 3.** (a) Experimental approach used to identify tumor archetypes in TNBC with AAnet. (b) Embedding of TNBC tumor cells (all tumors combined) with archetypes indicated by colored circles. (c) Heatmap of cosine similarity between archetypal expression vectors determined in each individual tumor (numbered) and all tumors combined (all). Orthologous archetypes between samples are indicated with colors. (d) Enriched HALLMARK genesets associated with each archetype in each tumor (colors indicate AT, numbers indicate tumor). P-value is false discovery rate corrected, gene sets where  $p > 0.05$  are shown in gray. Enrichment represents log<sub>2</sub> fold-change. (e) The percentage of cells in each sample committed to each archetype (colored) or uncommitted to an archetype (gray).

178 Having demonstrated the power of AAnet to deconstruct single cell data into meaningful archetypes (ATs), we  
179 sought to use AAnet to address a critical question in cancer biology — how does the cell state landscape change  
180 across primary and metastatic tumors? Resolving this question may lead to new strategies to prevent non-genetic  
181 adaptation that facilitates cancer progression.

182 To answer this, we generated a new scRNASeq dataset using an *in vivo* model of triple-negative breast cancer.  
183 Highly metastatic MDA-MB-231 breast cancer cells were injected into the mammary fat pad of NSG mice and  
184 left to grow (6-8 week, Female, n=4; Figure 1). At 12 weeks, mice underwent survival surgery to remove primary

185 tumors and allow metastases to develop further. Primary tumors were dissociated into single cells, sorted by flow  
186 cytometry to capture the human cancer cells only (CD298+ cells [26]) and immediately captured for scRNAseq.  
187 After an additional 4 weeks *in vivo*, lung, liver and lymph nodes were harvested and CD298+ human tumor  
188 cells were sorted and captured for scRNAseq. This hybrid human-mouse model has a number of important  
189 strengths that are well suited for this study: 1) the confident delineation of tumor cells from the surrounding  
190 microenvironment 2) the capture of matched tumors and metastases and 3) a homogeneous starting population  
191 to control for cell intrinsic differences such as genetic clonality. With this model, we used AAnet to deconvolute  
192 the archetypes underlying heterogeneity in primary tumors. A total of 28,478 cells were analyzed from four  
193 primary-tumors (5,118-8,163 cells) after quality control (Methods).

194 We first examined the archetypes contributing to cellular heterogeneity within primary tumors (Figure 3a).  
195 Each primary tumor was analyzed with AAnet individually (Tumor 1-4), as well as for all primary tumors  
196 combined (All). As hypothesized, each dataset showed a continuum of cellular expression states with multiple  
197 extrema, rather than discrete clusters or unidirectional trajectories (Figure 3b, Supplementary Figure 3, Supple-  
198 mentary Figure 4a). Five archetypes were identified in each dataset, representing distinct biological roles (Figure  
199 3b, Supplementary Figure 3). To elucidate the biology underlying these archetypes, marker genes were defined  
200 and used to identify hallmark genesets upregulated in each archetype and conserved between replicates (FDR <  
201 0.05) (Figure 3d, Supplementary Table 1, Methods). These genes and genesets summarize the biology of each  
202 archetype as follows:

203  
204 **Proliferative archetype (blue)** This archetype is enriched for hallmarks of cell proliferation (Hallmark  
205 genesets G2M checkpoint, E2F targets) and growth (MYC Targets V1/2, mTORC1 signaling). The top markers  
206 of this cluster include CDC20, CDK1 and CDK4, key regulators of phase transitions during the cell cycle [47].  
207 Concomitantly, analysis of cell cycle in the cells most strongly associated with this archetype revealed that >95%  
208 were in either the S or G2M phase (Supplementary Figure 5a-b).

209 **Oxidative/adipogenic archetype (yellow)** This archetype is associated with hallmarks of oxidative  
210 metabolism and stress. Oxidative phosphorylation (OXPHOS) is the most strongly enriched hallmark geneset  
211 across all replicates, driven by electron transport chain components which couple ATP-synthesis to oxygen  
212 availability in the mitochondria (Supplementary Table 1, [5]). MYC targets are also overrepresented, including  
213 many nuclear genes involved in mitogenesis. Genes in the reactive oxygen species (ROS) pathway are among the  
214 top markers of this archetype, including the 4/6 peroxiredoxin family of antioxidant enzymes (PRDX1/2/4/6),  
215 regulated by cancer cells in response to oxidative stress [36]. Genes involved in adipogenesis are also significantly  
216 overrepresented in all replicates, suggesting fatty acid synthesis may be important for this archetype.

217 **Hypoxic archetype (orange)** This archetype was significantly associated with the hypoxia hallmark in all  
218 replicates. While the canonical regulator of hypoxia HIF1A (hypoxia inducible factor 1) was not among the  
219 markers of this archetype, likely because it is degraded post-translationally in the presence of oxygen, CITED2, a  
220 HIF1A induced regulator of hypoxia known to promote both breast cancer, was one of the top markers associated  
221 with this archetype (Supplementary Table 1, [16]). Enrichment of this geneset was also driven by glycolytic  
222 enzymes among marker genes (the glycolysis hallmark enriched in 3/4 replicates). These included class 1 glucose  
223 transporters SLC2A1, SLC2A3 as well as genes linked to oxygen-independent energy production in cancer ENO2,  
224 HK2, LDHA, GAPDH [1]. Ribosomal subunits also featured prominently among marker genes, suggesting a  
225 relationship between ribosome biogenesis and hypoxia.

226 **Cell damage/death archetype (amber)** The top markers for this archetype were genes encoded on the  
227 mitochondrial genome, for which enrichment is associated with cell damage or death (Supplementary Table  
228 1, [30]). Indeed, analysis of the cells most strongly influenced by this archetype showed high expression of  
229 mitochondrially encoded components of the electron transport chain (ETC), yet limited expression of somatically  
230 encoded ETC components (Supplementary Figure 6). In addition to mitochondrial genes, TNF signaling via  
231 NF- $\kappa$ B and the epithelial-to-mesenchymal transition (EMT) hallmark genesets, both associated with cell stress,  
232 were also enriched. This associates this archetype with damaged and dying cells, potentially arising from technical  
233 variables.

234 **Immune-stimulatory archetype (purple)** Immune stimulatory proteins were among the top markers of  
235 this archetype, including CXCL1, IFITM2/3, BST2, HLA-A/B/C, B2M and ICAM1 (Supplementary Table 1).  
236 Concordantly, Hallmark analysis showed enriched genesets related to immune signaling (IFN- $\gamma$ , IFN- $\alpha$ , TNF- $\alpha$ ,

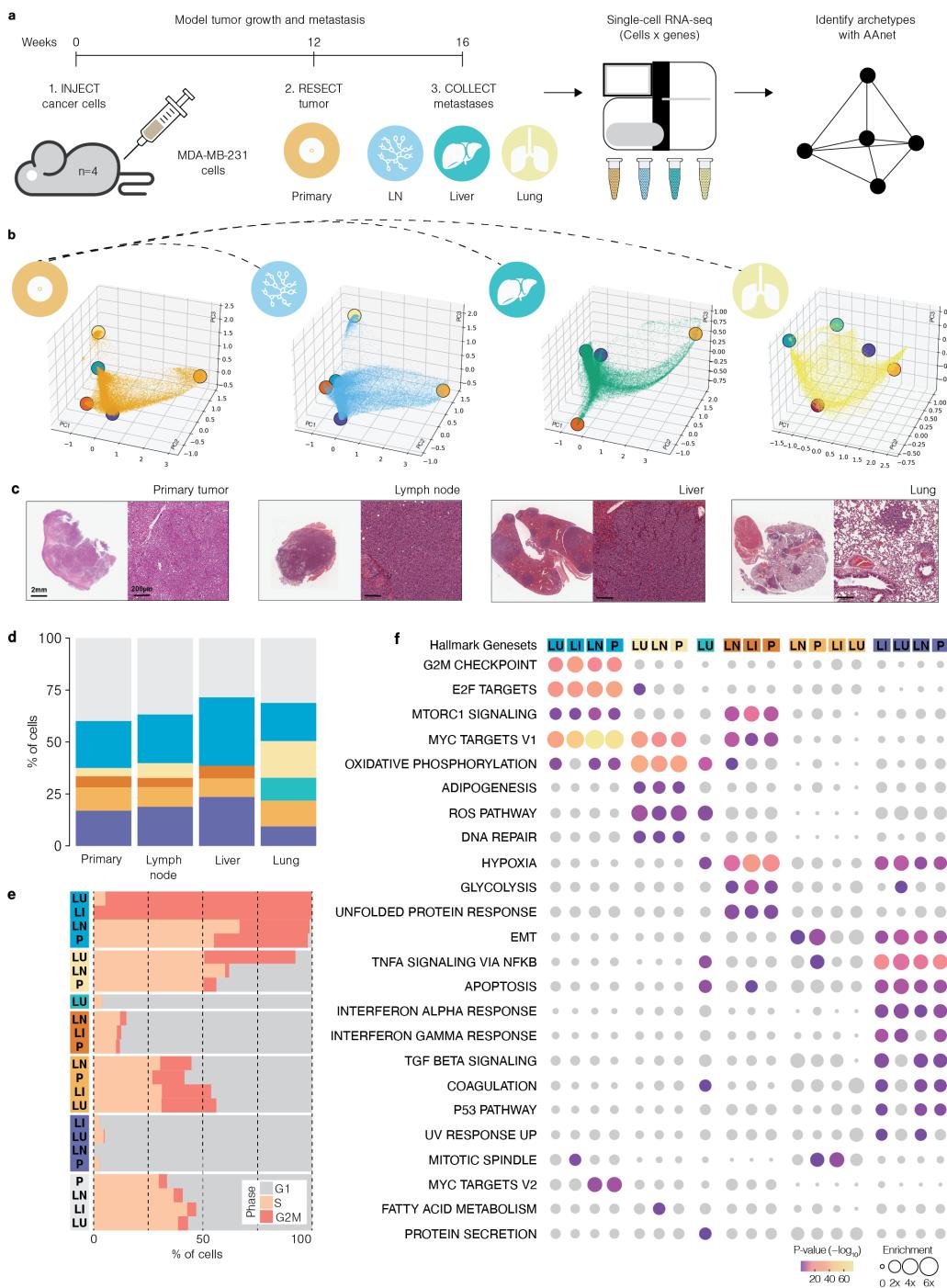
237 and TGF- $\beta$  pathways), and apoptosis (Apoptosis, p53). Analysis of cell cycle among the cells most strongly  
238 influenced by this archetype showed 96% percent of cells in G1, perhaps suggestive of a G1 arrest (Supplementary  
239 Figure 5a-b). These analyses indicate this archetype describes an apoptotic, immune-stimulatory expression  
240 pattern in tumor cells.

241  
242 Archetypes were highly conserved between primary tumors (Figure 3c). All archetypes had an orthologous  
243 archetype in each individual primary tumor and the combined dataset, with a mean cosine similarity 0.77-0.97.  
244 In contrast, individual archetypes showed modest pairwise cosine similarity (between -0.15 and -0.2), consistent  
245 with their classification of representing distinct cell states.

246 Having deconvoluted tumor cell heterogeneity into five biologically meaningful archetypes, we then determined  
247 the association of each cell with each archetype using AAnet. AAnet encodes a representation of each cell based  
248 on its relative association to each archetype, where the coordinates of each cell are non-negative and sum to  
249 one. We term this association "archetypal affinity" and define cells with affinity for one archetype greater than  
250 the sum of all other archetypal affinities (i.e. affinity > 0.5) as "committed". Cells that do not surpass this  
251 affinity threshold for any archetype can be considered "uncommitted". Using archetypal commitment to analyze  
252 the cellular composition of primary tumors, the abundance of cells committed to an archetype was consistently  
253 above 60% in all tumors (62%-67.8%), yet the abundance of committed cells varied between archetypes (Figure  
254 3e). The proliferative archetype was the most abundant archetype in the combined analysis and in two of four  
255 individual tumors, and the immune-stimulatory archetype was the second most abundant archetype in the  
256 combined analysis and in two of four individual tumors. The oxidative/adipogenic archetype represented a minor  
257 fraction of tumor cells across all samples, and the hypoxic archetype was among the most variable. Together,  
258 these analyses show that AAnet deconvolutes cancer cell heterogeneity into biologically meaningful archetypes  
259 that are reproducibly detected in discrete tumors.

260 **AAnet reveals conserved and *de novo* heterogeneity across distinct metastatic sites**

261 **Figure 4**



**Figure 4.** (a) Approach to define tumor archetypes in TNBC metastases from the lymph node (LN), liver and lung. (b) Embedding of tumor cells from each tissue with ATs indicated by colored circles. (c) Representative histological sections of tumors and metastases taken for single-cell sequencing. (d) Stacked barplot showing percentage of tumor cells in each tissue committed to each archetype (colored) or uncommitted to an archetype (gray). (e) Stacked bar plot showing the cell cycle phase of cells committed to each archetype. (f) Enriched HALLMARK genesets associated with ATs from each tissue (indicated by colors). P-value is false discovery rate corrected, genesets where  $p > 0.1$  are shown in gray. Enrichment represents log<sub>2</sub> fold-change. P = primary, LN = lymph node, LI = liver, LU = lung.

262 Next, we used AAnet to identify archetypes in metastases. scRNASeq was performed on matched lymph node  
263 (LN), liver, and lung metastases that were collected four weeks after resection of the primary tumor (Figure  
264 4a-c). A total of 42,250 metastatic cells were analyzed from LN ( $n = 4,660$ - $19,224$  cells per tumor), 42,647 cells  
265 from liver ( $n = 4,713$ - $18,671$  cells per tumor), and 17,687 cells from the lung ( $n = 3,559$ - $6,379$  cells per tumor)  
266 after removal of one outlier lung sample during QC (Methods, Supplementary Figure 7a). We combined data  
267 from metastases per tissue separately and defined archetypes in each site using AAnet.

268 Cellular heterogeneity in metastases followed a phenotypic continuum akin to that observed in primary  
269 tumors (Figure 4b, Supplementary Figure 4b-d). To characterize this continuum, AAnet defined five archetypes  
270 in the lymph node, four in the liver, and five in the lung, where all archetypes detected in the primary had an  
271 ortholog in at least two metastatic sites (Supplementary Figure 7b-c). Moreover, a highly conserved pattern  
272 of hallmark enrichment was observed in orthologous archetypes from different sites (Figure 4f, Supplementary  
273 Table 2). This demonstrates that the factors of variation contributing to cancer cell heterogeneity were largely  
274 recapitulated in primary tumors and metastases.

275 While orthologous archetypes were identified in many different metastases, analysis of archetypal affinity  
276 indicated their relative contribution to cellular heterogeneity in each tissue was somewhat distinct.

277 Lymph node metastases showed the strongest resemblance to primary tumors. All five archetypes identified  
278 in the lymph node had an ortholog in the primary tumor (Figure 4d-f). The proportions of cells committed to  
279 each archetype were also highly conserved, with the proliferative and immune-stimulatory archetypes the most  
280 abundant in both tissues (Figure 4d). Moreover, the changes in the relative abundance of the oxidative, hypoxic  
281 and cell death archetypes were marginal, collectively indicating cell heterogeneity is very similar between primary  
282 tumors and lymph node metastases.

283 Liver metastases differed from the primary and lymph node in the number and abundance of archetypes.  
284 Only four of the five archetypes identified in the primary tumor were detected in the liver, and more cells were  
285 committed to an archetype than in other tissues (Figure 4d-f). The absence of cells committed to the oxidative  
286 archetype was replaced by a greater proportion of cells committed to the proliferative and immune stimulatory  
287 archetypes, both nearly double their proportions in the primary tumor (32.94% and 23.54% respectively). Overall,  
288 while the archetypes were somewhat conserved with primary tumors and lymph node metastases, their relative  
289 contribution to cellular heterogeneity distinguished liver metastases from the other tissues.

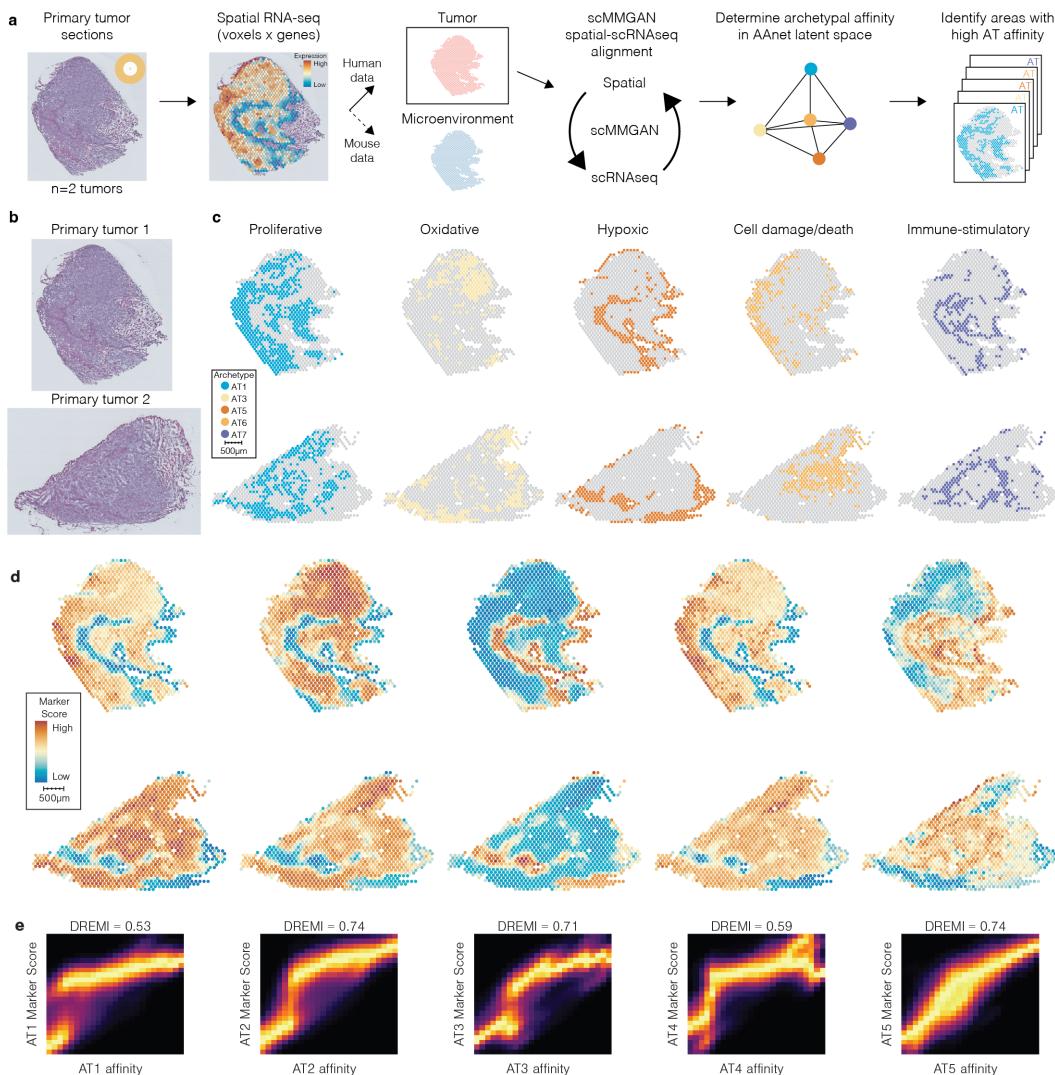
290 Lung metastases showed the greatest difference from primary tumors and the other metastases. Five archetypes  
291 were detected in the lung, four of which were orthologous to archetypes in the primary tumor (Figure 4d-f).  
292 While still the most abundant, the proportion of cells committed to the proliferative archetype were significantly  
293 lower than in other tissues (18.36%). The oxidative archetype, enriched for OXPHOS and ROS hallmarks, was  
294 the second most abundant and almost equal to the proliferative archetype (17.72%). Conversely, the hypoxic  
295 archetype was not detected in the lung metastases. This is consistent with cancer cellular heterogeneity in these  
296 metastases being shaped by a oxygen-rich lung microenvironment.

297 The remaining archetype was unique to the lung. Markers of this archetype were significantly enriched  
298 for the cancer hallmarks of an oxidative metabolism (OXHPOS and ROS pathway) driven by genes encoding  
299 ETC components (ATP5F1A,B, NDUFA4,8,9, COX7C) and antioxidant enzymes (PRDX2,4,6) respectively  
300 (Supplementary Table 2). Additionally, hallmarks of TNF signaling via NF- $\kappa$ B, apoptosis, hypoxia, protein  
301 secretion, and coagulation were overrepresented, a profile of pathways indicating similarity to the immune  
302 stimulatory archetype. Importantly, enrichment of hypoxia was driven by a separate subset of glycolytic enzymes  
303 to those associated with the hypoxic archetype in other tissues, and did not include markers of HIF1A induction  
304 such as CITED2 (Supplementary Table 2).

305 Together, these analyses show that AAnet-defined archetypes in primary tumors are also identified in  
306 metastases. This suggests that factors influencing cellular heterogeneity are highly conserved between primary  
307 tumors and metastases. Differences between tissues were largely driven by the differential influence of these  
308 archetypes, potentially due to interactions with the metastatic microenvironment.

309 Spatial Transcriptomics reveals organization and distinct cellular morphology of AAnet  
 310 archetypes

311 **Figure 5**



**Figure 5.** (a) Approach used to define a spatial map of TNBC tumor archetypes using scMMGAN and AAnet. (b) Histological sections of TNBC xenografts used for spatial transcriptome sequencing ( $n=2$ ). (c) Spatial voxels with a strong affinity ( $>0.3$ ) for each archetype (colored) as determined by scMMGAN and AAnet. (d) Spatial plots colored by marker geneset score for archetype. (e) Density resampled estimate of Mutual Information (DREMI) score between archetypal affinity and marker geneset scores.

312 To further validate the significance and biology of the archetypes identified by AAnet, we sought to investigate  
 313 their structural organization within tumors. Thus, we performed spatial transcriptomics (Visium 10X) on tissue  
 314 sections from two primary tumors that were grown *in vivo* for 8 weeks. These samples were collected from two  
 315 of the tumors used for scRNAseq, with parts of these tumors frozen in OCT for sectioning. Gene expression was  
 316 measured at 2,275 spatially distributed voxels in the two samples after QC (1,170 and 1,105 voxels respectively),  
 317 with each voxel assaying expression in an area approximately 3-10 cells in size.

318 To map the archetypes from the scRNAseq to the spatial transcriptomic data, we first had to overcome  
 319 the intrinsic differences in the data generated by these modalities. Raw data from these technologies is not  
 320 directly comparable as they have very different sample processing protocols (digestion and droplets vs freezing  
 321 and sectioning) and biological resolutions (single cells vs multiple cells). These differences create batch effect,  
 322 where noise introduced in the process of data generation dominates the biological differences between samples.

323 This batch effect was evident when spatial voxels were embedded with the scRNAseq data from the primary  
324 tumor, as there was little alignment between the modalities despite being generated from matched biological  
325 samples (Supplementary Figure 8a).

326 To address this, we leveraged our single-cell multi-modal generative adversarial network (scMMGAN) [4]  
327 (Methods). This approach uses adversarial learning to align data between modalities, enabling integrated  
328 downstream analysis while preserving data geometry (Figure 5a). We used scMMGAN to generate a scRNAseq  
329 measurement for each spatial voxel (scVoxels) that is aligned to our primary scRNAseq dataset. We can then  
330 input the aligned spatial data into the AAnet encoder trained on the primary scRNAseq data. This provides an  
331 archetypal representation for each spatial voxel based on our previously characterized primary archetypes (Figure  
332 5c, Supplementary Figure 8b-c). The ability to extend archetypal analysis to previously unseen data through the  
333 neural network framework is an important feature of AAnet versus existing archetypal analysis approaches.

334 With spatial transcriptomic data embedded as scVoxels in the AAnet latent space, we calculated their  
335 affinity to each archetype which was mapped back to their corresponding spatial location in the primary tumor  
336 (Methods, Supplementary Figure 8). Archetypal affinities and marker geneset expression scores shared high  
337 mutual information across voxels (Methods, Figure 5d-e). This indicates that biology that defined each archetype  
338 had been retained in the spatial mapping process.

339 Archetypes showed spatial organization and were associated with distinct cell morphologies in primary tumors:

340     **Proliferative archetype (blue)** Voxels with high affinity for the proliferative archetype (affinity > 0.3)  
341 formed the bulk of the tumor yet were markedly absent from the central areas of each tumor section. Areas with  
342 high affinity for this archetype were enriched for cycling cells (Supplementary Figure 5c). These findings were  
343 consistent with the scRNAseq commitment analysis (Figure 3e) indicating that proliferation was a dominant  
344 factor of cell heterogeneity in the primary samples.

345     **Oxidative/adipogenic archetype (yellow)** Areas of the tumor with high affinity for the oxidative  
346 archetype were located in close proximity to the proliferative archetype (Figure 5d). The affinity scores and  
347 expression of marker genesets were also significantly correlated (Figure 5e-f), indicating a relationship between  
348 oxidative metabolism and a proliferative cell state.

349     **Hypoxic archetype (orange)** Strikingly, areas with a strong affinity for the hypoxic archetype were  
350 localized to central and peripheral regions of primary tumors, devoid of the proliferative and oxidative archetypes  
351 (Figure 5c). These areas showed enriched expression of markers associated with oxygen-independent glycolysis  
352 and ribosomal subunits (Figure 5d-e).

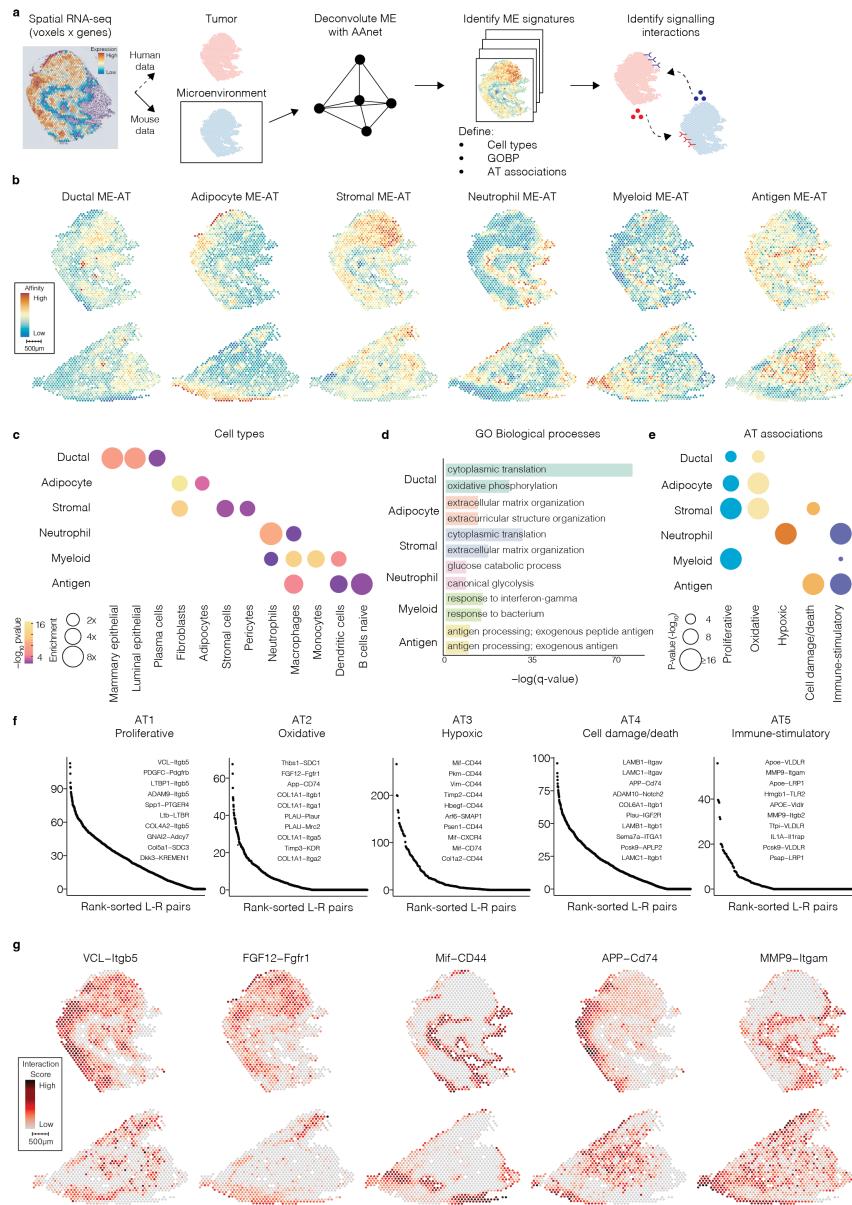
353     **Cell damage/death archetype (amber)** The cell death archetype localized to areas with high expression  
354 of mitochondrial genes and was strongly correlated with the proliferative and oxidative archetypes (Figure  
355 5c-e). While preferential enrichment of mitochondrially encoded genes is often used as a marker of cell death,  
356 mitochondrial genes encode critical components of the electron transport chain necessary for oxidative energy  
357 production. Consequently, the association between this archetype and the proliferative/oxidative archetypes may  
358 identify oxygen-rich areas within primary tumors rather than cell death.

359     **Immune-stimulatory archetype (purple)** Notably, affinity for the immune-stimulatory archetype was  
360 highest surrounding the hypoxic archetype within the tumor. These areas showed enriched expression (Figure 5c)  
361 of cytokines and antigen presenting proteins, such as CXCL1 and HLA-A/B/C/B2M (Figure 5d-e), consistent  
362 with the hallmark associations of this archetype. Interestingly, the immune-stimulatory archetype appears to  
363 demarcate the hypoxic and proliferative cancer cell archetypes.

364  
365     Together, this analysis shows identified archetypes have unique and distinct spatial organization within the  
366 tumor, further validating the ability of AAnet to identify unique cellular biology and structure of cancer cells  
367 within a phenotypic continuum of cell states. Further, the organization of the AAnet archetypes is consistent  
368 with a model whereby the local microenvironment may play a critical role in determining the organization of  
369 cellular heterogeneity in primary tumors.

371 Archetypes are associated with distinct cell types and metabolic niches in the microenvi-  
 372 ronment

373 **Figure 6**



**Figure 6.** (a) Approach to defining tumor microenvironments (ME) associated with key archetypes (AT) using AAnet. (b) Microenvironment-archetype (ME-AT) affinity scores for each spatial voxel. (c) Enrichment of cell-types associated with each ME-AT. P-value is false discovery rate corrected. Enrichment represents log2 fold-change. (d) Enrichment of top 2 biological processes from the gene ontology (GO) associated with each ME-AT. Processes are ranked by FDR (-log10). (e) ME-ATs (y-axis) significantly associated with each tumor AT (x-axis). Color reflects tumor AT, enrichment represents log2 fold-change. (f) Top ranked ligand-receptor pairs expressed in spatial voxels with a strong affinity for each archetype and their colocalized microenvironments. Capitalized gene symbols indicate genes with expression in tumor cells (human). Title case symbols indicate genes with expression in ME cells (mouse). Pairs are ranked by FDR (-log10). (g) Ligand-receptor interaction score across spatial voxels for top pairs associated with each archetype.

374 With clear spatial organization of cancer cell archetypes derived by AAnet within the tumor, we next sought  
 375 to determine if, beyond spatial location, the tumor microenvironment may also be playing a role in archetypal

development and commitment. We therefore assessed if microenvironmental cells were spatially structured into meaningful microenvironmental archetypes (ME-ATs) (Figure 6a). The xenograft model enabled separation of expression at each voxel into tumor and ME data based on alignment to the human or mouse genome, respectively. Then, we used AAnet to deconvolute the murine data into ME-ATs to investigate archetype-specific cell types and biological processes.

AAnet defined six ME-ATs with unique patterns of spatial organization (Figure 6b, Methods). Each ME-AT is enriched for specific cell types and biological processes (Figure 6c-d). These include a ductal ME-AT with high enrichment of mammary epithelial cell markers in voxels overlaying to breast ducts, as well as an adipocyte ME-AT enriched for adipocyte markers and most highly expressed on the margins of the tumor sections with residual mammary fat pad. This correspondence with underlying histological features provides orthogonal validation for AAnet. AAnet also identified a stromal ME-AT, enriched for fibroblasts, capillary-lining pericytes, and stromal cells, as well as biological processes related to ECM organization, oxidative phosphorylation and angiogenesis; a neutrophil ME-AT, enriched for neutrophils and biological processes related to glycolysis, leukocyte chemotaxis and hypoxia; a myeloid ME-AT enriched for markers of monocytes, macrophages and dendritic cells, with an enriched response to interferon-gamma, immune effector processes and leukocyte mediated cytotoxicity; and an antigen-response ME-AT enriched for macrophages and genes involved in antigen processing and presentation of exogenous peptide antigen via MHC class II, and immunoglobulin-mediated immune response (Figure 6c-d). Therefore, AAnet deconvolutes expression in the microenvironment into stromal and immune components.

Next, we explored the spatial relationship of the ME-ATs to the cancer cell ATs and uncovered specific associations (Figure 6e, Methods). Areas of the tumor with a high affinity for the proliferative AT were strongly associated with the myeloid and stromal ME-ATs, while the oxidative/adipogenic AT showed preferential enrichment for the stromal ME-AT. This suggests these archetypes colocalize with highly metabolic and vascularized microenvironments. Both of these archetypes were also associated with the ductal and fat-pad ME-ATs, concordant with the localization of these normal mammary gland structures in the histological sections. Notably, the hypoxic AT colocalized specifically with the neutrophil ME-AT. These interactions occurred at internal parts of the tumor section, indicating that these archetypes are associated with decreased oxygen availability and increased neutrophil chemotaxis. The cell death AT was strongly associated with the antigen ME-AT, indicating an active presentation of cancer cells to the immune system, cancer cell death and phagocytosis of dead and dying cells by macrophages. Interestingly, the immune-stimulatory cancer cell AT, defined by its enrichment for gene sets related to immune signaling, was associated with the antigen and neutrophil ME-ATs, indicating a strong overlap in signaling and cell phenotype between the cancer and microenvironmental cells. Thus, we observed multiple examples of spatially-localized phenotypic mimicry between cancer and microenvironmental cells, most notably in metabolic and immune signaling.

To investigate direct cell-cell signaling mechanisms for phenotypic mimicry, we analyzed ligand-receptor pairs (LR-pairs) for evidence of paracrine interactions (Methods). We again used our hybrid model system to delineate between tumor (human) and ME (mouse) expression, allowing us to establish the direction of signaling. Specifically, LR-pairs with a ligand expressed in the human data and its cognate receptor expressed in the mouse data indicate signaling from the tumor to the microenvironment, and vice versa (Figure 6f). We first determined the coexpression of annotated LR-pairs across all voxels and calculated their enrichment in areas with a high affinity for each archetype. For all archetypes, the proportion of ligands originating from the tumor and the microenvironment were approximately 50%. However, the strongest LR-pairs that localized to high affinity regions differed between archetypes.

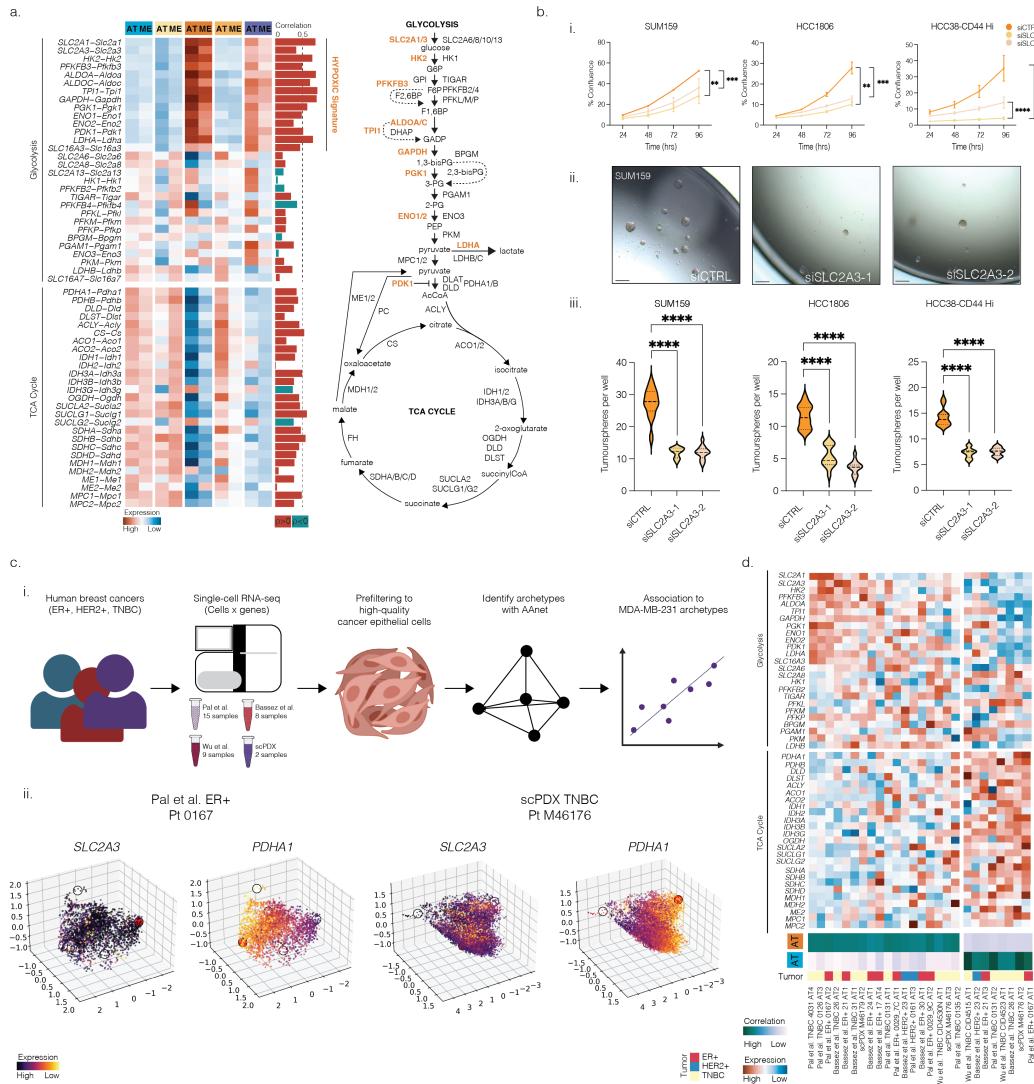
Tumor cells expressed stromal growth factors in areas with a high affinity for proliferative and oxidative/adipogenic archetypes. These included platelet-derived growth factor (AT Ligand: PDGFC, ME Receptor: Pdgfrb) and fibroblast growth factor (AT Ligand: FGF12, ME Receptor: Fgfr1) pathways that promote blood vessel formation [12]. The hypoxic archetype was characterized by interactions between CD44 expressed in tumor cells and a variety of ligands expressed in the microenvironment, including Mif, which may enhance neutrophil accumulation within the hypoxic AT. Examples of cancer-microenvironmental cross talk in the cell death AT are evidenced by App-Cd74, where Cd74 is an MHCII molecule, and sema7a-ITGA1, where sema7a is a potent immune modulator. We also observe strong overlap in LR-pairs regulating immune-cancer cell crosstalk in the immune-stimulatory AT including MMP-itgam, IL1A-Ilr1rap, and Hmgb1-TLR2. (Figure 6g).

Together, these results highlight the spatial colocalization of distinct cancer archetypes with unique microen-

428 environments, where paracrine interactions that may enhance phenotypic mimicry are an important determinant of  
 429 intratumoral heterogeneity.

## 430 Intratumoral metabolic heterogeneity in TNBC and alignment to human breast cancers

431 **Figure 7**



**Figure 7. (a)** Heatmap showing expression of metabolic genes (glycolysis and oxidative phosphorylation) across cancer archetypes and associated microenvironmental cells. **(b)** i. Cell growth of HCC1806, SUM159 and HCC38-CD44 Hi cells treated with either control siRNA (siCTRL) or siRNA targeting SLC2A3 (siSLC2A3-1, siSLC2A3-2) over a 96-hour period.  $n = 3$  independent experiments, triplicate wells analyzed per condition. Statistical significance defined by one-way ANOVA. ii. Representative images of SUM159 tumorspheres treated with control or SLC2A3 targeted siRNA. iii. Quantitation of tumorsphere-forming capacity in control (siCTRL) and SLC2A3 knockdown cells (siSLC2A3-1 or siSLC2A3-2) in SUM159, HCC1806 and HCC38-CD44 Hi cells.  $n = 3$  independent experiments, 12 wells analyzed per condition, statistical significance defined by ordinary one-way ANOVA with Tukey's multiple comparison post-hoc analysis. **(c)** i. Diagrammatic representation of human breast cancer sample single-cell analysis. ii. Visualization of cancer epithelial cells from two patients, colored by markers for oxygen-independent glycolysis (*SLC2A3*) and entry into the TCA cycle (*PDHA1*). **(d)** Heatmap showing expression of metabolic genes from (a) across 26 human cancer archetypes, 18 (left) associated with the hypoxic archetype and 8 (right) associated with the proliferative archetype. Archetypes are from breast cancer samples across three different breast cancer subtypes.

The discrete localization of distinct metabolic phenotypes within the tumor, including the oxidative phosphorylation-enriched proliferative AT versus the glycolytic-enriched hypoxic AT (Figure 3d and Figure 4f), as well as the metabolic mimicry of the microenvironmental cells in those regions (Figure 6d), led us to ask if targeting a specific metabolic program might impact tumor growth. To delve into this question, we first examined the metabolic state of the microenvironments colocalized with each archetype by comparing genes associated with glycolysis and the TCA cycle in each cancer AT and associated microenvironment. Indeed, the concordance between the cancer AT and the microenvironment was driven by the correlation in expression of metabolic enzymes (Figure 7a). Enzymes in the tricarboxylic acid cycle (TCA-cycle), a pathway which requires oxygen to generate energy, were uniformly highly expressed in the tumor and ME of the proliferative, oxidative and cell death archetypes. In contrast, glycolytic enzymes were clearly enriched in the hypoxic AT and associated ME cells (Figure 7a). Enzymes most highly enriched in the hypoxic area were associated with oxygen-independent glycolysis in tumors [1] and include PDK1, an inhibitor for the entry of pyruvate into the TCA cycle. This indicates that the metabolic heterogeneity in primary tumor archetypes is mirrored in their local microenvironments, and the intratumoral hypoxic regions of the tumor are driven by glycolysis ending in the accumulation of lactate. Interestingly, hypoxic niches are known to provide a permissive environment for maintenance of both pluripotent stem cell and cancer stem cell populations [34, 53]. AAnet identified that SLC2A1 (GLUT1) and SLC2A3 (GLUT3) are enriched in cancer cells that reside in that niche.

Given that GLUT1 is ubiquitously expressed in normal cells throughout the body, we asked if ablating GLUT3, whose expression is largely confined to the brain and sperm in normal tissues, could eradicate the aggressive phenotype of the cancer cells in the hypoxic niche. In addition, GLUT3 expression is increased in TNBC and associated with metastasis and poor prognosis [44]. Three TNBC cancer stem cell-enriched cell lines (SUM159, HCC1806 and HCC38-CD44Hi) were treated with a control siRNA (siCTRL) or two independent siRNAs targeting SLC2A3/GLUT3 (siSLC2A3-1 and siSLC2A3-2). Efficient knockdown of SLC2A3 was confirmed by qPCR (Supplementary Figure 9). SLC2A3 knockdown significantly inhibited cell proliferation in all cell lines tested (Figure 7b, Supplementary Figure 9). Excitingly, we confirm that SLC2A3 knockdown significantly inhibits tumorsphere formation, an *in vitro* surrogate assay for *in vivo* tumor-initiating ability (Figure 7b). Together, these results suggest that SLC2A3 is critical for maintenance of the cancer stem cell phenotype in TNBC and add to previous data indicating a role for GLUT3 in EMT and migration [44].

To examine the clinical relevance of the identified archetypes, we analyzed single-cell transcriptomes from human breast cancer cells across four distinct studies, corresponding to 34 samples from three major breast cancer subtypes (ER+, HER2+, TNBC) (Figure 7c) [6, 35, 51]. First, AAnet identifies 155 archetypes across the 34 samples and shows interesting similarities and differences across samples, cohorts, and breast cancer subtypes (Methods, Supplementary Figure 10). To further investigate the association of human archetypes with the proliferative AT (blue) and hypoxic AT (orange), we identify 18 archetypes across all human tumors that have transcriptomic profiles similar to AT3 (cosine similarity > 0.25) and dissimilar to AT1 (cosine similarity < -0.25), as well as 8 archetypes similar to AT1 and dissimilar to AT3 (Figure 7d). These metabolic archetypes are common in breast cancer, represented in 26 archetypes spanning 20 human tumors. Similar to the xenograft model, 6 tumors contain both a hypoxic AT and a proliferative AT within the same sample. Together, these results show correspondence between the metabolic profiles of archetypes identified by our model and human breast cancer tumors, and further, the identified archetypes are relevant across breast cancer subtypes.

Visualization of the metabolic markers from Figure 7a reveals, for human archetypes associated with hypoxia, enrichment for the hypoxic signature within the glycolysis pathway and low expression of genes related to the TCA cycle. Conversely, human archetypes similar to the proliferative archetype showed low expression for hypoxic genes and higher enrichment for TCA cycle genes.

Notably, there is no significant enrichment for a particular cancer subtype and association with AT1 or AT3 (KS test p>0.05 for all tests), nor significant difference between cancer subtypes in proportion of cells committed to hypoxic or proliferative archetypes (Wilcoxon rank sums test p>0.05 for all tests). These data show that AAnet can be used to identify phenotypic similarities across cancer subtypes, and thereby offers a functional method beyond hormone and molecular subtyping to classify breast cancers for therapeutic targeting.

## 481 Discussion

482 It is now recognized that non-genetic programs (e.g., epigenetic, transcriptional, translational) are a major driver  
483 of tumor heterogeneity. The dynamic and reversible nature of non-genetic heterogeneity likely favors rapid  
484 evolution of cancer cell states (e.g., seconds, minutes or hours) to enable survival in unfavorable microenvironments  
485 encountered throughout the metastatic cascade and in response to therapy. Thus, as single-cell technologies  
486 continue to resolve the breadth and structure of non-genetic heterogeneity in cancer and stromal cells within  
487 and across patient tumors, developing strategies to identify and validate the specific cell states and molecular  
488 mechanisms that fuel cancer progression remains a significant technological and biological challenge.

489 To address this knowledge gap, we developed AAnet, an archetypal analysis method to identify archetypal cell  
490 states within and between samples and their associated biological processes. Archetypal analysis is a framework  
491 to describe a dataset as a convex combination of extreme, or archetypal, observations. In contrast to other  
492 unsupervised approaches to characterize such data, archetypal analysis is aptly suited for both identifying key  
493 cell states reflecting distinct biological processes and analyzing the cellular state space as a continuum of cells  
494 committed to these processes. However, identifying the archetypal states remains a fundamental challenge of  
495 archetypal analysis. In particular, nonlinearities can worsen the performance of existing archetypal analysis tools,  
496 as the extreme states of the data geometry do not conform to the extreme states of the data space.

497 AAnet solves this problem by learning a transformation of the data into a simplex, rather than fitting a  
498 simplex on the data directly. The latent space of the autoencoder thus preserves relationships between cells and  
499 characterizes cells by their commitment to each archetype. We further regularize the latent space to initialize the  
500 archetypes to diffusion extrema (inferred from the cell-cell affinity graph) for improved accuracy and robustness.

501 Applied here to single-cell data from pre-clinical and clinical breast cancer samples, we have shown that  
502 AAnet enabled the discovery of biologically and functionally-distinct archetypes within a phenotypic continuum  
503 of cell states within the tumor and captured their associated molecular drivers. First, in a pre-clinical xenograft  
504 model comprising primary tumors matched with lung, liver, and lymph node metastases, we identify six unique  
505 archetypes across primary and metastatic tissues, with each archetype defining unique biology of cells committed  
506 to that extrema. Interestingly, we show that the number and distribution of cells committed to archetypes in  
507 primary tumors is remarkably similar to those found in lymph node metastases, yet liver metastases differ by  
508 the loss of one archetype, and the lung metastases deviate from primary, lymph node and liver metastases via  
509 the emergence of one new archetype. These analyses demonstrate that AAnet can reveal the emergence of new  
510 cell states, the number of cells committed to a specific archetype, and the underlying biology that facilitate  
511 site-specific metastatic adaptation.

512 Critically, we validate the significance of the archetypes identified by mapping the scRNAseq data to matched  
513 spatial transcriptomic data via scMMGAN [4]. These data confirm that AAnet-defined archetypes resolve  
514 into distinct spatially-localized regions within the tumor. Further, we show that AAnet has revealed a unique  
515 perspective on the organization of the associated microenvironments. Specifically, each archetype is enriched  
516 with distinct stromal cell types; for example, the proliferative archetype is enriched with fibroblasts, hypoxic  
517 archetype with neutrophils, and immune-stimulatory archetype with macrophages and dendritic cells. Thus,  
518 AAnet robustly identifies functional and spatially distinct cellular archetypes within a tumor.

519 Of note, we uncovered metabolic heterogeneity not seen before in TNBC, where cells in discrete archetypes  
520 utilize distinct metabolic programs. We have recently analyzed bulk RNA-seq data (METABRIC and TCGA)  
521 and shown that TNBC exhibit a unique highly metabolic gene expression phenotype, upregulating a range of  
522 pathways, including glycolysis, compared to other breast cancer subtypes such as Luminal A [37]. Our archetypal  
523 analysis provides further granularity to these data, showing the individual contributions of each archetype to  
524 this unique TNBC metabolic signature. For example, the hypoxic archetype clearly contributes to the high  
525 expression of SLC2A1, SLC2A3 (GLUT3), HK1, HK2, ALDOA, ALDOC, TPI, GAPDH, PGK1, ENO1, PDK1  
526 and LDHA in the bulk RNAseq data, with contributions also from the immune-stimulatory archetype but not  
527 from the most abundant proliferative archetype. By comparison, these findings reveal that different regions  
528 of the tumor uniquely invoke glycolysis or oxidative phosphorylation, which could not be determined when  
529 analyzing existing bulk RNA-seq data. Furthermore, we show that we could use our identified archetypes to  
530 predict novel therapeutic targets within distinct cellular subsets, such as the glucose transporter GLUT3 in stem  
531 cells within the hypoxic archetype. Interestingly, we also discovered that the distinct metabolic phenotypes of two

532 archetypes are strikingly reflected in the microenvironmental cells associated with those archetypes. Importantly,  
533 we also found these distinct archetypes are present in human breast cancer samples. GLUT3/SLC2A3 expression  
534 was present at highest levels in the hypoxic archetype, with high levels also seen in the surrounding immune  
535 archetype. A previous study has suggested a role for GLUT3 in regulating the inflammatory microenvironment  
536 in TNBC [44], which may also play a role in the matched metabolic phenotypes of the tumor and surrounding  
537 microenvironment cells.

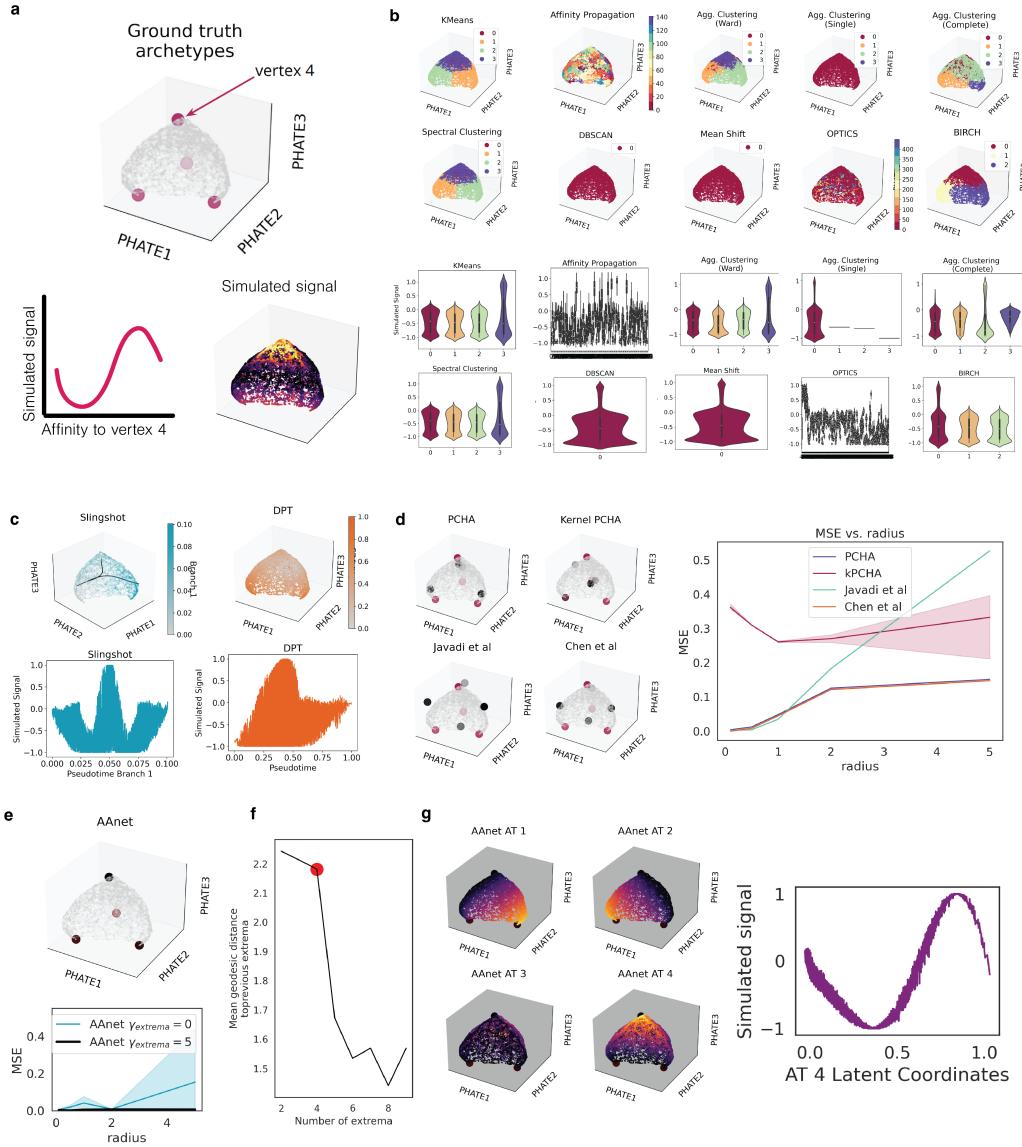
538 Tumor heterogeneity remains a significant clinical challenge for diagnosis and therapeutic management. The  
539 discovery of the AAnet archetypes in scRNAseq data has enabled segmentation of tumors into functionally  
540 distinct regions comprising cancer cell states associated with unique cellular microenvironments. Together,  
541 these data resolve tumor heterogeneity to a level not yet achieved with previous computational tools. In the  
542 future, classifying patients according to biological archetypes with tools like AAnet is likely to improve tumor  
543 sub-classifications. Applied to samples before and after specific treatment, we can begin to learn how archetypes  
544 change over time, in response to specific therapies, and in different metastatic sites. Moreover, this approach will  
545 reveal the molecular programs driving each cellular archetype, as well as when and how they emerge. Ultimately,  
546 these tools will deliver the knowledge to enable the development of improved and effective therapeutic strategies.

547 Importantly, AAnet is a flexible framework that can be used both independently and as a part of a large  
548 single-cell analysis pipeline to interpret the archetypal distribution underlying any single-cell dataset. Given its  
549 widespread utility and generalizability to characterize cells, AAnet is a valuable tool for the single-cell community.

## 550 Figures

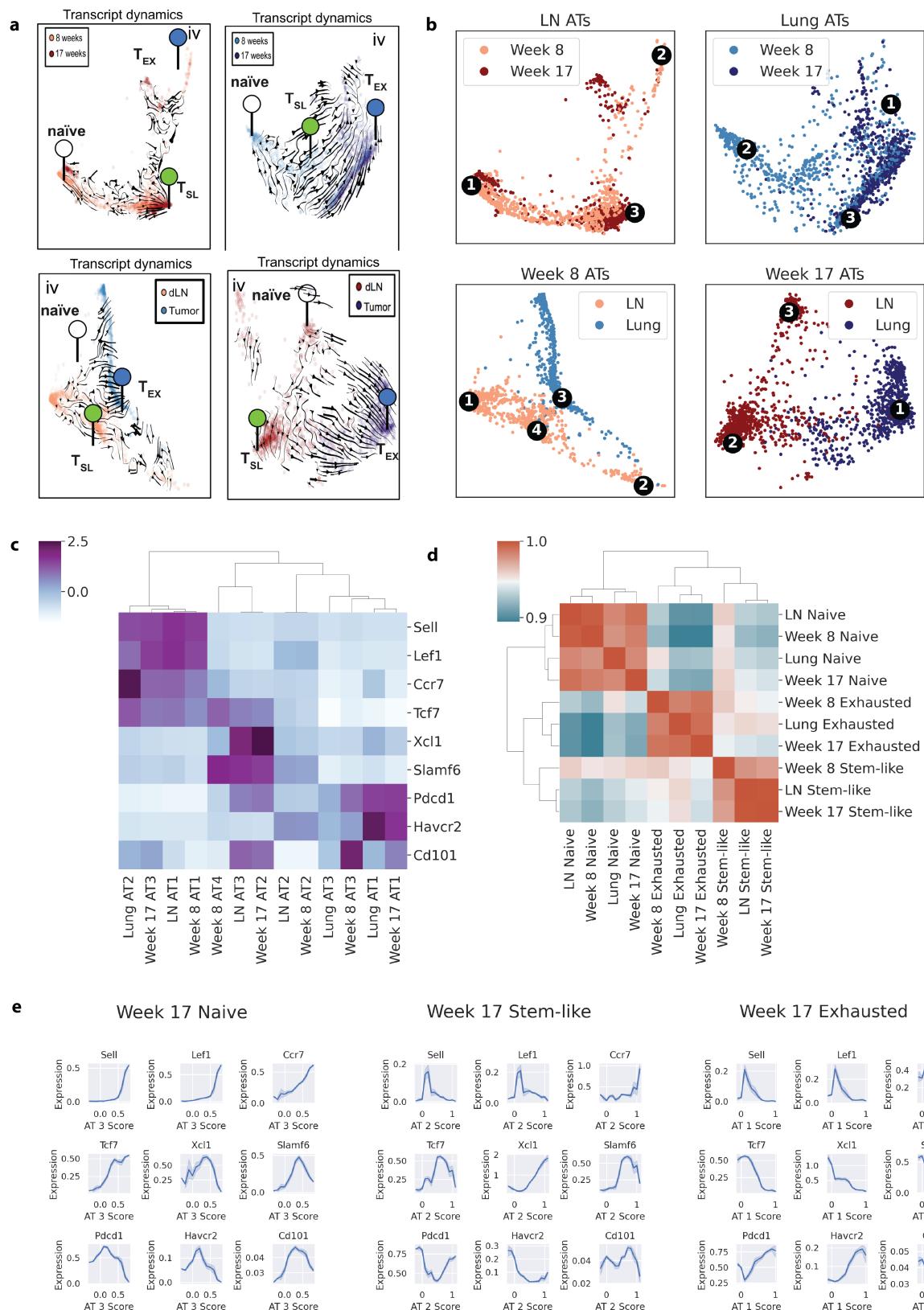
### 551 Supplemental Figures

#### 552 Supplementary Figure 1



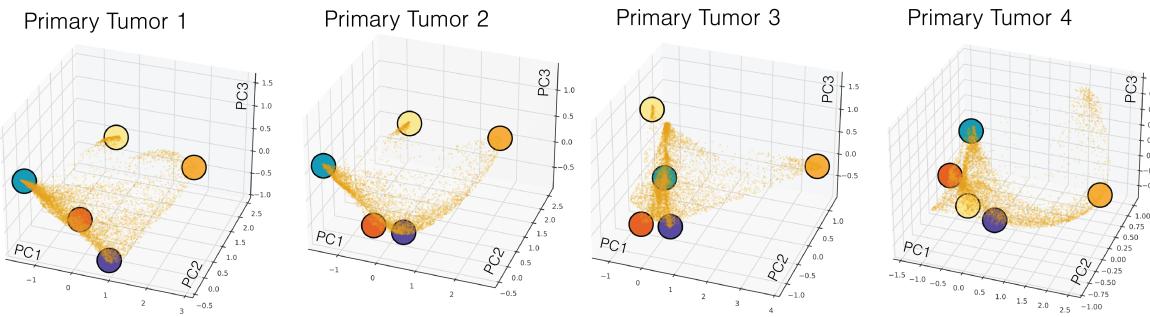
**Supplementary Figure 1. Curved Tetrahedron example.** (a) Simulated curved tetrahedron, colored by signal defined with respect to affinity to archetype 4 (vertex at top). (b) Cluster assignments and signal comparison for 10 clustering algorithms. (c) Trajectories and signal comparison from Slingshot with KMeans clusters as input (left) and diffusion pseudotime (right). (d) Inferred archetypes from each archetypal analysis method (black) versus ground truth archetypes (red) and mean squared error between real and ground truth archetypes over increasing curvature. (e) AANet-learned archetypes and mean squared error between real and ground truth archetypes over increasing curvature. (f) AANet-inferred number of archetypes. (g) AANet-learned archetypal coordinates. AANet recapitulates sine signal with respect to archetype 4's latent coordinates.

553 Supplementary Figure 2



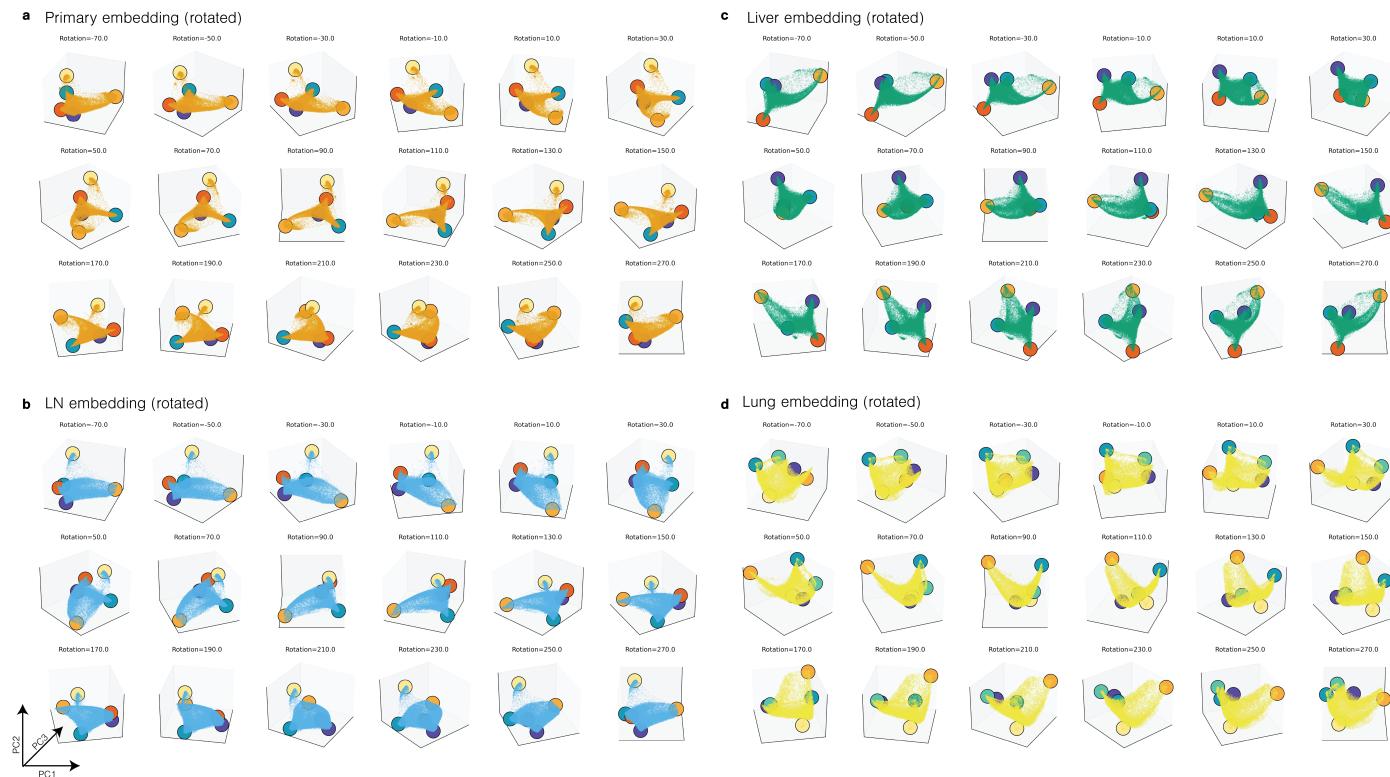
**Supplementary Figure 2. AAnet on CD8+ T cells.** (a) Hand-annotated archetypes from [10]. (b) AAnet-learned archetypes for each embedding. (c) Normalized expression for each archetype for key annotation genes. (d) Cosine similarity between archetypes for all measured genes. (e) Expression over archetypal coordinates for Week 17 embedding.

554 **Supplementary Figure 3**



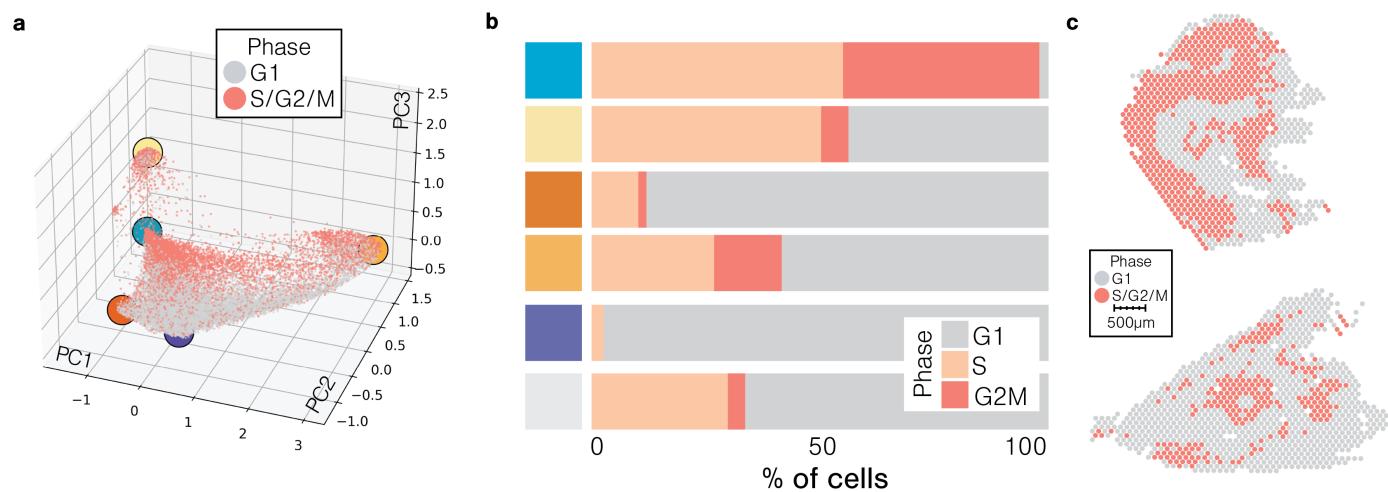
**Supplementary Figure 3. Primary tumor replicates independently characterized with AAnet.**  
Archetypes colored based on orthologous archetype in combined embedding (Figure 3).

555 **Supplementary Figure 4**



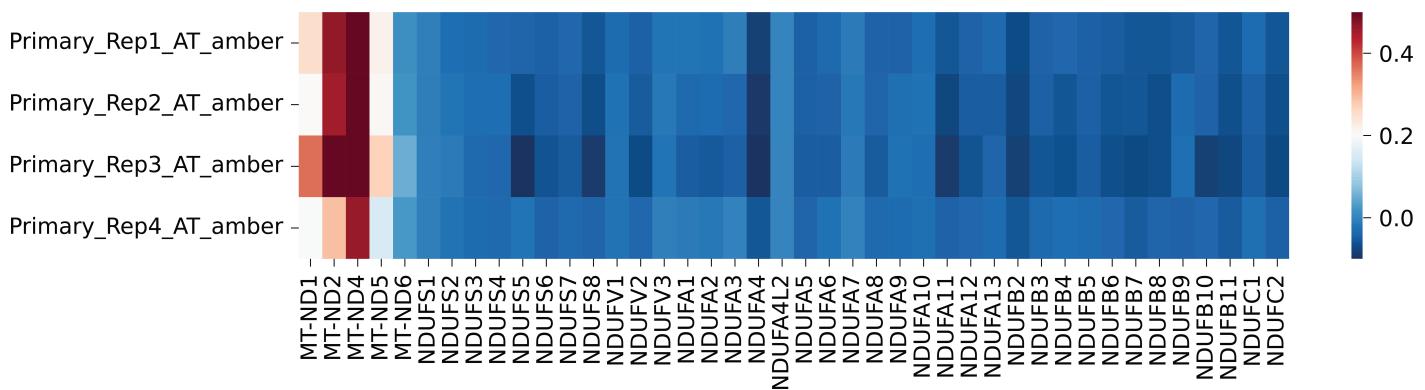
**Supplementary Figure 4. Cell embedding rotation.** PCA embedding rotated around PC3 for (a) Primary (b) LN (c) Liver (d) Lung.

556 Supplementary Figure 5



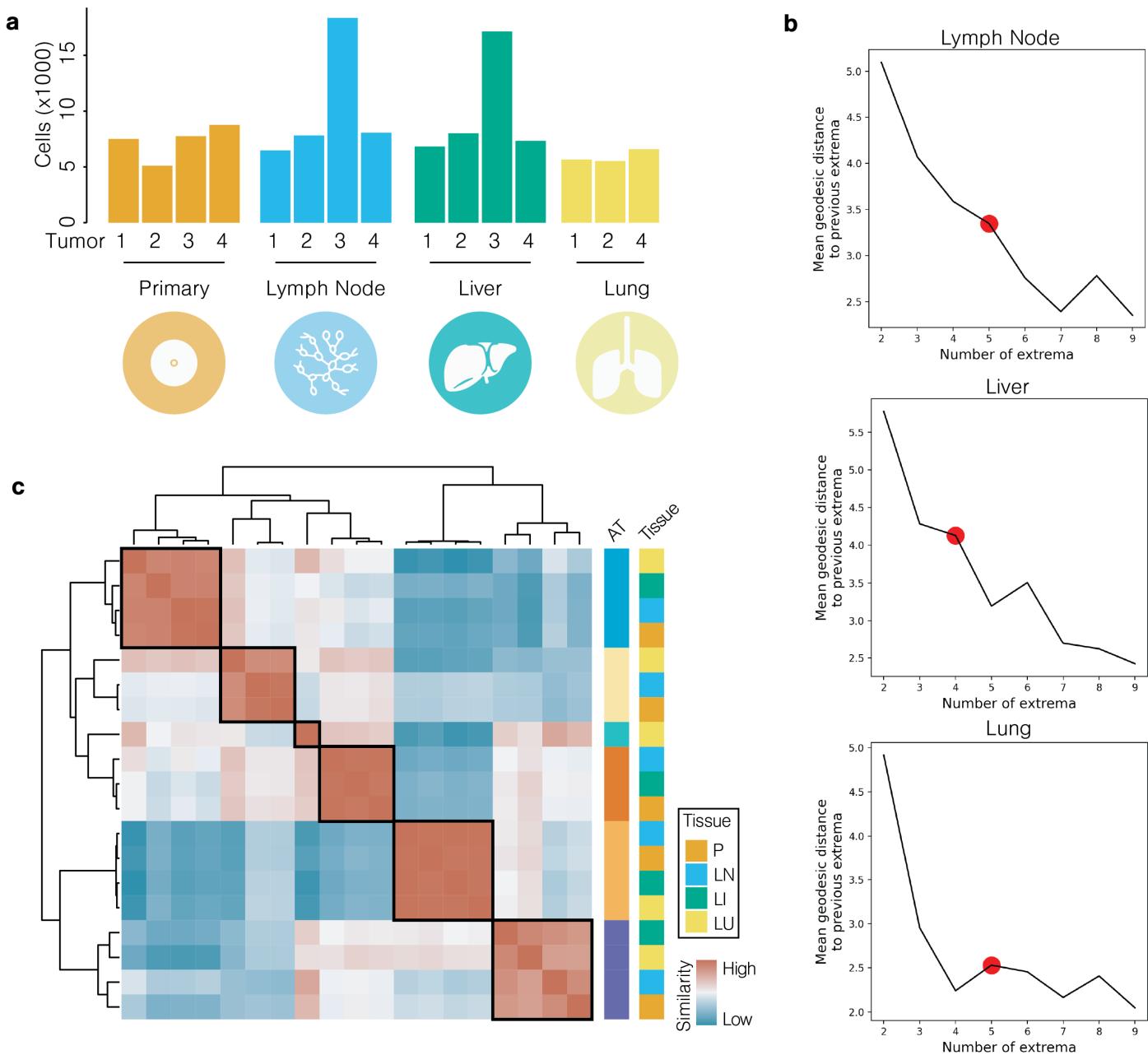
**Supplementary Figure 5. Cell cycle characterization.** (a) Cell cycle commitment based on competitive gene set enrichment from scRNASeq data. (b) Cell cycle commitment of cells closest to each archetype. (c) Commitment of each spatial voxel to each cell cycle phase.

557 Supplementary Figure 6



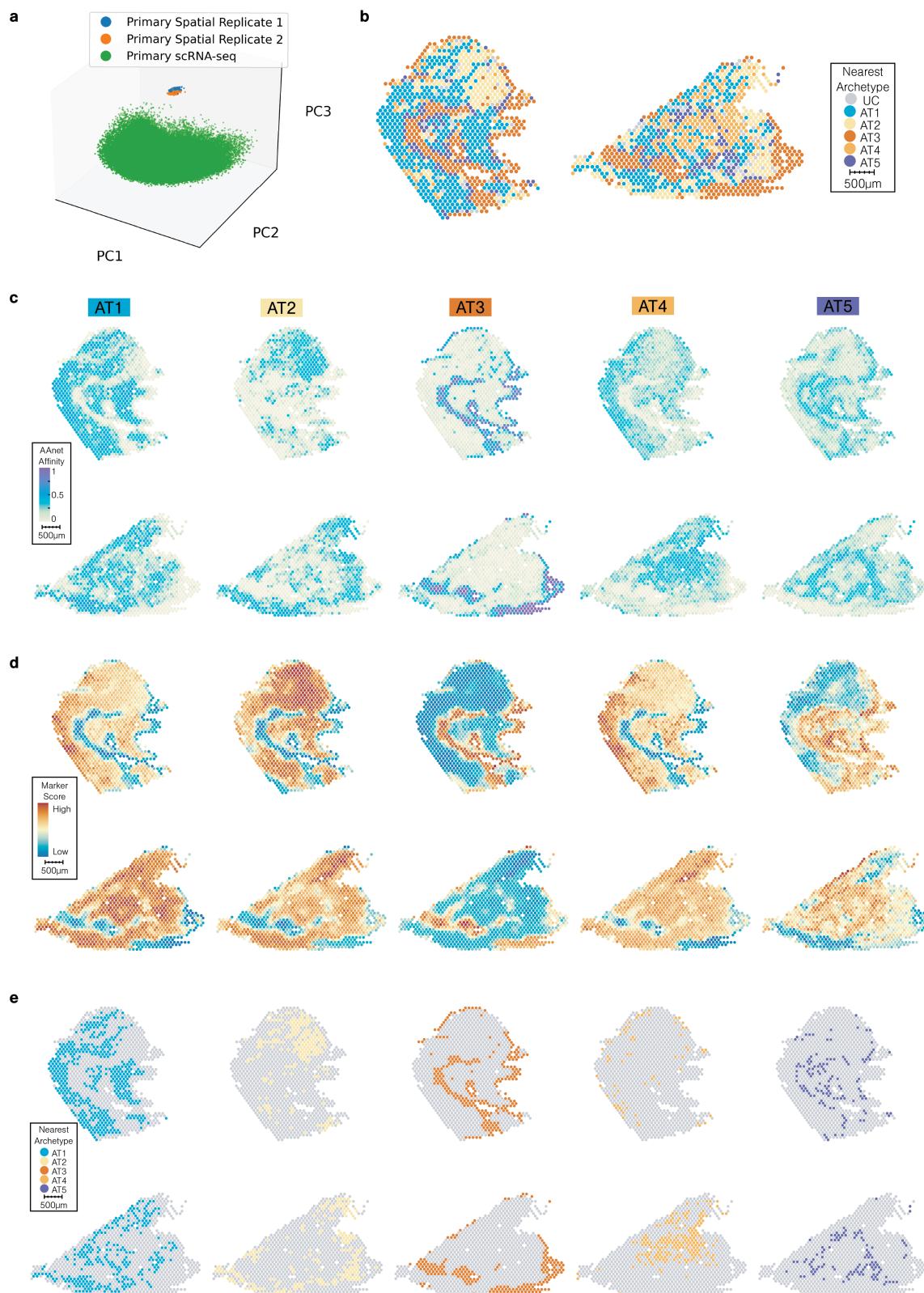
**Supplementary Figure 6. Mitochondrial expression.** Expression of mitochondrially and somatically-encoded electron transport chain genes in cell damage/death (amber) archetypes from primary replicates.

558 Supplementary Figure 7



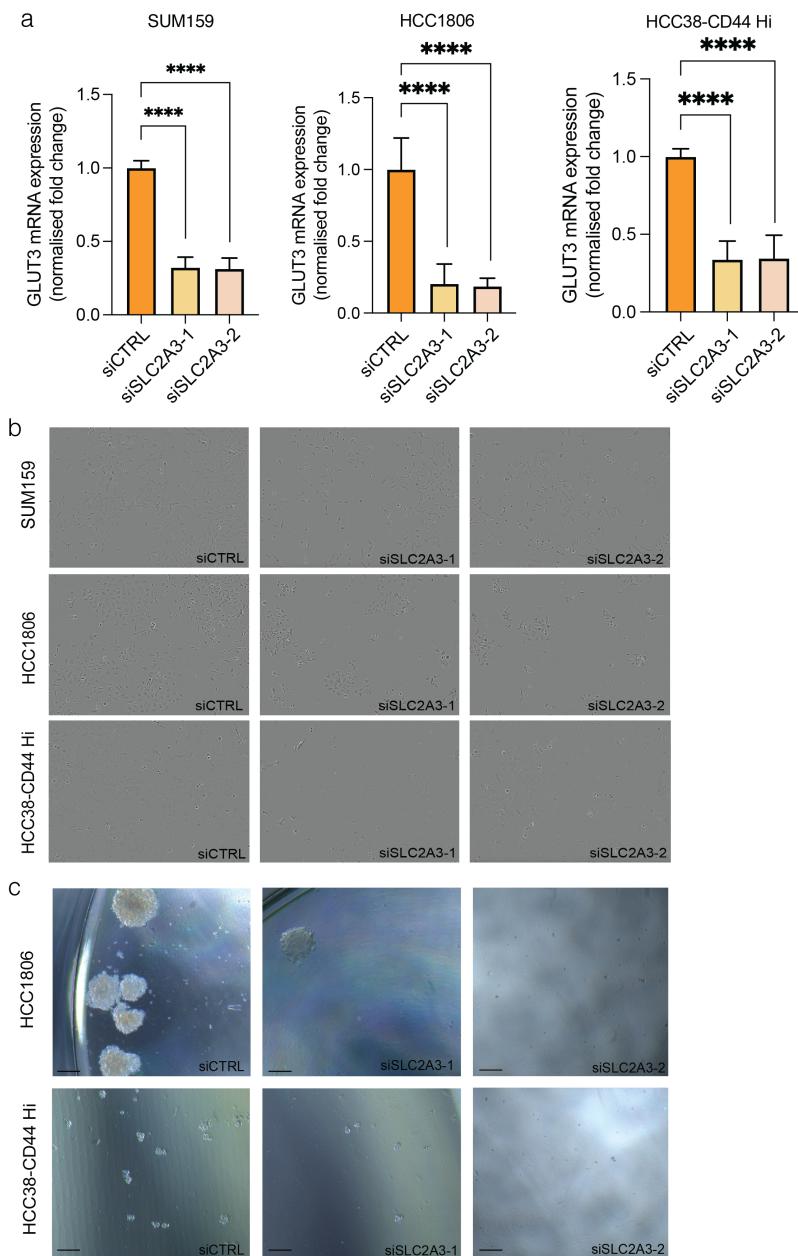
**Supplementary Figure 7. Comparison with metastatic tumor samples.** (a) Number of cells from each tumor for each tissue. (b) Number of archetypes for each tissue. (c) Cosine similarity characterizing relationships between archetypes.

559 **Supplementary Figure 8**



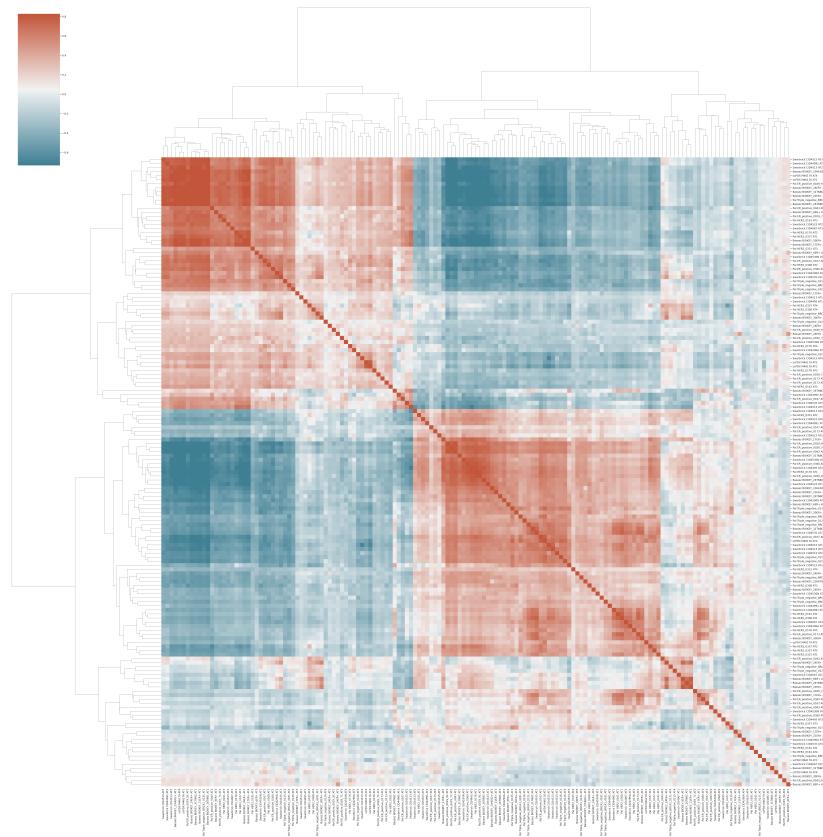
**Supplementary Figure 8. Spatial archetypal and gene set characterization.** (a) Embedding without scMMGAN alignment shows batch effect. (b) Overall commitment, where voxels that remained uncommitted colored gray. (c) Archetypal affinity for each archetype after scMMGAN alignment. (d) Core gene set enrichment for each archetype. (e) Commitment of each spatial voxel to each archetype.

560 Supplementary Figure 9



**Supplementary Figure 9. Metabolic heterogeneity and alignment to human breast cancers.** (a) Quantitative PCR analysis of SLC2A3 mRNA expression following control or SLC2A3 targeted siRNA treatment in SUM159, HCC1806 and HCC38-CD44 Hi cell lines.  $n=3$  independent experiments each with triplicate technical replicates analysed, significance measured using one-way ANOVA. (b) Representative images of proliferation analyses following control or SLC2A3 targeted siRNA treatment in SUM159, HCC1806 and HCC38-CD44 Hi cells. Images shown taken at experimental endpoint of 96 hours. (c) Representative images of tumorspheres derived from HCC1806 and HCC38-CD44 Hi cell lines following control or SLC2A3 targeted siRNA treatment. Images taken at experimental endpoint of 21 and 14 days for HCC1806 and HCC38-CD44 Hi respectively.

561 Supplementary Figure 10



Supplementary Figure 10. Cosine similarity across human breast cancer tumor archetypes.

## 562 Supplemental Table Legends

### 563 Supplementary Table 1.

564 Enriched marker genes and corresponding statistics from Wilcoxon rank sum test. Sheets correspond to genes  
565 associated with each archetype identified in each Primary tumor replicate.

### 566 Supplementary Table 2.

567 Enriched marker genes and corresponding statistics from Wilcoxon rank sum test. Sheets correspond to genes  
568 associated with each archetype identified in each tissue (Primary, LN, Liver, Lung).

## 569 Methods

### 570 Background on cell state heterogeneity analysis

#### 571 Clustering-based approaches unreliable identify clusters when data is a continuum of cells

572 Clustering is the most commonly used technique for characterizing cell state heterogeneity in single-cell data,  
573 and is considered a standard part of single-cell workflows and best practices [30]. However, clustering can be a  
574 nontrivial task, both with computational challenges and challenges with interpretation and annotation [23]. This  
575 is in part due to the fact that clustering assumes that data is composed of biologically distinct groups, such as  
576 discrete cell types.

577 After embedding the primary scRNA-seq data into 3-dimensions, it is evident that for the primary tumor,  
578 the cells are forming a connected manifold along the cellular state space, rather than separating into clusters  
579 (Figure 2a). After running Leiden clustering 100 times with default parameters, the cluster assignments changed  
580 at the boundaries, indicating that these cells are not strongly committed to one cluster.

581 We hypothesized that the cells at cluster boundaries are intermediate cells between the more distal extreme  
582 states. To test the ability of clustering-based analysis to characterize such datasets, we simulated a curved  
583 tetrahedron, where the datapoints are defined as a continuum between the vertices. We also defined a signal on  
584 the tetrahedron as a function of the affinity to one vertex (Supplementary Figure 2a). Clustering the simulated  
585 data with ten different clustering algorithms reveals (1) the lack of concordance across clustering methods when  
586 there is no latent cluster structure in the dataset and (2) the limitations of discretizing the cellular state space in  
587 characterizing continuous signals (Supplementary Figure 2b).

588

#### 589 Trajectory-based approaches enforce lineage structure that do not accurately capture simulated 590 signal

591 On the other hand, trajectory inference methods are commonly used to identify continuous paths in the datasets  
592 in order to define pseudotemporal ordering of cells, often for learning developmental decisions [20, 39, 42, 43].  
593 We show that, without clear lineage structure in the dataset, trajectory-based methods are not able to learn an  
594 intelligible ordering of cells or meaningfully characterize the defined signal (Supplementary Figure 2c).

595

### 596 Background on Archetypal Analysis

#### 597 Archetypal Analysis Overview

Archetypal analysis (AA) is an unsupervised learning method that aims to find extremal points, called *archetypes*, such that every point in a dataset can be approximated as a mixture of these archetypes [11]. Given a dataset  $X = \{x_1, \dots, x_N\} \subset \mathbb{R}^n$ , the archetypes  $\{z_1, \dots, z_k\} \subset \mathbb{R}^n$  are chosen so that for each data point  $x_i$  there exists  $\alpha_{i,1}, \dots, \alpha_{i,k} \in [0, 1]$  such that

$$\sum_{j=1}^k \alpha_{i,j} z_k \approx x_i \quad (1)$$

$$\sum_{j=1}^k \alpha_{i,j} = 1. \quad (2)$$

598 This type of linear combination where the coefficients are non-negative and sum to 1 is called a *convex combination*.  
599 The set of all such convex combinations of the archetypes  $\{z_1, \dots, z_k\}$  is a  $(k-1)$ -simplex. Note that the archetypes  
600 are not constrained to be points from the dataset.

## 601 Principal Convex Hull Analysis (PCHA)

602 One of the first AA algorithms was principal convex hull analysis (PCHA), proposed in [11]. PCHA constrains  
 603 the archetypes to be convex combinations of the input data points. It finds these archetypes through the following  
 604 optimization problem:

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{H}} \quad & \|\mathbf{X}^T - \mathbf{X}^T \mathbf{W} \mathbf{H}\|_F^2 \\ \text{s.t.} \quad & \mathbf{1}^T \mathbf{W} = \mathbf{1}, \mathbf{1}^T \mathbf{H} = \mathbf{1}, \mathbf{C} \geq 0, \mathbf{W} \geq 0. \end{aligned} \quad (3)$$

605 where  $\mathbf{W} \in \mathbb{R}^{N \times p}$  maps the data to the archetypes and  $\mathbf{H} \in \mathbb{R}^{p \times N}$  contains the coordinates of the archetypes in  
 606 the feature space. The constraints on  $\mathbf{W}$  guarantee that the archetypes are convex combinations of the data  
 607 points while the constraints on  $\mathbf{H}$  guarantee that the data points are convex combinations of the archetypes. [11]  
 608 also defines an optimization algorithm for the expression above using alternating convex least-squares. This  
 609 algorithm seeks solutions that satisfy the constraints  $\mathbf{1}^T \mathbf{W} = \mathbf{1}$  and  $\mathbf{1}^T \mathbf{H} = \mathbf{1}$  by adding auxiliary terms to the  
 610 objective function. [33] builds on this work by modifying the optimization algorithm to use projected gradient in  
 611 the alternating optimization steps to find solutions that satisfy the constraints.

## 612 Non-linear archetypal analysis variants (kPCHA, Javadi et al, Chen et al)

613 [33] also proposes kernel principal convex hull analysis (kPCHA), which is analogous to kernel principal  
 614 components analysis (kPCA). PCHA is rotation equivariant, meaning that applying a rotation to the input  
 615 dataset effectively results in the same rotation being applied the output features and archetypes. This implies that  
 616 the output of PCHA only depends on the kernel matrix  $\mathbf{X}^T \mathbf{X}$  rather than the actual input dataset  $\mathbf{X}$ . kPCHA  
 617 takes advantage of this fact by computing this kernel matrix in a possibly-infinite dimensional reproducing kernel  
 618 Hilbert space (RKHS), then running PCHA on this kernel matrix instead of  $\mathbf{X}^T \mathbf{X}$ . The mapping from the  
 619 input feature space into the RKHS is typically non-linear, allowing kPCHA to potentially take advantage of the  
 620 manifold geometry of the underlying dataset.

621 Several other works have also extended the algorithm proposed by [11]. In [22], the requirement that the  
 622 archetypes are convex combinations of input data is relaxed. The archetypes are found by optimizing an objective  
 623 function with two terms. The first term is similar to the objective function in 3 above and captures how well the  
 624 convex hull of the archetypes aligns with the dataset. The second term reflects how close the archetypes are to  
 625 the convex hull of the dataset. Meanwhile [9] proposes an algorithm to optimize 3 using an active-set approach.

## 626 Background on Machine Learning

### 627 Autoencoders

628 An autoencoder is a type of neural network that is used to learn compressed representations of data. Autoencoders  
 629 are comprised of two separate networks: an encoder and a decoder. The encoder network maps the input data  
 630 into a low-dimensional feature space or latent space, while the decoder tries to reconstruct the original data  
 631 from this low-dimensional representation. Through minimizing the error between the original data and the  
 632 reconstruction, termed *reconstruction loss*, autoencoders have been shown to successfully learn the structure of  
 633 data, and have had particular utility in capturing a meaningful representation of single-cell data [3, 14, 19, 29, 46].

### 634 Manifold learning

635 Manifold learning is a subfield of machine learning built around the manifold hypothesis, which asserts that  
 636 high-dimensional datasets are sampled from low-dimensional manifolds that lie in the high-dimensional space.  
 637 Here a *manifold* refers to a space that is locally isomorphic to a Euclidean space. Many methods in unsupervised  
 638 learning attempt to implicitly or explicitly capture the structure of the underlying data manifold. For instance,  
 639 the latent space of an autoencoder can be viewed as a parameterization of the data manifold.

## 640 AAnet Overview

### 641 AAnet architecture

642 AAnet is designed to have a flexible number and size of layers depending on the complexity for the task. For our  
 643 purposes, we found that a 2-layer (256 nodes, 128 nodes) encoder and (128 nodes, 256 nodes) 2-layer decoder  
 644 worked well. The batch size was 256, the optimizer was ADAM, the learning rate was set to 1e-3, and the weight  
 645 initialization was Xavier. All hidden layers contain Tanh nonlinear activations, besides layers directly before  
 646 and after archetypal layer which are linear so that each point is a linear combination of archetypes. The default  
 647 weight on *extrema loss*  $\gamma_{extrema}$  is set to 1. To encourage the archetypes to be tight, i.e. close to the data, we can  
 648 add Gaussian noise  $\sim N(0, 0.05)$  in the latent layer during training. For all datasets, we reduced dimensionality  
 649 using PCA before running AAnet and inverse-transformed the learned archetypes to the ambient space.

### 650 Reconstruction Loss

651 The main loss function for the autoencoder seeks to minimize the difference between the original input fed into  
 652 the encoder  $z = E(x)$  and the reconstructed input produced by the decoder  $\tilde{x} = D(z)$ , termed *reconstruction*  
 653 *loss*. Standardly, autoencoders use the mean squared difference of these two terms:

$$\text{Reconstruction MSE} = \mathbb{E}_{x \in X} [\|x - \tilde{x}\|^2] = \mathbb{E}_{x \in X} [\|x - D(E(x))\|^2]$$

### 654 Archetypal Loss

655 In addition to the reconstruction loss, we want to enforce the latent space of the autoencoder to learn the  
 656 structure of the data with respect to the archetypes. To this end, we convert the coordinates from Cartesian to  
 657 barycentric after the encoder learns the transformation. The barycentric coordinate system, related to Cartesian  
 658 coordinates, is a system in which each point is specified by reference to a simplex. When coordinates are  
 659 normalized to sum to 1, the vertices of the simplex are denoted by  $k+1$  one-hot vectors of length  $k+1$  for a  
 660  $k$ -simplex. For example, a triangle is a 2-simplex with 3 vertices, where the 3 vertices are (1,0,0), (0,1,0), and  
 661 (0,0,1). All coefficients of point  $P$  are positive if and only if  $P$  is inside the simplex.

662 As this coordinate system describes points with respect to a  $k$ -simplex, it is well-suited to be the latent space  
 663 for  $k$  archetypes.

664 To enable interpretation of points as convex combinations of archetypes, we enforce each point stays within  
 665 the simplex by adding an *archetypal loss* term, the mean squared error of the negative coefficients:

$$\text{Archetypal MSE} = \mathbb{E}_{x \in \text{neg. coefficients}} [\|x\|^2]$$

### 668 Extrema Loss

669 We developed a novel method to identify  $k$  plausible archetypes prior to model training. This method, explained  
 670 in detail below, builds a graph from the data and then uses the eigenvectors of the Laplacian matrix to find  
 671 extreme points in the datasets; these points will be referred to as *diffusion extrema*. We then include an *extrema*  
 672 *loss* term that penalizes large distances in the latent space between the diffusion extrema and the vertices of the  
 673 simplex. If the diffusion extrema and standard basis vectors in the latent space  $\mathbb{R}^k$  are labelled as  $\{\ell_i\}_{i=1}^k$  and  
 674  $\{e_i\}_{i=1}^k$ , respectively, then this loss term can be calculated as

$$\text{Diffusion extrema MSE} = \frac{1}{k} \sum_{i=1}^k \|E(\ell_i) - e_i\|_2^2$$

675 Let  $\{x_1, \dots, x_n\} \in \mathbb{R}^m$  be the points in the dataset. Then the procedure for finding these diffusion extrema  
 676 is as follows:

- 677 (1) Construct a graph  $G$  from the dataset. This can be done by computing a symmetrized  $k$ -nearest neighbors  
 678 graph from the dataset and then weighting the edges with a Gaussian kernel, as is done in [32].  
 679 (2) Let  $\psi_1, \dots, \psi_n$  denote the eigenvectors of the combinatorial Laplacian  $\mathbf{L}$  of  $G$  with corresponding eigenvalues  
 680  $\lambda_1 \leq \dots \leq \lambda_n$ . Compute

$$i_1 = \arg \max_{i \in \{1, \dots, n\}} |\psi_2(i)|,$$

681 where  $\psi_2(i)$  is the  $i$ th entry of  $\psi_2$ .

- 682 (3) Let  $\mathbf{L}^{(i_1)}$  denote  $\mathbf{L}$  with the entries in the  $i$ th row and  $i$ th column replaced by zeros. Likewise let  
 683  $\psi_1^{(i_1)}, \dots, \psi_n^{(i_1)}$  denote the eigenvectors of  $\mathbf{L}^{(i_1)}$ , again ordered in an ascending fashion by corresponding  
 684 eigenvalue. Compute

$$i_2 = \arg \max_{i \in \{1, \dots, n\}} |\psi_2^{(i_1)}(i)|.$$

- 685 (4) Let  $\mathbf{L}^{(i_1, \dots, i_m)}$  denote  $\mathbf{L}$  with the entries in the  $i_1$ th,  $\dots$ ,  $i_{m-1}$ st, and  $i_m$ th rows and columns replaced  
 686 by zeros. Likewise let  $\psi_1^{(i_1, \dots, i_m)}, \dots, \psi_n^{(i_1, \dots, i_m)}$  denote the eigenvectors of  $\mathbf{L}^{(i_1, \dots, i_m)}$ . Iteratively for each  
 687  $j = 3, \dots, k$  compute

$$i_j = \arg \max_{i \in \{1, \dots, n\}} |\psi_j^{(i_1, \dots, i_{j-1})}(i)|.$$

- 688 (5) The diffusion extrema are  $x_{i_1}, \dots, x_{i_k}$ .

689 The intuition behind this algorithm comes from an application of Courant-Fischer theorem for symmetric  
 690 matrices. Given an  $n \times n$  symmetric matrix  $\mathbf{A}$  with eigenvectors  $\mathbf{a}_1, \dots, \mathbf{a}_k$ , Courant-Fischer tells us that

$$\mathbf{a}_1 = \arg \min_{\|\mathbf{x}\|=1} \mathbf{x}^T \mathbf{A} \mathbf{x} \quad \text{and} \quad \mathbf{a}_2 = \arg \min_{\substack{\|\mathbf{x}\|=1 \\ \langle \mathbf{x}, \mathbf{a}_1 \rangle = 0}} \mathbf{x}^T \mathbf{A} \mathbf{x}$$

691 If  $\mathbf{A}$  is the Laplacian matrix  $\mathbf{L}$  for some weighted graph  $G = (V, E, w)$  then

$$\mathbf{x}^T \mathbf{L} \mathbf{x} = \sum_{(i,j) \in E} w_{i,j} (\mathbf{x}(i) - \mathbf{x}(j))^2$$

692 Intuitively the quadratic form  $\mathbf{x}^T \mathbf{L} \mathbf{x}$  captures how smoothly  $\mathbf{x}$  varies over the edges of  $G$ . Hence  $\psi_1$  is the  
 693 (normalized) constant vector, while  $\psi_1$  can be viewed as a smooth signal on  $G$  that is orthogonal to the constant  
 694 vector. Now if we consider the matrix  $\mathbf{L}^{(i)}$  we see that minimizing the quadratic form  $\mathbf{x}^T \mathbf{L}^{(i)} \mathbf{x}$  can be recast as  
 695 minimizing  $\mathbf{x}^T \mathbf{L} \mathbf{x}$  with the additional constraint that  $\mathbf{x}(i) = 0$ , i.e.

$$\psi_2^{(i)} = \min_{\substack{\|\mathbf{x}\|=1 \\ \langle \mathbf{x}, \mathbf{1} \rangle = 0}} \mathbf{x}^T \mathbf{L}^{(i)} \mathbf{x} = \min_{\substack{\|\mathbf{x}\|=1 \\ \mathbf{x}(i)=0}} \mathbf{x}^T \mathbf{L}^{(i)} \mathbf{x}$$

696 This is because  $\psi_1^{(i)} = e_i$ , the  $i$ th standard basis vector. Extending this reasoning to  $\mathbf{L}^{(i_1, \dots, i_m)}$ , where  $i_1, \dots, i_m \in$   
 697  $\{1, \dots, |V|\}$  are unique but arbitrary, we see that  $\psi_j^{(i_1, \dots, i_m)} = e_{i_j}$  for  $1 \leq j \leq m$ . Hence

$$\psi_{m+1}^{(i_1, \dots, i_m)} = \min_{\substack{\|\mathbf{x}\|=1 \\ \mathbf{x}(i_1) = \dots = \mathbf{x}(i_m) = 0}} \mathbf{x}^T \mathbf{L}^{(i_1, \dots, i_m)} \mathbf{x}.$$

698 This tells us that the  $m$ th eigenvector of  $\mathbf{L}^{(i_1, \dots, i_m)}$  is a signal on  $G$  that is as smooth as possible while also  
 699 having the value 0 on vertices  $i_1, \dots, i_m$ . Then we should expect  $\mathbf{L}^{(i_1, \dots, i_m)}$  to attain its largest absolute value  
 700 at a vertex that is far from the vertices  $i_1, \dots, i_m$ . This is the guiding principle behind the above diffusion  
 701 extrema method. Step (2) takes advantage of the fact that  $\psi_2$  is smooth on  $G$  and therefore is likely to attain  
 702 its maximum absolutely value at extremal points in the graph. The vertex at which  $\psi_2$  has its largest absolute

703 value is chosen to be the first diffusion extrema. Steps (3) and (4) then use the properties of  $\psi_{m+1}^{(i_1, \dots, i_m)}$  discussed  
704 above to iteratively find vertices in the graph that are far away from the extrema that have already been chosen.  
705

706 It has not been conclusively proven that the eigenvectors of the Laplacian  $\mathbf{L}$  and Laplacian variant  $\mathbf{L}^{(i_1, \dots, i_m)}$   
707 attain their maximum/minimal values at extremal points in the graph, nor is it entirely clear which vertices in a  
708 graph can be called in extremal. However through our experiments on both toy data and real biological data we  
709 believe this diffusion extrema method to be robust and effective in identify reasonable archetypes in a dataset.

## 709 Choosing number of archetypes

710 To choose the number of archetypes, we first define a range of possible archetype counts  $m = 2, 3, \dots, n_{at}$ . Then,  
711 we calculate  $n_{at}$  diffusion extrema. For the  $m^{th}$  diffusion extremum, we calculate the geodesic distance of that  
712 extremum to all previous extremum (first to  $m - 1^{th}$ ) and take the mean of all the distances. Intuitively, this  
713 score tells how transcriptionally distinct archetype  $m$  is from all existing archetypes, and if it is very similar,  
714 it is likely not a new archetype. Thus, we take the knee point of this score as the number of archetypes  $k$  for  
715 downstream analysis.

$$s_m = \frac{1}{m} \sum_{i=1}^{m-1} \text{geodesic distance}(\ell_m, \ell_i) \text{ for } m = 2, 3, \dots, n_{at}$$

## 716 Comparative Analysis

### 717 Simulated Datasets

718 To generate a simulated dataset for comparisons, we sample non-uniformly from a k-simplex projected onto a  
719 k-dimensional sphere using a stereographic projection. This enables comparisons on a dataset with nonlinear  
720 geometry and known vertices.

721 For comparisons in Supplementary Figure 1, we sample 10,000 points from a 3-simplex (a tetrahedron)  
722 projected onto a hypersphere. This results in a curved tetrahedron, where the ground truth archetypes do not  
723 correspond to the extrema in the data space (as they would for a classic tetrahedron). We embed the curved  
724 tetrahedron with PHATE [32], a dimensionality reduction tool designed to capture nonlinear local and global  
725 variation.

726 Next, we simulated a signal, defined by the *sine* of the affinity between vertex 4 and all other points in the  
727 tetrahedron based on data geometry. The ground truth signal is thus a sinusoidal function. This signal is defined  
728 by its relation to a ground truth archetype, and is designed to model biological processes that are enriched with  
729 respect to a cellular archetype.

### 730 Comparison to clustering methods

731 To compare archetypal analysis to clustering on the curved tetrahedron, we clustered the data using default  
732 parameters for 10 different clustering algorithms in `sklearn.cluster`. Five methods required the number of  
733 clusters to be specified, and for these we specified four clusters, the ground truth number of vertices for a  
734 tetrahedron.

735 We ran clustering using ten different clustering algorithms: five that require the number of clusters to be  
736 specified (KMeans, Agglomerative Clustering (Ward), Agglomerative Clustering (Single Linkage), Agglomerative  
737 Clustering (Complete Linkage), and Spectral Clustering), and five that infer the number of clusters from the  
738 data (Affinity Propagation, Mean Shift, DBSCAN, OPTICS, and BIRCH).

739

### 740 Comparison to trajectory-inference methods

741 For trajectory inference comparisons, we ran Slingshot and diffusion pseudotime (DPT). Slingshot requires  
742 cluster labels and a dimensionality-reduced representation. We used the PHATE embedding and the KMeans  
743 cluster labels as input for Figure 2, as KMeans is a popular method for clustering that produced stable results.

744 Besides these two inputs, all other parameters were default. DPT was also run with default parameters, with the  
745 required starting cell chosen randomly from cluster 3 of KMeans (the cluster containing extremum 4).

#### 746 Comparison to archetypal analysis methods

747 All archetypal analysis methods require input of the number of archetypes, so we specified four archetypes for  
748 each method, otherwise running with default parameters. To generate MSE calculations between the ground  
749 truth vertices and inferred archetypes, we ran each method 5 times and visualized the first run for each method.

750 By the definition of barycentric coordinates, a point has a value closer to 1 for dimension  $k$  if it has a high  
751 affinity to the  $k$ -th vertex, and a value closer to 0 if it has a low affinity to the  $k$ -th vertex. Therefore, we color  
752 the dataset with the latent coordinates for each dimension to determine if the latent space has semantic structure,  
753 and if AAnet is effectively learning affinity relative to each extremum.

#### 754 AAnet for published antigen-specific CD8+ T cells

755 For each dataset, we ran TruncatedSVD to reduce the dimensionality to 100, and then ran AAnet with default  
756 parameters. The expression levels in Supplementary Figure 2 were z-score normalized and archetypes were  
757 clustered based on cosine similarity.

### 758 Computational Methods for TNBC

#### 759 Single-cell RNA preprocessing

760 The CellRanger Analysis Pipeline (v3.0.2) was used to align the sequencing reads (fastq) to a pre-built reference  
761 genome (10x Genomics) containing both the human and mouse genomes (GRCh38 + mm10) and gene expression  
762 quantified in each cell. Cells from four primary tumors were sequenced (T1-T4) with a total of 33,938 cells  
763 sequenced (T1 = 8707, T2 = 5951, T3 = 8934, T4 = 10346). Gene expression data was loaded into Python  
764 (v3.8) and quality-control (QC) statistics were computed using the `scanpy.pp.calculate_qc_metrics` function  
765 (v1.9.1, [49]). To ensure only human tumor cells were taken for downstream archetypal analysis, cells with  
766 less than 99% of data aligning to the human genome were removed. Mouse genes were also removed prior to  
767 calculating library size. Cells with a total library size between 2000-50,000 UMI and expressing at least 1000  
768 genes were retained for downstream analysis. A total of 28,478 primary tumor cells passed QC filtering (T1 =  
769 7606, T2 = 5118, T3 = 8163, T4 = 7591). Remaining tumor cells were normalized to 10,000 reads per cell and  
770 square-root transformed using the scprep package (v1.2.3, [github.com/krishnaswamylab/scprep](https://github.com/krishnaswamylab/scprep)). To correct  
771 for dropout data were smoothed with the manifold smoothing method MAGIC (v3.0.0, [45]). Cell numbers  
772 remaining after QC were visualized using the R package ggplot2 (v3.4.2, [48]). Cell-cycle phase assignment  
773 was performed using the `scanpy.tl.score_genes_cell_cycle` function with previously defined S-phase and  
774 G2M-phase gene lists [41].

775 Highly variable genes were detected within each sample using the scprep function  
776 `scprep.select.highly_variable_genes` and a cellular graph constructed based on KNN and alpha decay  
777 kernel using graphtools  
778 (v1.5.3, [github.com/krishnaswamylab/graphtools](https://github.com/krishnaswamylab/graphtools)). For the combined analysis of all tumors we used an MNN  
779 kernel to build a cellular graph with batch correction between the replicates. We then used MAGIC (v3.0.0, [45])  
780 to transform each graph into the gene space, and ran TruncatedSVD to reduce the dimensionality to 100 and  
781 visualize samples in reduced dimensions.

#### 782 AT similarity, affinity and commitment

783 AAnet was run on both individual samples and a combined dataset using default parameters. Data archetypes  
784 were defined for each dataset using AAnet and archetypal expression vectors were generating by transforming  
785 archetype coordinates back into the gene space. Cosine similarity was calculated between expression vectors  
786 and hierarchical clustering used to identify five orthologous expression states between datasets. The affinity of  
787 each cell to each archetype was calculated based on the distance of each cell to each archetype in the AAnet  
788 latent embedding. As the combined affinity of each cell to all archetypes is regularized to 1, cells with an affinity

789 to single archetype greater than the sum of their affinity to all other archetypes (i.e. with an affinity >0.5 for  
790 any archetypes) were defined as committed to that archetype. The archetypal composition of each sample was  
791 determined by summing the number of committed cells per archetype and assigning cells with affinity for all  
792 archetypes <0.5 as uncommitted.

## 793 Biological characterization of archetypal states

794 Marker genes upregulated in each archetype were calculated by comparing gene expression between cells that  
795 were the most committed to each archetype. The top 125 cells with the strongest affinity for each archetype  
796 in the combined AAnet latent space were selected from each replicate (500 cells per archetype), as well as the  
797 125 cell per replicate furthest from any archetype as the uncommitted populations (500 cells uncommitted).  
798 Genes that were upregulated in any group relative to all other groups (FDR < 0.05) were determined using the  
799 FindMarkers function from the Seurat R package (v4.0, [40]). Cancer hallmark gene sets [28] overrepresented  
800 (FDR < 0.1) in markers associated with each archetype were determined by 1-sided Fisher's exact test using the  
801 clusterProfiler R-package (v3.16.1, [52]) and visualized using ggplot2.

## 802 Spatial RNA data generation

803 Spatial transcriptomics was conducted using 10X Visium Spatial Gene Expression Slide and Reagent Kit,  
804 16 rxns (PN-1000184), according to the protocol detailed in document CG000239RevD for the TNBCs and  
805 CG000239RevE for the Xenografts, available in 10X Genomics demonstrated protocols. Cryo-sectioning was  
806 done on OCT embedded and snap frozen tissue samples at 10um thickness and placed on cold Visium slide  
807 arrays. The sections were adhered by swiftly warming of the backside of the slide. The slides were then kept in  
808 -80°C less than 4 weeks before processing accordingly.

809 In short, the slides were first warmed at 37°C for one minute and then immersed in pre-chilled methanol  
810 (VWR EU, 20847.307) for 30 minutes at -20°C for fixation. Staining with hematoxylin and eosin was carried  
811 out by one min of drying of the tissues with 500uL of isopropanol (Fisher Scientific, A461-1) followed by air  
812 drying until sections turned white. Around 1 mL of with Mayer's hematoxylin (Agilent, S23309) was pipetted  
813 onto the slides and treated for four minutes. The slides were then washed in nuclease free water followed by a  
814 two-minute incubation with bluing buffer (Agilent, CS702), washed again and then counterstained with buffered  
815 eosin (Sigma-Aldrich, HT110216, 1:10 dilution in Tris-Acetic Acid Buffer). The slides were air dried for about 2  
816 minutes and then warmed for 5 minutes at 37°C before mounting using 85percent Glycerol (Merck, 104094) and  
817 a coverslip.

818 Bright field histology images were obtained using a 20X objective on Zeiss microscope using the Metafer  
819 VSlide system and the images were processed by the VSlide software. The images were extracted as jpgs for  
820 downstream analysis.

821 After imaging, the coverslip and the remaining glycerol was washed off in Milli-Q water and the slides were  
822 attached in plastic cassettes included in the reagent kit and first subjected to 20 minutes of permeabilization at  
823 37°C to let the mRNA reach the probes on the slide surface for binding.

824 The protocol was then followed without deviations to create amplified libraries which in the end were  
825 individually indexed using the Dual Index Kit TT SetA, (PN-1000215, 10X Genomics), quality controlled on a  
826 BioAnalyzer instrument and concentrations were measured using Qubit DNA HS. The libraries were pooled  
827 equimolarly (2nM) and sequenced on the Nextseq 500 (Illumina platform) for the tnbc and Nextseq 2000  
828 (Illumina platform) for the xenografts. To reach the appropriate read depth the recommended number of reads  
829 per ST spot were applied according to the protocol.

## 830 Spatial RNA preprocessing

831 The SpaceRanger Analysis Pipeline (v2.0.0, 10x Genomics) was used to align the sequencing reads (fastq) to a  
832 pre-built reference genome (10x Genomics) containing both the human and mouse genomes (GRCh38 + mm10)  
833 and gene expression quantified in each cell. Quality-control (QC) statistics were computed using the STUtility  
834 package [7]. Voxels with a library size < 3000 UMI were removed and remaining voxels manually curated to  
835 remove those that were disconnected from the main tissue section. Human and mouse genes were separated

836 to create two datasets per sample, one measuring the expression of human genes at each voxel, pertaining to  
837 tumor cells, and the other measuring expression of mouse genes, pertaining to cells from the microenvironment.  
838 Genes were removed if they were detected in less than 10 voxels. Finally each dataset was filtered to remove  
839 voxels with low library diversity for a given genome (<1000 human genes detected, <300 mouse genes detected).  
840 After QC filtering, the human dataset comprised of 14,819 genes measured at 1170 and 1105 voxels for each  
841 tumor, while the mouse dataset comprised of 12,119 genes measured at 1150 and 1079 voxels in TX and TX  
842 section respectively. Remaining datasets were each normalized to 10,000 reads per voxel and log transformed.  
843 To correct for dropout (human dataset = 69.8% zeros, mouse dataset = 90.6% zeros) data were smoothed with  
844 the manifold smoothing method MAGIC. Cell-cycle phase assignment was performed using the CellCycleScoring  
845 function in Seurat [40] with previously defined S-phase and G2M-phase gene lists [41]. Spatial plots throughout  
846 the manuscript were generated using STUtility.

#### 847 Mapping of archetypes to spatial transcriptomics with scMMGAN

848 To map archetypes identified in the single-cell data from primary tumors to voxels in the spatial transcriptomic  
849 samples we used scMMGAN [4]. Smoothed expression data were zero-centered and unit scaled each dimension,  
850 and reduce to 50 principal components used as input to the scMMGAN generator, with the combined scRNAseq  
851 data described above used as input for the discriminator layer. scMMGAN was run with a generator consisting  
852 of three internal layers of 128, 256, and 512 neurons with batch norm and leaky rectified linear unit activations  
853 after each layer, and a discriminator consisting of three internal layers with 1,024, 512, and 256 neurons with the  
854 same batch norm and activations except with minibatching after the first layer. We use the geometry-preserving  
855 correspondence loss with a coefficient of 10, cycle-loss coefficient of 1, learning rate of 0.0001, and batch size of  
856 256. This network was used to generate a single-cell-like representation of each spatial voxel. These generated  
857 single-cell values were then embedded into the AAnet latent space trained using the combined single-cell RNAseq  
858 dataset. The affinity of each voxel to each archetype was then determined based on the distance of its generated  
859 single-cell values to each archetype in the trained AAnet latent space. Each voxel was then assigned to an  
860 archetype to which it had the highest affinity, or uncommitted in the case the maximum affinity corresponded to  
861 more than one archetype. It is important to note that the resolution of spatial transcriptomics voxels above the  
862 single-cell level (estimated between 3-10 cells per voxel), thus archetypes represent the dominant expression state  
863 among cells in the voxel. We also scored voxels based on the expression of the marker gene sets associated with  
864 each archetype in the single-cell data using on the first principal component of gene set expression, analogous  
865 to the eigengene metric used to summarize gene coexpression networks [25]. Scores were calculated based on  
866 expression of the top 50 marker genes with the highest log fold enrichment per archetype, excluding mitochondrial  
867 and ribosomal genes.

#### 868 Microenvironment mapping and enrichment

869 We used data aligning to the mouse genome at each voxel to analyze the microenvironment associated with each  
870 archetype. Based on the archetypal assignment of voxels using tumor data, as described above, we identified  
871 differentially expressed genes between microenvironment spatially colocalized with each archetype using the  
872 FindMarkers function in Seurat (LFC > 0.1, FDR < 0.05) [40]. Enrichment of Gene Ontology Biological Processes  
873 and cell-types markers associated with "Connective tissue", "Immune system", "Smooth muscle", "Epithelium",  
874 "Vasculature", "Blood", "Mammary gland" or "Skeletal muscle" in the Pangloa database (v27/03/2020) [18]  
875 among differentially expressed genes was determined with a 1-sided Fisher's exact test using the clusterProfiler  
876 R-package [54]. Putative ligand-receptor interactions between archetypes were identified using CellPhoneDB  
877 (v2) [13]. Human orthologs to mouse genes were identified in biomaRt and counts matrices were merged based  
878 on gene id prior to running CellPhoneDB using parameters –iterations 1000 –threshold 0.2. and identifying  
879 significant interactions (FDR < 0.05). The metabolism between archetypes and microenvironment was compared  
880 based expression of key enzymes involved in glycolysis and the tricarboxylic acid (TCA) cycle in each archetype  
881 and associated microenvironment. Glycolytic enzymes were designated as either hypoxic based on their enriched  
882 expression (LFC > 0) in voxels assigned to the hypoxic AT5 archetype. Heatmaps were generated based on the  
883 mean of scaled values for voxels associated with each archetype (human genes) or microenvironment (orthologous  
884 mouse genes).

## 885 Cell Culture

886 SUM159 cells were cultured in Ham's F12 + 1% (v/v) hydrocortisone, HCC1806 and HCC38-CD44 Hi cells  
887 were cultured in RPMI media. Media were supplemented with 10% (v/v) fetal bovine serum (GE healthcare)  
888 and 1% (v/v) penicillin-streptomycin. To isolate CD44Hi cells from the parental HCC38 cell line, cells were  
889 expanded *in vitro* in 150 mm tissue culture-treated culture dish (Corning). For the initial FACS rounds, cells  
890 were trypsinized and 1-2 x 10<sup>7</sup> cells were stained for the membrane marker CD44 (BD anti-human CD44-PE-cy7  
891 (1:800)) for 25 min at 4C. Antibody titration was previously determined staining a battery of breast cancer cell  
892 lines. Pure CD44<sup>Hi</sup> cells were collected and replated for expansion in culture. Following sorting purification,  
893 cultures were supplemented with 0.1% (v/v) gentamicin and 1% (v/v) antibiotic-antimycotic for 2 passages to  
894 avoid contamination. Sequential rounds of FACS enrichment were performed until 100% pure populations were  
895 isolated. Data collection was performed using a BD Aria III and FACSDiva software (BD Biosciences). Flowjo  
896 X10.7.1 was used for data analysis.

## 897 siRNA treatment

898 Cells were plated, allowed to grow for 24h then transfected with either 20nM siRNA *Silencer*<sup>TM</sup> Select negative  
899 control (Invitrogen, 4390843) or SLC2A3 directed siRNA (Invitrogen, 4390824). siRNAs were transfected using  
900 Optimem Reduced Serum Media (ThermoFisher 31985070) and Lipofectamine 3000 (ThermoFisher scientific  
901 L3000015#) as per standard protocol. Cells were harvested at 24h for qPCR analysis or trypsinized and re-seeded  
902 for proliferation and/or tumorsphere assay.

## 903 qRT-PCR

904 Total RNA was extracted and purified using the RNeasy micro kit (Qiagen). Reverse transcription was performed  
905 from 1 µg of total RNA with Superscript IV reverse transcriptase (Life Technologies, CA, USA) according to the  
906 manufacturer's instructions. The reverse transcription product was diluted 1:10 with TaqMan Fast Advanced  
907 Master Mix (Invitrogen, 4444556) and used as a cDNA template for qPCR analysis. Real-time quantitative PCR  
908 was performed using the QuantStudio<sup>TM</sup> 7 Flex Real-Time PCR System (Applied Biosystems, CA, USA). Results  
909 are represented as mean values normalised to controls.

## 910 Proliferation

911 Following 24 hours of treatment of control or SLC2A3 targeted siRNA treatment, cells were seeded (1000 per  
912 well) in a 96-well plate. Proliferation was determined by confluence (%) per well, measured every 24 hours over a  
913 96-hour using an Incucyte (Sartorius).

## 914 Tumorsphere-forming assay

915 After 24 hours of treatment of control or SLC2A3 targeted siRNA treatment, cells were trypsinized and seeded  
916 in an ultra-low attachment 96 well plate (Corning) (300 cells per well) at 1000 cells per well. Tumorsphere  
917 media was comprised of methylcellulose (Sigma, m-7027) and basal media supplemented with 20ng/ml basic  
918 fibroblast growth factor (Millipore, GF003), 20ng/ml human epidermal growth factor (Sigma, E1264), B27 (Life  
919 Technologies, 17504-044) and 4µg/ml heparin (Sigma, H3149). 50ul additional tumorsphere media was added  
920 every 5 days and tumorspheres were counted and imaged at 14 or 21 days.

## 921 Human scRNA-seq data preprocessing

922 For experiments from [6,35,51], the datasets were first subset to cancer epithelial cells based on prior annotation,  
923 and only samples with  $\geq 1000$  cells were analyzed with AAnet. We then further preprocessed the data by  
924 removing genes expressed in fewer than 5 cells, filtering library size to between 1500 and 60000 UMI counts,  
925 and L1 normalizing for library size. We then square-root transformed the data, and filtered any remaining  
926 contaminating cells based on marker gene expression. Finally, we embedded the data with MAGIC and PCA  
927 before identifying archetypes with AAnet. For experiments based on the PDX model, the same pipeline was

928 followed, where additionally cells with less than 99% of reads aligning to human cells were removed, and mouse  
929 genes were removed.

## 930 Data Availability

931 The accession codes to the newly generated single-cell and spatial data will be provided before publication.  
932 Public CD8+ T cell data can be accessed at GEO under accession number GSE182509. Processed scRNA-seq  
933 data from [35] is available as GEO series GSE161529. Raw sequencing reads of all single-cell experiments from [6]  
934 have been deposited in the European Genome-phenome Archive (EGA) under study no. EGAS00001004809.  
935 Processed scRNA-seq data from [51] are also available through the Gene Expression Omnibus under accession  
936 number GSE176078.

## 937 Code Availability

938 The source code is available at <https://github.com/KrishnaswamyLab/AAnet>.

## 939 Author Contributions

940 A.V., D.B.B., A.B., D.V.D., and S.K. developed AAnet. C.L.C. and B.P.S.J. designed the biological experiment.  
941 B.P.S.J. performed biological experiments. M.A. aligned scRNA-seq and spatial transcriptomic data with  
942 scMMGAN. S.E.Y., A.V., and C.P. performed computational analyses. A.M. and J.K. performed and supported  
943 spatial transcriptomic experiments. J.H. performed analysis and interpretation of metabolomic data. L.D.G. and  
944 S.K. oversaw computational analyses done by S.E.Y., A.V., and C.P. A.V., S.E.Y., S.K. and C.L.C. wrote the  
945 manuscript.

## 946 Competing Interests

947 Dr. Smita Krishnaswamy is on the scientific advisory board of KovaDx and AI Therapeutics.  
948 Dr. Christine Chaffer is the Founder and Managing Director of Kembi Therapeutics.

## 949 References

- 950 [1] Ali F Abdel-Wahab, Waheed Mahmoud, and Randa M Al-Harizy. Targeting glucose metabolism to suppress  
951 cancer progression: prospective of anti-glycolytic cancer therapy. *Pharmacol. Res.*, 150(104511):104511,  
952 December 2019.
- 953 [2] Miri Adler, Yael Korem Kohanim, Avichai Tendler, Avi Mayo, and Uri Alon. Continuum of gene-expression  
954 profiles provides spatial division of labor within a differentiated cell type. *Cell Syst.*, 8(1):43–52.e5, January  
955 2019.
- 956 [3] Matthew Amodio, David van Dijk, Krishnan Srinivasan, William S Chen, Hussein Mohsen, Kevin R Moon,  
957 Allison Campbell, Yujiao Zhao, Xiaomei Wang, Manjunatha Venkataswamy, Anita Desai, V Ravi, Priti  
958 Kumar, Ruth Montgomery, Guy Wolf, and Smita Krishnaswamy. Exploring single-cell data with deep  
959 multitasking neural networks. *Nat. Methods*, 16(11):1139–1145, November 2019.
- 960 [4] Matthew Amodio, Scott E Youlten, Aarthi Venkat, Beatriz P San Juan, Christine L Chaffer, and Smita  
961 Krishnaswamy. Single-cell multi-modal GAN reveals spatial patterns in single-cell data from triple-negative  
962 breast cancer. *Patterns (N Y)*, 3(9):100577, September 2022.
- 963 [5] Thomas M. Ashton, W. Gillies McKenna, Leoni A. Kunz-Schughart, and Geoff S. Higgins. Oxidative  
964 phosphorylation as an emerging target in cancer therapy. *Clinical Cancer Research*, 24(11):2482–2490, May  
965 2018.
- 966 [6] Ayse Bassez, Hanne Vos, Laurien Van Dyck, Giuseppe Floris, Ingrid Arijs, Christine Desmedt, Bram Boeckx,  
967 Marlies Vanden Bempt, Ines Nevelsteen, Kathleen Lambein, Kevin Punie, Patrick Neven, Abhishek D Garg,  
968 Hans Wildiers, Junbin Qian, Ann Smeets, and Diether Lambrechts. A single-cell map of intratumoral  
969 changes during anti-PD1 treatment of patients with breast cancer. *Nat. Med.*, 27(5):820–832, May 2021.
- 970 [7] Joseph Bergensträhle, Ludvig Larsson, and Joakim Lundeberg. Seamless integration of image and molecular  
971 analysis for spatial transcriptomics workflows. *BMC Genomics*, 21(1):482, July 2020.
- 972 [8] Daniel B Burkhardt, Beatriz P San Juan, John G Lock, Smita Krishnaswamy, and Christine L Chaffer.  
973 Mapping phenotypic plasticity upon the cancer cell state landscape using manifold learning. *Cancer Discov.*,  
974 12(8):1847–1859, August 2022.
- 975 [9] Yuansi Chen, Julien Mairal, and Zaid Harchaoui. Fast and robust archetypal analysis for representation  
976 learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages  
977 1478–1485, 2014.
- 978 [10] Kelli A Connolly, Manik Kuchroo, Aarthi Venkat, Achia Khatun, Jiawei Wang, Ivana William, Noah I  
979 Hornick, Brittany L Fitzgerald, Martina Damo, Moujtaba Y Kasmani, Can Cui, Eric Fagerberg, Isabel  
980 Monroy, Amanda Hutchins, Julie F Cheung, Gena G Foster, Dylan L Mariuzza, Mursal Nader, Hongyu  
981 Zhao, Weiguo Cui, Smita Krishnaswamy, and Nikhil S Joshi. A reservoir of stem-like CD8+ T cells  
982 in the tumor-draining lymph node preserves the ongoing antitumor immune response. *Sci. Immunol.*,  
983 6(64):eabg7836, October 2021.
- 984 [11] Adele Cutler and Leo Breiman. Archetypal analysis. *Technometrics*, 36(4):338–347, 1994.
- 985 [12] Michele De Palma, Daniela Biziato, and Tatiana V Petrova. Microenvironmental regulation of tumour  
986 angiogenesis. *Nat. Rev. Cancer*, 17(8):457–474, August 2017.
- 987 [13] Mirjana Efremova, Miquel Vento-Tormo, Sarah A Teichmann, and Roser Vento-Tormo. CellPhoneDB:  
988 inferring cell-cell communication from combined expression of multi-subunit ligand-receptor complexes. *Nat.*  
989 *Protoc.*, 15(4):1484–1506, April 2020.
- 990 [14] Gökçen Eraslan, Lukas M Simon, Maria Mircea, Nikola S Mueller, and Fabian J Theis. Single-cell RNA-seq  
991 denoising using a deep count autoencoder. *Nat. Commun.*, 10(1):390, January 2019.

- 992 [15] Mark Esposito, Shridar Ganesan, and Yibin Kang. Emerging strategies for treating metastasis. *Nat. Cancer*,  
993 2(3):258–270, March 2021.
- 994 [16] Mónica T Fernandes, Sofia M Calado, Leonardo Mendes-Silva, and José Bragança. Cited2 and the modulation  
995 of the hypoxic response in cancer. *World Journal of Clinical Oncology*, 11(5):260–274, May 2020.
- 996 [17] David S. Fischer, Anna C. Schaar, and Fabian J. Theis. Modeling intercellular communication in tissues  
997 using spatial graphs of cells. *Nature Biotechnology*, October 2022.
- 998 [18] Oscar Franzén, Li-Ming Gan, and Johan L M Björkegren. PanglaoDB: a web server for exploration of mouse  
999 and human single-cell RNA sequencing data. *Database (Oxford)*, 2019, January 2019.
- 1000 [19] Christopher Heje Grønbech, Maximillian Fornitz Vording, Pascal N Timshel, Casper Kaae Sønderby, Tune H  
1001 Pers, and Ole Winther. scVAE: variational auto-encoders for single-cell gene expression data. *Bioinformatics*,  
1002 36(16):4415–4422, August 2020.
- 1003 [20] Laleh Haghverdi, Maren Büttner, F Alexander Wolf, Florian Buettner, and Fabian J Theis. Diffusion  
1004 pseudotime robustly reconstructs lineage branching. *Nat. Methods*, 13(10):845–848, October 2016.
- 1005 [21] Yuval Hart, Hila Sheftel, Jean Hausser, Pablo Szekely, Noa Bossel Ben-Moshe, Yael Korem, Avichai Tendler,  
1006 Avraham E Mayo, and Uri Alon. Inferring biological tasks using pareto analysis of high-dimensional data.  
1007 *Nat. Methods*, 12(3):233–5, 3 p following 235, March 2015.
- 1008 [22] Hamid Javadi and Andrea Montanari. Nonnegative matrix factorization via archetypal analysis. *Journal of*  
1009 *the American Statistical Association*, 115(530):896–907, 2020.
- 1010 [23] Vladimir Yu Kiselev, Tallulah S Andrews, and Martin Hemberg. Challenges in unsupervised clustering of  
1011 single-cell RNA-seq data. *Nat. Rev. Genet.*, 20(5):273–282, May 2019.
- 1012 [24] David Lähnemann, Johannes Köster, Ewa Szczurek, Davis J McCarthy, Stephanie C Hicks, Mark D Robinson,  
1013 Catalina A Vallejos, Kieran R Campbell, Niko Beerenwinkel, Ahmed Mahfouz, Luca Pinello, Pavel Skums,  
1014 Alexandros Stamatakis, Camille Stephan-Otto Attolini, Samuel Aparicio, Jasmijn Baaijens, Marleen Balvert,  
1015 Buys de Barbanson, Antonio Cappuccio, Giacomo Corleone, Bas E Dutilh, Maria Florescu, Victor Guriev,  
1016 Rens Holmer, Katharina Jahn, Thamar Jessurun Lobo, Emma M Keizer, Indu Khatri, Szymon M Kielbasa,  
1017 Jan O Korbel, Alexey M Kozlov, Tzu-Hao Kuo, Boudewijn P F Lelieveldt, Ion I Mandoiu, John C Marioni,  
1018 Tobias Marschall, Felix Mölder, Amir Niknejad, Lukasz Raczkowski, Marcel Reinders, Jeroen de Ridder,  
1019 Antoine-Emmanuel Saliba, Antonios Somarakis, Oliver Stegle, Fabian J Theis, Huan Yang, Alex Zelikovsky,  
1020 Alice C McHardy, Benjamin J Raphael, Sohrab P Shah, and Alexander Schönthuth. Eleven grand challenges  
1021 in single-cell data science. *Genome Biol.*, 21(1):31, February 2020.
- 1022 [25] Peter Langfelder and Steve Horvath. Eigengene networks for studying the relationships between co-expression  
1023 modules. *BMC Syst. Biol.*, 1(1):54, November 2007.
- 1024 [26] Devon A Lawson, Nirav R Bhakta, Kai Kessenbrock, Karin D Prummel, Ying Yu, Ken Takai, Alicia Zhou,  
1025 Henok Eyob, Sanjeev Balakrishnan, Chih-Yang Wang, Paul Yaswen, Andrei Goga, and Zena Werb. Single-cell  
1026 analysis reveals a stem-cell program in human metastatic breast cancer cells. *Nature*, 526(7571):131–135,  
1027 October 2015.
- 1028 [27] Michelle M Li and Marinka Zitnik. Deep contextual learners for protein networks. *ICML Computational*  
1029 *Biology*, 2021.
- 1030 [28] Arthur Liberzon, Chet Birger, Helga Thorvaldsdóttir, Mahmoud Ghandi, Jill P Mesirov, and Pablo Tamayo.  
1031 The molecular signatures database (MSigDB) hallmark gene set collection. *Cell Syst.*, 1(6):417–425, December  
1032 2015.
- 1033 [29] Romain Lopez, Jeffrey Regier, Michael B Cole, Michael I Jordan, and Nir Yosef. Deep generative modeling  
1034 for single-cell transcriptomics. *Nat. Methods*, 15(12):1053–1058, December 2018.

- 1035 [30] Malte D Luecken and Fabian J Theis. Current best practices in single-cell RNA-seq analysis: a tutorial.  
1036 *Mol. Syst. Biol.*, 15(6):e8746, June 2019.
- 1037 [31] Nicholas McGranahan and Charles Swanton. Clonal heterogeneity and tumor evolution: Past, present, and  
1038 the future. *Cell*, 168(4):613–628, February 2017.
- 1039 [32] Kevin R Moon, David van Dijk, Zheng Wang, Scott Gigante, Daniel B Burkhardt, William S Chen, Kristina  
1040 Yim, Antonia van den Elzen, Matthew J Hirn, Ronald R Coifman, Natalia B Ivanova, Guy Wolf, and Smita  
1041 Krishnaswamy. Visualizing structure and transitions in high-dimensional biological data. *Nat. Biotechnol.*,  
1042 37(12):1482–1492, December 2019.
- 1043 [33] Morten Mørup and Lars Kai Hansen. Archetypal analysis for machine learning and data mining. *Neurocom-*  
1044 *puting*, 80:54–63, 2012.
- 1045 [34] Adam Myszczyszyn, Anna M Czarnecka, Damian Matak, Lukasz Szymanski, Fei Lian, Anna Kornakiewicz,  
1046 Ewa Bartnik, Wojciech Kukwa, Claudine Kieda, and Cezary Szczylak. The role of hypoxia and cancer stem  
1047 cells in renal cell carcinoma pathogenesis. *Stem Cell Rev.*, 11(6):919–943, December 2015.
- 1048 [35] Bhupinder Pal, Yunshun Chen, François Vaillant, Bianca D Capaldo, Rachel Joyce, Xiaoyu Song, Vanessa L  
1049 Bryant, Jocelyn S Penington, Leon Di Stefano, Nina Tubau Ribera, Stephen Wilcox, Gregory B Mann,  
1050 kConFab, Anthony T Papenfuss, Geoffrey J Lindeman, Gordon K Smyth, and Jane E Visvader. A single-cell  
1051 RNA expression atlas of normal, preneoplastic and tumorigenic states in the human breast. *EMBO J.*,  
1052 40(11):e107333, June 2021.
- 1053 [36] Arden Perkins, Kimberly J Nelson, Derek Parsonage, Leslie B Poole, and P Andrew Karplus. Peroxiredoxins:  
1054 guardians against oxidative stress and modulators of peroxide signaling. *Trends Biochem. Sci.*, 40(8):435–445,  
1055 August 2015.
- 1056 [37] Lake-Ee Quek, Michelle van Geldermalsen, Yi Fang Guan, Kanu Wahi, Chelsea Mayoh, Seher Balaban, Angel  
1057 Pang, Qian Wang, Mark J Cowley, Kristin K Brown, Nigel Turner, Andrew J Hoy, and Jeff Holst. Glutamine  
1058 addiction promotes glucose oxidation in triple-negative breast cancer. *Oncogene*, 41(34):4066–4078, August  
1059 2022.
- 1060 [38] O Shoval, H Sheftel, G Shinar, Y Hart, O Ramote, A Mayo, E Dekel, K Kavanagh, and U Alon. Evolutionary  
1061 trade-offs, pareto optimality, and the geometry of phenotype space. *Science*, 336(6085):1157–1160, June  
1062 2012.
- 1063 [39] Kelly Street, Davide Risso, Russell B Fletcher, Diya Das, John Ngai, Nir Yosef, Elizabeth Purdom, and  
1064 Sandrine Dudoit. Slingshot: cell lineage and pseudotime inference for single-cell transcriptomics. *BMC*  
1065 *Genomics*, 19(1):477, June 2018.
- 1066 [40] Tim Stuart, Andrew Butler, Paul Hoffman, Christoph Hafemeister, Efthymia Papalex, William M Mauck  
1067 III, Yuhan Hao, Marlon Stoeckius, Peter Smibert, and Rahul Satija. Comprehensive integration of single-cell  
1068 data. *Cell*, 177:1888–1902, 2019.
- 1069 [41] Itay Tirosh, Benjamin Izar, Sanjay M Prakadan, Marc H Wadsworth, 2nd, Daniel Treacy, John J Trombetta,  
1070 Asaf Rotem, Christopher Rodman, Christine Lian, George Murphy, Mohammad Fallahi-Sichani, Ken  
1071 Dutton-Reeder, Jia-Ren Lin, Ofir Cohen, Parin Shah, Diana Lu, Alex S Genshaft, Travis K Hughes, Carly  
1072 G K Ziegler, Samuel W Kazer, Aleth Gaillard, Kellie E Kolb, Alexandra-Chloé Villani, Cory M Johannessen,  
1073 Aleksandr Y Andreev, Eliezer M Van Allen, Monica Bertagnolli, Peter K Sorger, Ryan J Sullivan, Keith T  
1074 Flaherty, Dennie T Frederick, Judit Jané-Valbuena, Charles H Yoon, Orit Rozenblatt-Rosen, Alex K Shalek,  
1075 Aviv Regev, and Levi A Garraway. Dissecting the multicellular ecosystem of metastatic melanoma by  
1076 single-cell RNA-seq. *Science*, 352(6282):189–196, April 2016.
- 1077 [42] Alexander Tong, Jessie Huang, Guy Wolf, David van Dijk, and Smita Krishnaswamy. TrajectoryNet: A  
1078 dynamic optimal transport network for modeling cellular dynamics. *Proc Mach Learn Res*, 119:9526–9536,  
1079 July 2020.

- 1080 [43] Cole Trapnell, Davide Cacchiarelli, Jonna Grimsby, Prapti Pokharel, Shuqiang Li, Michael Morse, Niall J  
1081 Lennon, Kenneth J Livak, Tarjei S Mikkelsen, and John L Rinn. The dynamics and regulators of cell fate  
1082 decisions are revealed by pseudotemporal ordering of single cells. *Nat. Biotechnol.*, 32(4):381–386, April  
1083 2014.
- 1084 [44] Tai-Hua Tsai, Ching-Chieh Yang, Tai-Chih Kou, Chang-En Yang, Jia-Zih Dai, Chia-Ling Chen, and Cheng-  
1085 Wei Lin. Overexpression of glut3 promotes metastasis of triple-negative breast cancer by modulating the  
1086 inflammatory tumor microenvironment. *Journal of Cellular Physiology*, 236(6):4669–4680, January 2021.
- 1087 [45] David van Dijk, Roshan Sharma, Juozas Nainys, Kristina Yim, Pooja Kathail, Ambrose J Carr, Cassandra  
1088 Burdziak, Kevin R Moon, Christine L Chaffer, Diwakar Pattabiraman, Brian Bierie, Linas Mazutis, Guy  
1089 Wolf, Smita Krishnaswamy, and Dana Pe'er. Recovering gene interactions from single-cell data using data  
1090 diffusion. *Cell*, 174(3):716–729.e27, July 2018.
- 1091 [46] Dongfang Wang and Jin Gu. VASC: Dimension reduction and visualization of single-cell RNA-seq data by  
1092 deep variational autoencoder. *Genomics Proteomics Bioinformatics*, 16(5):320–331, October 2018.
- 1093 [47] Qiushi Wang, Ann M. Bode, and Tianshun Zhang. Targeting cdk1 in cancer: mechanisms and implications.  
1094 *npj Precision Oncology*, 7(1), June 2023.
- 1095 [48] Hadley Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, 2016.
- 1096 [49] F Alexander Wolf, Philipp Angerer, and Fabian J Theis. SCANPY: large-scale single-cell gene expression  
1097 data analysis. *Genome Biol.*, 19(1), December 2018.
- 1098 [50] F Alexander Wolf, Fiona K Hamey, Mireya Plass, Jordi Solana, Joakim S Dahlin, Berthold Göttgens,  
1099 Nikolaus Rajewsky, Lukas Simon, and Fabian J Theis. PAGA: graph abstraction reconciles clustering with  
1100 trajectory inference through a topology preserving map of single cells. *Genome Biol.*, 20(1):59, March 2019.
- 1101 [51] Sunny Z Wu, Ghamdan Al-Eryani, Daniel Lee Roden, Simon Junankar, Kate Harvey, Alma Andersson,  
1102 Aatish Thennavan, Chenfei Wang, James R Torpy, Nenad Bartonicek, Taopeng Wang, Ludvig Larsson,  
1103 Dominik Kaczorowski, Neil I Weisenfeld, Cedric R Uytingco, Jennifer G Chew, Zachary W Bent, Chia-Ling  
1104 Chan, Vikkitharan Gnanasambandpillai, Charles-Antoine Dutertre, Laurence Gluch, Mun N Hui, Jane  
1105 Beith, Andrew Parker, Elizabeth Robbins, Davendra Segara, Caroline Cooper, Cindy Mak, Belinda Chan,  
1106 Sanjay Warrier, Florent Ginhoux, Ewan Millar, Joseph E Powell, Stephen R Williams, X Shirley Liu, Sandra  
1107 O'Toole, Elgene Lim, Joakim Lundeberg, Charles M Perou, and Alexander Swarbrick. A single-cell and  
1108 spatially resolved atlas of human breast cancers. *Nat. Genet.*, 53(9):1334–1347, September 2021.
- 1109 [52] Tianzhi Wu, Erqiang Hu, Shuangbin Xu, Meijun Chen, Pingfan Guo, Zehan Dai, Tingze Feng, Lang Zhou,  
1110 Wenli Tang, Li Zhan, Xiaocong Fu, Shanshan Liu, Xiaochen Bo, and Guangchuang Yu. clusterprofiler 4.0:  
1111 A universal enrichment tool for interpreting omics data. *Innovation (Camb.)*, 2(3):100141, August 2021.
- 1112 [53] Yoshinori Yoshida, Kazutoshi Takahashi, Keisuke Okita, Tomoko Ichisaka, and Shinya Yamanaka. Hypoxia  
1113 enhances the generation of induced pluripotent stem cells. *Cell Stem Cell*, 5(3):237–241, September 2009.
- 1114 [54] Guangchuang Yu, Li-Gen Wang, Yanyan Han, and Qing-Yu He. clusterprofiler: an R package for comparing  
1115 biological themes among gene clusters. *OMICS*, 16(5):284–287, May 2012.