# Video and 3D Generation

Kevin Lin

6/17/2024

# Evolution of Video Diffusion Model

- **Video-Gen Foundation Model**

- Faster training

- Faster inference
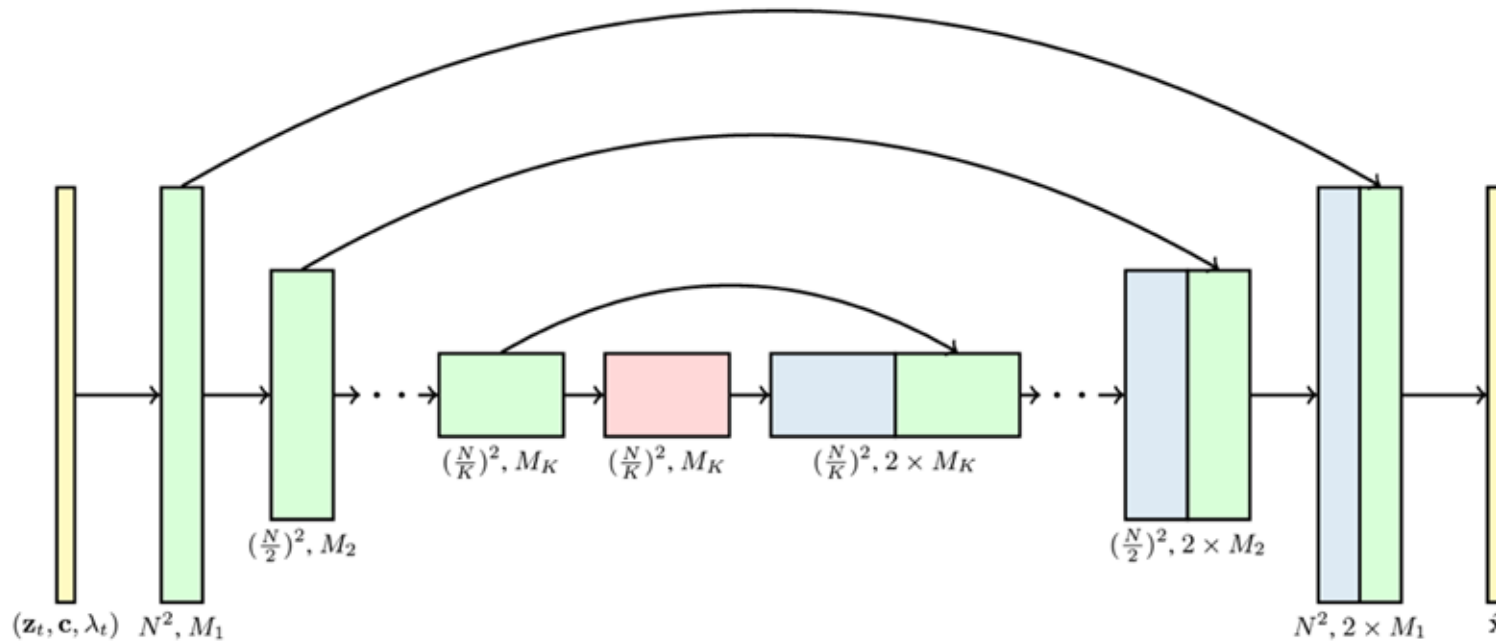
- Diverse content creation

# Creating Video from Text



Prompt: A flock of paper airplanes flutters through a dense jungle, weaving around trees as if they were migrating birds.
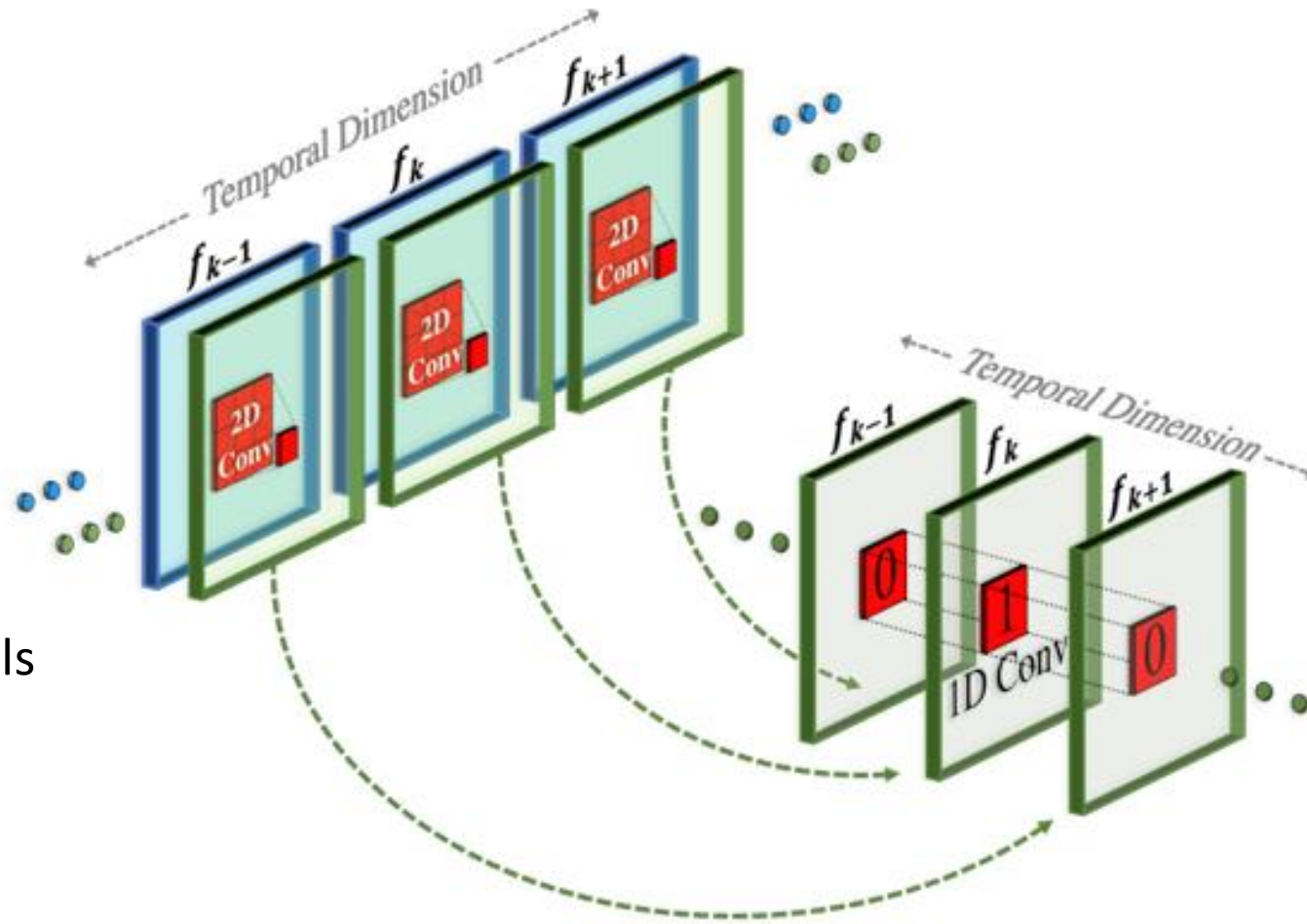
# Video Diffusion Models: *Pioneer Work*

- 3D UNet factorized over space and time
  - 2D conv is inflated to 3D
- Insert temporal attention layer that attends across the temporal dimension



Video Diffusion Models (Ho et al. NeurIPS'22)

# Make-A-Video



**Spatial Convolution**

Initialized from pre-trained T2I models

**Temporal Convolution**

Initialized with identity function

Make-A-Video (Singer et al. '22)

# Make-A-Video



**Spatial Attention**

Initialized from pre-trained T2I models

**Temporal Attention**

Initialized with zero projection (resulting in identity function)

Make-A-Video (Singer et al. '22)

# Preliminary Results

Prompt: Firework



Video Diffusion Models (Ho et al. '22)

# Curated Training Data Improves Performance

- Scaling training data from 10M to 577M video clips



Stable Video Diffusion (Blattmann et al. '23)

# Quality Improves as Training Compute Increases



Base compute            4x compute            32x compute

Sora: Video generation models as world simulators (OpenAI '24)

- Video-Gen Foundation Model

- Faster training

- Faster inference

- Diverse content creation

# Training a Video Diffusion Model is Expensive!

## SOTA Video Diffusion Model
**(Case Study: Stable Video Diffusion)**

577M Video Clips

1521M Parameters

Dedicated GPU clusters
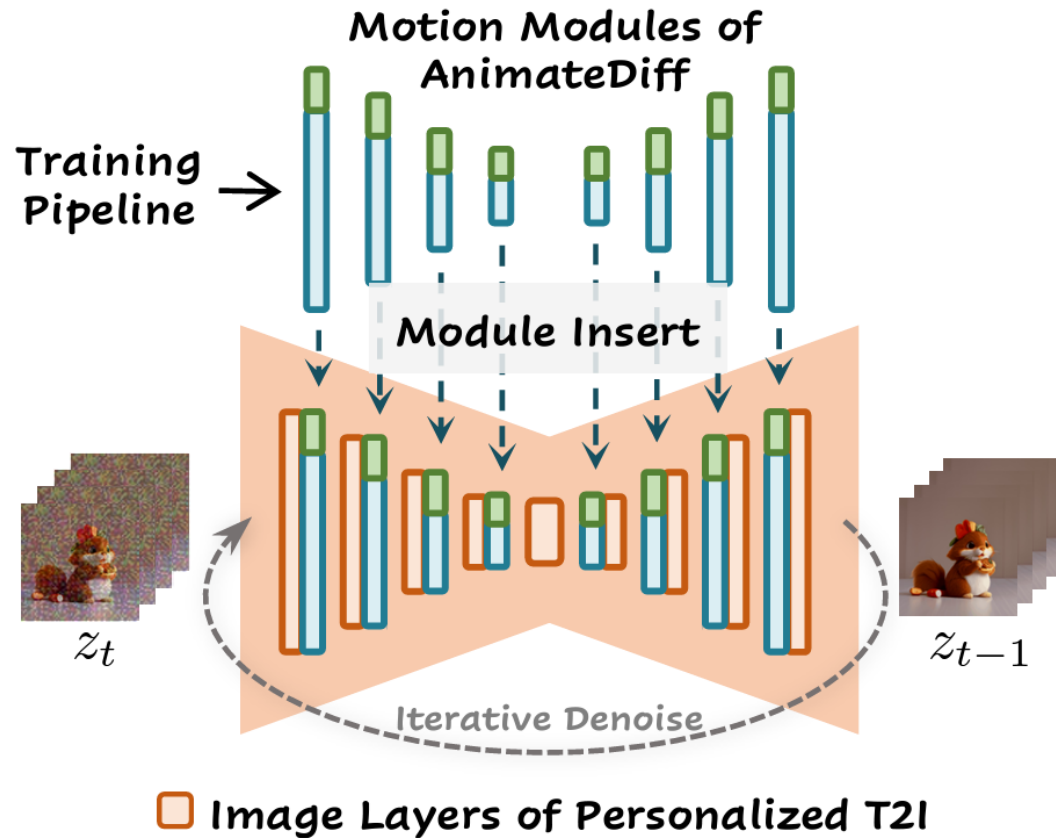
Slow and complex
training recipe

## Faster Training

20~50 Video Clips

+30M Parameters

~2 hours training

# Faster Training with Motion Modules



20~50 Video Clips

+30M Parameters

~2 hours training

AnimateDiff (Guo et al. ICLR'24)

# We have faster training now. How about inference?

**SOTA Video Diffusion Model**



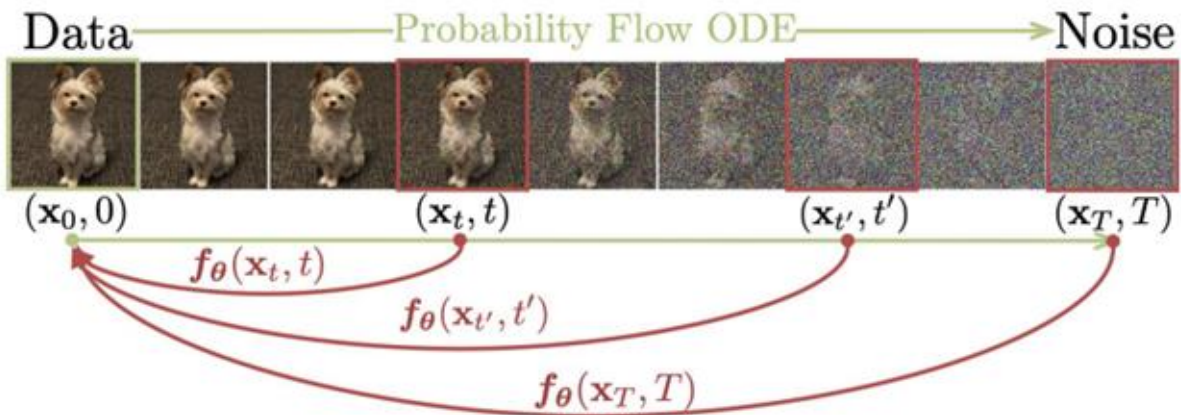250 Steps

**Faster Sampling**



4 Steps

# Challenges in Video Diffusion Distillation
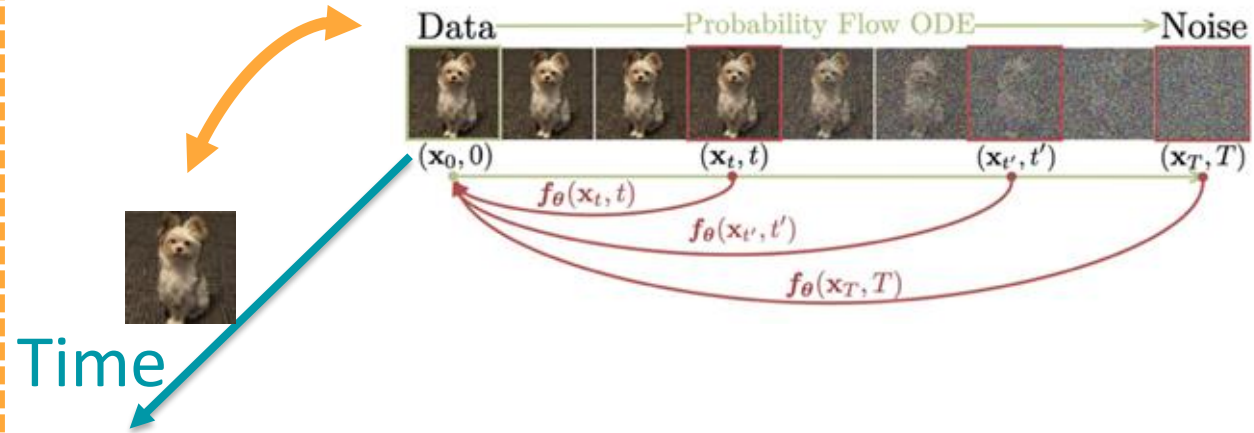
**Image** Diffusion Distillation



Consistency Model (Song et al. '23)

**Video** Diffusion Distillation

Appearance Consistency



Time

Motion Consistency
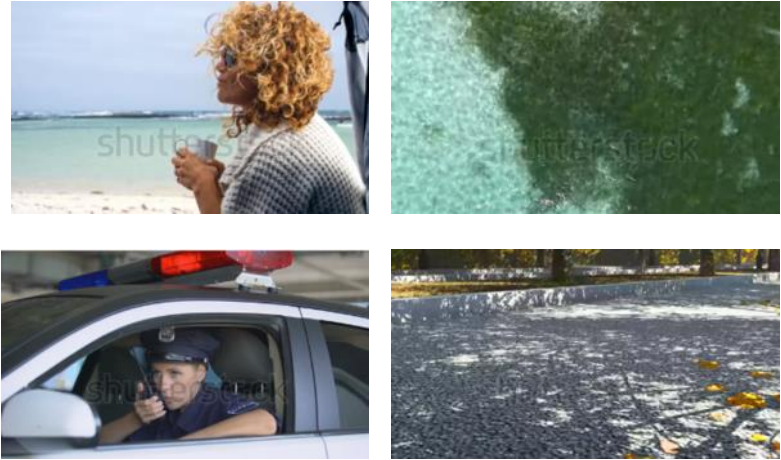
Motion Consistency Model (Zhai et al. '24)

# Challenges in Video Diffusion Distillation
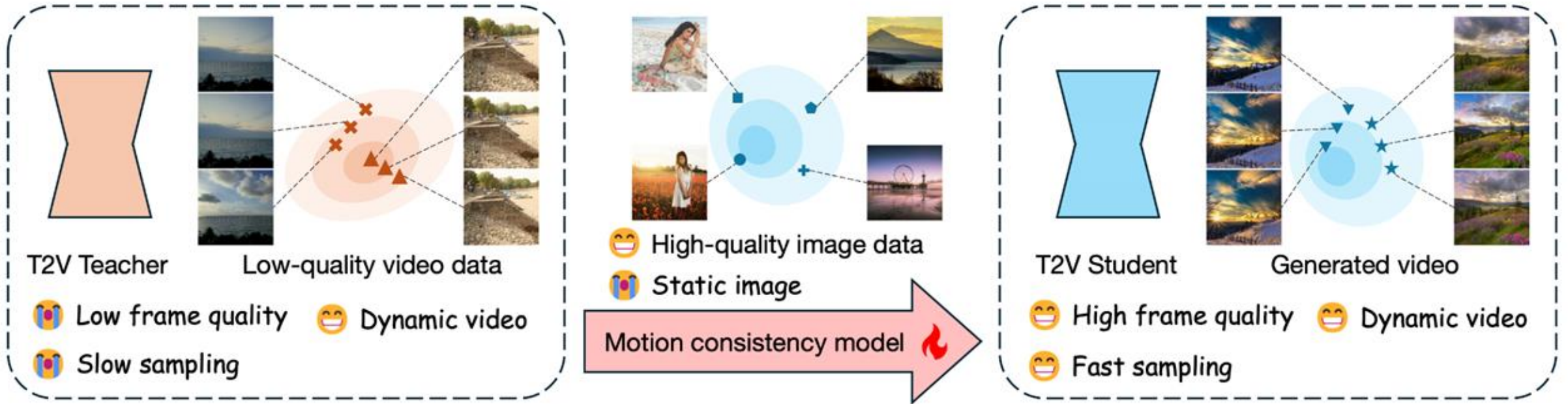
**High Quality Image** Data



LAION-Aesthetics

**Low Quality Video** Data



WebVid-10M

# Motion Consistency Model



Our motion consistency model not only distill the motion prior from the teacher to accelerate sampling, but also can benefit from an additional high-quality image dataset to improve the frame quality of generated videos.

Motion Consistency Model (Zhai et al. '24)

# Motion Consistency Model

| | Teacher<br>**50 steps** | MCM + WebVid<br>**4 steps** | MCM + LAION-aes<br>**4 steps** | MCM + Anime<br>**4 steps** | MCM + Realistic<br>**4 steps** | MCM + 3D Cartoon<br>**4 steps** |

Aerial uhd 4k view. mid-air flight over fresh and clean mountain river at sunny summer morning. Green trees and sun rays on horizon. Direct on sun.

Back of woman in shorts going near pure creek in beautiful mountains.

Misty mountain landscape



Motion Consistency Model (Zhai et al. '24)

✓ Video-Gen Foundation Model

✓ Faster training

✓ Faster inference

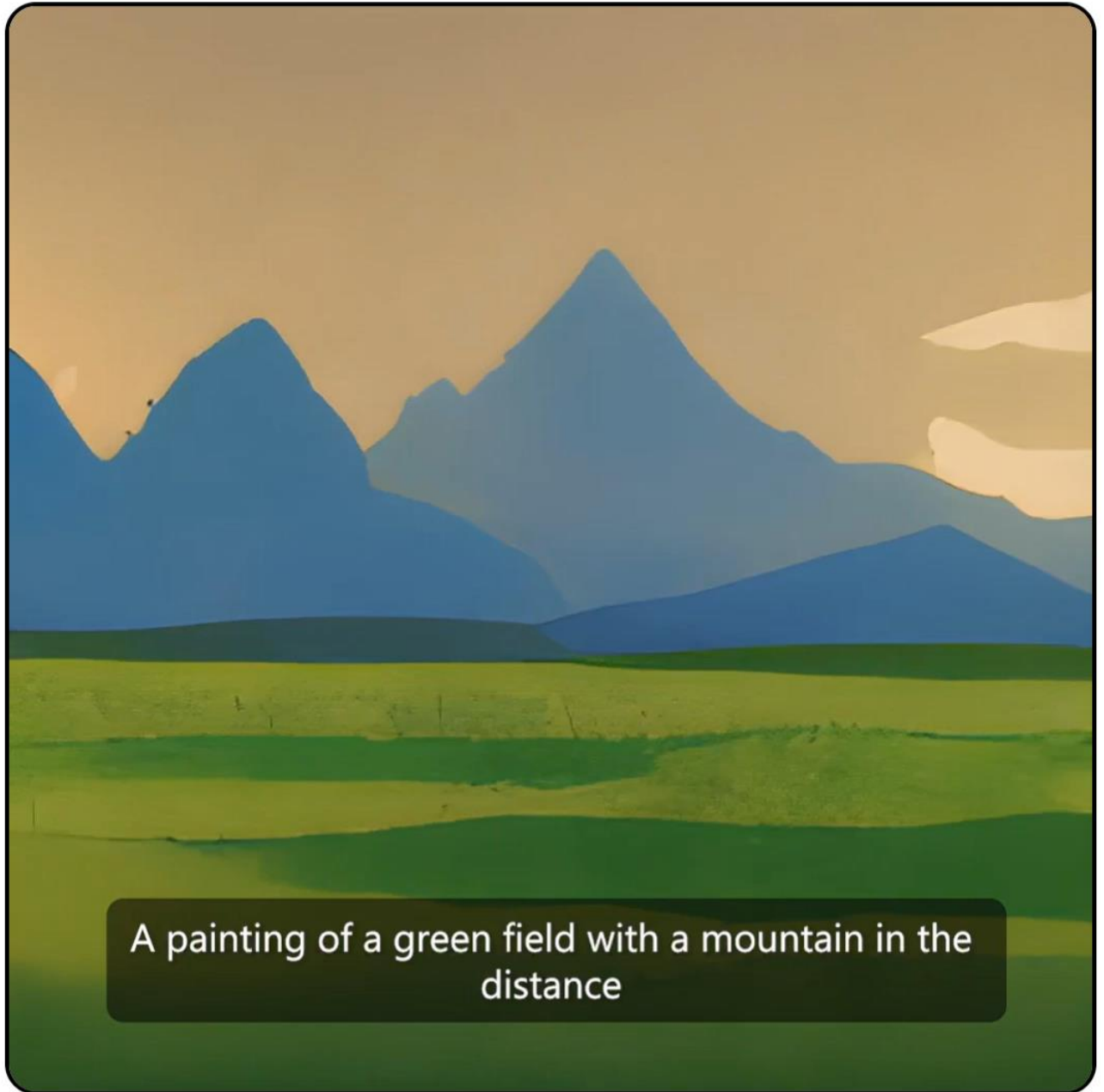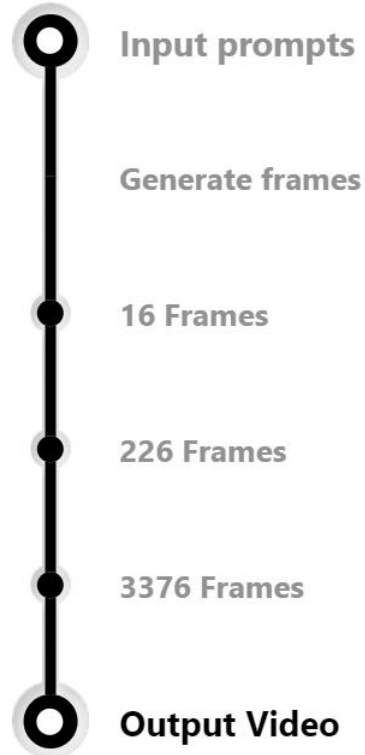✓ Diverse content creation:
   Variable durations   Controllability   Consistency

# LONG VIDEO

Given the prompts of a script, NUWA-XL can generate an extremely long video that conforms to it in a "coarse-to-fine" process.

- **Input prompts**
- **Generate frames**
- **16 Frames**
- **226 Frames**
- **3376 Frames**
- **Output Video**



A painting of a green field with a mountain in the distance

# Diffusion over diffusion



Figure 1: Overview of NUWA-XL for extremely long video generation in a "coarse-to-fine" process. A global diffusion model first generates $L$ keyframes which form a "coarse" storyline of the video, a series of local diffusion models are then applied to the adjacent frames, treated as the first and the last frames, to iteratively complete the middle frames resulting $O(L^m)$ "fine" frames in total.

NUWA-XL (Yu et al., ACL'23)

Simple Cartoon Videos

Rich and Diverse Contents?

# Human Dance Generation

- Different subject – Same pose



Target Pose    Dance#1    Dance#2    Dance#3    Dance#4    Dance#5

DisCo (Wang et al. CVPR'24)

# Human Dance Generation

- Same subject – Different pose



Reference Image     Dance #1     Dance #2     Dance #3     Dance #4

DisCo (Wang et al. CVPR'24)

# DisCo for Human Dance Generation

- By disentangling the control from all three conditions, DisCo enable arbitrary compositionality of human subjects, backgrounds, and dance-moves.



(a) Model Architecture with Disentangled Control

(b) Human Attribution Pre-training

DisCo (Wang et al. CVPR'24)

# Wonderjourney





Wonderjourney: Going from Anywhere to Everywhere (Yu et al. CVPR'24)

# CAT3D: Multi-View Latent Diffusion Model



Input Image(s) → Sample from multi-view diffusion model (5 seconds) → Generated Views → Optimize a NeRF (55 seconds) → 3D Model

NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. (Mildenhall et al., ECCV '22)
CAT3D: Create Anything in 3D with Multi-View Diffusion Models (Gao et al. '24)

✓ Video-Gen Foundation Model

✓ Faster training

✓ Faster inference

✓ <span style="color:red">Coherent video</span> + <span style="color:green">Realistic contents</span> + <span style="color:orange">3D consistency</span>

→ World Model!

# Video Generation Models as World Simulators



- ✓ 3D Consistency

- ✓ Coherence

Sora: Video generation models as world simulators (OpenAI '24)

# Discussion

- How to accurately model the physical and digital world?
  *Physics, object states, and things beyond languages*

- How to effectively evaluate the emerging capabilities?
  *Needs of exploration and new benchmark*

- Safety
  *Learning from real-world use is a critical component*

# MMWorld: World Model Evaluation in Videos



Q: What has been changed in the video?
A: The bottom drawer has been closed.

Q: How many animals appear in the video?
A: Two. There are a horse and a dog

Q: What is the reason that the lady decides to use the easy frost?
A: Because it has no-fuss frosting.

Attribution Understanding

Temporal Understanding

Domain Expertise

Procedure Understanding

Q: What was first added into the milk?
A: Cocoa powder.

5.5%
7.6%
18.0%

Multi-faceted Reasoning

8.8%

Future Prediction

Q: What will happen next as the price is below the blue and red lines?
A: The price will go down.

48.0%

10.9%

Explanation

Counterfactual Thinking

Q: How do the pulleys move when the hands are off the pulley system?
A: Two static and two moving upward.

Q: What would happen if the man skipped the step shown in the video?
A: The desktop of the coffee table will be upside down, which will make it impossible to mount the legs.

MMWorld: Towards Multi-discipline Multi-faceted World Model Evaluation in Videos (He et al. '24)

# Discussion: World Model



- Build internal representations of the 3D world

- Predict and simulate future events within the internal representation

- Reasoning and planning: governed by our brain's prediction of the future based on our internal world model

[1] Primary Visual Cortex Represents the Difference Between Past and Present. N. Nortmann et al. 2015
[2] Counterintuitive behavior of social systems. J.W. Forrester. 1971.
[3] Motion-Dependent Representation of Space in Area MT+. M. Gerrit et al. 2013

# Discussion: World Model



- Build internal representations of the 3D world

- Predict and simulate future events within the internal representation

- Reasoning and planning: governed by our brain's prediction of the future based on our internal world model

*The image of the world around us, which we carry in our head, is just a model. He selects concepts and relationships, and uses those to represent and simulate the real system.*

- Jay Wright Forrester, Father of System Dynamics talks about mental world models

VideoLCM
(Wang et al. '23)

AnimateLCM          MCM
(Wang et al. '24)   (Zhai et al. '24)

**Accelerating Video Diffusion**

AnimateDiff-Lightning
(Lin et al. '24)

AnimageDiff
(Guo et al. '23)

Text2Video-Zero
(Khachatryan et al. '23)

**Training-Efficient Techniques**

MagicVideo          SimDA
(Zhou et al. '22)   (Xing et al. '22)

DSDN
(Liu et al. '23)

Tune-A-Video
(Wu et al. '23)

Dreamix             Edit-A-Video
(Molad et al. '23)  (Shin et al. '23)

**Video Editing**

FateZero    MeDM         RAVE
(Qi et al. '23)(Chu et al. '23)(Kara et al. '24)

Ground-A-Video
(Jeong et al. '23)

Tokenflow
(Geyer et al. '23)

Stable Video Diffusion
(Blattmann et al. '23)

Sora                EmuVideo
(OpenAI '24)        (Girdhar et al. '2)

**Video Generation Foundation Model**

ModelScopeT2V       Veo
(Wang et al. '23)   (Google '24)

Show-1              Gen-2
( Zhang et al. '23) (Runway '23)

NUWA-XL
(Yin et al. '23)

**Long-From Video Generation**

LVDM
(He et al. '22)

4Real       CAT3D
(Yu et al. '24)(Gao et al. '24)

**3D and 4D Generation**

Animate124
(Zhao et al. '23)

Wonderjourney
(Yu et al. '24)

CameraCtrl          Pandora
(He et al. '24)     (Xiang et al. '24)

**World Model**

iVideoGPT           Genie
(Wu et al. '24)     (Bruce et al. '24)

3D-VLA
(Zhen et al. '24)

**Pioneering work in Video Generation**

VDM             Imagen Video        Make-A-Video        Align your latent
(Ho et al. '22) (Ho et al. '22)     (Singer et al. '22) (Blattmann et al. '23)

# Acknowledgment

Mike Shou's Youtube video

- [Tutorial: Video Diffusion Models](#)

Lilian Weng's blog

- [What are Diffusion Models?](#)

Hung-Yi Lee's Youtube video

- [Introduction to Diffusion Models](#)

Yuyang Zhao          Yining Hong          Yuanhao Zhai

Thank you!