

1 scGNN: a novel graph neural network framework for single-cell RNA-Seq analyses

2
3 Juexin Wang^{1, *}, Anjun Ma^{2, *}, Yuzhou Chang², Jianting Gong¹, Yuexu Jiang¹, Hongjun Fu³,
4 Cankun Wang², Ren Qi², Qin Ma^{2, §}, Dong Xu^{1, §}

5
6 ¹ Department of Electrical Engineering and Computer Science, and Christopher S. Bond Life
7 Sciences Center, University of Missouri, Columbia, MO 65211, USA

8 ² Department of Biomedical Informatics, College of Medicine, The Ohio State University,
9 Columbus, OH 43210, USA

10 ³ Department of Neuroscience, The Ohio State University, Columbus, OH 43210, USA

11 * These authors contributed equally to the paper as first authors

12 § To whom correspondence should be addressed

13 Dr. Qin Ma. Tel: (706)-254-4293; Email: qin.ma@osumc.edu

14 Dr. Dong Xu. Tel: (573)-882-2299; Email: xudong@missouri.edu.

15 16 17 **ABSTRACT**

18 Single-cell RNA-sequencing (scRNA-Seq) is widely used to reveal the heterogeneity and
19 dynamics of tissues, organisms, and complex diseases, but its analyses still suffer from multiple
20 grand challenges, including the sequencing sparsity and complex differential patterns in gene
21 expression. We introduce the scGNN (single-cell graph neural network) to provide a hypothesis-
22 free deep learning framework for scRNA-Seq analyses. This framework formulates and
23 aggregates cell-cell relationships with graph neural networks and models heterogeneous gene
24 expression patterns using a left-truncated mixture Gaussian model. scGNN integrates three
25 iterative multi-modal autoencoders and outperforms existing tools for gene imputation and cell
26 clustering on four benchmark scRNA-Seq datasets. In an Alzheimer's disease study with 13,214
27 single nuclei from postmortem brain tissues, scGNN successfully illustrated disease-related
28 neural development and the differential mechanism. scGNN provides an effective representation
29 of gene expression and cell-cell relationships. It is also a novel and powerful framework that can
30 be applied to scRNA-Seq analyses.

31 32 **BACKGROUND**

33 Single-cell RNA sequencing (scRNA-seq) techniques enable transcriptome-wide gene
34 expression measurement in individual cells, which are essential for identifying cell type clusters,
35 inferring the arrangement of cell populations according to trajectory topologies, and highlighting
36 somatic clonal structures while characterizing cellular heterogeneity in complex diseases^{1,2}.
37 scRNA-seq analysis for biological inference remains challenging due to its complex and un-
38 determined data distribution, which has an extremely large volume and high rate of dropout
39 events. Some pioneer methodologies, e.g., Phenograph³, MAGIC⁴, and Seurat⁵ use a k-nearest-
40 neighbor (KNN) graph to model the relationships between cells. However, such a graph
41 representation may over-simplify the complex cell and gene relationships of the global cell
42 population. Recently, the emerging graph neural network (GNN) has deconvoluted node
43 relationships in a graph through neighbor information propagation in a deep learning architecture⁶.
44 Compared with other autoencoders used in the scRNA-Seq analysis⁷⁻¹⁰ for revealing an effective
45 representation of scRNA-Seq data via recreating its own input, the unique feature of graph
46 autoencoder is in being able to learn a low dimensional representation of the graph topology and
47 train node relationships in a global view of the whole graph¹¹.

48
49 We introduce a multi-modal framework scGNN (single-cell graph neural network) for modeling
50 heterogeneous cell-cell relationships and their underlying complex gene expression patterns from
51 scRNA-Seq. scGNN trains low dimensional feature vectors (i.e., embedding) to represent

52 relationships among cells through topological abstraction based on both gene expression and
53 transcriptional regulation information. There are three unique features in scGNN: (i) scGNN
54 utilizes GNN with multi-modal autoencoders to formulate and aggregate cell-cell relationships,
55 providing a hypothesis-free framework to derive biologically meaningful relationships. The
56 framework does not need to assume any statistical distribution or relationships for gene
57 expression data or dropout events. (ii) Cell-type-specific regulatory signals are modeled in
58 building a cell graph, equipped with a left-truncated mixture Gaussian (LTMG) model for scRNA-
59 Seq data¹². This can improve the signal-to-noise ratio in terms of embedding biologically
60 meaningful information. (iii) Bottom-up cell relationships are formulated from a dynamically pruned
61 GNN cell graph. The entire graph can be represented by pooling on learned graph embedding of
62 all nodes in the graph. The graph embedding can be used as low-dimensional features with
63 tolerance to noises for the preservation of topological relationships in the cell graph. The derived
64 cell-cell relationships are adopted as regularizers in the autoencoder training to recover gene
65 expression values.

66
67 scGNN has great potential in capturing biological cell-cell relationships in terms of cell type
68 clustering, cell trajectory inference, cell lineages formation, and cells transitioning between states.
69 In this paper, we mainly focus on discovering its applicative power in two fundamental aspects
70 from scRNA-Seq data, i.e., gene imputation and cell clustering. Gene imputation aims to solve
71 the dropout issue which commonly exists in scRNA-Seq data where the expressions of a large
72 number of active genes are marked as zeros¹³⁻¹⁵. The excess of zero values often needs to be
73 recovered or handled to avoid the exaggeration of the dropout events in many downstream
74 biological analyses and interpretations. Existing imputation methods, such as MAGIC⁴ and
75 SAVER¹⁶, have an issue in generating biased estimates of gene expression and tend to induce
76 false-positive and biased gene correlations that could possibly eliminate some meaningful
77 biological variations^{17,18}. On the other hand, many studies, including Seurat⁵ and Phenograph³,
78 have explored the cell-cell relationships using raw scRNA-seq data, and built cell graphs with
79 reduced data dimensions and detected cell clusters by applying the Louvain modularity
80 optimization. Accurate cell-cell relationships obey the rule that cells are more homogeneous within
81 a cell type and more heterogeneous among different cell types¹⁹. The scGNN model provides a
82 global perspective in exploring cell relationships by integrating cell neighbors on the whole
83 population.

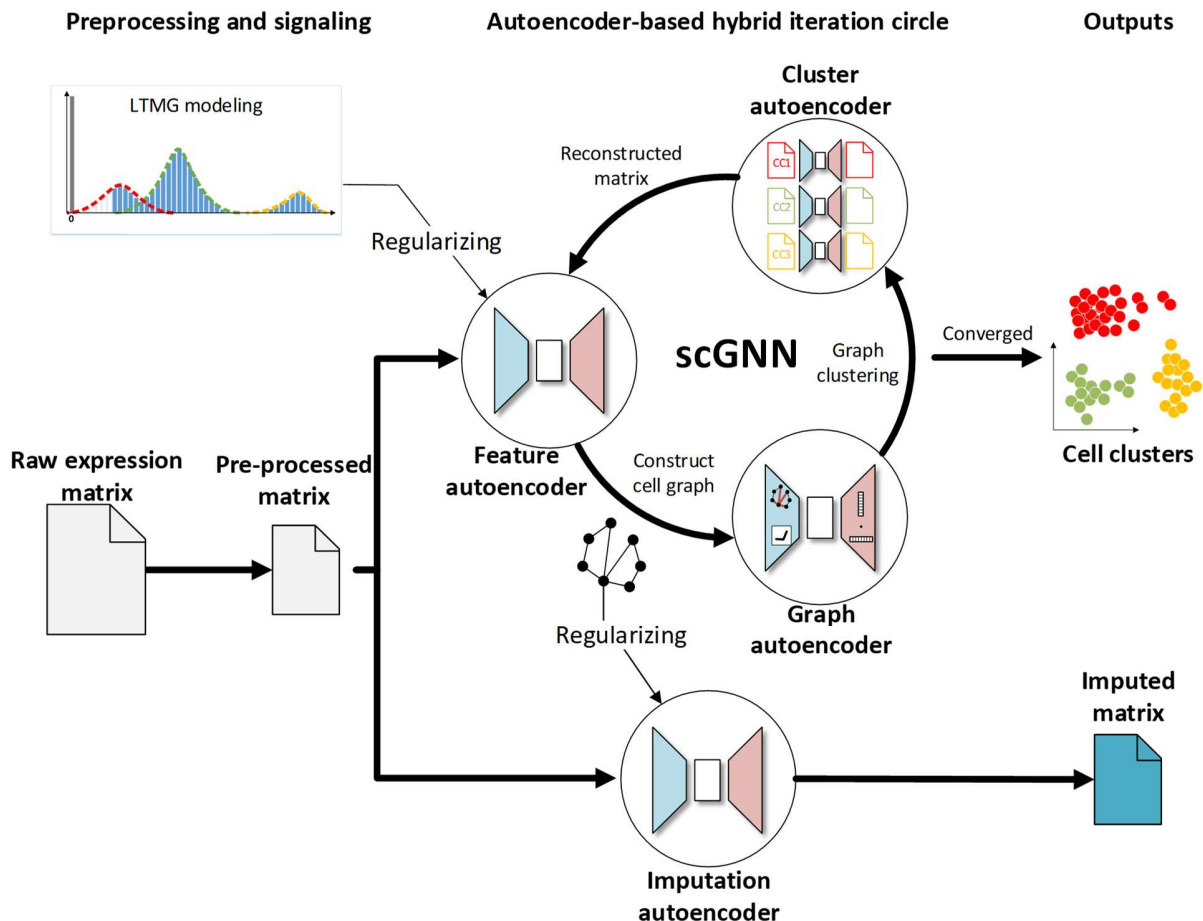
84
85 scGNN achieves promising performance in gene imputation and cell cluster prediction on four
86 scRNA-Seq datasets with gold-standard cell labels²⁰⁻²³, compared to seven existing imputation
87 and four clustering tools (**Supplementary Table S1**). We believe that the superior performance
88 in gene imputation and cell cluster prediction benefits from (i) our integrative autoencoder
89 framework, which synergistically determines cell clusters based on a bottom-up integration of
90 detailed pairwise cell-cell relationships and the convergence of predicted clusters, and (ii) the
91 integration of both gene regulatory signals and cell network representations in hidden layers as
92 regularizers of our autoencoders. To further demonstrate the power of scGNN in complex disease
93 studies, we applied it to an Alzheimer's disease (AD) dataset containing 13,214 single nuclei,
94 which elucidated its application power on cell-type identification and recovering gene expression
95 values²⁴. We claim that such a GNN-based framework is powerful and flexible enough to have
96 great potential in integrating scMulti-Omics data.

97 98 **RESULTS**

99 **The architecture of scGNN is comprised of stacked autoencoders**

100 The main architecture of scGNN is used to seek effective representations of cells and genes that
101 are useful for performing different tasks in scRNA-Seq data analyses (**Figure 1** and
102 **Supplementary Figure S1**). It has three comprehensive computational components in an

103 iteration process, including gene regulation integration in a feature autoencoder, cell graph
 104 representation in a graph autoencoder, gene expression updating in a set of parallel cell-type-
 105 specific cluster autoencoders, as well as the final gene expression recovery in an imputation
 106 autoencoder (**Figure 1**).
 107



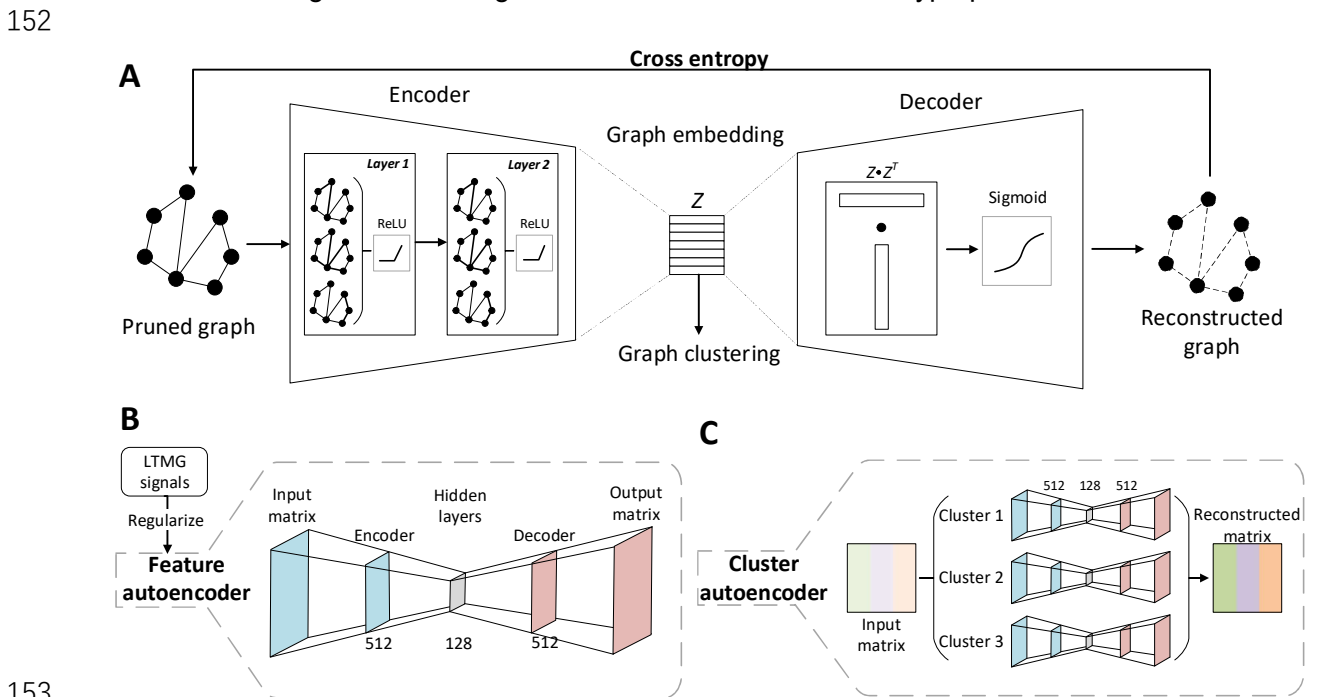
108 **Figure 1.** The architecture of scGNN. It takes the gene expression matrix generated from scRNA-Seq as
 109 the input. LTMG can translate the input gene expression data into a discretized regulatory signal as the
 110 regularizer for the feature autoencoder. The feature autoencoder learns a dimensional representation of
 111 the input as embedding, upon which a cell graph is constructed and pruned. The graph autoencoder
 112 learns a topological graph embedding of the cell graph, which is used for cell type clustering. The cells
 113 in each cell type have an individual cluster autoencoder to reconstruct gene expression values. The
 114 framework treats the reconstructed expression as a new input iteratively until converging. Finally,
 115 the imputed gene expression values are obtained by the feature autoencoder regularized by the cell-cell
 116 relationships in the learned cell graph on the original preprocessed raw expression matrix through the
 117 imputation autoencoder.

118 **Feature autoencoder.** This autoencoder intakes the pre-processed gene expression matrix after
 119 the removal of low-quality cells and genes, normalization, and variable gene ranking. First, the
 120 LTMG model^{12,25} is adopted to the top 2,000 variable genes to quantify gene regulatory signals
 121 encoded among diverse cell states in scRNA-Seq data (**Online Methods** and **Supplementary**
 122 **Figure S2**). This model was built based on the kinetic relationships between the transcriptional
 123 regulatory inputs and mRNA metabolism and abundance, which can infer the expression multi-
 124 modalities across single cells. The captured signals have a better signal-to-noise ratio to be used

125 as a high-order restraint to regularize the feature autoencoder. The aim of this regularization is to
 126 treat each gene differently based on their individual regulation status through a penalty in the loss
 127 function. The feature autoencoder learns a low dimensional embedding by the gene expression
 128 reconstruction together with the regularization. A cell-cell graph is generated from the learned
 129 embedding via the KNN graph, where nodes represent individual cells and the edges represent
 130 neighborhood relations among these cells^{26,27}. Then, the cell graph is pruned from selecting an
 131 adaptive number of neighbors for each node on the KNN graph by removing the noisy edges³.

132
 133 **Graph autoencoder.** Taking the pruned cell graph as input, the encoder of the graph autoencoder
 134 uses GNN to learn a low dimensional embedding of each node and then regenerates the whole
 135 graph structure through the decoder of the graph autoencoder (**Figure 2A**). Based on the
 136 topological properties of the cell graph, the graph autoencoder abstracts intrinsic high-order cell-
 137 cell relationships propagated on the global graph. The low dimensional graph embedding
 138 integrates the essential pairwise cell-cell relationships and the global cell-cell graph topology
 139 using a graph formulation by regenerating the topological structure of the input cell graph. Then
 140 the k-means clustering method is used to cluster cells on the learned graph embedding²⁸, where
 141 the number of clusters is determined by the Louvain algorithm on the cell graph.

142
 143 **Cluster autoencoder.** The expression matrix in each cell cluster from the feature autoencoder is
 144 reconstructed through the cluster autoencoder. Using the inferred cell type information from the
 145 graph autoencoder, the cluster autoencoder treats different cell types specifically and regenerates
 146 expression in the same cell cluster. The cluster autoencoder helps discover cell-type-specific
 147 information for each cell type in its individualized learning. Accompanied by the feature
 148 autoencoder, the cluster autoencoder leverages the inferences between global and cell-type-
 149 specific representation learning. Iteratively, the reconstructed matrix is fed back into the feature
 150 autoencoder. The iteration process stops until it converges with no change in cell clustering and
 151 this cell clustering result is recognized as the final results of cell type prediction.



153
 154
 155 **Figure 2. The architecture of scGNN Autoencoders.** (A) The graph autoencoder takes the adjacent
 156 matrix of the pruned graph as the input. The encoder consists of two layers of GNNs. In each layer, each

157 node of the graph aggregates information from its neighbors. The encoder learns a low dimensional
158 presentation (i.e., graph embedding) of the pruned cell graph. The decoder reconstructs the adjacent matrix
159 of the graph by dot products of the learned graph embedding followed by a sigmoid activation function. The
160 graph autoencoder is trained by minimizing the cross-entropy loss between the input and the reconstructed
161 graph. Cell clusters are obtained by applying k-means and Louvain on the graph embedding. (B) The
162 feature autoencoder takes the expression matrix as the input, regularized by LTMG signals. The dimensions
163 of the encoder and decoder layers are 512×128 and 128×512, respectively. The feature autoencoder is
164 trained by minimizing the difference between the input matrix and the output matrix. (C) The cluster
165 autoencoder takes a reconstructed expression matrix from the feature autoencoder as the input. An
166 individual encoder is built on the cells in each of the identified clusters, and each autoencoder is trained
167 individually. The concatenation of the results from all clusters is treated as the reconstructed matrix.

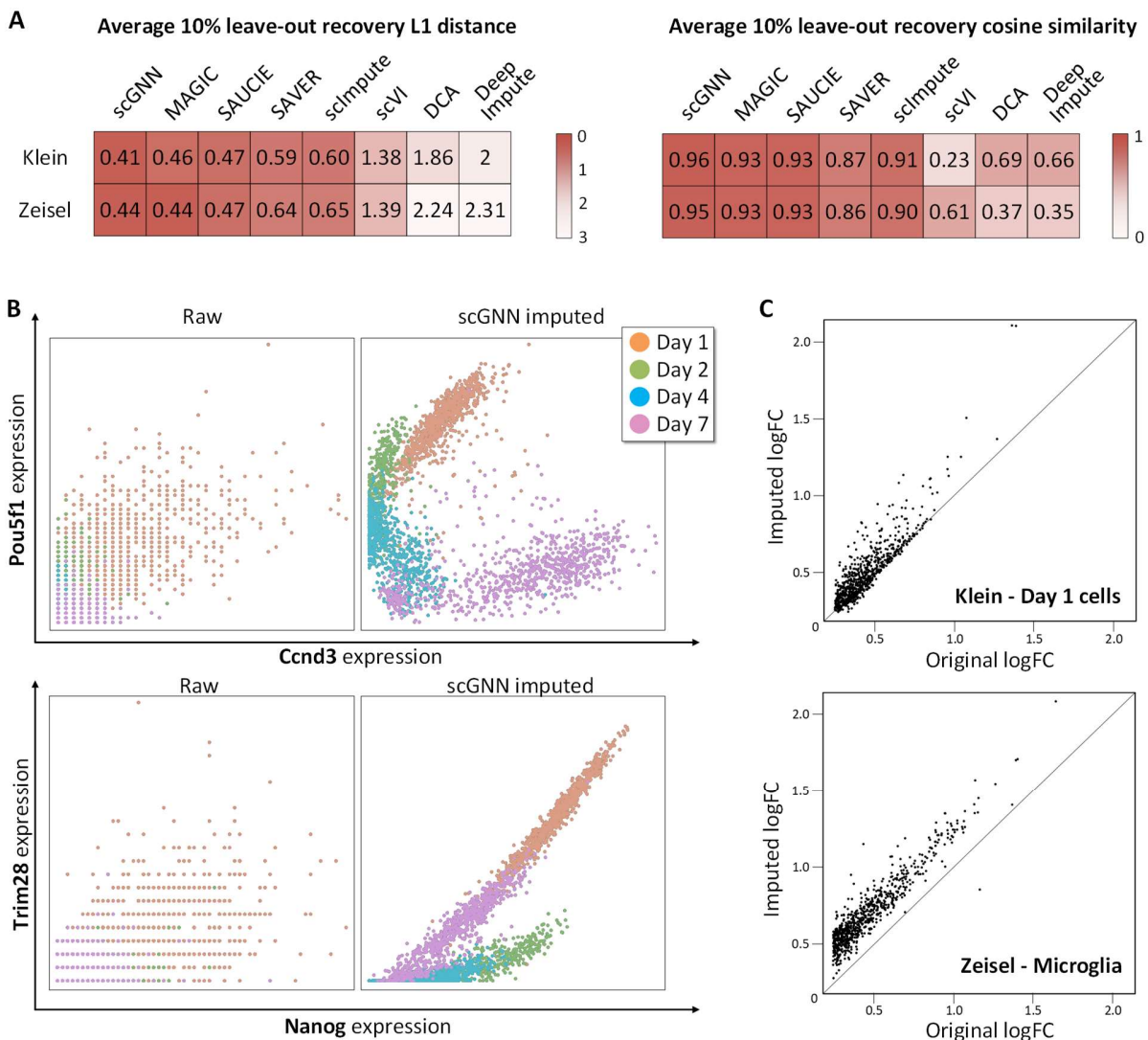
168
169 *Imputation autoencoder.* After the iteration stops, this imputation autoencoder takes the original
170 gene expression matrix as input and is trained with the additional L1 regularizer of the inferred
171 cell-cell relationships. The regularizers (see **Online Methods**) are generated based on edges in
172 the learned cell graph in the last iteration and their co-occurrences in the same predicted cell type.
173 Besides, the L1 penalty term is applied to increase the model generalization by squeezing more
174 zeroes into the autoencoder model weights. The sparsity brought by the L1 term benefits the
175 expression imputation in dropout effects. Finally, the reconstructed gene expression values are
176 used as the final imputation output.

177 **scGNN can effectively impute scRNA-Seq data and accurately predict cell clusters**

178 To assess the imputation and cell clustering performance of scGNN, four scRNA datasets (i.e.,
179 Chung²³, Kolodziejczyk²⁰, Klein²¹, and Zeisel²²) with gold-standard cell type labels are chosen as
180 the benchmarks (more performance evaluation on other datasets can be found in **Supplementary**
181 **Materials**). We manually simulated the dropout effects by randomly flipping 10% of the non-zero
182 entries to zeros. The median L1 distance between the original dataset and the imputed values for
183 these corrupted entries were evaluated to compare scGNN with MAGIC⁴, SAUCIE⁸, SAVER¹⁶,
184 scImpute²⁹, scVI³⁰, DCA⁹, and DeepImpute³¹. scGNN shows the lowest L1 distance and the
185 highest cosine similarity in recovering leave-out values, indicating that it can accurately capture
186 and restore true expression values (**Online Methods** and **Figure 3A**). Furthermore, scGNN
187 depicts the underlying gene-gene relationships missed due to the sparsity of scRNA-Seq. For
188 example, two pluripotency epiblast gene pairs, *Ccnd3* versus *Pou5f1* and *Nanog* versus *Trim28*,
189 are lowly correlated in the original raw data but show strong positive correlations, which are
190 differentiated by time points after scGNN imputation and, therefore, perform with a consistency
191 leading to the desired results sought in the original paper²¹ (**Figure 3B**). The relationships of four
192 more gene pairs are also enhanced (**Supplementary Figure S3**). In the Zeisel dataset, scGNN
193 amplifies differentially expressed genes (DEGs) signals with a higher fold change than the
194 original, using an imputed matrix to confidently depict the cluster heterogeneity (**Figure 3C** and
195 **Supplementary Figure S4**).

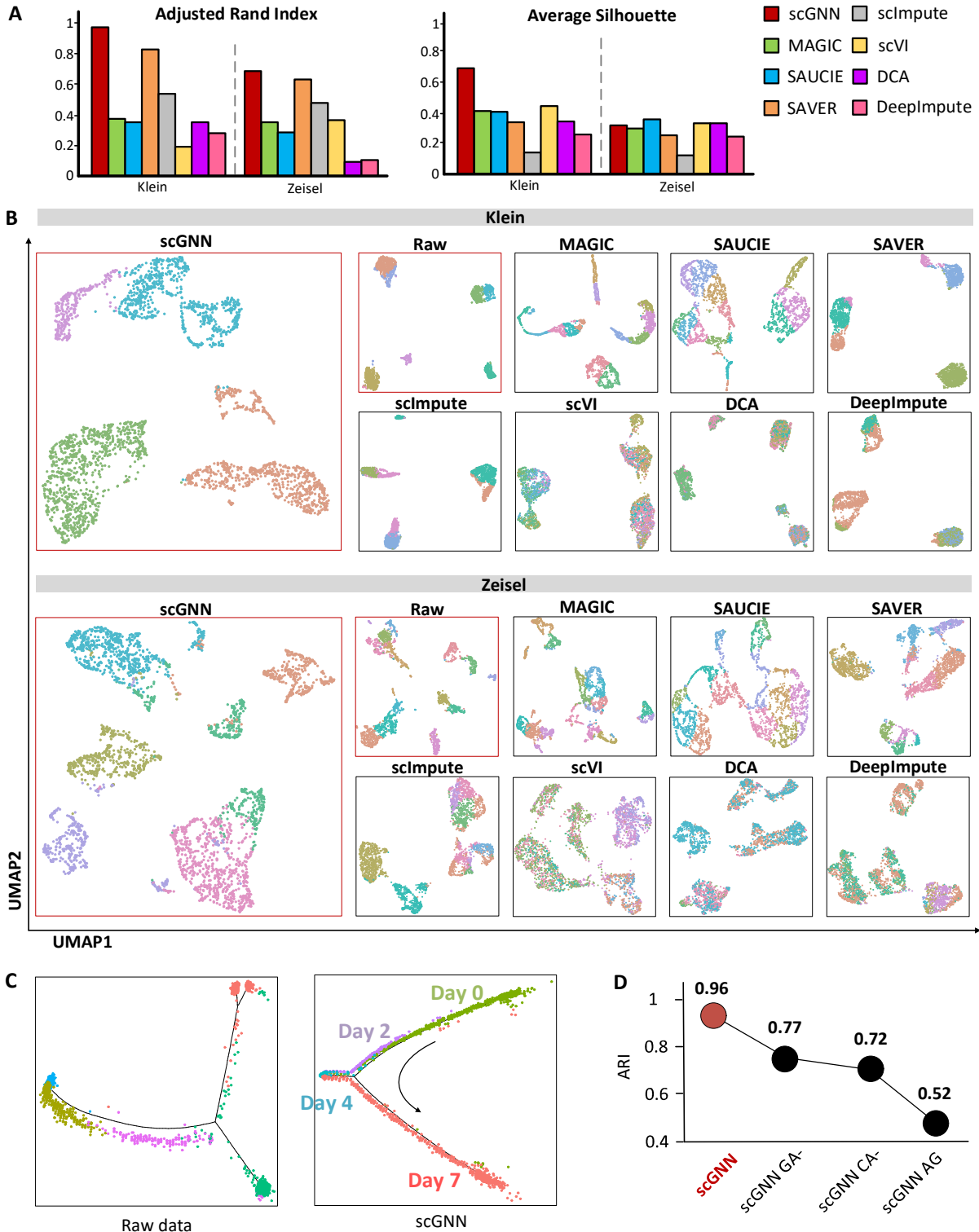
196
197 Besides the artificial dropout benchmarks, we continued to evaluate the clustering performance
198 of scGNN and the seven imputation tools on the same two datasets. The predicted cell labels
199 were systematically evaluated using 10 criteria including an adjusted Rand index (ARI)³²,
200 Silhouette³³, and eight other criteria (**Figure 4A**). By visualizing cell clustering results on UMAPs,
201 one can observe more apparent closeness of cells within the same cluster and separation among
202 different clusters when using scGNN embeddings compared to the other seven imputation tools
203 (**Figure 4B**). The expression patterns show heterogeneity along with embryonic stem cell
204 development. In the case of Klein's time-series data, scGNN recovered a complex structure that
205 was not well represented by the raw data, showing a well-aligned trajectory path of cell
206 development from Day 1 to Day 7 (**Figure 4C**). Moreover, scGNN showed significant
207 enhancement in cell clustering compared to the clustering tool (e.g., Seurat) when using the raw
208

209 data (**Supplementary Figure S5**). On top of that, to address the significance of using the graph
 210 autoencoder and cluster autoencoder in scGNN, we performed ablation tests to bypass each
 211 autoencoder and compare the ARI results on the Klein dataset (**Figure 4D**). The results showed
 212 that removing either of these two autoencoders dramatically decreased the performance of
 213 scGNN in terms of cell clustering accuracy. Another test using all genes rather than the top 2,000
 214 variable genes also showed poor performance in the results and doubled the runtime of scGNN,
 215 indicating that those low variable genes may reduce the signal-to-noise ratio and negatively affect
 216 the accuracy of scGNN. The design and comprehensive results of the ablation studies on both
 217 clustering and imputation are detailed in **Supplementary Method** and **Table S2-S7** and **S11**. We
 218 also extensively studied the parameter selection in **Supplementary Table S8-S10** and **S12**.
 219



220
 221 **Figure 3.** Comparison of the imputation performance. (A) The L1 distance (the lower the better) and cosine
 222 similarity (the higher the better) comparing a 10% leave-out test between scGNN and seven imputation
 223 tools on the Klein and Zeisel datasets. scGNN achieved the best scores in both datasets, indicating its
 224 superior performance in gene expression recovery. (B) Co-expression patterns can be addressed more
 225 explicitly after applying scGNN on the Klein data. No clear gene pair relationship of *Ccnd3* versus *Pou5f1*
 226 (*upper panel*) and *Nanog* versus *Trim28* (*lower panel*) is observed in the raw data (left) compared to the
 227 observation of unambiguous correlations within each cell type after scGNN imputation (right). (C)

228 Comparison of DEG logFC scores using the original expression value (x-axis) and the scGNN imputed
 229 expression values (y-axis) identified in Day 1 cells of the Klein data (up) and Microglia cells of the Zeisel
 230 data (bottom). The differentiation signals are amplified after imputation.
 231



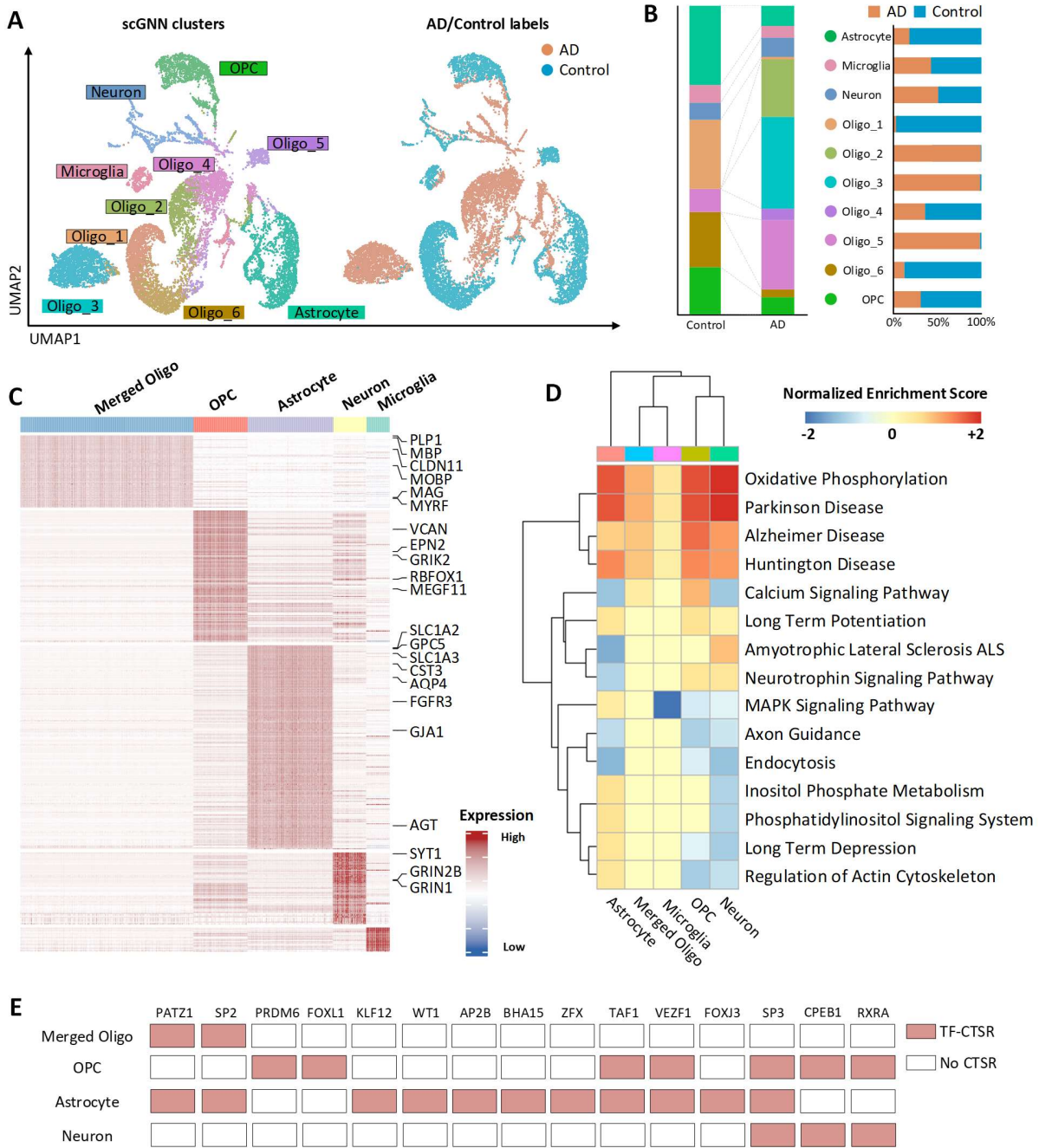
232

233 **Figure 4.** Cell clustering and trajectory evaluations. (A) Comparison of ARI and Silhouette scores among
234 scGNN and seven tools using Klein and Zeisel datasets. (B) Comparison of UMAP visualizations on the
235 same two datasets, indicating that when scGNN embeddings are utilized, cells are more closely grouped
236 within the same cluster but when other tools are used, cells are more separated between clusters. Raw
237 data is clustered and visualized using Seurat. (C) Pseudotime analysis using the raw expression matrix and
238 scGNN imputed matrix of the Klein dataset via Monocle2. (D) Justification of using the graph autoencoder,
239 the cluster autoencoder, and the top 2,000 variable genes on the Klein dataset in the scGNN framework,
240 in terms of ARI. scGNN CA- shows the results of the graph autoencoder's ablation, CA- shows the results
241 of the cluster autoencoder's ablation, and AG shows the results after using all genes in the framework.

242 243 **scGNN illustrates AD-related neural development and the underlying regulatory** 244 **mechanism**

245 To further demonstrate the applicative power of scGNN, we applied it to a scRNA-Seq dataset
246 (GEO accession number GSE138852) containing 13,214 single nuclei collected from six AD and
247 six control brains³⁴. scGNN identifies 10 cell clusters, including microglia, neurons,
248 oligodendrocyte progenitor cells (OPCs), astrocytes, and six sub-clusters of oligodendrocytes
249 (**Figure 5A**). Specifically, the proportions of these six oligodendrocyte sub-clusters differ between
250 AD patients (Oligos 2, 3, and 4) and healthy controls (Oligos 1, 5, and 6) (**Figure 5B**). Moreover,
251 the difference between AD and the control in the proportion of astrocyte and OPCs is observed,
252 indicating the change of cell population in AD patients compared to healthy controls (**Figure 5B**).
253 We then combined these six oligodendrocyte sub-clusters into one to discover DEGs. Since
254 scGNN can significantly increase true signals in the raw dataset, DEG patterns are more explicit
255 (**Supplementary Figure S6**). Among all DEGs, we confirmed 22 genes as cell-type-specific
256 markers for astrocytes, OPCs, oligodendrocytes, and neurons, in that order³⁵ (**Figure 5C**). A
257 biological pathway enrichment analysis shows several highly positive-enrichments in AD cells
258 compared to control cells among all five cell types. These enrichments include oxidative
259 phosphorylation and pathways associated with AD, Parkinson's disease, and Huntington
260 disease³⁶ (**Figure 5D** and **Supplementary Figure S7**). Interestingly, we observed a strong
261 negative enrichment of the MAPK (mitogen-activated protein kinase) signaling pathway in the
262 microglia cells, suggesting a relatively low MAPK regulation in microglia than other cells.

263
264 In order to investigate the regulatory mechanisms underlying the AD-related neural development,
265 we applied the imputed matrix of scGNN to IRIS3 (an integrated cell-type-specific regulon
266 inference server from single-cell RNA-Seq) and identified 21 cell-type-specific regulons (CTSR)
267 in five cell types³⁷ (**Figure 5E** and **Supplementary Table S13**; IRIS3 job ID: 20200626160833).
268 Not surprisingly, we identified several AD-related transcription factors (TFs) and target genes that
269 have been reported to be involved in the development of AD. SP2 is a common TF identified in
270 both oligodendrocytes and astrocytes. It has been shown to regulate the *ABCA7* gene, which is
271 an IGAP (International Genomics of Alzheimer's Project) gene that is highly associated with late-
272 onset AD³⁸. We also observed an SP2 CTSR in astrocytes that regulate *APOE*, *AQP4*, *SLC1A2*,
273 *GJA1*, and *FGFR3*. All of these five targeted genes are marker genes of astrocytes, which have
274 been reported to be associated with AD^{39,40}. In addition, the SP3 TF is identified in all cell clusters
275 which can regulate the synaptic function in neurons, and it is extremely activated in AD^{41,42}. We
276 identified CTSRs regulated by SP3 in OPCs, astrocytes, and neurons suggesting a significant
277 SP3 related regulation shifts in these three clusters. We observed 26, 60, and 22 genes that were
278 uniquely regulated in OPCs, astrocytes, and neurons, as well as 60 genes shared among the
279 three clusters (**Supplementary Table S14**). Such findings provide a direction for the discovery of
280 SP3 function in AD studies.



281
 282 **Figure 5.** Alzheimer's disease dataset (GSE138852) analysis based on scGNN. (A) Cell clustering UMAP.
 283 Labeled with scGNN clusters (left) and AD/control samples (right). (B) Comparison of cell proportions in
 284 AD/control samples (left) and each cluster (right). (C) Heatmap of DEGs (logFC > 0.25) in each cluster. Six
 285 oligodendrocyte sub-clusters are merged as one to compare with other cell types. Marker genes identified
 286 in DEGs are listed on the right. (D) Selected AD-related enrichment pathways in each cell type in the
 287 comparison between AD and control cells. (E) Underlying TFs are responsible for the cell-type-specific
 288 gene regulations identified by IRIS3.

289
 290 **DISCUSSION**

291 It is still a fundamental challenge to explore cellular heterogeneity in high-volume, high-sparsity,
292 and noisy scRNA-Seq data, where the high-order topological relationships of the whole-cell graph
293 are still not well explored and formulated. The key innovations of scGNN are incorporating global
294 propagated topological features of the cells through GNNs, together with integrating gene
295 regulatory signals in an iterative process for scRNA-Seq data analysis. The benefits of GNN is its
296 intrinsic learnable properties of propagating and aggregating attributes to capture relationships
297 across the whole cell-cell graph. Hence, the learned graph embedding can be treated as the high-
298 order representations of cell-cell relationships in scRNA-Seq data in the context of graph topology.
299 Unlike the previous autoencoder applications in scRNA-Seq data analysis, which only captures
300 the top-down distributions of the overall cells, scGNN can effectively aggregate detailed
301 relationships between similar cells using a bottom-up approach. Furthermore, scGNN integrates
302 gene regulatory signals efficiently by representing them discretely in LTMG in the feature
303 autoencoder regularization. These gene regulatory signals can help identify biologically
304 meaningful gene-gene relationships as they apply to our framework and eventually, they are
305 proven capable of enhancing performance. Technically, scGNN adopts multi-modal autoencoders
306 in an iterative manner to recover gene expression values and cell type prediction simultaneously.
307 Notably, scGNN is a hypothesis-free deep learning framework on a data-driven cell graph model,
308 and it is flexible to incorporate different statistical models (e.g. LTMG) to analyze complex scRNA-
309 Seq datasets.

310
311 Some limitations can still be found in scGNN. (i) It is prone to achieve better results with large
312 datasets, compared to relatively small datasets (e.g., less than 1,000 cells), as it is designed to
313 learn better representations with many cells from scRNA-Seq data, as shown in the benchmark
314 results, and (ii) Compared with statistics model-based methods, the iterative autoencoder
315 framework needs more computational resources, which is more time-consuming and less
316 interpretable. In the future, we will investigate creating a more efficient scGNN model with a lighter
317 and more compressed architecture.

318
319 In the future, we will continue to enhance scGNN by implementing heterogeneous graphs to
320 support the integration of single-cell multi-omics data (e.g., the intra-modality of Smart-Seq2 and
321 Droplet scRNA-Seq data; and the inter-modality integration of scRNA-Seq and scATAC-Seq
322 data). We will also incorporate attention mechanisms and graph transformer models⁴³ to make
323 the analyses more explainable. Specifically, by allowing the integration of scRNA-Seq and
324 scATAC-Seq data, scGNN has the potential to elucidate cell-type-specific gene regulatory
325 mechanisms⁴⁴. On the other hand, T cell receptor repertoires are considered as unique identifiers
326 of T cell ancestries that can improve both the accuracy and robustness of predictions regarding
327 cell-cell interactions⁴⁵. scGNN can also facilitate batch effects and build connections across
328 diverse sequencing technologies, experiments, and modalities. Moreover, scGNN can be applied
329 to analyze spatial transcription datasets regarding spatial coordinates as additional regularizers
330 to infer the cell neighborhood representation and better prune the cell graph. We plan to develop
331 a more user-friendly software system from our scGNN model, together with modularized analytical
332 functions in support of standardizing the data format, quality control, data integration, multi-
333 functional scMulti-seq analyses, performance evaluations, and interactive visualizations.

334

335 ONLINE METHODS

336 Dataset preprocessing

337 scGNN takes the scRNA-Seq gene expression profile as the input. Data filtering and quality
338 control are the first steps of data preprocessing. Due to the high dropout rate of scRNA-seq
339 expression data, only genes expressed as nonzero in more than 1% of cells, and cells expressed
340 as nonzero in more than 1% of genes are kept. Then, genes are ranked by standard deviation,
341 i.e., the top 2,000 genes in variances are used for the study. All the data are log-transformed.

342

343 Left Truncated Mixed Gaussian (LTMG) modeling

344 A mixed Gaussian model with left truncation assumption is used to explore the regulatory signals
345 from gene expression¹². The normalized expression values of gene X over N cells are denoted
346 as $X = \{x_1, \dots, x_N\}$, where $x_j \in X$ is assumed to follow a mixture of k Gaussian distributions,
347 corresponding to k possible gene regulatory signals (**TRSs**). The density function of X is:

348

$$349 \quad p(X; \theta) = \prod_{j=1}^N p(x_j; \theta) = \prod_{j=1}^N \sum_{i=1}^k \alpha_i p(x_j; \theta_i) = \prod_{j=1}^N \sum_{i=1}^k \alpha_i \frac{1}{\sqrt{2\pi}\sigma_i} e^{-\frac{(x_j - \mu_i)^2}{2\sigma_i^2}} = L(\theta; X) \quad (1)$$

350

351 where α_i is the mixing weight, μ_i and σ_i are the mean and standard deviation of the i^{th} Gaussian
352 distribution, which can be estimated by: $\theta^* = \underset{\theta}{\text{arg max}} L(\theta; X)$ to model the errors at zero and the low
353 expression values. With the left truncation assumption, the gene expression profile is split into
354 M , which is a truly measured expression of values, and $N - M$ representing left-censored gene
355 expressions for N conditions. The parameter θ maximizes the likelihood function and can be
356 estimated by an expectation-maximization algorithm. The number of Gaussian components is
357 selected by the Bayesian Information Criterion; then, the original gene expression values are
358 labeled to the most likely distribution under each cell. In detail, the probability that x_j belongs to
359 distribution i is formulated by:

$$360 \quad p(x_j \in \text{TRS } i | K, \theta^*) \propto \frac{\alpha_i}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(x_j - \mu_i)^2}{2\sigma_i^2}} \quad (2)$$

361

362 where x_j is labeled by TRS i if $p(x_j \in \text{TRS } i | K, \theta^*) = \max_{i=1, \dots, K} (p(x_j \in \text{TRS } i | K, \theta^*))$. Thus, the
363 discrete values $(1, 2, \dots, K)$ for each gene are generated.

364

365 Feature autoencoder

366 The feature autoencoder is proposed to learn the representative embedding of the scRNA
367 expression through stacked two layers of dense networks in both the encoder and decoder. The
368 encoder constructs the low dimensional embedding of X' from the input gene expression X , and
369 the encoder reconstructs the expression \hat{X} from the embedding; thus, $X, \hat{X} \in \mathbb{R}^{N \times M}$ and $X' \in$
370 $\mathbb{R}^{N \times M'}$, where M is the number of input genes, M' is the dimension of the learned embedding,
371 and $M' < M$. The objective of training the feature autoencoder is to achieve a maximum similarity
372 between the original and reconstructed through minimizing the loss function, in which $\sum (X - \hat{X})^2$
373 is the main term serving as the mean squared error (**MSE**) between the original and the
374 reconstructed expressions.

375

376 Regularization

377 Regularization is adopted to integrate gene regulation information during the feature autoencoder
378 training process. The aim of this regularization is to treat each gene differently based on their

379 individual gene regulation role through penalizing it in the loss function. In each cell, the MSE of
380 each gene is element-wise multiplication with discrete gene regulation signals from TRS, as
381 defined in Eq.(5).

$$382 \quad \alpha \sum (X - \hat{X})^2 \cdot TRS \quad (5)$$

383 where α is a parameter used to control the strength of gene regulation regularization; $\alpha \in [0,1]$.
384 Thus, the loss function of the feature autoencoder is shown as Eq.(6).

$$385 \quad Loss = (1 - \alpha) \sum (X - \hat{X})^2 + \alpha \sum (X - \hat{X})^2 \cdot TRS \quad (6)$$

387
388 In the encoder, the output dimensions of the first and second layers are set as 512 and 128,
389 respectively. Each layer is followed by the ReLU activation function. In the decoder, the output
390 dimensions of the first and second layers are 128 and 512, respectively. Each layer is followed
391 by a sigmoid activation function. The learning rate is set as 0.001. The cluster autoencoder has
392 the same architecture as the feature autoencoder, but without gene regulation regularization in
393 the loss function.

394

395 **Cell graph and pruning**

396 The cell graph formulates the cell-cell relationships using embedding learned from the feature
397 autoencoder. As done in previous works^{4,46}, the cell graph is built from a KNN graph, where nodes
398 are individual single-cells, and the edges are relationships between cells. K is the predefined
399 parameter used to control the scale of the captured interaction between cells. Each node finds its
400 neighbors within the K shortest distances and creates edges between them and itself. Euclidian
401 distance is calculated as the weights of the edges on the learned embedding vectors. The pruning
402 process selects an adaptive number of neighbors for each node on the original KNN graph and
403 keeps a more biologically meaningful cell graph. Here, Isolation Forest is applied to prune the
404 graph to detect the outlier in the K -neighbors of each node⁴⁷. Isolation Forest builds individual
405 random forest to check distances from the node to all K neighbors and only disconnects the
406 outliers.

407

408 **Graph autoencoder**

409 The graph autoencoder learns to embed and represent the topological information from the
410 pruned cell graph. For the input pruned cell graph, $G = (V, E)$ with $N = |V|$ nodes denoting the
411 cells and E representing the edges. A is its adjacency matrix and D is its degree matrix. The
412 node feature matrix of the graph autoencoder is the learned embedding X' from the feature
413 autoencoder.

414

415 The graph convolution network (GCN) is defined as $GCN(X', A) = ReLU(\tilde{A}X'W)$, and W is a
416 weight matrix learned from the training. $\tilde{A} = D^{-1/2}AD^{-1/2}$ is the symmetrically normalized
417 adjacency matrix and activation function $ReLU(\cdot) = \max(0, \cdot)$. The encoder of the graph
418 autoencoder is composed of two layers of GCN, and Z is the graph embedding learned through
419 the encoder in Eq.(7). W_1 and W_2 are learned weight matrices in the first and second layers, and
420 the output dimensions of the first and second layers are set at 32 and 16, respectively. The
421 learning rate is set at 0.001.

$$422 \quad Z = ReLU(\tilde{A}ReLU(\tilde{A}XW_1)W_2) \quad (7)$$

423

424 The decoder of the graph autoencoder is defined as an inner product between the embedding:

$$425 \quad \hat{A} = \text{sigmoid}(ZZ^T) \quad (8)$$

426 where \hat{A} is the reconstructed adjacent matrix of A . $\text{sigmoid}(\cdot) = 1/(1 + e^{-\cdot})$ is the sigmoid
 427 activation function.

428
 429 The goal of learning the graph autoencoder is to minimize the cross-entropy L between the input
 430 adjacent matrix A and the reconstructed matrix \hat{A} .

$$431 \quad L(A, \hat{A}) = -\frac{1}{N} \sum_{i=0}^N (A_i * \log(\hat{A}_i) + (1 - A_i) * \log(1 - \hat{A}_i)) \quad (9)$$

432
 433 where A_i and \hat{A}_i are the elements of adjacent matrix A and \hat{A} . N is the total number of
 434 elements in the adjacent matrix.

435 **Iterative process**

436
 437 The iterative process aims to build the single-cell graph iteratively until converging. The iterative
 438 process of the cell graph can be defined as:

$$439 \quad \tilde{A} = \lambda L_0 + (1 - \lambda) \frac{A_{ij}}{\sum_j A_{ij}} \quad (10)$$

440
 441 where L_0 is the normalized adjacency matrix of the initial pruned graph, and $L_0 = D_0^{-1/2} A_0 D_0^{-1/2}$,
 442 where D_0 is the degree matrix. λ is the parameter to control the converging speed, $\lambda \in [0,1]$.
 443 Each time in iteration t , two criteria are checked to determine whether to stop the iteration: (1)
 444 that is, to determine whether the adjacency matrix converges, i.e., $\tilde{A}_t - \tilde{A}_{t-1} < \gamma_1 \tilde{A}_0$, or (2)
 445 whether the inferred cell types are similar enough, i.e., $ARI < \gamma_2$. ARI is the similarity
 446 measurement, which is detailed in the next section. In our setting, $\lambda = 0.5$ and $\gamma_1, \gamma_2 = 0.99$. The
 447 cell type clustering results obtained in the last iteration are chosen as the final cell type results.

448 **Imputation autoencoder**

449
 450 After the iterative process stops, the imputation autoencoder imputes and denoises the raw
 451 expression matrix within the inferred cell-cell relationship. The imputation autoencoder shares the
 452 same architecture as the feature autoencoder, but it also uses three additional regularizers from
 453 the cell graph in Eq.(11), cell types in Eq.(12), and the L1 regularizer in Eq.(13).

$$454 \quad \gamma_1 \sum (X - \hat{X})^2 \cdot A \quad (11)$$

455
 456 where A is the adjacent matrix from the pruned cell graph in the last iteration. Cells within an
 457 edge in the pruned graph will be penalized in the training.

$$458 \quad \gamma_2 \sum (X - \hat{X})^2 \cdot B$$

$$459 \quad B_{ij} = \begin{cases} 1 & \text{where } i \text{ and } j \text{ in same cell type} \\ 0 & \text{else} \end{cases} \quad (12)$$

460
 461 where B is the relationship matrix between cells, and two cells in the same cell type have a B_{ij}
 462 value of 1. Cells within the same inferred cell type will be penalized in the training. γ_1, γ_2 are the
 463 intensities of the regularizers and $\gamma_1, \gamma_2 \in [0,1]$. The L1 regularizer is defined as

$$464 \quad \beta \sum |w| \quad (13)$$

465 which brings sparsity and increases the generalization performance of the autoencoder by
 466 reducing the number of non-zero w terms in $\sum |w|$, where β is a hyper-parameter controlling the
 467 intensity of the L1 term ($\beta \in [0,1]$). Therefore, the loss function of the imputation autoencoder is

468
$$Loss = (1 - \alpha) \sum (X - \hat{X})^2 + \beta \sum |w| + \sum (X - \hat{X})^2 (\alpha \cdot TRS + \gamma_1 A + \gamma_2 B) \quad (14)$$

469

470 **Benchmark evaluation compared to existing tools**

471 *Imputation evaluation.* For benchmarking imputation performance, we added noises by randomly
472 flipping 10% of the nonzero entries to zero to mimic the dropout effects. We evaluated both the
473 median L1 distance and cosine similarity between the original dataset and the imputed values for
474 these corrupted entries. For all the flipped entries, x is the row vector of the original expression,
475 and y is its corresponding row vector of the imputed expression. The L1 distance is the absolute
476 deviation between the value of the original and imputed expression. A lower L1 distance means
477 a higher similarity.

478
$$L1distance = |x - y|, \quad L1distance \in [0, +\infty) \quad (15)$$

479 The cosine similarity computes the dot products between original and imputed expression.

480
$$CosineSimilarity(x, y) = \frac{xy^T}{\|x\| \|y\|}, \quad CosineSimilarity \in [0, 1] \quad (16)$$

481 The process is repeated three times, and the mean and standard deviation were selected as a
482 comparison. The scores are compared between scGNN and seven imputation tools (i.e., MAGIC⁴,
483 SAUCIE⁸, SAVER¹⁶, scImpute²⁹, scVI³⁰, DCA⁹, and DeepImpute³¹), all using the default
484 parameters.

485

486 *Clustering evaluation.* We compared the cell clustering results of scGNN, the same seven
487 imputation tools, and four clustering tools (i.e., Seurat⁵, CIDR⁴⁸, Monocle⁴⁹, and RaceID⁵⁰), in
488 terms of ten clustering evaluation scores. The default parameters are applied in all test tools. ARI
489 ³² is used to compute similarities by considering all pairs of the samples that are assigned in
490 clusters in the current and previous clustering adjusted by random permutation.

491

492
$$ARI = \frac{RI - E[RI]}{\max(RI) - E[RI]} \quad (17)$$

493 where the unadjusted rand index (RI) is defined as

494
$$RI = \frac{a + b}{C_n^2} \quad (18)$$

495 where a is the number of pairs correctly labeled in the same sets, and b is the number of pairs
496 correctly labeled as not in the same dataset. C_n^2 is the total number of possible pairs. $E[RI]$ is
497 the expected RI of random labeling. More quantitative measurements are also used in the

498 **Supplemental Materials.**

499

500 **Case study of the AD database**

501 We applied scGNN on a public Alzheimer's disease (AD) scRNA-Seq data with 13,214 cells²⁴.
502 The resolution of scGNN was set to 1.0, KI was set to 20, and the remaining parameters were
503 kept as default. The AD patient and control labels were provided by the original paper and used
504 to color the cells on the same UMAP coordinates generated from scGNN. We simply combined
505 cells in six oligodendrocyte subpopulations into one cluster, referred to as merged oligo. The
506 DEGs were identified in each cell cluster via the Wilcoxon rank-sum test implemented in the
507 Seurat package along with adjusted p -values using the Benjamini-Hochberg procedure with a
508 nominal level of 0.05. DEGs with $\logFC > 0.25$ or < -0.25 were finally selected. We further
509 identified the DEGs between AD and control cells in each cluster using the same strategy and

510 applied GSEA for pathway enrichment analysis⁵¹. The imputed matrix, which resulted from
511 scGNN was then sent to IRIS3 for CTSR prediction, using the predicted cell clustering labels with
512 merged oligodendrocytes³⁷. The default parameters were served in regulatory analysis in IRIS3.

513

514 **Data availability**

515 Three benchmark and AD case datasets can be downloaded from GEO databases with accession
516 numbers of: GSE75688 (the Chung data); GSE65525 (the Klein data); GSE60361 (the Zeisel
517 data); and GSE138852 (AD case). The Kolodziejczy data can be accessed from EBI with an
518 accession number of E-MTAB-2600.

519

520 **Software Implementation**

521 Tools and packages used in this paper include: Python version 3.7.6, numpy version 1.18.1, torch
522 version 1.4.0, networkx version 2.4, pandas version 0.25.3, rpy2 version 3.2.4, matplotlib version
523 3.1.2, seaborn version 0.9.0, umap-learn version 0.3.10, munkres version 1.1.2, R version 3.6.1,
524 and igraph version 1.2.5. The IRIS3 website is at <https://bmbi.bmi.osumc.edu/iris3/index.php>.

525

526 **CODE AVAILABILITY**

527 Our tool is open source and publicly available at GitHub (<https://github.com/scgnn/scGNN>).

528

529 **ACKNOWLEDGEMENTS**

530 This work was supported by the National Institutes of Health's National Institute of General
531 Medical Sciences (awards R35-GM126985 and R01-GM131399).

532

533 **AUTHOR CONTRIBUTIONS**

534 Conceptualization: Q.M., D.X.; Methodology: J.W., A.M.; Software: J.W., C.Y.; Investigation: J.W.,
535 Q.R.; Formal Analysis: A.M., J.W., J.G., Y.C., J.Y. Resources and Reagents: J.W., J.G., R.Q.;
536 Writing, Review, and Editing: J.W., A.M., H.F., Q.M., D.X.

537

538 **COMPETING INTERESTS**

539 The authors declare no competing interests.

540

541 **REFERENCES**

- 542 1 Hwang, B., Lee, J. H. & Bang, D. Single-cell RNA sequencing technologies and bioinformatics
543 pipelines. *Exp Mol Med* **50**, 96, doi:10.1038/s12276-018-0071-8 (2018).
- 544 2 Gawel, D. R. *et al.* A validated single-cell-based strategy to identify diagnostic and therapeutic
545 targets in complex diseases. *Genome Med* **11**, 47, doi:10.1186/s13073-019-0657-3 (2019).
- 546 3 Levine, J. H. *et al.* Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells that
547 Correlate with Prognosis. *Cell* **162**, 184-197, doi:10.1016/j.cell.2015.05.047 (2015).
- 548 4 van Dijk, D. *et al.* Recovering Gene Interactions from Single-Cell Data Using Data Diffusion. *Cell*
549 **174**, 716-729 e727, doi:10.1016/j.cell.2018.05.061 (2018).
- 550 5 Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell transcriptomic
551 data across different conditions, technologies, and species. *Nat Biotechnol* **36**, 411-420,
552 doi:10.1038/nbt.4096 (2018).
- 553 6 Kipf, T. N. & Welling, M. Semi-Supervised Classification with Graph Convolutional Networks. *The*
554 *International Conference on Learning Representations (ICLR)* (2017).
- 555 7 Wang, W., Huang, Y., Wang, Y. & Wang, L. Generalized Autoencoder: A Neural Network
556 Framework for Dimensionality Reduction. *2014 IEEE Conference on Computer Vision and Pattern*
557 *Recognition Workshops*, 496-503, doi:10.1109/CVPRW.2014.79 (2014).
- 558 8 Amodio, M. *et al.* Exploring single-cell data with deep multitasking neural networks. *Nat*
559 *Methods* **16**, 1139-1145, doi:10.1038/s41592-019-0576-7 (2019).

- 560 9 Eraslan, G., Simon, L. M., Mircea, M., Mueller, N. S. & Theis, F. J. Single-cell RNA-seq denoising
561 using a deep count autoencoder. *Nat Commun* **10**, 390, doi:10.1038/s41467-018-07931-2
562 (2019).
- 563 10 Miao, Z. *et al.* Putative cell type discovery from single-cell gene expression data. *Nature*
564 *Methods* **17**, 621-628, doi:10.1038/s41592-020-0825-9 (2020).
- 565 11 Kipf, T. N. & Welling, M. Variational Graph Auto-Encoders. *arXiv e-prints*, arXiv:1611.07308
566 (2016). <<https://ui.adsabs.harvard.edu/abs/2016arXiv161107308K>>.
- 567 12 Wan, C. *et al.* LTMG: a novel statistical modeling of transcriptional expression states in single-
568 cell RNA-Seq data. *Nucleic Acids Res* **47**, e111, doi:10.1093/nar/gkz655 (2019).
- 569 13 Huang, M. *et al.* SAVER: gene expression recovery for single-cell RNA sequencing. *Nature*
570 *Methods* **15**, 539 (2018).
- 571 14 Risso, D., Perraudeau, F., Gribkova, S., Dudoit, S. & Vert, J. P. A general and flexible method for
572 signal extraction from single-cell RNA-seq data. *Nat Commun* **9**, 284, doi:10.1038/s41467-017-
573 02554-5 (2018).
- 574 15 Svensson, V. *et al.* Power analysis of single-cell RNA-sequencing experiments. *Nat Methods* **14**,
575 381-387, doi:10.1038/nmeth.4220 (2017).
- 576 16 Huang, M. *et al.* SAVER: gene expression recovery for single-cell RNA sequencing. *Nat Methods*
577 **15**, 539-542, doi:10.1038/s41592-018-0033-z (2018).
- 578 17 Wang, J. *et al.* Data denoising with transfer learning in single-cell transcriptomics. *Nat Methods*
579 **16**, 875-878, doi:10.1038/s41592-019-0537-1 (2019).
- 580 18 Zhang, L. & Zhang, S. Comparison of Computational Methods for Imputing Single-Cell RNA-
581 Sequencing Data. *IEEE/ACM Trans Comput Biol Bioinform* **17**, 376-389,
582 doi:10.1109/TCBB.2018.2848633 (2020).
- 583 19 Liu, B. *et al.* An entropy-based metric for assessing the purity of single cell populations. *Nature*
584 *Communications* **11**, 3155, doi:10.1038/s41467-020-16904-3 (2020).
- 585 20 Kolodziejczyk, A. A. *et al.* Single Cell RNA-Sequencing of Pluripotent States Unlocks Modular
586 Transcriptional Variation. *Cell Stem Cell* **17**, 471-485, doi:10.1016/j.stem.2015.09.011 (2015).
- 587 21 Klein, A. M. *et al.* Droplet barcoding for single-cell transcriptomics applied to embryonic stem
588 cells. *Cell* **161**, 1187-1201, doi:10.1016/j.cell.2015.04.044 (2015).
- 589 22 Zeisel, A. *et al.* Brain structure. Cell types in the mouse cortex and hippocampus revealed by
590 single-cell RNA-seq. *Science* **347**, 1138-1142, doi:10.1126/science.aaa1934 (2015).
- 591 23 Chung, W. *et al.* Single-cell RNA-seq enables comprehensive tumour and immune cell profiling in
592 primary breast cancer. *Nat Commun* **8**, 15081, doi:10.1038/ncomms15081 (2017).
- 593 24 Grubman, A. *et al.* A single-cell atlas of entorhinal cortex from individuals with Alzheimer's
594 disease reveals cell-type-specific gene expression regulation. *Nat Neurosci* **22**, 2087-2097,
595 doi:10.1038/s41593-019-0539-4 (2019).
- 596 25 Xie, J. *et al.* QUBIC2: a novel and robust biclustering algorithm for analyses and interpretation of
597 large-scale RNA-Seq data. *Bioinformatics* **36**, 1143-1149, doi:10.1093/bioinformatics/btz692
598 (2020).
- 599 26 Bendall, S. C. *et al.* Single-cell trajectory detection uncovers progression and regulatory
600 coordination in human B cell development. *Cell* **157**, 714-725, doi:10.1016/j.cell.2014.04.005
601 (2014).
- 602 27 Wolf, F. A. *et al.* PAGA: graph abstraction reconciles clustering with trajectory inference through
603 a topology preserving map of single cells. *Genome Biol* **20**, 59, doi:10.1186/s13059-019-1663-x
604 (2019).
- 605 28 Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities in
606 large networks. *Journal of Statistical Mechanics: Theory and Experiment* **2008**, P10008,
607 doi:10.1088/1742-5468/2008/10/p10008 (2008).

- 608 29 Li, W. V. & Li, J. J. An accurate and robust imputation method scImpute for single-cell RNA-seq
609 data. *Nat Commun* **9**, 997, doi:10.1038/s41467-018-03405-7 (2018).
- 610 30 Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling for single-
611 cell transcriptomics. *Nat Methods* **15**, 1053-1058, doi:10.1038/s41592-018-0229-2 (2018).
- 612 31 Arisdakessian, C., Poirion, O., Yunits, B., Zhu, X. & Garmire, L. X. DeepImpute: an accurate, fast,
613 and scalable deep neural network method to impute single-cell RNA-seq data. *Genome Biol* **20**,
614 211, doi:10.1186/s13059-019-1837-6 (2019).
- 615 32 Hubert, L. & Arabie, P. Comparing partitions. *J Classif* **2**, 193-218, doi:10.1007/bf01908075
616 (1985).
- 617 33 Rousseeuw, P. J. Silhouettes: A graphical aid to the interpretation and validation of cluster
618 analysis. *J Comput Appl Math* **20**, 53-65, doi:10.1016/0377-0427(87)90125-7 (1987).
- 619 34 Grubman, A. *et al.* A single-cell atlas of entorhinal cortex from individuals with Alzheimer's
620 disease reveals cell-type-specific gene expression regulation. *Nature Neuroscience* **22**, 2087-
621 2097, doi:10.1038/s41593-019-0539-4 (2019).
- 622 35 Tanzi, R. E. The genetics of Alzheimer disease. *Cold Spring Harb Perspect Med* **2**, a006296,
623 doi:10.1101/cshperspect.a006296 (2012).
- 624 36 Su, B. *et al.* Oxidative stress signaling in Alzheimer's disease. *Curr Alzheimer Res* **5**, 525-532,
625 doi:10.2174/156720508786898451 (2008).
- 626 37 Ma, A. *et al.* IRIS3: integrated cell-type-specific regulon inference server from single-cell RNA-
627 Seq. *Nucleic Acids Research*, doi:10.1093/nar/gkaa394 (2020).
- 628 38 Karch, C. M., Ezerskiy, L. A., Bertelsen, S., Alzheimer's Disease Genetics, C. & Goate, A. M.
629 Alzheimer's Disease Risk Polymorphisms Regulate Gene Expression in the ZCWPW1 and the
630 CELF1 Loci. *PLoS one* **11**, e0148717-e0148717, doi:10.1371/journal.pone.0148717 (2016).
- 631 39 Franzén, O., Gan, L.-M. & Björkegren, J. L. M. PanglaoDB: a web server for exploration of mouse
632 and human single-cell RNA sequencing data. *Database* **2019**, doi:10.1093/database/baz046
633 (2019).
- 634 40 Mathys, H. *et al.* Single-cell transcriptomic analysis of Alzheimer's disease. *Nature* **570**, 332-337,
635 doi:10.1038/s41586-019-1195-2 (2019).
- 636 41 Yamakawa, H. *et al.* The Transcription Factor Sp3 Cooperates with HDAC2 to Regulate Synaptic
637 Function and Plasticity in Neurons. *Cell Rep* **20**, 1319-1334, doi:10.1016/j.celrep.2017.07.044
638 (2017).
- 639 42 Boutillier, S. *et al.* Sp3 and sp4 transcription factor levels are increased in brains of patients with
640 Alzheimer's disease. *Neuro-degenerative diseases* **4**, 413-423, doi:10.1159/000107701 (2007).
- 641 43 Hu, Z., Dong, Y., Wang, K. & Sun, Y. Heterogeneous Graph Transformer. *Proceedings of The Web
642 Conference 2020* (2020).
- 643 44 Ma, A., McDermaid, A., Xu, J., Chang, Y. & Ma, Q. Integrative Methods and Practical Challenges
644 for Single-Cell Multi-omics. *Trends in Biotechnology*, doi:10.1016/j.tibtech.2020.02.013 (2020).
- 645 45 Han, A., Glanville, J., Hansmann, L. & Davis, M. M. Linking T-cell receptor sequence to functional
646 phenotype at the single-cell level. *Nature Biotechnology* **32**, 684-692, doi:10.1038/nbt.2938
647 (2014).
- 648 46 Grün, D. Revealing dynamics of gene expression variability in cell state space. *Nat Methods* **17**,
649 45-49, doi:10.1038/s41592-019-0632-3 (2020).
- 650 47 Liu, F. T., Ting, K. M. & Zhou, Z. in *2008 Eighth IEEE International Conference on Data Mining*.
651 413-422.
- 652 48 Lin, P., Troup, M. & Ho, J. W. J. G. b. CIDR: Ultrafast and accurate clustering through imputation
653 for single-cell RNA-seq data. **18**, 59 (2017).
- 654 49 Qiu, X. *et al.* Reversed graph embedding resolves complex single-cell trajectories. *Nat Methods*
655 **14**, 979-982, doi:10.1038/nmeth.4402 (2017).

656 50 Lin, P., Troup, M. & Ho, J. W. CIDR: Ultrafast and accurate clustering through imputation for
657 single-cell RNA-seq data. *Genome Biol* **18**, 59, doi:10.1186/s13059-017-1188-0 (2017).
658 51 Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for
659 interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* **102**, 15545-15550,
660 doi:10.1073/pnas.0506580102 (2005).
661