1   **Self-assembling Manifolds in Single-cell RNA Sequencing Data**

2

3   Alexander J. Tarashansky[1], Yuan Xue[1], Pengyang Li[1], Stephen R. Quake[1,2,3], Bo Wang[1,4,*]

4

5   [1]Department of Bioengineering, [2]Department of Applied Physics, Stanford University, Stanford,

6   CA, USA.

7   [3]Chan Zuckerberg Biohub, San Francisco, CA, USA.

8   [4]Department of Developmental Biology, Stanford University School of Medicine, Stanford, CA,

9   USA.

10

11   *Correspondence should be addressed to B.W. (wangbo@stanford.edu)

**Abstract**

Single-cell RNA sequencing has spurred the development of computational methods that enable researchers to classify cell types, delineate developmental trajectories, and measure molecular responses to external perturbations. Many of these technologies rely on their ability to detect genes whose cell-to-cell variations arise from the biological processes of interest rather than transcriptional or technical noise. However, for datasets in which the biologically relevant differences between cells are subtle, identifying these genes is a challenging task. We present the self-assembling manifold (SAM) algorithm, an iterative soft feature selection strategy to quantify gene relevance and improve dimensionality reduction. We demonstrate its advantages over other state-of-the-art methods with experimental validation in identifying novel stem cell populations of *Schistosoma*, a prevalent parasite that infects hundreds of millions of people. Extending our analysis to a total of 56 datasets, we show that SAM is generalizable and consistently outperforms other methods in a variety of biological and quantitative benchmarks.

**Introduction**

Single-cell RNA sequencing (scRNAseq) datasets typically contain tens of thousands of genes, although many of them may not be informative for differentiating between cell types or states. Feature selection is thus commonly used to select a subset of genes prior to downstream analyses, such as manifold reconstruction and cell clustering (Crow et al., 2018; Satija et al., 2015; Vallejos et al., 2015). However, current approaches have two major limitations.

First, feature selection methods filter genes based on arbitrarily or empirically chosen thresholds, small changes in which may result in different gene sets (Vallejos et al., 2017). In addition, the selection of features typically operates under the assumption that genes with highly variable expression between individual cells capture biologically meaningful variation. Because single-cell transcriptomes are inevitably contaminated by a combination of random transcriptional and technical noise (Grün et al., 2014), the variation in biologically relevant genes may be hard to distinguish from the background noise, especially when the differences between cell populations are subtle. Resolving these differences, or "signals", is essential to study a variety of biological problems, including identifying cell subtypes (Olsson et al., 2016; Treutlein et al., 2014; Lönnberg et al., 2017; Fincher et al., 2018; Baron et al., 2016) and quantifying the effects of molecular perturbations to otherwise homogeneous populations of cells (Lane et al., 2017). In such datasets, only a small fraction of the genes, and therefore only a small fraction of the total variation, may contain the signals relevant for distinguishing cell types. Choosing these features ahead of time without *a priori* knowledge remains an unmet computational challenge.

The second limitation is that existing methods have been almost exclusively benchmarked on well-annotated, gold standard datasets with clearly distinguishable cell types (Wang et al., 2017; Kiselev et al., 2017; Duò et al., 2019; Bahlo et al., 2018). These datasets are not informative

48    for distinguishing the performance between methods, because the differences between cell types

49    are relatively straightforward to detect. However, evaluating the performance of feature selection

50    and/or dimensionality reduction methods on datasets with more subtle signals is difficult as their

51    ground truth labels are typically ambiguous or nonexistent.

52         To overcome the shortcomings of current feature selection approaches, here, we introduce

53    the Self-Assembling Manifold (SAM) method, an unsupervised, "soft feature selection" algorithm

54    that iteratively rescales gene expressions to refine a nearest neighbor graph of cells until the graph

55    converges to a stable solution. At each iteration, SAM assigns more weight to genes that are

56    spatially variable across the constructed graph, and this weighted gene expression is then used to

57    improve the next nearest neighbor assignment. SAM presents two advantages: it rescales all genes

58    according to their weights, solving the problem of thresholding, and it prioritizes genes that are

59    variable across the intrinsic manifold of the data rather than selecting genes that are variable across

60    individual cells.

61         Second, in order to better distinguish the performance between methods, we define a

62    network sensitivity measure to identify datasets with subtle signals. With limited annotations in

63    most high-sensitivity datasets, we introduce unsupervised graph-based metrics to quantify the

64    degree of structure within the reconstructed manifolds for comparison between methods. In

65    addition, we perform benchmarking using known ground truth labels on simulated datasets

66    spanning a wide range of sensitivities by introducing increasing levels of noise to well-annotated

67    datasets. These analyses reveal that SAM consistently improves feature selection and cell

68    clustering.

69         To demonstrate the utility of SAM in practice, we provide an in-depth analysis of two

70    datasets that are challenging to analyze using existing methods: stem cells in a human parasitic

4

71    worm, *Schistosoma*, and activated macrophages (Lane et al., 2017). We show that SAM can

72    capture novel biology undetectable by other methods and validate these results with experimental

73    evidence.

## Results

### *The SAM algorithm*

The SAM algorithm begins with a random *k*-nearest neighbor (kNN) graph and averages the expression of each cell with its *k* nearest neighbors: $C = \frac{1}{k}NE$, where $N$ is the directed adjacency matrix and $E$ is the gene expression matrix (**Figure 1a**). For each gene *i*, SAM computes a spatial dispersion factor of the averaged expressions $C_i$, which measures variation across neighborhoods of cells rather than individual cells (**Methods**). These dispersions are used to calculate the gene weights, which then rescale the expression matrix: $\hat{E} = E\sqrt{W_D}$, where $W_D$ is a diagonal matrix with gene weights along the diagonal. Using the rescaled expressions $\hat{E}$, we compute a pairwise cell distance matrix and update the assignment of each cell's *k*-nearest neighbors accordingly. This cycle continues until the gene weights converge.

To demonstrate the implementation and utility of SAM, below we analyze a challenging dataset comprised of a few hundred relatively homogeneous stem cells isolated from *Schistosoma mansoni* (**Figure 1–Figure supplement 1**), a widespread human pathogen (Hoffman et al., 2014). Using a protocol we have established previously (Wang et al., 2018), these cells were collected by sorting dividing cells from juvenile parasites harvested from their mouse hosts at 2.5 weeks post infection. At this stage, the parasites use an abundant stem cell population (~15-20% of the total number of cells) for rapid organogenesis and growth (Wang et al., 2013; Wang et al., 2018). Testing several existing methods (Wang et al., 2017; Kiselev et al., 2017; Satija et al., 2015), we found that they were not able to identify distinct cell populations in this dataset. In contrast, SAM finds a stable solution independent of initial conditions (**Figure 1b**). A graph structure with clearly separated cell populations self-assembles through the iterative process (**Figure 1c**). In parallel, the gene weights converge onto the final weight vector. Eventually, only a small fraction of genes

6

97    (~1%) are strongly weighted and useful for separating cell clusters, reflecting the inherent

98    difficulty of analyzing this dataset.

99    **Figure 1d** shows that SAM iteratively improves a series of graph characteristics, including

100    the network-average clustering coefficient (NACC), modularity, and Euclidean norm of the spatial

101    dispersions (**Methods**). The NACC and modularity quantify the degree of structure within the

102    graphs – graphs with high NACC and modularity have regions of high density separated by regions

103    of low density. The dispersion quantifies the spatial organization of gene expression – the higher

104    the spatial dispersion the less uniformly distributed the gene expressions are along the graph.

105    Importantly, we verified that SAM does not artificially boost these metrics in data that lack

106    inherent structure: when applying SAM to a randomly shuffled expression matrix, none of these

107    metrics increased from the random initial conditions.

108

109    *SAM identifies novel subpopulations within schistosome stem cells*

110    Visualizing the converged graph in two dimensions using Uniform Manifold

111    Approximation and Projection (UMAP, Becht et al., 2019), we find that cells can be separated into

112    three main populations, with Louvain clustering (Blondel et al., 2008) further splitting one of these

113    clusters into two subpopulations (**Figure 2a**). In contrast, other commonly-used dimensionality

114    reduction methods, such as principal component analysis (PCA), Seurat (Satija et al., 2015), and

115    SIMLR (Wang et al., 2017), failed to distinguish these cell populations (see **Methods** for the

116    selection of algorithms for comparison). **Supplementary Table 1** lists genes with high SAM

117    weights, which includes most markers that were previously implicated to be enriched in subsets of

118    schistosome stem cells (Wang et al., 2013; Wang et al., 2018).

119    **Figure 2b** shows that the three populations include previously characterized δ′-cells, which

120     specifically express an RNA binding protein *nanos-2* (Smp_051920), and ε-cells, which are

121     marked by the expression of *eledh* (*eled*, Smp_041540) (Wang et al., 2018). More importantly,

122     SAM reveals a novel stem cell population, $\mu$, comprising ~30% of all sequenced cells ($\mu$ denotes

123     muscle progenitors as discussed below). While $\mu$-cells express ubiquitous stem cells markers (e.g.,

124     *ago2-1*, Smp_179320; *cyclin B*, Smp_082490) and cell cycle regulators (**Figure 2–Figure**

125     **supplement 1a**) (Collins et al., 2013; Wang et al., 2013; Wang et al., 2018), they are also strongly

126     enriched for a large set of genes, with a calcium binding protein (*cabp*, Smp_005350), an actin

127     protein (Smp_161920), an annexin homolog (Smp_074140), a helix-loop-helix transcription factor

128     (*dhand*, Smp_062490), and a phosphatase (*dusp10*, Smp_034500) as the most specific markers of

129     this population in comparison to other stem cells (**Figure 2–Figure supplement 1b**).

130         Fluorescent *in-situ* hybridization (FISH) in conjunction with EdU labeling of dividing cells

131     reveals that $\mu$-cells (*cabp*$^+$EdU$^+$) are distributed near the parasite surface right beneath a layer of

132     post-mitotic differentiated cells that also express *cabp* (**Figure 2c**). Close to the parasite surface,

133     there are two major cell types intertwined in space: the skin-like tegumental cells and the body

134     wall muscle cells. However, $\mu$-cells express none of the recently identified markers in tegumental

135     progenitors (Wendt et al., 2018), suggesting that they may be associated with the muscle lineage.

136     To test this idea, we performed double FISH experiments and observed in post-mitotic *cabp*$^+$ cells

137     the coexpression of a set of canonical muscle markers (Witchley et al., 2013), including

138     tropomyosin (Smp_031770), myosin (Smp_045220), troponin (Smp_018250), and collagen

139     (Smp_170340) (**Figure 2d**). These results suggest that *cabp* is a specific marker for parasite body

140     wall muscles and $\mu$-cells are the muscle progenitors. Why the juvenile parasites maintain such an

141     active pool of muscle progenitors will be an important question for future studies.

142         In addition, SAM identifies two subpopulations among ε-cells: $\varepsilon_a$-cells that are highly

8

143     enriched for an aschaete-scute transcription factor (*astf*, Smp_142120), and $\varepsilon_\beta$-cells that

144     abundantly express another basic helix-loop-helix protein (*bhlh*, Smp_087310) (**Figure 2b**, right

145     panels). FISH experiments confirm these cells to be in close spatial proximity but with no

146     coexpression of *astf* and *bhlh* (**Figure 2e**). Moreover, we observed with FISH that there are fewer

147     *astf*$^+$ cells in larger, more matured juveniles, suggesting $\varepsilon_\alpha$-cells are a dynamic population during

148     development. To verify this observation, we sequenced another ~370 stem cell from juveniles at a

149     later developmental time point (3.5 weeks post infection). After correcting for batch effects in the

150     combined 2.5- and 3.5-week datasets using the mutual nearest neighbors (MNN) algorithm

151     (Haghverdi et al., 2018), we find that $\delta'$-, $\mu$-, and $\varepsilon_\beta$-cells remain relatively constant throughout

152     both time points, whereas $\varepsilon_\alpha$-cells comprise a significantly smaller fraction of the stem cells at 3.5

153     weeks (7%) compared to 21% at 2.5 weeks (**Figure 2f**). Taken together, these analyses

154     demonstrate that SAM can identify experimentally validated stem cell populations that are

155     previously too subtle to separate using other methods but are closely associated with the

156     schistosome development.

157        The critical difference between SAM and other methods lies in how they select genes for

158     manifold reconstruction. SAM prioritizes genes with variable expressions across neighborhoods

159     of cells rather than individual cells as in other methods (e.g., Seurat). **Figure 2g** shows that genes

160     with high standardized dispersion across individual cells often have low SAM weights, indicating

161     that these highly variable genes (HVGs) are irrelevant to the topological relationships between

162     cells. Other methods (e.g. SC3, Kiselev et al., 2017) identify marker genes based on differential

163     gene expression between cell clusters, but this approach suffers when cell cluster assignment is

164     poor, especially when discrete cell groups are difficult to separate or absent. Indeed, SC3 failed in

165     the default mode as it incorrectly predicted there to be only one cluster in the schistosome dataset.

166    After we manually increased the number of clusters, SC3 could recover a few of the marker genes

167    associated with only one ($\mu$-cells, blue symbols in **Figure 2h**) of the populations detected by SAM.

168    Furthermore, changing the number of clusters resulted in different solutions and large variability

169    in SC3 scores for its top ranked genes.

170

171    ***SAM outperforms other state-of-the-art methods in extensive quantitative benchmarking***

172         Below, we assess the general applicability of SAM by benchmarking its performance

173    against state-of-the-art scRNAseq analysis methods on a large collection of datasets. We focus on

174    three methods, i.e., Seurat, SIMLR, and SC3, as they are mostly unsupervised, have been broadly

175    used, and were shown to outperform other methods through extensive benchmarking (Kiselev et

176    al., 2017; Wang et al., 2017; Duò et al., 2019; Bahlo et al., 2018; Tian et al., 2019, see **Methods**

177    for the selection of algorithms for comparison). We first benchmark against nine datasets

178    (**Supplementary Table 2**) that have high-confidence annotations to evaluate the accuracy of SAM

179    in assigning cell clusters. We find that SAM has the highest Adjusted Rand Index (ARI, a measure

180    of clustering accuracy) (Hubert and Arabie, 1985) on eight out of the nine datasets and does not

181    over cluster the data (**Figure 3a**). Furthermore, SAM converges to the same set of gene weights

182    for all datasets analyzed (**Figure 3b**, **Figure 3-Figure supplement 1a**) and its performance is

183    robust to the choice of parameters and random initial conditions (**Figure 3-Figure supplement**

184    **1b-c**). In contrast, applying SAM to randomly generated datasets (**Methods**), the resulting gene

185    weights are highly dissimilar across random initial conditions (**Figure 3b**), indicating that SAM

186    does not converge to a stable solution on datasets with no intrinsic structure. Finally, the scalability

187    of SAM is similar to that of Seurat, capable of analyzing hundreds of thousands of cells in minutes

188    (**Figure 3c**), whereas SIMLR and SC3 are orders of magnitudes slower and thus excluded from

10

189    further benchmarking which requires the analysis of many more datasets.

190         Because these nine datasets are all comprised of clearly distinguishable cell types, they

191    may not represent the performance of methods on other datasets that contain cell populations that

192    are only subtly different. To identify such datasets, we introduce a network sensitivity metric that

193    quantifies the changes in the cell-to-cell distances when randomly selecting a subset of features

194    from the gene expression matrices (**Methods**). High network sensitivity indicates that changes to

195    the selected features strongly alters the resulting topological network. Networks that are robust to

196    the selected features correspond to datasets that have many redundant signals or genes

197    corroborating that network structure. In the datasets we compiled (**Supplementary Table 2**), all

198    broadly-used benchmarking datasets have lower sensitivities whereas the schistosome dataset,

199    which we have shown to be challenging to analyze for other methods, has the highest sensitivity

200    (**Figure 4a**). The fraction of genes with large SAM weights (>0.5) is negatively correlated with

201    the network sensitivity, suggesting that the biologically relevant variation in datasets with high

202    sensitivity is captured by relatively fewer genes (**Figure 4b**). Analyzing all 56 datasets, we found

203    that SAM improves the clustering, modularity, and spatial organization of gene expression across

204    the graph in comparison to Seurat as the datasets become increasingly sensitive (**Figure 4c**).

205         Evaluating the clustering accuracy for the highly sensitive datasets, however, is

206    challenging, because many of them have incomplete or nonexistent cell type annotations.

207    Therefore, we use the nine well-annotated benchmarking datasets to simulate data across a wide

208    spectrum of sensitivities. For this, we corrupt the data by randomly permuting gradually increasing

209    fractions of the gene expressions. As illustrated by the Darmanis dataset (Darmanis et al., 2015),

210    **Figure 5a** shows that the sensitivity increases along with the corruption. SAM's ARI scores only

211    marginally decrease as the corruption (and thereby sensitivity) increases, whereas Seurat's

11

212    performance rapidly deteriorates. A similar contrast was observed with the NACC, modularity,

213    and spatial dispersion between SAM and Seurat. Importantly, passing the genes with high SAM

214    weights into Seurat rescued its performance across all metrics, indicating that SAM is able to

215    consistently capture the genes relevant to the underlying structure of the data even with increasing

216    levels of noise and illustrating the robustness of its feature selection strategy compared to the HVG

217    filtering approach of Seurat. These observations generalize to all nine benchmarking datasets,

218    quantified by the area under the curves (AUC) of the metrics with respect to corruption (**Figure**

219    **5b**).

220

221    *SAM clusters macrophages by their activation dynamics with proper temporal ordering*

222          We next highlight another dataset to show that SAM can recover biologically meaningful

223    information that other methods fail to capture. We chose this example, which contains ~600

224    macrophages treated with lipopolysaccharide (LPS) when individually trapped in microfluidic

225    channels (Lane et al., 2017), because it has high network sensitivity (**Figure 4a**) and has

226    accompanying single cell functional data of macrophage activation dynamics that may help

227    validate the results of our analysis. Applied to this dataset, SAM initially identifies two clusters

228    (**Figure 6a**, top). Performing gene set enrichment analysis (GSEA, **Methods**, Subramanian et al.,

229    2005), we find that genes with high SAM weights are dominated by cell cycle-related processes,

230    with one of the clusters heavily enriched for cell cycle genes (e.g., Top2a, Mki67, **Figure 6-Figure**

231    **supplement 1a**). After removing the cell cycle effects (**Methods**), SAM identifies two different

232    clusters in which cells are properly ordered by the time since LPS induction, with the highly

233    weighted genes being primarily involved in immune signaling (**Figure 6a**, bottom). These

234    observations demonstrate that, in conjunction with GSEA, the quantitative gene weights output by

235    SAM can be used to infer the biological pathways that drive the clustering of cells.

236         One of the two clusters is enriched for TNFα expression (**Figure 6b**). It is known that LPS

237    activates two independent pathways, one through the innate immune signal transduction adaptor

238    (Myd88) and the other through the TIR-domain-containing adapter-inducing interferon-β (TRIF)

239    (Lee et al., 2009). While the Myd88 pathway directly activates NF-κB, the TRIF pathway first

240    induces the secretion of TNFα, which subsequently binds to its receptor, TNFR, to prolong the

241    activation of NF-κB (**Figure 6c**). **Figure 6d** and **Figure 6-Figure supplement 1b** show examples

242    of genes that are highly enriched with TNFα, a number of which are inflammatory factors known

243    to accumulate due to prolonged NF-κB activation (Lane et al., 2017). These results suggest that

244    SAM grouped the cells based on their activated signaling pathways: one cluster is activated

245    through both Myd88 and TRIF pathways (MT) while the other is only activated through Myd88

246    (M).

247         To further verify that the separation between the MT and M clusters truly reflects the

248    dichotomy in cellular response to LPS induction, we noted that this dataset combines scRNAseq

249    with live-cell imaging of NF-κB activity in single cells. This allows us to directly test if the MT

250    and M clusters correspond to different signaling dynamics (**Methods**). We found that most of the

251    cells with prolonged NF-κB response (i.e., cells showing broad peaks of NF-κB activation in time)

252    are in fact in the MT cluster (**Figure 6e-f**, and **Figure 6-Figure supplement 2a**), consistent with

253    the expectation that TNFα signaling prolongs NF-κB activation. Although our interpretation of the

254    data matches that provided in the original study, we were able to analyze the dataset with almost

255    no *a priori* knowledge. In contrast, the original study required extensive manual curation, analyzed

256    only a subset of the dataset, and could not group cells by their NF-κB activation dynamics from

257    the gene expression data alone. Similarly, Seurat and SIMLR were unable to order the cells by the

13

258    time since LPS induction or group cells based on their activation dynamics after removing the cell

259    cycle effects (**Figure 6g**, and **Figure 6-Figure supplement 2b-c**).

**Discussion**

260

261        Here, we introduced a scRNAseq analysis method, SAM, that uses an unsupervised, robust,

262    and iterative strategy for feature selection and manifold reconstruction. As demonstrated by our

263    analysis of the schistosome stem cells and activated macrophages, SAM can capture biology that

264    is undetectable by other methods. While SAM has consistently higher clustering accuracy than

265    other state-of-the-art methods on datasets containing clearly distinct cell types, its advantages are

266    especially apparent on datasets in which cell states or types are only distinguishable through subtle

267    differences in gene expression.

268        The strength of SAM lies in the integration of three algorithmic components: spatial

269    dispersion to measure feature relevance, soft feature selection, and the iterative scheme. By

270    averaging the gene expression of a cell with that of its neighbors, the spatial dispersion quantifies

271    the variation across neighborhoods of cells rather than individual cells. Genes with high spatial

272    dispersion are more likely to be biologically relevant as they are capable of separating cells into

273    distinct topological locations. Soft feature selection includes all genes and weights their

274    contribution to the manifold reconstruction by their spatial dispersions. This mitigates the

275    shortcoming of existing approaches in which the selection of features is a binary decision: genes

276    are either included or not depending on arbitrarily chosen thresholds.

277        The conceptual challenge here is that calculating the gene weights requires the manifold,

278    but reconstructing the manifold requires the gene weights for feature selection. SAM thus uses an

279    iterative strategy to converge onto both the gene weights and the corresponding graph topology

280    from a random initial graph. Each successive iteration refines the gene weights and network

281    structure until the algorithm converges. Empirically, for all datasets analyzed we have shown that

282    SAM converges onto a stable solution and is robust to the random initial conditions. Practically,

283    we could initialize SAM using the graph output of another method such as Seurat, but using

284    random initial conditions avoids potential biases in the analysis and enables the evaluation of the

285    stability of SAM.

286          To demonstrate the strengths of SAM in practice, we analyzed the schistosome stem cells

287    and identified novel stem cell populations that were validated by FISH experiments (**Figure 2**). In

288    the analysis of activated macrophages, we showed that SAM can simultaneously order cells by the

289    time since LPS induction and group cells according to their respective activated signaling

290    pathways. We have validated this result using the live-cell imaging data presented in the original

291    study (**Figure 6**).

292          We expect the application of SAM is not limited to feature selection, cell clustering, and

293    manifold reconstruction; it can be readily integrated with existing analytical pipelines as its gene

294    weights and reconstructed manifolds can be used in downstream analyses. For example, we have

295    shown how the genes ranked by their SAM weights can be used as input to GSEA to determine

296    the biological processes enriched among the highly weighted genes (**Figure 6**), thus directly

297    testing if the weights reflect biologically relevant genes. Additionally, the manifold reconstructed

298    by SAM can be used as input to pseudotemporal ordering algorithms (Setty et al., 2016; Trapnell

299    et al., 2014).

300          Beyond the two example case studies, we have rigorously evaluated SAM on a total of 56

301    datasets. While previous studies benchmarked on datasets with clearly defined cell populations,

302    we defined a network sensitivity measure to rank the datasets based on the inherent difficulty of

303    their analysis (**Figure 4**). Using these datasets, we showed that SAM consistently outperforms

304    other methods in terms of both cell clustering accuracy measured by ground truth annotations, and

305    manifold reconstruction measured by quantitative graph characteristics. These improvements can

306    be attributed to the robust selection of features relevant for cell clustering and manifold

307    reconstruction even in the presence of significant amounts of random noise, as shown in the

308    corruption tests (**Figure 5**). Overall, the network sensitivity and quantitative benchmarking metrics

309    should help in characterizing the performance of future scRNAseq analysis methods across a wider

310    variety of datasets.

311    **Materials and Methods**

312

313    **Code and data availability.** The SAM source code and tutorials can be found at

314    https://github.com/atarashansky/self-assembling-manifold. The schistosome stem cell scRNAseq

315    data generated in this study is available through the Gene Expression Omnibus (GEO) under

316    accession number GSE116920.

317

318    **Data processing. Supplementary Table 2** summarizes all datasets used in this study as well as

319    the methods used to convert raw sequence read counts to gene expression, such as TPM (transcripts

320    per million), CPM (counts per million), RPKM (reads per kilobase per million), or FPKM

321    (fragments per kilobase per million). Datasets with asterisks next to their accession numbers are

322    sourced from the *conquer* database (Soneson and Robinson, 2018). The nine benchmarking

323    datasets used with high-confidence labels are marked by crosses. Gene expression is measured in

324    log space with a pseudocount of 1 (e.g., $\log_2(TPM+1)$). Genes expressed ($\log_2(TPM+1)>1$) in

325    fewer than 1% or more than 99% of cells are excluded from downstream analysis as these genes

326    lack statistical power. To reduce the influence of technical noise near the molecular detection limit,

327    we set gene expression to zero when $\log_2(TPM+1)<1$. We denote the resulting expression matrix

328    as $E$.

329            In the SAM algorithm (see below), we either standardize the gene expression matrix $E$ to

330    have zero mean and unit variance per gene (which corrects for differences in distributions between

331    genes) or normalize the expressions such that each cell has unit Euclidean (L2) norm (which

332    prevents cells with large variances in gene expressions from dominating downstream analyses)

333    prior to dimensionality reduction. In the below section, we denote the standardized or normalized

334    expression matrix as $\bar{E}$. Empirically, we have found that standardization performs well with large,

18

335    sparse datasets collected through droplet-based methods, whereas L2-normalization is more

336    suitable for smaller datasets with higher sequencing depth such as those prepared with the Smart-

337    Seq2 protocol (Picelli et al., 2013). This is likely due to the fact that standardization amplifies the

338    relative expression of genes specific to small populations in large datasets, thereby making them

339    easier to identify. In contrast, standardization decreases the relative expression of genes specific

340    to populations comprising larger fractions of the data, as is typically the case in smaller datasets,

341    thereby making distinct populations more difficult to identify. **Supplementary Table 2** documents

342    the preprocessing step used for each dataset.

343

344    **The SAM algorithm.** After first generating a random kNN adjacency matrix, the SAM algorithm

345    goes through three steps that are repeated until convergence.

346

347    *1) Calculate the gene weights*

348         First, the expression of each cell is averaged with its k-nearest neighbors:

$$C = \frac{1}{k}NE \tag{1}$$

349    where $N$ is the directed adjacency matrix for the kNN graph, and $E$ is the original $n \times m$ gene

350    expression matrix with rows as cells and columns as genes. Here, we do not use $\bar{E}$ as it may

351    contain negative values, for which dispersion is ill-defined. For each gene $i$, SAM computes the

352    Fano factor from the averaged expressions $C_i$:

$$\mu_{C_i} = \frac{1}{n}\sum_{j=1}^{n} C_{ji} \tag{2}$$

19

$$\sigma^2_{C_i} = \frac{1}{n}\sum_{j=1}^{n}(C_{ji} - \mu_{C_i})^2 \tag{3}$$

$$F_i = \frac{\sigma^2_{C_i}}{\mu_{C_i}} \tag{4}$$

353   where $\mu_{C_i}$ is the mean and $\sigma^2_{C_i}$ is the variance. We use the Fano factor to measure the gene

354   expression variance relative to the mean in order to account for the fact that genes with high mean

355   expressions tend to have higher variability. Computing the Fano factors based on the kNN-

356   averaged expressions links gene dispersion to the cellular topological structure: Genes that have

357   highly variable expressions among individual cells but are homogeneously distributed across the

358   topological representation should have small dispersions. $k$, set by default to 20, determines the

359   topological length scale over which variations in gene expression are quantified. **Figure 3-Figure**

360   **supplement 1b** shows that the downstream analysis is robust to the specific choice of $k$.

361   Additionally, the choice of $k$ does not significantly affect runtime complexity or scalability.

362   To compute the gene weights, we normalize the Fano factors to be between 0 and 1. First,

363   we saturate the Fano factors to ensure that outlier genes with large spatial dispersions do not skew

364   the distribution of weights:

$$\{F_i | F_i > z\} = z \tag{5}$$

365   where $z$ is by default the mean of the largest 50 dispersions. In other words, Fano factors exceeding

366   this number are saturated to be $z$. We then calculate the gene weights as:

$$W_i = \frac{F_i}{z} \tag{6}$$

367   *2) Rescale the expression matrix*

368   Having calculated the gene weights, SAM multiplies them into the preprocessed expression

369   matrix:

20

$$\hat{E} = \bar{E}\sqrt{W_D} \tag{7}$$

370    where $\bar{E}$ is the standardized or normalized expression matrix and $W_D$ is a diagonal matrix with $W_i$

371    along the diagonal. This matrix multiplication linearly rescales the gene expression variances and

372    gene-gene covariances by their respective weights, attenuating the influence of genes with low

373    dispersions across neighborhoods.

374

375    *3) Updating the kNN graph*

376         To compute pairwise cell-cell distances, we perform PCA on the rescaled expression

377    matrix $\hat{E}$. The variance-scaling operation in **Eq. 7** improves the robustness of PCA to variations in

378    genes that are uniformly distributed along the current graph (i.e., genes with low weights).

379    Furthermore, this weighting strategy eliminates the typical requirement of selecting a subset of

380    HVGs to feed into PCA, which often relies on arbitrary thresholds and heuristics. To perform PCA,

381    we first mean center $\hat{E}$ to form $\hat{E}_{\mu}$:

$$\hat{E}_{\mu} = \hat{E} - \frac{1}{n}ee^T\hat{E} \tag{8}$$

382    where $e$ is a column vector of ones with dimension $n$. We then compute the Singular Value

383    Decomposition (SVD) of $\hat{E}_{\mu}$:

$$\hat{E}_{\mu} = USV^T \tag{9}$$

384    with the principal components defined as

$$P = US \tag{10}$$

385    The eigenvalues corresponding to the eigendecomposition of the gene-gene covariance matrix

386    are defined in terms of the singular values as

$$\Lambda = \frac{S^2}{n-1} \tag{11}$$

387    where $S$ is a diagonal matrix with singular values along the diagonal. Using the PC matrix $P$, SAM

388    computes a pairwise cell-cell distance matrix. While typical dimension reduction approaches select

389    a subset of the PCs, which is often subjective or requires computationally intensive maximum-

390    likelihood approaches, we include all PCs and scale their variances by their corresponding

391    eigenvalues:

$$\hat{P} = P\sqrt{\Lambda} \tag{12}$$

392    As a result, PCs with small eigenvalues are weighted less in the calculation of the distance between

393    cells $i$ and $j$, $D_{P_i P_j}$. $D_{P_i P_j}$ is the Pearson correlation or Euclidean distance between rows $P_i$ and $P_j$

394    in the PC matrix. Pearson correlation distance is used by default, although **Figure 3-Figure**

395    **supplement 1c** shows that SAM is robust to the choice of distance metric. Using the distances to

396    define the $k$-nearest neighbors for each cell, SAM updates the kNN matrix and repeats steps 1-3.

397    The algorithm continues until convergence, defined as when the RMSE between gene weights in

398    adjacent iterations diminished as defined by:

$$\sqrt{\frac{1}{m} \sum_{j=1}^{m} (W_{i,j} - W_{i+1,j})^2} < 5 \times 10^{-3} \tag{13}$$

399    where $m$ is the number of genes and $W_{i,j}$ is the weight for gene $j$ at iteration $i$.

400

401    **Visualization.** To visualize the topological structure identified by SAM, we feed the final

402    weighted PCA matrix, $\hat{P}$, into UMAP (Becht et al., 2019) using Pearson correlation as the distance

403    metric by default. To directly visualize the final kNN adjacency matrix (**Figure 1c**), we used the

404    Fruchterman-Reingold force-directed layout algorithm and drawing tools implemented in the

405     Python package *graph-tool* (Peixoto, 2017).

406

407     **Choosing the benchmarking methods.** We used three main criteria for choosing the benchmark

408     scRNAseq analysis methods: they should be widely used, have done extensive benchmarking

409     against other methods, and be mostly unsupervised. We found on Web of Science that among the

410     highest cited scRNAseq analysis tools in 2017-2018 are Seurat, SC3, SIMLR, Reference

411     Component Analysis (RCA, Li et al., 2017), Monocle (Trapnell et al., 2014), zero-inflated factor

412     analysis (ZIFA, Pierson and Yau, 2015), and Wishbone (Setty et al., 2016), of which we chose

413     Seurat, SC3, and SIMLR.

414         SC3 is a consensus clustering algorithm that has done rigorous benchmarking against other

415     methods such as SINCERA (Guo et al., 2015), SNN-Cliq (Xu and Su, 2015) and pcaReduce

416     (Žurauskien and Yau, 2016) on 12 datasets with ground truth labels. SIMLR, a dimensionality

417     reduction and clustering algorithm, evaluated its clustering performance on four annotated datasets

418     against eight other dimensionality reduction methods, including PCA, Factor Analysis (FA), t-

419     SNE, multidimensional scaling (MDS), and (ZIFA). Both methods have demonstrated the highest

420     clustering accuracy across most of the tested datasets. Additionally, as both methods have built-in

421     functions to estimate the number of clusters present within the data, they are largely unsupervised.

422     We also selected Seurat as one of the benchmarking methods, because it is arguably the most

423     widely used tool for dimensionality reduction and clustering of scRNAseq data and has performed

424     well in rigorous benchmarking studies against various methods including SC3, SIMLR, RCA, and

425     pcaReduce (Duò et al., 2019; Bahlo et al., 2018).

426         We did not select Reference Component Analysis as it is primarily designed for cases in

427     which an atlas of bulk, cell-type specific, reference transcriptomes is present. Additionally, we did

23

428    not benchmark against Monocle and Wishbone, because they are pseudotime analysis methods

429    and are meant for datasets with continuous branching processes such as cell differentiation.

430    However, it is important to note that SAM can be used for dimensionality reduction upstream of

431    pseudotime algorithms for such datasets. Finally, we did not benchmark against ZIFA as it has

432    already been shown to have lower clustering accuracy than SIMLR.

433         In addition to measuring clustering accuracy, we also introduce the unsupervised NACC,

434    modularity, and spatial dispersion metrics to quantify both the degree of structure and spatial

435    organization of gene expression within a nearest-neighbor graph. As such, these metrics can only

436    be applied to dimensionality reduction methods that construct a graph representation of the dataset.

437    Consequently, we cannot use these metrics to evaluate SC3.

438         Although it does technically produce a graph representation of the data, SIMLR should be

439    considered as a hybrid between a clustering and dimensionality reduction method. Because its

440    similarity graph is assumed to have a block structure where the number of blocks is equal to the

441    prespecified number of clusters, the resulting nearest-neighbor graph will, by construction, tend to

442    have a higher degree of structure and therefore artificially inflated NACC and modularity.

443         Furthermore, the poor scalability of SC3 and SIMLR makes them difficult to run for many

444    trials across a large number of datasets. Although SIMLR, in particular, does provide an alternative

445    algorithm that can scale to much larger datasets, it has not been extensively benchmarked. Even

446    so, despite the improved speed of this large-scale implementation, estimating the number of

447    clusters using its built-in function remains a significant computational and memory bottleneck. For

448    example, when applied to the ~10,000 planarian neoblasts, neither implementations of SIMLR

449    could estimate the number of clusters within 2 hours.  As a result, we cannot run SIMLR in an

450    unsupervised manner on datasets significantly larger than ~3000 cells.

24

451   As there are few practical alternatives for manifold reconstruction that have been as

452   extensively benchmarked and widely used, we primarily compare SAM to Seurat in tests involving

453   the unsupervised, graph-based metrics to highlight the key, advantageous characteristics of SAM

454   as a manifold reconstruction and feature selection algorithm when applied to datasets with varying

455   sensitivities (**Figure 4a-c**).

456

457   **Benchmarking.** To generate the convergence curves in **Figure 1b**, we computed the root mean

458   square error (RMSE) of the gene weights averaged across all pairwise comparisons of ten

459   replicates starting from randomly generated initial graphs. In **Figure 3b**, we extend this analysis

460   to all datasets analyzed and report the final error. We use randomly generated datasets of varying

461   sizes (ranging from 200 to 5000 cells) as a negative control to show that SAM does not converge

462   onto the same solution across initial conditions when the data has no intrinsic structure. These

463   datasets were randomly generated by sampling gene expressions from a Poisson distribution with

464   mean drawn from a gamma distribution. To generate the convergence curves in **Figure 3-Figure**

465   **supplement 1a**, we computed the RMSEs, which are ensemble-averaged across ten replicate runs,

466   between the gene weights in adjacent iterations. We compute the adjacency error between kNN

467   adjacency matrices $N_i$ and $N_j$ (**Figure 1b**) as

$$A_{i,j} = \frac{e^T |N_i - N_j| e}{2 e^T N_i e} \tag{14}$$

468   where $e$ is a column vector of ones. This simply measures the fraction of total edges that are

469   different between the two graphs.

470   To compute the standardized dispersion factors in **Figure 2g**, we used Seurat's

471   methodology implemented in Scanpy (Wolf et al., 2018), which groups the genes into 20 bins

472   based on their mean expression values and computes the z-score of each gene's Fano factor with

25

473    respect to the mean and standard deviation of all Fano factors in its corresponding bin. To generate

474    the AUROC scores in **Figure 2h**, which quantify the likelihood of genes being cluster-specific

475    markers, we ran SC3 on the schistosome data with the number of clusters ranging from 2 to 12.

476    We used the AUROC scores corresponding to 4 clusters for the points on the scatter plot and the

477    standard deviations of the scores across all tested numbers of clusters for the error bars.

478         We evaluated each analysis method on nine gold standard datasets (**Figure 3a**) using the

479    Adjusted Rand Index (ARI), which measures the accuracy between two cluster assignments $X$ and

480    $Y$ while accounting for randomness in the clustering:

$$ARI = \frac{\sum \binom{n_{ij}}{2} - \left[\sum \binom{a_i}{2} \sum \binom{b_j}{2}\right]/\binom{n}{2}}{\frac{1}{2}\left[\sum \binom{a_i}{2} + \sum \binom{b_j}{2}\right] - \left[\sum \binom{a_i}{2} \sum \binom{b_j}{2}\right]/\binom{n}{2}} \tag{15}$$

481    where $n$ is the number of cells, and $n_{ij}$, $a_i$, and $b_j$ are elements from a contingency table that

482    summarizes the overlap between the assignments $X$ and $Y$ (Hubert and Arabie, 1985). $n_{ij}$ denotes

483    the number of cells assigned to $X_i$ that are also assigned to $Y_j$, while $a_i$ and $b_j$ are the sums of the

484    $i$th row $j$th column of the contingency table, respectively.

485         Seurat was implemented using the Scanpy package in Python (Wolf et al., 2018). For

486    Seurat, we selected the top 3000 variable genes according to their standardized dispersions and

487    chose the number of PCs (bounded between 6 and 50) which explain 30% of the variance for

488    dimensionality reduction. From these PCs, we calculated a cell-cell correlation distance matrix.

489    To keep the comparison between SAM and Seurat graphs consistent, this distance matrix was

490    converted into a kNN adjacency matrix with the value of $k$ used by SAM. To assign cluster labels

491    for SAM and Seurat, we applied HDBSCAN (Mcinnes et al., 2017), an unsupervised, density-

492    based clustering algorithm to their respective PCA outputs. As HDBSCAN does not cluster any

493    cell it deems an outlier, we assign the remaining outlier cells to clusters using kNN classification.

494    For each outlier cell, we identify its 20 nearest neighbors among the clustered cells. Outliers are

495    assigned to the same cluster as that of the majority of its neighbors. This minor extension to

496    HDBSCAN is available as the built-in function *hdbknn_clustering* in SAM. SC3 was run using

497    default parameters. The SIMLR package was implemented in R and run with the normalization

498    parameter set to "True", which mean-centers gene expressions after normalizing them to be

499    between 0 and 1. Both SC3 and SIMLR provide their own functions to estimate the number of

500    clusters and cluster assignments.

501        To compare the quality of graphs generated by different methods, we use the NACC,

502    modularity, and spatial dispersion. The NACC is the average of the local clustering coefficient for

503    each node of a graph and quantifies the degree of structure in the graph (Watts and Strogatz, 1998).

504    The local clustering coefficient is defined as

$$a_i = \frac{L_i}{k_i(k_i - 1)} \tag{16}$$

505    where $L_i$ is the number of edges between the $k_i$ neighbors of node $i$ and measures the degree of

506    connectedness in a particular node's local neighborhood. We calculate the NACC using the

507    implementation in *graph-tool* (Peixoto, 2017).

508        The modularity $Q$ of a graph is defined as

$$Q = \frac{1}{4m} \sum_{i,j}^{c} \left( A_{ij} - \frac{k_i k_j}{2m} \delta_{ij} \right) \tag{17}$$

509    where $A_{ij}$ is one if there is an edge from cell $i$ to cell $j$, $k_i$ is the degree of cell $i$, $k_j$ is the degree of

510    cell $j$, $m$ is the total number of edges, and $\delta_{ij}$ is one if cells $i$ and $j$ are in the same cluster. High

511    modularity indicates that clusters have on average many more edges within clusters than between

512    clusters. To find the optimal modularity for a particular graph, we used Louvain clustering, which

513    searches for a partition in which modularity is maximized.

27

514   To quantify the spatial organization of gene expression along the graph, we calculate the

515   Euclidean norm of the largest 500 spatial dispersions. Spatial dispersion is defined as before in the

516   SAM algorithm: $F_i = \dfrac{\sigma^2_{C_i}}{\mu_{C_i}}$, where $F_i$ is the Fano factor of the kNN-averaged expressions, $C_i =$

517   $\dfrac{1}{k}NE_i$. $N$ is the directed adjacency matrix output by SAM or Seurat and $E_i$ is a column vector of

518   expression values for gene $i$.

519   To measure the inherent sensitivity of each dataset, we randomly perturbed the gene

520   expression matrices of each dataset by randomly sampling 2000 genes and applied PCA to the

521   subsampled data. A correlation distance matrix was calculated from the top 15 PCs and

522   perturbations were repeated 20 times to generate distance matrix replicates. Sensitivity is then

523   defined as the average error across all pairwise comparisons between replicates. The error between

524   two distance matrices $j$ and $k$, $S_{jk}$, is defined as the average correlation distance between

525   corresponding pairs of rows in the distance matrices $d_j$ and $d_k$:

$$S_{jk} = \frac{1}{n} \sum_{i=1}^{n} D\{d_{j,i}, d_{k,i}\} \tag{18}$$

526   where $D\{d_{j,i}, d_{k,i}\}$ is the Pearson correlation distance between the distances from cell $i$ in distance

527   matrices $j$ and $k$.

528   We simulated datasets with increasing sensitivity by introducing increasing degrees of

529   corruption in each of the nine annotated datasets. To corrupt a dataset, we swapped random pairs

530   of elements in the expression matrix. The number of swaps, $p$, corresponds to the degree of

531   corruption, with $p$ varying from 0 to half of the total number of elements in the matrix. For each

532   annotated dataset, we simulated 10 replicates per value of $p$. SAM and Seurat were run with default

533   values on each corrupted dataset, clustering was performed using the *hdbknn_clustering* function

534   in SAM, and the ARI, NACC, modularity, and spatial dispersion metrics were recorded. The Area

28

535 Under the Curve (AUC) was calculated for each metric with respect to the fraction of elements

536 swapped, $p$, using the trapezoidal rule. Finally, to rescue the performance of Seurat, we used as

537 input to Seurat the top 3000 genes with the highest SAM weights.

538

539 **Gene set enrichment analysis (GSEA).** GSEA (Subramanian et al., 2005) is typically run on a

540 gene expression matrix with user-defined cluster assignments to quantify the differential

541 expression for each gene. By default, differential expression is quantified using a signal-to-noise

542 metric and the resulting scores are used to rank the genes in descending order.  However, GSEA

543 can also run in an alternative mode in which the user provides a predefined list of gene rankings.

544 Therefore, we used the genes ranked by their SAM weights as input to GSEA to determine the

545 biological processes enriched among the highly weighted genes. As shown in **Figure 6a**, we can

546 directly test if SAM captures the relevant biological processes. GSEA provides a number of

547 statistical measures to assess the significance of enriched gene sets, of which we use the False

548 Discovery Rate (FDR). The FDR quantifies the likelihood that a highly enriched gene set

549 represents a false positive. The significance threshold typically used with FDR is 25%, which

550 implies that the results are likely to be valid 75% of the time.

551

552 **Removal of cell cycle effects.** To remove cell cycle effects from the macrophage dataset, we

553 adopted a simpler version of the strategy used in *ccRemover* (Barron and Li, 2016), in which we

554 subtract from the data PCs that are significantly associated with known cell cycle genes. Letting $P$

555 represent the PCs and $L$ be the gene loadings, we quantify the association between the set of cell

556 cycle genes $G$ and PC $j$ as

$$A_j = \frac{1}{|G|} \sum_{i \in G} |L_{ji}| \qquad (19)$$

557    PC $j$ is selected if its association $A_j$ is at least two standard deviations above the mean of the

558    associations for all PCs. In the particular case of the macrophage data, we identified the set of PCs

559    $S = \{P_0, P_1, P_8\}$ as being significantly associated with the cell cycle genes. We next reconstruct the

560    data using these PCs, which thus captures the cell-cycle effects, and subtract the reconstructed data

561    from the expression matrix $E$:

$$E_{removed} = E - \sum_{j \in S} P_j L_j \sqrt{W} \qquad (20)$$

562    When reconstructing the data, we scale the gene loadings by the SAM weights $W$ so that only the

563    highly weighted SAM genes (which are initially dominated by cell cycle genes) contribute to the

564    cell cycle removal, as there may be genes involved in other biological processes that could also be

565    correlated with the PCs in $S$. To run SAM on the data with cell cycle effects removed, we use $E$

566    as opposed to $E_{removed}$ for the calculation of spatial dispersions, because the latter may contain

567    negative values, for which dispersion is ill-defined. This method is made available as a part of the

568    SAM package in the functions *calculate_regression_PCs* and *regress_genes*.

569

570    **Clustering the NF-κB activity time series.** In the original study,  Lane et al. combined imaging

571    and transcriptomics to link NF-κB nuclear translocation dynamics to changes in gene expression

572    within single cells. Macrophages stimulated with LPS were individually trapped in microfluidic

573    chambers and imaged for various lengths of time (75-300 min) prior to scRNAseq library

574    preparation. NF-κB was tagged with a fluorescent protein, and its activation was measured as the

575    nuclear-localized fluorescence intensity. Based on the imaging data, the authors identified three

576    main classes of NF-κB dynamics, the first with a transient initial response, the second with a

577    prolonged initial response, and the third with a recurrent response. Because the recurrent response

578    is found only in the 300 min time point and comprises only ~8% of these cells, we primarily

579    focused on clustering cells based on their initial dynamics. To do this, we used the *tslearn*

580    (Tavenard, 2017) python package to group cells based on their NF-κB activity time series. Because

581    these time series are quite noisy, we were conservative in labeling cells as having a prolonged

582    initial response in an effort to avoid false positives. As a result, these cells comprise only ~30% of

583    the dataset.

584         For the cells sampled at 75 and 150 min after LPS stimulation, we used the time series *k*-

585    means algorithm with the *softdtw* distance metric to cluster them each into three groups, which

586    resulted in representative time series with transient, intermediate, and prolonged responses.

587    Merging the cells with transient and intermediate responses into one cluster (which we simply

588    labeled as transient response), we obtained the 75 and 150 min (black and blue, respectively)

589    representative time series shown in **Figure 6e**. Because the cells sampled at 300 min displayed

590    much more variability in their NF-κB time series, we clustered them into 6 groups, labeling the

591    cluster whose representative time series had the broadest initial peak as the prolonged response

592    cluster (red in **Figure 6e**, right). The remaining groups were labeled as the transient response

593    cluster (red in **Figure 6e**, left).

594

595    **Mapping the schistosome datasets.** We used the Mutual Nearest Neighbors algorithm

596    (Haghverdi et al., 2018) with default values to generate an expression matrix $E_{corrected}$ in which

597    the batch effects between the 2.5-week and 3.5-week datasets were corrected for. To run SAM on

598    the batch-corrected data, we use $E$ for the calculation of spatial dispersions as opposed to

599    $E_{corrected}$.

600

601    **scRNAseq of schistosome stem cells.** Schistosome stem cells were isolated from juvenile

602    parasites retrieved from infected mice at 2.5 and 3.5 weeks post infection. We followed the

603    protocol as previously described (Wang et al., 2018). Briefly, we retrieved juvenile parasites from

604    schistosome-infected mice (Swiss Webster NR-21963) by hepatic portal vein perfusion. Parasites

605    were cultured at 37°C/5% $CO_2$ in Basch Medium 169 supplemented with 1X Antibiotic-

606    Antimycotic for 24-48 hr to allow complete digestions of host blood cell in parasite intestines. In

607    adherence to the Animal Welfare Act and the Public Health Service Policy on Humane Care and

608    Use of Laboratory Animals, all experiments with and care of mice were performed in accordance

609    with protocols approved by the Institutional Animal Care and Use Committees (IACUC) of

610    Stanford University (protocol approval number 30366).

611        Before dissociation, parasites were permeabilized in PBS containing 0.1% Triton X-100

612    and 0.1% NP-40 for 30 seconds and washed thoroughly to remove the surfactants. The

613    permeabilized parasites were dissociated in 0.25% trypsin for 20 min. Cell suspensions were

614    passed through a 100 μm nylon mesh (Falcon Cell Strainer) and centrifuged at 150 g for 5 min.

615    Cell pellets were gently resuspended, passed through a 30 μm nylon mesh, and stained with

616    Vybrant DyeCycle Violet (DCV; 5 μM, Invitrogen), and TOTO-3 (0.2 μM, Invitrogen) for 30–45

617    min. As the stem cells comprise the only proliferative population in schistosomes, we flow-sorted

618    cells at $G_2$/M phase of the cell cycle on a SONY SH800 cell sorter. Dead cells were excluded based

619    on TOTO-3 fluorescence. Single stem cells were gated using forward scattering (FSC), side

620    scattering (SSC), and DCV to isolate cells with doubled DNA content compared to the rest of the

621    population (Wang et al., 2018). Cells that passed these gates were sorted into 384-well lysis plates

622    containing Triton X-100, ERCC standards, oligo-dT, dNTP, and RNase inhibitor.

623     cDNA was reverse transcribed and amplified on 384-well plate following the Smart-Seq2

624     protocol (Picelli et al., 2013). For quality control, we quantified the histone *h2a* (Smp_086860)

625     levels using qPCR, as *h2a* is a ubiquitously expressed in all schistosomes stem cell (Collins et al.,

626     2013; Wang et al., 2013; Wang et al., 2018). We picked wells that generated $C_T$ values within 2.5

627     $C_T$ around the most probable values (~45% of total wells, **Figure 1-Figure supplement 1**). cDNA

628     was then diluted to 0.4 ng/µL for library preparation. Tagmentation and barcoding of wells were

629     prepared using Nextera XT DNA library preparation kit. Library fragments concentration and

630     purity were quantified by Agilent bioanalyzer and qPCR. Sequencing was performed on a NextSeq

631     500 using V2 150 cycles high-output kit at ~1 million reads depth per cell. Raw sequencing reads

632     were demultiplexed and converted to fastq files using bcl2fastq. Paired-end reads were mapped to

633     *S. mansoni* genome version WBPS9 (WormBase Parasite) using STAR. In 2.5-week dataset, 338

634     cells with more than 1700 transcripts expressed at >2 TPM were used for downstream analysis. In

635     the 3.5-weeks dataset, 338 cells with more than 1350 transcripts expressed at >2 TPM were used

636     for downstream analysis (**Figure 1-Figure supplement 1**).

637

638     **In situ hybridization and EdU labeling.** RNA FISH experiments were performed as detailed in

639     previous publications (Collins et al., 2013; Wang et al., 2013; Wang et al., 2018). Clones used

640     for riboprobe synthesis were generated as described previously, with oligonucleotide primers

641     listed in **Supplementary Table 3**. Juvenile parasites were cultured with 10 µM EdU overnight,

642     killed in 6 M $MgCl_2$ for 30s, and then fixed in 4% formaldehyde with 0.2% Triton X-100 and 1%

643     NP-40. Fixed parasites were sequentially dehydrated in methanol, bleached in 3% $H_2O_2$ for 30

644     min, and rehydrated. Parasites were permeabilized by 10 µg/mL proteinase K for 15 min and

645     post fixed with 4% formaldehyde. The hybridization was performed at 52°C with riboprobes

646     labeled with either digoxigenin-12-UTP (Roche) or fluorescein-12-UTP (Roche). For detection,

647     samples were blocked with 5% horse serum and 0.5% of Roche Western Blocking Reagent, and

648     then incubated with anti-digoxigenin-peroxidase (1:1000; Roche) or anti-fluorescein peroxidase

649     (1:1500; Roche) overnight at 4°C for tyramide signal amplification (TSA). For double FISH, the

650     first peroxidase was quenched for 30 min in 0.1% sodium azide solution before the detection of

651     the second gene. After FISH, EdU detection was performed by click reaction with 25μM Cy5-

652     azide conjugates (Click Chemistry Tools). Samples were mounted in *scale* solution (30%

653     glycerol, 0.1% Triton X-100, 4 M urea in PBS supplemented with 2 mg/mL sodium ascorbate)

654     and imaged on a Zeiss LSM 800 confocal microscope.

661    **References**

662    Bahlo, M., Tian, L., Lönnstedt, I., Ng, M. & Hicks, S. Comparison of clustering tools in R for
663    medium-sized 10x Genomics single-cell RNA-sequencing. *F1000Research* **7,** 1–26 (2018).
664
665    Baron, M. *et al*. A single-cell transcriptomic map of the human and mouse pancreas reveals inter-
666    and intra-cell population structure. *Cell syst.* **3,** 346–360 (2016).
667
668    Barron, M. & Li, J. Identifying and removing the cell- cycle effect from single-cell RNA-
669    Sequencing data. *Sci. Rep.* **6,** 33892 (2016).
670
671    Becht, E. *et al.* Dimensionality reduction for visualizing single-cell data using UMAP. *Nat.*
672    *Biotechnol.* **37,** 38–44 (2019).
673
674    Blondel, V. D., Guillaume, J., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities in
675    large networks. *J. Stat. Mech. Theory Exp.* **2008,** P10008 (2008).
676
677    Collins, J. J. *et al.* Adult somatic stem cells in the human parasite Schistosoma mansoni. *Nature*
678    **494,** 476–479 (2013).
679
680    Crow, M. *et al*. Characterizing the replicability of cell types defined by single cell RNA-
681    sequencing data using MetaNeighbor. *Nat. Commun.* **9,** 884 (2018).
682
683    Darmanis, S. *et al.* A survey of human brain transcriptome diversity at the single cell level. *Proc.*
684    *Natl. Acad. Sci. U.S.A.* **112,** 7285–7290 (2015).
685
686    Duò, A., Robinson, M. D. & Soneson, C. A systematic performance evaluation of clustering
687    methods for single-cell RNA-seq data. *F1000Research* **7,** 1–21 (2019).
688
689    Fincher, C. T. *et al*. Cell type transcriptome atlas for the planarian Schmidtea mediterranea.
690    *Science* **360,** eaaq1736 (2018).
691
692    Guo, M. *et al*. SINCERA: A Pipeline for Single-Cell RNA-Seq Profiling Analysis. *PLOS*
693    *Comput. Biol.* **11,** e1004575 (2015).
694
695    Grün, D., Kester, L. & Oudenaarden, A. V. Validation of noise models for single-cell
696    transcriptomics. *Nat. Methods* **11,** 637–640 (2014).
697
698    Haghverdi, L. *et al*. Batch effects in single-cell RNA-sequencing data are corrected by matching
699    mutual nearest neighbors. *Nat. Biotechnol.* **36,** 421–427 (2018).
700
701    Hoffmann, K. F., Brindley, P. J. & Berriman, M. Halting harmful helminths. *Nature* **168,** 168–
702    169 (2014).
703
704    Hubert, L. & Arabie, P. Comparing partitions. *J. Classif.* **2,** 193–218 (1985).
705

706    Kiselev, V. Y. *et al.* SC3: Consensus clustering of single-cell RNA-seq data. *Nat. Methods* **14,**
707    483–486 (2017).
708
709    Lane, K. *et al.* Measuring signaling and RNA-seq in the same cell links gene expression to
710    dynamic patterns of NF-κB activation. *Cell Syst.* **4,** 458–469 (2017).
711
712    Lee, T. K. *et al.* A Noisy Paracrine Signal Determines the Cellular NF- k B Response to
713    Lipopolysaccharide. *Sci. Immunol.* **2,** ra65 (2009).
714
715    Li, H. *et al.* Reference component analysis of single-cell transcriptomes elucidates cellular
716    heterogeneity in human colorectal tumors. *Nat. Genet.* **49,** 708–718 (2017).
717
718    Lönnberg, T. *et al.* Single-cell RNA-seq and computational analysis using temporal mixture
719    modeling resolves TH1/TFH fate bifurcation in malaria. *Sci. Immunol.* **2,** eaal2192 (2017).
720
721    Mcinnes, L., Healy, J. & Astels, S. hdbscan: Hierarchical density based clustering. *J. Open*
722    *Source Softw.* **2,** 205 (2017).
723
724    Olsson, A. *et al.* Single-cell analysis of mixed-lineage states leading to a binary cell fate choice.
725    *Nature* **537,** 698–702 (2016).
726
727    Peixoto, T. P. The graph-tool python library. (2017).
728    doi:https://doi.org/10.6084/m9.figshare.1164194.v14
729
730    Picelli, S. *et al.* Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat.*
731    *Methods* **10,** 1096–1100 (2013).
732
733    Pierson, E. & Yau, C. ZIFA: Dimensionality reduction for zero-inflated single-cell gene
734    expression analysis. *Genome Biol.* **16,** 241 (2015).
735
736    Qiu, X. *et al.* Reversed graph embedding resolves complex single-cell trajectories. *Nat. Methods*
737    **14,** 979–982 (2017).
738
739    Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-
740    cell gene expression data. *Nat. Biotechnol.* **33,** 495–502 (2015).
741
742    Schwalie, P. C. *et al.* A stromal cell population that inhibits adipogenesis in mammalian fat
743    depots. *Nature* **559,** 103–108 (2018).
744
745    Setty, M. *et al.* Wishbone identifies bifurcating developmental trajectories from single-cell data.
746    *Nat. Biotechnol.* **34,** 637–645 (2016).
747
748    Soneson, C. & Robinson, M. D. Bias, robustness and scalability in single-cell differential
749    expression analysis. *Nat. Methods* **15,** 255–261 (2018).
750

751    Subramanian, A. *et al*. Gene set enrichment analysis: A knowledge-based approach for
752    interpreting genome-wide. *Proc. Natl. Acad. Sci. U.S.A.* **102,** 15545–15550 (2005).
753
754    Tavenard, R. tslearn: A machine learning toolkit dedicated to time-series data. (2017).
755
756    Tian, L. *et al*. Benchmarking single cell RNA-sequencing analysis pipelines using mixture
757    control experiments. *Nat. Methods* **16,** 479–487 (2019).
758
759    Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by
760    pseudotemporal ordering of single cells. *Nat. Biotechnol.* **32,** 381–386 (2014).
761
762    Treutlein, B. *et al.* Reconstructing lineage hierarchies of the distal lung epithelium using single
763    cell RNA-seq. *Nature* **509,** 371–375 (2014).
764
765    Vallejos, C. A. *et al*. BASiCS: Bayesian Analysis of Single-Cell Sequencing Data. *PLOS*
766    *Comput. Biol.* **11.6,** e1004333 (2015).
767
768    Vallejos, C. A. *et al*. Normalizing single-cell RNA sequencing data: challenges and
769    opportunities. *Nat. Methods* **14,** 565–571 (2017).
770
771    Wang, B., Collins, J. J. & Newmark, P. A. Functional genomic characterization of neoblast-like
772    stem cells in larval Schistosoma mansoni. *eLife* **2,** e00768 (2013).
773
774    Wang, B., Zhu, J., Pierson, E., Ramazzotti, D. & Batzoglou, S. Visualization and analysis of
775    single-cell rna-seq data by kernel-based similarity learning. *Nat. Methods* **14,** 414–416 (2017).
776
777    Wang, B. *et al.* Stem cell heterogeneity drives the parasitic life cycle of Schistosoma mansoni.
778    *eLife* **7,** e35449 (2018).
779
780    Watts, D. J. & Strogatz, S. H. Collective dynamics of 'small-world' networks. *Nature* **393,** 440–
781    442 (1998).
782
783    Wendt, G.R. *et al.* Flatworm-specific transcriptional regulators promote the specification of
784    tegumental progenitors in Schistosoma mansoni. *eLife* **7,** e33221 (2018).
785
786    Witchley, J. N., Mayer, M., Wagner, D. E., Owen, J. H. & Reddien, P. W. Muscle cells provide
787    instructions for planarian regeneration. *Cell Rep.* **4,** 633–641 (2013).
788
789    Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: Large-scale single-cell gene expression data
790    analysis. *Genome Biol.* **19,** 1–5 (2018).
791
792    WormBase, https://www.wormbase.org, release WS268.
793
794    Xu, C. & Su, Z. Gene expression Identification of cell types from single-cell transcriptomes
795    using a novel clustering method. *Bioinformatics* **31,** 1974–1980 (2015).
796

797    Žurauskien, J. & Yau, C. pcaReduce: hierarchical clustering of single cell transcriptional
798    profiles. *BMC Bioinform.* **17,** 140 (2016).

799    **Figure legends**

800    **Figure 1. The SAM algorithm.** (**a**) SAM starts with a randomly initialized kNN adjacency matrix

801    and iterates to refine the adjacency matrix and gene weight vector until convergence. (**b**) Root

802    mean square error (RMSE) of the gene weights (top) and the fraction of different edges of the

803    nearest-neighbor adjacency matrices (bottom) between adjacent iterations (blue) and between

804    independent runs at the same iteration (orange) to show that SAM converges to the same solution

805    regardless of initial conditions. The differences between the gene weights and nearest-neighbor

806    graphs from independent runs are relatively small, indicating that SAM converges to the same

807    solution through similar paths. (**c**) Graph structures and gene weights of the schistosome stem cell

808    data converging to the final output over the course of 10 iterations (*i* denotes iteration number).

809    Top: nodes are cells and edges connect neighbors. Nodes are color-coded according to the final

810    clusters. Bottom: weights are sorted according to the final gene rankings. (**d**) Network properties

811    iteratively improve for the graphs reconstructed from the original data (red) but not on the

812    randomly shuffled data (blue). Dashed lines: metrics measured from the Seurat-reconstructed

813    graphs.

814

815    **Figure 2. SAM identifies novel subpopulations within schistosome stem cells.** (**a**) UMAP

816    projections of the manifolds reconstructed by SAM, PCA, and Seurat. SIMLR outputs its own 2D

817    projection based on its constructed similarity matrix using a modified version of t-SNE. The

818    schistosome cells are color-coded by the stem cell subpopulations $\mu$, $\delta$', $\varepsilon_\alpha$, and $\varepsilon_\beta$ determined by

819    Louvain clustering. (**b**) UMAP projections with gene expressions of subpopulation-specific

820    markers (*eledh*, *nanos-2 cabp*, *astf*, *bhlh*,) and a ubiquitous stem cell marker, *ago2-1*, overlaid.

821    Insets: magnified views of the expressing populations. (**c**) FISH of *cabp* and EdU labeling of

40

822     dividing stem cells in juvenile parasites at 2.5 weeks post infection show that μ-cells ($cabp^+EdU^+$,

823     arrowheads) are close to the parasite surface and beneath a layer of post-mitotic $cabp^+$ cells.

824     Dashed outline: parasite surface. Right: magnified views of the boxed region. (**d**) FISH of *cabp*

825     and a set of canonical muscle markers, *troponin*, *myosin*, *tropomyosin*, and *collagen*, shows

826     complete colocalization in post-mitotic $cabp^+$ cells. Images in (**c-d**) are single confocal slices. (**e**)

827     FISH of *astf* and *bhlh* shows their orthogonal expression in adjacent $EdU^+$ cells (arrowheads).

828     Bottom: magnified views of the boxed region. Image is a maximum intensity projection of a

829     confocal stack with a thickness of 12 µm. (**f**) UMAP projection of stem cells isolated from

830     juveniles at 2.5 and 3.5 weeks post infection. Cell subpopulation assignments based on marker

831     gene expressions are specified. Right: a magnified view to show the mapping of $\varepsilon_\alpha$- and $\varepsilon_\beta$-cells.

832     (**g**) Standardized dispersions as calculated by Seurat plotted vs. the SAM gene weights. (**h**) SC3

833     AUROC scores plotted vs. the SAM gene weights. Error bars indicate the standard deviation of

834     SC3 AUROC scores between trials using different chosen numbers of clusters. In (**g**) and (**h**), the

835     top 20 genes specific to each subpopulation are colored according to the color scheme used in (**a**).

836

837     **Figure 3. SAM improves clustering accuracy and runtime performance.** (**a**) Accuracy of

838     cluster assignment quantified by adjusted rand index (ARI) on nine annotated datasets (left). Right:

839     differences between the number of clusters found by each method and the number of annotated

840     clusters. Smaller differences indicate more accurate clustering. (**b**) RMSE of gene weights output

841     by SAM averaged across ten replicate runs with random initial conditions for 56 datasets (blue)

842     and simulated datasets with no intrinsic structure (green, **Methods**). (**c**) Runtime of SAM, SC3,

843     SIMLR, and Seurat as a function of the number of cells in each dataset. SC3 and SIMLR were not

844     run on datasets with >3000 cells as the run time exceeds 20 minutes.

845

846 **Figure 4. SAM improves the analysis of datasets with varying network sensitivities.** (**a**)

847 Network sensitivity of all 56 datasets ranked in descending order. Blue: the nine benchmarking

848 datasets used in **Figure 3a**. Sensitivity measures the robustness of a dataset to changes in the

849 selected features (**Methods**). (**b**) The network sensitivity plotted against the fraction of genes with

850 SAM weight greater than 0.5 (in log scale) with Spearman correlation coefficient specified in the

851 upper-right corner. (**c**) Fold improvement of SAM over Seurat for NACC, modularity, and spatial

852 dispersion with respect to sensitivity for all 56 datasets. These ratios are linearly correlated with

853 network sensitivity with Pearson correlations ($r^2$) specified in the upper-left corner of each plot.

854

855 **Figure 5. Robust feature selection improves cell clustering and manifold reconstruction.** (**a**)

856 Network sensitivity, ARI, NACC, modularity, and spatial dispersion with respect to corruption of

857 the Darmanis dataset, in which we swap random pairs of gene expressions with the number of

858 swaps ranging from 0-50% of the total number of elements in the data (**Methods**). Performance is

859 compared between SAM (blue), Seurat (red), and Seurat rescued with the top-ranked SAM genes

860 (indigo). Error bars indicate the standard deviations across 10 replicate runs. The errors for points

861 with no bars are too small to be seen. (**b**) Comparison of the area under curve (AUC) of the metrics

862 in (**a**) with respect to data corruption for all nine datasets. Error bars indicate the standard

863 deviations across 10 replicate runs. The errors for data with no error bars are too small to be seen.

864

865 **Figure 6. SAM captures the cellular activation dynamics in a stimulated macrophage dataset.**

866 (**a**) GSEA analysis (left) and SAM projections (right) of the activated macrophages[7] before (top)

867 and after (bottom) removing cell cycle effects. Teal: significantly enriched gene sets determined

868    by the significance threshold of 0.25 for the False Discovery Rate (FDR, dashed lines). Bottom:

869    the two clusters are denoted as MT and M with colors representing the time since LPS induction.

870    Arrows: evolution of time. (**b**) TNFα is enriched in the MT cluster. (**c**) Diagram of NF-κB

871    activation in response to LPS stimulation via the Myd88 and TRIF signaling pathways. (**d**) $Log_2$

872    fold changes of the average expressions of selected inflammatory genes in the MT cluster vs. the

873    M cluster. All genes are significantly differentially expressed between the two clusters according

874    to the Welch's two-sample t-test ($p < 5 \cdot 10^{-3}$). (**e**) Representative traces for transient (left) and

875    prolonged (right) NF-κB activation (**Methods**). (**f**) Cells with prolonged NF-κB response (denoted

876    as P) are primarily in the MT population. (**g**) Seurat and SIMLR projections show that they fail to

877    order the cells by time since LPS induction and do not identify cell clusters representing the

878    different modes of NF-κB activation.

879

880    **Figure 1 – Figure supplement 1. Quality control of library preparation and sequencing of**

881    **the schistosome stem cells.** (**a**) Histograms of *h2a* qPCR measurements in 2.5- (left) and 3.5-

882    (right) week samples. (**b**) Scatter plot of gene count (>2 TPM) vs mapped read count of individual

883    sequenced cells. Cells with low gene count or *h2a* expression are discarded and filtered from

884    analysis (red) and the remaining cells are analyzed (blue). The number of final cells kept for

885    analysis is annotated on the top left corner of each plot.

886

887    **Figure 2 – Figure supplement 1. $\mu$-cells express ubiquitous stem cell marker and population**

888    **specific genes.** UMAP projections with gene expressions of (**a**) stem cell markers and (**b**) $\mu$-cell-

889    specific genes overlaid.
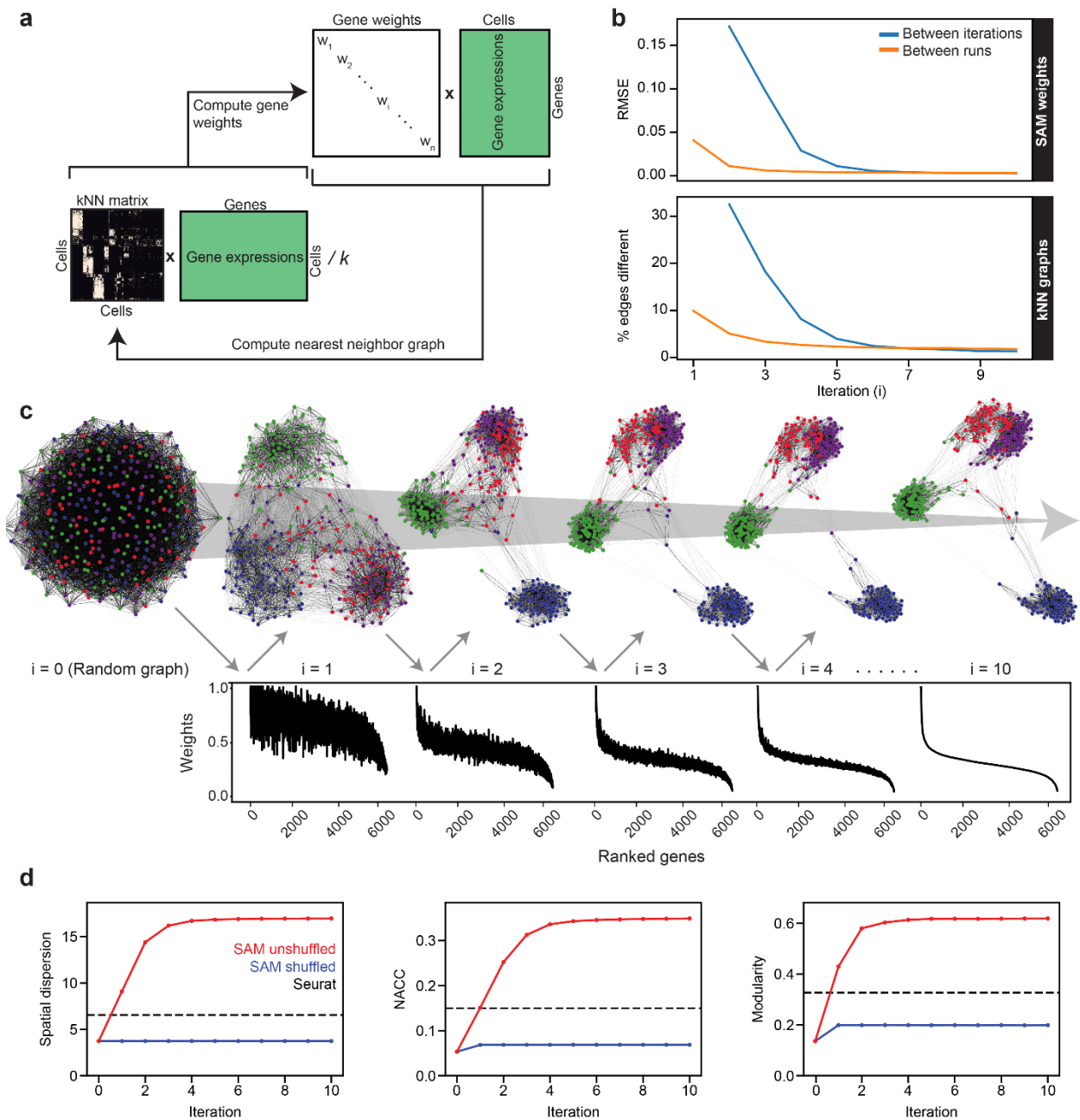
890

891    **Figure 3 – Figure supplement 1. SAM converges to a stable solution independent of random**

892    **initial conditions and is robust to the number of nearest neighbors and choice of distance**

893    **metric.** (**a**) RMSE of gene weights between adjacent iterations within a run, averaged across ten

894    replicate runs. (**b-c**) Average ARI scores for the nine annotated benchmarking datasets when

895    varying (**b**) the number of nearest neighbors, $k$, from 10 to 30 or (**c**) the choice of distance metric

896    (Euclidean or Pearson correlation). Error bars indicate standard deviations of ARI scores across

897    the different values of $k$ and distance metrics. The errors for data with no error bars are too small

898    to be seen.

899

900    **Figure 6 – Figure supplement 1. Cluster-specific marker genes before and after removing**

901    **cell cycle effects.** UMAP projections with marker genes specific to the dividing cells (**a**) and the

902    MT cluster (**b**) overlaid.

903

904    **Figure 6 – Figure supplement 2. SAM groups cells based on NF-κB activation dynamics while**

905    **other methods cannot.** (**a**) UMAP projection of the macrophage cells after the removal of cell

906    cycle effects. Cells with prolonged NF-κB dynamics are highlighted in red. (**b**) UMAP and t-SNE

907    projections for Seurat and SIMLR, respectively, after the removal of cell cycle effects. Cells with

908    prolonged NF-κB dynamics are highlighted in red. (**c**) UMAP projections with three MT-specific

909    marker gene expressions overlaid.

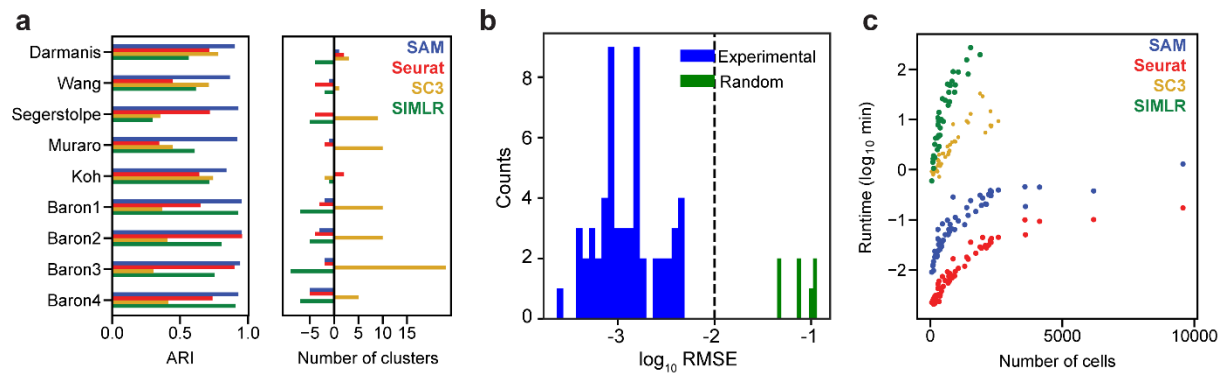910 **Figures**

911 **Figure 1.**



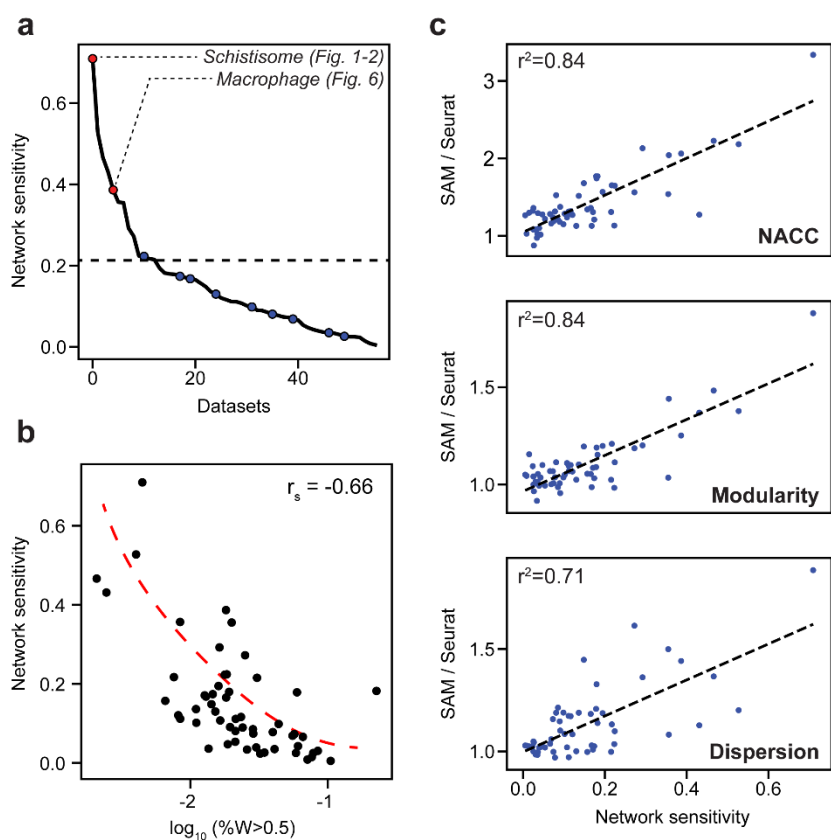912

913 **Figure 2.**



914
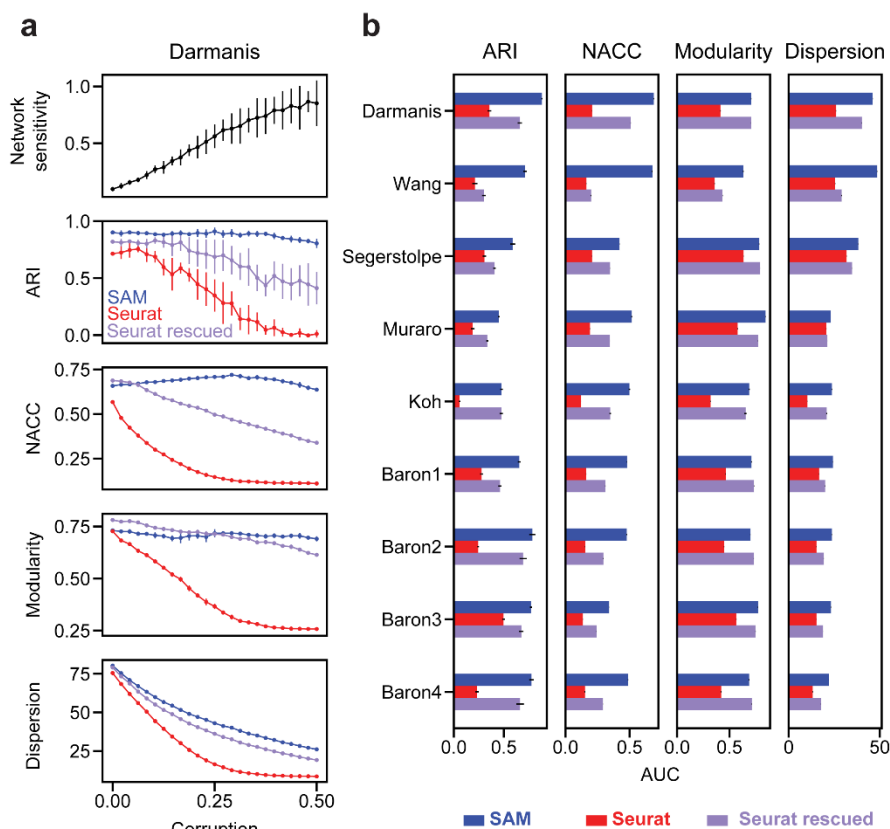
915    **Figure 3.**



916

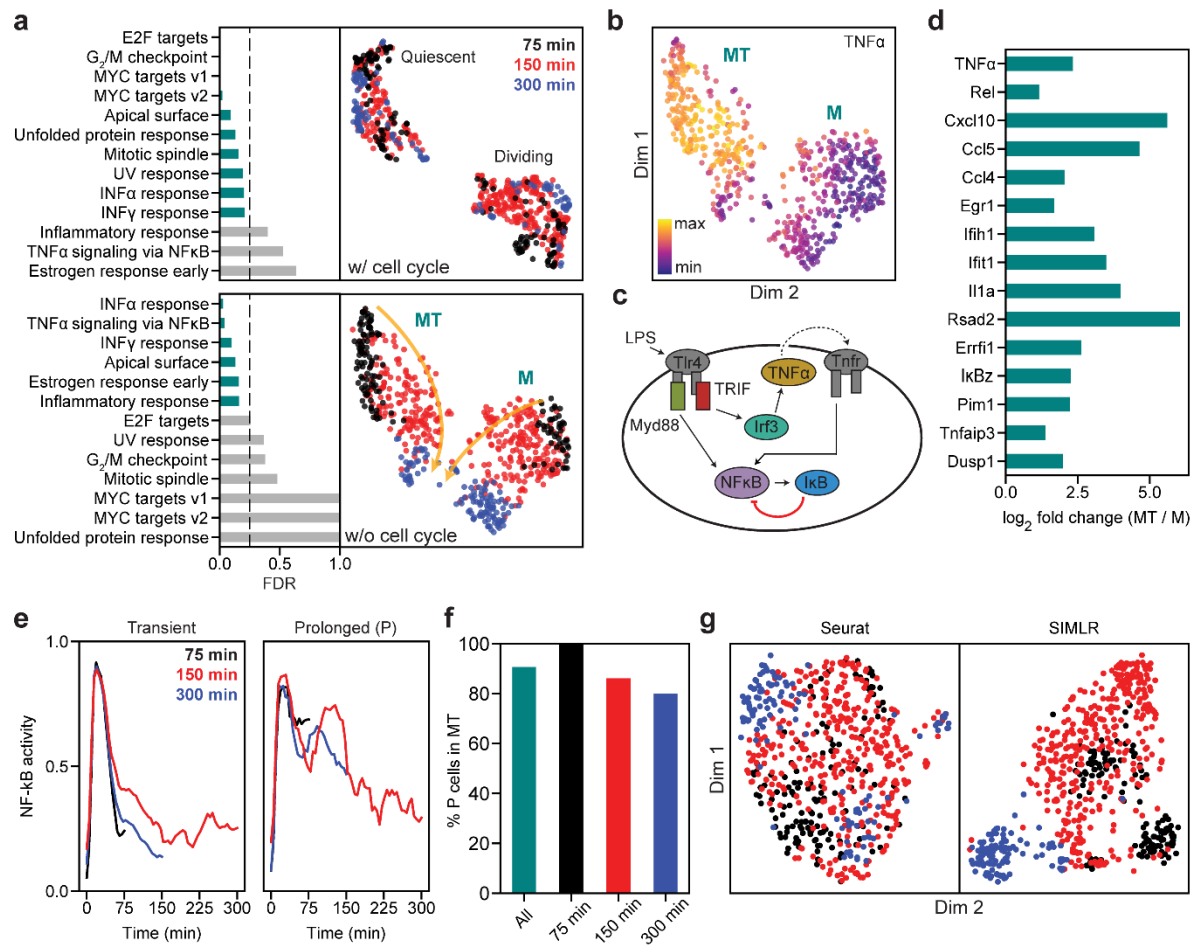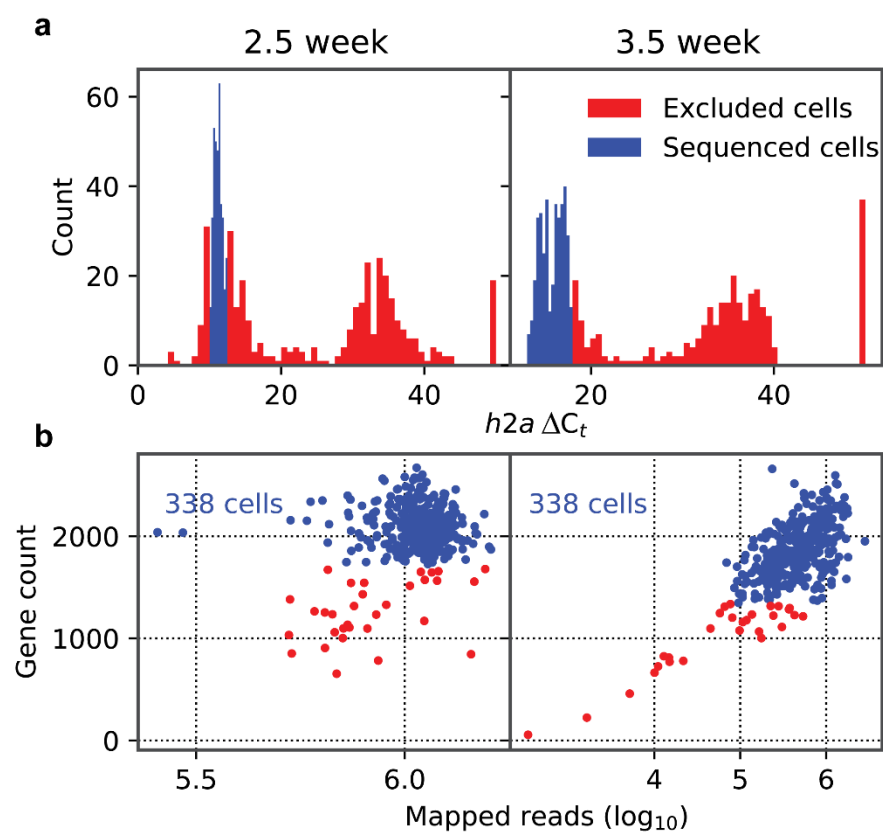917 **Figure 4.**



918

919    **Figure 5.**



920

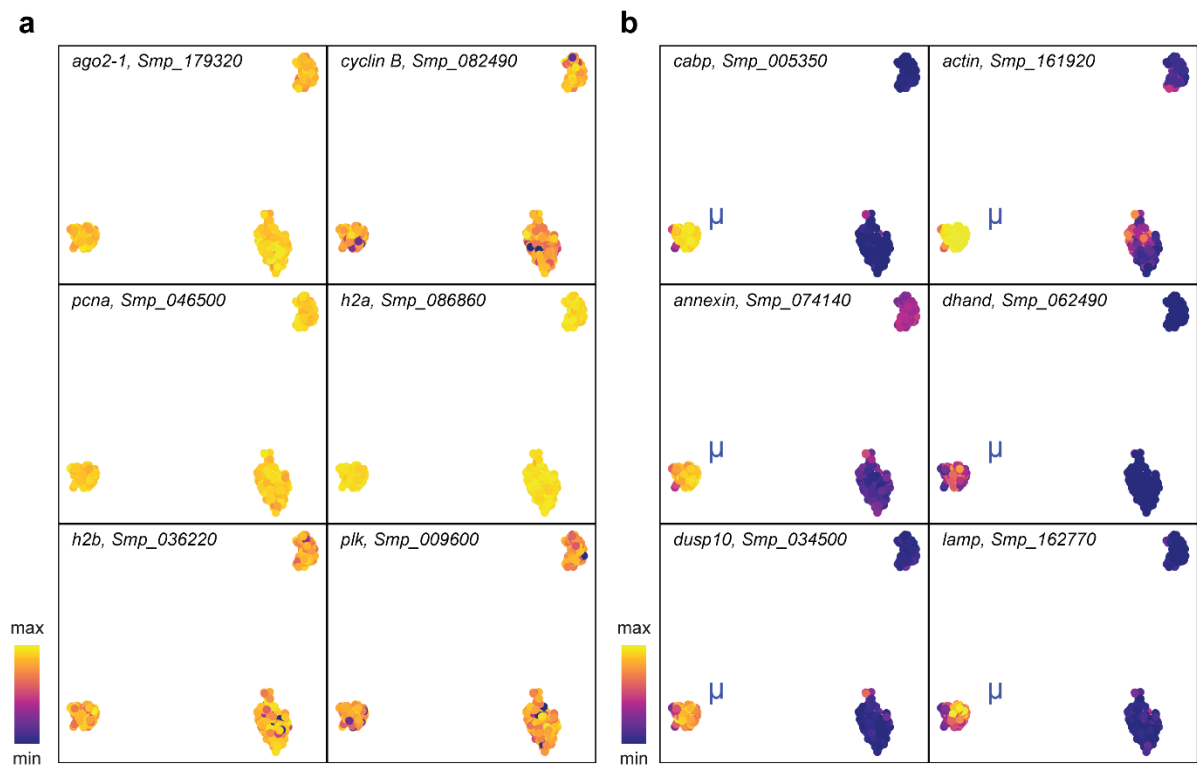921 **Figure 6.**



922

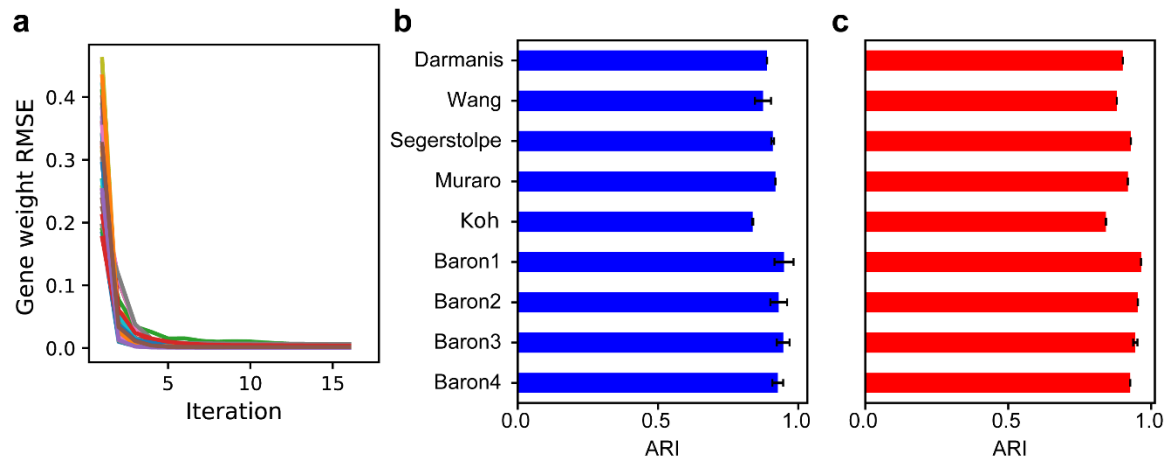923 **Figure 1 – Figure supplement 1.**



924

925 **Figure 2 – Figure supplement 1.**
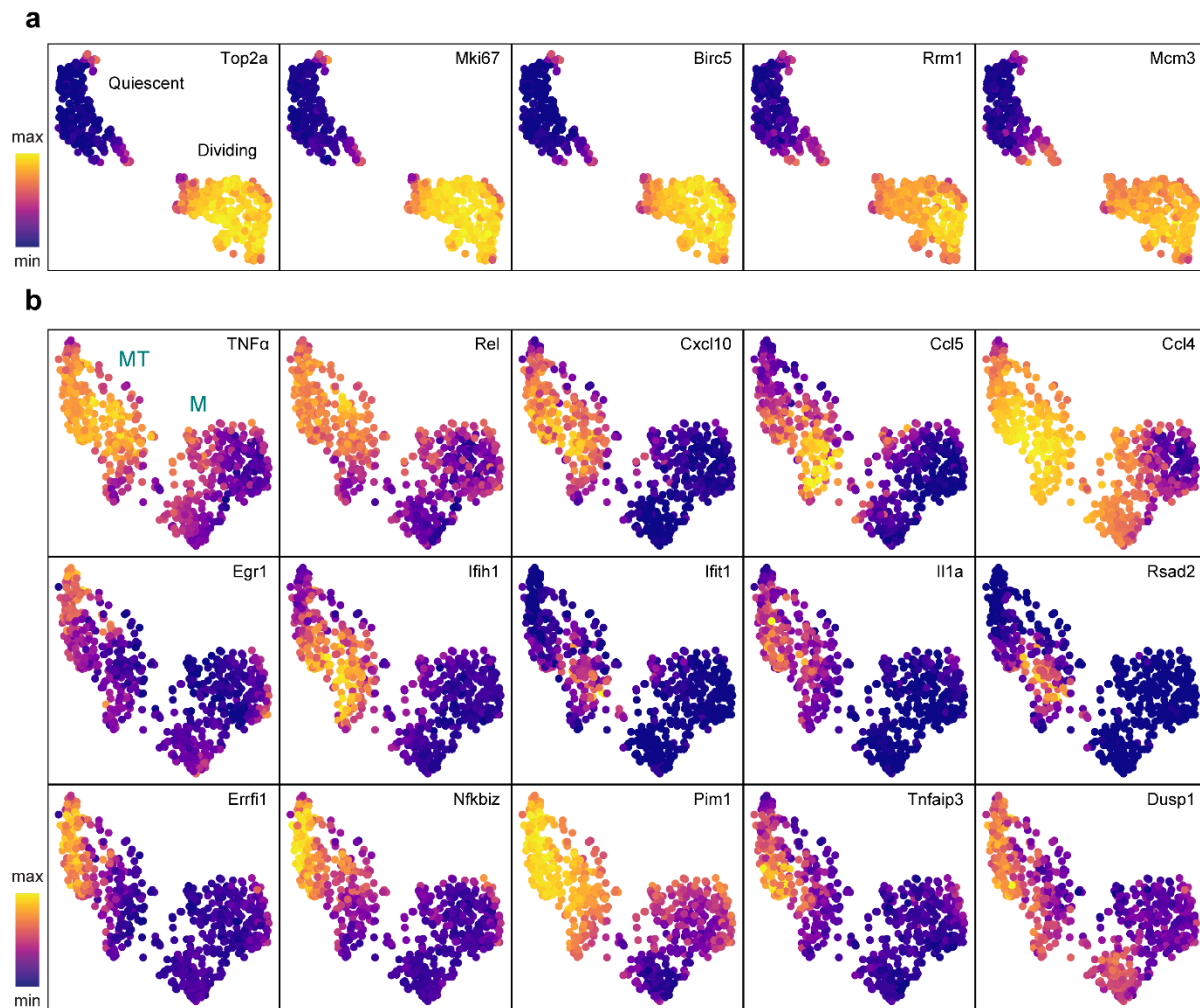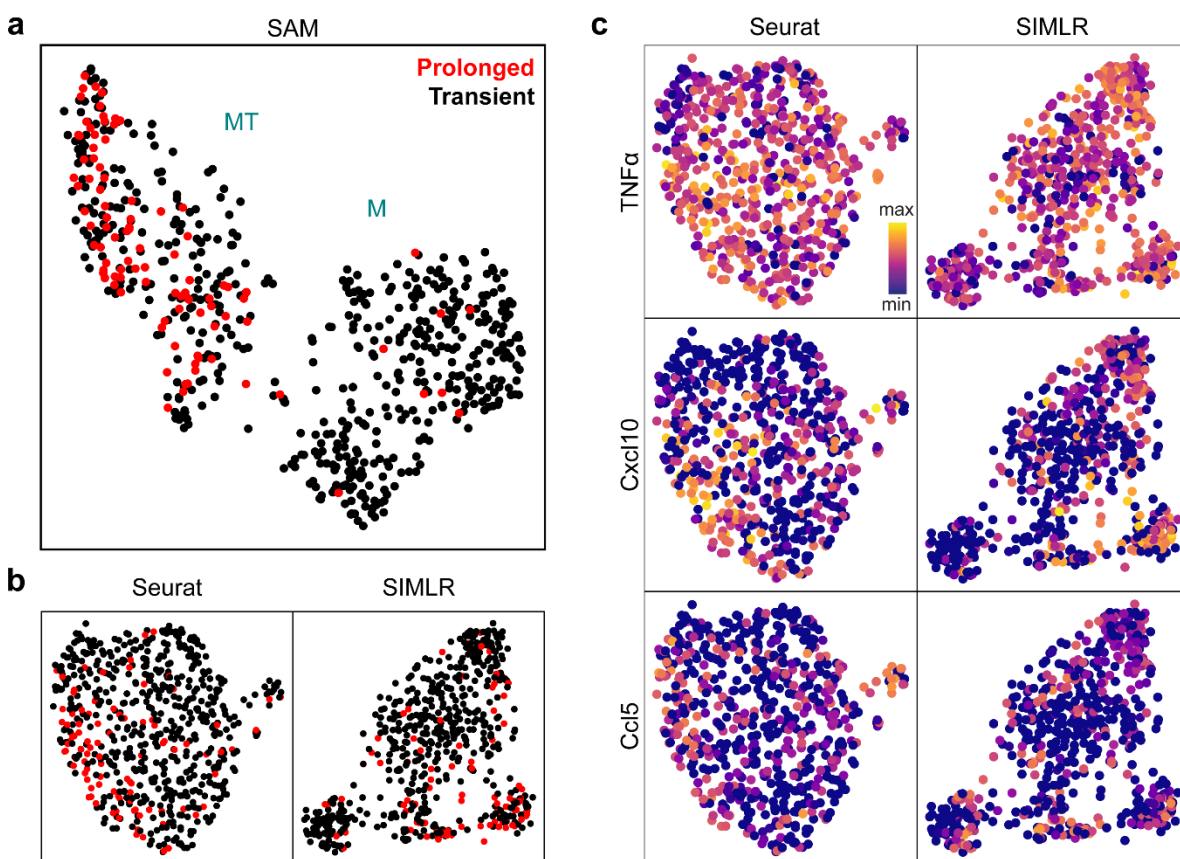


926

927    **Figure 3 – Figure supplement 1.**



928

929 **Figure 6 – Figure supplement 1.**



930

931  **Figure 6 – Figure supplement 2.**



932

933 **Supplementary Table legends:**

934

935 **Supplementary Table 1: Ranked gene list with high SAM weights in the schistosome stem**

936 **cell data**. Gene IDs and annotations are given in the *S. mansoni* genome version 9 (WormBase,

937 WS268). Genes are assigned to the cluster corresponding to the marker gene, *nanos-2*, *cabp*, *astf*,

938 or *bhlh*, that they have the highest correlation with. Genes found in our prior work[12] to be enriched

939 in subsets of stem cells are specified.

940

941 **Supplementary Table 2: Datasets used in this study.** Accession numbers, library size

942 normalization methods, data preprocessing methods, sensitivity scores, and corresponding

943 references are provided for each dataset. Accession numbers with asterisks indicate datasets that

944 are sourced from the *conquer* database (Soneson and Robinson, 2018). Accession numbers with

945 crosses indicate the nine well-annotated datasets that were used for benchmarking.

946

947 **Supplementary Table 3: Cloning primer sequences used for generating riboprobes for the**

948 **FISH experiments.** Functional annotations of the genes were given in the *S. mansoni* genome

949 version 9 (WormBase, WS268).