



ExtDM: Distribution Extrapolation Diffusion Model for Video Prediction

Zhicheng Zhang^{1,2,†} Junyao Hu^{1,2,†} Wentao Cheng^{1,‡} Danda Paudel^{3,4} Jufeng Yang^{1,2}

¹ VCIP & TMCC & DISec, College of Computer Science, Nankai University

² Nankai International Advanced Research Institute (SHENZHEN·FUTIAN)

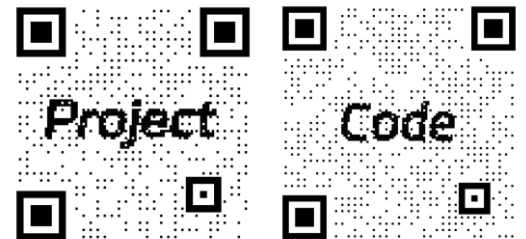
³ Computer Vision Lab, ETH Zurich ⁴ INSAIT, Sofia University



ETH zürich
INSAIT

<https://cv.nankai.edu.cn>

<https://github.com/nku-zhichengzhang/ExtDM>



- **Introduction**
- **Rethinking Previous Works**
- **ExtDM Architecture**
- **Experimental Results**
- **Conclusion**



Autonomous Driving



Sport Events

□ Video Prediction

□ Definition

It aims to capture the dynamic change from present x_c to future x_p .

□ Difference with Video Generation

building on existing video sequences v.s. creating from scratch

□ Application

Autonomous driving, sport events, video understanding, *etc.*



Prediction Performance of MCVD

Methods	<i>cond=10, pred = 40</i>				FPS↑
	SSIM↑	PSNR↑	LPIPS↓	FVD↓	
MCVD-c	0.793	26.20	0.124	276.6	6.35
MCVD-cpf	0.720	23.48	0.173	368.4	6.38
MCVD-s	0.744	26.40	0.115	331.6	2.29

Inference quality and speed of MCVD

Voleti V, Jolicoeur-Martineau A, Pal C. Mcvd-masked conditional video diffusion for prediction, generation, and interpolation[C]. NeurIPS, 2022.

□ Video Prediction

□ Challenges

□ Uncertainty and Complexity

especially in long-term video prediction

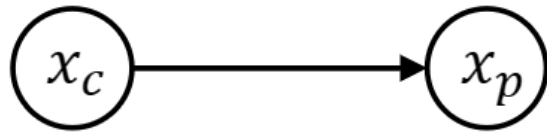
□ Modeling of Temporal Change

including dynamic variation and static background processing

□ Effectiveness and Usability

Trade-off between training computing cost and inference speed

Rethinking Previous Works



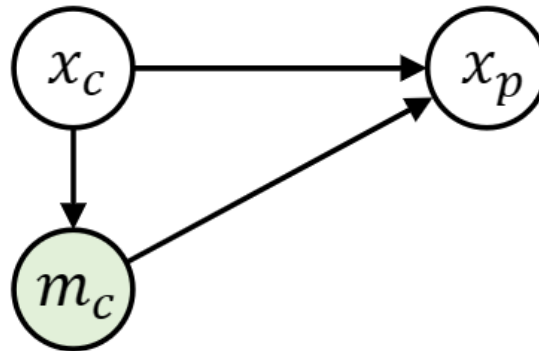
$$p(x_p|x_c)$$

Direct Method

- only **RGB**
- difficult to solve **complexity** in probability estimation

SRVP(ICML20)

SimVP(CVPR22)



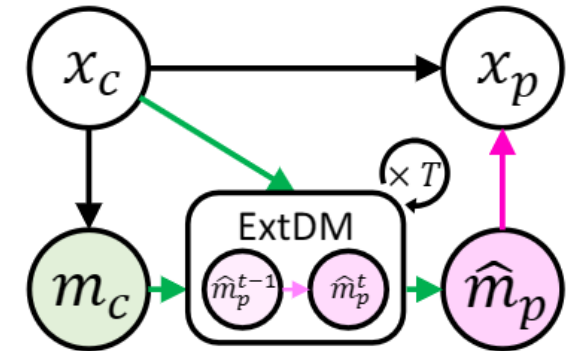
$$p(x_p|x_c, m_c)$$

In-context Learning Method

- **RGB + motion (implicit cues)**
- lack accuracy for **longer time**
- **counterfactual** results like fading, deformation, etc.

MCNet(ICLR17)

MOSO(CVPR23)

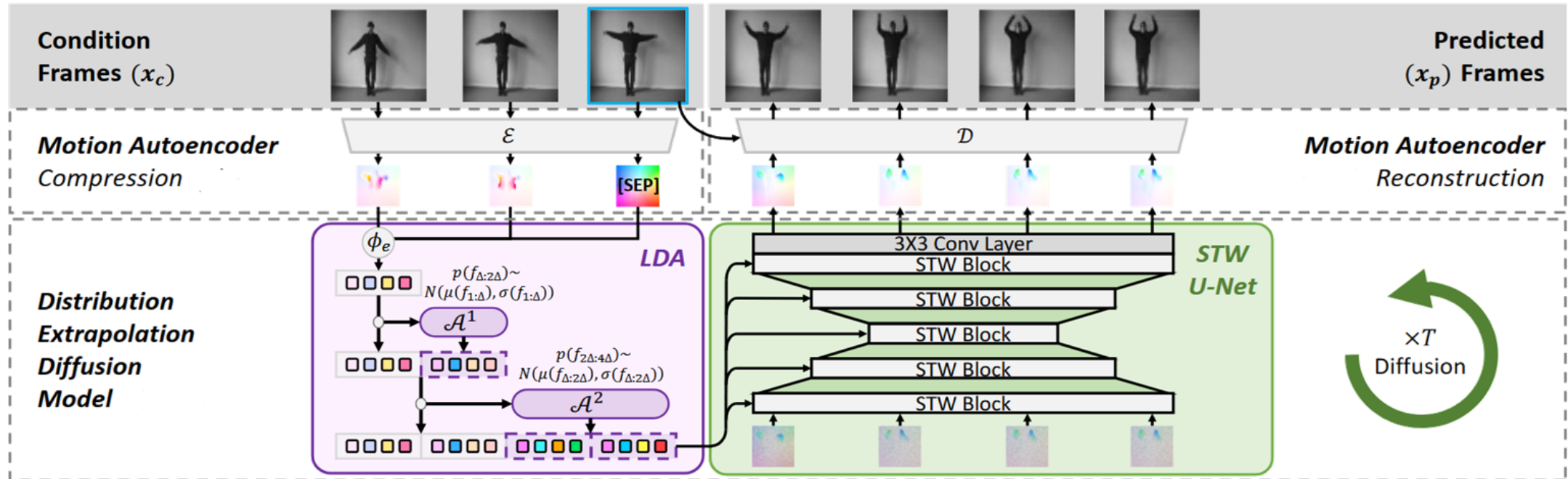


$$p(x_p|x_c, \zeta(x_c, m_c))$$

Extrapolation Method (Ours)

- **RGB + motion (explicit cues)**
- **Extrapolate** present deterministic motion cues into the future ones

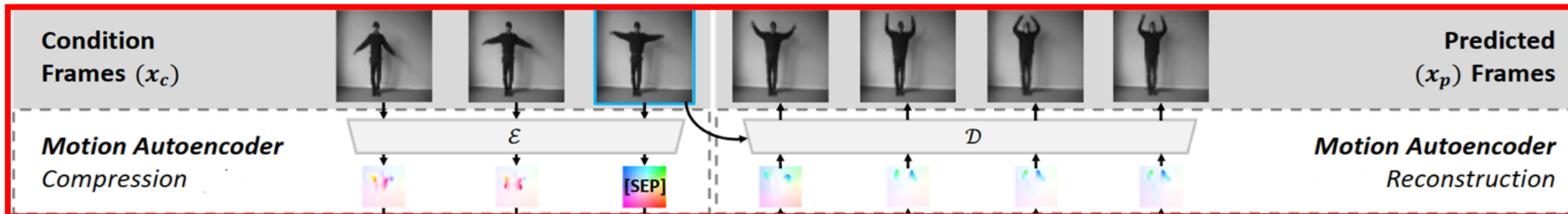
ExtDM Architecture



✓ Contributions

- ❑ A **distribution extrapolation** DM that predict future frames.
- ❑ An efficient VP 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 5 video prediction datasets.

ExtDM Architecture

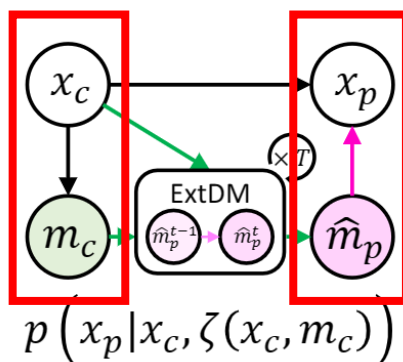


Compression

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

Reconstruction

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



Two Mapping Functions

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

- Motion Autoencoder Compression & Reconstruction

ExtDM Architecture

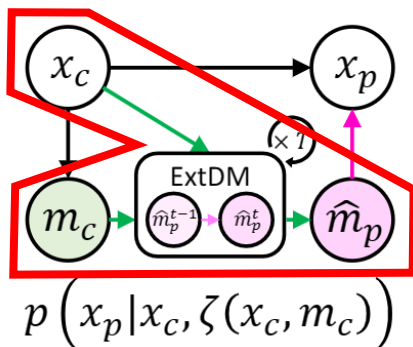
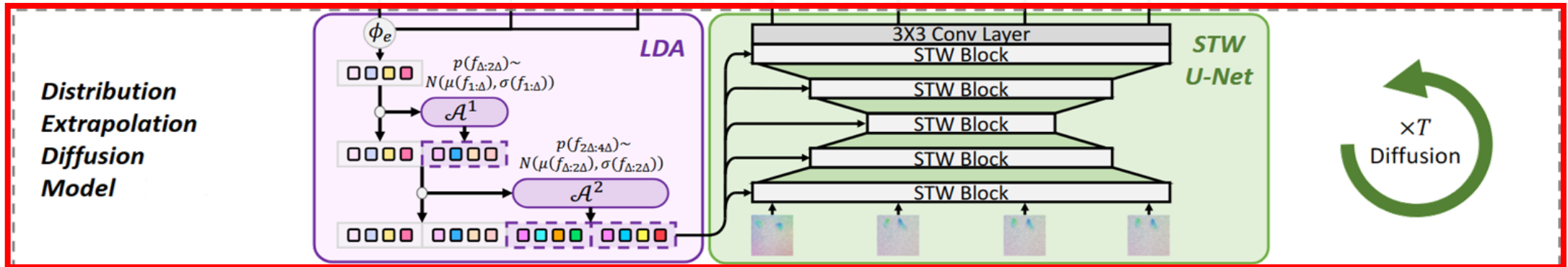
Condition
Frames (x_c)

Predicted
Frames (x_p)

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

Motion Autoencoder
Compression

Motion Autoencoder
Reconstruction

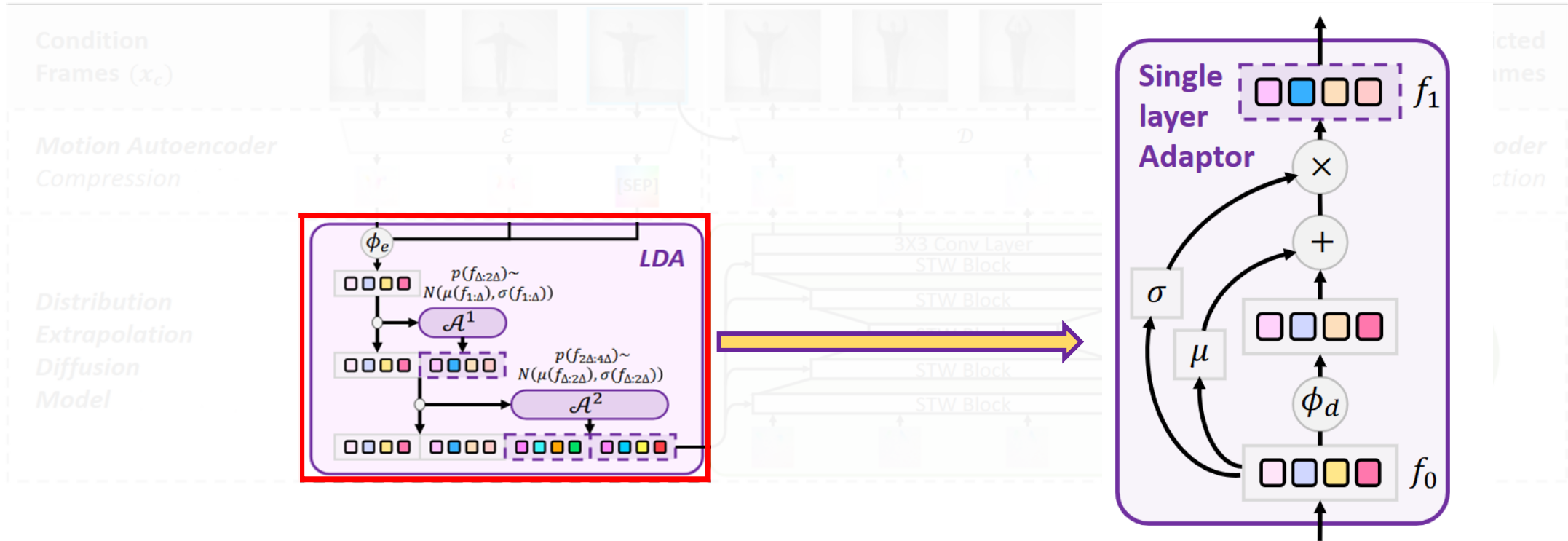


Two Mapping Functions

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

- Distribution Extrapolation Diffusion Model

ExtDM Architecture



Layered Distribution Adaptor

- estimate distribution params
- inference using distribution sampling

$$f_{1:\Delta} = \phi_e(f_c),$$

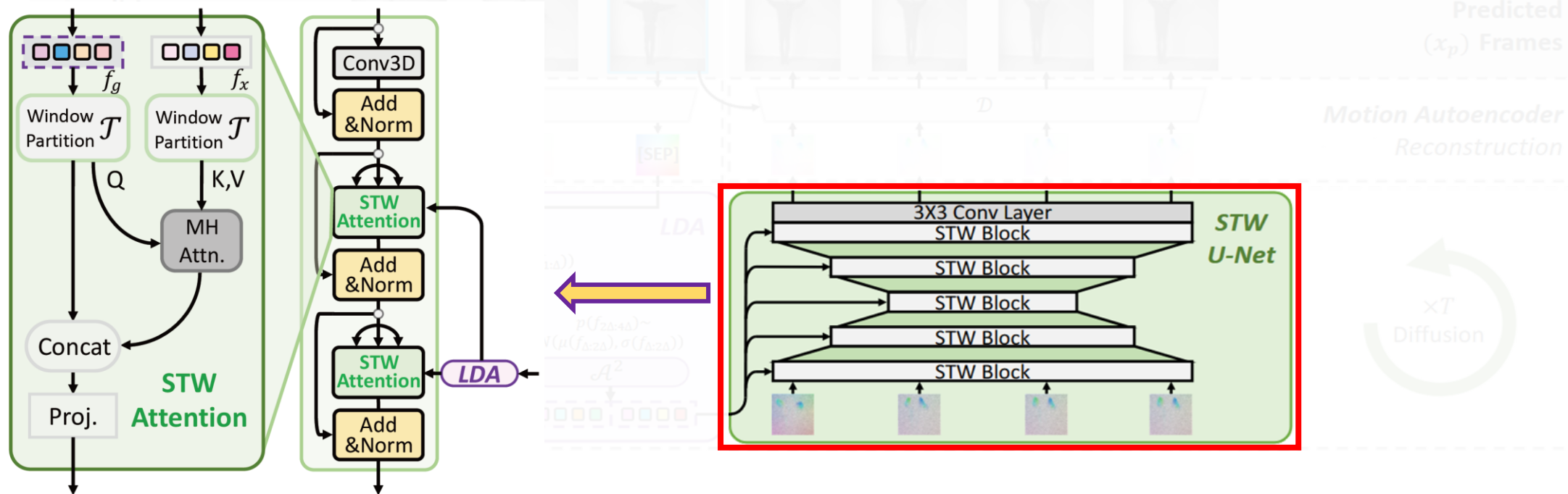
$$\hat{f}_{1:2^l\Delta} = (f_{1:2^{l-1}\Delta}, \mathcal{A}^{(l)}(f_{1:2^{l-1}\Delta})),$$

$$f_p = (\hat{f}_{1:\Delta}, \dots, \hat{f}_{2^{L-1}\Delta:2^L\Delta}).$$

$$f_b = \mathcal{A}(f_a)$$

$$= (\sigma(f_a) + \sigma')\phi_d \left(\frac{f_a - \mu(f_a)}{\sigma(f_a)} \right) + \mu(f_a) + \mu'$$

ExtDM Architecture

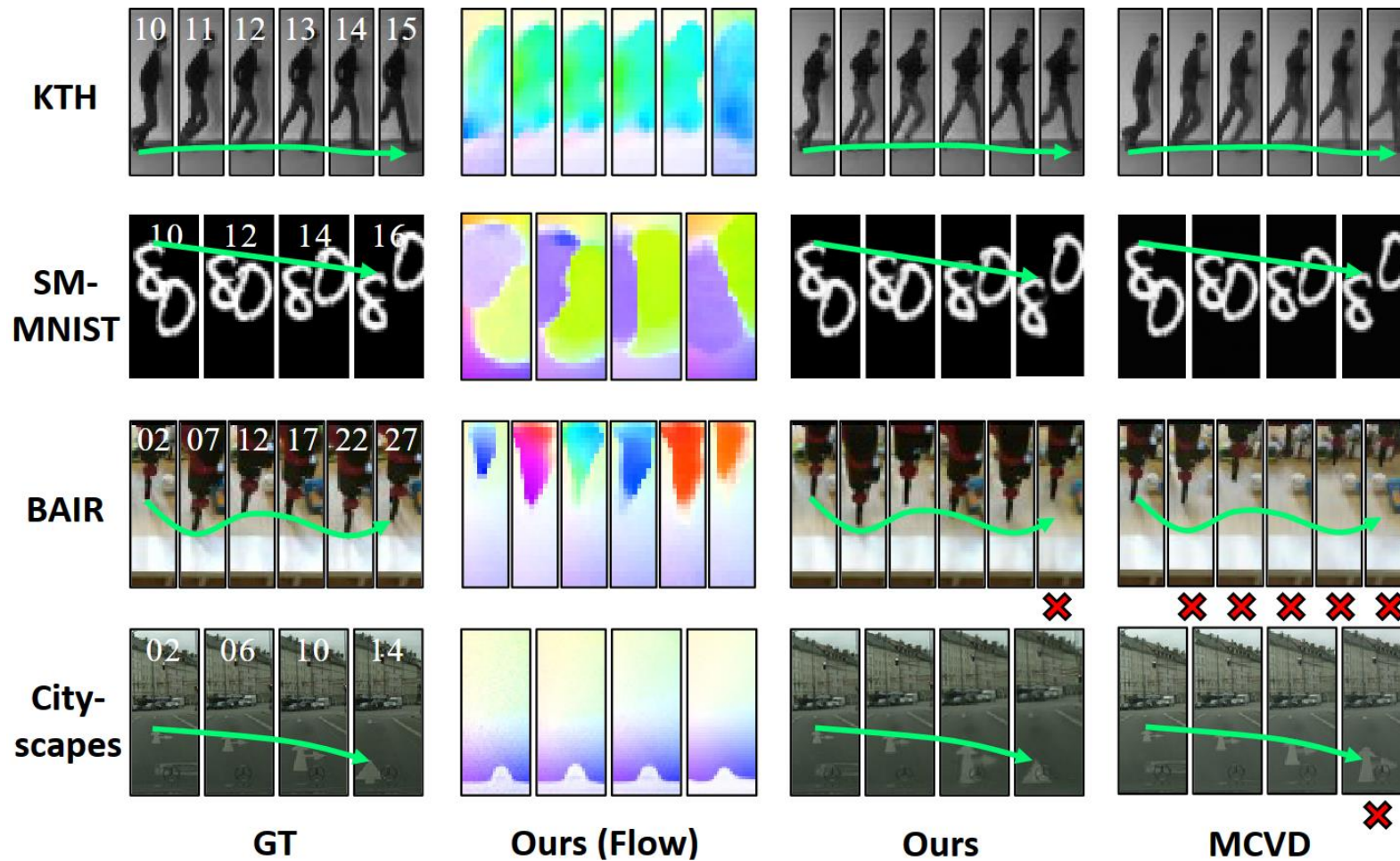


Spatiotemporal Window U-Net

- exploit the spatiotemporal coherence interaction via jointly conducting strided and grid window

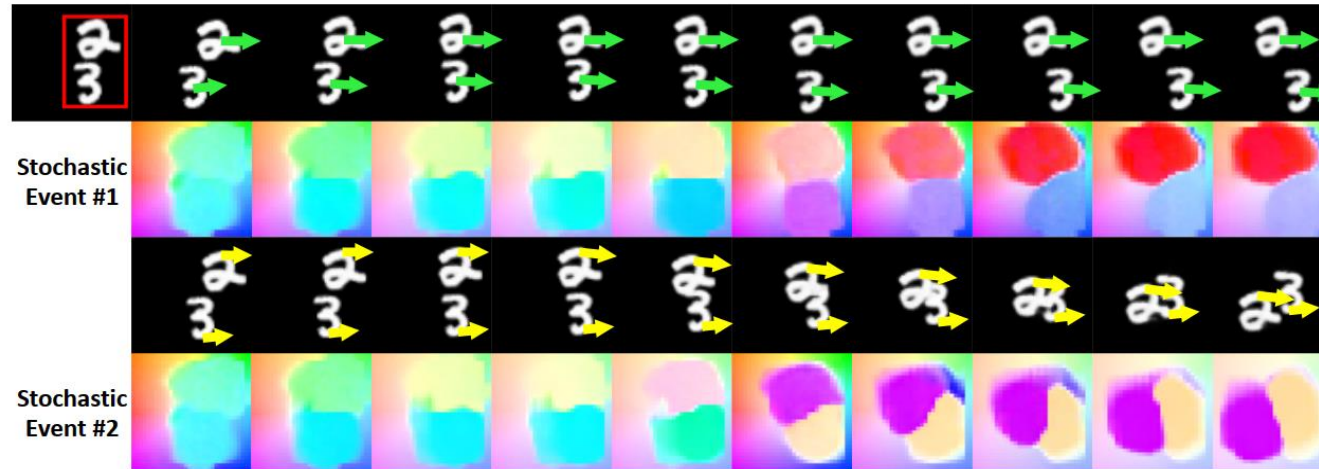
$$f_{x \rightarrow g} = \text{softmax}\left(\frac{[\mathcal{T}(f_x)\mathbf{W}^Q][\mathcal{T}(f_g)\mathbf{W}^K]^\top}{\sqrt{d}}\right)\mathcal{T}(f_x)\mathbf{W}^V$$

Experimental Results

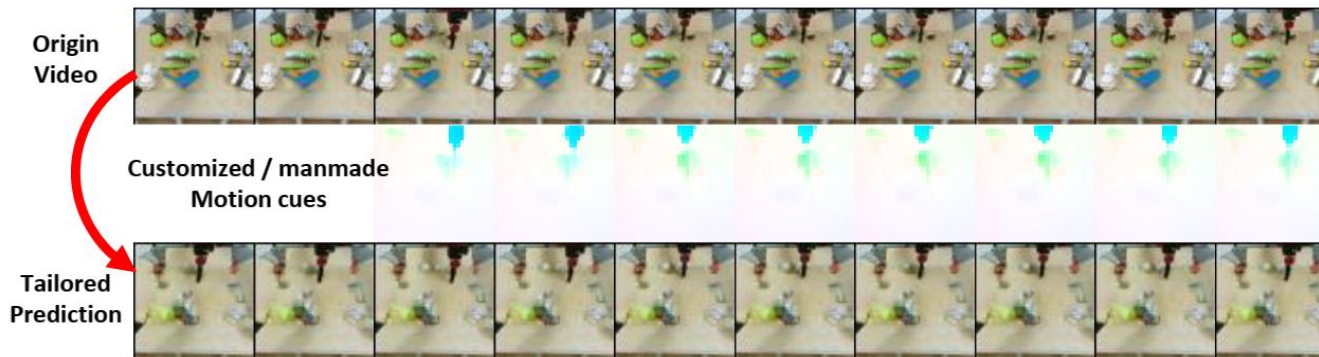


✓ It can predict the videos with **correct trajectories** of objects (**green curve** in the figure).

Experimental Results



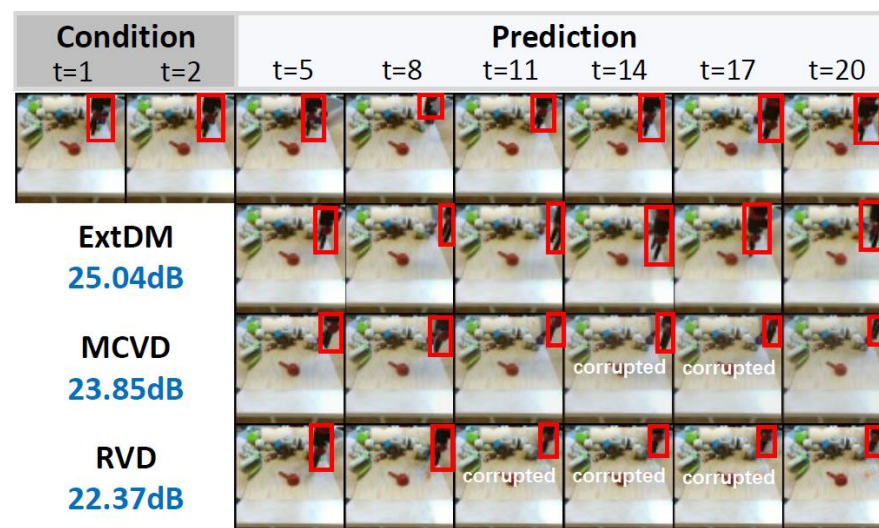
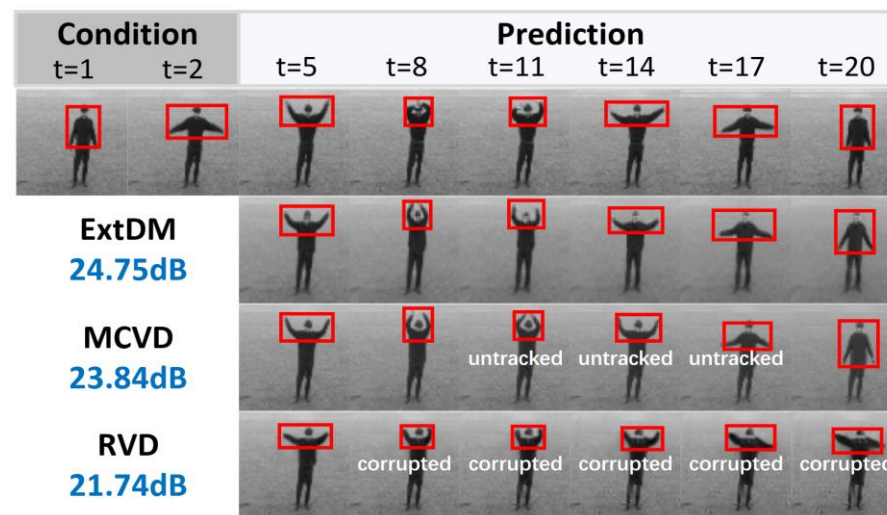
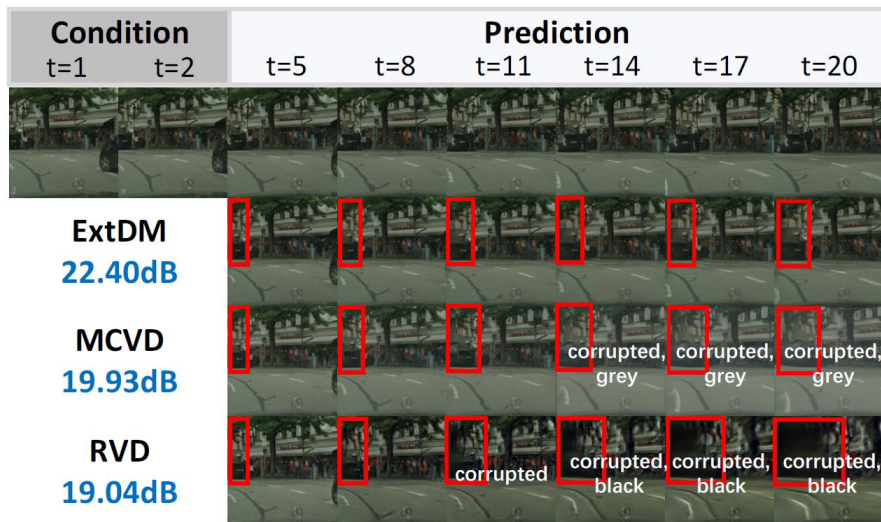
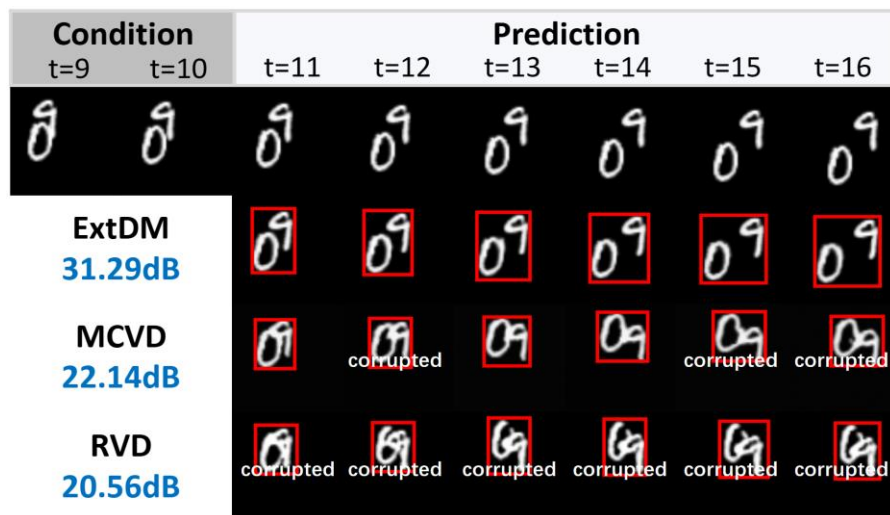
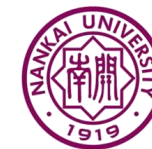
(a) Stochastic Events



(b) Tailored Predictions

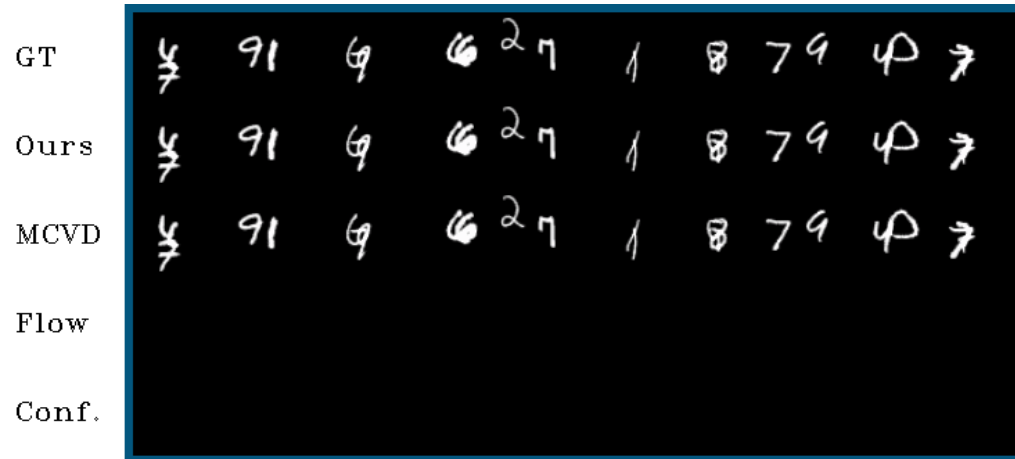
- ✓ Prediction results can be used to
 - (a) generate potential predictions
 - (b) customize a preferred trajectory.

Experimental Results



✓ Qualitative comparison on SMMNIST, KTH, Cityscapes and BAIR.

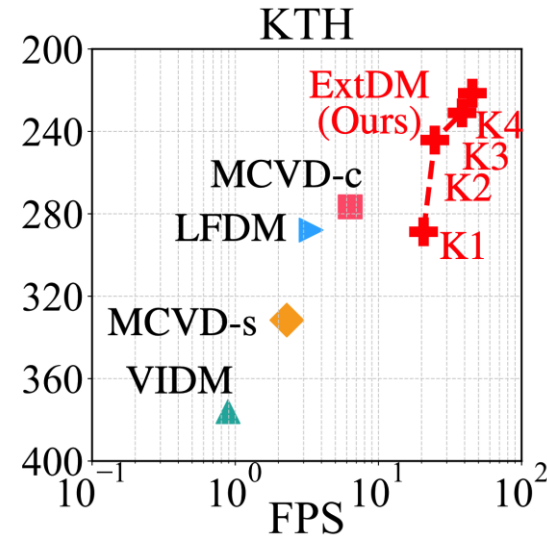
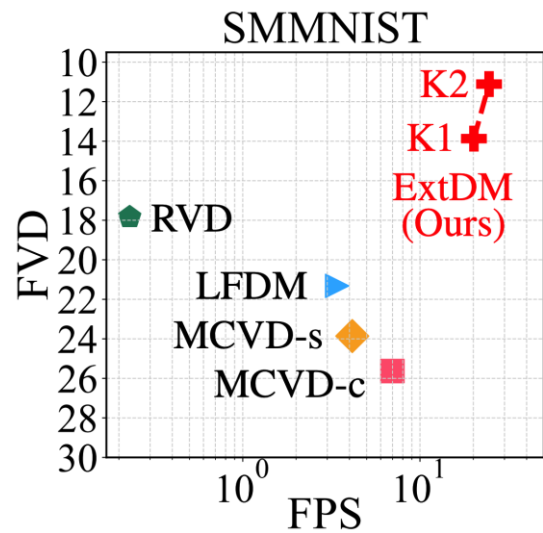
Experimental Results



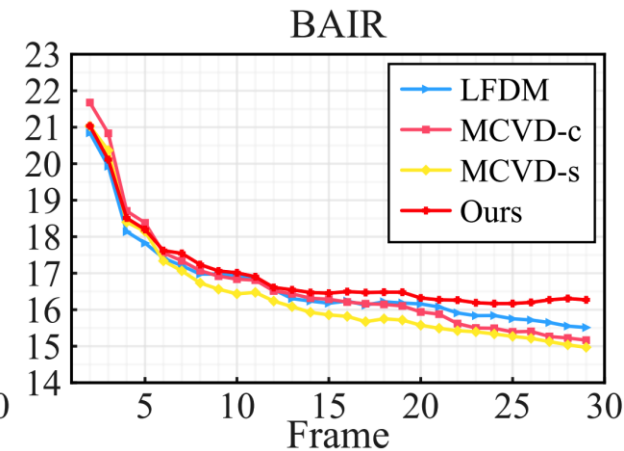
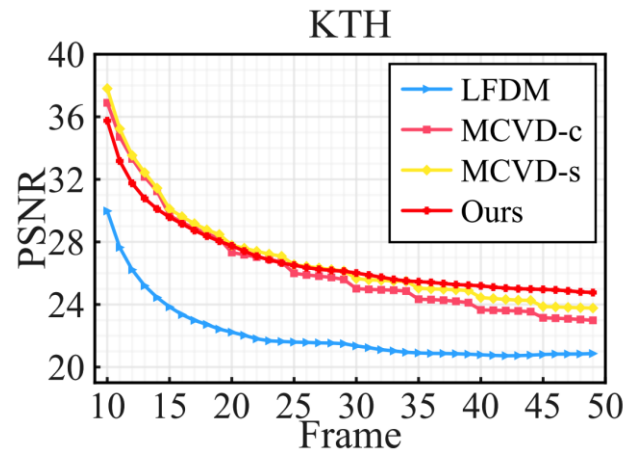
✓ Qualitative comparison on SMMNIST, KTH, Cityscapes and BAIR.



Experimental Results



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



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

Thank you



ExtDM: Distribution Extrapolation Diffusion Model for Video Prediction

<https://cv.nankai.edu.cn>

<https://github.com/nku-zhichengzhang/ExtDM>

