

---

# Hierarchical Quantized Diffusion Based Tree Generation Method for Hierarchical Representation and Lineage Analysis

---

Zelin Zang<sup>1,2</sup>, WenZhe Li<sup>1</sup>, Fei Chen<sup>1</sup>, Yongjie Xu<sup>1</sup>, Chang Yu<sup>1</sup>, Zhen Lei<sup>2,3,4</sup>, Stan Z. Li<sup>1\*</sup>

<sup>1</sup>School of Engineering, Westlake University, Hangzhou, China

<sup>2</sup>Centre for Artificial Intelligence and Robotics (CAIR), HKISI-CAS, Hong Kong, China

<sup>3</sup>State Key Laboratory of Multimodal Artificial Intelligence Systems, Institute of Automation, CAS, Beijing, China

<sup>4</sup>School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing, China

## Abstract

In single-cell research, tracing and analyzing high-throughput single-cell differentiation trajectories is crucial for understanding complex biological processes. Key to this is the modeling and generation of hierarchical data that represents the intrinsic structure within datasets. Traditional methods face limitations in terms of computational cost, performance, generative capacity, and stability. Recent VAEs based approaches have made strides in addressing these challenges but still require specialized network modules for each tree branch, limiting their stability and ability to capture deep hierarchical relationships. To overcome these challenges, we introduce diffusion-based approach called HDTree. HDTree captures tree relationships within a hierarchical latent space using a unified hierarchical codebook and quantized diffusion processes to model tree node transitions. This method improves stability by eliminating branch-specific modules and enhancing generative capacity through gradual hierarchical changes simulated by the diffusion process. HDTree’s effectiveness is demonstrated through comparisons on both general-purpose and single-cell datasets, where it outperforms existing methods in terms of accuracy and performance. These contributions provide a new tool for hierarchical lineage analysis, enabling more accurate and efficient modeling of cellular differentiation paths and offering insights for downstream biological tasks. The code of HDTree is available at anonymous link [https://anonymous.4open.science/r/code\\_HDTree\\_review-A8DB](https://anonymous.4open.science/r/code_HDTree_review-A8DB).

## 1 Introduction

In single-cell research, tracing and analyzing cellular differentiation trajectories is essential for understanding dynamic biological processes. This task requires not only effective modeling of hierarchical structures [Zeng et al., 2022], but also the ability to generate data that faithfully captures such hierarchies [Guo et al., 2024]. Accurately characterizing the hierarchical organization underlying cell differentiation facilitates the exploration of cellular systems and fate decisions, while conditional generation based on these hierarchies enables interpretable discovery of biological mechanisms. Importantly, hierarchical structures are not exclusive to biology—they also emerge in various domains such as recommendation systems, molecular design, and knowledge representation [Chehreghani and Chehreghani, 2024, Gyurek et al., 2024, Tian et al., 2024]. Therefore, developing models that can both represent and generate data along hierarchical relationships not only enhances performance

\*Corresponding author: Zhen Lei and Stan Z. Li.(zhen.lei@ia.ac.cn, Stan.ZQ.Li@westlake.edu.cn)

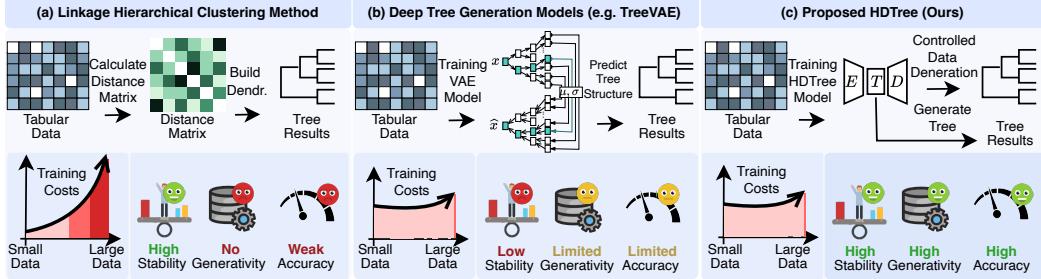


Figure 1: The motivations, base method and deep method cannot meet the requirements of hierarchical representation and lineage analysis in terms of stability, generalivity, accuracy, and training cost.

in downstream tasks such as classification and clustering, but also provides deeper insight into the intrinsic organization of complex data.

As shown in Fig. 1, Traditional methods [Murtagh and Legendre, 2014, Rokhlin and Tygert, 2017] often rely on a combination of dimension reduction [Jia et al., 2022], clustering [Oti and Olusola, 2024], and data regression [Ali and Younas, 2021] techniques to achieve hierarchical modeling and data generation. While these approaches can address the tasks to some extent, they face challenges regarding computation costs, performance, generative capacity, and stability [Zang et al., 2024b]. These limitations make them inadequate for handling the demands of large-scale, high-dimensional biological data. Recent state-of-the-art (SOTA) methods based on Variational Autoencoders (VAEs) [Manduchi et al., 2023, Xiao and Su, 2024, Majima et al., 2024] have unified generative tasks and hierarchical representation within a single modeling framework, achieving notable advancements. These VAE-based approaches effectively reduce computational costs when processing large-scale data while improving performance and data generation capabilities. However, a key limitation of existing SOTA methods lies in these methods often require a specialized network module for each tree branch [Manduchi et al., 2023]. This design not only reduces the stability but also constrains its ability to capture sufficiently deep and complex hierarchical tree relationships, limiting its applicability in scenarios requiring robust and deep hierarchical modeling.

To address these challenges, we propose a novel deep learning-based method, **Hierarchical vector quantized Diffusion Model (HDTREE)**, which captures the tree relationships in a hierarchical latent space with the hierarchical codebook [Huang et al., 2024] and modeling the transfer process of tree node with quantized diffusion model [Gu et al., 2022]. The core innovation of HDTREE lies in its integration of hierarchical latent space encoding with a quantized diffusion process, designed to systematically address limitations in computation costs, performance, generative capacity, and stability. First, stability is ensured by replacing branch-specific network modules with a unified hierarchical codebook, and conditional diffusion steps, eliminating fragmented architecture risks while maintaining adaptability to complex tree topologies. Second, generative capacity is strengthened by modeling branch transitions via a diffusion process [Liu et al., 2024], which simulates gradual hierarchical changes to produce diverse and biologically plausible outputs. Finally, performance gains arise from soft contrastive learning and multi-scale latent space regularization, which sharpen the representation of hierarchical dependencies and improve lineage analysis accuracy. The coordinated work of the above modules improves the performance of the entire model.

These advancements collectively enable HDTREE to capture deep hierarchical relationships robustly, while the learned tree-structured embeddings can be directly applied to downstream tasks—such as lineage analyses via computationally efficient graph-based algorithms. *The contributions of this work are as:* (a) We propose, HDTREE, a novel hierarchical (tree) embeddings & data generation method that captures complex hierarchical relationships and generates high-quality data. (b) We apply HDTREE to the task of lineage analyses. It analyzes the cell differentiation path through a pathfinding algorithm based on the generated tree structure. (c) Comparisons and visualizations of clustering performance, tree performance, generative performance, and lineage analyses performance on general-purpose datasets and single-cell datasets show that HDTREE surpasses existing methods in terms of accuracy and performance, providing new tools for hierarchical lineage analysis.

## 2 Related Work

**Tree-Structured Representation & Generative Models.** Tree-structured representations are crucial for modeling hierarchical relationships in data [Zang et al., 2024b]. Traditional methods like

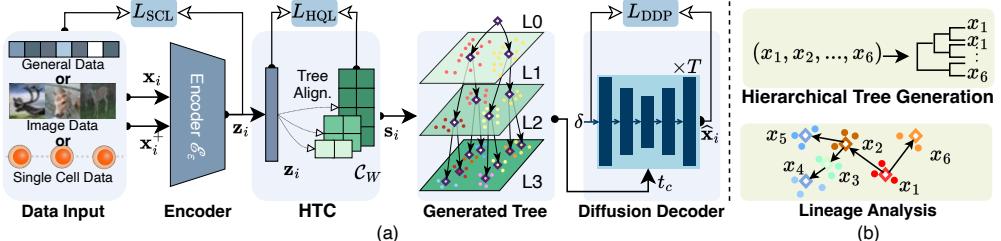


Figure 2: **Overview of the HDTTree framework & tasks.** (a) The framework of HDTTree, which consists of three main components: encoder for semantic representation, Hierarchical Tree Codebook (HTC) for tree-based structural modeling, and diffusion-based decoder for data generation. We use tree structures to model hierarchical relationships and generate data based on the hierarchical latent space. The soft contrastive loss (SCL), hierarchical quantization loss (HQL), and diffusion loss (DDP) are used to optimize the model. (b) The tasks of HDTTree include hierarchical tree generation and lineage analysis.

hierarchical clustering [Müllner, 2011] and distance-based techniques [Kaufman and Rousseeuw, 2009, Bouguettaya et al., 2015] rely on predefined metrics for tree construction. Recent deep learning advancements, such as NVAE [Vahdat and Kautz, 2020] and TreeVAE [Manduchi et al., 2024], leverage recursive and hierarchical latent structures. Hyperbolic geometry methods, including Poincaré Embeddings [Nickel and Kiela, 2017] and HGNNs [Zhou et al., 2023], offer efficient hierarchical representations. TreeVI [Xiao and Su, 2024] enhances variational inference by utilizing tree structures for scalable training and improved performance in tasks like clustering and link prediction. Tree-based generative models, such as GAN-Tree [Kundu et al., 2019], address mode-collapse issues with a hierarchical divisive strategy and enable incremental updates by modifying specific tree branches.

**Deep Learning Based Cell Lineage Analysis.** Cell lineage analysis is crucial for reconstructing developmental trajectories in single-cell genomics. Traditional methods like *Monocle*[Trapnell et al., 2014] and *Slingshot*[Street et al., 2018] infer pseudotime trajectories but are limited by predefined metrics and difficulty modeling unobserved progenitor states. Recent approaches such as *LineageVAE*[Majima et al., 2024] and *Waddington-OT*[Schiebinger et al., 2021] overcome some limitations with probabilistic models and optimal transport, though high dimensionality and sparsity remain challenges.

### 3 Methods

#### 3.1 Notation & Problem Definition

Given input data  $\mathbf{x}_i \in \mathbb{R}^D$  represents the original high-dimensional features of the dataset, where  $D$  is the number of input dimensions. To enhance the generalization capability and robustness of the model, an augmented version of the input data, denoted as  $\mathbf{x}_i^+ \in \mathbb{R}^D$ , is generated through a set of data augmentation techniques. In this work, we adopt knn-based augmentation [Zang et al., 2024a] to generate  $\mathbf{x}_i^+$ , which is used for contrastive learning.

**Definition 1 (Hierarchical Tree Generation Task).** Given a dataset of high-dimensional observations  $\{\mathbf{x}_i\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$ , the hierarchical tree generation task seeks to construct a tree-structured latent representation  $\mathcal{T}$  that captures multi-scale semantic relationships among the data points. Formally, the tree  $\mathcal{T}$  is defined as a rooted binary tree with maximum depth  $L$ , where each node at depth  $l$  is associated with a learnable code vector in a latent space.

**Definition 2 (Lineage Analysis Task).** Given the hierarchical tree  $\mathcal{T}$  generated by HDTTree, the lineage analysis task aims to infer developmental trajectories of cells by identifying the shortest path between a predefined origin and destination in the tree. The inferred trajectories provide a comprehensive representation of the underlying hierarchical relationships in the data, capturing the transitions between different cell states and the progression of cell differentiation processes.

#### 3.2 Model Design

HDTTree is achieved through three key modules: encoder  $\mathcal{E}_\epsilon$ , hierarchical VQ codebook  $\mathcal{C}_W$ , and diffusion decoder  $\mathcal{D}_\theta$ .

**Encoder  $\mathcal{E}_\epsilon$  & Hierarchical Tree Codebook (HTC)  $\mathcal{C}_W$ .** The encoder  $\mathcal{E}_\epsilon$  maps input data  $\mathbf{x}_i$  into a latent space  $\mathbf{z}_i \in \mathbb{R}^d$ , where  $\epsilon$  represents the learnable parameters of the encoder and  $d$  is the latent dimensions. The encoding process can be expressed as,  $\mathbf{z}_i = \mathcal{E}_\epsilon(\mathbf{x}_i)$ ,  $\mathbf{z}_i^+ = \mathcal{E}_\epsilon(\mathbf{x}_i^+)$ . To capture the

hierarchical relationships in the data, HTC  $\mathcal{C}_W$  is introduced, which is constructed as a binary tree, where each node represents a code vector in the latent space,

$$\mathcal{T}_j^l = \begin{cases} \mathbf{w}_j^l, & l = L \\ (\mathbf{w}_j^l, \mathcal{T}_{2j}^{(l+1)}, \mathcal{T}_{2j+1}^{(l+1)}), & l < L, \end{cases} \quad (1)$$

where  $\mathbf{w}_j^l$  is the learnable code vector at depth  $l$  and index  $j$ ,  $\mathcal{T}_0^0$  is the root node of the tree,  $L$  is the maximum depth of the tree. The  $\mathbf{w}_j^l \in \mathbb{R}^d$  is the node embeddings at level  $l$ , for level  $l$ , we have  $w^l$  nodes. The tree structure is optimized during training to capture the semantic and structural relationships in the data. For input data  $\mathbf{x}_i$ , the HTC is used to quantize the latent representation  $\mathbf{z}_i$  into a hierarchical sequence of code vectors  $\mathbf{s}_i$ ,

$$\mathbf{s}_i = [\Omega^{w1}(\mathbf{z}_i), \dots, \Omega^{wl}(\mathbf{z}_i), \dots, \Omega^{wL}(\mathbf{z}_i)], \Omega^{wl}(\mathbf{z}_i) = \operatorname{argmin}_{\mathbf{w}_j^l \in \text{Children}(\mathbf{w}_j^{l-1})} \|\mathbf{z}_i - \mathbf{w}_j^l\|_2, \quad (2)$$

where  $\Omega^l(\mathbf{z}_i)$  is the nearest code vector to  $\mathbf{z}_i$  at depth  $l$  of the tree. The  $\text{Children}(\mathbf{w}_j^{l-1})$  denotes the children of the code vector  $\mathbf{w}_j^{l-1}$  at level  $l - 1$  which defined in Eq. (1).

**Diffusion Decoder**  $\mathcal{D}_\theta$ . The diffusion decoder  $\mathcal{D}_\theta$  reconstructs data or generates new samples based on the hierarchical latent representations. It leverages a Denoising Diffusion Probabilistic Model (DDPM) [Ho et al., 2020b] to iteratively generate data starting from a noise distribution. Specifically, beginning with Gaussian noise  $\mathbf{x}_T \sim \Omega(0, \mathbf{I})$ , the model refines  $\mathbf{x}_T$  through  $T$  diffusion steps to produce the final data  $\mathbf{x}_0$ . The generation process  $\tilde{\mathbf{x}}_i = \text{Gen}(\delta, \mathbf{s}_i | \mathcal{D}_\theta(\cdot))$  is formulated as,

$$\text{Gen}(\delta, \mathbf{s}_i | \phi^*) = \left\{ \tilde{\mathbf{x}}^0 | \tilde{\mathbf{x}}^{t-1} = \frac{1}{\sqrt{\alpha_t}} (\tilde{\mathbf{x}}^t - \tilde{\alpha}) + \sigma_t \Omega(0, 1) \right\}, \tilde{\alpha} = \frac{1 - \alpha_t}{\sqrt{1 - \tilde{\alpha}_t}} \mathcal{D}_\theta(\tilde{\mathbf{x}}^t, t, \mathbf{s}_i), \quad (3)$$

where  $t \in \{T, \dots, 1\}$ , the  $\mathcal{D}_\theta(\cdot)$  is a neural network approximator intended to predict  $\delta$  with  $\tilde{\mathbf{x}}$ .

### 3.3 Loss Function Design

To optimize the HDT tree model, we design a composite loss function that integrates contrastive learning, vector quantization, and diffusion-based reconstruction.

**Soft Contrastive Learning Loss (SCL)** [ $\mathcal{L}_{\text{SCL}}(\cdot)$ ]. The SCL [Zang et al., 2024a] is redefined to preserve local similarities and hierarchical relationships. For a batch of data embeddings  $\{\mathbf{z}_i\}_{i=1}^{N_b}$  and their augmented counterparts  $\{\mathbf{z}_i^+\}_{i=1}^{N_b}$ , the loss is formulated as,

$$\mathcal{L}_{\text{SCL}} = \frac{1}{2N} \sum_{i_1=1}^{N_b} \left( \log \sum_{i_2=1}^{N_b} \mathbf{S}_{i_1 i_2}^{\mathbf{z} \mathbf{z}^+} + \log \sum_{i_2=1}^{N_b} \mathbf{S}_{i_1 i_2}^{\mathbf{z}^+ \mathbf{z}} \right) - \sum_{i_1=1}^{N_b} \log \operatorname{diag}(\mathbf{S}_{i_1 i_1}^{\mathbf{z} \mathbf{z}^+}), \quad (4)$$

where  $N_b$  is the batch size,  $\mathbf{S}_{ij}^{\mathbf{z} \mathbf{z}^+}$  represents the similarity matrix calculated using the t-distribution kernel:  $\mathbf{S}_{i_1 i_2} = (1 + (\mathbf{D}_{i_1 i_2}^2)/\nu)^{-\frac{\nu+1}{2}}$ , and  $\mathbf{D}_{i_1 i_2}$  is the pairwise distance between  $\mathbf{z}_{i_1}$  and  $\mathbf{z}_{i_1 i_2}$  in the hyperbolic space. Here,  $\nu = 0.1$  is the degrees of freedom of the t-distribution.

**Hierarchical Quantization Loss (HQL)** [ $\mathcal{L}_{\text{HQL}}(\cdot)$ ]. The HQL is designed to enable the learning of unified and robust hierarchical tree-structured representations within the framework of the HTC. Unlike vanilla VQ encoders [Van Den Oord et al., 2017], which rely heavily on initialization and local perceptual losses to achieve vector quantization, HQL approach allows the model to capture global structural relationships across multiple levels of hierarchy, ensuring consistency, thus better Lineage analysis. The HQL is expressed as,

$$\mathcal{L}_{\text{HQL}} = \sum_{l=1}^L \mathcal{A}(\mathbf{z}_i, c_{\mathbf{z}_i}^l) + \lambda \mathcal{A}(c_{\mathbf{z}_i}^l, \Psi^{\mathbf{z}_i}(c_{\mathbf{z}_i}^l)), \quad (5)$$

where  $c_{\mathbf{z}_i}^l$  is the nearest code vector to  $\mathbf{z}_i$  at level  $l$ ,  $\lambda = 2$  is hyperparameters controlling the contributions of the alignment and consistency terms. Similar to Eq. (2),  $\Psi^{\mathbf{z}_i}(\cdot)$  is the nearest code vector to  $\mathbf{z}_i$  in the latent space,

$$\Psi^{\mathbf{z}}(\mathbf{w}_j) = \operatorname{argmin}_{\mathbf{z}_i \in \mathbf{z}} \|\mathbf{z}_i - \mathbf{w}_j\|_2, \quad (6)$$

and  $\mathcal{A}(a, b) = \|\mathbf{sg}(a) - b\|_2^2 + \|a - \mathbf{sg}(b)\|_2^2$ , where  $\mathbf{sg}(\cdot)$  is the stop-gradient operation that prevents the gradients from flowing back to the encoder. Eq. 5 try to addresses two challenges inherent in hierarchical representation learning. The hierarchical constraint term,  $\mathcal{A}(\mathbf{z}_i, c_{\mathbf{z}_i}^l)$  ensures that the relationships between parent and child nodes in the hierarchical structure are explicitly preserved. The  $\mathcal{A}(\mathbf{z}_i, c_{\mathbf{z}_i}^l)$  resolves the instability issues often observed in vanilla VQ methods, where suboptimal initialization or limited local perception can hinder the model's performance.

**Diffusion Loss** [ $\mathcal{L}_{\text{DDP}}(\cdot)$ ]. The diffusion loss trains the decoder to iteratively reconstruct the data from noise while preserving hierarchical relationships. Following the Denoising Diffusion Probabilistic Model (DDPM) [Ho et al., 2020a],  $\mathcal{L}_{\text{DDP}}$

$$= \sum_{t=1}^T \left\{ \left\| \delta - \mathcal{D}_\theta \left( \sqrt{\bar{\alpha}_t} \tilde{\mathbf{x}}_c^t + \sqrt{1 - \bar{\alpha}_t}, t, \mathbf{s}_i \right) \right\|_2^2 \right\}, \quad (7)$$

where the conditional vector  $\mathbf{s}_i$  is generated from the semantic encoder in Eq. (2). The  $\mathcal{D}_\theta(\cdot)$  is the conditional diffusion neural network. The  $\bar{\alpha}_t = 1 - \alpha_t$ ,  $\alpha_t$  is defined in Eq. (3). The  $\tilde{\mathbf{x}}_c^t$  is the intermediate data in the diffusion process, and the  $\tilde{\mathbf{x}}_c^0 = \mathbf{x}_c$ .  $T$  is the time step of generation process.

**Overall Loss Function.** The overall loss is defined as,

$$\mathcal{L} = \mathcal{L}_{\text{SCL}} + \lambda_{\text{HQL}} \mathcal{L}_{\text{HQL}} + \lambda_{\text{DDP}} \mathcal{L}_{\text{DDP}}, \quad (8)$$

where  $\lambda_{\text{HQL}}$ , and  $\lambda_{\text{DDP}}$  are hyperparameters controlling the contributions of the contrastive learning loss, vector quantization loss, and diffusion-based reconstruction loss, respectively. This composite loss ensures balanced optimization across all components of the HDTREE model. The pseudocode for training the HDTREE model is provided in Algorithm 1.

### 3.4 Trajectory Analysis with HDTREE

**Graph Construction.** To infer developmental trajectories, the hierarchical tree structure generated by HDTREE is transformed into a weighted graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{W})$ . The graph  $\mathcal{G}$  inherits all nodes and edges from the hierarchical tree  $\mathcal{T}$ , while additional edges are added within the same depth using a  $k$ -nearest neighbors (KNN) approach. This augmentation enriches the connectivity of the graph by capturing local semantic relationships that are not explicitly represented in the original tree structure. The nodes in the graph correspond to the code vectors  $\mathbf{w}_j^l$  at each depth  $l$  and index  $j$ . The edges include both the hierarchical edges from  $\mathcal{T}$  and the newly introduced edges generated by the KNN process. For each edge  $(j_1, j_2)$  in the graph, the weight is,

$$\mathbf{w}(j_1, j_2) = \begin{cases} \|\mathbf{w}_{j_1} - \mathbf{w}_{j_2}\|_2, & j_1 = j_2 \\ \|\mathbf{w}_{j_1} - \mathbf{w}_{j_2}\|_2 + P^{L-l}, & j_1 \neq j_2, \end{cases} \quad (9)$$

where  $P^{L-l}$  is a penalty term that increases the weight of edges connecting nodes at different depths, ensuring that the developmental trajectories follow the hierarchical structure.

**Trajectory Inference.** Using the constructed graph  $\mathcal{G}$ , developmental trajectories are inferred by identifying the shortest path between a predefined origin  $\mathbf{w}_{\text{start}}$  and a destination  $\mathbf{w}_{\text{end}}$ . The shortest path is computed by minimizing the total edge weights along the trajectory,

$$\text{Path}_{\text{development}} = \arg \min_{\text{Path} \subseteq \mathcal{G}} \sum_{e \in \text{Path}} \mathbf{w}(j_1, j_2). \quad (10)$$

The inferred developmental trajectories provide a comprehensive representation of the underlying hierarchical relationships in the data, capturing the transitions between different cell states and the progression of cell differentiation processes.

## 4 Experiments

**Datasets & Baseline Methods.** To provide a comprehensive comparison of different methods, we use two types of datasets, general tabular, image, and text datasets (Mnist, Fashion-Mnist, 20news-groups, Cifar10) and single-cell datasets (Limb[Zhang et al., 2023], LHCO[He et al., 2022b], Weinreb[Weinreb et al., 2020b], ECL[Qiu et al., 2024b]). The scale and features of these datasets are detailed in Table. 1 and 2. Baseline methods include traditional approaches, such as Agglomerative Clustering (Agg) [Müllner, 2011], t-SNE [Linderman and Steinerberger, 2019], and UMAP [Dalmia

Table 1: **Comparison of tree performance, clustering performance, and reconstruction performance (Rec. Performance) on four general image and text datasets.** The <sup>A</sup> means directly use agglomerative clustering on the embeddings to get the tree performance. The -RL and LL are the reconstruction loss and negative log-likelihood. The best results are highlighted in **bold**. The number after/before  $\pm$  shows the mean/standard deviation with 10 different random seeds. ‘NG’ indicates these methods do not have the generation ability.

Dataset	Method	Tree Performance DP(↑) LP(↑)		Clustering Performance ACC(↑) NMI(↑)		Rec. Performance -RL(↑) LL(↑)		Average
Mnist (image, 70k × 784)	Agg	63.7 ± 0.0	78.6 ± 0.0	69.5 ± 0.0	71.1 ± 0.0	NG	NG	NG
	VAE <sup>A</sup>	79.9 ± 2.2	90.8 ± 1.4	86.6 ± 4.9	81.6 ± 2.0	-84.7 ± 2.6	-87.2 ± 2.0	27.8 ± 2.5
	LadderVAE <sup>A</sup>	81.6 ± 3.9	90.9 ± 2.5	80.3 ± 5.6	82.0 ± 2.1	-87.8 ± 0.7	-99.9 ± 0.3	24.5 ± 2.5
	DeepECT	74.6 ± 5.9	90.7 ± 3.2	74.9 ± 6.2	76.7 ± 4.2	NG	NG	NG
	TreeVAE	87.9 ± 4.9	96.0 ± 1.9	90.2 ± 7.5	90.0 ± 4.6	-80.3 ± 0.2	-92.9 ± 0.2	31.8 ± 3.2
	HDTTree <sup>A</sup>	<b>92.7 ± 0.3</b>	<b>97.1 ± 1.2</b>	<b>97.1 ± 0.1</b>	<b>92.8 ± 0.2</b>	NG	NG	NG
	HDTTree	<b>91.9 ± 2.8</b>	<b>96.6 ± 1.4</b>	<b>96.6 ± 1.4</b>	<b>92.4 ± 1.3</b>	<b>-77.9 ± 1.2</b>	<b>-85.4 ± 1.4</b>	<b>35.7 ± 1.6 (↑3.9)</b>
Fashion-Mnist (image, 70k × 784)	Agg	45.0 ± 0.0	67.6 ± 0.0	51.3 ± 0.0	52.6 ± 0.0	NG	NG	NG
	VAE <sup>A</sup>	44.3 ± 2.5	65.9 ± 2.3	54.9 ± 4.4	56.1 ± 3.2	-231 ± 3.2	-242 ± 3.2	-32.1 ± 3.1
	LadderVAE <sup>A</sup>	49.5 ± 2.3	67.6 ± 1.2	55.9 ± 3.0	60.7 ± 1.4	-231 ± 1.4	-239 ± 1.4	-39.5 ± 1.8
	DeepECT	44.9 ± 3.3	67.8 ± 1.4	51.8 ± 5.7	57.7 ± 3.7	NG	NG	NG
	TreeVAE	53.4 ± 2.4	70.4 ± 2.0	60.6 ± 3.3	64.7 ± 1.4	-226 ± 1.4	-234 ± 1.4	-35.4 ± 2.0
	HDTTree <sup>A</sup>	<b>47.7 ± 1.6</b>	<b>67.1 ± 1.5</b>	<b>64.6 ± 1.9</b>	<b>67.4 ± 1.2</b>	NG	NG	NG
	HDTTree	<b>57.4 ± 0.3</b>	<b>71.8 ± 0.3</b>	<b>71.1 ± 0.2</b>	<b>68.7 ± 0.2</b>	<b>-219 ± 0.1</b>	<b>-228 ± 0.1</b>	<b>-29.9 ± 0.2 (↑5.5)</b>
20news-groups (text, 19k × 2000)	Agg	13.1 ± 0.0	30.8 ± 0.0	26.1 ± 0.0	27.5 ± 0.0	NG	NG	NG
	VAE <sup>A</sup>	7.1 ± 0.3	18.1 ± 0.5	15.2 ± 0.4	11.6 ± 0.3	-45.5 ± 0.1	-44.2 ± 0.3	-6.3 ± 0.3
	LadderVAE <sup>A</sup>	9.0 ± 0.2	20.0 ± 0.7	17.4 ± 0.9	17.8 ± 0.6	-43.5 ± 0.1	-44.3 ± 0.6	-3.9 ± 0.5
	DeepECT	9.3 ± 1.8	17.2 ± 3.8	15.6 ± 3.0	18.1 ± 4.1	NG	NG	NG
	TreeVAE	17.5 ± 1.5	38.4 ± 1.6	32.8 ± 2.3	34.4 ± 1.5	-34.4 ± 1.5	-34.4 ± 1.5	9.1 ± 1.7
	HDTTree <sup>A</sup>	<b>22.0 ± 0.1</b>	<b>45.5 ± 0.4</b>	<b>44.6 ± 0.4</b>	<b>43.7 ± 0.2</b>	NG	NG	NG
	HDTTree	<b>23.7 ± 0.1</b>	<b>44.0 ± 0.2</b>	<b>41.8 ± 0.2</b>	<b>42.6 ± 0.2</b>	<b>-31.1 ± 0.3</b>	<b>-34.1 ± 1.5</b>	<b>19.0 ± 0.4 (↑9.9)</b>
Cifar10 (image, 50k × 32 × 32)	VAE <sup>A</sup>	10.5 ± 2.3	16.3 ± 2.3	16.3 ± 1.6	1.86 ± 4.2	<b>-31.7 ± 2.9</b>	<b>-39.2 ± 2.9</b>	-4.3 ± 2.7
	LadderVAE <sup>A</sup>	12.8 ± 3.9	25.3 ± 3.9	25.3 ± 2.0	7.41 ± 4.9	-41.8 ± 4.7	-40.2 ± 3.7	-1.9 ± 3.9
	DeepECT	10.5 ± 2.5	10.3 ± 2.5	10.3 ± 2.8	0.18 ± 4.2	NG	NG	NG
	TreeVAE	35.3 ± 4.0	53.8 ± 3.9	52.9 ± 7.0	41.4 ± 5.9	-47.0 ± 5.9	-48.3 ± 2.4	14.7 ± 4.9
	HDTTree <sup>A</sup>	<b>44.2 ± 1.5</b>	<b>55.2 ± 1.8</b>	<b>75.9 ± 4.3</b>	<b>55.3 ± 2.5</b>	NG	NG	NG
	HDTTree	<b>43.8 ± 1.7</b>	<b>55.1 ± 1.4</b>	<b>73.2 ± 2.7</b>	<b>53.9 ± 2.0</b>	-34.7 ± 1.9	-40.3 ± 3.6	<b>25.2 ± 2.2 (↑10.5)</b>

and Sia, 2021], as well as state-of-the-art (SOTA) deep learning methods, including VAE [Doersch, 2016, Lim et al., 2020], LadderVAE [Sønderby et al., 2016], DeepECT [Mautz et al., 2020], and TreeVAE [Manduchi et al., 2024]. Additionally, specialized models (Geneformer [Theodoris et al., 2023], LangCell [Zhao et al., 2024], CellPLM [Wen et al., 2024]) tailored for the single-cell domain are incorporated to ensure a thorough evaluation across diverse tasks. More details on the datasets and baselines are provided in the Appendix.

**Evaluation Metrics.** To comprehensively evaluate HDTTree and baseline methods, the testing protocol is divided into three parts: clustering performance, tree structure performance, and reconstruction performance. *Clustering performance* is measured using *Clustering Accuracy (ACC)* [Nazir et al., 2009] and *Normalized Mutual Information (NMI)* [Estévez et al., 2009]. To obtain these metrics, the input data is first mapped into a latent space using the respective method. The clustering results are then derived directly from this latent representation. For methods that do not inherently produce clustering results, hierarchical clustering is applied to the latent space to generate cluster labels. This ensures a fair and consistent comparison across all evaluated methods. **Tree structure performance (Tree performance)** is evaluated using *Leaf Purity (LP)* [Schütze et al., 2008] and *Depth Preservation (DP)* [Rokach and Maimon, 2005]. We predict the tree structures with different methods. **Reconstruction performance** is assessed using *Reconstruction Loss (RL)* and *Log-Likelihood (LL)*. These metrics quantify the ability of a method to recover the original input data from its latent space representation. Details on evaluation metrics are provided in the Appendix.

**Testing Protocol & Implementation.** For all experiments, the data is split into training, validation, and testing sets with an 8:1:1 ratio, ensuring unbiased evaluation. In testing, if the number of points in the dataset is greater than 10,000, we randomly sample 10,000 points from testing dataset. Details on downsampling and its rationale are provided in the Appendix. We implemented HDTTree using PyTorch and trained the model on a single NVIDIA A100 GPU. The model is trained using the AdamW optimizer with a learning rate of 1e-4 and a batch size of 128. The number of diffusion steps

Table 2: **Comparison of tree performance, clustering performance on three single cell datasets.**  
Since most of the methods are not generative models, we did not compare generative performance.

Dataset	Method	Year	Tree Performance		Clustering Performance		Average( $\uparrow$ )
			DP( $\uparrow$ )	LP( $\uparrow$ )	ACC( $\uparrow$ )	NMI( $\uparrow$ )	
Limb (cell lineage, 66,633 cells, celltype:10)	Geneformer <sup>A</sup>	2023	25.6 $\pm$ 5.4	35.9 $\pm$ 0.1	34.1 $\pm$ 0.1	34.9 $\pm$ 0.1	32.6 $\pm$ 1.4
	CellPLM <sup>A</sup>	2024	25.6 $\pm$ 0.1	39.9 $\pm$ 0.1	34.1 $\pm$ 0.2	32.9 $\pm$ 0.2	33.1 $\pm$ 0.2
	LangCell <sup>A</sup>	2024	25.3 $\pm$ 0.1	37.5 $\pm$ 0.1	33.9 $\pm$ 0.1	35.1 $\pm$ 0.1	33.0 $\pm$ 0.1
	TreeVAE <sup>A</sup>	2024	34.7 $\pm$ 1.7	55.6 $\pm$ 1.0	49.8 $\pm$ 0.1	50.0 $\pm$ 0.0	47.5 $\pm$ 0.7
	HDTTree <sup>A</sup>	Ours	<b>38.9<math>\pm</math>1.3</b>	<b>57.9<math>\pm</math>1.0</b>	<b>52.8<math>\pm</math>1.0</b>	<b>49.0<math>\pm</math>0.1</b>	<b>49.7<math>\pm</math>0.9</b>
	HDTTree	Ours	<b>41.0<math>\pm</math>0.4</b>	<b>57.2<math>\pm</math>1.4</b>	<b>55.0<math>\pm</math>1.4</b>	<b>46.6<math>\pm</math>0.4</b>	<b>50.0<math>\pm</math>0.9 (<math>\uparrow</math>2.5)</b>
LHCO (cell lineage, 10,628 cells, celltype:7)	CellPLM <sup>A</sup>	2024	27.0 $\pm$ 1.1	35.8 $\pm$ 2.7	16.8 $\pm$ 3.4	1.65 $\pm$ 5.2	20.3 $\pm$ 3.1
	LangCell <sup>A</sup>	2024	26.5 $\pm$ 1.2	35.2 $\pm$ 0.8	35.2 $\pm$ 0.6	0.02 $\pm$ 0.9	24.2 $\pm$ 0.9
	TreeVAE <sup>A</sup>	2024	38.3 $\pm$ 2.0	52.2 $\pm$ 0.1	37.9 $\pm$ 0.1	31.6 $\pm$ 0.0	40.0 $\pm$ 0.6
	HDTTree <sup>A</sup>	Ours	<b>38.8<math>\pm</math>0.3</b>	<b>52.1<math>\pm</math>0.4</b>	<b>46.4<math>\pm</math>0.3</b>	<b>34.7<math>\pm</math>0.5</b>	<b>43.0<math>\pm</math>0.4</b>
	HDTTree	Ours	<b>42.7<math>\pm</math>0.4</b>	<b>54.0<math>\pm</math>0.3</b>	<b>49.4<math>\pm</math>0.3</b>	<b>34.5<math>\pm</math>0.4</b>	<b>45.2<math>\pm</math>0.3 (<math>\uparrow</math>2.2)</b>
	LangCell <sup>A</sup>	2024	47.4 $\pm$ 0.1	54.8 $\pm$ 0.0	14.3 $\pm$ 0.5	34.3 $\pm$ 0.0	37.7 $\pm$ 0.2
Weinreb (cell lineage, 130,887 cells, celltype:11)	Geneformer <sup>A</sup>	2024	45.1 $\pm$ 0.4	55.3 $\pm$ 0.1	21.4 $\pm$ 0.1	32.3 $\pm$ 0.1	38.5 $\pm$ 0.2
	TreeVAE <sup>A</sup>	2024	60.4 $\pm$ 2.6	61.4 $\pm$ 0.5	41.0 $\pm$ 0.1	35.2 $\pm$ 0.0	49.5 $\pm$ 0.8
	HDTTree <sup>A</sup>	Ours	<b>63.3<math>\pm</math>2.6</b>	<b>78.2<math>\pm</math>1.1</b>	<b>50.6<math>\pm</math>1.0</b>	<b>45.2<math>\pm</math>1.2</b>	<b>59.3<math>\pm</math>1.5 (<math>\uparrow</math>7.5)</b>
	HDTTree	Ours	<b>61.0<math>\pm</math>0.4</b>	<b>67.0<math>\pm</math>0.3</b>	<b>62.6<math>\pm</math>0.3</b>	<b>42.6<math>\pm</math>0.3</b>	<b>58.3<math>\pm</math>0.4</b>

Table 3: **Comparisons on Lineage Ground Truth.** The ratio of observed time points (Lineage Ground Truth) in the  $k$ -neighborhood ( $k=30$ ). Top: LineageVAE dataset. Bottom: *C. elegans* dataset.

	Time	Waddington OT	LineageVAE (semi-supervised)	scVI+Agg (unsupervised)	TreeVAE (unsupervised)	HDTTree (unsupervised)
LineageVAE Dataset	Day 2	22.1%	2.2%	16.1%	12.1%	<b>23.2% (<math>\uparrow</math> +1.1%)</b>
	Day 4	21.4%	37.4%	28.2%	30.4%	<b>38.4% (<math>\uparrow</math> +1.0%)</b>
	Day 6	56.6%	60.3%	53.2%	56.4%	<b>62.0% (<math>\uparrow</math> +1.7%)</b>
C. Elegans Dataset	100–300	–	15.4%	10.8%	13.8%	<b>15.2% (<math>\uparrow</math> -0.2%)</b>
	300–500	–	41.5%	40.6%	32.5%	<b>45.8% (<math>\uparrow</math> +4.3%)</b>
	500–750	–	62.3%	51.8%	62.8%	<b>66.3% (<math>\uparrow</math> +4.0%)</b>

$T$  is set to 1000, and the tree depth  $L$  is set to 10. More details on the implementation are provided in the Appendix and code.

**Comparisons on General Datasets [better stability/accuracy/generativity].** The proposed HDTTree is both a tree generation and data generation method. To ensure a fair comparison of clustering, tree construction, and generation performance across different methods, we adopted the benchmarking strategy described in the benchmark of [Manduchi et al., 2024]. The results are shown in Table 1. The <sup>A</sup> means the methods directly use the agglomerative clustering method on the embeddings to calculate the Tree Performance. **Analysis:** (1) HDTTree achieves superior performance across all evaluated metrics, outperforming traditional and SOTA. This advantage stems from HDTTree’s explicit consideration of hierarchical tree structures, which enhances its ability to capture underlying data relationships. (2) Unlike TreeVAE, HDTTree uses a unified tree representation framework, leading to lower standard variance and enhanced stability. (3) The HDTTree exhibits smaller variance and generally better performance than HDTree<sup>A</sup>, demonstrating that its hierarchical tree representation effectively enhances modeling and generation capabilities. (4) HDTTree’s advantages become more pronounced with dataset complexity increases, highlighting its robustness and effectiveness.

**Comparisons on Single Cell Datasets [better stability/accuracy].** Due to the lack of established benchmarks in the single-cell domain, we incorporated data from high-impact published studies into the TreeVAE benchmark. The results are shown in Table 2. **Analysis:** (1) Similar to the general dataset, HDTTree consistently demonstrates superior performance across all evaluated metrics, including tree structure quality, clustering accuracy, and hierarchical integrity. (2) We observed that the zero-shot capabilities of single-cell large language models are often unsatisfactory and, in some cases, fail to surpass basic single-cell methods. This conclusion has also been validated in recent studies [Lan et al., 2024, He et al., 2024]. In comparison with foundational single-cell models, traditional single-cell tree analysis methods, and TreeVAE, HDTTree shows relative advantages in performance and achieves better stability. (3) These results establish HDTTree as a robust and reliable approach for single-cell data analysis.

**Comparisons on Lineage Ground Truth [better stability/accuracy].** We evaluate the alignment between latent space structure and true developmental progression using the ratio of observed time

Table 4: **Ablation Study on MNIST & ECL Datasets.** The best performance is highlighted in **bold**.

Ablation Setups	MNIST (General Dataset)				ECL (Single-Cell Dataset)			
	DP( $\uparrow$ )	LP( $\uparrow$ )	ACC( $\uparrow$ )	NMI( $\uparrow$ )	DP( $\uparrow$ )	LP( $\uparrow$ )	ACC( $\uparrow$ )	NMI( $\uparrow$ )
<b>A1. Full Model (Ours)</b>	<b>92.7</b>	<b>97.1</b>	<b>96.6</b>	<b>92.4</b>	<b>69.0</b>	<b>83.2</b>	<b>83.2</b>	<b>79.0</b>
A2. w/o HTC	87.4	85.3	84.1	75.2	58.7	71.4	70.8	66.5
A3. w/o SCL ( $\mathcal{L}_{SCL}$ )	78.9	82.1	81.5	73.1	55.6	68.9	68.3	63.1
A4. w/o HQL ( $\mathcal{L}_{HQL}$ )	84.1	89.3	86.8	81.7	61.4	74.2	73.7	69.4

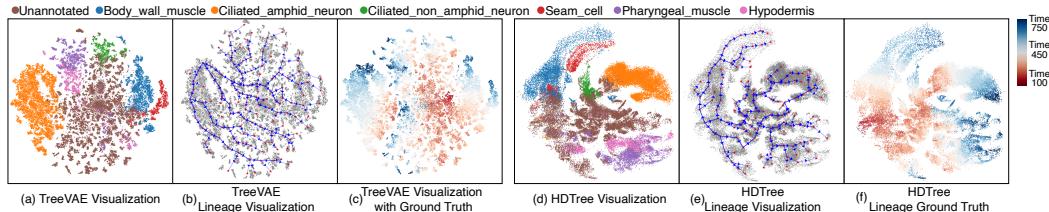


Figure 3: **Comparison of TreeVAE and HDTree methods for visualization and lineage inference.** (a) and (d) are the latent space visualization of TreeVAE and HDTree, color shows the cell type information. (b) and (e) are the lineage structure inferred by TreeVAE and HDTree, overlaid on the data distribution. The blue arrows indicate the inferred lineage relationships. (c) and (f) are the ground truth lineage visualization for comparison, the color shows the real-time information (from blue to red). The results show that HDTree captures more accurate lineage relationships and generates more realistic data than TreeVAE.

points (Lineage Ground Truth) in the  $k$ -nearest neighborhood ( $k = 30$ ). As shown in Table 3, HDTree consistently outperforms both classical and recent unsupervised methods across two benchmark datasets. On the LineageVAE dataset, HDTree achieves the highest local temporal consistency on all time points, even surpassing the semi-supervised LineageVAE by +1.7% at Day 6. This highlights the advantage of our diffusion-based hierarchical modeling in capturing temporal lineage progression without relying on labeled supervision. **Analysis:** On the *C. elegans* dataset, HDTree maintains strong performance across early, mid, and late developmental stages. It improves over TreeVAE by +4.3% and +4.0% in the 300-500 and 500-750 windows, respectively, and even slightly outperforms the supervised LineageVAE in the early stage. These results demonstrate that HDTree’s hierarchical latent structure provides a more faithful reflection of biological differentiation dynamics, offering robust generalization across varying levels of trajectory complexity.

**Case Study on HDTree Data Generation [better generativity].** The diffusion generation of HDTree can generate transformation processes between different tree branches according to the tree structure, which is very important in phylogenetic analysis because often people are interested in how different cell types are transformed under natural conditions (for example, from stem cells to somatic cells). We demonstrate how HDTree solves the above problem based on two datasets (MNIST and *C. elegans*).

The results are shown in Fig. 4. **Analysis:** (1) The results on the generic MNSIT data show the visualization of the result from the number 6 to three different numbers (3, 1, and 9). The result shows that the model can complete the generation process well, and the data changes slowly throughout the process. (2) The results on the *C. elegans* data show the visualization of the result from the stem cell to the somatic cell. Although we cannot visually display every gene, we have selected three iconic genes and then used pie charts to display the trend of data changes.

**Case Study on *C. elegans* Lineage [better accuracy].** To evaluate the performance of HDTree on real lineage labels, we utilized labeled data provided by [Packer et al., 2019]. The *C. elegans* dataset not only labels the type of cell but also the relative time at which the cell is collected, which can be regarded as the gold label for our analysis of the cell’s differentiation lineage. HDTree

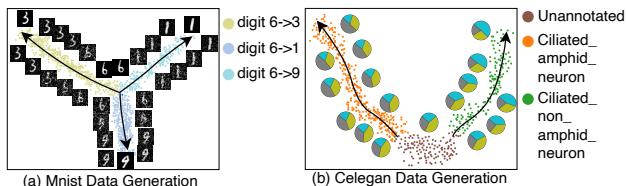


Figure 4: **Data generation processes of HDTree on MNIST and Celegan datasets.** Each scatter represents the generated data visualized by tSNE. Color indicates the label. For MNIST, data is generated from digit 6 to 3, 1, and 9. For Celegan, from stem cell to somatic cell. The pie charts show gene expression of iconic genes: **odr-10**, **osm-6**, **elt-5**.

demonstrates superior performance compared to TreeVAE in capturing lineage relationships and generating biologically meaningful results. A detailed introduction is in the caption of Fig. 3. **Analysis:** (1) By analyzing Fig. 3(a) and Fig. 3(d), both TreeVAE and HDTREE can distinguish different cell types well because they map cells of the same type to similar locations. (2) TreeVAE cannot accurately model the differentiation process of cells. This is because the differentiation lineage visualization (Fig. 3(b)) based on the TreeVAE representation does not match the real-time gold label (Fig. 3(c)). In contrast, the differentiation lineage inferred by HDTREE basically (Fig. 3(e)) overlaps with the time gold label (Fig. 3(f)).

### Comparisons on Computational Cost [better computational effectively].

To evaluate the computational efficiency of HDTREE, we compared the training time of HDTREE with TreeVAE and TreeVAE<sup>A</sup> on four dataset (in Table 5). We observe that traditional methods do have an advantage when dealing with small datasets. However, when the dataset size becomes large, traditional methods will become slow due to the their  $O^2$  complexity.

**Ablation Study [better accuracy].** To evaluate the contributions of HDTREE’s components, we conducted ablation experiments on MNIST and ECL datasets. The setups included **(A1)** Full Model (HDTREE), **(A2)** without the HCL and directly use vanilla codebook, **(A3)** without the SCL ( $\mathcal{L}_{SCL}$ ) and directly use the contrastive learning loss, and **(A4)** without the HQL loss and directly use the VQ Loss ( $\mathcal{L}_{hq}$ ). Performance is measured using tree structure ( $DP$ ,  $LP$ ) and clustering metrics ( $ACC$ ,  $NMI$ ). The results are shown in Table 4. **Analysis:** The full model consistently achieved the best performance. Removing the HCL (A2) caused the most significant performance drop across both datasets, highlighting its role in structural and clustering performance. The SCL (A3) is essential for maintaining tree depth and clustering interpretability. All components contribute to HDTREE’s success, with the HCL and contrastive loss being the most critical for optimal performance.

### Parameter Sensitivity Analysis [better stability].

To evaluate the impact of latent dimensionality on model performance, we conducted a sensitivity analysis using tree performance ( $DP$ ) as the primary metric. The baseline and SOTA methods are evaluated on the ECL dataset with varying latent dimensionalities. The results are presented in Fig. 6. **Analysis:** (1) HDTREE achieves its best performance at a latent dimensionality of 64-256, beyond which no further improvements are observed. This observation suggests that HDTREE effectively captures the essential structural information at moderate dimensionalities, avoiding over-parameterization. (2) Moreover, regardless of the choice of dimensionality, HDTREE maintains a consistent advantage over competing methods, underscoring its superiority in preserving tree structures and hierarchical relationships. The sensitivity of HDTREE to key hyperparameters is examined.

## 5 Conclusion

We introduce HDTREE, a unified diffusion-based framework for hierarchical representation and data generation. By combining a quantized diffusion process with a hierarchical codebook, HDTREE captures tree-structured relationships without relying on branch-specific modules, leading to enhanced stability, generative quality, and interpretability. Experimental results across both general and single-cell datasets demonstrate consistent improvements in clustering accuracy, tree structure fidelity, and lineage alignment. Notably, HDTREE outperforms existing unsupervised and even semi-supervised baselines in modeling complex differentiation trajectories. These findings highlight HDTREE as a robust and generalizable solution for structure-aware learning and biological discovery.

Table 5: Training time comparison on general and single-cell datasets. **Bold** denotes the best result. (mm:ss)

	tSNE+Agg	UMAP+Agg	TreeVAE	HDTREE
MNIST	854:10	<b>2:09</b>	192:09	42:23
F-MNIST	915:13	<b>2:22</b>	206:13	45:02
LHCO	1708:28	<b>13:51</b>	246:20	53:23
Weinreb	5879:27	340:18	361:12	<b>53:47</b>

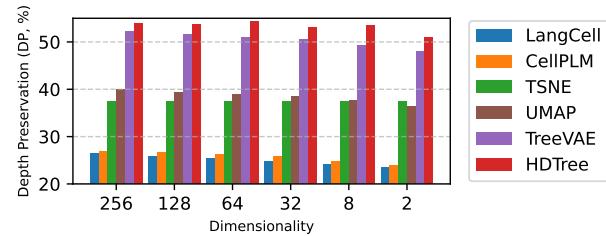


Table 6: Sensitivity Analysis of Dimensionality on Depth Preservation (DP) Across Methods. The performance of various methods on the ECL dataset with different latent dimensionalities.

## References

- Parveen Ali and Ahtisham Younas. Understanding and interpreting regression analysis. *Evidence-Based Nursing*, 24(4):116–118, 2021.
- Athman Bouguettaya, Qi Yu, Xumin Liu, Xiangmin Zhou, and Andy Song. Efficient agglomerative hierarchical clustering. *Expert Systems with Applications*, 42(5):2785–2797, 2015.
- Morteza Haghir Chehreghani and Mostafa Haghir Chehreghani. Hierarchical correlation clustering and tree preserving embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23083–23093, 2024.
- Ayush Dalmia and Suzanna Sia. Clustering with umap: Why and how connectivity matters. *arXiv preprint arXiv:2108.05525*, 2021.
- Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.
- Pablo A Estévez, Michel Tesmer, Claudio A Perez, and Jacek M Zurada. Normalized mutual information feature selection. *IEEE Transactions on neural networks*, 20(2):189–201, 2009.
- Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10696–10706, 2022.
- Zhiye Guo, Jian Liu, Yanli Wang, Mengrui Chen, Duolin Wang, Dong Xu, and Jianlin Cheng. Diffusion models in bioinformatics and computational biology. *Nature reviews bioengineering*, 2(2):136–154, 2024.
- Croix Gyurek, Niloy Talukder, and Mohammad Al Hasan. Binder: Hierarchical concept representation through order embedding of binary vectors. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 980–991, 2024.
- Laleh Haghverdi et al. Diffusion pseudotime robustly reconstructs lineage branching. *Nature Methods*, 13(10):845–848, 2016.
- Fei He, Ruixin Fei, Mingyue Gao, Li Su, Xinyu Zhang, and Dong Xu. Parameter-efficient fine-tuning enhances adaptation of single cell large language model for cell type identification. *bioRxiv*, 2024.
- Z. He, A. Maynard, A. Jain, et al. Lineage recording in human cerebral organoids. *Nat Methods*, 19:90–99, 2022a. doi: 10.1038/s41592-021-01344-8.
- Zhisong He, Ashley Maynard, Akanksha Jain, Tobias Gerber, Rebecca Petri, Hsiu-Chuan Lin, Malgorzata Santel, Kevin Ly, Jean-Samuel Dupré, Leila Sidow, et al. Lineage recording in human cerebral organoids. *Nature methods*, 19(1):90–99, 2022b.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *NeurIPS*, 2020a. URL <https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html>.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv preprint arXiv:2006.11239*, 2020b.
- Lun Huang, Qiang Qiu, and Guillermo Sapiro. Pq-vae: Learning hierarchical discrete representations with progressive quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7550–7558, 2024.
- Weikuan Jia, Meili Sun, Jian Lian, and Sujuan Hou. Feature dimensionality reduction: a review. *Complex & Intelligent Systems*, 8(3):2663–2693, 2022.
- Leonard Kaufman and Peter J Rousseeuw. *Finding groups in data: an introduction to cluster analysis*. John Wiley & Sons, 2009.

- Jogendra Nath Kundu, Maharshi Gor, Dakshit Agrawal, and R Venkatesh Babu. Gan-tree: An incrementally learned hierarchical generative framework for multi-modal data distributions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8191–8200, 2019.
- Wei Lan, Guohang He, Mingyang Liu, Qingfeng Chen, Junyue Cao, and Wei Peng. Transformer-based single-cell language model: A survey. *Big Data Mining and Analytics*, 7(4):1169–1186, 2024.
- Kart-Leong Lim, Xudong Jiang, and Chenyu Yi. Deep clustering with variational autoencoder. *IEEE Signal Processing Letters*, 27:231–235, 2020.
- George C Linderman and Stefan Steinerberger. Clustering with t-sne, provably. *SIAM journal on mathematics of data science*, 1(2):313–332, 2019.
- Qijiong Liu, Xiaoyu Dong, Jiaren Xiao, Nuo Chen, Hengchang Hu, Jieming Zhu, Chenxu Zhu, Tetsuya Sakai, and Xiao-Ming Wu. Vector quantization for recommender systems: A review and outlook. *arXiv preprint arXiv:2405.03110*, 2024.
- Koichiro Majima, Yasuhiro Kojima, Kodai Minoura, Ko Abe, Haruka Hirose, and Teppei Shimamura. Lineagevae: reconstructing historical cell states and transcriptomes toward unobserved progenitors. *Bioinformatics*, 40(10):btae520, 2024.
- Laura Manduchi, Moritz Vandenborth, Alain Ryser, and Julia Vogt. Tree variational autoencoders. *Advances in Neural Information Processing Systems*, 36:54952–54986, 2023.
- Laura Manduchi, Moritz Vandenborth, Alain Ryser, and Julia Vogt. Tree variational autoencoders. *Advances in Neural Information Processing Systems*, 36:54952–54986, 2024.
- Dominik Mautz, Claudia Plant, and Christian Böhm. Deepect: The deep embedded cluster tree. *Data Science and Engineering*, 5:419–432, 2020.
- Daniel Müllner. Modern hierarchical, agglomerative clustering algorithms. *arXiv preprint arXiv:1109.2378*, 2011.
- Fionn Murtagh and Pierre Legendre. Ward’s hierarchical agglomerative clustering method: which algorithms implement ward’s criterion? *Journal of classification*, 31:274–295, 2014.
- KA Abdul Nazeer, MP Sebastian, et al. Improving the accuracy and efficiency of the k-means clustering algorithm. In *Proceedings of the world congress on engineering*, volume 1, pages 1–3. Association of Engineers London, UK, 2009.
- Maximillian Nickel and Douwe Kiela. Poincaré embeddings for learning hierarchical representations. *Advances in neural information processing systems*, 30, 2017.
- E Oti and M Olusola. Overview of agglomerative hierarchical clustering methods. *British Journal of Computer, Networking and Information Technology*, 7:14–23, 2024.
- Jonathan S Packer, Qin Zhu, Chau Huynh, Priya Sivaramakrishnan, Elicia Preston, Hannah Dueck, Derek Stefanik, Kai Tan, Cole Trapnell, Junhyong Kim, et al. A lineage-resolved molecular atlas of *C. elegans* embryogenesis at single-cell resolution. *Science*, 365(6459):eaax1971, 2019.
- C. Qiu, B.K. Martin, I.C. Welsh, et al. A single-cell time-lapse of mouse prenatal development from gastrula to birth. *Nature*, 626:1084–1093, 2024a. doi: 10.1038/s41586-024-07069-w.
- Chengxiang Qiu, Beth K Martin, Ian C Welsh, Riza M Daza, Truc-Mai Le, Xingfan Huang, Eva K Nichols, Megan L Taylor, Olivia Fulton, Diana R O’Day, et al. A single-cell time-lapse of mouse prenatal development from gastrula to birth. *Nature*, 626(8001):1084–1093, 2024b.
- Lior Rokach and Oded Maimon. Clustering methods. *Data mining and knowledge discovery handbook*, pages 321–352, 2005.
- Vladimir Rokhlin and Mark Tygert. Hierarchical clustering for data sets and networks: Practical issues and adaptive algorithms. *Proceedings of the National Academy of Sciences*, 114(29):7535–7540, 2017.

- Geoffrey Schiebinger et al. Optimal transport for developmental trajectories. *Proceedings of NeurIPS*, 2021.
- Hinrich Schütze, Christopher D Manning, and Prabhakar Raghavan. *Introduction to information retrieval*, volume 39. Cambridge University Press Cambridge, 2008.
- Casper Kaae Sønderby, Tapani Raiko, Lars Maaløe, Søren Kaae Sønderby, and Ole Winther. Ladder variational autoencoders. *Advances in neural information processing systems*, 29, 2016.
- Kelly Street, Davide Risso, Robert B Fletcher, et al. Slingshot: cell lineage and pseudotime inference for single-cell transcriptomics. *BMC Genomics*, 19:477, 2018.
- Christina V Theodoris, Ling Xiao, Anant Chopra, Mark D Chaffin, Zeina R Al Sayed, Matthew C Hill, Helene Mantineo, Elizabeth M Brydon, Zexian Zeng, X Shirley Liu, et al. Transfer learning enables predictions in network biology. *Nature*, 618(7965):616–624, 2023.
- Feng Tian, Sen Lei, Yingbo Zhou, Jialin Cheng, Guohao Liang, Zhengxia Zou, Heng-Chao Li, and Zhenwei Shi. Hirennet: Hierarchical-relation network for few-shot remote sensing image scene classification. *IEEE Transactions on Geoscience and Remote Sensing*, 2024.
- Alex Tong and Xin Huang. Trajectorynet: Continuous modeling of cell trajectories with neural networks. *Proceedings of ISMB*, 2020.
- Cole Trapnell, Davide Cacchiarelli, James Grimsby, et al. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nature Biotechnology*, 32(4): 381–386, 2014.
- Arash Vahdat and Jan Kautz. Nvae: A deep hierarchical variational autoencoder. *Advances in neural information processing systems*, 33:19667–19679, 2020.
- Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.
- Caleb Weinreb, Alejo Rodriguez-Fraticelli, Fernando D. Camargo, and Allon M. Klein. Lineage tracing on transcriptional landscapes links state to fate during differentiation. *Science*, 367(6479): eaaw3381, 2020a. doi: 10.1126/science.aaw3381. URL <https://www.science.org/doi/abs/10.1126/science.aaw3381>.
- Caleb Weinreb, Alejo Rodriguez-Fraticelli, Fernando D Camargo, and Allon M Klein. Lineage tracing on transcriptional landscapes links state to fate during differentiation. *Science*, 367(6479): eaaw3381, 2020b.
- Hongzhi Wen, Wenzhuo Tang, Xinnan Dai, Jiayuan Ding, Wei Jin, Yuying Xie, and Jiliang Tang. CellPLM: Pre-training of cell language model beyond single cells. In *ICLR*, 2024. URL <https://openreview.net/forum?id=BKXvPDeKud>.
- Junxi Xiao and Qinliang Su. Treevi: Reparameterizable tree-structured variational inference for instance-level correlation capturing. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Zelin Zang, Yuhao Wang, Jinlin Wu, Hong Liu, Yue Shen, Stan Li, Zhen Lei, et al. Dmt-hi: Moe-based hyperbolic interpretable deep manifold transformation for unsupervised dimensionality reduction. *arXiv preprint arXiv:2410.19504*, 2024a.
- Zelin Zang, Yongjie Xu, Chenrui Duan, Jinlin Wu, Stan Z Li, and Zhen Lei. A review of artificial intelligence based biological-tree construction: Priorities, methods, applications and trends. *arXiv preprint arXiv:2410.04815*, 2024b.
- Andy GX Zeng, Suraj Bansal, Liqing Jin, Amanda Mitchell, Weihsu Claire Chen, Hussein A Abbas, Michelle Chan-Seng-Yue, Veronique Voisin, Peter van Galen, Anne Tierens, et al. A cellular hierarchy framework for understanding heterogeneity and predicting drug response in acute myeloid leukemia. *Nature medicine*, 28(6):1212–1223, 2022.
- B. Zhang, P. He, J.E.G. Lawrence, et al. A human embryonic limb cell atlas resolved in space and time. *Nature*, 635:668–678, 2024. doi: 10.1038/s41586-023-06806-x.

Bao Zhang, Peng He, John EG Lawrence, Shuaiyu Wang, Elizabeth Tuck, Brian A Williams, Kenny Roberts, Vitalii Kleshchevnikov, Lira Mamanova, Liam Bolt, et al. A human embryonic limb cell atlas resolved in space and time. *Nature*, pages 1–11, 2023.

Suyuan Zhao, Jiahuan Zhang, Yushuai Wu, YIZHEN LUO, and Zaiqing Nie. Langcell: Language-cell pre-training for cell identity understanding. In *Forty-first International Conference on Machine Learning*, 2024.

Min Zhou, Menglin Yang, Bo Xiong, Hui Xiong, and Irwin King. Hyperbolic graph neural networks: A tutorial on methods and applications. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 5843–5844, 2023.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Related Work</b>	<b>2</b>
<b>3</b>	<b>Methods</b>	<b>3</b>
<b>4</b>	<b>Experiments</b>	<b>5</b>
<b>5</b>	<b>Conclusion</b>	<b>9</b>
<b>A</b>	<b>Appendix: Details of Related Work</b>	<b>15</b>
<b>B</b>	<b>Appendix: Details of Dataset</b>	<b>15</b>
B.1	Datasets . . . . .	15
B.2	Dataset Organization . . . . .	16
B.3	Preprocessing . . . . .	17
B.3.1	Image Data . . . . .	18
B.3.2	Biological Data . . . . .	18
<b>C</b>	<b>Appendix: Details of Baseline Methods</b>	<b>20</b>
<b>D</b>	<b>Appendix: Details of Testing Protocol</b>	<b>21</b>
<b>E</b>	<b>Appendix: Details of Implementation</b>	<b>22</b>
<b>F</b>	<b>Appendix: Details of Downsampling in Testing</b>	<b>22</b>
<b>G</b>	<b>Appendix: Details of Experimental Environment</b>	<b>23</b>
<b>H</b>	<b>Appendix: Details of Experiments on Single Cell</b>	<b>23</b>
H.1	Detailed Comparisons on More Single Cell Datasets. . . . .	23
H.2	Stability of Model Performance Across Varying Training Data Sizes . . . . .	23

## A Appendix: Details of Related Work

**Tree-Structured Representation & Generative Models.** Tree-structured representations are essential for modeling hierarchical relationships within data [Zang et al., 2024b]. Traditional approaches, such as hierarchical clustering [Müllner, 2011] and distance-based methods [Kaufman and Rousseeuw, 2009, Bouguettaya et al., 2015], use predefined metrics to construct trees and have been foundational in many applications. Recent advancements in deep learning have introduced adaptive techniques for tree construction, such as the Nouveau VAE (NVAE)[Vahdat and Kautz, 2020], which captures hierarchical semantics through recursive structures, and the Tree Variational Autoencoder (Tree-VAE)[Manduchi et al., 2024], which encodes hierarchical latent structures in generative models. Furthermore, hyperbolic geometry has emerged as a powerful framework for representing hierarchical relationships, with methods like Poincaré Embeddings[Nickel and Kiela, 2017] and Hyperbolic Graph Neural Networks (HGNNS)[Zhou et al., 2023] offering efficient and expressive representations of such structures. TreeVI [Xiao and Su, 2024] extends variational inference by using a tree structure to efficiently capture correlations among latent variables in the posterior, enabling scalable reparameterization and training while improving performance in tasks like constrained clustering, user matching, and link prediction. Tree-based generative models offer powerful solutions for modeling hierarchical relationships and multi-modal data distributions. GAN-Tree[Kundu et al., 2019] introduces a hierarchical divisive strategy with a mode-splitting algorithm for unsupervised clustering, effectively addressing mode-collapse and discontinuities in data, while enabling incremental updates by modifying only specific tree branches.

**Deep Learning Based Cell Lineage Analysis.** Cell lineage analysis is a vital task in single-cell genomics, aiming to reconstruct developmental trajectories of cells. Traditional methods like *Monocle*[Trapnell et al., 2014] and *Slingshot*[Street et al., 2018] infer pseudotime trajectories but are limited by reliance on predefined metrics and inability to model unobserved progenitor states. Recent advances, such as *LineageVAE*[Majima et al., 2024] and *Waddington-OT*[Schiebinger et al., 2021], address these limitations using probabilistic models and optimal transport, yet often face challenges with the high dimensionality and sparsity of single-cell data. Generative models like *Diffusion Pseudotime Models*[Haghverdi et al., 2016] and *TrajectoryNet*[Tong and Huang, 2020] enhance scalability and interpretability but typically lack mechanisms to model hierarchical relationships in differentiation. Our proposed HDTree integrates hierarchical tree structures into the generative process, enabling accurate reconstruction of both observed and unobserved cell states while improving the interpretability of cell lineage trajectories.

## B Appendix: Details of Dataset

### B.1 Datasets

In this section, we provide an overview of the datasets used in our evaluation. We consider a diverse set of datasets, including image, text, and single-cell data, to assess the performance of hierarchical clustering methods across different domains. The datasets are selected based on their popularity, complexity, and relevance to real-world applications. We provide a brief description of each dataset, along with key statistics and preprocessing steps.

**MNIST:** A widely-used dataset consisting of 70,000 grayscale images of handwritten digits, each of size  $28 \times 28$ . Each image is flattened into a 784-dimensional vector. This dataset is primarily used for benchmarking tree structure modeling and hierarchical clustering methods. More details can be found at<sup>2</sup>.

**Fashion-MNIST:** A dataset of 70,000 grayscale images representing 10 categories of clothing items, each with a resolution of  $28 \times 28$ . The images are flattened into 784-dimensional vectors for analysis. This dataset is used to evaluate the robustness of methods on visual data with more complex patterns than MNIST. Dataset details are available at<sup>3</sup>.

**20News-Groups:** A text dataset with 18,846 newsgroup posts across 20 categories. The features are represented as TF-IDF vectors with dimensionality reduced to 2,000 features for computational

<sup>2</sup><http://yann.lecun.com/exdb/mnist/>

<sup>3</sup><https://github.com/zalandoresearch/fashion-mnist>

**Table 7: Summary of Dataset Statistics.**

Dataset	Type	Class	Samples	Features	Source URL or Reference
MNIST	Image	10	70,000	784	<a href="http://yann.lecun.com/exdb/mnist/">http://yann.lecun.com/exdb/mnist/</a>
Fashion-MNIST	Image	10	70,000	784	<a href="https://github.com/zalandoresearch/fashion-mnist">https://github.com/zalandoresearch/fashion-mnist</a>
20News-Groups	Text	20	18,846	2,000	<a href="http://qwone.com/~jason/20Newsgroups/">http://qwone.com/~jason/20Newsgroups/</a>
CIFAR-10	Image	10	60,000	3,072	<a href="https://www.cs.toronto.edu/~kriz/cifar.html">https://www.cs.toronto.edu/~kriz/cifar.html</a>
Limb	Single-cell	10	66,633	500	<a href="https://limb-dev.cellgeni.sanger.ac.uk/">https://limb-dev.cellgeni.sanger.ac.uk/</a>
LHCO	Single-cell	7	10,628	500	<a href="https://www.ebi.ac.uk/biostudies/arrayexpress/studies/E-MTAB-10973">https://www.ebi.ac.uk/biostudies/arrayexpress/studies/E-MTAB-10973</a>
Weinreb	Single-cell	11	130,887	500	<a href="https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSM4185642">https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSM4185642</a>
ECL	Single-cell	10	838,000	500	<a href="https://cellxgene.cziscience.com/collections/45d5d2c3-bc28-4814-aed6-0bb6f0e11c82">https://cellxgene.cziscience.com/collections/45d5d2c3-bc28-4814-aed6-0bb6f0e11c82</a>

feasibility. This dataset is employed to test methods on high-dimensional text data and hierarchical categorization tasks. More details can be found at<sup>4</sup>.

**CIFAR-10:** A dataset of 60,000 color images spanning 10 classes, with each image sized at  $32 \times 32$ . The pixel intensity values are flattened into a 3,072-dimensional vector. This dataset is utilized for evaluating hierarchical modeling methods on high-dimensional image data. Details are available at<sup>5</sup>.

**Limb:** A single-cell RNA-seq dataset collected from limb bud development experiments. It contains approximately 10,000 cells with 20,000 genes per cell after preprocessing. This dataset is used to study developmental trajectories in biological processes. Dataset details are available in Zhang et al. [2024].

**LHCO:** Single-cell transcriptomics data derived from lung and heart cell ontogeny studies. The features represent gene expression profiles across 15,000 cells with dimensionality reduced to 10,000 genes. This dataset is used to explore differentiation patterns in complex organ systems. Dataset details can be found in He et al. [2022a].

**Weinreb:** A lineage-specific dataset focused on Darwinian lineage inference from single-cell data. It contains 8,000 cells with high-dimensional transcriptomic data preprocessed to 12,000 genes. This dataset is employed for benchmarking methods for lineage tracing tasks. Dataset details are provided in Weinreb et al. [2020a].

**ECL:** Single-cell data from embryonic cell lineages, designed for studying early developmental processes. It comprises 12,000 cells with 15,000 genes per cell after quality filtering. This dataset is used to test hierarchical modeling on datasets with complex lineage structures. Details are available in Qiu et al. [2024a].

The key characteristics of the datasets are summarized in Table 7.

## B.2 Dataset Organization

In our experiment, the all dataset after download from our offered source link should be organized as follows: For Image data, the organization is:

<sup>4</sup><http://qwone.com/~jason/20Newsgroups/>

<sup>5</sup><https://www.cs.toronto.edu/~kriz/cifar.html>

### Image Datasets Directory Structure

```
datasets/
|-- MNIST/
|   |-- train-images-idx3-ubyte
|   |-- train-labels-idx1-ubyte
|   |-- t10k-images-idx3-ubyte
|   '-- t10k-labels-idx1-ubyte
|-- FashionMNIST/
|   |-- train-images-idx3-ubyte
|   |-- train-labels-idx1-ubyte
|   |-- t10k-images-idx3-ubyte
|   '-- t10k-labels-idx1-ubyte
|-- 20news/
|   |-- alt.atheism/
|   |   |-- 12345.txt
|   |   |-- 67890.txt
|   |   '-- ...
|   |-- comp.graphics/
|   |   |-- 12346.txt
|   |   |-- 67891.txt
|   |   '-- ...
|   '-- ...
\-- cifar-10-batches-py/
    |-- data_batch_1
    |-- data_batch_2
    |-- data_batch_3
    |-- data_batch_4
    |-- data_batch_5
    |-- test_batch
    '-- batches.meta
```

For biology data, the organization is:

### Biology Datasets Directory Structure

```
datasets_bio/
|- original/
| |- EpitheliaCell.h5ad
| |- LimbFilter.h5ad
| |- He_2022_NatureMethods_Day15.h5ad
| |- Weinreb_inVitro_clone_matrix.mtx
| |- Weinreb_inVitro_gene_names.txt
| |- Weinreb_inVitro_metadata.txt
| - Weinreb_inVitro_normed_counts.mtx
- processed/ (exists once the process is run)
|- EpitheliaCell_data_n.npy
|- EpitheliaCell_label.npy
|- LimbFilter_data_n.npy
|- LimbFilter_label.npy
|- LHC0.h5ad
- Weinreb.h5ad
```

## B.3 Preprocessing

Before training, datasets need to be preprocessed. Preprocessing steps differ depending on the type of dataset. Below are the detailed guidelines:

### B.3.1 Image Data

For image datasets (e.g., MNIST, FMINST), preprocessing is straightforward and can leverage the code from TreeVAE. The steps include: Initially, downloading and organizing the data ensures that all necessary dataset files, including images and corresponding labels, are retrieved from their sources and systematically placed into the appropriate directories within the project structure. This step is essential for maintaining data integrity and facilitating efficient access during subsequent processing stages. Next, converting raw data into NumPy arrays involves using provided scripts to load the raw image and label data. These scripts parse the binary or structured data formats and convert them into NumPy arrays, which are optimized for numerical computations in Python. This conversion facilitates further data manipulation and model training processes by providing a standardized format for handling large-scale datasets. Finally, normalization and formatting of the pixel values are performed. This typically includes scaling the pixel intensities to a range of [0, 1] to standardize the input features across different images. Additionally, any necessary adjustments are made to ensure the data format aligns with the requirements of the HDTree model, such as ensuring correct dimensionality and data type consistency. Proper normalization enhances model convergence and stability during training.

### B.3.2 Biological Data

For biological datasets (LHCO, Limb, Weinreb, ECL), preprocessing is more complex and tailored to each dataset. Below are the detailed preprocessing steps for each:

**LHCO:** Initially, data cleaning is performed to enhance the quality of the dataset. This includes the removal of duplicate entries and invalid samples that could introduce bias or noise into subsequent analyses. Additionally, strategies for handling missing values are implemented, which may involve imputation techniques or filtering out rows with an excessive amount of missing data. Following data cleaning, feature extraction is conducted to identify and extract relevant features from the raw data. For LHCO datasets, these features often include particle kinematics or event-level characteristics that are critical for analysis. Once extracted, feature scaling is applied to normalize or standardize the data, ensuring consistency across different scales and facilitating more efficient model training. Finally, the dataset is split into training, validation, and test sets according to predefined configurations.

This is the key code of our preprocessed methods:

```
1 adata = sc.read(f"{input_path}/He_2022_NatureMethods_Day15.h5ad")
2 sc.pp.highly_variable_genes(adata, n_top_genes=500)
3 adata = adata[:, adata.var['highly_variable']]
4 data = adata.X
5 data = adata.X.toarray()
6 data = np.array(data).astype(np.float32)
7 mean = data.mean(axis=0)
8 std = data.std(axis=0)
9 data = (data - mean) / std
```

Listing 1: Preprocess of Lhco

**Limb:** It begins with data loading, where the dataset is read from the provided files—typically in formats such as CSV or HDF5—into a structured representation suitable for further processing. Next, the data undergoes filtering and cleaning to enhance its quality and reliability. This includes the removal of noise or artifacts that may have been introduced during data acquisition, as well as deduplication to eliminate redundant entries. Missing values are addressed through appropriate strategies, such as imputation or selective removal of incomplete records. Following this, feature engineering is performed to transform raw biological signals into more interpretable and informative features. For instance, limb motion patterns or other domain-specific characteristics may be extracted to better capture the underlying structure of the data. These features are then normalized to ensure uniformity in scale, which is essential for many machine learning algorithms. Finally, the dataset is stratified and split into training, validation, and test subsets. This is the key code of our preprocessed methods:

```
1 adata = sc.read(f"{input_path}/LimbFilter.h5ad")
2 data_all = adata.X.toarray().astype(np.float32)
3 label_celltype = adata.obs['celltype'].to_list()
```

```

4  vars = np.var(data_all, axis=0) # HVG
5  mask_gene = np.argsort(vars)[-500:]
6  data_hvg = data_all[:, mask_gene]
7  label_count = {}
8  for i in list(set(label_celltype)):
9      label_count[i] = label_celltype.count(i)
10 label_count = sorted(label_count.items(), key=lambda x: x[1], reverse=True)
11 label_count = label_count[:10]
12 mask_top10 = np.zeros(len(label_celltype)).astype(np.bool_)
13 for str_label in label_count:
14     mask_top10[str_label[0]] == np.array(label_celltype)] = 1
15 data_n = np.array(data_hvg).astype(np.float32)[mask_top10]
16 mean = data_n.mean(axis=0)
17 std = data_n.std(axis=0)
18 data = (data_n - mean) / std

```

Listing 2: Preprocess of Limb

**Weinreb:** Initially, data transformation is performed to normalize the raw gene expression counts. This typically includes log-transformation or conversion into normalized values such as Counts Per Million (CPM), Transcripts Per Million (TPM), or Fragments Per Kilobase of transcript per Million mapped reads (FPKM). Low-expression genes, as well as cells with insufficient sequencing depth or missing data, are filtered out to improve signal-to-noise ratio and computational efficiency. Following normalization, dimensionality reduction techniques—such as PCA—reduce the feature space while preserving the major sources of variation in the data, which can improve model performance and reduce computational burden. When the dataset originates from multiple experimental batches or sources, batch effect correction is employed to mitigate technical variability that could confound biological signal detection. Various statistical or machine learning-based approaches may be used depending on the nature of the data and experimental design. Finally, the dataset is stratified and partitioned into training, validation, and test sets.

This is the key code of our preprocessed methods:

```

1 matrix_file = f'{input_path}Weinreb_inVitro_normed_counts.mtx'
2 genes_file = f'{input_path}Weinreb_inVitro_gene_names.txt'
3 metadata_file = f'{input_path}Weinreb_inVitro_metadata.txt'
4 mtx = mmread(matrix_file).tocsr()
5 genes = pd.read_csv(genes_file, header=None, names=['genes'])
6 adata = sc.AnnData(mtx, var=genes)
7 metadata = pd.read_csv(metadata_file, sep='\t')
8 adata.obs = metadata.set_index(adata.obs.index)
9 adata.write(f'{output_path}Weinreb.h5ad')
10 sc.pp.log1p(adata)
11 adata.obs['celltype']=adata.obs['Cell type annotation']
12 adata = adata[~adata.obs['celltype'].isna()]
13 sc.pp.highly_variable_genes(adata, n_top_genes=500)
14 adata = adata[:, adata.var['highly_variable']]
15 data = adata.X.toarray()
16 data = np.array(data).astype(np.float32)
17 mean = data.mean(axis=0)
18 std = data.std(axis=0)
19 data = (data - mean) / std

```

Listing 3: Preprocess of Weinreb

**ECL:** Initially, raw data files are parsed and converted into structured tabular formats that facilitate further computational processing. This is followed by a preprocessing stage that includes normalization—commonly achieved through z-score transformation or Min-Max scaling—and the handling of missing values, which may involve either imputation techniques or the removal of incomplete samples. Subsequently, feature selection is performed with an emphasis on retaining biologically meaningful attributes, often guided by domain-specific knowledge such as known molecular signatures. Finally, the dataset is partitioned into training, validation, and test subsets, ensuring that class distributions are preserved across splits to support unbiased model evaluation and generalization.

This is the key code of our preprocessed methods:

```

1 adata = sc.read(f"{input_path}/EpitheliaCell.h5ad")
2 adata.obs['celltype']=adata.obs['cell_type']
3 label_celltype = adata.obs['celltype'].to_list()
4 adata_sub = adata.copy()
5 sc.pp.subsample(adata_sub, fraction=0.1)
6 data_all = adata_sub.X.toarray().astype(np.float32)
7 vars = np.var(data_all, axis=0)
8 mask_gene = np.argsort(vars)[-500:]
9 adata = adata[:, mask_gene]
10 data = adata.X.toarray().astype(np.float32)
11 label_count = {}
12 for i in list(set(label_celltype)):
13     label_count[i] = label_celltype.count(i)
14 label_count = sorted(label_count.items(), key=lambda x: x[1], reverse=True)
15 label_count = label_count[:10]
16 mask_top10 = np.zeros(len(label_celltype)).astype(np.bool_)
17 for str_label in label_count:
18     mask_top10[str_label[0]] == np.array(label_celltype) = 1
19 data_n = np.array(data).astype(np.float32)[mask_top10]
20 label_train_str = np.array(list(np.squeeze(label_celltype))[
21     mask_top10])
# downsample the 10k data for every cell type
22 mask = np.zeros(len(label_train_str)).astype(np.bool_)
23 for i in range(10):
24     # random select 10k data for each cell type
25     random_index = np.random.choice(
26         np.where(label_train_str == label_count[i][0])[0],
27         10000, replace=False)
28     mask[random_index] = 1
29 data_n = data_n[mask]
30 mean = data_n.mean(axis=0)
31 std = data_n.std(axis=0)
32 data= (data_n - mean) / std

```

Listing 4: Preprocess of ECL

## C Appendix: Details of Baseline Methods

We used the TreeVAE, CellPLM, LangCell, Geneformer as the reference methods. The detail use should follow:

1. **TreeVAE** You should clone the project "treevae" from <https://github.com/lauramanduchi/treevae.git>, and install necessary package followed by "minimal\_requirements.txt".
2. **CellPLM** You should clone the project "CellPLM" from <https://github.com/OmicsML/CellPLM.git>, and install necessary package followed by "requirements.txt". The use of CellPLM is shown in the official tutorial in [https://github.com/OmicsML/CellPLM/blob/main/tutorials/cell\\_embedding.ipynb](https://github.com/OmicsML/CellPLM/blob/main/tutorials/cell_embedding.ipynb)
3. **LangCell** You shold clone the project "LangCell" from [gitclonehttps://github.com/PharMolix/LangCell.git](https://github.com/PharMolix/LangCell.git), and install necessary package followed by "requirements.txt". Then you should install the "geneformer\_001". The use of LangCell is shown in the official tutorial in <https://github.com/PharMolix/LangCell/blob/main/LangCell-annotation-zeroshot/zero-shot.ipynb>
4. **Geneformer** After you install LangCell, the Geneformer package has been installed. The use of Geneformer is shown in the official tutorial in [https://github.com/jkobject/geneformer/blob/main/examples/extract\\_and\\_plot\\_cell\\_embeddings.ipynb](https://github.com/jkobject/geneformer/blob/main/examples/extract_and_plot_cell_embeddings.ipynb)

## D Appendix: Details of Testing Protocol

To ensure comprehensive evaluation, the following metrics are used, each assessing specific aspects of clustering performance, tree structure performance, and reconstruction performance. Here, we provide detailed explanations of their computation methods.

**Clustering Accuracy (ACC):** Clustering accuracy is calculated by finding the best possible one-to-one mapping between predicted cluster labels and ground truth labels. Let  $y_i$  denote the ground truth label for the  $i$ -th sample, and  $\hat{y}_i$  denote the predicted label. ACC is defined as:

$$\text{ACC} = \frac{1}{n} \max_{\pi \in \text{Permutations}} \sum_{i=1}^n \mathbf{1}\{y_i = \pi(\hat{y}_i)\},$$

where  $n$  is the total number of samples,  $\pi$  is a permutation mapping of predicted labels to ground truth labels, and  $\mathbf{1}\{\cdot\}$  is the indicator function.

**Normalized Mutual Information (NMI):** NMI measures the amount of shared information between predicted clusters and ground truth labels. Given the predicted clustering  $C$  and ground truth clustering  $G$ , NMI is defined as:

$$\text{NMI}(C, G) = \frac{2 \cdot I(C; G)}{H(C) + H(G)},$$

where  $I(C; G)$  is the mutual information between  $C$  and  $G$ , and  $H(C)$  and  $H(G)$  are the entropies of  $C$  and  $G$ , respectively. NMI ranges from 0 to 1, with higher values indicating better clustering quality.

**Leaf Purity (LP):** Leaf Purity assesses the homogeneity of data within the leaf nodes of a tree. For a tree  $T$  with leaf nodes  $L_1, L_2, \dots, L_k$ , and each  $L_i$  containing samples with labels  $y_i$ , LP is computed as:

$$\text{LP} = \frac{1}{k} \sum_{i=1}^k \frac{\max_y |L_i^y|}{|L_i|},$$

where  $L_i^y$  represents the subset of samples in  $L_i$  with label  $y$ , and  $|L_i|$  is the total number of samples in  $L_i$ .

**Depth Preservation (DP):** Depth Preservation measures how well the hierarchical structure of the data is maintained in the generated tree. Let  $d_{\text{true}}(x_i, x_j)$  denote the true depth-based distance between two data points  $x_i$  and  $x_j$ , and  $d_{\text{tree}}(x_i, x_j)$  denote their distance in the tree. DP is calculated as:

$$\text{DP} = 1 - \frac{\sum_{i,j} |d_{\text{true}}(x_i, x_j) - d_{\text{tree}}(x_i, x_j)|}{\sum_{i,j} d_{\text{true}}(x_i, x_j)}.$$

This metric ranges from 0 to 1, with higher values indicating better depth preservation.

**Reconstruction Loss (RL):** Reconstruction Loss quantifies the error between the input data  $X$  and the reconstructed data  $\hat{X}$ . It is typically computed as the Mean Squared Error (MSE):

$$\text{RL} = \frac{1}{n} \sum_{i=1}^n \|x_i - \hat{x}_i\|^2,$$

where  $x_i$  is the original data point, and  $\hat{x}_i$  is its reconstruction. To ensure comparability with other metrics (higher is better), the negative value of RL is reported.

**Log-Likelihood (LL):** Log-Likelihood evaluates how well the model fits the observed data. For a probabilistic model with parameters  $\theta$ , the log-likelihood of data  $X$  is:

$$\text{LL} = \sum_{i=1}^n \log p(x_i|\theta),$$

where  $p(x_i|\theta)$  is the probability of observing  $x_i$  given the model parameters. Higher LL values indicate better model fit.

**Table 8: Time Efficiency Comparison Across Methods and Data Sizes:** We conducted an empirical evaluation by selecting varying numbers of highly expressed channels to assess the computational time cost of our model under different data sizes. For this experiment, we set the output dimensionality of the Embedding method to 512 dimensions and employed UMAP for dimensionality reduction to facilitate visualization. The resulting features were then subjected to clustering analysis. As the dataset size increased from 10,000 to 100,000 samples, we observed a dramatic increase in the computational time required for both the dimensionality reduction and clustering processes. Notably, the clustering process escalated from a response time measured in minutes (ranging from 30 to 41 seconds) to one that took several hours (ranging from 2201 to 20429 seconds). Such prolonged processing times are impractical for model evaluation purposes. Consequently, we opted to downsample the data to 10,000 samples for further processing.

Metod	PCA	TSNE	UMAP	CellPLM	LangCell	GeneFormer
<b>Cin=100, C=512, N=10000</b>						
Embedding	None	None	820s	6s	13s	27s
DR	None	None	120s	129s	133s	150s
Cluster	None	None	41s	37s	30s	25s
<b>Cin=1000, C=512, N=10000</b>						
Embedding	140s	None	800s	7s	13s	25s
Dimentional Reduction	140s	None	120s	130s	130s	105s
Cluster	45s	None	40s	37s	34s	30s
<b>Cin=1000, C=512, N=100000</b>						
Embedding	224s	None	3153s	13s	11s	220s
Dimentional Reduction	413s	None	825s	416s	430s	420s
Cluster	9637s	None	7083s	5037s	2032s	3285s
<b>Using K-means to 5000 clusters</b>						
Embedding	221s	None	1834s	14s	123s	110s
Dimentional Reduction	458s	None	411s	350s	427s	410s
Cluster	2925s	None	20429s	5500s	3432s	2200s

## E Appendix: Details of Implementation

For all experiments, the data is split into training, validation, and testing sets with an 8:1:1 ratio, ensuring unbiased evaluation. In testing, if the number of points in the dataset is greater than 10,000, we randomly sample 10,000 points from testing dataset. Details on downsampling and its rationale are provided in the Appendix. We implemented HDTTree using PyTorch and trained the model on a single NVIDIA A100 GPU. The model is trained using the AdamW optimizer with a learning rate of 1e-4 and a batch size of 128. The number of diffusion steps  $T$  is set to 1000, and the tree depth  $L$  is set to 10. The loss weights  $\lambda_{\text{tree}}$  and  $\lambda_{\text{vq}}$  were set to 1.0 and 0.25, respectively. The encoder and diffusion model are implemented using the multi-multilayer perceptron (MLP).

## F Appendix: Details of Downsampling in Testing

As shown in Table 8, the computational time required for encoding, dimensionality reduction (LowDim), and clustering significantly increases with larger data sizes. For instance, when the input size increases to N=100,000, the encoding time for UMAP rises from 820s to 3153s, and the clustering time for GeneFormer increases from 25s to 3285s. These trends are consistent across all tested methods.

Such exponential growth in computational overhead makes the evaluation process infeasible for large-scale datasets. To address this challenge, we uniformly downsample the data to 10,000 points for all methods during metric computation. This ensures consistent and fair comparisons while significantly improving testing efficiency. Additionally, the downsampling procedure preserves the overall data distribution and class proportions, ensuring that the evaluation results remain representative of the original dataset.

The detailed justification for selecting 10,000 points as the downsampling target is discussed in the Appendix.

## G Appendix: Details of Experimental Environment

The experiments were conducted on a high-performance computing system with robust hardware and software configurations to ensure efficient handling of large datasets and complex computations. The hardware setup included NVIDIA A100 GPUs with 40GB memory for accelerated computation, Intel Xeon Platinum 8260 CPUs with 24 cores operating at 2.40GHz for efficient multi-threaded processing, 512GB of RAM, and 10TB of NVMe SSD storage for fast input/output operations.

The software environment was carefully configured for compatibility and reproducibility. The operating system used was Ubuntu 20.04 LTS (64-bit), and the primary deep learning framework was PyTorch (version 2.4.1), supplemented by TorchVision and Torchaudio extensions. The experiments were conducted using Python 3.11.6, with Conda (version 24.9.2) employed as the package manager to handle dependencies and virtual environments.

Key Python packages essential for the experiments included NumPy (1.26.4), SciPy (1.14.1), and pandas (2.2.3) for data preprocessing and analysis. Visualization tasks were performed using Matplotlib (3.9.2), Seaborn (0.13.2), and Plotly (5.24.1). For deep learning tasks, PyTorch (2.4.1) served as the primary frameworks. Dimensionality reduction and clustering were supported by UMAP-learn (0.5.6) and scikit-learn (1.5.2). Single-cell data analysis relied on specialized packages like Scanpy (1.10.3) and anndata (0.10.9).

A comprehensive list of installed Python packages is available upon request, capturing all dependencies required for reproducing the experiments. This configuration ensures the reported results are reproducible and highlights the environment's compatibility with the experimental setup.

## H Appendix: Details of Experiments on Single Cell

### H.1 Detailed Comparisons on More Single Cell Datasets.

We selected four single-cell datasets: LHCO, Limb, Weinreb, and ECL. For LHCO, Limb, and Weinreb, we used the full data. For ECL, to reduce computational cost, we retained only the top 10 cell types by population size, each downsampled to 10,000 randomly selected cells.

We evaluated both traditional (t-SNE, UMAP) and deep learning-based methods (Geneformer, LangCell, CellPLM, TreeVAE). For traditional models, we first selected 500 highly variable genes (HVGs), applied log<sub>10</sub>p transformation where needed, and normalized the data using Z-score. For deep learning-based models, the preprocessing was handled automatically by their built-in pipelines. Thus, for fairness, we directly read the h5ad files, selected HVGs, and fed the results into the respective model pipelines to obtain cell embeddings.

After obtaining cell embeddings from all models, we performed unsupervised clustering using K-Means. Cluster assignments were compared with true cell type labels to compute clustering and tree-structure-based performance metrics. Due to the lack of established benchmarks in the single cell domain, we incorporated data from high-impact published studies into the TreeVAE benchmark. The results are shown in Table 10.

**Analysis:** (1) Similar to the general dataset, HDTTree consistently demonstrates superior performance across all evaluated metrics, including tree structure quality, clustering accuracy, and hierarchical integrity. (2) We observed that the zero-shot capabilities of single-cell large language models are often unsatisfactory and, in some cases, fail to surpass basic single-cell methods. This conclusion has also been validated in recent studies [Lan et al., 2024, He et al., 2024]. In comparison with foundational single-cell models, traditional single-cell tree analysis methods, and TreeVAE, HDTTree shows relative advantages in performance and achieves better stability. (3) These results establish HDTTree as a robust and reliable approach for single-cell data analysis.

### H.2 Stability of Model Performance Across Varying Training Data Sizes

To verify that downsampling did not significantly affect model performance, we conducted additional experiments on subsampled versions of the ECL dataset at varying scales. The results are shown in Table 9

Table 9: **Compare on different sample size**: HDTREE on different sample size of ECL. ECL-N k means we use ECL dataset and downsample it into N\*1000 data point for training, while keeping the same test dataset.

	ECL-10k	ECL-20k	ECL-30k	ECL-50k	ECL-100k	ECL-200k	ECL-300k
ACC	82.2 $\pm$ 1.5	82.2 $\pm$ 1.5	82.5 $\pm$ 0.3	82.6 $\pm$ 0.5	82.9 $\pm$ 0.3	83.1 $\pm$ 0.4	83.1 $\pm$ 0.3
DP	67.0 $\pm$ 1.2	67.6 $\pm$ 1.5	68.3 $\pm$ 0.9	69.1 $\pm$ 0.7	68.5 $\pm$ 0.9	68.9 $\pm$ 0.8	67.0 $\pm$ 0.8

Table 10: **Comparison of tree performance, clustering performance on four single cell datasets**. Since most of the methods are not generative models, we did not compare generative performance.

Dataset	Method	Year	Tree Performance		Clustering Performance		Average( $\uparrow$ )
			DP( $\uparrow$ )	LP( $\uparrow$ )	ACC( $\uparrow$ )	NMI( $\uparrow$ )	
Limb (cell lineage, 66,633 points, celltype:10)	tSNE+Agg <sup>A</sup>	2014	34.9 $\pm$ 1.4	55.4 $\pm$ 1.1	48.9 $\pm$ 0.7	47.5 $\pm$ 1.0	46.7 $\pm$ 1.0
	UMAP+Agg <sup>A</sup>	2018	30.9 $\pm$ 2.0	50.1 $\pm$ 1.0	49.1 $\pm$ 1.0	41.4 $\pm$ 1.4	42.9 $\pm$ 1.4
	Geneformer <sup>A</sup>	2023	25.6 $\pm$ 5.4	35.9 $\pm$ 0.1	34.1 $\pm$ 0.1	34.9 $\pm$ 0.1	32.6 $\pm$ 1.4
	CellPLM <sup>A</sup>	2024	25.6 $\pm$ 0.1	39.9 $\pm$ 0.1	34.1 $\pm$ 0.2	32.9 $\pm$ 0.2	33.1 $\pm$ 0.2
	LangCell <sup>A</sup>	2024	25.3 $\pm$ 0.1	37.5 $\pm$ 0.1	33.9 $\pm$ 0.1	35.1 $\pm$ 0.1	33.0 $\pm$ 0.1
	TreeVAE <sup>A</sup>	2024	34.7 $\pm$ 1.7	55.6 $\pm$ 1.0	49.8 $\pm$ 0.1	50.0 $\pm$ 0.0	47.5 $\pm$ 0.7
	HDTREE <sup>A</sup>	Ours	<b>38.9<math>\pm</math>1.3</b>	<b>57.9<math>\pm</math>1.0</b>	<b>52.8<math>\pm</math>1.0</b>	<b>49.0<math>\pm</math>0.1</b>	<b>49.7<math>\pm</math>0.9</b>
LHCO (cell lineage, 10,628 points, celltype:7)	HDTREE	Ours	<b>41.0<math>\pm</math>0.4</b>	<b>57.2<math>\pm</math>1.4</b>	<b>55.0<math>\pm</math>1.4</b>	<b>46.6<math>\pm</math>0.4</b>	<b>50.0<math>\pm</math>0.9 (<math>\uparrow</math>2.5)</b>
	tSNE+Agg <sup>A</sup>	2014	37.4 $\pm$ 1.6	52.8 $\pm$ 0.8	43.6 $\pm$ 1.0	29.8 $\pm$ 0.5	40.9 $\pm$ 1.0
	UMAP+Agg <sup>A</sup>	2018	40.0 $\pm$ 1.4	50.6 $\pm$ 0.2	46.2 $\pm$ 0.2	34.2 $\pm$ 0.2	43.0 $\pm$ 0.5
	CellPLM <sup>A</sup>	2024	27.0 $\pm$ 1.1	35.8 $\pm$ 2.7	16.8 $\pm$ 3.4	1.65 $\pm$ 5.2	20.3 $\pm$ 3.1
	LangCell <sup>A</sup>	2024	26.5 $\pm$ 1.2	35.2 $\pm$ 0.8	35.2 $\pm$ 0.6	0.02 $\pm$ 0.9	24.2 $\pm$ 0.9
	TreeVAE <sup>A</sup>	2024	38.3 $\pm$ 2.0	52.2 $\pm$ 0.1	37.9 $\pm$ 0.1	31.6 $\pm$ 0.0	40.0 $\pm$ 0.6
	HDTREE <sup>A</sup>	Ours	<b>38.8<math>\pm</math>0.3</b>	<b>52.1<math>\pm</math>0.4</b>	<b>46.4<math>\pm</math>0.3</b>	<b>34.7<math>\pm</math>0.5</b>	<b>43.0<math>\pm</math>0.4</b>
Weinreb (cell lineage, 130,887 points, celltype:11)	HDTREE	Ours	<b>42.7<math>\pm</math>0.4</b>	<b>54.0<math>\pm</math>0.3</b>	<b>49.4<math>\pm</math>0.3</b>	<b>34.5<math>\pm</math>0.4</b>	<b>45.2<math>\pm</math>0.3 (<math>\uparrow</math>2.2)</b>
	tSNE+Agg <sup>A</sup>	2014	57.9 $\pm$ 1.5	63.3 $\pm$ 1.1	35.3 $\pm$ 1.0	38.5 $\pm$ 0.6	48.8 $\pm$ 1.1
	UMAP+Agg <sup>A</sup>	2018	51.8 $\pm$ 0.1	62.1 $\pm$ 0.8	47.2 $\pm$ 4.9	46.1 $\pm$ 1.2	51.8 $\pm$ 1.8
	LangCell <sup>A</sup>	2024	47.4 $\pm$ 0.1	54.8 $\pm$ 0.0	14.3 $\pm$ 0.5	34.3 $\pm$ 0.0	37.7 $\pm$ 0.2
	Geneformer <sup>A</sup>	2024	45.1 $\pm$ 0.4	55.3 $\pm$ 0.1	21.4 $\pm$ 0.1	32.3 $\pm$ 0.1	38.5 $\pm$ 0.2
	TreeVAE <sup>A</sup>	2024	60.4 $\pm$ 2.6	61.4 $\pm$ 0.5	41.0 $\pm$ 0.1	35.2 $\pm$ 0.0	49.5 $\pm$ 0.8
	HDTREE <sup>A</sup>	Ours	<b>63.3<math>\pm</math>2.6</b>	<b>78.2<math>\pm</math>1.1</b>	<b>50.6<math>\pm</math>1.0</b>	<b>45.2<math>\pm</math>1.2</b>	<b>59.3<math>\pm</math>1.5 (<math>\uparrow</math>7.5)</b>
ECL (cell lineage, 100K points, celltype:10)	HDTREE	Ours	<b>61.0<math>\pm</math>0.4</b>	<b>67.0<math>\pm</math>0.3</b>	<b>62.6<math>\pm</math>0.3</b>	<b>42.6<math>\pm</math>0.3</b>	<b>58.3<math>\pm</math>0.4</b>
	tSNE+Agg <sup>A</sup>	2014	55.5 $\pm$ 5.4	73.7 $\pm$ 4.2	73.1 $\pm$ 5.4	70.9 $\pm$ 3.8	68.3 $\pm$ 4.7
	UMAP+Agg <sup>A</sup>	2018	53.2 $\pm$ 1.4	73.6 $\pm$ 1.0	71.2 $\pm$ 1.5	71.8 $\pm$ 0.8	67.4 $\pm$
	TreeVAE <sup>A</sup>	2024	41.86 $\pm$ 1.9	60.7 $\pm$ 2.0	57.0 $\pm$ 3.0	61.8 $\pm$ 1.8	55.3 $\pm$ 2.2
	HDTREE <sup>A</sup>	Ours	<b>60.1<math>\pm</math>0.1</b>	<b>74.7<math>\pm</math>0.5</b>	<b>70.9<math>\pm</math>0.4</b>	<b>78.9<math>\pm</math>0.4</b>	<b>71.2<math>\pm</math>0.4</b>
	HDTREE	Ours	<b>69.0<math>\pm</math>0.7</b>	<b>83.2<math>\pm</math>0.3</b>	<b>83.2<math>\pm</math>0.3</b>	<b>79.0<math>\pm</math>0.3</b>	<b>78.6<math>\pm</math>0.4 (<math>\uparrow</math>10.3)</b>

The experimental results showed that the model performance did not change significantly with variations in the size of the training data. This indicates that the model's performance remained relatively stable regardless of the dataset size, without notable improvements or declines. This finding suggests that the current amount of data is sufficient for the model to learn robust feature representations, or that data quantity is not the key factor influencing model performance in this task.