

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

AttentionFlow: Text-to-Video Editing Using Motion Map Injection Module

Anonymous WACV **Algorithms Track** submission

Paper ID 920

Abstract

Text-to-image diffusion, which has been trained with a large amount of text-image pair dataset, shows remarkable performance in generating high-quality images. Recent research using diffusion model has been expanded for text-guided video editing tasks by using text-guided image diffusion models as baseline. Existing video editing studies have devised an implicit method of adding cross-frame attention to estimate frame-frame attention to attention maps, resulting in temporal consistent editing. However, because these methods use generative models trained on text-image pair data, they do not take into account one of the most important characteristics of video: motion. When editing a video with prompts, the attention map of the prompt implying the motion of the video, such as ‘running’ or ‘moving’, is not clearly estimated and accurate editing cannot be performed. In this paper, we propose the ‘Motion Map Injection’ (MMI) module to perform accurate video editing by considering movement information explicitly. The MMI module provides a simple but effective way to convey video motion information to T2V models by performing three steps: 1) extracting motion map, 2) calculating the similarity between the motion map and the attention map of each prompt, and 3) injecting motion map into the attention maps. Considering experimental results, input video can be edited accurately and effectively with MMI module. To the best of our knowledge, our study is the first method that utilizes the motion in video for text-to-video editing.

1. Introduction

Recent research using text guided diffusion models and large-scale language models has led to unprecedented advances in image generation and image editing. In the field of image editing, various studies have been conducted, including fine-tuning a diffusion model for a few images [1, 2], zero-shot image editing [3, 4], guiding the region to be modified with mask [5, 6], or modifying an image using a text prompt [7–10]. Among them, Prompt-to-Prompt [7] offers

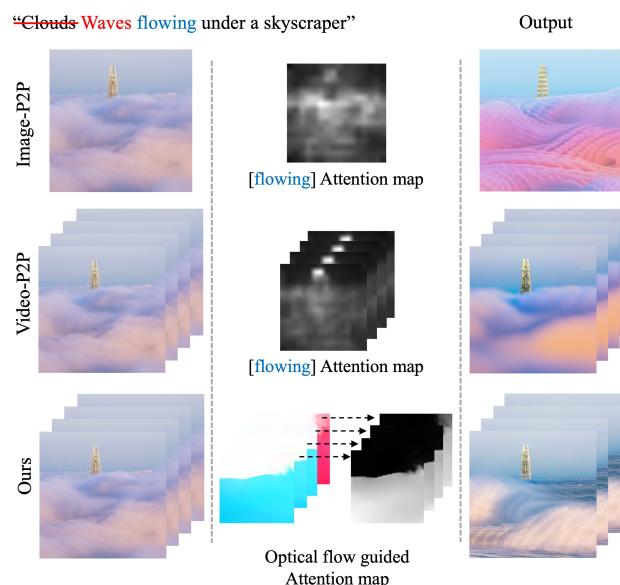


Figure 1. Comparison of video editing output and attention map compared to existing methods. Both Image-P2P and Video-P2P failed to accurately estimate the attention map, resulting in discrepancy with the prompt. Our study performed realistic video editing by enabling accurate attention maps through optical flow guided attention maps.

users intuitive image editing by proposing several methods and applications to control the attention map that identifies the semantic relationship between prompt token and input image.

Research using diffusion model [11–14] has been expanded for text-guided video editing tasks. Most of these studies use text-guided image diffusion models as baseline. However, since a video consists of several frames unlike images, it should be edited in consideration of temporal information. To this end, existing video editing studies [13, 15] have devised an implicit method of adding cross-frame attention to estimate frame-frame attention to attention maps, resulting in temporal consistent editing. Among them, video-P2P [14] is a study that enables video editing as

108 intuitive as image P2P [7] by applying image P2P to video
 109 editing.
 110

111 However, a vast amount of text-video pair dataset is re-
 112 quired to train the implicit structures that uses cross-frame
 113 attention to understand temporal information. However,
 114 there is no such amount of data enough to train those struc-
 115 ture. That is why the most of video editing algorithms have
 116 chosen the method of fine-tuning image diffusion model for
 117 the input video or zero-shot approach to edit video in con-
 118 sideration of temporal information. Because the Image dif-
 119 fusion model was trained on the text-still image pair dataset,
 120 it is not good at estimating the attention map for prompts
 121 containing motion as shown in Figure 1. Therefore, it is
 122 difficult to perform appropriate attention control for the
 123 prompt, which reduces the capability of editing videos. Motion
 124 information needs to be estimated and reinforced ex-
 125 plicitly through an additional network.
 126

127 Optical flow estimation is the task estimating the infor-
 128 mation (i.e. motion) of pixels that have moved between
 129 frames of video, and is the algorithm that has been used
 130 in a wide range of video tasks. In recent years, deep learn-
 131 ing has been applied to optical flow estimation, resulting in
 132 significant performance improvements. Starting with apply-
 133 ing convolution neural network to optical flow, a series of
 134 following studies proposed a coarse-to-fine and iterative es-
 135 timation methodology to reduce computational costs. Most
 136 recently, a methodology with Transformer improved accu-
 137 racy and efficiency of estimating motion, while addressing
 138 the long-standing challenges of motion estimation on oc-
 139 cluded pixels and large displacements. By exploiting optical
 140 flow estimation, which shows high performance and utility
 141 in motion extraction for video tasks, highly accurate motion
 142 can be estimated to be injected into the video editing model.
 143

144 In this paper, we propose a method to complement at-
 145 tention maps with optical flow to make video editing more
 146 effective and accurate. Our proposed framework consists of
 147 the following three steps. Firstly, existing optical flow al-
 148 gorithm estimates motion in input video. Secondly, attention
 149 maps identified by off-the-shelf video diffusion model
 150 are complemented with motion estimated by optical flow.
 151 Lastly, modify the video by applying the complemented at-
 152 tention map.
 153

154 The most important part of this framework is the Mo-
 155 tion Map Injection (MMI) module, which adds optical flow
 156 information to the existing attention map. First, the optical
 157 flow of the input video is estimated, and the magnitude of
 158 this value is used as the motion map. After that, the simi-
 159 larity through template matching between the attention map
 160 and the motion map of all prompts is calculated. At this
 161 point, Normalized Cross Correlation (NCC) is used for simi-
 162 larity. Finally, the motion map is multiplied by the simi-
 163 larity weight, and added to the every attention map of each
 164 prompt. However, the image diffusion model does not accu-
 165 rately estimate the attention map of the motion prompt to-
 166 ken such as 'flowing', so the motion map is directly inserted
 167 instead of template matching. Our proposed MMI module is
 168 a groundbreaking method that can supplement motion infor-
 169 mation in the prompt without learning a vast video dataset,
 170 which can dramatically improve the editing of the existing
 171 video editing framework.
 172

173 Another benefit of our module is additional capability of
 174 video editing with motion control. Using motion, selective
 175 editing of objects moving in a specific direction becomes
 176 possible when there are objects moving in various direc-
 177 tions. In this study, only contents moving in the direction
 178 specified by the user can be edited by rotating the motion
 179 map according to the direction provided by the user before
 180 injecting it. The motion can be estimated in various direc-
 181 tions depending on movements, but in this study, directions
 182 are typically considered for areas in eight directions.
 183

184 In essence, our main contributions include:
 185

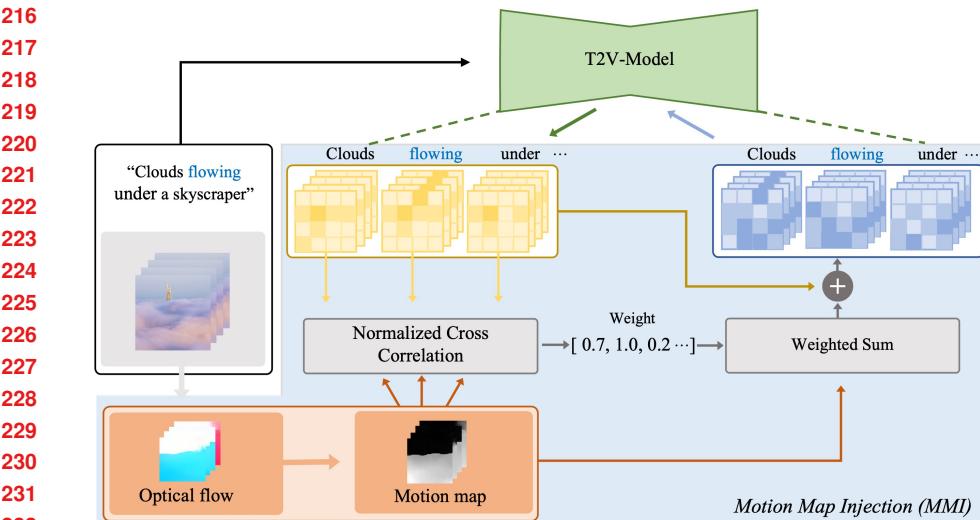
- To the best of our knowledge, our study is the first at-
 186 tempt to use optical flow information in the field of
 187 text-to-video editing.
 188
- Our proposed MMI module gives a way to effectively
 189 inject motion extracted from video into the attention
 190 map of the prompt.
 191
- We found that editing a video with the motion ex-
 192 tracted from the video improves general editing per-
 193 formance and enables selective editing according to the
 194 direction in which the object moves.
 195

2. Related Work

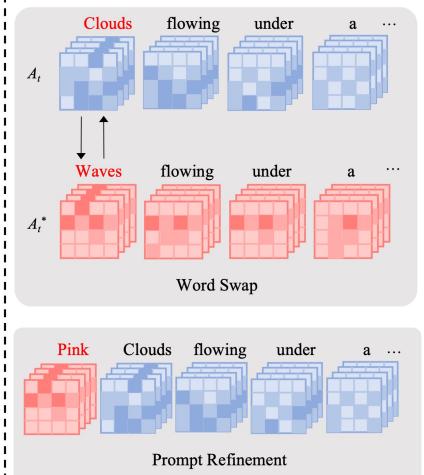
2.1. Text-Guided Editing

196 The diffusion model [16, 17], which has recently been
 197 actively studied, generates data from noise through the pro-
 198 cess of adding or removing noise. Based on this diffusion
 199 model, text-guided image editing models such as DALL-
 200 E2 [18], Imagen [19], and stable diffusion [20] show the
 201 results of high-quality image editing. In particular, Prompt-
 202 to-Prompt presents text-guided image editing that controls
 203 the relationship between the prompt text token and the cor-
 204 responding image pixel with the attention maps, enabling
 205 unprecedented semantic editing. In addition, subsequent pa-
 206 pers such as DreamBooth [1], EDICT [21], and Imagic [2]
 207 have been actively studied recently, showing impressive re-
 208 sults for text-guided image editing.
 209

210 Based on the significant progress of text-guided im-
 211 age editing, research has recently been expanded to text-
 212 guided video editing with the generative model. Dreamix
 213 [11] presents the first diffusion-based method of perform-
 214 ing text-guided motion and application editing of videos
 215 through fine-tuning, but there are difficulties with local-
 216 ized editing by replacing a word. Video-P2P [14] divides
 217



(a) Motion map Guided Text-to-Video Framework



(b) Text-to-Video editing

Figure 2. Overall framework of this study. First, the T2V-Model generates an attention map by receiving video and prompts as input. At the same time, the Motion Map Injection module receives the video frame, generates a motion map, and injects it into the attention map of the T2V-Model. After that, text-to-video editing is performed using the attention map that includes video motion information.

their framework into two branches for unchanged parts and edited parts, and incorporates each attention map to enable detailed editing.

Concurrent to above works, Vid2vid-zero [13] performs stable video reconstruction and editing by adding cross-frame attention to the U-Net structure of the existing diffusion model. In addition, FateZero [22], based on zero-shot, stores an attention map during the inversion process to maintain temporal consistency for structure and motion information. We attempt the first study to extract motion information directly from video and apply it to video editing.

2.2. Optical Flow Estimation

Optical flow estimation is a computer vision task that involves computing the motion of objects in a video sequence. Recently, this field is significantly advanced through the rise of deep neural networks. FlowNet [23] was the first fully convolutional neural network for estimating optical flow. Then, a series of works, represented by SpyNet [24], PWC-Net [25], LiteFlowNet [26], and RAFT [27] were proposed to reduce the computational costs through coarse-to-fine and iterative estimation methodology. Recently GM-Flow [28] were proposed to achieve highly accurate results without relying on a large number of refinements by performing global matching with a Transformer. Among the various optical flow estimation models, we employ recent GMFlow in this work.

Optical flow estimation is used in various video tasks.

First, video action recognition [29, 30] aims to automatically recognize the behavior of objects in video sequences, where optical flow is used as a useful motion representation in video motion representation. Using spatio-temporal information from surrounding scenes to fill in new content for damaged areas, video inpainting [31–33] enables spatio-temporally stable synthesis between frames of video through optical flow. Video super resolution [34–36] is the field of generating high-resolution video frames from low-resolution video frames, and generally maintains temporal consistency between video frames by using optical flow as motion compensation. Video frame interpolation (VFI) [37, 38], a technology that generates an intermediate frame between two consecutive frames, also effectively extracts motion and shape information between frames by utilizing optical flow to estimate motion information between frames. Our work is the first attempt to apply optical flow to text-guided video editing where motion information is important based on the proven validity of the optical flow estimation in various video fields.

3. Proposed Method

Let be input video \mathcal{V} , consists of n frames. We define the source prompt as \mathcal{P} just like the Prompt-to-Prompt setting. In \mathcal{P} , the prompt containing motion information of the video (for example, running, moving, etc.) is called motion prompt \mathcal{P}_M .

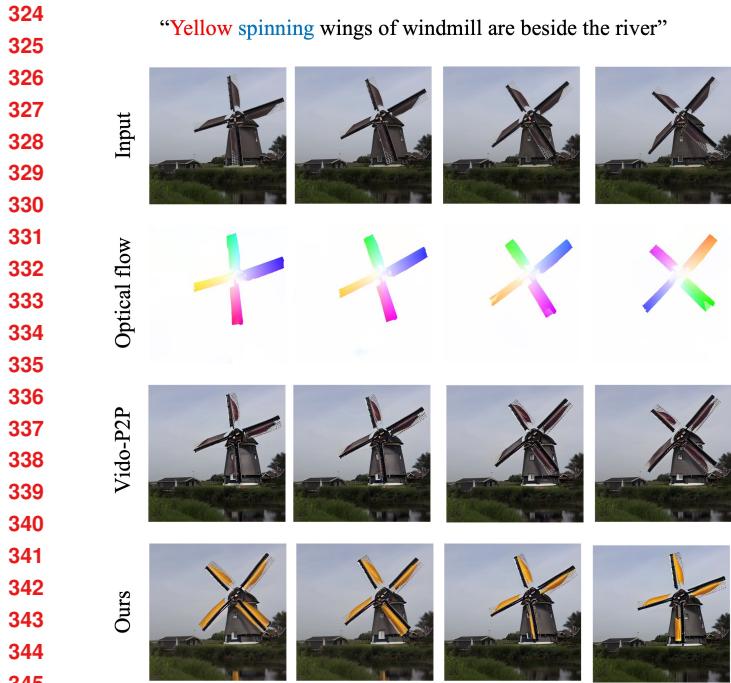


Figure 3. Examples that the existing video editing model could not edit object having motion. However, editing was performed by directly injecting optical flow into the object prompt.

In this section, we provide an overview of our framework, which is illustrated in Fig. 2. Our framework is structured to complement the text-to-video model through the Motion Map Injection (MMI) module. The attention map $\mathcal{A}_{\mathcal{P}}$ representing the correlation between the frame of video \mathcal{V} and the prompt \mathcal{P} is estimated through the cross-attention layer inside the network. The estimated motion map is injected into $\mathcal{A}_{\mathcal{P}}$ through the MMI Module, and then edited video \mathcal{V} with n frames through the editing process of the Text-to-Video Editing network. Our proposed MMI Module consists of two steps for effective video editing. At the first step, pre-trained optical flow estimation network [28] estimates optical flow V_{flow} of input video \mathcal{V} . At the second step, a motion map is obtained by applying L2 norm to the optical flow. By calculating the correlation between the attention map of all prompts and the motion map, it complements the existing attention map according to the correlation. Lastly inject the motion map into the attention map of motion prompt. This framework is a method that enables effective video editing even for motion prompts.

In the following sections, we describe the Prompt-to-Prompt [7] method of editing an image using an attention map through Preliminary (Sec 3.1). In addition, we propose a new method of injecting motion maps into attention maps in MMI module (Sec 3.2.)

3.1. Preliminary

Prompt-to-Prompt P2P [7] based on text-guided diffusion model edits image by modifying source prompt \mathcal{P} . At each diffusion step t , diffusion model predicts noise ϵ from noise image z_t and text embedding $\psi(\mathcal{P})$ through U-Net. Cross-attention layer produces cross-attention map indicating spatial correlation between visual and textual feature. Spatial features of noise image $\phi(z_t)$ are projected on query matrix $Q = \ell_Q(\phi(z_t))$ and key matrix $K = \ell_K(\psi(\mathcal{P}))$ through learned linear projections ℓ_Q and ℓ_K . Attention map \mathcal{A} can be written as

$$A = \text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right), \quad (1)$$

where d denotes latent projection dimension of key and query. Because the attention map A equals matrix product of spatial information of visual feature Q and transpose of textual embedding's key K^T , we can see where the textual feature affects the image. Because P2P can edit images by controlling these attention, it gives intuitive image editing by modifying only the prompt. Recent studies extend Prompt-to-Prompt to video domain from image domain.

3.2. Motion Map Injection Module

Motion Map Extraction In this paper, we use optical flow estimation algorithm to estimate motion of video. Optical flow refers to the amount of movement that appears between the previous frame F_{t-1} and the subsequent frame F_t . Among the various optical flow estimation algorithm, we exploit Unimatch [28]. Firstly, we calculate the correlation for the pixels of the two frames by matrix product and then normalize for the last two dimensions using the softmax function. Then, matching distribution D_{flow} is obtained for each pixel location in F_{t-1} with respect to all pixel locations in F_t .

$$D_{flow} = \text{softmax} \left(\frac{F_{t-1}F_t^T}{\sqrt{N}} \right) \in \mathbb{R}^{H \times W \times H \times W}, \quad (2)$$

where N denotes normalization coefficient to prevent the value from increasing after internal operation. We calculate the weighted average of the matching distribution D_{flow} on the 2D coordinates of the pixel grid G_{2D} to obtain correspondence \hat{G}_{2D} . Finally, the optical flow V_{flow} can be obtained by calculating the difference of corresponding pixel coordinates as follow

$$V_{flow} = \hat{G}_{2D} - G_{2D} \in \mathbb{R}^{H \times W \times 2}. \quad (3)$$

We then apply L2 norm to the generated optical flow V_{flow} to obtain motion map V_{flow}^* . Before proposing our method, we would like to first show the feasibility of our method. Even if the prompt is written more specifically

432 or accurately to be given P2P-based video editing models,
 433 there are times when editing is not accurate. To solve this
 434 problem, we utilized the optical flow estimation task, which
 435 estimates motion information of the video. As a result of
 436 injecting optical flow directly into the attention map of the
 437 motion prompt, it was confirmed that accurate editing was
 438 performed, as shown in Fig 3. Inspired by the fact that sim-
 439 ply injecting motion information can do an edit that existing
 440 video editing models cannot, we construct an algorithm to
 441 edit video by injecting the motion information of the video
 442 into the attention map. Specifically, motion information is
 443 injected into motion prompts (e.g. running, moving, etc.),
 444 which are essential to describe movement, but are not accu-
 445 rately estimated. We also propose a method of propagating
 446 the motion information of the video to the entire prompt in
 447 consideration of the correlation between all prompts except
 448 the motion prompt.
 449

450 **Motion Map Injection using NCC** It is necessary to cal-
 451 culate the similarity of each other to effectively supplement
 452 the existing attention map using the estimated motion map.
 453 To this end, we use the template matching algorithm. We
 454 set the source image to attention map and the target image
 455 to motion map and calculate the Normalized Cross Corre-
 456 lation (NCC) as below

$$457 C_k = \frac{1}{n} \sum_{x,y} \frac{(V_{flow}^*(x,y) - \bar{V}_{flow}^*)(A_k(x,y) - \bar{A}_k)}{\sigma_V \sigma_A}, \quad (4)$$

461 where n denotes the number of pixels, x, y denotes pixel
 462 of source image and target image. V_{flow}^* and \bar{A}_k denote the
 463 average of pixels of V_{flow}^* and A_k . σ_A denotes the standard
 464 deviation of pixels of V_{flow}^* and A_k . k denotes the index
 465 of a particular word in the entire prompt. We calculate the
 466 correlation score by performing template matching at every
 467 denoising step t of the diffusion model, making the sizes of
 468 V_{flow}^* and A_k alike for all k .

469 Finally, as shown in (5), a new attention map is identified
 470 by injecting $C_k \cdot V_{flow}^*$ to attention map of each prompt

$$472 A_k^* = A_k + \lambda \cdot \frac{(C_k * V_{flow}^*)}{t} \quad (5)$$

474 where λ is a hyperparameter for motion map injection rate,
 475 t denotes denoising step. The effect of each hyperparameter
 476 is shown in ablation study of experimental result (Sec 4.5)
 477 for all k .

4. Experiments

4.1. Experimental Setup

482 **Baseline Model** We used the Video-P2P model as a base-
 483 line model, which edits video with text. Video-P2P is a
 484 model that can intuitively edits only with text without ad-
 485 dditional information by manipulating a cross attention map

Table 1. A user study evaluating the performance of Video-P2P and the output of our model conducted on three items: Structure Preserving, Text Alignment, Realism and Quality.

User Preference	Video-P2P (%)	Ours (%)
Structure Preserving	31.25	68.75
Text Alignment	21.87	78.12
Realism & Quality	43.75	56.25

Table 2. Quantitative comparison results. Higher performance was recorded for CLIP Score and Masked PSNR compared to Video-P2P.

CLIP Score		Masked PSNR	
Video-P2P	Ours	Video-P2P	Ours
28.98	29.63	22.09	22.60

representing the relationship between text and video. We added the proposed MMI module to Video-P2P to verify the effectiveness of our framework.

Dataset Experiments were conducted using Davis video dataset, which has been used for various video tasks, and YouTube videos were collected and used to make more experiments.

Implementation details We experimented with RTX 3090 GPUs and set the image resolution to 512×512 just like existing Video-P2P. The number of video frames was fixed to four, and the UniMatch [28] model was used to extract the optical flow.

Evaluation Metrics CLIP score [39] and masked PSNR [14] were used to evaluate textual similarity and region preservation. CLIP Score is a metric which can be used to evaluate the correlation between a prompt and image. It has been found to be highly correlated with human judgement. Masked PSNR is a metric which evaluates the quality of structure preservation only for the masked region. In this experiment, we exploited non-moving region as a mask.

4.2. User Study

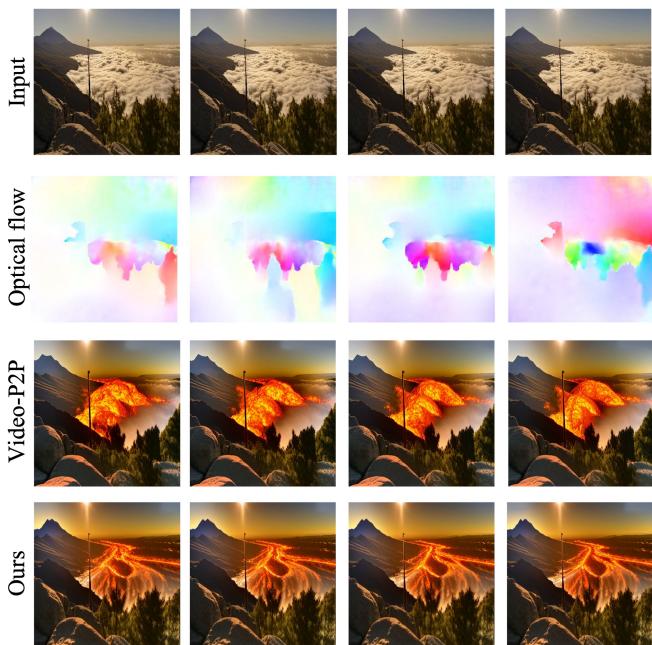
To compare our proposed model with Video-P2P, we conducted user study on three criteria : structure preserving, text alignment, and realism & quality. We presented input video, target prompt, and videos edited with each model to users. Between two edited videos, the evaluation of outputs was conducted by making them choose more faithful videos for the presented criteria. The results filled in Table 1 show that the output of our module is more preferable.

4.3. Quantitative Results

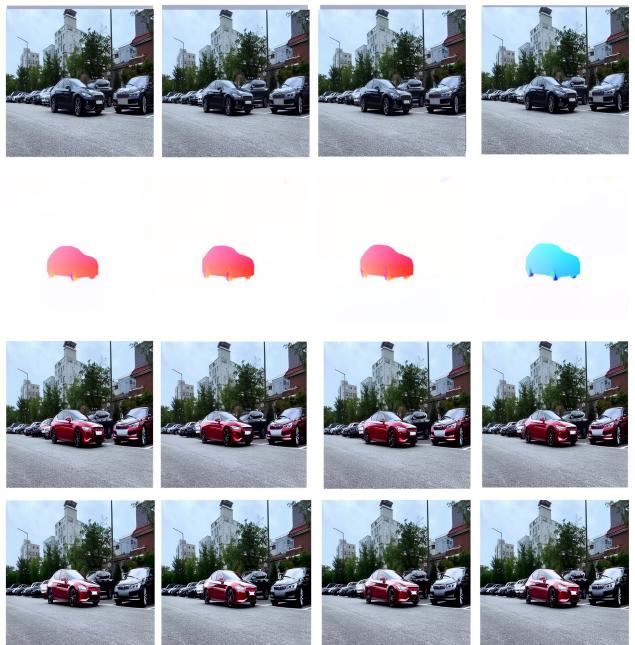
The quantitative results of this study can be found in Table 2. Clip Score was measured by comparing target

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

“There is flowing clouds lava under the mountains”

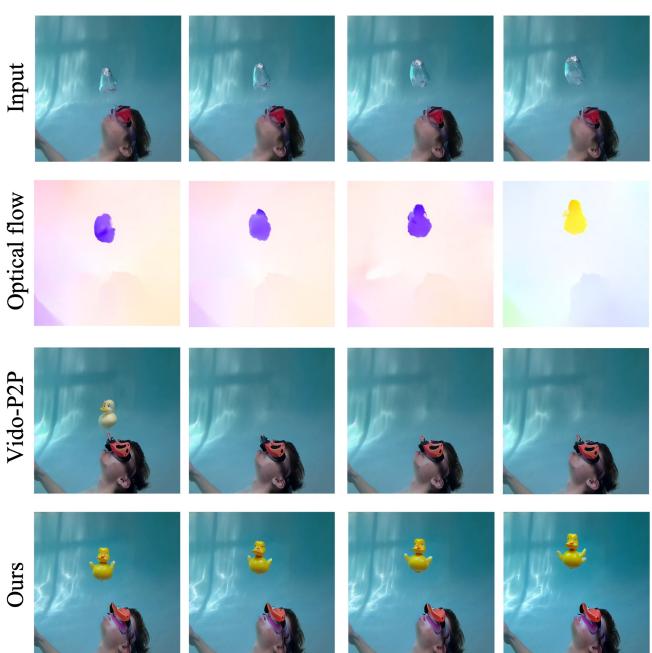


“There is a red driving car in the parking lot”



594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

“Floating bubble duck on the face”



“There is colorful spreading water on the man”

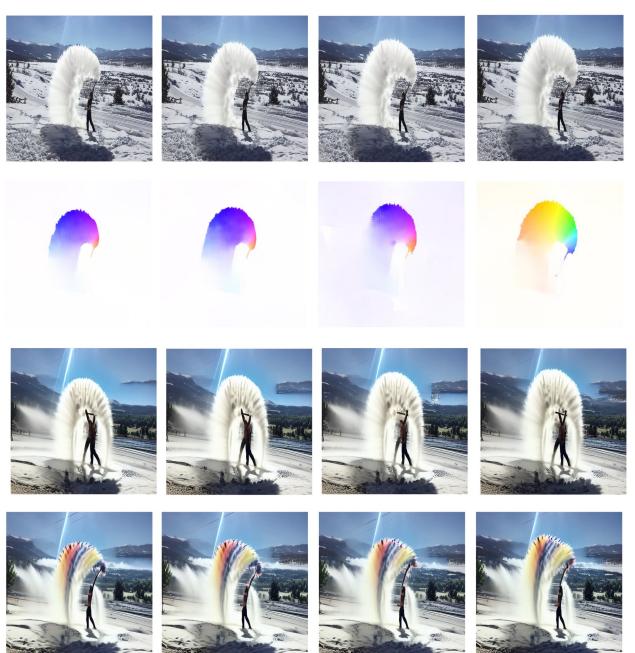


Figure 4. First experimental results of our study. It shows the result of four examples corresponding to each target prompt. The top row of each example is the input video frame, the second row is the optical flow extracted through the MMI module, the third row is the output of Video-P2P, and the last row is the result of our research. It can be seen that the result of the model receiving the motion information of the image through Optical Flow is better. It also shows good performance in editing not only fixed-shaped objects but also irregular objects (water droplets, snow, etc.).

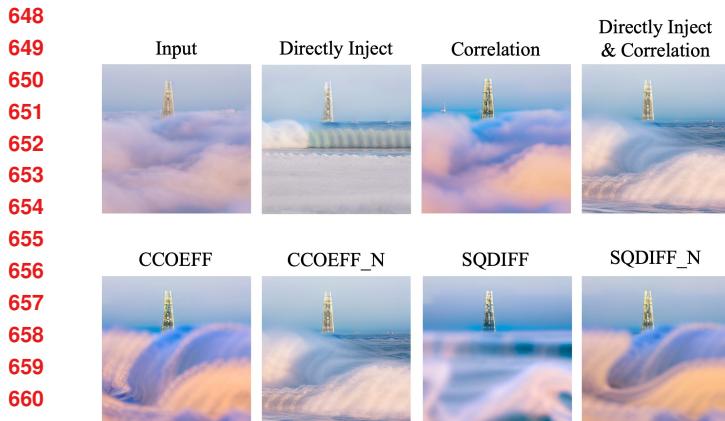


Figure 5. Output of various methods to measure the correlation between the motion map and the attention map.



Figure 6. Editing method for objects moving in the direction specified by the user. Before editing, the user first selects one of the 8 directions.

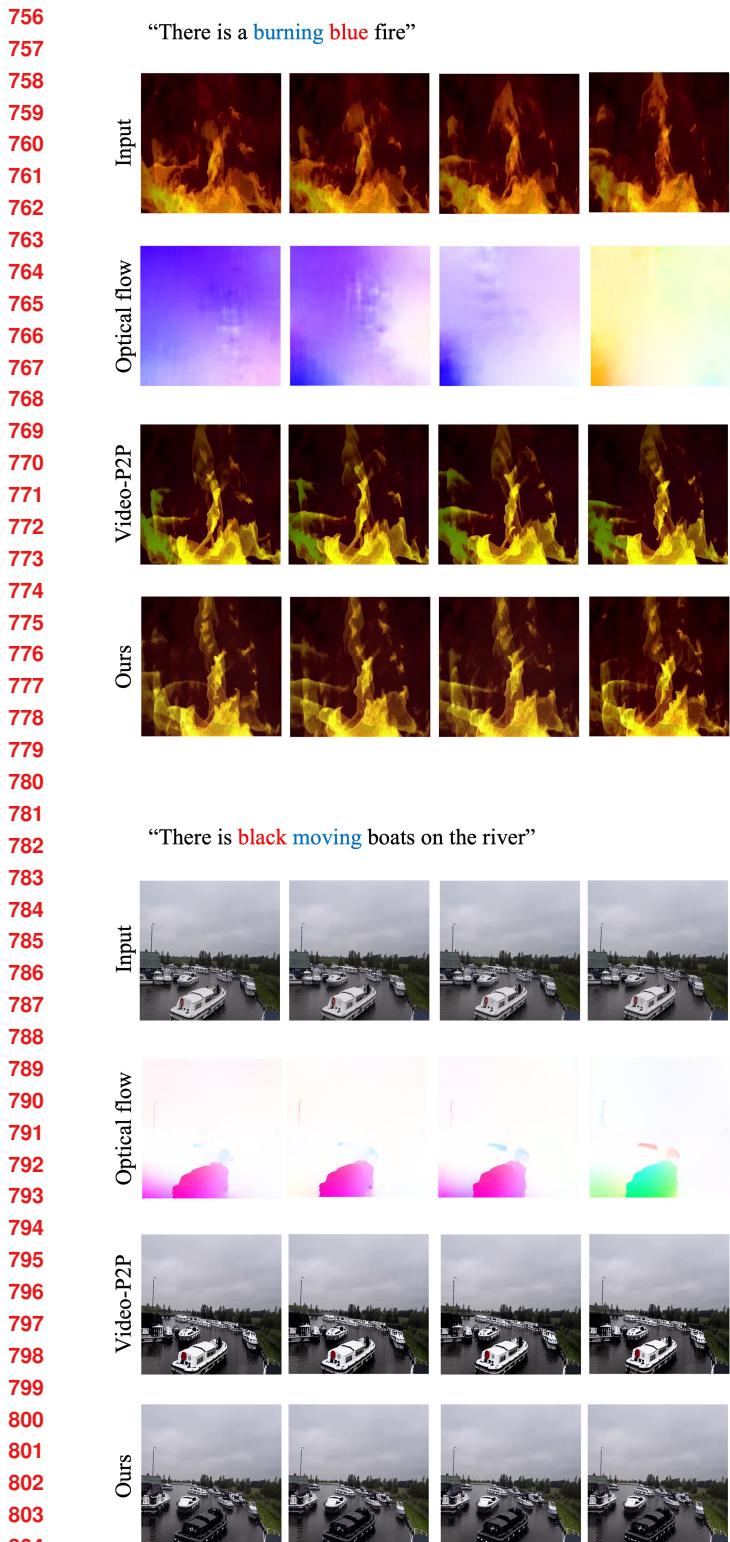
prompt with edited image. It measures how well edit image representative the meaning of edit prompt. The results of this study confirmed that our method was better understood meaning of prompt compared to other comparison model. For Masked PSNR, the performance of our framework scored slightly higher. This means that the existing model did more edit in areas that were not intended than the model we propose. It means that it can accurately edit to the object that you want to edit, this is measured by applying a mask to the edited part.

4.4. Qualitative Results

Fig. 4 shows qualitative results of Video-P2P algorithm and our proposed MMI algorithm. Previous methods tend to produce structurally unnatural images, but our method produces structurally coherent and content-preserving frames. In ‘Lava’ video, only part of the cloud area of the Input frame is activated, and when edited with the video-P2P, it remains as a cloud without changing to lava. When applying ours, the corresponding area could be estimated through the optical flow of the cloud and edited the entire cloud to lava. In ‘Car’ video, when edited through the Video-P2P model, even though prompt contain “driving”, the parked car is also caught in attention map of car, so it can be seen that some other car is edited red. Using our model, editing was applied only to moving cars that we want to refer to by estimating optical flow for moving cars through our model. In ‘Bubble’ video, the attention map of the bubble did not accurately represent the area because the information that “floating” was not apprehend. The editing was done only in first frame and the duck disappeared in the other frames. Through our model, the entire frame was edited through enhanced attention by estimating and injecting bubble’s motion map. Editing was not made in Video-P2P due to inaccurate attention map of “spreading” not being grasped. In our model, it changed when we injected motion map.

4.5. Ablation Study

Fig. 5 is the results for different injection methods of motion map. “Directly Inject” injects the motion map directly into the attention map of the motion prompt. “correlation” injects motion map to entire prompt’s attention maps through the correlation between the each word attention maps and the motion map. “Directly Inject & Correlation” is combination of above described two methods. Specifically, “Correlation” is firstly applied, and then, “Directly Injection” is applied to the motion prompt. To inject the motion map into all words, the correlation between the attention maps of the input prompt and motion map is calculated. We applied various functions of template matching that represent a correlation between the two images. The “CCOEFF” applied to the first image used correlation coefficient, and the second image is the result of “CCOEFF_N” that normalizes it. The third image used “SQDIFF”, a sum of squared differences, and the fourth image used “SQDIFF_N”, which was normalized. As you can see in Fig. 5, the most suitable function to reinforce semantically editable is seen as “CCOEFF_N”, which helps to edit semantically by varying weights depending on the degree of association of the prompter word. Therefore, we choose “CCOEFF_N” in method which is same as NCC.



805 Figure 7. Video-p2p and our video editing results and estimated
806 optical flow in inaccurate motion estimation

807

808

809

4.6. Application

Optical Flow Rotation Since V_{flow} has information on the magnitude of pixel movement between frames, as well as the direction in which pixel moved between frames, the user can select and edit the motion value in the desired direction. Our model allows the user to edit contents in a specific direction by rotating the optical flow V_{flow} according to the direction D provided by the user before injecting it. You can find the result in Fig. 6. Optical flow can be estimated to having various directions depending on movement, but in this study, D were dealt with only eight directions. The specific method would be explained written in the supplementary material.

4.7. Limitations

Accurate motion estimation of input video is essential for editing using optical flow. Therefore, even if optical flow is used, the bad results as shown in Fig. 7 may be obtained when it is difficult to estimate motion information from an image. The optical flow for the movement of fire could not be estimated, so there was no difference from the existing Video-P2P. In addition, in the example of boat, motion was estimated only for ships that occupy a large area of the image, and small ships were not estimated. If the estimated motion of the optical flow for input video is not accurate, it is confirmed that our model, like existing Video-P2P, is difficult to perform accurate editing.

5. Conclusion

There have been a problem that existing video editing algorithms applying an implicit method for the image diffusion model to see the relationship between frames do not understand the information between frames, ‘motion’. We propose a Motion Map Injection (MMI) module to inject estimated motion map into the attention map of the image diffusion model.

Our proposed MMI module injects the motion extracted from input video by the additional network into the every attention map of each prompt estimated by the image diffusion model. To this end, the motion map is multiplied by the similarity (defined as NCC) between the attention map and the motion map. Only the attention map of the motion prompt, which is difficult to estimate due to the limitations of the image diffusion model, is replaced entirely by the estimated motion map.

Injecting motion map into attention map with our proposed MMI module makes general video editing performance improved because attention map of motion prompt become apparent. This is proven by improved evaluation metrics. In addition, selective editing according to direction of motion is enabled.

864

References

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- [1] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023.
- [2] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6007–6017, 2023.
- [3] Chen Henry Wu and Fernando De la Torre. Making text-to-image diffusion models zero-shot image-to-image editors by inferring “random seeds”. In *NeurIPS 2022 Workshop on Score-Based Methods*.
- [4] Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun Li, Jingwan Lu, and Jun-Yan Zhu. Zero-shot image-to-image translation. *arXiv preprint arXiv:2302.03027*, 2023.
- [5] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18208–18218, 2022.
- [6] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance. *arXiv preprint arXiv:2210.11427*, 2022.
- [7] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022.
- [8] Jooyoung Choi, Yunjey Choi, Yunji Kim, Junho Kim, and Sungroh Yoon. Custom-edit: Text-guided image editing with customized diffusion models. *arXiv preprint arXiv:2305.15779*, 2023.
- [9] Zhixing Zhang, Ligong Han, Arnab Ghosh, Dimitris N Metaxas, and Jian Ren. Sine: Single image editing with text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6027–6037, 2023.
- [10] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023.
- [11] Eyal Molad, Eliahu Horwitz, Dani Valevski, Alex Rav Acha, Yossi Matias, Yael Pritch, Yaniv Leviathan, and Yedid Hoshen. Dreamix: Video diffusion models are general video editors. *arXiv preprint arXiv:2302.01329*, 2023.
- [12] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Weixian Lei, Yuchao Gu, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. *arXiv preprint arXiv:2212.11565*, 2022.

- [13] Wen Wang, Kangyang Xie, Zide Liu, Hao Chen, Yue Cao, Xinlong Wang, and Chunhua Shen. Zero-shot video editing using off-the-shelf image diffusion models. *arXiv preprint arXiv:2303.17599*, 2023.
- [14] Shaoteng Liu, Yuechen Zhang, Wenbo Li, Zhe Lin, and Jiaya Jia. Video-p2p: Video editing with cross-attention control. *arXiv preprint arXiv:2303.04761*, 2023.
- [15] Shuai Yang, Yifan Zhou, Ziwei Liu, and Chen Change Loy. Rerender a video: Zero-shot text-guided video-to-video translation. *arXiv preprint arXiv:2306.07954*, 2023.
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [17] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [18] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [19] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022.
- [20] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.
- [21] Bram Wallace, Akash Gokul, and Nikhil Naik. Edict: Exact diffusion inversion via coupled transformations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22532–22541, 2023.
- [22] Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei, Xintao Wang, Ying Shan, and Qifeng Chen. Fatezero: Fusing attentions for zero-shot text-based video editing. *arXiv preprint arXiv:2303.09535*, 2023.
- [23] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2758–2766, 2015.
- [24] Anurag Ranjan and Michael J Black. Optical flow estimation using a spatial pyramid network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4161–4170, 2017.
- [25] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8934–8943, 2018.

- 972 [26] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. Lite- 1026
973 flownet: A lightweight convolutional neural network for 1027
974 optical flow estimation. In *Proceedings of the IEEE conference 1028
975 on computer vision and pattern recognition*, pages 8981– 1029
976 8989, 2018. 1030
- 977 [27] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field 1031
978 transforms for optical flow. In *Computer Vision–ECCV 1032
979 2020: 16th European Conference, Glasgow, UK, August 23– 1033
980 28, 2020, Proceedings, Part II 16*, pages 402–419. Springer, 1034
981 2020. 1035
- 982 [28] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezatofighi, 1036
983 Fisher Yu, Dacheng Tao, and Andreas Geiger. Unify- 1037
984 ing flow, stereo and depth estimation. *arXiv preprint arXiv: 1038
985 2211.05783*, 2022. 1039
- 986 [29] Shuyang Sun, Zhanghui Kuang, Lu Sheng, Wanli Ouyang, 1040
987 and Wei Zhang. Optical flow guided feature: A fast and 1041
988 robust motion representation for video action recognition. In 1042
989 *Proceedings of the IEEE conference on computer vision and 1043
990 pattern recognition*, pages 1390–1399, 2018. 1044
- 991 [30] Laura Sevilla-Lara, Yiyi Liao, Fatma Güney, Varun Jampani, 1045
992 Andreas Geiger, and Michael J Black. On the integration 1046
993 of optical flow and action recognition. In *Pattern Recogni- 1047
994 tion: 40th German Conference, GCPR 2018, Stuttgart, Ger- 1048
995 many, October 9–12, 2018, Proceedings 40*, pages 281–297. 1049
996 Springer, 2019. 1050
- 997 [31] Rui Xu, Xiaoxiao Li, Bolei Zhou, and Chen Change Loy. 1051
998 Deep flow-guided video inpainting. In *Proceedings of the 1052
999 IEEE/CVF Conference on Computer Vision and Pattern 1053
1000 Recognition*, pages 3723–3732, 2019. 1054
- 1001 [32] Dahun Kim, Sanghyun Woo, Joon-Young Lee, and In So 1055
1002 Kweon. Deep video inpainting. In *Proceedings of the 1056
1003 IEEE/CVF Conference on Computer Vision and Pattern 1057
1004 Recognition*, pages 5792–5801, 2019. 1058
- 1005 [33] Kaidong Zhang, Jingjing Fu, and Dong Liu. Flow-guided 1059
1006 transformer for video inpainting. In *Computer Vision–ECCV 1060
1007 2022: 17th European Conference, Tel Aviv, Israel, Octo- 1061
1008 ber 23–27, 2022, Proceedings, Part XVIII*, pages 74–90. 1062
1009 Springer, 2022. 1063
- 1010 [34] Zhigang Tu, Hongyan Li, Wei Xie, Yuanzhong Liu, Shifu 1064
1011 Zhang, Baoxin Li, and Junsong Yuan. Optical flow for video 1065
1012 super-resolution: a survey. *Artificial Intelligence Review*, 1066
1013 55(8):6505–6546, 2022. 1067
- 1014 [35] Mehdi SM Sajjadi, Raviteja Vemulapalli, and Matthew 1068
1015 Brown. Frame-recurrent video super-resolution. In *Pro- 1069
1016 ceedings of the IEEE conference on computer vision and pattern 1070
1017 recognition*, pages 6626–6634, 2018. 1071
- 1018 [36] Longguang Wang, Yulan Guo, Li Liu, Zaiping Lin, Xinpu 1072
1019 Deng, and Wei An. Deep video super-resolution using hr 1073
1020 optical flow estimation. *IEEE Transactions on Image Pro- 1074
1021 cessing*, 29:4323–4336, 2020. 1075
- 1022 [37] Zhen Li, Zuo-Liang Zhu, Ling-Hao Han, Qibin Hou, Chun- 1076
1023 Le Guo, and Ming-Ming Cheng. Amt: All-pairs multi-field 1077
1024 transforms for efficient frame interpolation. In *Proceedings 1078
1025 of the IEEE/CVF Conference on Computer Vision and Pat- 1079
tern Recognition*, pages 9801–9810, 2023.