scPaLM: Pre-training of Single-cell Language Models through Genetic Pathway Learning

Anonymous ACL submission

Abstract

As scRNA-seq expands dataset richness and 001 advances cellular biology, transformer-based single-cell foundation models have emerged. Despite their efficacy, they face two key challenges: (1) Biological Perspective: They overlook the unordered nature of gene expression and fail to incorporate pathway priors essen-007 tial for capturing functional interactions. (2) Computational Perspective: The high dimensionality of gene tokens leads to excessive tokenization and a lack of an efficient mechanism for handling long-context sequences. Driven by these challenges, we propose a novel Singlecell Pre-trained Language Model via Genetic **Pa**thway Learning, named scPaLM¹, that addresses these challenges through three key innovations: **O Permutation-Invariant Embed-**017 deding where we handle high numer of genes via patching technique at the same time keeping permuation-invaraince; **2** a **Genetic Pathway** Learning module that is designed to learn discrete representations, enabling the modeling of collective gene behaviors in a data-driven way;
Cell-level Information Aggregation that progressively aggregates cell representations into a designated token during the training phase, with a tailored masking strategy and a token-level contrastive regularizer. Comprehensive evaluations across four biological benchmarks demonstrate scPaLM's superiority: it achieves average 10.1% improvement in cell type annotation compared to scGPT across all datasets, 5.15% increase in drug response prediction correlation compared to scFoundation.

1 Introduction

040

Single-cell RNA sequencing has emerged as the state-of-the-art method for elucidating the intricacies and diversity inherent in RNA transcripts at the individual cell level and providing insights into the composition of distinct cell types and

Token Length = N_{gene}



Figure 1: Compared to existing works where N_{gene} often reaches tens of thousands, we leverage the inherent functional similarity among gene groups, such as pathway-level organization (e.g., TREM2, HLA-DRA, and CD86 within the Immune Response Pathway), to enable latent pathway learning with an efficiently reduced token count $K << N_{gene}$ in the latent space.

their respective functions within tissues, organs, and organisms (Jovic et al., 2022). The massive amount of data generated by scRNA-seq techniques has provided massive information on various cells, enabling a better understanding of them, and hence benefit diverse research areas such as development (Semrau et al., 2017), auto-immune diseases (Gaublomme et al., 2015) and cancer diagnosis or prognosis (Patel et al., 2014).

To effectively model the scRNA-seq data, various computational methods (Cui et al., 2024a; Hao et al., 2024a) with different architecture designs have been proposed. Recent progress in deep learning has inspired the usage of advanced machine learning techniques, such as transformers (Vaswani, 2017), to scRNA-seq data analysis. These models, particularly those inspired by the success of

¹Source code is provided in supplementary.

BERT (Devlin et al., 2018) GPT (Radford et al., 058 2019), adopt a token-based approach, treating the 059 expression count of each gene as a "token", a con-060 cept widely embraced in natural language processing (NLP). These tokens are then assembled into a "sentence" representing the genetic expression 063 profile of a cell. Such models have demonstrated 064 exceptional performance in capturing the underlying structures inherent to scRNA-seq data and 066 have consistently outperformed conventional algo-067 rithms across various downstream tasks, including cell type annotation and bulk drug response prediction (Cui et al., 2024b; Hao et al., 2024b; Theodoris et al., 2023; Yang et al., 2022). Despite its effectiveness, several major bottlenecks remain:

Biological Perspective. Unlike NLP domain, where the order of words carries semantic meaning, gene expression in scRNA-seq does not follow a strict order. Traditional transformer models assume sequential order (like in sentences), but gene expression is inherently unordered-what matters is the presence, absence, and expression levels of genes (Cui et al., 2024b). This naturally raises the challenge of efficiently handling the vast number of genes in single-cell data. Moreover, existing approaches lack explicit mechanisms to incorporate pathway prior knowledge - the well-established bi-084 ological principle that genes operate in coordinated functional units to execute cellular processes (Cui et al., 2024b; Liang et al., 2023). This oversight limits their ability to model the hierarchical organization of genetic regulation and discover biologically interpretable patterns.

Computational Perspective. Considering that each cell is expressed with vast number of genes, treating each individual gene as a distinct token leads to a substantial increase in the overall token count, resulting in significant computational demands. Traditional approaches to address this issue is to select a subset of highly variable genes, such as the top 2,000 genes, which can notably reduce the overall gene count (Cui et al., 2024b). However, this gene exclusion inevitably entails the loss of valuable biological information, as certain essential genes, like housekeeping genes, may not exhibit high variability while maintaining pivotal regulatory functions (Joshi et al., 2022).

101

102

103

104 105

106

107

109

In this paper, we propose scPaLM, a transformerbased model that effectively harnesses massive scRNA-seq data. scPaLM incorporates multiple innovative elements: **0** We propose an **Permutation-Invariant Embedding**, an efficient embedding process that condenses the information from all 110 the genes into a reduced number of tokens by 111 leveraging a symmetric encoder-decoder architec-112 ture while ensuring the nature of permutation-113 invariance, substantially mitigates the computa-114 tional expenses and facilitates rapid training and 115 inference: ⁽²⁾ Acknowledging the collective nature 116 of gene functionality, we introduce a genetic path-117 way learning, an encoder which seeks to encode 118 gene tokens into genetically related pathway to-119 kens and acquire discrete representations for them 120 to capture the collective yet distinct functionality 121 of genes; Finally, 3 to aggregate cell-specific infor-122 mation, we establish a tailored training framework, 123 Cell-level Information Aggregation to learn a des-124 ignated token to represent cells, with the establish-125 ment of a masking strategy and a token-level con-126 trastive regularizer. Extensive downstream tasks 127 including cell type annotation and drug response 128 prediction, scPaLM achieves greater performance 129 compared to baseline methods, including Gene-130 former (Theodoris et al., 2023), scGPT (Cui et al., 131 2023) and scFoundation (Hao et al., 2023). 132

Our contributions are summarized as follows:

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

- * We propose scPaLM, a pretrained foundation model which incorporates biological pathway with efficient tokenziation on massive scRNAseq data that achieves state-of-the-art performance on various downstream tasks.
- * We devise multiple innovative elements that contribute to the success of scPaLM: (1) a novel permutation-invariant embedding process that efficiently maps gene expression values into representations for subsequent modeling; (2) a genetic pathway encoder that is designed to model the collective behaviors of genes by learning discrete representations for their tokens; and (3) a training scheme that aggregates cell-specific information into a designated token, and two techniques that augment the aggregation process.
- * We demonstrate performance superiority of scPaLM on various downstream tasks. For the cell type annotation task, we outperform Geneformer and scGPT by {8.4%,4.9%}/{8.4%,3.0%} and {6.5%,0.9%}/{5.4%,0.2%} in terms of the ARI and NMI scores of two scRNA-seq datasets. For the drug response prediction task, we outperform scFoundation by up to 5.15% in terms of the correlation of IC50 values. We also achieve

254

255

208

159 160

162

163

164

165

167

168

169

170

171

172

173

174

175

176

177

179

180

181

182

184

185

186

188

189

191

192

193

194

197

198

201

204

207

state-of-the-art performance on the drug response perturbation prediction task.

2 Related Works

Single-cell Data Analysis. Several methods have been proposed to model transcriptome measurements at the single-cell level that reflect biological diversity (Lopez et al., 2018). MAGIC (Dijk et al., 2017) aims to predict the missing measurements, often referred to as "dropouts", by propagating in a graph constructed based on cell-cell similarity. scImpute (Li and Li, 2018) learns to accurately and robustly identify dropouts and perform imputation on these identified positions. SAVER (Huang et al., 2018) leverages gene-to-gene relationships to recover the expression level of each individual cell. Lopez et al. (Lopez et al., 2018) developed a scalable framework called scVI for probabilistic representation and analysis of gene expression in single cells. More recently, there has been a notable utilization of pretrained transformers in the context of modeling single-cell RNA sequencing (scRNAseq) data. In the realm of encoder-only transformers, scBERT (Yang et al., 2022) embeds each gene into a token and leverages an efficient transformer to model over 16,000 genes for each individual cell. Subsequently, scFoundation (Hao et al., 2023) has made advancements in the embedding process introduced by scBERT, which has resulted in enhanced performance. Geneformer (Theodoris et al., 2023) discards the original measurement of transcriptome and constructs input sequences that account for the ranking of measurements across the entirety of the dataset, thereby creating a representation that encapsulates the relative expression levels of all genes within each cell. For decoderonly models, scGPT (Cui et al., 2023) leverages the concept of next token prediction in NLP to iteratively predict the masked genes, creating a novel path for scRNA-seq data modeling. A concurrent work CellPLM (Wen et al., 2023) encodes cell-cell relations by leveraging spatially-resolved transcriptomic data in pre-training. In this work, our objective is to deploy a vector quantization technique to learn discrete genetic pathway representation.

3 Methodology

Overview. A scRNA-seq dataset is usually stored as a matrix $X_{raw} \in \mathbb{N}^{N_c \times N_g}$, where N_c is the number of cells and N_g is the number of genes. In X, each row represents the expression values of genes in a cell. A transformation (*e.g.*, log1p) is usually applied on X_{raw} to obtain X, *i.e.* data with normalized scales (Yeo and Johnson, 2000).

The main components of scPaLM are transformers (Vaswani, 2017). To facilitate the learning process, we propose a novel embedding process (§3.1), referred to as $embed(\cdot)$, that maps each row of X into N tokens. Note that our embedding process can produce a reduced amount of tokens, which is more memory- and computation-efficient compared to existing works which usually need to construct a large number of tokens (typically equals to N_q). The tokens are then fed into the transformers in scPaLM, which leverage the attention mechanism to capture the interaction between tokens and learn the biological implications behind the scRNAseq data. The transformers have multiple layers and process the token representations sequentially with the following formula: $h_i = \text{Layer}_i(h_{i-1})$, where h_i indicates the token representations generated by the *i*-th transformer layer.

The training process for scPaLM consists of two distinct stages. *During the initial stage*, we train both an encoder and a decoder using a reconstruction loss. Specifically, the encoder is trained to map the raw gene tokens to tokens that represent genetic pathways (as discussed in §3.2), while the decoder's role is to reverse this mapping, converting genetic pathway tokens back to the original expression levels. *In the second stage*, we train another encoder with an additional token designed to capture cell-specific information (as discussed in §3.3). The two encoders we have trained in these two stages collectively empower various downstream tasks, exhibiting superior performance.

3.1 Permutation-Invariant Embedding

The Tokenization Dilemma in scRNA Modeling Current tokenization strategies for scRNAseq data face a fundamental trade-off: full-gene tokenization preserves molecular granularity but incurs prohibitive $O(N_g^2)$ computational complexity (20,000 $\leq N_g \leq 60,000$), while gene filtering discards critical biological signals. Even advanced methods like HVG selection (Cui et al., 2024b) risk eliminating housekeeping genes with low variability but essential functions (Joshi et al., 2022). This dilemma severely limits Transformer-based models' scalability and biological fidelity.

Pathway-Inspired Patch ConstructionDraw-256ing inspiration from Vision Transformers (Dosovit-257



Figure 2: The overview of our framework. Three main innovative components are introduced in our frameworks: an efficient embedding process that is permutation-invariant; a module that captures genetic pathways; and a training framework for cell information aggregation.

skiy et al., 2020), we reconceptualize gene patches as functional units analogous to biological pathways. However, unlike image pixels with spatial coherence, genes lack inherent positional relationships. Direct application of ViT-style patching introduces order sensitivity—permuting gene order alters patch composition, contradicting biological reality where expression vectors are permutationinvariant. To resolve this, we develop a symmetric embedding architecture that dynamically clusters genes into N pathway-aligned patches ($N \ll N_g$) through order-agnostic operations.

260

261

262

264

Symmetric Embedding Pipeline As detailed in Algorithm 1, our pipeline enforces permutation 271 invariance through three key stages. First, gene 272 expressions are projected into d-dimensional vectors via learnable matrices $P_1 \in \mathbb{R}^{N_g \times d}$, scaled 274 by gene-specific coefficients α . A hierarchical 275 pooling strategy then aggregates these embeddings: zero-expression genes are replaced by a learnable 277 embedding and pooled separately, then combined with active genes through concatenation and linear 279 fusion. This two-stage process ensures balanced representation of both silent and active genomic regions. The final embeddings $\bar{E} \in \mathbb{R}^{N imes d}$ are computed through linear interpolation of pooled features with learnable basis vectors P', eliminat-284 ing positional dependencies. This architecture reduces token count by 99.6% (from $N_q = 60,664$ to N = 256) while achieving 16.1% higher clustering accuracy than gene-wise baselines (Table 6),

demonstrating superior biological fidelity.

3.2 Genetic Pathway Learning

Biologically, most genes do not function in isolation; instead, they function in concert to perform biological functions (Alexa et al., 2006; Shastry, 2009; Mi et al., 2013). Here, groups of biologically related genes that demonstrate substantial associations with specific biological processes are commonly referred to as *pathways*. Recognizing the activated pathways within a cell holds paramount importance in comprehending its characteristics (Wang and Sherwood, 2011). Despite such significance of pathways current methodologies frequently overlook this aspect. 290

291

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

Pathway Modeling through Discrete Representations Learning. We propose to learn distinct "pathway" tokens, represented as discrete codes, by training an encoder and a vector quantizer. To be more specifically, the encoder, implemented as a transformer, maps \bar{E} into hidden representations, denoted as $\mathcal{E} = \{e_1, e_2, \dots, e_N\}$. Subsequently, the quantizer learns a codebook \mathcal{V} = $\{v_1, v_2, \ldots, v_K\}$, and associates each e_i with the closest entry in \mathcal{V} in terms of distance. More precisely, for each embedding with index *i*, we derive the corresponding embedding from the codebook with the following formula: $oldsymbol{z}_i = rg \min_{v \in \mathcal{V}} \|oldsymbol{v} - oldsymbol{v}\|$ $e_i \parallel_2$. After acquiring $\mathcal{Z} = \{z_1, z_2, \dots, z_N\}$, we input it into an additional decoder, which is implemented as another transformer model, and obtain the output vector $\boldsymbol{o} \in \mathbb{R}_q^N$. These modules are

325

327

329

331

333

334

3.3

tion.

336 337

- 338
- 340

341 342

344

346 347

357

Token Scaling and Reconstruction Despite 361 pathway tokens being fewer than genes, our scaling algorithm (Algorithm 3) effectively aligns their dimensions through linear interpolation. This approach preserves biological fidelity while reducing computational complexity, as evidenced by the 366

this modified embedding using a decoder.

trained using reconstruction tasks, where a mean-

squared-error (MSE) loss is employed to minimize

the dissimilarity between x and o. Furthermore,

we introduce a commitment loss (Huh et al., 2023)

to minimize the distance between each pair of e_i

 $\mathcal{L}_{\text{cmt}} = \sum_{i=1}^{N} \beta \| \text{sg}(\boldsymbol{z}_i) - \boldsymbol{e}_i \| + (1 - \beta) \| \boldsymbol{z}_i - \text{sg}(\boldsymbol{e}_i) \|,$

where sg indicates the stop-gradient operation. To

mitigate the issue of index collapses in VQ tech-

niques, we follow Huh et al. (2023) to regularly

replace the unused tokens in codebooks with ran-

domly re-initialized tokens, and leverage an affine

parameterization to minimize interval covariate

While the pathway encoder encodes gene expression values into pathway tokens, it does not pro-

vide a cellular-level representation. One possible

approach to constructing cell representations in-

volves concatenating the representation of genetic

pathway tokens. However, this results in a pro-

hibitively high-dimensional representation, leading

to increased computational costs in downstream

tasks. Alternatively, using the average representa-

tion of pathway tokens, while simpler, yields infe-

Additional tokens to aggregate cell informa-

tively and efficiently at the cellular level, we intro-

duce a learnable token, e_C , designed to encapsu-

late cell-specific information. This token is asso-

ciated with a subset, denoted as $z_{i_1}, z_{i_2}, \ldots, z_{i_{N'}}$,

randomly selected with monotonically increasing

indices from the set \mathcal{Z} , therefore information is

transferred to e_C . We employ an encoder to con-

vert these tokens into representations, which we denote as $\mathcal{H} = \{\boldsymbol{h}_C, \boldsymbol{h}_{i_1}, \boldsymbol{h}_{i_2}, \dots, \boldsymbol{h}_{i_{N'}}\}$. Subse-

quently, we exclude h_C , replace the positions of

the previously omitted pathway tokens with a com-

mon learnable token e_M , and proceed to decode

To aggregate gene representations effec-

rior performance as shown in §5.

shifts. More details are provided in Appendix C.

Cell Information Aggregation

(1)

and z_i , which has the following form:

superior clustering performance of h_C over geneaveraged baselines in Table 6. The architecture freezes the encoder and quantizer introduced in §3.2, minimizing reconstruction loss between original and decoded expressions, ensuring stable training while maintaining pathway-aware representations.

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

384

386

387

388

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

Contrastive Regularization with Dynamic **Pseudo-Labels** To enhance cell-specific discriminability, we design a contrastive framework leveraging K-means-derived pseudo-labels. During training, two augmented views $(\mathcal{H}_i, \mathcal{H}'_i)$ of each cell are generated by masking different pathway token subsets. Embeddings $h_{i,C}$ and $h'_{i,C}$ from the same cell cluster form positive pairs, while crosscluster cells serve as negatives (limited to K pairs per batch). The contrastive loss

$$\mathcal{L}_{CL} = \sum_{i} -\log \frac{s(\boldsymbol{h}_{i,C}, \boldsymbol{h}_{i,C}')/\tau}{s(\boldsymbol{h}_{i,C}, \boldsymbol{h}_{i,C}') + \sum_{j \in \mathsf{neg}(i)} s(\boldsymbol{h}_{j,C}, \boldsymbol{h}_{i,C})}$$
(2)

maximizes similarity within pseudo-classes through cosine similarity $s(\cdot, \cdot)$, with temperature τ sharpening distinctions (Chen et al., 2020). A fixed-length queue Q stores recent h_C embeddings, enabling periodic K-means updates to adapt pseudo-labels to evolving representations. This implicit modeling of cell-type relationships enforces topological consistency in the latent space, as visualized in Figures 3-4. The detailed pipeline is in Algorithm 2.

4 **Experiments**

In this section, we detail our experimental setting and demonstrate the superior performance of scPaLM. A series of experiments are conducted on various downstream tasks, and we have provided ablation studies in §5 to validate the importance of our proposed techniques.

4.1 Implementation Details

Pretraining Data. scPaLM is pretrained on singlecell RNA-seq data covering different types of cells, having in total 43, 312, 189 cells and 60, 664 genes collected from CELLxGENE platform (Megill et al., 2021). The statistics and description of the pre-training data are in Appendix B.

Architectures. The overall architecture of scPaLM can be split into three parts: (1) embedding layers, where we use mainly MLPs as described in Algorithm 1. The N is set to 256 in our experiment;



Figure 4: Unsupervised clustering performance on the COVID dataset.

(2) genetic pathway learning, where we introduce an encoder, a decoder, and a quantizer. The en-coder and the decoder are developed based on the transformer architecture, which contains 12 layers and a hidden dimension of 768. The quantizer has the same hidden dimension with a codebook size of 128; (3) cell information aggregation, where we introduce another encoder and decoder that have the same architectures and configuration. More details on the architectures can be found in Table 4.

Training Settings. The optimizer we use in our experiments is AdamW (Loshchilov and Hutter, 2017). In the first stage of training, where we train the encoder for pathway tokens and the vector quantizer, we adopt a learning rate of 0.001 and a batch size of 128. In the second stage, we train other components with a learning rate of 5×10^{-4} and using the same batch size of 128.

Benchmarks. We compare against baseline methods on various benchmark datasets that are manually excluded from the pretraining data: (1) the CLL (GEO: GSE111014) dataset (Rendeiro et al., 2020), which originally contains 48016 cells with 33694 genes and 6 types of cells. We further filter out cells without type annotations, resulting in 30K cells; (2) the COVID (GEO: GSE150861) dataset (Guo et al., 2020), which contains 11931 cells; (3) the Jurkat from 10x Genomics dataset², which contains 3258 cells. We filter out zero-count genes and retain 17753 genes. (4) the PBMC-5k dataset that also comes from 10x Genomics³ which

contains around 5K cells; and finally the Cancer Cell Line Encyclopedia (CCLE) (Barretina et al., 2012) and Genomics of Cancer Drug Sensitivity (GDSC) (Iorio et al., 2016) datasets which are leveraged in the drug response prediction task (§4.3). **Baselines.** We compare scPaLM with various baseline methods on different tasks. For cell-type annotation, we compare with PCA (where we derive the first 256 principal components on the log-normalized expression values), Gene-former (Theodoris et al., 2023), scGPT (Cui et al., 2023). The latter two are current state-of-the-art algorithms for this task. For imputation, we compare with SAVER (Huang et al., 2018), scImpute (Li and Li, 2018), DCA (Eraslan et al., 2019), which are widely used methods for this task. For drug response prediction, we compare with Deep-CDR (Liu et al., 2020) and another transformer-based algorithm, scFoundation (Hao et al., 2023).

4.2 Unsupervised Cell Type Annotation

Our first set of experiments involves applying computational methods on unseen scRNA-seq data and providing type annotations to those unseen cells in an unsupervised manner. We compare the performance of scPaLM with three baselines, namely PCA, Geneformer, and scGPT. These experiments are conducted on the CLL and the COVID dataset. Figure 3 and 4 display UMAP visualizations created from the cell representations, *i.e.*, h_C . We use the Leiden (Traag et al., 2019) algorithm with a resolution of 1.0 to cluster the embeddings and assess the clustering performance with the adjusted rand index (ARI) and normalized mutual information (NMI) scores. Qualitatively speaking, PCA

 $^{^2\}mbox{Avaiable}$ at https://www.10xgenomics.com/datasets/jurkat-cells-1-standard-1-1-0

³Avaiable at https://www.10xgenomics.com/datasets/5khuman-pbmcs-3-v3-1-chromium-controller-3-1-standard



Figure 5: Comparison of Pearson correlation coefficient (PCC \uparrow) between the predicted and the ground-truth IC50 values using different settings of feature . We compare scPaLM's performance with two baseline algorithms, DeepCDR and scFoundation.

479

480 481

482

483

484

485

486

487

488

489

491

492

493

494

495

496

497

498

499

503

504

505

507

509

511

512

513

514

515

exhibits the poorest results, whereas the UMAPs generated by the other three models demonstrate significantly superior clustering quality. We can also observe that scPaLM possesses a smoother and more clustered latent space with respect to the ground-truth cell type labels. Quantitative results also confirm scPaLM's ability to annotate types of cells. Notably, our method achieves higher ARI and NMI scores compared to the baselines by clear margins on both datasets. On CLL, scPaLM outperforms the baselines by $0.084 \sim 0.169$ in terms of the ARI score and $0.054 \sim 0.158$ in terms of the NMI score. Similarly, on COVID, it outperforms the baselines by $0.030 \sim 0.356$ in terms of the ARI score and $0.002 \sim 0.247$ in terms of the NMI score. These improvements indicate the high quality of produced embeddings.

4.3 Cancer Drug Response Prediction

Cancer Drug Responses is an important task that can help guide the design of anti-cancer drugs and also understand the cancer biology (Unger et al., 2015). Following the setting in scFoundation (Hao et al., 2023), we combine scPaLM with a CDR prediction framework, DeepCDR (Liu et al., 2020), to provide prediction of the IC50 values (i.e., halfmaximal inhibitory concentrations) of drugs across different cells. We adopt the settings from scFoundation (Hao et al., 2023) to fuse the extracted representations from gene expression values with the representations of drugs and fit a graph convolution network (GCN) to learn representations that encompass information from multiple sources and modalities. We follow the settings of DeepCDR and experiment with different options: (1) Use Mut, which indicates the usage of genomic mutation information; and (2) Use Methy, which indicates the usage of DNA methylation data.

From Figure 5, we can observe that both scFoun-

dation and our method outperform the baseline framework DeepCDR significantly and achieve a stronger correlation between the prediction and the IC50 values. Notably, when using no additional information from the mutation and methylation, our method significantly outperforms scFoundation by 5% in terms of the Pearson Correlation Coefficient (PCC). When having additional mutation and methylation information, all methods demonstrate higher PCCs, yet our method remains the top performance among them all.

To have a better understanding of the performance gain, we provide pairwise visualization and case study of the correlation achieved by our method and scFoundation in Figure 6 and 7. Detailed analysis can be found in Appendix E.

4.4 Imputation

Imputation is an important task where the model is asked to recover the expression value of genes within individual cells. It has real-world implications because the measurement of expression levels often exhibits noise (Grün et al., 2014) and suffer from dropout events (Kharchenko et al., 2014). We conduct a series of simulated experiments on the Jurkat and the PBMC dataset to assess scPaLM's ability to accurately predict the expression levels of the missing genes. We randomly sample 10% of genes from each cell with a probability that is proportional to the exponent of negative gene expression values and mask them as 0.

Table 1: Imputation performance of various methods on the Jurkat and the PBMC dataset. We use the rooted mean square error (RMSE) and the mean absolute error (MAE) as the measurements.

Method	Jurkat		PBMC	
	$RMSE\downarrow$	$MAE\downarrow$	$RMSE\downarrow$	$MAE\downarrow$
SAVER	0.841	0.664	0.779	0.594
scImpute	1.178	0.838	1.528	1.132
DCA	0.937	0.629	0.833	0.638
scPaLM (Zero Shot)	0.494	0.397	0.674	0.539

Table 1 presents the rooted mean square error (RMSE) and the mean absolute error (MAE) between the ground truth and the predicted expression values on the masked genes across different cells. Note that these metrics are calculated based on the log-normalized expression values. Even under a zero-shot setting, scPaLM achieves superior performance compared to most baselines, which estimate their parameters on the downstream datasets. This experiment confirms scPaLM's ability in denoising the expression data and capturing the interactions

516

517

545

546

547

555

556

583

between cells and genes.

557

559

560

563

564

565

566

571

574

578

579

4.5 Genetic Pathway Identification

We finally conducted an experiment to understand the obtained pathway tokens from scRNA-seq datasets. We follow the setting from scGPT (Cui et al., 2023) where we aim to identify genetic pathways on the Immune Human dataset. To associate a gene with a certain pathway token, we first derive the \mathcal{V} for each gene, and for every v, we calculate the associated gene expression vector weighted by the occurrence percentage of v for each cell. Finally, for every v, we obtain the list of associated genes by calculating their relative prevalence. A more detailed algorithm is deferred in §D. We associate 10 genes to each pathway token and run the gene set enrichment analysis (GSEA) algorithm to search for pathways in Reactome Pathway Database (Fabregat et al., 2018). Note that this dataset is not included in our training set; therefore, it constitutes a zero-shot setting. Nevertheless, our method identifies two significant pathways related to the immune system, as shown in Table 7. Particularly, it identifies and clusters the CD1 gene family (CD1E and CD1B), which is involved in antigen presentation that is related to immune reaction.

5 Ablation Studies

Table 2: Clustering performance of different variants for cell representations on the CLL dataset (Rendeiro et al., 2020). We compare the adjusted rand index (ARI), normalized mutual information (NMI), silhouette score (S-score), and clustering time between models.

Method	$\text{ARI} \uparrow$	$\rm NMI\uparrow$	S-score \uparrow
Mean	0.015	0.059	-0.133
Concatenated	0.181	0.478	0.310
$oldsymbol{h}_C$ (No CL)	0.275	0.573	0.361
h_C (Ours)	0.292	0.593	0.376

The effectiveness of the cell information aggregation process. We conduct a series of experiments on the CLL dataset to compare two alternatives for building cell representations, where we use the average and the concatenated representations of pathway tokens to represent cells. Table 5 presents the performance on the cell type annotation task, where we can observe that using our cell information aggregation technique yields the best performance among all the variants. We have also conducted an experiment where we do not use the token-level contrastive learning framework to train the embedding of h_C . The decreased scores of these experiments demonstrate the importance of the token-level contrastive learning regularizer.

596

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

The effectiveness of the pathway encoder. We conduct experiments excluding the genetic pathway learning module discussed in §3.2. Instead of our proposed method, we train the embedding layers and also aggregate information directly from the embeddings of genes on a small subset of training data that have around 100K cells. Table 3 demonstrates the importance of the encoder and the quantizer. We see that the introduced genetic pathway encoder helps improve the clustering performance, improving the metrics by 0.14 and 0.16, respectively. The usage of the quantizer also further improves the performance by an additional 1%.

Table 3: Comparisons on COVID dataset (Guo et al., 2020) with different configurations. The models are trained on a small subset of the pre-training data.

Configuration		ARI↑	NMI ↑	
Encoder	Quantizer		1 1111	
X	×	0.050	0.120	
1	×	0.191	0.286	
1	1	0.201	0.291	

The effectiveness of the embedding process. To evaluate the effectiveness of our embedding process, we explore several alternatives and compare them to our proposed embedding process: (1) per-gene, where we directly employs gene-specific embeddings E. This is a widely adopted option in various methods such as scFoundation (Hao et al., 2023) and Geneformer (Theodoris et al., 2023); (2) shared-first-layer, where we employ only a shared **P** for all the genes. The results are presented in Table 6, where we can observe that these alternatives demonstrate either degraded performance, or suffer from overly high computational cost. The per-gene variant results in out-of-memory (OOM) error even using a batch size of 1. Using a shared first layer requires less amount of GPU memory but yields inferior performance compares to our method.

6 Conclusion

This work presents scPaLM, a foundation model pre-trained on single-cell RNA-seq data. We devise several novel techniques that efficiently represent gene expression values into tokens, model the collective function of genes, and effectively aggregate cell-specific information into a single token. We evaluate scPaLM on a wide range of downstream tasks, and demonstrate it reaches SoTA.

637 Limitations

Implicit Pathway Modeling. While our genetic pathway learning module captures collective gene behaviors through discrete representations, it currently relies on data-driven discovery rather than explicit integration of established pathway databases (e.g., Reactome or KEGG). Future work could enhance biological interpretability by incorporating curated pathway knowledge through hybrid architectures that combine learned codebooks with prior biological constraints.

648Human-Centric Data Bias.Our pre-training649dataset primarily focuses on human single-cell650transcriptomes from Tabula Sapiens.651enables strong performance on human biological652tasks, the cross-species generalizability of our path-653way representations remains uncertain.654lutionary divergence of gene regulatory networks655across species may require specialized adaptation656mechanisms when applying scPaLM to model non-657human organisms.

Ethics Statement

659This work adheres to ethical research practices660in computational biology. All datasets used for661pre-training and evaluation are publicly available662through GEO/SRA archives or 10x Genomics, with663proper ethical approvals obtained in their origi-664nal studies. Our framework processes only de-665identified genomic data, containing no protected666health information. While foundation models like667scPaLM could theoretically accelerate therapeutic668development, we emphasize that any clinical appli-669cation requires rigorous validation through estab-670lished biomedical research protocols.

References

671

675

677

679

- Adrian Alexa, Jörg Rahnenführer, and Thomas Lengauer. 2006. Improved scoring of functional groups from gene expression data by decorrelating go graph structure. *Bioinformatics*, 22(13):1600–1607.
- Alexei Baevski, Steffen Schneider, and Michael Auli. 2019. vq-wav2vec: Self-supervised learning of discrete speech representations. *arXiv preprint arXiv:1910.05453*.
- Jordi Barretina, Giordano Caponigro, Nicolas Stransky, Kavitha Venkatesan, Adam A Margolin, Sungjoon Kim, Christopher J Wilson, Joseph Lehár, Gregory V Kryukov, Dmitriy Sonkin, et al. 2012. The cancer cell line encyclopedia enables predictive modelling of anticancer drug sensitivity. *Nature*, 483(7391):603– 607.

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR. 687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709 710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

- Marco Colonna. 2023. The biology of trem receptors. *Nature Reviews Immunology*, 23(9):580–594.
- Haotian Cui, Chloe Wang, Hassaan Maan, Kuan Pang, Fengning Luo, Nan Duan, and Bo Wang. 2024a. scgpt: toward building a foundation model for singlecell multi-omics using generative ai. *Nature Methods*, pages 1–11.
- Haotian Cui, Chloe Wang, Hassaan Maan, Kuan Pang, Fengning Luo, Nan Duan, and Bo Wang. 2024b. scgpt: toward building a foundation model for singlecell multi-omics using generative ai. *Nature Methods*, pages 1–11.
- Haotian Cui, Chloe Wang, Hassaan Maan, Kuan Pang, Fengning Luo, and Bo Wang. 2023. scgpt: Towards building a foundation model for single-cell multiomics using generative ai. *bioRxiv*, pages 2023–04.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- David van Dijk, Juozas Nainys, Roshan Sharma, Pooja Kaithail, Ambrose J Carr, Kevin R Moon, Linas Mazutis, Guy Wolf, Smita Krishnaswamy, and Dana Pe'er. 2017. Magic: A diffusion-based imputation method reveals gene-gene interactions in single-cell rna-sequencing data. *BioRxiv*, page 111591.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.
- Gökcen Eraslan, Lukas M Simon, Maria Mircea, Nikola S Mueller, and Fabian J Theis. 2019. Singlecell rna-seq denoising using a deep count autoencoder. *Nature communications*, 10(1):390.
- Patrick Esser, Robin Rombach, and Bjorn Ommer. 2021. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883.
- Antonio Fabregat, Steven Jupe, Lisa Matthews, Konstantinos Sidiropoulos, Marc Gillespie, Phani Garapati, Robin Haw, Bijay Jassal, Florian Korninger, Bruce May, et al. 2018. The reactome pathway knowledgebase. *Nucleic acids research*, 46(D1):D649–D655.

- 740 741 742
- 744
- 745
- 747 748
- 751
- 756
- 757
- 759 760
- 761
- 762 763

- 772 775 776
- 781
- 782

790

794

- Jellert T Gaublomme, Nir Yosef, Youjin Lee, Rona S Gertner, Li V Yang, Chuan Wu, Pier Paolo Pandolfi, Tak Mak, Rahul Satija, Alex K Shalek, et al. 2015. Single-cell genomics unveils critical regulators of th17 cell pathogenicity. Cell, 163(6):1400–1412.
- Robert Gray. 1984. Vector quantization. IEEE Assp Magazine, 1(2):4–29.
- Dominic Grün, Lennart Kester, and Alexander Van Oudenaarden. 2014. Validation of noise models for single-cell transcriptomics. Nature methods, 11(6):637-640.
- Chuang Guo, Bin Li, Huan Ma, Xiaofang Wang, Pengfei Cai, Qiaoni Yu, Lin Zhu, Liying Jin, Chen Jiang, Jingwen Fang, et al. 2020. Single-cell analysis of two severe covid-19 patients reveals a monocyte-associated and tocilizumab-responding cytokine storm. Nature communications, 11(1):3924.
- Minsheng Hao, Jing Gong, Xin Zeng, Chiming Liu, Yucheng Guo, Xingyi Cheng, Taifeng Wang, Jianzhu Ma, Le Song, and Xuegong Zhang. 2023. Large scale foundation model on single-cell transcriptomics. bioRxiv, pages 2023-05.
- Minsheng Hao, Jing Gong, Xin Zeng, Chiming Liu, Yucheng Guo, Xingyi Cheng, Taifeng Wang, Jianzhu Ma, Xuegong Zhang, and Le Song. 2024a. Largescale foundation model on single-cell transcriptomics. Nature Methods, pages 1–11.
- Minsheng Hao, Jing Gong, Xin Zeng, Chiming Liu, Yucheng Guo, Xingyi Cheng, Taifeng Wang, Jianzhu Ma, Xuegong Zhang, and Le Song. 2024b. Largescale foundation model on single-cell transcriptomics. *Nature Methods*, pages 1–11.
- Mo Huang, Jingshu Wang, Eduardo Torre, Hannah Dueck, Sydney Shaffer, Roberto Bonasio, John I Murray, Arjun Raj, Mingyao Li, and Nancy R Zhang. 2018. Saver: gene expression recovery for single-cell rna sequencing. Nature methods, 15(7):539-542.
- Minyoung Huh, Brian Cheung, Pulkit Agrawal, and Phillip Isola. 2023. Straightening out the straightthrough estimator: Overcoming optimization challenges in vector quantized networks. arXiv preprint arXiv:2305.08842.
- Francesco Iorio, Theo A Knijnenburg, Daniel J Vis, Graham R Bignell, Michael P Menden, Michael Schubert, Nanne Aben, Emanuel Gonçalves, Syd Barthorpe, Howard Lightfoot, et al. 2016. A landscape of pharmacogenomic interactions in cancer. Cell, 166(3):740-754.
- Chintan J Joshi, Wenfan Ke, Anna Drangowska-Way, Eyleen J O'Rourke, and Nathan E Lewis. 2022. What are housekeeping genes? PLoS computational biology, 18(7):e1010295.
- Dragomirka Jovic, Xue Liang, Hua Zeng, Lin Lin, Fengping Xu, and Yonglun Luo. 2022. Single-cell rna sequencing technologies and applications: A

brief overview. Clinical and Translational Medicine, 12(3):e694.

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

- Peter V Kharchenko, Lev Silberstein, and David T Scadden. 2014. Bayesian approach to single-cell differential expression analysis. *Nature methods*, 11(7):740– 742.
- Wei Vivian Li and Jingyi Jessica Li. 2018. An accurate and robust imputation method scimpute for singlecell rna-seq data. Nature communications, 9(1):997.
- Qingnan Liang, Yuefan Huang, Shan He, and Ken Chen. 2023. Pathway centric analysis for single-cell rnaseq and spatial transcriptomics data with gsdensity. Nature communications, 14(1):8416.
- Qiao Liu, Zhiqiang Hu, Rui Jiang, and Mu Zhou. 2020. Deepcdr: a hybrid graph convolutional network for predicting cancer drug response. Bioinformatics, 36(Supplement_2):i911-i918.
- Romain Lopez, Jeffrey Regier, Michael B Cole, Michael I Jordan, and Nir Yosef. 2018. Deep generative modeling for single-cell transcriptomics. Nature methods, 15(12):1053-1058.
- Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101.
- Colin Megill, Bruce Martin, Charlotte Weaver, Sidney Bell, Lia Prins, Seve Badajoz, Brian McCandless, Angela Oliveira Pisco, Marcus Kinsella, Fiona Griffin, et al. 2021. Cellxgene: a performant, scalable exploration platform for high dimensional sparse matrices. bioRxiv, pages 2021-04.
- Huaiyu Mi, Anushya Muruganujan, John T Casagrande, and Paul D Thomas. 2013. Large-scale gene function analysis with the panther classification system. Nature protocols, 8(8):1551–1566.
- Anoop P Patel, Itay Tirosh, John J Trombetta, Alex K Shalek, Shawn M Gillespie, Hiroaki Wakimoto, Daniel P Cahill, Brian V Nahed, William T Curry, Robert L Martuza, et al. 2014. Single-cell rnaseq highlights intratumoral heterogeneity in primary glioblastoma. Science, 344(6190):1396-1401.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. OpenAI blog, 1(8):9.
- Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. 2019. Generating diverse high-fidelity images with vq-vae-2. Advances in neural information processing systems, 32.
- André F Rendeiro, Thomas Krausgruber, Nikolaus Fortelny, Fangwen Zhao, Thomas Penz, Matthias Farlik, Linda C Schuster, Amelie Nemc, Szabolcs Tasnády, Marienn Réti, et al. 2020. Chromatin mapping and single-cell immune profiling define the temporal dynamics of ibrutinib response in cll. Nature Communications, 11(1):577.

- 851 852 853 854
- 855 856
- 8! 8!
- 0
- 0 8
- 8 8 8

- 869 870
- 871 872
- 873 874 875
- 81
- 88
- 8

- 887 888
- 890
- 891
- 892 893
- 894 895

- 899 900

- Alex Rogozhnikov. 2022. Einops: Clear and reliable tensor manipulations with einstein-like notation. In International Conference on Learning Representations.
- Stefan Semrau, Johanna E Goldmann, Magali Soumillon, Tarjei S Mikkelsen, Rudolf Jaenisch, and Alexander Van Oudenaarden. 2017. Dynamics of lineage commitment revealed by single-cell transcriptomics of differentiating embryonic stem cells. *Nature communications*, 8(1):1096.
- Barkur S Shastry. 2009. Snps: impact on gene function and phenotype. *Single nucleotide polymorphisms: methods and protocols*, pages 3–22.
- Christina V Theodoris, Ling Xiao, Anant Chopra, Mark D Chaffin, Zeina R Al Sayed, Matthew C Hill, Helene Mantineo, Elizabeth M Brydon, Zexian Zeng, X Shirley Liu, et al. 2023. Transfer learning enables predictions in network biology. *Nature*, pages 1–9.
- Vincent A Traag, Ludo Waltman, and Nees Jan Van Eck. 2019. From louvain to leiden: guaranteeing well-connected communities. *Scientific reports*, 9(1):5233.
- Florian T Unger, Irene Witte, and Kerstin A David. 2015. Prediction of individual response to anticancer therapy: historical and future perspectives. *Cellular and molecular life sciences*, 72:729–757.
- Aaron Van Den Oord, Oriol Vinyals, et al. 2017. Neural discrete representation learning. *Advances in neural information processing systems*, 30.
- A Vaswani. 2017. Attention is all you need. Advances in Neural Information Processing Systems.
- Zheng Wang and David R Sherwood. 2011. Dissection of genetic pathways in c. elegans. *Methods in cell biology*, 106:113–157.
- Hongzhi Wen, Wenzhuo Tang, Xinnan Dai, Jiayuan Ding, Wei Jin, Yuying Xie, and Jiliang Tang. 2023.Cellplm: Pre-training of cell language model beyond single cells. *bioRxiv*, pages 2023–10.
- Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. 2021. Videogpt: Video generation using vq-vae and transformers. *arXiv preprint arXiv:2104.10157*.
- Fan Yang, Wenchuan Wang, Fang Wang, Yuan Fang, Duyu Tang, Junzhou Huang, Hui Lu, and Jianhua Yao. 2022. scbert as a large-scale pretrained deep language model for cell type annotation of singlecell rna-seq data. *Nature Machine Intelligence*, 4(10):852–866.
- In-Kwon Yeo and Richard A Johnson. 2000. A new family of power transformations to improve normality or symmetry. *Biometrika*, 87(4):954–959.

A Related Work

901

902

903

904

905

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

924

926

927

928

930

Vector Quantization (VQ). In terms of tokenizaiton, VQ (Gray, 1984) is a classical quantization technique in various fields. VQ operates by utilizing a *codebook*, which consists of multiple representations referred to as codes, and associating an input vector with the code within the codebook that is closest in proximity. This approach can be viewed as a form of discrete representation learning since typically only a single token is activated. Researchers have demonstrated that the usage of discrete representation in computer vision can improve the robustness of models. VQ techniques have found application in diverse fields, such as image generation (Van Den Oord et al., 2017; Razavi et al., 2019; Esser et al., 2021), video generation (Yan et al., 2021), and speech recognition (Baevski et al., 2019). In this work, we take a pioneering step by applying the concept of VQ to the field of biosciences, to learn discrete representations capturing genetic pathways.

B Dataset Description

Our pretraining data comes from CELLx-GENE (Megill et al., 2021), which has 43, 312, 189 cells from more than 8 tissues.

C More Details on Methodology

Algorithm 1 Permutation-Invariant Embedding

- Input: A batch of normalized expression value vectors X ∈ ℝ^{B×Ng} where B is the batch size, learnable gene embeddings E_{gene}, learnable projection matrix P₁ ∈ ℝ^{Ng×d}, P₂ ∈ ℝ^{d×d}, P' ∈ ℝ^{d×h}.
- 2: **Output**: Embedding \bar{E} of X, $f(\cdot)$ is a symmetric pooling function.
- 3: $X \leftarrow \text{einsum}(\text{"bi,ij->bij"}, X, P_1)$
- 4: $X \leftarrow \text{leaky_relu}(X)$
- 5: $E_{\text{value}} \leftarrow ext{einsum}(" ext{ij,bij->bij}", m{lpha}, m{X}) + ext{MLP}(m{X})$
- 6: Fill positions of zeros in E with a learnable embedding e_{zero}
- 7: $E \leftarrow \mathsf{concat}(E_{\mathsf{value}}, E_{\mathsf{gene}})$

8:
$$E \leftarrow einsum("bj, jl->bjl", f(E), P')$$

9: return $oldsymbol{E}$

Permutation-Invariant Embedding. Algorithm 1 transforms gene expression vectors into permutation-invariant embeddings. It projects input values into latent features via learnable

matrices, combines them with gene embeddings, applies symmetric pooling to aggregate features, and produces final embeddings through linear transformation. This approach reduces computational costs while maintaining invariance to gene order. The einsum operation denotes the Einstein summation convention, performing element-wise operations and summation along specified axes indicated by the letters (Rogozhnikov, 2022).

Genetic Pathway Learning. Algorithm 2 outlines the training pipeline for cell representation learning. It first encodes gene expressions into pathway tokens \mathcal{Z} and aggregates cell-level information via a learnable token e_C appended to masked representation. When contrastive learning is enabled, K-Means periodically updates pseudo-labels using stored embeddings in queue Q, and Equation 2 optimizes cluster consistency. The model trains e_C , e_M , and networks using reconstruction loss for masked tokens while updating h_C in Q.

Algorithm 2 Training Pipeline (One Step)

- 1: **Input**: A batch of gene expression vectors $X \in \mathbb{N}^{B \times N_g}$, current time step T, an interval for re-fit T_r , a K-Means classifier K, a queue Q.
- 2: Obtain $\bar{E} \in \mathbb{R}^{B \times N \times d}$ according to §3.1.
- 3: Obtain $\mathcal{Z} \in \mathbb{R}^{B \times N \times d}$ according to §3.2.
- 4: Randomly select p% tokens from Z for each sample and prepend a token e_C. Obtain H and assign clusters.
- 5: if use token-level contrastive learning then
- 6: **if** $T\%T_r == 0$ **then**
- 7: Run K-means based on embeddings in Q.
- 8: **end if**
- 9: Randomly select p% tokens from \mathcal{Z} for each sample and prepend e_C . Obtain \mathcal{H}' .
- 10: Calculate the contrastive learning loss according to Equation 2.
- 11: end if
- 12: Fill e_M into \mathcal{H} at positions previously excluded during sampling.
- 13: Train e_C , e_M , and the two networks with reconstruction loss for excluded tokens (*i.e.*, e_M).
- 14: Store h_C in Q.

Mask Construction and Output Reshaping. In

Algorithm 3, we introduce a simple way to make the shape of the output from the decoder introduced in §3.3 consistent with the original gene expression 931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

vector. Essentially, we flatten the hidden represen-955 tations, and we assign a region according to the 956 indices of leave-out tokens in which we calculate 957 the MSE loss.

959

961

962

963

964

965

966

967

969

970

971

972

973

974

VQ-Techniques For Stable Training. To address the potential index collapse when applying VQ techniques in training neural networks, we follow the pipeline introduced in Huh et al. (Huh et al., 2023). Firstly, they introduce an affine transformation to reparameterize the representation in the codebook with the following formula:

$$\boldsymbol{v}_i = \boldsymbol{c}_{ ext{mean}} + \boldsymbol{c}_{ ext{std}} imes \boldsymbol{c}_i$$

where the c_i represents the original code vector, and c_{mean} and c_{std} indicate the shared affine parameters. Moreover, they introduce several minor modifications to the codebook update process to enhance the stability.

More Details on Experiments D

Model Configurations. Table 4 presents the hyperparameters of our scPaLM.

Table 4: Configurations of our scPaLM.

Hyperparameters	Value
Hidden Size	768
Intermediate Size	3072
Number of Layers	12
Number of Attention Heads	8
Dropout Probability	0.0
Attention Dropout Probability	0.0

Algorithm 3 Reshaping Masks and Outputs For Loss Calculation.

- 1: Input: a gene expression vector $\boldsymbol{x} \in \mathbb{R}^{N_g}$, the hold-out indices $I = \{i_1, \ldots, i_m\}$, number of tokens N.
- 2: Calculate the scaling factor $s \leftarrow \lceil N_q/N \rceil$.
- 3: Initialize a mask vector $\boldsymbol{m} \leftarrow \mathbf{0}^{N_g}$.
- 4: for j = 1, 2, ..., m do

5:
$$\boldsymbol{m}_{s \times i_j: s \times (i_j+1)} \leftarrow \mathbf{1}^s$$
.

- 6: end for
- 7: Flatten the output from the decoder which also has the shape of $N \times s$ to $1 \times (N \times s)$, and store it as o.
- 8: Crop both o and m to have the length of N_g . Calculate the MSE loss as \mathcal{L}_{MSE} = $\|\boldsymbol{o} - \boldsymbol{x}\|_2^2 / \|\boldsymbol{m}\|_1.$

Association of Genes with Pathway Tokens. In this section, we describe the methodology employed for associating genes with specific pathway identifiers through an algorithmic approach. The process involves the utilization of a matrix with a dimension of K by N_g , where K represents the number of tokens and N_q the number of genes. For each vector of gene expression, we obtain the set of tokens that are activated within the codebook. Upon activation of the K_i -th token, the corresponding raw gene expression vector is scaled by the frequency of K_i token occurrences among the activated tokens and subsequently aggregated to the K_i -th row of the matrix. This procedure is iterated across the entire gene dataset. Subsequent to the completion of this iterative process, we perform a normalization step on each column, which correlates to individual genes. Following normalization, for each row, we identify and select the genes that exhibit the most significant values (Colonna, 2023).

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

Table 5: Clustering performance of different variants for cell representations on the CLL dataset. We compare the adjusted rand index (ARI), normalized mutual information (NMI), silhouette score (S-score), and clustering time between models.

Method	ARI ↑	NMI \uparrow	S-score ↑
Mean	0.015	0.059	-0.133
Concatenated	0.181	0.478	0.310
$oldsymbol{h}_C$ (No CL)	0.275	0.573	0.361
h_C (Ours)	0.292	0.593	0.376

Ε **More Experimental Results**

Additional Drug Response Prediction Result Analysises. In these experiments, we follow the setting of scFoundation and disable the mutation and the methylation features, to focus on the benefit brought by the incorporation of embeddings from 1000 gene expression values. From Figure 6 and 7, we 1001 can observe that scPaLM achieves better PCCs on 1002 all but one cancer type and improves the metrics on a majority of cell lines. Following the analysis, we 1004 further visualize the best prediction case of the can-1005 cer type, namely the low-grade gliomas (LGG) in 1006 Figure 7, where we observe both methods achieve 1007 high PCC values despite that the IC50 values have 1008 a large range from -6 to 6. scPaLM outperforms 1009 scFoundation by 2% and 4% in terms of the PCC 1010 and the Spearman correlation coefficient. These 1011 results showcase the effectiveness of scPaLM. It 1012



Figure 6: Pairwise visualization of the Pearson correlation coefficient of scFoundation and scPaLM based on different grouping strategies. Left: grouping with respect to the cell lines; Middle: grouping with respect to the cancer type; Right: grouping with respect to the drug type. The red lines indicate the relationship of y = x.

is also noteworthy that the embeddings generated by scPaLM are smaller in dimension compared to those of scFoundation, which implies that scPaLM is more efficient in modeling scRNA-seq data.

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1030



Figure 7: Scatter plots of the predicted and the observed IC50 values on the samples with the cancer type of low-grade gliomas.

Genetic Pathway Results. We present the identified significant genetic pathway in Table 7. Two pathways with p-values that are significantly smaller than 1×10^{-5} are identified, which are all related to immunological reactions.

> Table 6: Comparison between different embedding processes. The models are trained on a subset of the pretraining data and evaluated on the COVID dataset.

Embedding Algorithm	Memory Usage \downarrow	ARI↑	NMI↑
Per-gene Shared-first-layer	>80G 42378MiB	- 0.167	0.251
Permutation-invariant (Ours)	43678MiB	0.201	0.291

Comparison of Different Embedding Processes. We conduct an experiment to compare the memory usage and the subsequent clustering performance of models with different embedding processes on the COVID dataset. The results are presented in Table 6.

F Potential Risks

While scPaLM demonstrates promising capabilities for therapeutic discovery and cellular analy-

Table 7: Significant pathways identified by GSEA (p-value $<1\times10^{-5}$). Two pathways related to the immune system are selected.

Gene Lists (To-	Term	P-value
ken ID)		
CD1E,	Immunoregulatory Interac-	$2.8 imes 10^{-7}$
TREM2,	tions Between A Lymphoid	
ICAM5,	And A non-Lymphoid Cell	
CD1B (25)		
BTN1A1,	Adaptive Immune System	4.4×10^{-7}
MRC1, CD1E,		
TREM2,		
ICAM5,		
CD1B (25)		

sis, we acknowledge the dual-use potential inherent in any foundational biomedical AI technology. The model's ability to predict cell-drug interactions could theoretically be misapplied to screen compounds with harmful biological activity. To mitigate this risk, we emphasize that any clinical translation of our method must occur within established regulatory frameworks requiring rigorous safety evaluation and ethical oversight. Researchers utilizing this technology should adhere to institutional bio-safety review processes and existing chemical/biological weapons conventions to prevent misuse.

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041