
Cell ontology guided transcriptome foundation model

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Transcriptome foundation models (TFMs) hold great promises of deciphering the
2 transcriptomic language that dictate diverse cell functions by self-supervised learn-
3 ing on large-scale single-cell gene expression data, and ultimately unraveling the
4 complex mechanisms of human diseases. However, current TFMs treat cells as
5 independent samples and ignore the taxonomic relationships between cell types,
6 which are available in cell ontology graphs. We argue that effectively leveraging
7 this ontology information during the TFM pre-training can improve learning biolog-
8 ically meaningful gene co-expression patterns while preserving TFM as a general
9 purpose foundation model for downstream zero-shot and fine-tuning tasks. To this
10 end, we present single cell, Cell-ontology guided TFM (scCello). We introduce
11 cell-type coherence loss and ontology alignment loss, which are minimized along
12 with the masked gene expression prediction loss during the pre-training. The novel
13 loss component guide scCello to learn the cell-type-specific representation and the
14 structural relation between cell types from the cell ontology graph, respectively.
15 We pre-trained scCello on 22 million cells from CellxGene database leveraging
16 their cell-type labels mapped to the cell ontology graph from Open Biological
17 and Biomedical Ontology Foundry. Our TFM demonstrates competitive general-
18 ization and transferability performance over the existing TFMs on biologically
19 important tasks including identifying novel cell types of unseen cells, prediction of
20 cell-type-specific marker genes, and cancer drug responses.

21 1 Introduction

22 Cells are basic units of all living organisms. Deciphering diverse cell functions through gene
23 expression is a long-standing challenge in life science and yet the essential path towards precision
24 and personalized medicine. In this context, single-cell RNA sequencing (scRNA-seq) has emerged as
25 a pivotal technique to measure the gene expression in individual cells. The vast amount of publicly
26 available scRNA-seq data offers a rich transcriptomic data source [44] for learning cell representations
27 towards various research applications, such as cancer therapy [56] and drug discovery [4].

28 Recently, several *Transcriptome Foundation Models* (TFMs) were developed to improve cell represen-
29 tation learning. They mainly utilize pre-training methods analogous to natural language processing
30 like masked token prediction, treating genes as “tokens” and cells as “sentences” [14, 55, 63, 51].
31 However, the existing TFMs treat cells as independent samples and ignore their cell-type lineages. On
32 the other hand, prior knowledge of the taxonomic relationships of cell types has been made available
33 through the cell ontology graph by Open Biological and Biomedical Ontology Foundry [3]. Effectively
34 leveraging the ontology knowledge can improve the quality of the pre-training on large-scale
35 scRNA-seq atlases, which are heterogeneous and encompass hundreds of cell types. This can be
36 done by training the TFM to recognize the inherent ontology relationships among cell types, thereby
37 refining the cell representations. For instance, “mature α - β T cell” should be closer to “mature T

38 cells” compared to more general term “T cells” and farther from neurons and astrocytes from the
39 brain (e.g., Tab. 7).

40 To capture this intuition, we propose scCello, a single cell, **Cell-ontology** guided TFM. scCello
41 learns cell representation by integrating cell type information and cellular ontology relationships
42 into its pre-training framework. scCello’s pre-training framework is structured with three levels
43 of objectives: (1) **gene level**: a masked token prediction loss to learn gene co-expression patterns,
44 enriching the understanding of gene interactions (Sec. 2.2); (2) **intra-cellular level**: an ontology-
45 based cell-type coherence loss to encourage cell representations of the same cell type to aggregate,
46 prompting consistency between cells and their types (Sec. 2.3); and (3) **inter-cellular level**: a
47 relational alignment loss to guide the cell representation learning by consulting the cell-type lineage
48 from the cell ontology graph (Sec. 2.4)..

49 We demonstrate the generalizability and transferability of scCello on 22 million cells from CellxGene.
50 For model generalization, we observe that scCello excels on cell type identification across all datasets
51 in both zero-shot setting (i.e., directly using the pre-trained model) (Sec. 4.2.1) and fine-tuning setting
52 (Sec. 4.2.2). In particular, scCello accurately classifies novel cell types by leveraging the ontology
53 graph structure (Sec. 4.3). For transferability, scCello demonstrates competitive performances in pre-
54 dicting cell-type-specific marker genes (Sec. 4.4) and cancer drug responses (Sec. 4.5). Additionally,
55 scCello is robust against batch effects (Sec. 4.6). Finally, we validate our contribution via ablation
56 study (Sec. 4.7).

57 2 Method

58 Fig. 1 illustrates an overview of scCello. We present the details of individual components below.

59 2.1 Data Preprocessing

60 **Cell ontology graph.** Cell ontology is a widely used metadata schema for standard cell type
61 annotations [16]. We downloaded the ontology from Open Biological and Biomedical Ontology
62 Foundry (<https://obofoundry.org/>). It is structured as an unweighted directed acyclic graph
63 $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where each node $v \in \mathcal{V}$ corresponds to a distinct cell type and each directed edge
64 $(u, v) \in \mathcal{E}$ denotes a hierarchical lineage relationship of the form "is a subtype of" between cell types
65 (Fig. 1a). To accurately represent the inherently symmetric "being biologically similar" relationship
66 between cell types, the directed graph was transformed into an undirected one for subsequent
67 calculation of cellular ontology relationships in Sec. 2.4.

68 **scRNA-seq data.** The scRNA-seq data were downloaded from CellxGene. After the preprocessing
69 (App. B), we obtained 22 million cells. Each single-cell transcriptome is represented by a sequence
70 of tuples, each containing genes and their expression counts.¹ Each sequence was then ordered by the
71 rank of the gene expression values [55], akin to the sequential ordering of natural languages. Given a
72 batch of B cells, each cell $i \in \{1, \dots, B\}$ was assigned a cell type ontology identifier $c_i \in \mathcal{V}$ from
73 the CellxGene database, to enable mapping between cell and cell ontology.

74 2.2 Masked Gene Prediction

75 Same as BERT [15], scCello predicts a randomly masked gene token in each cell based on its
76 surrounding context in the sequence. This objective \mathcal{L}_{MGP} aims to learn the dynamic gene co-
77 expression network.

78 2.3 Intra-Cellular Ontology Coherence

79 A straightforward approach to encourage learning the cell representations that are coherent to the
80 cell type labels is to apply cross-entropy loss for supervised cell type classification. However, this
81 approach is limited in learning cell representation for the foundation model. Instead, we employed
82 a supervised contrastive loss as our objective $\mathcal{L}_{\text{Intra}}$, which directly optimizes the TFM rather than

¹scRNA-seq data was from CellxGene database <https://cellxgene.cziscience.com/>.

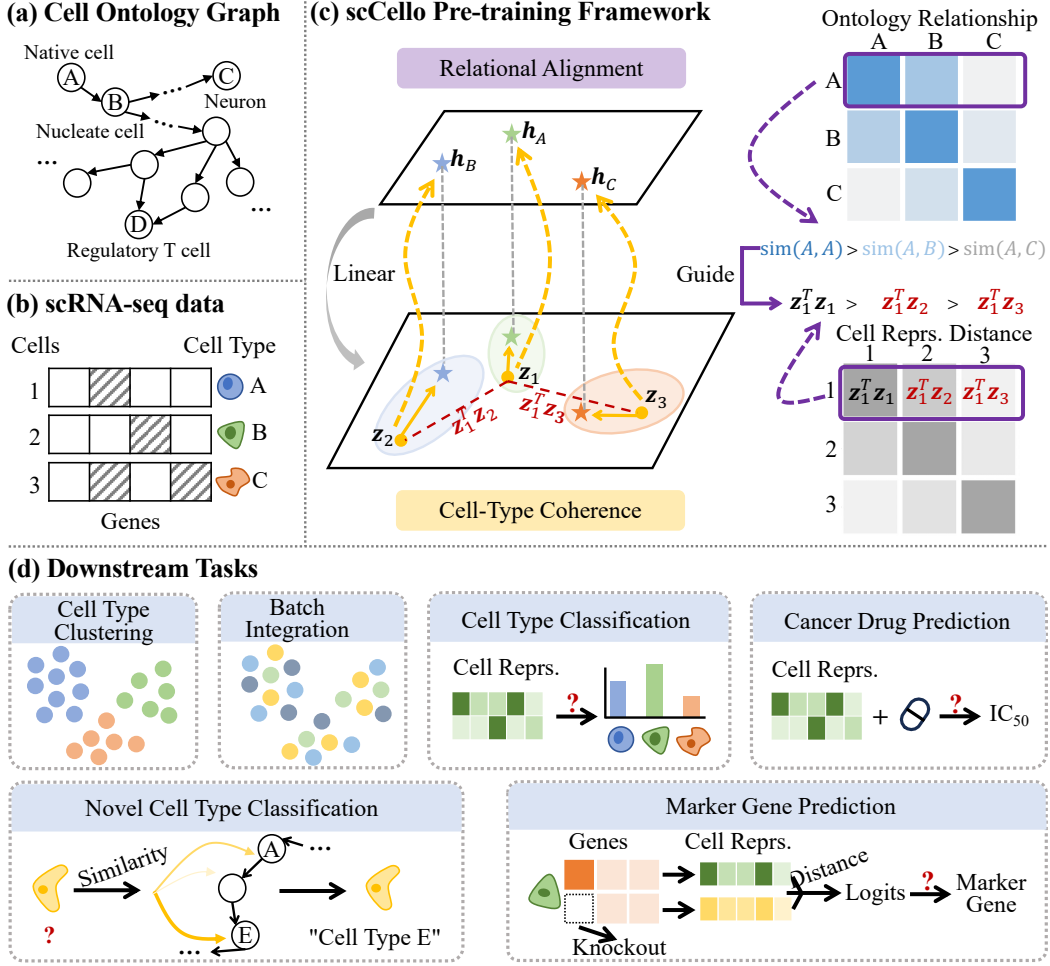


Figure 1: (a) Cell ontology graph describes taxonomic relationships between cell types. (b) Each cell in scRNA-seq data is represented by gene sequences, and associated with a cell type ontology identifier. (c) The pre-training framework of scCello is structured with three levels of objectives: gene-level masked gene prediction, intra-cellular level cell type coherence and inter-cellular level ontology alignment. For example, as shown in panel b, cells 1, 2, and 3 are labelled with cell type A, B and C. The intra-cellular cell type coherence loss encourages alignment of embedding z_1 with h_A , z_2 with h_B , and z_3 with h_C . The inter-cellular level ontology alignment loss encourages representational learning of cell similarities $z_i^\top z_j$ between cell i and j to be consistent to the similarity of their corresponding cell types $sim(c_i, c_j)$ based on the ontology relationships. (d) Downstream tasks enabled by scCello and demonstrated in the study.

83 merely learning through the linear classifier:

$$\mathcal{L}_{\text{Intra}} = - \sum_{i=1}^B \log \left(\frac{\exp(z_i^\top h_{c_i} / \tau)}{\exp(z_i^\top h_{c_i} / \tau) + \sum_{j=1, j \neq i}^B \exp(z_i^\top h_{c_j} / \tau)} \right). \quad (1)$$

84 where z_i and h_{c_i} denote the latent representation of cell i and cell type c_i , respectively.

85 This supervised contrastive loss pulls representations of the same class (positives) and repels rep-
 86 presentations of different classes (negatives). It often leads to representations that are at least as
 87 discriminative as the cross-entropy loss [22]. However, both cross entropy and contrastive loss are
 88 prone to class collapse, where all samples in a class are mapped to the same representation [30, 11].
 89 The resulting model may produce simplistic representations that perform well on similar training
 90 tasks like cell type clustering or classification but generalize poorly to new tasks. This defeats the

91 purpose of pre-training a versatile and general-purpose TFM. To tackle this limitation, we introduce a
 92 regularization term \mathcal{L}_{Reg} :

$$\mathcal{L}_{\text{Reg}} = \sum_{i=1}^B \|\text{Linear}(\mathbf{h}_{c_i}) - \mathbf{z}_i\|_2^2, \quad (2)$$

93 where the linear layer is shared across all cells and cell types. Thereby, it constrains the cell type
 94 representation space to be an affine transformation of the cell representation space, thus reducing the
 95 degrees of freedom available for TFM optimization and the chance for class collapse.

96 2.4 Inter-Cellular Relational Alignment

97 To encourage TFMs to learn inter-cellular ontology relationships, scCello forces cell representations
 98 to truthfully reflect the pairwise node structural similarity derived from the cell ontology graph, using
 99 a relational alignment objective. This objective constitutes the most important part of scCello.

100 **Ontology relationships.** To effectively quantify ontology relationships between cell types from the
 101 ontology graph, scCello estimates pairwise node structural similarities as proxies using Personalized
 102 PageRank (PPR) [20]. PPR is a graph learning algorithm. The PPR score $\text{PPR}(u, v)$ estimates the
 103 probability for a random walk. It starts from a given target node $u \in \mathcal{V}$ and terminates at another node
 104 $v \in \mathcal{V}$. Importantly, this is a context-sensitive structural similarity measure that accounts both direct
 105 connections and broader subgraph patterns [60]. It also provides robustness against variations in
 106 global network structures, such as variable node degrees and clustering coefficients [10]. To improve
 107 robustness (as justified in App. A), we transform $\text{PPR}(\cdot)$ through a non-linear function to derive the
 108 structural similarities $\text{sim}(\cdot)$ as ontology relationships tunable by a hyper-parameter threshold s :

$$\text{sim}(u, v) = \begin{cases} \lfloor \log_2(\frac{\text{PPR}(u, v)}{s} + 1) \rfloor, & \text{if } \text{PPR}(u, v) \geq s \\ 1, & \text{otherwise} \end{cases}. \quad (3)$$

109 **Relational alignment.** Cells with closely related cell types tend to be more similar than those with
 110 distinct cell types. This observation guides scCello to align the distances between cell representations
 111 *w.r.t.* a target cell, with their structural similarities $\text{sim}(\cdot)$ (as shown in Fig. 1c). Specifically, given a
 112 batch of B cells, if we consider a target cell i and another cell in the batch $j \neq i$, the representation
 113 distance $\mathbf{z}_i^T \mathbf{z}_j$ should reflect their structural similarity $\text{sim}(c_i, c_j)$. Accordingly, a negative sample
 114 set $\Omega_{i, j} = \{k | \text{sim}(c_i, c_j) > \text{sim}(c_i, c_k), 1 \leq k \leq B\}$ can be produced, where cell pair (i, k) are
 115 considered less similar to the cell pair (i, j) and should be contrasted against in the representation
 116 space using the objective $\mathcal{L}_{\text{Inter}}$:

$$\mathcal{L}_{\text{Inter}} = - \sum_{i=1}^B \sum_{j=1, j \neq i}^B \log \left(\frac{\exp(\mathbf{z}_i^T \mathbf{z}_j / \tau)}{\exp(\mathbf{z}_i^T \mathbf{z}_j / \tau) + \sum_{k \in \Omega_{i, j}} \exp(\mathbf{z}_i^T \mathbf{z}_k / \tau)} \right). \quad (4)$$

117 Notably, ancestor cell types, which can reach the target cell type via the directed "is a subtype of"
 118 edge on the ontology graph, are structurally distant from the target cell type. Despite being distant,
 119 they fall into the same, broader cell type category. Contrasting cells associated with these distant
 120 ancestor cell types with the target cell is counter-intuitive. Therefore, scCello explicitly excludes such
 121 cells from the negative sample set, avoiding inappropriately pushing away biologically similar cells.
 122 This enhances scCello's capability to discern subtle similarities and differences within the cell types.

123 2.5 Overall Pre-training Objective

124 During pre-training, we seek to minimize the loss functions of all pre-training tasks simultaneously:

$$\theta^* \leftarrow \arg \min_{\theta} \mathcal{L}_{\text{MGP}} + \mathcal{L}_{\text{Inter}} + \mathcal{L}_{\text{Intra}} + \mathcal{L}_{\text{Reg}} \quad (5)$$

125 where θ denotes all learnable parameters in scCello, which adopts transformer stacks as model
 126 backbones. We state the detailed information of model architectures in App. C.

127 3 Related Work

128 The rapid growth of scRNA-seq datasets has opened new avenues for constructing TFMs, enabling
 129 transfer learning across various biological downstream tasks. Initial efforts, such as scBERT [65],

130 Exceiver [13] and Geneformer [55], borrows the concept of masked language modeling [15] from
 131 natural language processing (NLP) domain for pre-training, by treating cells as sentences and genes as
 132 tokens. Concurrently, tGPT [54] and scGPT [14] explored generative modeling [49], and CellLM [67]
 133 adapted the idea of contrastive learning [38]. Following the concept of “scaling” towards emergent
 134 ability [63] in NLP, scFoundation [24] proposes the largest foundation model at the time in terms
 135 of model size and pre-training data size; scHyena [45] scales modeling context window size to the
 136 full length of scRNA-seq data with Hyena operator [47] instead of conventionally used transformers.
 137 scTab [18] is the first to explore large-scale supervised learning mechanism for scRNA-seq pre-
 138 training, and is capable of annotating unseen tissue cells for real-world applications. Moreover,
 139 SCimilarity [27] and UCE [51] focus on developing a unified latent space as a large-scale reference
 140 atlas for querying new cells. Yet, these TFMs mainly treat cells as independent samples during
 141 training and ignore their biological ontology relationships. scCello bridges this gap by incorporating
 142 cell type relationships derived from the cell ontology graph into TFM pre-training. This strengthens
 143 TFMs’ model generalization and transferability capability, as shown in Sec. 4.

144 4 Experiments

145 As an overview, the following experiments show that, **(1)** scCello can generalize to unseen cells,
 146 and to more difficult settings, such as cells of unseen cell types, tissues, and donors (Sec. 4.2.1);
 147 **(2)** scCello can benefit from fine-tuning on target datasets (Sec. 4.2.2); **(3)** the structural similarity
 148 embedded in scCello helps to classify novel cell types in a zero-shot manner (Sec. 4.3); **(4)** scCello
 149 effectively transfers to different downstream tasks (Sec. 4.4 and Sec. 4.5); **(5)** scCello is robust to
 150 batch effects that arise from different experimental conditions (Sec. 4.6); **(6)** Each loss component in
 151 Eqn. 5 is beneficial to scCello (Sec. 4.7). For every table reported, we used **bold** to highlight the best
 152 performance and results within 0.005 difference from the best. We used underlining to denote the
 153 second-best performances. For all metrics, \uparrow indicates the higher the better.

154 4.1 Setups

155 **Pre-training and downstream datasets.** We collected a large pre-training dataset consisting of 22
 156 million cells along with downstream datasets. In particular, we generated one in-distribution (ID)
 157 and six out-of-distribution (OOD) datasets (App. B). The ID dataset is denoted as D^{id} . For the OOD
 158 setting, we introduced three scenarios: unseen cell types ($\{D_i^{ct}\}_{i=1}^2$), unseen cell tissues ($\{D_i^{ts}\}_{i=1}^2$),
 159 and unseen donors ($\{D_i^{dn}\}_{i=1}^2$). Each scenario has two datasets. Notably, the OOD donor setting
 160 presents more realistic challenges than ID and other OOD settings because of the potential batch
 161 effects in the test donors.

162 **Pre-training configurations.** An Adam optimizer [35] (learning rate: 0.001, weight decay: 0.001,
 163 warm-up steps: 3, 333) was used to train the scCello for 40, 000 steps on 4 NVIDIA A100 GPUs on
 164 Compute Canada. We used 192 for batch size. More details are introduced in App. C.

165 **Baselines.** Across all downstream tasks, scCello is benchmarked with leading open-source large-
 166 scale TFMs: Geneformer [55], scGPT [14], scTab [18], UCE [51], and three TFM ablations. We also
 167 implemented ablated versions of scCello that only differ in the pre-training objectives from scCello:
 168 scCello using only the masked gene prediction loss (denoted as MGP), scCello using only the cell
 169 type supervised classification (denoted as Sup), and scCello using only the two losses (denoted as
 170 MGP+Sup). The three ablated TFMs provide a reference to isolate the effect of implementation
 171 details and training configurations. For each task, we also selected state-of-the-art non-TFM methods
 172 for fair comparison.

173 **Downstream metrics.** We evaluated the 3 tasks by the following metrics. (1) Clustering metrics
 174 include normalized mutual information (NMI), adjusted rand index (ARI), average silhouette width
 175 (ASW), and the average of the 3 scores (AvgBio) to assess both between-cluster separation and
 176 within-cluster closeness [14]. The batch integration task (Sec. 4.6) is evaluated by ASW_b , graph
 177 connectivity (GraphConn) and their average (AvgBatch), along with an overall score (Overall =
 178 $0.6 \times \text{AvgBio} + 0.4 \times \text{AvgBatch}$) to balance biological relevance and batch consistency following [14].
 179 (2) Classification metrics include accuracy (Acc), Macro F1 and area under the ROC curve (AU-
 180 ROC) [46]. (3) Regression task metrics include Pearson correlation coefficient score (PCC) [46].
 181 Details for each metric were provided in App. D.1.

Table 1: Zero-shot cell type clustering on the curated ID and OOD datasets.

Method	In-Distribution (ID)				Out-of-Distribution (OOD)						
	D^{id}				D_1^{ct}	D_2^{ct}	D_1^{ts}	D_2^{ts}	D_1^{dn}	D_2^{dn}	OOD Avg.↑
	NMI↑	ARI↑	ASW↑	AvgBio↑	AvgBio↑	AvgBio↑	AvgBio↑	AvgBio↑	AvgBio↑	AvgBio↑	
Non-TFM Methods											
Raw Data	0.566	0.237	0.453	0.419	0.703	0.629	0.540	0.631	0.458	0.460	0.570
Seurat	0.648	0.270	0.407	0.442	0.752	0.737	0.587	0.636	0.466	0.489	0.611
Harmony ¹	0.621	0.261	0.382	0.421	0.432	0.417	0.462	0.515	0.456	0.474	0.459
scVI	0.660	0.297	0.464	0.474	0.760	0.725	0.577	0.634	0.478	0.502	0.613
Ontology-Agnostic TFMs											
Geneformer	0.616	0.261	0.418	0.432	0.689	0.668	0.539	0.597	0.468	0.482	0.574
scGPT	0.615	0.258	0.442	0.438	0.707	0.720	0.544	0.627	0.456	0.477	0.589
scTab	<u>0.707</u>	<u>0.479</u>	0.544	<u>0.577</u>	<u>0.759</u>	0.726	0.515	0.657	OOM	OOM	/
UCE	0.670	0.304	0.494	0.489	0.772	0.741	0.598	0.670	0.485	0.506	0.629
MGP	0.662	0.306	0.451	0.473	0.714	0.740	0.576	0.628	0.488	0.518	0.611
Sup	0.703	0.393	<u>0.569</u>	0.555	0.767	<u>0.775</u>	<u>0.605</u>	<u>0.680</u>	0.552	<u>0.573</u>	<u>0.659</u>
MGP+Sup	0.661	0.337	0.550	0.516	0.758	0.764	0.610	0.672	<u>0.553</u>	0.570	0.655
Ontology-Enhanced TFMs											
scCello	0.785	0.558	0.667	0.670	0.769	0.786	0.612	0.705	0.608	0.643	0.687

¹ Harmony could be over-corrected *w.r.t.* batch labels for datasets with many batches [9].

182 4.2 Cell Type Identification

183 4.2.1 Zero-shot Cell Clustering Results

184 **Setup.** For the cell type clustering task, TFM baselines and four non-TFM methods were evaluated:
 185 (1) raw data expressions of highly variable genes (*abbr.*, Raw Data) [33]; (2) Seurat [26]; (3)
 186 Harmony [36] (4) scVI [40]. Cell representations were extracted from the baselines and clustered by
 187 Louvain algorithm [5]. We evaluated the clustering performance of each method on both ID dataset
 188 D^{id} and OOD datasets D_i^{cond} ($cond \in \{ct, ts, dn\}$, $i \in \{1, 2\}$).

189 **ID and OOD generalization.** We reported zero-shot cell type clustering performance in Tab. 1,
 190 and included all the metrics for all datasets in App. D.2.1 due to space constraint. For both the ID
 191 and OOD settings, scCello consistently outperforms all baselines, achieving a 16.1% improvement in
 192 AvgBio on the ID dataset and a 12.1% improvement in average AvgBio across the six OOD datasets.
 193 Interestingly, while scCello outperforms non-TFM methods by a large margin, Geneformers and
 194 scGPT barely surpass these methods. The latter is consistent with previous observations [66].

195 In the OOD experiments, scCello confers strong generalization capability across unseen cell types
 196 tissue, and donors. In cell type clustering, scCello is the second best only trailing UCE by 0.03
 197 and the best method for dataset 1 and 2. The OOD tissue setting highlights scCello’s ability to
 198 transfer its learned knowledge to different unseen tissues. Specifically, scCello achieve 0.6 and 0.7
 199 while most methods conferred below 0.6 and 0.7 for the two datasets, respectively. For the unseen
 200 OOD donor scenario, most methods perform poorly with AvgBio ranging between 0.45 and 0.55.
 201 scCello led the chart achieving AvgBio above 0.6 in both datasets. Overall, scCello showcases strong
 202 model generalization capabilities across a range of biological conditions, which is attributable to the
 203 integration of cell ontology priors during its TFM pre-training. Indeed, the ablated models namely
 204 MGP, Sup, and MGP+Sup conferred lower scores compared to the full model.

205 4.2.2 Fine-tuning Results

206 **Setup.** We benchmarked all TFM baselines except UCE for its lack of fine-tuning support. These
 207 TFMs were fine-tuned on a subset of our pre-training data with supervised classification loss (details
 208 in App. D.2.2). We assessed both classification and clustering performance on the ID dataset D^{id} .

209 **Improvement with fine-tuning.** In Tab. 2, The fine-tuned scCello outperforms other TFMs on
 210 both classification and clustering metrics, achieving up to 25.9% improvement in Macro F1 over
 211 the best baseline. Moreover, scCello without fine-tuning still surpasses the performance of the other
 212 fine-tuned methods, further highlighting its superior transferability.

213 4.3 Novel Cell Type Classification

214 Novel cell type classification aims to label cells of unseen cell types without further fine-tuning.
 215 This task is useful for annotating completely new scRNA-seq datasets but infeasible for most of the

Table 2: Cell type identification using fine-tuned TFMs. Both the classification and clustering performances on the ID dataset D^{id} are reported.

Method	Classification		Clustering
	Acc \uparrow	Macro F1 \uparrow	AvgBio \uparrow
Scratch	0.621	0.223	0.544
Ontology-Agnostic TFMs			
Geneformer	0.747	0.440	0.439
scGPT	0.712	0.344	0.477
scTab	0.778	0.373	0.606
MGP	0.722	0.287	0.607
Sup	0.812	0.363	<u>0.659</u>
MGP+Sup	<u>0.820</u>	<u>0.406</u>	<u>0.607</u>
Ontology-Enhanced TFMs			
scCello	0.867	0.511	0.694

Table 3: Marker gene prediction, a binary classification task to identify cell-type-specific marker genes.

Method	D_1^{mk}	D_2^{mk}	Avg. \uparrow
	AUROC \uparrow	AUROC \uparrow	
Ontology-Agnostic TFMs			
Geneformer	0.452	0.470	0.461
scGPT	0.385	0.387	0.386
scTab	0.672	0.727	0.700
UCE	0.500	0.500	0.500
MGP	0.579	0.629	0.604
Sup	0.699	<u>0.693</u>	0.696
MGP+Sup	<u>0.730</u>	0.730	<u>0.730</u>
Ontology-Enhanced TFMs			
scCello	0.756	0.729	0.743

Table 4: Cancer drug response prediction: a regression task to predict the IC_{50} values of drugs.

Method	Non-TFM Methods		Ontology-Agnostic TFMs						Ontology-Enhanced TFMs	
	DeepCDR	scFoundation	Geneformer	scGPT	scTab	UCE	MGP	Sup	MGP+Sup	scCello
PCC \uparrow	0.854	0.882	0.911	0.919	0.913	0.922	0.872	0.915	<u>0.916</u>	0.917

216 supervised methods that solely rely on the labels observed in the training data [7, 29, 61]. Leveraging
 217 the cell ontology graph that comprises the lineage relations among all of the known cell types, scCello
 218 makes this task feasible.

219 **Setup.** Our goal is to classify new query cells into "novel cell types" not seen during pre-training. To
 220 do this, we generate representations for both query cells and novel cell types, using similarity measures
 221 for classification. This process involves utilizing similarities between TFM-derived representations
 222 for the former and biological relationships from the cell ontology graph for the later. Details were
 223 described in App. D.3.

224 We benchmarked all TFMs and evaluated them on OOD cell type datasets D_1^{ct} and D_2^{ct} . We increased
 225 the difficulty of this task by the number of novel cell types (#Cell Types) that exist among the query
 226 cells. Specifically, we simulated five difficulty levels, with the number of novel cell types ranging
 227 from 10% to 100% of the total cell types. To assess the variance of the performance, we randomly
 228 sampled cell type combinations 20 times at each level.

229 **OOD generalization.** In Fig. 2, scCello led other TFMs by a large margin, achieving up to 76.8%
 230 Acc to classify 9 novel cell types (i.e., 10% of the total heldout cell types) and 33.5% Acc to classify
 231 up to 87 novel cell types (i.e., 100% of the total heldout cell types) (Tab. 16 and Tab. 17). These
 232 results show a significant leap from the existing TFMs, which either do not work or only work for
 233 annotating a handful of novel types [61, 41, 59].

234 4.4 Marker Gene Prediction

235 Cell-type-specific genes, or marker genes, are highly expressed in a specific cell type but exhibit low
 236 expression in others. These genes play a crucial role in delineating cell functions in diverse tissue
 237 contexts. Identifying marker genes in less characterized cell types is an ongoing challenge [48].

238 **Setup.** We sought to assess whether the pre-trained TFMs can discriminate marker from non-marker
 239 genes for any cell type without any supervised fine-tuning. This zero-shot experiment evaluates
 240 whether the TFM is able to learn biologically meaningful gene co-expression patterns without
 241 supervision. For each cell, we quantified the marker gene potential of each gene by the changes in
 242 TFM-generated cell representations after *in-silico* knockout of the target gene (details in App. D.4).
 243 Here we assume that the larger the change the higher the marker gene potential. We discussed the
 244 caveat of this approach in Sec. 5. As test data, we used GSE96583 [31] (D_1^{mk}) and GSE130148 [58]
 245 (D_2^{mk}). We obtained the marker gene labels from CellMarker2 [28] and PanglaoDB [21].

246 **Zero-shot transferability.** In Tab. 3, scCello outperforms other TFMs, improving upon the second-
 247 best method by 1.8% in average AUROC. The inclusion of cell label information during pre-training
 248 boosts TFM performance, as evidenced by the strong results of scTab, Sup, MGP+Sup and scCello.

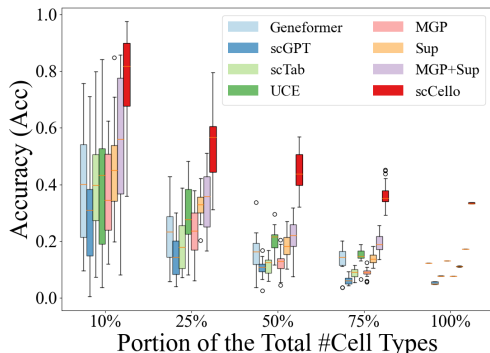


Figure 2: Novel cell type classification on OOD cell type dataset D_1^{ct} for increasing difficulties.

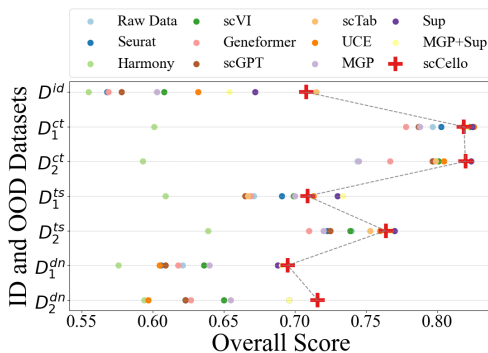


Figure 3: Batch integration on the curated ID and OOD datasets.

249 This is due to the biological correlation between marker genes and cell types. Furthermore, employing
 250 cell ontology graphs further improves the prediction accuracy over MGP+Sup.

251 4.5 Cancer Drug Response Prediction

252 Developing effective drugs for cancer treatment is challenging due to individual variability in drug
 253 responses. Accurately predicting cancer drug responses (CDR) can greatly aid anti-cancer drug
 254 development and improve our understanding of cancer biology [39].

255 **Setup.** Following the approach of scFoundation [25], cell representations were extracted from
 256 fixed TFMs and integrated into the DeepCDR [39] pipeline to estimate the half-maximal inhibitory
 257 concentration (IC_{50}) values of drugs (details in App. D.5). We benchmarked our method against
 258 DeepCDR, scFoundation, and other TFM baselines, using the same pre-processed data as DeepCDR.

259 **Zero-shot transferability.** In Tab. 4, scCello is among the top 3 along with scGPT and UCE,
 260 achieving 7.4% improvement in PCC over the base method DeepCDR. This highlights scCello’s
 261 transferability in enhancing specialized task-oriented methods. In particular, it can be used as an
 262 powerful feature extractor for diverse downstream tasks.

263 4.6 Batch Integration

264 The scRNA-seq atlases, assembled from datasets across various labs and conditions, are prone to
 265 unwanted technical variations known as batch effects [42]. These effects can significantly affect the
 266 generalization ability of TFMs especially because they require pre-training on a massive amount of
 267 heterogeneous scRNA-seq data pooled from many studies. Here we sought to evaluate scCello’s
 268 robustness to batch effects without fine-tuning.

269 **Setup.** We adopted the same baselines as in zero-shot cell type clustering (Sec. 4.2.1), and followed
 270 the evaluation protocol of scGPT [14]. We evaluated on one ID dataset D^{id} and six OOD datasets
 271 D_i^{cond} ($cond \in \{ct, ts, dn\}, i \in \{1, 2\}$) (see complete results of all metrics in App. D.6).

272 **Robustness to data noise.** Fig. 3 shows that scCello excels in 3 out of 7 datasets, and achieves
 273 comparable performance on another 3 datasets. The performance is attributable to the use of cell type
 274 information as the ablated baseline MGP conferred much lower batch integration score compared to
 275 Sup and scCello.

276 4.7 Ablation Study

277 **Ablation of pre-training losses.** Tab. 5 reports the cell type clustering (Sec. 4.2.1) and novel
 278 cell type classification (Sec. 4.3) performance of scCello by using full or partial pre-training losses.
 279 Removing any of the four losses in Eqn. 5 resulted in decreased performance, corroborating the
 280 benefits of the proposed pre-training losses. Notably, removing the inter-cellular ontology relation
 281 loss \mathcal{L}_{Inter} led to 56.1% and 65.3% decrease in terms of Acc. and Macro F1 on novel cell type
 282 classification task, respectively. This shows the upmost importance of the structurally induced loss
 283 and ultimately the use of cell ontology graph information.

Table 5: Pre-training loss ablation on the cell type clustering and novel cell type classification (*abbr.*, "clf.") tasks.

Config	Cell Type Clustering		Novel Cell Type Clf.	
	D_2^{ct}	D_2^{dn}	D_1^{ct}	
	AvgBio \uparrow	AvgBio \uparrow	Acc \uparrow	Macro F1 \uparrow
Full Loss	0.786	0.643	0.335	0.150
<i>w/o</i> \mathcal{L}_{MGP}	0.774 ($\downarrow 1.5\%$)	0.640 ($\downarrow 0.5\%$)	0.287 ($\downarrow 14.3\%$)	0.131 ($\downarrow 12.7\%$)
<i>w/o</i> \mathcal{L}_{Inter}	0.778 ($\downarrow 1.0\%$)	0.620 ($\downarrow 3.6\%$)	0.147 ($\downarrow 56.1\%$)	0.052 ($\downarrow 65.3\%$)
<i>w/o</i> \mathcal{L}_{Intra}	0.730 ($\downarrow 7.1\%$)	0.626 ($\downarrow 2.6\%$)	0.280 ($\downarrow 16.4\%$)	0.118 ($\downarrow 21.3\%$)
<i>w/o</i> \mathcal{L}_{Reg}	0.764 ($\downarrow 2.8\%$)	0.638 ($\downarrow 0.8\%$)	0.296 ($\downarrow 11.6\%$)	0.134 ($\downarrow 10.7\%$)

Table 6: Overall performance *v.s.* the number of parameters.

Method	Perf. Rank	#Params (M)
Geneformer	6.3	<u>10.3</u>
scGPT	6.2	51.3
scTab	4.2	9.7
UCE	4.8	674.7
MGP	6.7	<u>10.3</u>
Sup	3.3	10.4
MGP+Sup	<u>3.2</u>	10.9
scCello	1.3	10.7

284 **Parameter efficiency.** Tab. 6 demonstrates that scCello is highly parameter-efficient, utilizing up to
 285 60 times fewer parameters than the largest existing TFM, UCE, while still achieving the best average
 286 performance rankings across all downstream tasks. With an average performance rank of 1.3, scCello
 287 consistently ranks first or near the top in nearly every task.

288 **Visualization.** Visualization and analysis of scCello’s learned cell representations were presented
 289 in App. D.7. In short, biologically similar cell types are closer to each other and farther from those
 290 dissimilar ones in the t-SNE 2D space (Fig. 11).

291 5 Discussion and Conclusion

292 **Limitation and future work.** The cell ontology is constantly revised and expanded. In the future,
 293 we plan to investigate more efficient methods for fine-tuning scCello to enable continual learning of
 294 updated ontology, rather than retraining the entire model. Additionally, we aim to scale up the model
 295 size of scCello to increase its expressiveness and capacity. For the zero-shot marker gene prediction
 296 experiments (Sec. 4.4), one caveat is that our in-silico gene knockout approach also detects essential
 297 genes such as housekeeping genes [17] and transcription factors that are master regulators [8], which
 298 may not necessarily be marker genes. Nonetheless, deletion of these influential genes will also lead
 299 to large change of the transcriptome landscape of the cell. We will explore this in future study.

300 **Societal impact.** This work proposes a novel cell ontology-guided TFM, scCello, to enhance cell
 301 representation learning. On the positive side, once pre-trained, scCello can serve as a foundational
 302 model capable of facilitating scientific discoveries across various downstream tasks related to cells
 303 and cellular processes. However, on the negative side, the pre-training of scCello requires significant
 304 computational resources, potentially resulting in substantial carbon dioxide emissions that could
 305 contribute to environmental harm.

306 **Conclusion.** The proposed scCello incorporates cell ontology knowledge into its pre-training
 307 process by simultaneously modeling at the gene level, intra-cellular level, and inter-cellular level. We
 308 constructed a large-scale cell type identification benchmark to evaluate the model’s generalization
 309 capabilities, both in-distribution and out-of-distribution. Our evaluation demonstrates that scCello
 310 also exhibits strong transferability, as evidenced by its performance on other biologically meaningful
 311 downstream tasks such as zero-shot novel cell type classification and cell-type-specific marker gene
 312 prediction. Foundational models are typically heavy on the parameters for them to have sufficient
 313 capacity to learn from unlabeled data from scratch. This limits their usage to only fine-tuning tasks as
 314 pre-training them is prohibitive without large compute. Our proposed approach provides an efficient
 315 way of leveraging the prior knowledge at the pre-training, which led to much smaller parameter
 316 size while achieving performance comparable of the TFMs that are 5-60 times bigger. Together,
 317 scCello is a knowledge-informed and general purpose deep learning model that can be fine-tuned
 318 for a wide array of downstream applications, aiding in the rapid identification of novel cell types,
 319 disease-associated genes, and effective cancer drugs.

320 References

321 [1] Shibli Abdulla, Brian D. Aevertmann, Pedro Assis, Seve Badajoz, Sidney M. Bell, Emanuele
 322 Bezzi, Batuhan Cakir, Jim Chaffer, Signe Chambers, J. Michael Cherry, Tiffany Chi, Jennifer
 323 Chien, Leah Dorman, Pablo Garcia-Nieto, Nayib Gloria, Mim Hastie, Daniel Hegeman, Jason
 324 Hilton, Timmy Huang, Amanda Infeld, Ana-Maria Istrate, Ivana Jelic, Kuni Katsuya, Yang-Joon

- 325 Kim, Karen Liang, Mike Lin, Maximilian Lombardo, Bailey Marshall, Bruce Martin, Fran
326 McDade, Colin Megill, Nikhil Patel, Alexander V. Predeus, Brian Raymor, Behnam Robotmili,
327 Dave Rogers, Erica Rutherford, Dana Sadgat, Andrew Shin, Corinn Small, Trent Smith, Prathap
328 Sridharan, Alexander Tarashansky, Norbert Tavares, Harley Thomas, Andrew Tolopko, Meg
329 Urisko, Joyce Yan, Garabet Yeretssian, Jennifer Zamanian, Arathi Mani, Jonah Cool, and
330 Ambrose J. Carr. Cz cell×gene discover: A single-cell data platform for scalable exploration,
331 analysis and modeling of aggregated data. *bioRxiv*, 2023.
- 332 [2] Sercan Ö Arik and Tomas Pfister. Tabnet: Attentive interpretable tabular learning. In *Proceed-*
333 *ings of the AAAI conference on artificial intelligence*, volume 35, pages 6679–6687, 2021.
- 334 [3] Jonathan Bard, Seung Y Rhee, and Michael Ashburner. An ontology for cell types. *Genome*
335 *biology*, 6:1–5, 2005.
- 336 [4] Nurken Berdigaliyev and Mohamad Aljofan. An overview of drug discovery and development.
337 *Future medicinal chemistry*, 12(10):939–947, 2020.
- 338 [5] Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. Fast
339 unfolding of communities in large networks. *Journal of statistical mechanics: theory and*
340 *experiment*, 2008(10):P10008, 2008.
- 341 [6] Felipe A. Vieira Braga, Gozde Kar, Marijn Berg, Orestes A. Carpaij, Krzysztof Polański,
342 Lukas M. Simon, Sharon Brouwer, Tomás Gomes, Laura Hesse, Jian Jiang, Eirini Sofia Fasouli,
343 Mirjana Efremova, Roser Vento-Tormo, Carlos Talavera-López, Marnix R. Jonker, Karen
344 Affleck, Subarna Palit, Paulina M. Strzelecka, Helen V. Firth, Krishnaa T. Mahbubani, Ana
345 Cvejic, Kerstin B. Meyer, Kourosh Saeb-Parsy, Marjan A. Luinge, Corry-Anke Brandsma,
346 Wim Timens, Ilias Angelidis, Maximilian Strunz, Gerard H. Koppelman, Antoon J. M. van
347 Oosterhout, Herbert B. Schiller, Fabian J Theis, Maarten van den Berge, Martijn C. Nawijn, and
348 Sarah A. Teichmann. A cellular census of human lungs identifies novel cell states in health and
349 in asthma. *Nature Medicine*, 25:1153 – 1163, 2019.
- 350 [7] Maria Brbic, Marinka Zitnik, Sheng Wang, Angela Oliveira Pisco, Russ B. Altman, Spyros
351 Darmanis, and Jure Leskovec. Mars: discovering novel cell types across heterogeneous single-
352 cell experiments. *Nature Methods*, 17:1200 – 1206, 2020.
- 353 [8] Sunny Sun-Kin Chan and Michael Kyba. What is a master regulator? *Journal of stem cell*
354 *research & therapy*, 3, 2013.
- 355 [9] Ruben Chazarra-Gil, Stijn van Dongen, Vladimir Yu Kiselev, and Martin Hemberg. Flexible
356 comparison of batch correction methods for single-cell rna-seq using batchbench. *Nucleic acids*
357 *research*, 49(7):e42–e42, 2021.
- 358 [10] Fan Chen, Yini Zhang, and Karl Rohe. Targeted sampling from massive block model graphs
359 with personalized pagerank. *Journal of the Royal Statistical Society Series B: Statistical*
360 *Methodology*, 82(1):99–126, 2020.
- 361 [11] Mayee Chen, Daniel Y Fu, Avaniika Narayan, Michael Zhang, Zhao Song, Kayvon Fatahalian,
362 and Christopher Ré. Perfectly balanced: Improving transfer and robustness of supervised
363 contrastive learning. In *International Conference on Machine Learning*, pages 3090–3122.
364 PMLR, 2022.
- 365 [12] David Combe, Christine Largeron, Mathias Géry, and Előd Egyed-Zsigmond. I-louvain:
366 An attributed graph clustering method. In *Advances in Intelligent Data Analysis XIV: 14th*
367 *International Symposium, IDA 2015, Saint Etienne, France, October 22-24, 2015. Proceedings*
368 *14*, pages 181–192. Springer, 2015.
- 369 [13] William Connell, Umair Khan, and Michael J Keiser. A single-cell gene expression language
370 model. *arXiv preprint arXiv:2210.14330*, 2022.
- 371 [14] Haotian Cui, Chloe X. Wang, Hassaan Maan, Kuan Pang, Fengning Luo, Nan Duan, and
372 Bo Wang. scgpt: toward building a foundation model for single-cell multi-omics using genera-
373 tive ai. *Nature methods*, 2024.

- 374 [15] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
375 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,
376 2018.
- 377 [16] Alexander D. Diehl, Terrence F. Meehan, Yvonne M. Bradford, Matthew H. Brush, Wasila M.
378 Dahdul, David S. Dougall, Yongqun He, David Osumi-Sutherland, Alan Ruttenberg, Sirarat
379 Sarntivijai, Ceri E. Van Slyke, Nicole A. Vasilevsky, Melissa A. Haendel, Judith A. Blake,
380 and Christopher J. Mungall. The cell ontology 2016: enhanced content, modularization, and
381 ontology interoperability. *Journal of biomedical semantics*, 7(44), 2016.
- 382 [17] Eli Eisenberg and Erez Y Levanon. Human housekeeping genes, revisited. *TRENDS in Genetics*,
383 29(10):569–574, 2013.
- 384 [18] Felix Fischer, David S Fischer, Evan Biederstedt, Alexandra-Chloé Villani, and Fabian J Theis.
385 Scaling cross-tissue single-cell annotation models. *bioRxiv*, 2023.
- 386 [19] Felix Fischer, David S. Fischer, Evan Biederstedt, Alexandra-Chloé Villani, and Fabian J. Theis.
387 Scaling cross-tissue single-cell annotation models. *bioRxiv*, 2023.
- 388 [20] Dániel Fogaras, Balázs Rácz, Károly Csalogány, and Tamás Sarlós. Towards scaling fully
389 personalized pagerank: Algorithms, lower bounds, and experiments. *Internet Mathematics*,
390 2(3):333–358, 2005.
- 391 [21] Oscar Franzén, Li-Ming Gan, and Johan LM Björkegren. Panglaodb: a web server for explo-
392 ration of mouse and human single-cell rna sequencing data. *Database*, 2019:baz046, 2019.
- 393 [22] Florian Graf, Christoph Hofer, Marc Niethammer, and Roland Kwitt. Dissecting supervised
394 contrastive learning. In *International Conference on Machine Learning*, pages 3821–3830.
395 PMLR, 2021.
- 396 [23] Aric Hagberg, Pieter Swart, and Daniel S Chult. Exploring network structure, dynamics, and
397 function using networkx. Technical report, Los Alamos National Lab.(LANL), Los Alamos,
398 NM (United States), 2008.
- 399 [24] Minsheng Hao, Jing Gong, Xin Zeng, Chiming Liu, Yucheng Guo, Xingyi Cheng, Taifeng
400 Wang, Jianzhu Ma, Le Song, and Xuegong Zhang. Large scale foundation model on single-cell
401 transcriptomics. *bioRxiv*, pages 2023–05, 2023.
- 402 [25] Minsheng Hao, Jing Gong, Xin Zeng, Chiming Liu, Yucheng Guo, Xingyi Cheng, Taifeng
403 Wang, Jianzhu Ma, Le Song, and Xuegong Zhang. Large scale foundation model on single-cell
404 transcriptomics. *bioRxiv*, 2023.
- 405 [26] Yuhan Hao, Tim Stuart, Madeline H Kowalski, Saket Choudhary, Paul Hoffman, Austin
406 Hartman, Avi Srivastava, Gesmira Molla, Shaista Madad, Carlos Fernandez-Granda, et al.
407 Dictionary learning for integrative, multimodal and scalable single-cell analysis. *Nature*
408 *biotechnology*, 42(2):293–304, 2024.
- 409 [27] Graham Heimberg, Tony Kuo, Daryle DePianto, Tobias Heigl, Nathaniel Diamant, Omar Salem,
410 Gabriele Scalia, Tommaso Biancalani, Shannon Turley, Jason Rock, et al. Scalable querying of
411 human cell atlases via a foundational model reveals commonalities across fibrosis-associated
412 macrophages. *bioRxiv*, pages 2023–07, 2023.
- 413 [28] Congxue Hu, Tengyue Li, Yingqi Xu, Xinxin Zhang, Feng Li, Jing Bai, Jing Chen, Wenqi
414 Jiang, Kaiyue Yang, Qi Ou, et al. Cellmarker 2.0: an updated database of manually curated
415 cell markers in human/mouse and web tools based on scrna-seq data. *Nucleic Acids Research*,
416 51(D1):D870–D876, 2023.
- 417 [29] Aleksandr Ianevski, Anil K. Giri, and Tero Aittokallio. Fully-automated and ultra-fast cell-type
418 identification using specific marker combinations from single-cell transcriptomic data. *Nature*
419 *Communications*, 13, 2022.
- 420 [30] Jiachen Jiang, Jinxin Zhou, Peng Wang, Qing Qu, Dustin Mixon, Chong You, and Zhihui Zhu.
421 Generalized neural collapse for a large number of classes. *ArXiv*, abs/2310.05351, 2023.

- 422 [31] Hyun Min Kang, Meena Subramaniam, Sasha Targ, Michelle Nguyen, Lenka Maliskova,
423 Elizabeth McCarthy, Eunice Wan, Simon Wong, Lauren Byrnes, Cristina M Lanata, et al. Multi-
424 plexed droplet single-cell rna-sequencing using natural genetic variation. *Nature biotechnology*,
425 36(1):89–94, 2018.
- 426 [32] Hyun Min Kang, Meena Subramaniam, Sasha Targ, Michelle Ly Thai Nguyen, Lenka Maliskova,
427 Elizabeth E. McCarthy, Eunice Wan, Simon Wong, Lauren E. Byrnes, Cristina M. Lanata,
428 Rachel E. Gate, Sara Mostafavi, Alexander Marson, Noah A. Zaitlen, Lindsey A. Criswell, and
429 Chun Jimmie Ye. Multiplexed droplet single-cell rna-sequencing using natural genetic variation.
430 *Nature biotechnology*, 36:89 – 94, 2017.
- 431 [33] Kasia Zofia Kedzierska, Lorin Crawford, Ava Pardis Amini, and Alex X Lu. Assessing the
432 limits of zero-shot foundation models in single-cell biology. *bioRxiv*, pages 2023–10, 2023.
- 433 [34] Hadas Keren-Shaul, Ephraim Kenigsberg, Diego Adhemar Jaitin, Eyal David, Franziska Paul,
434 Amos Tanay, and Ido Amit. Mars-seq2. 0: an experimental and analytical pipeline for indexed
435 sorting combined with single-cell rna sequencing. *Nature protocols*, 14(6):1841–1862, 2019.
- 436 [35] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
437 *arXiv:1412.6980*, 2014.
- 438 [36] Ilya Korsunsky, Nghia Millard, Jean Fan, Kamil Slowikowski, Fan Zhang, Kevin Wei, Yuriy
439 Baglaenko, Michael Brenner, Po-ru Loh, and Soumya Raychaudhuri. Fast, sensitive and
440 accurate integration of single-cell data with harmony. *Nature methods*, 16(12):1289–1296,
441 2019.
- 442 [37] Samuel A Lambert, Arttu Jolma, Laura F Campitelli, Pratyush K Das, Yimeng Yin, Mihai
443 Albu, Xiaoting Chen, Jussi Taipale, Timothy R Hughes, and Matthew T Weirauch. The human
444 transcription factors. *Cell*, 172(4):650–665, 2018.
- 445 [38] Phuc H Le-Khac, Graham Healy, and Alan F Smeaton. Contrastive representation learning: A
446 framework and review. *Ieee Access*, 8:193907–193934, 2020.
- 447 [39] Qiao Liu, Zhiqiang Hu, Rui Jiang, and Mu Zhou. Deepcdr: a hybrid graph convolutional
448 network for predicting cancer drug response. *bioRxiv*, 2020.
- 449 [40] Romain Lopez, Jeffrey Regier, Michael B Cole, Michael I Jordan, and Nir Yosef. Deep
450 generative modeling for single-cell transcriptomics. *Nature methods*, 15(12):1053–1058, 2018.
- 451 [41] Mohammad Lotfollahi, Mohsen Naghipourfar, Malte D. Luecken, Matin Khajavi, Maren
452 Büttner, Marco Wagenstetter, Žiga Avsec, Adam Gayoso, Nir Yosef, Marta Interlandi, Sergei
453 Rybakov, Alexander V. Misharin, and Fabian J Theis. Mapping single-cell data to reference
454 atlases by transfer learning. *Nature Biotechnology*, 40:121 – 130, 2021.
- 455 [42] Malte D. Luecken, Maren Büttner, Kridsakorn Chaichoompu, Anna Danese, Marta Inter-
456 landi, MF Mueller, D Strobl, Luke Zappia, Martin Dugas, Maria Colomé-Tatché, and F Theis.
457 Benchmarking atlas-level data integration in single-cell genomics. *Nature Methods*, 19:41 – 50,
458 2020.
- 459 [43] Malte D Luecken, Maren Büttner, Kridsakorn Chaichoompu, Anna Danese, Marta Interlandi,
460 Michaela F Müller, Daniel C Strobl, Luke Zappia, Martin Dugas, Maria Colomé-Tatché, et al.
461 Benchmarking atlas-level data integration in single-cell genomics. *Nature methods*, 19(1):41–50,
462 2022.
- 463 [44] Colin Megill, Bruce Martin, Charlotte Weaver, Sidney Bell, Lia Prins, Seve Badajoz, Brian
464 McCandless, Angela Oliveira Pisco, Marcus Kinsella, Fiona Griffin, et al. Cellxgene: a
465 performant, scalable exploration platform for high dimensional sparse matrices. *bioRxiv*, pages
466 2021–04, 2021.
- 467 [45] Gyutaek Oh, Baekgyu Choi, Inkyung Jung, and Jong Chul Ye. schyena: Foundation model for
468 full-length single-cell rna-seq analysis in brain. *arXiv preprint arXiv:2310.02713*, 2023.

- 469 [46] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,
470 Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-
471 learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830,
472 2011.
- 473 [47] Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y Fu, Tri Dao, Stephen Baccus, Yoshua
474 Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional
475 language models. In *International Conference on Machine Learning*, pages 28043–28078.
476 PMLR, 2023.
- 477 [48] Yixuan Qiu, Jiebiao Wang, Jing Lei, and Kathryn Roeder. Identification of cell-type-specific
478 marker genes from co-expression patterns in tissue samples. *Bioinformatics*, 37(19):3228–3234,
479 2021.
- 480 [49] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language
481 understanding by generative pre-training. 2018.
- 482 [50] William M Rand. Objective criteria for the evaluation of clustering methods. *Journal of the*
483 *American Statistical association*, 66(336):846–850, 1971.
- 484 [51] Yanay Rosen, Yusuf Roohani, Ayush Agrawal, Leon Samotorcan, Tabula Sapiens Consortium,
485 Stephen R Quake, and Jure Leskovec. Universal cell embeddings: A foundation model for cell
486 biology. *bioRxiv*, pages 2023–11, 2023.
- 487 [52] Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster
488 analysis. *Journal of computational and applied mathematics*, 20:53–65, 1987.
- 489 [53] Robert Salomon, Dominik Kaczorowski, Fatima Valdes-Mora, Robert E Nordon, Adrian Neild,
490 Nona Farbehi, Nenad Bartonicek, and David Gallego-Ortega. Droplet-based single cell rnaseq
491 tools: a practical guide. *Lab on a Chip*, 19(10):1706–1727, 2019.
- 492 [54] Hongru Shen, Jilei Liu, Jiani Hu, Xilin Shen, Chao Zhang, Dan Wu, Mengyao Feng, Meng
493 Yang, Yang Li, Yichen Yang, et al. Generative pretraining from large-scale transcriptomes for
494 single-cell deciphering. *Iscience*, 26(5), 2023.
- 495 [55] Christina V. Theodoris, Ling Xiao, Anant Chopra, Mark D. Chaffin, Zeina R Al Sayed,
496 Matthew C. Hill, Helene Mantineo, Elizabeth M Brydon, Zexian Zeng, X. Shirley Liu, and
497 Patrick T. Ellinor. Transfer learning enables predictions in network biology. *Nature*, 618:616–
498 624, 2023.
- 499 [56] Ander Urruticoechea, Ramon Alemany, J Balart, Alberto Villanueva, Francesc Vinals, and
500 Gabriel Capella. Recent advances in cancer therapy: an overview. *Current pharmaceutical*
501 *design*, 16(1):3–10, 2010.
- 502 [57] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
503 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
504 *processing systems*, 30, 2017.
- 505 [58] Felipe A Vieira Braga, Gozde Kar, Marijn Berg, Orestes A Carpaij, Krzysztof Polanski, Lukas M
506 Simon, Sharon Brouwer, Tomás Gomes, Laura Hesse, Jian Jiang, et al. A cellular census of
507 human lungs identifies novel cell states in health and in asthma. *Nature medicine*, 25(7):1153–
508 1163, 2019.
- 509 [59] Hui Wan, Liang Chen, and Min Deng. scemail: Universal and source-free annotation method
510 for scrna-seq data with novel cell-type perception. *Genomics, Proteomics & Bioinformatics*,
511 20:939 – 958, 2022.
- 512 [60] Hanzhi Wang, Zhewei Wei, Junhao Gan, Sibow Wang, and Zengfeng Huang. Personalized
513 pagerank to a target node, revisited. In *Proceedings of the 26th ACM SIGKDD International*
514 *Conference on Knowledge Discovery & Data Mining*, pages 657–667, 2020.
- 515 [61] Wenchuan Wang, Fan Yang, Yuejing Fang, Duyu Tang, Junzhou Huang, Hui Lu, and Jianhua
516 Yao. scbert as a large-scale pretrained deep language model for cell type annotation of single-cell
517 rna-seq data. *Nature Machine Intelligence*, 4:852 – 866, 2022.

- 518 [62] Xiliang Wang, Yao He, Qiming Zhang, Xianwen Ren, and Zemin Zhang. Direct comparative
519 analyses of 10x genomics chromium and smart-seq2. *Genomics, Proteomics and Bioinformatics*,
520 19(2):253–266, 2021.
- 521 [63] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani
522 Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large
523 language models. *arXiv preprint arXiv:2206.07682*, 2022.
- 524 [64] F Alexander Wolf, Philipp Angerer, and Fabian J Theis. Scanpy: large-scale single-cell gene
525 expression data analysis. *Genome biology*, 19:1–5, 2018.
- 526 [65] Fan Yang, Wenchuan Wang, Fang Wang, Yuan Fang, Duyu Tang, Junzhou Huang, Hui Lu, and
527 Jianhua Yao. scbert as a large-scale pretrained deep language model for cell type annotation of
528 single-cell rna-seq data. *Nature Machine Intelligence*, 4(10):852–866, 2022.
- 529 [66] Hongyu Zhao, Tianyu Liu, Kexing Li, Yuge Wang, and Hongyu Li. Evaluating the utilities of
530 large language models in single-cell data analysis. 2023.
- 531 [67] Suyuan Zhao, Jiahuan Zhang, and Zaiqing Nie. Large-scale cell representation learning via
532 divide-and-conquer contrastive learning. *arXiv preprint arXiv:2306.04371*, 2023.

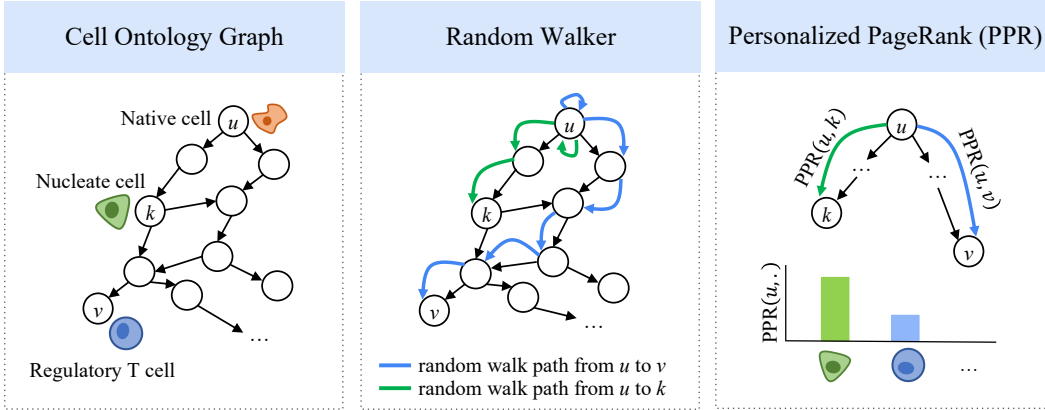


Figure 4: Graphical illustration of applying the Personalized PageRank (PPR) algorithm to cell ontology graph. As explained in App. A, PPR conducts random walks over the ontology graph with respect to a target cell type u , and converges to a steady state when the likelihood of terminating on each node stabilizes into a steady distribution. This likelihood distribution determines the final PPR score $\text{PPR}(\cdot)$ and reflects the structural similarity between cell types.

533 A PPR Transformation

534 **Personalized PageRank (PPR).** Personalized PageRank (PPR) extends the classic PageRank algo-
 535 rithm, which Google originally developed to rank web pages in search engines. PageRank conducts
 536 this by analyzing large-scale hyperlinked graphs on the web using random walker simulations. Unlike
 537 traditional PageRank that assigns a universal score to each web page, PPR customizes these scores.
 538 Specifically, individual user preferences during searches are incorporated, so that PPR can focus on
 539 web pages particularly relevant to each user. Due to its flexibility and effectiveness, PPR has been
 540 widely applied in graph learning across various fields, such as social networks, recommendation
 541 systems, and biological data analysis.

542 As illustrated in Fig. 4, this algorithm starts with a predefined preference node (or target node), which
 543 is emphasized according to the user’s interests. Subsequently, a random walk is conducted on the
 544 graph to facilitate graph traversal. At each step of the walk, there is a fixed probability α that the
 545 walker will jump back to the target node from the current node instead of moving to an adjacent
 546 node chosen at random. This process of jumping, commonly referred to as "teleportation", biases the
 547 walk towards subgraphs that are of particular importance to the target node, thus personalizing the
 548 results according to user preferences. The walk continues until it reaches a steady state, at which
 549 point the likelihood of being on each node stabilizes into a steady-state distribution. These stabilized

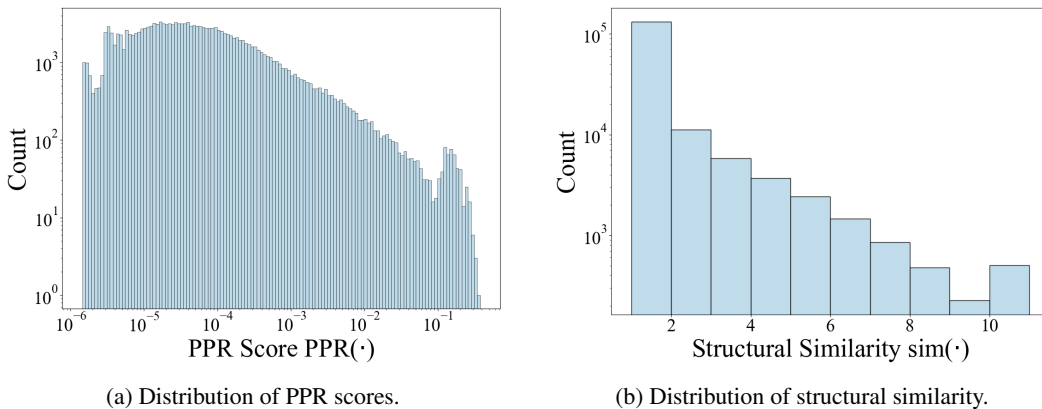


Figure 5: Comparison of the distributions for the PPR scores $\text{PPR}(\cdot)$ and the structural similarity $\text{sim}(\cdot)$ after the transformation.

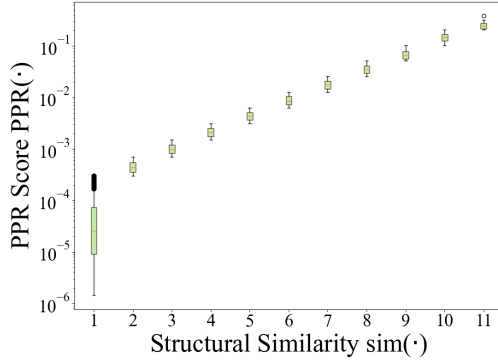


Figure 6: Relationships between the structural similarity $\text{sim}(\cdot)$ after PPR transformation and the original PPR scores $\text{PPR}(\cdot)$.

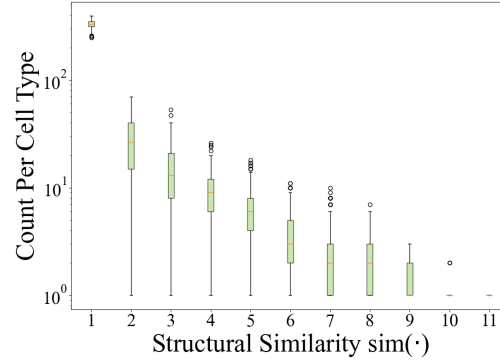


Figure 7: Frequency for each target cell type to be associated with other cell types that is at specific levels of structural similarity.

550 probabilities, reflecting both the graph’s structure and the user’s preferences, determine the PPR
 551 scores. These scores effectively evaluate each node’s structural similarities and rank them according
 552 to their relevance and importance from a personalized perspective.

553 **PPR transformation.** In scCello, the PPR algorithm is applied to the cell ontology graph to
 554 assess the structural similarities among cell types, or to measure their importance relative to a
 555 specified target cell type. We implemented PPR using the "pagerank" function in NetworkX [23]
 556 with "personalization" as arguments.

557 However, modification is needed to integrate PPR into TFM pre-training. The PPR scores are in
 558 real-number format and susceptible to numerical noise. Also, as shown in Fig. 5a, these scores
 559 typically exhibit a skewed distribution, concentrated around lower magnitudes. Consequently, setting
 560 precise thresholds to differentiate between node similarity and dissimilarity is challenging. Moreover,
 561 the vast amount small PPR values may be indistinguishable from noise.

562 To mitigate the effects of numerical noise and skewed magnitudes for the PPR scores, we employ
 563 truncation, logarithmic scaling, and discretization as outlined in Eqn. 3. Note that Eqn. 3 defines a
 564 monotonic, non-decreasing function that preserves the relative order between nodes. Its minimum
 565 value is set to 1 for the least similar cell types.

566 This equation transforms the raw PPR score, $\text{PPR}(\cdot)$, into the final structural similarity, $\text{sim}(\cdot)$. This
 567 transformation ensures that $\text{sim}(\cdot)$ accurately reflects pronounced similarities as defined by the cell
 568 ontology and avoids emphasizing minor dissimilarities that could mislead during TFM pre-training.

569 **Analyses.** In Fig. 5, we present a comparison of the distributions for the PPR score, $\text{PPR}(\cdot)$,
 570 and the transformed structural similarity, $\text{sim}(\cdot)$. After transformation, the distribution of $\text{sim}(\cdot)$ is
 571 less skewed and exhibits clear discretization. This facilitates the setting of definitive thresholds for
 572 distinguishing between similarity and dissimilarity among cell types, thereby enabling the effective
 573 incorporation of the cell ontology graph in scCello’s pre-training.

574 In addition, we provide detailed insights into the scale of structural similarity, the distribution of these
 575 similarities for each cell type, and examples of cell types associated with various levels of structural
 576 similarity:

577 (1) Fig. 6 illustrates the correspondence between the structural similarity after PPR transfor-
 578 mation and the original PPR scores, showcasing a log-linear relationship as expected. This
 579 helps clarify the scaling of structural similarity, which is discretized into integer levels
 580 ranging from 1 to 11.

581 (2) Fig. 7 demonstrates how frequently each target cell type is associated with other cell types
 582 at specific levels of structural similarity. Consequently, during scCello’s pre-training, a
 583 substantial number of negative samples are expected to be utilized in the inter-cellular
 584 relational alignment objective, as outlined in Sec. 2.4.

Table 7: Examples of cell types associated with various levels of structural similarity, $\text{sim}(\cdot)$, for specified target cell types. Cell types demonstrated in the cell ontology graph in Fig. 1 are underlined.

Target Type	$\text{sim}(\cdot)$	Corresponding Cell Types
T Cell	8	"gamma-delta T cell", " <u>mature T cell</u> ", "lymphocyte"
	7	"mature gamma-delta T cell", " α - β T cell", " <u>mature α-β T cell</u> ", "thymocyte"
	6	"B cell", "double-positive, α - β thymocyte", "CD8-positive, α - β thymocyte", "CD4-positive, α - β T cell", " <u>CD8-positive, α-β T cell</u> ", "double negative thymocyte"
	5	"dendritic cell", "innate lymphoid cell", "plasmablast", "mononuclear cell", " <u>regulatory T cell</u> ", "memory T cell", "myeloid leukocyte", "naive T cell", "mature B cell", "CD4-positive, CD25-positive, α - β regulatory T cell"

	1	"renal intercalated cell", "smooth muscle fiber of ileum", "type II pneumocyte", "hematopoietic cell", " <u>neuron</u> ", "common lymphoid progenitor", ...
Neuron	7	"secretory cell"
	6	"glutamatergic neuron", "GABAergic neuron", "motor neuron", "neural cell", "peripheral nervous system neuron"
	5	"glycinergic neuron", "retinal bipolar neuron", " <u>native cell</u> ", "enteric neuron", "retina horizontal cell", "amacrine cell", "neuronal receptor cell"
	4	"retinal ganglion cell", "endocrine cell", "neuroendocrine cell", "cerebral cortex GABAergic interneuron", "muscle cell", "somatic cell"

	1	"germ cell", " <u>T cell</u> ", "tracheal goblet cell", "DN3 thymocyte", "promonocyte", "cerebral cortex endothelial cell", ...

585 (3) Tab. 7 displays examples of highly similar and dissimilar cell types categorized into various
586 levels of structural similarity, specifically targeting "T cell" and "neuron" types.

587 B Data Preprocessing Details

588 **Download and Preprocessing.** We downloaded from CellxGene [1] census version 2023-7-25. We
589 focused on 291 datasets for human scRNA-seq. We preprocessed the dataset by the following steps:

- 590 (1) **Remove non-primary cells.** Some data on CellxGene was duplicated due to multiple
591 submissions of the same dataset from different research groups, therefore cells marked as
592 "non-primary" were filtered out to prevent label leakage between pre-training and down-
593 stream.
- 594 (2) **Filter out cells not produced by 10x-based [62] sequencing protocols.** There are numerous
595 sequencing protocols in CellxGene database besides 10x-based sequencing [62], such as
596 Drop-seq [53] and MARS-seq [34]. Only sequencing data from 10x-based sequencing
597 protocols was kept to avoid large variation of data signals [42].
- 598 (3) **Exclude cancer cells.** Cancer cells were highly dissimilar to normal cells and even occupied
599 a large amount in the CellxGene database (nearly 12%). These cells could bring unexpected
600 signals and skew the data, therefore we excluded these cancer cells.

601 To build downstream datasets for out-of-distribution (OOD) generalization evaluation, we first held
602 out two category sets for each of the three settings: unseen cell types, unseen tissues and unseen

Table 8: Data statistics for our curated pre-training and downstream datasets, where the downstream datasets encompass one ID dataset and six OOD datasets under three different OOD scenarios, including unseen cell types, unseen tissues and unseen donors (Sec. 4.1). The blue colored numbers represent disjoint categories of that column. For example, in the "cell type" column, the cell type set in the pre-training data, and the cell type set in the OOD cell type dataset D_1^{ct} and D_2^{ct} are disjoint.

Dataset	#Total Cells	#Cell Types	#Tissues	#Donors	#Conditions	#Batches
Pre-training data	22,293,755	398	140	4,103	55	267
ID dataset D^{id}	22,317	318	132	3,447	54	261
OOD cell type dataset D_1^{ct}	486,810	87	125	122	35	90
OOD cell type dataset D_2^{ct}	435,791	87	128	106	40	117
OOD tissue dataset D_1^{ts}	335,675	186	32	1,801	10	28
OOD tissue dataset D_2^{ts}	341,681	205	32	2,052	7	25
OOD donor dataset D_1^{dn}	2,528,134	439	91	525	36	127
OOD donor dataset D_2^{dn}	2,521,868	404	101	525	33	123
In Total	/	572	204	5153	/	/

603 donors. Each category set were randomly selected with selection ratios 15%, 15% and 10% for
604 the three OOD settings respectively. During the selection, we prohibited any category associated
605 with more than 0.1% of the total pre-processed cells from being selected. This avoids losing too
606 much data for pre-training. After the selection, cells associated with each held category set are
607 collected, resulting in two OOD downstream datasets for each of the three OOD settings. These
608 datasets are denoted as $\{D_i^{ct}\}_{i=1}^2$ for the OOD cell type setting, $\{D_i^{ts}\}_{i=1}^2$ for the OOD tissue setting,
609 and $\{D_i^{dn}\}_{i=1}^2$ for the OOD donor setting.

610 By excluding cells with at least one property belong to any of the six held category sets, the remaining
611 data is further split into 99.9% as our pre-training data and 0.1% as the in-distribution (ID) downstream
612 dataset D^{id} . This way, our pre-training data and the ID dataset D^{id} share similar data distributions.

613 **Data Statistics.** We summarize the data statistics for our curated pre-training dataset, one ID dataset
614 and six OOD datasets in Tab. 8.

615 C Implementation Details

616 **scRNA-seq Data.** scRNA-seq can enable the quantification of gene expression profiles of
617 individual cells. Each cell’s gene expression profile can be described by the set $\hat{X} =$
618 $\{(e_1, g_1), (e_2, g_2), \dots, (e_M, g_M)\}$, where e_k denotes the expression count of gene g_k , with $e_k \geq 0$.
619 A value of $e_k = 0$ indicates that the gene g_k is not expressed or not detected by the sequencing
620 experiment. We use the same gene vocabulary set as [55], with the number of genes $M=25,424$.

621 **Gene Token Vocabulary.** The gene vocabulary set contains both protein-coding genes and miRNA
622 genes. M , the number of genes, is not the same as the number of all tokens in the model vocabulary.
623 scCello has M gene tokens plus three more special tokens [MASK] for masking, [CLS] for the start of
624 a sentence and [PAD] for padding.

625 **Rank Value Encoding.** Unlike natural languages, which inherently follow a sequential order,
626 scRNA-seq data presents a unique challenge due to the lack of intrinsic order among gene tokens.
627 Therefore, we employ Rank Value Encoding [55] approach to rank genes based on their normalized
628 expression set $\{(\tilde{e}_i, g_i)\}_{i=1}^M$. Specifically, gene expressions are first normalized by the total count
629 within a cell [64] in a cell-wise manner, and then normalized through gene-specific weighting factors
630 in a gene-wise manner. These factors are adopted from [55], which calculates the non-zero median
631 value of expression of each detected gene across all cells. By design, these factors are assigned to
632 emphasize lowly-expressed but essential genes, such as transcription factors [37], while deprioritizing
633 ubiquitously expressed housekeeping genes [17].

634 After the normalization and ranking, it results in an ordered sequence of gene identities $X =$
635 $[g_{\pi(1)}, g_{\pi(2)}, \dots, g_{\pi(M)}]$ with an index permutation $\pi(\cdot)$, satisfying $\tilde{e}_{\pi(1)} \geq \tilde{e}_{\pi(2)} \geq \dots \geq \tilde{e}_{\pi(M)}$.
636 To mitigate memory consumption, zero-expressed genes are removed and the gene sequence is

Table 9: Hyper-parameters comparison between TFM baselines (introduced in Sec. 4.1) and our TFM scCello. "The number of" is denoted with the symbol #.

Configuration	Geneformer [55]	scGPT [14]	scTab [19]	UCE [51]	scCello
#Parameters	10,316,196	51,330,049	9,655,628	674,745,857	10,683,654
Total GPUs	12 * V100 (32G)	4 * A100	1 * A100	24 * A100 (80G)	4 * A100 (40G)
Training Time	3 days	3 days	/	43.5 days	2 days
Sequence Length	2,048	1,200	19,331	1,024 [N]	2,048
Gene Mask Ratio	15%	/	/	20%	15%
Batch Size Per GPU	12	32	2,048	6	12
Gradient Accumulation Steps	1	1	1	4	4
Effective Batch Size	144	128	2048	576	192
Cell Reprs.	Avg. pooling	CLS	/	CLS	CLS
#Genes in Token Vocabulary	25,424	48,292	19,331	Any protein-coding genes	25,424
#Transformer Layers	6	12	/	33	6
Transformer Layer Hidden Dimension	512	512	/	5,120	512
Transformer Layer Embedding Size	256	512	/	1,280	256
#Transformer Heads	4	8	/	20	4
Transformer Layer Activation Function	GeLU	ReLU	/	ReLU	ReLU
MLP Layer Activation Function	ReLU	ReLU	/	GeLU	ReLU
Dropout	0.02	0.2	/	0.05	0.02

637 further truncated with a context length $L=2,048$ in practice. This rank-based approach offers better
638 robustness against technical artifacts than directly using the original numerical expressions, which
639 can vary significantly in magnitude across different experimental assays [42].

640 **Cell and Cell Type Representations.** Given a pre-training dataset with N cells $\mathcal{X} =$
641 $\{X_1, X_2, \dots, X_N\}$, each cell X_i can be mapped to a specific cell type ontology identifier $c_i \in \mathcal{V}$.
642 For analyzing, scCello denotes cell X_i 's representation as z_i and cell type c_i 's representation as h_{c_i} .

643 **Masked Gene Prediction.** Given a batch of cells $\{X_i\}_{i=1}^B$, scCello predicts a gene token g_k based
644 on the ordered gene sequence context $X_{i,\setminus k}=[g_1, \dots, g_{k-1}, [\text{MASK}], g_{k+1}, \dots, g_M]$ after replacing
645 the token with a special [MASK]. This objective (term as \mathcal{L}_{MGP}) aims to capture complex but
646 important gene-gene interactions within one cell, like regulatory mechanisms between transcription
647 factors and other genes:

$$\mathcal{L}_{\text{MGP}} = - \sum_{k=1}^B \mathbb{E}_{i \sim \Psi} - \log p(x_i | X_{k,\setminus i}) \quad (6)$$

648 where tokens are masked by a pre-defined distribution Ψ , same as that in BERT [15]. Specifically,
649 80% selected genes are replaced with [MASK], 10% selected genes are kept the same as its original,
650 and 10% selected genes are replaced with random gene tokens.

651 **Model Architecture.** scCello utilizes a stack of self-attention transformer encoder layers [57],
652 each composed of a self-attention and feedforward neural networks. The self-attention mechanism
653 processes the input sequence, effectively capturing interactions between gene tokens.

654 **Configuration Hyper-parameters.** Besides scCello, we also summarize essential hyper-parameters
655 for TFM baselines in Tab. 9 for comparison. It includes pre-training configurations like batch size,

Table 10: Metrics used in downstream tasks.

Task	Metrics
Cell Type Clustering (Sec. 4.2.1)	NMI, ARI, ASW, AvgBio
Cell Type Classification (Sec. 4.2.2)	Acc, Macro F1, AvgBio, Δ_{AvgBio}
Novel Cell Type Classification (Sec. 4.3)	Acc, Macro F1
Marker Gene Prediction (Sec. 4.4)	AUROC
Cancer Drug Response Prediction (Sec. 4.5)	PCC
Batch Integration (Sec. 4.6)	NMI, ARI, ASW, AvgBio, ASW _b , GraphConn, AvgBatch, Overall

656 sequence length, and training time consumed. It also includes architecture configurations for the
657 transformer model backbone, such as the number of transformer layers and the embedding size of
658 transformer layers. Note that scTab uses TabNet [2] instead of transformer layers as model backbone,
659 therefore its architecture configurations are not recorded in the table.

660 D Downstream Experiment Details

661 D.1 Evaluation Metrics

662 All metrics used in downstream tasks are summarized in Tab. 10 and introduced below.

663 **Normalized Mutual Info Score (NMI).** The NMI is a metric that quantifies the similarity between
664 two different clustering assignments or labelings of the same set of samples. We use NMI to
665 compare the cell-type labels, with the cluster indices obtained from applying the Louvain clustering
666 algorithm [12] on the target dataset.

667 We denote the two label assignments of the same N cell samples as C and K , representing the cell-type
668 labels and the Louvain cluster indices, respectively. The entropy of a label assignment, say C , is a
669 measure of the uncertainty associated with that assignment set. It’s calculated as:

$$H(C) = - \sum_{i=1}^{|C|} P(i) \log P(i) \quad (7)$$

670 where $|C|$ is the number of unique cell types and $P(i) = \frac{|C_i|}{N}$ is the probability that a randomly
671 selected sample belongs to the class C_i . The entropy $H(K)$ for the cluster indices K is computed
672 similarly, with $Q(j) = \frac{|K_j|}{N}$ being the probability of a sample belonging to the cluster K_j :

$$H(K) = - \sum_{j=1}^{|K|} Q(j) \log Q(j) \quad (8)$$

673 The mutual information (MI) between C and K quantifies the amount of information shared between
674 the two label assignments. It is calculated by:

$$\text{MI}(C, K) = \sum_i^{|C|} \sum_j^{|K|} R(i, j) \log \frac{R(i, j)}{P(i)Q(j)} \quad (9)$$

675 where $R(i, j) = \frac{|C_i \cap K_j|}{N}$ is the probability that a randomly selected sample belongs to both the class
676 C_i and the cluster K_j .

677 The normalized mutual information (NMI) is defined as:

$$\text{NMI}(C, K) = \frac{\text{MI}(C, K)}{\text{mean}(H(C), H(K))} \quad (10)$$

678 NMI is a normalized version of MI, scaled by the mean of the entropy terms for cell-type labels and
 679 cluster indices. This normalization ensures that NMI values range from 0 to 1, where 0 indicates no
 680 correlation between the two label assignments, and 1 represents a perfect match.

681 To obtain the best match between the clusters and the cell-type labels, we performed optimized
 682 Louvain clustering over a range of resolutions from 0.1 to 2, in steps of 0.1. The clustering output
 683 with the highest NMI score, when compared to the cell-type label set, was selected as the optimal
 684 clustering result. The implementation of NMI used in this study was from the scib python library [43].

685 **Adjusted Rand Index Score (ARI).** The ARI is another metric used to evaluate the similarity
 686 between the clustering assignment and the cell type labels of the same set of samples, similar to
 687 the NMI metric. In this context, we similarly denote the cell-type labels as C and the Louvain [12]
 688 cluster indices computed on the target dataset as K .

689 The Rand Index (RI) is a measure of the overlap between the two clusterings, C and K . It considers
 690 both the correct clustering overlaps and the correct disagreements between the two clusterings [50].
 691 Formally, if we define a as the number of pairs of elements that belong to the same set in both C and
 692 K , and b as the number of pairs of elements that are in different sets in C and in different sets in K ,
 693 the unadjusted RI is given by:

$$\text{RI} = \frac{a + b}{C_2^N} \quad (11)$$

694 where N is the total number of cell samples and C_2^N represents the total number of possible pairs in
 695 the dataset.

696 However, the unadjusted RI does not account for the possibility of random label assignments leading
 697 to correct overlaps by chance. To address this issue, the adjusted RI (ARI) is introduced, which
 698 corrects for randomly correct labels by discounting the expected RI of random labelings:

$$\text{ARI} = \frac{\text{RI} - \mathbb{E}[\text{RI}]}{\max(\text{RI}) - \mathbb{E}[\text{RI}]} \quad (12)$$

699 The ARI ranges from 0 to 1, where 0 corresponds to a random labeling, and 1 indicates a perfect
 700 match between the two clustering assignments.

701 Similar to NMI, we performed NMI-optimized Louvain clustering to obtain the best match between
 702 the clusters and the cell-type labels. Specifically, we executed Louvain clustering over a range of
 703 resolutions and selected the clustering output with the highest NMI score when compared to the cell
 704 type label set. The implementation of ARI used in this study was from the scib python library [43].

705 **Average Silhouette Width Score (ASW).** The silhouette width [52] is a metric that evaluates
 706 the quality of a clustering solution by quantifying the relationship between the within-clustering
 707 distances and the between-cluster distances for each data point. Like the NMI and the ARI, the
 708 silhouette calculates the similarity between the clustering assignment and the cell type labels of the
 709 same set of samples.

710 For each cell sample, the silhouette width is computed based on two scores: (1) a : the mean distance
 711 between a sample and all other samples in the same cluster; and (2) b the mean distance between a
 712 sample and all samples in the nearest neighboring cluster. The silhouette score s_i for each sample i is
 713 defined as

$$s_i = \frac{b - a}{\max(a, b)} \quad (13)$$

714 The silhouette score ranges from -1 to 1, with higher values indicating that the sample is well-matched
 715 to its own cluster and dissimilar to the nearest neighboring cluster.

716 To obtain an overall assessment of the clustering quality, the average silhouette width (ASW) is
 717 calculated by averaging the silhouette scores s_i across all samples. This overall ASW, denoted as
 718 ASW_o , ranges between -1 and 1, with the following interpretations:

- 719 • ASW_o close to 1: The clusters are dense and well-separated.

- 720 • ASW_o around 0: The clusters overlap, and the between-cluster and within-cluster variability
721 are approximately equal.
- 722 • ASW_o near -1: Strong misclassification has occurred, where the within-cluster variability is
723 greater than the between-cluster variability.

724 To ensure that the final ASW metric falls within the range of 0 to 1, a scaling operation is often
725 applied:

$$ASW = \frac{ASW_o + 1}{2} \quad (14)$$

726 This scaled ASW value, ranging from 0 to 1, provides a convenient measure for evaluating the quality
727 of the clustering solution, with higher values indicating better separation and cohesion of the clusters.

728 **AvgBio.** This score combines the three clustering metrics: NMI, ARI and ASW.

$$AvgBio = \frac{1}{3}(NMI + ARI + ASW) \quad (15)$$

729 **Silhouette Variant Score (ASW_b).** To evaluate the effectiveness of the batch integration task
730 (Sec. 4.6), a variant of the average silhouette width score (ASW) is employed, referred to as the
731 ASW_b . Unlike ASW based on cell type labels, ASW_b considers batch labels. This score is designed
732 to assess the degree of batch mixing, where a score of 0 indicates well-mixed batches, and deviations
733 from 0 suggest the presence of a batch effect.

734 We take the absolute value of the original silhouette width score \tilde{s}_i for sample i based on batch labels:

$$s'_i = |\tilde{s}_i| \quad (16)$$

735 To ensure higher scores indicate better batch mixing, these scores are scaled by subtracting them
736 from 1. As we expect batches to integrate within cell identity clusters, we compute the $ASW_{b,j}$ score
737 for each cell label j separately, using the following equation:

$$ASW_{b,j} = \frac{1}{|C_j|} \sum_{i \in C_j} 1 - s(i)' \quad (17)$$

738 where $C_j = \{i | c_i = j\}_{i=1}^N$ is the set of cell indices whose cell type label is exactly j .

739 To obtain the final ASW_b score, the label-specific $ASW_{b,j}$ scores are averaged across the set of
740 unique cell type labels:

$$ASW_b = \frac{1}{|\mathcal{V}|} \sum_{j \in \mathcal{V}} ASW_{b,j} \quad (18)$$

741 where \mathcal{V} represents the set of unique cell type labels.

742 **Graph Connectivity (GraphConn).** The GraphConn metric is designed to assess whether the
743 k -nearest neighbor (k NN) graph representation of the integrated data directly connects all cells
744 with the same cell type label. This metric operates on the k NN graph, denoted as G_{kNN} , which is
745 pre-processed by the Scanpy library using the "scanpy.pp.neighbors" function.

746 For each cell type label $v \in \mathcal{V}$, where \mathcal{V} represents the set of cell type labels (Sec. 2), a subset k NN
747 graph $G_{kNN}(\mathcal{V}_v; \mathcal{E}_v)$ is created. This subset graph contains only cells from the given label v .

748 Using these subset k NN graphs, the GraphConn score is computed as follows:

$$GraphConn = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \frac{|LCC(G_{kNN}(\mathcal{V}_v, \mathcal{E}_v))|}{|\mathcal{V}_v|} \quad (19)$$

749 Here, $|LCC(\cdot)|$ is the number of nodes in the largest connected component of the graph and $|\mathcal{V}_v|$ is
750 the number of nodes with cell type v .

751 The resultant GraphConn score has a range of $(0; 1]$, where a score of 1 indicates that all cells with
752 the same cell type are connected in the integrated k NN graph. The lowest possible score indicates a
753 graph where no cell is connected to any other cell.

754 It's important to note that the GraphConn score is computed directly on the k NN graph representation
755 of the integrated data. As a result, this metric can be used to evaluate the quality of any integration
756 output, regardless of the specific integration method used.

757 **AvgBatch.** This score combines two metrics: ASW_b and GraphConn.

$$\text{AvgBatch} = \frac{1}{2}(ASW_b + \text{GraphConn}) \quad (20)$$

758 **Overall.** We follow scGPT [14] to calculate a weighted average score of both the batch removal
759 score ASW_b and the bio-conservation score AvgBio to balance biological relevance and batch
760 consistency, following the equation:

$$\text{Overall} = 0.6 * \text{AvgBio} + 0.4 * \text{AvgBatch} \quad (21)$$

761 **Accuracy (Acc).** In classification tasks like cell type classification (Sec. 4.2.2) and novel cell type
762 classification (Sec. 4.3), we denote the predicted values of the i -th sample as \hat{y}_i and the corresponding
763 true label as y_i . Then the accuracy metric is defined as

$$\text{Acc}(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}[\hat{y}_i == y_i] \quad (22)$$

764 where the $\mathbb{1}[\cdot]$ is the indicator function.

765 **Macro F1 Score (Macro F1).** The F1 Score is essentially defined for binary classification tasks.

$$F_1 = \frac{2}{\text{Recall}^{-1} + \text{Precision}^{-1}} \quad (23)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (24)$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (25)$$

766 where TP is the number of true positives, FN the number of false negatives, and FP the number of
767 false positives. The recall is intuitively the ability of the classifier to find all the positive samples; The
768 precision is intuitively the ability of the classifier not to label as positive a sample that is negative.
769 For multi-class classification, macro F1 is defined as the average F1 taken over all different classes.

770 **ROC AUC Score (AUROC).** The Area Under the Receiver Operating Characteristic (AUROC)
771 curve is a metric commonly used to evaluate the performance of binary classification models. It
772 provides a comprehensive measure of the trade-off between the true positive rate (sensitivity) and the
773 false positive rate (1 - specificity) across different classification thresholds.

774 In a binary classification task, the model's output is typically a probability or score that represents
775 the likelihood of a sample belonging to the positive class. By varying the classification threshold,
776 different operating points on the ROC curve can be obtained, where each point represents a specific
777 combination of true positive rate (TPR) and false positive rate (FPR).

778 The ROC curve is created by plotting the TPR (y-axis) against the FPR (x-axis) for different classifi-
779 cation thresholds. The AUROC is then calculated as the area under this ROC curve, providing a single
780 scalar value that summarizes the overall performance of the binary classifier. The AUROC ranges
781 from 0 to 1, with the following interpretations: (1) AUROC=1 indicates perfect classification, where
782 the classifier can perfectly distinguish between the positive and negative classes; (2) AUROC=0.5
783 indicates random guessing, indicating that the classifier performs no better than a random prediction.

784 The AUROC is a widely used metric because it provides a comprehensive evaluation of the classifier's
785 performance across all possible classification thresholds. It is invariant to class imbalance and does
786 not require choosing a specific threshold, making it a robust and threshold-agnostic measure.

787 Furthermore, the AUROC has a statistical interpretation as the probability that a randomly chosen
788 positive instance will have a higher predicted probability than a randomly chosen negative instance,
789 which provides a clear interpretation of the metric's value.

Pearson correlation coefficient score (PCC). The PCC is a widely used measure of the linear
relationship between two variables. It quantifies the strength and direction of the linear association

between the variables, ranging from -1 to 1. The formula for the PCC between two variables, A and B, is given by:

$$r_{AB} = \frac{\sum_{i=1}^n (A_i - \bar{A})(B_i - \bar{B})}{\sqrt{\sum_{i=1}^n (A_i - \bar{A})^2} \sqrt{\sum_{i=1}^n (B_i - \bar{B})^2}}$$

790 where A_i and B_i are the individual observations of variables A and B, respectively. \bar{A} and \bar{B} are the
791 sample means of A and B, respectively. n is the number of observations.

792 The numerator represents the covariance between A and B, which measures how much A and B vary
793 together from their respective means. The denominator normalizes the covariance by the product of
794 the standard deviations of A and B, ensuring that the correlation coefficient falls within the range of -1
795 to 1. The interpretation of this PPC metric is as follows: (1) $r_{AB}=1$ indicates perfect positive linear
796 correlation (as A increases, B increases proportionally); (2) $r_{AB}=-1$ indicates perfect negative linear
797 correlation (as A increases, B decreases proportionally); (3) $r_{AB}=0$ indicates no linear correlation
798 between A and B; (4) $0 < |r_{AB}| < 1$ indicates that the strength of the linear correlation increases as
799 the value approaches 1 (either positive or negative).

800 In the context of regression analysis, computing the PCC between each regressor (independent
801 variable) and the target variable can provide insights into the linear relationships between the
802 predictors and the response variable.

803 D.2 Cell Type Identification

804 D.2.1 Zero-shot Identification (*i.e.*, Cell Type Clustering)

805 **Method.** We here discuss the experimental details for Sec. 4.2.1. Cell representations extracted
806 from each baseline model are used to compute the k nearest neighbor (k NN) graph using Scanpy’s
807 standard protocols [64]. These representations and the k NN graph are then processed with Louvain
808 clustering algorithms at various resolutions, ranging from 0.1 to 2 in steps of 0.1. The optimized
809 clustering result is determined by the highest gained NMI score achieved across all the resolutions.

810 For implementation, we accelerated Louvain clustering by adopting RAPIDS, a software library that
811 enhances data science pipelines by entirely utilizing NVIDIA GPUs instead of traditional CPUs.
812 Additionally, we conducted ten iterations of dataset down-sampling and reported the averaged NMI,
813 ARI, ASW, and AvgBio scores. This approach significantly reduced the time required to evaluate a
814 dataset, such as D^{id} , from days to just a few minutes.

815 **Datasets.** As introduced in Sec. 4.2.1, we evaluate one ID dataset (D^{id}) and six OOD datasets
816 (D_i^{cond} with $cond \in \{ct, ts, dn\}$ and $i \in \{1, 2\}$) to demonstrate our model’s generalization capabili-
817 ties. These evaluations address various scenarios involving unseen cells for comprehensive testing,
818 including cells with distributions similar to our pre-training dataset, as well as those associated with
819 unseen cell types, tissues, and donors.

820 **Hyper-parameters.** We used $k = 15$ neighbors to compute the k NN graph, with node distances
821 calculated using the euclidean distance between cell representations. The Louvain clustering used
822 seed 0 as the random state and treated the k NN graph as unweighted and directed.

823 **Performance.** In Sec. 4.2.1, Tab. 1 reports only the AvgBio metric for six OOD datasets due to
824 space constraints. Full metrics, including NMI, ARI, and ASW, are detailed in: (1) Tab. 11 for the
825 two OOD cell type datasets (D_1^{ct} and D_2^{ct}); (2) Tab. 12 for the two OOD tissue datasets (D_1^{ts} and
826 D_2^{ts}); and (3) Tab. 13 for the two OOD donor datasets (D_1^{dn} and D_2^{dn}).

827 D.2.2 Identification with Fine-tuning (*i.e.*, Cell Type Classification)

828 **Method.** In this setting, the TFMs are further fine-tuned by adding a simple linear layer atop
829 their model backbones, which transforms the hidden representations into prediction logits. The
830 dimensions of these logits correspond to the number of cell type classes predicted. Importantly,
831 all model parameters, including those of the TFM backbone and the newly added linear layer, are
832 trainable during fine-tuning. The model checkpoint that achieves the highest Macro F1 score on the
833 validation data is then selected for final testing.

Table 11: Full results for the OOD unseen cell type datasets D_1^{ct} and D_2^{ct} in the cell type clustering.

Method	OOD CellType Data (D_1^{ct})				OOD CellType Data (D_2^{ct})			
	NMI↑	ARI↑	ASW↑	AvgBio↑	NMI↑	ARI↑	ASW↑	AvgBio↑
Non-TFM Methods								
Raw Data	0.864	0.718	0.529	0.703	0.823	0.557	0.505	0.629
Seurat	<u>0.893</u>	0.773	0.590	0.752	0.884	0.723	0.605	0.737
Harmony	0.553	0.241	0.432	0.432	0.594	0.248	0.411	0.417
scVI	0.905	<u>0.797</u>	0.577	<u>0.760</u>	0.889	0.709	0.577	0.725
Ontology-Agnostic TFMs								
Geneformer	0.846	0.697	0.525	0.689	0.846	0.629	0.530	0.668
scGPT	0.866	0.705	0.551	0.707	0.873	0.724	0.564	0.720
scTab	0.886	0.807	0.584	0.759	0.867	0.754	0.557	0.726
UCE	0.902	0.802	0.612	0.772	0.892	0.695	0.635	0.741
MGP	0.860	0.710	0.573	0.714	0.881	0.745	0.595	0.740
Sup	0.892	0.787	<u>0.621</u>	0.767	0.910	<u>0.793</u>	<u>0.622</u>	<u>0.775</u>
MGP+Sup	0.888	0.775	0.611	0.758	<u>0.901</u>	0.779	0.611	0.764
Ontology-Enhanced TFMs								
scCello	0.887	0.781	0.640	0.769	0.909	0.817	0.632	0.786

Table 12: Full results for the OOD unseen tissue datasets D_1^{ts} and D_2^{ts} in the cell type clustering.

Method	OOD Tissue Data (D_1^{ts})				OOD Tissue Data (D_2^{ts})			
	NMI↑	ARI↑	ASW↑	AvgBio↑	NMI↑	ARI↑	ASW↑	AvgBio↑
Non-TFM Methods								
Raw Data	0.733	0.405	0.481	0.540	0.800	0.585	0.508	0.631
Seurat	<u>0.777</u>	0.497	0.488	0.587	0.813	0.560	0.535	0.636
Harmony	0.649	0.302	0.436	0.462	0.684	0.400	0.460	0.515
scVI	0.774	0.443	0.516	0.577	0.816	0.550	0.537	0.634
Ontology-Agnostic TFMs								
Geneformer	0.736	0.412	0.468	0.539	0.787	0.499	0.505	0.597
scGPT	0.739	0.407	0.486	0.544	0.794	0.556	0.531	0.627
scTab	0.754	0.492	0.515	0.515	0.815	0.616	0.541	0.657
UCE	0.787	0.476	0.531	0.598	0.836	0.610	0.562	0.670
MGP	0.766	0.472	0.491	0.576	0.802	0.544	0.537	0.628
Sup	0.788	<u>0.502</u>	<u>0.527</u>	<u>0.605</u>	0.838	<u>0.621</u>	<u>0.580</u>	<u>0.680</u>
MGP+Sup	0.789	0.518	0.524	0.610	<u>0.833</u>	0.612	0.573	0.672
Ontology-Enhanced TFMs								
scCello	0.784	0.519	0.534	0.612	0.839	0.675	0.601	0.705

834 **Datasets.** We fine-tuned TFMs on a subset of our curated pre-training data, randomly selecting
835 90% for training and using the remaining 10% for validation. The final performance was tested on
836 the ID dataset D^{id} , which consists of cell samples never seen during scCello’s pre-training. We
837 explored two subset sizes, 0.1% and 1% of the pre-training data, to simulate scenarios where $10\times$
838 more annotated data becomes available. This exploration is meaningful for real-world applications,
839 where annotating data is both costly and time-consuming.

840 **Hyper-parameters.** For scCello, we set the following hyper-parameters for fine-tuning: a learning
841 rate of 5.0×10^{-5} , a linear learning rate scheduler with 500 warmup steps, a weight decay of 0.001,
842 and a batch size of 24. The same fine-tuning configuration was applied to the three ablation TFMs

Table 13: Full results for the OOD unseen donor datasets D_1^{dn} and D_2^{dn} in the Cell Type Clustering. Note that scTab is OOM on these two datasets.

Method	OOD Donor Data (D_1^{dn})				OOD Donor Data (D_2^{dn})			
	NMI↑	ARI↑	ASW↑	AvgBio↑	NMI↑	ARI↑	ASW↑	AvgBio↑
Non-TFM Methods								
Raw Data	0.665	0.247	0.462	0.458	0.665	0.251	0.462	0.460
Seurat	0.691	0.294	0.413	0.466	0.711	0.335	0.420	0.489
Harmony	0.679	0.286	0.405	0.456	0.690	0.324	0.408	0.474
scVI	0.699	0.269	0.466	0.478	0.722	0.311	0.471	0.502
Ontology-Agnostic TFMs								
Geneformer	0.666	0.303	0.434	0.468	0.686	0.327	0.433	0.482
scGPT	0.656	0.259	0.452	0.456	0.677	0.298	0.456	0.477
scTab	/	/	/	OOM	/	/	/	OOM
UCE	0.718	0.245	0.491	0.485	0.737	0.284	0.496	0.506
MGP	0.713	0.294	0.457	0.488	0.734	0.359	0.462	0.518
Sup	<u>0.754</u>	<u>0.357</u>	<u>0.545</u>	<u>0.552</u>	<u>0.768</u>	<u>0.395</u>	<u>0.556</u>	<u>0.573</u>
MGP+Sup	<u>0.754</u>	<u>0.373</u>	<u>0.532</u>	<u>0.553</u>	<u>0.768</u>	<u>0.398</u>	<u>0.544</u>	<u>0.570</u>
Ontology-Enhanced TFMs								
scCello	0.774	0.426	0.625	0.608	0.794	0.486	0.649	0.643

Table 14: Cell type identification with fine-tuning evaluated on the ID dataset D^{id} , as the pre-training subset data size for fine-tuning increases from 0.1% to 1% for the subset selection ratio.

Methods	Cell Type Classification		Cell Type Clustering
	Acc↑ (0.1% → 1%)	Macro F1↑ (0.1% → 1%)	AvgBio↑ (0.1% → 1%)
Ontology-Agnostic TFMs			
Geneformer	0.747 → 0.872	0.440 → 0.664	0.439 → 0.469
scGPT	0.712 → 0.862	0.344 → 0.636	0.477 → 0.481
scTab	0.778 → 0.773	0.373 → 0.455	0.606 → 0.589
MGP	0.722 → 0.861	0.287 → 0.639	0.607 → 0.631
Sup	0.812 → <u>0.902</u>	0.363 → 0.718	<u>0.659</u> → <u>0.668</u>
MGP+Sup	<u>0.820</u> → <u>0.902</u>	<u>0.406</u> → <u>0.735</u>	0.607 → 0.667
Ontology-Enhanced TFMs			
scCello	0.867 → 0.910	0.511 → 0.761	0.694 → 0.699

843 pre-trained using scCello’s codebase (MGP, Sup, and MGP+Sup). For other TFM baselines, we
844 searched for the optimal learning rate to report the final performance.

845 **Performance.** In Sec. 4.2.2, we reported classification and clustering metrics for TFMs fine-tuned
846 with the 0.1% subset of the pre-training data. Here, we extend our reporting to TFMs fine-tuned
847 with 1% of a pre-training subset that is $10 \times$ larger. We compare performances at these two subset
848 selection ratios in Tab. 14. We observe that,

- 849 (1) As the size of fine-tuning data increases, all TFMs except scTab show benefits and scCello
850 achieves 48.9% improvement in Macro F1 when the data size gets $10 \times$ larger. scTab’s
851 underperformance may be related to its model capacity, as it employs a TabNet architec-
852 ture [2]—unlike others that use the powerful standard Transformers [57].
- 853 (2) Across both the classification and clustering metrics, scCello’s prevails other TFM baselines
854 by a large margin. Remarkably, even when fine-tuned with a smaller 0.1% pre-training
855 subset, scCello surpasses TFMs fine-tuned with a much larger 1% subset, achieving a 3.9%

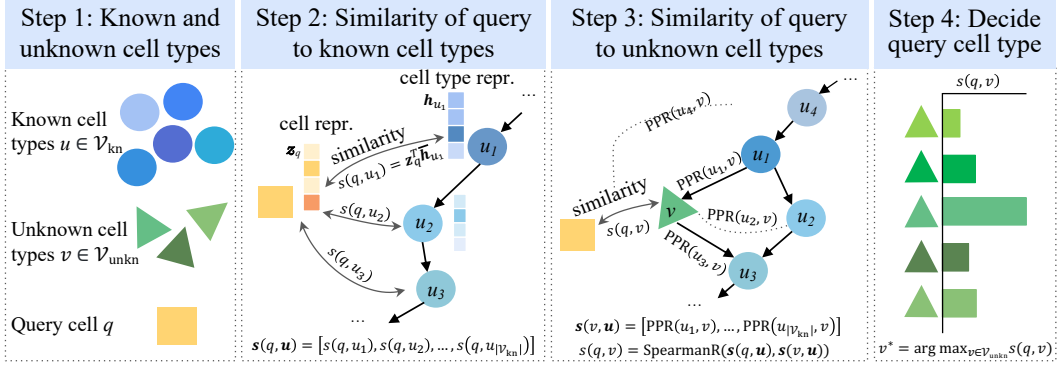


Figure 8: Graphical illustration of our approach for classifying novel cell types (*i.e.*, unknown cell types) (introduced in App. D.3).

856 improvement over the best baseline. This underscores scCello’s superiority, attributed to its
 857 cell ontology-guided pre-training.

858 (3) Interestingly, clustering performance does not necessarily correlate directly with classifica-
 859 tion performance. For instance, while MGP+Sup outperforms Sup in classification metrics,
 860 it does not do so in clustering metrics. This observation underscores the importance to
 861 evaluate both the clustering and classification performances for cell type identification with
 862 model fine-tuning, which can make the evaluation setting more comprehensive and rigorous.

863 D.3 Novel Cell Type Classification

864 **Method.** In this task, we define "known cell types" $\mathcal{V}_{\text{kn}} \subseteq \mathcal{V}$ as the 398 cell types from our
 865 labeled pre-training dataset (see dataset statistics Tab. 8). "Novel cell types", or "unknown cell types"
 866 $\mathcal{V}_{\text{unkn}} \subseteq \mathcal{V}$, are those present only in the target downstream dataset and not observed during TFM
 867 pre-training ($\mathcal{V}_{\text{unkn}} = \mathcal{V} \setminus \mathcal{V}_{\text{kn}}$).

868 Given a new query cell q , we aim to classify it to one of the unknown cell types $\mathcal{V}_{\text{unkn}}$. To solve this
 869 problem, we choose to first calculate representations for both the query cell sample and the unknown
 870 cell types. And then we measure the similarity between the two representations to determine the
 871 prediction results $v_q \in \mathcal{V}_{\text{unkn}}$.

872 Since unknown cell types are absent from the pre-training dataset, their representations cannot be
 873 directly obtained from any TFM baselines or our model, despite its ability to learn representations
 874 for known cell types. To address this problem, we leverage the known cell types \mathcal{V}_{kn} as a bridge to
 875 represent the query cells through the similarity between the cell and cell type representations produced
 876 by TFMs, and also represent the unknown cell types using the structural similarity relationships
 877 between the known and unknown ones derived from the cell ontology graph.

878 Specifically, our approach is illustrated in Fig. 8 and involves the following steps:

879 (1) **Representations for known cell types.** Although scCello inherently learns cell type repre-
 880 sentations during pre-training, most existing TFMs do not output cell type representations
 881 directly. For benchmarking, we propose a protocol to calculate known cell type repre-
 882 sentations for general TFMs. Specifically, the representation for each known cell type is
 883 calculated by averaging cell representations derived from TFMs across cells belonging to
 884 this cell type. We used cell samples from a subset (10%) of our curated pre-training dataset,
 885 because the whole 22 million dataset is too large to fit.

886 We denote the known cell type representations as $\{\bar{h}_u\}_{u \in \mathcal{V}_{\text{kn}}}$, to differentiate with the nota-
 887 tion of scCello’s learned cell type representations $\{h_u\}_{u \in \mathcal{V}_{\text{kn}}}$ introduced in Sec. 2. For fair
 888 comparison, scCello also follows this protocol to generate known cell type representations,
 889 instead of using its learned ones. Nevertheless, we emphasize scCello’s capability to conduct
 890 this task alone without further accessing reference databases like our pre-training dataset.

891 (2) **Similarity vector for a query cell to known cell types.** We first derive the cell repre-
 892 sentations for the query cell q from TFMs. Then, we estimate the similarity between the
 893 query cell q and any known cell type $u \in \mathcal{V}_{\text{kn}}$ using the cosine similarity between their
 894 representations $s(q, u) = \mathbf{z}_q^T \bar{\mathbf{h}}_u$. For all known cell types, this results in a similarity vector:

$$\mathbf{s}(q, \mathbf{u}) = [d(q, u_1), d(q, u_2), \dots, d(q, u_{|\mathcal{V}_{\text{kn}}|})] \quad (26)$$

895 where we define the order of vector indices as $\mathbf{u} = [u_1, u_2, \dots, u_{|\mathcal{V}_{\text{kn}}|}]$ satisfying $u_1 <$
 896 $u_2 < \dots < u_{|\mathcal{V}_{\text{kn}}|}$.

897 (3) **Similarity vector for unknown cell types to known cell types.** For each unknown
 898 cell type $v \in \mathcal{V}_{\text{unkn}}$, we estimate the similarity $s(v, \mathbf{u})$ between the unknown v and the
 899 known cell types \mathbf{u} . To achieve this, we leverage the cell ontology graph to calculate
 900 structural proximities as proxies. The proximities are measured using the raw PPR score
 901 $\text{PPR}(u, v)$, $u \in \mathcal{V}_{\text{kn}}, v \in \mathcal{V}_{\text{unkn}}$, which is introduced in Sec. 2.4. Therefore, the similarity
 902 vector can be represented as:

$$\mathbf{s}(v, \mathbf{u}) = [\text{PPR}(u_1, v), \text{PPR}(u_2, v), \dots, \text{PPR}(u_{|\mathcal{V}_{\text{kn}}|}, v)], \quad (27)$$

903 (4) **Align the similarity vectors for the query cell and the unknown cell types.** Intuitively,
 904 the similarity vector $\mathbf{s}(q, \mathbf{u})$ indicates a profiling for the query cell q , with known cell types
 905 \mathbf{u} as a frame of reference; and the similarity vector $\mathbf{s}(v, \mathbf{u})$ conveys similar profiling for
 906 an unknown cell type v . Therefore, the more similar the two similarity vectors $\mathbf{s}(q, \mathbf{u})$ and
 907 $\mathbf{s}(v, \mathbf{u})$ is, the higher possibility for the query cell to be alike this unknown cell type. We
 908 derive it using Spearman Ratio [46] $\text{SpearmanR}(\cdot)$ as the similarity measure:

$$s(q, v) = \text{SpearmanR}(\mathbf{s}(q, \mathbf{u}), \mathbf{s}(v, \mathbf{u})). \quad (28)$$

909 Other formulas for the vector similarity function are available, like the commonly used
 910 cosine similarity (i.e., $d(q, v) = \mathbf{d}(q, \mathbf{u})^T \mathbf{s}(q, \mathbf{u})$). Our approach is not sensitive to the
 911 choice of the similarity metric. As shown in Fig. 10, using the dot product as the similarity
 912 score led to similar relative performance as in Fig. 2, where scCello generally performs
 913 better or on par with other TFMs. Therefore, we used Spearman Ratio throughout the
 914 experiments.

915 (5) **Select the final answer.** The unknown cell type v^* with the largest distance is selected as
 916 the prediction for novel cell type classification:

$$v^* = \arg \max_{v \in \mathcal{V}_{\text{unkn}}} s(q, v) \quad (29)$$

917 In real-world applications, our approach is still applicable since almost all cell types are included in
 918 the cell ontology graph. But we won't be able to know whether the newly coming query cells are
 919 from unknown cell types $\mathcal{V}_{\text{unkn}}$ or known cell types \mathcal{V}_{kn} . Therefore, we can expand the unknown cell
 920 type set $\mathcal{V}_{\text{unkn}}$ to all the cell type defined in the ontology graph \mathcal{V} , and conduct similar processes in
 921 our approach.

922 **Datasets.** We evaluate on OOD cell type datasets D_1^{ct} and D_2^{ct} . The cell types in D_1^{ct} and D_2^{ct}
 923 are already aligned to the cell ontology graph using the ontology identifiers provided by CellxGene
 924 database, and are a subset of all the unknown cell types $\mathcal{V}_{\text{unkn}}$. We recognize that the prediction task
 925 becomes more challenging as the number of novel cell types increases. Therefore, we constrain the
 926 complete unknown cell type set to the cell types occurred in the datasets we used.

927 To further reflect the challenge, we created five difficulty levels, where the number of cell types
 928 spanned from 10%, 25%, 50%, 75% to 100% of the total cell type count. For example, if we use 25%
 929 cell types in the OOD cell type dataset D_1^{ct} with a total 87 cell types, the unknown cell types include
 930 $(87 \times 25\% \approx 22)$ randomly selected cell types from the complete set $\{c_i | X_i \in D_1^{ct}\}$. To account for
 931 potential biases, we randomly sampled 20 distinct combinations of cell types for each difficulty level.

932 **Hyper-parameters.** The $\text{PPR}(\cdot)$ score is calculated using the "nx.pagerank" function with alpha
 933 hyper-parameter set to 0.9.

934 **Performance.** The full metrics for both accuracy and macro f1 score on the two OOD cell type
 935 datasets D_1^{ct} and D_2^{ct} are reported in Fig. 9. Besides plots, the numerical results are also summarized
 936 in Tab. 16 and Tab. 17 for reference.

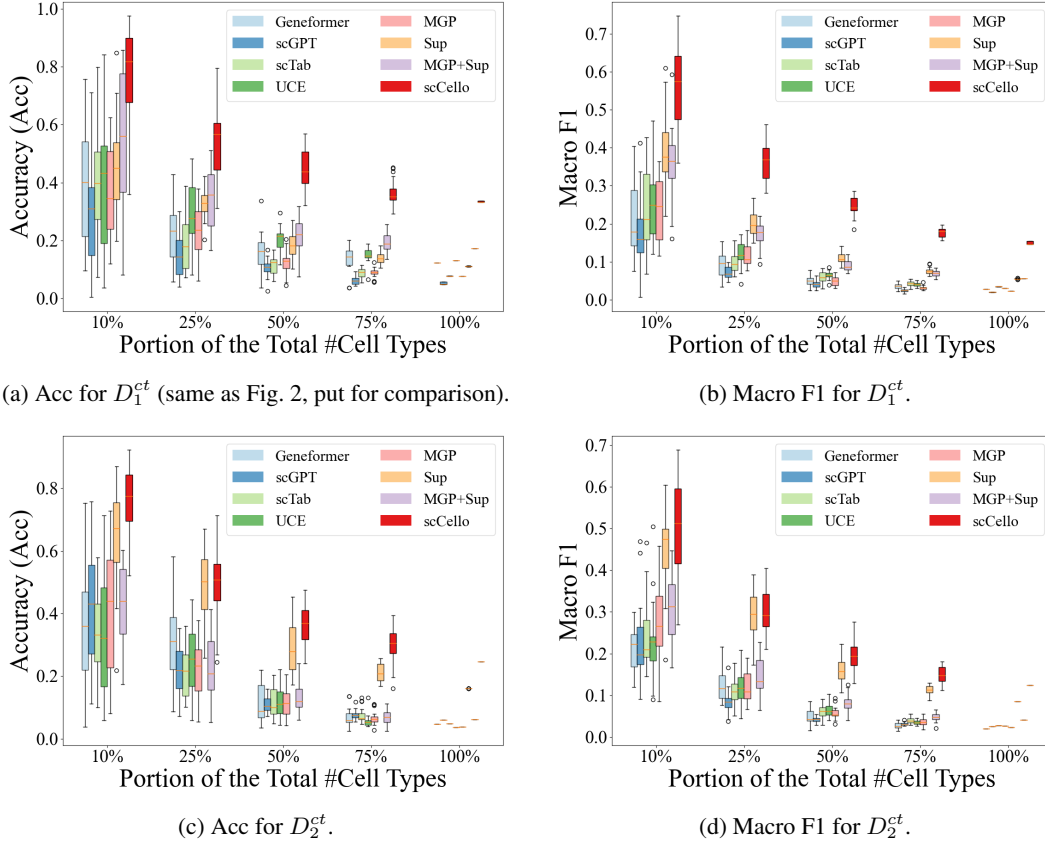


Figure 9: Novel cell type classification on two OOD cell type datasets D_1^{ct} and D_2^{ct} , using the Spearman Ratio similarity measure to compare the representations of the query cells and the novel cell types (App. D.3). Two metrics Acc and Macro F1 are reported.

Table 16: Novel cell type classification results on OOD cell type dataset D_1^{ct} .

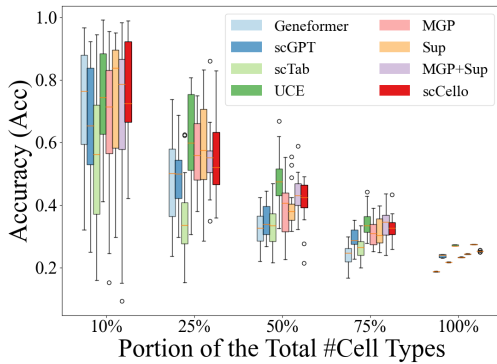
Method	10% cell types		25% cell types		50% cell types		75% cell types		100% cell types	
	Acc \uparrow	F1 \uparrow	Acc \uparrow	F1 \uparrow	Acc \uparrow	F1 \uparrow	Acc \uparrow	F1 \uparrow	Acc \uparrow	F1 \uparrow
Ontology-Agnostic TFMs										
Geneformer	0.392	0.207	0.226	0.095	0.157	0.050	0.135	0.036	0.123	0.027
scGPT	0.291	0.178	0.148	0.072	0.105	0.041	0.062	0.024	0.052	0.020
scTab	0.380	0.248	0.191	0.096	0.114	0.058	0.088	0.042	0.077	0.035
UCE	0.399	0.253	0.289	0.120	0.205	0.064	0.149	0.040	0.131	0.030
MGP	0.361	0.243	0.233	0.119	0.125	0.048	0.089	0.032	0.076	0.022
Sup	0.464	<u>0.389</u>	0.329	<u>0.200</u>	0.187	<u>0.109</u>	0.139	<u>0.075</u>	0.111	0.055
MGP+Sup	<u>0.556</u>	0.358	<u>0.341</u>	0.172	<u>0.217</u>	0.089	<u>0.193</u>	0.069	<u>0.172</u>	<u>0.056</u>
Ontology-Enhanced TFMs										
scCello	0.768	0.559	0.547	0.365	0.442	0.246	0.364	0.177	0.335	0.150

937 D.4 Marker Gene Prediction

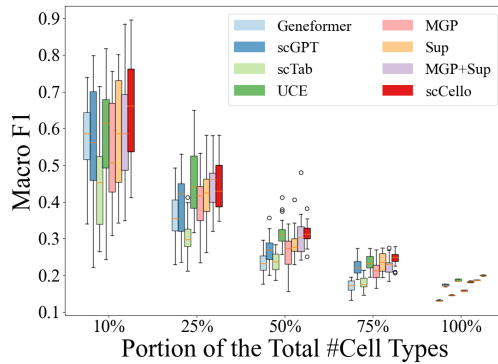
938 **Method.** We here explain our approach for this task in details. Given a cell’s gene expression profile,
 939 we enumerate each gene and attempt to knock it out, either by replacing it with a special [MASK]
 940 token or by reducing its expression to zero. The former method is used for Geneformer, MGP, Sup,
 941 MGP+Sup, and scCello, while the latter is applied to scGPT, scTab, and UCE. By comparing the
 942 cell representations of the mutated expression and those of the original expression, we assess the
 943 impact of each gene’s knockout. A greater impact suggests a higher likelihood of the gene being a

Table 17: Novel cell type classification results on OOD cell type dataset D_2^{ct} .

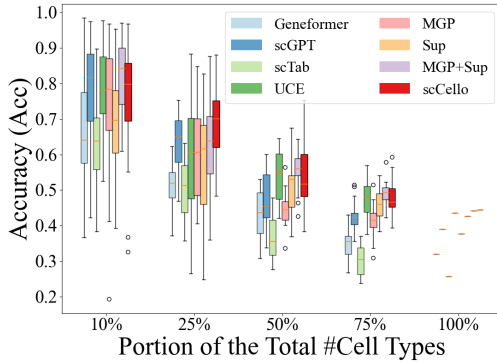
Method	10% cell types		25% cell types		50% cell types		75% cell types		100% cell types	
	Acc \uparrow	F1 \uparrow	Acc \uparrow	F1 \uparrow	Acc \uparrow	F1 \uparrow	Acc \uparrow	F1 \uparrow	Acc \uparrow	F1 \uparrow
Ontology-Agnostic TFMs										
Geneformer	0.367	0.213	0.310	0.125	0.107	0.048	0.069	0.027	0.047	0.020
scGPT	0.411	0.223	0.217	0.085	0.108	0.041	0.077	0.031	0.061	0.025
scTab	0.338	0.245	0.175	0.102	0.113	0.053	0.074	0.038	0.049	0.027
UCE	0.339	0.227	0.244	0.119	0.119	0.064	0.056	0.035	0.037	0.027
MGP	0.411	0.270	0.225	0.120	0.114	0.057	0.066	0.036	0.038	0.023
Sup	<u>0.581</u>	<u>0.372</u>	<u>0.325</u>	<u>0.204</u>	<u>0.199</u>	<u>0.112</u>	<u>0.140</u>	<u>0.081</u>	<u>0.108</u>	<u>0.063</u>
MGP+Sup	0.428	0.315	0.228	0.143	0.131	0.082	0.069	0.047	0.061	0.041
Ontology-Enhanced TFMs										
scCello	0.763	0.500	0.498	0.304	0.364	0.196	0.297	0.149	0.247	0.124



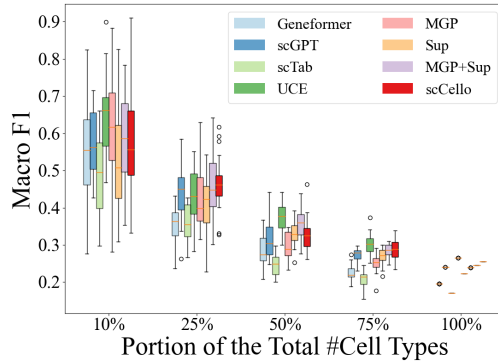
(a) Acc for D_1^{ct} .



(b) Macro F1 for D_1^{ct} .



(c) Acc for D_2^{ct} .



(d) Macro F1 for D_2^{ct} .

Figure 10: Novel cell type classification on two OOD cell type datasets D_1^{ct} and D_2^{ct} , using the cosine similarity measure to compare the representations of the query cells and the novel cell types (App. D.3). Two metrics Acc and Macro F1 are reported.

944 marker gene. This zero-shot approach requires no further fine-tuning and is particularly useful when
 945 additional computational resources or annotated datasets for fine-tuning are unavailable.

946 Notably, we acknowledge the shortage of our method: for house keeping genes (*i.e.*, non-marker
 947 genes), knocking out these genes will also have large impact on the cell because the cell would
 948 die [17]. Therefore, a high impact from gene knockout does not necessarily indicate a marker gene,
 949 but rather an "important" gene. However, this issue is not critical empirically, as the number of
 950 well-documented housekeeping genes is about 400, which is small compared to the extensive gene
 951 token vocabulary of $M = 25,424$.

Table 19: Full results for the five data subsets from GSE96583 (D_1^{mk}) and one dataset from GSE130148 (D_2^{mk}) in the marker gene prediction task (Sec. 4.4).

Method	GSE96583 (D_1^{mk})					Avg.↑	GSE130148 (D_2^{mk})	Avg.↑ of D_1^{mk} and D_2^{mk}
	GSE96583_1 AUROC↑	GSE96583_2 AUROC↑	GSE96583_3 AUROC↑	GSE96583_4 AUROC↑	GSE96583_5 AUROC↑		AUROC↑	
Ontology-Agnostic TFMs								
Geneformer	0.445	0.447	0.478	0.484	0.408	0.452	0.470	0.461
scGPT	0.423	0.387	0.344	0.385	0.388	0.385	0.387	0.386
scTab	0.666	0.654	0.689	0.693	0.660	0.672	0.727	0.700
UCE	0.502	0.499	0.500	0.499	0.500	0.500	0.500	0.500
MGP	0.572	0.560	0.606	0.589	0.567	0.579	0.629	0.604
Sup	0.707	0.697	0.694	0.699	0.700	0.699	0.693	0.696
MGP+Sup	0.734	0.720	0.739	0.734	0.724	0.730	0.730	0.730
Ontology-Enhanced TFMs								
scCello	0.767	0.753	0.754	0.748	0.760	0.756	0.729	0.743

Table 20: Cell types for the two marker gene prediction datasets GSE96583 (D_1^{mk}) and GSE130148 (D_2^{mk}).

Dataset	Cell Types
GSE96583	"Dendritic cells", "CD8 T cells", "NK cells", "B cells", "Megakaryocytes", "FCGR3A+ Monocytes", "CD14+ Monocytes", "CD4 T cells", "Not Known"
GSE130148	"Macrophages", "T cell", "NK cell", "Mast cell", "Endothelium", "Lymphatic", "Pulmonary Alveolar Type II", "Transformed epithelium", "Ciliated", "Pulmonary Alveolar Type I", "B cell", "Fibroblast", "Secretory"

952 **Datasets.** As introduced in Sec. 4.4, we used the datasets from GSE96583 [32] and GSE130148 [6].
953 One the one hand, the GSE96583 dataset D_1^{mk} inherently contains five cell subsets associated
954 with 9 cell type classes. The five cell subsets are denoted as "GSE96583_1", "GSE96583_2",
955 "GSE96583_3", "GSE96583_4", "GSE96583_5", respectively. On the other hand, the GSE130148
956 dataset D_2^{mk} contains 13 cell type classes. The size of these two datasets are summarized in Tab. 21,
957 and their associated cell types are recorded in Tab. 20 for demonstration. Additionally, the ground
958 truth cell-type-specific marker genes are originally sourced from two databases: CellMarker2 [28]
959 and PanglaoDB [21].

960 **Performance.** In Sec. 4.4, we only reported the average performance across the 5 subsets of
961 GSE96583 (D_1^{mk}) and the individual performance of GSE130148 (D_2^{mk}) in Tab 3. Here, we provide
962 complete results for all five subsets in Tab. 19.

963 D.5 Cancer Drug Response Prediction

964 **Method.** In this task, we first compute cell line level representations from scRNA-seq data and drug
965 representations for associated drugs. Both these two representations are then input into the DeepCDR
966 framework for training. Finally, we calculate the PCC between the predicted and actual IC50 values
967 for each drug across all cell lines and report the average performance across all tested drugs.

968 Specifically, for TFMs, single-cell gene expression data are inputted into each model to generate cell-
969 specific representations for each gene. These are then aggregated into cell line-level representations
970 through max-pooling across all genes for each dimension. Conversely, the DeepCDR method uses raw
971 gene expressions, aggregating them directly before max-pooling. Additionally, drugs are represented
972 as graphs and encoded using graph neural networks to obtain drug representations.

973 **Datasets.** In our experiments, we utilized cell line and drug-paired data pre-processed by Deep-
974 CDR [39], including 223 drugs and 561 cell line bulk gene expression profiles for 697 genes from 31

Table 21: The number of cell samples (#Cells) for the marker gene prediction datasets GSE96583 (D_1^{mk}) and GSE130148 (D_2^{mk}).

Dataset	GSE96583_1	GSE96583_2	GSE96583_3	GSE96583_4	GSE96583_5	GSE130148
#Cells	4,246	3,639	14,619	14,446	6,145	10,360

Table 22: The correlation of the ontology structure and the pairwise similarity of known cell type representations

Method	Spearman R \uparrow
Non-TFM Methods	
Raw Data	0.212
Seurat	<u>0.316</u>
Harmony	0.262
Ontology-Agnostic TFMs	
Geneformer	0.284
scGPT	0.037
scTab	0.209
UCE	0.285
MGP	0.275
Sup	0.229
MGP+Sup	0.238
Ontology-Enhanced TFMs	
scCello	0.506

975 different cancer types. Among the dataset, 89,585 cell line-drug samples were used for training and
 976 4,729 for testing [25].

977 **Hyper-parameters.** We following scFoundation’s implementation to set the parameters in the
 978 DeepCDR framework, like “-use_gexp” as True, and both “-use_mut” and “-use_methy” as False.

979 **Performance.** Results are already reported in Tab. 4 in Sec. 4.5.

980 D.6 Batch Integration

981 **Method.** This batch integration task aims to seamlessly integrate scRNA-seq data from different
 982 batches, which can be conducted using the same protocol as cell type clustering. After clustering,
 983 model performance is evaluated. Besides using cell type labels and clustering indices from the
 984 optimized Louvain algorithm to calculate the preservation of biological signals (NMI, ARI, ASW
 985 and AvgBio), this task also use batch labels to measure the removal of batch effects (ASW_b and
 986 AvgBatch). See App. D.1 for metric calculation details.

987 **Datasets.** As introduced in Sec. 4.6, all datasets used in the cell type clustering task (Sec. 4.2.1)
 988 are evaluated, including one ID dataset D^{id} and six OOD datasets D_i^{cond} ($cond \in \{ct, ts, dn\}$,
 989 $i \in \{1, 2\}$).

990 **Hyper-parameters.** We use the same hyper-parameters as that in cell type clustering.

991 **Performance.** In Sec. 4.6, the Overall score, a weighted average of AvgBio and AvgBatch, is
 992 already reported in Fig. 3. Complete results for all metrics are included in Tab. 23 for the ID dataset
 993 D^{id} , Tab. 24 for the OOD cell type datasets D_1^{ct} and D_2^{ct} , Tab. 25 for the OOD tissue datasets D_1^{ts}
 994 and D_2^{ts} , and Tab. 26 for the OOD donor datasets D_1^{dn} and D_2^{dn} .

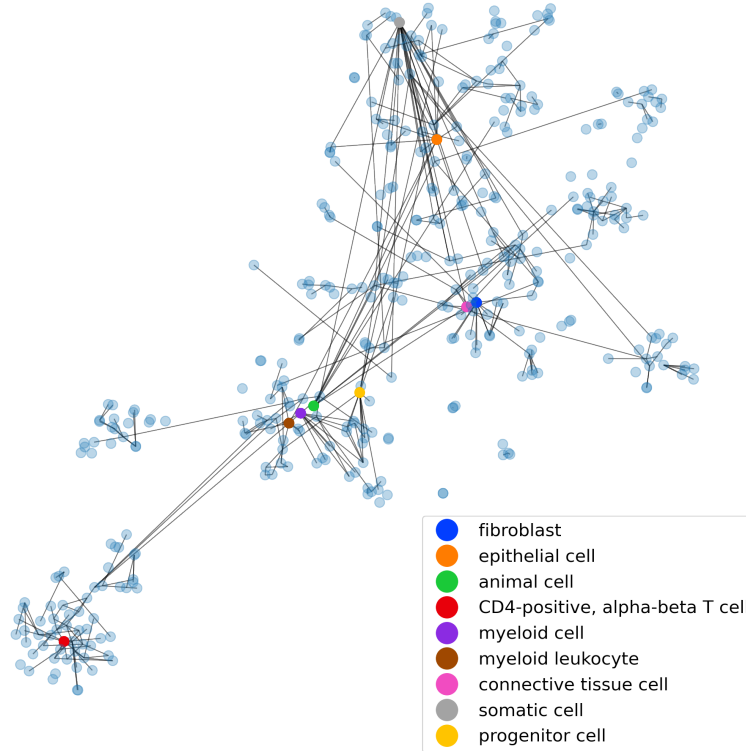


Figure 11: Visualization for learned cell representations of scCello (introduced in App. D.7). The nodes are different cell types in the pre-training dataset and the edges denote "is a subtype of" relationships in cell ontology \mathcal{G} . The coordinates of nodes are calculated using tSNE dimensional reduction for cell type representations derived from scCello. As expected, highly ontology-correlated cell type pairs are very close in the latent space, such as myeloid leukocyte and myeloid cell, as well as fibroblast and connective tissue cell. Meanwhile, dissimilar cell type pairs remain distant, such as CD4-positive, alpha-beta T cell and epithelial cell. The highly biologically informative representation space implies scCello's potential generalization ability to other cell-type-related downstream tasks.

995 D.7 Visualization for Learned Cell Representations

996 We calculate known cell type representation as introduced in Sec. D.3, by averaging cell representa-
 997 tions for each type on 10% of the pre-training data. Then we apply tSNE to project the known cell
 998 type representations to 2D space and visualize in Fig. 11. Highly correlated cell types are clustered
 999 together as expected, and dissimilar cell types are distant.

1000 We also calculate the Spearman R correlation of the pairwise similarity of known cell type representa-
 1001 tions and the ontology structure (1 for an edge between two cell types and 0 for no edge between
 1002 them) in Tab. 22. As expected, scCello learned a biologically informative representation space that
 1003 is much more correlated to the true ontology structure than other methods. This implies scCello's
 1004 potential generalization ability to other cell-type-related downstream tasks.

Table 23: Batch integration on ID dataset D^{id} .

Method	ID Unseen Data (D^{in})				
	ASW _b ↑	GraphConn↑	AvgBatch↑	AvgBio↑	Overall↑
Non-TFM Methods					
Raw Data	0.951	0.806	<u>0.878</u>	0.419	0.603
Seurat	0.829	0.686	<u>0.757</u>	0.442	0.568
Harmony	0.824	0.688	0.756	0.421	0.555
scVI	0.880	0.738	0.809	0.474	0.608
Ontology-Agnostic TFMs					
Geneformer	0.875	0.676	0.775	0.432	0.569
scGPT	0.887	0.691	0.789	0.438	0.578
scTab	<u>0.917</u>	0.925	0.921	<u>0.577</u>	0.715
UCE	0.906	0.788	0.847	0.489	0.632
MGP	0.870	0.728	0.799	0.473	0.603
Sup	0.885	0.809	0.847	0.555	0.672
MGP+Sup	0.892	<u>0.829</u>	0.860	0.516	0.654
Ontology-Enhanced TFMs					
scCello	0.834	0.697	0.766	0.670	<u>0.708</u>

Table 24: Batch integration on OOD cell type datasets D_1^{ct} and D_2^{ct} .

Method	OOD CellType Data (D_1^{ct})					OOD CellType Data (D_2^{ct})				
	ASW _b ↑	GraphConn↑	AvgBatch↑	AvgBio↑	Overall↑	ASW _b ↑	GraphConn↑	AvgBatch↑	AvgBio↑	Overall↑
Non-TFM Methods										
Raw Data	0.934	0.940	0.937	0.703	0.797	0.939	0.895	0.917	0.629	0.744
Seurat	0.831	0.928	0.880	0.752	0.803	0.844	0.932	0.888	0.737	0.797
Harmony	0.909	0.800	0.855	0.432	0.601	0.898	0.817	0.858	0.417	0.593
scVI	0.875	0.959	<u>0.917</u>	<u>0.760</u>	0.823	0.880	0.952	0.916	0.725	0.801
Ontology-Agnostic TFMs										
Geneformer	<u>0.915</u>	0.907	0.911	0.689	0.778	0.915	0.917	0.916	0.668	0.767
scGPT	<u>0.903</u>	0.913	0.908	0.707	0.787	0.896	0.927	0.912	0.720	0.797
scTab	0.908	0.904	0.906	0.759	0.818	0.910	0.905	0.908	0.726	0.799
UCE	0.867	0.947	0.907	0.772	0.826	0.854	<u>0.946</u>	0.900	0.741	<u>0.805</u>
MGP	0.894	0.903	0.899	0.714	0.788	<u>0.925</u>	0.580	0.753	0.740	0.745
Sup	0.879	0.944	0.912	0.767	0.825	0.879	0.914	0.897	<u>0.775</u>	0.824
MGP+Sup	0.885	<u>0.946</u>	0.916	0.758	0.821	0.885	0.925	0.905	0.764	0.820
Ontology-Enhanced TFMs										
scCello	0.877	0.911	0.894	0.769	<u>0.819</u>	0.858	0.884	0.871	0.786	0.820

Table 25: Batch integration on OOD tissue datasets D_1^{ts} and D_2^{ts} .

Method	OOD Tissue Data (D_1^{ts})					OOD Tissue Data (D_2^{ts})				
	ASW _b ↑	GraphConn↑	AvgBatch↑	AvgBio↑	Overall↑	ASW _b ↑	GraphConn↑	AvgBatch↑	AvgBio↑	Overall↑
Non-TFM Methods										
Raw Data	0.941	0.792	0.867	0.540	0.671	0.946	0.862	0.904	0.631	0.740
Seurat	0.865	0.830	0.847	0.587	0.691	0.867	0.841	0.854	0.636	0.723
Harmony	0.905	0.755	0.830	0.462	0.609	0.908	0.744	0.826	0.515	0.639
scVI	0.901	0.861	0.881	0.577	0.699	0.910	0.881	0.896	0.634	0.739
Ontology-Agnostic TFMs										
Geneformer	<u>0.925</u>	0.804	0.865	0.539	0.669	<u>0.924</u>	0.835	0.880	0.597	0.710
scGPT	0.916	0.776	0.846	0.544	0.665	0.920	0.826	0.873	0.627	0.725
scTab	0.916	0.872	<u>0.894</u>	0.515	0.667	0.917	0.874	0.896	0.657	0.753
UCE	0.905	0.864	0.885	0.598	<u>0.713</u>	0.911	0.879	0.895	0.670	0.760
MGP	0.887	0.887	0.887	0.576	0.700	0.901	0.815	0.858	0.628	0.720
Sup	0.903	<u>0.932</u>	0.918	<u>0.605</u>	0.730	0.899	<u>0.911</u>	0.905	<u>0.680</u>	0.770
MGP+Sup	0.900	0.941	0.921	0.610	0.734	0.898	0.922	0.910	0.672	0.767
Ontology-Enhanced TFMs										
scCello	0.868	0.841	0.855	0.612	0.709	0.884	0.819	0.852	0.705	<u>0.764</u>

Table 26: Batch integration on OOD donor datasets D_1^{dn} and D_2^{dn} .

Method	OOD Donor Data (D_1^{dn})					OOD Donor Data (D_2^{dn})				
	ASW $_b$ ↑	GraphConn↑	AvgBatch↑	AvgBio↑	Overall↑	ASW $_b$ ↑	GraphConn↑	AvgBatch↑	AvgBio↑	Overall↑
Non-TFM Methods										
Raw Data	0.945	0.785	0.865	0.458	0.621	0.946	0.787	0.867	0.460	0.623
Seurat	0.875	0.759	0.817	0.466	0.606	0.876	0.771	0.824	0.489	0.623
Harmony	0.893	0.618	0.756	0.456	0.576	0.891	0.655	0.773	0.474	0.594
scVI	0.914	0.831	0.872	0.478	0.636	0.909	0.837	0.873	0.502	0.650
Ontology-Agnostic TFMs										
Geneformer	<u>0.921</u>	0.763	0.842	0.468	0.618	0.919	0.768	0.844	0.482	0.627
scGPT	0.920	0.757	0.839	0.456	0.609	<u>0.920</u>	0.763	0.842	0.477	0.623
scTab	/	/	/	OOM	OOM	/	/	/	OOM	OOM
UCE	0.904	0.665	0.784	0.485	0.605	0.907	0.558	0.733	0.506	0.597
MGP	0.910	0.824	0.867	0.488	0.640	0.906	0.814	0.860	0.518	0.655
Sup	0.909	<u>0.877</u>	<u>0.893</u>	0.552	0.688	0.902	<u>0.857</u>	<u>0.880</u>	<u>0.573</u>	0.696
MGP+Sup	0.910	0.888	0.899	<u>0.553</u>	0.691	0.903	0.869	0.886	0.570	0.696
Ontology-Enhanced TFMs										
scCello	0.845	0.805	0.825	0.608	0.695	0.849	0.802	0.826	0.643	0.716

1005 **NeurIPS Paper Checklist**

1006 **1. Claims**

1007 Question: Do the main claims made in the abstract and introduction accurately reflect the
1008 paper's contributions and scope?

1009 Answer: [\[Yes\]](#)

1010 Justification: In this work, we propose scCello, a novel TFM designed to leverage cell
1011 ontology priors for enhancing cell representation and understanding. This main claim is
1012 accurately and clearly stated and emphasized in the abstract and introduction. It reflects our
1013 contribution and scopes.

1014 Guidelines:

- 1015 • The answer NA means that the abstract and introduction do not include the claims
1016 made in the paper.
- 1017 • The abstract and/or introduction should clearly state the claims made, including the
1018 contributions made in the paper and important assumptions and limitations. A No or
1019 NA answer to this question will not be perceived well by the reviewers.
- 1020 • The claims made should match theoretical and experimental results, and reflect how
1021 much the results can be expected to generalize to other settings.
- 1022 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1023 are not attained by the paper.

1024 **2. Limitations**

1025 Question: Does the paper discuss the limitations of the work performed by the authors?

1026 Answer: [\[Yes\]](#)

1027 Justification: We state several limitations of our work in Sec. 5, such as the lack of continue
1028 learning capability for our proposed model scCello, the relative small model scale for
1029 scCello, and our downstream approach for the zero-shot marker gene experiment.

1030 Guidelines:

- 1031 • The answer NA means that the paper has no limitation while the answer No means that
1032 the paper has limitations, but those are not discussed in the paper.
- 1033 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1034 • The paper should point out any strong assumptions and how robust the results are to
1035 violations of these assumptions (e.g., independence assumptions, noiseless settings,
1036 model well-specification, asymptotic approximations only holding locally). The authors
1037 should reflect on how these assumptions might be violated in practice and what the
1038 implications would be.
- 1039 • The authors should reflect on the scope of the claims made, e.g., if the approach was
1040 only tested on a few datasets or with a few runs. In general, empirical results often
1041 depend on implicit assumptions, which should be articulated.
- 1042 • The authors should reflect on the factors that influence the performance of the approach.
1043 For example, a facial recognition algorithm may perform poorly when image resolution
1044 is low or images are taken in low lighting. Or a speech-to-text system might not be
1045 used reliably to provide closed captions for online lectures because it fails to handle
1046 technical jargon.
- 1047 • The authors should discuss the computational efficiency of the proposed algorithms
1048 and how they scale with dataset size.
- 1049 • If applicable, the authors should discuss possible limitations of their approach to
1050 address problems of privacy and fairness.
- 1051 • While the authors might fear that complete honesty about limitations might be used by
1052 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
1053 limitations that aren't acknowledged in the paper. The authors should use their best
1054 judgment and recognize that individual actions in favor of transparency play an impor-
1055 tant role in developing norms that preserve the integrity of the community. Reviewers
1056 will be specifically instructed to not penalize honesty concerning limitations.

1057 **3. Theory Assumptions and Proofs**

1058 Question: For each theoretical result, does the paper provide the full set of assumptions and
1059 a complete (and correct) proof?

1060 Answer: [NA]

1061 Justification: We propose a novel cell-ontology guided transcriptomic foundation model sc-
1062 Cello in this work. It's a biological insight driven method and focuses on strong downstream
1063 applications. We do not establish theoretical results in this paper.

1064 Guidelines:

- 1065 • The answer NA means that the paper does not include theoretical results.
- 1066 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
1067 referenced.
- 1068 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 1069 • The proofs can either appear in the main paper or the supplemental material, but if
1070 they appear in the supplemental material, the authors are encouraged to provide a short
1071 proof sketch to provide intuition.
- 1072 • Inversely, any informal proof provided in the core of the paper should be complemented
1073 by formal proofs provided in appendix or supplemental material.
- 1074 • Theorems and Lemmas that the proof relies upon should be properly referenced.

1075 4. Experimental Result Reproducibility

1076 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
1077 perimental results of the paper to the extent that it affects the main claims and/or conclusions
1078 of the paper (regardless of whether the code and data are provided or not)?

1079 Answer: [Yes]

1080 Justification: We provide all the details needed to reproduce our work, such as all experi-
1081 mental results for every metric and every dataset, pre-training setups in Sec. 4.1 and App. C,
1082 and downstream task settings in App. D.

1083 Guidelines:

- 1084 • The answer NA means that the paper does not include experiments.
- 1085 • If the paper includes experiments, a No answer to this question will not be perceived
1086 well by the reviewers: Making the paper reproducible is important, regardless of
1087 whether the code and data are provided or not.
- 1088 • If the contribution is a dataset and/or model, the authors should describe the steps taken
1089 to make their results reproducible or verifiable.
- 1090 • Depending on the contribution, reproducibility can be accomplished in various ways.
1091 For example, if the contribution is a novel architecture, describing the architecture fully
1092 might suffice, or if the contribution is a specific model and empirical evaluation, it may
1093 be necessary to either make it possible for others to replicate the model with the same
1094 dataset, or provide access to the model. In general, releasing code and data is often
1095 one good way to accomplish this, but reproducibility can also be provided via detailed
1096 instructions for how to replicate the results, access to a hosted model (e.g., in the case
1097 of a large language model), releasing of a model checkpoint, or other means that are
1098 appropriate to the research performed.
- 1099 • While NeurIPS does not require releasing code, the conference does require all submis-
1100 sions to provide some reasonable avenue for reproducibility, which may depend on the
1101 nature of the contribution. For example
 - 1102 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
1103 to reproduce that algorithm.
 - 1104 (b) If the contribution is primarily a new model architecture, the paper should describe
1105 the architecture clearly and fully.
 - 1106 (c) If the contribution is a new model (e.g., a large language model), then there should
1107 either be a way to access this model for reproducing the results or a way to reproduce
1108 the model (e.g., with an open-source dataset or instructions for how to construct
1109 the dataset).

1110 (d) We recognize that reproducibility may be tricky in some cases, in which case
1111 authors are welcome to describe the particular way they provide for reproducibility.
1112 In the case of closed-source models, it may be that access to the model is limited in
1113 some way (e.g., to registered users), but it should be possible for other researchers
1114 to have some path to reproducing or verifying the results.

1115 5. Open access to data and code

1116 Question: Does the paper provide open access to the data and code, with sufficient instruc-
1117 tions to faithfully reproduce the main experimental results, as described in supplemental
1118 material?

1119 Answer: [No]

1120 Justification: Our code and datasets will be released upon acceptance.

1121 Guidelines:

- 1122 • The answer NA means that paper does not include experiments requiring code.
- 1123 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
1124 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1125 • While we encourage the release of code and data, we understand that this might not be
1126 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
1127 including code, unless this is central to the contribution (e.g., for a new open-source
1128 benchmark).
- 1129 • The instructions should contain the exact command and environment needed to run to
1130 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
1131 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1132 • The authors should provide instructions on data access and preparation, including how
1133 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1134 • The authors should provide scripts to reproduce all experimental results for the new
1135 proposed method and baselines. If only a subset of experiments are reproducible, they
1136 should state which ones are omitted from the script and why.
- 1137 • At submission time, to preserve anonymity, the authors should release anonymized
1138 versions (if applicable).
- 1139 • Providing as much information as possible in supplemental material (appended to the
1140 paper) is recommended, but including URLs to data and code is permitted.

1141 6. Experimental Setting/Details

1142 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
1143 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
1144 results?

1145 Answer: [Yes]

1146 Justification: All the training and test details are justified in App. C for pre-training, and in
1147 App. D for downstreams.

1148 Guidelines:

- 1149 • The answer NA means that the paper does not include experiments.
- 1150 • The experimental setting should be presented in the core of the paper to a level of detail
1151 that is necessary to appreciate the results and make sense of them.
- 1152 • The full details can be provided either with the code, in appendix, or as supplemental
1153 material.

1154 7. Experiment Statistical Significance

1155 Question: Does the paper report error bars suitably and correctly defined or other appropriate
1156 information about the statistical significance of the experiments?

1157 Answer: [Yes]

1158 Justification: In this work, we provide error bars whenever random sampling on key exper-
1159 imental factors is involved. For example, in the novel cell type prediction task, we have
1160 random sampling procedures for cell type combinations and we report box plots with error
1161 bars to demonstrate the statistical significance of this experiments.

1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: As indicated in Tab. 9, the pre-training our proposed model requires training for 2 days on $4 \times$ A100 NVIDIA GPUs, each with 40G GPU memory.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: We follow the NeurIPS Code of Ethics throughout the entire project.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266

Answer: [Yes]

Justification: We discuss the social impact in Sec. 5 for both positive and negative influences.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]

Justification: All the datasets we utilized for both pre-training and downstream tasks are publicly available, and we did not implement explicit safeguards for these datasets. While our proposed model, scCello, aims to advance our understanding of cell representation learning for scientific discovery purposes and carries relatively little inherent risk for misuse, we acknowledge that as an open-source model, we cannot guarantee zero potential for misuse if the methods were to fall into malicious hands. Despite our intentions for beneficial applications, the open availability of scCello introduces a degree of uncertainty regarding potential mishandling or nefarious exploitation that we cannot definitively preclude.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

1267 Justification: All baseline methods and datasets are properly credited. Their license and
1268 terms of use are properly respected.

1269 Guidelines:

- 1270 • The answer NA means that the paper does not use existing assets.
- 1271 • The authors should cite the original paper that produced the code package or dataset.
- 1272 • The authors should state which version of the asset is used and, if possible, include a
1273 URL.
- 1274 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1275 • For scraped data from a particular source (e.g., website), the copyright and terms of
1276 service of that source should be provided.
- 1277 • If assets are released, the license, copyright information, and terms of use in the
1278 package should be provided. For popular datasets, `paperswithcode.com/datasets`
1279 has curated licenses for some datasets. Their licensing guide can help determine the
1280 license of a dataset.
- 1281 • For existing datasets that are re-packaged, both the original license and the license of
1282 the derived asset (if it has changed) should be provided.
- 1283 • If this information is not available online, the authors are encouraged to reach out to
1284 the asset’s creators.

1285 13. New Assets

1286 Question: Are new assets introduced in the paper well documented and is the documentation
1287 provided alongside the assets?

1288 Answer: [No]

1289 Justification: Code and datasets will be released upon acceptance. Clear documentations
1290 will be provided along.

1291 Guidelines:

- 1292 • The answer NA means that the paper does not release new assets.
- 1293 • Researchers should communicate the details of the dataset/code/model as part of their
1294 submissions via structured templates. This includes details about training, license,
1295 limitations, etc.
- 1296 • The paper should discuss whether and how consent was obtained from people whose
1297 asset is used.
- 1298 • At submission time, remember to anonymize your assets (if applicable). You can either
1299 create an anonymized URL or include an anonymized zip file.

1300 14. Crowdsourcing and Research with Human Subjects

1301 Question: For crowdsourcing experiments and research with human subjects, does the paper
1302 include the full text of instructions given to participants and screenshots, if applicable, as
1303 well as details about compensation (if any)?

1304 Answer: [NA]

1305 Justification: We did not perform crowdsourcing experiments and research with human
1306 subjects.

1307 Guidelines:

- 1308 • The answer NA means that the paper does not involve crowdsourcing nor research with
1309 human subjects.
- 1310 • Including this information in the supplemental material is fine, but if the main contribu-
1311 tion of the paper involves human subjects, then as much detail as possible should be
1312 included in the main paper.
- 1313 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1314 or other labor should be paid at least the minimum wage in the country of the data
1315 collector.

1316 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 1317 Subjects

1318 Question: Does the paper describe potential risks incurred by study participants, whether
1319 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1320 approvals (or an equivalent approval/review based on the requirements of your country or
1321 institution) were obtained?

1322 Answer: [NA]

1323 Justification: This paper does not involve crowdsourcing nor research with human subjects.

1324 Guidelines:

- 1325 • The answer NA means that the paper does not involve crowdsourcing nor research with
1326 human subjects.
- 1327 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1328 may be required for any human subjects research. If you obtained IRB approval, you
1329 should clearly state this in the paper.
- 1330 • We recognize that the procedures for this may vary significantly between institutions
1331 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1332 guidelines for their institution.
- 1333 • For initial submissions, do not include any information that would break anonymity (if
1334 applicable), such as the institution conducting the review.