# F2FLDM: Latent Diffusion Models with Histopathology Pre-Trained Embeddings for Unpaired Frozen Section to FFPE Translation

Man M. Ho[1]    Shikha Dubey [1]    Yosep Chong[2,3]
Beatrice Knudsen[3,4]    Tolga Tasdizen[1,5]

[1] Scientific Computing and Imaging Institute, University of Utah, USA
[2] The Catholic University of Korea College of Medicine, Korea
[3] Departmant of Pathology, University of Utah, USA
[4] Huntsman Cancer Institute, University of Utah Health, USA
[5] Department of Electrical and Computer Engineering, University of Utah, USA

**Abstract.** The Frozen Section (FS) technique is a rapid and efficient method, taking only 15-30 minutes to prepare slides for pathologists' evaluation during surgery, enabling immediate decisions on further surgical interventions. However, FS process often introduces artifacts and distortions like folds and ice-crystal effects. In contrast, these artifacts and distortions are absent in the higher-quality formalin-fixed paraffin-embedded (FFPE) slides, which require 2-3 days to prepare. While Generative Adversarial Network (GAN)-based methods have been used to translate FS to FFPE images (F2F), they may leave morphological inaccuracies with remaining FS artifacts or introduce new artifacts, reducing the quality of these translations for clinical assessments. In this study, we benchmark recent generative models, focusing on GANs and Latent Diffusion Models (LDMs), to overcome these limitations. We introduce a novel approach that combines LDMs with Histopathology Pre-Trained Embeddings to enhance restoration of FS images. Our framework leverages LDMs conditioned by both text and pre-trained embeddings to learn meaningful features of FS and FFPE histopathology images. Through diffusion and denoising techniques, our approach not only preserves essential diagnostic attributes like color staining and tissue morphology but also proposes an embedding translation mechanism to better predict the targeted FFPE representation of input FS images. As a result, this work achieves a significant improvement in classification performance, with the Area Under the Curve rising from 81.99% to 94.64%, accompanied by an advantageous CaseFD. This work establishes a new benchmark for FS to FFPE image translation quality, promising enhanced reliability and accuracy in histopathology FS image analysis. Our work is available at https://minhmanho.github.io/f2f_ldm/.

**Keywords:** Frozen Section to FFPE Translation · Generative Models · Latent Diffusion Models · Histopathology Image Analysis.
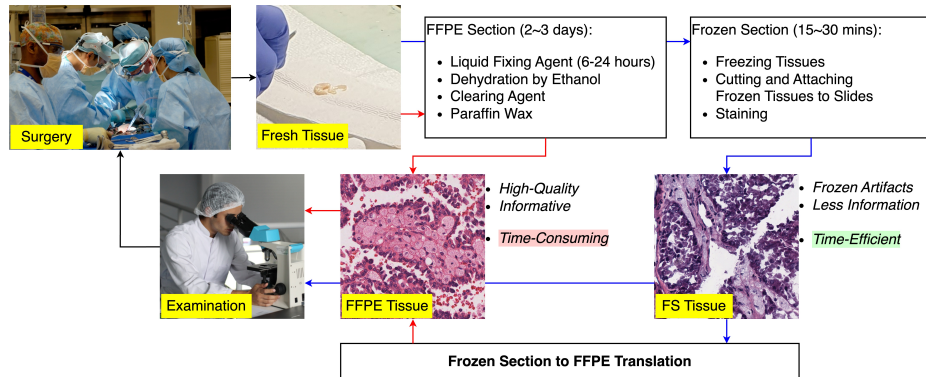
**Fig. 1.** Overview of FS and FFPE processes and our motivation.

## 1  Introduction

Histopathology, the microscopic examination of tissues to diagnose diseases, relies heavily on the Formalin-Fixed Paraffin-Embedded (FFPE) technique for producing high-quality tissue slides. Considered the gold standard, FFPE offers detailed, artifact-free slides crucial for accurate diagnosis. Yet, the FFPE process is notably slow, often requiring two to three days to prepare slides, making it unsuitable for surgeries that necessitate immediate decisions. Therefore, Frozen Section (FS) procedure is preferred due to its quick processing time. FS allows pathologists to examine tissue samples within minutes, providing surgeons with instant information that can influence surgical decisions and potentially improve patient outcomes. Despite its speed, FS image quality suffers from the introduction of artifacts such as tissue folds and ice crystals. These can obscure crucial histological details and complicate diagnoses, as described in Figure 1.

The integration of Artificial Intelligence (AI) into histological analysis has been revolutionized by Generative Adversarial Networks (GANs) [2, 8, 10, 12, 13, 17, 18, 20] and Latent Diffusion Models (LDMs) [5, 11, 15, 19], facilitating domain-to-domain image translation without paired samples. These technologies have made significant strides, yet occasionally struggle to maintain the detailed accuracy crucial for disease diagnosis in histological images. Challenges include not fully removing artifacts from Frozen Sections (FS) or introducing new artifacts. Specifically, the approach described in [12] aims to improve the quality of FS images towards FFPE standards by employing Unpaired Contrastive Translation. This method is primarily focused on correcting artifacts using an attention mechanism and self-regulation to preserve clinically relevant features. Nonetheless, it does not entirely overcome the challenge of maintaining the complex tissue structure characteristic of FFPE images, which is essential for accurate diagnoses. Moreover, diffusion models employing a Gaussian latent space formulation [19] present a novel approach to image translation. By unifying the latent space and exploiting its cyclic nature, this method facilitates image translation
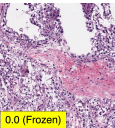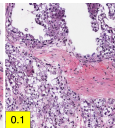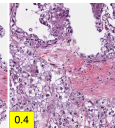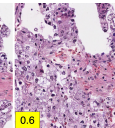
| Strength | AUC | Accuracy | CaseFD w/ HIPT-256 |
|---|---|---|---|
| 0.0 (Frozen) | 79.65 ± 0.03 | 58.12 ± 0.04 | 1139.48 |
| 0.1 | 72.75 ± 0.05 | 55.98 ± 0.03 | 1098.43 |
| 0.2 | 73.26 ± 0.03 | 55.56 ± 0.04 | 1082.73 |
| 0.3 | **73.45 ± 0.03** | **56.41 ± 0.04** | 1060.59 |
| 0.4 | 71.63 ± 0.03 | 55.56 ± 0.04 | 1030.55 |
| 0.5 | 69.77 ± 0.03 | 57.26 ± 0.04 | 986.39 |
| 0.6 | 67.29 ± 0.03 | 55.56 ± 0.06 | 952.91 |
| 0.7 | 67.23 ± 0.04 | 52.56 ± 0.05 | 938.49 |
| 0.8 | 65.93 ± 0.06 | 52.14 ± 0.03 | 923.92 |
| 0.9 | 66.04 ± 0.05 | 48.29 ± 0.03 | 885.19 |
| 1.0 | 60.09 ± 0.04 | 43.59 ± 0.03 | **858.26** |



**Fig. 2.** Performance on CycleDiffusion-restored slides on kidney subtype classification, including Area Under the Curve (AUC) and accuracy, alongside case-wise Fréchet Distance in the HIPT-256 feature space (FD-HIPT256). The higher **strength**, the more added noise and denoising timesteps, the closer to FFPE domain.

across different domains while preserving content and adjusting domain-specific features. While this process improves Fréchet Distance (FD) scores, it decreases the accuracy of downstream tasks, such as classification, as shown in Figure 2.

Given these challenges, our study introduces a novel framework for unpaired FS to FFPE image translation, leveraging the capabilities of LDMs enhanced with Histopathology Pre-Trained Embeddings. This approach aims to produce high-fidelity FS and FFPE images by leveraging text descriptions and pre-trained embeddings. To address the challenge of unpaired translation - where the target FFPE embeddings are absent, we employ a GAN-based U-style Fully Connected Network. This network effectively converts FS embeddings into their FFPE versions, offering improved guidance for generating authentic FFPE images. This novel strategy promises to significantly improve the translation and artifact restoration in FS images, addressing common issues such as FS artifact presence and morphological inaccuracies, thereby improving the accuracy and reliability of histological analysis for clinical assessments.

Our contributions are as follows: 1) Benchmarking latest generative models to address the FS to FFPE image translation challenge. 2) Developing a FS to FFPE image translation framework that improves artifact restoration in FS images, utilizing LDMs with Histopathology Pre-Trained Embeddings, including a mechanism for FS to FFPE embedding translation to overcome the absence of direct FFPE embeddings. 3) Establishing robust evaluation metrics for FS to FFPE image translation , utilizing case-wise Fréchet Distance (CaseFD) within a histopathology pre-trained latent space, alongside downstream classification tasks, evaluating how translation enhances classification accuracy. 4) Significantly outperforming existing state-of-the-art solutions, such as AI-FFPE, UVCGAN2, and CycleDiffusion in downstream classification performance with favorable CaseFD.
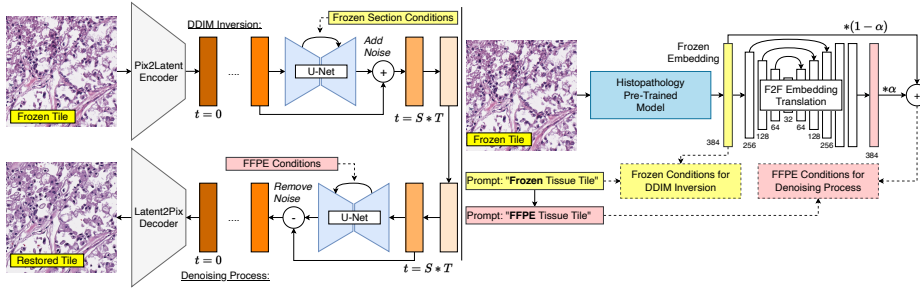
**Fig. 3.** Overview of our FS to FFPE image translation framework.

## 2   Methods

In this study, we present a novel FS to FFPE image translation framework that leverages Latent Diffusion Models (LDMs) and Histopathology Pre-Trained Embeddings. By conditioning LDMs with text descriptions and embeddings, our denoising model effectively captures and translates the unique characteristics of FS and FFPE images with high fidelity. The process begins with transforming FS images into noisy versions via DDIM inversion, tailored with FS-specific data. A denoising U-Net then restores these images, which are further refined into FFPE equivalents using FFPE-specific data and a GAN-based network for embedding translation, as shown in Figure 3. This method not only maximizes the LDMs' potential in producing authentic FFPE images but also focuses on artifact restoration, crucial for improving histopathology image analysis.

### 2.1   Latent Diffusion Models with Histopathology Pre-Trained Embeddings

Our research utilizes Hematoxylin and Eosin (H&E) stained images $x \in \mathbb{R}^{H \times W \times 3}$, applying Latent Diffusion Models (LDMs) enhanced with histopathology pre-trained embeddings for accurate FS to FFPE image translation. FS images $(x_{\text{fs}})$ and Formalin-Fixed Paraffin-Embedded (FFPE) images $(x_{\text{ffpe}})$ are associated with respective text descriptions $p$ and embeddings $e = \mathcal{V}(x)$, where $\mathcal{V}$ denotes models pre-trained on histopathology data, facilitating precise feature capture and translation. Through a pix2latent encoder $(\mathcal{E})$, images are converted into simpler latent representations $(z)$, which are then refined back into detailed images by a latent2pix decoder $(\mathcal{D})$, as described in Figure 3.

The translation process introduces Gaussian noise to these latent representations, incrementally adjusted over $T$ timesteps. A specialized denoising U-Net $(\epsilon_\theta)$, trained to predict and remove the added noise guided by $p$ and $e$, then accurately restores FS/FFPE images. The training loss function for the denoising U-Net is simplified as: $\mathcal{L}_{LDM} = \mathbb{E}_{z,\epsilon \sim \mathcal{N}(0,I),t,p,e} \left[ ||\epsilon - \epsilon_\theta(z_t, t, p, e)||_2^2 \right]$, where $z_t$ denotes the noisy version of $z_0$ at the timestep $t$, sampled uniformly

**Table 1.** Our train, validation, and test sets with three kidney subtypes equally distributed over cases.

| | Train | | Validation | | Test |
|---|---|---|---|---|---|
| #cases | #patches (1024x1024) | #cases | #patches (1024x1024) | #cases | #patches (4096x4096) |
| 180 | 314,114 | 39 | 70,484 | 39 | 636 |

from the set $\{1, ..., T\}$, where $T = 1000$. To ensure efficient training and effective FS/FFPE image generation, we adopt the pre-trained Stable Diffusion XL (SDXL) [14], incorporating extra linear layers for processing pre-trained embeddings and fine-tuning their text encoders and denoising U-Net with Low-Rank Adaptation (LoRA) [7].

### 2.2 Frozen Section to FFPE Image Translation

Leveraging the refined SDXL framework, we developed a FS to FFPE image translation method, $F2F : x_{fs} \rightarrow x_{ffpe}$, using histopathology-specific text descriptions $p$ and embeddings $e$. The process starts with encoding FS images into latent representations, $z_{fs\_0} = \mathcal{E}(x_{fs})$, which are then noised by $\epsilon_\theta(z_t, t, p_{fs}, e_{fs})$ via DDIM inversion with a strength $S \in [0.0, 1.0]$, $T = 50$, and $t = 0, 1, 2, \ldots, S * T$. Next, a U-style fully connected network $G$, trained with WGAN-GP [3] in cycle fashion [20], translates FS embeddings $e_{fs}$ to FFPE equivalents $\hat{e}_{ffpe}$, blending them using an interpolation weight $\alpha$ to maintain FS image identity. The final step involves progressively eliminating the added noise $\epsilon_\theta(z_t, t, p_{ffpe}, \hat{e}_{ffpe})$ from $z_{S*T}$ to $t = 0$, culminating in the latent base $\hat{z}_{ffpe\_0}$, which is then transformed into the FFPE image $\hat{x}_{ffpe} = \mathcal{D}(\hat{z}_{ffpe\_0})$, as shown in Figure 3.

**Enhancing Translation with Classifier-Free Guidance (CFG) and L0 Regularization**. Our exploration into LDMs' capacity for artifact removal in FS images, while maintaining and refining morphology, led us to innovate in guiding the translation process. Increasing the Guidance Scale (GS) enhances the FFPE feature presence and minimize FD between the translated and real FFPE image distributions in latent space. However, an unmoderated increase in GS risks distorting FS image morphology, potentially reducing downstream classification performance, as discussed in Figure 2. Furthermore, reliance on unconditional predicted noise in CFG [6] could introduce artifacts. In pursuit of a balanced translation that respects both FFPE conditions and the inherent creativity of the model, we replace the unconditional noise with noise conditioned by FFPE-translated embeddings and apply L0 Regularization [4] to harmonize the conditional noise with embedding-guided noise as follows:

$$\hat{\epsilon}_t = \epsilon_\theta(z_t, t, \varnothing, \hat{e}_{ffpe}) + GS * \text{prox}_\lambda\big(\epsilon_\theta(z_t, t, p_{ffpe}, \hat{e}_{ffpe}) - \epsilon_\theta(z_t, t, \varnothing, \hat{e}_{ffpe})\big) \quad (1)$$

where $\text{prox}_\lambda(d) = d$ if $|d| > \sqrt{2\lambda}$, and 0 otherwise. $\lambda$ is set as the 70% quantiles of the absolute values of the noise difference, in line with [4]. Note that our model without $\text{prox}_\lambda$ will apply all noise changes.
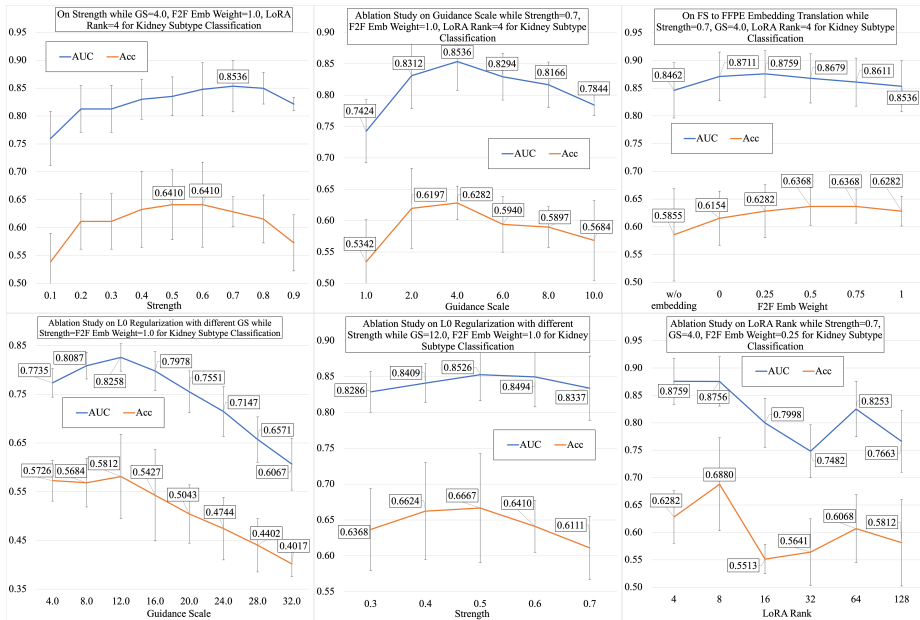
**Fig. 4.** Ablation Studies on classifier-free Guidance Scale (GS), Strength, FS to FFPE Embedding Translation, LoRA rank, and L0 Regularization in Restoration of Artifacts in FS images, evaluated on downstream kidney subtype classfication in macro-averaged Area Under the Curve (AUC) and sample-wise Accuracy (Acc).

## 3  Experiments

**Datasets**. We use the TCGA-Kidney dataset, including 258 cases evenly spread across three kidney subtypes (ccRCC, ChRCC, PRCC) with 516 slides in total. The division is 180 cases for training (70%), 39 for validation (15%), and 39 for testing (15%) (details described in Table 1). Training and validation use $1024 \times 1024$ patches, while $4096 \times 4096$ patches are for slide-level test evaluation [1]. Notably, our evaluation focuses on restoring one-by-one $1024 \times 1024$ patches from these larger test set patches, and only approximate 200 test set patches are utilized for efficiently evaluating our ablation models.

**Training Details**. The FS to FFPE image translation is fine-tuned using the AdamW optimizer with a learning rate of 0.0001, $\beta_1 = 0.9$, $\beta_2 = 0.999$, and a weight decay of 0.01. Training is conducted with a batch size of 1 on $512 \times 512$ images. Ablation models are fine-tuned over 50,000 iterations (approximately 7 hours), while the final models including ours and previous works are refined across 150,000 iterations (21 hours) leveraging NVIDIA RTX A6000 GPUs.

**Evaluation Metrics**. While distribution-based metrics like the Fréchet Distance (FD) are standard for assessing generated images, they do not account for the morphological accuracy crucial for clinical evaluations. To address this, we introduce the Case-wise Fréchet Distance (CaseFD) to compute FD between
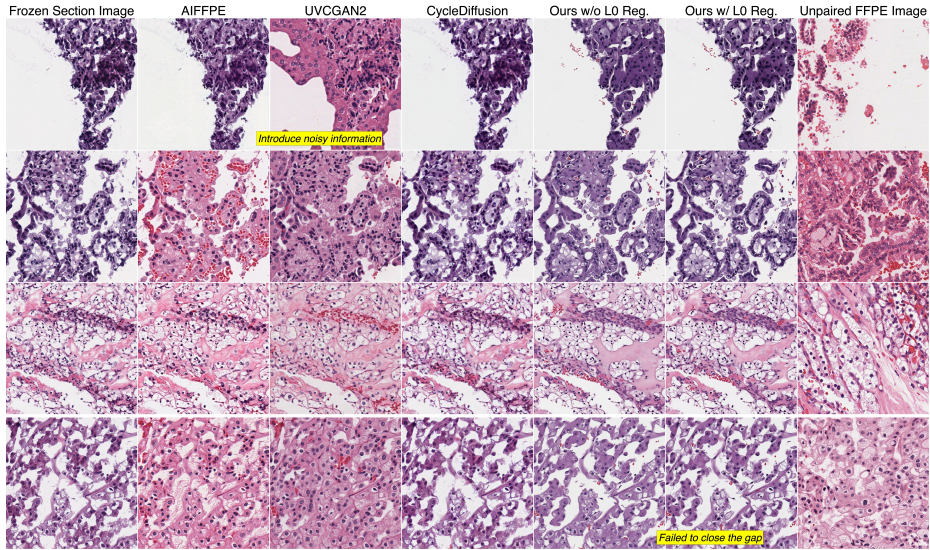
**Fig. 5.** A qualitative comparison between AIFFPE [12], UVCGAN2 [17], CycleDiffusion [19], and ours. Unpaired FFPE images are from a different tissue.

translated and real images within the same case in the latent space by pre-trained models such as DINOv2 ViT-L14 [16], HIPT-256 [1], and ViT-DINO [9]. These models are pre-trained on general images, TCGA FFPE slides, and both TCGA and TULIP datasets, respectively. Moreover, we implement a Fully-Connected Multiple Instance Learning (MIL) network based on the HIPT [1], specifically for kidney subtyping on FFPE slides, using a 6-fold leave-one-out cross-validation from our dataset. This approach assesses whether FS to FFPE image translation enhances kidney subtype classification accuracy, providing an in-depth evaluation of its clinical relevance.

**On Strength, Guidance Scale (GS), and LoRA Rank**. Tuning these hyperparameters is crucial for optimal translation. As shown in the top-left, top-middle, and bottom-right charts of Figure 4, adjusting Strength and GS demonstrates a bell-shaped relationship with AUC and accuracy performance, highlighting their impact on downstream classification tasks. These adjustments reveal a trade-off: while increasing Strength and GS can enhance the FFPE-like appearance of images (see Figure 2), it can also lead to the introduction of new artifacts, negatively affecting performance. A higher LoRA rank enhances domain adaptation but requires more training time, affecting efficiency. After considering both AUC and accuracy, we identify the optimal settings for our final models as $S = 0.7$, $GS = 4.0$, and LoRA rank of 8. A qualitative result can be found in Supplementary Document.

**On Histopathology Pre-Trained Embeddings and Translation**. Integrating histopathology pre-trained embeddings improves FS to FFPE image translation, boosting AUC to 0.8711 from 0.8462. However, FS embeddings alone

**Table 2.** A quantitative comparison between AIFFPE [12], UVCGAN2 [17] and Cy-cleDiffusion [19], and ours. This work outperforms previous works on the downstream kidney subtype classfication in macro-averaged Area Under the Curve (AUC) and sample-wise Accuracy (Acc), while obtaining the favorable Case-wise Fréchet Distances (CaseFD). **Bold**/underlined values indicate **best**/second-best performance.

| Method | AUC | Acc | CaseFD w/ DINOv2 ViT-L14 [16] | CaseFD w/ HIPT-256 [1] | CaseFD w/ ViT-DINO [9] |
|---|---|---|---|---|---|
| FFPE | $94.63 \pm 0.02$ | $88.89 \pm 0.03$ | $\infty$ | $\infty$ | $\infty$ |
| Frozen Section | $81.99 \pm 0.03$ | $61.97 \pm 0.08$ | 546.86 | 1044.24 | 1581.22 |
| AIFFPE [12] | $75.46 \pm 0.04$ | $62.82 \pm 0.03$ | 554.43 | 887.47 | 1243.67 |
| UVCGAN2 [17] | $84.89 \pm 0.01$ | $70.09 \pm 0.03$ | **513.34** | **808.47** | **1205.63** |
| CycleDiffusion [19] | $70.55 \pm 0.02$ | $53.42 \pm 0.05$ | 621.69 | 972.96 | 1486.58 |
| Ours w/o L0 Reg. | $\mathbf{94.64 \pm 0.01}$ | $73.5 \pm 0.06$ | 544.85 | 839.67 | 1235.65 |
| Ours w/ L0 Reg. | $94.26 \pm 0.03$ | $\mathbf{80.34 \pm 0.07}$ | 546.65 | 822.66 | 1240.07 |

do not fully match FFPE image characteristics. To address this, we introduced a translation mechanism using U-style fully-connected layers based on [3, 20], predicting FFPE embeddings while maintaining FS identity through an inter-polation weight $\alpha$. Choosing $\alpha$ based on AUC and accuracy, we found that a $\alpha = 0.25$ further raised AUC to 0.8759 and accuracy to 0.6282. See Supplementary Document for qualitative outcomes.

**On L0 Regularization**. To mitigate new artifacts from increasing Strength and GS, we adopt L0 Regularization on noise difference between conditional noise and embedding-guided noise for image generation guidance, as detailed in Equation 1. This ensures significant FFPE-conditioned changes align with the latent representation, disregarding minor and noisy alterations. Consequently, the translation is more robust, resulting lower $S = 0.5$ - faster translation and higher $GS = 12.0$ - more robust FFPE patterns, as shown in the bottom-left and middle-left charts of Figure 4. Refer to Supplementary Document for a qualitative result.

**Qualitative and quantitative results**. Following our ablation studies, we subsequently train our model until 150,000 iterations with a LoRA rank of 8, applying the translation to all test patches with and without L0 Regulariza-tion for comprehensive evaluation. We conduct a comparative analysis against AIFFPE [12], UVCGAN2 [17], and CycleDiffusion [19], across qualitative results, CaseFB, and downstream kidney classification accuracy (exclusively trained on FFPE images). Qualitatively, UVCGAN2 showed a better performance in mim-icking FFPE image characteristics, notably in filling white gaps with tissue tex-tures. However, this was accompanied by the introduction of additional artifacts that reduced its effectiveness for clinical assessments. Conversely, our models sub-stantially improved histological details without significant artifact introduction. Quantitatively, while UVCGAN2 led in CaseFD (**513.34**, **808.47**, and **1205.63**) and demonstrated improvements of $+\mathbf{2.9}$ in AUC and $+\mathbf{8.3}$ in accuracy, our model with L0 Regularization outperformed in enhancing AUC by $+\mathbf{12.27}$ and accuracy by $+\mathbf{18.55}$, achieving favorable CaseFD results (**546.65**, **822.66**, and **1240.07**), as described in Table 2. More results and a comparison to CycleDiffu-
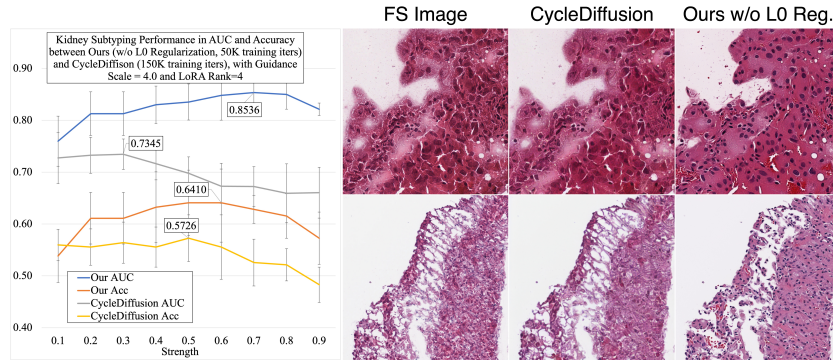
**Fig. 6.** A comparison to CycleDiffusion with the same hyper-parameters.

sion with similar hyper-parameters are in Supplementary Document and Figure 6, respectively.

## 4    Conclusion

We benchmarked the latest Generative Adversarial Networks (GANs) and Latent Diffusion Models (LDMs) to tackle the issues of translating Frozen Section (FS) images to Formalin-Fixed Paraffin-Embedded (FFPE) images and restoring FS artifacts. To address the remaining issues, we introduce an innovative framework utilizes LDMs, enhanced with histopathology pre-trained embeddings and a FS to FFPE embedding translation mechanism, to deliver high-quality image translations that preserve essential histological details for precise clinical assessments. Our model surpasses state-of-the-art methods like AIFFPE [12], UVCGAN2 [17], and CycleDiffusion [19], establishing new standards in FS to FFPE image translation and artifact restoration.

## References

1. Chen, R.J., Chen, C., Li, Y., Chen, T.Y., Trister, A.D., Krishnan, R.G., Mahmood, F.: Scaling vision transformers to gigapixel images via hierarchical self-supervised learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 16144–16155 (June 2022) 6, 7, 8
2. Dubey, S., Kataria, T., Knudsen, B., Elhabian, S.Y.: Structural cycle gan for virtual immunohistochemistry staining of gland markers in the colon. In: Machine Learning in Medical Imaging. pp. 447–456. Springer Nature Switzerland (2023) 2
3. Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., Courville, A.C.: Improved training of wasserstein gans. Advances in neural information processing systems **30** (2017) 5, 8
4. Han, L., Wen, S., Chen, Q., Zhang, Z., Song, K., Ren, M., Gao, R., Stathopoulos, A., He, X., Chen, Y., et al.: Proxedit: Improving tuning-free real image editing with proximal guidance. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 4291–4301 (2024) 5

5. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems **33**, 6840–6851 (2020) 2

6. Ho, J., Salimans, T.: Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598 (2022) 5

7. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685 (2021) 5

8. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1125–1134 (2017) 2

9. Kang, M., Song, H., Park, S., Yoo, D., Pereira, S.: Benchmarking self-supervised learning on diverse pathology datasets. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3344–3354 (June 2023) 7, 8

10. Kang, M., Zhu, J.Y., Zhang, R., Park, J., Shechtman, E., Paris, S., Park, T.: Scaling up gans for text-to-image synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10124–10134 (2023) 2

11. Moghadam, P.A., Van Dalen, S., Martin, K.C., Lennerz, J., Yip, S., Farahani, H., Bashashati, A.: A morphology focused diffusion probabilistic model for synthesis of histopathology images. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 2000–2009 (2023) 2

12. Ozyoruk, K., Can, S., Darbaz, B., Başak, K., Demir, D., Gokceler, I., Serin, G., Hacısalihoglu, P., Kurtuluş, E., Lu, M., Chen, T., Williamson, D., Yılmaz, F., Mahmood, F., Turan, M.: A deep-learning model for transforming the style of tissue images from cryosectioned to formalin-fixed and paraffin-embedded. Nature Biomedical Engineering **6** (12 2022). https://doi.org/10.1038/s41551-022-00952-9 2, 7, 8, 9

13. Park, T., Efros, A.A., Zhang, R., Zhu, J.Y.: 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. pp. 319–345. Springer (2020) 2

14. Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., Rombach, R.: Sdxl: Improving latent diffusion models for high-resolution image synthesis. arXiv preprint arXiv:2307.01952 (2023) 5

15. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10684–10695 (2022) 2

16. Stein, G., Cresswell, J., Hosseinzadeh, R., Sui, Y., Ross, B., Villecroze, V., Liu, Z., Caterini, A.L., Taylor, E., Loaiza-Ganem, G.: Exposing flaws of generative model evaluation metrics and their unfair treatment of diffusion models. Advances in Neural Information Processing Systems **36** (2024) 7, 8

17. Torbunov, D., Huang, Y., Tseng, H.H., Yu, H., Huang, J., Yoo, S., Lin, M., Viren, B., Ren, Y.: Rethinking cyclegan: Improving quality of gans for unpaired image-to-image translation. arXiv preprint arXiv:2303.16280 (2023) 2, 7, 8, 9

18. Torbunov, D., Huang, Y., Yu, H., Huang, J., Yoo, S., Lin, M., Viren, B., Ren, Y.: Uvcgan: Unet vision transformer cycle-consistent gan for unpaired image-to-image translation. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 702–712 (2023) 2

19. Wu, C.H., De la Torre, F.: Unifying diffusion models' latent space, with applications to cyclediffusion and guidance. arXiv preprint arXiv:2210.05559 (2022) 2, 7, 8, 9

20. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: Computer Vision (ICCV), 2017 IEEE International Conference on (2017) 2, 5, 8