

# Towards Intelligent Interactive 3D-aware Video World Model

**Weijie Wang**

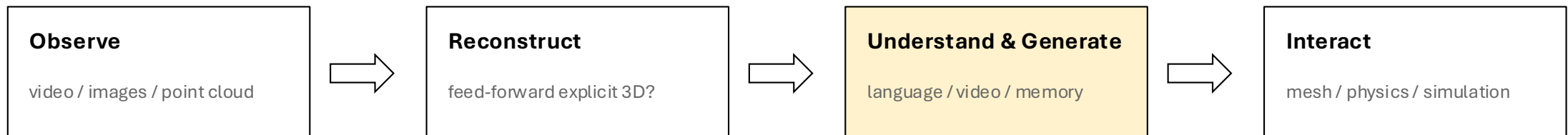
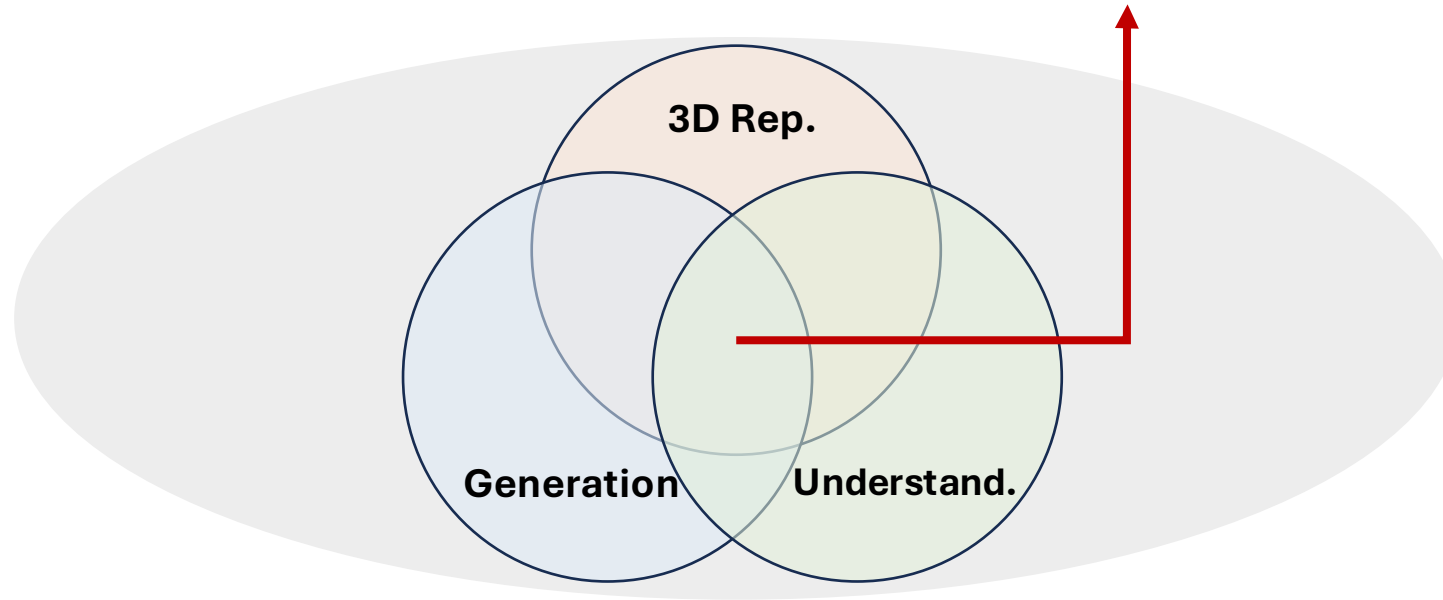
<http://lhmd.top>

College of Computer Science and Technology, Zhejiang University

Ph.D. student @ ZIP Lab / Intern @ ByteDance Seed

2026/06/17

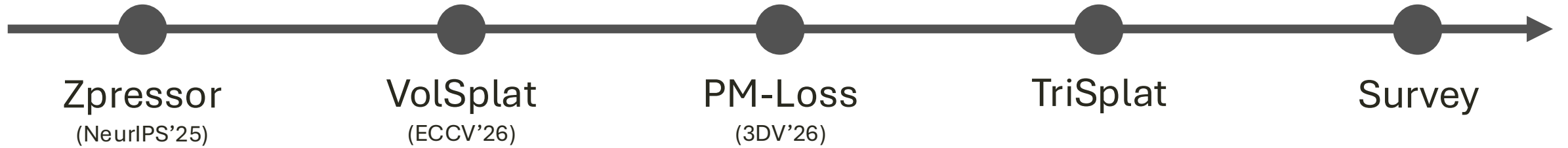
## Build an intelligent, interactive, 3D-aware world model



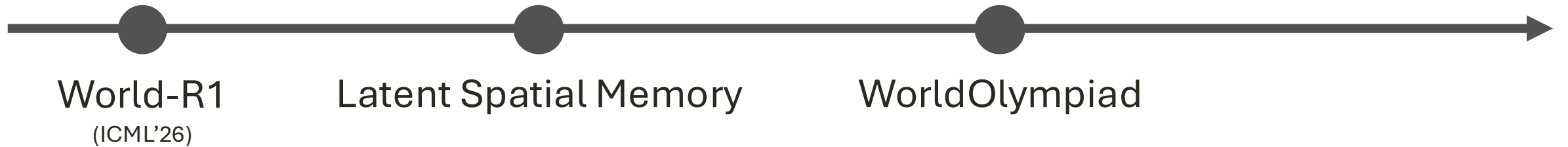
**3D is the interface that can connect reconstruction, generation, understanding, and physical interaction.**

# Two Key Directions

## Key Direction 1: Breaking the limits of Feed-Forward 3D Reconstruction



## Key Direction 2: How to construct a 3D-aware video world model



# Enhance 3D Consistency for Video World Models via **Post-Training**

# Video Foundation Models

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

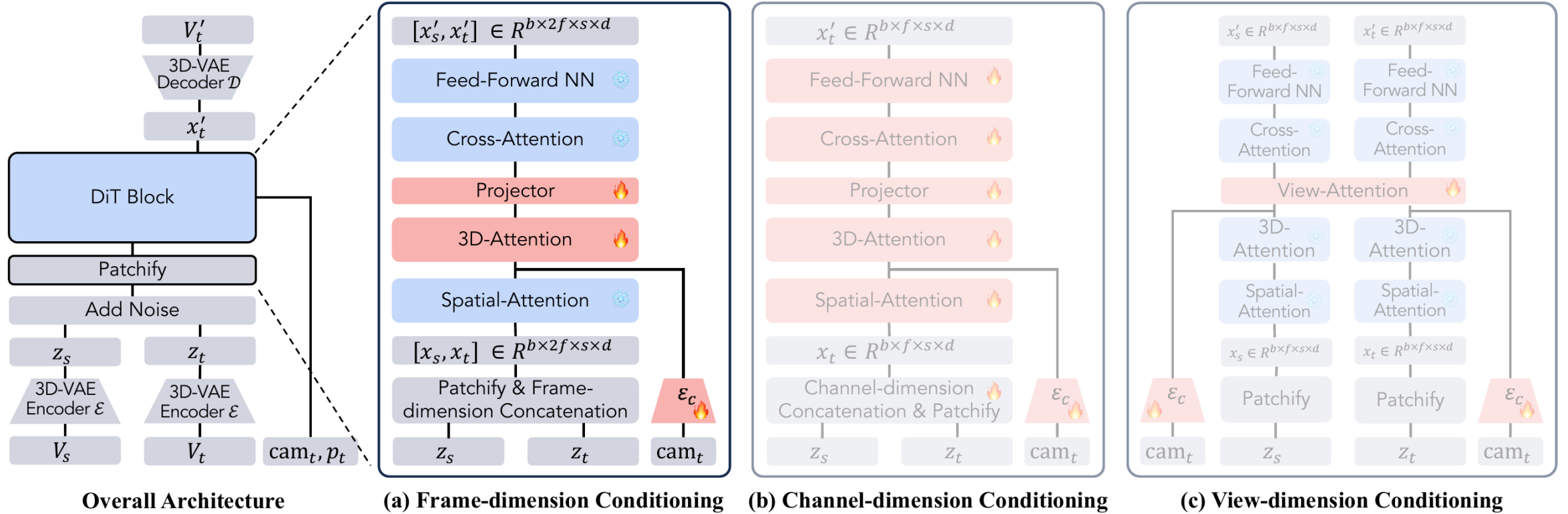


**Poor controllability**



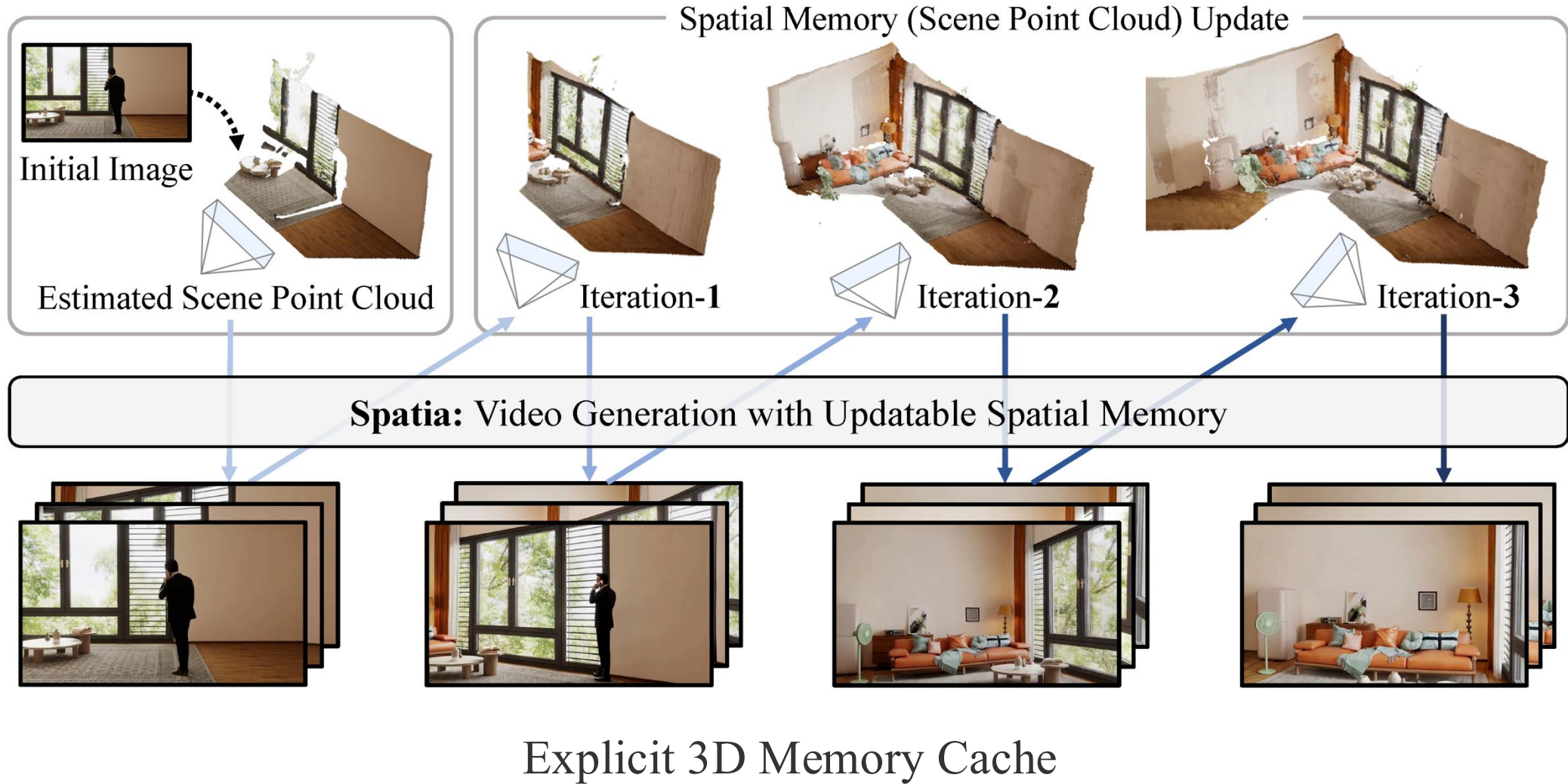
**Poor 3D consistency**

# Existing Methods for Camera Control



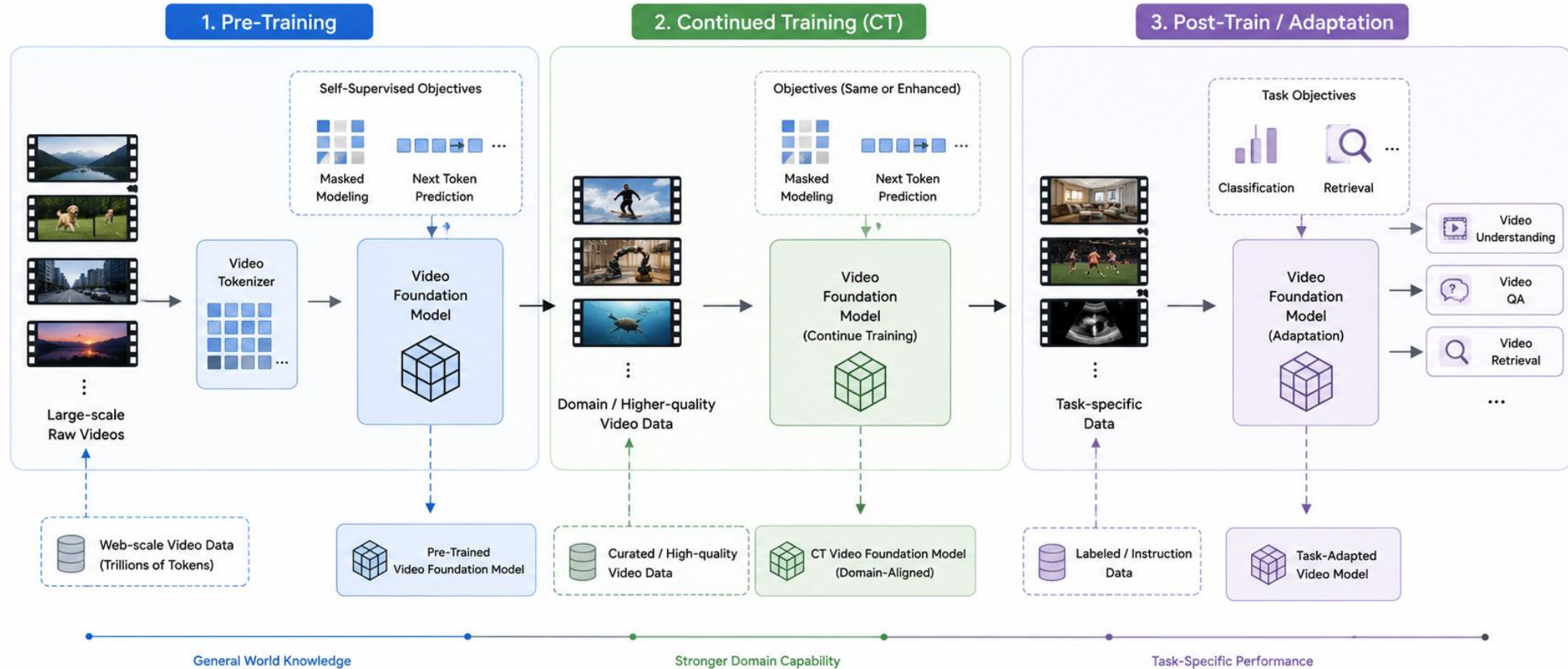
New Modules for Controller

# Existing Methods for 3D Consistency



# Video Foundation Models

## Video Foundation Model Training Pipeline



**Where to inject 3D information? Pretrain (Data)? CT? or PT?**

# 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**

# Our Answer: World-R1 (Post-Train)

Camera push in, An orange-hued 3D game world with an industrial aesthetic, featuring skyscrapers in the distance and railways alongside in-game weapons in the foreground.



Camera push in. A delicate afternoon tea set with macarons and pastries arranged on a tiered stand.



Camera move right. Giant spheres of shimmering liquid metal hovering silently over a mirrored lake.



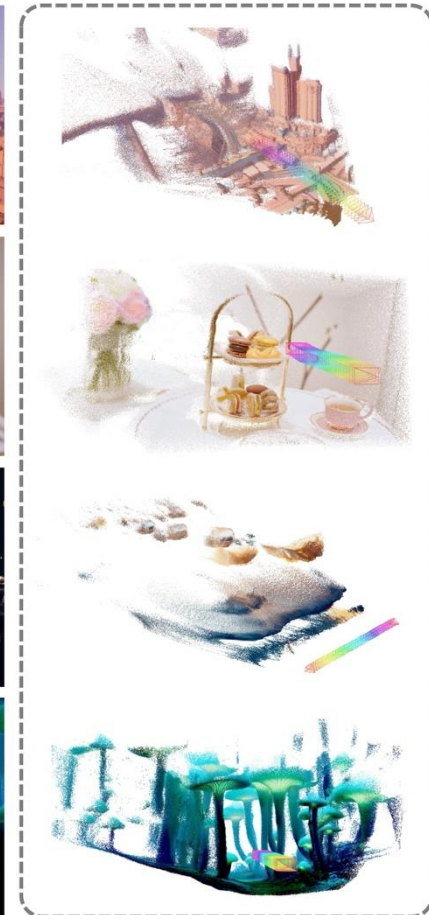
Camera push in and turn right. A forest of towering fungi that glow with soft blue and green light in the dark.



Text Prompt

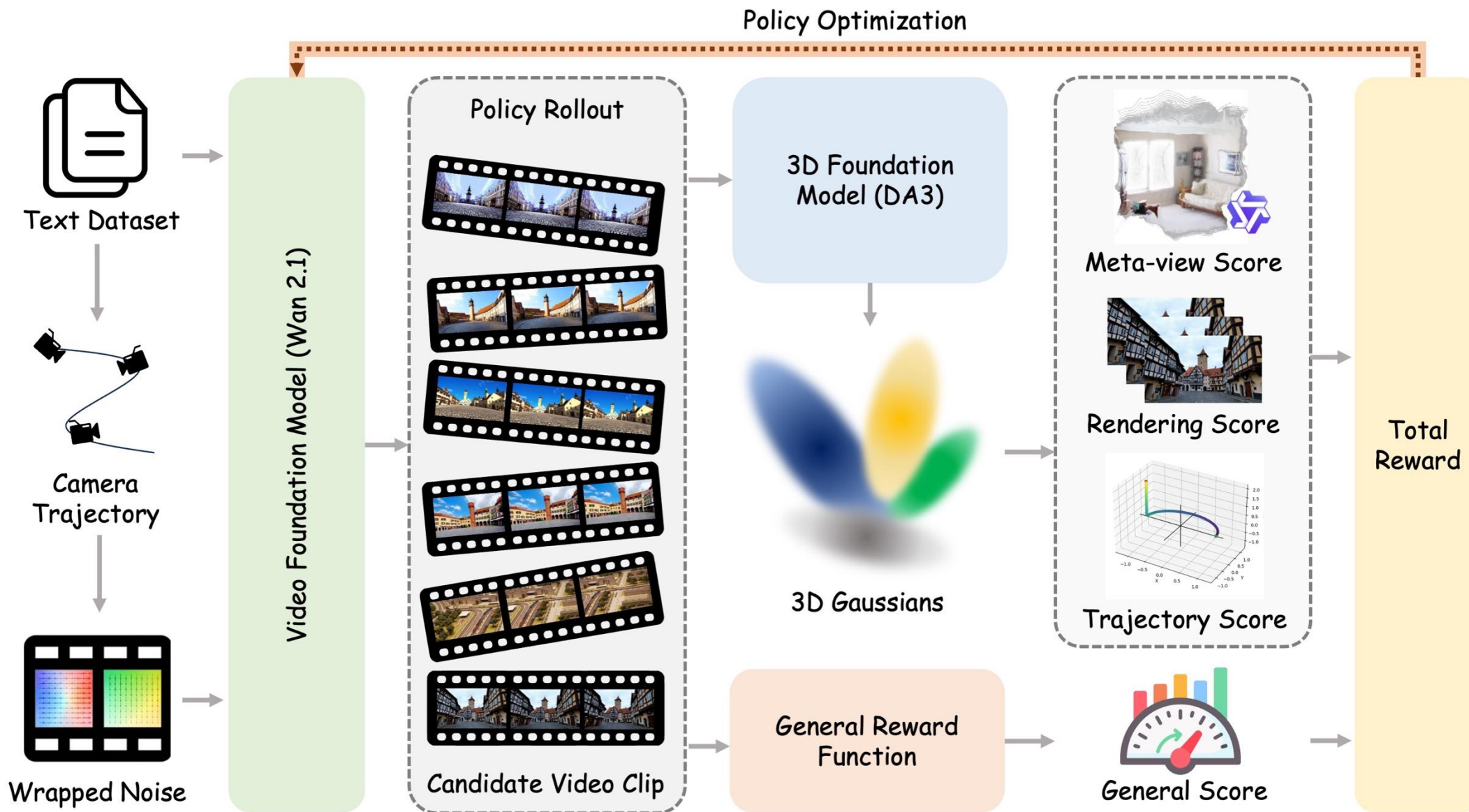
Generated Video

3D World Visualization

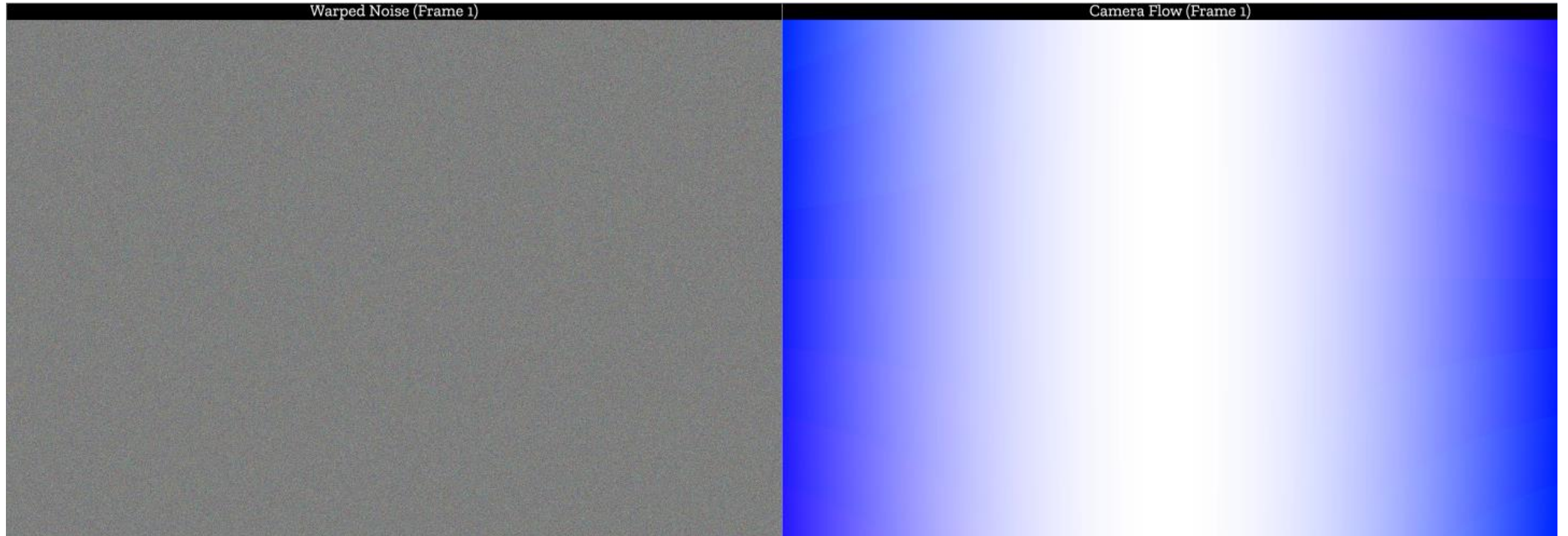


**3D consistency can be converted into RL rewards for text-to-video generation.**

# World-R1



# Implicit Camera Conditioning



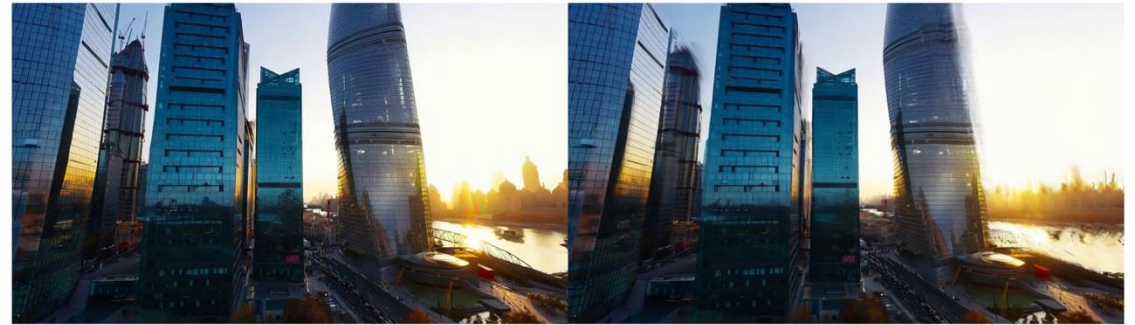
Instead of using Gaussian noise, using warped noise

# 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.



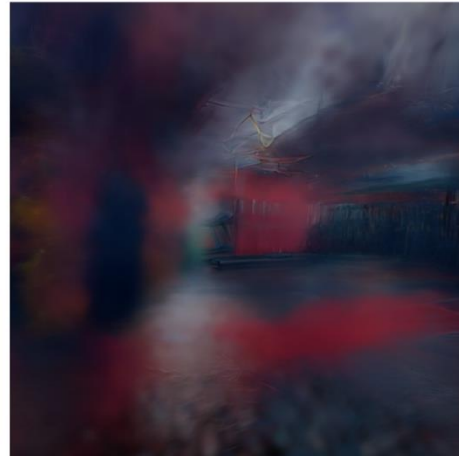
Meta score: 2

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



Meta score: 7

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



Meta score: 1

Camera push in.  
Windmills on the  
grassland.



Meta score: 8

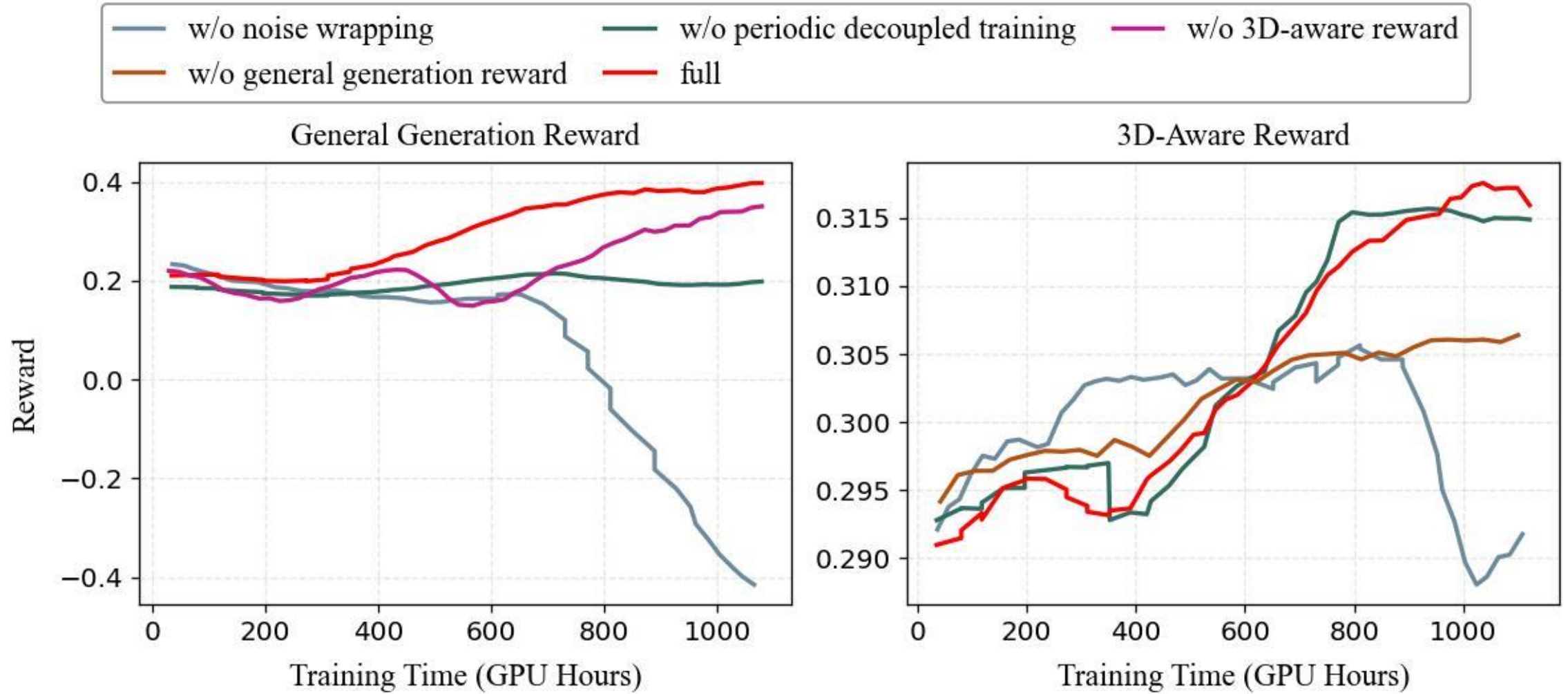
VLM determines the quality of the generated scene

# 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>	<b>65.74</b>	<b>67.53</b>	<b>98.55</b>	<b>97.58</b>	<b>96.67</b>

Even higher general video generation metrics

# Comparison with Baselines

Wan2.2-T2V-14B

Weijie Wang

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



Wan2.1-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

## World-R1-Large (Ours)

A lion roaring with its mane shaking in the wind.



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



# Results of Dynamic Scene

## World-R1-Large (Ours)

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

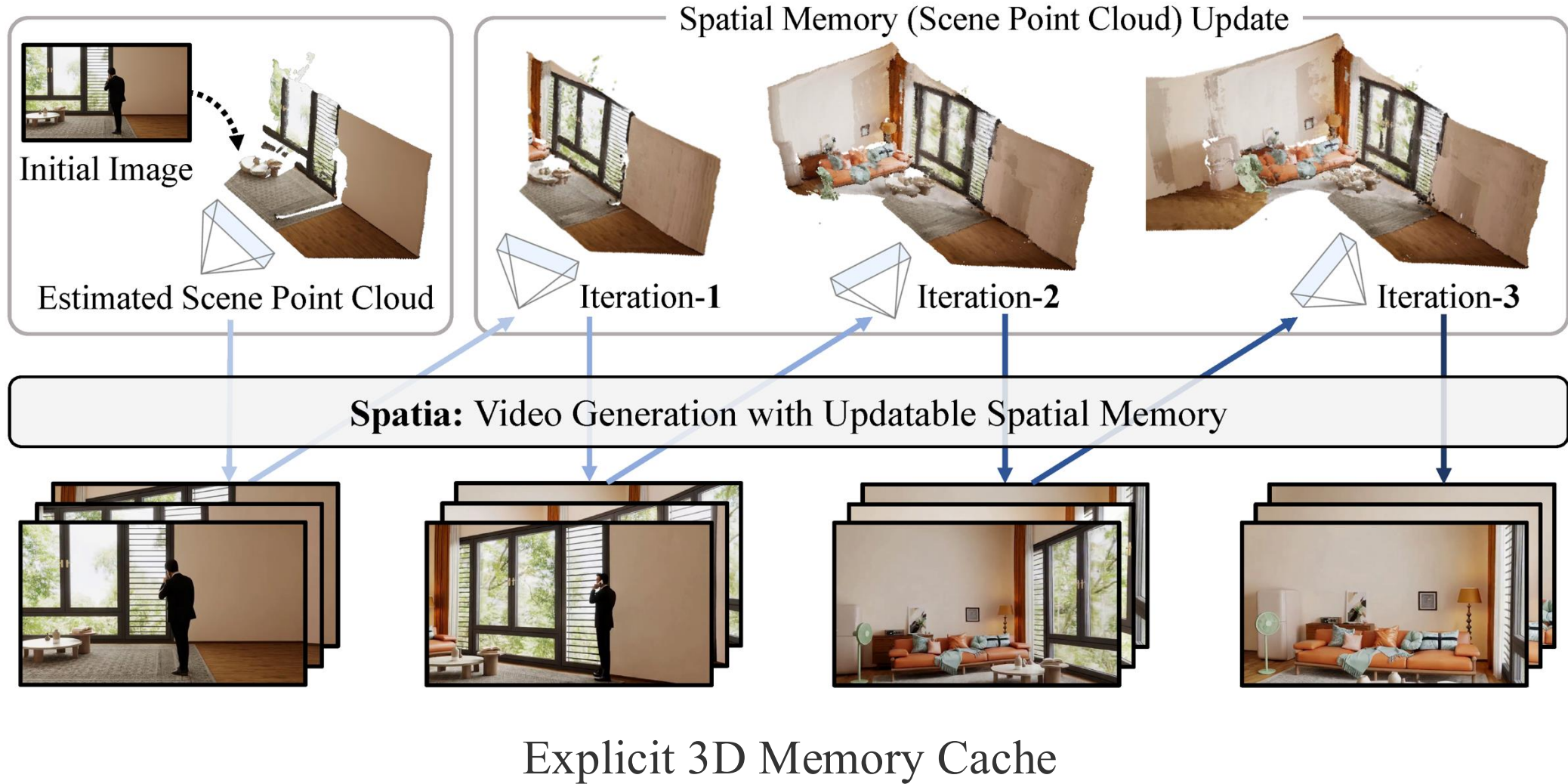


Camera move left. A fighter jet performing an aileron roll.

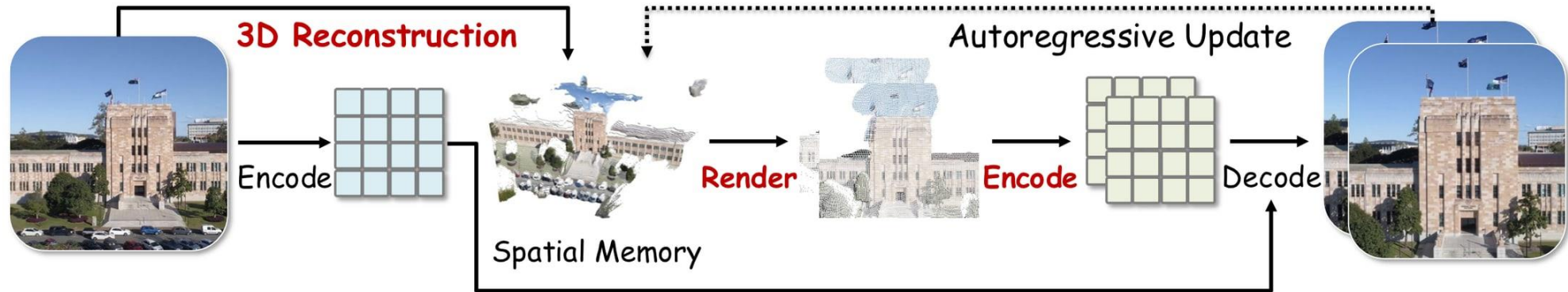


**But it's so expensive! We can do much more with **3D memory!****

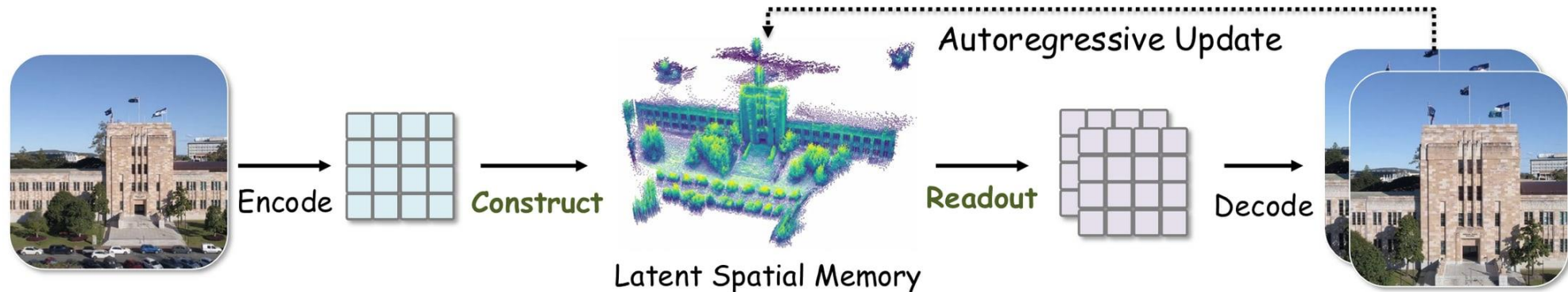
# Review: Existing Methods for 3D Consistency



# Latent Spatial Memory



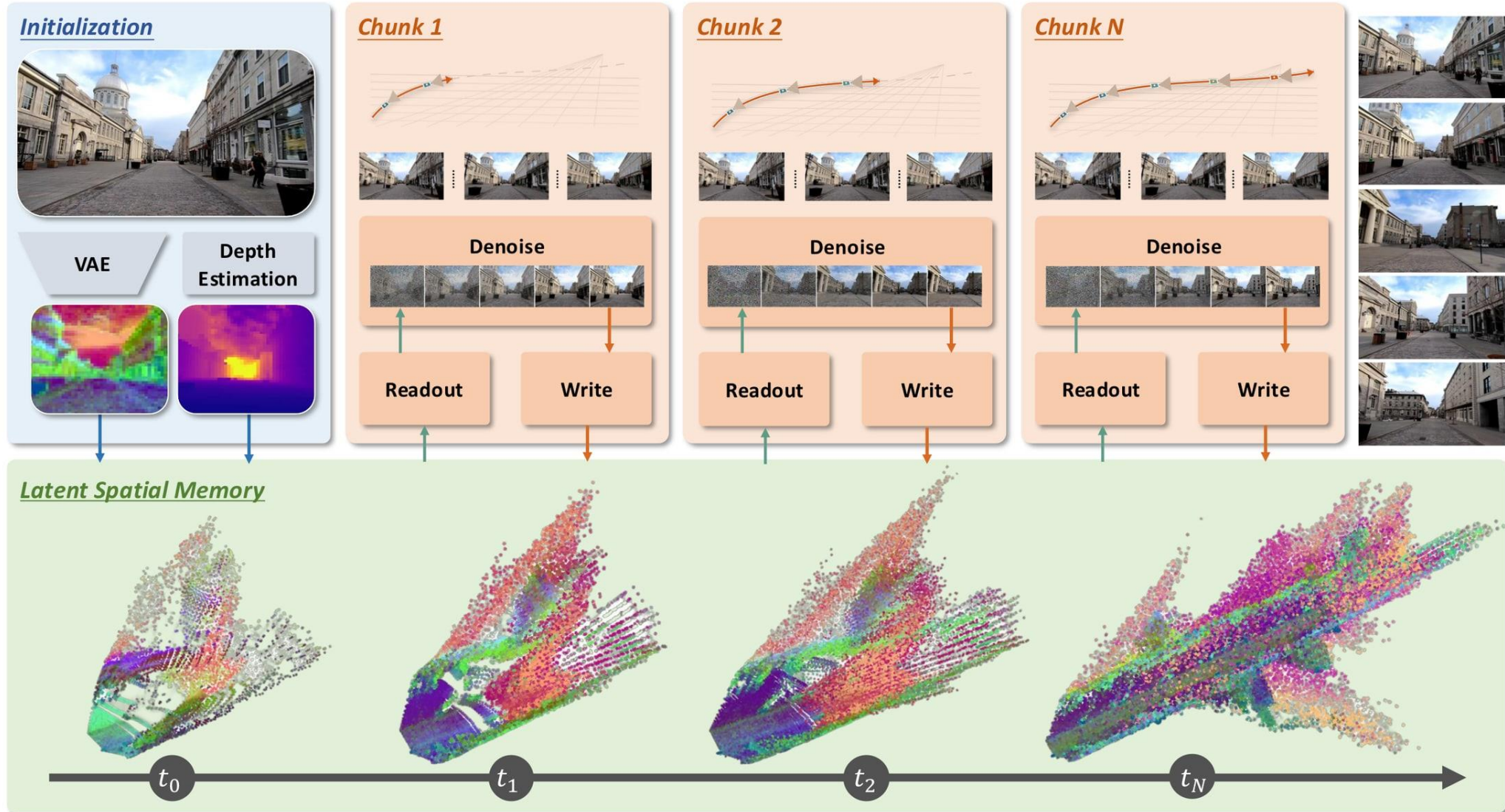
Spatial Memory for Video World Models



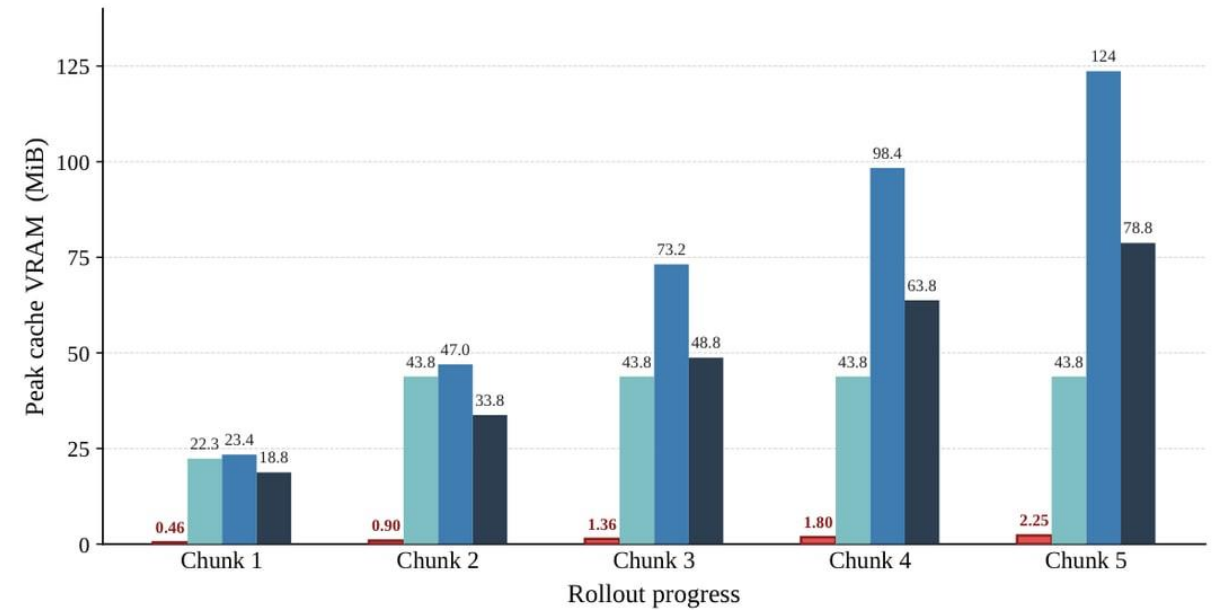
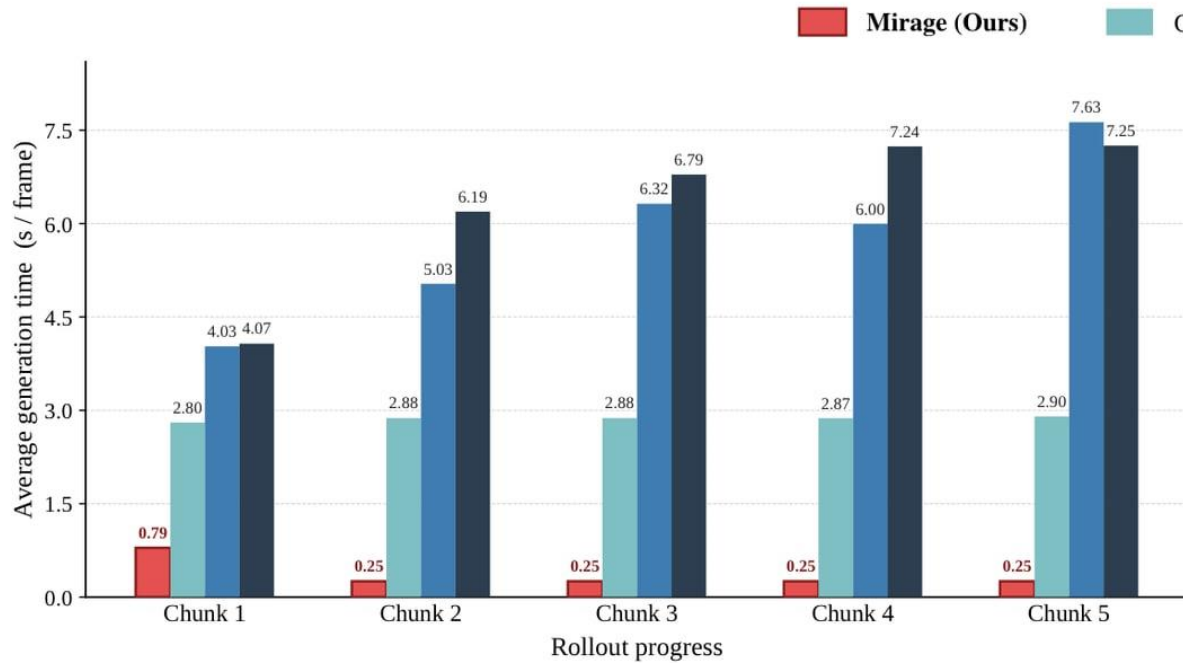
Latent Spatial Memory for Video World Models (Ours)

**Video World Models can use Latent Spatial Memory to maintain 3D Consistency**

# Mirage: Pipeline with Latent Spatial Memory



# Latent Spatial Memory

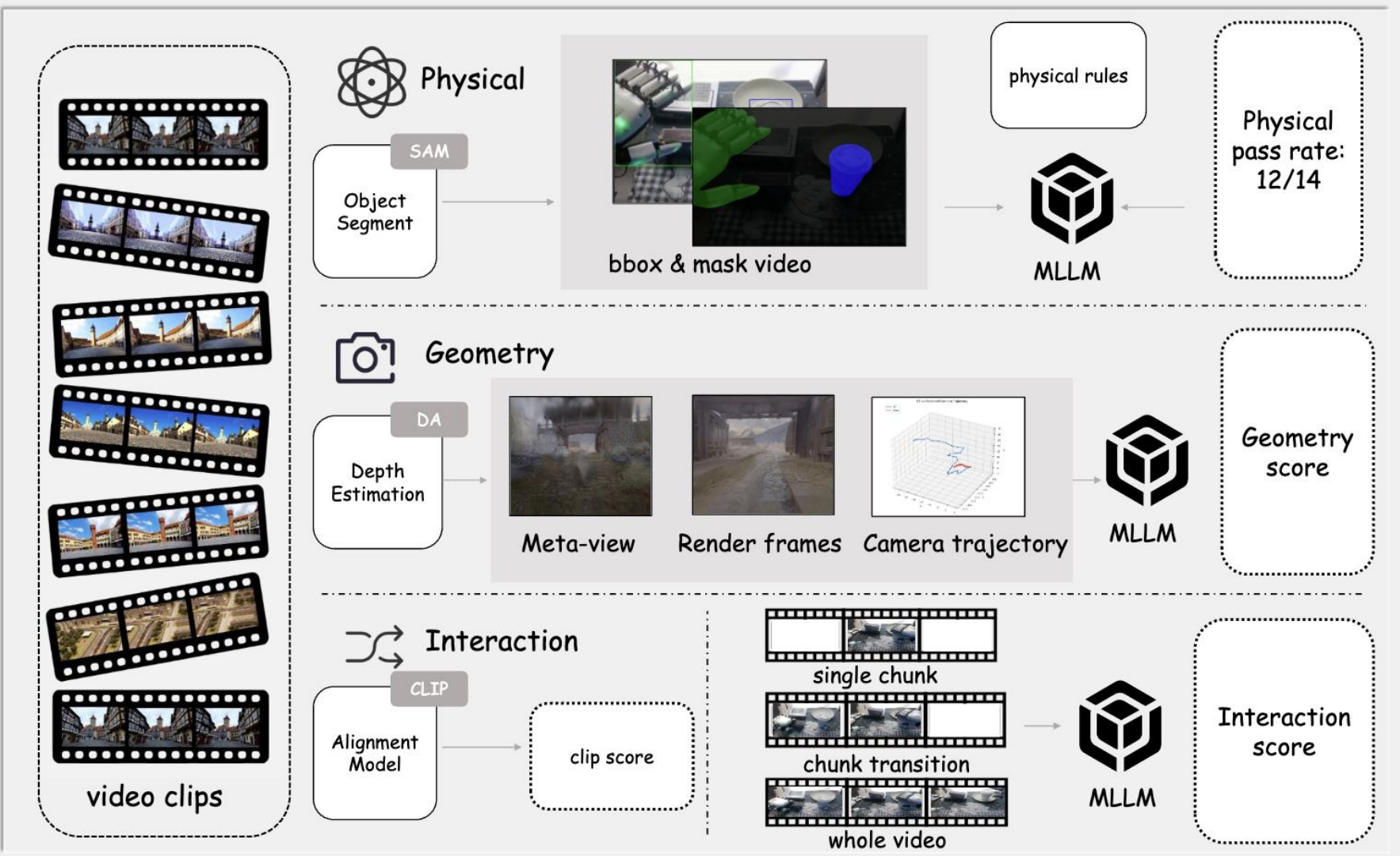
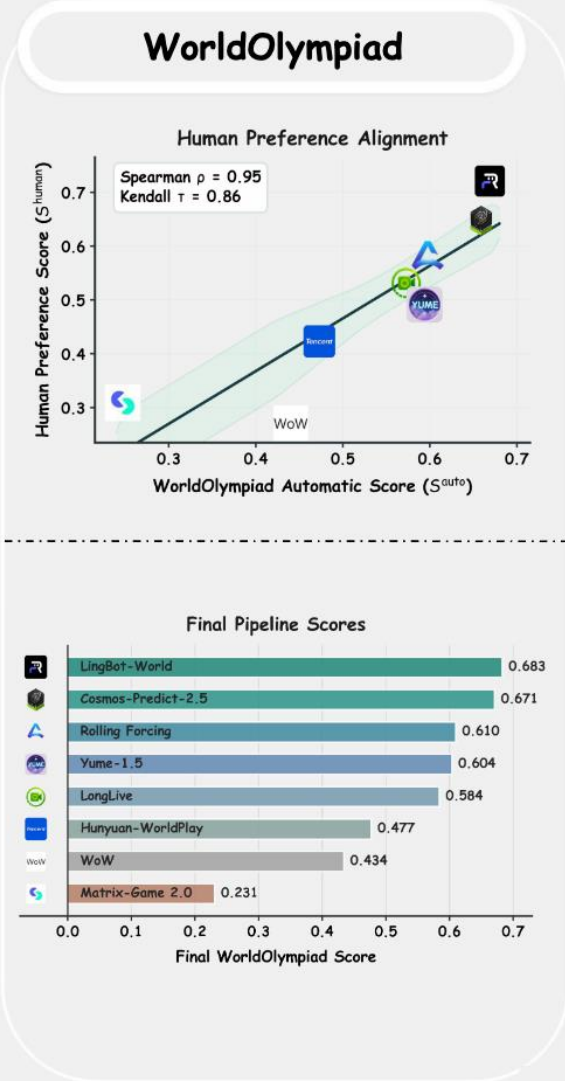


# Latent Spatial Memory

for Video World Models

# How to Benchmark these works?

# WorldOlympiad



Benchmark	Eval Metrics				Video Tasks		
	Long Video	Physical	Geometry	Interaction	Gaming	Robotics	Real-world
VBench <a href="#">[15]</a>	X	X	X	X	X	X	✓
VBench++ <a href="#">[16]</a>	✓	X	X	X	X	X	✓
VBench 2.0 <a href="#">[53]</a>	X	✓	X	X	X	X	✓
MIND <a href="#">[47]</a>	✓	X	X	✓	✓	X	X
EWMBench <a href="#">[12]</a>	X	X	X	✓	X	✓	X
WorldEval <a href="#">[20]</a>	X	X	X	✓	X	✓	X
WorldArena <a href="#">[29]</a>	X	✓	✓	✓	X	✓	X
<b>WorldOlympiad</b>	✓	✓	✓	✓	✓	✓	✓

# Summary

- Enhance the 3D capabilities of video foundation models – ***World-R1***
- Use 3D latent to represent the entire world – ***Latent Spatial Memory***
- Benchmarking the video world model – ***WorldOlympiad***

# Thank You

Towards Intelligent Interactive 3D-aware Video World Model