



EUROPEAN CONFERENCE ON COMPUTER VISION

MILANO
2024

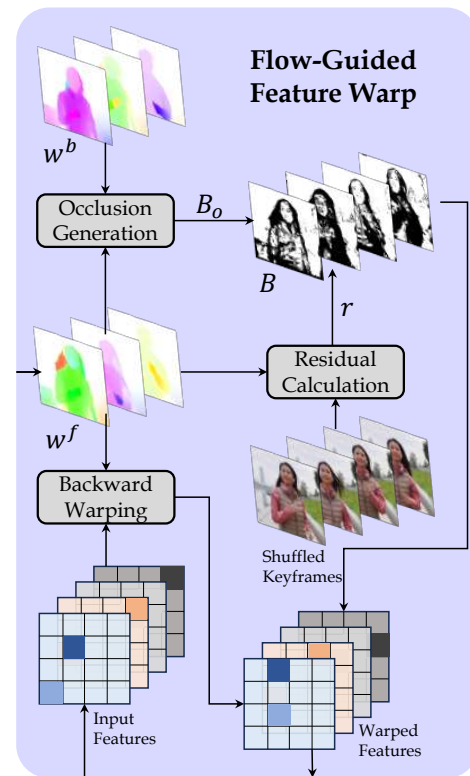
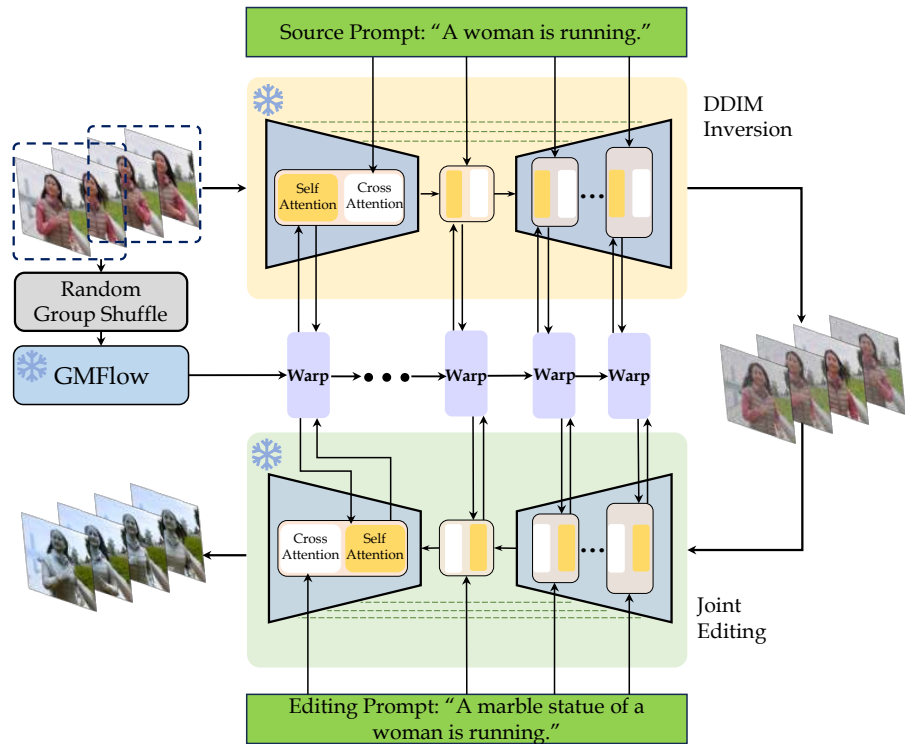
WAVE: Warping DDIM Inversion Features for Zero-shot Text-to-Video Editing

Yutang Feng*, Sicheng Gao*, Yuxiang Bao*, Xiaodi Wang, Shumin Han[†],
Juan Zhang[†], Baochang Zhang, and Angela Yao

Given a source video and a target textual prompt, WAVE employs a flow-guided feature warping strategy during both DDIM inversion and joint editing processes for powerful video editing.



Our framework



$$z_{t-1} = \sqrt{\alpha_{t-1}} \frac{z_t - \sqrt{1 - \alpha_t} \varepsilon_\theta}{\sqrt{\alpha_t}} + \sqrt{1 - \alpha_{t-1}} \varepsilon_\theta, \quad (4)$$

$$\hat{z}_t = \sqrt{\alpha_t} \frac{\hat{z}_{t-1} - \sqrt{1 - \alpha_{t-1}} \varepsilon_\theta}{\sqrt{\alpha_{t-1}}} + \sqrt{1 - \alpha_t} \varepsilon_\theta, \quad (5)$$

$$\hat{z}_{n+1}^i = B * w^b(\hat{z}_n^i) + (1 - B) z_{n+1}^i, \quad (8)$$

Visualization

Input video



A marble statue of a woman is running



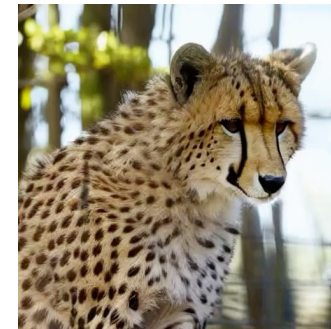
Bronze statue of a woman



Input video



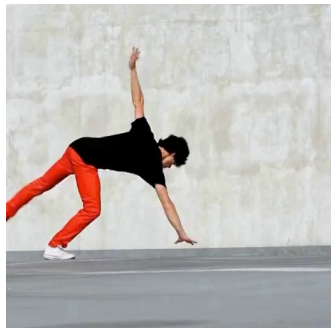
A cheetah



A husky



Input video



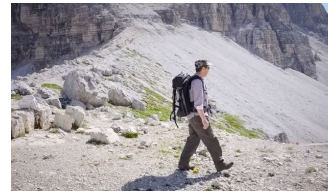
A man wearing yellow trousers



Makoto Shinkai Style



Input video



Makoto Shinkai Style



Sunset ambiance





EUROPEAN CONFERENCE ON COMPUTER VISION

MILANO
2024

Thank you