

CVPR
JUNE 3-7, 2026



DENVER
COLORADO

Enable Explicit 3D/4D Controls for Pre-trained Generative Models

Yaoyao Liu

University of Illinois Urbana-Champaign

CVPR 2026 Workshop on Long-form Video Understanding, Generation, and Action

June 4, 2026



UNIVERSITY OF
ILLINOIS
URBANA-CHAMPAIGN

Humans recognize objects by understanding their 3D structure



AI models:



Input image

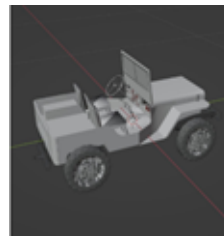


Prediction
e.g., vehicle

Humans:



Input image



3D structure



Prediction
e.g., vehicle

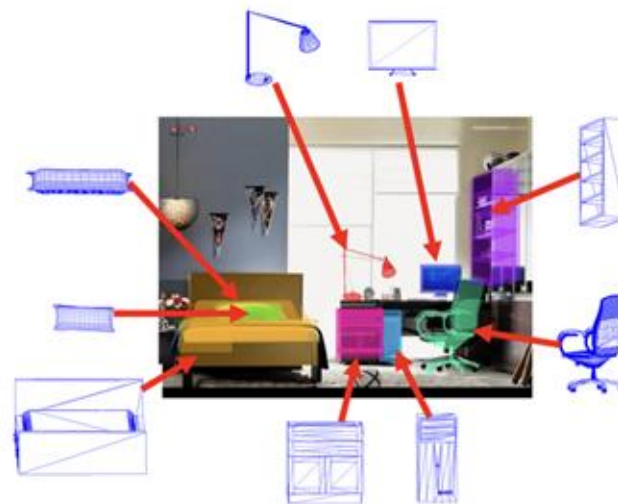
Advantages: more robust to occlusions, out-of-distribution cases, etc.

Can we learn AI models with the knowledge of the 3D world?

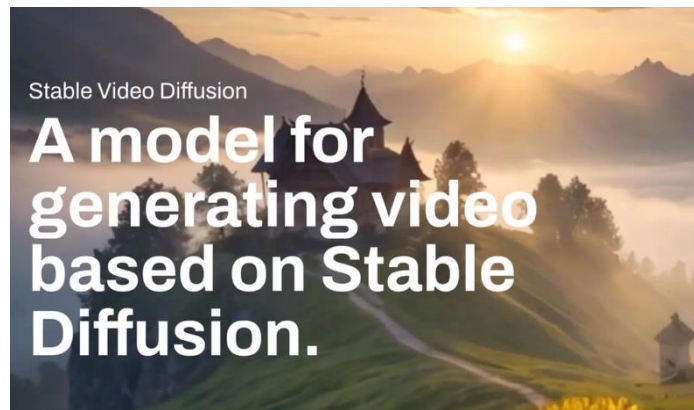


Major challenge: we don't have a large-scale dataset with 3D annotations

Dataset	# categories	# instances
PASCAL3D+ (2014) [13]	12	12,000
ObjectNet3D (2016) [24]	100	57,000
CAMERA25 (2019) [27]	6	1,000
REAL275 (2019) [27]	6	24
Objectron (2021) [28]	9	18,000
Wild6D (2022) [14]	5	2,000



**We now have powerful generation models,
but they still lack explicit 3D/4D control.**



**Can we use pre-trained generative models to
create data with 3D/4D annotations?**

Video generation with 3D (4D) controls

SIGGRAPH 2026 Conference Papers

go.illinois.edu/4dvideo



FreeOrbit4D: Training-Free Arbitrary Camera Redirection for Monocular Videos via Foreground-Complete 4D Reconstruction

WEI CAO, University of Illinois Urbana-Champaign, USA

HAO ZHANG, University of Illinois Urbana-Champaign, USA

FENGRUI TIAN, University of Pennsylvania, USA

YULUN WU, University of Illinois Urbana-Champaign, USA

YINGYING LI, University of Illinois Urbana-Champaign, USA

SHENLONG WANG, University of Illinois Urbana-Champaign, USA

NING YU, Eyeline Labs, USA and Netflix, USA

YAOYAO LIU*, University of Illinois Urbana-Champaign, USA



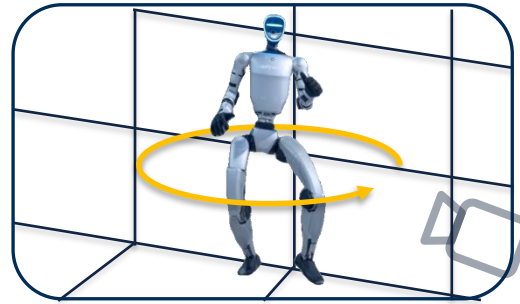
Poster session today: 5:15 pm, Exhibit Hall A, No. 303



Can we enable 4D controls for video redirection?



Monocular Video



Bullet
Time



Existing methods: large-angle redirection remains challenging

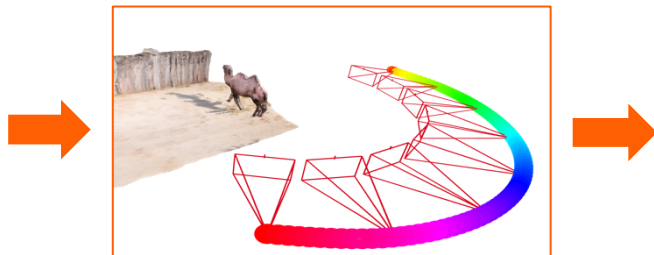
go.illinois.edu/4dvideo



Input Video



Target Trajectory



Explicit Warping
e.g., using depth maps

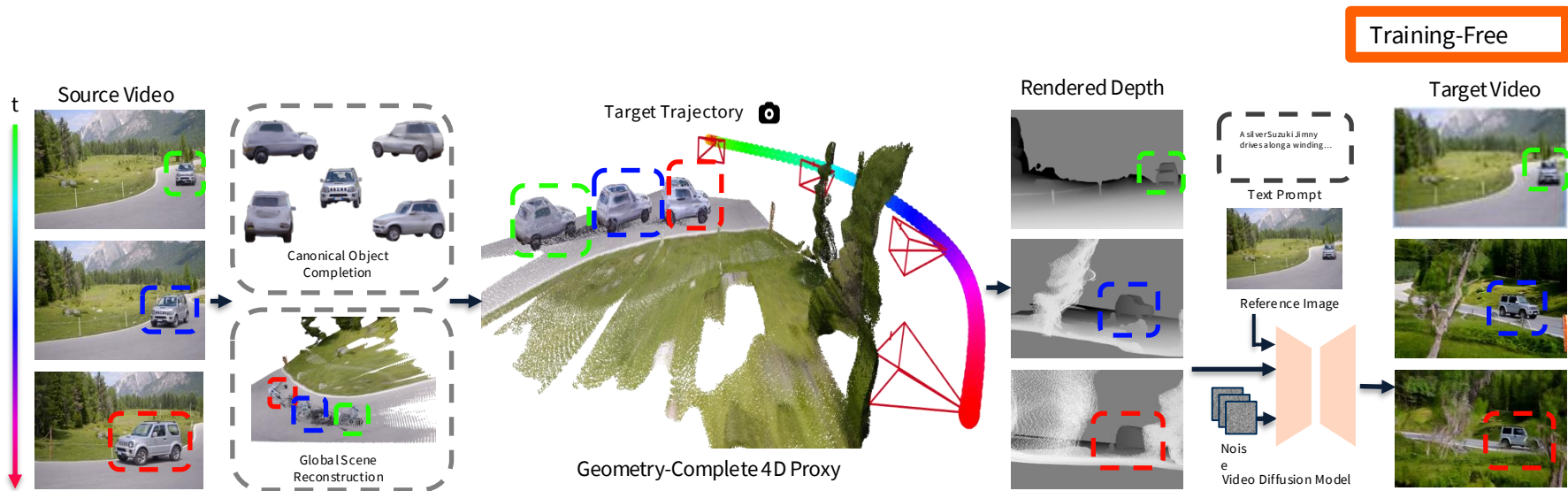


Implicit Control
e.g., text prompts



Our pipeline: geometry-complete 4D proxy for synthesis

go.illinois.edu/4dvideo



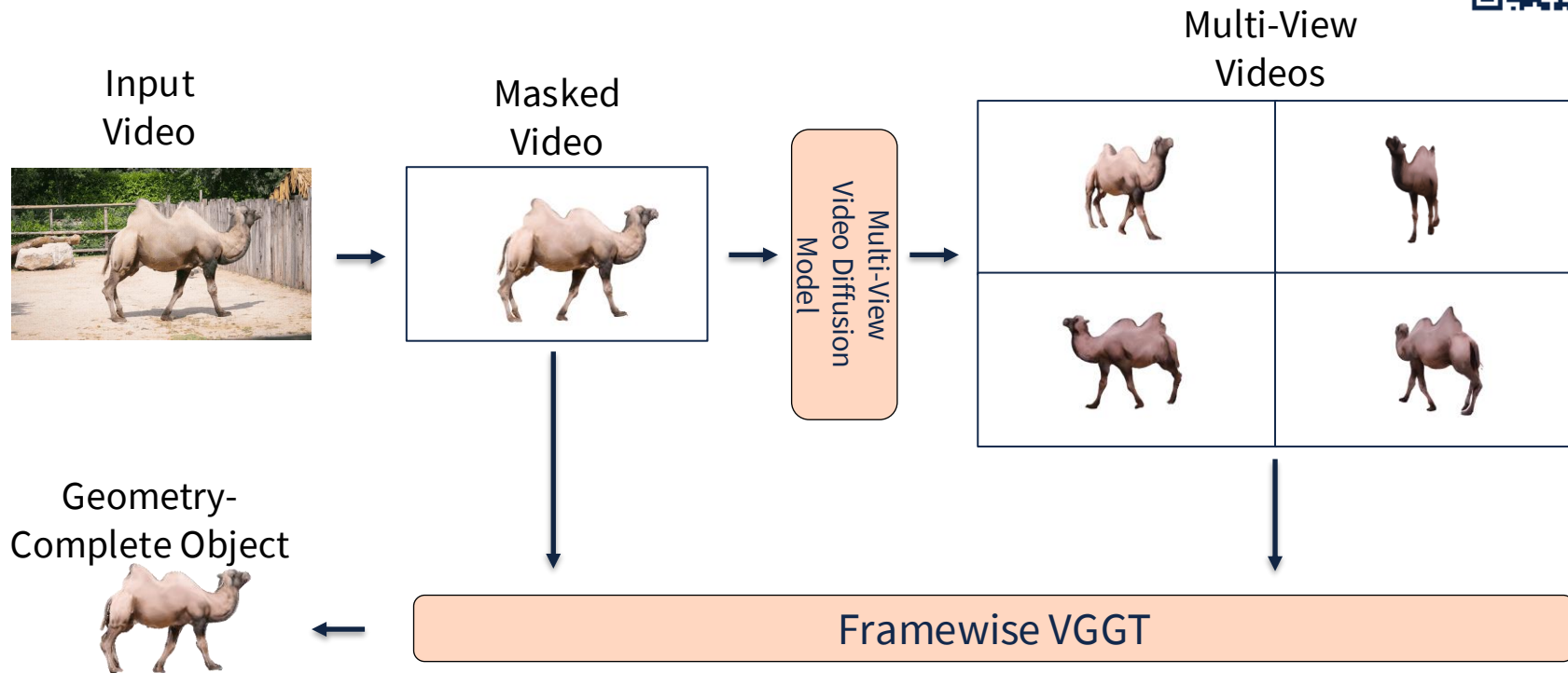
Decoupled Reconstruction

4D Proxy Alignment

Geometry-Guided Synthesis



Canonical object completion: from masked monocular video to geometry-complete object





Global scene reconstruction: recover background and visible foreground in world space

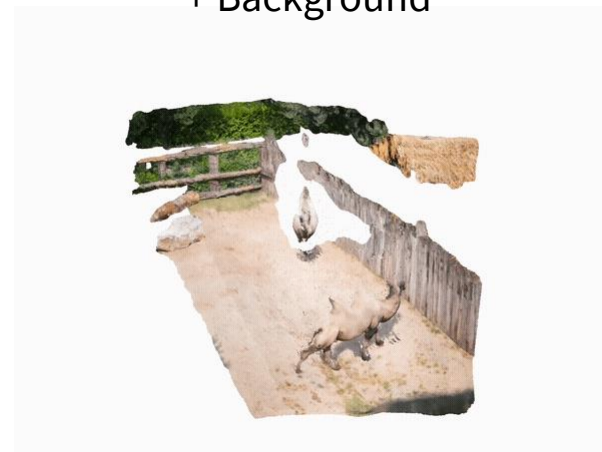
Input
Video



Dynamic-aware
Feed-Forward Network

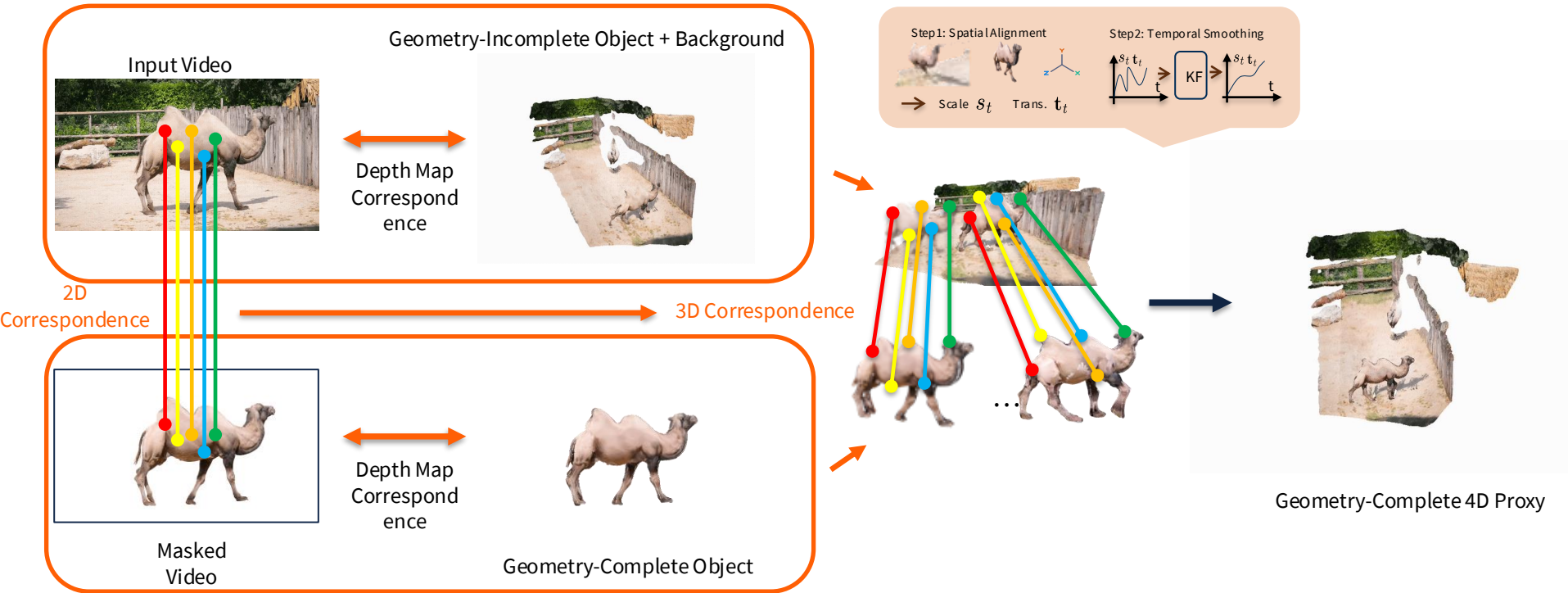


Geometry-Incomplete Object
+ Background





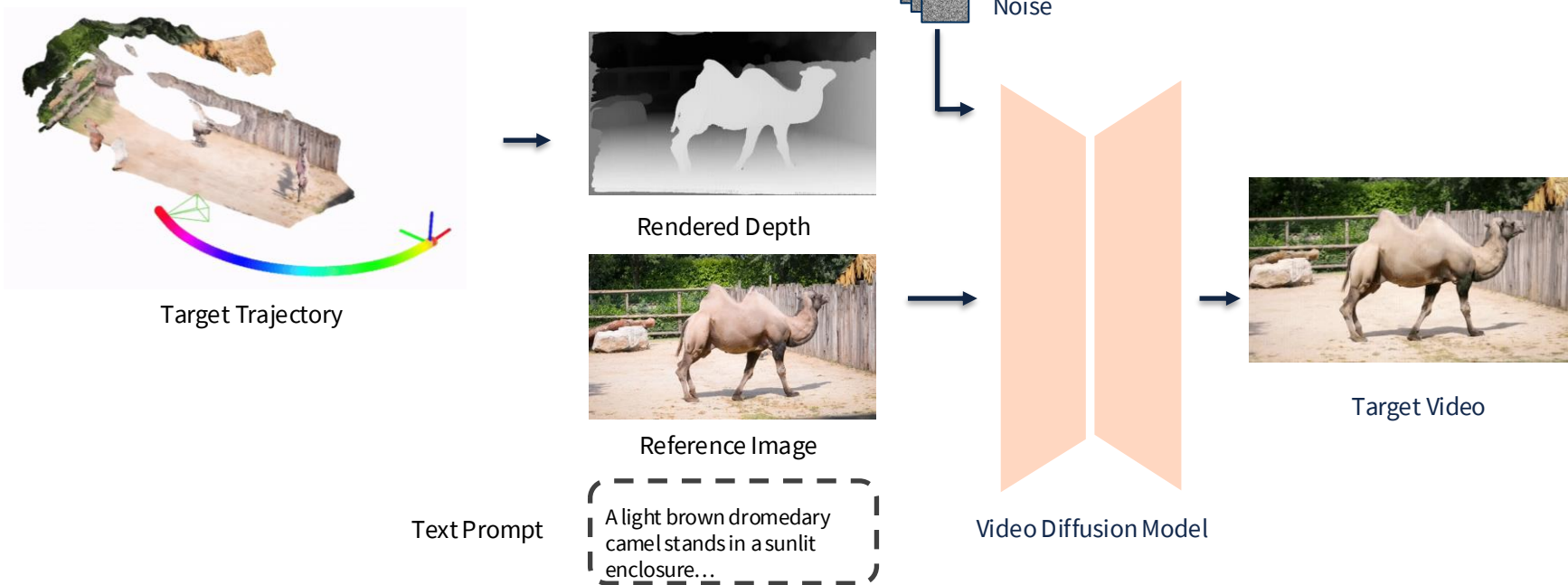
3D correspondence-aware alignment: complete geometry + correct scene placement → unified 4D proxy





Geometry-guided video synthesis: rendered depth from the 4D proxy guides target-view generation

Geometry-conditioned Video Synthesis



Fast motion and thin structures



Input Video



ReCamMaster



TrajectoryCrafter



Ours



EX4D



Gen3C



Large camera displacement



Input Video



ReCamMaster



TrajectoryCrafter



Ours



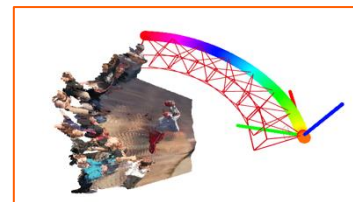
EX4D



Gen3C



Complex human motion



Input Video



ReCamMaster



TrajectoryCrafter



Ours



EX4D



Gen3C



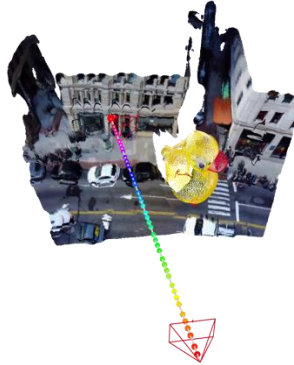


Flexible camera control

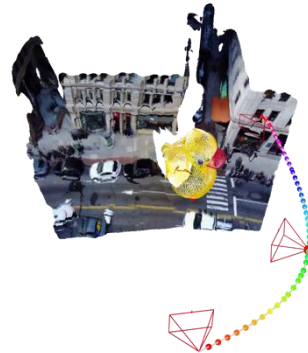
Input Video



Target Trajectory



Target Video



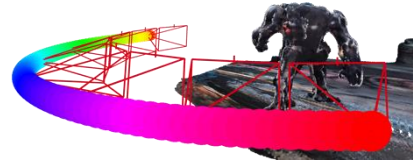


Flexible camera control

Input Video

Target Trajectory

Target Video



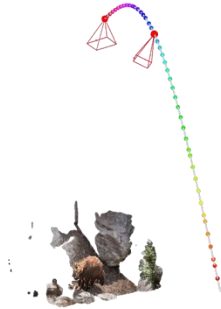
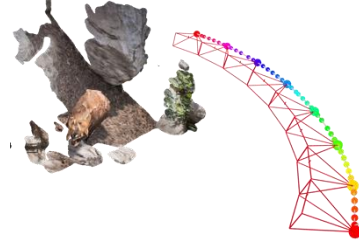


Flexible camera control

Input Video

Target Trajectory

Target Video

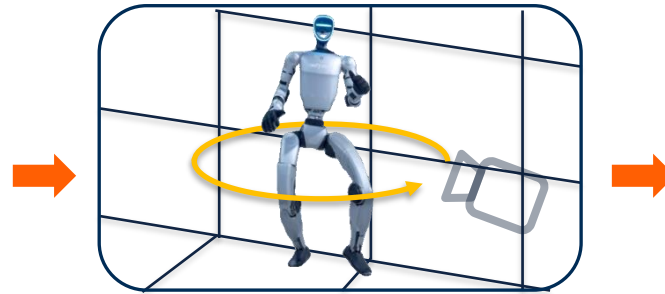




Back to our motivating example



Monocular Video



Bullet Time



Target Video



Consistency across random seeds

Input Video



Target Video

Seed #1



Seed #2





Appearance propagation from a single edited frame

Input Video



Reference Image



Target Video





Direct geometry editing in explicit 4D space

Input Video



Reference Image



Target Video





Quantitative results

Table 1. **Quantitative comparison and User Study.** We report VBench for perceptual video quality, DINO/CLIP-SIM for semantic similarity, FID-V/FVD-V for distributional fidelity, and user ratings (1–5 scale). **Bold**: best; underline: second-best.

Method	VBench \uparrow						Similarity & Fidelity				User Study		
	Subject Consis.	BG Consis.	Motion Smooth.	Overall Consis.	Aesth. Qual.	Imaging Qual.	DINO-SIM (\uparrow)	CLIP-SIM (\uparrow)	FID-V ($\downarrow \times 10^2$)	FVD-V ($\downarrow \times 10^3$)	Overall (\uparrow)	Motion (\uparrow)	Stab. (\uparrow)
ReCamMaster	<u>0.84</u>	<u>0.92</u>	0.98	0.16	0.39	43	0.37	0.75	2.6	3.9	2.0	2.5	2.0
TrajectoryCrafter	0.80	0.91	0.94	<u>0.19</u>	<u>0.47</u>	<u>53</u>	<u>0.47</u>	<u>0.79</u>	<u>2.0</u>	<u>3.6</u>	<u>2.8</u>	<u>3.2</u>	<u>2.9</u>
EX-4D	0.76	0.89	0.94	0.16	0.42	46	0.28	0.69	3.2	3.8	2.0	2.5	2.0
GEN3C	0.79	0.88	0.95	0.18	0.42	49	0.43	0.75	2.3	3.3	2.4	<u>3.5</u>	2.3
Ours	0.88	0.94	<u>0.96</u>	0.24	0.52	64	0.65	0.84	1.7	<u>3.6</u>	4.6	4.5	4.5

Best in 5/6 VBench metrics $\frac{3}{4}$ Similarity & Fidelity metrics and all user-study axes



Extension: we can create a video dataset with 4D annotations

Redirect4D-Bench: A Scalable Benchmark for Camera Redirection of Monocular Dynamic Videos with Pseudo-4D Ground Truth




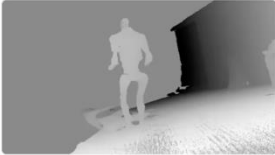








Wei Cao¹, Hao Zhang¹, Jiapeng Tang², Yulun Wu¹,
Yingying Li¹, Ning Yu³, Shenlong Wang¹, Yaoyao Liu¹

¹ University of Illinois Urbana-Champaign

² Technical University of Munich ³ Netflix



Extension: we can create a video dataset with 4D annotations

 <p>Source RGB - Robot #2</p>	 <p>4D reconstruction (interactive)</p>	 <p>Target pseudo-GT mask</p>	 <p>Target rendered depth</p>
 <p>Source RGB - Tiger</p>	 <p>4D reconstruction (interactive)</p>	 <p>Target pseudo-GT mask</p>	 <p>Target rendered depth</p>
 <p>Source RGB - Cow</p>	 <p>4D reconstruction (interactive)</p>	 <p>Target pseudo-GT mask</p>	 <p>Target rendered depth</p>

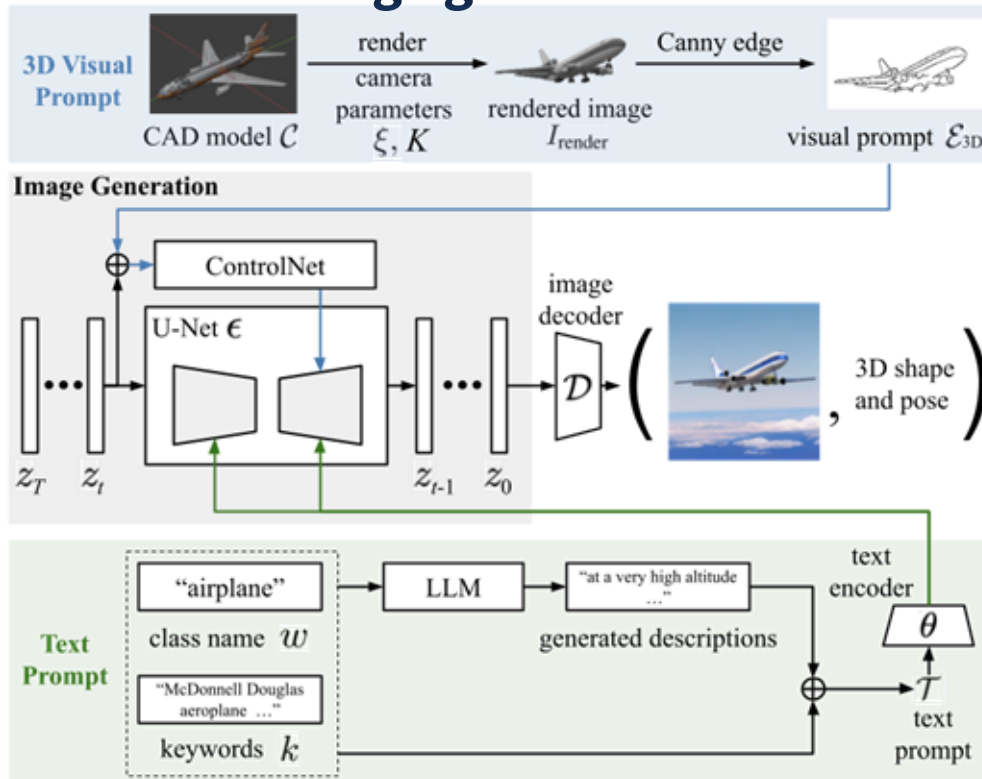


Extension: we can create a video dataset with 4D annotations

<p>Source RGB - Robot #3</p>	<p>3.3</p>	<p>Target pseudo-GT mask</p>	<p>Target rendered depth</p>
<p>4D reconstruction (interactive)</p>			
<p>Source RGB - Robot #4</p>	<p>3.5</p>	<p>Target pseudo-GT mask</p>	<p>Target rendered depth</p>
<p>4D reconstruction (interactive)</p>			
<p>Source RGB - Robot #5</p>	<p>2.2</p>	<p>Target pseudo-GT mask</p>	<p>Target rendered depth</p>



For image generation: 3D diffusion style transfer (3D-DST) enable 3D controls for image generation





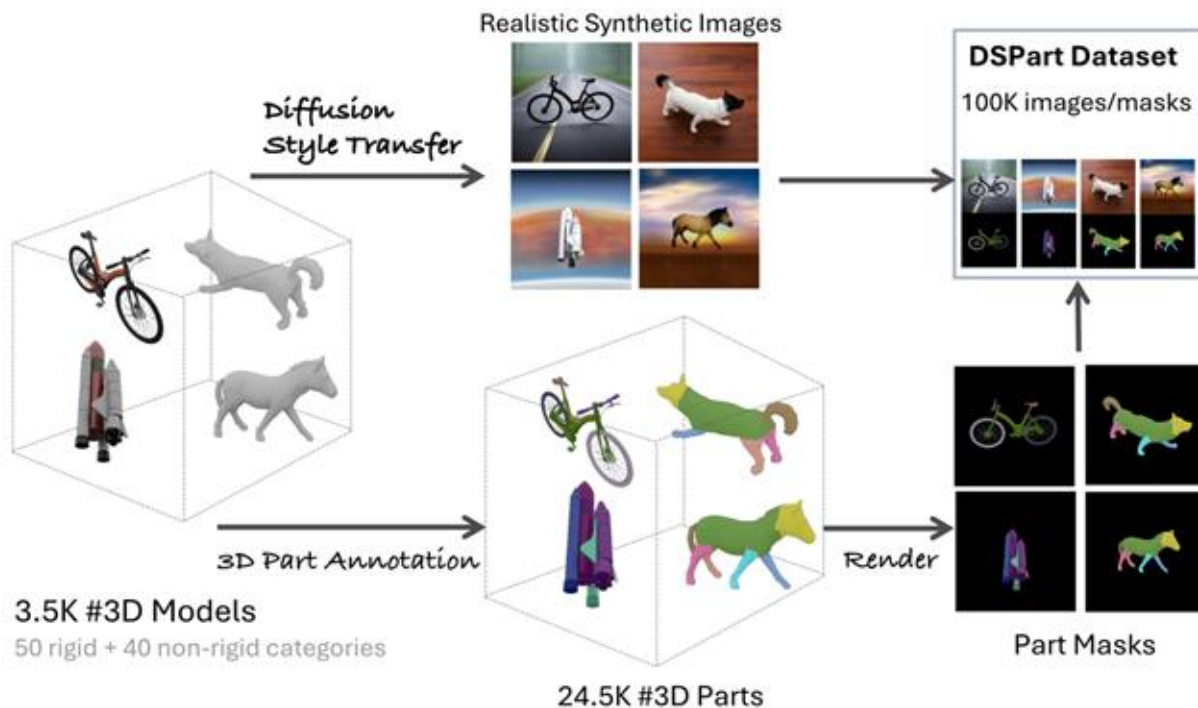
For image generation: 3D diffusion style transfer (3D-DST) enable 3D controls for image generation



CAD Model

Generated Images

Extension: generating 3D annotations based on articulated models





Enable 3D/4D controls for pre-trained generative models

Videos with 4D controls
(SIGGRAPH 2026 Conference)

go.illinois.edu/4dvideo

Videos with 4D annotations
(Under review)

go.illinois.edu/4dbenchmark

Images with 3D controls
(ICLR 2023 Spotlight)

go.illinois.edu/3dimage

- Cao, W., ..., **Liu, Y.** *FreeOrbit4D: Training-free Arbitrary Camera Redirection for Monocular Videos via Foreground-Complete 4D Reconstruction*. SIGGRAPH 2026.
- Cao, W., ..., **Liu, Y.** *Redirect4D-Bench: A Scalable Benchmark for Camera Redirection of Monocular Dynamic Videos with Pseudo-4D Ground Truth*. Under review.
- Ma, W., ..., **Liu, Y.* (corresponding author)**, Yuille, A. *Generating Images with 3D Annotations Using Diffusion Models*. ICLR 2023.

Thanks!
Any questions?

lyy@illinois.edu
<https://yaoyaoliu.web.illinois.edu/>

