

IN2OUT: FINE-TUNING VIDEO INPAINTING MODEL FOR VIDEO OUTPAINTING USING HIERARCHICAL DISCRIMINATOR

Sangwoo Youn¹, Minji Lee², Nokap Tony Park³, Yeonggyoo Jeon³, Taeyoung Na³

¹KAIST, ²Columbia University, ³SK Telecom

ABSTRACT

Video outpainting presents a unique challenge of extending the borders while maintaining consistency with the given content. In this paper, we suggest the use of video inpainting models that excel in object flow learning and reconstruction in outpainting rather than solely generating the background as in existing methods. However, directly applying or fine-tuning inpainting models to outpainting has shown to be ineffective, often leading to blurry results. Our extensive experiments on discriminator designs reveal that a critical component missing in the outpainting fine-tuning process is a discriminator capable of effectively assessing the perceptual quality of the extended areas. To tackle this limitation, we differentiate the objectives of adversarial training into global and local goals and introduce a hierarchical discriminator that meets both objectives. Additionally, we develop a specialized outpainting loss function that leverages both local and global features of the discriminator. Fine-tuning on this adversarial loss function enhances the generator’s ability to produce both visually appealing and globally coherent outpainted scenes. Our proposed method outperforms state-of-the-art methods both quantitatively and qualitatively. Supplementary materials including the demo video and the code are available in [SigPort](#).

Index Terms— Video Outpainting, Hierarchical Discriminator

1. INTRODUCTION

Despite the advances of diffusion models and generative adversarial networks, video outpainting has not been as extensively studied as image outpainting. Image and video outpainting are inherently distinct due to the possible existence of information about the extended region in the other frames of the video.

In contrast, video inpainting, which involves filling in objects or free-form masks within a video in a contextually and temporally consistent manner, has been extensively studied

in [1, 2, 3]. Notably, recent advancements like ProPainter [3] and E²FGVI[2], which propagate features using completed flow and reconstruct the frames through a novel architecture, have shown excellent results in both inpainting background and foreground. The ability of these models to estimate flow and reconstruct objects underscores their potential for outpainting applications. However, directly applying the inpainting model for outpainting is infeasible, producing blurry results, as pointed out in [4] and our results (Fig. 4). While Dehan *et al.* [4] attribute this failure to the inherent problem of outpainting, where less surrounding information is available compared to inpainting, we, *however*, attribute this to the current adversarial loss used in inpainting training. We argue that adequate fine-tuning with discriminators that assess intermediate features can successfully adapt the video inpainting model for outpainting.

In this paper, we propose a novel approach to video outpainting, termed IN2OUT. In order to focus on the challenge of achieving both local perceptual quality and global consistency in generated video regions, we introduce a hierarchical discriminator that leverages the properties of convolutional layers. The early layers assess the local quality of the video, while the deeper layers evaluate global consistency by contrasting various patches of frames. To tailor the layers to the purpose, we introduce an outpainting loss function that operates on local and global features derived from real and generated videos. On the whole, our proposed adversarial framework optimizes the generator’s performance in both local detail and global scene coherence. Since our idea is orthogonal to the generator architecture, our discriminator and generative loss can be used with any video inpainting model.

Our main contributions are as follows:

1. An investigation, through extensive comparisons, into the failures of commonly-used discriminators in video outpainting, highlighting the critical role of the discriminator during fine-tuning;
2. A novel adversarial objective specifically tailored for video outpainting that reduces blur in outpainted regions;
3. The first successful adaptation of a video inpainting model to the outpainting task;
4. Achievement of state-of-the-art performance compared to previous outpainting methods and inpainting baselines.

This research was supported by Culture, Sports, and Tourism R&D Program through the Korea Creative Content Agency(KOCCA) grant funded by the Ministry of Culture, Sports, and Tourism(MCST) in 2024(Project Name: Development of Technology for Convergence Performance Planning and Production Platform to Revitalize the Production of Convergence Performance by Traditional Artist Dance Music, Project Number: RS-2024-00398536, Contribution Rate: 100%)

2. RELATED WORKS

2.1. Video Outpainting

Video outpainting extends the contents of a video frame beyond its original boundaries while preserving the consistency of contents across neighboring frames. There is comparatively less work on video outpainting, and they mostly provide incomplete solutions. Lee *et al.* [5] warps and blends neighboring frames to extend the region based on the observed pixel, but regions that were never visible are left blank. While some video inpainting methods [6] evaluate their methods in video outpainting as well, they perform worse than inpainting.

Background-only methods. Dehan *et al.* [4] use a video object segmentation (VOS) network to detach objects from the background and then employ a flow-based video completion network to generate the background. Jin *et al.* [7] similarly employ a VOS network but stretch the background rather than generating content.

Generative methods. Fan *et al.* [8] propose a masked 3D diffusion model with classifier-free guidance [9] to tackle video outpainting. Recently, Wang *et al.* [10] introduce a diffusion based pipeline comprises input-specific adaptation and pattern-aware outpainting for video outpainting. However, diffusion based methods [8, 10] are limited to process only particular sizes of videos.

2.2. Discriminators in Image/Video Inpainting

Discriminator and generative loss are widely used in image and video inpainting to enhance the perceptual quality of the generated results. Pathak *et al.* [11] first propose to use adversarial loss to alleviate blurry results caused by the pixel-wise reconstruction loss in image inpainting. They use *global* discriminator that looks at an entire image to evaluate the consistency between generated features and real features. To further focus on the perceptual quality of the generated region, Iizuka *et al.* [12] propose to use *partial* discriminator that looks only at the inpainted region together with the global discriminator. Due to the inapplicability of partial discriminator in free-form inpainting, where mask can exist anywhere in any shape, Yu *et al.* [13] propose to apply adversarial loss on the feature maps of the discriminator, instead of the single predicted log likelihood value. Chang *et al.* [1] extends this *global feature* loss to temporal dimension, by using 3-dimensional convolution. This T-PatchGAN discriminator and loss is widely used in video inpainting [14, 2, 15, 16, 3].

3. PROPOSED METHOD

In this section, we propose *hierarchical discriminator* driven from the failures of existing discriminators, and formulate the video outpainting loss.

The spatio-temporal feature discriminator \mathcal{D} learns to classify each patch of given video as real or fake. Given real video,

$x \sim P_Y(x)$ and video generated by generator \mathcal{G} , $z \sim P_{\hat{Y}}(z)$, the general training objective of discriminator is a hinge loss on model output,

$$\mathcal{L}_{\mathcal{D}} = \mathbb{E}_{x \sim P_Y(x)} [\text{ReLU}(1 - \mathcal{D}(x))] + \mathbb{E}_{z \sim P_{\hat{Y}}(z)} [\text{ReLU}(1 + \mathcal{D}(z))].$$

This objective aims to maximize the margin between the real patches and fake patches. The inpainting generator \mathcal{G} is typically trained on multiple objectives including reconstruction loss \mathcal{L}_{rec} and adversarial loss \mathcal{L}_{adv} ,

$$\begin{aligned} \mathcal{L}_{\mathcal{G}} &= \lambda_{\text{rec}} \mathcal{L}_{\text{rec}} + \lambda_{\text{adv}} \mathcal{L}_{\text{adv}} + \lambda_{\text{flow}} \mathcal{L}_{\text{flow}}, \\ \mathcal{L}_{\text{rec}} &= ||\hat{Y} - Y||_1, \\ \mathcal{L}_{\text{adv}} &= -\mathbb{E}_{z \sim P_{\hat{Y}}(z)} [\mathcal{D}(z)]. \end{aligned}$$

Our work focuses on training the discriminator \mathcal{D} and defining the loss functions $\mathcal{L}_{\mathcal{D}}$ and \mathcal{L}_{adv} in a way that effectively adapts the *inpainting* generator $\mathcal{G}_{\theta_{\text{in}}}$ to the *outpainting* generator $\mathcal{G}_{\theta_{\text{out}}}$.

3.1. Hierarchical Discriminator

The purpose of the discriminator can be divided into two: (i) ensuring **global** consistency of the scene, and (ii) ensuring **local** perceptual quality of generated region. The T-PatchGAN discriminator primarily targets the former, a global objective, as per the design in which the receptive field of the last convolutional layer covers an entire video. However, the discriminator that effectively evaluates the local quality of the outpainted region is necessary, especially in an outpainting setting.

While the approach proposed by Iizuka *et al.* [12], which involves the training of two distinct discriminators for each of these objectives, might seem viable, it is fraught with its own set of challenges. Managing multiple GAN losses introduces a delicate balance and sensitivity to hyperparameters, often resulting in training failures (See Sec. 4.5.2).

We believe that a single discriminator can effectively capture both local and global objectives in outpainting scenarios by cleverly leveraging the properties of convolutional layers. As layers progress deeper, the receptive field of features expands, enabling the local features to have a smaller receptive field while the global features perceive the entire video. Thus, the receptive field of earlier layers is restricted to the outpainted region for the pixels at the side, and to the generated region for the pixels in the center, as illustrated in Fig. 1. Motivated by this, we propose a *hierarchical discriminator* $\mathcal{D}_{\text{hierarchical}}$ where the initial layers, termed feature extraction module (FEM), focus on assessing the local quality of video, whereas deeper layers, termed feature comparison module (FCM), focus on comparing the different patches of frames and assess global consistency. (See Fig. 1)

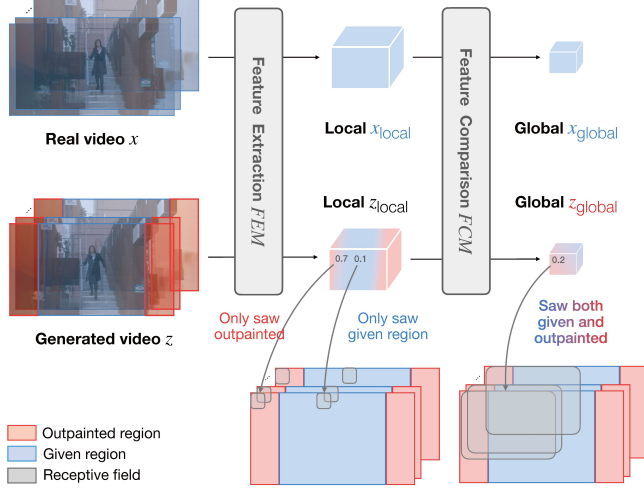


Fig. 1. Outpainting loss calculation using hierarchical discriminator. We employ the local and global features, which are the output of the feature extraction module (FEM) and feature comparison module (FCM), respectively. $\text{Out}(z_{\text{local}})$ in Eq. 1 only saw outpainted region (red).

3.2. Outpainting Loss

We enforce these unique roles of modules via our proposed *outpainting loss* (See Eq. 1) that operates on both local features $x_{\text{local}}, z_{\text{local}}$ and global features $x_{\text{global}}, z_{\text{global}}$. Given real video, $x \sim P_Y(x)$ and video generated by generator \mathcal{G} , $z \sim P_{\hat{Y}}(z)$, the hierarchical discriminator compute features,

$$\begin{aligned} x_{\text{local}} &= \text{FEM}(x), z_{\text{local}} = \text{FEM}(z), \\ x_{\text{global}} &= \text{FCM}(x_{\text{local}}), z_{\text{global}} = \text{FCM}(z_{\text{local}}). \end{aligned}$$

The outpainting loss is defined as

$$\begin{aligned} \mathcal{L}_{\text{out}} &= \mathbb{E}_{x \sim P_Y(x)} [\alpha_{\text{local}} \cdot \text{ReLU}(1 - x_{\text{local}}) \\ &\quad + \alpha_{\text{global}} \cdot \text{ReLU}(1 - x_{\text{global}})] + \\ &\quad \mathbb{E}_{z \sim P_{\hat{Y}}(z)} [\alpha_{\text{local}} \cdot \text{ReLU}(1 + \text{Out}(z_{\text{local}})) \\ &\quad + \alpha_{\text{global}} \cdot \text{ReLU}(1 + z_{\text{global}})]. \end{aligned} \quad (1)$$

Let mask ratio m . For simplicity, we define procedure Out which indicates the outpainted region of video, *i.e.* $\text{Out}(x) = x(i, j)$ such that $i < (m/2) \cdot \text{width}(x)$ or $i > (1 - m/2) \cdot \text{width}(x)$. Here, note that we use $\text{Out}(z_{\text{local}})$ instead of z_{local} . Since the FEM has a small receptive field, z_{local} contains local information of generated inputs. Thus, the discriminator should not be trained to classify the center of z_{local} , which is the feature of given region, as fake. The mapping Out only reflects calculations for the outpainted regions for the local feature. Specifically, FEM is designed to have a receptive field size identical to the size of the outpainted region.

For the video inpainting generator, the adversarial loss is defined as:

$$\mathcal{L}_{\text{adv}} = \mathbb{E}_{z \sim P_{\hat{Y}}(z)} [\text{FCM}(\text{FEM}(z))]. \quad (2)$$

4. RESULTS

In this section, we compare our method with several state-of-the-art video outpainting methods and demonstrate the impact of discriminator design on outpainting performance.

4.1. Settings

Models. Our method is agnostic to the specific generator used. To demonstrate its effectiveness, we fine-tune two video inpainting models with our discriminator and loss: ProPainter [3] and E²FGVI [2], and report the performance improvements. The training details are included in the supplementary material. We compare our method with the background-only [4] and diffusion-based [8, 10] video outpainting models. Additionally, we present the baseline performance of the video inpainting model FuseFormer [15]. Note that FuseFormer, M3DDM, and MOTIA can only process videos with a resolution of 240p, 144p, and 256p, respectively. Therefore, we downsampled the input videos and upsampled the generated videos when evaluating their performance.

Datasets. To assess the performance of our proposed approach, we conduct evaluations on two recognized video datasets: YouTube-VOS [17] and DAVIS [18]. For DAVIS, following Liu *et al.* [15], we evaluate in 50 video clips from the test set. The videos of DAVIS dataset we evaluated are 480p. We fine-tuned our model on the train set of YouTube-VOS dataset resized to 240p. During the evaluation, we used the test set of Youtube-VOS dataset resized to 360p.

Metrics. We choose Peak Signal To Noise Ratio (PSNR), structural similarity index measure (SSIM) [19], and Video Frchet Inception Distance (VFID) [20] to evaluate the quality of the outpainted videos. Note that we compute the metrics in whole video, not only the outpainted region. VFID measures the perceptual similarity between two input videos using a pretrained I3D [21] model and has been widely used in recent video inpainting works.

4.2. Quantitative Results

As shown in Tab. 1, our method demonstrates superior reconstruction performance on both Youtube-VOS and DAVIS datasets compared to SOTA models. Notably, on Youtube-VOS dataset, E²FGVI[2] adapted to outpainting using our IN2OUT method attains a PSNR 1.9dB higher than the baseline and 3.7dB higher than Dehan *et al.*. We also outperform M3DDM [8] and MOTIA [10] with a large margin. These results demonstrate that our method excels in outpainting and successfully adapts the inpainting model for outpainting.

Model		Youtube-VOS			DAVIS		
		PSNR \uparrow	SSIM \uparrow	VFID \downarrow	PSNR \uparrow	SSIM \uparrow	VFID \downarrow
BACKGROUND-ONLY	Dehan <i>et al.</i> [4]	21.99	0.8632	0.085	25.78	0.8901	0.104
DIFFUSION	M3DDM[8]	24.16	0.8862	0.091	24.64	0.8641	0.187
	MOTIA[10]	22.95	0.8795	0.208	24.51	0.8624	0.177
INPAINTING	FuseFormer [15]	23.78	0.7899	0.098	25.55	0.7861	0.193
	ProPainter[3]	22.74	0.9292	0.097	25.14	0.9353	0.144
	E ² FGVI[2]	23.81	0.9378	0.093	24.73	0.9290	0.158
IN2OUT (Ours)	ProPainter	25.18 ($\blacktriangle 2.4$)	0.9399	0.075	27.33 ($\blacktriangle 2.2$)	0.9431	0.115
	E ² FGVI	25.71 ($\blacktriangle 1.9$)	0.9464	0.096	26.61 ($\blacktriangle 1.9$)	0.9385	0.139

Table 1. Quantitative comparisons on Youtube-VOS and DAVIS datasets. Mask ratio is set to 1/4. \uparrow indicates higher is better, and \downarrow indicates lower is better. The **value** in the parentheses indicate the increase in PSNR by fine-tuning inpainting models using our discriminator.

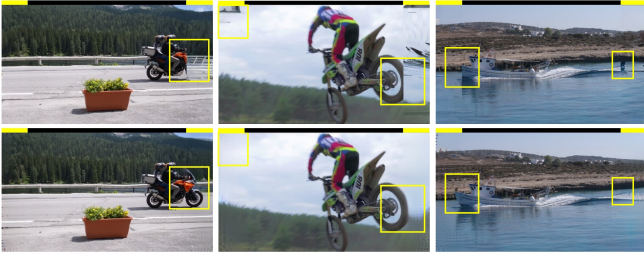


Fig. 2. Qualitative comparisons of Dehan *et al.* [4] (top) and our IN2OUT fine-tuned E²FGVI (bottom) on 480p DAVIS dataset.

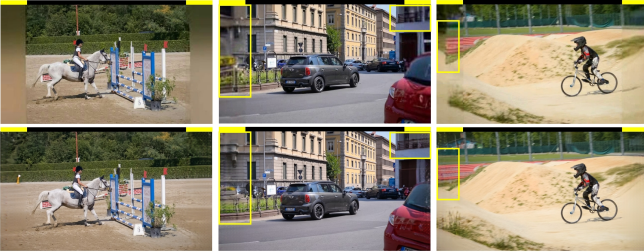


Fig. 3. Qualitative comparisons of M3DDM [8] (top) and our IN2OUT fine-tuned E²FGVI (bottom) on 480p DAVIS dataset.

4.3. Qualitative Results

The yellow line on the top of the video indicates the horizontally outpainted region. Figure 2 visually compares the three outpainted videos of Dehan *et al.* [4] and our method. In the outpainted video by Dehan *et al.*, objects are truncated, whereas our method seamlessly completes objects moving into the outpainted area.

M3DDM fails to generate content and produces blurry results in 15 out of 50 videos from the DAVIS test dataset. Examples of these failures are shown in the leftmost image in Fig. 3. Even in videos where M3DDM successfully outpaints,

our method provides more accurate and complete results, as shown in the two right images in Fig. 3.

As shown in Fig. 4, the baseline inpainting model produces blurry results, especially at the boundaries of the frames. Our method allows the generator to adapt to outpainting, and significantly reduce blurry artifacts.

4.4. Efficiency

Dehan *et al.* takes about 21s/frame to outpaint a 480p video, due to their iterative outpainting scheme. In contrast, our approach employing the end-to-end inpainting model takes about 0.4s/frame over 52 times faster than theirs. Additionally, M3DDM and MOTIA takes about 5s/frame and 24.3s/frame, respectively, even though M3DDM operates on 144p video and MOTIA operates on 256p video, which is order of magnitude slower than inpainting models. This underscores that employing a video inpainting approach for outpainting is a promising and effective strategy.

4.5. Discriminator Designs

In this section, we extensively study the effect of discriminator design on the outpainting adaption of video inpainting model. We used E²FGVI for the inpainting generator. Table 2 compares the designs of each discriminator, starting with *None* where no discriminator is employed during fine-tuning. Other designs include: *Global* where the discriminator processes the entire video and utilizes only the final output to calculate the loss (equivalent to T-PatchGAN discriminator); *Partial-only* that exclusively processes the outpainted region and considers only the final output of the discriminator to evaluate the loss; *Global & partial* that averages the losses from both discriminators to compute the total discriminator loss (equivalent to the discriminator proposed by Iizuka *et al.* [12]); and lastly, *Local-only* that takes the full video as input but exclusively utilizes the local features x_{local} and $\text{Out}(z_{\text{local}})$ of the discriminator to determine the loss.

Discriminator	Youtube-VOS			DAVIS		
	PSNR \uparrow	SSIM \uparrow	VFID \downarrow	PSNR \uparrow	SSIM \uparrow	VFID \downarrow
None	24.53	0.9256	0.115	24.51	0.8984	0.220
Global (T-PatchGAN [1])	24.28	0.9237	0.107	24.04	0.8958	0.166
Partial-only	24.04	0.9083	0.086	25.67	0.9272	0.164
Global & partial [12]	13.11	0.7869	0.181	12.79	0.7709	0.317
Local-only	24.47	0.9179	0.082	25.74	0.9322	0.162
Hierarchical (Ours)	25.71	0.9464	0.096	26.61	0.9385	0.139

Table 2. Quantitative comparisons of discriminator designs on Youtube-VOS and DAVIS datasets. Mask ratio is set to 1/4. \uparrow indicates higher is better, and \downarrow indicates lower is better.

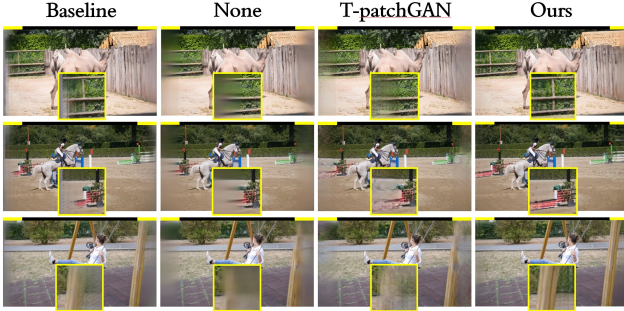


Fig. 4. Qualitative comparisons of discriminator designs on 480p DAVIS dataset.

4.5.1. Quantitative Results

As shown in Tab. 2, all other discriminator designs show a decrease in performance compared to fine-tuning without discriminator, except our proposed design, on Youtube-VOS dataset. The trend is similar on DAVIS dataset, but partial-only and local-only discriminators outperform the fine-tuning without discriminator. The success of discriminators focusing on the quality of outpainted regions highlights the importance of assessing local quality in the outpainting task. Global & partial discriminator shows the most severe performance degradation. We observed from the log that the loss fails to converge, and thus the training fails, underscoring the infeasibility of balancing multiple discriminators targeting different objectives.

4.5.2. Qualitative Results

We also show the visual comparison of the results of the different discriminator designs in Fig. 4. The fine-tuning without a discriminator and the fine-tuning with the global discriminator led to more blurry results compared to the baseline. Compared to other designs, our hierarchical discriminator achieves the most accurate and consistent restoration of foreground and least blurry artifacts, demonstrating the effectiveness of our proposed outpainting loss that considers both local and global objectives in increasing the perceptual quality of the outpainted region.

4.6. Ablation Study on Mask Ratio

Method	Ratio	
	1/3	1/6
Dehan <i>et al.</i> [4]	23.34 / 0.8234	29.46 / 0.9359
M3DDM [8]	21.92 / 0.8146	27.53 / 0.9133
E ² FGVI[2]	22.13 / 0.8790	28.54 / 0.9385
ProPainter	22.61 / 0.8792	28.30 / 0.9736
Ours (E²FGVI)	23.94 / 0.9251	30.05 / 0.9765
Ours (ProPainter)	24.72 / 0.8916	30.70 / 0.9780

Table 3. Comparison of PSNR/SSIM by mask ratios on DAVIS dataset. See Supplementary Sec. 2 for the full comparison including VFID metrics.

Table 3 presents the outpainting performance at varying mask ratios. As the mask ratio increases, the task becomes more challenging, resulting in generally lower performance at a ratio of 1/3 and improved performance at 1/6. In both scenarios, our method achieves higher PSNR and SSIM compared to Dehan *et al.*, M3DDM, E²FGVI baseline and ProPainter baseline. This consistently high performance highlights the robustness of our fine-tuning approach.

5. CONCLUSION

In this work, we propose a novel adversarial framework to fine-tune video inpainting generators to outpainting, paving the effective way to exploit the powerful priors for a task that is comparatively less studied. The contribution also lies in that we extensively compare different discriminator designs, and suggest that the discriminator enforcing both local and global objective may be the missing piece of the successful adaptation of inpainting to outpainting. Our experiments demonstrate that our proposed method outperforms existing video outpainting models in terms of quantitative and qualitative measures. Notably, our discriminator can be integrated into any existing video inpainting model, providing a solid starting point for future research in this domain.

6. REFERENCES

- [1] Ya-Liang Chang, Zhe Yu Liu, Kuan-Ying Lee, and Winston Hsu, “Free-form video inpainting with 3d gated convolution and temporal patchgan,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9066–9075. 1, 2, 5
- [2] Zhen Li, Cheng-Ze Lu, Jianhua Qin, Chun-Le Guo, and Ming-Ming Cheng, “Towards an end-to-end framework for flow-guided video inpainting,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17562–17571. 1, 2, 3, 4, 5
- [3] Shangchen Zhou, Chongyi Li, Kelvin CK Chan, and Chen Change Loy, “Propainter: Improving propagation and transformer for video inpainting,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 10477–10486. 1, 2, 3, 4
- [4] Loïc Dehan, Wiebe Van Ranst, Patrick Vandewalle, and Toon Goedemé, “Complete and temporally consistent video outpainting,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2022, pp. 687–695. 1, 2, 3, 4, 5
- [5] Sangwoo Lee, Jungjin Lee, Bumki Kim, Kyehyun Kim, and Junyong Noh, “Video extrapolation using neighboring frames,” *ACM Transactions on Graphics (TOG)*, vol. 38, no. 3, pp. 1–13, 2019. 2
- [6] Chen Gao, Ayush Saraf, Jia-Bin Huang, and Johannes Kopf, “Flow-edge guided video completion,” in *European Conference on Computer Vision*, 2020, pp. 713–729. 2
- [7] Jun-Gyu Jin, Jaehyun Bae, Han-gyul Baek, and Sanghyo Park, “Object-ratio-preserving video retargeting framework based on segmentation and inpainting,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 497–503. 2
- [8] Fanda Fan, Chaoxu Guo, Litong Gong, Biao Wang, Tiezheng Ge, Yuning Jiang, Chunjie Luo, and Jianfeng Zhan, “Hierarchical masked 3d diffusion model for video outpainting,” in *Proceedings of the 31st ACM International Conference on Multimedia*, 2023, pp. 7890–7900. 2, 3, 4, 5, 1
- [9] Jonathan Ho and Tim Salimans, “Classifier-free diffusion guidance,” *arXiv preprint arXiv:2207.12598*, 2022. 2
- [10] Fu-Yun Wang, Xiaoshi Wu, Zhaoyang Huang, Xiaoyu Shi, Dazhong Shen, Guanglu Song, Yu Liu, and Hongsheng Li, “Be-your-outpainter: Mastering video outpainting through input-specific adaptation,” in *European Conference on Computer Vision*, 2024, pp. 153–168. 2, 3, 4
- [11] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros, “Context encoders: Feature learning by inpainting,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2536–2544. 2
- [12] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa, “Globally and locally consistent image completion,” *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, pp. 1–14, 2017. 2, 4, 5
- [13] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang, “Free-form image inpainting with gated convolution,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4471–4480. 2
- [14] Yanhong Zeng, Jianlong Fu, and Hongyang Chao, “Learning joint spatial-temporal transformations for video inpainting,” in *European Conference on Computer Vision*, 2020, pp. 528–543. 2
- [15] Rui Liu, Hanming Deng, Yangyi Huang, Xiaoyu Shi, Lewei Lu, Wenxiu Sun, Xiaogang Wang, Jifeng Dai, and Hongsheng Li, “Fuseformer: Fusing fine-grained information in transformers for video inpainting,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14040–14049. 2, 3, 4, 1
- [16] Kaidong Zhang, Jingjing Fu, and Dong Liu, “Flow-guided transformer for video inpainting,” in *European Conference on Computer Vision*, 2022, pp. 74–90. 2
- [17] Ning Xu, Linjie Yang, Yuchen Fan, Dingcheng Yue, Yuchen Liang, Jianchao Yang, and Thomas Huang, “Youtube-vos: A large-scale video object segmentation benchmark,” *arXiv preprint arXiv:1809.03327*, 2018. 3
- [18] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung, “A benchmark dataset and evaluation methodology for video object segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 724–732. 3
- [19] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli, “Image quality assessment: from error visibility to structural similarity,” *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004. 3
- [20] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro, “Video-to-video synthesis,” *arXiv preprint arXiv:1808.06601*, 2018. 3
- [21] Joao Carreira and Andrew Zisserman, “Quo vadis, action recognition? a new model and the kinetics dataset,” in *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6299–6308. 3

1. TRAINING DETAILS

We used E²FGVI HQ [2] and ProPainter [3] for a baseline pre-trained generator. The generator and discriminator are trained simultaneously using Adam optimizer for $5 \cdot 10^4$ iterations. Learning rate is set to $4 \cdot 10^{-5}$ for both models. For E²FGVI, we set $\lambda_{\text{rec}} = \lambda_{\text{valid}} = 1$, $\lambda_{\text{flow}} = 0.01$, $\lambda_{\text{adv}} = 0.04$, and $\alpha_{\text{local}} = \alpha_{\text{global}} = 0.5$. For ProPainter, we set the values same as E²FGVI except $\lambda_{\text{flow}} = 1$ and $\lambda_{\text{adv}} = 0.01$. During training, all frames are resized into 432×240 and the number of local frames and non-local frames (See E²FGVI [2]) are set to 5 and 3, respectively. Training took approximately 390 hours on one RTX 4090 GPU when fine-tuning E²FGVI. During evaluation and test, following the previous practices, we use sliding window with the size of 10.

Masks. While our primary target is outpainting 4:3 videos to 16:9 videos ($m = 1/4$), we fine-tuned the generator to mask ratio of minimum 1/12 to maximum 1/3 to increase robustness of the model.

Model architecture. For FEM, we stack three 3D convolutional layers with a spatial stride size of 2. The receptive field is $\approx 2^3 \cdot 7 = 56$ which is similar to the width of the outpainted region when mask ratio $m = 1/4$, 54. For FCM, we also stack three 3D convolutional layers with a spatial stride size of 2. The receptive field is $\approx 2^6 \cdot 7 = 448$ which is larger than the width of the training data, 432.

2. EXTENDED RESULTS

Here we present the VFID results of Tab. 3.

Method	1/3	1/6
Dehan <i>et al.</i> [4]	0.130	0.071
M3DDM [8]	0.277	0.120
E ² FGVI[2]	0.217	0.095
ProPainter[3]	0.193	0.105
Ours (E ² FGVI)	0.204	0.092
Ours (ProPainter)	0.156	0.075

Table 4. VFID by the outpainting ratios on the DAVIS dataset.

3. EXTENDED ABLATION STUDIES

3.1. Ablation on Additional Generator

As shown in Tab. 5, our fine-tuning framework increases the performance of FuseFormer [15] in both PSNR and SSIM metrics, compared to the T-PatchGAN discriminator. Thus, effectiveness of our method is not restricted to E²FGVI[2] and ProPainter[3], and can be used with any video inpainting model.

Discriminator	PSNR	SSIM	VFID
w/o Fine-tuning	25.55	0.7861	0.193
T-PatchGAN [1]	26.06	0.7907	0.167
Ours	26.24	0.7916	0.177

Table 5. Quantitative comparison of discriminator design on DAVIS dataset and FuseFormer [15] generator.

3.2. Flow loss weight

λ_{gen}	λ_{flow}	PSNR \uparrow	SSIM \uparrow	VFID \downarrow
1	0.01	26.61	0.9385	0.139
1	0.1	26.43	0.9375	0.146
1	1.0	26.26	0.9363	0.147

Table 6. Ablation study on the flow loss weight on the DAVIS dataset. Note that E²FGVI baseline is trained to $\lambda_{\text{flow}} = 1$.

As shown in Tab. 6, lower flow weight in generator loss led to a slight increase in all metrics. This is expected since the inpainting task that incorporates object mask during training is better for learning the flow estimation.

3.3. Generative loss weight

α_{inter}	α_{global}	PSNR \uparrow	SSIM \uparrow	VFID \downarrow
0.9	0.1	26.50	0.9383	0.149
0.1	0.9	26.31	0.9365	0.137
0.5	0.5	26.61	0.9385	0.139

Table 7. Ablation study on the local and global loss weight on the DAVIS dataset.

As shown in Tab. 7, different configurations of hyperparameters do not markedly affect the performance in all metrics, highlighting the robustness of our method to hyperparameters.

4. LIMITATION

Fig. 5 shows the failure case when outpainting static video. Our method sometimes blurs (left) or omits (right) the foreground that is never seen in a given region. This shows the continuing challenge of static videos in video outpainting.



Fig. 5. Failure cases when outpainting static videos in 480p DAVIS dataset. The yellow line on the top of the video indicates the horizontally outpainted region.