



Zhicheng Zhang<sup>1,2,†</sup> Junyao Hu<sup>1,2,†</sup> Wentao Cheng<sup>1,‡</sup> Danda Paudel<sup>3,4</sup> Jufeng Yang<sup>1,2</sup>

<sup>1</sup>VCIP & TMCC & DISec, College of Computer Science, Nankai University

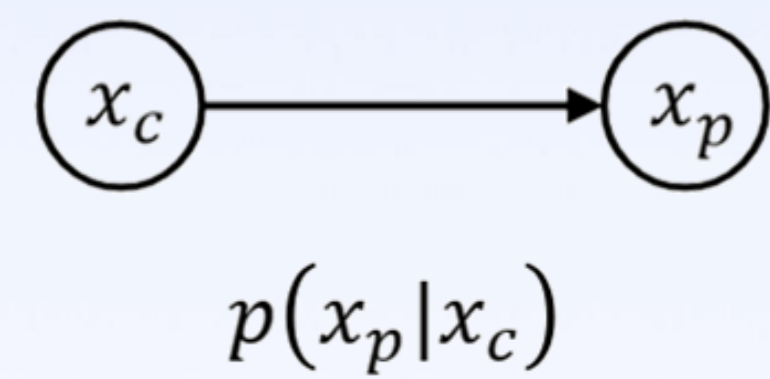
<sup>2</sup>Nankai International Advanced Research Institute (SHENZHEN·FUTIAN)

<sup>3</sup>Computer Vision Lab, ETH Zurich <sup>4</sup>INSAIT, Sofia University



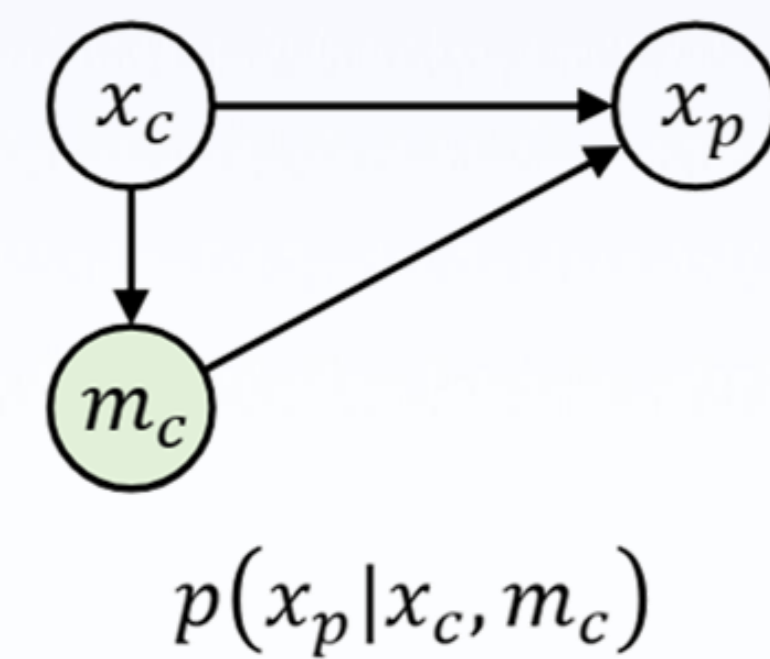
# Introduction

Video prediction aims to capture the dynamic change from present  $x_c$  to future  $x_p$ .



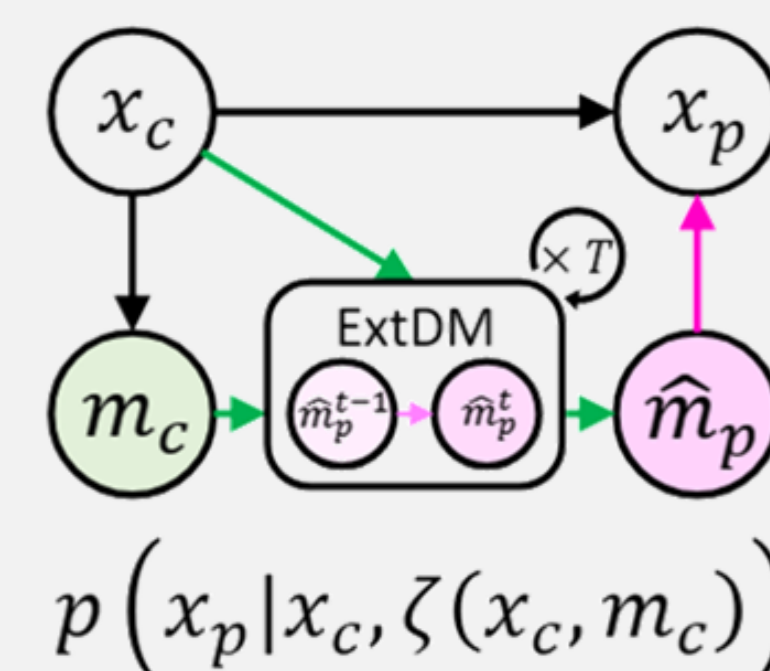
## Direct Method

- only RGB
  - difficult to solve complexity in probability estimation
- SRVP(ICML20) SimVP(CVPR22)



## In-context Learning Method

- RGB + motion as implicit cues
  - lack accuracy for longer time
  - counterfactual results like fading, deformation, etc.
- MCNet(ICLR17) MOSO(CVPR23)



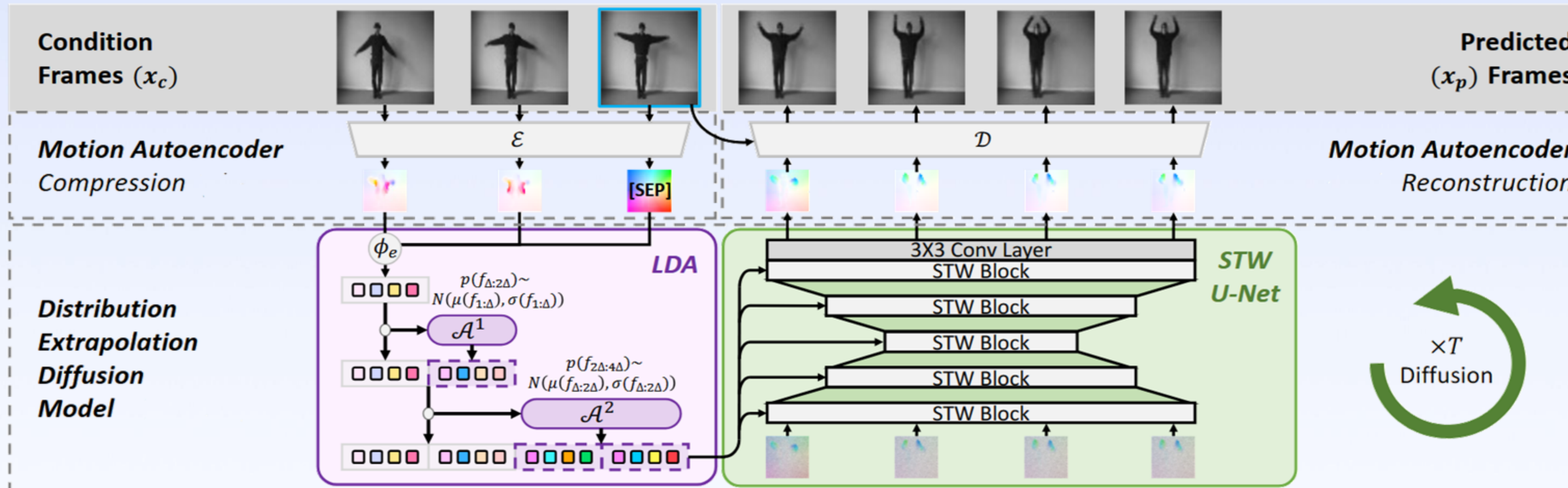
## Extrapolation Method (Ours)

- RGB + motion as explicit cues
- Extrapolate present deterministic motion cues into the future ones

## Contributions

- A distribution extrapolation diffusion model that forecasts future frames by extrapolating from the present ones.
- An efficient video prediction method includes compression and reconstruction, which can create multiple tailored proposals for stochastic events by imitating motion cues.
- Effectiveness for short/long-term videos in KTH, BAIR, Cityscapes, UCF, and SMMNIST.

# Methodology



## Two Mapping Functions

Step 1:  $x_c \rightarrow m_c$  &  $\hat{m}_p \rightarrow x_p$

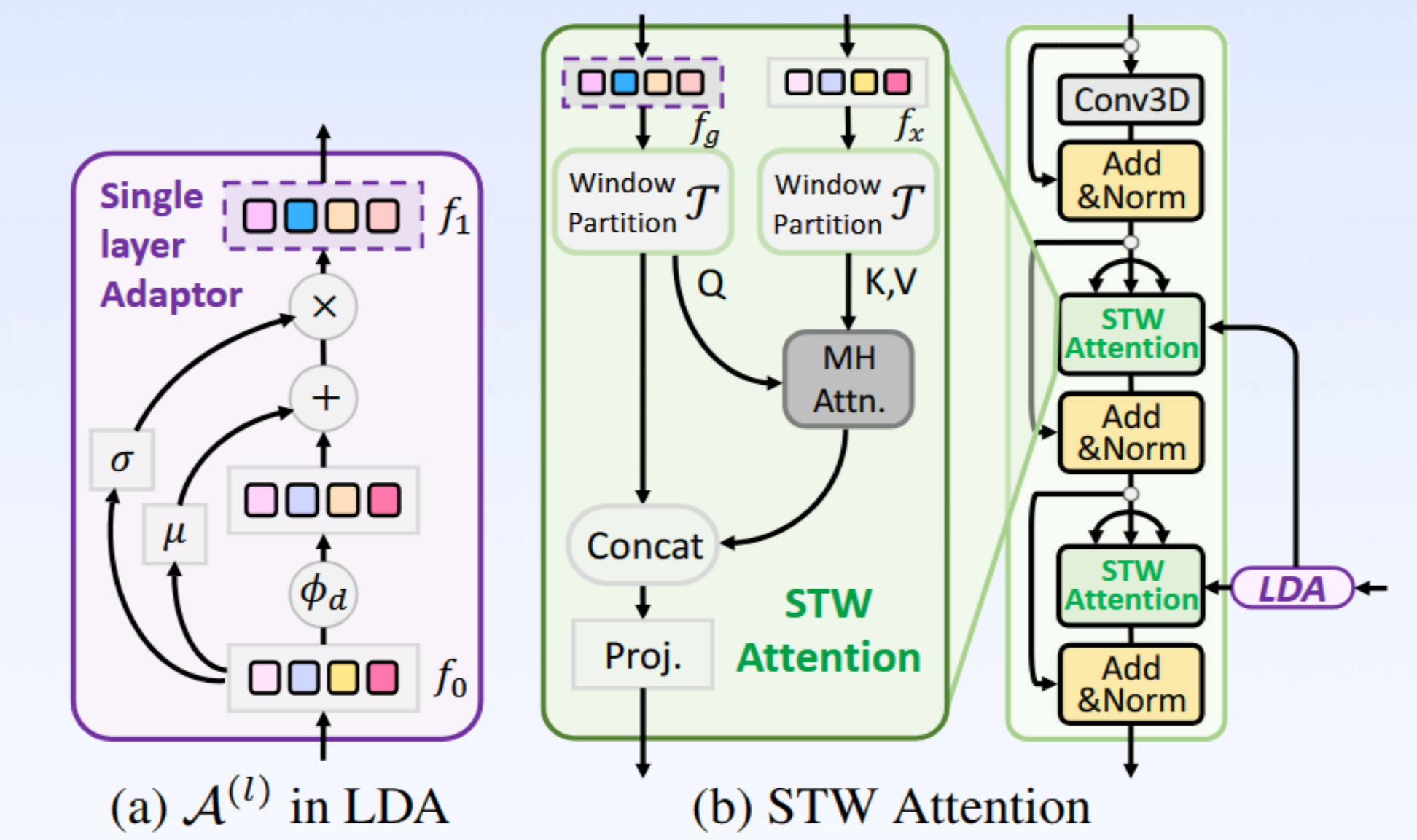
$$m_c = \{m_i \in \mathbb{R}^{3hw} \mid m_i = \mathcal{E}(x_i, x_u) = \begin{bmatrix} w_i \\ o_i \end{bmatrix}\}$$

$$x_p = \{x_j \in \mathbb{R}^{3HW} \mid x_j = \mathcal{D}(\hat{m}_j, x_u) = \mathcal{G}(o_j \odot \mathcal{W}(z_u, w_j))\}$$

Step 2:  $m_c, x_c \rightarrow \hat{m}_p$

$$\mathcal{L}_{ext} = \mathbb{E}_{t \sim \mathcal{U}(0,1), \hat{m}_p^0 \sim q(\hat{m}_p^0), \epsilon \sim \mathcal{N}(0,1)} [\|\epsilon - \epsilon_\theta(\hat{m}_p^t, t, c)\|^2]$$

• Distribution Extrapolation Diffusion Model



## (a) Layered Distribution Adaptor

- estimate distribution params
- inference using distribution sampling

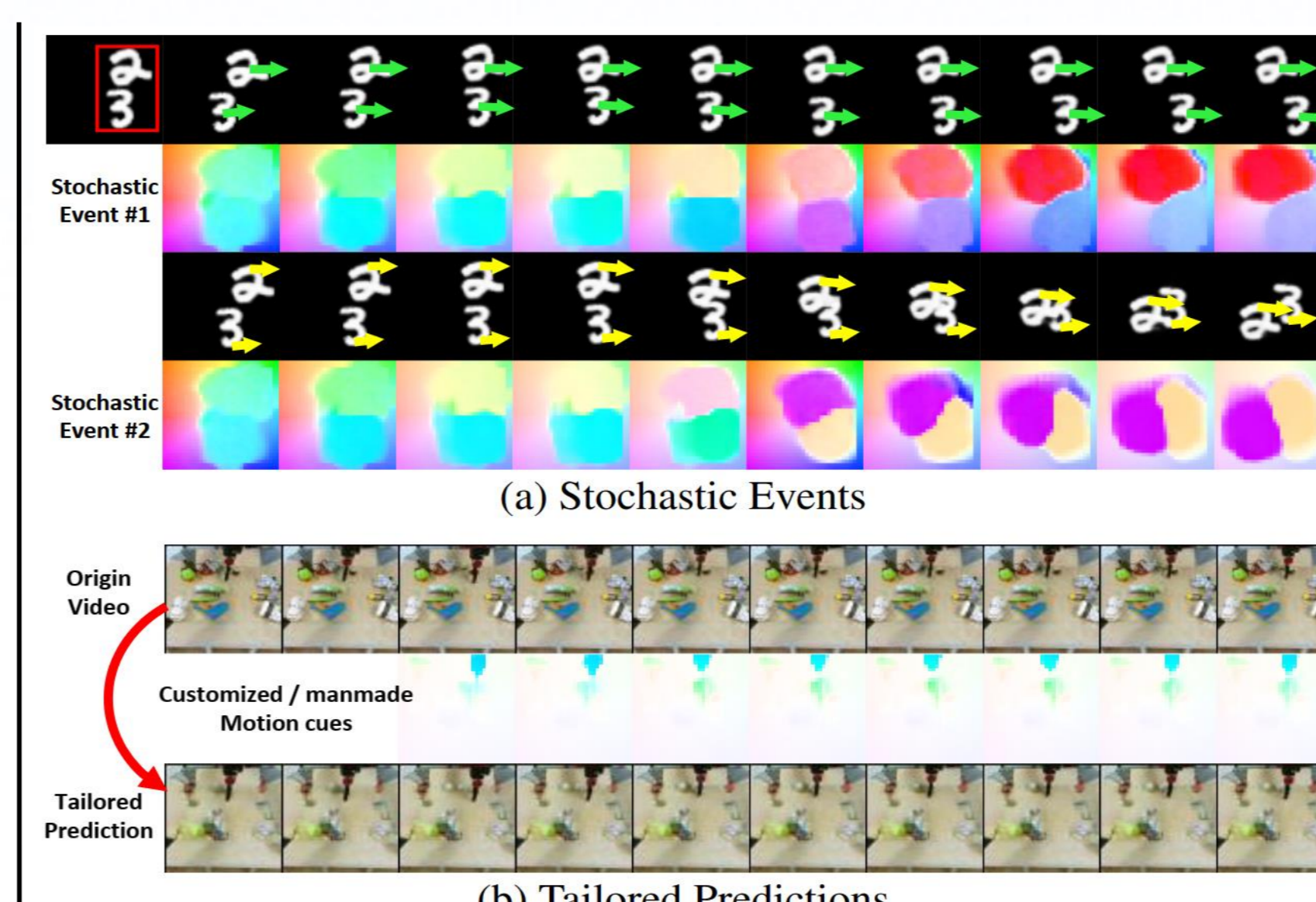
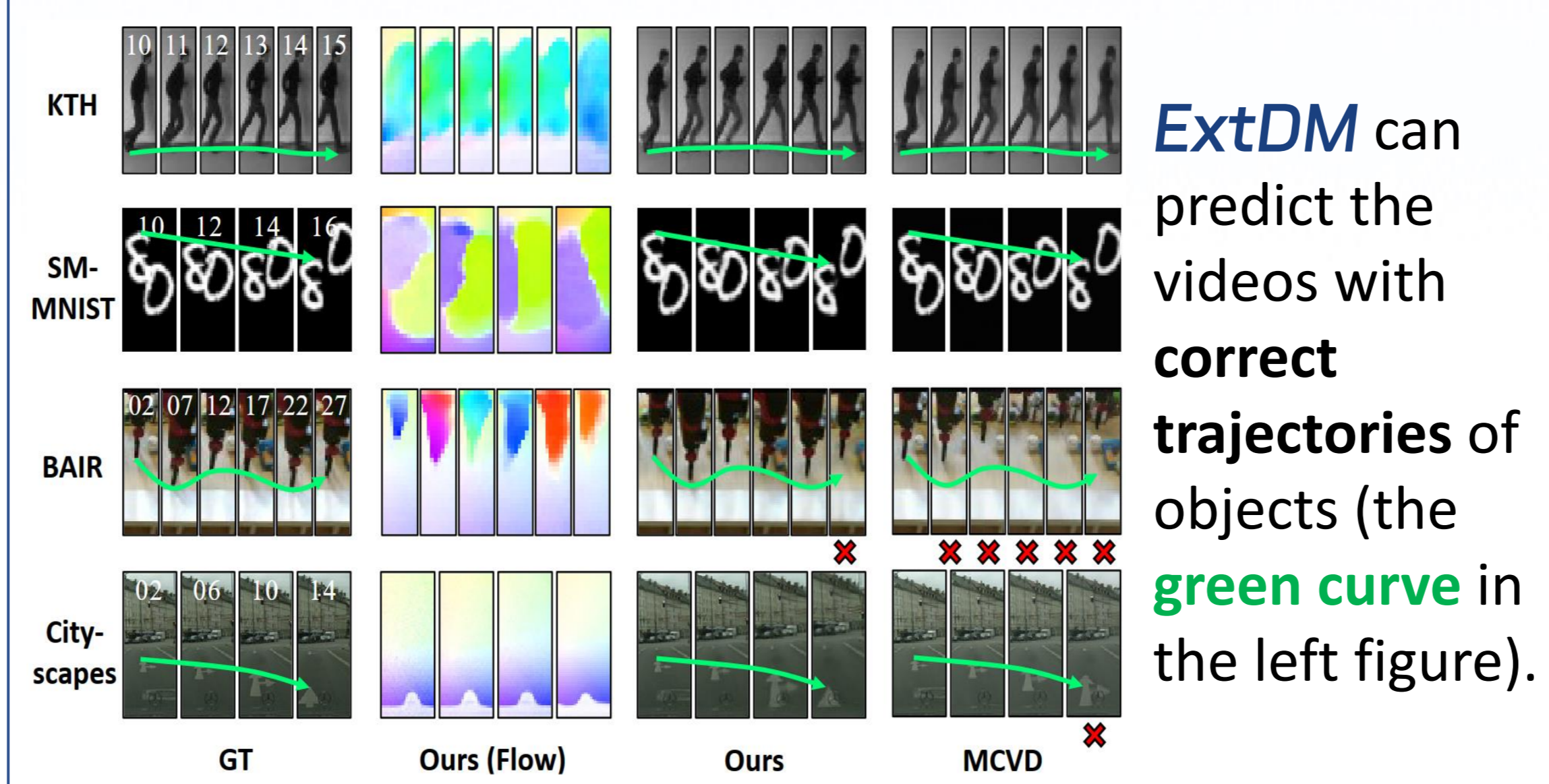
$$f_{1:\Delta} = \phi_e(f_c), \quad f_b = \mathcal{A}(f_a)$$

$$\hat{f}_{1:2^l\Delta} = (f_{1:2^{l-1}\Delta}, \mathcal{A}^{(l)}(f_{1:2^{l-1}\Delta})), \quad = (\sigma(f_a) + \sigma')\phi_d \left( \frac{f_a - \mu(f_a)}{\sigma(f_a)} \right)$$

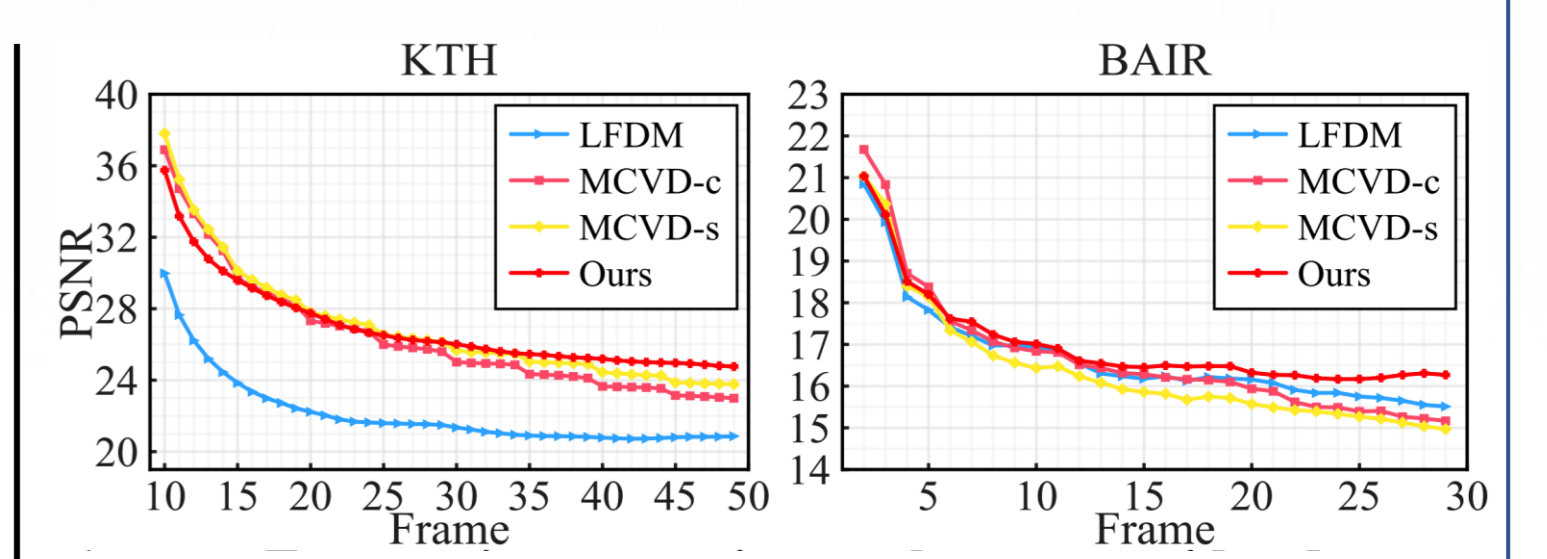
$$f_p = (\hat{f}_{1:\Delta}, \dots, \hat{f}_{2^{L-1}\Delta:2^L\Delta}), \quad + \mu(f_a) + \mu'$$

## (b) Spatiotemporal Window U-Net

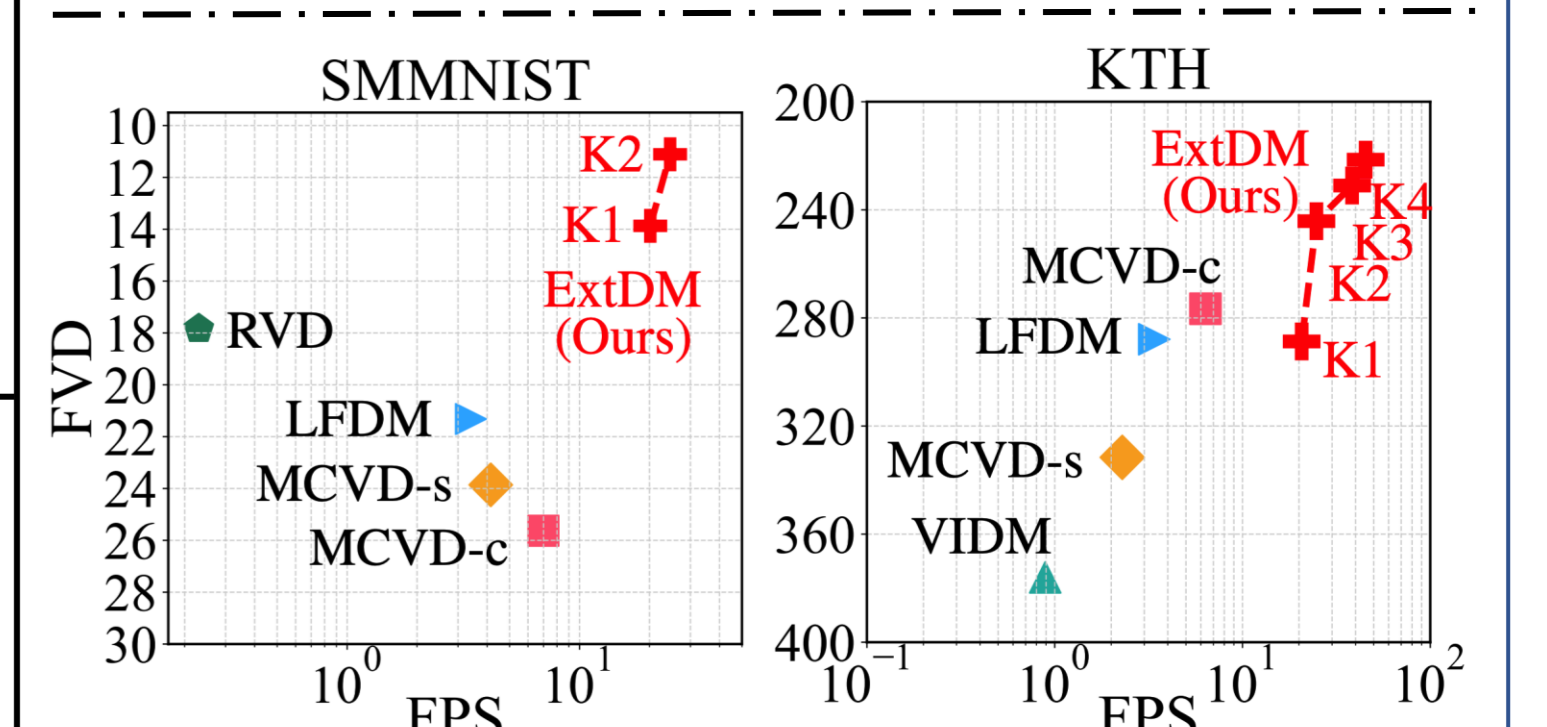
# Performance



Prediction results can be used to generate potential predictions and customize a preferred trajectory.



Frame-wise PSNR comparison on long-term video datasets.



Comparison of quality and speed of SOTA DMs for short- and long-term video prediction.

