



浙江大學  
ZHEJIANG UNIVERSITY



ICML  
International Conference  
On Machine Learning

# World-R1: Reinforcing 3D Constraints for Text-to-Video Generation

Weijie Wang<sup>1,2,\*†</sup>, Xiaoxuan He<sup>1,\*</sup>, Youping Gu<sup>1,\*</sup>, Yifan Yang<sup>2,‡</sup>  
Zeyu Zhang<sup>3</sup>, Yefei He<sup>1</sup>, Yanbo Ding<sup>2</sup>, Xirui Hu<sup>3</sup>  
Donny Y. Chen<sup>3</sup>, Zhiyuan He<sup>2</sup>, Yuqing Yang<sup>2,‡</sup>, Bohan Zhuang<sup>1,‡</sup>

<sup>1</sup>Zhejiang University   <sup>2</sup>Microsoft Research   <sup>3</sup>Independent Researcher

\* Equal contribution

† Work done during an internship at MSRA

‡ Corresponding authors

# Video Foundation Models

Prompt: Camera move left. Modernist glass skyscrapers reflecting the Shanghai Bund waterfront during golden hour.

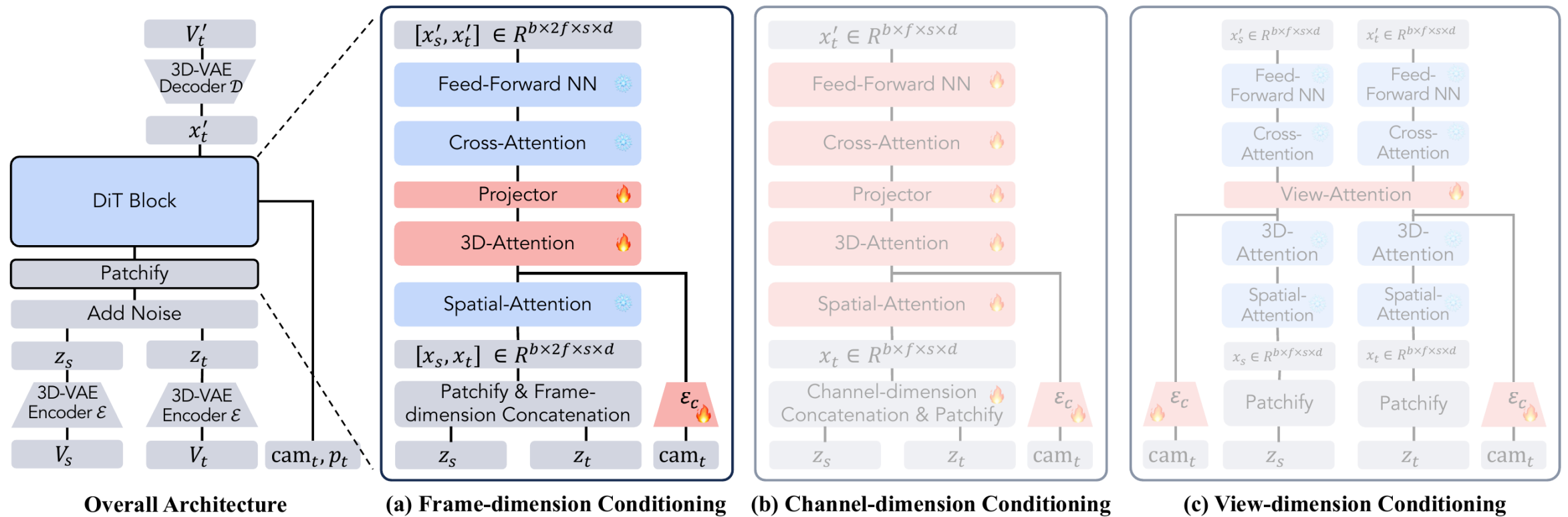


**Poor controllability**



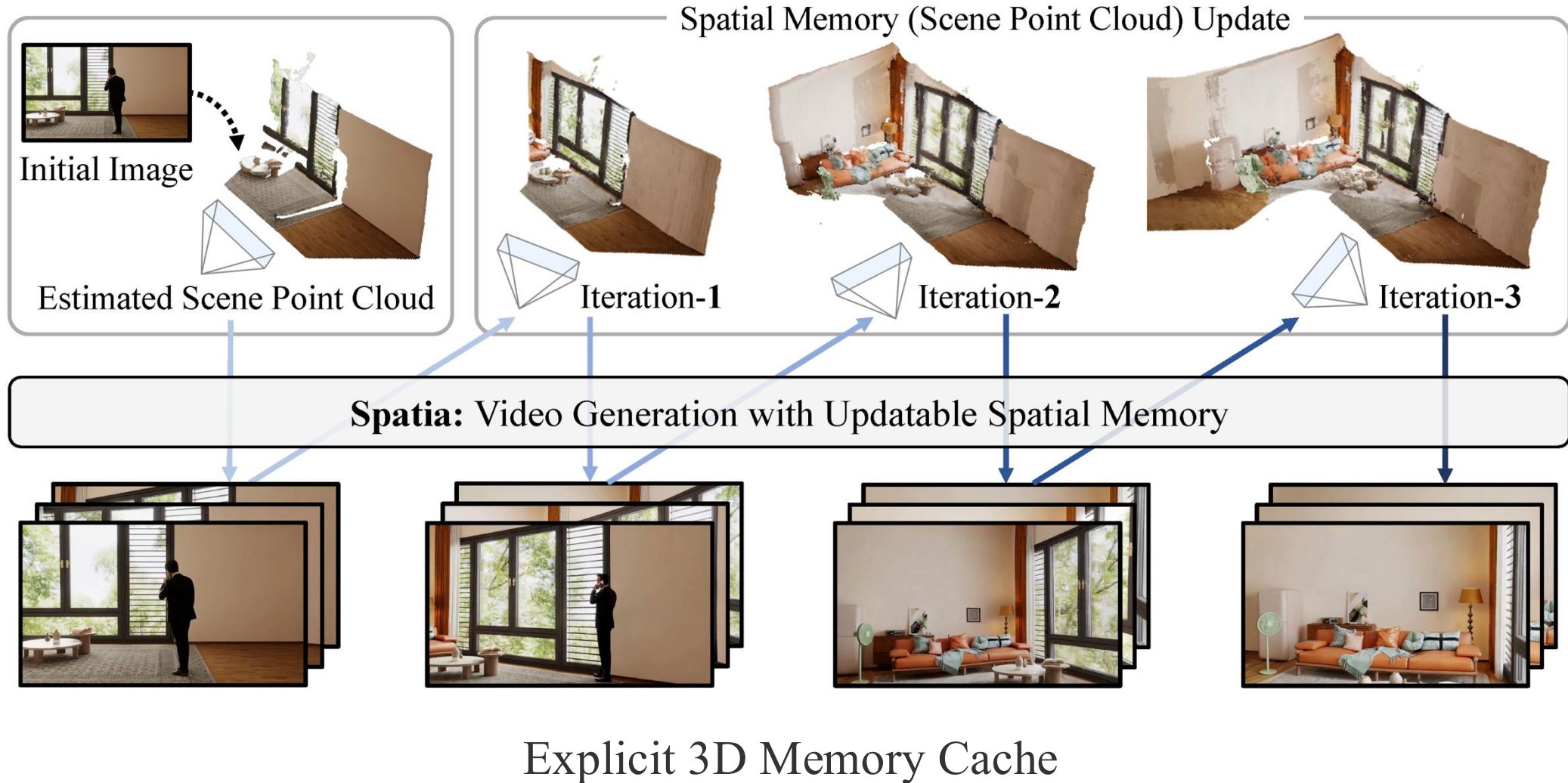
**Poor 3D consistency**

# Existing Methods for Camera Control



New Modules for Controller

# Existing Methods for 3D Consistency



# Motivation

No Explicit Memory Module

No Control Module

No architectural  
modifications

No Specific Video Training Data

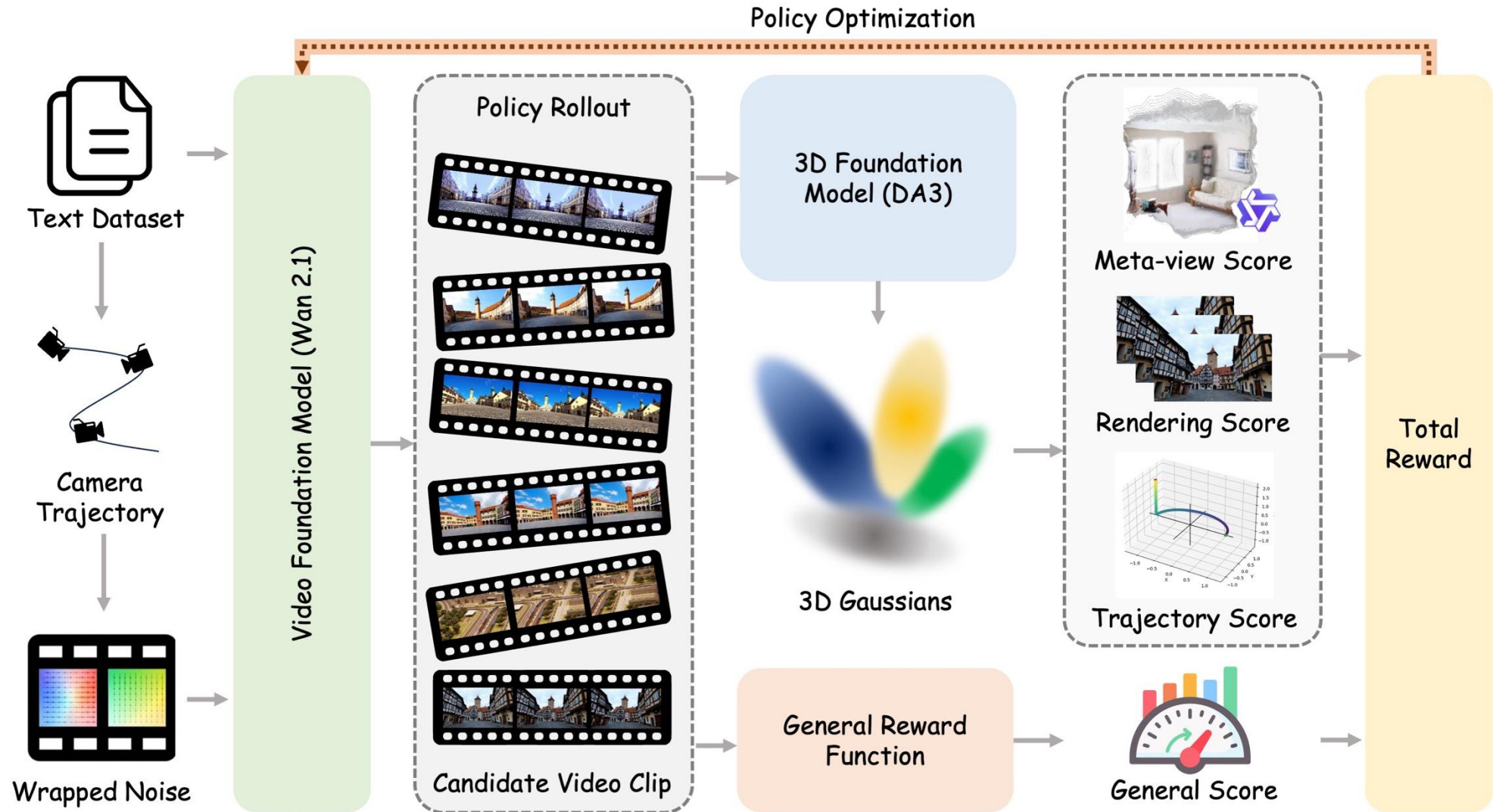
No extra inference cost

Not only applicable to I2V

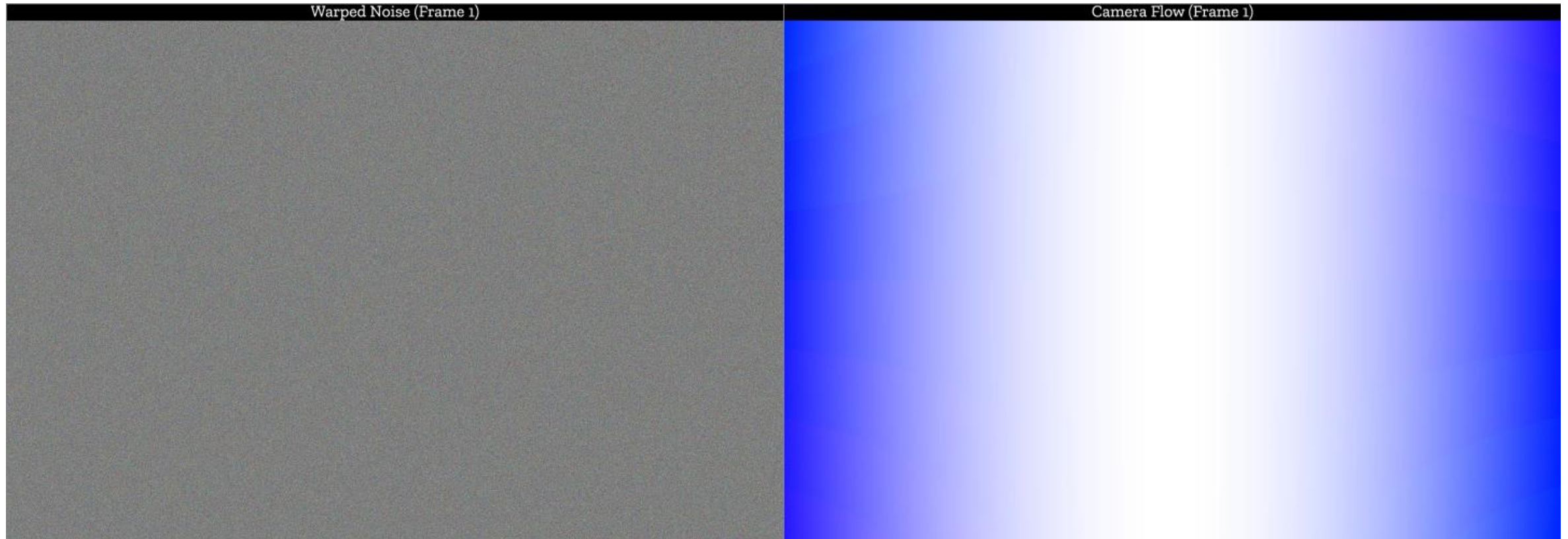


**ALL-IN-ONE-METHOD**

# Overview



# Implicit Camera Conditioning



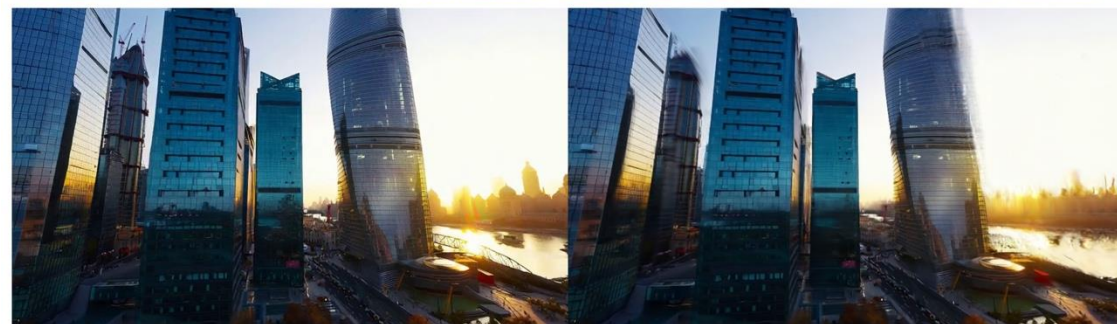
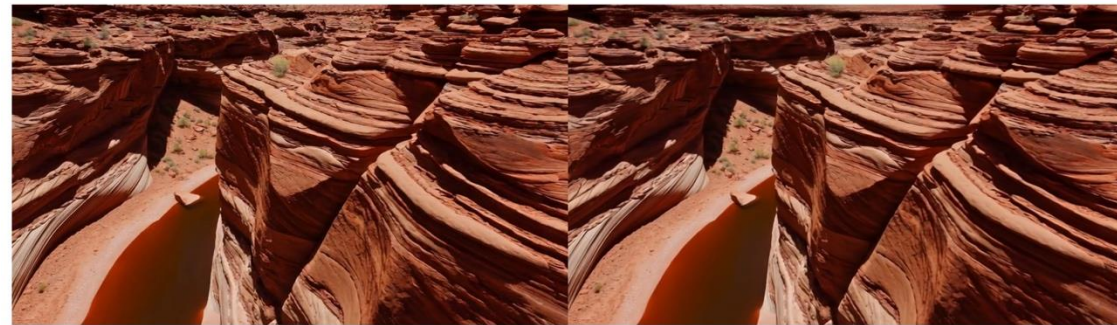
Instead of using Gaussian noise, using warped noise

# Reward Design: Rendering Score



Generated Video Frame

Reconstructed Video Frame



Generated Video Frame

Reconstructed Video Frame

Re-rendered visual alignment metrics

# Reward Design: Meta-view Score

Camera pan right.  
Over a messy  
teenager's room.



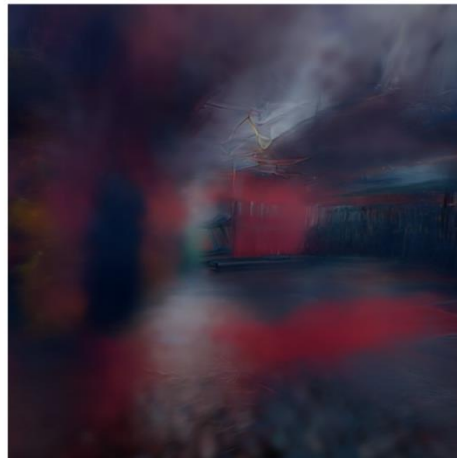
Meta score: 2



Camera push in, then  
pan left. A  
basketball court  
next to a red sports  
arena.



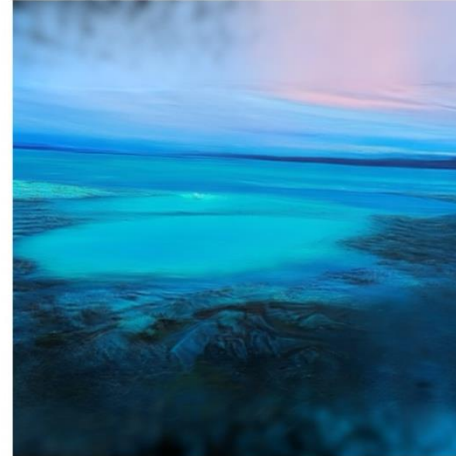
Meta score: 1



Camera pull out, then  
orbit right. A  
dramatic reveal of a  
glacial lake, circling  
the turquoise water.



Meta score: 7



Camera push in.  
Windmills on the  
grassland.



Meta score: 8



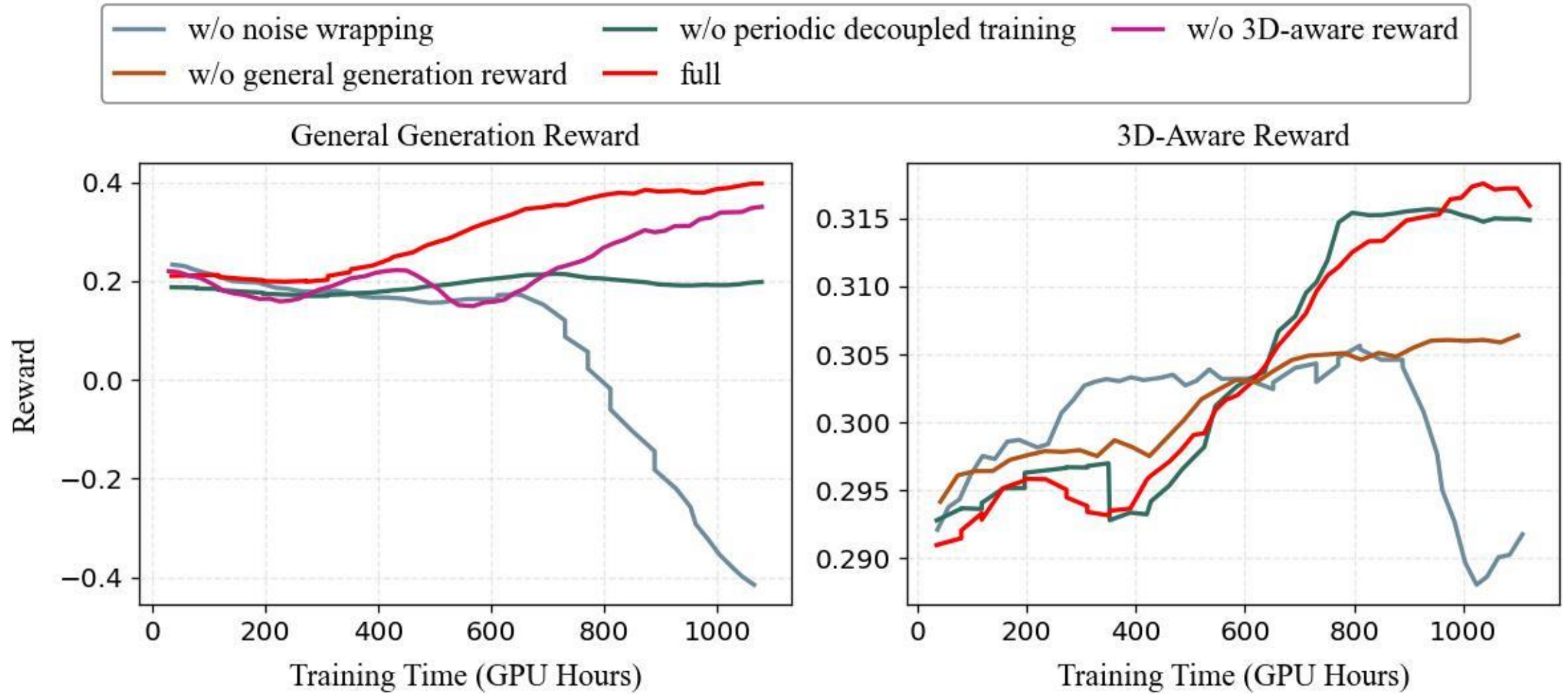
VLM determines the quality of the generated scene

# Reward Design: Trajectory Score



Error in calculating the injected camera trajectory

# Training Strategy



Periodic Decoupled Training ensures the quality of dynamic scenes.

# Quantitative Results

Method	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
CogVideoX-1.5-5B (Yang et al., 2025)	24.44	0.783	0.242
Wan2.2-T2V-14B (Wan et al., 2025)	23.47	0.779	0.253
Wan2.2-T2V-5B (Wan et al., 2025)	22.36	0.716	0.303
Wan2.1-T2V-14B (Wan et al., 2025)	19.76	0.629	0.405
Wan2.1-T2V-1.3B (Wan et al., 2025)	17.40	0.550	0.467
<b>World-R1-Small (Ours)</b>	27.63	0.858	0.201
<b>World-R1-Large (Ours)</b>	<b>27.67</b>	<b>0.865</b>	<b>0.162</b>

**+10dB PSNR, +0.3 SSIM, -0.3 LPIPS**

# Quantitative Results

Method	Aesthetic Quality ↑	Imaging Quality ↑	Motion Smoothness ↑	Subject Consistency ↑	Background Consistency ↑
CogVideoX-1.5-5B (Yang et al., 2025)	62.07	65.34	98.15	96.56	96.81
Wan2.1-T2V-1.3B (Wan et al., 2025)	62.43	66.51	97.44	96.34	97.29
GCD (Van Hoorick et al., 2024)	38.21	41.56	98.37	88.94	92.00
Trajectory-Attention (Xiao et al., 2024)	38.50	51.00	98.21	90.60	92.83
DAS (Gu et al., 2025)	39.86	51.55	99.14	90.34	92.03
ReCamMaster (Bai et al., 2025a)	42.70	53.97	99.28	92.05	93.83
<b>World-R1-Small (Ours)</b>	65.74	67.53	98.55	97.58	96.67

Even higher general video generation metrics

# Comparison with Baselines

Camera push in. Deep canyon walls made of layered red rock, with a winding river at the bottom.



**Wan2.1-T2V-14B**

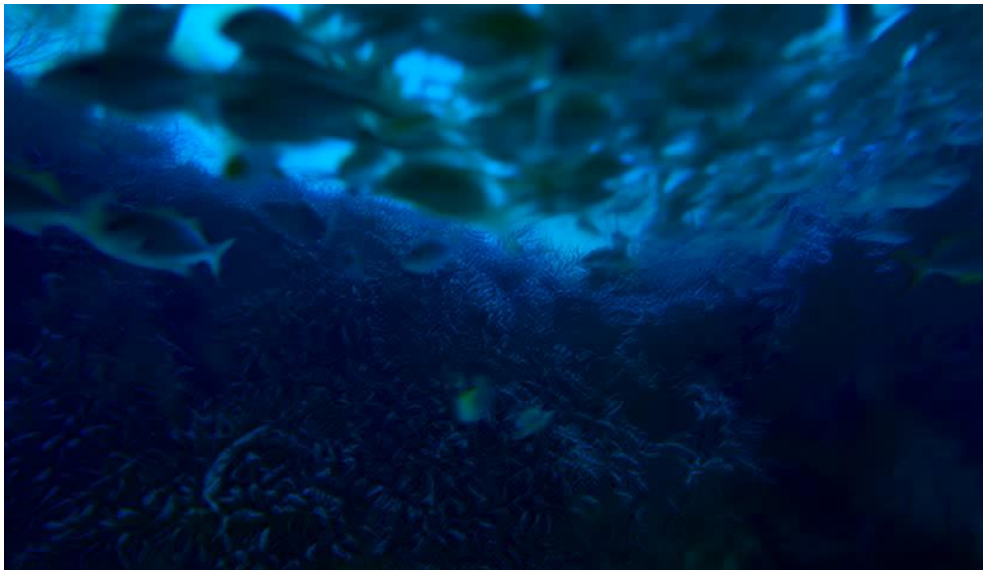
**Wan2.2-T2V-14B**



**World-R1-Large (Ours)**

# Comparison with Baselines

Camera orbit left, then push in.  
Deep sea coral reefs teeming  
with colorful fish and  
bioluminescent plant life.



**Wan2.1-T2V-14B**



**Wan2.2-T2V-14B**



**World-R1-Large (Ours)**

# Comparison with Baselines

Camera move left. Modernist glass skyscrapers reflecting the Shanghai Bund waterfront during golden hour.



**Wan2.1-T2V-1.3B**

**CogVideoX-1.5-5B**



**World-R1-Small (Ours)**

# Results of Dynamic Scene

A lion roaring with its mane shaking in the wind.



Camera pan right. A drone flying through a complex obstacle course.



World-R1-Large (Ours)

# Results of Dynamic Scene

World-R1-Large (Ours)

Camera pan left. Soldiers marching in synchronization across a dusty field.



Camera move left. A fighter jet performing an aileron roll.



# More Information



Paper, code and more visualizations are available on our project page.



Weijie Wang's homepage.  
Actively seeking cooperation opportunities.

## Conclusion:

- **World-R1** reinforces 3D constraints for text-to-video generation without extra inference-time modules.
- **Improving** camera control and 3D coherence across static and dynamic scenes.
- **Can** be scaled up to the video foundation model.