

1 **scEMB: Learning context representation of genes based on large-scale single-cell**
2 **transcriptomics**

3 Kang-Lin Hsieh^{1,†}, Yan Chu^{2,†}, Xiaoyang Li^{3,4}, Patrick G. Pilié¹, Yulin Dai^{3,*}

4 ¹Department of Genitourinary Medical Oncology, Division of Cancer Medicine, UT MD Anderson
5 Cancer Center, Houston, TX 77030, USA

6 ²Department of Radiation Physics, Division of Radiation Oncology, UT MD Anderson Cancer
7 Center, Houston, TX 77030, USA

8 ³Center for Precision Health, McWilliams School of Biomedical Informatics, The University of
9 Texas Health Science Center at Houston, Houston, TX 77030, USA

10 ⁴Department of Biostatistics and Data Science, School of Public Health, The University of Texas
11 Health Science Center at Houston, Houston, TX 77030, USA

12

13 *To whom correspondence should be addressed:

14 Yulin Dai, Ph.D.

15 Center for Precision Health

16 McWilliams School of Biomedical Informatics

17 The University of Texas Health Science Center at Houston

18 7000 Fannin St. Suite E760A Houston, TX 77030

19 Phone: 713-500-3462

20 Email: Yulin.Dai@uth.tmc.edu

21 **ABSTRACT**

22 Background: The rapid advancement of single-cell transcriptomic technologies has led to the
23 curation of millions of cellular profiles, providing unprecedented insights into cellular
24 heterogeneity across various tissues and developmental stages. This growing wealth of data
25 presents an opportunity to uncover complex gene-gene relationships, yet also poses significant
26 computational challenges.

27 Results: We present scEMB, a transformer-based deep learning model developed to capture
28 context-aware gene embeddings from large-scale single-cell transcriptomics data. Trained on
29 over 30 million single-cell transcriptomes, scEMB utilizes an innovative binning strategy that
30 integrates data across multiple platforms, effectively preserving both gene expression
31 hierarchies and cell-type specificity. In downstream tasks such as batch integration, clustering,
32 and cell type annotation, scEMB demonstrates superior performance compared to existing
33 models like scGPT and Geneformer. Notably, scEMB excels *in silico* correlation analysis,
34 accurately predicting gene perturbation effects in CRISPR-edited datasets and microglia state
35 transition, identifying a few known Alzheimer's disease (AD) risks genes in top gene list.
36 Additionally, scEMB offers robust fine-tuning capabilities for domain-specific applications,
37 making it a versatile tool for tackling diverse biological problems such as therapeutic target
38 discovery and disease modeling.

39 Conclusions: scEMB represents a powerful tool for extracting biologically meaningful insights
40 from complex gene expression data. Its ability to model *in silico* perturbation effects and conduct
41 correlation analyses in the embedding space highlights its potential to accelerate discoveries in
42 precision medicine and therapeutic development.

43 **Keywords:** single-cell transcriptomics, Transformer, gene-gene relationship, *in silico*
44 perturbation analysis, *in silico* correlation analysis

45

46 INTRODUCTION

47 The rapid advancement of single-cell technologies has enabled international consortia, including
48 the Human Cell Atlas (HCA)¹, Tabula Sapiens², Human Cell Landscape³, and the more recent
49 95 million CELLxGENE corpus⁴, to characterize the cellular heterogeneity in major human
50 tissues and organ systems in different developmental stages. However, managing and
51 extracting meaning and value from such large-scale data is highly challenging and this technical
52 gap creates opportunities for researchers to innovate and develop advanced computational
53 methods capable of harnessing the full potential of these vast datasets, enabling deeper
54 biological understanding and practical applications.

55 Clustered regularly interspaced short palindromic repeats (CRISPR) system is a genetic
56 engineering tool for gene editing that can activate (CRISPRa), repress (CRISPRi), or modify
57 gene expressions for specific genes⁵. This technology allows us to understand the systematic
58 effects of perturbing one or multiple genes in different cell types, which provides ideal data to
59 learn the relationship between different perturbations^{6,7}. However, genome-scale gene
60 perturbations and their potential combinations have been impeded by the cost, emphasizing the
61 need for novel cost-efficient methods to prioritize their effects.

62 Recent advances in artificial intelligence (AI), such as variational autoencoder (VAE) models⁸⁻¹⁰,
63 which explicitly learn low-dimensional embeddings from single-cell transcriptome to capture the
64 biological meaning embedding. Moreover, these embedding-based methods¹¹ [Autoencoder¹²,
65 VAE¹³⁻¹⁶, GNN¹⁷], has been widely used to along with the cellular perturbation tasks along with
66 single-cell CRISPR data, to compress the transcriptomics data (mainly single-cell) into the latent
67 space and compare the effects of different perturbations and possess a more accurate
68 molecular measurement as well as in silico-induced perturbations^{18,19}. However, these
69 methods^{18,19} were mainly task-specific and the perturbation effect were encoded on different
70 embedding systems⁹, hindering the transfer learning performance.

71 More recently, large-language models (LLM) such as GPT^{20,21}, have revolutionized various fields
72 by leveraging deep neural networks trained on vast text datasets. These models generalize
73 knowledge from pretraining, enabling efficient task transfer with minimal data. The self-attention
74 mechanism enhances their ability to focus on relevant information, improving predictions across
75 diverse applications.

76 In single-cell data analysis, transformer-based models offer a promising approach to
77 overcoming batch effects through batch-unaware pretraining, which has proven resilient to
78 some technical artifacts. For instance, models like Universal Cell Embeddings (UCE)²² and
79 GeneCompass²³ have integrated molecular profiles across studies, tissues, and species,
80 enabling cell type annotations to transfer across previously unseen species. More recent
81 Geneformer²⁵, scGPT²⁴, and scFoundation²⁴ trained the foundation model from large-scale
82 single-cell corpus and offered transfer learning features²⁵ to makes the transformer an ideal
83 backbone for inferring unseen perturbations^{17,26,27}, showing great promise to mitigate previous
84 limitations.

85 Despite significant advances, a key gap remains in applying transformer models specifically to
86 single-cell data. While transformers have shown promise in handling large-scale transcriptomic
87 data, optimization for single-cell applications is still needed. To address this, we aim to develop
88 the scEMB model to explore broader application scenarios in single-cell foundation models and
89 establish best practices across diverse datasets. Additionally, we will investigate the potential of
90 in silico perturbation prediction using both zero-shot and fine-tuned models, enabling more
91 accurate biological insights and expanding applications in disease modeling and therapeutic
92 discovery.

93

94 **RESULT**

95 **Overview of scEMB**

96 scEMB is an attention-based deep learning model pretrained on large-scale transcriptomic data
97 to capture complex biological networks. It uses self-attention to focus on the most critical genes
98 expressed in each single cell, optimizing predictive accuracy through various learning objectives.
99 scEMB is designed to capture complex biologically meaning alteration underlying cell type state
100 transition and perturbation response (Fig. 1).

101 scEMB represents each single cell's transcriptome as a rank-binned gene expression encoding,
102 ranking genes by their expression within individual cells. While this rank-based approach has
103 limitations, such as not fully utilizing precise gene expression measurements from transcript
104 counts, it provides a non-parametric representation of each cell's transcriptome. This method
105 leverages the vast observations of gene expression across the cell-x-gene 30M dataset to
106 highlight genes that characterize cell states by composing cellular biological networks.

107 The rank value encoding of each single cell's transcriptome then passes through twelve
108 transformer encoder units, each consisting of a self-attention layer and a feed-forward neural
109 network layer. Pretraining utilized a masked learning objective, a technique proven in other
110 information domains to enhance the generalizability of foundational knowledge acquired during
111 pretraining, for a broad spectrum of downstream fine-tuning tasks such as batch integration, cell
112 type annotation, *in silico* perturbation, and correlation analysis.

113

114 **scEMB shows comparable performance in clustering and batch integration**

115 As the standard tasking in single-cell foundation model, we tested the performance in single-cell
116 clustering using scEMB, scGPT, Geneformer, and unintegrated methods on PBMC 10k
117 dataset²⁸. Fig. 2 shows the PBMC 10k dataset before and after dimensionality reduction using
118 UMAP (Uniform Manifold Approximation and Projection), with the legend indicating various cell
119 types, including B cells, CD4 T cells, CD8 T cells, CD14+ Monocytes, Dendritic Cells,

120 FCGR3A+ Monocytes, Megakaryocytes, NK cells, and other. For batch integration and cell type
121 clustering tasks, we present UMAP plots that illustrate cell clustering across different datasets
122 using four methods, with distinct colors indicating various cell types. Below the UMAP plots, a
123 performance comparison table presents various metrics for each method, including isolation
124 score, kBET score, and batch mixing entropy. The table also includes aggregate scores for
125 metrics such as batch correction, biological conservation, and overall performance. This
126 comprehensive visualization enables a direct comparison of how well each method integrates
127 data across batches while preserving biological information.

128

129 **High consistent performance in cell type annotation task**

130 One of the main advantages of foundation models is their ability to leverage pretraining on
131 large-scale general datasets, enabling fine-tuning for a wide range of downstream tasks, even
132 when the available task-specific data is sparse to generate meaningful predictions. We
133 evaluated scEMB's performance on cell-type annotation tasks by training a supervised learning
134 model on cell-type annotations from a reference dataset, then predicting cell types in an
135 independent, unseen dataset using two batches of data from MTG brain region dataset²⁹.

136 After inputting the transcriptome into scEMB Encoder, the cell-level embedding generated by
137 scEMB, representing a specific cellular state, can be used to infer cell type. The cell type
138 annotation was trained in one batch of the brain MTG dataset, and test on the other batch (Fig.
139 3a). To evaluate cell type prediction accuracy, we used confusion matrices for scEMB, scGPT,
140 and Geneformer (Fig. 3b-3d). While all methods performed well, we observed slight differences,
141 especially for less common cell types. This visualization offers both qualitative (UMAP) and
142 quantitative (confusion matrix) insights into scEMB's effectiveness in predicting cell types from
143 single-cell RNA sequencing data. Overall, all foundation models demonstrated strong potential

144 in accurately annotating cell types using well-annotated reference data, a critical step in single-
145 cell data analysis.

146

147 **scEMB *in silico* perturbation analysis identified genes highly consistent with CRISPRi**
148 **data**

149 As demonstrated by Geneformer and scGPT, single-cell foundation models show great potential
150 in predicting cellular gene expression responses to specific gene perturbations. This highlights a
151 major advantage of transfer learning, which leverages biological knowledge from millions of
152 human cells to improve accuracy and scalability. To investigate potential deleterious effects in
153 response to cellular perturbation, we designed an *in silico* perturbation response prediction task,
154 which was validated using a single-cell CRISPRi dataset with known perturbation outcomes (Fig.
155 4a). Specifically, we measured these effects using a cosine similarity score calculated from the
156 scEMB-encoded embeddings in an iPSC-derived microglia dataset³⁰. Among the 39 perturbed
157 gene conditions (Fig. 4b), we identified the top 10 most accurately predicted *in silico* perturbed
158 genes, which were validated against ground truth CRISPRi data. Among them, we identified a
159 few well know microglia function genes, such CSF1R (Colony Stimulating Factor 1 Receptor),
160 which is critical for the survival, differentiation, and function of microglia³¹. TGFBR1/2
161 (Transforming Growth Factor Beta Receptor 1/2) TGFBR1/2 are involved in the TGF- β signaling
162 pathway, which is essential for microglial homeostasis and modulation of their inflammatory
163 response. TGF- β helps maintain the quiescent state of microglia in the healthy brain but also
164 plays a role in activating microglia during injury or neurodegenerative processes³². Interestingly,
165 the housekeeping gene like AARS (Alanyl-tRNA Synthetase)³³ demonstrated the highest cosine
166 similarity in predictions, indicating that scEMB may perform well in predicting housekeeping
167 genes. However, this result also suggests that the universal expression of housekeeping genes

168 across cells could be more complex, as highlighted by findings from the single-cell foundation
169 model²⁷ .

170

171 **scEMB *in silico* correlation analysis highlighted AD risk genes that might contribute to**
172 **the AD pathogenesis**

173 As there is still limited real world perturbation data existing, not to mention most perturbation
174 response data or conditions is not included in training in current foundation models^{24,27,34} .

175 Therefore, perturbation response predicted from zero-shot learning still requires careful use. We
176 designed a finetune task based alternative approach as revealed in Fig. 5. scEMB *in silico*
177 correlation analysis could be applied to identify potential reversed relationship at cellular
178 embedding level between cellular state alteration effect and perturbation effect among single-
179 cell transcriptome data. Here, we first project the cellular state alteration embeddings using
180 scEMB Encode, focusing on microglia from individuals with Alzheimer's Disease (AD) and
181 cognitively normal (CN) individuals from the Religious Orders Study/Memory and Aging Project
182 (ROS/MAP) cohort³⁶. Then, we followed our previously perturbation strategy^{10,17} to integrate
183 pretrained GO Gene GNN node embeddings, allowing us to propagate the impact of
184 perturbations from drug treatment or iPSC-derived CRISPRi data³⁰. As illustrated in Fig. 5b, the
185 violin plot presents the top 15 absolute cosine similarity scores, highlighting cellular state
186 alterations in microglia by comparing cells from individuals with Alzheimer's Disease (AD) to
187 those from cognitively normal controls. Additionally, the plot reveals the effects of gene
188 perturbations on iPSC-derived microglia. Notably, several well-known AD risk genes from
189 genome-wide association study³⁶, including *PLCG2*, *SORL1*, and *TREM2*, emerged among the
190 top genes. These findings suggest that the scEMB *in silico* correlation analysis may offer
191 valuable insights into genes involved in disease pathogenesis, providing clues about underlying
192 molecular mechanisms. Furthermore, it points to potential therapeutic targets, as the identified

193 genes could play a pivotal role in modulating microglial function and influencing disease
194 progression in AD.

195

196

197 **METHODS AND MATERIALS**

198 **scEMB architecture**

199 scEMB comprises twelve transformer encoder units, each containing a self-attention layer and a
200 feed-forward neural network layer. The model's key parameters include an input size of 2,048,
201 embedding dimensions of 768, 12 attention heads, and a feed-forward network size of 3,072.
202 The input size was maximized to capture the most context possible through full attention, based
203 on the typical number of genes detected in each cell in the pretraining dataset. To speed up the
204 training process for this large dataset, scEMB employed Scaled Dot-Product Attention (SDPA)
205 across the entire 2,048 input size. The model's depth was determined by the maximum level for
206 which sufficient pretraining data was available. Other minor parameters include the Gaussian
207 Error Linear Unit as the activation function, a dropout probability of 0.1 for fully connected layers,
208 and a dropout ratio of 0.1 for attention probabilities. The weight matrices were initialized with a
209 standard deviation of 0.02, and an epsilon value of 1×10^{-12} was used for layer normalization.
210 The model was built in PyTorch and utilized the Huggingface Transformers library for
211 configuration, data loading, and training.

212

213 **scEMB pretraining**

214 The primary pretraining objective of scEMB is masked language modeling. This approach,
215 proven effective in various fields, enhances the generalizability of foundational knowledge
216 acquired during pretraining, benefiting a wide range of downstream fine-tuning objectives.

217 During pretraining, 15% of the genes within each transcriptome are masked. The model then
218 learns to predict which gene should occupy each masked position in that specific cell state,
219 using the context provided by the remaining unmasked genes. This self-supervised approach
220 can be accomplished on completely unlabeled data, allowing the inclusion of large amounts of
221 training data without being restricted to samples with accompanying labels.

222 The pretraining process utilized optimized hyperparameters to enhance model performance.
223 These included a maximum learning rate of 1×10^{-4} , a linear learning scheduler with warmup,
224 and the Adam optimizer with weight decay fix. Additionally, 10,000 warmup steps were
225 employed, along with a weight decay of 0.001 and a batch size of 8. These carefully selected
226 parameters contributed to the model's effective training. During model training, we adapted a
227 custom tokenizer from the Huggingface Transformers library to implement dynamic, length-
228 grouped padding. This approach minimized computation on padding and achieved a 29.4x
229 speedup in pretraining. The process involves randomly sampling a megabatch, then ordering
230 minibatches by their length in descending order. These minibatches are dynamically padded,
231 reducing wasted computation on padding due to their grouped lengths. The distributed training
232 is implemented by Deepspeed, which partitions parameters, gradients, and optimizer states
233 across available GPUs. Overall, pretraining was achieved in approximately 10 days, distributed
234 across one node with 8 Nvidia H100 96GB GPUs.

235

236 **scEMB finetuning**

237 Fine-tuning scEMB involves initializing the model with pretrained scEMB weights and adding a
238 final task-specific transformer layer. This process aims to recalibrate the distribution between
239 the pre-trained model and the task-specific dataset. The recalibration can be divided into two
240 types: 1) using a small dataset to adjust the layers within the foundation model, or 2) fine-tuning
241 external knowledge or embeddings, generated either from a task-specific dataset or expert

242 knowledge, to feed into the foundation model. In the first type of fine-tuning, the number of
243 frozen layers must be carefully considered. Based on our experience and suggestions from
244 relevant literature, applications more closely related to the pretraining objective benefit from
245 freezing more layers. This prevents overfitting to limited task-specific data. Conversely,
246 applications that diverge more from the pretraining objective benefit from fine-tuning more layers,
247 optimizing performance for the new task. In the second type of fine-tuning, the extraction of
248 features and appropriate use of expert knowledge still need comprehensive discussion. The
249 task-specific dataset underwent the same preprocessing pipeline as scEMB. To demonstrate
250 the effectiveness of the pretrained scEMB model in enhancing the predictive performance of
251 downstream fine-tuning tasks, we employed consistent fine-tuning hyperparameters across all
252 applications. Special attention must be given to determining the appropriate number of frozen
253 layers during the fine-tuning process.

254 Several applications of scEMB have been demonstrated, including cell type annotations, batch
255 corrections, and perturbation effects. These applications further test the two types of
256 recalibration. For cell type annotations and batch corrections, we used the PBMC dataset. In
257 both subtasks, we tested LoRA (Low-Rank Adaptation) for the first type of fine-tuning and Prefix
258 Tuning for the second type. The key difference is that LoRA aligns the model with a new dataset,
259 while Prefix Tuning aligns the external knowledge vector, typically called an embedding, with a
260 new dataset. Consequently, the cell representation from LoRA fine-tuning uses cell gene
261 expression, whereas Prefix Tuning combines external knowledge and cell gene expression. The
262 cell representation is the average of hidden layer outputs for all 12 layers. The benchmarks for
263 cell type annotation were AUROC, F1 score, and accuracy. For batch corrections, we used the
264 following metrics: Isolated labels, KMeans NMI, KMeans ARI, Silhouette label, cLISI, Silhouette
265 batch, iLISI, KBET, Graph connectivity, and PCR comparison.

266

267 ***in silico* perturbation analysis**

268 To evaluate the model's ability to capture gene perturbation effects, we conducted an *in silico*
269 perturbation analysis, focusing on cosine similarity between cell embeddings of *in silico*
270 perturbed control cells and CRISPR-edited cells. Cell embeddings, representing each cell as a
271 768-dimensional vector, were obtained by averaging the encoder output. For control cells, we
272 masked the CRISPR-targeted gene, then calculated the cosine similarity between control and
273 perturbed groups. A high cosine similarity suggests a minimal effect on the gene, indicating less
274 significant perturbation. Conversely, lower similarity highlights more impactful gene changes.
275 This evaluation method provides insight into the genes most significantly affected by CRISPR
276 perturbations. The results, as shown in Fig. 4, demonstrate how cosine similarity can effectively
277 reflect gene impact, where high similarity indicates low gene perturbation impact and vice versa.
278 This approach showcases the model's precision in detecting subtle yet significant gene
279 expression changes after perturbations.

280

281 ***in silico* correlation analysis**

282 To further investigate the model's capability to recognize the similarity between disease state
283 transitions and gene perturbation effects, we conducted an *in silico* correlation analysis. Single-
284 cell transcriptomes from non-targeting controls (NTC) in the Alzheimer's Disease (AD) dataset
285 were tokenized and processed through the scEMB encoder, producing 768-dimensional cell
286 embeddings. These embeddings were combined with gene perturbation embeddings (emb_p),
287 which were derived from a protein-protein interaction (PPI)-guided gene relationship graph,
288 informed by gene ontology (GO). The graph utilized a 3-layer graph attention network (GAT) to
289 conduct self-supervised learning, enabling the extraction of embeddings for each gene^{10,17}.
290 These combined embeddings were then passed through a two-layer Multilayer Perceptron
291 (MLP), which linked the cell embeddings to gene features, allowing the model to predict altered
292 gene expression profiles. Additionally, the *in silico* correlation analysis was employed to assess

293 the relationship between cellular state alterations and gene perturbation effects. This analysis
294 aimed to identify potential therapeutic targets by capturing the similarity or reversed
295 relationships between disease-related cellular changes and the impact of specific gene
296 perturbations. To evaluate these relationships, we performed a cosine similarity analysis,
297 quantifying the similarity between control and AD cellular states and the effects of gene
298 perturbations on induced pluripotent stem cell (iPSC)-derived microglia. The results, presented
299 in Fig. 5, highlight the AD risk genes, allowing for a better understanding of how perturbations
300 influence the disease state at the cellular level. This approach offers a robust framework for
301 studying perturbation effects and identifying key genes involved in disease mechanisms.

302

303 **Benchmark with other models**

304 To comprehensively evaluate the performance of scEMB, we compared its clustering and batch
305 integration capabilities against two other single-cell foundation models, Geneformer²⁷ and
306 scGPT³⁴, as well as the standard single-cell data analysis approach in Scanpy³⁷. These
307 comparisons were conducted in a zero-shot setting using the PBMC 10k dataset²⁸ with default
308 configurations. Additionally, we assessed cell type annotation performance in a fine-tuning
309 scenario by conducting a classification task on a brain microglia dataset³⁶, comparing scEMB to
310 Geneformer and scGPT.

311

312 **Data preprocess**

313 ***30M single-cell transcriptome corpus curation***

314 For pretraining the whole-human foundation model, we sourced data through the Census API
315 from the CELLxGENE portal³⁸ (<https://cellxgene.cziscience.com/>), which provides regularly
316 updated datasets (accessed on May 9, 2024). We included both scRNA-seq and snRNA-seq
317 sequencing protocols, focusing on samples from healthy conditions. To ensure data quality, we

318 filtered out cells expressing fewer than 200 genes or mitochondria gene expression percent >
319 10 using python library Scanpy³⁷. After applying these filters, the final dataset comprised
320 sequencing data from 30 million cells.

321

322 ***PBMC 10k dataset***

323 The PBMC 10k dataset consists of two single-cell RNA sequencing batches of peripheral blood
324 mononuclear cells (PBMCs) obtained from a healthy donor. This dataset was reanalyzed by
325 Gayoso et al.²⁸, identifying 3,346 differentially expressed genes. The first batch consists of
326 7,982 cells, and the second batch contains 4,008 cells. Cell type annotations were conducted
327 using the R package Seurat³⁹, categorized the cells into nine distinct groups. The preprocessed
328 data was adopted from scGPT³⁴ (<https://github.com/bowang-lab/scGPT>, accessed on June 16,
329 2024).

330

331 ***MTG Brain dataset***

332 We incorporated two brain samples from the middle temporal gyrus (MTG) region, provided by
333 the Seattle Alzheimer's Disease Brain Cell Atlas (SEA-AD) consortium via Amazon AWS Bucket
334 (accessed on August 9, 2024). These two samples (H19.33.004 and H19.30.001) were profiled
335 from two different batches. Cell type annotations were derived from the original study²⁹.

336

337 ***Microglia from ROS/MAP snRNA-seq and polygenic risk score process***

338 We used snRNA-seq data from the Synapse portal (syn2580853, accessed on April 15, 2023),
339 which includes 454 participants from the Religious Orders Study/Memory and Aging Project
340 (ROS/MAP) cohort³⁶. Matched whole-genome sequencing (WGS) data were obtained from
341 Synapse (syn11724057, accessed on November 10, 2022). Individual genetic risk was
342 estimated using LDpred2⁴⁰, based on variant effect sizes from Wightman et al.'s GWAS
343 summary statistics⁴¹. This provided a dataset of 407 individuals with both snRNA-seq and WGS
344 data. To focus on individuals with high polygenic risk scores (PRS) for Alzheimer's disease, we
345 selected the top 20% PRS, yielding 53 high-risk AD cases and 15 high-risk cognitively normal
346 individuals⁴². Data was processed using R package Seurat³⁹ with filters applied for “percent.mt
347 <= 50” and “nFeature_RNA > 200,” and cell type annotations were adopted from the original
348 study.

349

350 ***CRISPRi iPSC-derived microglia dataset***

351 The CRISPR iPSC-derived microglia dataset was downloaded from GSE178317³⁰ (accessed on
352 March 10, 2023) and processed using the CROP-seq pipeline⁴³. The gene feature matrix was
353 obtained by alignment and gene expression quantification for scRNA-seq libraries and sgRNA-
354 enriched libraries using Cell Ranger and STAR⁴⁴. 2) The sgRNA matrix was assigned based on
355 the demuxEM algorithm⁴⁵ and a modified z-score cut-off method⁴⁶. We categorized the cells into
356 two groups: those with a single sgRNA and those with two sgRNAs. In total, we identified 39
357 single sgRNA-targeted genes, each captured by at least 200 cells.

358

359 **DISCUSSION**

360 In this study, we introduce scEMB, a novel single-cell transcriptome foundation model
361 developed to extract complex gene-gene interaction information from a large corpus of 30

362 million human single-cell transcriptomes. A key feature of scEMB is its specially designed
363 tokenization mechanism, which was engineered to generate gene embeddings that are both
364 stable and sensitive. This design allows scEMB to be highly effective across a broad spectrum
365 of downstream tasks, from gene interaction analysis to large-scale biological modeling.

366 scEMB's performance has been rigorously evaluated in tasks such as clustering and batch
367 integration, and it demonstrated competitive zero-shot performance across multiple datasets.
368 Remarkably, its performance was on par with leading foundation models like Geneformer and
369 scGPT, reinforcing scEMB's utility as a powerful tool for a variety of biological applications.

370 One of the novel contributions of this work is the exploration of scEMB's potential in perturbation
371 response prediction, using an Alzheimer's Disease (AD) dataset as a case study. We estimated
372 the cosine similarity between perturbation effects and cellular state transitions, demonstrating
373 how scEMB can model complex cellular responses to external stimuli or genetic modifications.
374 This ability to predict perturbation effects holds significant promise for the study of disease
375 mechanisms and therapeutic interventions.

376 In our evaluation of fine-tuning approaches, we identified a critical gap in current practices: the
377 lack of a standardized method for fine-tuning large models in biological contexts, particularly in
378 classification tasks. Existing classification tasks often fail to capture the complexity of biological
379 research needs, limiting the practical utility of these models. To address this challenge, we
380 introduced alternative fine-tuning strategies, including Prefix-tuning and LoRA (Low-Rank
381 Adaptation), which are specifically tailored for downstream applications such as cell-type
382 annotation and *in silico* perturbation analysis. These methods offer greater flexibility and more
383 accurate biological insights, especially in tasks that involve nuanced gene expression changes
384 across different conditions.

385 Furthermore, based on observations from models like Geneformer and scGPT, we confirmed
386 that scEMB's performance adheres to the scaling law principle. This principle suggests that
387 pretraining on larger and more diverse corpora consistently enhances predictive power. As
388 scEMB was pretrained on hundreds of experimental datasets, it encountered various batch
389 effects, technical artifacts, and individual variability during training, which ultimately improved its
390 robustness and generalization across datasets. Larger pretraining corpora allowed scEMB to
391 develop deeper, more predictive models capable of addressing the complexity of real-world
392 biological data.

393 Finally, the introduction of *in silico* perturbation analysis opens new avenues for using scEMB in
394 drug response prediction and cellular behavior modeling. scEMB's ability to predict how cells
395 will respond to various perturbations positions it as a powerful tool for future therapeutic
396 discovery and precision medicine applications.

397 Moving forward, we plan to expand scEMB's capabilities further. Leveraging generative
398 modeling, scEMB can implicitly capture gene-gene interactions through its embeddings and
399 attention maps, enabling the exploration of Gene Regulatory Networks (GRNs). We propose
400 GRN inference workflows that utilize both pretrained and fine-tuned versions of scEMB, where
401 the gene embeddings reflect dataset-level interactions, and attention maps reveal specific gene
402 activation patterns across diverse cell states. By validating these inferred networks against
403 established biological data, we demonstrated scEMB's potential for gene program discovery
404 and its ability to uncover previously unknown regulatory pathways.

405

406 **DATA AVAILABILITY**

407 All the data generated or analyzed in this study is available from the authors upon reasonable
408 request.

409

410 **ACKNOWLEDGEMENTS**

411 **FUNDING**

412 This research was partially supported by National Institutes of Health grants awarded to Y.D

413 R21AG087299.

414

415

416 REFERENCES

- 417 1. Regev, A. *et al.* The Human Cell Atlas. *Elife* **6**, (2017).
- 418 2. Tabula Sapiens Consortium* *et al.* The Tabula Sapiens: A multiple-organ, single-cell
419 transcriptomic atlas of humans. *Science* **376**, eabl4896 (2022).
- 420 3. Han, X. *et al.* Construction of a human cell landscape at single-cell level. *Nature* **581**, 303–
421 309 (2020).
- 422 4. Chanzuckerberg Initiative. CZ CELLxGENE Discover. assessed by 8/6/2023.
- 423 5. Adli, M. The CRISPR tool kit for genome editing and beyond. *Nat. Commun.* **9**, 1911 (2018).
- 424 6. Faial, T. Single-cell CRISPR screen for GWAS loci. *Nature genetics* vol. 55 904 (2023).
- 425 7. Bock, C. *et al.* High-content CRISPR screening. *Nat. Rev. Methods Primers* **2**, (2022).
- 426 8. Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling for
427 single-cell transcriptomics. *Nat. Methods* **15**, 1053–1058 (2018).
- 428 9. Lotfollahi, M. *et al.* Mapping single-cell data to reference atlases by transfer learning. *Nat.*
429 *Biotechnol.* **40**, 121–130 (2022).
- 430 10. Hsieh, K.-L., Chu, Y., Pilié, P. G., Zhang, K. & Dai, Y. Learning interpretable cellular
431 embedding for inferring biological mechanisms underlying single-cell transcriptomics.
432 *bioRxiv* (2024) doi:10.1101/2024.03.29.24305092.
- 433 11. Pan, X. *et al.* Deep learning for drug repurposing: methods, databases, and applications.
434 *arXiv [q-bio.BM]* (2022).
- 435 12. Jia, P. *et al.* Deep generative neural network for accurate drug response imputation. *Nat.*
436 *Commun.* **12**, 1740 (2021).
- 437 13. Lotfollahi, M., Wolf, F. A. & Theis, F. J. scGen predicts single-cell perturbation responses.
438 *Nat. Methods* **16**, 715–721 (2019).

- 439 14. Rampášek, L., Hidru, D., Smirnov, P., Haibe-Kains, B. & Goldenberg, A. Dr.VAE: improving
440 drug response prediction via modeling of drug perturbation effects. *Bioinformatics* **35**,
441 3743–3751 (2019).
- 442 15. Chen, J. *et al.* Deep transfer learning of cancer drug responses by integrating bulk and
443 single-cell RNA-seq data. *Nat. Commun.* **13**, 6494 (2022).
- 444 16. Lotfollahi, M., Naghipourfar, M., Theis, F. J. & Wolf, F. A. Conditional out-of-distribution
445 generation for unpaired data using transfer VAE. *Bioinformatics* **36**, i610–i617 (2020).
- 446 17. Roohani, Y., Huang, K. & Leskovec, J. Predicting transcriptional outcomes of novel
447 multigene perturbations with GEARS. *Nat. Biotechnol.* (2023) doi:10.1038/s41587-023-
448 01905-6.
- 449 18. Lotfollahi, M. *et al.* Predicting cellular responses to complex perturbations in high-
450 throughput screens. *Mol. Syst. Biol.* **19**, e11517 (2023).
- 451 19. Hetzel, L. *et al.* Predicting cellular responses to novel drug perturbations at a single-cell
452 resolution. *arXiv [cs.LG]* (2022).
- 453 20. Brown, T. B. *et al.* Language Models are Few-Shot Learners. *arXiv [cs.CL]* (2020).
- 454 21. Bubeck, S. *et al.* Sparks of Artificial General Intelligence: Early experiments with GPT-4.
455 *arXiv [cs.CL]* (2023).
- 456 22. Rosen, Y. *et al.* Universal cell embeddings: A foundation model for cell biology. *bioRxiv*
457 (2023) doi:10.1101/2023.11.28.568918.
- 458 23. Yang, X. *et al.* GeneCompass: Deciphering universal gene regulatory mechanisms with
459 knowledge-informed cross-species foundation model. (2023)
460 doi:10.1101/2023.09.26.559542.
- 461 24. Hao, M. *et al.* Large-scale foundation model on single-cell transcriptomics. *Nat. Methods* **21**,
462 1481–1491 (2024).
- 463 25. Hosna, A. *et al.* Transfer learning: a friendly introduction. *J. Big Data* **9**, 102 (2022).

- 464 26. Cui, H. *et al.* ScGPT: Towards building a foundation model for single-cell multi-omics using
465 generative AI. *bioRxiv* 2023.04.30.538439 (2023) doi:10.1101/2023.04.30.538439.
- 466 27. Theodoris, C. V. *et al.* Transfer learning enables predictions in network biology. *Nature* **618**,
467 616–624 (2023).
- 468 28. Gayoso, A. *et al.* A Python library for probabilistic analysis of single-cell omics data. *Nat.*
469 *Biotechnol.* **40**, 163–166 (2022).
- 470 29. Gabitto, M. I. *et al.* Integrated multimodal cell atlas of Alzheimer’s disease. *Res. Sq.* (2023)
471 doi:10.21203/rs.3.rs-2921860/v1.
- 472 30. Dräger, N. M. *et al.* A CRISPRi/a platform in human iPSC-derived microglia uncovers
473 regulators of disease states. *Nat. Neurosci.* **25**, 1149–1162 (2022).
- 474 31. Elmore, M. R. P. *et al.* Colony-stimulating factor 1 receptor signaling is necessary for
475 microglia viability, unmasking a microglia progenitor cell in the adult brain. *Neuron* **82**, 380–
476 397 (2014).
- 477 32. Zöller, T. *et al.* Silencing of TGF β signalling in microglia results in impaired homeostasis.
478 *Nat. Commun.* **9**, 1–13 (2018).
- 479 33. Hounkpe, B. W., Chenou, F., de Lima, F. & De Paula, E. V. HRT Atlas v1.0 database:
480 redefining human and mouse housekeeping genes and candidate reference transcripts by
481 mining massive RNA-seq datasets. *Nucleic Acids Res.* **49**, D947–D955 (2021).
- 482 34. Cui, H. *et al.* scGPT: toward building a foundation model for single-cell multi-omics using
483 generative AI. *Nat. Methods* **21**, 1470–1480 (2024).
- 484 35. Fujita, M. *et al.* Cell subtype-specific effects of genetic variation in the Alzheimer’s disease
485 brain. *Nat. Genet.* **56**, 605–614 (2024).
- 486 36. Schwartzenuber, J. *et al.* Genome-wide meta-analysis, fine-mapping and integrative
487 prioritization implicate new Alzheimer’s disease risk genes. *Nat. Genet.* **53**, 392–402 (2021).
- 488 37. Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data
489 analysis. *Genome Biol.* **19**, 15 (2018).

- 490 38. CZI Single-Cell Biology Program *et al.* CZ CELLxGENE Discover: A single-cell data
491 platform for scalable exploration, analysis and modeling of aggregated data. (2023)
492 doi:10.1101/2023.10.30.563174.
- 493 39. Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573-3587.e29
494 (2021).
- 495 40. Privé, F., Arbel, J. & Vilhjálmsson, B. J. LDpred2: better, faster, stronger. *Bioinformatics* **36**,
496 5424–5431 (2021).
- 497 41. Wightman, D. P. *et al.* A genome-wide association study with 1,126,563 individuals
498 identifies new risk loci for Alzheimer’s disease. *Nat. Genet.* **53**, 1276–1282 (2021).
- 499 42. Li, X. *et al.* Genetically-regulated pathway-polygenic risk score (GRPa-PRS): A risk
500 stratification method to identify genetically regulated pathways in polygenic diseases.
501 *medRxiv* 2023.06.19.23291621 (2023) doi:10.1101/2023.06.19.23291621.
- 502 43. Hill, A. J. *et al.* On the design of CRISPR-based single-cell molecular screens. *Nat.*
503 *Methods* **15**, 271–274 (2018).
- 504 44. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
- 505 45. Gaublomme, J. T. *et al.* Nuclei multiplexing with barcoded antibodies for single-nucleus
506 genomics. *Nat. Commun.* **10**, 2907 (2019).
- 507 46. Tian, R. *et al.* CRISPR interference-based platform for multimodal genetic screens in
508 human iPSC-derived neurons. *Neuron* **104**, 239-255.e12 (2019).

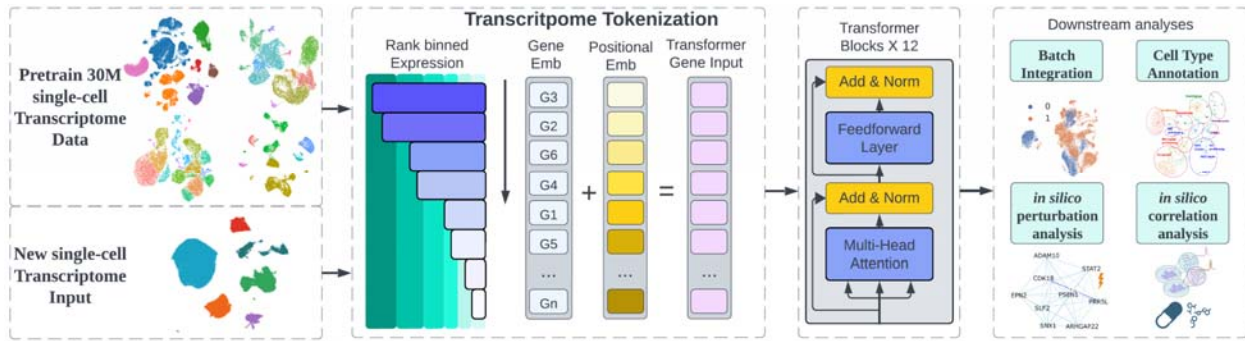
509

510

511

512

513 **FIGURES AND FIGURE LEGENDS**

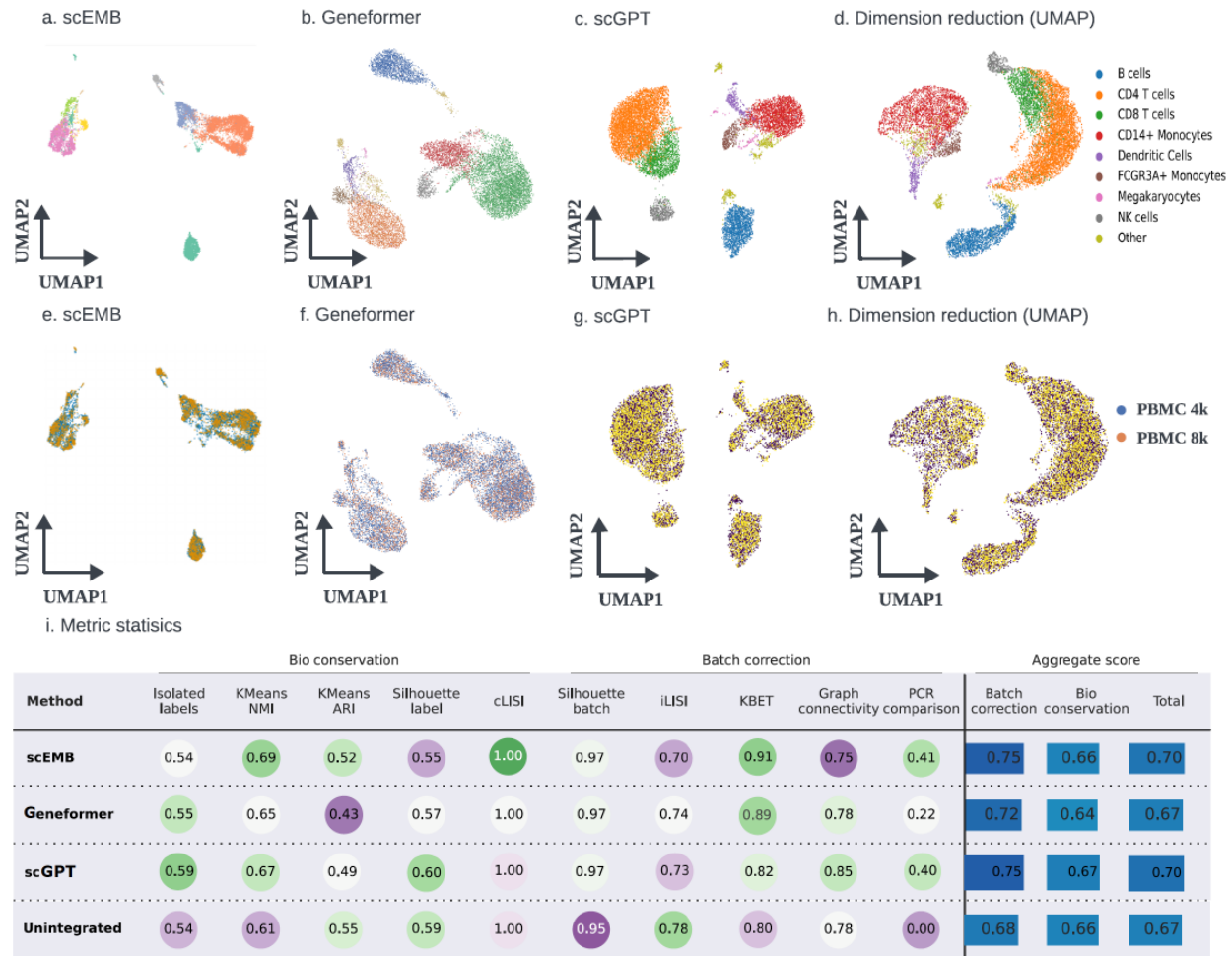


514

515 **Fig. 1 Framework of scEMB.** The 30M single-cell transcriptome dataset, curated by CZ
516 CELLxGENE Discover, was processed by normalizing gene expression to counts per million
517 and applying a log_{1p} transformation. We employed a binned expression strategy adapted from
518 scGPT, binning genes into 100 intervals and ranking them based on their real expression values.
519 Leveraging a BERT model as the backbone, we engineered the positional embeddings to
520 preserve the order of real expression values. The resulting concatenated embeddings were fed
521 into 12 transformer blocks, training the model to capture gene order and generate a gene-
522 attention map to represent cells. During inference, the gene expression data for new cells were
523 tokenized and inputted into the pretrained model to obtain both gene and cell embeddings.
524 These embeddings were then used in various downstream tasks to streamline conventional
525 single-cell analyses. This figure was created using materials adapted from Biorender.com.

526

527 **Fig. 2 Clustering performance of PBMC10k data.** a-d, UMAP plots illustrating clustering
528 performance for scEMB, Geneformer, scGPT, and conventional dimensionality reduction
529 methods on the PBMC10k dataset. e-h, UMAP plots illustrating batch integration performance
530 for scEMB, Geneformer, scGPT, and conventional dimensionality reduction methods on the
531 PBMC10k dataset. i, Five clustering metrics are compared across the four benchmark methods,
532 with scores ranging from 0 to 1, where higher values indicate better performance.

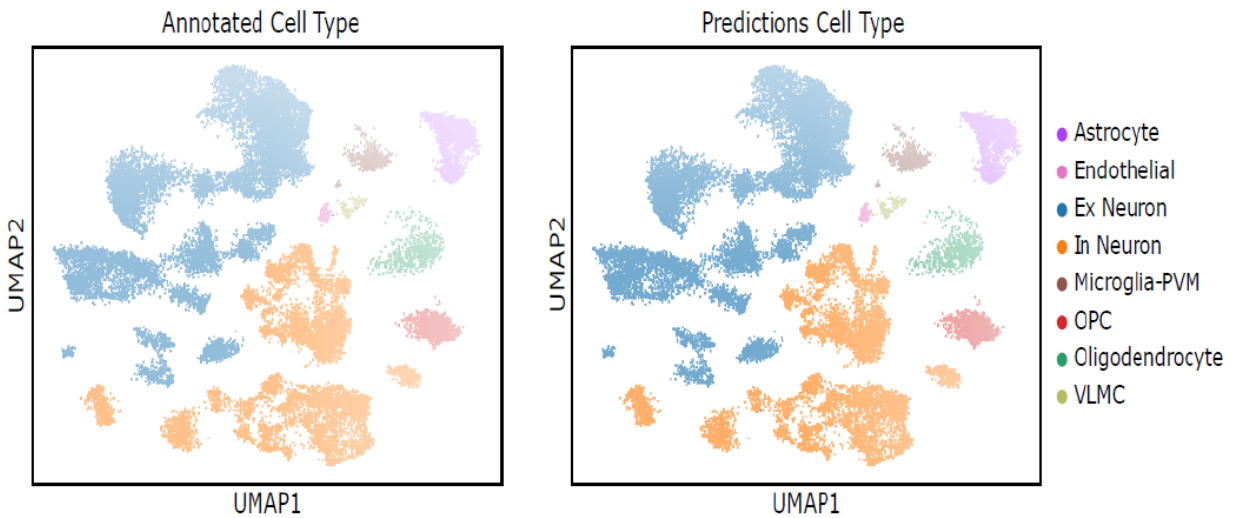


533

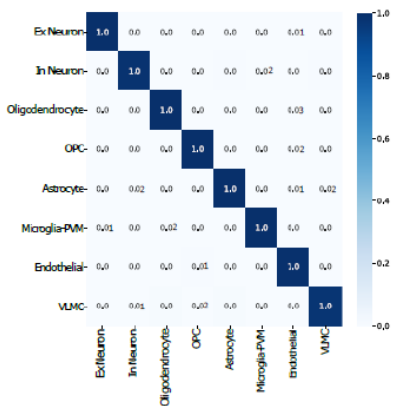
534

535 **Fig.3 Cell type annotation performance after fine-tuning on brain MTG datasets. We**
 536 **adapted two samples from SEA-AD cohort.** a. One sample was used for training the cell type
 537 classifier, while the other sample was reserved for testing the classification performance. b-d.
 538 We benchmarked the cell type annotation performance of scEMB against scGPT and
 539 Geneformer, and presented the results using confusion matrix. Darker blue along the diagonal
 540 indicates higher accuracy.

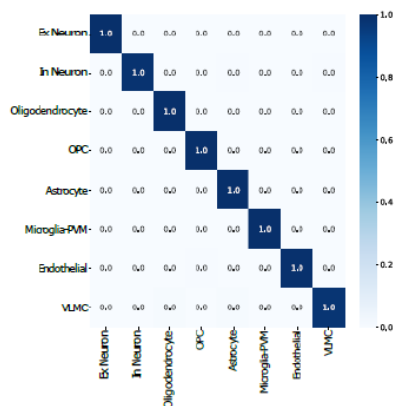
a. Cell type annotation in training and testing



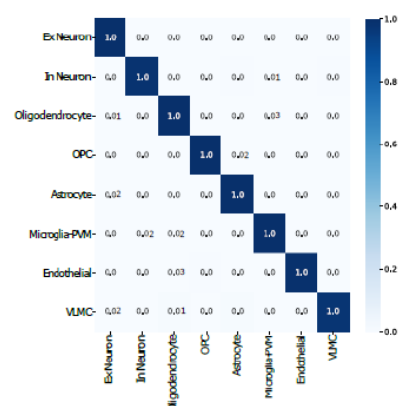
b. scEMB



c. scGPT



d. Geneformer



541

542

543

544

545

546

547 **Fig.4 scEMB *in silico* perturbation analysis.** a. Diagram of the scEMB *in silico* perturbation

548 analysis framework, from input data processing to the comparison of cellular embedding

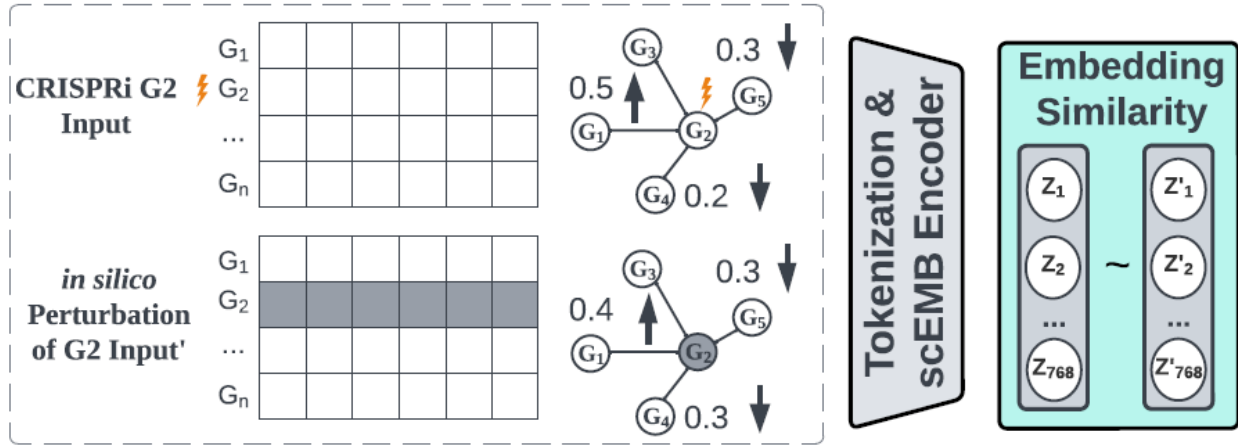
549 similarities. To evaluate the performance of the *in silico* perturbation, we used a CRISPRi

550 dataset as the ground truth and measured prediction accuracy by calculating the cosine

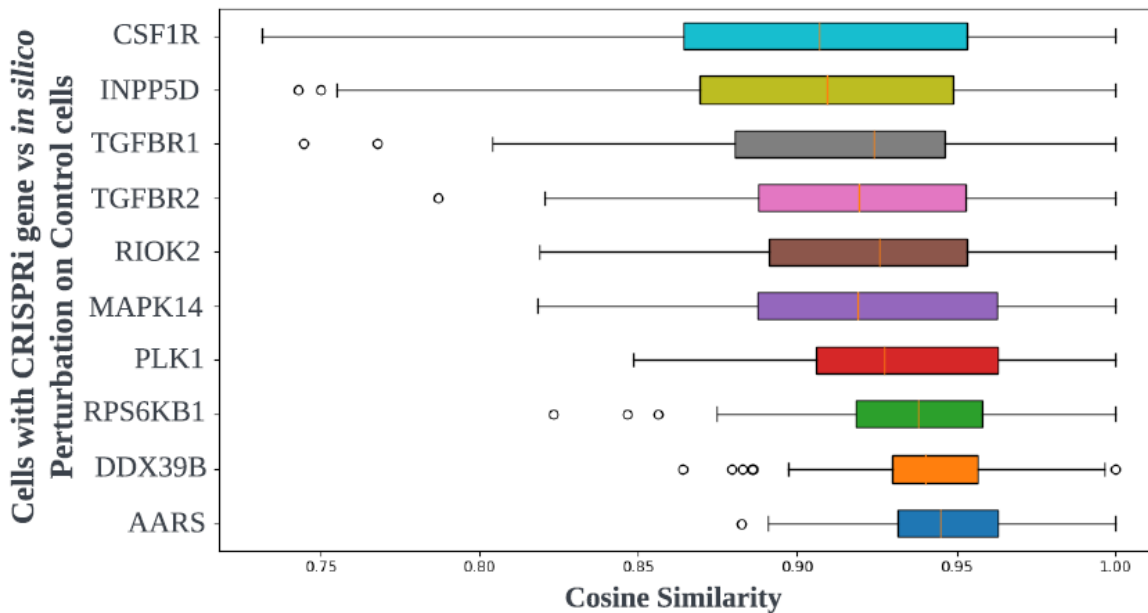
551 similarity of the cell embeddings. b. Bar plot showing the cosine similarity of cell embeddings for

552 10 CRISPR-edited genes, comparing their *in silico* perturbations on controls with the
 553 corresponding CRISPRi gene perturbed cells.

a. *in silico* perturbation response of gene expression on CRISPRi data.



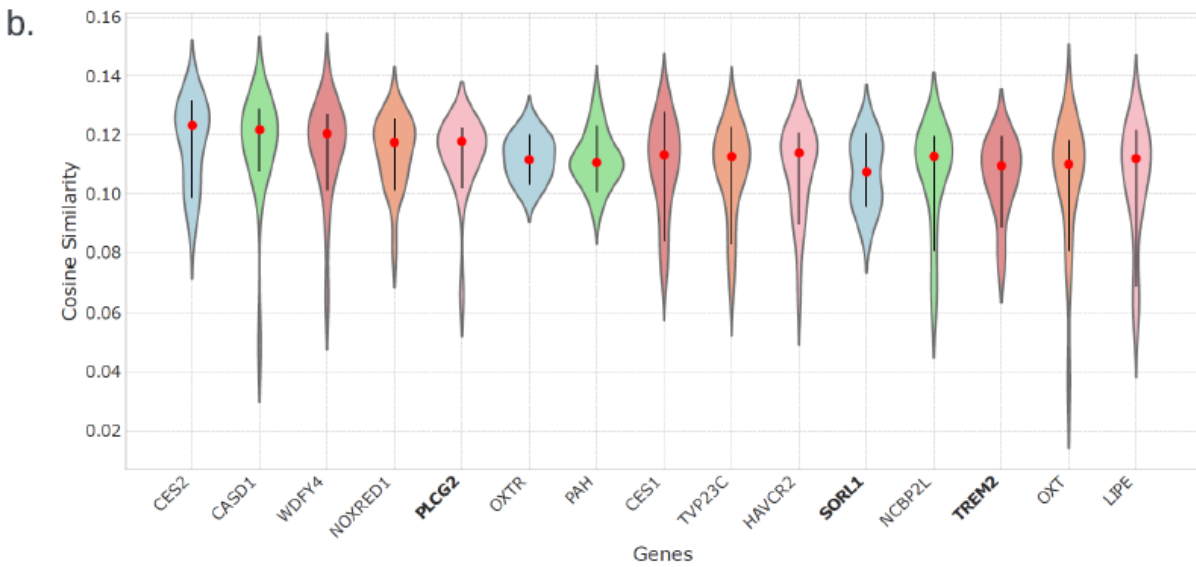
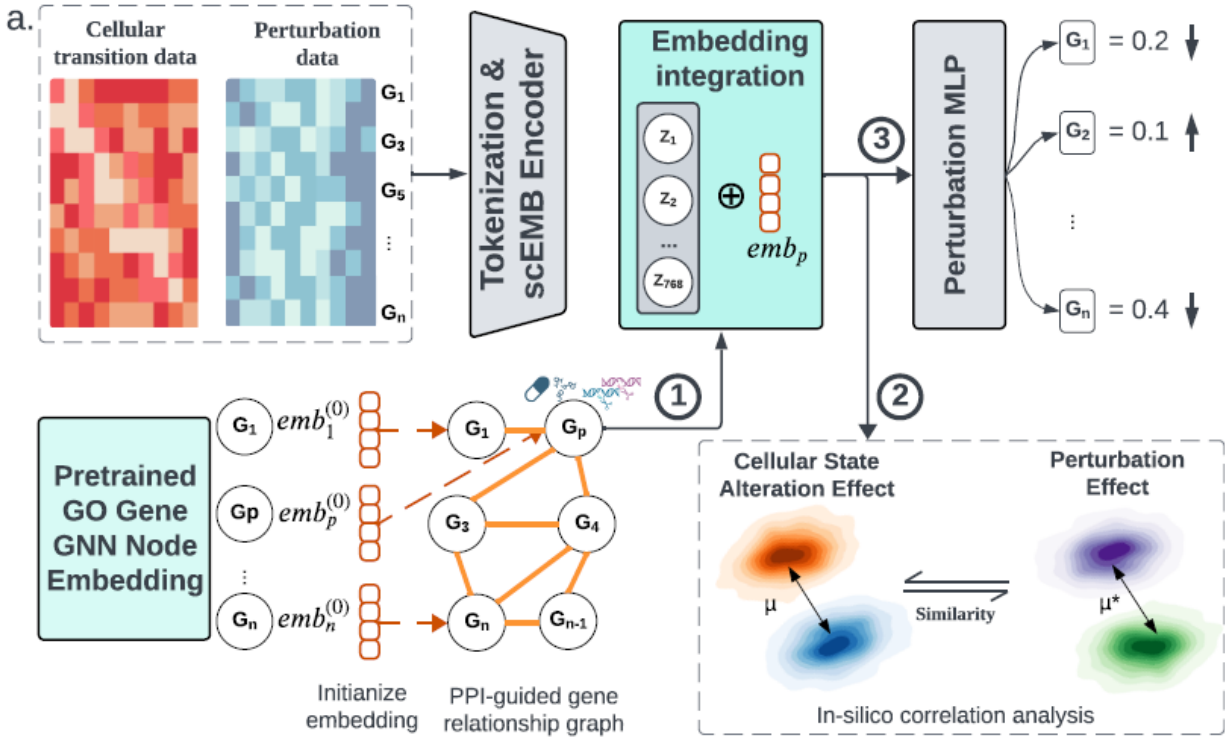
b. *in silico* perturbation analysis measure by cosine similarity.

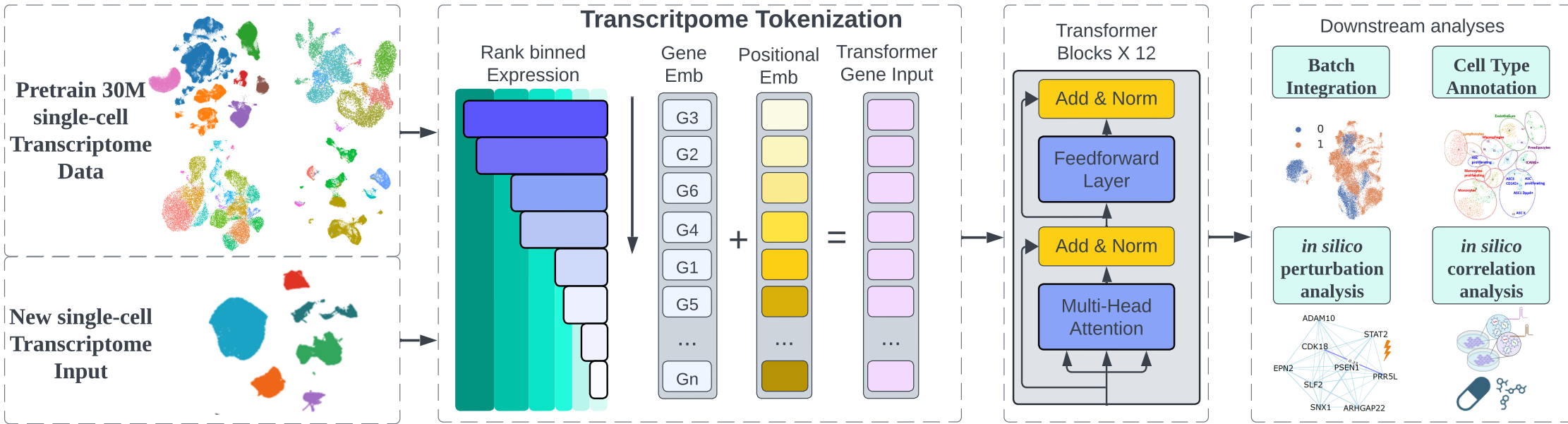


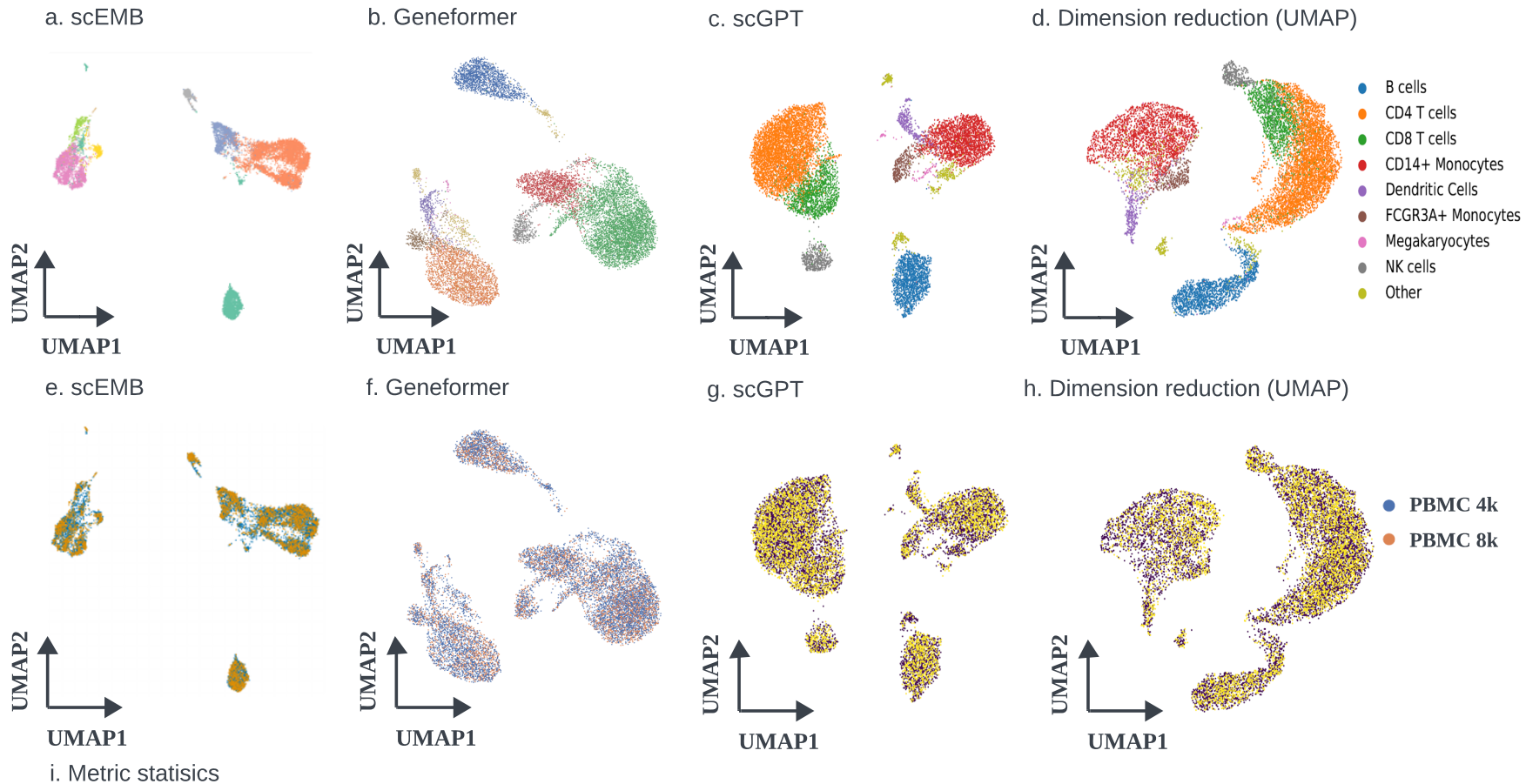
554

555 **Fig. 5 scEMB perturbation response and *in silico* correlation analysis.** a. All the single-cell
 556 transcriptome were first tokenized and went through scEMB Encoder. The 786-dimensional
 557 cellular embedding will be concatenated with gene perturbation embeddings (emb_p), which are
 558 derived from a PPI-guided gene relationship graph (step 1) following our previous method¹⁰.
 559 This graph is built using genetic perturbation data and propagates the impact of perturbations (P

560 on gene G_p) across gene-gene relationships informed by gene ontology (GO). scEMB provides
561 two downstream tasks to analyze perturbation effects at both the gene and cellular levels,
562 respectively. For gene-level analysis in step 2, the concatenated embeddings are processed
563 through a two-layer Multilayer Perceptron (MLP), designed to link cell embeddings to gene
564 features. The model outputs altered gene expression profiles, representing the predicted overall
565 impact of perturbations on other genes. For cellular level *in silico* correlation analysis, we
566 designed to test the correlation analysis between cellular state alteration effect and Perturbation
567 effect, and therefore capture to a potential reversing effect might provide the insight of therapeutic
568 targets (step 3). b. The violin plot shows the top 15 absolute cosine similarity scores for cellular
569 state alterations in microglia, comparing cells from Alzheimer's Disease (AD) with cells from
570 controls, as well as the effects of gene perturbations on iPSC-derived microglia. The absolute
571 cosine similarity measures the similarity between these two effects in 768 dimensions. The x-
572 axis represents the perturbed genes, with AD risk genes highlighted in bold. This figure was
573 created using materials adapted from Biorender.com.



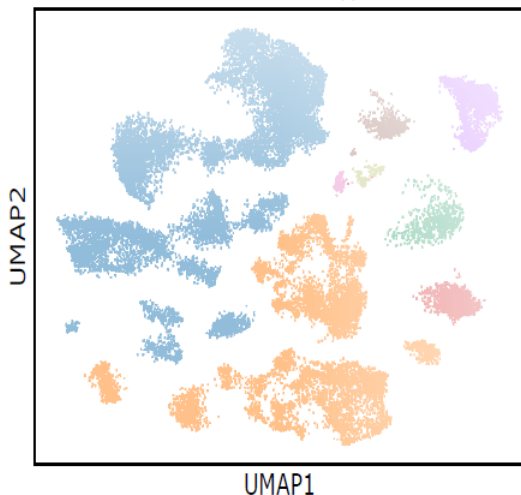




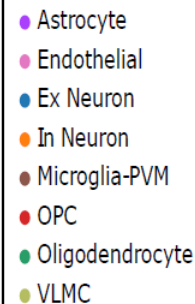
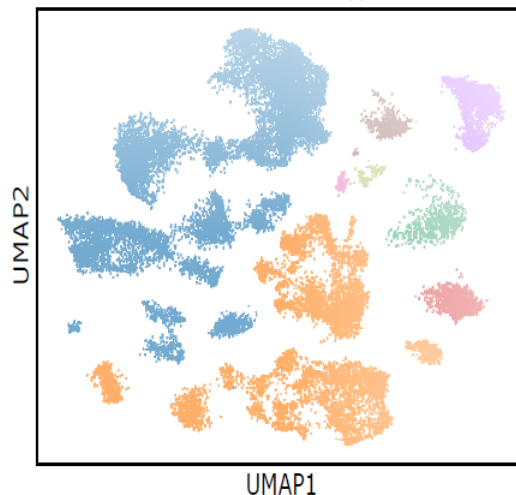
Method	Bio conservation					Batch correction					Aggregate score		
	Isolated labels	KMeans NMI	KMeans ARI	Silhouette label	cLISI	Silhouette batch	iLISI	KBET	Graph connectivity comparison	PCR	Batch correction	Bio conservation	Total
scEMB	0.54	0.69	0.52	0.55	1.00	0.97	0.70	0.91	0.75	0.41	0.75	0.66	0.70
Geneformer	0.55	0.65	0.43	0.57	1.00	0.97	0.74	0.89	0.78	0.22	0.72	0.64	0.67
scGPT	0.59	0.67	0.49	0.60	1.00	0.97	0.73	0.82	0.85	0.40	0.75	0.67	0.70
Unintegrated	0.54	0.61	0.55	0.59	1.00	0.95	0.78	0.80	0.78	0.00	0.68	0.66	0.67

a. Cell type annotation in training and testing

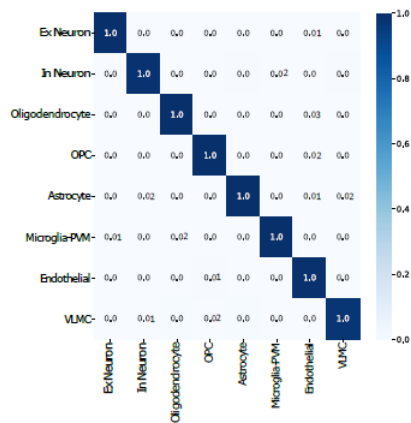
Annotated Cell Type



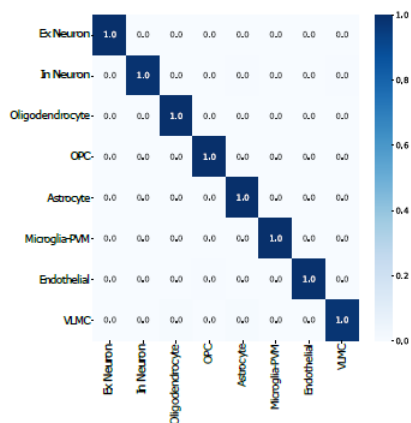
Predictions Cell Type



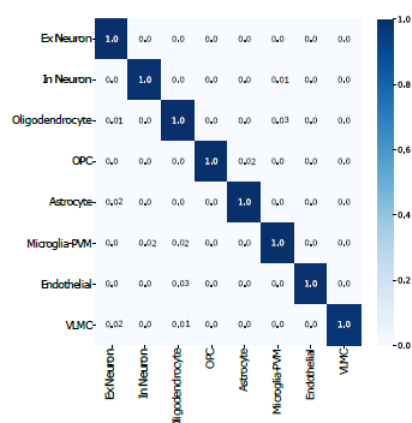
b. scEMB



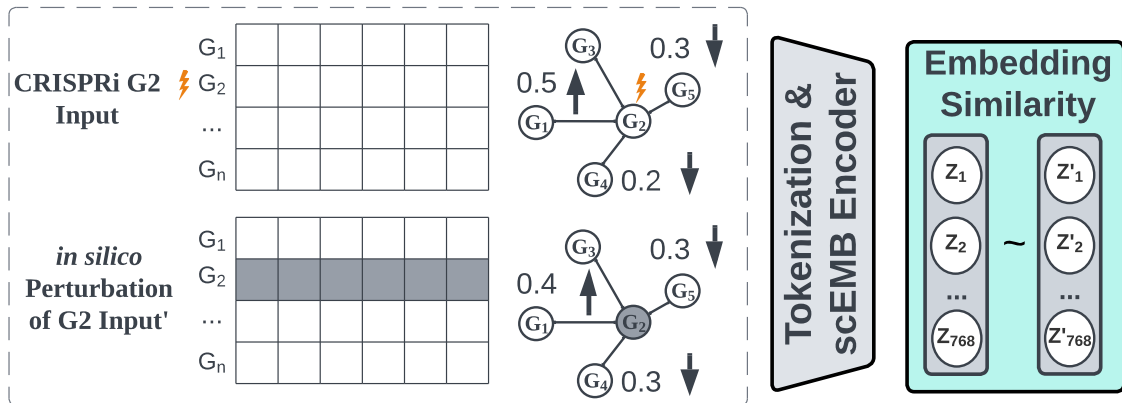
c. scGPT



d. Geneformer



a. *in silico* perturbation response of gene expression on CRISPRi data.



b. *in silico* perturbation analysis measure by cosine similarity.

