



spRefine denoises and imputes spatial transcriptomic data with a reference-free framework powered by genomic language model

Tianyu Liu, Tinglin Huang, Wengong Jin, et al.

Genome Res. 2026 36: 754-768 originally published online February 3, 2026

Access the most recent version at doi:[10.1101/gr.281001.125](https://doi.org/10.1101/gr.281001.125)

References This article cites 73 articles, 6 of which can be accessed free at:
<http://genome.cshlp.org/content/36/4/754.full.html#ref-list-1>

Open Access Freely available online through the *Genome Research* Open Access option.

Creative Commons License This article, published in *Genome Research*, is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

Email Alerting Service Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Method

spRefine denoises and imputes spatial transcriptomic data with a reference-free framework powered by genomic language model

Tianyu Liu,^{1,2} Tinglin Huang,³ Wengong Jin,^{4,5} Tinyi Chu,² Rex Ying,³ and Hongyu Zhao^{1,2}

¹Interdepartmental Program in Computational Biology & Bioinformatics, Yale University, New Haven, Connecticut 06511, USA;

²Department of Biostatistics, ³Department of Computer Science, Yale University, New Haven, Connecticut 06511, USA;

⁴Department of Computer Science, Northeastern University, Boston, Massachusetts 02115, USA; ⁵Broad Institute, Cambridge, Massachusetts 02142, USA

The analysis of spatial transcriptomic data is hindered by high noise levels and missing gene measurements, challenges that are further compounded by the higher cost of spatial data compared to traditional single-cell data. To overcome this challenge, we introduce spRefine, a deep learning framework that leverages genomic language models to jointly denoise and impute spatial transcriptomic data. Our results demonstrate that spRefine yields more robust cell- and spot-level representations after denoising and imputation, substantially improving data integration. In addition, spRefine serves as a strong framework for model pretraining and the discovery of novel biological signals, as highlighted by multiple downstream applications across data sets of varying scales. Notably, spRefine enhances the accuracy of spatial aging clock estimations and uncovers new aging-related relationships associated with key biological processes, such as neuronal function loss, which offers new insights for analyzing aging effect with spatial transcriptomics.

[Supplemental material is available for this article.]

Spatially-resolved transcriptomic (SRT) technologies have enabled the investigation of the cellular functions under the spatial context (Williams et al. 2022), which cannot be directly accessed through single-cell RNA sequencing (scRNA-seq) technology. SRT technologies fall into two primary categories: (1) imaging-based methods, including smFISH (Shah et al. 2016; Codeluppi et al. 2018), MERFISH (Chen et al. 2015), seqFISH (Eng et al. 2019), and Xenium (Janesick et al. 2023; Marco Salas et al. 2025) and (2) sequencing-based methods, including Visium, STARmap (Wang et al. 2018), and Slide-seq (Rodriques et al. 2019). Imaging-based methods can measure gene expression at the subcellular resolution and generate single-cell gene expression profiles through aggregation. However, the small area covered as well as the limited number of genes pose challenges in characterizing the spatial context of the target tissue comprehensively (Biancalani et al. 2021; You et al. 2024). On the other hand, although sequencing-based methods can profile at the genome-wide level, they tend to have lower resolution and higher noise (Xu et al. 2024). These challenges presented in these data sets hidden biological signals such as aging-associated or survival-associated features (Sun et al. 2024). Therefore, there is a need to develop computational tools to address the limitations of different platforms and improve the reliability of the interpretation of spatial transcriptomics.

Efforts have been made to address the limitations of these two platforms. For imaging-based technologies, researchers utilize reference scRNA-seq data sets with a larger number of measured genes to impute the unmeasured genes in the spatial transcriptomics,

whose related methods were already benchmarked (Li et al. 2022). As examples, methods have been proposed to predict unmeasured gene expression by aggregating nearest neighbors of cells for regression (Stuart et al. 2019; Welch et al. 2019; Abdelaal et al. 2020; Shengquan et al. 2021), joint probabilistic modeling (Lopez et al. 2019; Wan et al. 2023; Haviv et al. 2024; Liu et al. 2026), and optimal transport (Cang and Nie 2020; Biancalani et al. 2021; Moriel et al. 2021). However, imputation results are limited by the quality of scRNA-seq data, which may lead to challenges in interpreting the imputed profiles. For transcriptomics from sequencing-based methods, Sprod (Wang et al. 2022) was developed to integrate both image features and gene expression levels to improve data quality by graph smoothing. SEDR (Xu et al. 2024) improves data quality and generates better representations through an unsupervised learning framework. However, these methods are restricted to sequencing-based methods and cannot uncover the missing information from unmeasured genes, and they are not able to process large-scale SRT data sets. Therefore, it is critical to have a flexible and efficient method for processing data sets from different platforms as well for driving new biology discoveries.

Recently, a method has been proposed for imputing imaging-based SRT data sets from gene networks based on protein language models (PLMs) (Zeng et al. 2024), which play an important role in the unmeasured gene expression prediction task according to their ablation tests. PLMs are foundation models trained with large-scale amino-acid sequences to model protein data as a language (Ofer et al. 2021; Zhang et al. 2024), and the success of applying PLMs to this task suggests the utility of PLMs to build expression

Corresponding author: hongyu.zhao@yale.edu

Article published online before print. Article, supplemental material, and publication date are at <https://www.genome.org/cgi/doi/10.1101/gr.281001.125>. Freely available online through the *Genome Research* Open Access option.

© 2026 Liu et al. This article, published in *Genome Research*, is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

relationships among genes. However, this method is limited to the imputation of protein-coding genes. To simultaneously impute and denoise gene expression profiles, we may consider genomic language models (gLMs) (Avsec et al. 2021; de Almeida and Pierrot 2022; Benegas et al. 2025), which model DNA sequence as a language, to derive reliable gene–gene correlations in the expression levels. There exist gLMs pretrained with either DNA sequence information purely or jointly between DNA sequences and expression profiles, and the latter setting is more related to our problem. Furthermore, the imputation process does not require external information such as scRNA-seq data, which will be more flexible for SRT data without paired scRNA-seq data or SRT data on a large scale. In principle, leveraging gene–gene correlations from gLMs pretrained with different gene expression profiles will also improve data quality.

In this article, we investigate the data noise and measured gene expression profiles in spatial transcriptomic data analysis. We define noise as imprecise gene expression measurement, which can reduce the identification of biological signals, such as cell types. We are also interested in the possibility of extending these tasks for model pretraining with enriched spatial transcriptomic data. Here, we present spRefine, a novel method to refine the quality of SRT data collected from different platforms. spRefine utilizes both measured gene expression profiles as well as gene embeddings from gLMs as inputs and learns the hidden gene–gene interactions in SRT data with a coupled Auto-Encoder, which generates spot embeddings and data set–specific gene embeddings simultaneously. It then imputes and denoises the given data by multiplying the learned embeddings to output a new matrix with improved expression levels.

Results

Overview of spRefine

We design spRefine with a decomposed auto-encoder for both imputing gene expression and denoising the measured gene expression of spatial transcriptomics. spRefine contains two modules for learning cell embeddings and gene embeddings separately, while the cell embeddings are generated based on measured expression profiles, and the gene embeddings are generated based on a pretrained sequence-to-function model. To generate the final imputed data matrix, we multiply the cell embeddings and gene embeddings optimized by reconstructing the denoised expression profiles, shown in Figure 1A. Therefore, spRefine is capable of performing data set–specific analysis for various downstream applications including spatial domain identification, cell–cell communication (CCC) inference, batch effect correction, and spatial aging clock improvement, with the help of task-specific tools shown in Figure 1B. spRefine can also work as a pretraining architecture, which treats imputation + denoising as a pretraining task and learns cell embeddings from large-scale spatial transcriptomics. spRefine performs better than baselines and can facilitate several downstream analyses with pretraining design, including survival prediction, phenotype identification, and disease-state diagnosis, shown in Figure 1C. We demonstrate better performance of spRefine through its application to SRT data of different scales from various platforms. With a larger number of inferred genes as well as better data quality, spRefine can reduce batch effect, discover novel CCCs, and improve spatial aging clock modeling. Moreover, spRefine can also work as a pretraining framework to unify cell and spot representations. In addition, we show that

spRefine can leverage the disease-associated marker genes selected based on the imputed HEST-1k data (Jaume et al. 2024) to better predict survival or disease states of patients from The Cancer Genome Atlas (TCGA) database (Weinstein et al. 2013). Our results show the strong capacity of spRefine in transferring the learned knowledge to biomedical data with different resolutions and improving our understanding of aging and cancer processes. Details of model development are discussed in the Methods section.

spRefine improves data quality through imputing and denoising profiles

spRefine performs imputation and denoising by predicting the unmeasured genes' expressions and improving the observed ones. We explain the difference between imputation (predicting the expression levels of unmeasured genes across cells or spots) and denoising (enhancing gene expression levels with known measured profiles) in Supplemental Figure 1A and B. To unify our contribution, we call this process as “refine” in this article. Inspired by Tang et al. (2024), we designed a biology-driven framework for evaluating the refined gene expression profiles by assessing the discovery of biological signals by comparing our method with different baselines. By collecting observed cell-type labels, we tested the clustering performances based on the expression profiles and utilized normalized mutual information (NMI), adjusted rand index (ARI), average silhouette width (ASW), and the averaged score (Avg) as metrics for the evaluation of signal discovery (Pedregosa et al. 2011; Luecken et al. 2022). Furthermore, we also considered demonstrating the contribution of refined profiles for correcting batch effect in the spatial transcriptomics, which can be evaluated by the metrics focusing on biological signal preservation (S_{bio}) and batch effect removal (S_{batch}) from scIB (Luecken et al. 2022). These metrics include Isolated Label Score, NMI, ARI, cell-type ASW (cASW), cell-type Local Inverse Simpson's Index (cLISI), batch LISI (bLISI), batch ASW (bASW), k -nearest-neighbor batch-effect test (kBET), Graph connectivity score (GC), and PCR comparison score. All of the metrics are in (0,1) and higher values mean better model performance. Details of our evaluation framework are summarized in the Methods section.

To measure the effectiveness of imputation, we considered some existing imputation methods, including VISTA (Liu et al. 2026), ENVI (Haviv et al. 2024), Tangram (Biancalani et al. 2021), gimVI (Lopez et al. 2019), TransImp (Qiao and Huang 2024), and SpaGE (Abdelal et al. 2020). Most of these methods were found to have relatively good performance in benchmark studies (Li et al. 2022; Marco Salas et al. 2025). Moreover, to test if incorporating spatial information can improve cell-type annotation, we also aggregated the spots into niches with major-voting labels as a baseline for spatial effect (named as Niches). Furthermore, we included a foundation model trained with spatial transcriptomics, known as Novae (Blampey et al. 2025), in our baseline to test if our imputed profiles can generate better representations than those pretrained models. We considered four data sets with subcellular spatial transcriptomics (SST), named as xenium_breast (Lin et al. 2020), xenium_brain (Lin et al. 2020), seqfish_embryo (Lohoff et al. 2022), and osmfish_brain (Zeisel et al. 2015). According to Figure 2A,B, spRefine has the best overall performance across four data sets, followed by VISTA. Moreover, Novae and Niches performed worse than spRefine, suggesting that spRefine can better capture information than the pretrained models or generated niches based on spatial location. Moreover,

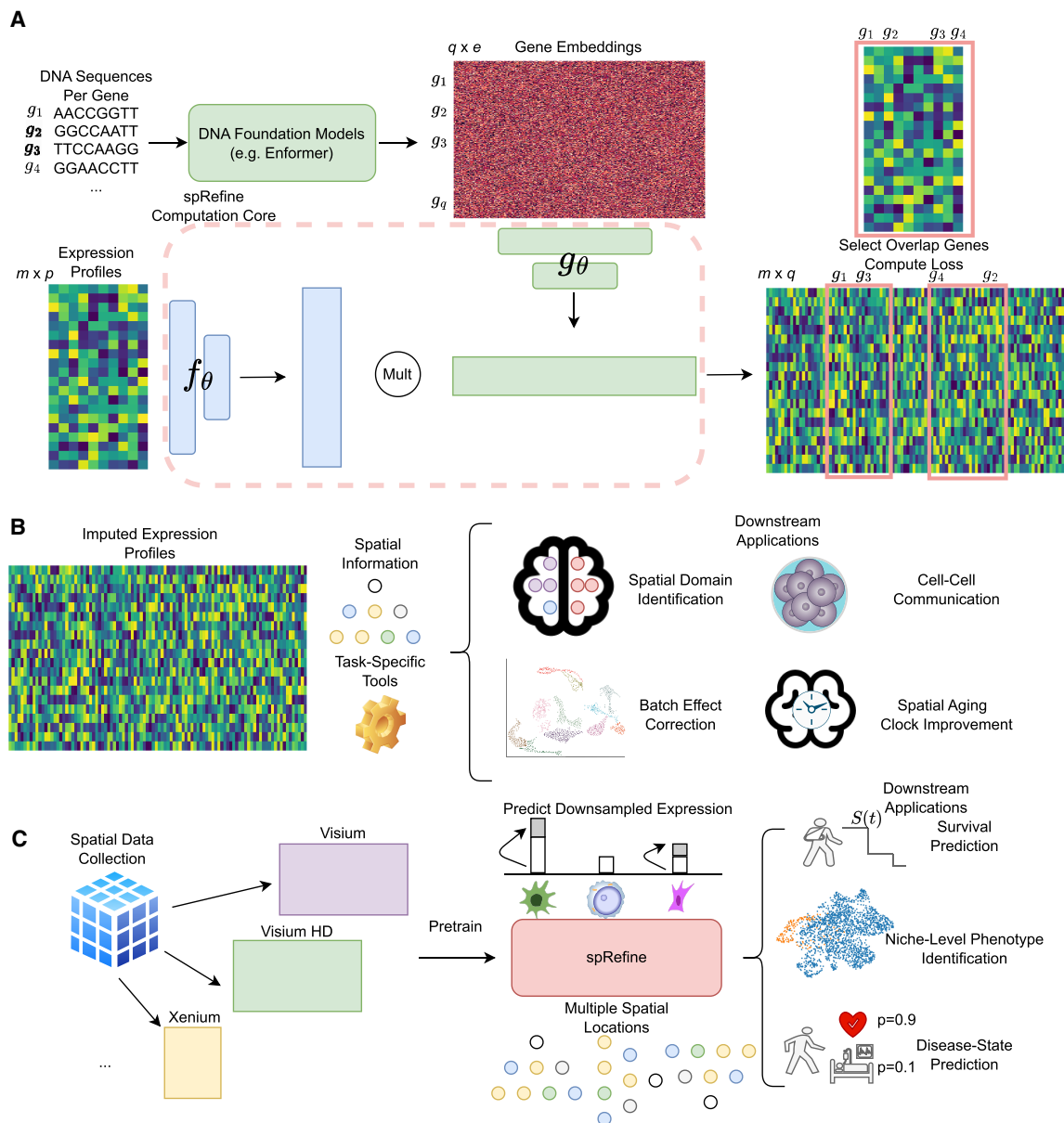


Figure 1. The landscape of spRefine and related applications. (A) The computation framework of spRefine as a tool for imputation and denoising. spRefine utilizes prior information from DNA foundation models to help uncover missing gene–gene interactions and gene expressions unmeasured by the raw transcriptomic profiles. (B) The applications of spRefine including spatial domain identification, CCC inference, batch effect correction, and spatial aging clock improvement. (C) spRefine as a pretraining framework for representing phenotype information. spRefine is capable of performing survival prediction, identifying spot-level phenotype information, and predicting disease states.

there was better cell type separation after imputation for the xenium_breast sample, supported by the cell-type-level similarity change shown in Supplemental Figure 2A. spRefine is also robust to the change of random seeds, shown in Supplemental Figure 2B.

We also performed several ablation tests, including the choices of initial gene embeddings, the choices of encoder, and the choices of loss function design, shown in Supplemental Figure 3A–C. In our comparison for gene embeddings, we found that using embeddings from Enformer gave us the best cell-type resolution, whereas other embeddings such as text-based and random-based did not help much. In our comparison for the choices of encoder, we found that modeling the spatial information with

graph neural network did not directly contribute to this task. This is explainable as different cells and cell types might have different spatial variation. In our comparison for the loss function design, we found that modeling the data with Poisson distribution loss or evidence lower bound (ELBO) did not improve the model performances. Details are discussed in the Methods section. The details of hyper-parameter tuning can also be found in the Methods section.

To directly evaluate the performance of spRefine in denoising measured gene expression profiles, we masked 30% expression levels of genes in different cells, and trained different methods to recover the masked information. Our baseline methods include

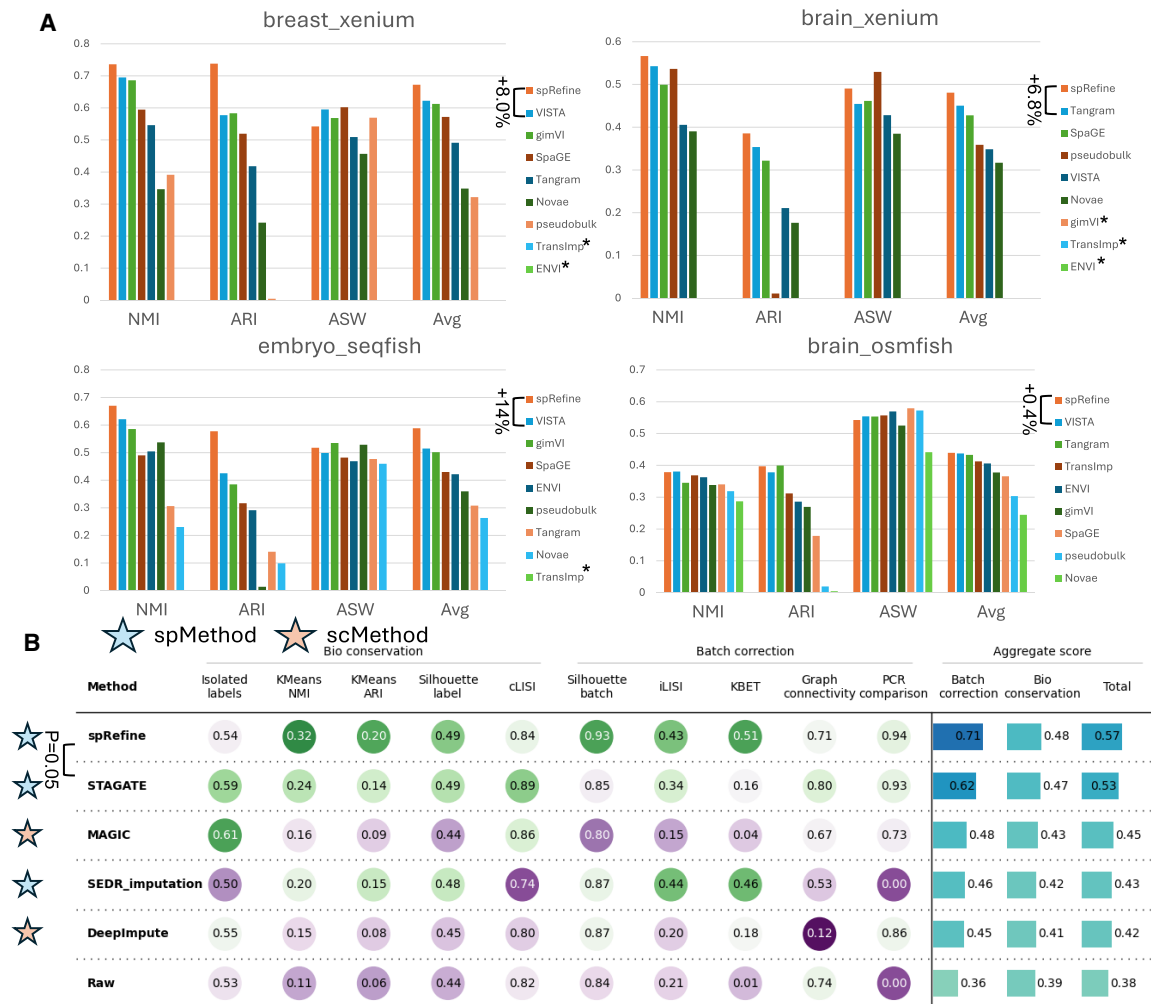


Figure 2. Comparisons between spRefine and other baselines for cell clustering and batch effect correction. (A) The clustering performance based on different imputation methods or spatial foundation models across four sub-cellular-resolved spatial transcriptomic data. The star means that the selected method met out-of-memory (OOM) errors in the specific data set. The methods are ranked by the Avg score for each data set, and we highlight the increment made by spRefine compared with the second-best baseline method. (B) The performance of batch effect correction based on different methods for processing the SpatialLIBD data set. Here, *spMethod* represents the method designed for spatial transcriptomics, whereas *scMethod* represents the method designed for single-cell transcriptomics. We performed the Wilcoxon rank-sum test (one-side) between spRefine and the second-best baseline method, and annotated the p -value (P) in this panel.

commonly used tools such as scVI (Lopez et al. 2018), DeepImpute (Arisdakessian et al. 2019), MAGIC (Van Dijk et al. 2018), and DCA (Eraslan et al. 2019), for denoising transcriptomic profiles. Here, we select Mean Squared Error (MSE) score as the evaluation metric, and report the averaged score as well as scaled standard deviation with 0.1, as some methods have very large standard deviations, which hurt the visualization. According to Supplemental Figure 4A–D, in our evaluation based on all of the four spatial transcriptomic data sets, spRefine outperforms the rest of baseline methods for recovering the masked information. Therefore, spRefine shows advantages in denoising spatial transcriptomic profiles by leveraging the prior information provided by Enformer, and it also does not suffer from the problems caused by oversmoothing.

To evaluate the performance of spRefine in denoising measured gene expression profiles and correcting batch effect, we followed the idea from Xu et al. (2024), and utilized the batch effect correction task as an indicator to evaluate model performances in

reducing the noise level of the given data set. The baselines we considered include SEDR (Xu et al. 2024) and STAGATE (Dong and Zhang 2022), which were designed for spatial transcriptomics denoising and tested for reducing batch effect. We also considered representative denoising methods designed for single-cell transcriptomics, including MAGIC and DeepImpute according to their performances in benchmarking analysis (Ding et al. 2024). We used the 10x Visium data set from SpatialLIBD (Maynard et al. 2021) for evaluation. This data set contains multiple slides and also offers expert annotation of spot types to validate the preservation of biological signals. Our selected baselines provide recommended hyper-parameter settings based on this data set for our evaluation. To perform a fair comparison, we imputed the SpatialLIBD data set with different methods and utilized Harmony (Korsunsky et al. 2019) for data integration. Harmony accepts cell embeddings generated by the encoder part of spRefine as inputs. According to Figure 2B, spRefine outperformed

the baselines in reducing the batch effect significantly, as reflected in S_{bio} , S_{batch} , and S_{total} . Again, including a GNN-based encoder does not reduce the noise of tested spatial transcriptomic data sets, shown in Supplemental Figure 3D. We further visualized the Uniform Manifold Approximation and Projection (UMAP) (McInnes et al. 2018; Becht et al. 2019) before and after integration, shown in Supplemental Figure 5. This figure also demonstrates that spRefine with Harmony can help with data integration.

Discovering novel CCCs for breast cancer

Imputing spatial transcriptomics can facilitate CCC inference in the spatial context (Almet et al. 2021) by expanding the list of genes. CCC is a technique used to capture the cell–cell interactions grouped by cell types from gene–gene interactions. Here, we considered the breast cancer data set and imputed gene expressions using spRefine. We compared the number of significant CCCs based on the raw and imputed profiles using COMMOT (Cang et al. 2023), which utilizes optimal transport (Gabriel and Marco 2019) across spatial distributions to capture the dynamic process of gene expression changes and infer CCCs. According

to Figure 3A, a higher number of CCCs (19 vs. 3) were inferred based on the imputed profiles, as expected. By detecting the CCCs based on scRNA-seq data from the same tissue with CellChat (Jin et al. 2021) as external validation, we found that the raw mode can identify only two overlapped CCCs whereas the imputed mode can identify six CCCs, which means that spRefine can also produce more confident CCCs with external validation compared with raw expression profiles. We further visualized the cell-type composition of the given sample in Figure 3B, which contains the location information of malignant cells as well as immune cells. It has been shown that chemokines play an important role in the development and migration of cancer cells, and thus we focused on the interaction of chemokines-related genes to understand the landscape of the microenvironment in breast cancer. According to Figure 3C,D, there exists a strong interaction between the sender and receiver at the junction of immune cells and ECM 1+ /CRABP2+ Malignant cells, and the signal flow further validated this discovery. The role of CXCL12 and CXCR4 has been widely discussed in the mechanism and treatment of cancer, and thus this interaction is an important marker for microenvironment (Li et al. 2012; Zhuo et al. 2012). Moreover, spRefine

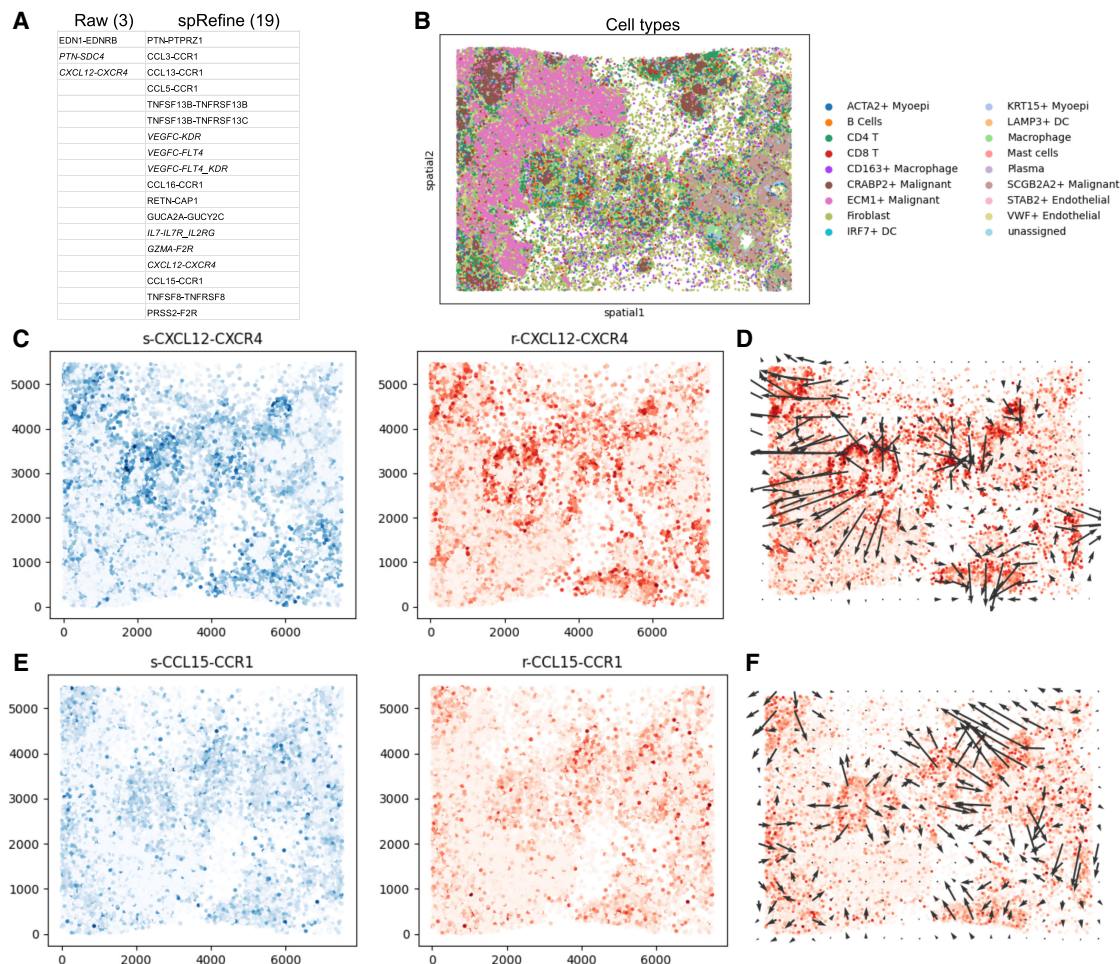


Figure 3. Using spRefine for discovering novel cell–cell communications (CCCs). (A) The comparison of discovered signals between the results of COMMOT based on the raw expression profiles (Raw) and imputed expression profiles (spRefine). In (C) and (E), “s” represents sender and “r” represents receiver. (B) The visualization of cells in xenium_breast data set colored by cell types based on spatial location. (C) The plots containing sender and receiver strength for the CCC CXCL12-CXCR4. (D) The signal direction of the CCC CXCL12-CXCR4. (E) The plots containing sender and receiver strength for the CCC CCL15-CCR1. (F) The signal direction of the CCC CCL15-CCR1.

identified a novel CCC between CCL15 and CCR1, which was not inferred based on the raw data set. According to Figure 3E,F, this CCC shows a strong interaction in the junction between immune cells and SCGB2A2⁺ malignant cells, as another key signal for breast cancer microenvironment (Li et al. 2016; Korbecki et al. 2020). Therefore, spRefine can offer additional insights through more identified CCCs. The full list of cell–cell interactions inferred by spRefine is summarized in Supplemental Figures 6 and 7.

Utilizing imputed profiles as a pretraining framework for identifying disease-associated clusters

spRefine may work as an effective pretraining framework to preserve biological signals with the help of intrinsic gene–gene interaction supported by the embeddings from DNA foundation models. To validate this assumption, we collected large-scale spatial transcriptomics from two publicly available databases, including HEST-1k (Jaume et al. 2024) and STImage-1K4M (Chen et al. 2024). These two data sets contain data from different sequencing technologies with little overlap. There are more than 500 samples and some are annotated with spot identification, for example, tumor cells versus normal cells from the patient sample. With such information, we could validate learned representations based on identifying the disease-level (or phenotype-level) information based on cell embeddings. Phenotype information (e.g., cancer cells) is usually treated as a different covariate compared with cell types (Liu et al. 2025a), and thus we intend to examine if we can utilize pretrained model to identify such covariates by clustering. The workflow in this section is summarized in Figure 4A.

Here, we pretrained spRefine based on the combined data sets and using the expression enhancement and imputation approaches as the pretraining strategies. Different from our previous experiments, we focused on more diverse data and phenotype-level differences. We first visualize representative samples based on the UMAP colored by their cell-state labels, shown in Figure 4B. This figure contains samples whose clusters are either easy or difficult to identify after imputation. The easier sample contains cells with or without tertiary lymphoid structures (TLS), whereas the difficult sample contains a more complicated tissue structure, including five different annotations. We can see that imputation does not always identify specific cells or spots, and the more diverse the cells are, the more difficult it is to distinguish the phenotypes of cells by the profiles after imputation. We computed the clustering metrics NMI score and ARI score, as well as the win rate after imputation for these two scores, where win rate = $S_{\text{imputed}} - S_{\text{raw}}$. The relationship between the number of the clusters in the tested samples and the win rates of two scores are presented in Figure 4C,D, which show a clear negative correlation between them (PCC = -0.55 , p -value = $3.91E5$ for NMI, and PCC = -0.44 , p -value = $1.8e-3$ for ARI). Therefore, imputation profiles are better suited to help identify cell types in samples with fewer predefined categories, for example, a sample with only tumor cells and normal cells. However, we also found that for samples with more cell types, clustering is not necessarily reduced after imputation, as many samples cluster around win rate = 0, and the score correlation coefficients between imputed and raw profiles are positive and significant (PCC = 0.90, p -value = $1.19E18$ for NMI, and PCC = 0.78, p -value = $7.19E11$ for ARI), which demonstrate that the imputation profile can also preserve most of the biological signals (Fig. 4E,F). As to the reason why spRefine did not work well in multiple cell-state differentiation, we speculated that the smoothing and denoising effects of impu-

tation on gene expression may have masked some very strong signals that are supposed to be associated with the classification of cell type. Therefore, we recommend checking the predefined number of clusters in the testing data sets before using spRefine as a pretraining framework.

Predicting survival and disease states for cancer patients from imputed profiles

The imputed gene expression profiles not only provide enriched gene expression information but also help select important features, for example, marker genes (Liu et al. 2025b) for different conditions. The robust features can also transfer the knowledge from the spatial transcriptomic domain to other omics-type data, for example, predicting the survival information based on the RNA-seq data from The Cancer Genome Atlas (TCGA) (Weinstein et al. 2013). Here, we utilized the imputed profiles to select marker genes across the samples from the HEST-1k data set with different disease states and included the overlapped genes between the marker genes and genes from RNA-seq to predict patient-level survival from the TCGA data set. Previous research on marker genes has shown that selecting important genes may increase efficiency and improve statistical power (Dumitrascu et al. 2021; Liu et al. 2025b). We considered nine cancer types in TCGA, and partitioned the data with RNA-seq and survival information into training, validation, and testing data sets. Here, we selected DeepSurv (Katzman et al. 2018) as the baseline method and compared the prediction results between DeepSurv based on marker genes from spRefine and that based on all the genes (Raw) or selected from other sources, including marker genes from scRNA-seq data set (Tyler et al. 2025) (sc_marker) and variable genes (Pedregosa et al. 2011). We followed the default setting of DeepSurv and considered three different metrics to evaluate the prediction performance, including concordance index (CI), integrated Brier score (IBS), and integrated negative binomial log-likelihood score (INBLLS) (Kvamme et al. 2019). Details of these metrics can be found in the Methods section, and higher CI, as well as lower IBS and INBLLS indicate better results.

The results for the Raw and spRefine settings for COAD samples are summarized in Figure 5A, where spRefine outperformed the Raw setting based on all three metrics, and made significant improvement evaluated by CI and IBS. The running time based on spRefine was also 37.2% faster than the Raw mode. Furthermore, we visualize the scores of CI based on three different gene lists based on the GBM samples because we can only access the cancer markers from GBM samples at the single-cell resolution in Figure 5B. In this figure, the performance based on spRefine is comparable with the Raw setting while outperforming the other two feature selection approaches. This suggests that spRefine offers better marker selections than the other methods. As selecting marker genes might not always improve prediction performance, we plot the prediction performances across nine diseases in Figure 5C, choosing marker genes improves prediction in five (COAD, COADREAD, GBM, PRAD, READ) of the nine conditions.

Furthermore, the imputed gene expression profiles can improve the prediction of disease states, leading to a better understanding of diseases. Here, we compared classification performance using either the raw expression profiles or imputed expression profiles. We used logistic regression and traditional classification metrics for evaluating (Pedregosa et al. 2011), and we kept the same parameters for a fair comparison. According to Supplemental Figure 8, using the imputed profiles can improve

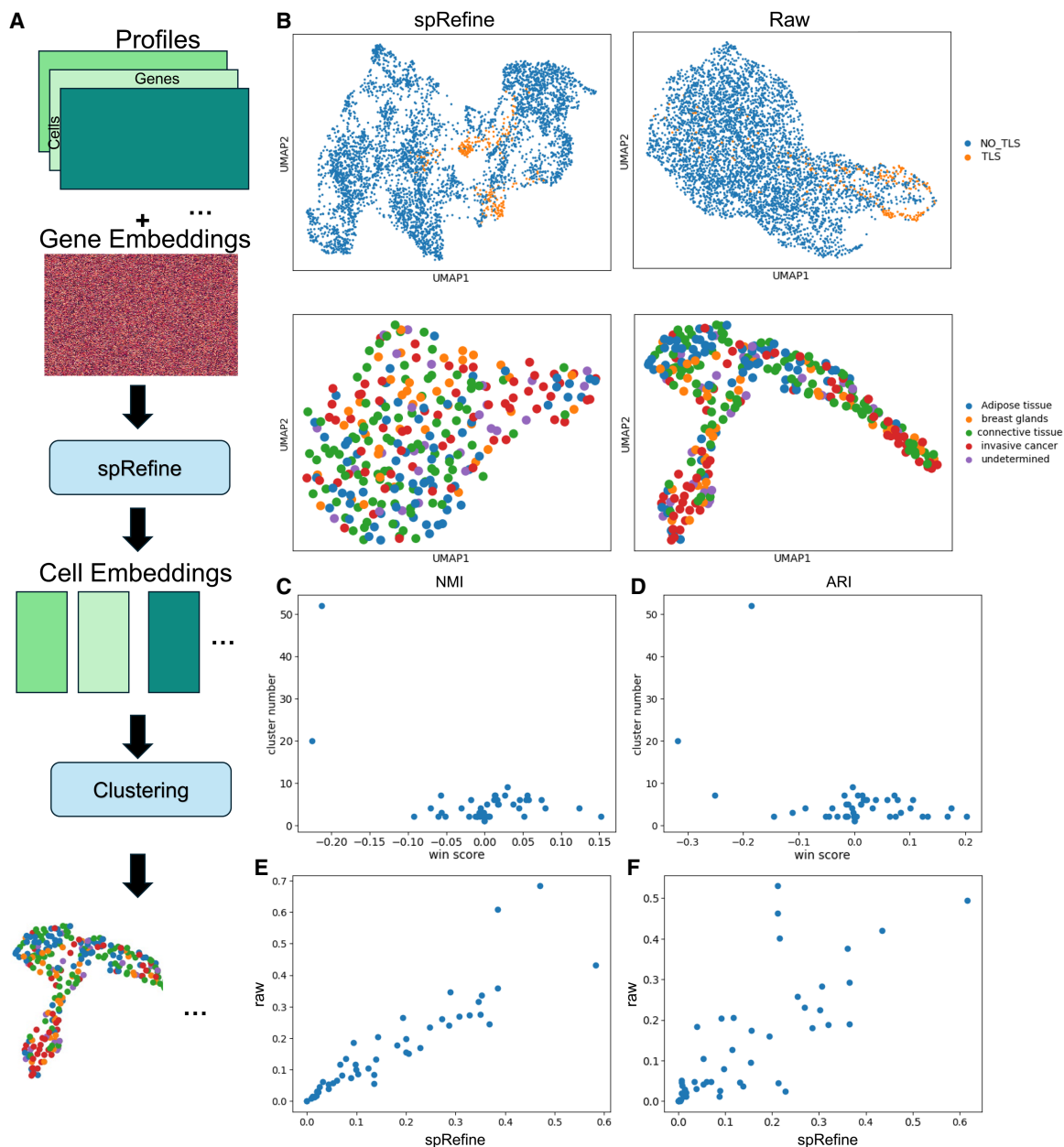


Figure 4. Analyzing the results of spRefine pretraining with large-scale spatial transcriptomics. (A) The workflow of pretraining spRefine for phenotype-level identification. (B) The UMAPs for visually comparing spot representations before (*right*) and after (*left*) imputation based on the two selected samples with different cell states. (C) The relationship between win rate computed based on NMI and cluster number. (D) The relationship between win rate computed based on ARI and cluster number. (E) The relationship between the NMI scores of raw data and data imputed by spRefine. (F) The relationship between the ARI scores of raw data and data imputed by spRefine.

classification accuracy by achieving higher precision and weighted F1 scores, with similar accuracy.

Imputing spatial transcriptomic data with spRefine can better characterize the spatial aging clock

Understanding brain aging can help us reveal the cause of neurodegenerative diseases (Cole and Franke 2017; Buckley et al. 2023; Sun et al. 2024), whose risk increases with aging. The spatial transcriptomic data measured in the brain stratified by age groups provide valuable resources to analyze the effect of spatial context

toward aging and local cell–cell interaction in the aging process. However, one challenge of building an aging model is the lack of measured genes, especially for genes associated with aging effects. The subcellular-level platforms such as MERFISH (Zhang et al. 2021) or STARMAP (Wang et al. 2018) can only measure expressions from a gene panel and we need to fill the gene expression levels of unmeasured genes. Here, we refer to the pipeline of the spatial aging clock (Sun et al. 2024), which is a model trained with spatial transcriptomics to predict the age of testing data. We have modified the pipeline by first imputing the gene expression profiles based on spRefine and then estimating the aging

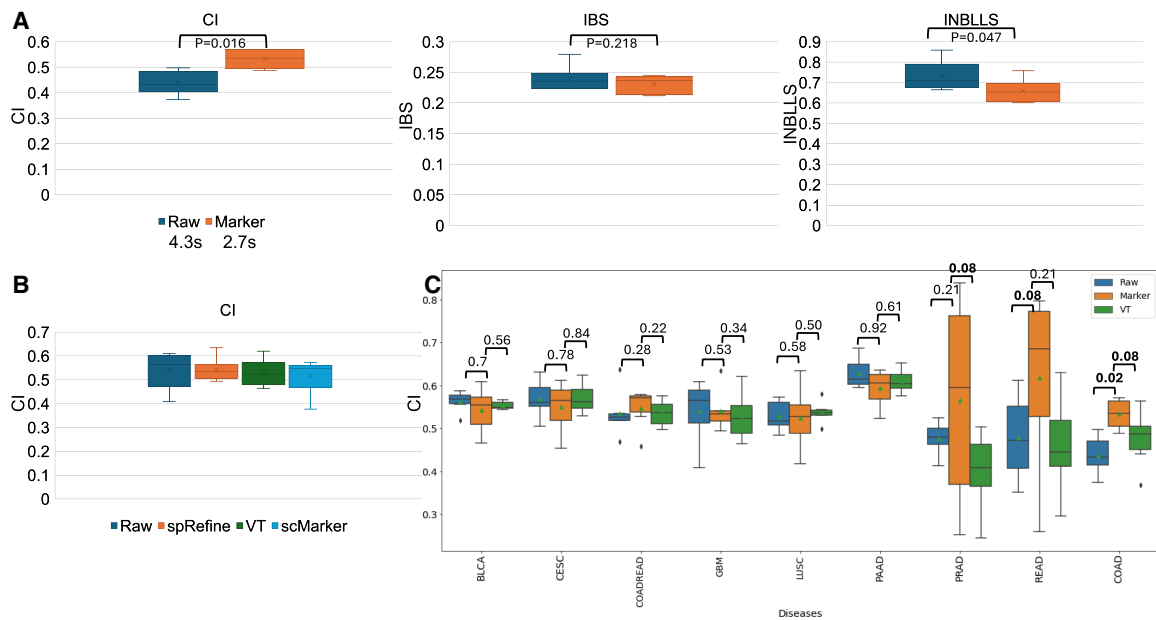


Figure 5. Applications of spRefine for disease modeling. (A) Comparison of survival prediction performances based on three different metrics between raw expression profiles and expression profiles only containing genes selected by spRefine. The cancer type of selected data is COAD. We highlighted the running time comparison and the Wilcoxon rank-sum test (one-side) p -value (P) in this panel. (B) Comparison of survival prediction performances across marker genes from different sources. (C) Comparison of survival prediction performances across different gene sets based on the cancer types from different cancer types. We annotated the Wilcoxon rank-sum test (one-side) p -value (P) in this panel, and highlighted the significant result (p -value < 0.1).

effect by considering the cellular interaction in the spatial context with a graph neural network. The modified pipeline is illustrated in Figure 6A. We considered the mouse brain data set from the original paper on the spatial aging clock. We first demonstrated the contribution of spRefine by showing the larger number of overlapped genes associated with aging after imputation (~1000 genes), shown in Figure 6B. The model of the spatial aging clock has a default gene set used for training and inference, with 100 genes not included in the gene panel, and only two genes not included after imputation. We further evaluated the prediction performances by showing the distribution differences in Figure 6C, with the difference statistically significant by the Mann-Whitney U test. Therefore, after imputation with spRefine, the refined spatial aging clock identified most cells in their correct age group. Our imputed genes can also reflect the differences of cells in the two age groups, shown in Figure 6D. Here, we visualize the differentially expressed genes (DEGs) identified by the Wilcoxon rank-sum test between the young (AC124742.1) and old group (Insc), and we can observe specific expression patterns at the spatial level. The earlier expressed genes tend to express in the embryo tissues or early-stage neuron tissue <https://www.bgee.org/gene/ENSMUSG00000118552>, whereas the later expressed genes tend to express in mature tissues <https://www.bgee.org/gene/ENSMUSG00000487822>.

In Figure 6E, we show the spatial context of the given sample colored by different information, including cell types (upper panel), age groups (middle panel), and normalized aging acceleration effects (lower panel). Higher acceleration effects mean faster aging speed rate, which is highlighted in red, which contains cells from both young and old groups, and the cellular composition of this region is also more complex. Therefore, we need to analyze the effects of cell type on senescence by analyzing cellular interactions. Following the aging effect analysis in the spatial aging clock pipe-

line, we calculated and tested the aging effect of local cell-cell interaction based on the comparison between the nearby cells and distant cells, defined as cell proximity and shown in Figure 6F. We used the cross to mark cell-cell interaction with p -value < 0.05 , where a higher proximity value represents a stronger aging effect. Compared with the proximity estimation based on the raw expression profiles, the spatial aging clock based on spRefine can identify more significant cell-cell interactions. Cell types such as Neuron-Inhibitory and Neuron-MSN have the largest number of significant interactions with other cells, which is consistent with recent experimental results for the aging effect such as neuron function loss in the mouse brain (Stanley et al. 2012; Oh et al. 2022). The third-ranked cell type is NSC, which was reported by the original spatial aging clock paper as an informative cell type. Moreover, we also discovered several less-explored cell-cell interactions with significant aging effect, including Endothelial + Ependymal, Ependymal + OPC, Microglia + VLMC, and VLMC + Neuron-Inhibitory. Considering the important roles of these cell types in brain functions, these interactions are worth future investigation.

Finally, we performed external validation based on the scRNA-seq data (Buckley et al. 2023) with mouse brains of different ages, summarized in Figure 7. As shown in Figure 7A, the selected single-cell data set has clear cell-type patterns with nearly no batch effect. We also computed the spatial aging clock information based on this data set, including predicted age distribution (shown in Fig. 7B) and predicted normalized age acceleration effect (shown in Fig. 7C). By comparing the estimation results from spatial data and single-cell data, we found that the spatial aging effect estimation is more precise based on spatial data, which is supported by the clearer age distribution difference based on the imputed spatial data. Furthermore, we computed the cell proximity effects based on the principal components of scRNA-seq data, shown in Figure

7D, and after comparing it with effects estimated based on spatial data, we highlighted the cell–cell interactions with the same significance and aging effect in the same figure. We found overlapped patterns between cell types in these two data sets, such as Neuron cells as effectors and Microglia cells as targets. We visualize the overlapped DEGs based on ages between two data sets, shown in Figure 7E for genes associated with a younger age and Figure 7F for genes associated with a younger age. These selected genes show similar expression patterns in overlapped cell types such as Neuroblast and Oligodendrocyte. Overall, spatial data have higher-resolved cell types and we discovered more cell–cell interactions with experimental support. As a result, the summary based on spatial transcriptomics is more comprehensive, which is also supported by the higher PCCs between predicted and recorded ages estimated from spatial transcriptomics shown in Figure 2 of Sun et al. (2024).

Therefore, the spatial aging clock estimated from the imputed expression profiles not only preserves the biological signals validated by prior knowledge but also inspires new directions to further investigate the roles of certain cell types in the brain aging process.

Discussion

Spatial transcriptomics can provide new insights into understanding biological processes and disease mechanisms in the spatial context. However, there are many challenges in processing spatial transcriptomics, including the imputation of unmeasured genes and the possible denoising of raw gene expression profiles. With imputation and denoising, we can obtain spatial transcriptomics with better coverage and more accurate gene expressions, which can facilitate downstream analyses.

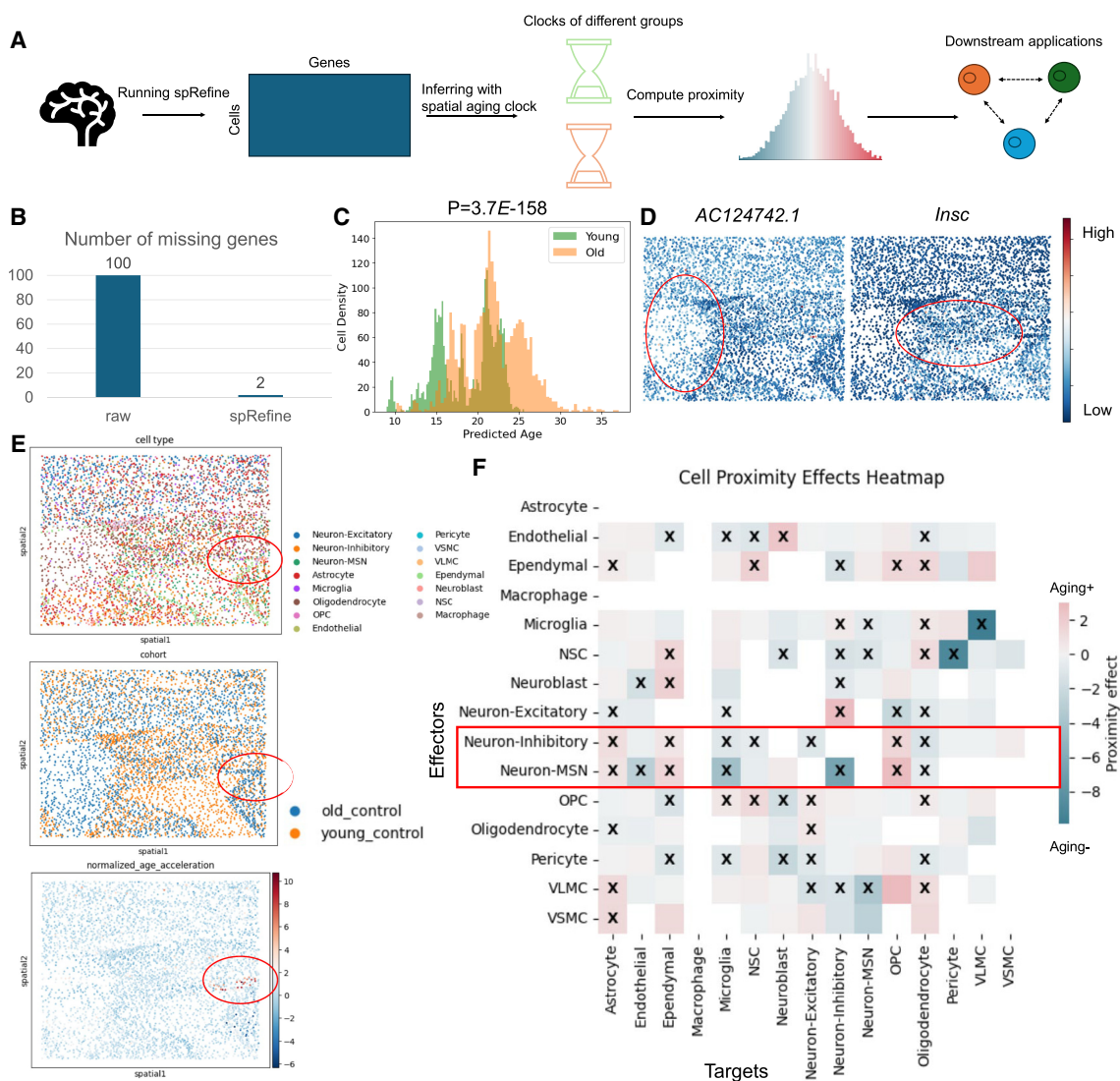


Figure 6. Leveraging contributions of spRefine to understand the biological structure of aging at the spatial level. (A) The pipeline of aging effect estimation with spRefine. (B) Comparison of the missing genes between the raw gene set and the imputed gene set. (C) The density of cells predicted with different ages, and the colors represent the measured age group. (D) The spatial-level expression patterns of two selected DEGs from different age groups. Highly expressed regions are highlighted. (E) Visualization of spots colored by cell types, age groups, and normalized age acceleration rate based on spatial location. The area with a strong aging association was highlighted in a red circle. (F) Heatmap of cell proximity effects by measuring the interactions of different cell types.

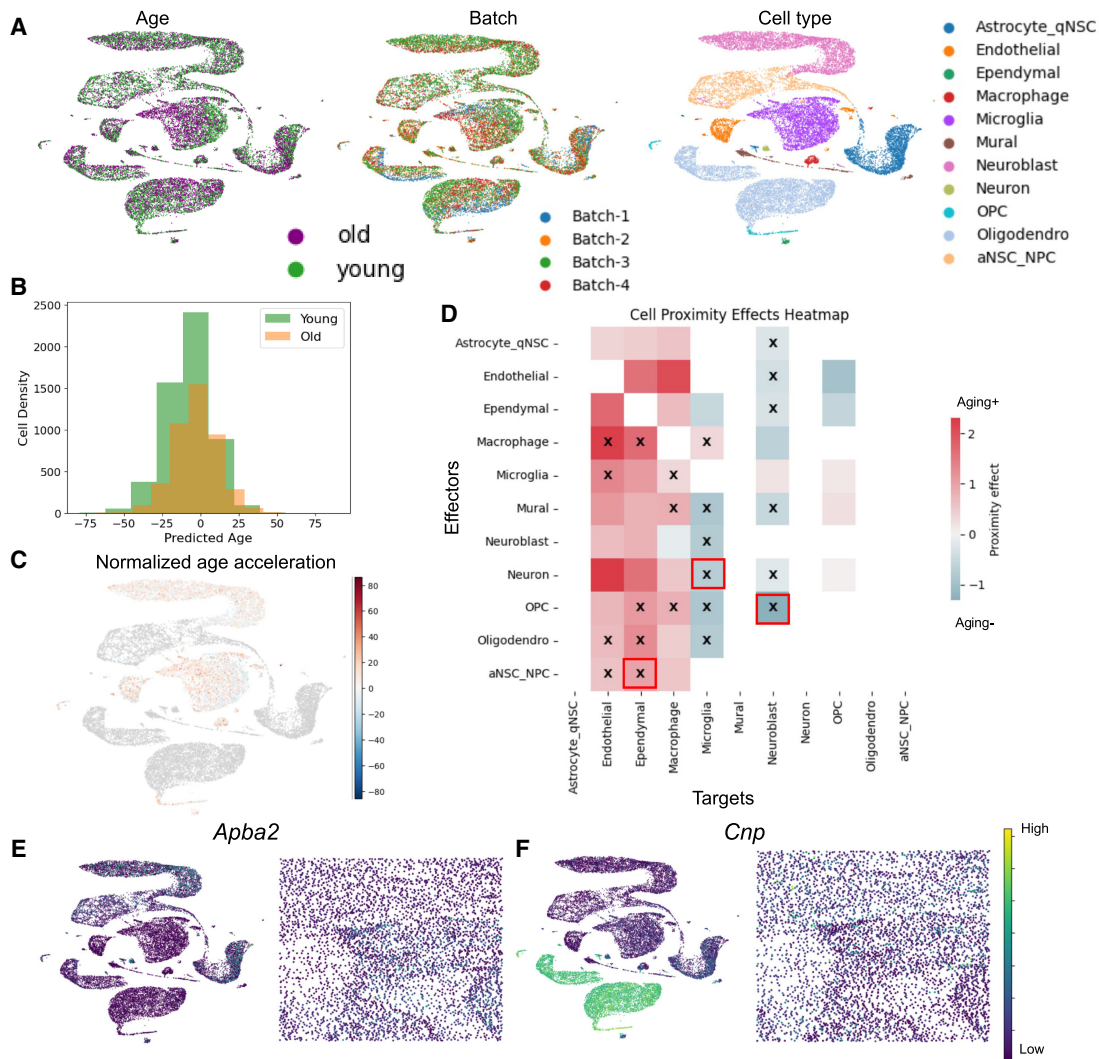


Figure 7. The aging effect analysis based on scRNA-seq data. (A) The UMAP visualization of measured scRNA-seq data colored by age (left panel), batch label (middle panel), and cell type (right panel). (B) The cell density of predicted cell age colored by observed age information. (C) The UMAP visualization of measured scRNA-seq data colored by normalized age acceleration effect. (D) The cell proximity effect of cell types based on the scRNA-seq data. (E) and (F) represent examples of expression patterns of age-specific marker genes shared by both scRNA-seq data and spatial transcriptomics. *Apba2* (E) is a marker from young group and *Cnp* (F) is a marker from old group. For each gene, the left panel represents UMAPs of scRNA-seq data and the right panel represents spatial location of spatial transcriptomics. Each panel is colored by gene expression levels.

In this article, we have proposed spRefine, a reference-free imputation and denoising pipeline for spatial transcriptomic data analysis, for imputation and denoising. spRefine takes measured spatial transcriptomics as inputs, and utilizes gene embeddings from DNA foundation models to establish gene–gene interactions. We then integrated these two data types to impute unmeasured genes and reduced the noise in the measured genes. Through the analyses of several spatial transcriptomic data, we showed that the results produced by spRefine can better cluster cells and enable more effective downstream analyses.

Given the growing volume of spatial transcriptomic data, there is the potential to collect a large amount of data for pretraining a model to analyze these data. We demonstrated that spRefine can be formatted as a pretraining framework for large-scale spatial transcriptomics collected from different platforms. Our results showed that the imputed expression profiles generated using the pretrained model can enhance cell classification for some data

sets, especially for data with fewer target cell states. Meanwhile, the gene expression profiles obtained after training on large-scale data can also be used for several downstream analyses related to diseases, including the extraction of marker genes to improve survival prediction, and better prediction of sample-level disease state. Therefore, spRefine can serve as a useful pipeline for building foundation models for spatial transcriptomics.

spRefine can also identify CCCs and spatial aging effects, which may be highly associated with certain cell states or niches. By imputing gene expression levels, we can obtain a profile with more genes to infer CCCs in the given sample. Moreover, imputing the aging-associated genes can identify more aging proximity CCCs and thus have a better model to describe the spatial aging clock. All these novel biological findings demonstrate the potential of the reference-free imputation framework.

Taken together, our framework should be understood as reference-free with respect to external single-cell or spatial atlases, but

not devoid of priors. By incorporating Enformer embeddings, we leverage sequence-trained models that encode regulatory features learned from DNA. Recognizing this distinction is critical for interpreting the reported gains: they may reflect both the utility of Enformer-derived priors and the capacity of spatial data to complement them. We view this integration as a promising paradigm for combining sequence-level priors with reference-free spatial modeling, supported by ablation studies.

Finally, we note that the performance of spRefine may be affected by the choices of models to generate the gene embeddings, and thus a better model can improve the performance of spRefine. In addition, considering that there are only a few publicly available large-scale spatial transcriptomic data sets, while some of them are not well-annotated, collecting and organizing spatial transcriptomic data is also an important research direction. The improvements of both fields may enhance the performance of spRefine to the next level.

Methods

Problem definition

In this article, we consider a data set containing multiple (m) spatial transcriptomic profiles, which can be denoted as $\mathcal{D} = \{[X^{n_1, p_1}, S^{n_1, 2}], [X^{n_2, p_2}, S^{n_2, 2}], \dots, [X^{n_m, p_m}, S^{n_m, 2}]\}$, where X represents the expression profile and S represents the location information. Our target is to utilize the reference gene panel $G^{q, e}$ with q genes and e dimensions, where $q \gg p$, to impute the unmeasured gene expression in X with the final imputed and denoised expression profiles denoted as $X^{n, q}$. Formally, our target is to learn a function \mathcal{M} , which can perform:

$$X^{n, q} = \mathcal{M}(X^{n, p}, G^{q, e}).$$

Model architecture

spRefine considers two components as inputs, including a cell embedding encoder CE_e and a gene embedding encoder GE_e . The former takes raw spatial data as input, and the latter takes raw gene embeddings as input. Our gene embeddings are computed based on the pretrained Enformer (Avsec et al. 2021). We tested other choices and they did not work as well as the setting based on Enformer. The loss is computed based on the measured genes between the raw input data and imputed outputs, and thus this process can be represented as:

$$\begin{aligned} e_c &= CE_e(X), \\ g_c &= GE_e(G), \\ X' &= F(e_c, g_c^T), \\ \mathcal{L} &= \text{MSE}(X, X', g_o), \end{aligned}$$

where g_o represents the set of overlapped genes and MSE represents the mean squared error loss function modified for a selected gene set. F is a function used to constrain the output format. If we expect to have imputed profiles equal to or larger than 0, then $F = \text{Softplus}()$; otherwise $F = \text{Identity}()$. Considering the density of measured gene expression profiles in the selected gene panel design of spatial transcriptomics, we do not model the distribution of outputs in other distributions, but keep on using Normal distribution to compute the loss.

To use spRefine as a pretraining framework, our idea is to impute the missing gene expression profiles proportion to the magnitude of expression levels, which can be defined as a downsampling strategy used in Kalfon et al. (2025) and Hao et al. (2024). The mod-

el architecture is the same for the imputation mode, while we first mask the input profiles based on a Poisson distribution and fill the masked value based on the imputed results; that is:

$$\begin{aligned} p_i &\sim \text{Poisson}(x_i \times r), \\ \pi_i &= I(u \geq r), \\ u_i &\sim U(0, 1), \\ x_{i, \text{update}} &= \max(x_i p_i \times \pi_i, 0). \end{aligned}$$

Here, $r = 0.605$ represents the dropout rate (Kalfon et al. 2025), Poisson represents the Poisson distribution, I represents the indicator function, and U represents the uniform distribution. We use x_i , which represents the raw data from spot i , and $x_{i, \text{update}}$ represents the data after downsampling selection, which is the final input for pretraining. Therefore, for the given cell, the final pretraining loss is:

$$\begin{aligned} e_c &= CE_e(x_{i, \text{update}}), \\ g_c &= GE_e(G), \\ x' &= F(e_c, g_c^T), \\ \mathcal{L}_p &= \text{MSE}(x_i, x', g_o). \end{aligned}$$

Obviously, this design can impute the missing genes by introducing a set of gene embeddings with unmeasured genes. Moreover, the gene-gene interaction as well as the target of learning a mean value of the measured gene expression profiles can help us reduce the noise level of the given data set, discussed in Liu et al. (2024).

Model training and hyper-parameter

We utilize PyTorch-lightning (<https://github.com/Lightning-AI/pytorch-lightning>) to construct and train the model. Our default optimizer is Adam (Kinga and Adam 2015), with a learning rate of $1E-4$. Our optimal batch size and latent dimensions are tuned for data set-specific settings, while the batch size is set as large as possible by default, which is shown in Supplemental Figure 9A–C. We also use the early-stopping method to reduce overfitting and split the original data into training, validation, and testing data sets to evaluate the performance of unsupervised learning.

Ablation test

When we design spRefine, we consider various ablation tests for the model training stage, including different choices of gene embeddings (Liu et al. 2025c), whether using GNN to encode spot embeddings, whether using variation auto-encoder, whether modeling the decoder output with distribution other than Normal distribution. The results of these settings are shown in Supplemental Figure 3A–C. We found that making the current model more complicated does not improve its performance, and thus our current design is efficient and optimized.

Applications

spRefine is capable of various downstream applications after imputation and denoising, including CCC discovery, cell-state stratification, survival prediction, and disease-state prediction.

1. CCC discovery. By using COMMOT (Cang et al. 2023), we can identify the cell-cell interactions with spatial context and visualize the communication scores for sender genes and receiver genes. We can also visualize the signal flow existing in the given tissue sample, which can help us describe the microenvironment of the disease in a higher resolution.
2. Cell-state stratification. By pretraining spRefine with large-scale spatial transcriptomics, we can generate spot embeddings

within the context of diverse disease states. Thus, the provided spot embeddings might work better for identifying disease-associated spots and regions for a given slide. Here, we demonstrate that the embeddings from spRefine can be directly used to identify cells with different states and also explain the performance difference affected by the cell-state resolutions.

- Survival prediction. Traditional methods for survival prediction start from transcriptomic data and train models to predict survival based on the expression profiles of all measured genes. Here, we demonstrate that by using the marker genes extracted from the imputed gene expression profiles, we can better predict the survival function of the given samples across different cancer types, and using a smaller number of genes can also improve the model efficiency.
- Disease-state prediction. We show that using the imputed gene expression profiles from large-scale spatial transcriptomics can help infer the disease state for a whole-slide sample, which further demonstrates the potential of unsupervised learning.
- Spatial aging clock construction. We use spRefine to impute the MERFISH sample with aging information and offer new insights for analyzing the spatial aging clock model by identifying aging-associated marker genes, aging-associated spatial context, and aging-associated cell–cell interactions. The pipeline we used is the same as the default setting, while the input data are imputed.

Metrics

Here, we discuss the metrics we used to evaluate the performances of spRefine and other baselines in the tasks we performed in this article.

- For imputation and denoising, we considered NMI, ARI, and ASW as metrics. Details are discussed below:
 - Normalized Mutual Information (NMI): We calculate NMI score based on computing the mutual information between the optimal Louvain clusters and the known cell-type labels and then take the normalization. Therefore, $NMI \in (0, 1)$ and higher NMI means better performance.
 - Adjusted Rand Index (ARI): We calculate ARI score by measuring the agreement between optimal Louvain clusters and cell-type labels. Therefore, $ARI \in (0, 1)$ and higher ARI means better performance.
 - cell-type Average Silhouette Width (cASW): Here, we only consider ASW for cell types. For one cell point, ASW calculates the ratio between the inner cluster distance and the intra cluster distance for this cell. Therefore, higher ASW_{cell} means better biological information preservation. For ASW_{cell} , we take the normalization; that is:

$$ASW_{cell} = \frac{ASW_{cell}^{raw} + 1}{2}.$$

- Avg: This metric represents the average score among these three metrics. All of the metrics are in (0,1) and higher values mean a better model.
- For batch effect correction, in addition to the three metrics above, we have added more metrics. For metrics covering the evaluation for S_{bio} , including NMI, ARI, cell-type ASW (cASW), and cell-type Local Inverse Simpson's Index (cLISI). For metrics covering the evaluation for S_{batch} , including bLISI, bASW, kBET, GC, and PCR comparison score.
 - Isolated Label Score: Isolated Label computes the batch correction score for the labels within the least number of batches, which offers weights for isolated clusters. Here, we utilize the

default ASW mode of isolated label score, which means we compute the ASW score for spots from the given label as the final score.

- Local Inverse Simpson's Index (LISI): This metric is used to evaluate LISI is a metric to evaluate whether data sets are well-mixed under batch labels (*bLISI*) or can be discerned with different cell types (*cLISI*). We first compute the k -nearest-neighbor list of one cell and count the number of cells that can be extracted from the neighbors before one label is observed twice.
- batch Average Silhouette Width (bASW): For one cell, ASW calculates the ratio between the inner cluster distance and the intra cluster distance for this cell. The bASW is computed by treating cluster labels as batch labels, and we took the inverse of the computed value to obtain bASW.
- kBET (Büttner et al. 2019): the kBET algorithm is used to determine if the label composition of the k -nearest-neighbors of a cell is similar to the expected label composition. For the batch label mixture of cells in the same cell types, the proportion of cells from different batches for the neighbors of one cell should match the global level distribution.
- Graph Connectivity (GC): GC measures the connectivity of cells in different cell types. If the batch effect is substantially removed, the connectivity of cells of the same cell type from different batches will have a higher connectivity score based on the k -NN neighbor graph. Therefore, we can compute the GC score for each cell type and take the average.
- PCR Comparison score: PCR is a metric to evaluate the performance of batch effect correction. We calculate the R^2 for a linear regression of the covariate of interest onto each principal component. The variance contribution of the batch effect for all the PCs is based on the sum of the product between the variance of each PC and the R^2 of each PC across all the PCs. Therefore, the score can be represented as:

$$PCR = \sum_{i=1}^G Var(C | PC_i) \times R^2(PC_i | B),$$

where G denotes the number of PCs and B denotes the batch information.

For evaluating batch effect correction, all metrics are in (0,1) and higher scores represent better preservation of biological variation or better correction of batch effect. To compute the final score S_{final} , we compute the weighted average: $S_{final} = 0.6S_{bio} + 0.4S_{batch}$. This paper (cite) has demonstrated that the weighted setting will not affect the ranking of models and thus we keep the default setting.

- Cell-state stratification. Here, we directly evaluate the ability of using the embeddings from pretrained spRefine to identify the cell states based on different samples. Here, we compare the clustering metrics (NMI, ARI, and ASW) between the clustering scores from the raw data set and the imputed data set.
- Survival prediction. To evaluate the model performance for survival prediction, we consider three metrics to evaluate the prediction accuracy, including concordance_td score, integrated_brier_score, and integrated_nbl score (<https://github.com/havakv/pycox>).
 - concordance_td score: The time-dependent concordance index (C-index) is a measure used to evaluate the predictive accuracy of survival models. It quantifies how well the model predicts the ordering of individuals' event times. Higher score represents better model performance.
 - integrated_brier_score: The integrated Brier score measures the mean squared difference between the observed outcomes and

the predicted probabilities at various times. Lower score represents better model performance.

- integrated_nbl score: The Negative Binomial Log-Likelihood score measures the probability of the observed outcomes given the model predictions. Lower scores represent better model performance.
5. Disease-state prediction. To evaluate model performance for this task, we utilize well-known metrics for classification performance evaluation, including Accuracy and Weighted F1-score, based on scikit-learn (Pedregosa et al. 2011). Higher score represents better model performance.

Baselines

Here, we consider different baseline models for different tasks. The methods are ranked in alphabet (A–Z) order.

1. Imputation and denoising.

- ENVI (Haviv et al. 2024) is a model based on conditional auto-encoder. It learns the embeddings of scRNA-seq data and multiplexed spatial transcriptomics simultaneously and decodes the embeddings to expression space for imputation. However, ENVI meets OOM issues in our benchmarking process.
- gimVI (Lopez et al. 2019) also models the gene expression from these two modalities into a joint latent space and uses variational inference to generate the output distribution and impute the expression levels of missing genes. However, gimVI does not consider the neighborhood relation in the spatial data and its implementation is not efficient, as shown in our results. Moreover, gimVI meets OOM errors in our benchmarking process.
- SpaGE (Abdelaal et al. 2020) is a model based on dimension reduction and regression. It firstly reduces the high dimensions of the input data, and in the joint low-dimensional space, it trains a regression model to impute the value of missing genes. However, SpaGE is not efficient for Xenium-based data with moderate performance.
- Tangram (Biancalani et al. 2021) is a model based on optimal matching. Tangram learns the best match relation between scRNA-seq data and spatial data by learning a mapping function and then performs the imputation based on minimizing the loss of such mapping function. However, Tangram is not scalable and the batch version of Tangram does not consider the difference between local optimal solutions and global optimal solutions.
- TransImp (Qiao and Huang 2024) is a model based on regression and spatial information regularization. It also relies on dimension reduction for the first step, and in the regression step, it considers both minimizing the loss between predicted data and ground truth data and minimizing the difference of spatial information. However, TransImp meets OOM errors in our benchmarking process and its imputation results lead to poor performance for some downstream applications, for example, RNA velocity inference.
- VISTA (Liu et al. 2026) models the gene expression profiles based on a coupled graph variational auto encoder. VISTA learns the representations of spots and cells based on encoder and then reduce the distance between two different domains and impute the gene expression profiles measured in the spatial context based on a decoder. VISTA is capable of performing various downstream applications, including CCC discovery, spatial RNA velocity inference, and spatial perturbation simulation.

- scVI (Lopez et al. 2019) utilizes negative binomial distribution to model single-cell transcriptomic data, and can learn cell embeddings as well as denoised gene expression profiles with variational auto encoder.
- DCA (Eraslan et al. 2019) models single-cell transcriptomic data based on count-data distribution, and performs data denoising based on latent space correction with decoding.

2. Batch effect correction.

- DeepImpute (Arisdakessian et al. 2019) uses a denoising auto-encoder to impute the single-cell transcriptomics. It only considers filling zero-expressed genes in cells.
- MAGIC (Van Dijk et al. 2018) uses a graph-diffusion map approach to impute the single-cell transcriptomics. It considers imputing both unmeasured gene expression levels as well as enhancing known gene expression levels.
- SEDR (Xu et al. 2024) utilizes a graph auto encoder to encode the gene expression of spots into a latent space and then reconstruct the masked gene expression profiles to perform unsupervised training. It then utilizes embeddings to enhance other downstream applications, including batch effect correction and clustering.
- STAGATE (Dong and Zhang 2022) considers a cell-type-level graph auto encoder as well as a spatial-level graph auto encoder to learn the representation with a weighted-sum approach. It also takes the representation as inputs and decodes them into expression profiles, which have been denoised after training. This method can also identify spatial domains and help extract 3D spatial architecture.

3. Cell-state stratification.

- Raw expression profiles. Here, we take the gene expression profiles without imputing and denoising as the baseline.

4. Survival prediction.

- All gene set. Here, we first consider using all of the genes to predict the survival function of the given sample.
- Features from RNA-seq. Here, we utilize variance threshold approach to select marker genes from the training data set, and then test the model performance.
- Marker gene from scRNA-seq. Here, we collect the disease marker genes from 3CA database (Tyler et al. 2025) of diseases and select these markers for testing

5. Disease-state prediction.

- Raw expression profiles. Here, we take the gene expression profiles without imputing and denoising as the baseline.

6. Spatial aging clock construction.

- Raw expression profiles. Here, we take the gene expression profiles without imputing and denoising as the baseline.

Data sets

Details of used data in this paper are summarized in [Supplemental Table 1](#). These data include (<https://linnarssonlab.org/osmFISH/availability/>), (<https://www.embopress.org/doi/full/10.15252/msb.20199389>), (<https://www.biorxiv.org/content/10.1101/2023.02.13.528102v1.abstract>), (<https://zenodo.org/records/7556184>), (<https://www.cancer.gov/ccg/research/genome-sequencing/tcga>), (<https://huggingface.co/datasets/MahmoodLab/hest>), (<https://huggingface.co/datasets/jiawennnn/STimage-1K4M>).

Code availability

Our computation resources are based on Yale High-performance Computing Center (YCRC). For data preprocessing, we reply on one CPU with up to 800 GB. For model training, we utilize one NVIDIA A100 (H100) GPU with up to 100 GB RAM. Details of running statistics and hyperparameter settings can be found in Supplemental Table 2. The codes of this project can be found at GitHub ([https://github.com/HelloWorldLTY/sprefine](https://github.com>HelloWorldLTY/sprefine)) and as Supplemental Code. The license is MIT license.

Competing interest statement

The authors declare no competing interests.

Acknowledgments

We thank Dr. Zhang Le for suggesting methods to analyze aging effect in the brain region. This project is partially supported by National Human Genome Research Institute grant K99HG013429 (T.C.), National Institutes of Health grants U24 HG012108 and U01 HG013840 (H.Z.), and National Science Foundation IIS Div Of Information & Intelligent Systems 2,403,317 (R.Y.).

Author contributions: T.L. proposed the study. T.L., T.H., W.J., and T.C. designed the model. T.L. ran all the experiments. R.Y. provided the computation resources. T.L., T.H., W.J., R.Y., and H.Z. wrote the manuscript. H.Z. supervised this project.

References

- Abdelaal T, Mourragui S, Mahfouz A, Reinders MJT. 2020. SpaGE: spatial gene enhancement using scRNA-seq. *Nucleic Acids Res* **48**: e107. doi:10.1093/nar/gkaa740
- Almet AA, Cang Z, Jin S, Nie Q. 2021. The landscape of cell–cell communication through single-cell transcriptomics. *Curr Opin Syst Biol* **26**: 12–23. doi:10.1016/j.coisb.2021.03.007
- Arisdakessian C, Poirion O, Yunits B, Zhu X, Garmire LX. 2019. DeepImpute: an accurate, fast, and scalable deep neural network method to impute single-cell RNA-seq data. *Genome Biol* **20**: 211. doi:10.1186/s13059-019-1837-6
- Avsec Z, Agarwal V, Visentin D, Ledsam JR, Grabska-Barwinska A, Taylor KR, Assael Y, Jumper J, Kohli P, Kelley DR. 2021. Effective gene expression prediction from sequence by integrating long-range interactions. *Nat Methods* **18**: 1196–1203. doi:10.1038/s41592-021-01252-x
- Becht E, McInnes L, Healy J, Dutertre C-A, Kwok IWH, Ng LG, Ginhoux F, Newell EW. 2019. Dimensionality reduction for visualizing single-cell data using UMAP. *Nat Biotechnol* **37**: 38–44. doi:10.1038/nbt.4314
- Benegas G, Ye C, Albers C, Li JC, Song YS. 2025. Genomic language models: opportunities and challenges. *Trends Genet* **41**: 286–302. doi:10.1016/j.tig.2024.11.013
- Biancalani T, Scalia G, Buffoni L, Avasthi R, Lu Z, Sanger A, Tokcan N, Vanderburg CR, Segerstolpe Å, Zhang M, et al. 2021. Deep learning and alignment of spatially resolved single-cell transcriptomes with tanger. *Nat Methods* **18**: 1352–1362. doi:10.1038/s41592-021-01264-7
- Blampey Q, Benkirane H, Bercovici N, Mulder K, Gessain G, Ginhoux F, André F, Courmède PH. 2025. Novae: a graph-based foundation model for spatial transcriptomics data. *Nat Methods* **22**: 2539–2550. doi:10.1038/s41592-025-02899-6
- Buckley MF, Sun ED, George BM, Liu L, Schaub N, Xu L, Reyes JM, Goodell MA, Weissman IL, Wyss-Coray T, et al. 2023. Cell-type-specific aging clocks to quantify aging and rejuvenation in neurogenic regions of the brain. *Nat Aging* **3**: 121–137. doi:10.1038/s43587-022-00335-4
- Büttner M, Miao Z, Wolf FA, Teichmann SA, Theis FJ. 2019. A test metric for assessing single-cell RNA-seq batch correction. *Nat Methods* **16**: 43–49. doi:10.1038/s41592-018-0254-1
- Cang Z, Nie Q. 2020. Inferring spatial and signaling relationships between cells from single cell transcriptomic data. *Nat Commun* **11**: 2084. doi:10.1038/s41467-020-15968-5
- Cang Z, Zhao Y, Almet AA, Stabell A, Ramos R, Plikus MV, Atwood SX, Nie Q. 2023. Screening cell–cell communication in spatial transcriptomics via collective optimal transport. *Nat Methods* **20**: 218–228. doi:10.1038/s41592-022-01728-4
- Chen KH, Boettiger AN, Moffitt JR, Wang S, Zhuang X. 2015. Spatially resolved, highly multiplexed RNA profiling in single cells. *Science* **348**: aaa6090. doi:10.1126/science.aaa6090
- Chen J, Zhou M, Wu W, Zhang J, Li Y, Li D. 2024. STImage-1K4M: A histopathology image-gene expression dataset for spatial transcriptomics. *Adv Neural Inf Process Syst* **37**: 35796–35823. doi:10.52202/079017-1129
- Codeluppi S, Borm LE, Zeisel A, La Manno G, van Lunteren JA, Svensson CI, Linnarsson S. 2018. Spatial organization of the somatosensory cortex revealed by osmFISH. *Nat Methods* **15**: 932–935. doi:10.1038/s41592-018-0175-z
- Cole JH, Franke K. 2017. Predicting age using neuroimaging: innovative brain ageing biomarkers. *Trends Neurosci* **40**: 681–690. doi:10.1016/j.tins.2017.10.001
- de Almeida BP, Pierrot T. 2022. Large language models for genomics. In *Proceedings of the LLMs4Bio Workshop at AAAI*. Virtual conference.
- Ding J, Liu R, Wen H, Tang W, Li Z, Venegas J, Su R, Molho D, Jin W, Wang Y, et al. 2024. Dance: a deep learning library and benchmark platform for single-cell analysis. *Genome Biol* **25**: 72. doi:10.1186/s13059-024-03211-z
- Dong K, Zhang S. 2022. Deciphering spatial domains from spatially resolved transcriptomics with an adaptive graph attention auto-encoder. *Nat Commun* **13**: 1739. doi:10.1038/s41467-022-29439-6
- Dumitrascu B, Villar S, Mixon DG, Engelhardt BE. 2021. Optimal marker gene selection for cell type discrimination in single cell analyses. *Nat Commun* **12**: 1186. doi:10.1038/s41467-021-21453-4
- Eng C-HL, Lawson M, Zhu Q, Dries R, Koulena N, Takei Y, Yun J, Cronin C, Karp C, Yuan G-C, et al. 2019. Transcriptome-scale super-resolved imaging in tissues by RNA seqFISH+. *Nature* **568**: 235–239. doi:10.1038/s41586-019-1049-y
- Eraslan G, Simon LM, Mircea M, Mueller NS, Theis FJ. 2019. Single-cell RNA-seq denoising using a deep count autoencoder. *Nat Commun* **10**: 390. doi:10.1038/s41467-018-07931-2
- Gabriel P, Marco C. 2019. Computational optimal transport: with applications to data science. *Found Trends Mach Learn* **11**(5–6): 355–607. doi:10.1561/MAL
- Hao M, Gong J, Zeng X, Liu C, Guo Y, Cheng X, Wang T, Ma J, Zhang X, Song L. 2024. Large-scale foundation model on single-cell transcriptomics. *Nat Methods* **21**: 1481–1491. doi:10.1038/s41592-024-02305-7
- Haviv D, Remšik J, Gatie M, Snopkowski C, Takizawa M, Pereira N, Bashkin J, Jovanovich S, Nawy T, Chaligne R, et al. 2024. The covariance environment defines cellular niches for spatial inference. *Nat Biotechnol* **43**: 269–280. doi:10.1038/s41587-024-02193-4
- Janesick A, Shelansky R, Gottscho AD, Wagner F, Williams SR, Rouault M, Beliakoff G, Morrison CA, Oliveira MF, Sicherman JT, et al. 2023. High resolution mapping of the tumor microenvironment using integrated single-cell, spatial and in situ analysis. *Nat Commun* **14**: 8353. doi:10.1038/s41467-023-43458-x
- Jaume G, Doucet P, Song AH, Lu MY, Almagro-Pérez C, Wagner SJ, Vaidya AJ, Chen RJ, Williamson DFK, Kim A, et al. 2024. HEST-1k: a dataset for spatial transcriptomics and histology image analysis. *Adv Neural Inf Process Syst* **37**: 53798–53833. doi:10.52202/079017-1704
- Jin S, Guerrero-Juarez CF, Zhang L, Chang I, Ramos R, Kuan C-H, Myung P, Plikus MV, Nie Q. 2021. Inference and analysis of cell–cell communication using cellchat. *Nat Commun* **12**: 1088. doi:10.1038/s41467-021-21246-9
- Kalfon J, Samaran J, Peyré G, Cantini L. 2025. scPRINT: pre-training on 50 million cells allows robust gene network predictions. *Nat Commun* **16**: 3607. doi:10.1038/s41467-025-58699-1
- Katzman JL, Shaham U, Cloninger A, Bates J, Jiang T, Kluger Y. 2018. DeepSurv: personalized treatment recommender system using a cox proportional hazards deep neural network. *BMC Med Res Methodol* **18**: 24. doi:10.1186/s12874-018-0482-1
- Kinga D, Adam JB. 2015. A method for stochastic optimization. In *International conference on learning representations (ICLR)*, San Diego, Vol. 5, no. 6.
- Korbecki J, Kojder K, Simińska D, Bohatyrewicz R, Gutowska I, Chlubek D, Baranowska-Bosiacka I. 2020. CC chemokines in a tumor: a review of pro-cancer and anti-cancer properties of the ligands of receptors CCR1, CCR2, CCR3, and CCR4. *Int J Mol Sci* **21**: 8412. doi:10.3390/ijms21218412
- Korsunsky I, Millard N, Fan J, Slowikowski K, Zhang F, Wei K, Baglaenko Y, Brenner M, Loh P-R, Raychaudhuri S. 2019. Fast, sensitive and accurate integration of single-cell data with harmony. *Nat Methods* **16**: 1289–1296. doi:10.1038/s41592-019-0619-0
- Kvamme H, Borgun Ø, Scheel I. 2019. Time-to-event prediction with neural networks and cox regression. *J Mach Learn Res* **20**: 1–30.
- Li J-Y, Ou Z-L, Yu S-J, Gu X-L, Yang C, Chen A-X, Di G-H, Shen Z-Z, Shao Z-M. 2012. The chemokine receptor CCR4 promotes tumor growth and lung metastasis in breast cancer. *Breast Cancer Res Treat* **131**: 837–848. doi:10.1007/s10549-011-1502-6

- Li Y, Wu J, Zhang P. 2016. CCL15/CCR1 axis is involved in hepatocellular carcinoma cells migration and invasion. *Tumor Biol* **37**: 4501–4507. doi:10.1007/s13277-015-4287-0
- Li B, Zhang W, Guo C, Xu H, Li L, Fang M, Hu Y, Zhang X, Yao X, Tang M, et al. 2022. Benchmarking spatial and single-cell transcriptomics integration methods for transcript distribution prediction and cell type deconvolution. *Nat Methods* **19**: 662–670. doi:10.1038/s41592-022-01480-9
- Lin Y, Cao Y, Kim HJ, Salim A, Speed TP, Lin DM, Yang P, Yang JYH. 2020. scClassify: sample size estimation and multiscale classification of cells using single and multiple reference. *Mol Syst Biol* **16**: e9389. doi:10.15252/msb.20199389
- Liu T, Li K, Wang Y, Li H, Zhao H. 2024. Evaluating the utilities of foundation models in single-cell data analysis. bioRxiv doi:10.1101/2023.09.08.555192
- Liu T, De Brouwer E, Kuo T, Diamant N, Missarova A, Wang H, Hao M, Corrado Bravo H, Scalia G, Regev A, et al. 2025a. Learning multi-cellular representations of single-cell transcriptomics data enables characterization of patient-level disease states. In *International Conference on Research in Computational Molecular Biology*, pp. 303–306. Springer Nature, Cham, Switzerland.
- Liu T, Long W, Cao Z, Wang Y, He CH, Zhang L, Strittmatter SM, Zhao H. 2025b. Cosgenegate selects multi-functional and credible biomarkers for single-cell analysis. *Brief Bioinform* **26**: bbae626. doi:10.1093/bib/bbae626
- Liu T, Huang T, Wang L, Lin Y, Ying R, Zhao H. 2025c. UNICORN: towards universal cellular expression prediction with a multi-task learning framework. *Nat Commun* **16**: 9455. doi:10.1038/s41467-025-64506-8
- Liu T, Lin Y, Luo X, Sun Y, Zhao H. 2026. VISTA uncovers missing gene expression and spatial-induced information for spatial transcriptomic data analysis. *Commun Biol* **9**: 203. doi:10.1038/s42003-025-09479-6
- Lohoff T, Ghazanfar S, Missarova A, Koulana N, Pierson N, Griffiths JA, Bardot ES, Eng C-HL, Tyser RCV, Argelaguet R, et al. 2022. Integration of spatial and single-cell transcriptomic data elucidates mouse organogenesis. *Nat Biotechnol* **40**: 74–85. doi:10.1038/s41587-021-01006-2
- Lopez R, Regier J, Cole MB, Jordan MI, Yosef N. 2018. Deep generative modeling for single-cell transcriptomics. *Nat Methods* **15**: 1053–1058. doi:10.1038/s41592-018-0229-2
- Lopez R, Nazaret A, Langevin M, Samaran J, Regier J, Jordan MI, Yosef N. 2019. A joint model of unpaired data from scRNA-seq and spatial transcriptomics for imputing missing gene expression measurements. arXiv:1905.02269 [cs.LG].
- Luecken MD, Büttner M, Chaichoompu K, Danese A, Interlandi M, Mueller MF, Strobl DC, Zappia L, Dugas M, Colomé-Tatché M, et al. 2022. Benchmarking atlas-level data integration in single-cell genomics. *Nat Methods* **19**: 41–50. doi:10.1038/s41592-021-01336-8
- Marco Salas S, Kuemmerle LB, Mattsson-Langseth C, Tismeyer S, Avenel C, Hu T, Rehman H, Grillo M, Czarnewski P, Helgadottir S, et al. 2025. Optimizing xenium in situ data utility by quality assessment and best-practice analysis workflows. *Nat Methods* **22**: 813–823. doi:10.1038/s41592-025-02617-2
- Maynard KR, Collado-Torres L, Weber LM, Uyttingco C, Barry BK, Williams SR, Catalini JL, Tran MN, Besich Z, Tippani M, et al. 2021. Transcriptome-scale spatial gene expression in the human dorsolateral prefrontal cortex. *Nat Neurosci* **24**: 425–436. doi:10.1038/s41593-020-00787-0
- McInnes L, Healy J, Saul N, Großberger L. 2018. UMAP: uniform manifold approximation and projection. *J Open Source Softw* **3**: 861. doi:10.21105/joss
- Moriel N, Senel E, Friedman N, Rajewsky N, Karaiskos N, Nitzan M. 2021. NovoSpaRc: flexible spatial reconstruction of single-cell gene expression with optimal transport. *Nat Protoc* **16**: 4177–4200. doi:10.1038/s41596-021-00573-7
- Ofer D, Brandes N, Linial M. 2021. The language of proteins: NLP, machine learning & protein sequences. *Comput Struct Biotechnol J* **19**: 1750–1758. doi:10.1016/j.csbj.2021.03.022
- Oh YM, Lee SW, Kim WK, Chen S, Church VA, Cates K, Li T, Zhang B, Dolle RE, Dahiya S, et al. 2022. Age-related Huntington's disease progression modeled in directly reprogrammed patient-derived striatal neurons highlights impaired autophagy. *Nat Neurosci* **25**: 1420–1433. doi:10.1038/s41593-022-01185-4
- Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, et al. 2011. Scikit-learn: Machine learning in Python. *J Mach Learn Res* **12**: 2825–2830.
- Qiao C, Huang Y. 2024. Reliable imputation of spatial transcriptomes with uncertainty estimation and spatial regularization. *Patterns* **5**: 101021. doi:10.1016/j.patter.2024.101021
- Rodrigues SG, Stickels RR, Goeva A, Martin CA, Murray E, Vanderburg CR, Welch J, Chen LM, Chen F, Macosko EZ. 2019. Slide-seq: a scalable technology for measuring genome-wide expression at high spatial resolution. *Science* **363**: 1463–1467. doi:10.1126/science.aaw1219
- Shah S, Lubeck E, Zhou W, Cai L. 2016. In situ transcription profiling of single cells reveals spatial organization of cells in the mouse hippocampus. *Neuron* **92**: 342–357. doi:10.1016/j.neuron.2016.10.001
- Shengquan C, Boheng Z, Xiaoyang C, Xuegong Z, Rui J. 2021. stPlus: a reference-based method for the accurate enhancement of spatial transcriptomics. *Bioinformatics* **37**(Supplement_1): i299–i307. doi:10.1093/bioinformatics/btab298
- Stanley EM, Fadel JR, Mott DD. 2012. Interneuron loss reduces dendritic inhibition and GABA release in hippocampus of aged rats. *Neurobiol Aging* **33**: 431.e1–431.e13. doi:10.1016/j.neurobiolaging.2010.12.014
- Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM, Hao Y, Stoeckius M, Smibert P, Satija R. 2019. Comprehensive integration of single-cell data. *Cell* **177**: 1888–1902.e21. doi:10.1016/j.cell.2019.05.031
- Sun ED, Zhou OY, Hauptschein M, Rappoport N, Xu L, Navarro Negredo P, Liu L, Rando TA, Zou J, Brunet A. 2024. Spatial transcriptomic clocks reveal cell proximity effects in brain ageing. *Nature* **638**: 160–171. doi:10.1038/s41586-024-08334-8
- Tang Z, Chen G, Chen S, Yao J, You L, Chen CY-C. 2024. Modal-nexus auto-encoder for multi-modality cellular data integration and imputation. *Nat Commun* **15**: 9021. doi:10.1038/s41467-024-53355-6
- Tyler M, Gavish A, Barbolin C, Tschernichovsky R, Hoefflin R, Mints M, Puram SV, Tirosh I. 2025. The curated cancer cell atlas provides a comprehensive characterization of tumors at single-cell resolution. *Nat Cancer* **6**: 1088–1101. doi:10.1038/s43018-025-00957-8
- van Dijk D, Sharma R, Nainys J, Yin K, Kathail P, Carr AJ, Burdziaik C, Moon KR, Chaffer CL, Pattabiraman D, et al. 2018. Recovering gene interactions from single-cell data using data diffusion. *Cell* **174**: 716–729.e27. doi:10.1016/j.cell.2018.05.061
- Wan X, Xiao J, Tam SST, Cai M, Sugimura R, Wang Y, Wan X, Lin Z, Wu AR, Yang C. 2023. Integrating spatial and single-cell transcriptomics data using deep generative models with spatialscope. *Nat Commun* **14**: 7848. doi:10.1038/s41467-023-43629-w
- Wang X, Allen WE, Wright MA, Sylwestrak EL, Samusik N, Vesuna S, Evans K, Liu C, Ramakrishnan C, Liu J, et al. 2018. Three-dimensional intact-tissue sequencing of single-cell transcriptional states. *Science* **361**: eaat5691. doi:10.1126/science.aat5691
- Wang Y, Song B, Wang S, Chen M, Xie Y, Xiao G, Wang L, Wang T. 2022. Spro for de-noising spatially resolved transcriptomics data based on position and image information. *Nat Methods* **19**: 950–958. doi:10.1038/s41592-022-01560-w
- Weinstein JN, Collisson EA, Mills GB, Shaw KRM, Ozenberger BA, Ellrott K, Shmulevich I, Sander C, Stuart JM. 2013. The cancer genome atlas pan-cancer analysis project. *Nat Genet* **45**: 1113–1120. doi:10.1038/ng.2764
- Welch JD, Kozareva V, Ferreira A, Vanderburg C, Martin C, Macosko EZ. 2019. Single-cell multi-omic integration compares and contrasts features of brain cell identity. *Cell* **177**: 1873–1887.e17. doi:10.1016/j.cell.2019.05.006
- Williams CG, Lee HJ, Asatsuma T, Vento-Tormo R, Haque A. 2022. An introduction to spatial transcriptomics for biomedical research. *Genome Med* **14**: 68. doi:10.1186/s13073-022-01075-1
- Xu H, Fu H, Long Y, Ang KS, Sethi R, Chong K, Li M, Uddamvathanak R, Lee HK, Ling J, et al. 2024. Unsupervised spatially embedded deep representation of spatial transcriptomics. *Genome Med* **16**: 12. doi:10.1186/s13073-024-01283-x
- You Y, Fu Y, Li L, Zhang Z, Jia S, Lu S, Ren W, Liu Y, Xu Y, Liu X, et al. 2024. Systematic comparison of sequencing-based spatial transcriptomic methods. *Nat Methods* **21**: 1743–1754. doi:10.1038/s41592-024-02325-3
- Zeisel A, Muñoz-Manchado AB, Codeluppi S, Lönnerberg P, La Manno G, Jureus A, Marques S, Munguba H, He L, Bethsholtz C, et al. 2015. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science* **347**: 1138–1142. doi:10.1126/science.aaa1934
- Zeng Y, Song Y, Zhang C, Li H, Zhao Y, Yu W, Zhang S, Zhang H, Dai Z, Yang Y. 2024. Imputing spatial transcriptomics through gene network constructed from protein language model. *Commun Biol* **7**: 1271. doi:10.1038/s42003-024-06964-2
- Zhang M, Eichhorn SW, Zingg B, Yao Z, Cotter K, Zeng H, Dong H, Zhuang X. 2021. Spatially resolved cell atlas of the mouse primary motor cortex by MERFISH. *Nature* **598**: 137–143. doi:10.1038/s41586-021-03705-x
- Zhang Z, Wayment-Steele H, Brixi G, Wang H, Kern D, Ovchinnikov S. 2024. Protein language models learn evolutionary statistics of interacting sequence motifs. *Proc Natl Acad Sci* **121**: e2406285121. doi:10.1073/pnas.2406285121
- Zhuo W, Jia L, Song N, Lu X-A, Ding Y, Wang X, Song X, Fu Y, Luo Y. 2012. The CXCL12–CXCR4 chemokine pathway: a novel axis regulates lymphangiogenesis. *Clin Cancer Res* **18**: 5387–5398. doi:10.1158/1078-0432.CCR-12-0708

Received June 2, 2025; accepted in revised form January 21, 2026.