# VLM-based Prompts as the Optimal Assistant for Unpaired Histopathology Virtual Staining

Zizhi Chen*
Fudan University
China
chenzz24@m.fudan.edu.cn

Xinyu Zhang*
Central South University
China
8208220519@csu.edu.cn

Minghao Han*
Fudan University
China
mhhan22@m.fudan.edu.cn

Yizhou Liu
Harbin Institute of Technology
China
kevinclaint.liu@gmail.com

Ziyun Qian
Fudan University
China
zyqian22@m.fudan.edu.cn

Weifeng Zhang
Central South University
China
8203211723@csu.edu.cn

Xukun Zhang
Fudan University
China
zhangxk21@m.fudan.edu.cn

Jingwei Wei†
Chinese Academy of Sciences
China
weijingwei2014@ia.ac.cn

Lihua Zhang†
Fudan University
China
lihuazhang@fudan.edu.cn

## Abstract

In histopathology, tissue sections are typically stained using common H&E staining or special stains (MAS, PAS, PASM, *etc.*) to clearly visualize specific tissue structures. The rapid advancement of deep learning offers an effective solution for generating virtually stained images, significantly reducing the time and labor costs associated with traditional histochemical staining. However, a new challenge arises in separating the fundamental visual characteristics of tissue sections from the visual differences induced by staining agents. Additionally, virtual staining often overlooks essential pathological knowledge and the physical properties of staining, resulting in only style-level transfer. To address these issues, we introduce, for the first time in virtual staining tasks, a pathological vision-language large model (VLM) as an auxiliary tool. We integrate contrastive learnable prompts, foundational concept anchors for tissue sections, and staining-specific concept anchors to leverage the extensive knowledge of the pathological VLM. This approach is designed to describe, frame, and enhance the direction of virtual staining. Furthermore, we have developed a data augmentation method based on the constraints of the VLM. This method utilizes the VLM's powerful image interpretation capabilities to further integrate image style and structural information, proving beneficial in high-precision pathological diagnostics. Extensive evaluations on publicly available multi-domain unpaired staining datasets demonstrate that our method can generate highly realistic images and enhance the accuracy of downstream tasks, such as glomerular detection and segmentation. Our code is available at: https://github.com/CZZZZZZZZZZZZZZZZZ/VPGAN-HARBOR
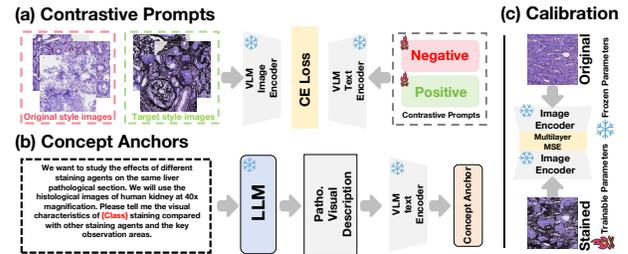
## CCS Concepts

• **Computing methodologies** → **Reconstruction**.

## Keywords

Digital Pathology, Virtual Staining, Vision Language Model

---

*Equal contribution.
†Corresponding author



Figure 1: Three auxiliary methods for virtual staining tasks proposed by us using the VLM: (a) Learnable contrastive prompts based on the classification task. (b) Concept anchors design based on LLM. (c) Visual calibration based on the VLM.

## 1 Introduction

Histopathological examination is widely regarded as the gold standard for the clinical diagnosis of diseases [5, 27, 45]. This process typically involves histochemical staining, which differentiates various tissue components through distinct colors to aid in pathological diagnosis. Routine pathological examination commonly employs Hematoxylin and Eosin (H&E) staining to highlight tissue morphology for initial diagnostic purposes. However, H&E staining often fails to provide sufficient diagnostic information for many diseases. Consequently, the use of special stains [32, 61] offers critical diagnostic insights across multiple dimensions [39]. For instance, in renal pathology, Masson's Trichrome (MAS) staining is utilized to distinguish collagen fibers from muscle fibers, Periodic Acid-Schiff (PAS) staining is employed to better visualize the Glomerular basement membrane, Tubular basement membrane, and Mesangial matrix (GTM) [3, 47], and Periodic Acid-Schiff with Methenamine Silver (PASM) staining can more clearly delineate the GTM on the basis of PAS [8, 42]. Nevertheless, the application of special stains generally requires more time and incurs higher labor costs. Moreover, when patients suffer from non-neoplastic

kidney diseases [62], liver cirrhosis [31], or other conditions, pathologists may necessitate multiple types of special stains to achieve a more accurate diagnosis. Therefore, the development of virtual staining technology, which reduces the costs of special stains for both pathologists and patients while addressing the need for multi-staining on the same tissue section, holds significant importance in clinical practice [66, 68].

Recent advancements in generative model [11, 21, 56] technology have also spurred progress in the image-to-image (I2I) domain within pathology [23, 33, 39, 41]. These advances enable stain style transfer for color normalization [7, 60] and serve as a feature enhancement strategy [36]. They also train effective feature extractors, boosting performance in subtype classification [24, 43]. However, current pathological I2I methods largely follow the technical frameworks established in the natural image domain, focusing primarily on style transformation while neglecting the texture and cytological structures of pathological sections, as well as the physical and chemical properties of staining agents. This limitation undermines the reliability of virtually stained sections as diagnostic tools. In our view, expecting current generative models to possess the extensive domain knowledge and cellular-level visual discernment required in pathology is overly demanding. Such models urgently require a pathology expert-level "assistant" to aid them in more effectively accomplishing virtual staining tasks.

The advent of the VLM in pathology [25, 26, 44, 57–59, 64] has made this endeavor feasible. Empowered by millions of pathology image-caption pairs, it possesses extensive pathological knowledge and robust capabilities in pathological image recognition. It has achieved state-of-the-art performance in various tasks, including pathology image classification, segmentation, caption generation, text-to-image synthesis, image and text retrieval. Given its role as an expert-level assistant in clinical decision-making, providing comprehensive support to pathologists, we aspire to extend its all-round excellence to the field of virtual staining. The provision of advanced pathological knowledge by the VLM is expected to potentially ensure that virtual staining results meet medical and chemical standards. Furthermore, intermediate staining processes could be characterized and fine-grained visual details might be captured through the leveraging of the VLM's powerful multimodal information extraction capability, potentially leading to enhanced performance of virtual staining.

In this paper, we present three attempts to leverage VLMs for guiding virtual staining tasks. As shown in Figure 1(a), the **contrastive prompting tuning** employs contrastive learning strategies and binary classification tasks to decode and extract the rich information embedded in pathology VLMs. This enables the system to articulate stain differentiation and staining processes that are typically challenging to describe in human language. The **conceptual anchoring method** as presented in Figure 1(b) generates foundational and stain-specific concept anchors by leveraging the rich corpora produced by Large Language Model (LLM) [1, 14] and the information compression capabilities of pathology VLM, guiding the "variation and invariance" during the staining process. For these two prompting strategies, we designed the **C**ontrastive **P**rompt **T**ransfer (CPT), **C**onstant **C**oncept **A**nchoring (CCA), and **I**ndependent **C**oncept **R**einforcement (ICR) modules, respectively. Together with a unpaired I2I model as baseline, these components

form our method, the **V**LM-based **P**rompts **G**enerative **A**dversarial **N**etwork (VPGAN), which, to the best of our knowledge, represents the first attempt to bridge GANs and pathology VLM. Additionally, inspired by Xiong *et al.* [65], we argue that inference enhancement based on DDIM [48] can effectively meet the demands of high-resolution diagnostic tasks. However, in practice, existing methods risk visual domain collapse (*e.g.* H&E2PASM). To address this, as illustrated in Figure 1(c), we adopted a **multi-level calibration strategy** based on the VLM, successfully improving the stability and performance of inference enhancement. Built on CLIP [52] with ResNet101 [19], our method dynamically adjusts texture and color details across network layers, significantly improving inference robustness. The inference enhancement method empowered by VPGAN and the multi-level calibration strategy together constitute our **H**istopathology st**A**ining expe**R**t **B**ased **O**n p**R**ompts (HARBOR). The contribution of this paper is summarized as follows:

- As far as we know, our proposed VPGAN is the first GAN based on diversified prompts from the pathological VLM, which serves as a super assistant for virtual staining.
- We designed a VLM-based multi-level visual calibration module to tackle data staining domain disintegration and enhance data augmentation stability and performance.
- Our proposed method produced satisfactory images in three virtual staining tasks and showed optimal performance under different inference strategies, indicating it can meet diverse scenario, cost and detection task requirements.
- Segmentation and detection accuracy across diverse glomerulus datasets are improved by our method, and strong clinical potential is demonstrated.

## 2 Related Work

### 2.1 Virtual Staining in Pathology Analysis

Virtual staining originates from the image-to-image (I2I) task in the domain of natural images, aiming to accomplish the transfer of image styles. With the recent surge in generative model technologies, I2I tasks have also seen significant advancements by leveraging baseline such as GAN [11] and DDPM [21]. Due to the scarcity and high production costs of paired data, unpaired I2I models have demonstrated greater potential. Zhu *et al.* [73]designed parallel sets of generators and discriminators to perform staining and restoration on unpaired data separately, achieving commendable results. CUT [49] introduced contrastive learning methods, successfully simplifying the style transfer process by eliminating the need for paired generator-discriminator combinations. Building upon CUT, Jung *et al.* [28]utilized graph neural networks to identify closely related patches, optimizing the contrastive learning algorithm. Zhao *et al.* [70]incorporated stochastic differential equations into I2I tasks, achieving a breakthrough in the application of diffusion architectures for unpaired image style transfer. Kim *et al.* [29]attempted to leverage Schrödinger bridges to compute the optimal path for domain transfer.

In the field of pathology, virtual staining has demonstrated significant clinical value. A series of H&E-to-IHC staining workflows facilitate the detection of tumor markers, aiding in tumor classification and diagnosis [4, 33, 35, 41, 50, 67]. Virtual staining is employed for tasks such as normalization [7, 20, 60] and tailing
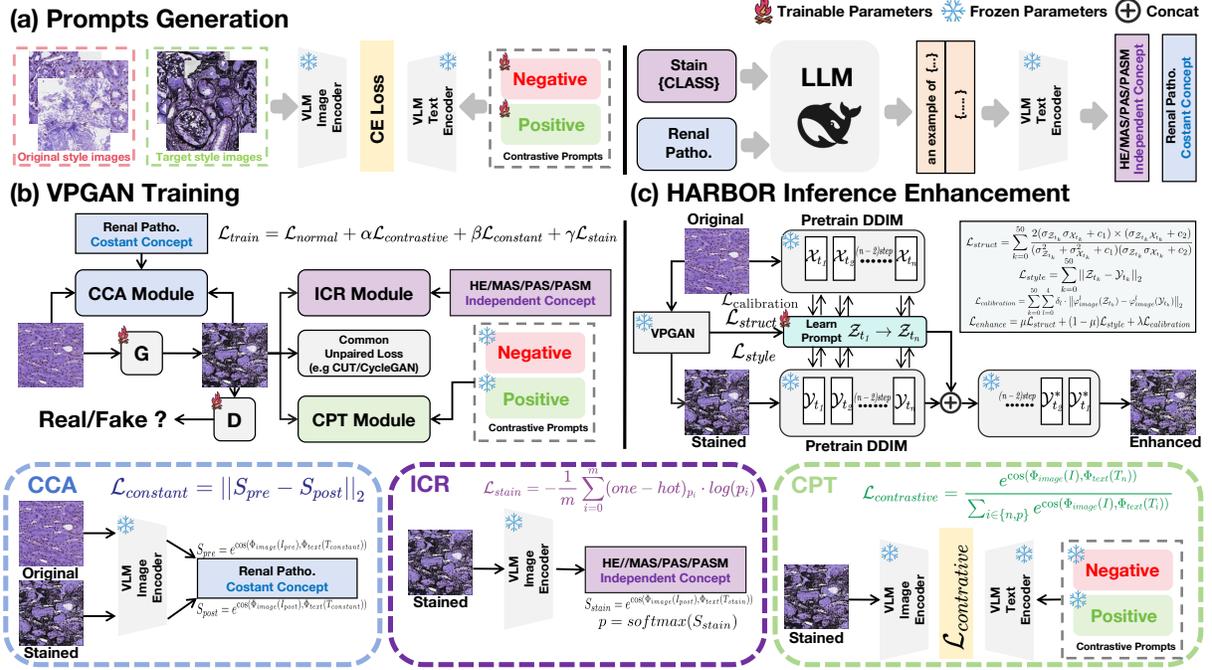
**Figure 2: Overview of the proposed VPGAN and HARBOR. In the prompt generation phase, we employed prompt tuning method to generate contrastive prompts based on a binary classification task. Utilizing the DeepSeek-R1, we created constant concept anchors and independent concept anchors of different staining agents. During the training phase, we leveraged three types of prompts and a pathological VLM to achieve the description, framing, and reinforcement of the virtual staining direction, thereby optimizing the original virtual staining model. In the inference enhancement phase, we trained learnable denoising prompt blocks based on structural and stylistic constraints, further improving the performance of virtual staining.**

artifact reduction [22, 23, 38], enhancing the clarity of tissue sections by mitigating tailing effects and reducing batch effects in pathological slides. Additionally, virtual staining exhibits potential as a proxy task to improve the visual perception capabilities of models [24, 43]. Focusing on virtual staining tasks for kidney tissue sections, Lin *et al.* proposed the UMDST [39] model, which leverages multi-task learning based on stain classification and virtual staining. This model can achieve style transfer across different stains and even simulate virtual scenarios involving multiple stains. Guan *et al.* [13]designed a style-guided module, showcasing exceptional performance and validating the medical significance of virtual staining in glomerular segmentation and detection tasks. Xiong *et al.* [65]concentrated on inference enhancement strategies, utilizing the DDIM [48] architecture to achieve a leap in structural consistency and clarity in virtual staining, thereby addressing the challenges of high-precision medical diagnosis. However, none of these methods incorporate the pathological VLM as a robust assistant to provide enhanced visual recognition capabilities and pathological knowledge, thereby improving the effectiveness of staining.

## 2.2 Prompt Tuning in Vision-Language Models

As of now, the CLIP [52] architecture remains the mainstream in the VLM and continues to play an irreplaceable role across various visual domains. prompt tuning based on CLIP has also been

demonstrated to be an exceedingly simple yet effective strategy. CoOp [72] and CoCoOp [71] have proven that learnable prompts possess even greater fitting capabilities. Additionally, prompt tuning has achieved exceptional results in tasks such as lighting adjustment [37], rain and haze removal [63], image restoration [46] and artistic style transfer [69].

The advancement of prompt tuning methods in the field of pathology is closely intertwined with the maturation of the VLM specialized in pathology. Works such as CONCH [44], MUSK [64], and PLIP [25] have constructed diverse million-scale pathology text-image pairs, demonstrating exceptional performance. Qu *et al.* [51] pioneered the application of prompt tuning on pathological images, enhancing the few-shot classification performance of pathological slides and catalyzing the emergence of outstanding works in subtype classification tasks using prompt Tuning [10, 12, 15, 16, 34, 54, 55]. Liu *et al.* [40]achieved the first application of prompt tuning in survival analysis by transforming continuous survival labels into textual prompts for ordinal survival learning. To the best of our knowledge, our proposed VPGAN and HARBOR represent the first application of prompt tuning in virtual staining and, more broadly, in medical I2I tasks. Our remarkable performance has been highly encouraging, suggesting that prompt tuning based on the pathology-specific VLM has the potential to become a "super assistant" across all tasks in the field of pathology.

## 3 Method

### 3.1 Prompts Generation

Our method generates different prompts to leverage the pathological knowledge of the VLM for more precise guidance on the staining domain. Moreover, it uses learnable prompts to capture the key information during the virtual staining intermediate process.

**Learnable contrastive prompts.** Contrastive learning methods [6, 17] often excel in learning the fine-grained features of data such as images. Meanwhile, prompt tuning methods [37, 51, 71, 72] enhance the accuracy of image text descriptions through a learnable process. Inspired by both, It is proven by us that it is feasible to use the contrastive learning method to capture important information from the intermediate steps of virtual staining and convert it into learnable text prompts during the virtual staining process.

As shown in Figure 1(a) and Figure 2(a) for the training process of contrastive text prompts, we use the training set data from both the source domain images $I_s \in \mathbb{R}^{H \times W \times 3}$ and the target domain images $I_t \in \mathbb{R}^{H \times W \times 3}$ of the subsequent staining task as the overall training set. We randomly initialize a positive prompt $T_p \in \mathbb{R}^{N \times 512}$ and a negative prompt $T_n \in \mathbb{R}^{N \times 512}$. N represents the number of embedded tokens in each prompt. Then, we feed the source and target images to the image encoder $\phi_{image}$ of the VLM to obtain their latent code. Meanwhile, we also extract the latent code of the positive and negative prompts by feeding them to the text encoder $\phi_{text}$. Based on the text-image similarity in the VLM latent space, we use the binary cross entropy loss of classifying the source and target images to learn the contrastive prompt pair:

$$\mathcal{L}_{prompt} = -(a * log(\hat{a}) + ((1 - a) * log(1 - \hat{a}))), \quad (1)$$

$$\hat{a} = \frac{e^{\cos(\Phi_{image}(I), \Phi_{text}(T_p))}}{\sum_{i \in \{n,p\}} e^{\cos(\Phi_{image}(I), \Phi_{text}(T_i))}}, \quad (2)$$

where $I \in \{I_s, I_t\}$ and $a$ is the label of the current image, 0 is for negative sample $I_s$ and 1 is for positive sample $I_t$.

**Concept Anchors.** In this section, we aim to generate pathological visual descriptions to serve as prior linguistic knowledge for guiding the virtual staining task. To minimize manual effort, large language models (LLMs) [1, 14] are employed to produce descriptions related to different staining agents. Specifically, we input the following query into the LLM: *"We want to study the effects of different staining agents on the same liver pathological section. We will use the histological images of human kidney at 40x magnification. Please tell me the visual characteristics of Class staining compared with other staining agents and the key observation areas."* Following a similar approach, we obtain the intrinsic feature description sets for kidney tissue sections, thereby deriving a total of five concept knowledge sets corresponding to the four staining classes and the intrinsic features. Ultimately, we utilized the text encoder $\phi_{text}$ of VLM to generate the final concept anchors.

### 3.2 VLM-based Prompts GAN

Building upon the aforementioned learnable contrastive prompts, the invariant concept anchors of kidney tissue sections, and the independent concept anchors of different staining agents, we propose the **V**LM-based **P**rompts **G**enerative **A**dversarial **N**etwork

(VPGAN). This framework enhances the original GAN [11] architecture for unpaired data by meticulously characterizing the intermediate staining processes, the fixed concepts of kidney tissue, and the staining agent-specific concepts, thereby improving the overall image generation quality. Furthermore, our VPGAN is adaptable to any GAN architecture as an optimization method. After extensive experimentation, we selected CycleGAN [73] as the baseline for this paper, thereby achieving the optimal generation results. Apart from our design, the configurations for the generator, discriminator, and learning rate are all aligned with the original settings of CycleGAN.

**Contrastive Prompt Transfer Module.** Inspired by CLIP-LIT [37], we prove that learnable contrastive prompts can achieve effective image enhancement. Building on this, we introduce the Contrastive Prompt Transfer Module (CPT), which decodes pathological information from the VLM to describe the intermediate staining processes and subtle domain-specific differences in staining. To the best of our knowledge, this represents the first application of learnable contrastive prompts in GAN models.

Given the learnable contrastive prompts obtained from the prompts generation step, we can train the CPT module with VLM-aware loss. This loss is based on the contrastive differences between staining domains and depicts the staining transfer process, thereby improving the quality of virtual staining.

$$\mathcal{L}_{contrastive} = \frac{e^{\cos(\Phi_{image}(I), \Phi_{text}(T_n))}}{\sum_{i \in \{n,p\}} e^{\cos(\Phi_{image}(I), \Phi_{text}(T_i))}}. \quad (3)$$
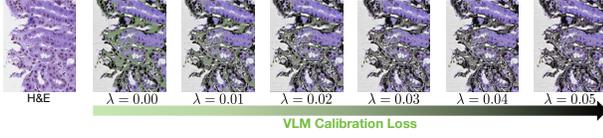
**Constant Concept Anchoring Module.** In fact, the most significant difference between virtual staining and natural image style transfer lies in the need to consider the preservation of texture and shape, the authenticity of pathological features, and the real physicochemical properties of staining agents. Merely achieving the highest degree of fitting in terms of style will lead to a substantial reduction in clinical efficacy. This is also the crucial reason why a large number of evaluation metrics for virtual staining tasks place greater emphasis on structural preservation [23, 39, 65]. In addition, the inherent properties in imaging, such as the magnification of the slices and the imaging equipment used for the dataset, are also taken into account. As described in the prompt generation, we leverage the powerful capabilities of DeepSeek-R1 [14] and online searches to obtain relatively accurate descriptions, and manually remove parts with factual errors caused by LLM hallucinations.

Next, the goal of Constant Concept Anchoring Module (CCA) is to quantify the invariance between the pre-staining images $I_{pre} \in \mathbb{R}^{H \times W \times 3}$ and the post-staining images $I_{post} \in \mathbb{R}^{H \times W \times 3}$. To this end, we generate the concept of renal slice invariance, denoted as $\mathbb{R}^{1 \times 512}$, which is derived by first generating conceptual descriptions via the LLM and subsequently transforming them into textual embeddings through the VLM. We use CPT's method to measure image-text correspondence with cosine similarity, obtaining cosine similarities $S_{pre}$ and $S_{post}$ for subsequent concept invariance analysis. The formula is as follows:

$$S_{pre} = e^{\cos(\Phi_{image}(I_{pre}), \Phi_{text}(T_{constant}))}, \quad (4)$$

$$S_{post} = e^{\cos(\Phi_{image}(I_{post}), \Phi_{text}(T_{constant}))}. \quad (5)$$

Subsequently, to ensure the invariance of the constant concept before and after staining, we employed the Mean Squared Error (MSE) loss function to calculate the mean of the sum of squared

**Figure 3: We demonstrate a fine-grained verification process based on the VLM on the H&E2PASM task, enabling the progressive and successful verification of the staining domains.**

differences between the cosine similarities $S_{pre}$ and $S_{post}$ before and after staining, thereby completing the delineation of the virtual staining range. The formula for the constant concept loss function $\mathcal{L}_{constant}$ is as follows:

$$\mathcal{L}_{constant} = \|S_{pre} - S_{post}\|_2 . \tag{6}$$

**Independent Concept Reinforcement Module.** The staining effects are aimed to be further enhanced by us by leveraging the independent concept anchors of staining agents generated by LLM. Based on textual prompts from VLMs, more microscopic details that are difficult to capture in conventional GANs can be obtained, such as the staining effects of specific agents on cell nuclei in pathological slides. It is important to emphasize that certain "shortcut" prompts that quickly deceive the discriminator (*e.g.* black streaks produced by PASM staining) will be eliminated, as their use would undermine our goal of achieving fine-grained textual prompts and instead exacerbate overfitting in style transfer. The final selected textual prompts for the four staining agents are collectively represented as $\mathbb{R}^{4 \times 512}$ after being encoded into text embeddings by the VLM. Subsequently, we compute the cosine similarity $S_{stain}$ between the stained image $I_{post}$ and the textual prompts $T_{stain}$ containing information about the four types of staining agents (H&E, MAS, PAS, PASM). Based on the inherent preprocessing method and computational rules of the CE loss, we perform a four-class proxy task to ensure that the stained image $I_{post}$ approximates the target staining domain as closely as possible. The formula is as follows:

$$S_{stain} = e^{\cos(\Phi_{image}(I_{post}), \Phi_{text}(T_{stain}))}, \tag{7}$$

$$p = softmax(S_{stain}), \tag{8}$$

$$\mathcal{L}_{stain} = -\frac{1}{m} \sum_{i=0}^{m} (one - hot)_{p_i} \cdot log(p_i), \tag{9}$$

where $p$ represents the result of normalizing the similarity $S_{stain}$. We use $one - hot$ encoding to describe the category of the virtually stained image, which is consistent with the category of the target staining domain. In summary, the loss function of VPGAN is as follows, where $\alpha$, $\beta$, and $\gamma$ are the hyperparameters therein:

$$\mathcal{L}_{train} = \mathcal{L}_{normal} + \alpha \mathcal{L}_{contrastive} + \beta \mathcal{L}_{constant} + \gamma \mathcal{L}_{stain}. \tag{10}$$

## 3.3 Inference Enhancement

Dual-Path Inference (DPI) [65], based on DDIM [48], achieves inference enhancement. It has been demonstrated that this method improves the integrity of pathological tissue structures and significantly reduces the distortion of virtually stained images. However, it carries a substantial risk of staining domain collapse. To address this issue, we introduce fine-grained structural verification based

on the VLM, which successfully resolves the problem of staining domain collapse and further enhances performance. The progressive correction effects are illustrated in Figure 3.

**Inference Enhancement Baseline.** We adopted the same settings as DPI and pre-trained a DDIM with a step size of 50. As a result, we can obtain the noisy image $\mathcal{X}_{t_k}$ at the $k$-th step based on DDIM for the original image $I_{pre}$. Similarly, we can get the intermediate noisy image $\mathcal{Y}_{t_k}$ for the image $I_{post}$ after virtual staining by VPGAN. The recurrence formulas for $\mathcal{X} = \{X_{t_1}, X_{t_2}, \cdots, X_{t_k}, \cdots, X_{t_{50}}\}$ and $\mathcal{Y} = \{\mathcal{Y}_{t_1}, \mathcal{Y}_{t_2}, \cdots, \mathcal{Y}_{t_k}, \cdots, \mathcal{Y}_{t_{50}}\}$ in DDIM are as follows:

$$X_{t_{k+1}} = \sqrt{\alpha_{k+1}} \left( \frac{X_{t_k} - \sqrt{1 - \alpha_k}\epsilon_\theta(X_{t_k}, k, C_S)}{\sqrt{\alpha_k}} \right) + \sqrt{1 - \alpha_{k+1}}\epsilon_\theta(X_{t_k}, k, C_S), \tag{11}$$

$$\mathcal{Y}_{t_{k+1}} = \sqrt{\alpha_{k+1}} \left( \frac{\mathcal{Y}_{t_k} - \sqrt{1 - \alpha_k}\epsilon_\theta(\mathcal{Y}_{t_k}, k, C_E)}{\sqrt{\alpha_k}} \right) + \sqrt{1 - \alpha_{k+1}}\epsilon_\theta(\mathcal{Y}_{t_k}, k, C_E). \tag{12}$$

The diffusion model's conditional variables are denoted as $C_S$ and $C_T$, where $C_S$ represents the source domain category conditional variable, and $C_E$ represents signifies the absence of additional conditional variables to mitigate errors in the style trajectory. $\epsilon_\theta$ is a neural network controlled by the parameter $\theta$, and the main function of this network is to predict noise.

Our objective is to augment each stained image by in proper order training $\mathcal{Z} = \{\mathcal{Z}_{t_1}, \mathcal{Z}_{t_2}, \cdots, \mathcal{Z}_{t_k}, \cdots, \mathcal{Z}_{t_{50}}\}$, an initially zero-initialized empty prompt map, which serves as an additional noise prompt to achieve the effect of data augmentation. $\mathcal{Z}$ is trained based on the SSIM structural constraint of $\mathcal{X}$ and the MSE stylization constraint of $\mathcal{Y}$. The formula is as follows:

$$\mathcal{L}_{struct} = \sum_{k=0}^{50} \frac{2(\sigma_{\mathcal{Z}_{t_k}} \sigma_{X_{t_k}} + c_1) \times (\sigma_{\mathcal{Z}_{t_k}X_{t_k}} + c_2)}{(\sigma_{\mathcal{Z}_{t_k}}^2 + \sigma_{X_{t_k}}^2 + c_1)(\sigma_{\mathcal{Z}_{t_k}} \sigma_{X_{t_k}} + c_2)}, \tag{13}$$

$$\mathcal{L}_{style} = \sum_{k=0}^{50} \|\mathcal{Z}_{t_k} - \mathcal{Y}_{t_k}\|_2 . \tag{14}$$

Among them, $\sigma_{X_{t_k}}$ and $\sigma_{\mathcal{Z}_{t_k}}$ are the variances of $X_{t_k}$ and $\mathcal{Z}_{t_k}$, $\sigma_{\mathcal{Z}_{t_k}X_{t_k}}$ is their covariance.

We denoise and restore $\mathcal{Y}$ based on the cues from $\mathcal{Z}$ to obtain the enhanced image $I_{enhance} \in \mathbb{R}^{H \times W \times 3}$. The formula is as follows:

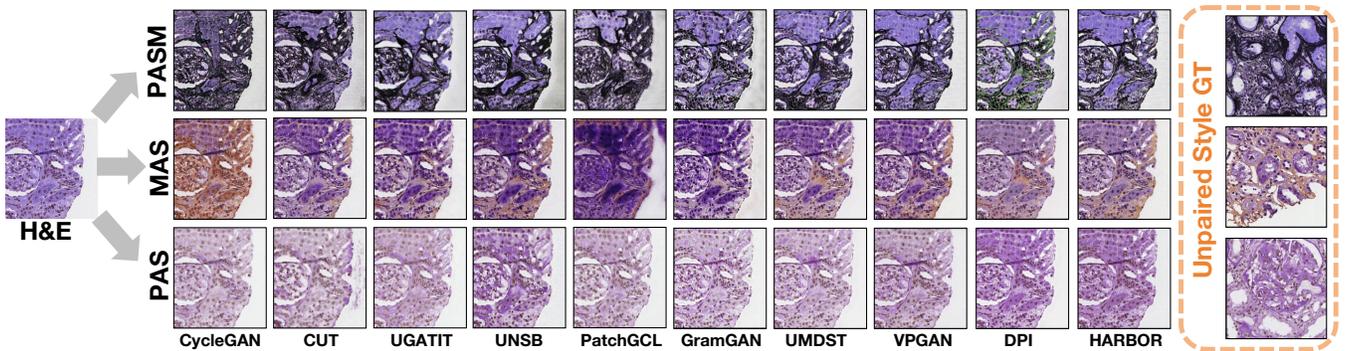$$\mathcal{Y}_{t_k}^* = \mathcal{Y}_{t_k} + \mathcal{Z}_{t_k}, \tag{15}$$

$$\psi(\mathcal{Y}_{t_{k-1}}^*, C_T) = \sqrt{\alpha_{k-1}} \left( \frac{\mathcal{Y}_{t_k}^* - \sqrt{1 - \alpha_k}\epsilon_\theta(\mathcal{Y}_{t_k}^*, k, C_T)}{\sqrt{\alpha_k}} \right) + \sqrt{1 - \alpha_{k-1}}\epsilon_\theta(\mathcal{Y}_{t_k}^*, k, C_T), \tag{16}$$

where the $\psi$ function represents conditional sampling, and $C_T$ represents the target domain label. Finally, $I_{enhance}$ is obtained.

**Structural Verification based on the VLM.** Without altering the DDIM-based inference enhancement framework, we aim to leverage the powerful capabilities of VLMs to address the issue of staining domain degradation. Since the focus here is primarily on correcting structural and color-related aspects, we employ CLIP [52], a general-purpose VLM, rather than pathology-specialized VLMs such as CONCH. For the visual encoder $\varphi_{image}$, we utilize ResNet101 [19], which allows us to extract intermediate layer features for detailed corrections, thereby resolving staining domain degradation and

**Table 1: Comparison of different methods on H&E2MAS,H&E2PAS and H&E2PASM datasets. The parts with a red background represent zero-cost inference methods, and the best results are marked in red. The parts with a blue background represent inference enhancement methods, and the best results are marked in blue.**

| Method | H&E2MAS | | | | | H&E2PAS | | | | | H&E2PASM | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | SSIM↑ | CSS↑ | MS-SSIM↑ | PSNR↑ | FID | SSIM↑ | CSS↑ | MS-SSIM↑ | PSNR↑ | FID | SSIM↑ | CSS↑ | MS-SSIM↑ | PSNR↑ | FID |
| CycleGAN [73] | 0.7824 | 0.8336 | 0.7843 | 15.71 | 133.17 | 0.8654 | 0.8888 | 0.8914 | 17.12 | 149.81 | 0.5490 | 0.6391 | 0.4013 | 10.54 | 140.44 |
| CUT [49] | 0.7075 | 0.7519 | 0.8090 | 16.03 | 125.64 | 0.7157 | 0.7337 | 0.8496 | 16.65 | 103.02 | 0.2914 | 0.3180 | 0.3293 | 10.14 | 103.96 |
| UGATIT [30] | 0.6091 | 0.6322 | 0.8412 | 16.24 | 152.48 | 0.5658 | 0.5793 | 0.8542 | 16.70 | 111.04 | 0.3680 | 0.4013 | 0.3024 | 10.96 | 136.72 |
| UNSB [29] | 0.6224 | 0.6773 | 0.7823 | 14.38 | 120.13 | 0.6648 | 0.6752 | 0.8701 | 18.3 | 112.46 | 0.2849 | 0.3119 | 0.2938 | 10.39 | 87.54 |
| PatchGCL [28] | 0.3677 | 0.4199 | 0.6777 | 11.99 | 123.73 | 0.4940 | 0.5039 | 0.8011 | 16.46 | 91.99 | 0.2346 | 0.2573 | 0.2623 | 10.04 | 95.37 |
| GramGAN [13] | 0.6260 | 0.6767 | 0.8073 | 14.77 | 175.83 | 0.6938 | 0.7096 | 0.8739 | 17.12 | 154.06 | 0.5088 | 0.5653 | 0.5327 | 12.17 | 174.81 |
| UMDST [39] | 0.7514 | 0.7864 | 0.8571 | 17.16 | 187.59 | 0.7762 | 0.7958 | 0.9254 | 17.12 | 154.06 | 0.5845 | 0.6276 | 0.5432 | 12.23 | 129.59 |
| VPGAN(Ours) | 0.8158 | 0.8648 | 0.8526 | 16.49 | 112.06 | 0.9173 | 0.9339 | 0.9457 | 19.06 | 132.95 | 0.6650 | 0.7372 | 0.5997 | 12.65 | 125.28 |
| DPI [65] | 0.8971 | 0.9040 | 0.9278 | 20.86 | 193.94 | 0.8935 | 0.8883 | 0.9508 | 22.25 | 157.47 | —— | —— | —— | —— | —— |
| HARBOR(Ours) | 0.9063 | 0.9149 | 0.9312 | 21.02 | 152.94 | 0.9302 | 0.9343 | 0.9643 | 23.64 | 154.09 | 0.6736 | 0.7323 | 0.6498 | 13.30 | 132.77 |



**Figure 4: The performance comparison of various existing methods and our proposed method for multiple stain transfer of the same H&E-stained image.**

further enhancing the robustness and performance of the inference enhancement. The loss function is as follows:

$$\mathcal{L}_{calibration} = \sum_{k=0}^{50} \sum_{l=0}^{4} \delta_l \cdot \left\| \varphi_{image}^l(\mathcal{Z}_{t_k}) - \varphi_{image}^l(\mathcal{Y}_{t_k}) \right\|_2, \quad (17)$$

where $\delta_l$ is the weight of the $l$-th layer of the image encoder in the ResNet101 CLIP model. Finally, our inference-enhanced loss function is as follows:
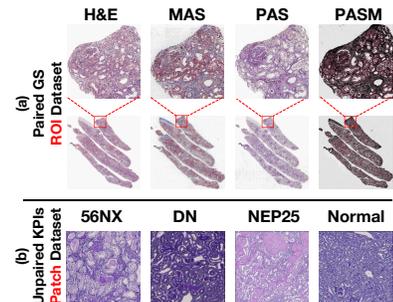
$$\mathcal{L}_{enhance} = \mu \mathcal{L}_{struct} + (1 - \mu) \mathcal{L}_{style} + \lambda \mathcal{L}_{calibration}, \quad (18)$$

$\mu$ and $\lambda$ all adjustable parameters. Details of hyperparameter settings in our method can be found in the **supplementary materials**.

## 4 Experiments and Results

### 4.1 Datasets and Experiment Setup

**Datasets.** As depicted in Figure 5, we evaluated the performance of our method on three open-source datasets. Following the settings of UMDST [39] and DPI [65], we obtained and partitioned ANHIR dataset [2] slices stained with H&E, MAS, PAS, and PASM, enabling model training and performance validation. On the GS [13] dataset, the virtual staining effect on glomerular detection and segmentation at the ROI level was examined by us adhering to the GramGAN setup. For the KPIs dataset [9], we normalized PAS slices across different disease categories according to its data settings, thereby



**Figure 5: Overview of Downstream Task Datasets**

enhancing the performance of glomerular object detection at the patch level. Please refer to the **supplementary materials** for all specific details of dataset division and preprocessing.

**Evaluation Metrics.** In this experiment, we employed five metrics to comprehensively evaluate the performance of pathological image translation. First, the Structural Similarity Index **(SSIM)** was used to measure the similarity in luminance, contrast, and structure between images. Second, the Contrast Structural Similarity **(CSS)** focused on assessing the preservation of contrast and structural details. Additionally, the Multi-Scale Structural Similarity **(MS-SSIM)**

evaluated both fine details and global structures through multiscale analysis. Meanwhile, the Peak Signal-to-Noise Ratio **(PSNR)** quantified the noise and distortion levels in the images. Finally, the Fréchet Inception Distance **(FID)** assessed the distribution consistency between generated and real images in the feature space. These metrics ensured a thorough and reliable evaluation of the results from multiple perspectives. Our approach achieved state-of-the-art (SOTA) performance across multiple metrics, validating the effectiveness of the proposed method.

**Implementation Details.** Our method is implemented in PyTorch and trained on a workstation with 8 NVIDIA H100 GPUs. We adopted CycleGAN [73] as the baseline for training VPGAN, setting the batch size to 1 and training for 50 epochs. During the inference enhancement phase, we utilized DPI [65] as the baseline. The data preprocessing methods, optimizer, and learning rate settings were kept consistent with the baseline. The remaining hyperparameter configurations are detailed in the **supplementary materials**.

## 4.2 Comparison Results

VPGAN and HARBOR were compared with previous unpaired image translation and virtual staining methods by us. Given the substantial computational cost and time overhead of the DDIM-based inference enhancement technique (requiring 5-10 minutes for inference enhancement on a single $256 \times 256$ image), we deemed it necessary to separately evaluate zero-cost inference methods and inference enhancement methods. This approach further validates the versatility of our method across diverse scenarios.

**Zero-Cost Inference Methods.** We selected **CycleGAN** as the baseline for **VPGAN** and compared it with various two-domain unpaired image translation methods such as **CUT** [49], **UGATIT** [30], **UNSB** [29], and **PatchGCL** [28], as well as multi-domain unpaired image translation methods like **GramGAN** and **UMDST**. As shown in Table 1 and Figure 4, VPGAN achieved SOTA performance across 12 metrics on three datasets compared to other zero-cost inference methods, and also demonstrated the best visual quality. It is clearly evident that compared to the baseline CycleGAN, our method effectively addresses its limitations in the transitional style transfer for the H&E2PASM and H&E2MAS tasks. By leveraging VLM-based prompt constraints, VPGAN achieves remarkable results that faithfully adhere to the staining characteristics of pathological images. Both CUT and PatchGCL demonstrate structural artifacts, and the PatchGCL method entirely collapses in the H&E2MAS task. GramGAN exhibits noticeable blurring at the edges of the images and irregular stains in the H&E2PASM task, which contradicts the physical properties of the staining agents. In contrast, UNSB and UMDST perform slightly worse in terms of image clarity and structural preservation. The diversity of issues encountered with other methods underscores the versatility of VLMs as a virtual staining assistant, playing a significant role in various aspects such as pathological knowledge guidance and structural preservation.

**Inference Enhancement Methods.** Our method, HARBOR, in comparison to the baseline DPI, has further rectified the deviation in the style domain. In the H&E2PASM task, it successfully repaired the complete collapse of the DPI staining domain (manifested by the emergence of irrelevant green colors), and in the H&E2MAS task, it also addressed the issue of green residual shadows in the cell

nuclei. In addition, visual calibration based on the VLM also enables HARBOR to achieve the SOTA performance in all indicators among inference enhancement methods. Compared to VPGAN, HARBOR better demonstrates the texture and veins of the images. This is due to the original images' texture cues and VLMs' strong visual discrimination. This may be the underlying reason for its superior performance over VPGAN across multiple metrics.

## 4.3 Ablation Study on VPGAN

As shown in Table 2 and Table 3, we conducted a series of ablation studies on each module of VPGAN and the underlying VLM assistants. These studies demonstrated the necessity of each module and enabled the selection of the Best VLM for the staining task.

**Module Ablation of VPGAN.** We investigated the functions and necessity of the CPT, CCA, and ICR modules in VPGAN, as shown in Table 2. Excitingly, the results on three datasets indicate that the functions of these modules are complementary rather than simply a superposition of performance. It can be observed that the enhancement of fixed style domains is achieved by the ICR module, which is also reflected in the general improvement of the FID metric. The CPT module is capable of representing complex intermediate coloring processes, leading to performance enhancements across multiple metrics. The CCA module, when used alone for structural invariance correction, may even yield results inferior to the baseline. However, when combined with other modules that describe the coloring process and reinforce specific coloring domains, it collectively achieves superior outcomes. This conclusion is quite intriguing. Drawing an analogy to our daily lives, an assistant who only points out what cannot be done might be quite frustrating. Yet, after summarizing the specific workflow and key tasks, appropriate regulations and reminders can further enhance work efficiency.

**Performance Differences across VLMs.** Due to differences in data sources, data volume, and training methods, VLMs may exhibit performance variations in virtual staining tasks. We conducted a comparative analysis of the effectiveness of CLIP [52] on natural images and pathology-specialized models such as PLIP [25], MUSK [64], and CONCH [44] on VPGAN across three datasets: H&E2MAS, H&E2PAS, and H&E2PASM. As shown in Table 3, the experimental results demonstrate that CONCH achieves the optimal performance on VPGAN. Based on the experimental results, we observed that CLIP even leads to a performance degradation compared to the baseline, which may stem from its inherent incompatibility with medical texts and staining tasks. The performance differences between PLIP and other pathology-specialized VLMs are attributed to its slightly inferior data volume and quality, which is also reflected in VLM-based subtype classification and survival analysis tasks. CONCH and MUSK demonstrated comparable results, but CONCH exhibited more balanced outcomes, likely due to its training data sources. Consequently, we selected CONCH as the VLM for our method, as it achieved the best average performance.

## 4.4 The Optimal Interval of the Calibration

In contrast to the reasoning enhancement component in H&E2PASM, we observed that on H&E2MAS, reasoning enhancement exhibits a more pronounced effect within a certain range of the correction

**Table 2: An ablation study was conducted on the CPT, CCA, and ICR modules in the VPGAN. All the research was carried out on H&E2MAS, H&E2PAS, and H&E2PASM. The best values are highlighted.**
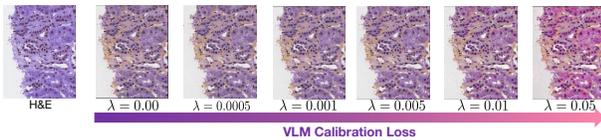
| CPT | CCA | ICR | SSIM↑ | CSS↑ | MS-SSIM↑ | PSNR↑ | FID | SSIM↑ | CSS↑ | MS-SSIM↑ | PSNR↑ | FID | SSIM↑ | CSS↑ | MS-SSIM↑ | PSNR↑ | FID |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | H&E2MAS | | | | | H&E2PAS | | | | | H&E2PASM | | |
| - | - | - | 0.7824 | 0.8336 | 0.7843 | 15.71 | 133.17 | 0.8654 | 0.8888 | 0.8914 | 17.12 | 149.81 | 0.5490 | 0.6391 | 0.4013 | 10.54 | 140.44 |
| ✓ | - | - | 0.7964 | 0.8574 | 0.8335 | 15.35 | 121.83 | 0.8729 | 0.8926 | 0.8907 | 18.52 | 143.20 | 0.5447 | 0.6422 | 0.4342 | 10.27 | 138.88 |
| - | ✓ | - | 0.7686 | 0.8375 | 0.8321 | 14.80 | 124.10 | 0.8601 | 0.8813 | 0.8931 | 17.92 | 147.93 | 0.5236 | 0.6144 | 0.3732 | 10.50 | 140.06 |
| - | - | ✓ | 0.7599 | 0.8310 | 0.8180 | 14.66 | **109.41** | 0.9031 | 0.9234 | 0.9219 | 18.27 | 134.02 | 0.5858 | 0.6541 | 0.4214 | 11.51 | 143.79 |
| ✓ | ✓ | - | 0.7763 | 0.8618 | 0.8411 | 15.80 | 135.45 | 0.8839 | 0.8812 | 0.8923 | 18.84 | 153.31 | 0.6278 | 0.6972 | **0.6069** | 11.97 | 151.22 |
| ✓ | - | ✓ | 0.7993 | 0.8520 | 0.8416 | 15.59 | 122.43 | 0.9102 | 0.9081 | 0.9414 | 18.02 | 140.16 | 0.6495 | 0.7031 | 0.5734 | 12.57 | 127.65 |
| - | ✓ | ✓ | 0.7450 | 0.8061 | 0.7983 | 15.28 | 113.73 | 0.8592 | 0.8932 | 0.8864 | 17.31 | 137.84 | 0.5901 | 0.6528 | 0.4370 | 11.52 | 132.18 |
| ✓ | ✓ | ✓ | **0.8158** | **0.8648** | **0.8526** | **16.49** | 112.06 | **0.9173** | **0.9339** | **0.9457** | **19.06** | 132.95 | **0.6650** | **0.7372** | 0.5997 | **12.65** | **125.28** |

**Table 3: Comparison of different VLMs' effect on H&E2MAS, H&E2PAS and H&E2PASM datasets. The best values are highlighted.**

| VLM | SSIM↑ | CSS↑ | MS-SSIM↑ | PSNR↑ | FID | SSIM↑ | CSS↑ | MS-SSIM↑ | PSNR↑ | FID | SSIM↑ | CSS↑ | MS-SSIM↑ | PSNR↑ | FID |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | H&E2MAS | | | | | H&E2PAS | | | | | H&E2PASM | | |
| CLIP [52] | 0.7215 | 0.7531 | 0.7309 | 16.20 | 155.34 | 0.8336 | 0.8528 | 0.8817 | 17.04 | 145.98 | 0.5827 | 0.6465 | 0.4835 | 11.89 | 143.08 |
| PLIP [25] | 0.7776 | 0.8567 | 0.8217 | 15.12 | 114.03 | 0.8745 | 0.8866 | 0.9000 | 18.23 | 141.29 | 0.6591 | 0.6922 | 0.5075 | 10.93 | 139.67 |
| MUSK [64] | **0.8259** | 0.8624 | **0.8670** | **16.71** | 130.43 | 0.8854 | 0.8871 | 0.9202 | 18.95 | 153.22 | 0.6437 | 0.7052 | **0.6214** | 12.42 | 127.42 |
| CONCH [44] | 0.8158 | **0.8648** | 0.8526 | 16.49 | **112.06** | **0.9173** | **0.9339** | **0.9457** | **19.06** | **132.95** | **0.6650** | **0.7372** | 0.5997 | **12.65** | **125.28** |

**Table 4: Search for the optimal value of the hyperparameter $\lambda$ for calibration in the H&E2MAS task.**

| $\lambda$ setting | SSIM↑ | CSS↑ | MS - SSIM↑ | PSNR↑ | FID |
|---|---|---|---|---|---|
| $\lambda = 0.00$ | 0.8794 | 0.8942 | 0.9254 | 19.86 | 197.52 |
| $\lambda = 0.0005$ | 0.8797 | 0.9023 | 0.9231 | 20.12 | 173.20 |
| $\lambda = 0.001$ | **0.9063** | **0.9149** | **0.9312** | **21.02** | 152.94 |
| $\lambda = 0.005$ | 0.8925 | 0.9010 | 0.9219 | 20.83 | 147.94 |
| $\lambda = 0.01$ | 0.8738 | 0.9048 | 0.9283 | 20.07 | **142.83** |
| $\lambda = 0.05$ | 0.8682 | 0.8700 | 0.8928 | 19.53 | 230.78 |

**Table 5: We use MAP@[0.50:0.95] to evaluate the accuracy of ROI-level glomeruli detection and segmentation.**

| Tasks | H&E (real) | PASM (generated) | PAS (generated) | MAS (generated) |
|---|---|---|---|---|
| Detection | 0.543 | **0.559** | 0.548 | 0.546 |
| Segmentation | 0.567 | **0.581** | 0.574 | 0.528 |

**Table 6: Patch-level glomerulus segmentation accuracy.**

| Method | Average | Merge | Normalization by VPGAN |
|---|---|---|---|
| Unet [53] | 87.93 | 87.12 | **88.57** |



**Figure 6: Unlike H&E2PASM, H&E2MAS correction shows post-staining domain re-collapse from over-correction, necessitating exploration of the optimal correction interval.**

loss function compared to H&E2PAS and H&E2PASM tasks. However, when the constraint properties become excessively strong, it leads to a secondary collapse in the staining domain, resulting in an unrealistic bright pink color. This intriguing phenomenon necessitates the search for the optimal hyperparameter $\lambda$ in the correction loss function on the H&E2MAS dataset, guiding the reasoning enhancement to the optimal staining range.

Table 4 demonstrates the effects of reasoning enhancement under different $\lambda$ settings in the H&E2MAS task. It can be observed that, within a certain parameter range, the calibration loss function effectively corrects the collapse in the style domain, which is also reflected in the continuous decrease of the FID score. At $\lambda = 0.001$, the optimal structural correction result is achieved, along with a

style correction result that approximates the optimal outcome. The optimal range for correction is derived from our experiments. Interestingly, when the $\lambda$ used for correcting the color domain is greater than or equal to 0.05, another type of collapse in the color domain occurs. This indicates that the correction process based on the VLM is dynamic, and granting the "assistant" too much power can lead to disastrous results.

## 4.5 Downstream Tasks

Our model's superior performance in multi-scale glomerular detection and segmentation on ROI and patch-level datasets was validated. Due to significant color and pathological differences among DS, KPIs, and ANINR datasets, we trained and tested each separately, maintaining consistent preprocessing and training. Given DDIM's size constraints, we conducted downstream experiments only on VPGAN, showcasing our method's clinical value.

**ROI-level glomerular detection and segmentation.** The GS dataset, a paired human kidney slice dataset with four virtual stain registrations and manually annotated glomerular masks, was utilized to validate the ROI-level performance of VPGAN. Following

the methodology of GramGAN, H&E was virtually stained into MAS, PAS, and PASM, and glomerular detection and segmentation were implemented using Mask R-CNN [18]. As shown in Table 5, VPGAN-based virtual staining significantly enhances detection and segmentation performance, particularly in the H&E2PASM task.

**Patch-level glomerular segmentation.** We performed patch-level glomerulus segmentation following KPIs, using four PAS-stained viral sub-datasets with distinct color variations, as shown in Figure 5(b). We hypothesized VPGAN normalization could improve segmentation, validated in Table 6. Average reflects mean test results from separate sub-dataset training, merge represents combined dataset training results, and normalization denotes pre-merge VPGAN data normalization, highlighting our method's potential.

## 5 Conclusion

In summary, we have introduced a novel unpaired slice virtual staining model designed for the virtual staining of pathological image slices. Our approach employs multiple Vision Language Model (VLM)-based prompts to achieve staining domain delineation and enhancement that aligns with actual pathological characteristics. It is crucial to highlight that we are the first in the pathological field to utilize contrastive learning methods to describe the complexity information in the staining process, thereby enhancing the staining effects. Additionally, we have provided a VLM-revised inference enhancement scheme to mitigate the risk of staining domain collapse. Our method has demonstrated the effectiveness of VLM-assisted virtual staining tasks and has been proven to serve as a data augmentation method for downstream tasks, such as glomerulus detection and segmentation. Importantly, the boundaries of VLM-based pathological prompt tuning tasks have been expanded by us, and more prompt schemes in virtual staining have been showcased.

## Acknowledgments

## References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[2] Jiří Borovec, Jan Kybic, Ignacio Arganda-Carreras, Dmitry V Sorokin, Gloria Bueno, Alexander V Khvostikov, Spyridon Bakas, Eric I-Chao Chang, Stefan Heldmann, Kimmo Kartasalo, et al. 2020. ANHIR: automatic non-rigid histological image registration challenge. *IEEE transactions on medical imaging* 39, 10 (2020), 3042–3052.

[3] Gloria Bueno, Lucia Gonzalez-Lopez, Marcial Garcia-Rojo, Arvydas Laurinavicius, and Oscar Deniz. 2020. Data for glomeruli characterization in histopathological images. *Data in brief* 29 (2020), 105314.

[4] Fuqiang Chen, Ranran Zhang, Boyun Zheng, Yiwen Sun, Jiahui He, and Wenjian Qin. 2024. Pathological semantics-preserving learning for H&E-to-IHC virtual staining. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 384–394.

[5] Richard J Chen, Tong Ding, Ming Y Lu, Drew FK Williamson, Guillaume Jaume, Bowen Chen, Andrew Zhang, Daniel Shao, Andrew H Song, Muhammad Shaban, et al. 2024. Towards a General-Purpose Foundation Model for Computational Pathology. *Nature Medicine* (2024).

[6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*. PmLR, 1597–1607.

[7] Cong Cong, Sidong Liu, Antonio Di Ieva, Maurice Pagnucco, Shlomo Berkovsky, and Yang Song. 2022. Colour adaptive generative networks for stain normalisation of histopathology images. *Medical Image Analysis* 82 (2022), 102580.

[8] Kevin de Haan, Yijie Zhang, Jonathan E Zuckerman, Tairan Liu, Anthony E Sisk, Miguel FP Diaz, Kuang-Yu Jen, Alexander Nobori, Sofia Liou, Sarah Zhang, et al. 2021. Deep learning-based transformation of H&E stained tissues into special stains. *Nature communications* 12, 1 (2021), 4884.

[9] Ruining Deng, Tianyuan Yao, Yucheng Tang, Junlin Guo, Siqi Lu, Juming Xiong, Lining Yu, Quan Huu Cap, Pengzhou Cai, Libin Lan, et al. 2025. KPIs 2024 Challenge: Advancing Glomerular Segmentation from Patch-to Slide-Level. *arXiv preprint arXiv:2502.07288* (2025).

[10] Kexue Fu, Linhao Qu, Shuo Wang, Ying Xiong, Ilias Maglogiannis, Longxiang Gao, Manning Wang, et al. 2024. Fast: A dual-tier few-shot learning paradigm for whole slide image classification. *Advances in Neural Information Processing Systems* 37 (2024), 105090–105113.

[11] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. *Advances in neural information processing systems* 27 (2014).

[12] Jiaxiang Gou, Luping Ji, Pei Liu, and Mao Ye. 2024. Queryable Prototype Multiple Instance Learning with Vision-Language Models for Incremental Whole Slide Image Classification. *arXiv preprint arXiv:2410.10573* (2024).

[13] Xianchao Guan, Yifeng Wang, Yiyang Lin, Xi Li, and Yongbing Zhang. 2024. Unsupervised multi-domain progressive stain transfer guided by style encoding dictionary. *IEEE Transactions on Image Processing* 33 (2024), 767–779.

[14] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025).

[15] Zhengrui Guo, Conghao Xiong, Jiabo Ma, Qichen Sun, Lishuang Feng, Jinzhuo Wang, and Hao Chen. 2024. FOCUS: Knowledge-enhanced Adaptive Visual Compression for Few-shot Whole Slide Image Classification. *arXiv preprint arXiv:2411.14743* (2024).

[16] Minghao Han, Linhao Qu, Dingkang Yang, Xukun Zhang, Xiaoying Wang, and Lihua Zhang. 2024. MSCPT: Few-shot Whole Slide Image Classification with Multi-scale and Context-focused Prompt Tuning. *arXiv preprint arXiv:2408.11505* (2024).

[17] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 9729–9738.

[18] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*. 2961–2969.

[19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.

[20] Martin J Hetz, Tabea-Clara Bucher, and Titus J Brinker. 2024. Multi-domain stain normalization for digital pathology: A cycle-consistent adversarial network for whole slide images. *Medical Image Analysis* 94 (2024), 103149.

[21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems* 33 (2020), 6840–6851.

[22] Ming-Yang Ho, Che-Ming Wu, Min-Sheng Wu, and Yufeng Jane Tseng. 2024. Every pixel has its moments: Ultra-high-resolution unpaired image-to-image translation via dense normalization. In *European Conference on Computer Vision*. Springer, 312–328.

[23] Ming-Yang Ho, Min-Sheng Wu, and Che-Ming Wu. 2022. Ultra-high-resolution unpaired stain transformation via kernelized instance normalization. In *European Conference on Computer Vision*. Springer, 490–505.

[24] Shengyi Hua, Fang Yan, Tianle Shen, Lei Ma, and Xiaofan Zhang. 2024. Pathoduet: Foundation models for pathological slide analysis of H&E and IHC stains. *Medical Image Analysis* 97 (2024), 103289.

[25] Zhi Huang, Federico Bianchi, Mert Yuksekgonul, Thomas J Montine, and James Zou. 2023. A visual–language foundation model for pathology image analysis using medical twitter. *Nature medicine* 29, 9 (2023), 2307–2316.

[26] Wisdom Ikezogwo, Saygin Seyfioglu, Fatemeh Ghezloo, Dylan Geva, Fatwir Sheikh Mohammed, Pavan Kumar Anand, Ranjay Krishna, and Linda Shapiro. 2023. Quilt-1m: One million image-text pairs for histopathology. *Advances in neural information processing systems* 36 (2023), 37995–38017.

[27] Maximilian Ilse, Jakub Tomczak, and Max Welling. 2018. Attention-based deep multiple instance learning. In *International conference on machine learning*. PMLR, 2127–2136.

[28] Chanyong Jung, Gihyun Kwon, and Jong Chul Ye. 2024. Patch-wise graph contrastive learning for image translation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 13013–13021.

[29] Beomsu Kim, Gihyun Kwon, Kwanyoung Kim, and Jong Chul Ye. [n. d.]. Unpaired Image-to-Image Translation via Neural Schrödinger Bridge. In *The Twelfth International Conference on Learning Representations*.

[30] Junho Kim, Minjae Kim, Hyeonwoo Kang, and Kwang Hee Lee. 2020. U-GAT-IT: Unsupervised Generative Attentional Networks with Adaptive Layer-Instance Normalization for Image-to-Image Translation. In *International Conference on Learning Representations*.

[31] Murli Krishna. 2013. Role of special stains in diagnostic liver pathology. *Clinical liver disease* 2, S1 (2013), S8–S10.

[32] Thomas Lampert, Odyssée Merveille, Jessica Schmitz, Germain Forestier, Friedrich Feuerhake, and Cédric Wemmert. 2019. Strategies for training stain invariant CNNs. In *2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019)*. IEEE, 905–909.

[33] Fangda Li, Zhiqiang Hu, Wen Chen, and Avinash Kak. 2023. Adaptive supervised patchnce loss for learning h&e-to-ihc stain translation with inconsistent groundtruth image pairs. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 632–641.

[34] Hao Li, Ying Chen, Yifei Chen, Rongshan Yu, Wenxian Yang, Liansheng Wang, Bowen Ding, and Yuchen Han. 2024. Generalizable whole slide image classification with fine-grained visual-semantic interaction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11398–11407.

[35] Jiahan Li, Jiuyang Dong, Shenjin Huang, Xi Li, Junjun Jiang, Xiaopeng Fan, and Yongbing Zhang. 2024. Virtual immunohistochemistry staining for histological images assisted by weakly-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11259–11268.

[36] Jingxiong Li, Sunyi Zheng, Chenglu Zhu, Yuxuan Sun, Pingyi Chen, Zhongyi Shui, Yunlong Zhang, Honglin Li, and Lin Yang. 2024. PathUp: Patch-wise Timestep Tracking for Multi-class Large Pathology Image Synthesising Diffusion Model. In *Proceedings of the 32nd ACM International Conference on Multimedia*. 3984–3993.

[37] Zhexin Liang, Chongyi Li, Shangchen Zhou, Ruicheng Feng, and Chen Change Loy. 2023. Iterative prompt learning for unsupervised backlit image enhancement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 8094–8103.

[38] Yiyang Lin, Yifeng Wang, Zijie Fang, Zexin Li, Xianchao Guan, Danling Jiang, and Yongbing Zhang. 2024. A Multi-Perspective Self-Supervised Generative Adversarial Network for FS to FFPE Stain Transfer. *IEEE Transactions on Medical Imaging* (2024).

[39] Yiyang Lin, Bowei Zeng, Yifeng Wang, Yang Chen, Zijie Fang, Jian Zhang, Xiangyang Ji, Haoqian Wang, and Yongbing Zhang. 2022. Unpaired multi-domain stain transfer for kidney histopathological images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 36. 1630–1637.

[40] Pei Liu, Luping Ji, Jiaxiang Gou, Bo Fu, and Mao Ye. 2025. Interpretable Vision-Language Survival Analysis with Ordinal Inductive Bias for Computational Pathology. In *The Thirteenth International Conference on Learning Representations*. https://openreview.net/forum?id=trj2Jq8riA

[41] Shengjie Liu, Chuang Zhu, Feng Xu, Xinyu Jia, Zhongyue Shi, and Mulan Jin. 2022. Bci: Breast cancer immunohistochemical image generation through pyramid pix2pix. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 1815–1824.

[42] Ying-Chih Lo, I-Fang Chung, Shin-Ning Guo, Mei-Chin Wen, and Chia-Feng Juang. 2021. Cycle-consistent GAN-based stain translation of renal pathology images with glomerulus detection application. *Applied Soft Computing* 98 (2021), 106822.

[43] Wei Lou, Guanbin Li, Xiang Wan, and Haofeng Li. 2024. Multi-modal Denoising Diffusion Pre-training for Whole-Slide Image Classification. In *Proceedings of the 32nd ACM International Conference on Multimedia*. 10804–10813.

[44] Ming Y Lu, Bowen Chen, Drew FK Williamson, Richard J Chen, Ivy Liang, Tong Ding, Guillaume Jaume, Igor Odintsov, Long Phi Le, Georg Gerber, et al. 2024. A visual-language foundation model for computational pathology. *Nature Medicine* 30, 3 (2024), 863–874.

[45] Ming Y Lu, Drew FK Williamson, Tiffany Y Chen, Richard J Chen, Matteo Barbieri, and Faisal Mahmood. 2021. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature Biomedical Engineering* 5, 6 (2021), 555–570.

[46] Ziwei Luo, Fredrik K Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B Schön. 2023. Controlling Vision-Language Models for Universal Image Restoration. *arXiv preprint arXiv:2310.01018* (2023).

[47] Jeffrey H Miner. 2012. The glomerular basement membrane. *Experimental cell research* 318, 9 (2012), 973–978.

[48] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. 2023. Null-text inversion for editing real images using guided diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 6038–6047.

[49] Taesung Park, Alexei A Efros, Richard Zhang, and Jun-Yan Zhu. 2020. Contrastive learning for unpaired image-to-image translation. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX 16*. Springer, 319–345.

[50] Qiong Peng, Weiping Lin, Yihuang Hu, Ailisi Bao, Chenyu Lian, Weiwei Wei, Meng Yue, Jingxin Liu, Lequan Yu, and Liansheng Wang. 2024. Advancing H&E-to-IHC Virtual Staining with Task-Specific Domain Knowledge for HER2 Scoring. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 3–13.

[51] Linhao Qu, Kexue Fu, Manning Wang, Zhijian Song, et al. 2023. The rise of ai language pathologists: Exploring two-level prompt learning for few-shot weakly-supervised whole slide image classification. *Advances in Neural Information Processing Systems* 36 (2023), 67551–67564.

[52] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PmLR, 8748–8763.

[53] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*. Springer, 234–241.

[54] Jiangbo Shi, Chen Li, Tieliang Gong, Chunbao Wang, and Huazhu Fu. 2024. CoD-MIL: Chain-of-Diagnosis Prompting Multiple Instance Learning for Whole Slide Image Classification. *IEEE Transactions on Medical Imaging* (2024).

[55] Jiangbo Shi, Chen Li, Tieliang Gong, Yefeng Zheng, and Huazhu Fu. 2024. ViLa-MIL: Dual-scale Vision-Language Multiple Instance Learning for Whole Slide Image Classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11248–11258.

[56] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. [n. d.]. Score-Based Generative Modeling through Stochastic Differential Equations. In *International Conference on Learning Representations*.

[57] Yuxuan Sun, Yixuan Si, Chenglu Zhu, Xuan Gong, Kai Zhang, Pingyi Chen, Ye Zhang, Zhongyi Shui, Tao Lin, and Lin Yang. 2024. CPath-Omni: A Unified Multimodal Foundation Model for Patch and Whole Slide Image Analysis in Computational Pathology. *arXiv preprint arXiv:2412.12077* (2024).

[58] Yuxuan Sun, Yunlong Zhang, Yixuan Si, Chenglu Zhu, Kai Zhang, Zhongyi Shui, Jingxiong Li, Xuan Gong, XINHENG LYU, Tao Lin, et al. [n. d.]. PathGen-1.6 M: 1.6 Million Pathology Image-text Pairs Generation through Multi-agent Collaboration. In *The Thirteenth International Conference on Learning Representations*.

[59] Yuxuan Sun, Chenglu Zhu, Sunyi Zheng, Kai Zhang, Lin Sun, Zhongyi Shui, Yunlong Zhang, Honglin Li, and Lin Yang. 2024. Pathasst: A generative foundation ai assistant towards artificial general intelligence of pathology. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 5034–5042.

[60] Cheng-Chang Tsai, Yuan-Chih Chen, and Chun-Shien Lu. 2024. Test-Time Stain Adaptation with Diffusion Models for Histopathology Image Classification. In *European Conference on Computer Vision*. Springer, 257–275.

[61] Jelica Vasiljević, Zeeshan Nisar, Friedrich Feuerhake, Cédric Wemmert, and Thomas Lampert. 2022. CycleGAN for virtual stain transfer: Is seeing really believing? *Artificial Intelligence in Medicine* 133 (2022), 102420.

[62] Patrick D Walker, Tito Cavallo, and Stephen M Bonsib. 2004. Practice guidelines for the renal biopsy. *Modern Pathology* 17, 12 (2004), 1555–1563.

[63] Yuanbo Wen, Tao Gao, and Ting Chen. 2024. Unpaired Photo-realistic Image Deraining with Energy-informed Diffusion Model. In *Proceedings of the 32nd ACM International Conference on Multimedia*. 360–369.

[64] Jinxi Xiang, Xiyue Wang, Xiaoming Zhang, Yinghua Xi, Feyisope Eweje, Yijiang Chen, Yuchen Li, Colin Bergstrom, Matthew Gopaulchan, Ted Kim, et al. 2025. A vision–language foundation model for precision oncology. *Nature* (2025), 1–10.

[65] Bing Xiong, Yue Peng, RanRan Zhang, Fuqiang Chen, JiaYe He, and Wenjian Qin. 2024. Unpaired Multi-Domain Histopathology Virtual Staining using Dual Path Prompted Inversion. *arXiv preprint arXiv:2412.11106* (2024).

[66] Renao Yan, Qiming He, Yiqing Liu, Peng Ye, Lianghui Zhu, Shanshan Shi, Jizhou Gou, Yonghong He, Tian Guan, and Guangde Zhou. 2023. Unpaired virtual histological staining using prior-guided generative adversarial networks. *Computerized Medical Imaging and Graphics* 105 (2023), 102185.

[67] Wei Zhang, Tik Ho Hui, Pui Ying Tse, Fraser Hill, Condon Lau, and Xinyue Li. 2024. High-Resolution Medical Image Translation via Patch Alignment-Based Bidirectional Contrastive Learning. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 178–188.

[68] Yijie Zhang, Kevin de Haan, Yair Rivenson, Jingxi Li, Apostolos Delis, and Aydogan Ozcan. 2020. Digital synthesis of histological stains using micro-structured and multiplexed virtual staining of label-free tissue. *Light: Science & Applications* 9, 1 (2020), 78.

[69] Zhanjie Zhang, Quanwei Zhang, Wei Xing, Guangyuan Li, Lei Zhao, Jiakai Sun, Zehua Lan, Junsheng Luan, Yiling Huang, and Huaizhong Lin. 2024. Artbank: Artistic style transfer with pre-trained diffusion model and implicit style prompt bank. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 7396–7404.

[70] Min Zhao, Fan Bao, Chongxuan Li, and Jun Zhu. 2022. Egsde: Unpaired image-to-image translation via energy-guided stochastic differential equations. *Advances in Neural Information Processing Systems* 35 (2022), 3609–3623.

[71] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Conditional prompt learning for vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 16816–16825.

[72] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Learning to prompt for vision-language models. *International Journal of Computer Vision* 130, 9 (2022), 2337–2348.

[73] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*. 2223–2232.

## A Dataset Description

**Stain Dataset.** In the ANHIR dataset [2], there are five sets of high-resolution tissue slides of the human kidney. Each set has four successive slides of the same tissue stained with different types of stain: H&E, MAS, PAS, and PASM. We use four sets (Patient 1, Patient 2, Patient 3, and Patient 4) as the training set and one set (Patient 5) as the testing set. Following the settings in UMDST [39] and DPI [65], H&E stained samples from Patient 1 were excluded from training due to large color and staining differences in the slices. Each slide is cropped into a series of $256 \times 256$ patches with an overlap of 192, where the background regions (saturation<15) are discarded, and all remaining patches are used to train and test our model. There are 40,258 patches in the training set (7,688 for H&E, 12,132 for MAS, 1,1458 for PAS, and 8,980 for PASM), and 8,070 patches in the testing set (1,989 for H&E, 2,062 for MAS, 1,900 for PAS, and 2,119 for PASM). All methods are trained on the training set and tested on the test set. This also ensures that the epochs and iterations are the same, enabling a more fair comparison.

**Downstream Tasks Dataset.** We tested the performance of our method on two glomerular detection and segmentation datasets, GS [13] and KPIs [9], under different levels of images and various tasks. On the GS dataset, referring to the dataset settings in Gram-GAN, we divided the dataset into a training set and a test set at a ratio of 4:1. We randomly extracted patches of different sizes, ranging from 200×200 to 1,000×1,000, to test the performance of virtual staining in the object detection task at the ROI level. On the KPIs dataset, we adhered to the initial settings of the dataset. We achieved color difference normalization of PAS-stained mouse slices under four states: 56NX, DN, NEP25, and Normal, to improve the performance of detection and segmentation.

## B Baseline Detail

### B.1 DeepSeek-R1

Based on comprehensive considerations of performance, cost efficiency, and other factors, we selected DeepSeek-R1 [14] as the large language model (LLM) for concept anchoring generation. We attribute the success of our approach to its several key advantages:

First, it demonstrates exceptional role-playing capabilities, effectively adopting the professional perspective of a pathologist in its text outputs. Second, its web search functionality not only addresses the prevalent hallucination issues in the medical domain but also enables a more contextually appropriate grasp of data and task-specific nuances. In fact, medical data, and even clinical scenarios, often exhibit strong regional specificity, shaped by complex factors such as geographical conditions, climate, and sociocultural environments. The web search feature allows the model to approximate the localized expertise of physicians. Finally, we must commend the model's robust reasoning capabilities, which empower DeepSeek-R1 to deliver more comprehensive descriptions of stains or pathologies, encompassing a wealth of nuanced details.

### B.2 CycleGAN

The following constraints are satisfied by two mappings, $G_{AB}$ : $A \rightarrow B$ and $G_{BA}$ : $B \rightarrow A$, parameterized by neural networks, which are used by the CycleGAN model [73] to estimate these conditionals. First, the output of each mapping should match the empirical distribution of the target domain, when marginalized over the source domain. Then, mapping an element from one domain to the other, and then back, should produce a sample close to the original element. The former technique serves as the cornerstone of all generative adversarial networks (GAN) [11]. Mappings $G_{AB}$ and $G_{BA}$ are given by neural networks trained to fool adversarial discriminators $D_B$ and $D_A$, respectively. Enforcing marginal matching on target domain $B$, marginalized over source domain $A$, involves minimizing an adversarial objective with respect to $G_{AB}$:

$$\mathcal{L}_{GAN}^B(G_{AB}, D_B) = \mathbb{E}_{b \sim p_d(b)}[\log D_B(b)] + \mathbb{E}_{a \sim p_d(a)}[\log(1 - D_B(G_{AB}(a)))],$$
(19)

while the discriminator $D_B$ is trained to maximize it. A similar adversarial loss $\mathcal{L}_{GAN}^A(G_{BA}, D_A)$ is defined for marginal matching in the reverse direction.

Cycle-consistency enforces that, when starting from a sample $a$ from $A$, the reconstruction $a^\star = G_{BA}(G_{AB}(a))$ remains close to the original $a$. For image domains, closeness between $a$ and $a^\star$ is typically measured with $L_1$ or $L_2$ norms. When using the $L_1$ norm, cycle-consistency starting from $A$ can be formulated as:

$$\mathcal{L}_{CYC}^A(G_{AB}, G_{BA}) = \mathbb{E}_{a \sim p_d(a)}\|G_{BA}(G_{AB}(a)) - a\|_1.$$
(20)

And similarly for cycle-consistency starting from B. The full Cycle-GAN objective is given by:

$$\mathcal{L}_{GAN}^A(G_{BA}, D_A) + \mathcal{L}_{GAN}^B(G_{AB}, D_B) + \\ \nu \mathcal{L}_{CYC}^A(G_{AB}, G_{BA}) + \nu \mathcal{L}_{CYC}^B(G_{AB}, G_{BA}),$$
(21)

where $\nu$ is a hyper-parameter that balances between marginal matching and cycle-consistency.

CycleGAN's success can be attributed to the complementary roles of marginal matching and cycle-consistency in its objective. Marginal matching enforces realism per domain.

### B.3 DDIM Inversion

DDIM [48] utilizes an implicit non-Markovian process for sample generation, differing from DDPM [21] which relies on a Markovian chain. This non-Markovian approach enables accelerated sampling through step-skipping in the reverse diffusion process. The core reverse-process equation of DDIM is expressed as:

$$\mathbf{x}_{t-1} = \sqrt{\alpha_{t-1}}\left(\frac{\mathbf{x}_t - \sqrt{1 - \alpha_t}\epsilon_\theta(\mathbf{x}_t, t)}{\sqrt{\alpha_t}}\right) \\ + \sqrt{1 - \alpha_{t-1}}\epsilon_\theta(\mathbf{x}_t, t),$$
(22)

$\epsilon_\theta$ is the network predicting noise at each step.

The inversion process involves running the DDIM sampling process in reverse, which can be formulated as:

$$\mathbf{x}_{t+1} = \sqrt{\alpha_{t+1}}\left(\frac{\mathbf{x}_t - \sqrt{1 - \alpha_t}\epsilon_\theta(\mathbf{x}_t, t)}{\sqrt{\alpha_t}}\right) \\ + \sqrt{1 - \alpha_{t+1}}\epsilon_\theta(\mathbf{x}_t, t).$$
(23)

We denote $\epsilon_t^*$ as the groundtruth of prediction. To enhance the inversion process, consider the following modifications:

$$\mathbf{x}_0 = \mathbf{x}_t^* + \left(\frac{1}{\alpha_t} - 1\right)\sigma, \quad \sigma > 0,$$
(24)

$$\epsilon_\theta(\mathbf{x}_t, t) = \epsilon_\theta^*(\mathbf{x}_t, t) + \sigma, \quad \sigma > 0.$$
(25)

In non-conditional DDIM inversion, larger time steps help reduce error since error decreases as $t$ increases. The process works by iteratively applying the equations to trace the sample back to its original noise vector.

## C Experimental parameter settings

Most parameter settings have already been configured in the code. Please refer to the code content for details. To ensure fair comparison and demonstrate optimal performance, we follow all parameter settings of the baseline [65, 73]. Some hyperparameters in our method may have different numerical values depending on the dataset, as shown in the table 7 below.

**Table 7: The hyperparameter values of our method on various datasets.**

| hyperparameter setting | H&E2MAS | H&E2PAS | H&E2PASM |
|:---:|:---:|:---:|:---:|
| $\alpha$ | 30 | 50 | 30 |
| $\beta$ | 0.1 | 0.1 | 0.1 |
| $\gamma$ | 0.1 | 0.1 | 0.1 |
| $\mu$ | 0.05 | 0.55 | 0.8 |
| $\lambda$ | 0.001 | 0.001 | 0.05 |

## D Limitation and Future Work

### D.1 Limitation

**Inference Time Burden.** Although our inference enhancement achieves strong performance, DDIM still incurs substantial computational overhead even with just 50 sampling steps, requiring 5–10 minutes to enhance a single 256×256 image. This imposes a prohibitive time burden for whole-slide image (WSI) staining enhancement. We experimented with several acceleration methods for diffusion models, but all led to degraded inference quality.

**ResNet CLIP not Specialized for Pathology.** Our stain normalization method leverages multi-level MSE based on ResNet101 [19] CLIP [52], an approach that is difficult to replicate with ViT-based architectures like CLIP, CONCH [44], or MUSK [64]. We attribute this advantage to ResNet's hierarchical structure, which effectively captures diverse visual features across different scales. However, a notable limitation is that all current pathology-specific vision-language model (VLM) adopt ViT architectures. This prevents direct comparison between natural image optimized ResNet CLIP and pathology specialized ResNet encoders, where the former excels in color perception while the latter better understands histopathology.

### D.2 Future Work

We aim to explore the greater potential of VLMs and Diffusion models in pathological downstream tasks, extending to more diverse staining or other tasks, while seeking tighter integration methods between the two technologies.

On the other hand, we also look forward to implementing more diverse multimodal data fusion approaches, where images, text, bulk RNA-seq data, spatial transcriptomics, and other information can be synergistically integrated to achieve superior results.