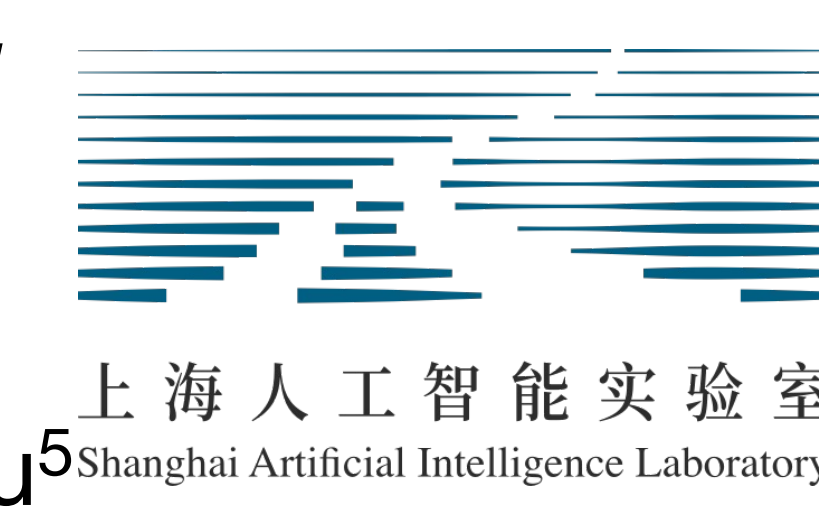


# DCM: Dual-Expert Consistency Model for Efficient and High-Quality Video Generation

Zhengyao Lv<sup>2,3\*</sup> Chenyang Si<sup>1\*</sup> Tianlin Pan<sup>1,4</sup> Zhaoxi Chen<sup>5</sup> Kwan-Yee K. Wong<sup>2</sup> Yu Qiao<sup>3</sup> Ziwei Liu<sup>5</sup>  
 NJU<sup>1</sup> HKU<sup>2</sup> Shanghai Artificial Intelligence Laboratory<sup>3</sup> UCAS<sup>4</sup> NTU<sup>5</sup>

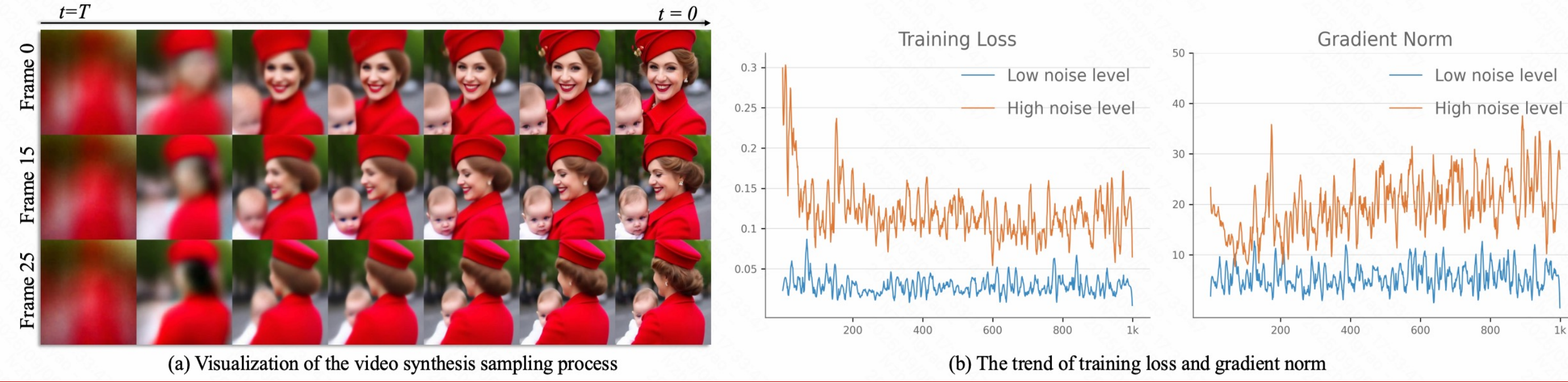


## Motivation

□ **Task:** In this paper, we propose a **parameter-efficient Dual-Expert Consistency Model (DCM)**, maintain visual quality with significantly reduced sampling steps, **down to as few as four steps**.

□ **Motivation:**

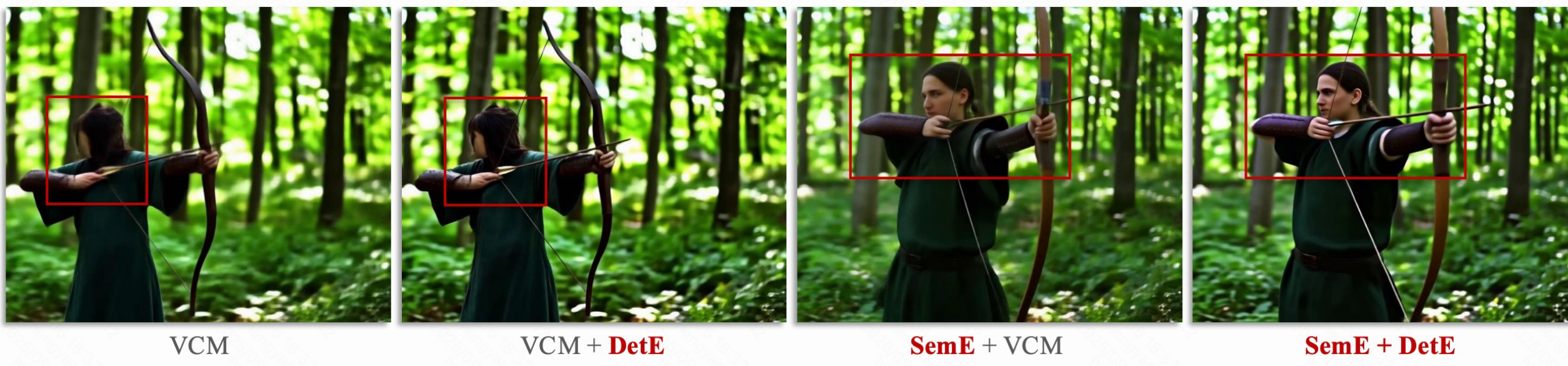
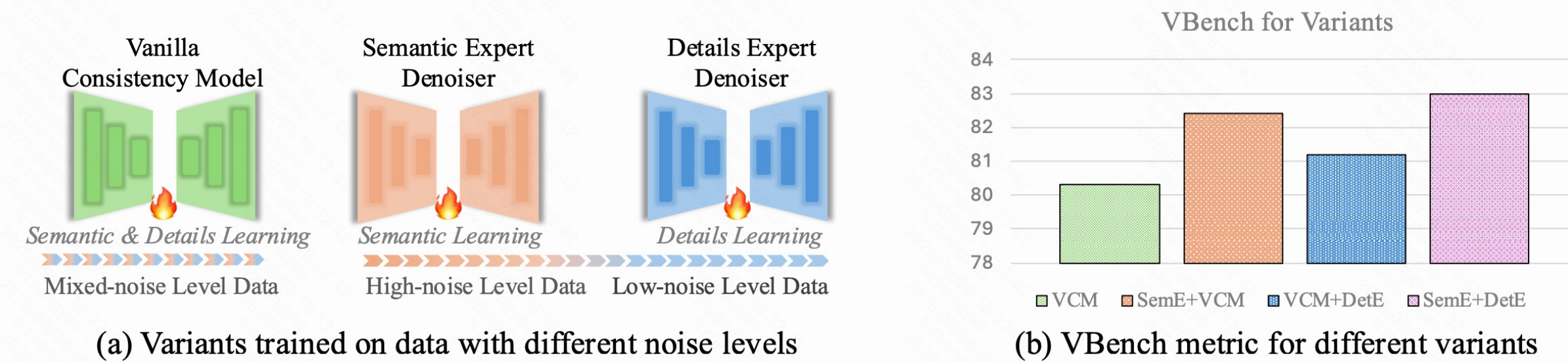
- Differences between adjacent steps are **substantial in early sampling stages**, while the changes become **more gradual in later stages**.
- The **magnitude of loss and gradient during distillation differ markedly between different noise levels**, suggesting that distilling a single student to capture both semantic layout and fine-detail synthesis may cause optimization interference and yield suboptimal results.



## Investigation

□ **Suboptimal Solution in Consistency Distillation**

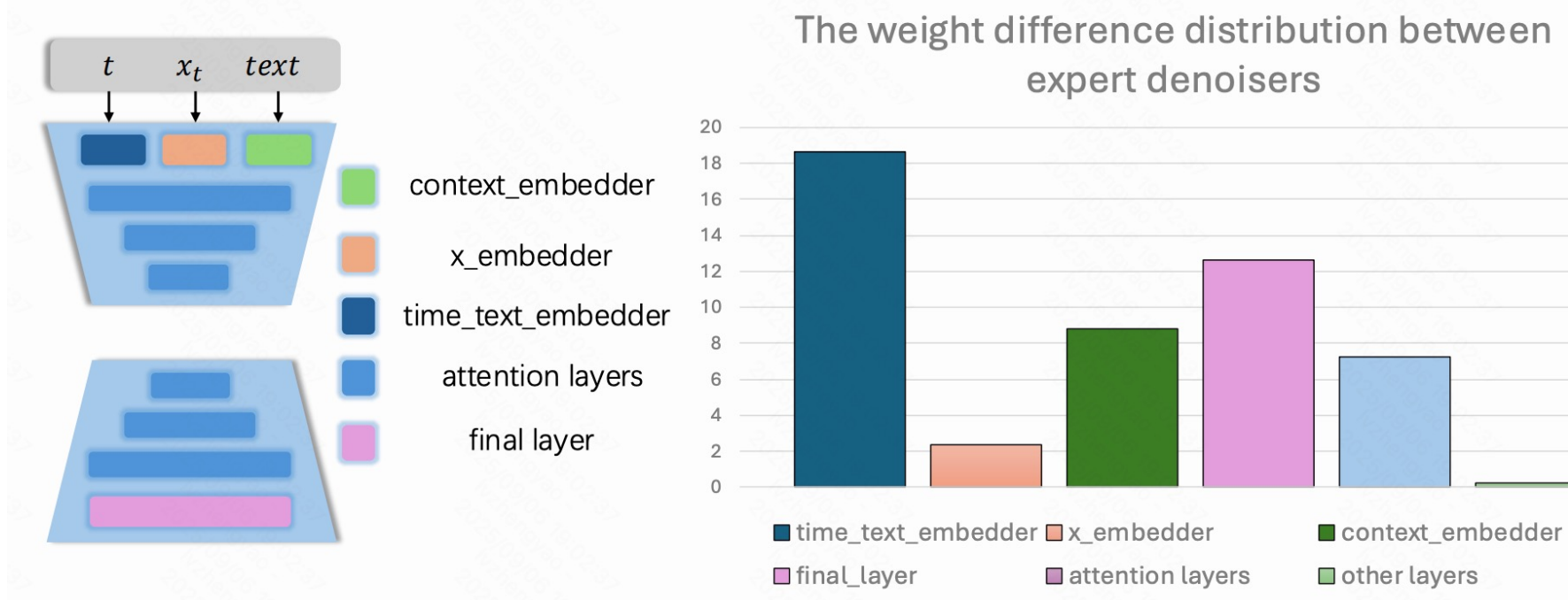
- Comparison of results of denoiser variants trained at different noise level samples.
- By **optimizing two expert denoisers to decouple the distillation process into semantic learning and detail learning**, and combining them during inference, we achieve the best quantitative and qualitative visual results.



## Method

□ **Parameter-efficient Dual-Expert Distillation**

- While training two expert denoisers improves video quality, it significantly increases model parameters and GPU memory consumption during inference.
- We found that the primary differences in model parameters lie in the embedding layers and the linear layers within the attention layers.



- Based on the above observations, we propose the parameter-efficient Dual-Expert distillation strategy.

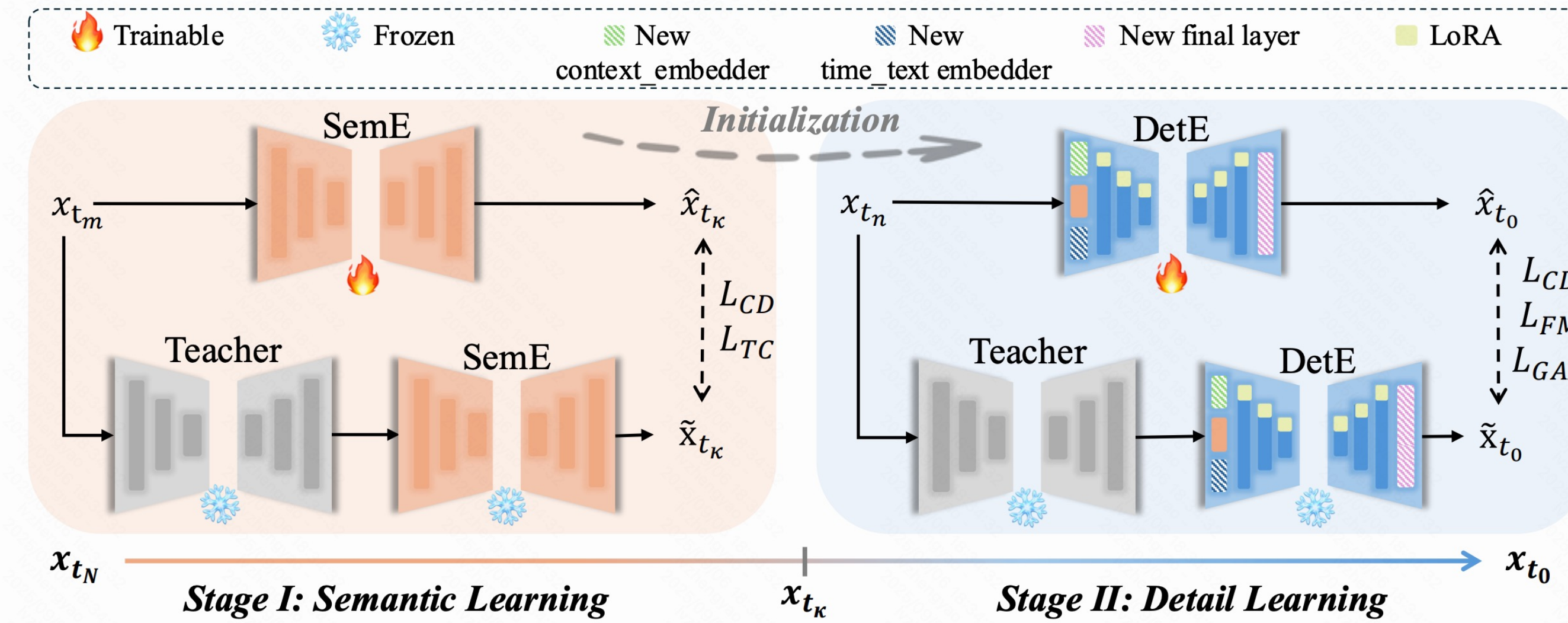


Figure 4. The training process of DCM consists of two stages. In the semantic learning stage, we train SemE on high-noise samples with consistency loss and temporal coherence loss as the learning objectives. In the detail learning stage, we initialize DetE with the weights of SemE and introduce a set of time-dependent layers and LoRA. DetE is then trained on low-noise samples, where **only the newly added**

□ **Expert-specific Optimization Objective**

- Temporal Coherence Loss for SemE

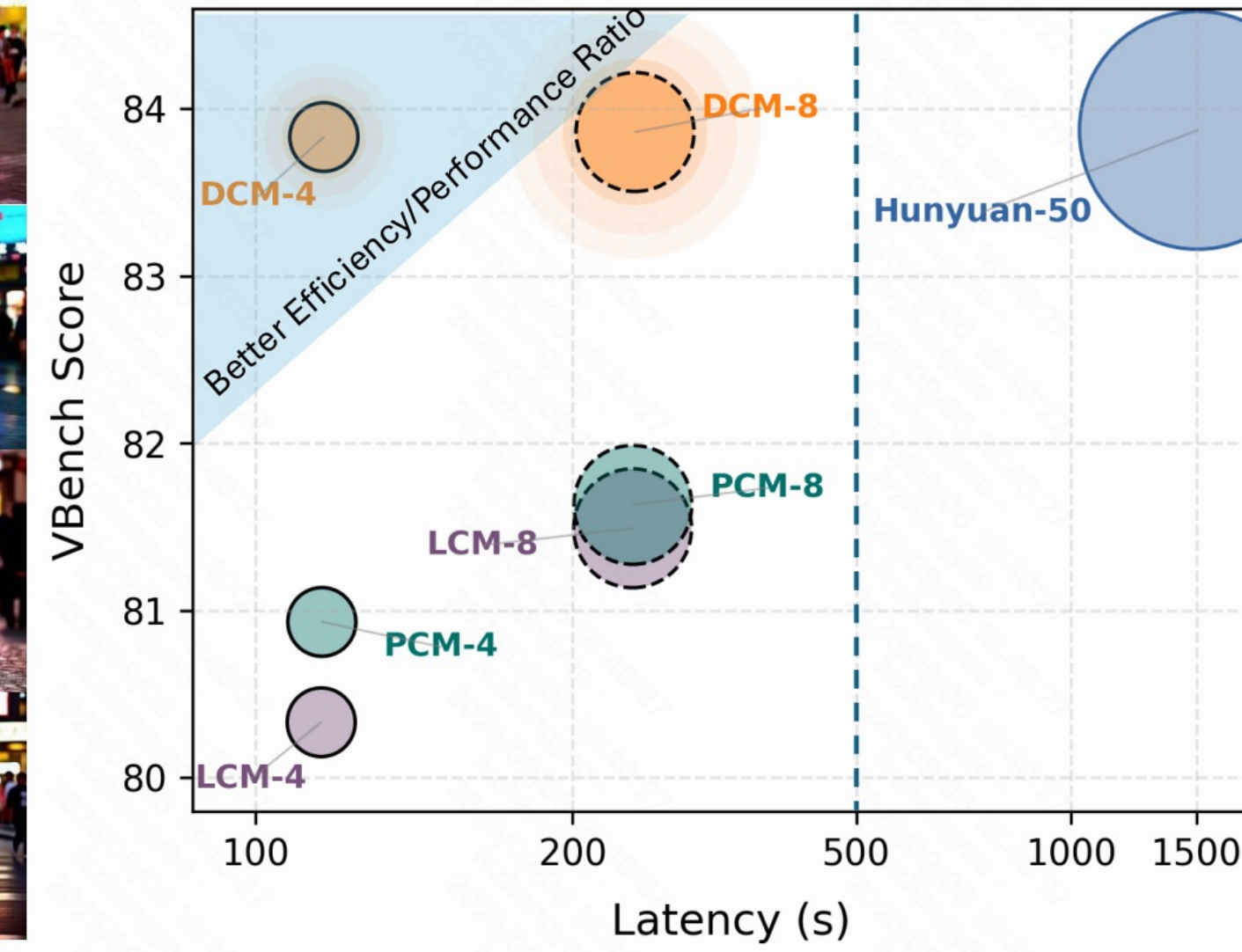
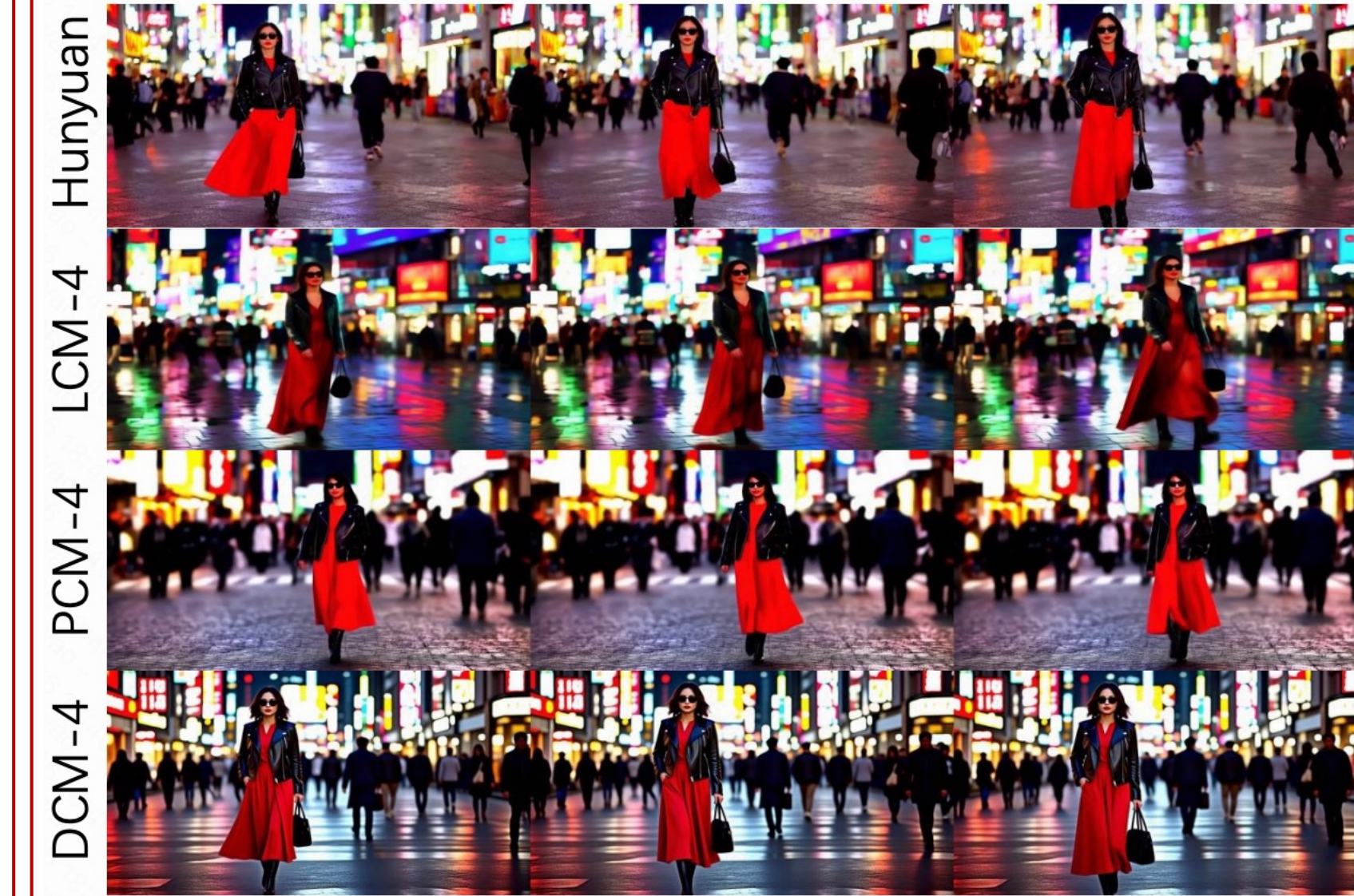
$$\begin{aligned} \mathbf{x}_{t_\kappa} &= \Phi(\mathbf{x}_{t_m}, F_{\text{SemE}}(\mathbf{x}_{t_m}, t_m, c), t_\kappa), \\ \hat{\mathbf{x}}_{t_\kappa} &= \Phi(\hat{\mathbf{x}}_{t_{m-1}}, F_{\text{SemE}}^-(\hat{\mathbf{x}}_{t_{m-1}}, t_{m-1}, c), t_\kappa), \\ \mathcal{L}_{TC} &= \|(\mathbf{x}_{t_\kappa}^{t_\kappa} - \mathbf{x}_{0:L-l}^{t_\kappa}) - (\hat{\mathbf{x}}_{t_\kappa}^{t_\kappa} - \hat{\mathbf{x}}_{0:L-l}^{t_\kappa})\|_2^2. \end{aligned}$$

- Generative Adversarial and Feature Matching Loss for DetE

$$\begin{aligned} \mathcal{L}_{FM} &= \mathbb{E}_{\mathbf{x}, t_n} \|\Omega(\mathbf{x}_{fake}) - \Omega(\mathbf{x}_{real})\|_2^2, \\ \mathcal{L}_G &= \mathbb{E}_{\mathbf{x}, t_n} [1 - f_D(\Omega(\mathbf{x}_{fake}))] + \mathcal{L}_{FM}, \\ \mathcal{L}_D &= \mathbb{E}_{\mathbf{x}, t_n} [f_D(\Omega(\mathbf{x}_{fake}))] + \mathbb{E}_{\mathbf{x}, t_n} [1 - f_D(\Omega(\mathbf{x}_{real}))]. \end{aligned}$$

## Results

□ **Quantitative Comparison & Qualitative Comparison**



□ **Ablation Study**



Figure 8. Impact of temporal coherence loss.

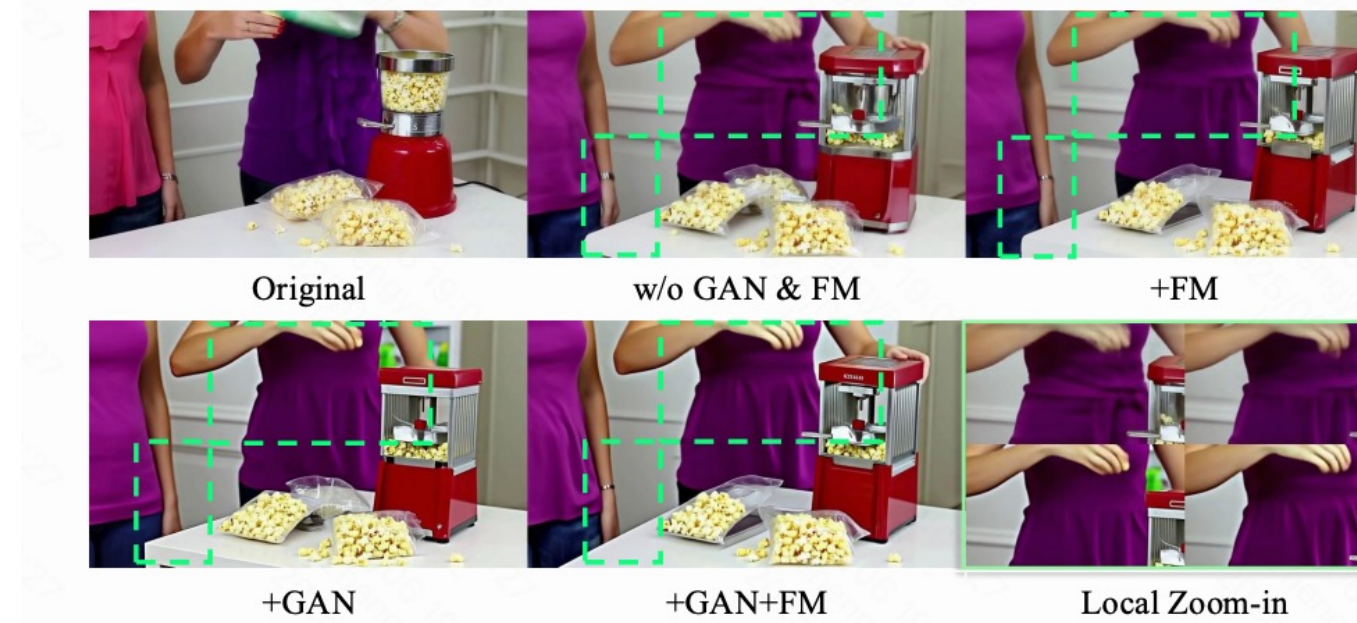


Figure 9. Impact of the GAN loss and Feature Matching term.

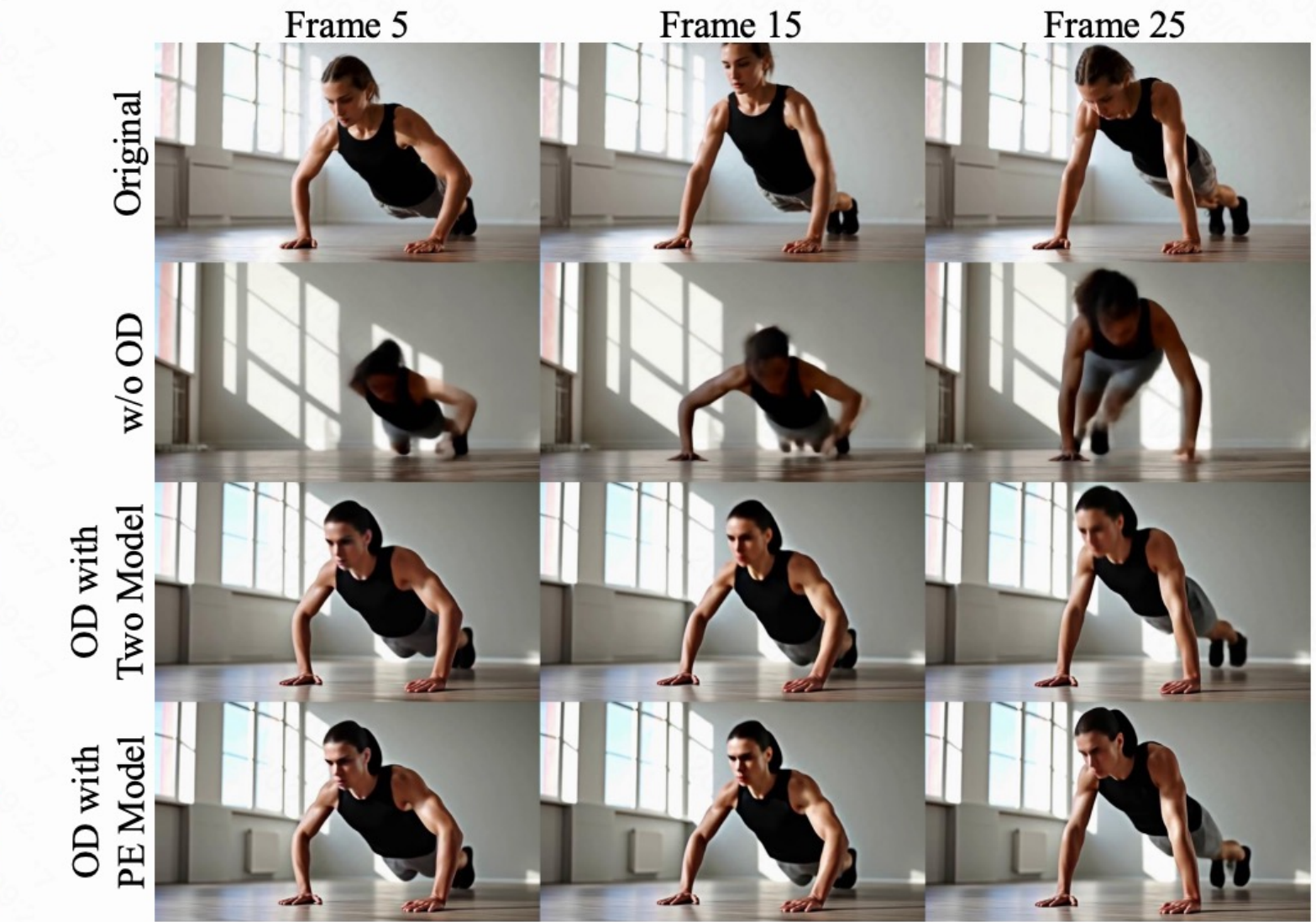


Figure 7. Impact of optimization decoupling and parameter-efficient distillation.

□ **For More Detail**



Project Page



GitHub Repo



Paper