

# Generative Video Diffusion for Unseen Novel Semantic Video Moment Retrieval

Dezhao Luo<sup>1</sup>, Shaogang Gong<sup>1</sup>, Jiabo Huang<sup>2</sup>, Hailin Jin<sup>3</sup>, Yang Liu<sup>4,5</sup>

<sup>1</sup>Queen Mary University of London, <sup>2</sup>Sony AI, <sup>3</sup>Adobe Research,

<sup>4</sup>WICT, Peking University, <sup>5</sup>State Key Laboratory of General Artificial Intelligence, Peking University,

## 1. Problem Definition

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

## 2. Solution

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

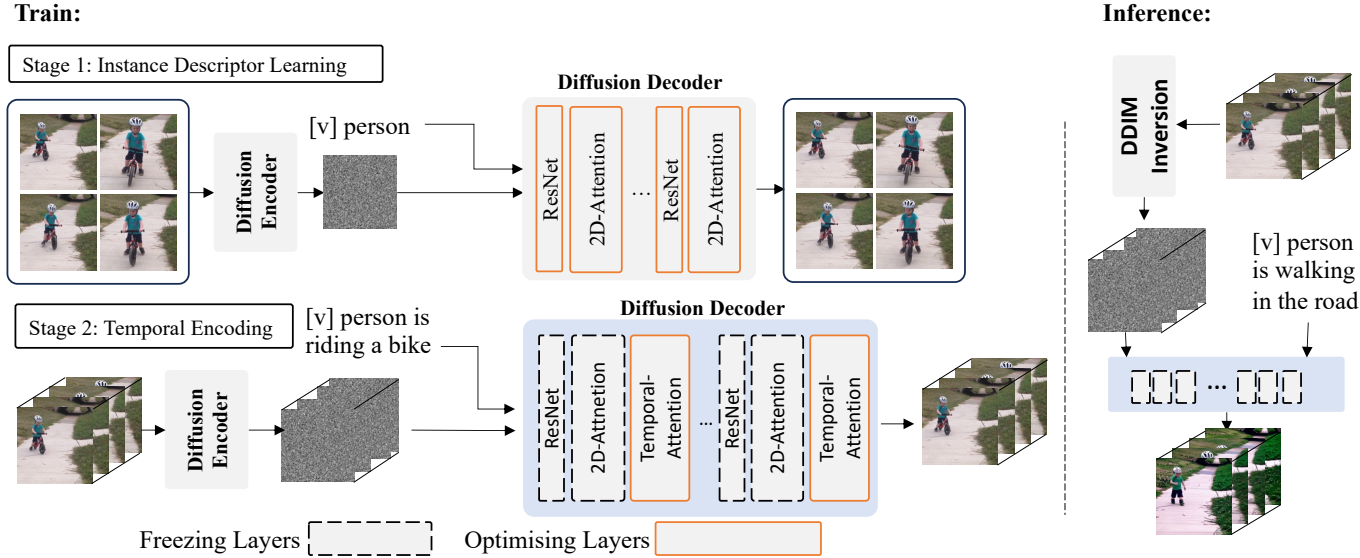
## 5. Data Selection

Cross-modal relevance:  $s_c(p_e, m_e) = \frac{1}{N} \sum_{i=1}^N \cos(\text{VLM}(p_e), \text{VLM}(f_{m_e}^i))$

Uni-modal structure:  $s_u(m_s, m_e) = \frac{1}{N} \sum_{i=1}^N \cos(\text{VM}(f_{m_s}^i), \text{VM}(f_{m_e}^i))$

Model performance:  $D_{\text{mpd}} = \text{TOP}_l(\{(d, -\text{VMR}(d)) \mid d \in D_{cu}\})$

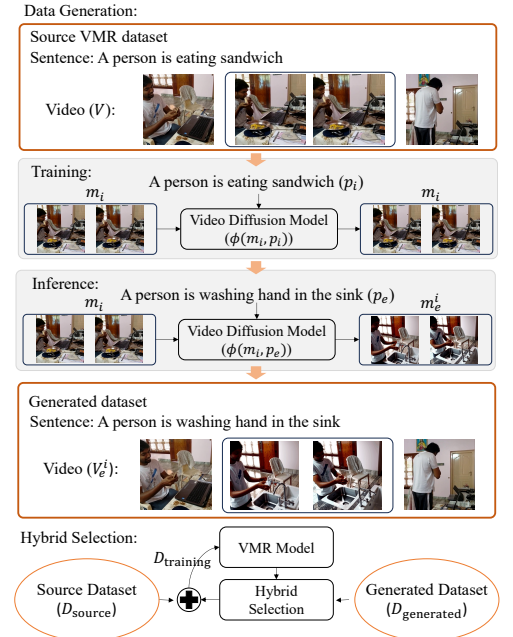
## 3. Video Diffusion Model



## 6. Video Editing Ability



## 4. Data Generation



## 7. Conclusion

- FVE is able to generate high-quality training data that benefits the VMR task (44.89%vs 44.01%)
- FVE is able to change the action in a video and maintain other details.