

Lumiere: A Space-Time Diffusion Model for Video Generation

Omer Bar-Tal*^{1 2} Hila Chefer*^{1 3} Omer Tov*¹ Charles Herrmann^{†1} Roni Paiss^{†1} Shiran Zada^{†1}
Ariel Ephrat^{†1} Junhwa Hur^{†1} Guanghui Liu¹ Amit Raj¹ Yuanzhen Li¹ Michael Rubinstein¹
Tomer Michaeli^{1 4} Oliver Wang¹ Deqing Sun¹ Tali Dekel^{1 2} Inbar Mosseri^{†1}

Google Research

날짜: 2024-11-26

발표자: 홍진욱

9th-together-RL

Contents

1

Introduction

2

Method

3

Contribution

4

Experiments

5

Conclusion

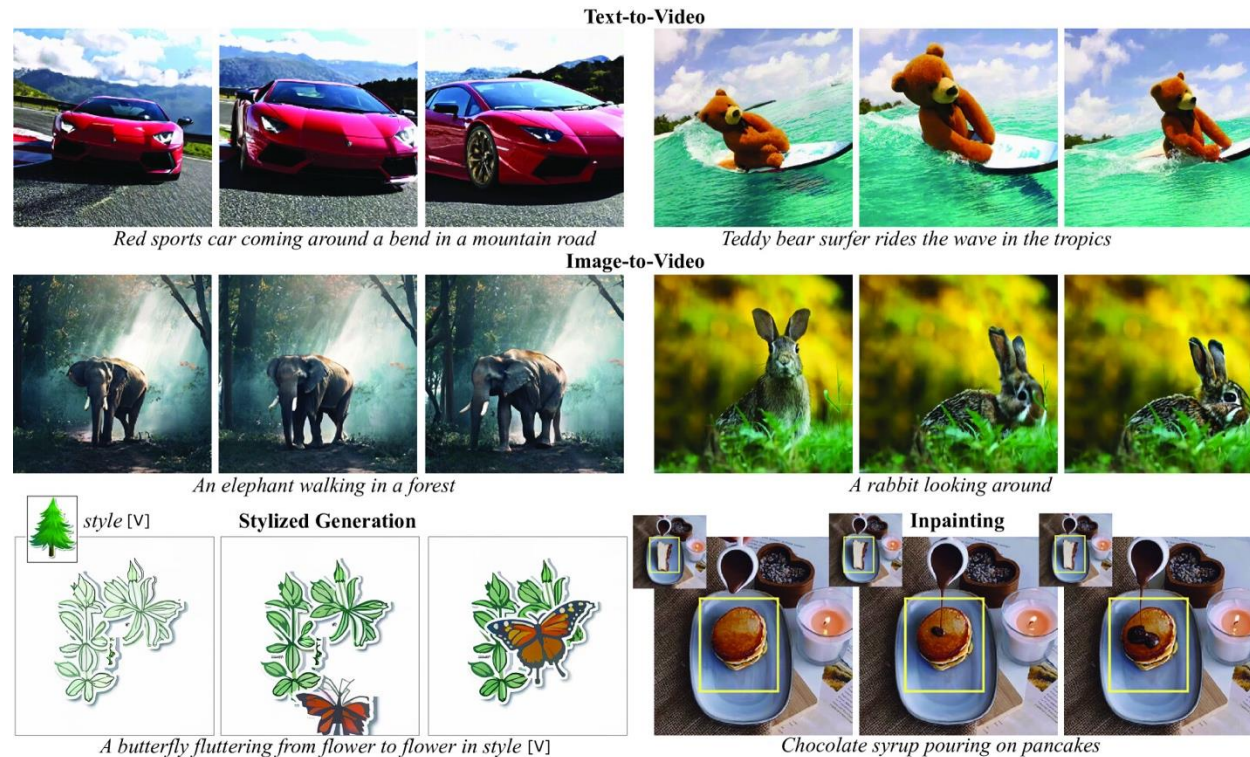
6

Q&A

Introduction

Background

- While deep generative models have shown astonishing results in the visual domain.
- Recent progress achieves promising results of video generation such as text-to-video, image-to-video, and video-to-video.



Introduction

Summary

- Lumiere successfully extend the capability of a pretrained text-to-image model (Imagen) into video generation for various applications.
- Different from common approaches, which adopt a cascade of temporal super-resolution, Lumiere can generate full frames at once (80 frames, 16 fps, 5s).
- Space-Time UNet (STUNet) + Temporal super-resolution w/ multi-diffusion
- Competitive results on benchmark (UCF101) and user-study.

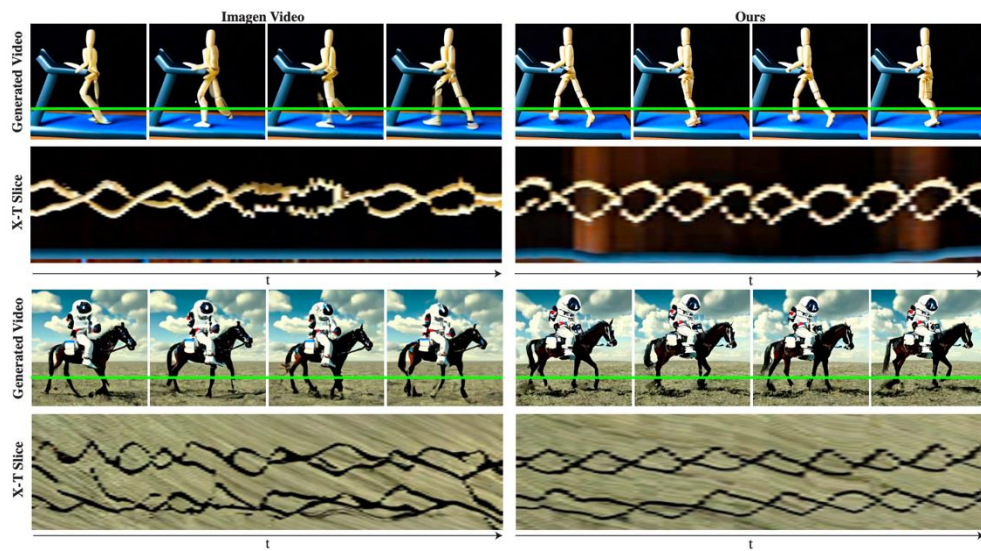
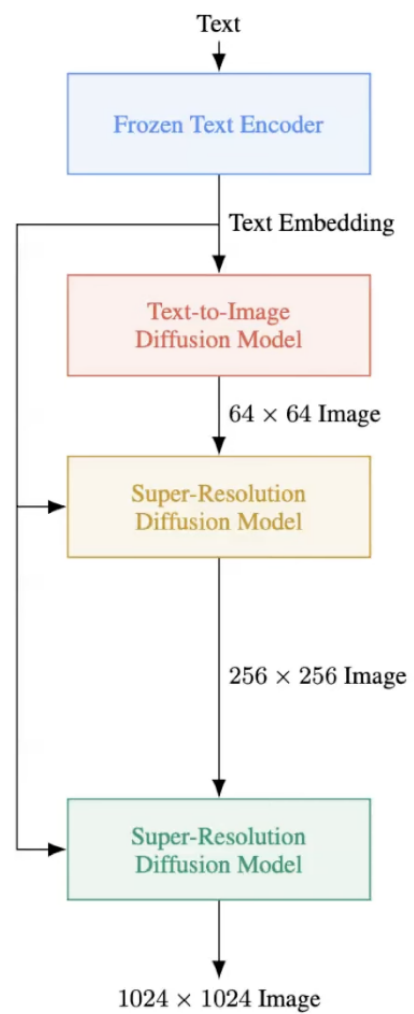


Imagen (NeurIPS '22)



"A Golden Retriever dog wearing a blue checkered beret and red dotted turtleneck."

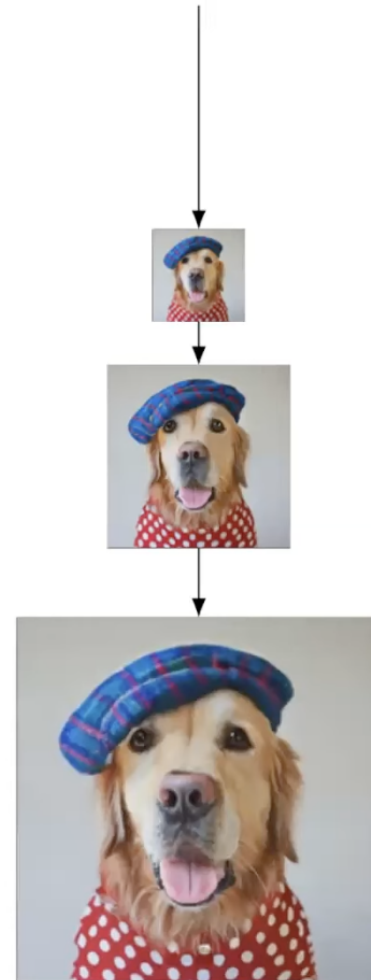
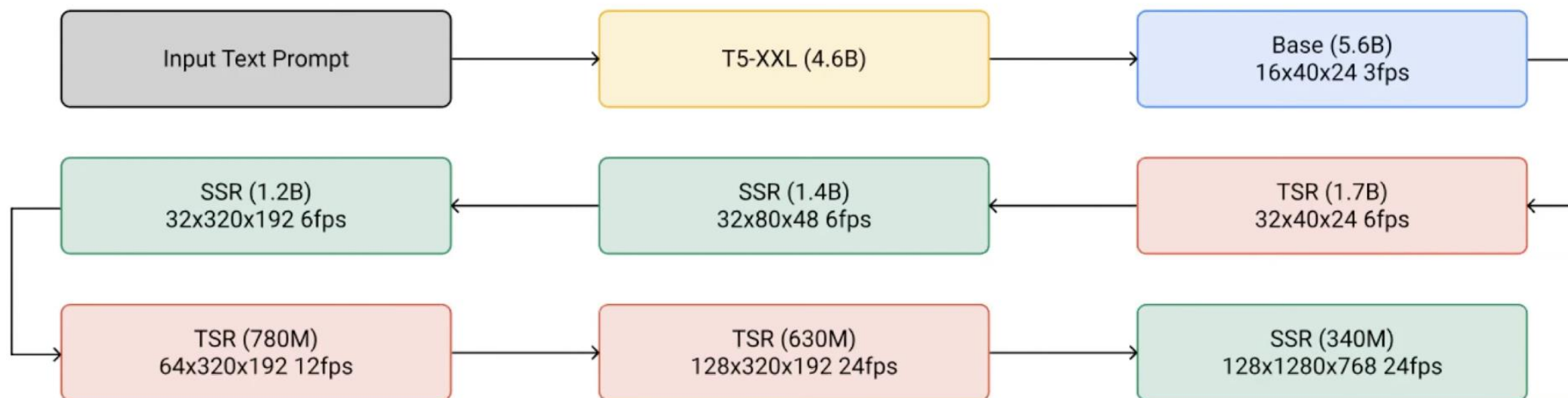


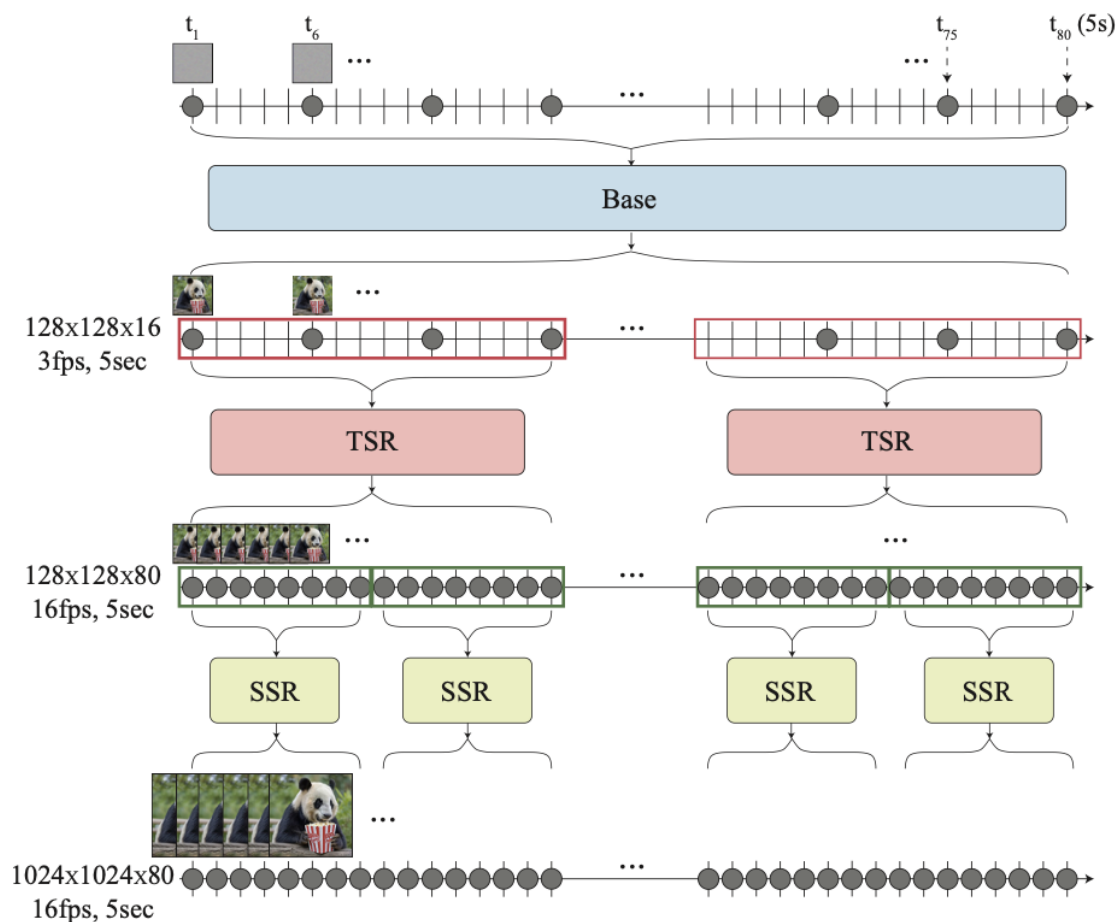
Imagen Video for Text-to-Video Generation (ArXiv '22)

- Imagen video generates the contents of videos from a user-provided text prompt, sequentially increasing the spatial and temporal resolutions.
- Here, spatial and temporal super-resolutions are separately processed.



Common Approaches of T2V vs. Lumiere

(a) Common Approach with TSR model(s)



First, generate key frames at a low resolution and a low frame ratio (fps)

Divide the total duration into some chunks, and conduct temporal super-resolution (TSR).

(-) base model: aliasing on fast motions

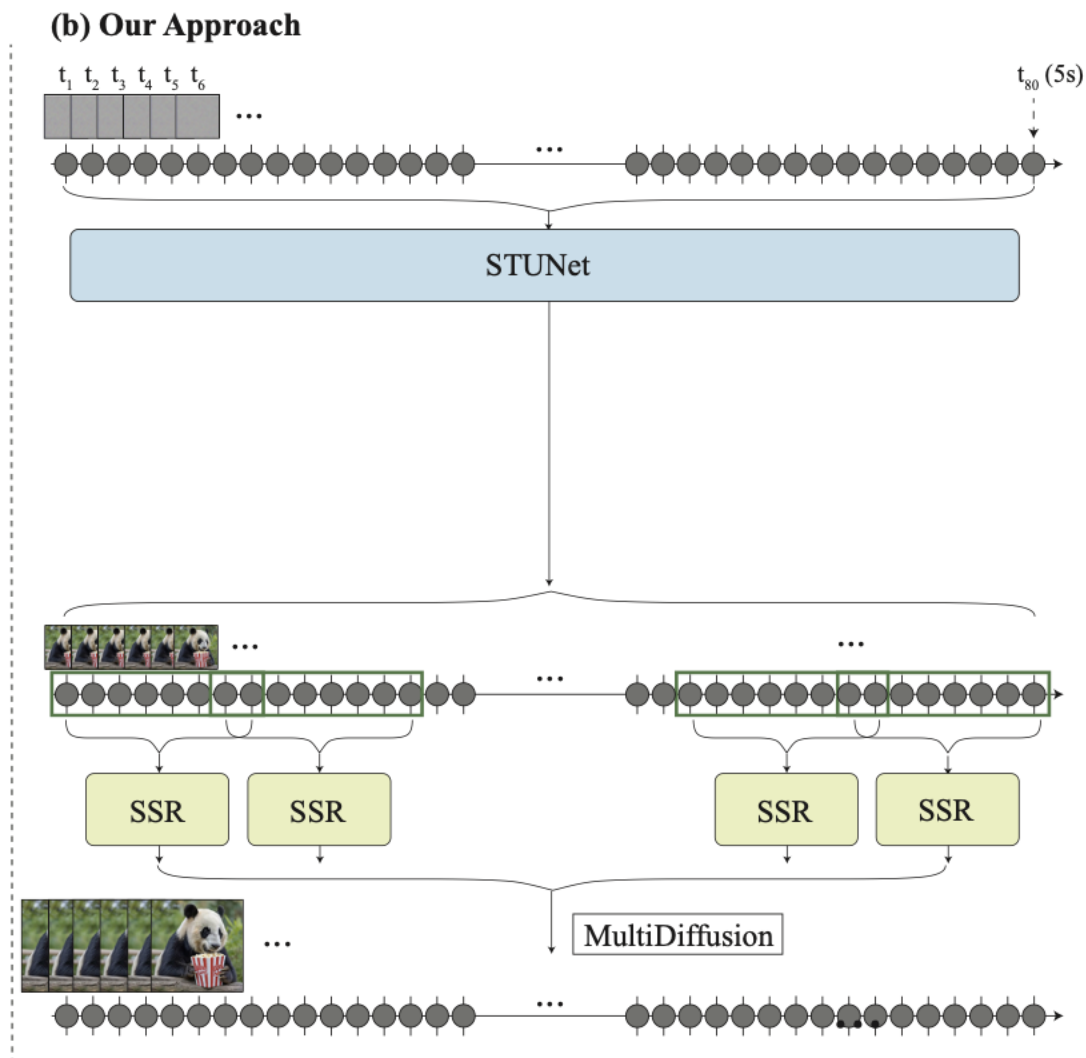
(-) TSR model: inconsistent motions, covariate shifts

For each chunk at the high frame ratio, conduct spatial super-resolutions.

Common Approaches of T2V vs. Lumiere

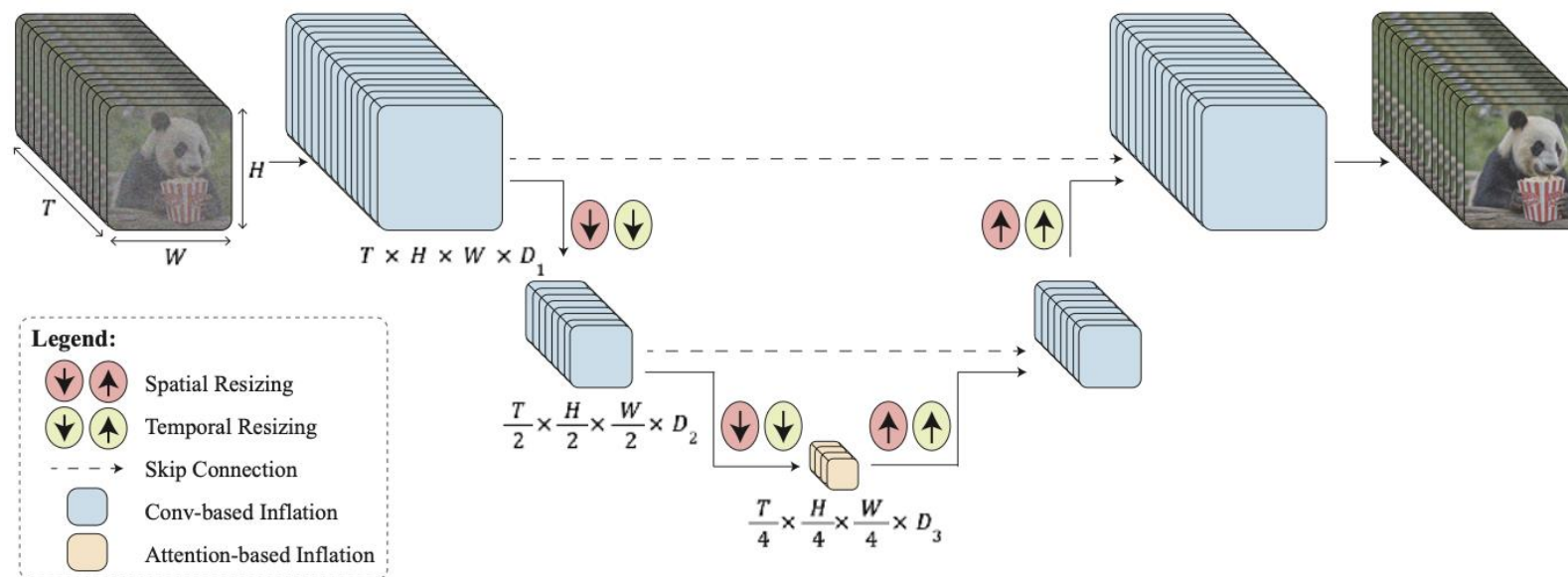
Generate whole frames at a low spatial resolution (128x128, 80 frames, 16 fps, 5 seconds)

Conduct spatial super-resolutions for overlapped video clips using the technique of multi-diffusions.

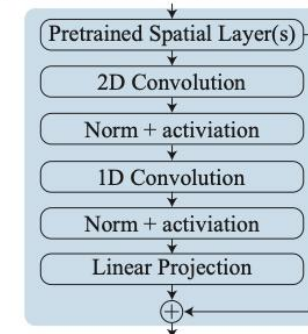


Space-Time UNet (STUNet)

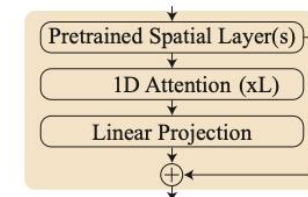
(a) Space-Time UNet (STUNet)



(b) Convolution-based Inflation Block

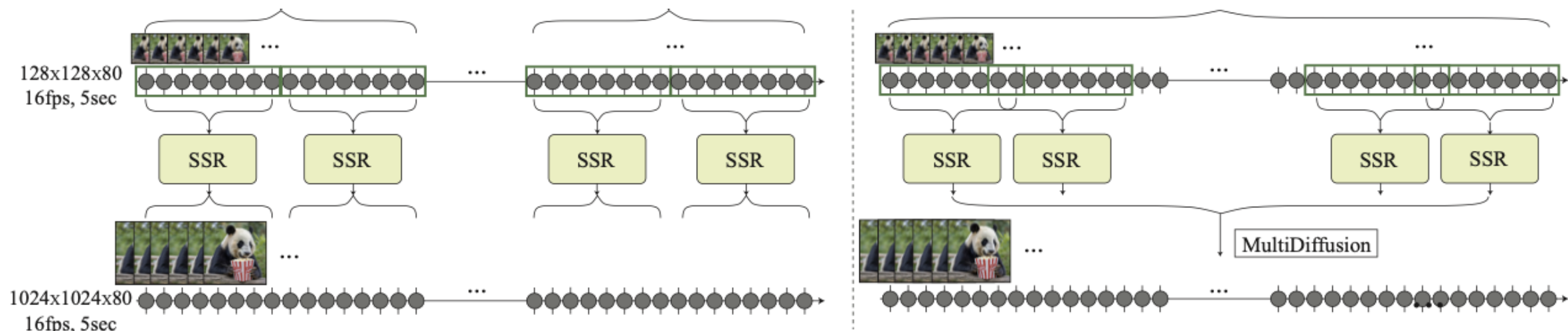


(c) Attention-based Inflation Block



- Temporal layers are added into the base model of Imagen that freeze its weights.
- For the computational efficiency, temporal attentions are used only for coarsest level.

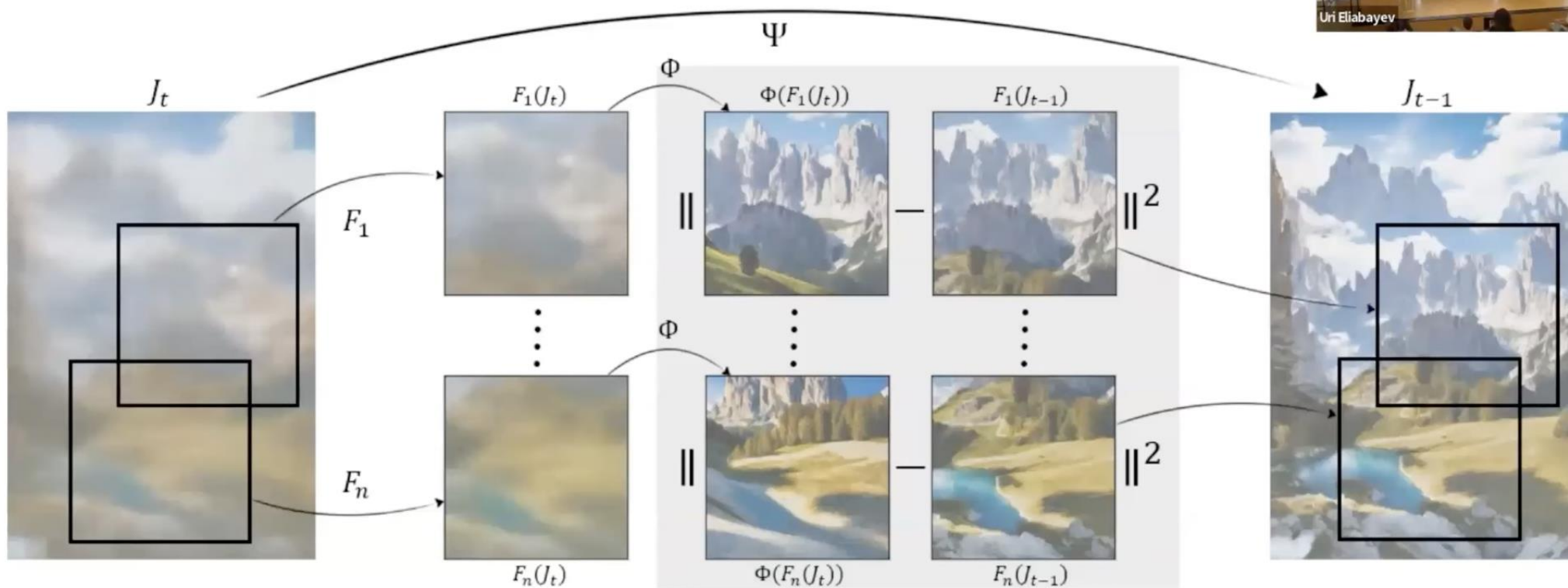
Common Approaches of T2V vs. Lumiere



- Due to memory constraints, the inflated SSR network can operate only on short segments of the video.
- To avoid temporal boundary artifacts, Lumiere employs Multidiffusion along the temporal axis for smooth transitions.

Space-Time UNet (STUNet)

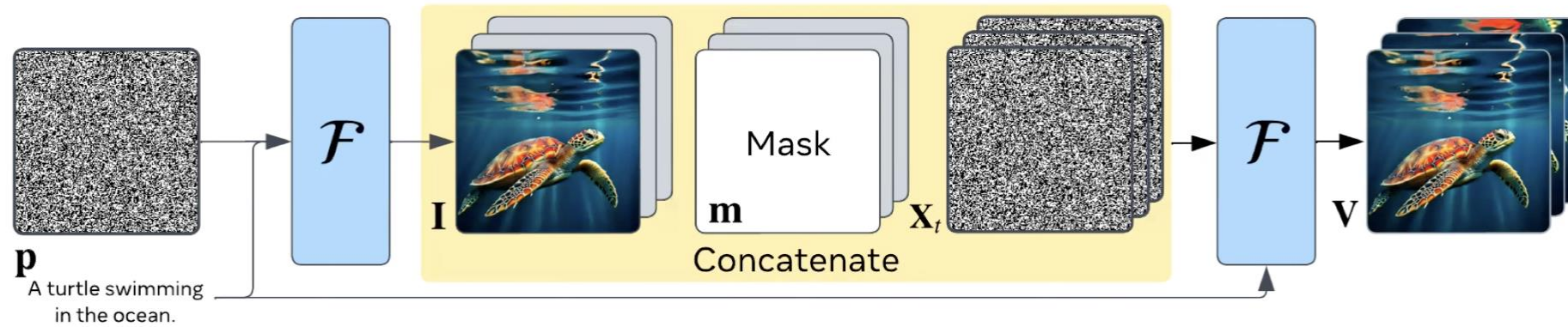
MultiDiffusion process: Panorama



Ψ as close as possible to Φ

$$\Psi(J_t) = \operatorname{argmin}_J \sum_{i=1}^n \|F_i(J) - \Phi(F_i(J_t))\|^2$$

Conditional Video Generation



- For conditional video generation, an input condition and a binary mask map are concatenated into the noised input of STUNet.
- Various applications can be implemented with different making patterns : Image-to-Video, Inpainting, Cinemagraphs

Conditional Video Generation

Text-to-Video



Golden retriever puppy in the park, autumn



Astronaut walking on the planet Mars



Beer pouring into glass



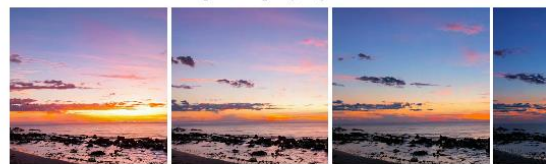
Epic tracking of a gorilla walking gracefully



Chocolate syrup pouring on a vanilla ice cream



A panda playing a ukulele at home



Sunset time lapse at the beach



Flying through a temple in ruins, epic, mist

Image-to-Video



A girl winking and smiling



Bigfoot walking through the woods in Northern California



Flying through an intense battle between pirate ships in a stormy ocean



Bee buzzing busily around a field of blooming wildflowers



A cat playing piano



A teddy bear running in New York City



Jack russell terrier dog snowboarding. GoPro shot

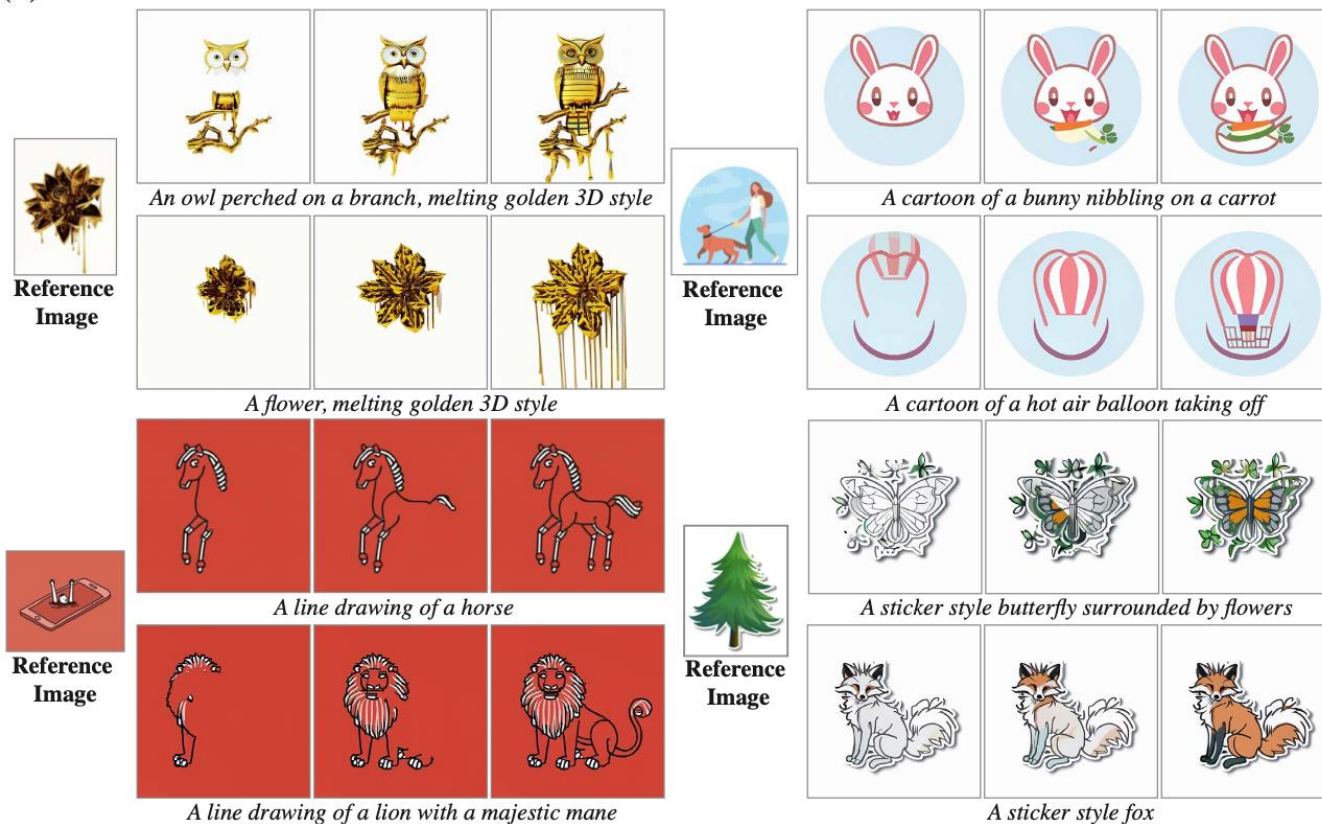


Ancient pharaoh singing and smiling and shaking his head like

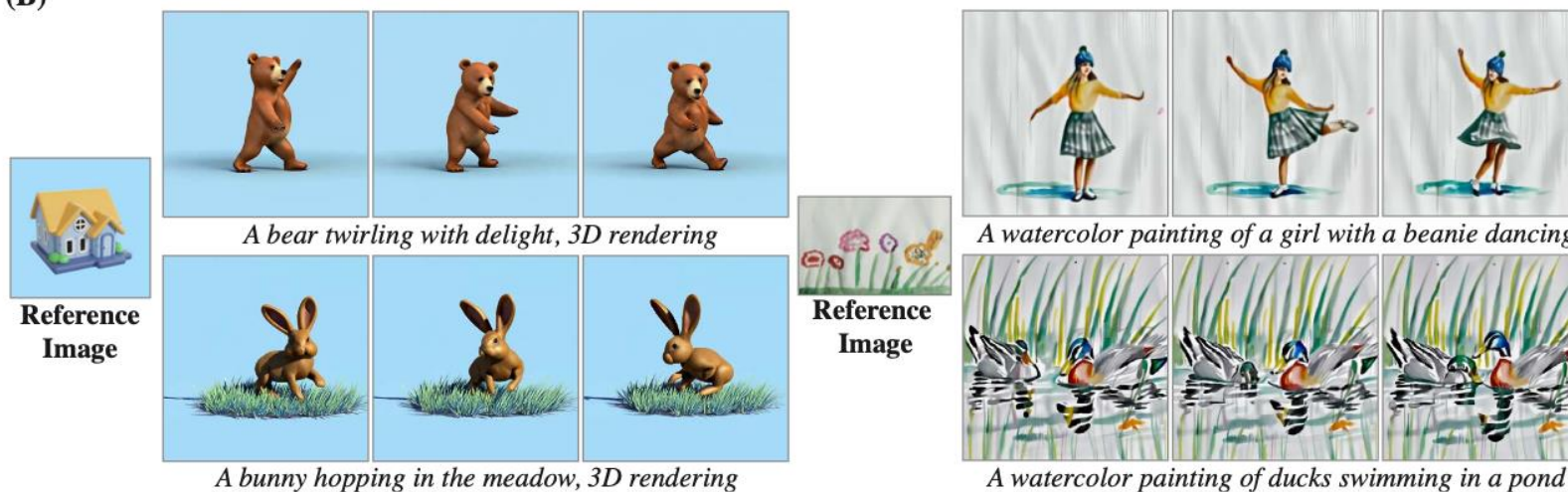
Stylized Generation

$$W_{\text{interpolate}} = \alpha \cdot W_{\text{style}} + (1 - \alpha) \cdot W_{\text{orig}}.$$

(A)



(B)



Video-to-Video via SDEdit

Input



Output



Made of colorful toy bricks



Sculpture made of flowers



Ultra high detail mech robot

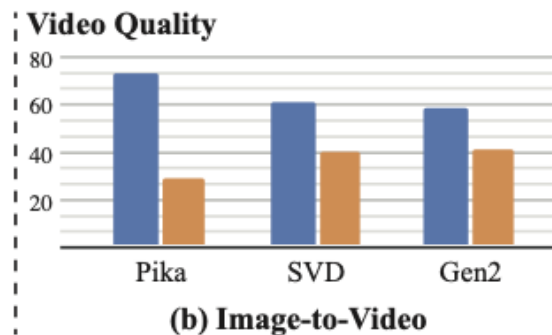
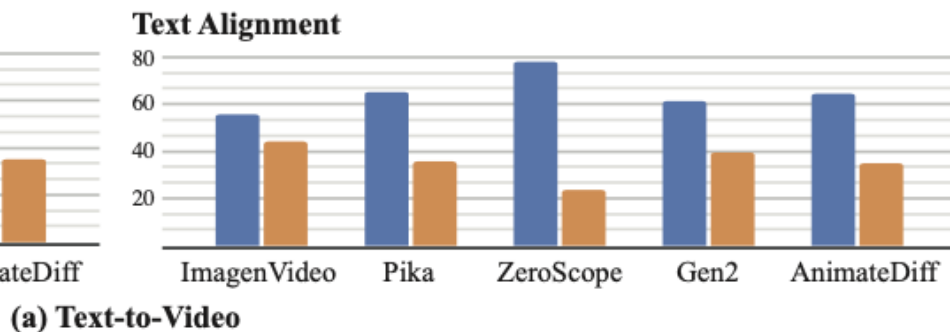
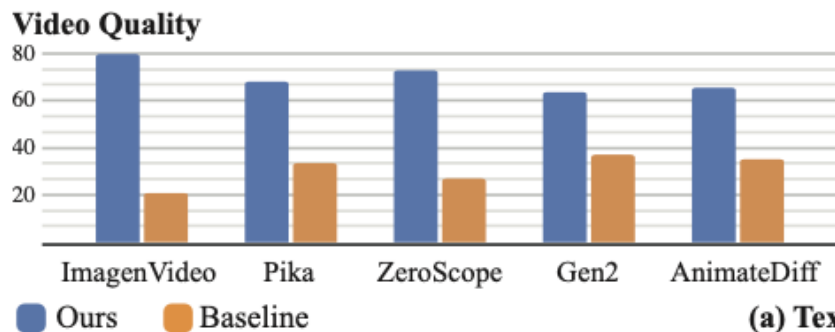


Made of stacked wooden blocks

Results

- Zero-shot text-to-video on UCF101.
- User study on 113 text prompts (~ 400 user judgements).

Method	FVD ↓	IS ↑
MagicVideo (Zhou et al., 2022)	655.00	-
Emu Video (Girdhar et al., 2023)	606.20	42.70
Video LDM (Blattmann et al., 2023b)	550.61	33.45
Show-1 (Zhang et al., 2023a)	394.46	35.42
Make-A-Video (Singer et al., 2022)	367.23	33.00
PYoCo (Ge et al., 2023)	355.19	47.76
SVD (Blattmann et al., 2023a)	242.02	-
Lumiere (Ours)	332.49	37.54





Results

Pika			
Gen2			
AnimateDiff			
ImagenVideo			
ZeroScope			
Ours			

A sheep to the right of a wine glass

Teddy bear skating in Times Square



Conclusion

- A new text-to-video generation framework is presented using a pretrained T2I model.
- To learn globally-coherent motion, a space-time U-Net architecture design enables to generate the whole frames at once with a low resolution.
- Lumiere shows State-of-the-art results and is applicable for various applications.