

Supplementary Material for “Enhancing Text-to-Video Editing with Motion Map Injection”

Seong-Hun Jeong^{*1} In-Hwan Jin^{*1} Haesoo Choo^{*1} Hyeonjun Na¹ Kyeongbo Kong²
Pukyong National University¹, Pusan National University²
tlqwkrk915@pukyong.ac.kr {bds06081, tndi28, wewelst}@naver.com, kbkong@pknu.ac.kr

In this supplementary material, we describe related works, a method for calculating optical flow, an ablation study, a method for optical flow rotation, further explanation of the metric, and a github link to the code that could not be included in the extended abstract due to the lack of pages.

1. Code Descriptions

Our code is based on PyTorch version of Video-P2P [1]. We use Video-P2P [1] 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://github.com/currycurry915/Motion_Map_Injection.

And extensive experimental results are in <https://currycurry915.github.io/MMI/>.

2. Related Work

2.1. Text-Guided Editing

The diffusion model [2, 3], which has recently been actively studied, generates data from noise through the process of adding or removing noise. Based on this diffusion model, text-guided image editing models such as DALL-E2 [4], Imagen [5], and stable diffusion [6] show the results of high-quality image editing. In particular, Prompt-to-Prompt [7] presents text-guided image editing that controls the relationship between the prompt text token and the corresponding image pixel with the attention maps, enabling unprecedented semantic editing. In addition, subsequent papers such as DreamBooth [8], EDICT [9], and Imagic [10] have been actively studied recently, showing impressive results for text-guided image editing.

Based on the significant progress of text-guided image editing, research has recently been expanded to text-guided video editing with the generative model. Dreamix [11] presents the first diffusion-based method of performing text-guided motion and application editing of videos through fine-tuning, but there are difficulties with local-

ized editing by replacing a word. Video-P2P [1] divides 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 [12] performs stable video reconstruction and editing by adding cross-frame attention to the U-Net structure of the existing diffusion model. In addition, FateZero [13], 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 [14] was the first fully convolutional neural network for estimating optical flow. Then, a series of works, represented by SpyNet [15], PWC-Net [16], LiteFlowNet [17], and RAFT [18] were proposed to reduce the computational costs through coarse-to-fine and iterative estimation methodology. Recently GMFlow [19] were proposed to achieve highly accurate results without relying on a large number of refinements by performing global matching with a Transformer.

Optical flow estimation is used in various video tasks. First, video action recognition [20, 21] 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 [22, 23, 24] enables spatio-temporally stable synthesis between frames of video through optical flow. Video super resolution [25, 26, 27] 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



Figure 1. Outputs of various methods to measure the correlation between the motion map and the attention maps.

Pouring ~~water~~ beer into glass

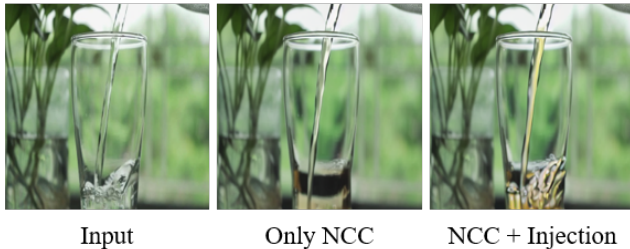


Figure 2. Comparison results of directly injecting motion maps and not injecting them.

Table 1. Evaluation on 4 template matching operations using BRISQUE and NIQE

	BRISQUE	NIQE
TM_CCOEFF	47.87	14.92
TM_CCOEFF_NORMED	27.91	11.52
TM_SQDIFF	64.45	15.80
TM_SQDIFF_NORMED	52.48	13.91

(VFI) [28, 29], 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. Ablation Study

In Fig. 1, 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 “CCOEF” applied to the first image used correlation coefficient, and the second image is the result of “CCOEF_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. 1, 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 prompt word. In Table 1, template

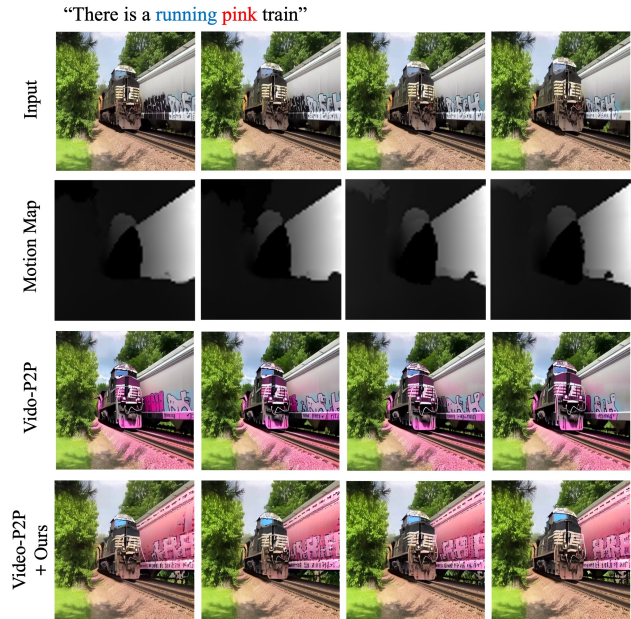


Figure 3. Editing method for objects moving in the direction specified by the user. Before editing, the user first selects one of the 8 directions.

matching operations were also measured through NIQE and BRISQUE scores. “CCOEFF_N” was demonstrated to be the most appropriate template matching operation by showing the highest quality in both metrics. Therefore, we choose “CCOEFF_N” in method which is same as NCC.

The result of our proposed framework using the fixed parameter of $\lambda = 0.3$ is shown in Fig. 2. First, editing was difficult when only the motion map was injected into the motion prompt due to the difference in scale from other attention values. Such as, in image P2P[7], the weight for the attention value for a single word is concentrated and increased, making it impossible to edit. With NCC, the scale of the attention was normalized and matched. In addition, because motion information was injected by calculating the correlation between the attention map of the entire prompt and the motion map, the entire attention was enhanced and showed better performance than existing T2V models [1, 12, 13]. Meanwhile, due to the inherent limitation of image diffusion model, the NCC score between motion map and attention map of motion prompt could not be attained. Consequently, we set the motion prompt’s NCC score to 1 to increase its level of attention, so that the editing goes desirable.

4. 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 current 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. 3.

5. Experimental Details

5.1. Evaluation Metrics

CLIP Score [30] is calculated in the CLIP model [31], generating embedding vectors for input images and prompts. CLIP Score [30] 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 [30]. We measured the CLIP Score [30] 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 [30] 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 denotes 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 [1]. It indicates 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 [1], $B(V, M) = V_M$ is defined as a reversed mask binary function, so only regions not to be changed are involved in measuring masked PSNR.

BRISQUE With a No-Reference Image Quality Assessment, the quality of the image is evaluated with only the input image without any comparison image. Using this metric, We evaluated the edited video without the original one.

5.2. Dataset and Implementation Details

Dataset Experiments were conducted with DAVIS video dataset [32] and Youtube videos that we collected. We chose 20 videos having motion, and successive 4 frames from those. They were cropped and resized to 512×512 .

Implementation details We experimented with NVIDIA RTX A6000 GPUs and set the resolution to 512×512 like Video-P2P [1], vid2vid-zero [33], and FateZero [13]. The number of video frames was set to 4 because this number is sufficient to demonstrate how well our method accomplishes our goal. The UniMatch [19] model was used to extract the optical flow.

6. Details about User Study

In this section, we describe the details of the user study. A total of 60 participants were asked to choose the more preferable video for 20 videos. Among the 20 videos, 10 were output of Video-P2P [1], and vid2vid-zero [12] and FateZero [13] were included with 5 each. The users were asked to answer following questions: ‘‘Edits well, reflecting the target prompt accurately’’, ‘‘Maintains overall structure well after editing.’’, and ‘‘Edited result is realistic and of high quality’’. The order of videos were randomized.

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. 4 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 [1]. 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 [1], is difficult to perform accurate editing.

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

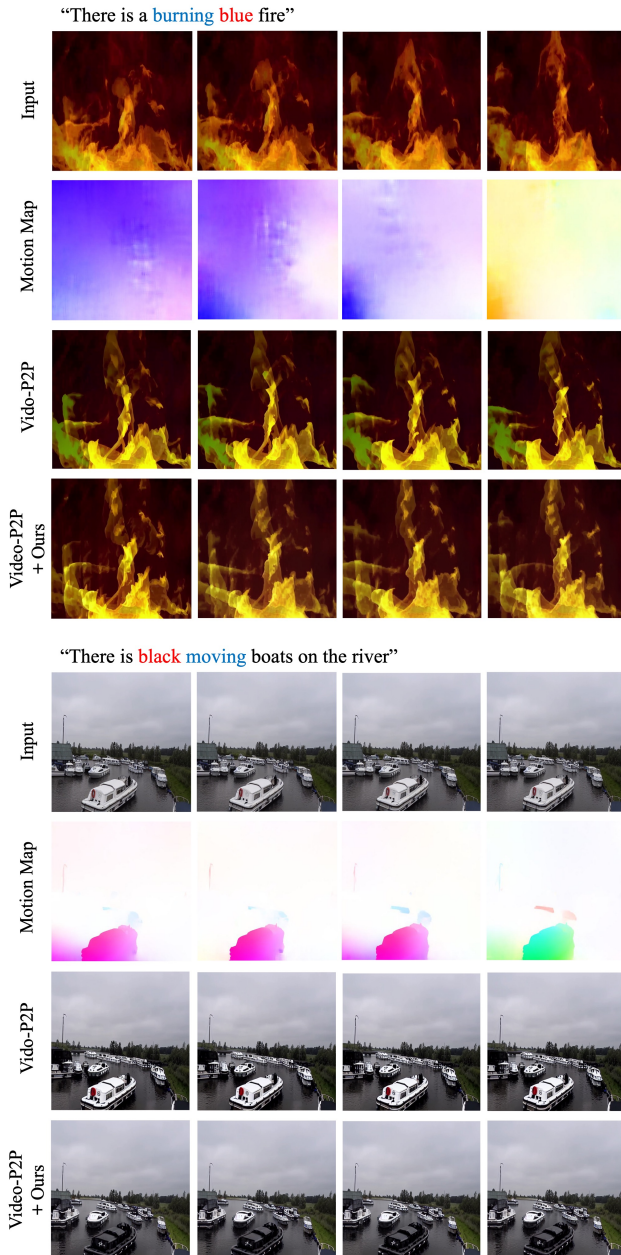


Figure 4. Results of Video-p2p [1] and our video editing model with inaccurately estimated optical flow

References

- [1] 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](#), [2](#), [3](#), [4](#)
- [2] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. [1](#)
- [3] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021. [1](#)
- [4] Aditya Ramesh, Prfulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. [1](#)
- [5] 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. [1](#)
- [6] 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. [1](#)
- [7] 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](#)
- [8] 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](#)
- [9] 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. [1](#)
- [10] 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](#)
- [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. [1](#)
- [12] 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. [1](#), [2](#), [3](#)
- [13] 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. [1](#), [2](#), [3](#)
- [14] 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. [1](#)
- [15] 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. 1
- [16] 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. 1
- [17] Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. Lite-flownet: A lightweight convolutional neural network for optical flow estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8981–8989, 2018. 1
- [18] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 402–419. Springer, 2020. 1
- [19] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezatofghi, Fisher Yu, Dacheng Tao, and Andreas Geiger. Unifying flow, stereo and depth estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023. 1, 3
- [20] Shuyang Sun, Zhanghui Kuang, Lu Sheng, Wanli Ouyang, and Wei Zhang. Optical flow guided feature: A fast and robust motion representation for video action recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1390–1399, 2018. 1
- [21] Laura Sevilla-Lara, Yiyi Liao, Fatma Güneç, Varun Jampani, Andreas Geiger, and Michael J Black. On the integration of optical flow and action recognition. In *Pattern Recognition: 40th German Conference, GCPR 2018, Stuttgart, Germany, October 9–12, 2018, Proceedings 40*, pages 281–297. Springer, 2019. 1
- [22] Rui Xu, Xiaoxiao Li, Bolei Zhou, and Chen Change Loy. Deep flow-guided video inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3723–3732, 2019. 1
- [23] Dahun Kim, Sanghyun Woo, Joon-Young Lee, and In So Kweon. Deep video inpainting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5792–5801, 2019. 1
- [24] Kaidong Zhang, Jingjing Fu, and Dong Liu. Flow-guided transformer for video inpainting. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVIII*, pages 74–90. Springer, 2022. 1
- [25] Zhigang Tu, Hongyan Li, Wei Xie, Yuanzhong Liu, Shifu Zhang, Baoxin Li, and Junsong Yuan. Optical flow for video super-resolution: a survey. *Artificial Intelligence Review*, 55(8):6505–6546, 2022. 1
- [26] Mehdi SM Sajjadi, Raviteja Vemulapalli, and Matthew Brown. Frame-recurrent video super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6626–6634, 2018. 1
- [27] Longguang Wang, Yulan Guo, Li Liu, Zaiping Lin, Xinpu Deng, and Wei An. Deep video super-resolution using hr optical flow estimation. *IEEE Transactions on Image Processing*, 29:4323–4336, 2020. 1
- [28] Zhen Li, Zuo-Liang Zhu, Ling-Hao Han, Qibin Hou, Chun-Le Guo, and Ming-Ming Cheng. Amt: All-pairs multi-field transforms for efficient frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9801–9810, 2023. 2
- [29] Wenbo Bao, Wei-Sheng Lai, Chao Ma, Xiaoyun Zhang, Zhiyong Gao, and Ming-Hsuan Yang. Depth-aware video frame interpolation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3703–3712, 2019. 2
- [30] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7514–7528. Association for Computational Linguistics, nov 2021. 3
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 3
- [32] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 724–732, 2016. 3
- [33] 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. 3