

# UC Irvine

## UC Irvine Previously Published Works

### Title

Impeller: a path-based heterogeneous graph learning method for spatial transcriptomic data imputation.

### Permalink

<https://escholarship.org/uc/item/3zp1k3q4>

### Journal

Computer applications in the biosciences : CABIOS, 40(6)

### Authors

Duan, Ziheng

Riffle, Dylan

Li, Ren

et al.

### Publication Date

2024-06-03

### DOI

10.1093/bioinformatics/btae339

Peer reviewed

## Gene expression

# Impeller: a path-based heterogeneous graph learning method for spatial transcriptomic data imputation

Ziheng Duan <sup>1</sup>, Dylan Riffle<sup>1</sup>, Ren Li<sup>2</sup>, Junhao Liu<sup>1</sup>, Martin Renqiang Min<sup>3</sup>, Jing Zhang <sup>1,\*</sup>

<sup>1</sup>Department of Computer Science, University of California, Irvine, Irvine, CA 92697, United States

<sup>2</sup>Mathematical, Computational, and Systems Biology, University of California, Irvine, Irvine, CA 92697, United States

<sup>3</sup>Department of Machine Learning, NEC Labs America, Princeton, NJ 08540, United States

\*Corresponding author. Department of Computer Science, University of California, 4084 Bren Hall University of California, Irvine (UCI) Irvine, CA 92697-3435, United States. E-mail: zhang.jing@uci.edu (J.Z.)

Associate Editor: Anthony Mathelier

## Abstract

**Motivation:** Recent advances in spatial transcriptomics allow spatially resolved gene expression measurements with cellular or even sub-cellular resolution, directly characterizing the complex spatiotemporal gene expression landscape and cell-to-cell interactions in their native microenvironments. Due to technology limitations, most spatial transcriptomic technologies still yield incomplete expression measurements with excessive missing values. Therefore, gene imputation is critical to filling in missing data, enhancing resolution, and improving overall interpretability. However, existing methods either require additional matched single-cell RNA-seq data, which is rarely available, or ignore spatial proximity or expression similarity information.

**Results:** To address these issues, we introduce Impeller, a path-based heterogeneous graph learning method for spatial transcriptomic data imputation. Impeller has two unique characteristics distinct from existing approaches. First, it builds a heterogeneous graph with two types of edges representing spatial proximity and expression similarity. Therefore, Impeller can simultaneously model smooth gene expression changes across spatial dimensions and capture similar gene expression signatures of faraway cells from the same type. Moreover, Impeller incorporates both short- and long-range cell-to-cell interactions (e.g. via paracrine and endocrine) by stacking multiple GNN layers. We use a learnable path operator in Impeller to avoid the over-smoothing issue of the traditional Laplacian matrices. Extensive experiments on diverse datasets from three popular platforms and two species demonstrate the superiority of Impeller over various state-of-the-art imputation methods.

**Availability and implementation:** The code and preprocessed data used in this study are available at <https://github.com/aicb-ZhangLabs/Impeller> and <https://zenodo.org/records/11212604>.

## 1 Introduction

The orchestration of cellular life hinges on the precise control of when and where genes are activated or silenced. Characterizing such spatiotemporal gene expression patterns is crucial for a better understanding of life, from development to disease to adaptation (Mantri *et al.* 2021). While single-cell RNA sequencing (scRNA-seq) is a revolutionary and widely available technology that enables simultaneous gene expression profiling over thousands of cells, it usually needs to dissociate cells from their native tissue and thus loses the spatial context (Lähnemann *et al.* 2020). Recent advances in spatial transcriptomics (Stahl *et al.* 2016) allow spatially resolved gene expression measurements at a single-cell or even sub-cellular resolution, providing unprecedented opportunities to characterize the complex landscape of spatiotemporal gene expression and understand the intricate interplay between cells in their native microenvironments (Strell *et al.* 2019). However, due to technical and biological limitations, most spatial transcriptomic profiling technologies still yield incomplete datasets with excessive missing gene expression values, hindering our biological interpretation of such valuable datasets (Choe *et al.* 2023). Therefore, gene imputation

is a critical task to enrich spatial transcriptomics by filling in missing data, enhancing resolution, and improving the overall quality and interpretability of the datasets.

Several methods have been successfully developed for gene imputation in spatial transcriptomics, which can be broadly summarized into two categories—reference-based and reference-free approaches. Since scRNA-seq data usually offer a deeper dive into transcriptome profiling, reference-based methods integrated spatial transcriptomic data with matched scRNA-seq data from the same sample for accurate imputation. While promising, these referenced-based methods usually suffer from two limitations. First, most studies do not always have matched scRNA-seq data, especially those using valuable and rare samples. Second, even with matched data, there can be significant gene expression distribution shifts due to sequencing protocol differences (e.g. single nuclei RNA-seq versus whole cell spatial transcriptomics) (Zeng *et al.* 2022).

Researchers also used reference-free methods for direct gene expression imputation. For instance, traditional gene imputation methods designed for scRNA-seq data, such as scVI (Lopez *et al.* 2018), ALRA (Linderman *et al.* 2018),

Magic (van Dijk *et al.* 2018), and scGNN (Wang *et al.* 2021), have been adapted for spatial transcriptomic data imputation. While effectively capturing cell-type-specific gene expression signatures, these methods completely ignored the rich spatial information, resulting in suboptimal results. Later, scientists emphasized the importance of spatial context for cell-to-cell interaction (CCI) in modulating expression changes in response to external stimuli (Armingol *et al.* 2021). Therefore, Graph Neural Network (GNN) based methods have been developed to mimic CCIs for imputation tasks with improved performance. However, different types of CCI involve distinct cell signaling mechanisms with varying interaction ranges. Existing GNN-based methods used very shallow convolutional layers for computational convenience, successfully modeling short-range CCI (e.g. via autocrine and juxtacrine) but ignoring long-range interactions (e.g. via paracrine and endocrine). As a result, they cannot fully exploit the spatial information for gene expression imputation.

To address the abovementioned issues, we propose Impeller, a path-based heterogeneous graph learning method for accurate spatial transcriptomic data imputation. Impeller contains two unique components to exploit both transcriptomic and spatial information. First, it builds a heterogeneous graph with nodes representing cells and two types of edges describing expression similarity and spatial proximity. Therefore, the expression-based edges allow it to capture cell-type-specific expression signatures of faraway cells from the same type, and the proximity-based edges incorporate CCI effects in the spatial context. Second, Impeller models long-range CCI by stacking multiple GNN layers and uses a learnable path operator instead of the traditional Laplacian matrices to avoid the over-smoothing problem. Extensive experiments on diverse datasets from three popular platforms and two species demonstrate the superiority of Impeller over various state-of-the-art imputation methods.

Our main contributions are summarized below:

- We propose a graph neural network, Impeller, for reference-free spatial transcriptomic data imputation. Impeller incorporates cell-type-specific expression signatures and CCI via a heterogeneous graph with edges representing transcriptomic similarity and spatial proximity.
- Impeller stacks multiple GNN layers to include both short- and long-range cell-to-cell interactions in the spatial context. Moreover, it uses a learnable path-based operator to avoid over-smoothing.
- To the best of our knowledge, this is the first paper to combine cell-type-specific expression signatures with spatial short- and long-range CCI for gene expression imputation.
- We extensively evaluate Impeller alongside state-of-the-art competitive methods on datasets from three sequencing platforms and two species. The results demonstrate that Impeller outperforms all of the baselines.

## 2 Related work

### 2.1 Imputation methods ignoring spatial information

Earlier spatial transcriptomic data imputation methods adapted the computational strategies originally developed for scRNA-seq data, overlooking the spatial coordinate

information of each spot. For instance, eKNN (expression-based K nearest neighbor), and eSNN (expression-based Shared nearest neighbor) are methods implemented using the Seurat R-package that rely on gene expressions of nearest neighbors. MAGIC adopted data diffusion across similar cells to impute missing transcriptomic data. ALARA used low-rank approximation to distinguish genuine nonexpression from technical dropouts, thus preserving true gene absence in samples. scVI used a deep variational autoencoder for gene imputation by assuming the read counts per gene follow a zero-inflated negative binomial distribution. However, these methods completely ignored the rich spatial information, resulting in sub-optimal performance.

### 2.1 Imputation methods utilizing spatial information

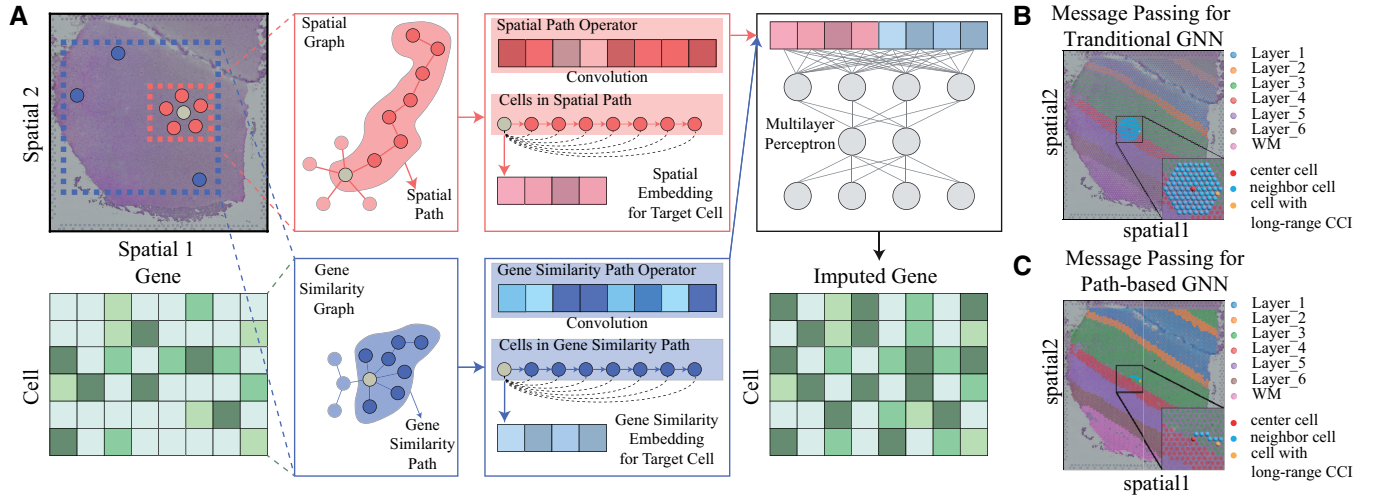
Later on, several methods were developed to exploit the spatial coordinate information to improve imputation accuracy. Since scRNA-seq data are usually sequenced deeper to provide more accurate expression measurements, several methods incorporated additional scRNA-seq data during the imputation process. For instance, gimVI used a low-rank approximation and included scRNA reference (Lopez *et al.* 2019). Tangram mapped scRNA-seq data onto spatial transcriptomics data to facilitate imputation by fitting expression values on the shared genes (Biancalani *et al.* 2021). STLearn used gene expression data, spatial distance, and tissue morphology data for imputing absent gene reads (Pham *et al.* 2020). However, additional scRNA-seq data are not always available and there can be large gene expression distribution shifts between these datasets due to differences in sequencing protocols (e.g. single-cell versus single-nuclei), resulting in limited applications for reference-based methods.

On the other hand, several reference-free methods have been developed for more generalized settings. For example, the seKNN (spatial-expression-based K nearest neighbor) and seSNN (spatial-expression-based shared nearest neighbor) models (Satija *et al.* 2015, Butler *et al.* 2018, Stuart *et al.* 2019, Hao *et al.* 2021) incorporated cell-to-cell distance when defining the KNN for imputation tasks. Recently, STAGATE (Dong and Zhang 2022) is a graph attention auto-encoder framework that effectively imputes genes by integrating spatial data and cell type labels. Overall, these methods did not deeply integrate and exploit the full potential of combining expression and spatial data.

## 3 Materials and methods

### 3.1 Problem definition

Here, we aim to impute the excessive missing gene expression values in spatial transcriptomics data without matched reference scRNA-seq data. Formally, given a sparse cell-by-gene count matrix  $\mathbf{X}_{\text{obs}} \in \mathbb{R}^{n \times m}$  which represents observations for  $n$  cells across  $m$  genes, and the spatial coordinates  $\mathbf{C} \in \mathbb{R}^{n \times 2}$  of these cells, our goal is to impute the gene expression matrix  $\hat{\mathbf{X}} \in \mathbb{R}^{n \times m}$ .  $\mathbf{X}_{\text{obs}}$  is derived from the ground truth matrix  $\mathbf{X}_{\text{gt}} \in \mathbb{R}^{n \times m}$ , which contains the observed nonzero entries pre-masking. To simulate real-world data conditions, 10% of the nonzero entries in  $\mathbf{X}_{\text{gt}}$  are masked to form a test set and another 10% for validation, thus creating  $\mathbf{X}_{\text{obs}}$ . This matrix serves as the input for our imputation model. The major challenge is to generate  $\hat{\mathbf{X}}$  that is as close as possible to the ground truth gene expression  $\mathbf{X}_{\text{gt}}$ , using both the observed gene expressions in  $\mathbf{X}_{\text{obs}}$  and the spatial information in  $\mathbf{C}$ .



**Figure 1.** The overview of Impeller. (A) Given the observed matrix  $\mathbf{X}_{\text{obs}} \in \mathbb{R}^{n \times m}$  of  $n$  cells and  $m$  genes, and the cells’ spatial coordinates  $\mathbf{C} \in \mathbb{R}^{n \times 2}$ , we build the spatial graph  $\mathbf{G}_s$  and the gene similarity graph  $\mathbf{G}_g$ . The learned spatial and gene similarity path operators  $\text{op}_s$  and  $\text{op}_g$  are obtained through  $\text{path}_s$  and  $\text{path}_g$ , respectively. Convoluting cell features with path operators yields spatial/gene similarity embeddings, which are concatenated and fed into a multilayer perceptron for final gene imputation. (B) and (C) Comparison of neighbor aggregation methods in GNNs. (B) Traditional GNN stacks multiple layers to gather information from distant nodes. (C) The path-based GNN, Impeller samples a path to the target node.

### 3.2 Heterogeneous graph construction

As shown in Fig. 1, we build our Impeller model based on two widely accepted biological insights—(i) gene expression can be modulated by surrounding cells via CCI; (ii) faraway cells of the same cell type may share stable gene expression signatures. Therefore, Impeller first builds a heterogeneous graph  $\mathbf{G}$  to fully exploit both spatial and cell-type information, with nodes and edges representing cells and their relationships.

Specifically,  $\mathbf{G}$  contains two complementary graphs: a spatial graph ( $\mathbf{G}_s$ ) and a gene similarity graph ( $\mathbf{G}_g$ ). Edges in  $\mathbf{G}_s$  represent the cell’s spatial proximity to model CCI, while edges in  $\mathbf{G}_g$  denote the cell’s transcriptomic similarity to capture the cell-type-specific expression signatures.

#### 3.2.1 Spatial graph construction

The spatial graph  $\mathbf{G}_s(\mathbf{V}_s, \mathbf{E}_s)$  is created based on the spatial distance between cells, with nodes  $\mathbf{V}_s$  representing the cells and edges in  $\mathbf{E}_s$  connecting nearby cells. Specifically, an edge  $e_{s,\{ij\}}$  in  $\mathbf{G}_s$  is established between  $v_i, v_j \in \mathbf{V}_s$  if and only if their Euclidean distance  $d_{ij}$  is less than a predefined threshold  $d_{\text{thr}}$ , which can be represented as:

$$e_{s,\{ij\}} = \text{if } \|\mathbf{C}_i - \mathbf{C}_j\|_2 \leq d_{\text{thr}} \text{ else } 0, \quad (1)$$

where  $\mathbf{C}_i = [C_{i,0}, C_{i,1}]$  and  $\mathbf{C}_j = [C_{j,0}, C_{j,1}]$  are 2D spatial coordinates of cell  $i$  and  $j$ , respectively.

#### 3.2.2 Gene similarity graph construction

Impeller also builds a gene expression similarity graph  $\mathbf{G}_g$  similar to that in scRNA-seq analysis. Specifically, we first extract the highly variable genes (default 3100). Then, for each target cell, we select its top  $K$  most similar cells. Mathematically,

$$e_{g,\{ij\}} = 1 \text{ if } j \in \mathbf{K}_g(\mathbf{X}_i^h) \text{ else } 0, \quad (2)$$

where  $\mathbf{X}_i^h$  is the expression vector of highly variable genes in cell  $i$ ,  $\mathbf{K}_g(\mathbf{X}_i^h)$  returns the top  $k_g$  cells most similar to cell  $i$

(e.g. using the Euclidean distance as the similarity metric), and  $e_{g,\{ij\}}$  is the edge between cells  $i$  and  $j$  in  $\mathbf{G}_g$ .

### 3.3 GNN model on heterogeneous graph

With the heterogeneous graph built, Impeller uses a path-based heterogeneous GNN to synthesize the impacts of spatial CCI ( $\mathbf{G}_s$ ) and cell-type-specific expression signatures ( $\mathbf{G}_g$ ) for the imputation task. We introduce the problem of traditional GNN, our learnable path operator, and the overall architecture of Impeller as follows.

#### 3.3.1 Problem of traditional GNN

We aim to impute the missing gene expression values in spatial transcriptomics data by incorporating its physical and transcriptional neighbors via a heterogeneous graph. By treating expression profiles as initial cell embeddings ( $\mathbf{f}^{(0)} = \mathbf{X}_{\text{obs}}$ ), the  $l$ th ( $l \in \{1, 2, \dots, L-1\}$ ) GNN layer follows a message passing form (Duan *et al.* 2022a<sup>b,c</sup>, Wang *et al.* 2019, 2022, Xu *et al.* 2020, 2021, Duan *et al.* 2023, 2024) to generate cell  $i$ ’s embedding in layer  $l$  as follows:

$$\mathbf{f}_i^{(l)} = \gamma_{\theta}(\mathbf{f}_i^{(l-1)}) \oplus_{j \in \mathbf{N}_s(i)} \phi_{\theta}(\mathbf{f}_i^{(l-1)}, \mathbf{f}_j^{(l-1)}, e_{s,\{ij\}}) \oplus_{j \in \mathbf{N}_g(i)} \psi_{\theta}(\mathbf{f}_i^{(l-1)}, \mathbf{f}_j^{(l-1)}, e_{g,\{ij\}}), \quad (3)$$

where  $\mathbf{f}_i^{(l)} \in \mathbb{R}^{d_{\text{emb}}^{(l)}}$  is the embedding of cell  $i$  at  $l$ th layer,  $d_{\text{emb}}^{(l)}$  is the embedding dimension at  $l$ th layer, and  $\mathbf{N}_s(i)$  and  $\mathbf{N}_g(i)$  are neighboring cell  $i$  in  $\mathbf{G}_s$  and  $\mathbf{G}_g$ .  $\oplus$  denotes a differentiable, permutation invariant function, e.g. sum, mean, and  $\gamma_{\theta}$ ,  $\phi_{\theta}$ , and  $\psi_{\theta}$  denote differentiable functions such as MLPs. After  $L$  layers, we obtain the imputed gene expressions, denoted as  $\hat{\mathbf{X}} = \mathbf{f}^{(L)} \in \mathbb{R}^{n \times m}$ .

In order to capture long-range CCI interaction, we have to include relatively far away cells by stacking multiple GNN layers via a larger  $L$ . Traditional Laplacian matrices-based GNN suffers from over-smoothing, resulting in deteriorated performance as  $L$  increases (Eliasof *et al.* 2022). Therefore, we introduce a learnable path operator to overcome this issue and better capture the long-range CCI.

### 3.3.2 Learnable path operator

We first define path  $\mathbf{P}_s = (s_1, s_2, \dots, s_{k_s})$  on  $\mathbf{G}_s$  of length  $k_s$  and path  $\mathbf{P}_g = (g_1, g_2, \dots, g_{k_g})$  on  $\mathbf{G}_g$  of length  $k_g$ , where  $s_i$  and  $g_i$  are node (cell) indexes. Node embeddings at  $l$ th layer are denoted by  $\mathbf{f}_{s_i}^{(l)} \in \mathbb{R}^{d_{emb}^{(l)}}$  and  $\mathbf{f}_{g_i}^{(l)} \in \mathbb{R}^{d_{emb}^{(l)}}$ . Then,  $\mathbf{op}_s^{(l)} \in \mathbb{R}^{k_s \times d_{emb}^{(l)}}$  and  $\mathbf{op}_g^{(l)} \in \mathbb{R}^{k_g \times d_{emb}^{(l)}}$  are two learnable path operators which allow us to convolve node embeddings along paths:

$$\begin{aligned} \mathbf{op}_s^{(l)}(\mathbf{P}_s) * \mathbf{f}^{(l)} &= \sum_{i=1}^{k_s} \mathbf{op}_{s,i}^{(l)} * \mathbf{f}_{s_i}^{(l)} = \sum_{i=1}^{k_s} \sum_{j=1}^{d_{emb}^{(l)}} \mathit{op}_{s,i}^{(l)}[j] \cdot f_{s_i}^{(l)}[j], \\ \mathbf{op}_g^{(l)}(\mathbf{P}_g) * \mathbf{f}^{(l)} &= \sum_{i=1}^{k_g} \mathbf{op}_{g,i}^{(l)} * \mathbf{f}_{g_i}^{(l)} = \sum_{i=1}^{k_g} \sum_{j=1}^{d_{emb}^{(l)}} \mathit{op}_{g,i}^{(l)}[j] \cdot f_{g_i}^{(l)}[j], \end{aligned} \quad (4)$$

where  $*$  denotes the convolution operation, and  $\cdot$  symbol is the multiplication operation between two scalars. Here  $\mathit{op}_{s,i}^{(l)}[j]$ ,  $\mathit{op}_{g,i}^{(l)}[j]$ ,  $f_{s_i}^{(l)}[j]$  and  $f_{g_i}^{(l)}[j]$  represent the  $j$ th scalars of the  $d_{emb}^{(l)}$ -dimensional vector  $\mathbf{op}_{s,i}^{(l)}$ ,  $\mathbf{op}_{g,i}^{(l)}$ ,  $\mathbf{f}_{s_i}^{(l)}$  and  $\mathbf{f}_{g_i}^{(l)}$ , respectively. Starting from each node, we generate multiple paths on  $\mathbf{G}_s$  and  $\mathbf{G}_g$  and aggregate results for a more expressive representation:

$$\begin{aligned} \mathbf{op}_s^{(l)}(\mathcal{P}_s) * \mathbf{f}^{(l)} &= \frac{1}{T_s} \sum_{\mathbf{P}_s \in \mathcal{P}_s} \mathbf{op}_s^{(l)}(\mathbf{P}_s) * \mathbf{f}^{(l)}, \\ \mathbf{op}_g^{(l)}(\mathcal{P}_g) * \mathbf{f}^{(l)} &= \frac{1}{T_g} \sum_{\mathbf{P}_g \in \mathcal{P}_g} \mathbf{op}_g^{(l)}(\mathbf{P}_g) * \mathbf{f}^{(l)}, \end{aligned} \quad (5)$$

where  $\mathcal{P}_s$  and  $\mathcal{P}_g$  are sets of paths sampled from the  $\mathbf{G}_s$  and  $\mathbf{G}_g$ , each containing  $T_s$  and  $T_g$  paths. Each path  $\mathbf{P}_s \in \mathcal{P}_s$  and  $\mathbf{P}_g \in \mathcal{P}_g$  are separately convolved using  $\mathbf{op}_s^{(l)}$  or  $\mathbf{op}_g^{(l)}$ , and the results are averaged to acquire the node embeddings.

### 3.3.3 The overall architecture of impeller

After convolving both spatial and gene similarity paths, we concatenate their embeddings to form the overall node embeddings, as in

$$\mathbf{f}^{(l+1)} = \sigma(\mathbf{W}_1^{(l)} [\mathbf{op}_s^{(l)}(\mathcal{P}_s) * \mathbf{f}^{(l)}, \mathbf{op}_g^{(l)}(\mathcal{P}_g) * \mathbf{f}^{(l)}]), \quad (6)$$

where  $\sigma(\cdot)$  denotes the ReLU activation function,  $\mathbf{W}_1^{(l)} \in \mathbb{R}^{d_{emb}^{(l+1)} \times 2 \cdot d_{emb}^{(l)}}$  is the learnable weight matrix,  $d_{emb}^{(l)}$  is the embedding dimension at  $l$ th layer, and  $[\cdot, \cdot]$  denotes concatenation operation. Then, Impeller tries to minimize the Mean Squared Error (MSE) between  $\hat{\mathbf{X}}$  and  $\mathbf{X}_{gt}$ :

$$\mathcal{L} = \frac{\sum_{i=1}^n \sum_{j=1}^m \mathbb{1}[X_{gt,(i,j)} \neq 0] (\hat{X}_{i,j} - X_{gt,(i,j)})^2}{\sum_{i=1}^n \sum_{j=1}^m \mathbb{1}[X_{gt,(i,j)} \neq 0]}, \quad (7)$$

where  $\mathbb{1}[\cdot]$  is an indicator function that equals 1 if the condition inside brackets is met ( $X_{gt,(i,j)} \neq 0$ ), and 0 otherwise. The loss is computed only over nonzero entries of  $\mathbf{X}_{gt}$ .

## 3.4 Computational complexity analysis

### 3.4.1 $k$ -hop complexity analysis

Traditional GNNs need to gather information from  $k$ -hop neighbor nodes after stacking of  $k$  layers. Given the complexity of each layer as  $O(n \times d_l)$ , where  $n$  is the number of nodes

and  $d_l$  is the average node degree, the overall complexity becomes  $O(n \times d_l \times k)$ . In contrast, Impeller can directly access neighbors up to  $k$ -hop distance via a single layer by setting  $k_s = k_g = k$ . The computational complexity per layer for Impeller is  $O(n \times (T_s \times k_s + T_g \times k_g))$ , with  $T_s$  and  $T_g$  representing the number of paths in  $\mathbf{G}_s$  and  $\mathbf{G}_g$ ,  $k_s$  and  $k_g$  denoting path lengths. As a result, when  $T_s < d_l$  (a condition satisfied in our task), Impeller offers superior computational efficiency.

### 3.4.2 The number of parameters

Traditional GNNs have  $O(d_{emb}^{(l)} \times d_{emb}^{(l+1)})$  parameters per layer while the path operator of Impeller adds  $(k_s + k_g) \times d_{emb}^{(l)}$  parameters. Since  $k_s + k_g$  is typically much smaller than  $d_{emb}^{(l+1)}$ , Impeller's number of parameters remains on par with traditional GNNs.

## 4 Experiments

### 4.1 Detailed experimental setup

#### 4.1.1 Data sources and preprocessing

In our study, we tested Impeller using diverse datasets from three popular sequencing platforms and two organisms. Specifically, we included 10X Visium datasets from the human dorsolateral prefrontal cortex (DLPFC), (Maynard *et al.* 2021), Steroseq datasets from mouse olfactory bulb (Chen *et al.* 2021), and Slide-seqV2 from mouse olfactory bulb (Stickels *et al.* 2021) in our analyses. Detailed attributes of these datasets are summarized in Table 1 (for filter details and visualizations, see the Supplementary Appendix). After standard pre-processing and normalization procedures, we downsampled the data according to scGNN, where 10% of nonzero entries in the dataset were used as a test set, and another 10% of nonzero entries were reserved for validation. For a fair comparison, we repeat ten times with different mask configurations.

#### 4.1.2 Baseline methods for benchmarking

We conducted a comparative study utilizing 12 state-of-the-art methods, including reference-free and reference-based methods that originally required additional scRNA-seq data. However, in our analysis, we did not use any additional scRNA-seq data for a fair comparison.

First, we included methods directly adapted from scRNA-seq data imputation and completely ignored the rich spatial information, including a deep generative model scVI, a low-rank approximation model ALRA, nearest neighbors-based models eKNN and eSNN, a diffusion-based model MAGIC, and a GNN-based model scGNN. Furthermore, we used several imputation methods specifically designed for spatial transcriptomic data, such as seKNN (spatial-expression-based K nearest neighbor), and seSNN (spatial-expression-based shared nearest neighbor). gimVI and Tangram need additional scRNA-seq from matched samples, so we used a reference-free implementation available through their website for a fair comparison. Lastly, we included STAGATE a graph attention auto-encoder framework by amalgamating spatial data and gene expression profiles. We use default parameters in most baseline methods (for details, see the Supplementary Appendix).

**Table 1.** Summary of datasets.

Platform	Organism	Sample ID	Raw matrix (cell, gene)	Raw density	Filter matrix (cell, gene)	Filter density	No. of imputed entries
10xVisium	Human dorsolateral prefrontal cortex (DLPFC)	151507	4226, 33 538	0.042	4117, 4028	0.261	437 240
		151508	4384, 33 538	0.036	4148, 3342	0.258	358 184
		151509	4789, 33 538	0.043	4700, 4188	0.258	508 186
		151510	4643, 33 538	0.041	4547, 3908	0.259	461 112
		151669	3661, 33 538	0.054	3617, 5246	0.277	525 930
		151670	3498, 33 538	0.050	3433, 4909	0.272	457 770
		151671	4110, 33 538	0.055	3988, 5539	0.278	615 111
		151672	4015, 33 538	0.052	3809, 5273	0.279	561 166
		151673	3639, 33 538	0.066	3628, 6538	0.286	677 473
		151674	3673, 33 538	0.080	3668, 7796	0.305	871 032
		151675	3592, 33 538	0.054	3565, 5454	0.267	518 515
		151676	3460, 33 538	0.058	3449, 5784	0.274	545 920
		Stereoseq	Mouse	/	19 109, 14 376	0.024	4036, 1581
SlideSeqV2	Mouse	/	20 139, 11 750	0.031	5161, 2611	0.217	292 418

### 4.1.3 Evaluation metrics

We first define a test mask  $\mathbf{M} \in \mathbb{R}^{n \times m}$  where the entries to be imputed are marked as 1 and the others as 0. Then we extract the relevant entries from both the imputed matrix  $\hat{\mathbf{X}}$  and the ground truth matrix  $\mathbf{X}_{\text{gt}}$  to form two vectors:  $\hat{\mathbf{x}}$  (from  $\hat{\mathbf{X}}$ ) and  $\mathbf{x}_{\text{gt}}$  (from  $\mathbf{X}_{\text{gt}}$ ), each of length  $N$ , where  $N$  is the total number of entries to be imputed. Following scGNN settings, we use L1 Distance, Cosine Similarity, and Root-Mean-Square Error (RMSE) to compare imputed gene expressions  $\hat{\mathbf{x}}$  with the ground truth  $\mathbf{x}_{\text{gt}}$ . Mathematically:

$$\text{L1 Distance} = |\hat{\mathbf{x}} - \mathbf{x}_{\text{gt}}|, \quad (8)$$

$$\text{Cosine Similarity}(\hat{\mathbf{x}}, \mathbf{x}_{\text{gt}}) = \frac{\hat{\mathbf{x}} \mathbf{x}_{\text{gt}}^T}{\|\hat{\mathbf{x}}\| * \|\mathbf{x}_{\text{gt}}\|}, \quad (9)$$

$$\text{RMSE}(\hat{\mathbf{x}}, \mathbf{x}_{\text{gt}}) = \sqrt{\frac{\sum_{i=1}^N (\hat{x}_i - x_{\text{gt}i})^2}{N}}. \quad (10)$$

## 4.2 Experimental results

### 4.2.1 Improved imputation accuracy

We benchmarked our performance against 12 leading methods by assessing imputation accuracy across 14 datasets. These datasets span three prominent sequencing platforms (10x Visium, Stereoseq, and SlideSeq) and two species (human and mouse). [Table 2](#) summarizes the performance of Impellers and other baselines (for results of the other six samples of the DLPFC dataset, please see the [Supplementary Appendix](#)). For a fair comparison, we did not include any additional scRNA-seq data to facilitate the imputation task. Overall, Impeller consistently outperforms others in all datasets using L1 distance, Cosine Similarity, and RMSE, indicating the effectiveness and robustness of our strategy.

In addition, we found that most methods utilizing spatial information (w\* group in [Table 2](#)) demonstrated higher imputation accuracy than those ignoring spatial information (wo\* group in [Table 2](#)), validating the presence of rich information in the spatial context. Notably, Impeller surpasses even the best gene expression-only method, eKNN, with improvements of 11.32% on 10xVisium DLPFC, 31.09% on Stereoseq, and 6.01% on SlideSeqV2 Mouse. Furthermore, compared to uniform averaging using KNN, GNN allows for more flexible neighbor information aggregation for better imputation accuracy, as reflected by the noticeably improved performance of Impeller and STAGATE.

### 4.2.2 Impact of long-range CCI

To probe disparities between Impeller and traditional GNNs in capturing long-range cell dependencies, we examined several models—Impeller, GCN ([Kipf and Welling 2016](#)), GraphSAGE ([Hamilton et al. 2017](#)), GAT ([Veličković et al. 2017](#)), and GraphTransformer ([Shi et al. 2020](#))—across varying receptive fields in the Stereoseq dataset.

In [Table 3](#), GAT and GraphSAGE suffer from gradient vanishing/exploding issues as more layers are added to capture long-range CCI, resulting in quickly degraded performance. GCN works best initially, but its performance drops with more layers added. This could be because the number of neighbors grows fast as we increase the receptive field, leaving it difficult for the target cell to understand the influence of each neighbor. Furthermore, GraphTransformer starts with high errors at a receptive field of 2. It works best at a receptive field of 8, but the error goes up again at 32. This increase in error is similar to the problem of GCN, as all cells start to look too similar to make useful representations. On the other hand, Impeller effectively tackles these challenges by the path operator, as reflected by the consistently improved results until the receptive field of 32. As the receptive field continues to grow, Impeller’s performance slightly declines, likely because distant information becomes less relevant for the target cell’s gene imputation. An additional perturbation study, demonstrating the effectiveness of Impeller in capturing CCI, is shown in the [Supplementary Appendix](#).

### 4.2.3 Advantage of heterogeneous graph

In our study, we explored the influence of graph modalities on imputation accuracy by assessing three key variants:  $var_s$ , using solely the spatial graph;  $var_g$ , utilizing only the gene similarity graph;  $var_b$ , integrating both graphs. We then calculated the performance improvement from adding  $\mathbf{G}_g$  by comparing  $var_b$  with  $var_s$ , and the improvement from adding  $\mathbf{G}_s$  by comparing  $var_b$  with  $var_g$ . As shown in [Fig. 2](#), the majority of the cases (22 out of 24) exhibit positive improvements. Specifically, in the DLPFC sample 151674, the inclusion of the gene similarity graph yields a 17.3% improvement, and the 13.6% enhancement is achieved by adding the spatial graph alone. Similarly, in sample 151508, the gene similarity graph and the spatial graph contribute to improvements of 3.6% and 9.9%, respectively. These results underscore the efficacy of our approach, particularly in scenarios where the complex interaction between spatial and

Table 2. Gene imputation benchmark.<sup>a</sup>

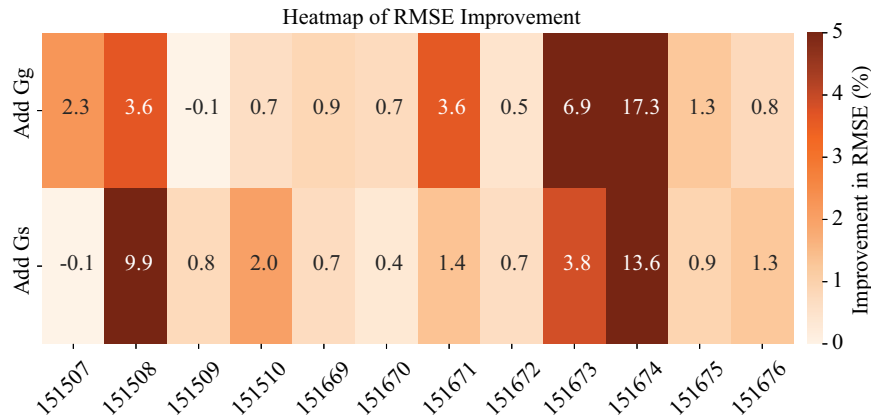
Metric	Method	Platform and dataset										
		10xVisium					SlideSeqV2					
		151507	151508	151509	151510	151669	151670	Mouse	Mouse	Mouse	Mouse	
L1 distance	wo	scVI	0.794 ± 0.004	0.838 ± 0.006	0.800 ± 0.002	0.670 ± 0.003	0.810 ± 0.003	0.696 ± 0.005	1.442 ± 0.005	1.127 ± 0.006		
		ALRA	0.499 ± 0.003	0.512 ± 0.001	0.490 ± 0.001	0.496 ± 0.001	0.467 ± 0.002	0.472 ± 0.002	0.406 ± 0.013	0.406 ± 0.013		
		eKNN	0.274 ± 0.001	0.281 ± 0.001	0.266 ± 0.000	0.275 ± 0.001	0.269 ± 0.000	0.275 ± 0.000	0.205 ± 0.001	0.205 ± 0.001		
		eSNN	1.254 ± 0.001	1.373 ± 0.001	1.266 ± 0.001	1.294 ± 0.000	1.017 ± 0.001	1.071 ± 0.001	2.802 ± 0.002	2.802 ± 0.002		
		Magic	0.779 ± 0.001	0.825 ± 0.001	0.787 ± 0.000	0.664 ± 0.001	0.795 ± 0.001	0.692 ± 0.000	1.324 ± 0.001	1.324 ± 0.001		
		scGNN	0.583 ± 0.011	0.665 ± 0.085	0.589 ± 0.011	0.584 ± 0.004	0.550 ± 0.006	0.532 ± 0.009	0.819 ± 0.240	0.819 ± 0.240		
		gimVI	0.838 ± 0.003	0.890 ± 0.003	0.835 ± 0.001	0.737 ± 0.002	0.863 ± 0.003	0.765 ± 0.001	1.325 ± 0.001	1.325 ± 0.001		
		seKNN	0.306 ± 0.000	0.309 ± 0.001	0.307 ± 0.000	0.307 ± 0.000	0.281 ± 0.000	0.289 ± 0.000	0.263 ± 0.001	0.263 ± 0.001		
		seSNN	1.254 ± 0.001	1.371 ± 0.001	1.266 ± 0.000	1.294 ± 0.000	1.017 ± 0.001	1.072 ± 0.001	2.775 ± 0.002	2.775 ± 0.002		
		Tangram	1.691 ± 0.001	1.811 ± 0.001	1.689 ± 0.000	1.420 ± 0.000	1.728 ± 0.001	1.474 ± 0.000	2.899 ± 0.001	2.899 ± 0.001		
		STLearn	1.333 ± 0.001	1.423 ± 0.001	1.332 ± 0.001	1.148 ± 0.001	1.369 ± 0.002	1.206 ± 0.001	NA	NA		
		STAGATE	0.297 ± 0.001	0.300 ± 0.002	0.295 ± 0.005	0.294 ± 0.004	0.274 ± 0.005	0.278 ± 0.002	0.190 ± 0.005	0.190 ± 0.005		
		Impeller	<b>0.248 ± 0.001</b>	<b>0.252 ± 0.003</b>	<b>0.247 ± 0.003</b>	<b>0.254 ± 0.003</b>	<b>0.242 ± 0.004</b>	<b>0.237 ± 0.001</b>	<b>0.242 ± 0.004</b>	<b>0.242 ± 0.004</b>		
		scVI	0.907 ± 0.001	0.913 ± 0.001	0.906 ± 0.001	0.903 ± 0.001	0.909 ± 0.001	0.904 ± 0.001	0.941 ± 0.001	0.941 ± 0.001		
		ALRA	0.948 ± 0.002	0.952 ± 0.002	0.952 ± 0.001	0.952 ± 0.001	0.938 ± 0.006	0.944 ± 0.003	0.980 ± 0.002	0.980 ± 0.002		
		eKNN	0.983 ± 0.000	0.984 ± 0.000	0.983 ± 0.000	0.984 ± 0.000	0.979 ± 0.000	0.979 ± 0.000	0.993 ± 0.000	0.993 ± 0.000		
eSNN	0.842 ± 0.000	0.841 ± 0.000	0.839 ± 0.000	0.840 ± 0.000	0.846 ± 0.000	0.843 ± 0.000	0.777 ± 0.001	0.777 ± 0.001				
Magic	0.915 ± 0.000	0.920 ± 0.000	0.914 ± 0.000	0.909 ± 0.000	0.916 ± 0.000	0.910 ± 0.000	0.968 ± 0.002	0.968 ± 0.002				
scGNN	0.933 ± 0.004	0.927 ± 0.016	0.932 ± 0.002	0.932 ± 0.000	0.917 ± 0.002	0.929 ± 0.002	0.948 ± 0.002	0.948 ± 0.002				
gimVI	0.957 ± 0.000	0.965 ± 0.001	0.955 ± 0.001	0.947 ± 0.001	0.962 ± 0.001	0.948 ± 0.002	0.964 ± 0.000	0.964 ± 0.000				
seKNN	0.982 ± 0.000	0.985 ± 0.000	0.982 ± 0.000	0.983 ± 0.000	0.979 ± 0.000	0.980 ± 0.000	0.995 ± 0.000	0.995 ± 0.000				
seSNN	0.843 ± 0.000	0.841 ± 0.000	0.840 ± 0.000	0.841 ± 0.000	0.851 ± 0.000	0.847 ± 0.000	0.768 ± 0.000	0.768 ± 0.000				
Tangram	0.713 ± 0.001	0.725 ± 0.001	0.717 ± 0.001	0.716 ± 0.001	0.717 ± 0.001	0.715 ± 0.000	0.772 ± 0.001	0.772 ± 0.001				
STLearn	0.718 ± 0.000	0.718 ± 0.000	0.715 ± 0.001	0.724 ± 0.000	0.715 ± 0.001	0.717 ± 0.000	NA	NA				
STAGATE	0.983 ± 0.000	0.985 ± 0.000	0.983 ± 0.001	0.984 ± 0.001	0.980 ± 0.001	0.980 ± 0.000	0.990 ± 0.000	0.990 ± 0.000				
Impeller	<b>0.987 ± 0.000</b>	<b>0.988 ± 0.000</b>	<b>0.987 ± 0.000</b>	<b>0.987 ± 0.000</b>	<b>0.983 ± 0.001</b>	<b>0.985 ± 0.000</b>	<b>0.987 ± 0.000</b>	<b>0.987 ± 0.000</b>				
scVI	0.940 ± 0.005	0.993 ± 0.006	0.949 ± 0.003	0.803 ± 0.003	0.735 ± 0.004	0.743 ± 0.003	1.628 ± 0.005	1.628 ± 0.005				
ALRA	0.784 ± 0.003	0.810 ± 0.005	0.766 ± 0.001	0.777 ± 0.001	0.735 ± 0.004	0.743 ± 0.003	0.402 ± 0.008	0.402 ± 0.008				
eKNN	0.380 ± 0.001	0.395 ± 0.002	0.382 ± 0.001	0.384 ± 0.001	0.368 ± 0.001	0.374 ± 0.001	0.416 ± 0.003	0.416 ± 0.003				
eSNN	1.378 ± 0.001	1.503 ± 0.000	1.393 ± 0.000	1.419 ± 0.001	1.143 ± 0.002	1.199 ± 0.001	2.778 ± 0.001	2.778 ± 0.001				
Magic	0.917 ± 0.001	0.972 ± 0.001	0.929 ± 0.000	0.792 ± 0.000	0.936 ± 0.001	0.824 ± 0.000	1.453 ± 0.001	1.453 ± 0.001				
scGNN	0.755 ± 0.016	0.850 ± 0.096	0.762 ± 0.011	0.755 ± 0.002	0.717 ± 0.007	0.686 ± 0.010	1.051 ± 0.307	1.051 ± 0.307				
gimVI	0.955 ± 0.002	1.002 ± 0.001	0.957 ± 0.001	0.858 ± 0.001	0.970 ± 0.002	0.890 ± 0.002	1.448 ± 0.001	1.448 ± 0.001				
seKNN	0.392 ± 0.001	0.395 ± 0.000	0.392 ± 0.000	0.392 ± 0.000	0.361 ± 0.001	0.370 ± 0.000	0.361 ± 0.001	0.361 ± 0.001				
seSNN	1.354 ± 0.001	1.474 ± 0.000	1.370 ± 0.000	1.395 ± 0.001	1.119 ± 0.001	1.175 ± 0.001	2.770 ± 0.002	2.770 ± 0.002				
Tangram	1.768 ± 0.001	1.889 ± 0.001	1.767 ± 0.000	1.503 ± 0.000	1.504 ± 0.001	1.557 ± 0.000	2.970 ± 0.001	2.970 ± 0.001				
STLearn	1.516 ± 0.001	1.629 ± 0.001	1.521 ± 0.001	1.300 ± 0.001	1.556 ± 0.002	1.362 ± 0.001	NA	NA				
STAGATE	0.384 ± 0.002	0.393 ± 0.002	0.379 ± 0.007	0.380 ± 0.007	0.357 ± 0.007	0.365 ± 0.004	0.485 ± 0.008	0.485 ± 0.008				
Impeller	<b>0.337 ± 0.001</b>	<b>0.341 ± 0.000</b>	<b>0.336 ± 0.000</b>	<b>0.340 ± 0.004</b>	<b>0.327 ± 0.007</b>	<b>0.323 ± 0.000</b>	<b>0.277 ± 0.002</b>	<b>0.277 ± 0.002</b>				

<sup>a</sup> The best results are bolded. Results marked “NA” for stLearn indicate unavailable HE stained images required by the method.

**Table 3.** Performance of different receptive fields (RMSE).

Receptive field	GCN	GraphSAGE	GAT	GraphTransformer	Impeller
2	0.339 ± 0.000	0.352 ± 0.008	0.360 ± 0.005	1.058 ± 0.463	<b>0.310 ± 0.016</b>
4	0.348 ± 0.001	0.362 ± 0.013	0.372 ± 0.022	0.424 ± 0.031	<b>0.286 ± 0.009</b>
8	0.386 ± 0.012	0.496 ± 0.033	0.454 ± 0.024	0.351 ± 0.002	<b>0.279 ± 0.001</b>
16	0.403 ± 0.001	0.617 ± 0.054	0.506 ± 0.061	0.435 ± 0.017	<b>0.286 ± 0.010</b>
32	0.418 ± 0.015	1.466 ± 0.031	1.458 ± 0.037	0.420 ± 0.000	<b>0.277 ± 0.002</b>
64	0.430 ± 0.002	1.615 ± 0.010	1.621 ± 0.004	0.420 ± 0.001	<b>0.302 ± 0.028</b>
128	0.429 ± 0.001	1.629 ± 0.003	1.614 ± 0.022	0.420 ± 0.001	<b>0.357 ± 0.001</b>

The best imputation performance is highlighted in bold.

**Figure 2.** RMSE improvement by adding different graph modalities.

gene expression data is pivotal for enhancing gene imputation accuracy.

### 4.3 Ablation study

We conducted an ablation study to evaluate the performance of four primary path operator variants of Impeller:  $\text{op}_{\text{glo}}$ , where all Impeller layers and channels (each channel representing 1D of  $\mathbf{f}^{(l)}$ ) share one path operator;  $\text{op}_{\text{cha}}$ , where channels share an operator but layers have distinct ones;  $\text{op}_{\text{lay}}$ , where all layers share one, but channels have individual operators; and  $\text{op}_{\text{ind}}$  where every layer and channel possesses an independent path operator. As depicted in Fig. 3, both  $\text{op}_{\text{glo}}$  and  $\text{op}_{\text{cha}}$  performed poorly on the DLPFC dataset, indicating the importance of distinct operators for each channel. Notably,  $\text{op}_{\text{lay}}$  and  $\text{op}_{\text{ind}}$  showed comparable results, suggesting that layer-specific operators might be optional, depending on the specific application. Another ablation study regarding different path construction and graph construction methods is shown in the [Supplementary Appendix](#).

### 4.4 Parameter analysis

To investigate the influence of Impeller’s various hyperparameters, we conducted extensive experiments using the DLPFC dataset (Sample ID: 151507) and reported the mean and standard deviation of the imputation accuracy over ten repetitions.

First, we studied the impact of  $q_s$  and  $q_g$  on the RMSE of a random walk on  $\mathbf{G}_s$  and  $\mathbf{G}_g$  following the Node2Vec mechanism. Higher values of  $q$  (i.e.  $q_s$  and  $q_g$ ) encourage the walk to sample more distant nodes, enhancing the exploration of the global graph structure, while lower values bias the walk toward neighboring nodes, facilitating local exploration. As shown in Fig. 4A, Impeller exhibits strong robustness with RMSE from 0.33 to 0.36 when  $q_s$  and  $q_g$  varied from 0.1 to

5. However, higher values of  $q_s$  and  $q_g$  tend to induce larger errors. For generality, we selected 1 as the default value for  $q_s$  and  $q_g$ .

We investigated the impact of random walk length ( $k_s$  and  $k_g$ ) and layer number ( $L$ ), shown in Fig. 4B. A path length of 2 with 10 layers results in maximum errors, reducing our model to a standard ten-layer GCN. This is because, at this path length, the model focuses on immediate neighbors, akin to how traditional GCNs operate. Such a setup, while deep, limits neighborhood exploration and increases over-smoothing risk. Conversely, a path length of 8 with 4 layers allows for capturing broader interactions (up to 28 hops), balancing extended reach and computational efficiency, thus avoiding over-smoothing and optimizing long-range CCI capture.

Next, we examined the impact of the number of random walks ( $T_s$  and  $T_g$ ). As shown in Fig. 4C,  $T_s$  and  $T_g$  appeared to have a minimal effect on results, due to the robustness of Impeller which resamples at each epoch during training. We chose 8 as the default number of random walks.

Lastly, we evaluated how the embedding dimension  $d_{\text{emb}}^{(l)}$  affects Impeller’s performance. As shown in Fig. 4D, smaller  $d_{\text{emb}}^{(l)}$  (such as 2, 4, 8) leads to limited expressive power and larger imputation errors. As  $d_{\text{emb}}^{(l)}$  increases to 16, 32, 64, or 128, Impeller’s expressive power improves and operations converge well in each run. Due to our early stopping criterion, we cease training if the validation RMSE does not improve for 50 consecutive epochs. When  $d_{\text{emb}}^{(l)}$  was set to 256 or 512, it’s hard for Impeller to converge quickly at these dimensions. To strike a balance between complexity and representational power, we opted for  $d_{\text{emb}}^{(l)}$  of 64.

In summary, these comprehensive parameter analyses reveal that Impeller is robust across a wide range of parameter settings, while still providing tunable options for balancing computational efficiency and prediction accuracy. These



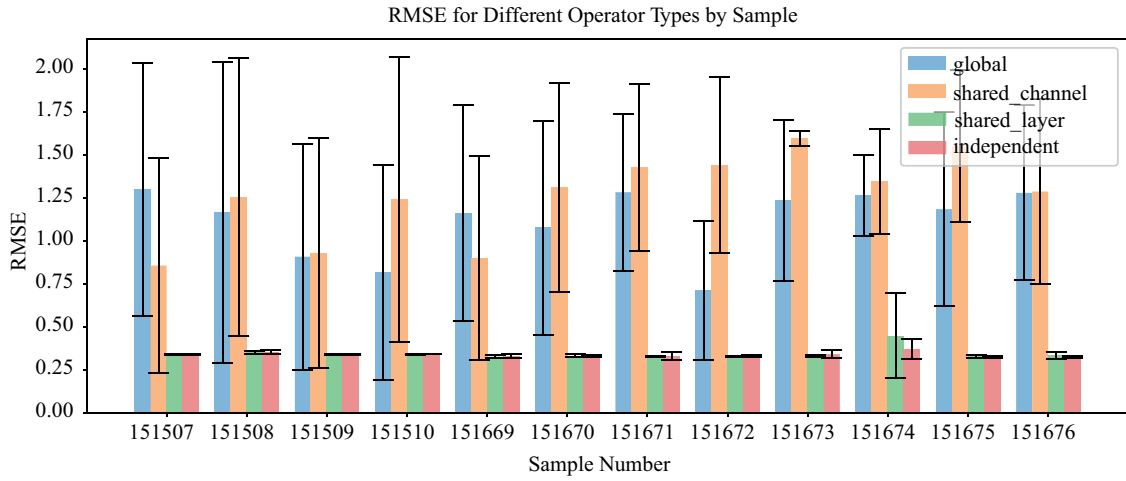


Figure 3. RMSE w.r.t. different path operators.

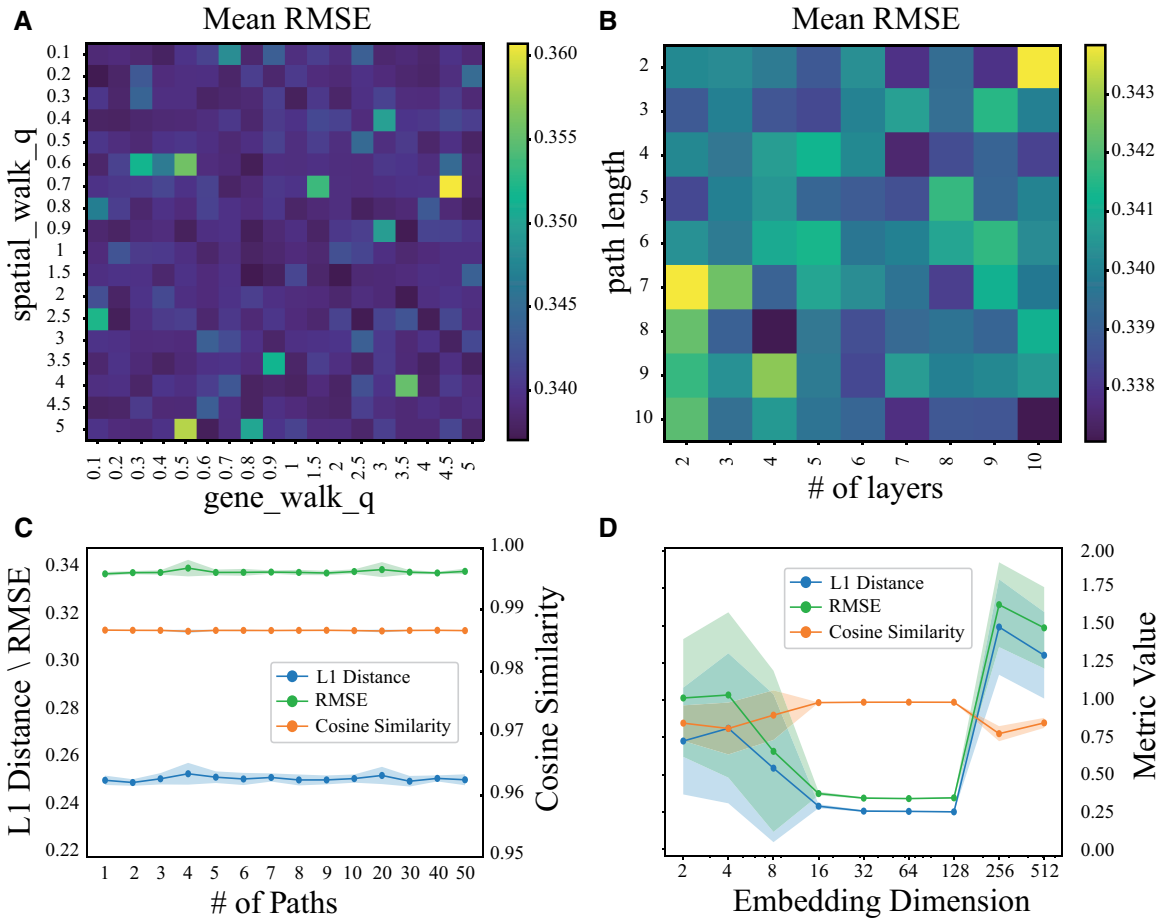


Figure 4. Parameter analysis. (A) The mean RMSE w.r.t. different  $q_s$  and  $q_g$  for generating random walks in  $G_s$  and  $G_g$ . (B) The mean RMSE w.r.t. different path lengths  $k_s$  and  $k_g$ , and the number of Impeller layers  $L$ . (C) The L1 distance, cosine similarity, and RMSE w.r.t. different number of paths  $T_s$  and  $T_g$ . (D) The L1 distance, cosine similarity and RMSE w.r.t. different embedding dimensions  $d_{emb}^{(l)}$ .

results further substantiate the effectiveness and practicality of our proposed model for gene imputation tasks.

#### 4.5 Neighbor visualization

To better understand the differences between traditional GNNs and our path-based GNN, Impeller, we turned to a

visual example (sample 151507 from the DLPFC dataset). Figure 1B shows how the typical GNN gathers information from far-away neighbors. The center node (red sphere) stacks five GNN layers to gather information from distant nodes like the one shown in yellow. But this method sometimes pulls in extra information from different tissue layers that is

**Table 4.** Running time summary of graph-based models.

Model	No. of parameters	Path (ms)	Training (ms)	Inference (ms)	RMSE
GCN	519 676		21.73 ± 10.44	21.04 ± 10.84	0.37 ± 0.02
GAT	523 768		23.84 ± 13.25	24.99 ± 11.41	0.36 ± 0.02
GraphSAGE	1 035 260		18.77 ± 11.06	16.97 ± 12.16	0.38 ± 0.01
Transformer	2 078 704		33.94 ± 16.16	38.13 ± 8.71	0.36 ± 0.01
Impeller	538 108	1.61 ± 0.60	<b>6.35 ± 0.30</b>	<b>8.43 ± 0.25</b>	<b>0.34 ± 0.00</b>

Training and inference times with the fastest performance, as well as the best imputation performance (RMSE), are highlighted in bold.

not needed. On the other hand, Fig. 1C shows our Impeller model. Instead of stacking GNN layers, Impeller samples a direct path from the center node to the target node. While using this direct path method, Impeller offers better gene imputation performance by capturing the relevant long-range CCI.

#### 4.6 Running time analysis

As shown in Table 4, we conducted a comparative model parameter and runtime analysis with popular graph-based models (GCN, GAT, GraphSAGE, and Transformer) on the DLPFC dataset. As discussed in Section 3.4.2, our model maintains a parameter count comparable to traditional GNNs, with the complexity per layer defined as  $O(d_{emb}^{(l)} \times d_{emb}^{(l+1)})$ . Specifically, our model introduces only a 3.5% increase in parameters for GCN and a 2.7% increase for GAT. In contrast, it achieves a 48.0% reduction in parameters for GraphSAGE and a 74.1% reduction for GraphTransformer (Table 4). Despite its additional path sampling step, Impeller remarkably outperformed the others in training and inference efficiency. This can be partially credited to leveraging the DGL library's optimized implementation for path sampling (<https://docs.dgl.ai/en/0.8.x/api/python/dgl.sampling.html>) and the inherently faster multiplication process used in path-based convolution compared to edge-wise information aggregation in traditional GNNs. In addition, Impeller showed the lowest RMSE, indicating superior prediction accuracy. Hence, Impeller offers a balanced blend of efficiency and precision for spatial transcriptomic data imputation, outperforming other graph-based models.

## 5 Conclusion

In this study, we introduced Impeller, a path-based heterogeneous graph learning approach tailored for spatial transcriptomic data imputation. By constructing a heterogeneous graph capturing both spatial proximity and gene expression similarity, Impeller offers a refined representation of cellular landscapes. Further, its integration of multiple GNN layers, coupled with a learnable path operator, ensures comprehensive modeling of both short and long-range cellular interactions while effectively averting over-smoothing issues. Benchmark tests across diverse datasets spanning various platforms and species underscore Impeller's superior performance compared to state-of-the-art imputation methods. This work not only establishes Impeller's prowess in spatial transcriptomic imputation but also underscores its potential to model both short- and long-range cell-cell interactions.

## Acknowledgements

The authors thank the anonymous reviewers for their valuable suggestions and the UCI OIT department for GPU resources.

## Supplementary data

Supplementary data are available at *Bioinformatics* online.

## Conflict of interest

None declared.

## Funding

This work was supported by the National Institutes of Health [R01HG012572, R01NS128523].

## Data availability

The code and preprocessed data used in this study are available at <https://github.com/aicb-ZhangLabs/Impeller> and <https://zenodo.org/records/11212604>.

## References

- Armingol E, Officer A, Harismendy O *et al.* Deciphering cell–cell interactions and communication from gene expression. *Nat Rev Genet* 2021;22:71–88.
- Biancalani T, Scalia G, Buffoni L *et al.* Deep learning and alignment of spatially resolved single-cell transcriptomes with tangram. *Nat Methods* 2021;18:1352–62. <https://doi.org/10.1038/s41592-021-01264-7>.
- Butler A, Hoffman P, Smibert P *et al.* Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat Biotechnol* 2018;36:411–20. <https://doi.org/10.1038/nbt.4096>.
- Chen A, Liao S, Cheng M *et al.* Large field of view-spatially resolved transcriptomics at nanoscale resolution. bioRxiv, <https://doi.org/10.1101/2021.01.17.427004>, 2021, preprint: not peer reviewed.
- Choe K, Pak U, Pang Y *et al.* Advances and challenges in spatial transcriptomics for developmental biology. *Biomolecules* 2023;13:156.
- Dong K, Zhang S. Deciphering spatial domains from spatially resolved transcriptomics with an adaptive graph attention auto-encoder. *Nat Commun* 2022;13:1739.
- Duan Z, Lee C, Zhang J. EXAD-GNN: explainable graph neural network for Alzheimer's disease state prediction from single-cell data. *SIP* 2023;12:e201.
- Duan Z, Wang Y, Ye W *et al.* Connecting latent relationships over heterogeneous attributed network for recommendation. *Appl Intell* 2022a;52:16214–32.
- Duan Z, Xu H, Huang Y *et al.* Multivariate time series forecasting with transfer entropy graph. *Tsinghua Sci Technol* 2022b;28:141–9.
- Duan Z, Xu H, Wang Y *et al.* Multivariate time-series classification with hierarchical variational graph pooling. *Neural Netw* 2022c;154:481–90.
- Duan Z, Dai Y, Hwang A *et al.* iherd: an integrative hierarchical graph representation learning framework to quantify network changes and prioritize risk genes in disease. *PLoS Comput Biol* 2023;19:e1011444.
- Duan Z, Xu S, Sai Srinivasan S *et al.* scenscore: leveraging single-cell epigenetic data to predict chromatin conformation using graph embedding. *Brief Bioinform* 2024;25:bbae096.

- Eliasof M, Haber E, Treister E. pathgcn: learning general graph spatial operators from paths. In: *International Conference on Machine Learning*. PMLR, 2022, 5878–91.
- Hamilton W, Ying Z, Leskovec J. Inductive representation learning on large graphs. *Adv Neural Inf Process Syst* 2017;30.
- Hao Y, Hao S, Andersen-Nissen E *et al.* Integrated analysis of multimodal single-cell data. *Cell* 2021;184:3573–87.e29. <https://doi.org/10.1016/j.cell.2021.04.048>.
- Kipf TN, Welling M. Semi-supervised classification with graph convolutional networks. arXiv, arXiv:1609.02907, 2016, preprint: not peer reviewed.
- Lähnemann D, Köster J, Szczurek E *et al.* Eleven grand challenges in single-cell data science. *Genome Biol* 2020;21:31–5.
- Linderman GC, Zhao J, Roulis M *et al.* Zero-preserving imputation of single-cell RNA-seq data[J]. *Nature communications*, 2022, 13 (1): 192.
- Lopez R, Regier J, Cole MB *et al.* Deep generative modeling for single-cell transcriptomics. *Nat Methods* 2018;15:1053–8. <https://doi.org/10.1038/s41592-018-0229-2>.
- Lopez R, Nazaret A, Langevin M *et al.* A joint model of unpaired data from scRNA-seq and spatial transcriptomics for imputing missing gene expression measurements. CoRR, CoRR:abs/1905.02269, 2019, preprint: not peer reviewed.
- Mantri M, Scuderi GJ, Abedini-Nassab R *et al.* Spatiotemporal single-cell RNA sequencing of developing chicken hearts identifies interplay between cellular differentiation and morphogenesis. *Nat Commun* 2021; 12:1771. <https://doi.org/10.1038/s41467-021-21892-z>.
- Maynard KR, Collado-Torres L, Weber LM *et al.* Transcriptome-scale spatial gene expression in the human dorsolateral prefrontal cortex. *Nat Neurosci* 2021;24:425–36.
- Pham D, Tan X, Balderson B *et al.* Robust mapping of spatiotemporal trajectories and cell–cell interactions in healthy and diseased tissues [J]. *Nature communications*, 2023, 14(1): 7739.
- Satiya R, Farrell JA, Gennert D *et al.* Spatial reconstruction of single-cell gene expression data. *Nat Biotechnol* 2015;33:495–502. <https://doi.org/10.1038/nbt.3192>.
- Shi Y, Huang Z, Feng S *et al.* Masked label prediction: Unified message passing model for semi-supervised classification. arXiv, arXiv:2009.03509, 2020, preprint: not peer reviewed.
- Stähl PL, Salmén F, Vickovic S *et al.* Visualization and analysis of gene expression in tissue sections by spatial transcriptomics. *Science* 2016;353:78–82.
- Stickels RR, Murray E, Kumar P *et al.* Highly sensitive spatial transcriptomics at near-cellular resolution with slide-seq2. *Nat Biotechnol* 2021;39:313–9.
- Strell C, Hilscher MM, Laxman N *et al.* Placing RNA in context and space—methods for spatially resolved transcriptomics. *FEBS J* 2019; 286:1468–81.
- Stuart T, Butler A, Hoffman P III, *et al.* Comprehensive integration of single-cell data. *Cell* 2019;177:1888–902.e21. <https://doi.org/10.1016/j.cell.2019.05.031>.
- van Dijk D, Sharma R, Nainys J *et al.* Recovering gene interactions from single-cell data using data diffusion. *Cell* 2018;174:716–29. e27. <https://doi.org/10.1016/j.cell.2018.05.061>.
- Veličković P, Cucurull G, Casanova A *et al.* Graph attention networks. arXiv, arXiv:1710.10903, 2017, preprint: not peer reviewed.
- Wang J, Ma A, Chang Y *et al.* scgcn is a novel graph neural network framework for single-cell RNA-seq analyses. *Nat Commun* 2021; 12:1882.
- Wang Y, Duan Z, Liao B *et al.* Heterogeneous attributed network embedding with graph convolutional networks. *AAAI* 2019;33:10061–2.
- Wang Y, Duan Z, Huang Y *et al.* Mthetgcn: a heterogeneous graph embedding framework for multivariate time series forecasting. *Pattern Recognition Letters* 2022;153:151–8.
- Xu H, Chen R, Wang Y *et al.* Cosimgcn: Towards large-scale graph similarity computation. arXiv, arXiv:2005.07115, 2020, preprint: not peer reviewed.
- Xu H, Duan Z, Wang Y *et al.* Graph partitioning and graph neural network based hierarchical graph matching for graph similarity computation. *Neurocomputing* 2021;439:348–62.
- Zeng Z, Li Y, Li Y *et al.* Statistical and machine learning methods for spatially resolved transcriptomics data analysis. *Genome Biol* 2022; 23:83–23.