

MotionV2V: Editing Motion in a Video

Supplementary Material

775

6. Human Interaction



Figure 8. Interface for creating motion edits. The red arrow shows the transformation from source trajectory (line) to target trajectory (triangle).

776 Our motion editing interface (Figure 8) provides an
 777 intuitive way to specify complex motion changes without
 778 requiring any laborious segmentation or rotoscoping.
 779 Users simply click points on the video to initialize
 780 points which are then tracked bidirectionally to create
 781 source trajectories. Then, the user manipulates these
 782 trajectories using Bezier splines to define the desired
 783 target motion, with arrows indicating the transforma-
 784 tion from source to target.

785 The red arrow in the figure illustrates a typical edit,
 786 showing the transformation from the source trajectory
 787 (line) to the target trajectory (triangle). The interface
 788 allows users to scrub through video frames and place
 789 trajectory points as needed. Different tracking points
 790 are represented by different colors: red, green, blue,
 791 cyan, magenta, yellow, and white.

7. User Study

793 As described in Section 4, we conducted a user study
 794 with 41 participants evaluating 20 test videos to com-
 795 pare our method against state-of-the-art baselines.

796 **Evaluation Protocol** Participants used the inter-
 797 face shown in Figure 10 to compare our method against
 798 three baselines: ATI [37], ReVideo [24], and Go-with-
 799 the-Flow [3]. The interface presents side-by-side com-
 800 parisons of the original input video alongside results

from our method and all baselines, with colored track-
 ing dots visualizing the motion edits so users can clearly
 see how each method interprets the desired motion
 changes. For each test case, participants were asked
 three questions:

- **Q1:** “Which video better preserves the input video’s content?”
- **Q2:** “Which video better reflects the desired motion?”
- **Q3:** “Which video is overall a better edit of the input video?”

801
 802
 803
 804
 805
 806
 807
 808
 809
 810
 811

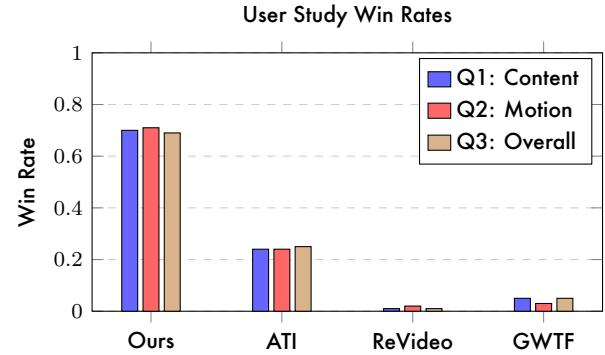


Figure 9. User study win rates per question (see Table 1 for values).

The results, visualized in Figure 9 and detailed in Table 1, show users greatly preferred our method, with rates around 70% compared to approximately 25% for ATI and less than 5% for ReVideo and Go-with-the-Flow.

812
 813
 814
 815
 816

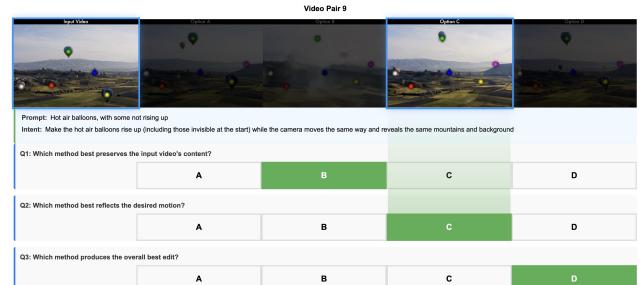
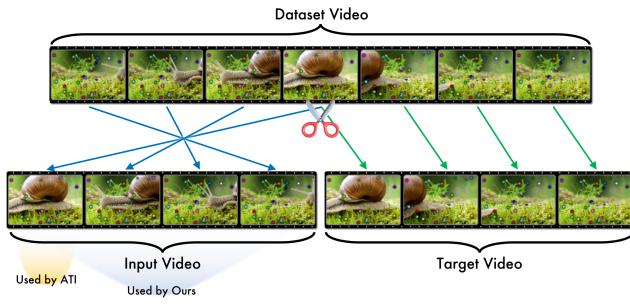


Figure 10. User study evaluation interface. Side-by-side video comparisons with visualized motion edits.

817

8. Quantitative Evaluation Dataset

818 Our quantitative evaluation dataset construction is
 819 fully described in Section 4 (Dataset Construction sub-
 820 section). We created $N_{\text{test}} = 100$ test videos by split-
 821 ting videos at their temporal midpoint and reversing
 822 one half to create video pairs with common starting
 823 frames (Figure 11). This protocol ensures that image-
 824 to-video baselines, which require first-frame correspon-
 825 dences, receive inputs they can properly handle since
 826 the tracking points match the first frame. The dataset
 827 specifically includes videos where significant content
 828 appears in middle frames but not in the first frame,
 829 quantified by tracking $N_{\text{points}} = 25$ points bidirection-
 830 ally from the midpoint.

Quantitative Evaluation: Video Preparation

831 **Figure 11. Test Data Generation.** A video is separated
 832 at the middle, and then one half is reversed. This results
 833 in two videos with a common starting frame.

834

9. Baseline Implementation Details

835 All baseline methods are image-to-video (I2V) algo-
 836 rithms with motion control, fundamentally different
 837 from our video-to-video (V2V) approach:

838 **ATI** [37] ATI is based on Wan2.1. It is a point-
 839 based image-to-video algorithm with controllable motion.
 840 For our baseline, we take the target tracks and
 841 apply it to the first frame of our counterfactual videos,
 842 using the target prompts for text guidance. We use the
 843 default number of diffusion steps and CFG as provided
 844 by their public code repository.

845 **ReVideo** [24] ReVideo is based on Stable Video
 846 Diffusion. It takes no prompt as an input. It is
 847 an image-to-video point-based motion-controllable al-
 848 gorithm, with the ability to specify editable regions.
 849 Since we are avoiding manual labor such as rotoscoping
 850 we designate the entire video as an editable region.

Go-with-the-Flow [3] Go-with-the-Flow is not a
 849 point-based motion control algorithm, but is instead an
 850 image-to-video motion-controllable algorithm driven
 851 by warped noise, which often comes from optical flow.
 852 To make this baseline work, we run the rasterized tar-
 853 get tracks through RAFT to get optical flows, and from
 854 that create warped noise that is used to generate the
 855 output videos. We use a more recent Wan2.2-based ver-
 856 sion of Go-with-the-Flow to test with, as it is a tougher
 857 baseline than their original CogVideoX-5B model.

The fundamental limitation of all these baselines is
 858 their I2V formulation: they can only access informa-
 859 tion from the first frame, preventing them from han-
 860 dling content that appears later in the video, camera
 861 viewpoint changes, or complex temporal reordering.
 862 Our V2V approach, in contrast, can leverage informa-
 863 tion from any frame of the input video.

864

10. Ablations



865 **Figure 12. The effects of trajectory jitter on motion**
 866 **editing.** Top: without jitter, a second basketball appears.
 867 Bottom: with 1-2 pixel jitter, the edit follows correctly.

868 We discovered an interesting phenomenon during
 869 inference: when tracking points are pixel-perfectly
 870 aligned with the input video trajectories across mul-
 871 tiple frames, the model exhibits a strong bias toward
 872 reproducing the original video’s semantics rather than
 873 following the edited motion.

874 Figure 12 illustrates this effect. In the top
 875 row (without jitter), when tracking points are pixel-
 876 perfectly aligned with the input video, the model ex-
 877 hibits a bias toward reproducing the original video’s
 878 semantics. Although the basketball successfully goes
 879 through the hoop following the edited trajectory, a sec-
 880 ond basketball mysteriously appears behind the hoop
 881 to match the original video where the basketball passes
 882 in front. This occurs because the pixel-perfect align-
 883 ment of other tracking points signals to the model
 884 that it should preserve the content of the entire input
 885 video, which is often the only case where the points

884 are aligned that perfectly during training. In the bot-
885 tom row (with jitter), this identity-copying behavior
886 is eliminated and the edit follows the intended motion
887 correctly.

888 To address this, we introduce a simple but effective
889 inference-time technique: “Jitter”. We add small ran-
890 dom noise $\epsilon \sim \mathcal{U}(-2, 2)$ pixels to the (x, y) positions of
891 all tracking points at each frame. Importantly, this is
892 an inference-time modification only—the model is not
893 trained with this jitter. This minimal perturbation (1-
894 2 pixels) is imperceptible in the rendered tracks but
895 sufficient to break the model’s tendency to copy the
896 input video’s identity, allowing it to follow the edited
897 trajectories more faithfully.

898 11. Future Work

899 In future work, we consider creating large-scale syn-
900 thetic datasets with precise motion counterfactuals
901 made with 3d software. While our current approach
902 leverages real videos paired with diffusion-generated
903 conuterfactuals, synthetic 3d data would provide per-
904 fect ground truth motion-edit pairs, enabling exact
905 control over individual object trajectories, physical
906 interactions, and the resulting lighting and shading
907 changes. This would improve the precision of our train-
908 ing dataset, possibly allowing even less points to be
909 used for control.