



3D Vision and  
Robotics Lab

# [CVPR 2024] Generative Rendering: Controllable 4D-Guided Video Generation with 2D Diffusion Models

Shengqu Cai<sup>1,2</sup>, Duygu Ceylan<sup>\*2</sup>, Matheus Gadelha<sup>\*2</sup>, Chun-Hao Huang<sup>2</sup>, Tuanfeng Wang<sup>2</sup>, Gordon Wetzstein<sup>1</sup>

<sup>1</sup>Stanford University, <sup>2</sup>Adobe Research

\*equal contribution

Gyeongsu Cho

@UNIST

Lab Seminar 2024.05.09. (Thu)

# Contents

---

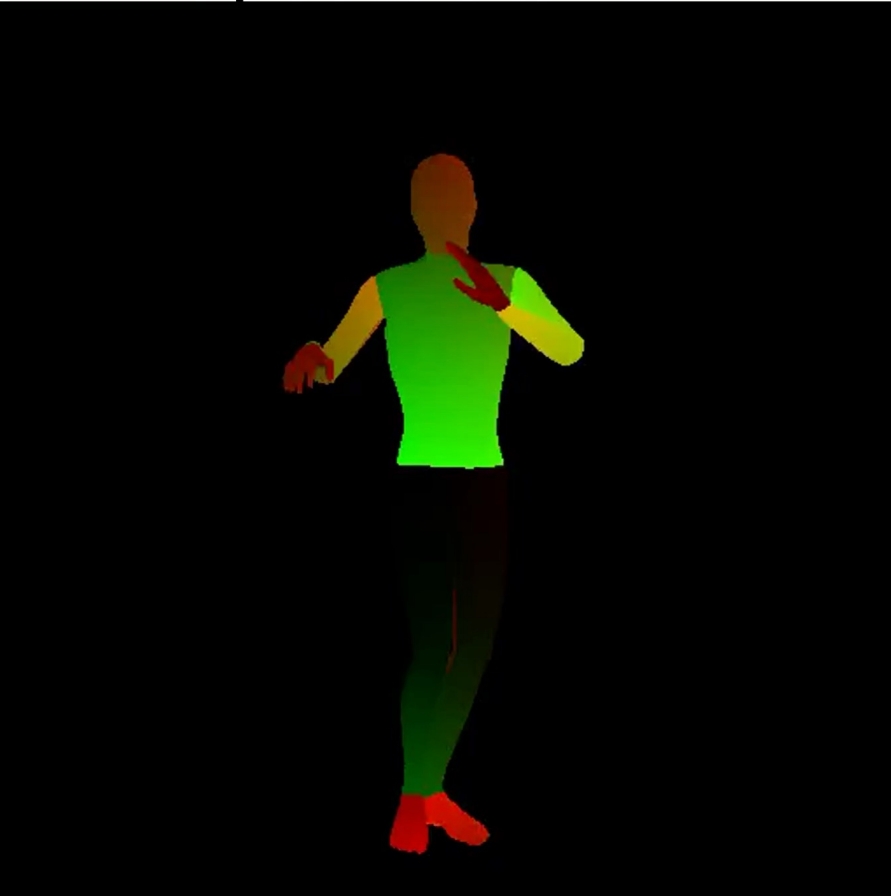
- Introduction
- Method
- Experiments
- Conclusion

# Introduction

---

Controllable video generation?

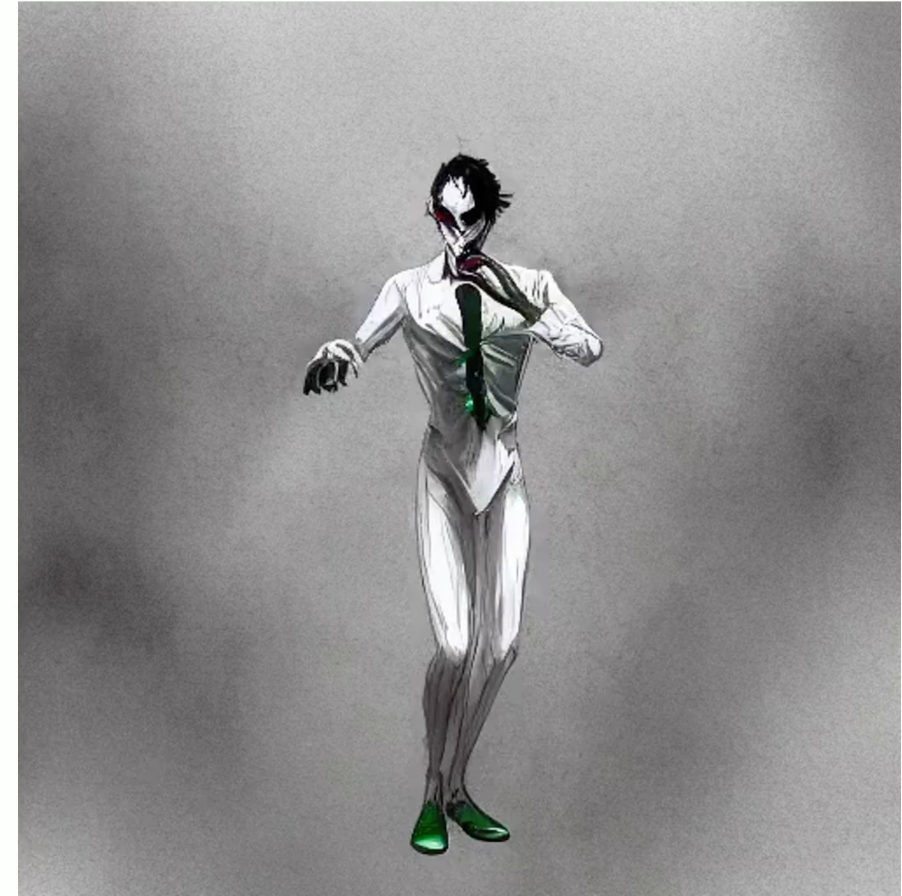
**input animated mesh**



**rendering prompt**

a manga of  
the Joker dancing

**rendered video**



[CVPR 2024] Generative Rendering

# Introduction

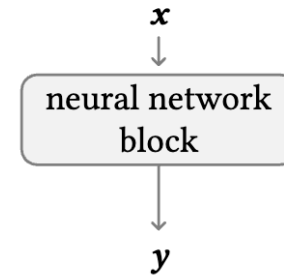
## Video Generation



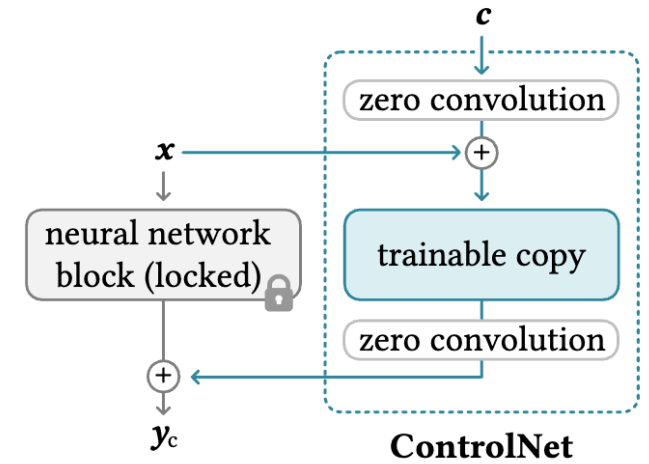
OpenAI Sora

# Introduction

## Controllable image generation



(a) Before



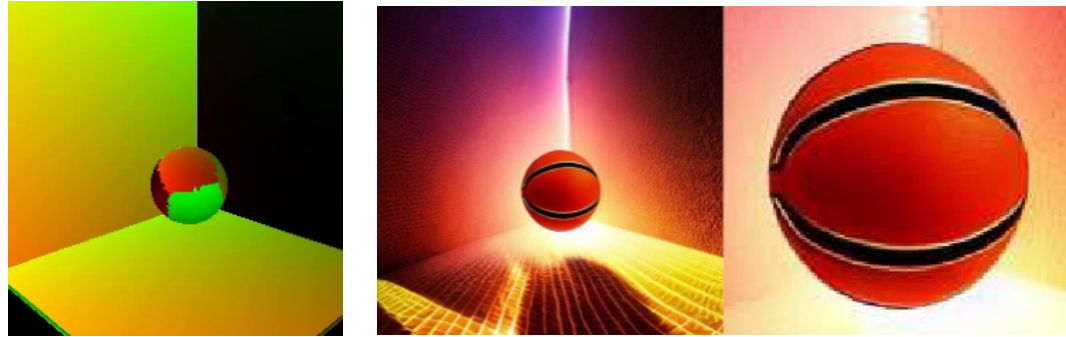
(b) After

[ICCV 2023] ControlNet

# Introduction

## Controllable image generation

Frame N



Frame M



Input

output

[ICCV 2023] ControlNet

ControlNet doesn't have consistency between frames.

# Introduction

## Video-to-video diffusion model



[ICCV 2023] Pix2Video



[ICLR 2024] TokenFlow

They require a high-fidelity video as input, which is not always available!

# Introduction

## Controllable Image-to-Video Synthesis



[arXiv 2023] Animate Anyone

We need a lot of training data to train Animate Anyone!

# Introduction

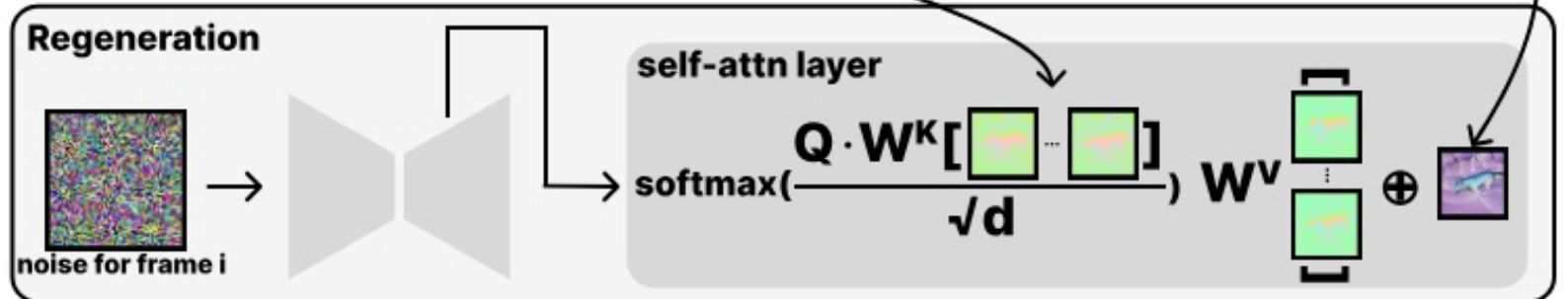
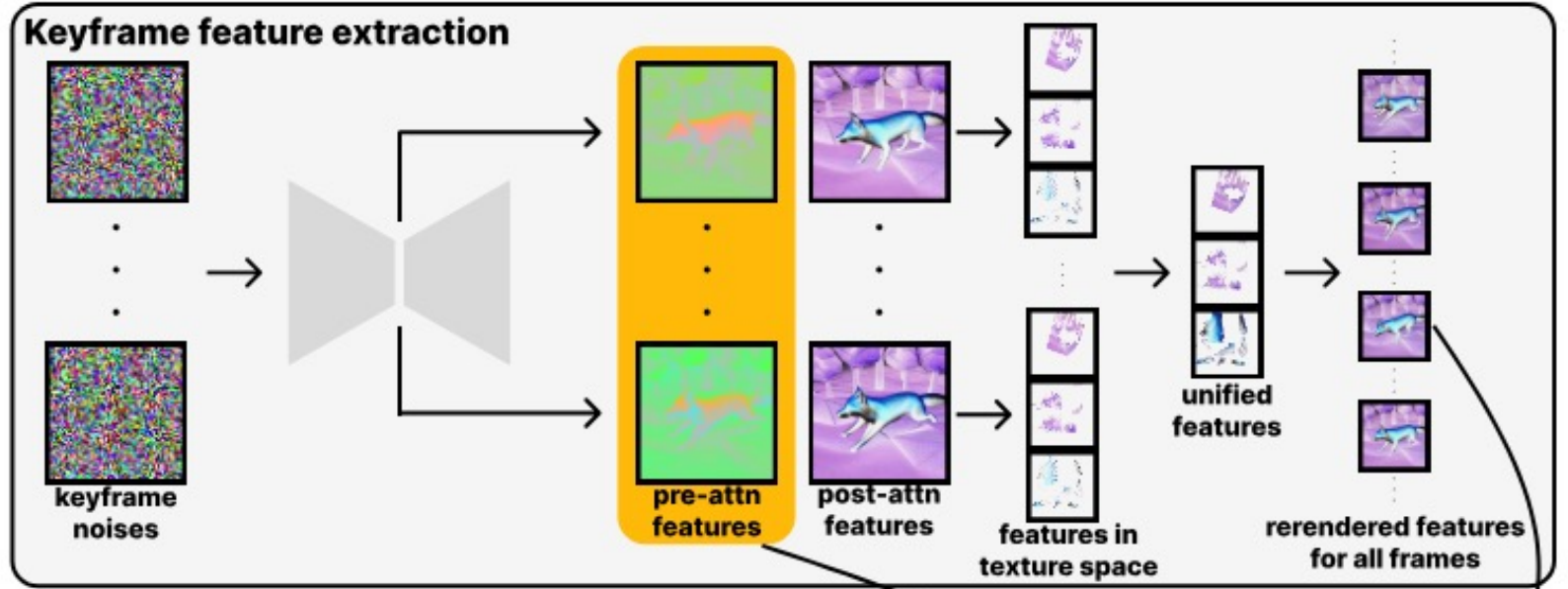
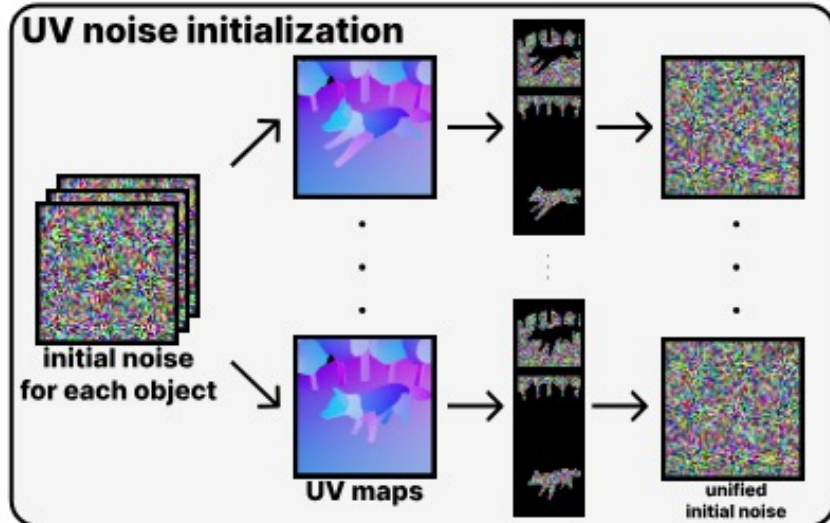
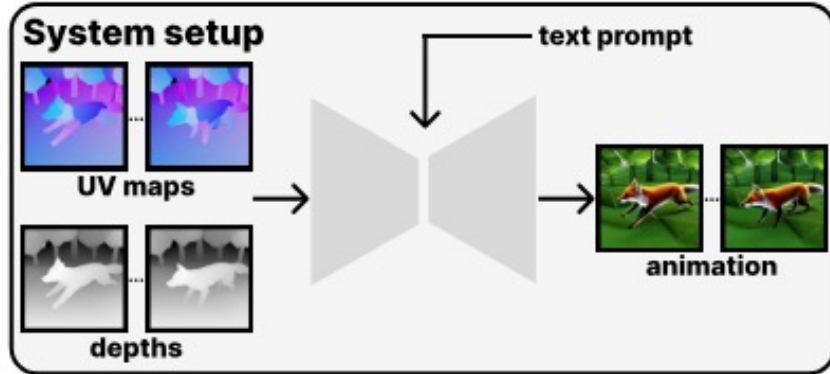
---

## Contributions

- They present a framework for 4D-guided animation synthesis utilizing the pre-trained T2I generation models as a multi-frame renderer.
- They enhance the self-attention layers of the image generation model by performing correspondence-aware blending of both input and output features to enforce consistent appearance synthesis.
- They introduce a UV-space noise initialization mechanism.

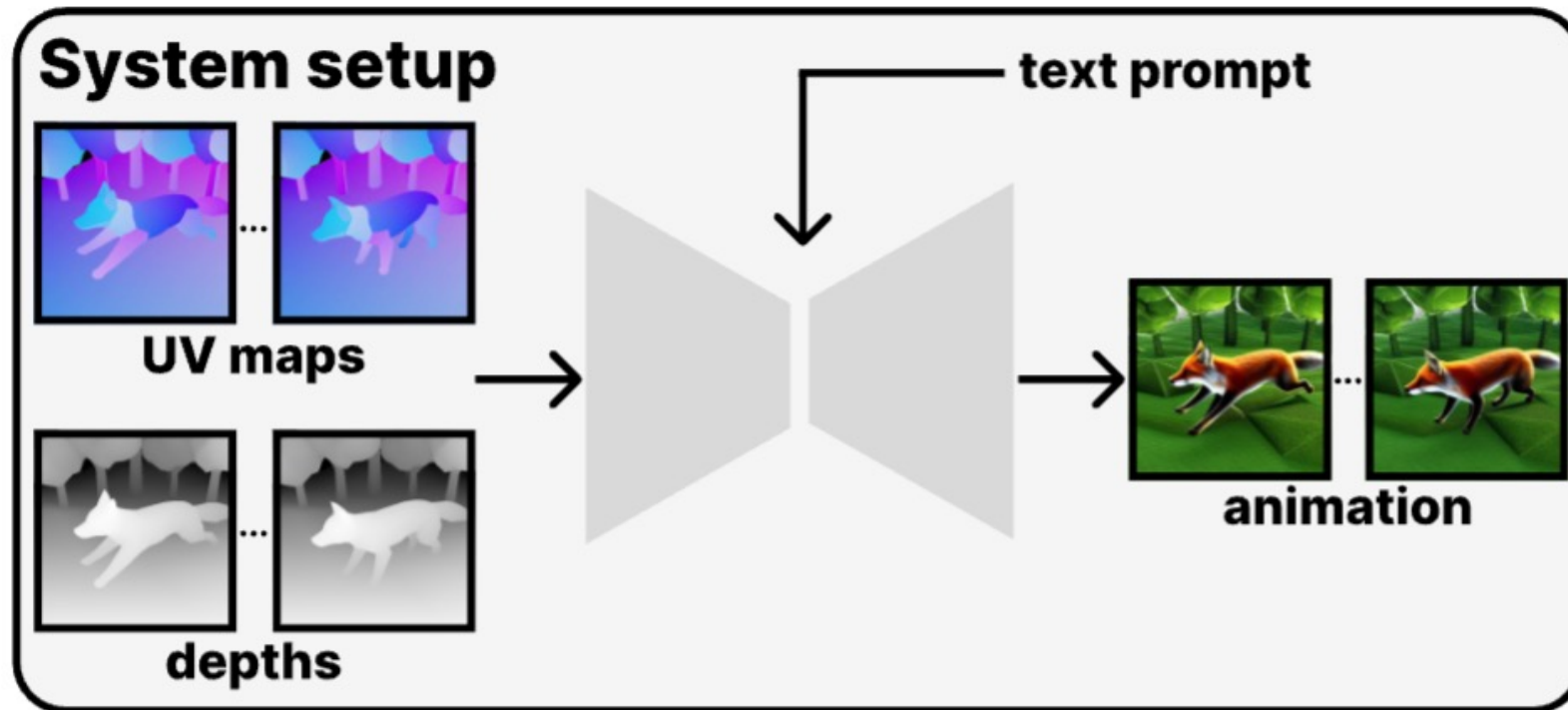
# Method

## Overview



# Method

## System setup



Input: ground truth UV maps, ground truth depth maps, text prompt

Output: animation video

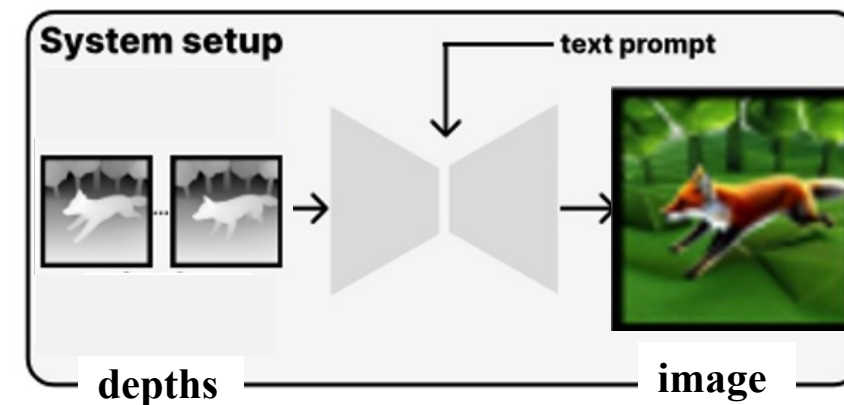
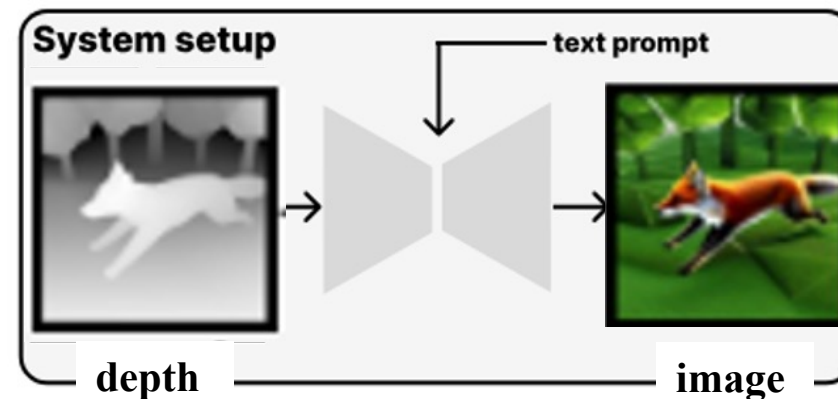
# Method

## 1. Self-attention feature injection

Naïve self-attention

$$\mathbf{F} = \text{Attn}(\mathbf{Q}; \mathbf{K}; \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right) \cdot \mathbf{V}, \quad (1)$$

$$\mathbf{Q}^{(i,l)} = \mathbf{W}^Q \mathbf{f}^{(i,l)}, \quad \mathbf{K}^{(i,l)} = \mathbf{W}^K \mathbf{f}^{(i,l)}, \quad \mathbf{V}^{(i,l)} = \mathbf{W}^V \mathbf{f}^{(i,l)} \quad (2)$$



Self-attention feature injection

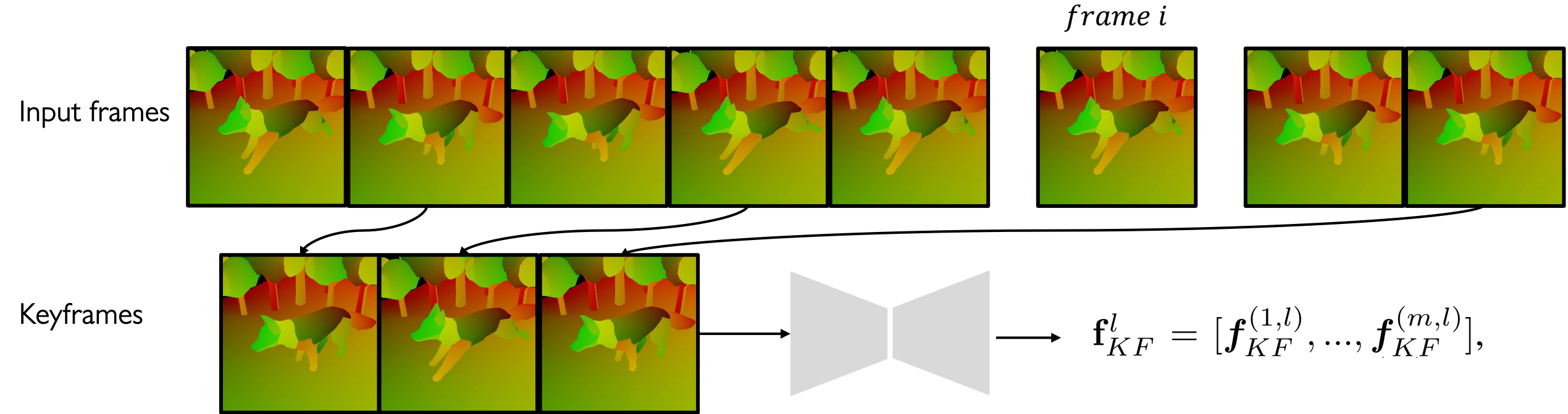
$$\mathbf{K}^{(i,l)} = \mathbf{W}^K [\mathbf{f}^{(1,l)}, \dots, \mathbf{f}^{(N,l)}], \quad \text{pre-attn} \quad (3)$$

$$\mathbf{V}^{(i,l)} = \mathbf{W}^V [\mathbf{f}^{(1,l)}, \dots, \mathbf{f}^{(N,l)}].$$

$$\mathbf{F}^{(j,l)} = \pi_{i,j}(\mathbf{F}^{(i,l)}), \quad \text{post-attn} \quad (4)$$

# Method

## 2. UV-space feature injection



$$\begin{aligned} \mathbf{K}^{(i,l)} &= \mathbf{W}^K [\mathbf{f}^{(1,l)}, \dots, \mathbf{f}^{(N,l)}], \\ \mathbf{V}^{(i,l)} &= \mathbf{W}^V [\mathbf{f}^{(1,l)}, \dots, \mathbf{f}^{(N,l)}]. \end{aligned} \quad (3)$$

$$\mathbf{F}^{(j,l)} = \pi_{i,j}(\mathbf{F}^{(i,l)}), \quad (4)$$

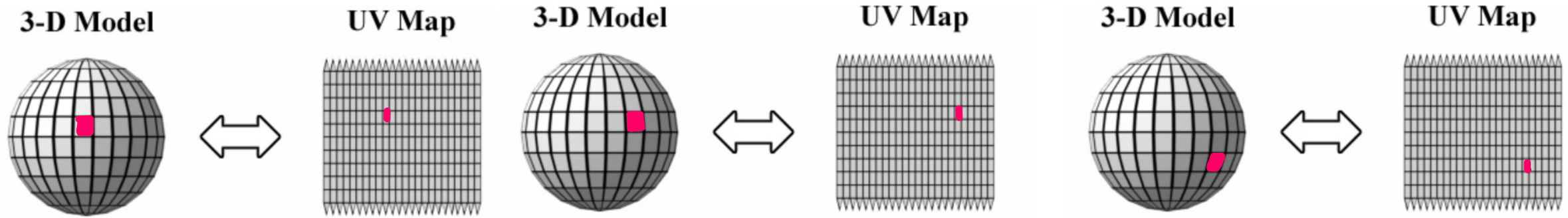
$$\begin{aligned} \mathbf{K}^{(i,l)} &= \mathbf{W}^K [\mathbf{f}^{(i,l)}, \mathbf{f}_{KF}^l], \\ \mathbf{V}^{(i,l)} &= \mathbf{W}^V [\mathbf{f}^{(i,l)}, \mathbf{f}_{KF}^l]. \end{aligned} \quad (5)$$

$$\mathbf{F}_{out}^{(i,l)} = \alpha \cdot \bar{\mathbf{F}}^{(i,l)} + (1 - \alpha) \cdot \hat{\mathbf{F}}^{(i,l)} \quad (6)$$

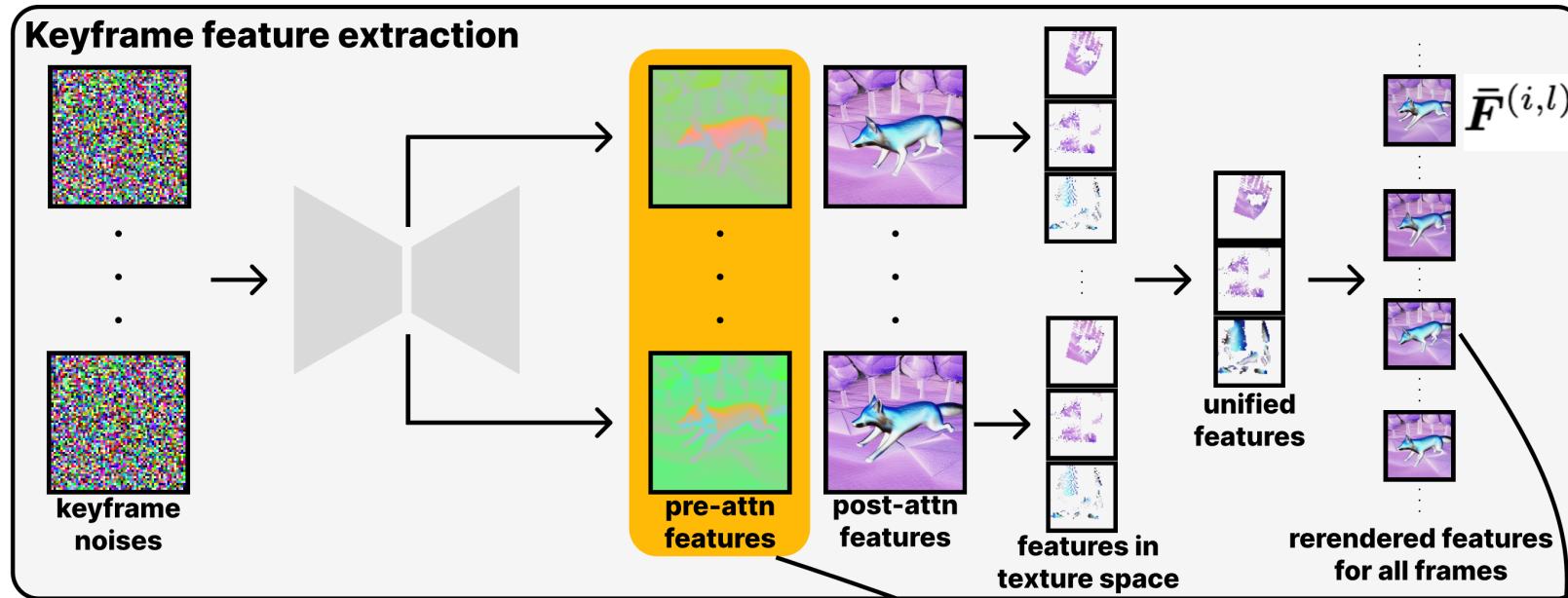
# Method

## 2. UV-space feature injection

How to inject the UV-space feature?

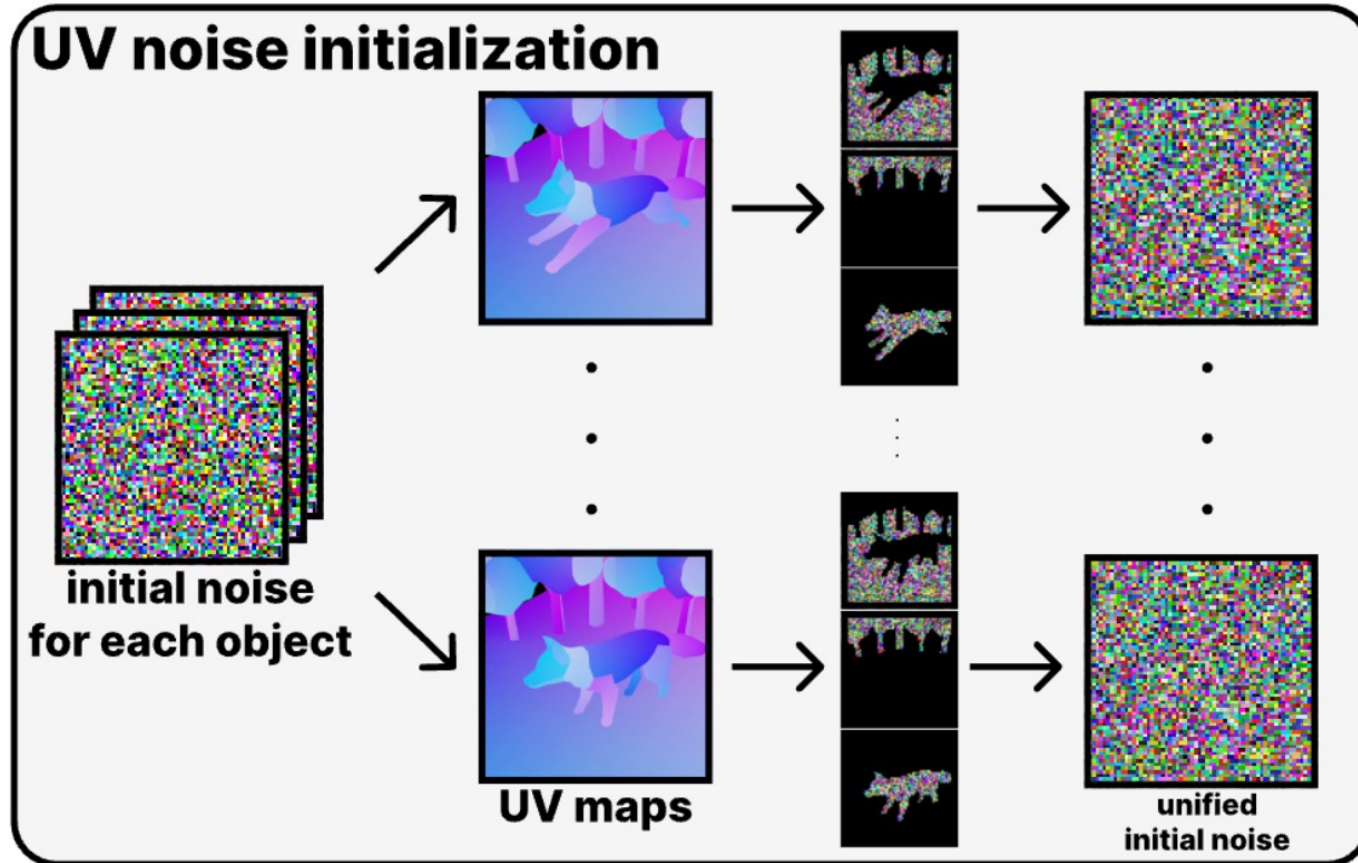


We can get UV map correspondences between frames from 3D assets!



# Method

## 3. UV noise initialization



- Initialize the noise for each object by creating a Gaussian noise texture
- Project this noise to each frame by utilizing frame-UV correspondences

# Experiments

---

- Dataset
  1. Camera rotations
  2. Physical simulation
  3. Character animation
- Baseline
  1. [ICCV 2023] Pix2Video
  2. [ICLR 2024] TokenFlow
- Evaluation metrics
  1. Frame consistency: CLIP cosine similarity between all pairs
  2. Prompt fidelity: mean CLIP embedding similarity score between video and prompt

# Experiments

## Object Rotation



input



a black sneaker



a blue sneaker



input



a forest house



a cyberpunk house



input



in a sunny forest



in an antique room

## Camera Rotation (constant background)



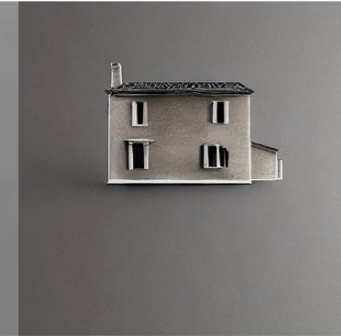
a black sneaker



a blue sneaker



a mordern house



an antique house



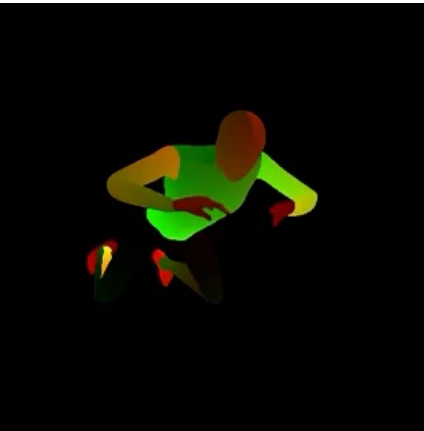
a clean bed



a pink bed

Camera rotations

# Experiments



animated mesh



a stormtrooper swimming



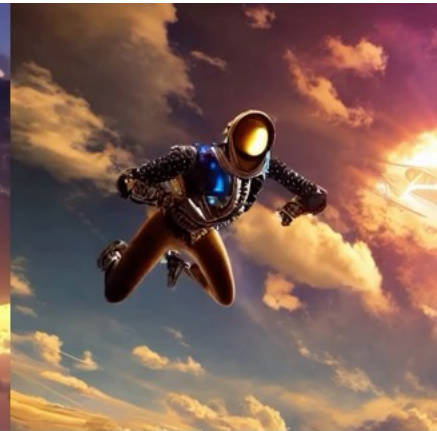
a stormtrooper swimming



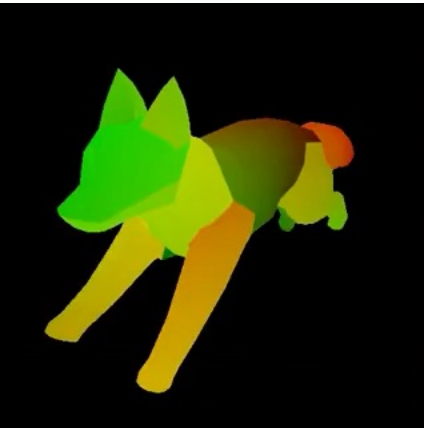
a stormtrooper swimming



machinery swimming



machinery swimming



animated mesh



a fox running in the forest



a fox running under warm sunlight



dreamlike fox running



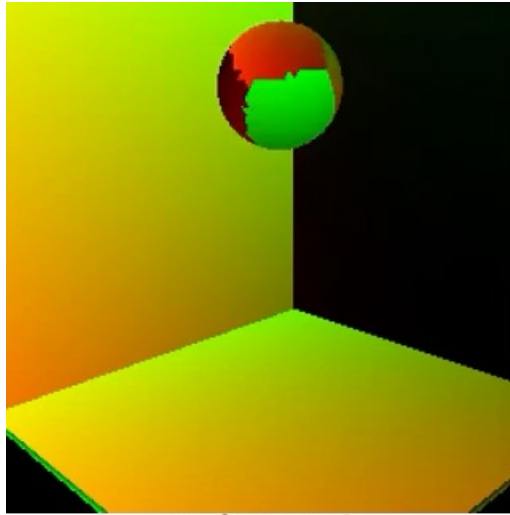
an 80s' anime fox running



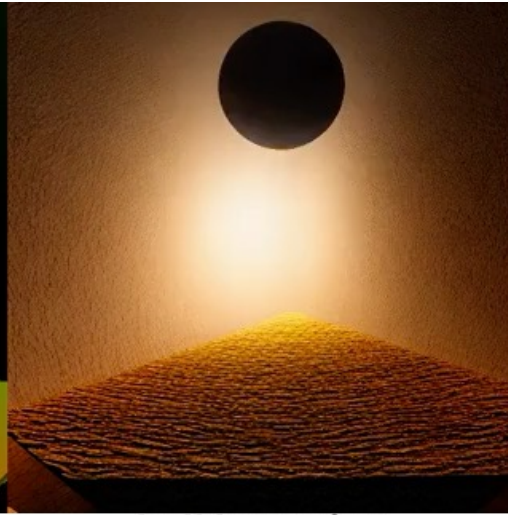
a lineart fox running

Character animation

# Experiments



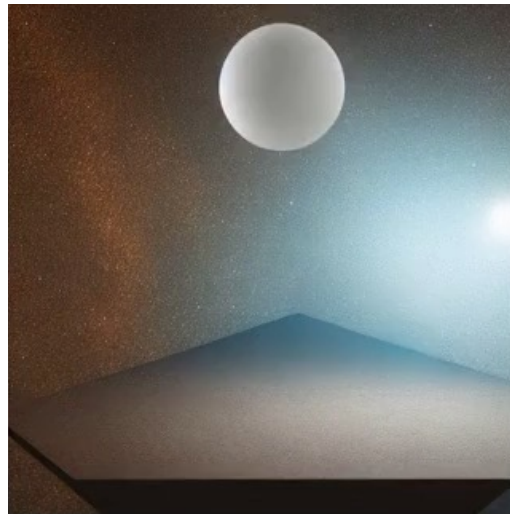
**animated  
mesh**



**ball bouncing  
in a cave**



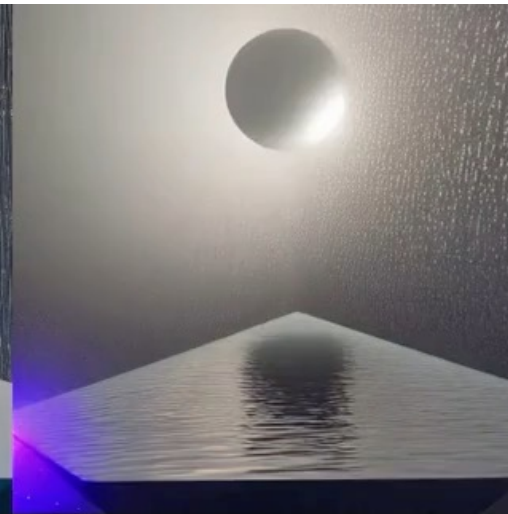
**ball bouncing  
in a cave**



**ball bouncing  
in a chamber**



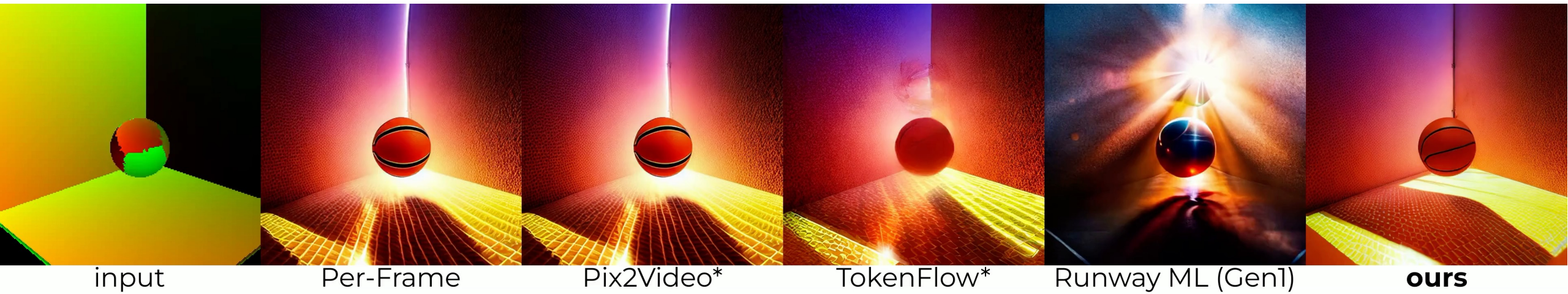
**ball bouncing  
in a chamber**



**ball bouncing  
in a chamber**

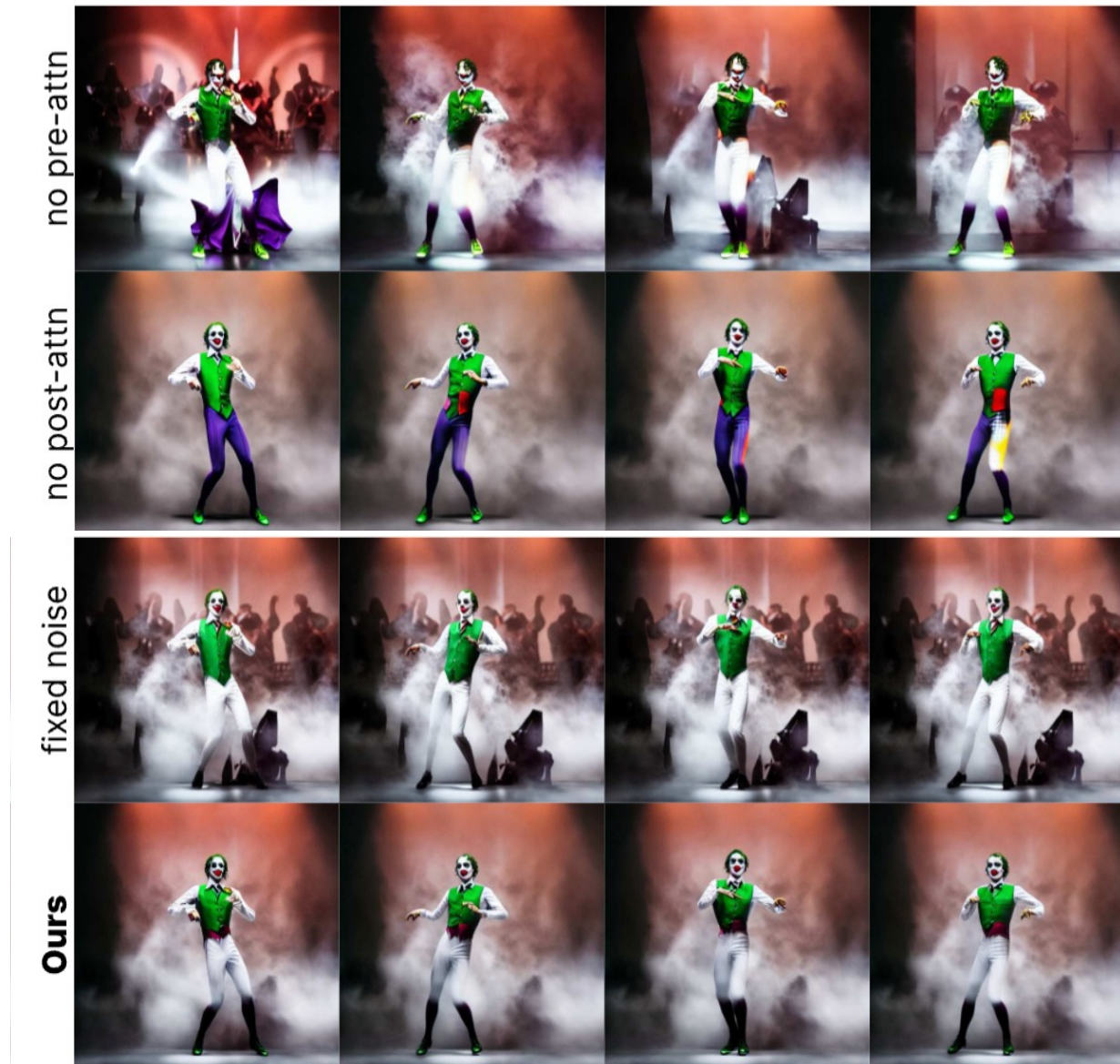
Physical simulation

# Experiments



Qualitative comparisons

# Experiments



Ablation study

# Experiments

Method	Frame Consistency $\uparrow$	Prompt Fidelity $\uparrow$
Per-Frame	0.9547	<b>0.3233</b>
Pix2Video*	0.9630	0.2983
TokenFlow*	0.9822	0.2737
Gen1	<b>0.9907</b>	0.3029
<b>Ours</b>	<u>0.9845</u>	<u>0.3227</u>

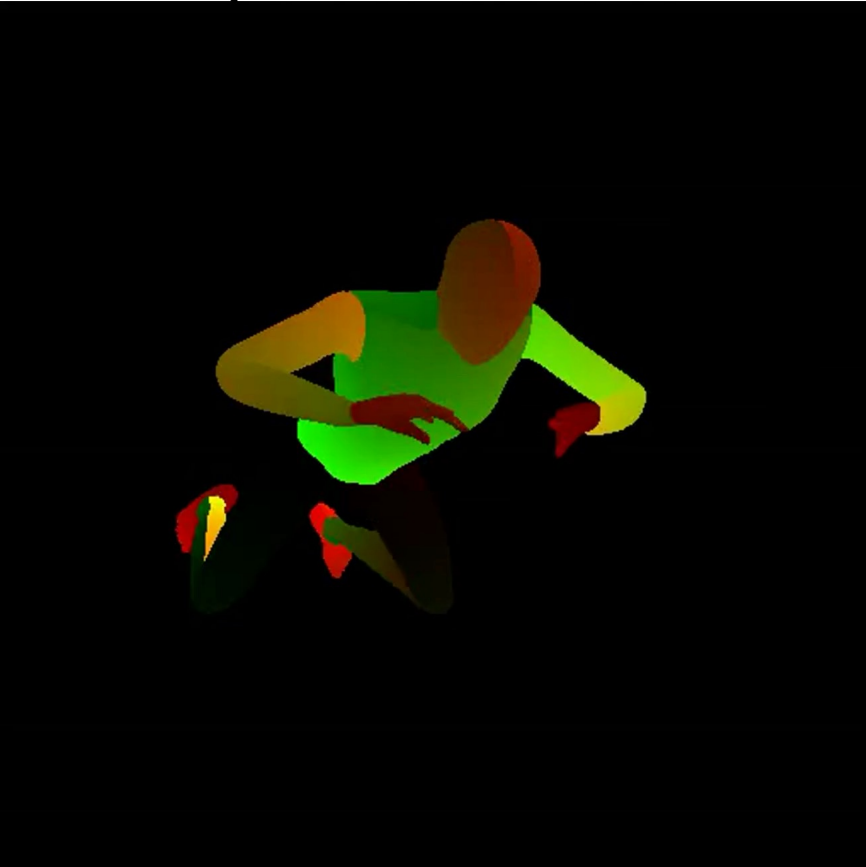
Method	Pairwise Frame Interval				
	1	5	10	15	20
Per-Frame	0.9547	0.9173	0.9109	0.9145	0.9062
Pix2Video*	0.9630	0.9503	0.9494	0.9415	0.9471
TokenFlow*	0.9822	<u>0.9754</u>	<u>0.9728</u>	<u>0.9706</u>	<u>0.9712</u>
Gen1	<b>0.9907</b>	0.9715	0.9582	0.9624	0.9601
<b>Ours</b>	<u>0.9845</u>	<b>0.9815</b>	<b>0.9738</b>	<b>0.9749</b>	<b>0.9727</b>

Quantitative evaluation

# Conclusion

---

input animated mesh



rendering prompt

a stormtrooper  
swimming

rendered video



**Generative Rendering** can animate creator-defined low-fidelity meshes and motion sequences, thereby bypassing steps requiring significant manual labor such as detailing, texturing, physical simulation.

# Q&A