

# FantasyWorld: Geometry-Consistent World Modeling via Unified Video and 3D Prediction



ICLR 2026 · Poster

Yixiang Dai<sup>\*1</sup> Fan Jiang<sup>\*†‡1</sup> Chiyu Wang<sup>\*1</sup> Mu Xu<sup>1</sup> YongGang Qi<sup>‡2</sup>

<sup>1</sup>AMAP, Alibaba Group

<sup>2</sup>Beijing University of Posts and Telecommunications

# Motivation

## Video Generation Models

- Strength: Strong 2D visual priors.
- Limitation: Weak 3D consistency.

Imaginative



## Feed-Forward 3D Reconstruction Models

- Strength: Fast and High-quality reconstruction.
- Limitation: Reliant on acquiring high-quality 2D source data.

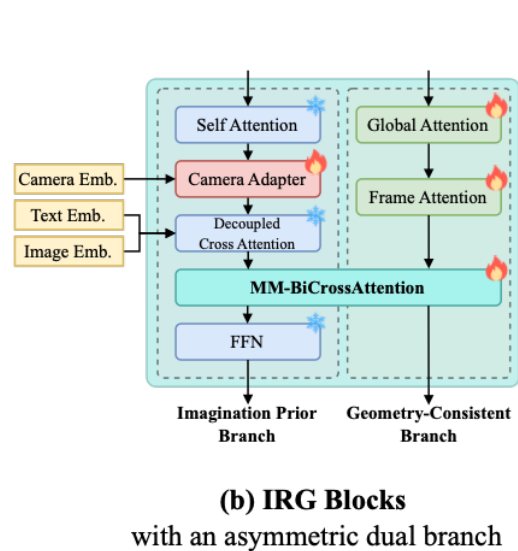
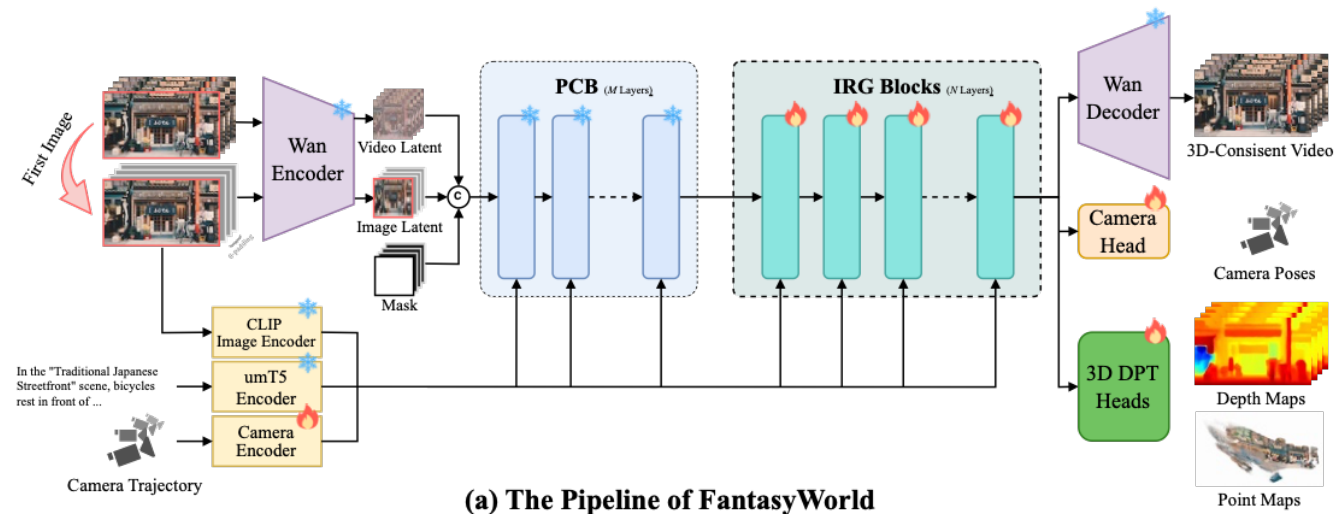
Geometric

**Our goal:** Bidirectional mutual promotion in unified generation framework.

**3D** → **2D**: Improve the geometric consistency of the 2D branch.

**2D** → **3D**: Provide visual priors for the 3D branch.

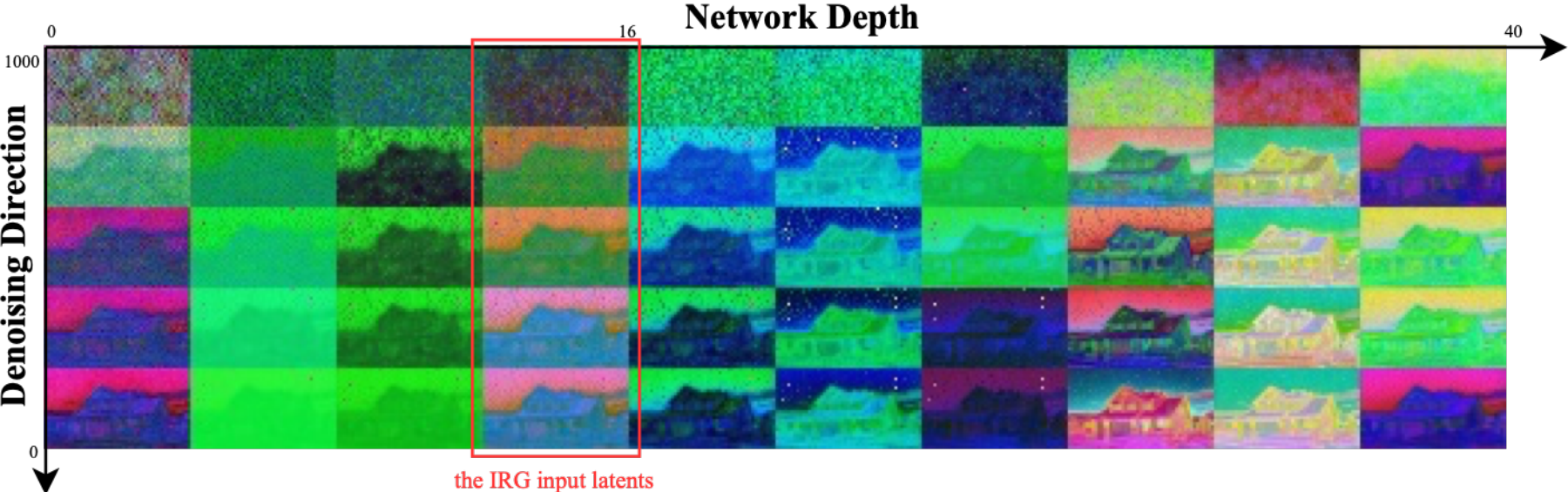
# Overview



Given image + text + camera trajectory  $\rightarrow$  video + implicit 3D field (single forward pass)

# Preconditioning Blocks (PCB)

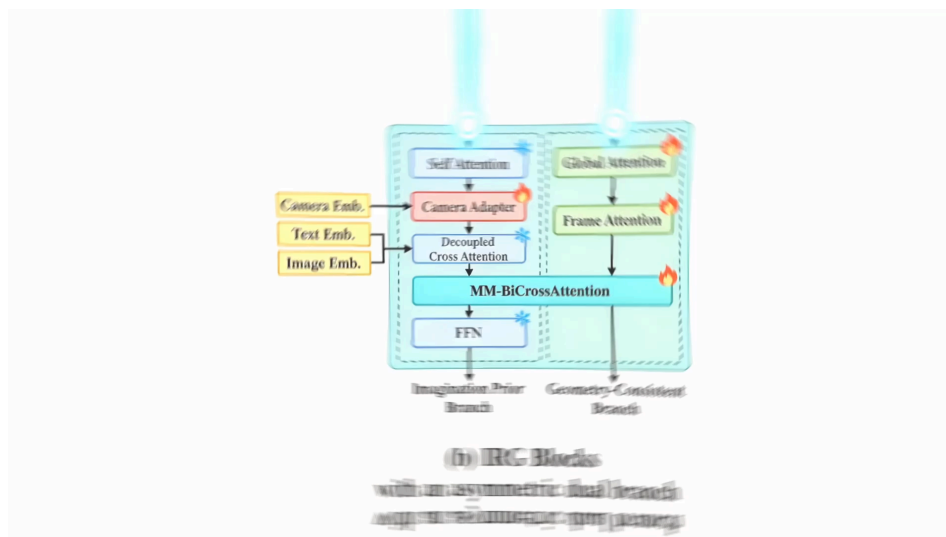
The features extracted by the PCB module ensures that the geometry branch is guided by meaningful signals rather than pure noise



# Integrated Reconstruction & Generation Blocks (IRG)

## Asymmetric dual-branch structure:

- Imaginative Branch: Frozen weights, Appearance synthesis
- Geometric Branch: Trainable, Explicit 3D reasoning
- Coupled via lightweight adapters + MM-BiCrossAttention



# Results on WorldScore Leaderboard

Method	Motion	3D Consist.	Photo Consist.	Style Consist.	Camera Ctrl.	Object Ctrl.	Content Align.	Subjective Qual.
WonderWorld	Small	82.85 $\pm$ 19.69	67.86 $\pm$ 23.56	55.79 $\pm$ 34.89	<b>92.32</b>	47.63	<b>79.09</b>	<b>69.03</b>
AETHER	Small	79.84 $\pm$ 14.68	58.68 $\pm$ 38.59	72.09 $\pm$ 32.62	57.44	52.26	28.06	41.11
Uni3C	Small	78.59 $\pm$ 21.08	85.48 $\pm$ 20.98	88.32 $\pm$ 18.47	62.94	45.83	47.40	57.00
Voyager	Small	56.00 $\pm$ 26.32	80.68 $\pm$ 16.32	72.89 $\pm$ 29.78	45.92	<b>57.69</b>	48.36	44.74
<b>Ours w/o 3D</b>	Small	79.77 $\pm$ 16.06	83.86 $\pm$ 8.73	92.54 $\pm$ 12.90	57.94	37.33	43.31	55.85
<b>Ours w/ 3D</b>	Small	<b>83.31</b> $\pm$ 14.24	<b>86.11</b> $\pm$ 7.97	<b>94.22</b> $\pm$ 9.11	57.05	34.46	38.45	57.40
WonderWorld	Large	63.70 $\pm$ 24.37	3.22 $\pm$ 8.47	35.95 $\pm$ 33.47	<b>96.28</b>	38.61	<b>97.10</b>	<b>72.46</b>
AETHER	Large	63.97 $\pm$ 17.39	33.11 $\pm$ 23.99	61.99 $\pm$ 32.24	4.43	34.78	33.69	35.09
Uni3C	Large	73.95 $\pm$ 17.55	46.78 $\pm$ 32.64	71.43 $\pm$ 29.38	8.69	34.28	77.88	51.12
Voyager	Large	13.82 $\pm$ 19.96	9.52 $\pm$ 17.17	61.34 $\pm$ 35.29	0.00	<b>49.23</b>	64.10	39.21
<b>Ours w/o 3D</b>	Large	72.06 $\pm$ 20.14	56.98 $\pm$ 23.60	81.59 $\pm$ 22.23	9.32	34.44	75.85	46.96
<b>Ours w/ 3D</b>	Large	<b>74.83</b> $\pm$ 16.31	<b>60.61</b> $\pm$ 21.39	<b>82.02</b> $\pm$ 19.56	11.24	31.96	77.20	50.46

# Qualitative Results

**Indoor**



Voyager



WonderWorld



AETHER

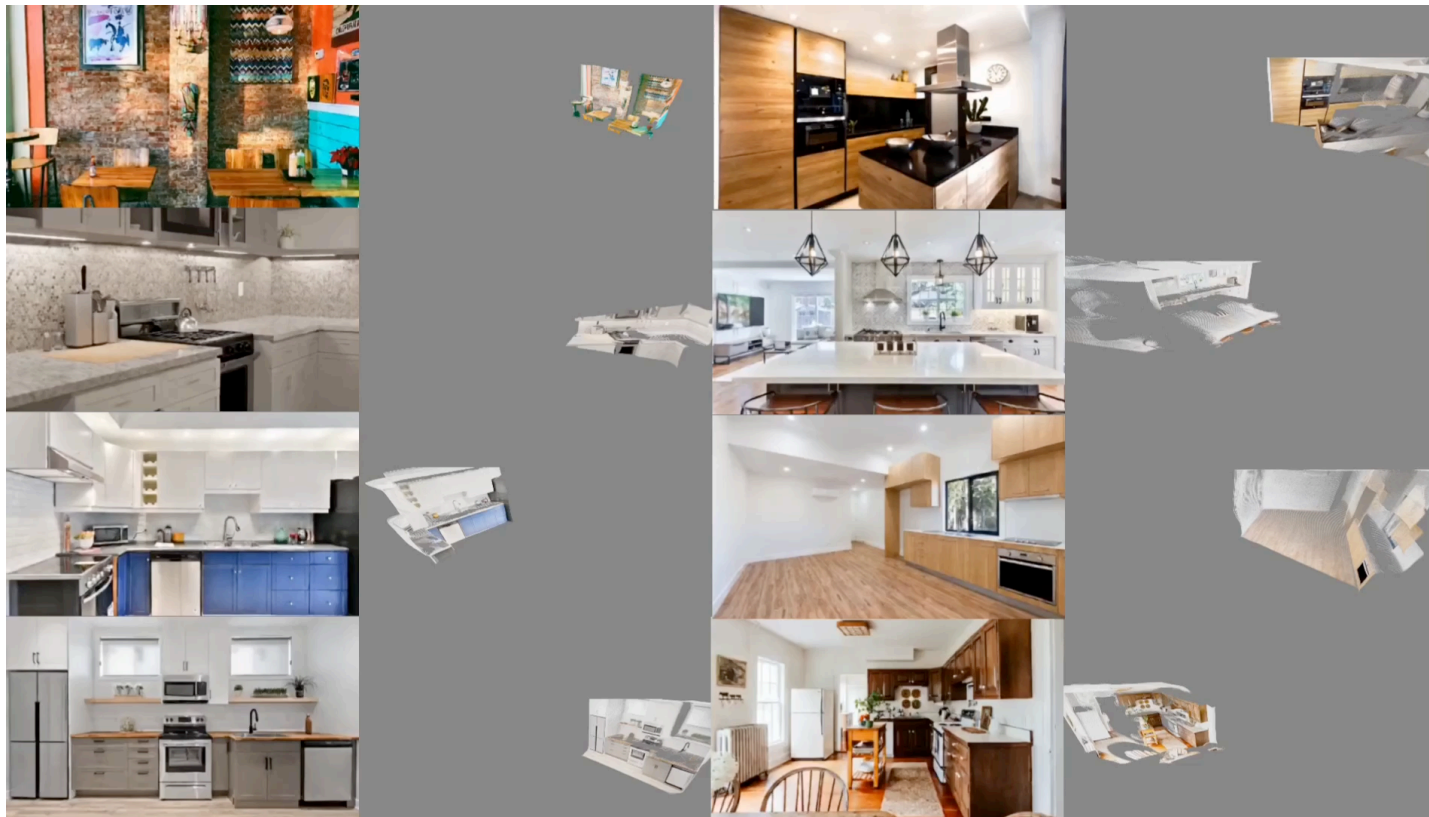


Uni3C



FantasyWorld

# More Results



# Summary

- FantasyWorld: Unified feed-forward model for joint video + 3D generation
- Cross-branch supervision between video priors and geometry cues
- Ranked 1st on WorldScore Leaderboard
- No per-scene optimization — generalizable 3D features for downstream tasks

---

Code & Models: <https://github.com/Fantasy-AMAP/fantasy-world>

Paper: [arxiv.org/abs/2509.21657](https://arxiv.org/abs/2509.21657)