

Generative Video Diffusion for Unseen Novel Semantic Video Moment Retrieval

Dezhao Luo¹, Shaogang Gong¹, Jiabo Huang², Hailin Jin³, Yang Liu^{4,5}

¹Queen Mary University of London, ²Sony AI, ³Adobe Research,

⁴WICT, Peking University, ⁵State Key Laboratory of General Artificial Intelligence, Peking University,

Problem Definition

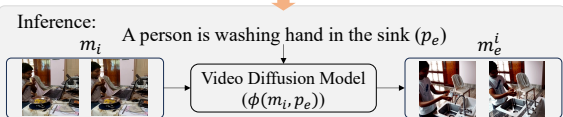
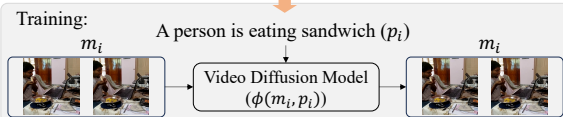
The correlation between video moments and text is crucial for **video moment retrieval**, yet there is a scarcity of large-scale datasets.

Solution

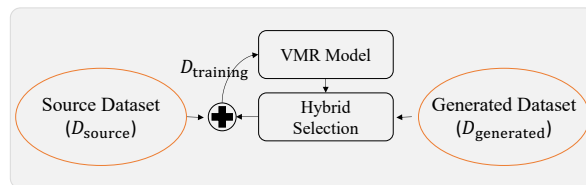
- A video diffusion model that synthesises training data
- A data selection module that selects beneficial data for the VMR task

Data Generation Pipeline

Data Generation

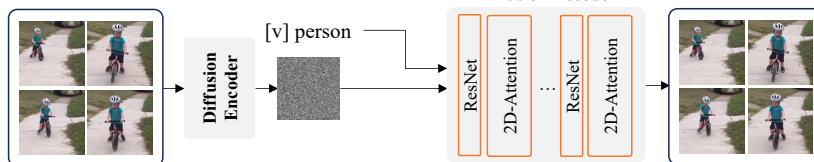


Hybrid Selection:

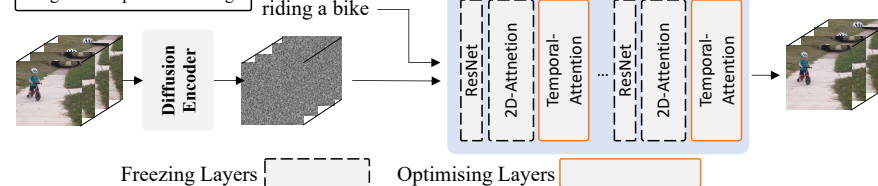


Train:

Stage 1: Instance Descriptor Learning

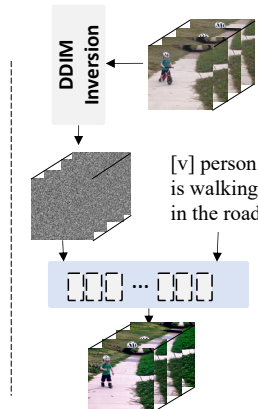


Stage 2: Temporal Encoding



Model Framework

Inference:



Video Editing Ability



Conclusion

- FVE is able to generate training data that benefits the VMR task
- FVE is able to change the action in a video and maintain other details.