# Auxiliary Gene Learning:
# Spatial Gene Expression Estimation by Auxiliary Gene Selection

**Kaito Shiku[1], Kazuya Nishimura[2], Shinnosuke Matsuo[1], Yasuhiro Kojima[2], Ryoma Bise[1]**

[1] Department of Advanced Information Technology, Kyushu University, Japan
[2] Laboratory of Computational Life Science, National Cancer Center Japan
kaito.shiku@human.ait.kyushu-u.ac.jp

## Abstract

Spatial transcriptomics (ST) is a novel technology that enables the observation of gene expression at the resolution of individual spots within pathological tissues. ST quantifies the expression of tens of thousands of genes in a tissue section; however, heavy observational noise is often introduced during measurement. In prior studies, to ensure meaningful assessment, both training and evaluation have been restricted to only a small subset of highly variable genes, and genes outside this subset have also been excluded from the training process. However, since there are likely co-expression relationships between genes, low-expression genes may still contribute to the estimation of the evaluation target. In this paper, we propose *Auxiliary Gene Learning* (AGL) that utilizes the benefit of the ignored genes by reformulating their expression estimation as auxiliary tasks and training them jointly with the primary tasks. To effectively leverage auxiliary genes, we must select a subset of auxiliary genes that positively influence the prediction of the target genes. However, this is a challenging optimization problem due to the vast number of possible combinations. To overcome this challenge, we propose Prior-Knowledge-Based Differentiable Top-$k$ Gene Selection via Bi-level Optimization (DkGSB), a method that ranks genes by leveraging prior knowledge and relaxes the combinatorial selection problem into a differentiable top-$k$ selection problem. The experiments confirm the effectiveness of incorporating auxiliary genes and show that the proposed method outperforms conventional auxiliary task learning approaches.

**Code** — https://github.com/Shiku-Kaito/AGL

## Introduction

Spatial transcriptomics (ST) is a novel technology that enables the observation of gene expression at the resolution of individual spots within pathological tissues, playing a crucial role in evaluating disease progression and drug efficacy (Ståhl et al. 2016). While spatial gene expression data obtained through ST provides detailed insights into molecular activity within pathological tissue, its high observational cost has prompted the development of neural network models that aim to estimate spatial gene expression from pathological images (He et al. 2020; Pang, Su, and Li 2021; Yuansong et al. 2022; Yang et al. 2023, 2024).

ST quantifies the expression of tens of thousands of genes in a tissue section. Yet, for most genes, the observed counts are extremely low and often contain dropout noise (Mejia et al. 2024, 2023) (zero or near-zero measurements that arise from technical limitations rather than true absence of expression). Because these noisy readings are statistically unreliable, they cannot serve as dependable reference values for benchmarking model performance. To ensure meaningful assessment, earlier studies have therefore restricted both training and evaluation to a narrow set of highly variable genes (Zheng et al. 2017; Satija et al. 2015; Stuart et al. 2019) whose expression consistently rises above the noise floor and is considered trustworthy. In practice, the evaluation target genes consist of only a small subset of the total genes, while the remaining tens of thousands are entirely excluded from performance evaluation (Pang, Su, and Li 2021; Yuansong et al. 2022; Chung et al. 2024; He et al. 2020).

Although the specific relationships have yet to be thoroughly investigated, many of the genes excluded from training may still share regulatory (M. Ribeiro, Ziyani, and Delaneau 2022) or co-expression patterns (Wang, Maletic-Savatic, and Liu 2022) with the target genes used for evaluation. Even when their individual measurements are noisy or sparse, these genes can still provide useful contextual signals that help the model learn richer representations for the evaluation target genes, thereby improving prediction accuracy.

We propose *Auxiliary Gene Learning* (AGL) that utilizes the benefit of the ignored genes by reformulating their expression estimation as auxiliary tasks and training them jointly with the primary tasks, as illustrated in Figure 1 (a). An auxiliary task is not evaluated directly, yet learning it together with the primary task can improve the primary-task accuracy (Figure 1 (b)). In our case, the primary tasks are the expression predictions for the small subset of highly variable genes, whereas the auxiliary tasks correspond to the remaining genes that were previously discarded. Introducing this auxiliary task framework into the problem of predicting gene expression from pathological images is, to our knowledge, the first attempt of its kind. Consequently, our study is the first to make practical use of information that earlier work had always thrown away.

However, using all remaining genes as auxiliary targets may be counter-productive, because many of them have very low counts or are dominated by measurement noise. There-
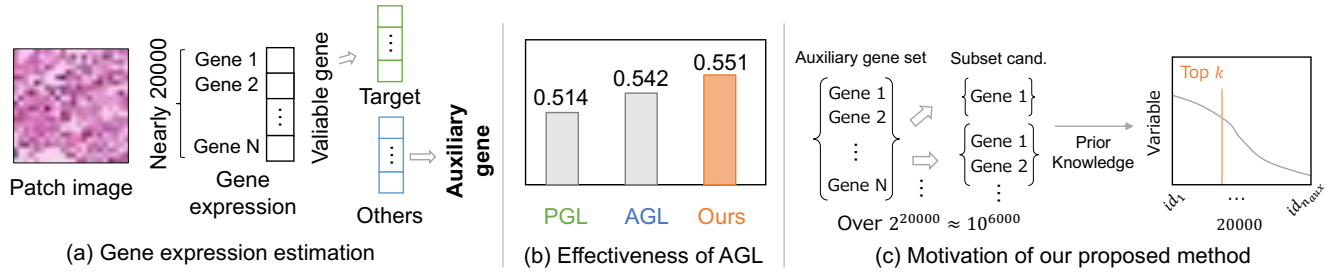
Figure 1: (a) Conventional gene expression estimation focuses solely on predicting primary genes, typically ignoring the remaining ones. In this study, we treat these overlooked genes as auxiliary genes. (b) Effectiveness of *Auxiliary Gene Learning* (AGL). PGL denotes primary gene learning, which uses only the target gene for training. AGL represents our auxiliary gene learning, which jointly estimates primary genes and previously ignored auxiliary genes, selecting auxiliaries via a differentiable cut-off. (c) Illustration of our top-k gene selection approach. As the number of possible subsets exceeds $10^{6000}$, we relax this combinatorial selection into a top-$k$ problem by leveraging prior knowledge of gene-expression signal quality.

fore, we must select a subset of auxiliary genes that positively influence the prediction of the evaluation genes.

As Figure 1 (c) shows, nearly twenty thousand genes are available as auxiliary-task candidates. Selecting an appropriate subset from such a vast pool leads to a combinatorial explosion and is far from trivial. Existing work on selective auxiliary-task learning assumes at most five to ten candidate tasks and is designed to pick one or a few from that small set (Kendall, Gal, and Cipolla 2018; Aviv et al. 2021; Ko et al. 2023; Jiang et al. 2023; Sivasubramanian et al. 2023). Approaches developed under those conditions do not transfer well to the much larger search space we face here, a limitation that our experiments later confirm.

To overcome this large-scale selection challenge, we introduce Prior-Knowledge-Based Differentiable Top-$k$ Gene Selection via Bi-level Optimization (DkGSB). DkGSB first ranks all auxiliary-gene candidates by their expression variance, which serves as a simple yet effective proxy for signal quality. Instead of searching over every subset, the method learns a single scalar $k$ that determines how many of the top-ranked genes will remain as auxiliary tasks. A soft and differentiable relaxation makes this top-$k$ operator compatible with gradient descent, and a bi-level objective updates $k$ together with the network weights so that the selected auxiliary genes maximize the accuracy of the target genes.

In experiments conducted on publicly available datasets (Jaume et al. 2024), the proposed AGL approach was shown to outperform conventional methods that discard lowly expressed genes. Furthermore, under the setting where auxiliary genes are used during training, we demonstrate that the proposed method, which incorporates soft parameterization and bi-level optimization with prior knowledge, outperforms conventional auxiliary task learning approaches and also achieves better performance than using all auxiliary genes without selection.

## Related work

### Gene expression estimation for pathological images

With recent advances in observation technologies enabling the measurement of gene expression at the resolution of individual spots within pathological tissue, several studies have

accordingly explored approaches to estimate gene expression from spot-level pathological images (He et al. 2020; Dawood et al. 2021; Xie et al. 2023; Pang, Su, and Li 2021; Yang et al. 2023, 2024; Chung et al. 2024; Nishimura et al. 2025). Research on gene expression estimation from spot images has primarily followed two main streams: one in which gene expression is predicted independently for each spot image, and another that incorporates surrounding contextual information into the prediction. Among the methods that predict gene expression independently for each spot image, the most popular is ST-Net (He et al. 2020), which employs a CNN-based backbone and formulates gene expression estimation as a multi-output regression task. Other approaches that have been proposed include stain-aware prediction methods (Dawood et al. 2021) that perform stain deconvolution on H&E images, and joint embedding methods (Xie et al. 2023) that learn a shared representation of spot images and gene expression, similar to CLIP (Radford et al. 2021). In studies that incorporate surrounding context into prediction, methods employing Transformers (Pang, Su, and Li 2021), graph neural networks (Yang et al. 2023, 2024), and multi-scale feature embeddings (Chung et al. 2024) have achieved state-of-the-art performance by modeling interactions between spots within the WSI, outperforming approaches that treat spot-level predictions independently.

However, these studies commonly train models using only the top-expressing genes with stable expression levels, while discarding low-expressing genes, thereby overlooking potentially informative signals that these discarded genes may provide for learning. To the best of our knowledge, this is the first attempt to utilize low-expressing genes as auxiliary supervision in the training process. Furthermore, our proposed auxiliary gene selection method is model-agnostic, making it easily pluggable into existing approaches.

### Auxiliary task learning

Auxiliary Task Learning (ATL) aims to improve the performance of a primary task by leveraging information from related auxiliary tasks. In ATL, various methods have been proposed to select beneficial tasks from multiple auxiliary tasks, such as weighting based on gradient similarity (Chen et al. 2018), task-wise uncertainty (Kendall, Gal, and Cipolla

2018), and techniques like ForkMerge (Jiang et al. 2023), which create a branch for each auxiliary task combined with the primary task, and aggregate the model parameters through weighted averaging based on the performance improvement rate on the primary task. In recent years, methods that aggregate auxiliary task losses have become mainstream in the field, including approaches that assign loss weights using learnable parameters (Sivasubramanian et al. 2023), methods that integrate losses non-linearly using MLPs (Aviv et al. 2021), and techniques that leverage Transformers to combine losses (Ko et al. 2023).

In typical ATL settings, the number of auxiliary tasks is assumed to range from 1 to 10, or up to around 300 in large-scale cases (Aviv et al. 2021). In contrast, our proposed auxiliary gene learning (AGL) involves approximately 20,000 auxiliary tasks, making the selection of optimal auxiliary tasks significantly difficult. Moreover, due to computational cost constraints, methods such as training on all pairs of the primary task and each auxiliary task (Jiang et al. 2023), or using Transformers whose computational complexity increases exponentially with the number of tokens (Ko et al. 2023), are practically infeasible to apply to AGL.

## Preliminary: general formulation of subset auxiliary task selection

We first summarize the general formulation of subset selection. The aim is to choose, from $n_{\mathrm{aux}}$ auxiliary-task candidates, the subset that most improves the performance of the primary tasks evaluated at test time. The problem can be written as:

$$\boldsymbol{\lambda}^{\star} = \operatorname*{arg\,min}_{\boldsymbol{\lambda} \in \{0,1\}^{n_{\mathrm{aux}}}} \sum_{j=1}^{n_{\mathrm{pri}}} L_{\mathrm{pri},j}\big(\theta^{\star}(\boldsymbol{\lambda})\big),$$

$$\text{s.t.} \quad \theta^{\star}(\boldsymbol{\lambda}) = \operatorname*{arg\,min}_{\theta}\Big[\sum_{j=1}^{n_{\mathrm{pri}}} L_{\mathrm{pri},j}(\theta) + \sum_{j=1}^{n_{\mathrm{aux}}} \lambda_j \, L_{\mathrm{aux},j}(\theta)\Big],$$

$$(1)$$

where $L_{\mathrm{pri},j}$ and $L_{\mathrm{aux},j}$ denote the losses for the $j$-th primary and auxiliary task, respectively; $n_{\mathrm{pri}}$ and $n_{\mathrm{aux}}$ are the numbers of primary and auxiliary tasks. The binary mask $\boldsymbol{\lambda} \in \{0,1\}^{n_{\mathrm{aux}}}$ specifies which auxiliary genes are selected ($\lambda_j = 1$) or discarded ($\lambda_j = 0$).

However, as noted in the introduction, the number of auxiliary tasks in our setting is enormous ($n_{\mathrm{aux}} \approx 20{,}000$). The corresponding search space $\binom{n_{\mathrm{aux}}}{k}$ grows exponentially, so solving Eq. (1) exhaustively is infeasible.

## Auxiliary Gene Learning with Differentiable Top-$k$ Gene Selection

Figure 2 outlines our *Auxiliary Gene Learning* (AGL) framework with the Prior-Knowledge-Based Differentiable Top-$k$ Gene Selection (DkGSB) module. The procedure has three steps: (i) all $n_{\mathrm{aux}} \approx 20{,}000$ auxiliary genes are ranked once by a variance-based score (Section 4.2); (ii) a single learnable scalar $k$ defines a soft top-$k$ mask $\boldsymbol{\lambda}(k)$, obtained through a differentiable relaxation of the hard cut-off; (iii)
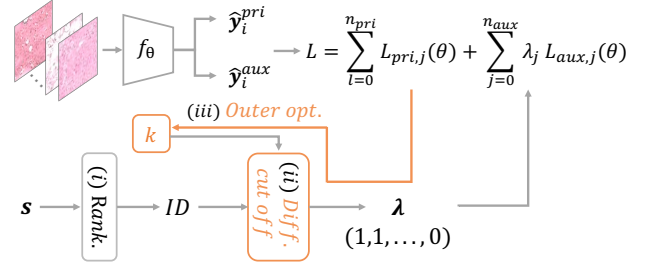


Figure 2: Overview of Proposed **DkGSB**. The procedure has three steps: (i) auxiliary genes are ranked based on a variance-based score $\boldsymbol{s}$; (ii) a single learnable scalar $k$ defines a soft top-$k$ mask $\boldsymbol{\lambda}(k)$, obtained through a differentiable relaxation of the hard cut-off; (iii) $k$ is optimized together with the network weights by a bi-level scheme.

$k$ is optimized together with the network weights by a bi-level scheme (Section 4.3). This reduces the combinatorial search to a one-dimensional optimization and retains only auxiliaries that improve primary gene prediction.

### Problem setup

Let $\mathcal{D} = \big\{(\mathbf{x}_i, \mathbf{y}_i^{\mathrm{pri}}, \mathbf{y}_i^{\mathrm{aux}})\big\}_{i=1}^{N}$ be the spatial gene expression dataset, where $\mathbf{x}_i$ is the $i$-th spot image, $\mathbf{y}_i^{\mathrm{pri}} \in \mathbb{R}^{n_{\mathrm{pri}}}$ contains the expressions of the $n_{\mathrm{pri}}$ *primary* genes, and $\mathbf{y}_i^{\mathrm{aux}} \in \mathbb{R}^{n_{\mathrm{aux}}}$ stores the remaining $n_{\mathrm{aux}} \approx 20{,}000$ *auxiliary* genes. Note that $n_{\mathrm{pri}} \ll n_{\mathrm{aux}}$ in this setting. A neural network $f_\theta$ predicts both sets of expressions from an image, $(\hat{\mathbf{y}}_i^{\mathrm{pri}}, \hat{\mathbf{y}}_i^{\mathrm{aux}}) = f_\theta(\mathbf{x}_i)$.

After ranking, only the top $k$ auxiliary genes are kept; the soft mask $\boldsymbol{\lambda}(k) \in [0,1]^{n_{\mathrm{aux}}}$ returned by the relaxation assigns high weights to genes with rank $\leq k$ and small weights otherwise. During training, the network weights $\theta$ are updated with losses from both the primary genes and the masked auxiliary genes, while the scalar $k$ is adjusted so that the validation loss of the primary genes decreases. The formal bi-level objective is given in Section 4.3.

### Prior knowledge-based auxiliary gene ranking

In gene expression analysis, it is widely accepted that genes exhibiting high variability relative to their mean expression within biological tissues provide more informative signals (Satija et al. 2015). We therefore rank all auxiliary genes by a *highly variable gene* (HVG) score, which corrects for the dependence of dispersion on mean expression.

To compute the HVG score, we calculate the mean expression $\mu_j$ and the raw dispersion $\delta_j$ for the $j$-th gene across the entire dataset, where the raw dispersion $\delta_j$ is obtained by dividing the variance by the mean expression. All genes are divided into twenty bins based on their mean expression $\mu_j$, and within each bin, the dispersions are $z$-score normalized to yield scale-independent HVG scores

$$s_j = \text{z-score}\big(\delta_j\big), \qquad j = 1, \ldots, n_{\mathrm{aux}}. \qquad (2)$$

We then sort the genes in descending order of their HVG

Algorithm 1: Prior Knowledge-Based Differentiable Top-$k$ Gene Selection via Bi-level Optimization.

---

**Require:** Training data $\mathcal{D}^t$, Validation data $\mathcal{D}^v$, Temperature $\tau$, Learning rates $(\alpha, \beta)$, Inner steps $H$.

1: **while** not Converge **do**
2:    **while** not Converge **do**
3:      $\{\mathbf{x}_i, \mathbf{y}_i\} \leftarrow SampleMiniBatch(\mathcal{D}^t)$
4:      Obtain $\tilde{\boldsymbol{\lambda}}(k, \tau)$ by Eq.( 6).
5:      $(\hat{\mathbf{y}}_i^{\mathrm{pri}}, \hat{\mathbf{y}}_i^{\mathrm{aux}}) = f_\theta(\mathbf{x}_i)$
6:      Update: $\theta \leftarrow \theta - \alpha \nabla_\theta \left( \sum_{j=1}^{n_{\mathrm{pri}}} L_{\mathrm{pri},j}(\theta) + \sum_{j=1}^{n_{\mathrm{aux}}} \tilde{\lambda}_j(k) L_{\mathrm{aux},j}(\theta) \right)$
7:    **end while**
8:    **while** not Converge **do**
9:      $\{\mathbf{x}_i, \mathbf{y}_i\} \leftarrow SampleMiniBatch(\mathcal{D}^t)$
10:     Obtain $\tilde{\boldsymbol{\lambda}}(k, \tau)$ by Eq.( 6).
11:     **for** $h = 1, ..., H$ **do**
12:       Update: $\theta^+ \leftarrow \theta - \alpha \nabla_\theta \left( \sum_{j=1}^{n_{\mathrm{pri}}} L_{\mathrm{pri},j}(\theta) + \sum_{j=1}^{n_{\mathrm{aux}}} \tilde{\lambda}_j(k) L_{\mathrm{aux},j}(\theta) \right)$
13:     **end for**
14:     $\{\mathbf{x}_i, \mathbf{y}_i\} \leftarrow SampleMiniBatch(\mathcal{D}^v)$
15:     Update: $k \leftarrow k - \beta \nabla_k \sum_{j=1}^{n_{\mathrm{pri}}} L_{\mathrm{pri},j}(\theta^+)$
16:    **end while**
17: **end while**

---

scores,
$$ID = (id_1, id_2, \ldots, id_{n_{\mathrm{aux}}}), s_{id_1} > s_{id_2} > \cdots > s_{id_{n_{\mathrm{aux}}}}, \tag{3}$$
and pass this ranked list to the differentiable top-$k$ selection module described in the next subsection.

## Differentiable top-$k$ selection via bi-level optimization

Based on gene ranking results, we find the optimal cut-off $k$ that selects the top-ranked auxiliary genes most helpful for predicting the primary genes. The ranking itself is fixed by the HVG scores; only $k$ is optimized. Our AGL optimization is reformulated as:

$$\underset{k \in \{1, \ldots, n_{\mathrm{aux}}\}}{\arg\min} \sum_{j=1}^{n_{\mathrm{pri}}} L_{\mathrm{pri},j}(\theta^\star(k)),$$

$$\text{s.t.} \quad \theta^\star(k) = \underset{\theta}{\arg\min} \left[ \sum_{j=1}^{n_{\mathrm{pri}}} L_{\mathrm{pri},j}(\theta) + \sum_{j=1}^{n_{\mathrm{aux}}} \lambda_j(k) L_{\mathrm{aux},j}(\theta) \right], \tag{4}$$

where $\lambda_j$ is 1 if $id_j \le k$ and 0 otherwise.

In this study, we use the Pearson correlation coefficient loss function for $L_{\mathrm{pri}}, L_{\mathrm{aux}}$, which is known to reduce batch-effect–related scaling bias.

$$L = 1 - \frac{\sum_{i=1}^{N} (\hat{y}_i - \hat{Y})(y_i - Y)}{\sqrt{\sum_{i=1}^{N} (\hat{y}_i - \hat{Y})^2} \sqrt{\sum_{i=1}^{N} (y_i - Y)^2}}, \tag{5}$$

where $N$ denotes the mini-batch [1] size, and $Y$ and $\hat{Y}$ represent the mean ground-truth and predicted expression values within the mini-batch, respectively.

---

[1] In this paper, we refer to a patient as a "batch," and a subset of data used for a single model-update step as a "mini-batch."

Because the ordinary top-$k$ operator is not differentiable, we replace it with a temperature-controlled soft mask so that the cut-off can be updated by gradient descent. The resulting mask $\tilde{\boldsymbol{\lambda}}(k, \tau)$ is given by

$$\tilde{\lambda}_j = \frac{\exp(k/\tau)}{\exp(id_j/\tau) + \exp(k/\tau)}, \qquad j = 1, \ldots, n_{\mathrm{aux}}, \tag{6}$$

where $id_j$ is the rank of the $j$-th auxiliary gene, $k$ is the learnable threshold, and $\tau$ controls the softness of the step. As $\tau \to 0$, the mask approaches the hard top-$k$ selection $\mathbb{1}[\, id_j \le k\,]$.

Algorithm 1 learns the network weights $\theta$ and the cut-off $k$ in a bi-level manner. At the start of each outer iteration, $k$ is fixed and an inner loop runs on the training data $\mathcal{D}^t$. For each training mini-batch, the current $k$ is converted into a soft mask $\tilde{\boldsymbol{\lambda}}(k, \tau)$, the mini-batch is forward-propagated, and $\theta$ is updated with the total loss that combines the primary terms with the auxiliary terms weighted by $\tilde{\boldsymbol{\lambda}}$. This sequence is repeated a fixed number of times, denoted $H$, so that $\theta$ is fitted to the current auxiliary subset.

After these $H$ weight updates, a validation mini-batch from $\mathcal{D}^v$ is used to refine $k$. The primary-gene validation loss is back-propagated, and a single gradient step is taken on $k$ through the differentiable mask $\tilde{\boldsymbol{\lambda}}$. Because $\tilde{\boldsymbol{\lambda}}$ is a smooth function of $k$, standard optimizers with learning rate $\beta$ can be applied directly. By alternating training-data updates of $\theta$ ($H$ steps) with validation-data updates of $k$ (one step), the algorithm converges to the value of $k$ that minimises the primary-gene validation error, reducing an otherwise combinatorial subset search to a one-dimensional optimization learned end-to-end with the model.

# Experiments

Spatial gene expression data are frequently affected by batch effects, which arise from differences in experimental conditions or technical variation introduced by the operator during data acquisition. These effects often lead to significant biases in gene expression measurements across different slides (Lopez et al. 2018; Shaham et al. 2017). In this experiment, our objective is to analyze the impact of auxiliary genes on the estimation performance of target genes; however, in the presence of batch effects, there is a risk that the effects of auxiliary genes may not be accurately evaluated. Therefore, we conduct experiments based on two settings: an intra-batch experiment, where batch effects are absent, and an inter-batch experiment, where training and evaluation are performed across different batches.

## Dataset

For both intra-batch and inter-batch experimental settings, we used data from the Hest-1k dataset (Jaume et al. 2024), a large-scale public dataset comprising paired spatial gene expression and pathological images. The details of each experimental setting are as follows.
**Intra-batch experiment.** We used slides from two bowels (**BOWEL A**, **BOWEL B**) and one ovary (**OVARY**) organs, along with spatial gene expression data acquired using Visium technology (Williams et al. 2022). After quality control,

Table 1: **Comparison with conventional method** includes both the baseline setting without auxiliary genes and comparisons with existing auxiliary task learning methods under the "**AGL**" setting. The performance of each method is evaluated through cross-validation, and the reported values represent the mean and standard deviation. Best performances are bold.

| Method | | Intra-batch | | | Inter-batch | Average |
|---|---|---|---|---|---|---|
| | | BOWEL A | BOWEL B | OVARY | HEART | |
| PGL | | 0.514±0.009 | 0.419±0.004 | 0.448±0.008 | 0.245±0.034 | 0.407 |
| **AGL** | + ALL | 0.542±0.008 | 0.430±0.006 | 0.451±0.005 | 0.248±0.038 | 0.418 |
| **AGL** | + Uncertainty | 0.543±0.009 | 0.431±0.005 | 0.451±0.006 | 0.252±0.038 | 0.419 |
| **AGL** | + AuxLearn | 0.541±0.007 | 0.411±0.012 | 0.445±0.008 | 0.251±0.037 | 0.412 |
| **AGL** | + AMAL | 0.535±0.010 | 0.416±0.005 | 0.443±0.007 | 0.248±0.039 | 0.411 |
| **AGL** | **+ DkGSB** | **0.551**±0.009 | **0.440**±0.008 | **0.458**±0.006 | **0.256**±0.039 | **0.426** |

**BOWEL A**, **BOWEL B**, and **OVARY** contain 4,096, 4,617, and 5,774 patch images of size 224×224, respectively, each cropped at the center of a spot with a width of 55 $\mu$m. The corresponding expression data include 18,066, 18,054, and 18,043 gene types, respectively. As a preprocessing step, the expression values for each spot were log-normalized. Following the experimental setting described in the Hest-1k dataset paper and other studies using this dataset (Jaume et al. 2024; Cho et al. 2025), we also use the top 50 highly variable genes as prediction targets. We split the patch images into five folds using a 3:1:1 ratio for the training, validation, and test sets on each slide, and performed 5-fold cross-validation.[2]

**Inter-batch experiment.** We used four slides from the heart organ (**HEART**), along with spatial gene expression data acquired using Visium technology. Following (Xie et al. 2023), we performed leave-one-batch-out cross-validation by splitting the data into training and evaluation sets based on batch identity. After quality control, each slide contained 2,857, 3,400, 3,236, and 3,042 patch images of size 224×224, respectively. We selected genes of the 15,904 type that were commonly present across all batches as training targets.

### Implementation details and evaluation metric

We implemented the proposed method using PyTorch (Adam et al. 2019). As the feature extractor for $f_\theta$, we employed a ResNet18 (He et al. 2016) model pretrained on ImageNet (Deng et al. 2009). To train the network, we used the Adam optimizer (Kingma and Ba 2014) with a learning rate of $\alpha = 3 \times 10^{-5}$, a mini-batch size of 128, and trained the model for up to 1,000 epochs with early stopping set to 20 epochs. In the bi-level optimization, the number of inner steps was set to $H = 2$, and the learning rate $\beta$ was set to $3 \times 10^{-3}$. The temperature parameter $\tau$ of the differentiable cut-off is set to 0.01. The proposed method is trained on a system equipped with an Intel Xeon Gold 5122 CPU and an NVIDIA RTX A6000 GPU.

We evaluated the performance of the proposed method based on the Pearson correlation coefficient (**PCC**) (Chung et al. 2024), which is calculated as the correlation between the predicted and ground-truth gene expression values. The

reported value is the average computed over all target genes after calculating the PCC for each gene within the test dataset.

### Comparison

To demonstrate the effectiveness of our *Auxiliary Gene Learning* (**AGL**) framework, we compared the performance of the following methods: 1) "Primary gene learning (PGL)" trains the model using only the primary genes, without incorporating any auxiliary genes. 2) "**AGL**+ All" trains the model under the proposed AGL setting, but without any auxiliary–task selection: every gene other than the primary is included as an auxiliary task. This baseline tests whether simply adding all remaining genes can improve primary-gene prediction without task selection. 3) "**AGL**+ Uncertainty" (Kendall, Gal, and Cipolla 2018): Instead of selecting a subset of auxiliaries, this variant keeps all auxiliary genes but assigns each task a weight that is learned from uncertainty: tasks with lower predictive uncertainty receive higher weights, while highly uncertain tasks are down-weighted. 4) "**AGL**+ AuxLearn" (Aviv et al. 2021): This variant retains all auxiliary genes and combines their losses through *AuxLearn*, which feeds the individual losses into a small neural network that predicts a non-linear weighting for each task according to its estimated contribution to the target task. 5) "**AGL**+ AMAL": Based on the adaptive multi-task weighting strategy AMAL (Sivasubramanian et al. 2023), this variant attaches one learnable scalar to every auxiliary task and multiplies that scalar by the task's loss. A temperature-controlled sigmoid keeps each scalar in the interval $[0, 1]$; values close to 1 retain the task, whereas values near 0 effectively drop it, so the mechanism performs soft task selection driven by the contribution of each auxiliary gene to the primary loss. 6) "**AGL**+DkGSB(ours)": Uses the proposed differentiable top-$k$ mask on the HVG ranking; $k$ is learned via the bi-level scheme in Section 4.3.

For a fair comparison, we adopt ST-Net (He et al. 2020), which is the most widely used architecture for gene-expression prediction, as the common backbone in all experiments. Our framework is model-agnostic and could be applied to any ST predictor, but fixing the backbone isolates the effect of the auxiliary-loss strategies. Apart from the way auxiliary losses are combined, every competing method

---

[2]No data leakage occurs when updating the parameter $k$, as the validation set used is completely independent of the test set.
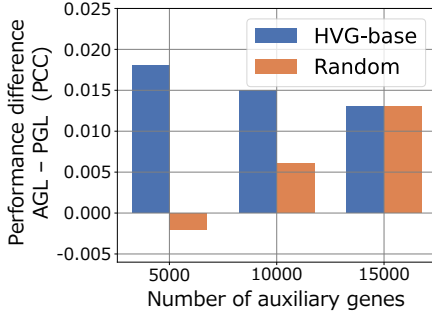
Figure 3: **Reasonability of HVG score-based selection.** Performance of primary gene expression estimation using HVG score–based selection (blue) and random selection (orange). The vertical axis shows the performance difference between models trained with auxiliary genes selected by each method and those trained using only primary genes ("PGL"), while the horizontal axis indicates the number of selected auxiliary genes. The experiments were conducted using **BOWEL B**.

shares the same network architecture, optimization schedule, and hyper-parameters used in AGL.

**Intra-batch experimentu.** BOWEL A, **BOWEL B**, and **OVARY** in Table 1 presents the PCC scores for the three intra-batch tissues. Training on the 50 primary genes alone ("PGL") yields the lowest accuracy throughout. Adding all auxiliary genes without selection ("**AGL** + All") improves performance across all datasets, indicating that auxiliary signals are broadly beneficial. Weighting losses by predictive uncertainty ("**AGL** + Uncertainty") yields only a modest performance gain, while the adaptive weighting schemes "**AGL** + AuxLearn" and "**AGL** + AMAL" offer no consistent improvement: optimizing weights for nearly 20,000 tasks proves difficult without additional guidance. The proposed "**AGL+DkGSB**", which replaces exhaustive subset search with a single learnable HVG-based cut-off, achieves superior performance over all other methods across all tissue types.

**Inter-batch experiment.** In the inter-batch experiment on the **HEART** dataset in Table 1, all methods showed decreased performance compared to those in the intra-batch datasets. This result also suggests that batch effects have a substantial impact on the performance of spatial gene expression prediction. Even under this setting, applying "**AGL**" successfully outperforms "PGL", and further performance improvement is achieved by selecting auxiliary genes using the proposed "**DkGSB**". This result confirms that our approach retains its effectiveness even in the presence of batch effects.

## Analysis

**Reasonability of HVG score-based selection.** To demonstrate the reasonability of selecting auxiliary genes based on HVG scores, we compared the performance of primary gene estimation using HVG score-based selection and random selection, with 5,000, 10,000, and 15,000 auxiliary genes, respectively.

Figure 3 shows, on the vertical axis, the performance difference between models trained with a subset of selected auxiliary genes and those trained using only primary genes
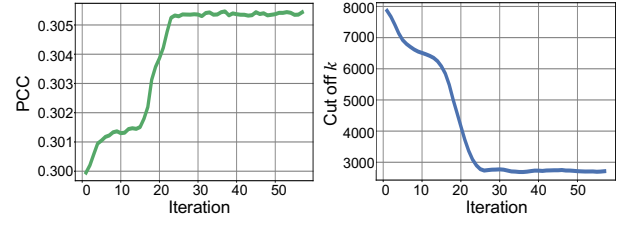


Figure 4: **Behavior during cut-off $k$ optimization.** The left panel shows the changes in validation performance during the optimization process of the *outer* loop, while the right panel shows the changes in the cut-off $k$. The experiments were conducted using **HEART** dataset.

("PGL"), and on the horizontal axis, the number of selected auxiliary genes. Values greater than 0.0 indicate an improvement in performance, while values below 0.0 represent a decrease in performance. The blue bars represent the HVG score-based selection method, while the orange bars represent the random selection method. The experiments were conducted using **BOWEL B**.

Experimental results show that when the number of auxiliary genes is 5,000 or 10,000, the HVG score-based selection yields a significantly greater performance improvement compared to the random selection method. In particular, at 5,000 auxiliary genes, the random selection method even leads to a decrease in performance. When the number of auxiliary genes approaches the total number of genes (15,000), both methods yield similar levels of performance improvement. These results suggest that the HVG score-based selection method is capable of preferentially selecting auxiliary genes that contribute to improving the prediction performance of the target genes.

**Behavior during cut-off $k$ optimization.** Figure 4 illustrates how the bi-level algorithm optimizes the cut-off $k$ on the **HEART** dataset. The *left graph* traces the validation Pearson correlation (PCC) of the primary genes over successive outer iterations. Beginning with the full mask ($k = n_{\text{aux}}$), the PCC rises monotonically for about 15–20 outer steps and then stabilises at $\approx 0.305$. The *right graph* shows the simultaneous evolution of the cut-off $k$. Over the same 20 iterations the algorithm lowers $k$ from the full auxiliary pool to $k = 2,698$, after which both $k$ and the PCC remain essentially flat. Since the **HEART** dataset contains roughly 15,000 auxiliary genes, the final threshold retains only $\sim 18\%$ of the candidates while discarding the remaining $82\%$. The fact that the accuracy curve has already reached its plateau when $k$ converges confirms that the optimizer has settled on a subset that is near-optimal for this tissue.

**Visualization of the expression levels for the selected auxiliary genes.** To compare the expression patterns of auxiliary genes selected by the proposed "**DkGSB**" and the conventional auxiliary task learning method "**AGL** +AMAL", Figure 5 visualizes the expression of the selected genes on the **HEART** dataset. The top row shows genes selected by "**AGL+DkGSB**" but not by "**AGL+AMAL**," and the bottom row shows the reverse. The name of each visualized gene appears above its slide.

Table 2: **Robustness to the number of primary genes.** The performance of "PGL", "**AGL + All**", and "**AGL + DkGSB**" when the number of primary genes is varied among 25, 50, 75, and 100. The performance of each method is evaluated via cross-validation on the **BOWEL A** dataset. Reported values indicate the mean and standard deviation, with the best performances highlighted in bold.

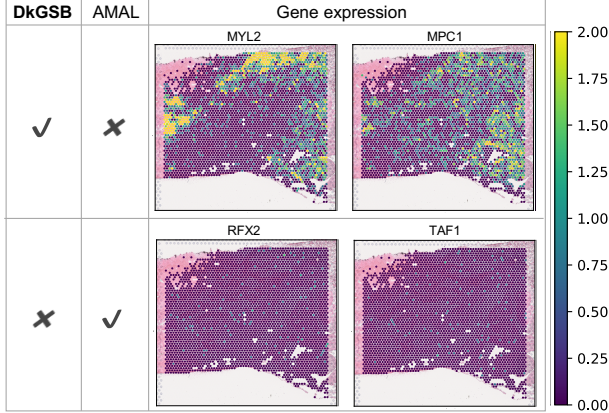| Method | 25 | 50 | 75 | 100 | Average |
|---|---|---|---|---|---|
| PGL | 0.445±0.009 | 0.514±0.009 | 0.542±0.008 | 0.561±0.009 | 0.516 |
| **AGL** + ALL | 0.481±0.009 | 0.542±0.008 | 0.557±0.007 | 0.571±0.008 | 0.538 |
| **AGL** + **DkGSB** | **0.487**±0.008 | **0.551**±0.009 | **0.564**±0.008 | **0.579**±0.008 | **0.545** |



Figure 5: **Visualization of the expression levels for the selected auxiliary genes.** In the top row, the expression patterns of genes that were selected by the proposed method but not by "**AGL+AMAL**," and in the bottom row, the opposite: genes that were not selected by the proposed method but were selected by "**AGL+AMAL**." The name of the visualized gene is shown at the top of each slide.

As shown in the top row, the proposed "**AGL+DkGSB**" tends to select highly expressed genes as auxiliary genes, whereas, as shown in the bottom row, "**AGL+AMAL**" selects genes that are barely expressed across the entire slide and are unlikely to contribute effectively to the learning process. This result suggests that "**AGL+AMAL**" is unable to effectively select appropriate gene combinations because it directly tackles a high-dimensional combinatorial optimization problem. **Robustness to the number of primary genes.** Although the Hest-1k dataset is commonly used with the top 50 highly variable genes as prediction targets (Jaume et al. 2024; Cho et al. 2025), whose expression levels consistently rise above the noise floor and are considered reliable, we conducted experiments by varying the number of primary genes to demonstrate the robustness of the proposed "**AGL+DkGSB**" to such changes. Table 2 presents the performance of "PGL", "**AGL + All**", and "**AGL + DkGSB**" on the **BOWEL A** dataset when the number of primary genes is set to 25, 50, 75, and 100. These results indicate that even when the number of primary genes varies, leveraging previously discarded genes as auxiliary genes remains effective. Moreover, the proposed "**DkGSB**" consistently improves performance regardless of the number of primary genes. These findings demonstrate that the proposed "**AGL+DkGSB**" is robust to variations in the number of primary genes.

## Limitations

While formulating the auxiliary gene selection problem as a top-$k$ choice based on HVG score ranking improves primary gene prediction, several limitations remain. First, since the approach relies solely on HVG scores, it is challenging to fully capture the underlying biological functional relationships among genes. When incorporating the biological relationships, the problem involves both determining the set of informative genes to select and identifying which biological functions are activated, and it is complicated. In the present study, we restricted our analysis to the informativeness of genes and did not model pathway activation effects as a first step. As a potential direction for future work, one could consider predicting the contribution of each gene based on both the strength of its expression signal and its functional relevance to the primary genes, and using this information to guide auxiliary gene selection.

Second, the contribution of an auxiliary gene to the estimation of primary genes may vary depending on the spatial location within the tissue; a gene may be informative in one region but act as noise in another. A spatially adaptive mechanism that detects and down-weights unhelpful image–gene pairs during training, for example by monitoring per-sample losses or employing co-teaching strategies (Han et al. 2018) from the noisy-label literature, could further improve robustness. In such a case, the addition of data selection alongside gene selection would increase the overall complexity of the problem, potentially requiring the development of an alternative framework.

## Conclusion

In this study, we proposed *Auxiliary Gene Learning* (AGL), which leverages the benefits of previously ignored genes by reformulating their expression estimation as auxiliary tasks and jointly training them with the primary tasks. In AGL, it is necessary to select an appropriate subset of auxiliary genes from a larger set that is often affected by observational noise. However, this becomes a high-dimensional combinatorial optimization problem, making it challenging to solve. To overcome this challenge, we proposed Prior-Knowledge-Based Differentiable Top-$k$ Gene Selection via Bi-level Optimization (DkGSB), which ranks all auxiliary gene candidates and performs top-$k$ selection in a differentiable manner. The experimental results demonstrate the effectiveness of incorporating previously ignored genes into the learning process as auxiliary tasks, and show that the proposed DkGSB method outperforms conventional auxiliary task learning approaches.

## Acknowledgments

## References

Adam, P.; Sam, G.; Francisco, M.; Adam, L.; James, B.; Gregory, C.; Trevor, K.; Zeming, L.; Natalia, G.; Luca, A.; Alban, D.; Andreas, K.; Edward, Y.; Zach, D.; Martin, R.; Alykhan, T.; Sasank, C.; Benoit, S.; Lu, F.; Junjie, B.; and Soumith, C. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Neural Information Processing Systems*, 8026 – 8037.

Aviv, N.; Idan, A.; Haggai, M.; Gal, C.; and Ethan, F. 2021. Auxiliary Learning by Implicit Differentiation. In *International Conference on Learning Representations*.

Chen, Z.; Badrinarayanan, V.; Lee, C.-Y.; and Rabinovich, A. 2018. Gradnorm: Gradient Normalization for Adaptive Loss Balancing in Deep Multitask Networks. In *International Conference on Machine Learning*, 794–803.

Cho, W. J.; Yoon, H.; Jeong, D.; Lim, H.; and Chong, Y. 2025. $MV_{Hybrid}$: Improving Spatial Transcriptomics Prediction with Hybrid State Space-Vision Transformer Backbone in Pathology Vision Foundation Models. In *MICCAI Workshop on Computational Pathology with Multimodal Data (COMPAYL)*.

Chung, Y.; Ha, J. H.; Im, K. C.; and Lee, J. S. 2024. Accurate Spatial Gene Expression Prediction by Integrating Multi-resolution Features. In *Computer Vision and Pattern Recognition*, 11591–11600.

Dawood, M.; Branson, K.; Rajpoot, N. M.; and Minhas, F. u. A. A. 2021. All You Need is Color: Image based Spatial Gene Expression Prediction Using Neural Stain Learning. In *Machine Learning and Knowledge Discovery in Databases*, 437–450.

Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A Large-scale Hierarchical Image Database. In *Computer Vision and Pattern Recognition*, 248–255.

Han, B.; Yao, Q.; Yu, X.; Niu, G.; Xu, M.; Hu, W.; Tsang, I.; and Sugiyama, M. 2018. Co-teaching: Robust Training of Deep Neural Networks with Extremely Noisy Labels. In *Neural Information Processing Systems*.

He, B.; Bergenstråhle, L.; Stenbeck, L.; Abid, A.; Andersson, A.; Borg, Å.; Maaskola, J.; Lundeberg, J.; and Zou, J. 2020. Integrating Spatial Gene Expression and Breast Tumour Morphology via Deep Learning. *Nature Biomedical Engineering*, 4(8): 827–834.

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep Residual Learning for Image Recognition. In *Computer Vision and Pattern Recognition*, 770–778.

Jaume, G.; Doucet, P.; Song, A.; Lu, M. Y.; Almagro Pérez, C.; Wagner, S.; Vaidya, A.; Chen, R.; Williamson, D.; Kim, A.; et al. 2024. Hest-1k: A dataset for Spatial Transcriptomics and Histology Image Analysis. *Neural Information Processing Systems*, 37: 53798–53833.

Jiang, J.; Chen, B.; Pan, J.; Wang, X.; Liu, D.; Jiang, J.; and Long, M. 2023. Forkmerge: Mitigating Negative Transfer in Auxiliary-task Learning. *Neural Information Processing Systems*, 36: 30367–30389.

Kendall, A.; Gal, Y.; and Cipolla, R. 2018. Multi-task Learning Using Uncertainty to Weigh Losses for Scene Geometry and Semantics. In *Computer Vision and Pattern Recognition*, 7482–7491.

Kingma, D. P.; and Ba, J. 2014. Adam: A Method for Stochastic Optimization. *arXiv preprint arXiv:1412.6980*.

Ko, D.; Choi, J.; Choi, H. K.; On, K.-W.; Roh, B.; and Kim, H. J. 2023. Meltr: Meta Loss Transformer for Learning to Fine-tune Video Foundation Models. In *Computer Vision and Pattern Recognition*, 20105–20115.

Lopez, R.; Regier, J.; Cole, M. B.; Jordan, M. I.; and Yosef, N. 2018. Deep Generative Modeling for Single-cell Transcriptomics. *Nature methods*, 15(12): 1053–1058.

M. Ribeiro, D.; Ziyani, C.; and Delaneau, O. 2022. Shared Regulation and Functional Relevance of Local Gene Co-expression Revealed by Single Cell Analysis. *Communications Biology*, 5(1): 876.

Mejia, G.; Cárdenas, P.; Ruiz, D.; Castillo, A.; and Arbeláez, P. 2023. SEPAL: Spatial Gene Expression Prediction from Local Graphs. In *International Conference on Computer Vision*, 2294–2303.

Mejia, G.; Ruiz, D.; Cárdenas, P.; Manrique, L.; Vega, D.; and Arbeláez, P. 2024. Enhancing Gene Expression Prediction from Histology Images with Spatial Transcriptomics Completion. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 91–101.

Nishimura, K.; Hirose, H.; Bise, R.; Shiku, K.; and Kojima, Y. 2025. Learning Relative Gene Expression Trends from Pathology Images in Spatial Transcriptomics. In *Neural Information Processing Systems*.

Pang, M.; Su, K.; and Li, M. 2021. Leveraging Information in Spatial Transcriptomics to Predict Super-resolution Gene Expression from Histology Images in Tumors. *BioRxiv*, 2021–11.

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning Transferable Visual Models from Natural Language Supervision. In *International Conference on Machine Learning*.

Satija, R.; Farrell, J. A.; Gennert, D.; Schier, A. F.; and Regev, A. 2015. Spatial Reconstruction of Single-cell Gene Expression Data. *Nature Biotechnology*, 33(5): 495–502.

Shaham, U.; Stanton, K. P.; Zhao, J.; Li, H.; Raddassi, K.; Montgomery, R.; and Kluger, Y. 2017. Removal of Batch Effects Using Distribution-matching Residual Networks. *Bioinformatics*, 33(16): 2539–2546.

Sivasubramanian, D.; Maheshwari, A.; Shenoy, P.; Ramakrishnan, G.; et al. 2023. Adaptive Mixing of Auxiliary Losses in Supervised Learning. In *Conference on Artificial Intelligence*, 9855–9863.

Ståhl, P. L.; Salmén, F.; Vickovic, S.; Lundmark, A.; Navarro, J. F.; Magnusson, J.; Giacomello, S.; Asp, M.; Westholm, J. O.; Huss, M.; et al. 2016. Visualization and Analysis of Gene Expression in Tissue Sections by Spatial Transcriptomics. *Science*, 353(6294): 78–82.

Stuart, T.; Butler, A.; Hoffman, P.; Hafemeister, C.; Papalexi, E.; Mauck, W. M.; Hao, Y.; Stoeckius, M.; Smibert, P.; and Satija, R. 2019. Comprehensive Integration of Single-cell Data. *cell*, 177(7): 1888–1902.

Wang, L.; Maletic-Savatic, M.; and Liu, Z. 2022. Region-specific Denoising Identifies Spatial Co-expression Patterns and Intra-tissue Heterogeneity in Spatially Resolved Transcriptomics Data. *Nature Communications*, 13(1): 6912.

Williams, C. G.; Lee, H. J.; Asatsuma, T.; Vento-Tormo, R.; and Haque, A. 2022. An Introduction to Spatial Transcriptomics for Biomedical Research. *Genome medicine*, 14(1): 68.

Xie, R.; Pang, K.; Chung, S.; Perciani, C.; MacParland, S.; Wang, B.; and Bader, G. 2023. Spatially Resolved Gene Expression Prediction from Histology Images via Bi-modal Contrastive Learning. *Neural Information Processing Systems*, 36: 70626–70637.

Yang, Y.; Hossain, M. Z.; Stone, E.; and Rahman, S. 2024. Spatial Transcriptomics Analysis of Gene Expression Prediction Using Exemplar Guided Graph Neural Network. *Pattern Recognition*, 145: 109966.

Yang, Y.; Hossain, M. Z.; Stone, E. A.; and Rahman, S. 2023. Exemplar Guided Deep Neural Network for Spatial Transcriptomics Analysis of Gene Expression Prediction. In *Winter Conference on Applications of Computer Vision*, 5039–5048.

Yuansong, Z.; Zhuoyi, W.; Weijiang, Y.; Rui, Y.; Bingling, L.; Zhonghui, T.; Yutong, L.; and Yuedong, Y. 2022. Spatial Transcriptomics Prediction from Histology Jointly through Transformer and Graph Neural Networks. *bioRxiv*.

Zheng, G. X.; Terry, J. M.; Belgrader, P.; Ryvkin, P.; Bent, Z. W.; Wilson, R.; Ziraldo, S. B.; Wheeler, T. D.; McDermott, G. P.; Zhu, J.; et al. 2017. Massively Parallel Digital Transcriptional Profiling of Single Cells. *Nature communications*, 8(1): 14049.