

000
001
002
003
004
005
006
007
008
009
010
011
012
013054
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

Enhancing Text-to-Video Editing with Motion Map Injection

Anonymous ICCV submission

Paper ID 27

Abstract

Based on the remarkable performance of text-to-image diffusion models, text-guided video editing studies recently have been expanded. Existing video editing studies have introduced an implicit method of adding cross-frame attention to estimate inter-frame attention, resulting in temporally consistent videos. However, because these methods use models pre-trained on text-image pair data, they do not handle unique property of video: motion. When editing a video with prompts, the attention map of the prompt implying the motion of the video (e.g. ‘running’, ‘moving’) is prone to be poorly estimated, which causes inaccurate video editing. To address this problem, we propose the ‘Motion Map Injection’ (MMI) module to consider motion explicitly. The MMI module provides text-to-video (T2V) models a simple but effective way to convey motion in 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. Given experimental results, input video can be edited accurately 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. Extensive experimental results are in <https://currycurry915.github.io/Attention-Flow/>

1. Introduction

Recent research on 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 [1, 3, 4, 5, 6, 7, 8, 9, 10, 11] have been conducted. Among them, Prompt-to-Prompt (P2P) [1] offers 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 on diffusion model has been expanded to text-guided video editing tasks. Most of these studies [2, 12, 13, 14] use text-guided image diffusion models for baseline. However, since a video consists of several frames unlike images, it should be edited in consideration of tempo-

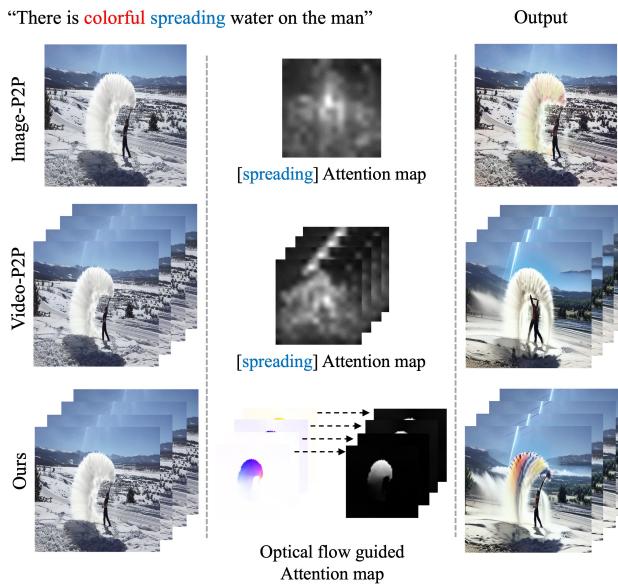


Figure 1: Comparison of attention map and edited video output derived from each existing method. Image-P2P [1] and Video-P2P [2] failed to estimate the attention map, resulting in discrepancy between prompt and video. Our method edited video realistically by exploiting optical flow guided attention maps.

ral information. To this end, existing video editing studies [14, 15, 16] have devised an implicit method of adding cross-frame attention to estimate inter-frame attention, resulting in temporal consistent editing. Recently, starting with Video-P2P [2] and vid2vid-zero [16], research on image based P2P [1] have been expanded to video editing.

A vast amount of text-video pair dataset is required to train the implicit structures that use cross-frame attention to understand temporal information. However, there are not enough of data to train them. That is why the most of video editing algorithms have chosen the fine-tuning method of image diffusion model [12] for the input video or zero-shot approach [14] to edit video in consideration of temporal information. Because the image diffusion model was trained on the text-image pair dataset, it is not good at estimating the attention map for prompts containing motion as shown

108 in Fig 1. As a result, it is difficult to perform appropriate
 109 attention control for the prompt, which reduces the capability
 110 of editing videos. One solution is to estimate motion
 111 information explicitly.
 112

113 Optical flow is the information of pixels that have moved
 114 between frames of video (i.e. ‘motion’). There are a lot of
 115 optical flow estimation methods that have been used in wide
 116 range of video tasks. With existing state-of-the-art optical
 117 flow estimation network [17], we can obtain highly accurate
 118 motion which can be injected into video editing model.
 119

120 In this paper, we propose a framework to complement at-
 121 tention maps with optical flow to make video editing more
 122 effective and accurate. It consists of the following three
 123 steps. Firstly, existing optical flow estimation algorithm
 124 [17] estimates motion from the input video. Secondly, with
 125 the motion, attention maps of off-the-shelf video diffusion
 126 model are complemented. Lastly, modify the video by ap-
 127 plying the improved attention map.
 128

129 The most important part of this framework is first and
 130 second steps. We propose ‘Motion Map Injection (MMI)’
 131 module for these steps. Specifically, MMI module has also
 132 three steps. Firstly, module estimates optical flow of the
 133 input video and derives motion map from the flow. The
 134 motion map is magnitude of the optical flow. Secondly,
 135 the similarities between the motion map and attention maps
 136 of each prompt are computed with template matching al-
 137 gorithm [18]. For template matching, Normalized Cross
 138 Correlation (NCC) [19] is used. Finally, the motion map is
 139 multiplied by the similarity weights, and added to the every
 140 attention map of each prompt. However, the image diffu-
 141 sion model does not accurately estimate the attention map of
 142 the motion prompt, so the attention map of motion prompt
 143 such as ‘moving’ is directly replaced with the motion map
 144 without template matching. Our proposed MMI module is
 145 a efficient method that can supplement motion information
 146 in the prompt without learning a vast video dataset, which
 147 can dramatically improve the editing of the existing video
 148 editing framework.
 149

2. Proposed Method

150 Let be input video \mathcal{V} which consists of n frames. We de-
 151 fine the source prompt as \mathcal{P} just like the Prompt-to-Prompt
 152 [1] setting. In \mathcal{P} , the prompt containing motion information
 153 of the video is called motion prompt \mathcal{P}_M (e.g. ‘running’,
 154 ‘moving’).
 155

156 In this section, we provide an overview of our frame-
 157 work, which is illustrated in Fig. 2. The attention map $\mathcal{A}_{\mathcal{P}}$
 158 representing the correlation between the frame of \mathcal{V} and \mathcal{P}
 159 is estimated through the cross-attention layer in the T2V-
 160 Model. Our proposed MMI module estimates the optical
 161 flow V_{flow} of \mathcal{V} using the pre-trained optical flow estima-
 162 tion network [17] in the first step. In the second step, a
 163 motion map is obtained by applying L2 norm to the optical
 164 flow. After calculating the correlation between the attention
 165 map of all prompts and the motion map using NCC, it com-
 166plements the existing attention map according to the corre-
 167 lation. Lastly, the motion map is inserted into the Attention
 168 Map of the motion prompt \mathcal{P}_M . Using this approach, the
 169 MMI module provides the optical flow information of the
 170 video to the T2V-Model, enabling improved video editing
 171 capabilities.
 172

173 **2.1. Preliminary**

174 **Prompt-to-Prompt** P2P [1] based on text-guided diffusion
 175 model edits image by modifying source prompt \mathcal{P} . Cross-
 176 attention layer inside the diffusion model produces cross-
 177 attention map indicating spatial correlation between visual
 178 and textual feature. Spatial features of noise image $\phi(z_t)$
 179 are projected on query matrix $Q = \ell_Q(\phi(z_t))$ and key ma-
 180 trix $K = \ell_K(\psi(\mathcal{P}))$ through learned linear projections ℓ_Q
 181 and ℓ_K . Attention map A can be written as
 182

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

183 where d denotes latent projection dimension of K and Q .
 184 Because P2P [1] can edit images by controlling these at-
 185 tentions, it gives intuitive image editing by modifying only
 186 \mathcal{P} . P2P [1] can edit images by replacing A of the word to
 187 be edited in \mathcal{P} with the attention of the target prompt or by
 188 providing an additional A .
 189

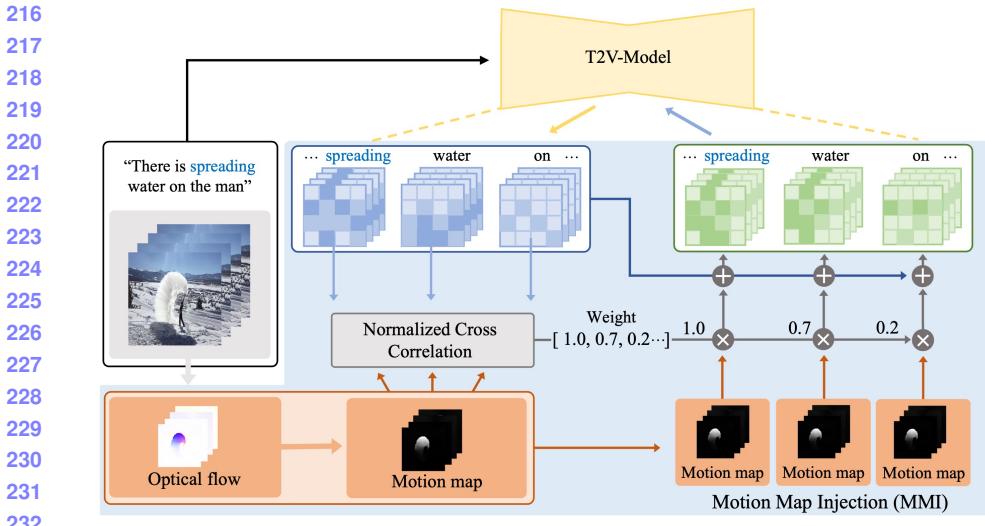
2.2. Motion Map Injection Module

190 **Motion Map Extraction** In this paper, we use pre-trained
 191 optical flow estimation algorithm, UniMatch [17], to esti-
 192 mate motion of video. Firstly, we calculate the correlation
 193 for the pixels of the two frames by matrix product and then
 194 normalize for the last two dimensions using the softmax
 195 function. Then, matching distribution D_{flow} is obtained for
 196 each pixel location in F_{t-1} with respect to all pixel loca-
 197 tions in F_t

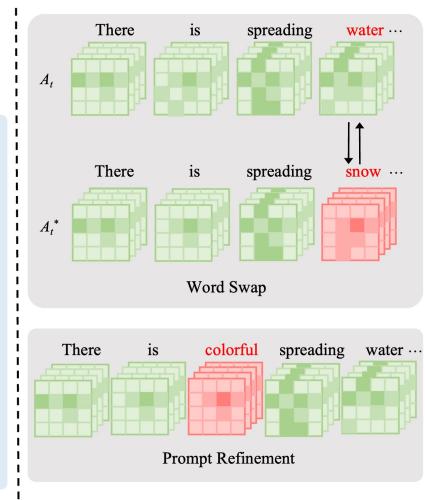
$$D_{flow}(i, j, k, l) = \mathcal{S} \left(\frac{\sum_{d'=1}^d F_{t-1}(i, j, d') \cdot F_t(k, l, d')}{\sqrt{N}} \right), \quad (2)$$

198 where N denotes normalization coefficient to prevent the
 199 value from increasing after internal operation. \mathcal{S} represents
 200 the softmax function. We calculate the weighted average of
 201 the matching distribution D_{flow} on the 2D coordinates of
 202 the pixel grid G_{2D} to obtain correspondence \hat{G}_{2D} . Finally,
 203 the optical flow V_{flow} can be obtained by subtracting the
 204 two grids. We obtain the motion map V_{flow}^* by applying the
 205 L2 norm to the optical flow V_{flow} .
 206

207 **Motion Map Injection using NCC** It is necessary to cal-
 208 culate the similarity of each other to effectively supplement
 209 the existing attention map using the estimated motion map.
 210 To this end, we use the template matching algorithm [19].
 211



(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.

We set the source image to attention map and the target image to motion map and calculate the Normalized Cross Correlation (NCC) as below

$$C_k = \frac{1}{n} \sum_{x,y} \frac{(V_{flow}^*(x,y) - \bar{V}_{flow}^*)(A_k(x,y) - \bar{A}_k)}{\sigma_V \sigma_A}, \quad (3)$$

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

Finally, as shown in (4), a new attention map is identified by injecting $C_k \cdot V_{flow}^*$ to attention map of each prompt

$$A_k^* = A_k + \lambda \cdot \frac{(C_k \cdot V_{flow}^*)}{t}, \quad (4)$$

where λ is a hyperparameter for motion map injection rate, t denotes denoising step. Here, the image diffusion model does not accurately estimate the attention map of the motion prompt token such as 'flowing', so the motion map is directly inserted instead of template matching.

3. Experiments

3.1. Experimental Setup

Baseline model We used the Video-P2P [2] and vid2vid-zero [16] for a baseline model, which can intuitively edits a video with text by manipulating a cross attention map representing the relationship between text and video. We added

Table 1: User study to evaluate the performance of Video-P2P [2] and our model on three criteria: 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

MMI module to both video editing models to verify the effectiveness of our framework.

Dataset Experiments were conducted with Davis video dataset [20]. Also, we collected and used YouTube videos to make more experiments.

Implementation details We experimented with NVIDIA RTX A6000 GPUs and set the resolution to 512×512 like Video-P2P [2] and vid2vid-zero [16]. The number of video frames was fixed to four, and the UniMatch [17] model was used to extract the optical flow.

Evaluation Metrics CLIP Score [21] and masked PSNR [2] were used to evaluate textual similarity and region preservation. CLIP Score [21] is a metric 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 to evaluate structure preservation only for the masked region. In this experiment, we exploited non-moving region as a mask.

3.2. User Study

To compare our proposed model with Video-P2P [2] and vid2vid-zero [16], we conducted user study on three

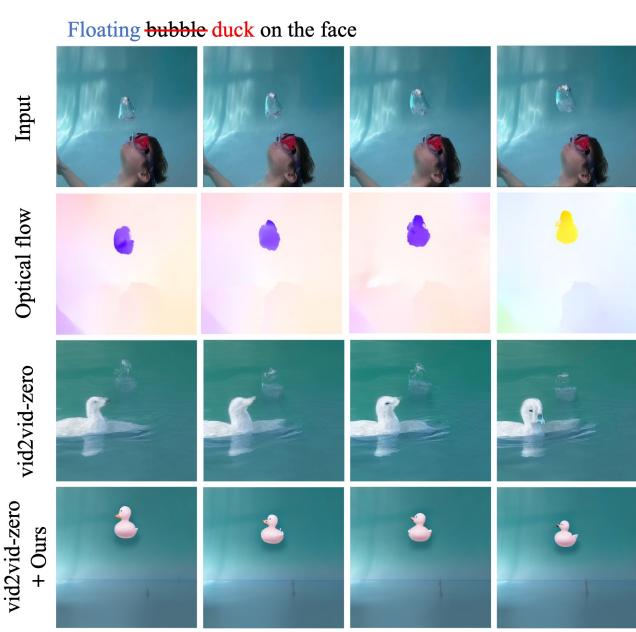
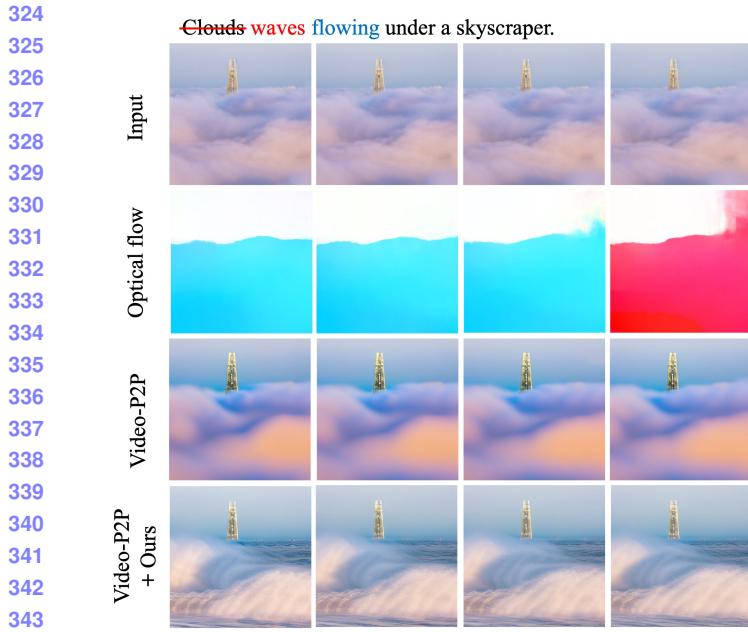


Figure 3: Experimental results of our study. This shows the result of two examples corresponding to each different model [2, 16]. The top row is the input video frame, the second row is the optical flow extracted through the MMI module, the third row is the output of Video-P2P [2] and vid2vid-zero [16], 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.

Table 2: Quantitative comparison results. Higher performance was recorded for CLIPScore [21] and Masked PSNR compared to Video-P2P [2] and vid2vid-zero [16].

	Video-P2P	Video-P2P + Ours
CLIP Score	28.98	29.63
Masked-PSNR	28.88	29.38
	vid2vid-zero	vid2vid-zero + Ours
CLIP Score	22.09	22.60
Masked-PSNR	21.80	25.30

361 criteria: structure preserving, text alignment, and realism
362 & quality. We presented input video, target prompt, and
363 videos edited with each model to users. Between two edited
364 videos, the outputs were evaluated by making users choose
365 more faithful videos for the presented criteria. The results
366 filled in Table 1 show that the output of our module is more
367 preferable.

3.3. Quantitative Results

371 The quantitative results of this study are in Table 2. CLIP
372 Score [21] was measured by comparing target prompt with
373 edited image. The results confirmed that our method was
374 better at understanding meaning of prompt compared to the
375 other models. For masked PSNR, the performance of our
376 framework scored slightly higher. This means that the ex-
377 isting models edited more in areas not intended than our

proposed model.

3.4. Qualitative Results

Figure 3 shows the qualitative results of the Video-P2P [2] algorithm, the vid2vid-zero algorithm [16], and the proposed MMI algorithm. Previous methods tend to produce structurally unnatural images, but our method produces structurally coherent and content-preserving frames. In Video-P2P [2], due to inaccurate estimation of the attention map for the motion prompt, it was unable to perform editing that aligns with the target prompt. However, with the injection of the motion map using our proposed method, we achieve better performance in editing. In vid2vid-zero [16], the performance is lower compared to Video-P2P model [2] that fine-tuning on input video. Nevertheless our model enhances attention and enables control via target prompts. You can see more results on the project page attached to the abstract and in the zipped file of the supplementary.

4. Conclusion

We propose a Motion Map Injection (MMI) module to inject estimated motion map into the attention map of the image diffusion model. 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. As a future work, we will implicitly apply optical flow to the video editing module.

432 In this supplementary material, we describe related
 433 works, a method for calculating optical flow, an ablation
 434 study, a method for optical flow rotation, further explanation
 435 of the metric, and a github link to the code.
 436

437 A. Related Work

438 A.1. Text-Guided Editing

439 The diffusion model [22, 23], which has recently been
 440 actively studied, generates data from noise through the process
 441 of adding or removing noise. Based on this diffusion
 442 model, text-guided image editing models such as DALL-E2 [24],
 443 Imagen [25], and stable diffusion [26] show the
 444 results of high-quality image editing. In particular, Prompt-
 445 to-Prompt [1] presents text-guided image editing that
 446 controls the relationship between the prompt text token and
 447 the corresponding image pixel with the attention maps, en-
 448 abling unprecedented semantic editing. In addition, sub-
 449 sequent papers such as DreamBooth [3], EDICT [27], and
 450 Imagic [4] have been actively studied recently, showing im-
 451 pressive results for text-guided image editing.
 452

453 Based on the significant progress of text-guided im-
 454 age editing, research has recently been expanded to text-
 455 guided video editing with the generative model. Dreamix
 456 [12] presents the first diffusion-based method of perform-
 457 ing text-guided motion and application editing of videos
 458 through fine-tuning, but there are difficulties with local-
 459 ized editing by replacing a word. Video-P2P [2] divides
 460 their framework into two branches for unchanged parts and
 461 edited parts, and incorporates each attention map to enable
 462 detailed editing.
 463

464 Concurrent to above works, vid2vid-zero [14] performs
 465 stable video reconstruction and editing by adding cross-
 466 frame attention to the U-Net structure of the existing dif-
 467 fusion model. In addition, FateZero [28], based on zero-
 468 shot, stores an attention map during the inversion process
 469 to maintain temporal consistency for structure and motion
 470 information. We attempt the first study to extract motion
 471 information directly from video and apply it to video editing.
 472

473 A.2. Optical Flow Estimation

474 Optical flow estimation is a computer vision task that in-
 475 volves computing the motion of objects in a video sequence.
 476 Recently, this field is significantly advanced through the rise
 477 of deep neural networks. FlowNet [29] was the first fully
 478 convolutional neural network for estimating optical flow.
 479 Then, a series of works, represented by SpyNet [30], PWC-
 480 Net [31], LiteFlowNet [32], and RAFT [33] were proposed
 481 to reduce the computational costs through coarse-to-fine
 482 and iterative estimation methodology. Recently GMFlow
 483 [17] were proposed to achieve highly accurate results with-
 484 out relying on a large number of refinements by performing
 485 global matching with a Transformer.



486
 487
 488
 489
 490
 491
 492
 493
 494
 495
 496
 497
 498
 499
 500
 501
 502
 503
 504
 505
 506
 507
 508
 509
 510
 511
 512
 513
 514
 515
 516
 517
 518
 519
 520
 521
 522
 523
 524
 525
 526
 527
 528
 529
 530
 531
 532
 533
 534
 535
 536
 537
 538
 539

Figure 4: Outputs of various methods to measure the correlation between the motion map and the attention maps.

503 Optical flow estimation is used in various video tasks.
 504 First, video action recognition [34, 35] aims to automatically
 505 recognize the behavior of objects in video sequences,
 506 where optical flow is used as a useful motion representa-
 507 tion in video motion representation. Using spatio-temporal
 508 information from surrounding scenes to fill in new content
 509 for damaged areas, video inpainting [36, 37, 38] enables
 510 spatio-temporally stable synthesis between frames of video
 511 through optical flow. Video super resolution [39, 40, 41]
 512 is the field of generating high-resolution video frames from
 513 low-resolution video frames, and generally maintains tem-
 514 poral consistency between video frames by using optical
 515 flow as motion compensation. Video frame interpolation
 516 (VFI) [42, 43], a technology that generates an intermediate
 517 frame between two consecutive frames, also effectively ex-
 518 tracts motion and shape information between frames by uti-
 519 lizing optical flow to estimate motion information between
 520 frames. Our work is the first attempt to apply optical flow
 521 to text-guided video editing where motion information is
 522 important based on the proven validity of the optical flow
 523 estimation in various video fields.

524 B. Ablation Study

525 Fig. 4 shows the results for different injection meth-
 526 ods of motion map. “Directly Inject” injects the motion
 527 map directly into the attention map of the motion prompt.
 528 “correlation” injects motion map to entire prompt’s atten-
 529 tion maps through the correlation between the each word
 530 attention maps and the motion map. “Directly Inject &
 531 Correlation” is combination of above described two methods.
 532 Specifically, “Correlation” is firstly applied, and then, “Di-
 533 rectly Injection” is applied to the motion prompt. To inject
 534 the motion map into all words, the correlation between the
 535 attention maps of the input prompt and motion map is cal-
 536 culated. We applied various functions of template matching
 537 that represent a correlation between the two images. The

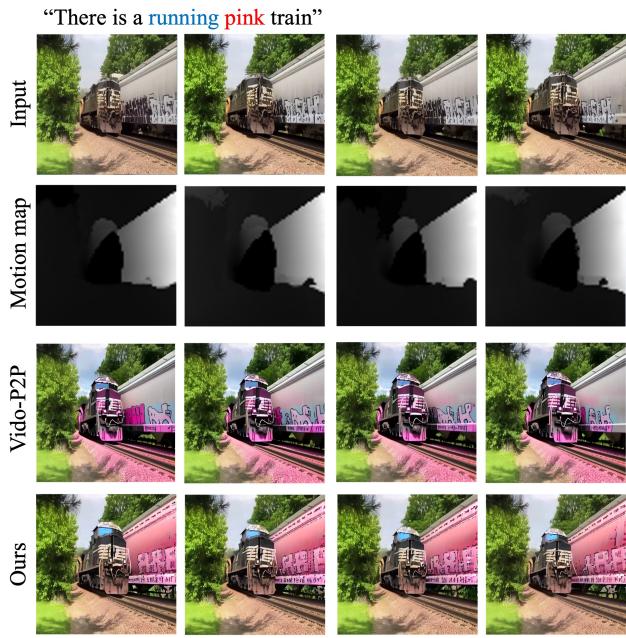


Figure 5: Editing method for objects moving in the direction specified by the user. Before editing, the user first selects one of the 8 directions.

“CCOE” applied to the first image used correlation coefficient, and the second image is the result of “CCOE_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. 4, the most suitable function to reinforce semantically editable is seen as “CCOE_N”, which helps to edit semantically by varying weights depending on the degree of association of the prompter word. Therefore, we choose “CCOE_N” in method which is same as NCC.

C. Method for 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. We propose a method to edit using information on the direction in which pixels move in optical flow representing information on pixel motion between previous frame F_{t-1} and present frame F_t . Our proposed model receives one of eight directions from the user, including Northeast (NE), Southeast (SE), Southwest (SW), and Northwest (NW), which are made from a combination of four directions in the 2D coordinate system.

$$\begin{bmatrix} X' \\ Y' \end{bmatrix} = \begin{bmatrix} X \cos \theta_D & -Y \sin \theta_D \\ X \sin \theta_D & Y \cos \theta_D \end{bmatrix}, \quad (1)$$

where θ_D denotes the angle between axis in frame and user provided direction D . X, Y denotes each motion vector in axis X , and Y . X', Y' denote the motion vectors rotated by θ in each axis.

After rotating the optical flow for the user-provided direction D , only the region of pixels with positive directional motion of the x axis in the rotated coordinate system is specified. The value of the motion map is extracted from the specific region and video edit is performed on the corresponding area. The results can be seen in Fig. 5.

D. Detailed Description of Evaluation Metrics

CLIPS core [21] is calculated in the CLIP model [44], generating embedding vectors for input images and prompts. CLIP Score [21] is measured by computing the cosine similarity between image and caption embedding. We measured how close the target prompt and the edited video frames are semantically in the CLIP Score[21]. We measured the CLIP Score [21] between target prompt and each edited video frame, and quantitatively compared the performance of our model and other models by the average of the scores measured per each frame. The CLIP Score [21] is calculated with the following equation.

$$extCLIPScore(F_t, \mathcal{P}^*) = max(100 * cos(E_{F_t}, E_{\mathcal{P}^*}), 0), \quad (2)$$

where F_t denote the t th edited frame, and \mathcal{P}^* denotes the target prompt. We use official ViT-Base-Patch16 CLIP model.

Masked PSNR To evaluate whether our proposed model performs undesired edit out of target region to be edited, we measured masked PSNR (M.PSNR) proposed by Video-P2P [2]. On this purpose, we measure how much the external region of the target region has changed from the frame of the original video. In consideration of the averaged attention mask sequence M of the changed object, we measure masked PSNR by computing the pixel distance in the out-of-target regions of the edited video V^* and the input video V ,

$$M.PSNR(V^*, V) = PSNR(B(V^*, M), B(V, M)), \quad (3)$$

according to Video-P2P [2], $B(V, M) = V_M$ is defined as a reversed mask binary function, so only regions not to be changed are involved in measuring masked PSNR.

“There is a **burning blue** fire”

Input

Optical flow

Video-P2P

Ours

“There is **black** moving boats on the river”

Figure 6: Results of Video-p2p [2] and our video editing model with inaccurately estimated optical flow

D.1. Ablation Study

D.2. 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. 6 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 Video-P2P [2]. 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 [2], is difficult to perform accurate editing.

E. Code Descriptions

Our code is based on PyTorch version of Video-P2P [2]. We use Video-P2P [2] to edit videos. We set the parameters as follows: frame_size_h = 512, frame_size_w = 512, number of frames = 4,

Code is available at <https://anonymous.4open.science/r/AttentionFlow-197C/README.md>

F. Additional Qualitative Results

Additional experimental results and code can be found in the supplementary archive zipped with the supplementary paper.

References

- [1] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. Prompt-to-prompt image editing with cross-attention control. In *The Eleventh International Conference on Learning Representations*, 2023. [1](#), [2](#)
 - [2] 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. [1](#), [3](#), [4](#), [2](#)
 - [3] 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. [1](#)
 - [4] 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. [1](#)
 - [5] 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*. [1](#)

- 756 [6] Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun
757 Li, Jingwan Lu, and Jun-Yan Zhu. Zero-shot image-to-image
758 translation. In *ACM SIGGRAPH 2023 Conference Proceedings*, SIGGRAPH '23. Association for Computing Machinery,
759 2023. 1
- 760 [7] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended
761 diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision
762 and Pattern Recognition*, pages 18208–18218, 2022. 1
- 763 [8] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and
764 Matthieu Cord. Diffedit: Diffusion-based semantic image
765 editing with mask guidance. In *The Eleventh International
766 Conference on Learning Representations*, 2023. 1
- 767 [9] Jooyoung Choi, Yunjey Choi, Yunji Kim, Junho Kim, and
768 Sungroh Yoon. Custom-edit: Text-guided image editing
769 with customized diffusion models. *arXiv preprint arXiv:2305.15779*, 2023. 1
- 770 [10] Zhixing Zhang, Ligong Han, Arnab Ghosh, Dimitris N
771 Metaxas, and Jian Ren. Sine: Single image editing with
772 text-to-image diffusion models. In *Proceedings of the IEEE/CVF
773 Conference on Computer Vision and Pattern Recognition*,
774 pages 6027–6037, 2023. 1
- 775 [11] Tim Brooks, Aleksander Holynski, and Alexei A Efros. In-
776 structpix2pix: Learning to follow image editing instructions.
777 In *Proceedings of the IEEE/CVF Conference on Computer
778 Vision and Pattern Recognition*, pages 18392–18402, 2023.
779 1
- 780 [12] Eyal Molad, Eliahu Horwitz, Dani Valevski, Alex Rav
781 Acha, Yossi Matias, Yael Pritch, Yaniv Leviathan, and Yedid
782 Hoshen. Dreamix: Video diffusion models are general video
783 editors. *arXiv preprint arXiv:2302.01329*, 2023. 1
- 784 [13] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Weixian Lei,
785 Yuchao Gu, Wynne Hsu, Ying Shan, Xiaohu Qie, and
786 Mike Zheng Shou. Tune-a-video: One-shot tuning of image
787 diffusion models for text-to-video generation. *arXiv preprint arXiv:2212.11565*, 2022. 1
- 788 [14] Wen Wang, Kangyang Xie, Zide Liu, Hao Chen, Yue Cao,
789 Xinlong Wang, and Chunhua Shen. Zero-shot video editing
790 using off-the-shelf image diffusion models. *arXiv preprint arXiv:2303.17599*, 2023. 1
- 791 [15] Shuai Yang, Yifan Zhou, Ziwei Liu, and Chen Change
792 Loy. Rerender a video: Zero-shot text-guided video-to-video
793 translation. *arXiv preprint arXiv:2306.07954*, 2023. 1
- 794 [16] Wen Wang, kangyang Xie, Zide Liu, Hao Chen, Yue Cao,
795 Xinlong Wang, and Chunhua Shen. Zero-shot video editing
796 using off-the-shelf image diffusion models. *arXiv preprint arXiv:2303.17599*, 2023. 1, 3, 4
- 797 [17] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezatofighi,
798 Fisher Yu, Dacheng Tao, and Andreas Geiger. Unifying flow,
799 stereo and depth estimation. *IEEE Transactions on Pattern
800 Analysis and Machine Intelligence*, 2023. 2, 3, 1
- 801 [18] Artiom Basulto-Lantsova, José Alfredo Padilla-Medina,
802 Francisco Javier Pérez-Pinal, and Alejandro Israel Barranco
803 Gutiérrez. Performance comparative of opencv template
804 matching method on jetson tx2 and jetson nano developer
805 kits. *2020 10th Annual Computing and Communication
806 Workshop and Conference (CCWC)*, pages 0812–0816,
807 2020. 2
- 808 [19] John P Lewis. Fast template matching. In *Vision interface*,
809 volume 95, pages 15–19. Quebec City, QC, Canada, 1995. 2
- 810 [20] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc
811 Van Gool, Markus Gross, and Alexander Sorkine-Hornung.
812 A benchmark dataset and evaluation methodology for video
813 object segmentation. In *Proceedings of the IEEE conference
814 on computer vision and pattern recognition*, pages 724–732,
815 2016. 3
- 816 [21] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras,
817 and Yejin Choi. Clipscore: A reference-free evaluation met-
818 ric for image captioning. In *Proceedings of the 2021 Confer-
819 ence on Empirical Methods in Natural Language Processing*,
820 pages 7514–7528. Association for Computational Linguis-
821 tics, nov 2021. 3, 4, 2
- 822 [22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffu-
823 sion probabilistic models. *Advances in Neural Information
824 Processing Systems*, 33:6840–6851, 2020. 1
- 825 [23] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denois-
826 ing diffusion implicit models. In *International Conference
827 on Learning Representations*, 2021. 1
- 828 [24] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu,
829 and Mark Chen. Hierarchical text-conditional image gen-
830 eration with clip latents. *arXiv preprint arXiv:2204.06125*,
831 2022. 1
- 832 [25] Chitwan Saharia, William Chan, Saurabh Saxena, Lala
833 Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour,
834 Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans,
835 et al. Photorealistic text-to-image diffusion models with deep
836 language understanding. *Advances in Neural Information
837 Processing Systems*, 35:36479–36494, 2022. 1
- 838 [26] Robin Rombach, Andreas Blattmann, Dominik Lorenz,
839 Patrick Esser, and Björn Ommer. High-resolution image
840 synthesis with latent diffusion models. In *Proceedings of
841 the IEEE/CVF Conference on Computer Vision and Pattern
842 Recognition*, pages 10684–10695, 2022. 1
- 843 [27] Bram Wallace, Akash Gokul, and Nikhil Naik. Edict: Exact
844 diffusion inversion via coupled transformations. In *Proceed-
845 ings of the IEEE/CVF Conference on Computer Vision and
846 Pattern Recognition*, pages 22532–22541, 2023. 1
- 847 [28] Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei,
848 Xintao Wang, Ying Shan, and Qifeng Chen. Fatezero: Fus-
849 ing attentions for zero-shot text-based video editing. *arXiv
850 preprint arXiv:2303.09535*, 2023. 1
- 851 [29] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip
852 Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van
853 Der Smagt, Daniel Cremers, and Thomas Brox. Flownet:
854 Learning optical flow with convolutional networks. In *Pro-
855 ceedings of the IEEE international conference on computer
856 vision*, pages 2758–2766, 2015. 1

- 864 [30] Anurag Ranjan and Michael J Black. Optical flow estimation
865 using a spatial pyramid network. In *Proceedings of the*
866 *IEEE conference on computer vision and pattern recognition*, pages
867 4161–4170, 2017. 1 918
868 [31] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz.
869 Pwc-net: Cnns for optical flow using pyramid, warping, and
870 cost volume. In *Proceedings of the IEEE conference on*
871 *computer vision and pattern recognition*, pages 8934–8943,
872 2018. 1 919
873 [32] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. Lite-
874 flownet: A lightweight convolutional neural network for
875 optical flow estimation. In *Proceedings of the IEEE conference*
876 *on computer vision and pattern recognition*, pages 8981–
877 8989, 2018. 1 920
878 [33] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field
879 transforms for optical flow. In *Computer Vision–ECCV*
880 *2020: 16th European Conference, Glasgow, UK, August 23–*
881 *28, 2020, Proceedings, Part II 16*, pages 402–419. Springer,
882 2020. 1 921
883 [34] Shuyang Sun, Zhanghui Kuang, Lu Sheng, Wanli Ouyang,
884 and Wei Zhang. Optical flow guided feature: A fast and
885 robust motion representation for video action recognition. In
886 *Proceedings of the IEEE conference on computer vision and*
887 *pattern recognition*, pages 1390–1399, 2018. 1 922
888 [35] Laura Sevilla-Lara, Yiyi Liao, Fatma Güney, Varun Jampani,
889 Andreas Geiger, and Michael J Black. On the integration
890 of optical flow and action recognition. In *Pattern Recognition:*
891 *40th German Conference, GCPR 2018, Stuttgart, Germany,*
892 *October 9–12, 2018, Proceedings 40*, pages 281–297.
893 Springer, 2019. 1 923
894 [36] Rui Xu, Xiaoxiao Li, Bolei Zhou, and Chen Change Loy.
895 Deep flow-guided video inpainting. In *Proceedings of the*
896 *IEEE/CVF Conference on Computer Vision and Pattern*
897 *Recognition*, pages 3723–3732, 2019. 1 924
898 [37] Dahun Kim, Sanghyun Woo, Joon-Young Lee, and In So
899 Kweon. Deep video inpainting. In *Proceedings of the*
900 *IEEE/CVF Conference on Computer Vision and Pattern*
901 *Recognition*, pages 5792–5801, 2019. 1 925
902 [38] Kaidong Zhang, Jingjing Fu, and Dong Liu. Flow-guided
903 transformer for video inpainting. In *Computer Vision–ECCV*
904 *2022: 17th European Conference, Tel Aviv, Israel, October*
905 *23–27, 2022, Proceedings, Part XVIII*, pages 74–90.
906 Springer, 2022. 1 926
907 [39] Zhigang Tu, Hongyan Li, Wei Xie, Yuanzhong Liu, Shifu
908 Zhang, Baoxin Li, and Junsong Yuan. Optical flow for video
909 super-resolution: a survey. *Artificial Intelligence Review*,
910 55(8):6505–6546, 2022. 1 927
911 [40] Mehdi SM Sajjadi, Raviteja Vemulapalli, and Matthew
912 Brown. Frame-recurrent video super-resolution. In *Proceed-
913 ings of the IEEE conference on computer vision and pattern*
914 *recognition*, pages 6626–6634, 2018. 1 928
915 [41] Longguang Wang, Yulan Guo, Li Liu, Zaiping Lin, Xinpu
916 Deng, and Wei An. Deep video super-resolution using hr
917 optical flow estimation. *IEEE Transactions on Image Pro-
cessing*, 29:4323–4336, 2020. 1 929
918 [42] Zhen Li, Zuo-Liang Zhu, Ling-Hao Han, Qibin Hou, Chun-
919 Le Guo, and Ming-Ming Cheng. Amt: All-pairs multi-field
920 transforms for efficient frame interpolation. In *Proceedings*
921 *of the IEEE/CVF Conference on Computer Vision and Pat-
tern Recognition*, pages 9801–9810, 2023. 1 930
922 [43] Wenbo Bao, Wei-Sheng Lai, Chao Ma, Xiaoyun Zhang,
923 Zhiyong Gao, and Ming-Hsuan Yang. Depth-aware video
924 frame interpolation. In *Proceedings of the IEEE/CVF Con-
ference on Computer Vision and Pattern Recognition*, pages
925 3703–3712, 2019. 1 931
926 [44] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
927 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
928 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning
929 transferable visual models from natural language supervi-
930 sion. In *International conference on machine learning*, pages
931 8748–8763. PMLR, 2021. 2 932
933 [45] ... 933
934 [46] ... 934
935 [47] ... 935
936 [48] ... 936
937 [49] ... 937
938 [50] ... 938
939 [51] ... 939
940 [52] ... 940
941 [53] ... 941
942 [54] ... 942
943 [55] ... 943
944 [56] ... 944
945 [57] ... 945
946 [58] ... 946
947 [59] ... 947
948 [60] ... 948
949 [61] ... 949
950 [62] ... 950
951 [63] ... 951
952 [64] ... 952
953 [65] ... 953
954 [66] ... 954
955 [67] ... 955
956 [68] ... 956
957 [69] ... 957
958 [70] ... 958
959 [71] ... 959
960 [72] ... 960
961 [73] ... 961
962 [74] ... 962
963 [75] ... 963
964 [76] ... 964
965 [77] ... 965
966 [78] ... 966
967 [79] ... 967
968 [80] ... 968
969 [81] ... 969
970 [82] ... 970
971 [83] ... 971