# Lumiere: A Space-Time Diffusion Model for Video Generation

Authors:
Omer Bar-Tal, Hila Chefer, Omer Tov,  Charles Herrmann, Roni Paiss, Shiran Zada, Ariel Ephrat, Junhwa Hur, Guanghui Liu, Amit Raj, Yuanzhen Li, Michael Rubinstein, Tomer Michaeli, Oliver Wang, Deqing Sun, Tali Dekel, Inbar Mosseri

Presented by:
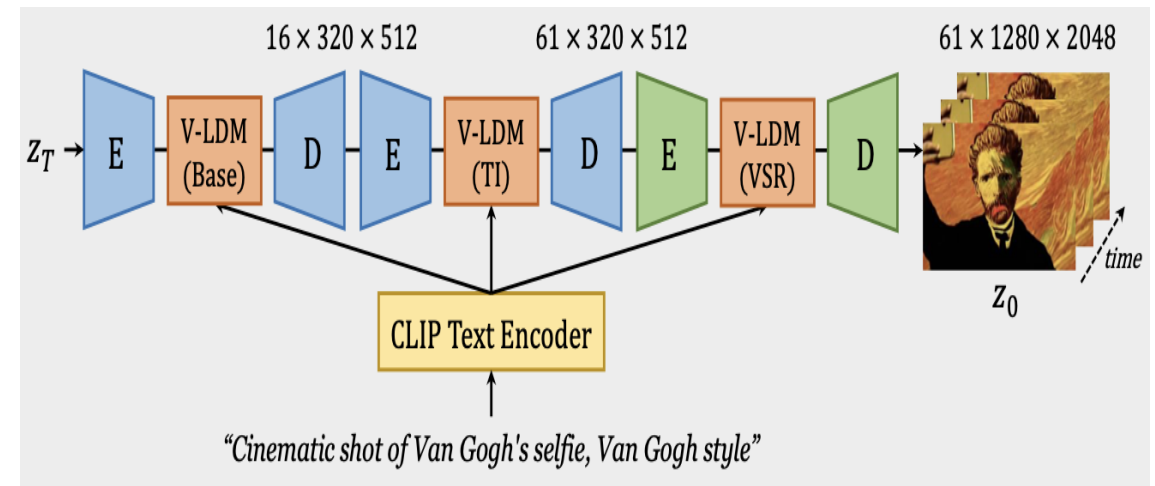Anantapadmanaabha Prasannakumar

# Outline

- Motivation
- Method
- Evaluations
- Applications
- Societal Impact
- Limitations
- Conclusion

# Motivation

▶ Restricted capability of existing models

  ▶ Sensitive to error

  ▶ Suffers from memory and computing constraints

  ▶ Obtaining large-scale data is cumbersome

  ▶ Training large-scale T2V is challenging

# Motivation

▶ Employing temporal cascade design is hindersome

▶ Generates aggressively sub-sampled set of keyframes

▶ TSR modules are constrained to fixed, small temporal context

▶ Cascaded training suffers from domain gap



LaVie: High-Quality Video Generation with Cascaded Latent Diffusion Models, IJCV 2024, https://arxiv.org/pdf/2309.15103

# Method - Lumiere

▶ Utilizes Diffusion Probabilistic Models

   ▶ Through denoising steps, trained to approximate
      data distribution

   ▶ Starting from noise, a clean sample is drawn from the
      targeted distribution

▶ Incorporates additional guiding signals

# Common T2V Framework

- Base Model

- Temporal Super-resolution Model (TSR)

- Spatial Super-resolution Model (SSR)

# Lumiere Framework:

- Base Model

- Spatial Super-resolution Model (SSR)

- Multidiffusion

Lumiere: A Space-Time Diffusion Model for Video Generation, SIGGRAPH 2024, https://arxiv.org/pdf/2401.12945
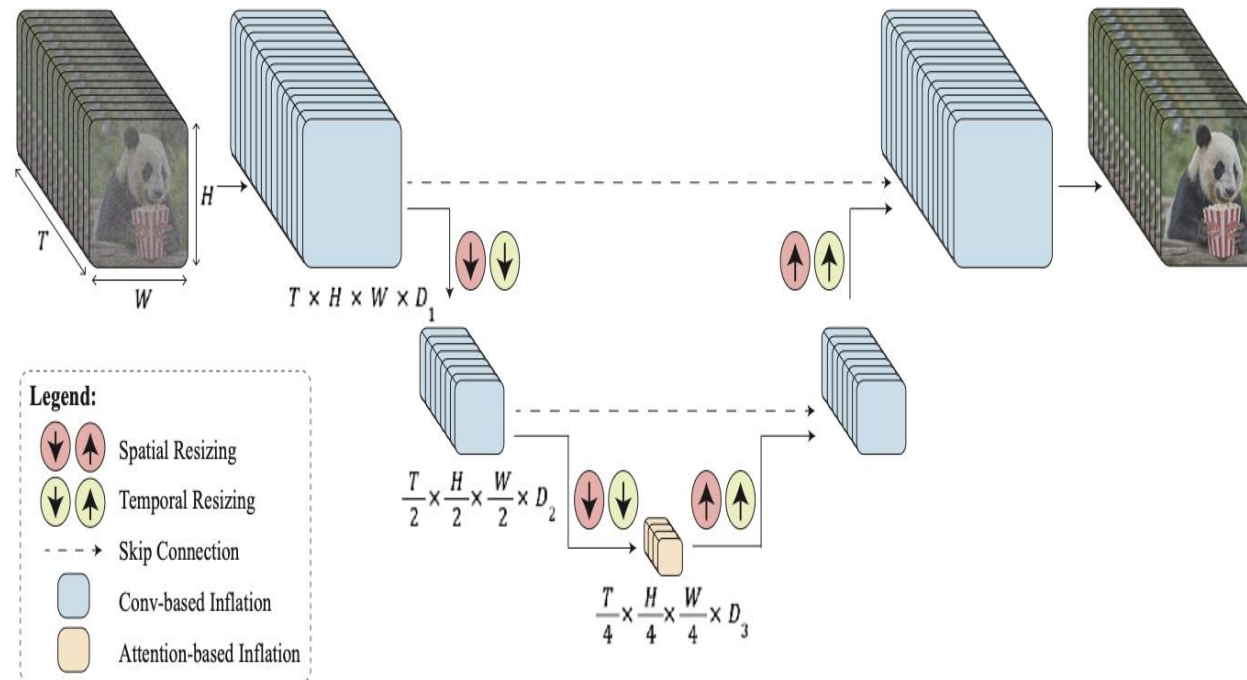
# Method – U-Net

- Encoder
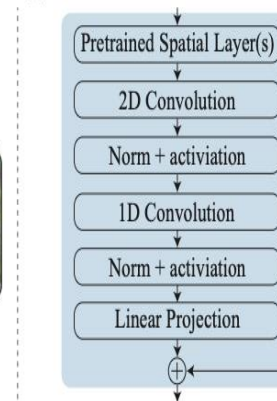
- Decoder

# Method - STUNet

- Employs the U-Net architecture

- Consists of 2 inflation blocks

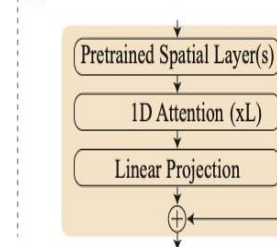- Interleave temporal blocks to T2I Architecture

# Method - STUNet



(a) Space-Time UNet (STUNet)

Legend:
- Spatial Resizing
- Temporal Resizing
- Skip Connection
- Conv-based Inflation
- Attention-based Inflation

(b) Convolution-based Inflation Block
- Pretrained Spatial Layer(s)
- 2D Convolution
- Norm + activiation
- 1D Convolution
- Norm + activiation
- Linear Projection

(c) Attention-based Inflation Block
- Pretrained Spatial Layer(s)
- 1D Attention (xL)
- Linear Projection

- Trains only new parameters

- Performs identity Initialization

- Low computational overhead

# Method - MultiDiffusion

- New generation process

- Employs one global denoising step

MultiDiffusion: Fusing Diffusion Paths for Controlled Image Generation, ICML 2023, https://arxiv.org/pdf/2302.08113

# Method - MultiDiffusion



Generation with independent diffusion paths



Generation with fused diffusion paths using MultiDiffusion

# SSR with Multidiffusion

- ▶ An inflated SSR network can only operate on short videos

- ▶ Employ multidiffusion for smooth temporal transition

- ▶ Multidiffusion prevents temporal artifacts

  - ▶ Resolved by linearly combining video segments

# SSR with Multidiffusion

▶ At each generation step:

    ▶ split noisy input video $J \in \mathbb{R}^{H \times W \times T \times 3}$ into $1 \ldots N$ overlapping segments

    ▶ Where $J_i \in \mathbb{R}^{H \times W \times T' \times 3}$ is the $i^{th}$ segment

    ▶ Temporal duration: $T' < T$

▶ To reconcile per-segment SSR predictions:

$$\arg\min_{J'} \sum_{i=1}^{n} \left\| J' - \Phi(J_i) \right\|^2 .$$

# Evaluation Setup

- Train T2V model on 30M videos with text prompts

  - Videos are 80 frames long at 16 fps

  - 109 text prompts

  - Base model dimension: 128 x 128 frames

  - SSR dimension: 1024 x 1024 frames
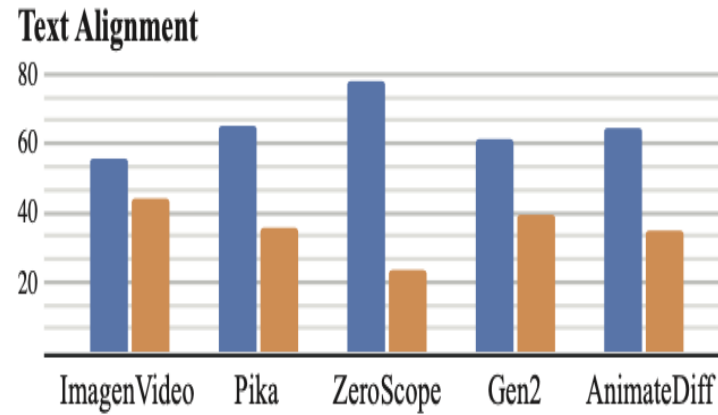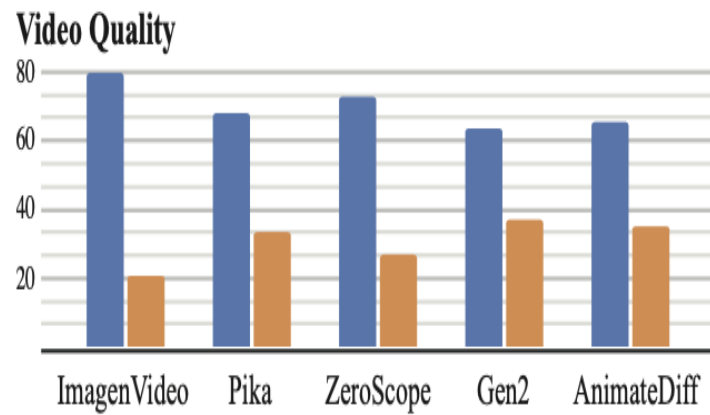
# Zero-shot on UCF-101

| Method | FVD ↓ | IS ↑ |
|---|---|---|
| MagicVideo (Zhou et al., 2022) | | |
| Emu Video (Girdhar et al., 2023) | | |
| Video LDM (Blattmann et al., 2023b) | | |
| Show-1 (Zhang et al., 2023a) | | |
| Make-A-Video (Singer et al., 2022) | | |
| PYoCo (Ge et al., 2023) | | |
| SVD (Blattmann et al., 2023a) | | |
| **Lumiere (Ours)** | | |

# User Study

- Two-alternative Forced Choice protocol Adopted
  - Randomly ordered pairs of videos are provided
  - 400 user judgments obtained
  - 109 prompts were utilized
  - Fixed random seed
  - Spatial and Temporal alignment

# User Study



Text-to-video
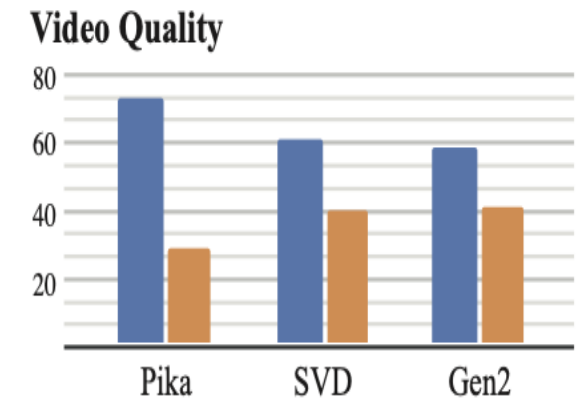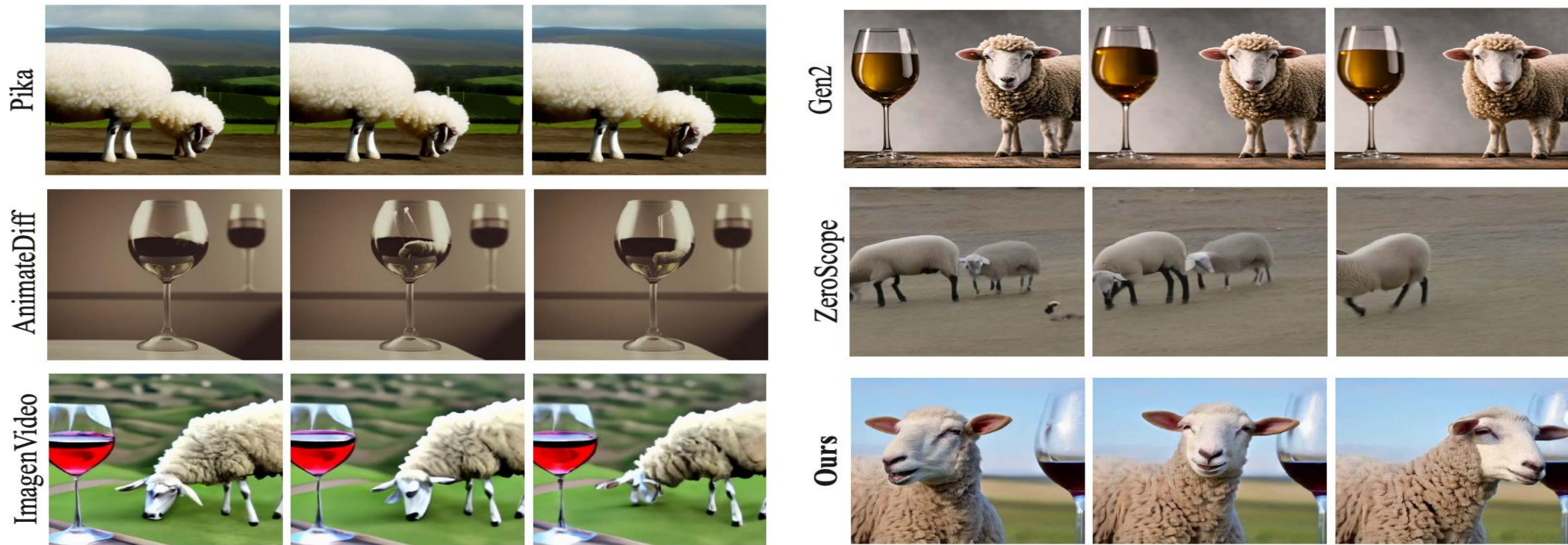
Image-to-video

# User Study

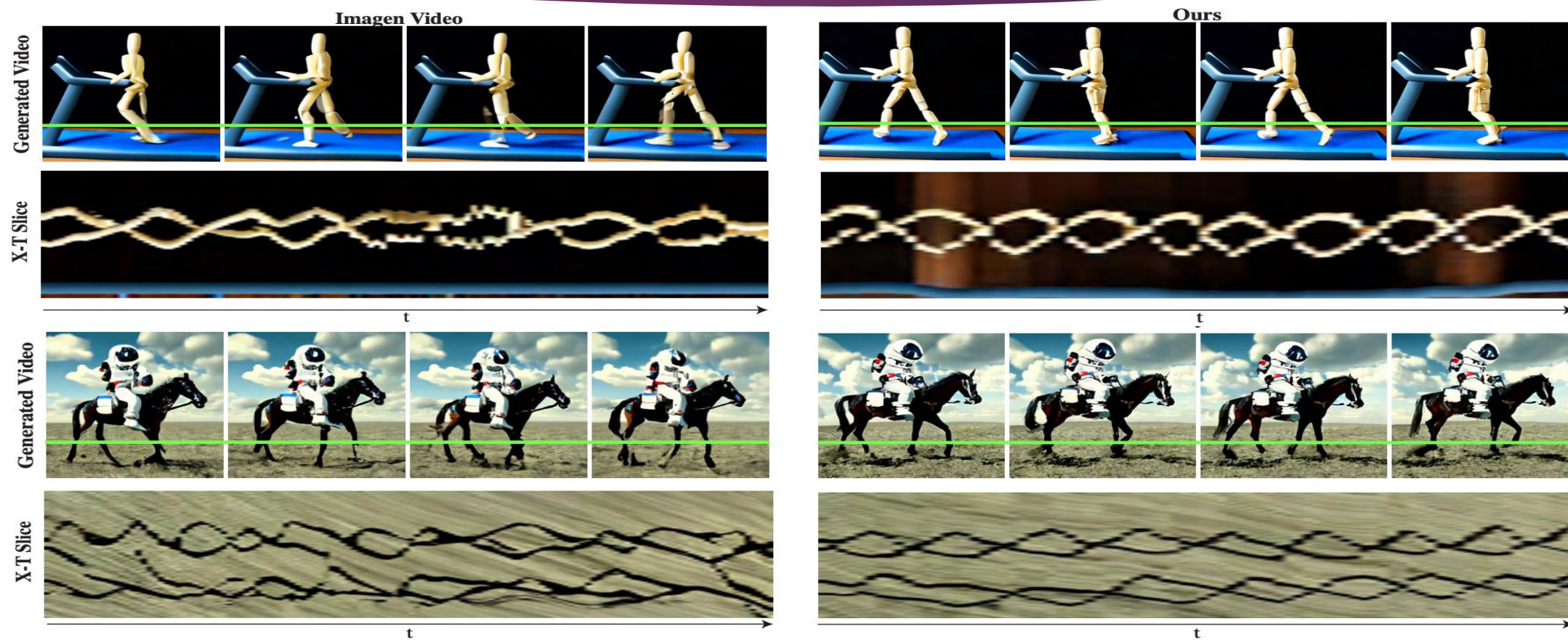# Qualitative Evaluation



A sheep to the right of the wine glass

# Temporal Consistency

# Ablation - Initialization



Loss Functions Over Epochs (UCF-101)

# Visualize Initialization Schemes

# Ablation - Multidiffusion



Ours

Xt-Slice

Ours

Without MultiDiffusion

# Applications – Video Editing



Original Video

Generated Video

SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations, Jan 2022, https://arxiv.org/pdf/2108.01073

# Application – Stylized Generation

▶ Pre-trained T2I weights remain fixed

▶ Newly added temporal layers are trained

▶ Linear interpolation between fixed and fine-tuned T2I weights
  ▶ $W_{interpolate} = \alpha \cdot W_{style} + (1 - \alpha) \cdot W_{orig}$
  ▶ Where $\alpha \in [0.5, 1]$

# Application – Stylized Generation
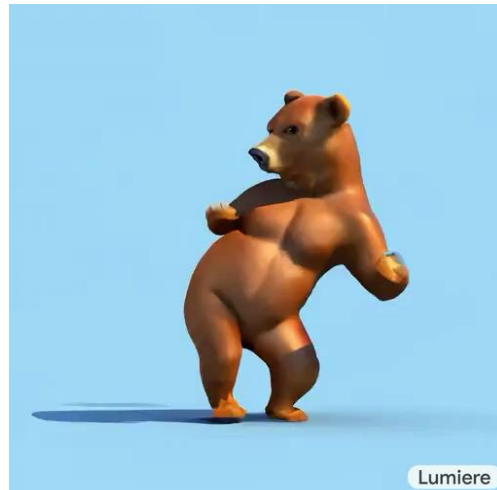
Vector art styles



Reference Image                                                Output

# Application – Stylized Generation

Realistic styles



Reference Image                                    Output

# Application – Conditional Generation

▶ Model conditioned on additional input signals

   ▶ Noisy video $J \in \mathbb{R}^{H \times W \times T \times 3}$

   ▶ Text prompt

   ▶ Masked conditioning video $C \in \mathbb{R}^{H \times W \times T \times 3}$

   ▶ Binary Mask $M \in \mathbb{R}^{H \times W \times T \times 1}$


▶ Concatenated Tensor $\langle J, C, M \rangle = \mathbb{R}^{T \times H \times W \times 7}$
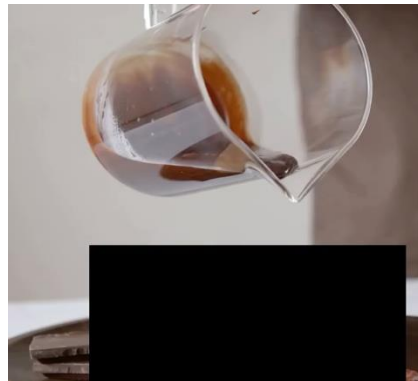
# Application – Image to Video
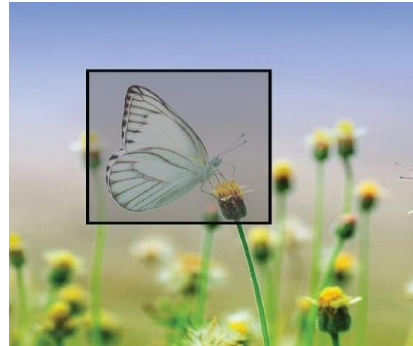
# Application - Inpainting



Video + Mask



Output

# Application - Cinemagraphs



Source Image + Mask



Output

# Societal Impact

▶ Risk of misuse

   ▶ Tools for detecting biases and malicious use cases

   ▶ To ensure safe and fair use

# Limitations

▶ The model cannot generate videos

  ▶ Multiple shots

  ▶ Transition between scenes

▶ The model operates in pixel space

# Conclusion

▶ Presents a novel T2V framework

    ▶ Built on a pre-trained T2I model

    ▶ Introduces space-time U-Net Architecture

    ▶ Utilizes Multidiffusion framework

▶ Demonstrates state-of-the-art generation results

▶ Showcases applicability to various downstream tasks

Thank you