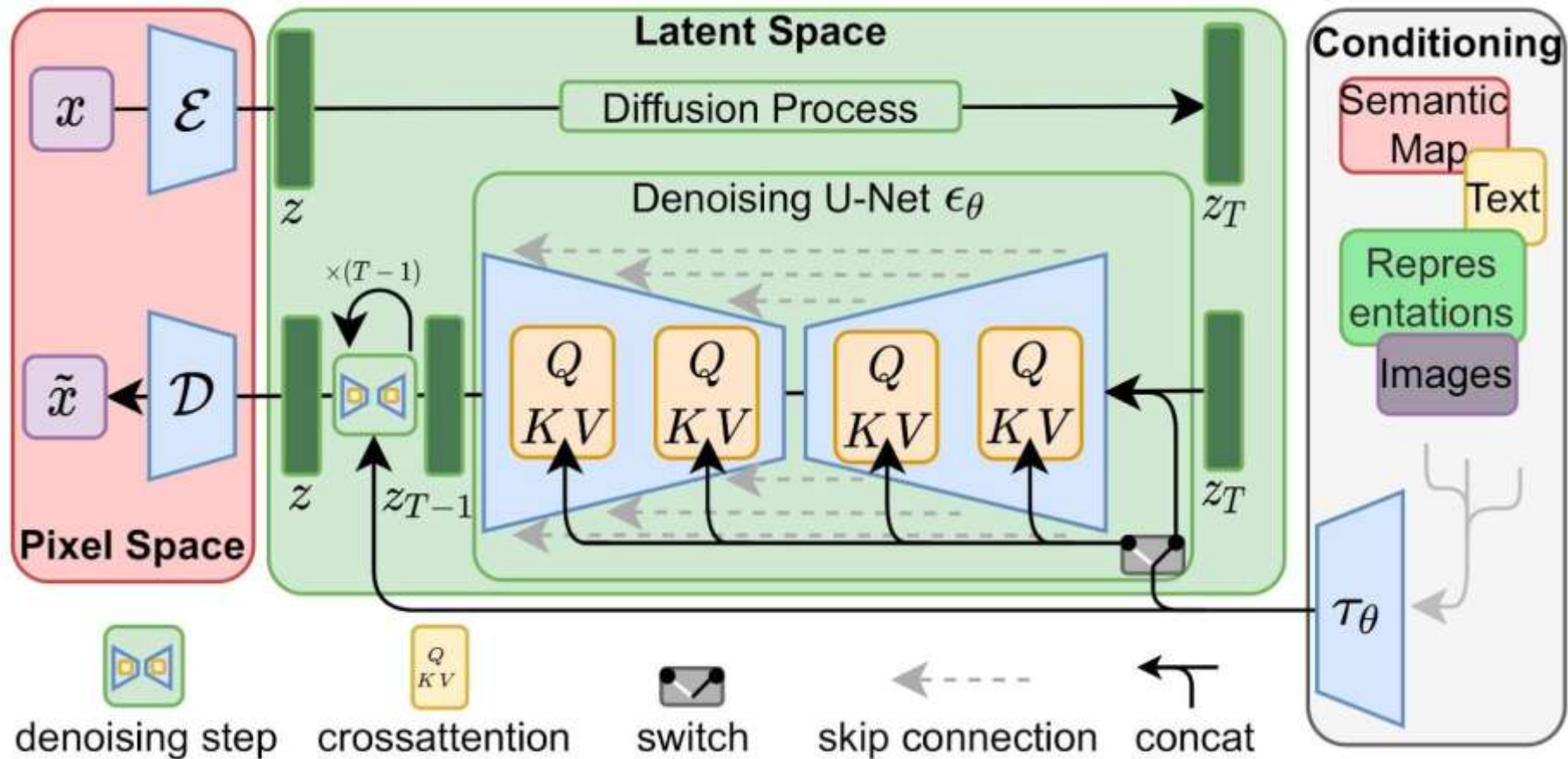
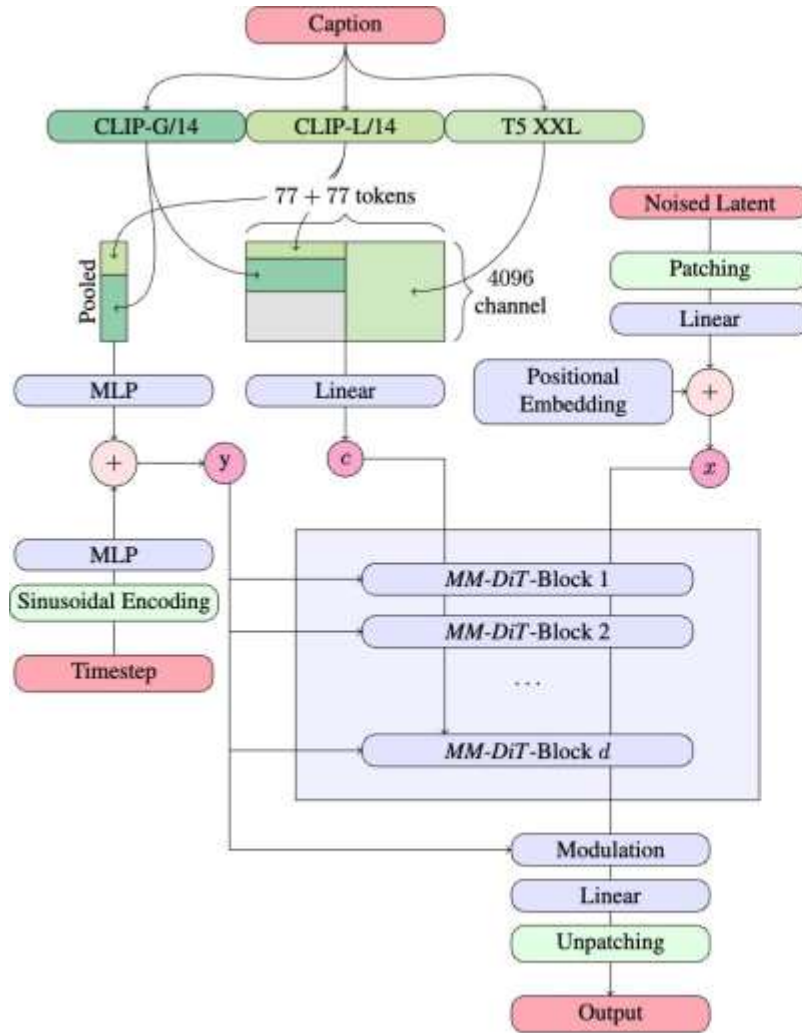


**Recap.**

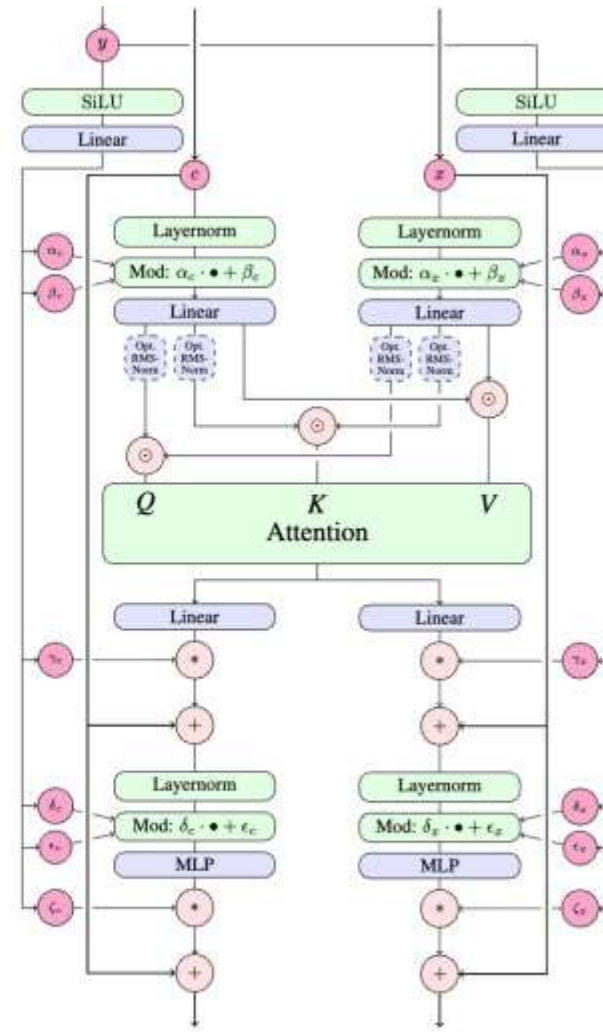
# Attempt 1: Cross Attention



# Attempt 2: Double Stream Multimodal-DiT

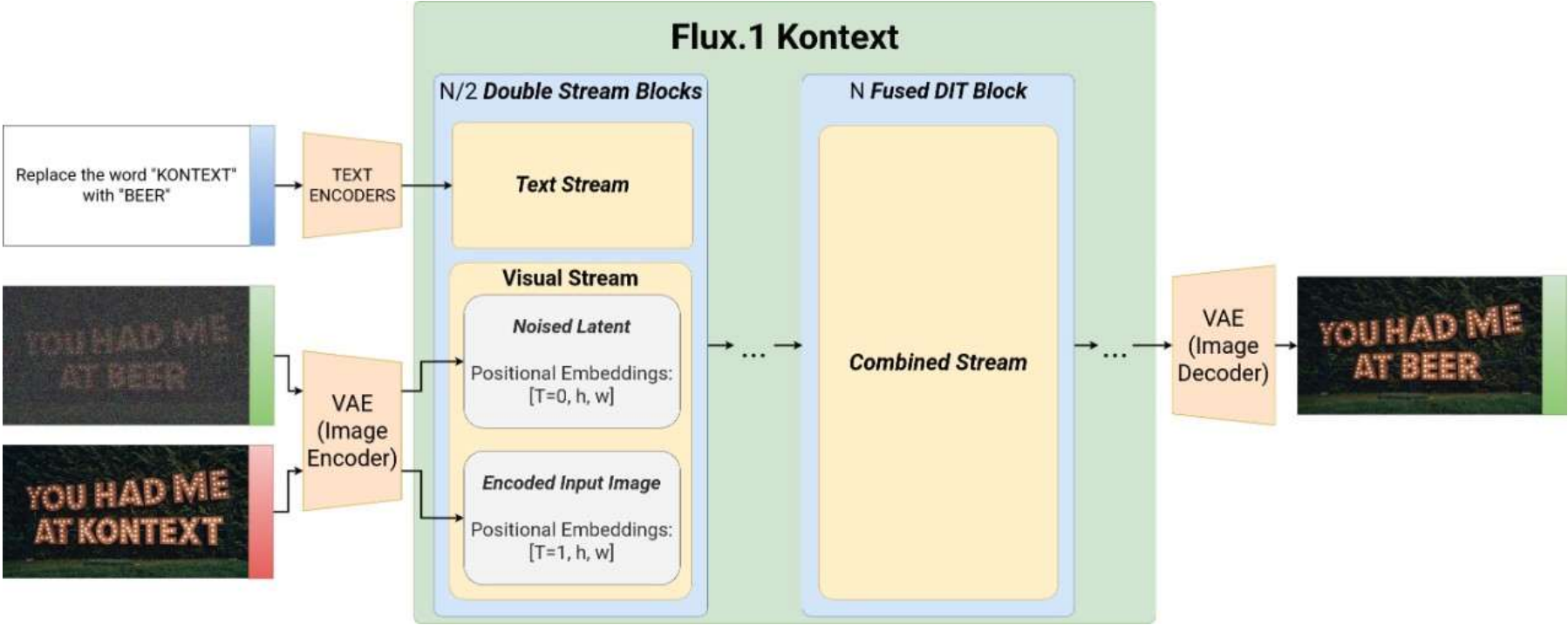


(a) Overview of all components.

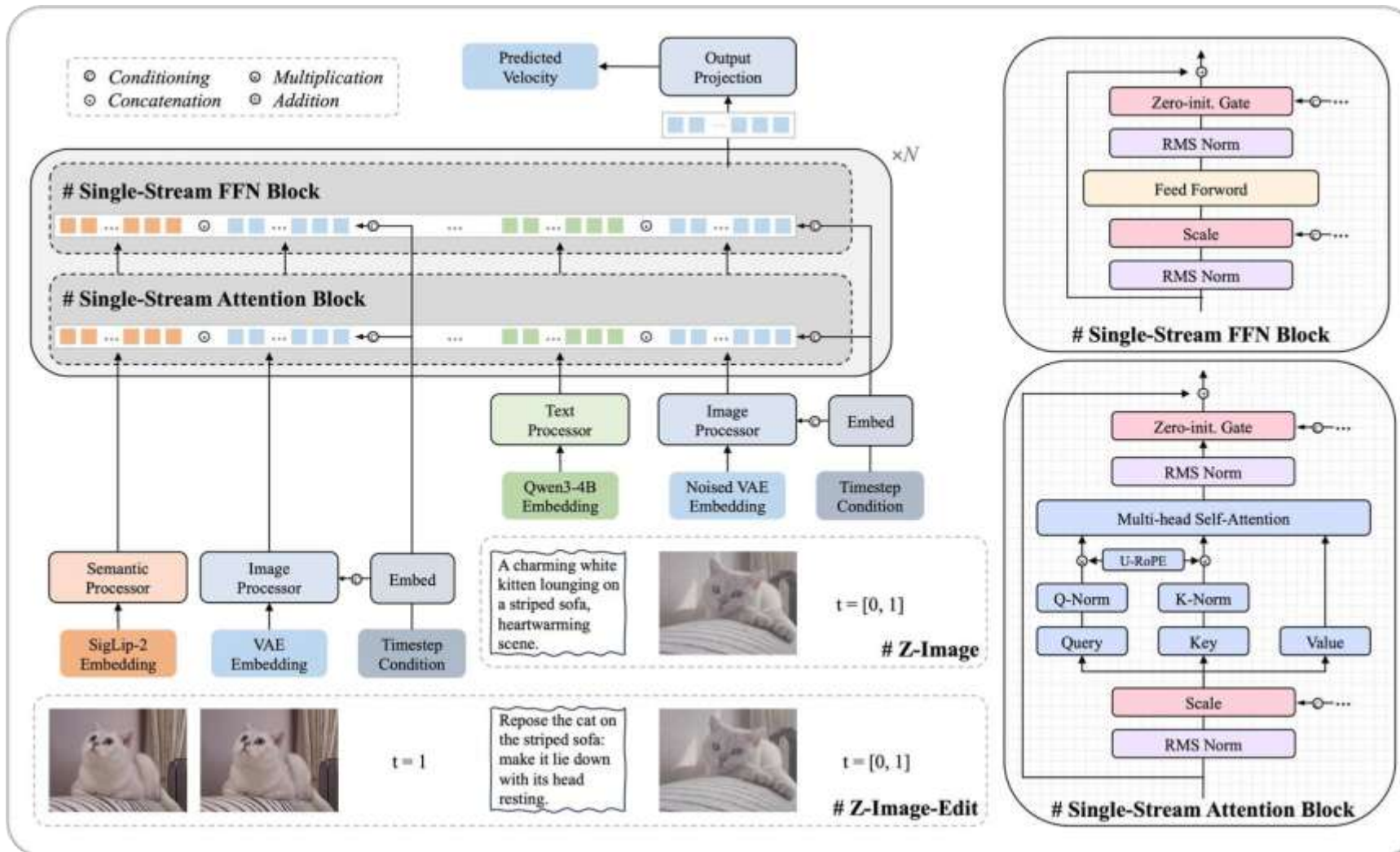


(b) One *MM-DiT* block

# Attempt 2: Double stream -> merged stream MM-DiT



# Attempt 3: Single stream MM-DiT



# Attempt 3.5: “Native multimodal model”

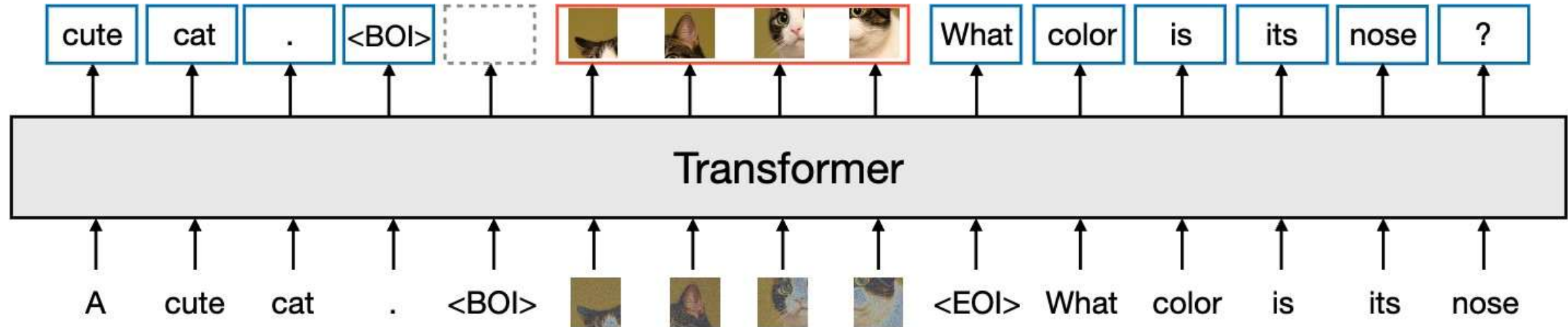
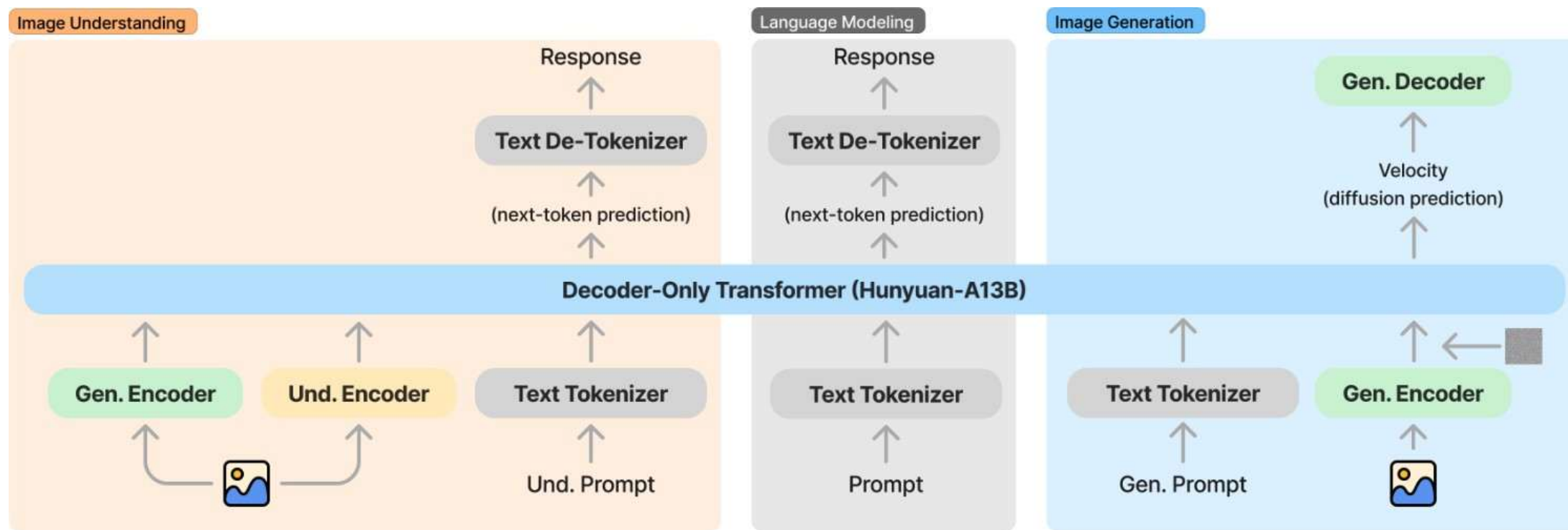
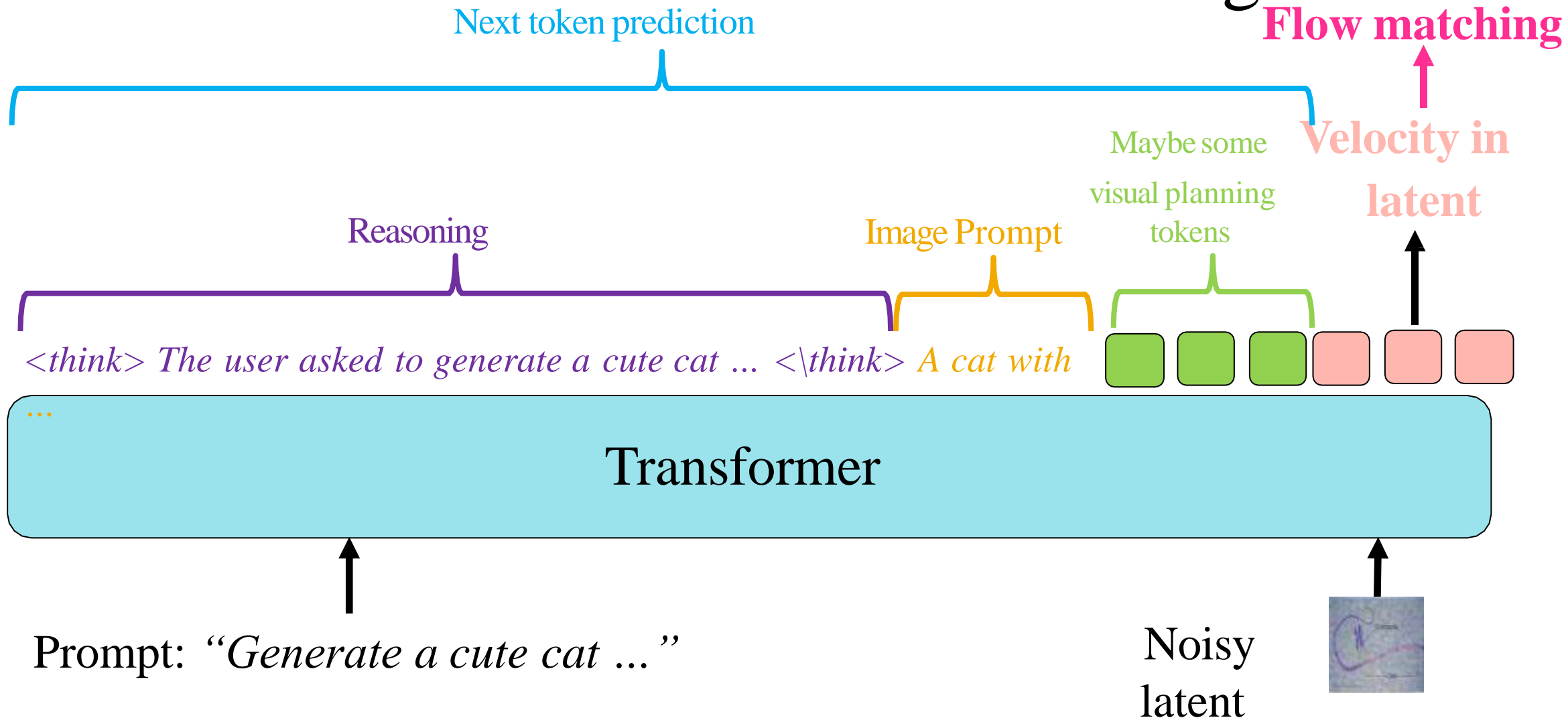


Figure 1: A high-level illustration of Transfusion. A single transformer perceives, processes, and produces data of every modality. Discrete (text) tokens are processed autoregressively and trained on the **next token prediction** objective. Continuous (image) vectors are processed together in parallel and trained on the **diffusion** objective. Marker BOI and EOI tokens separate the modalities.

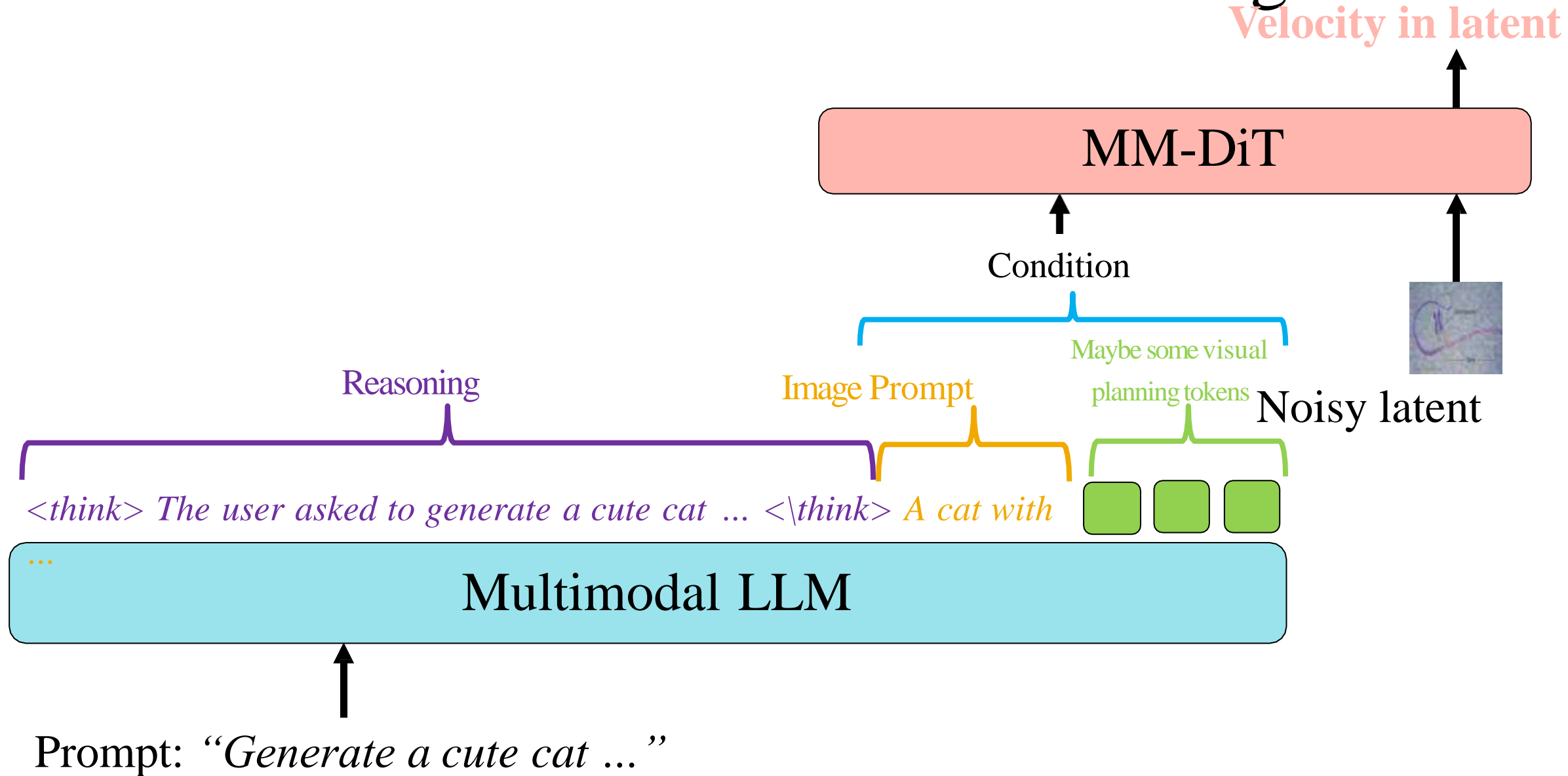
# Attempt 3.5: “Native multimodal model”



# Guess on how Nano Banana & GPT-4o Image work



# Guess on how Nano Banana & GPT-4o Image work



## The design space of text-to-image generation

### Training

- Training paradigm
  - DDPM
  - Flow matching

### Model

- Latent Space
  - VQ-VAE/VQGAN
  - Advanced VAE
- Model architecture
  - U-Net
  - DiT

### Text Encoding

- Text Encoder
  - CLIP
  - CLIP + T5
  - LLM/VLM/MLLM
- Text Conditioning
  - Cross Attention
  - MM-DiT
  - Native MM

## The design space of text-to-image generation

### Training

- Training paradigm
  - DDPM
  - Flow matching

### Model

- Latent Space
  - VQ-VAE/VQGAN
  - Advanced VAE
- Model architecture
  - U-Net
  - DiT

### Text Encoding

- Text Encoder
  - CLIP
  - CLIP + T5
  - LLM/VLM/MLLM
- Text Conditioning
  - Cross Attention
  - MM-DiT
  - Native MM

**This is Stable Diffusion 1 & 2!**

## The design space of text-to-image generation

### Training

- Training paradigm
  - DDPM
  - Flow matching

### Model

- Latent Space
  - VQ-VAE/VQGAN
  - Advanced VAE
- Model architecture
  - U-Net
  - DiT

### Text Encoding

- Text Encoder
  - CLIP
  - CLIP + T5
  - LLM/VLM/MLLM
- Text Conditioning
  - Cross Attention
  - MM-DiT
  - Native MM

This is Stable Diffusion 3 & Flux 1!

## The design space of text-to-image generation

### Training

- Training paradigm
  - DDPM
  - Flow matching

### Model

- Latent Space
  - VQ-VAE/VQGAN
  - Advanced VAE
- Model architecture
  - U-Net
  - DiT

### Text Encoding

- Text Encoder
  - CLIP
  - CLIP+ T5
  - LLM/VLM/MLLM
- Text Conditioning
  - Cross Attention
  - MM-DiT
  - Native MM

This is Flux 2, Z-Image, Qwen-Image, etc!

## The design space of text-to-image generation

### Training

- Training paradigm
  - DDPM
  - Flow matching

### Model

- Latent Space
  - VQ-VAE/VQGAN
  - Advanced VAE
- Model architecture
  - U-Net
  - DiT

### Text Encoding

- Text Encoder
  - CLIP
  - CLIP + T5
  - LLM/VLM/MLLM
- Text Conditioning
  - Cross Attention
  - MM-DiT
  - Native MM

This is Transfusion, Hunyuan 3.0, etc!

## The design space of text-to-image generation

### Training

- Training paradigm
  - DDPM
  - Flow matching

### Model

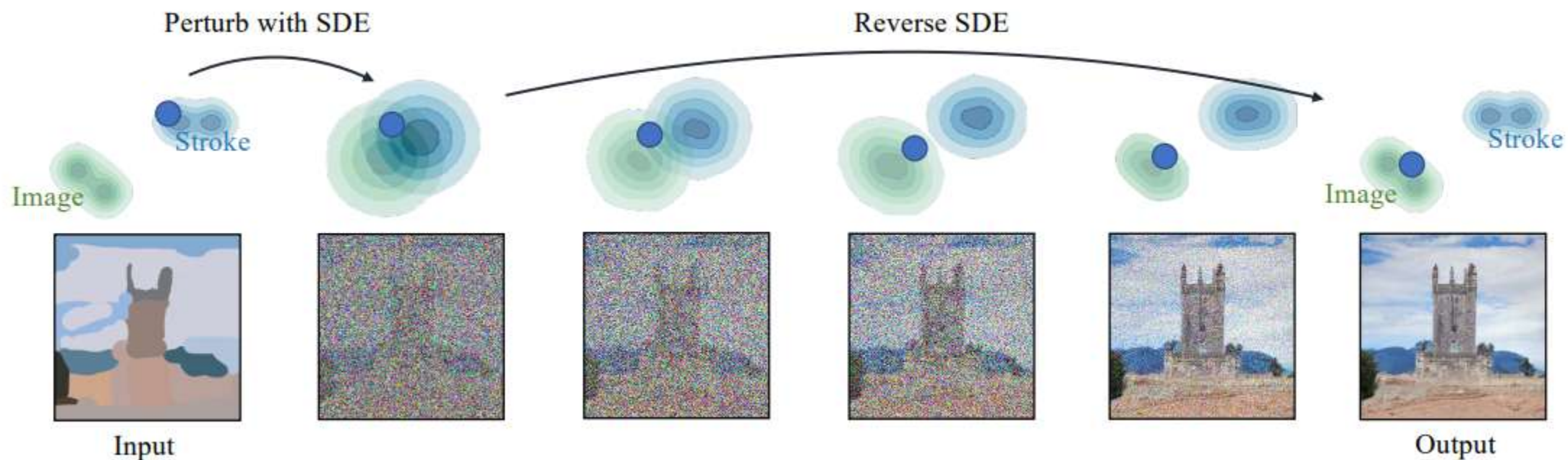
- Latent Space
  - VQ-VAE/VQGAN
  - Advanced VAE
- Model architecture
  - U-Net
  - DiT

### Text Encoding

- Text Encoder
  - CLIP
  - CLIP + T5
  - LLM/VLM/MLLM
- Text Conditioning
  - Cross Attention
  - MM-DiT
  - Native MM

(Probably also Nano Banana & GPT-4o Image)

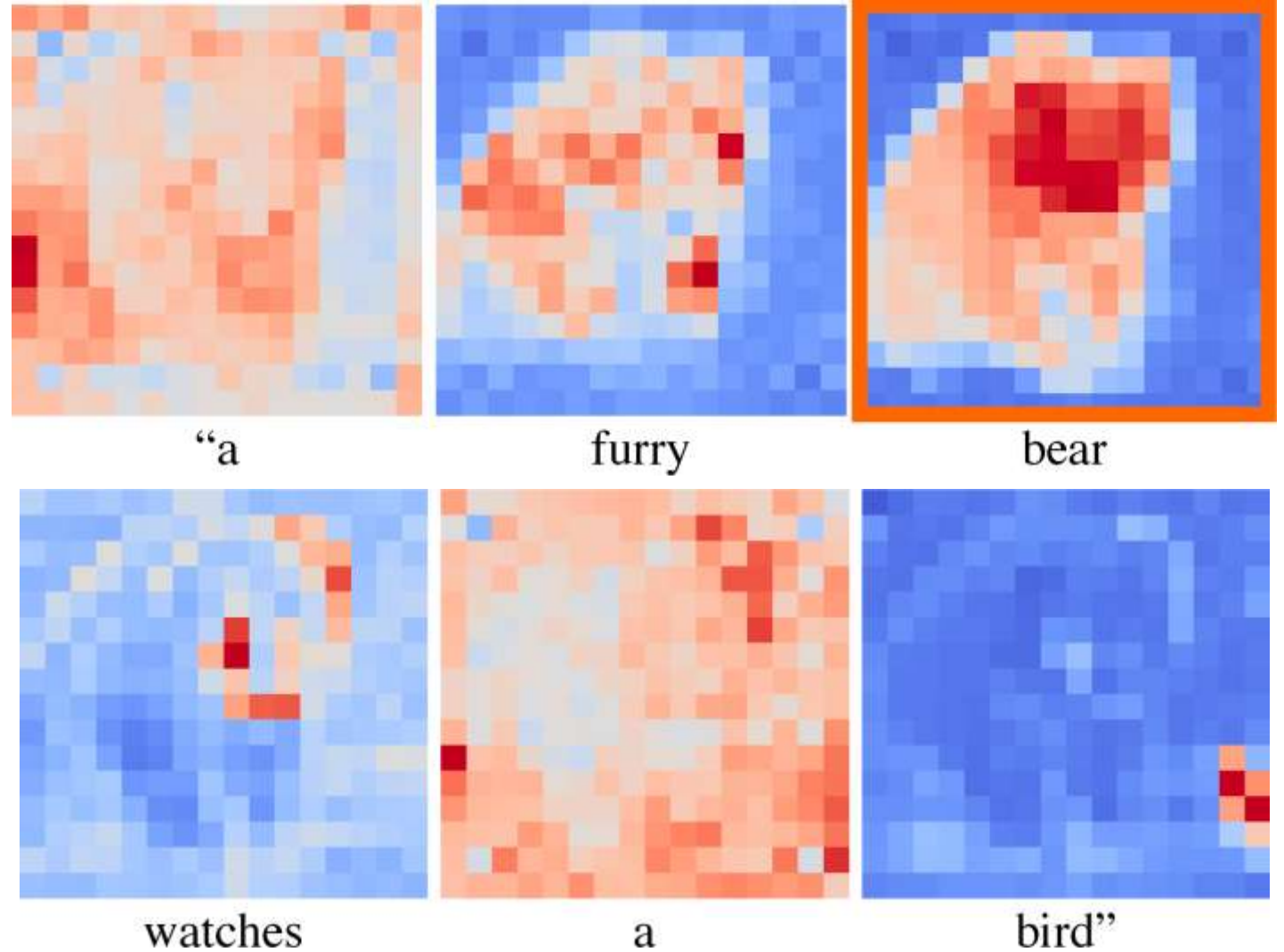
# SDEdit: Guided Image Synthesis with Diffusion



# Attention visualization



“a furry bear  
watches a bird”



# Prompt-to-Prompt

“Photo of a cat riding on a bicycle.”



source image



cat → dog



cat → chicken

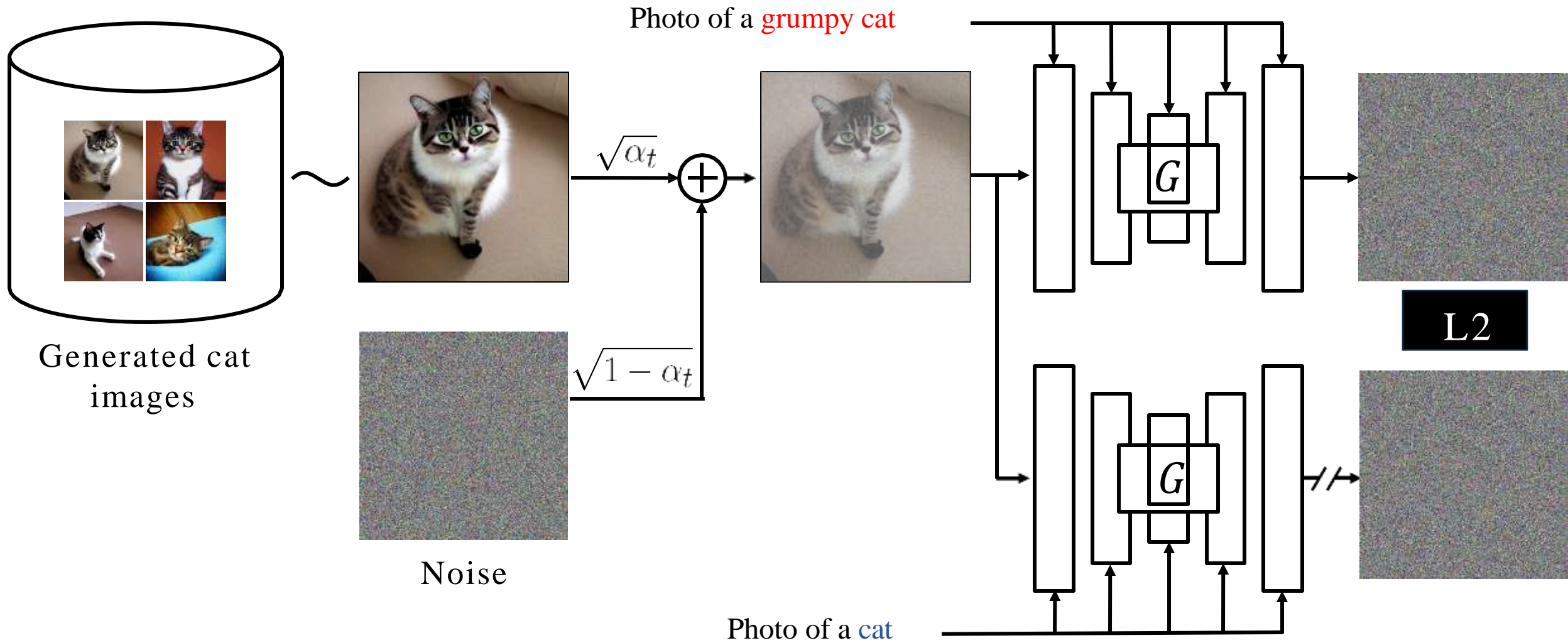


cat → squirrel



cat → elephant

# Final Method



Ablated



Target removed

Photo of a grumpy cat  
Target concept

Ablated



Nearby preserved

Photo of a british shorthair cat  
Nearby concept

# **Video Generation**

# Sora

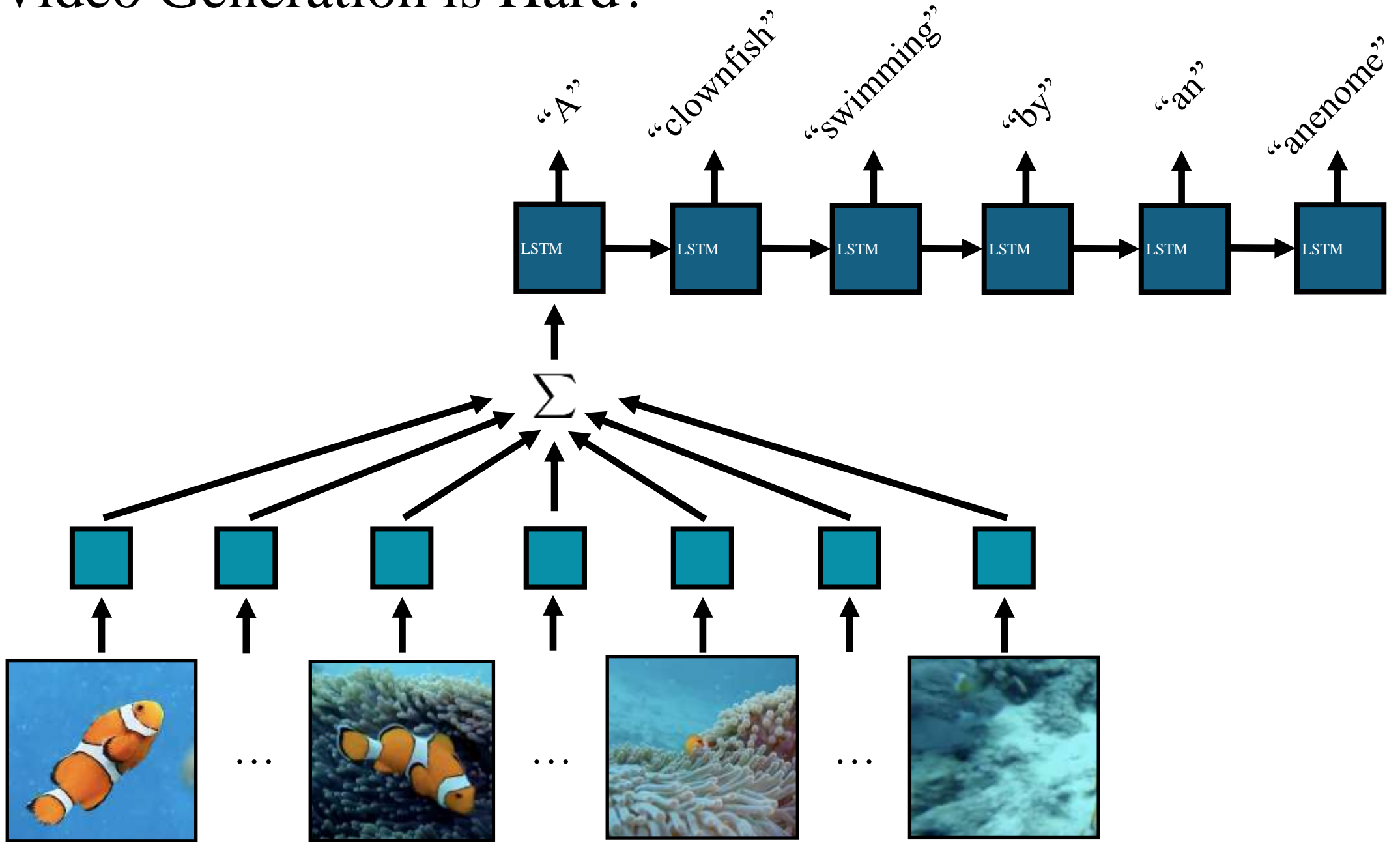


# Why Video Generation is Hard?

Outputs

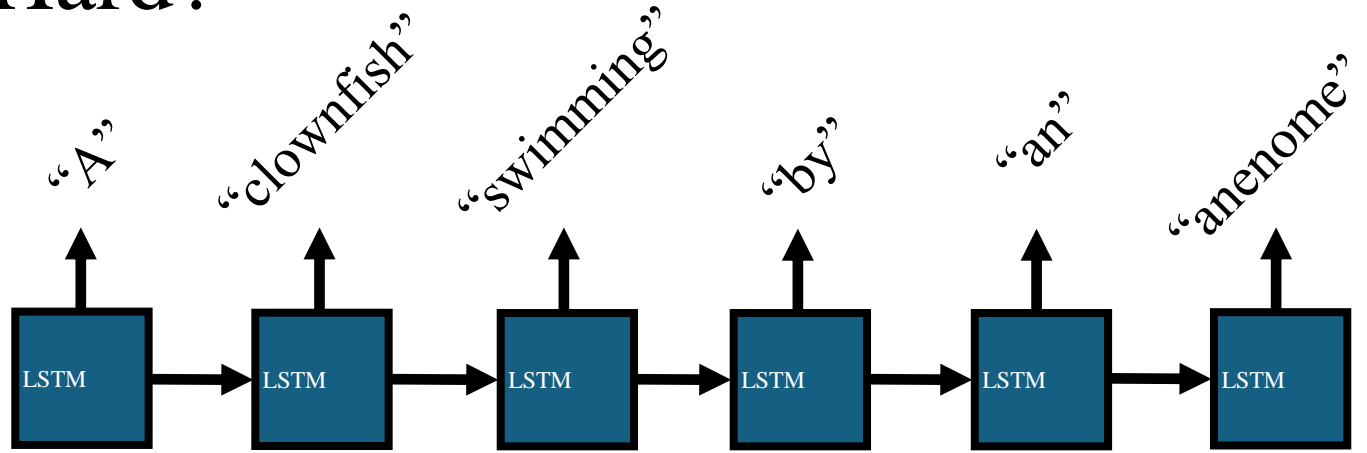
Hidden

Input



# Why Video Generation is Hard?

Outputs



$\Sigma$

Hidden

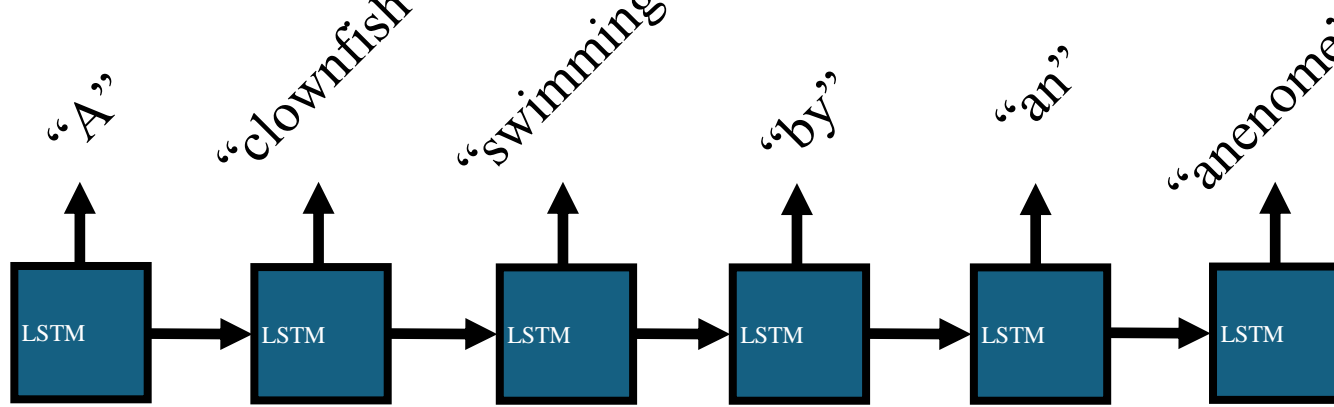


Input



# Why Video Generation is Hard?

Outputs



$\Sigma$

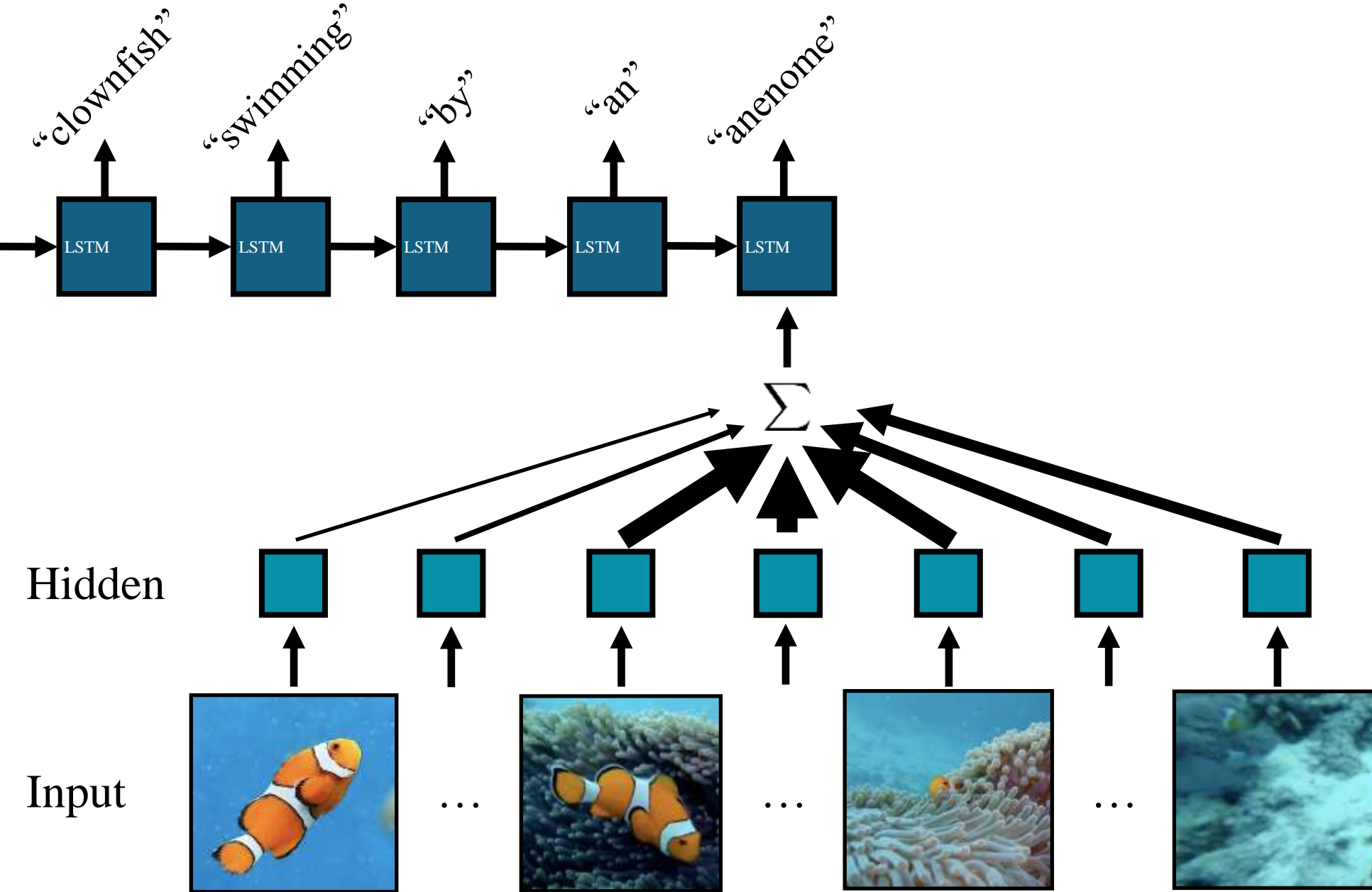
Hidden



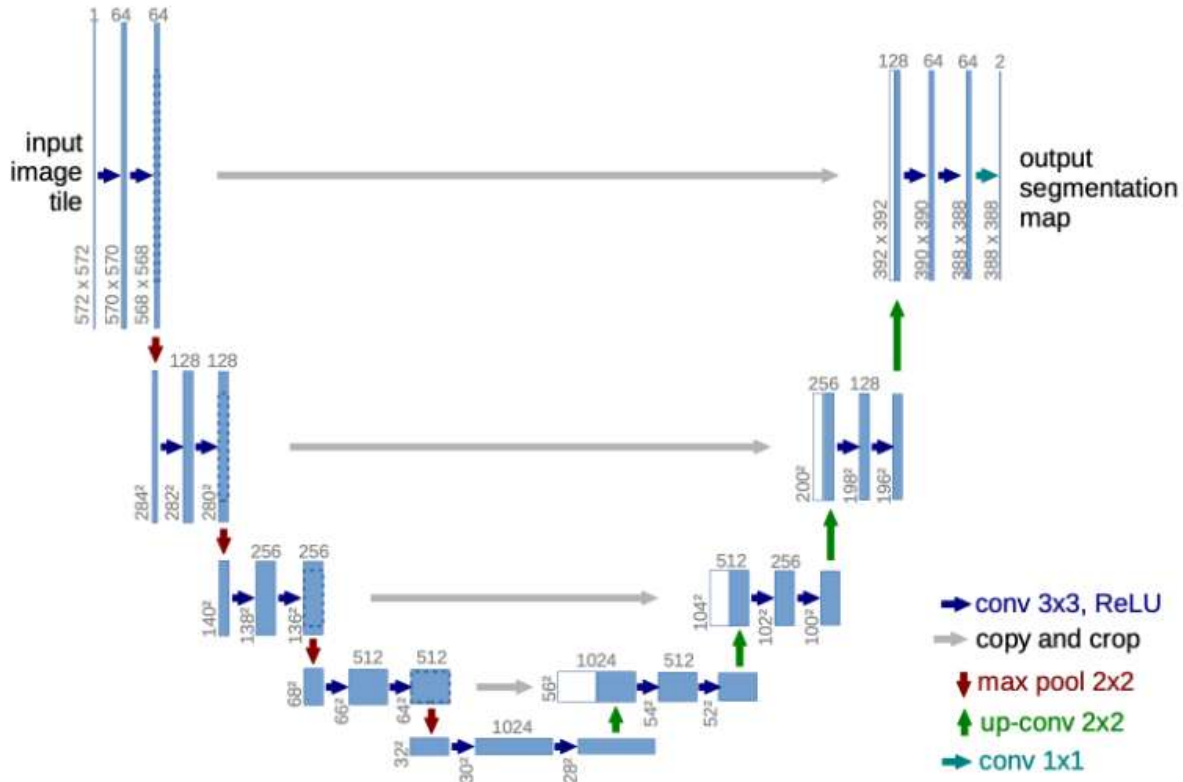
Input



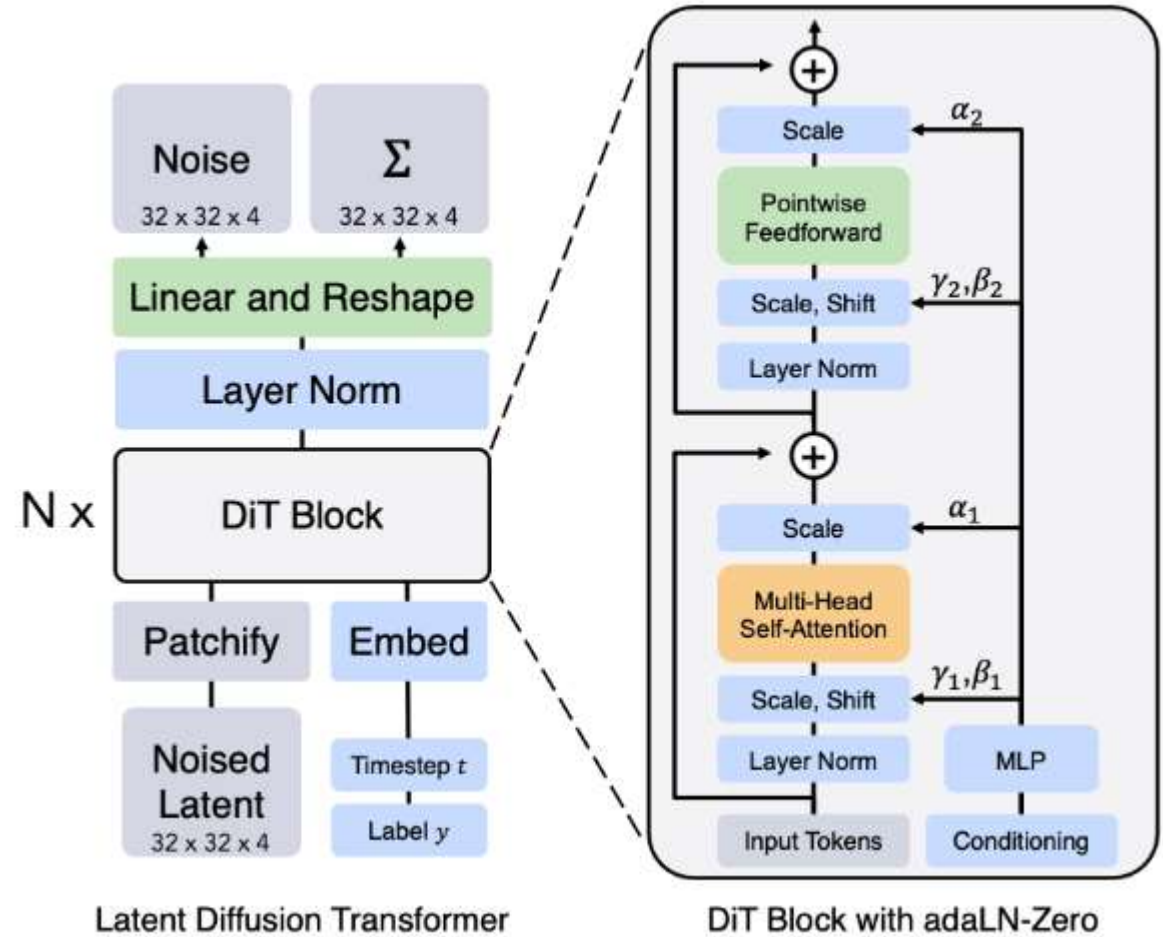
# Why Video Generation is Hard?



# Model Architecture



U-Net: <https://arxiv.org/abs/1505.04597>

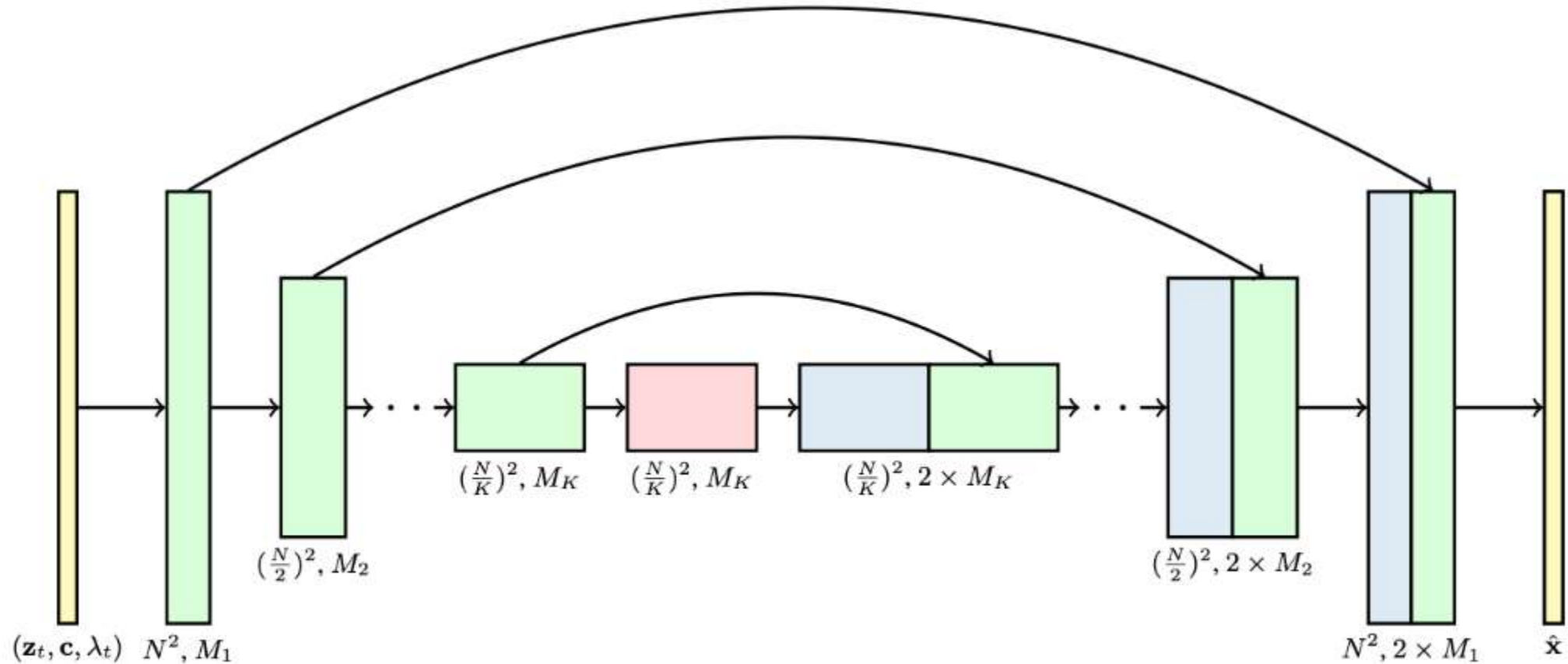


Latent Diffusion Transformer

DiT Block with adaLN-Zero

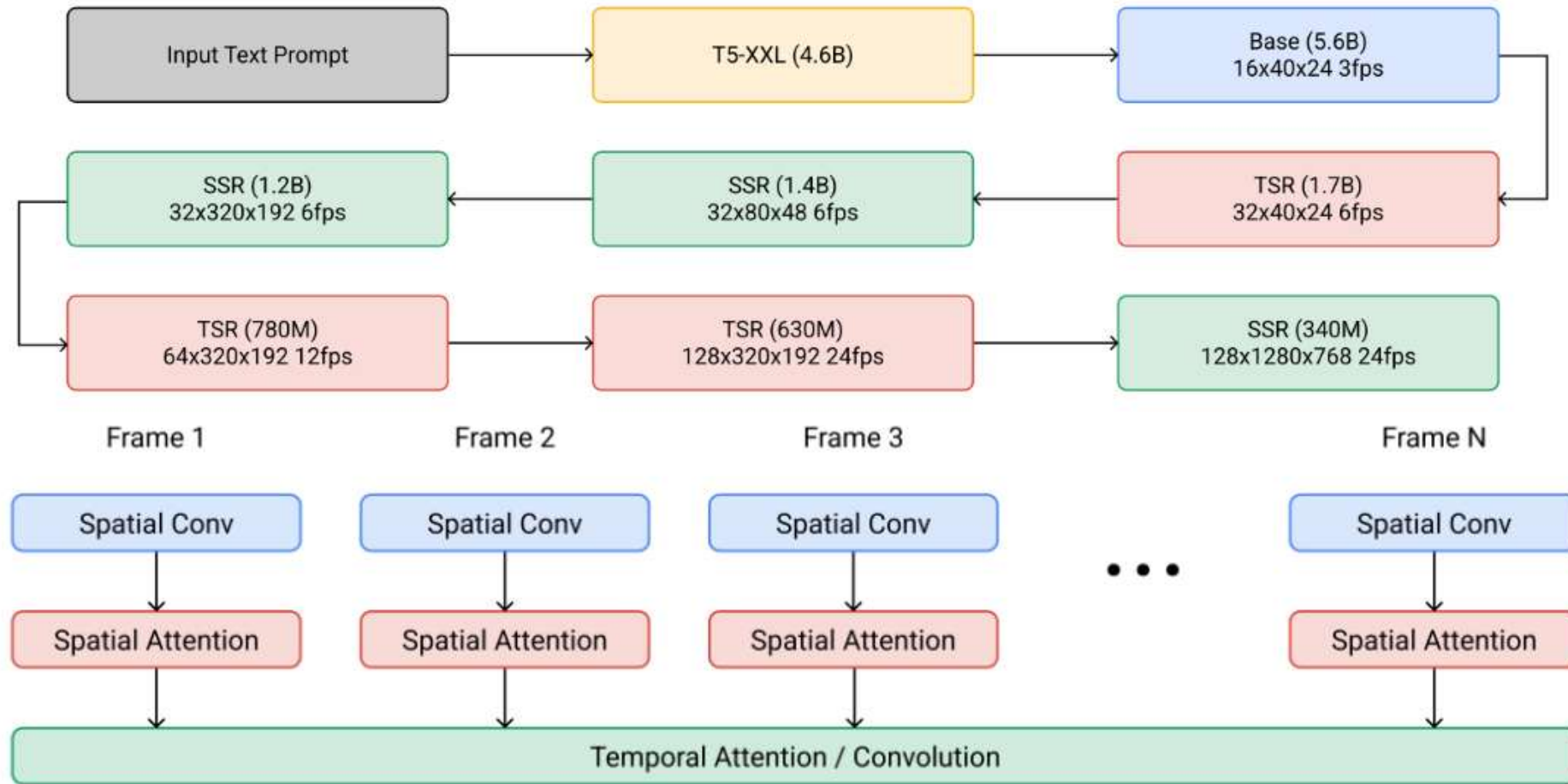
DiT: <https://arxiv.org/abs/2212.09748>

# Video Diffusion Models



The 3D U-net architecture. The noisy video  $\mathbf{z}_t$ , conditioning information  $\mathbf{c}$  and the log signal-to-noise ratio (log-SNR)  $\lambda_t$  are inputs to the network. The channel multipliers  $M_1, \dots, M_K$  represent the channel counts across layers.

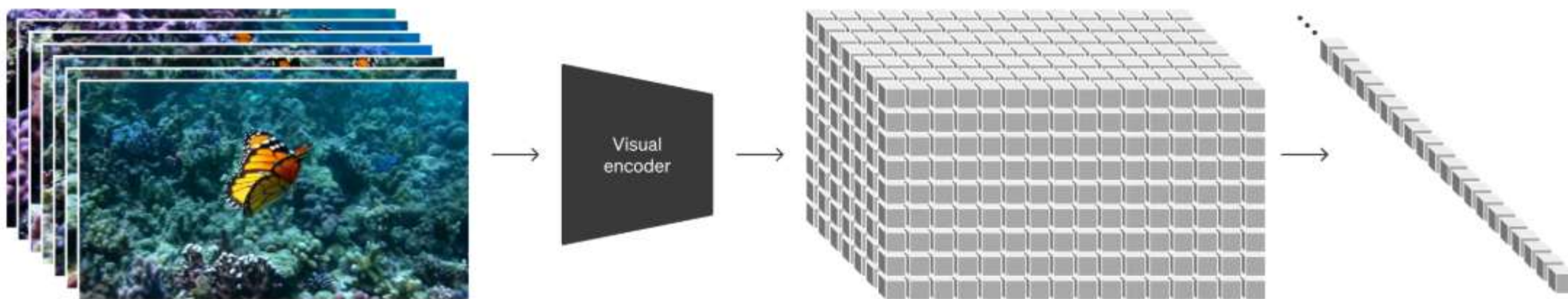
# Imagen Video



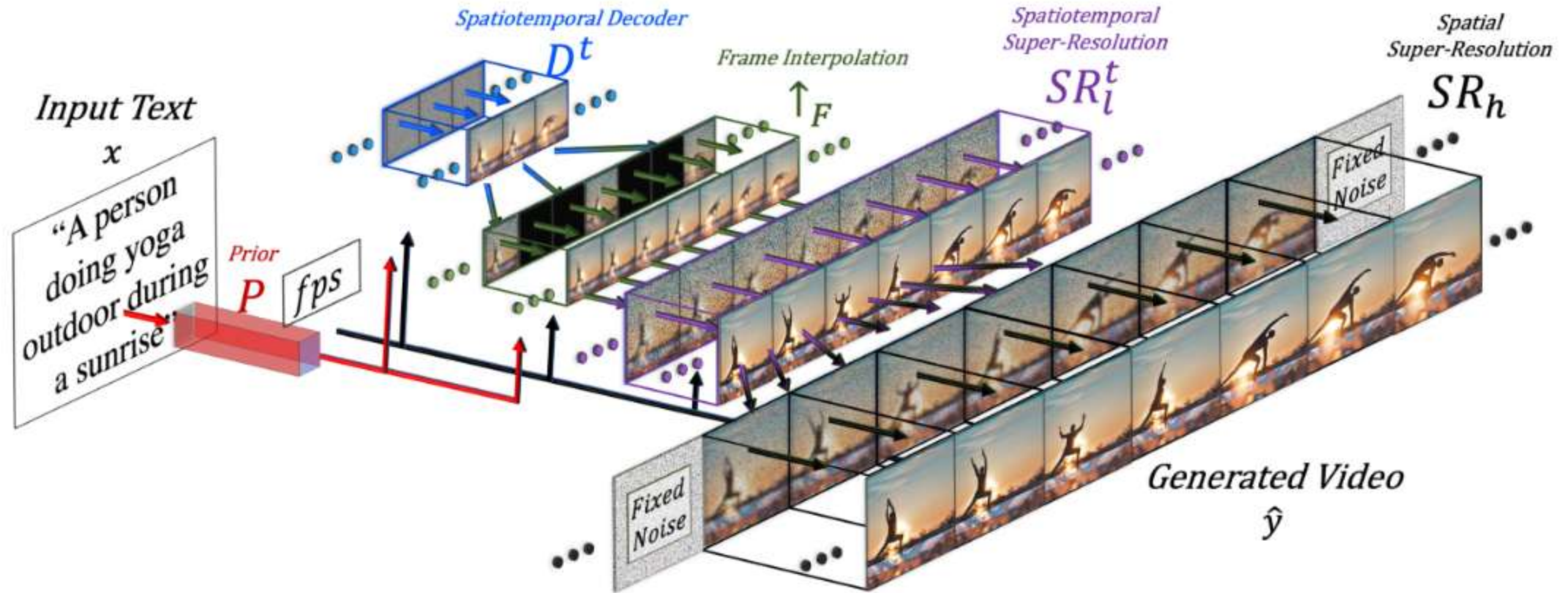
Upper: The cascaded sampling pipeline in Imagen Video. In practice, the text embeddings are injected into all components, not just the base model.

Below: The architecture of one space-time separable block in the Imagen Video diffusion model.

# Sora: DiT for Video Generation



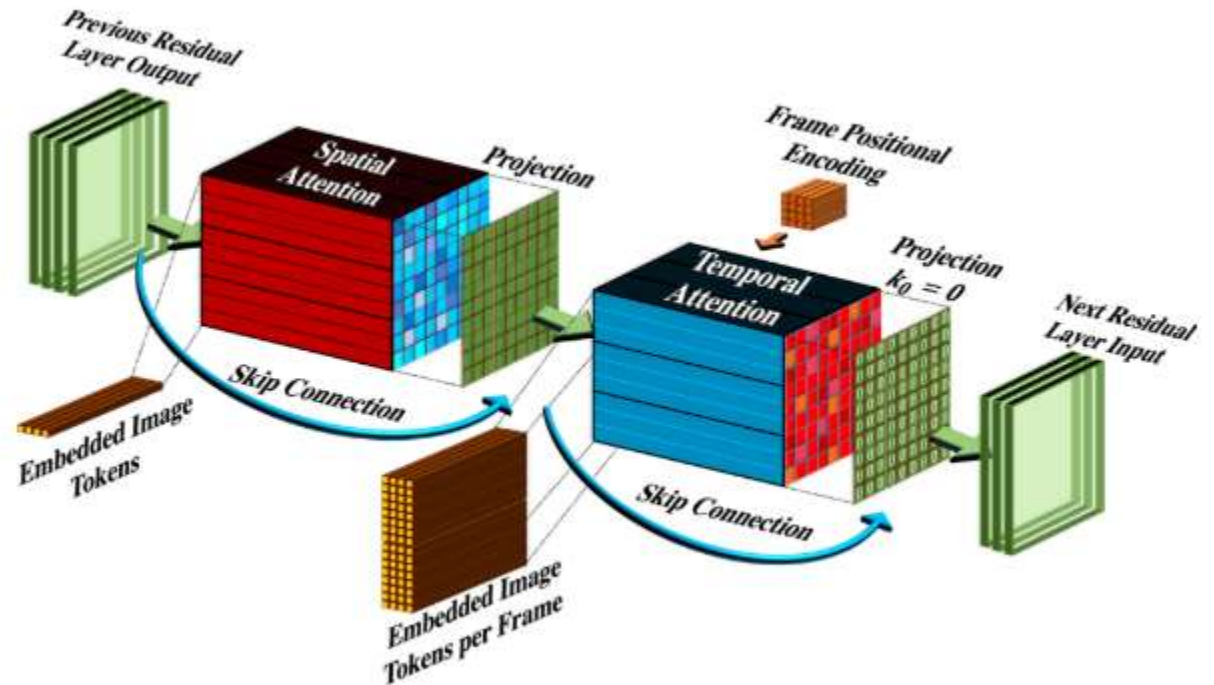
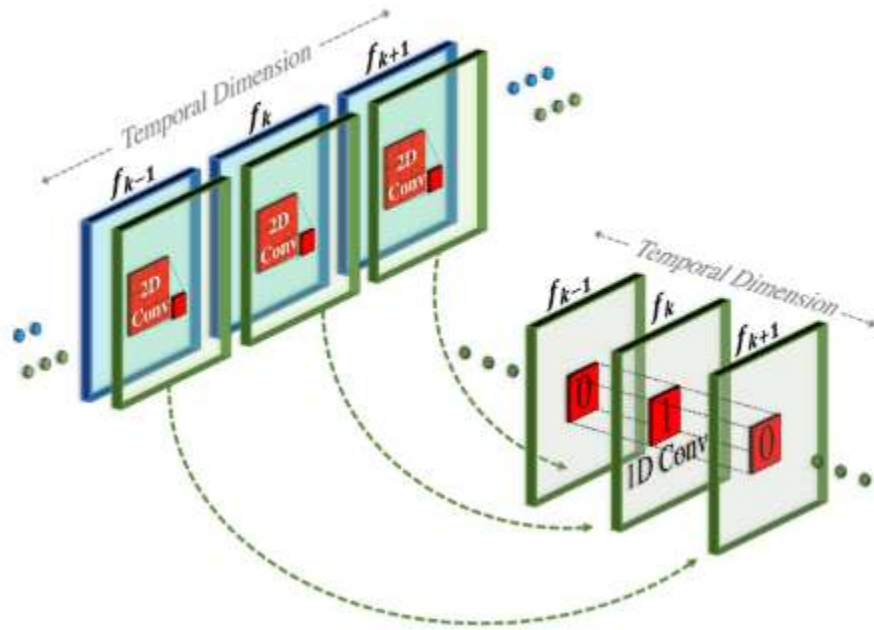
# Adapting Image Models to Generate Videos



$$\hat{y}_t = SR_h \circ SR_l^t \circ \uparrow_F \circ D^t \circ P \circ (\hat{x}, CLIP_{\text{text}}(x))$$

# Adapting Image Models to Generate Videos

pseudo-3D convo layers and pseudo-3D attention layers



$$\text{Conv}_{P3D} = \text{Conv}_{1D}(\text{Conv}_{2D}(\mathbf{h}) \circ T) \circ T$$

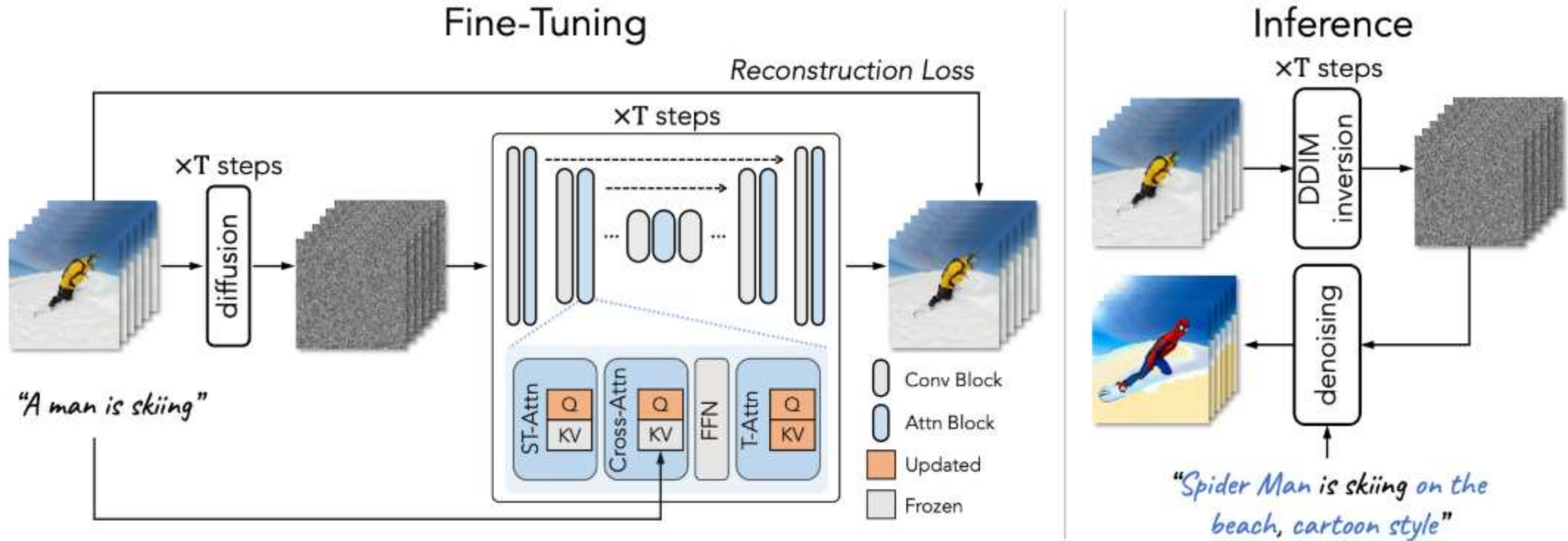
$$\text{Attn}_{P3D} = \text{flatten}^{-1}(\text{Attn}_{1D}(\text{Attn}_{2D}(\text{flatten}(\mathbf{h})) \circ T) \circ T)$$

# Adapting Image Models to Generate Videos



<https://ai.meta.com/blog/generative-ai-text-to-video/>

# Adapting Image Models to Generate Videos



$$\mathbf{Q} = \mathbf{W}^Q \mathbf{z}_{v_i}, \quad \mathbf{K} = \mathbf{W}^K [\mathbf{z}_{v_1}, \mathbf{z}_{v_{i-1}}], \quad \mathbf{V} = \mathbf{W}^V [\mathbf{z}_{v_1}, \mathbf{z}_{v_{i-1}}]$$

$$\mathbf{O} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right) \cdot \mathbf{V}$$

# Adapting Image Models to Generate Videos



"A man is skiing"

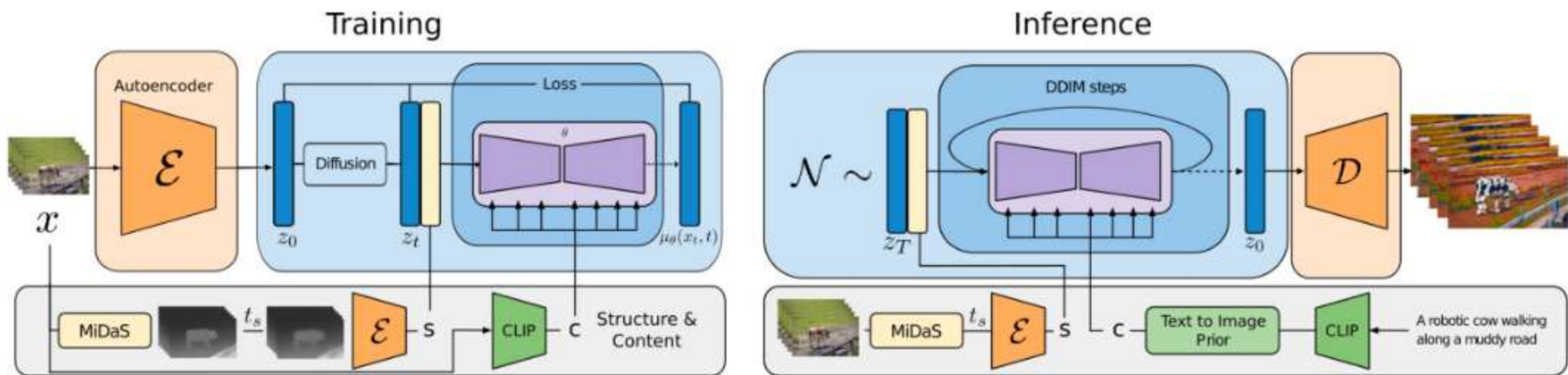


"Spider Man is skiing on the beach, cartoon style"



"A man, wearing pink clothes, is skiing at sunset"

# Adapting Image Models to Generate Videos



# Adapting Image Models to Generate Videos



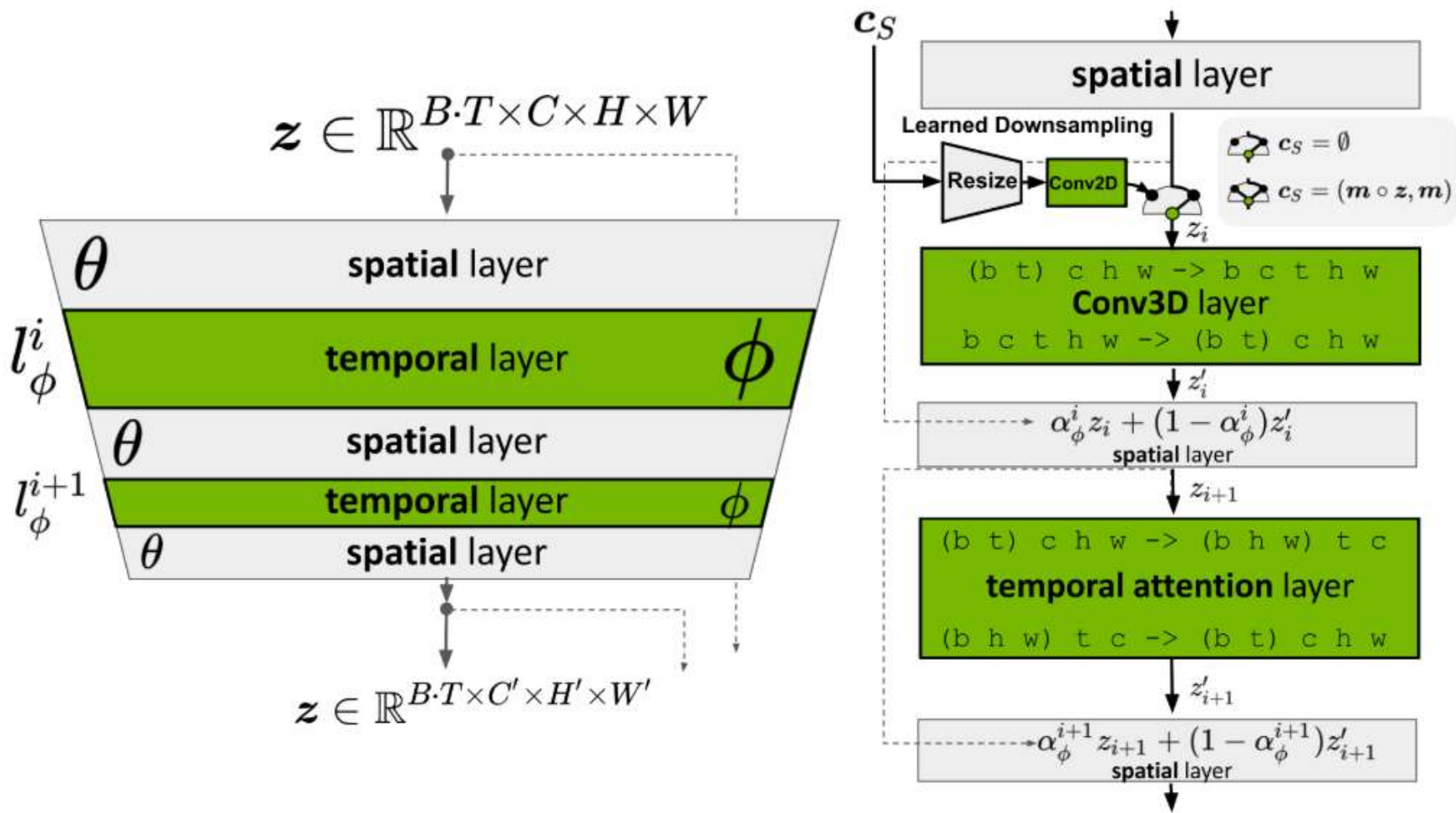
# Recap.



Bilibili BV1Lg4y1t797 with ControlNet

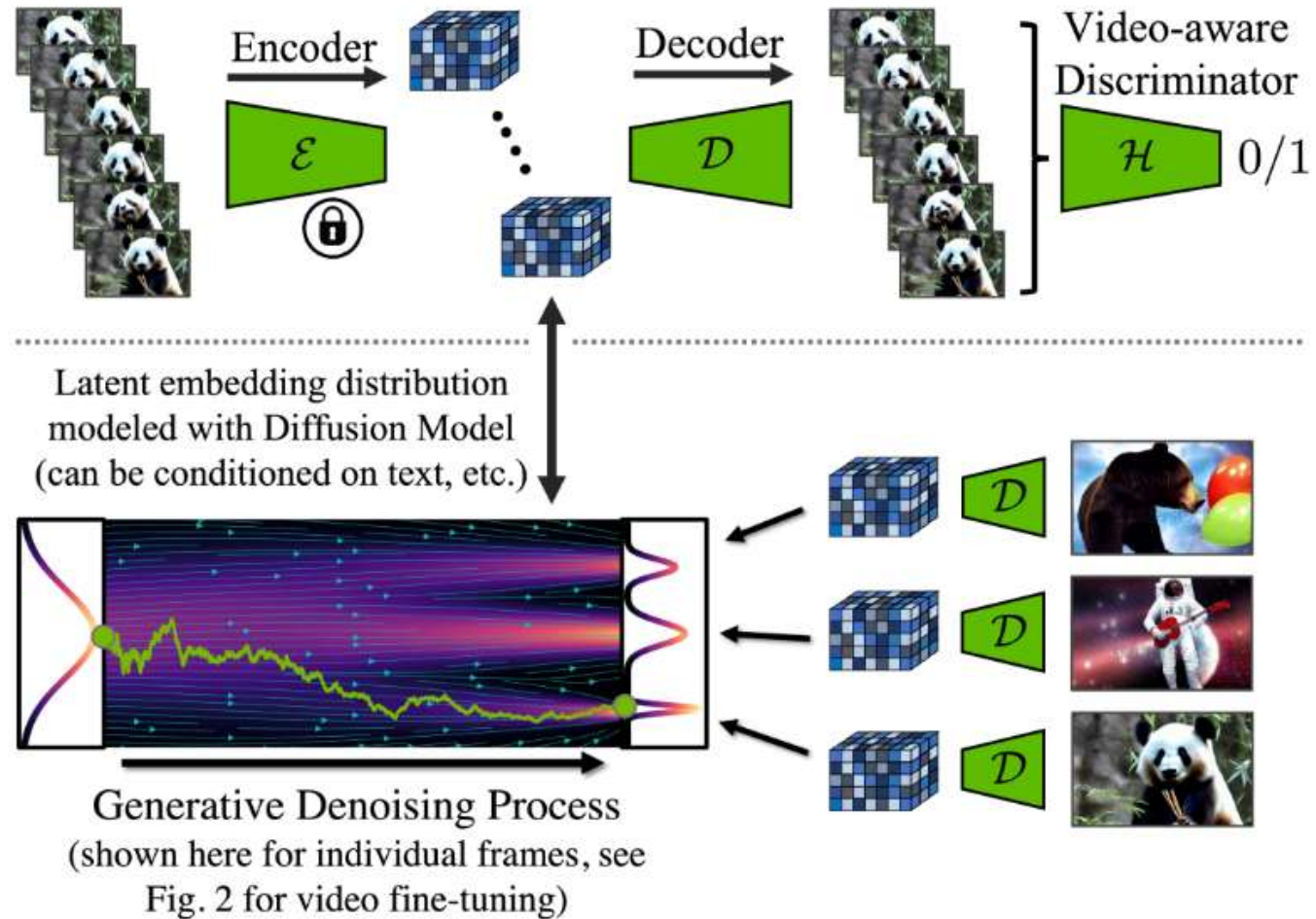
<https://arxiv.org/abs/2302.05543>

# Adapting Image Models to Generate Videos



Base image model  $\theta$

# Adapting Image Models to Generate Videos



# Adapting Image Models to Generate Videos

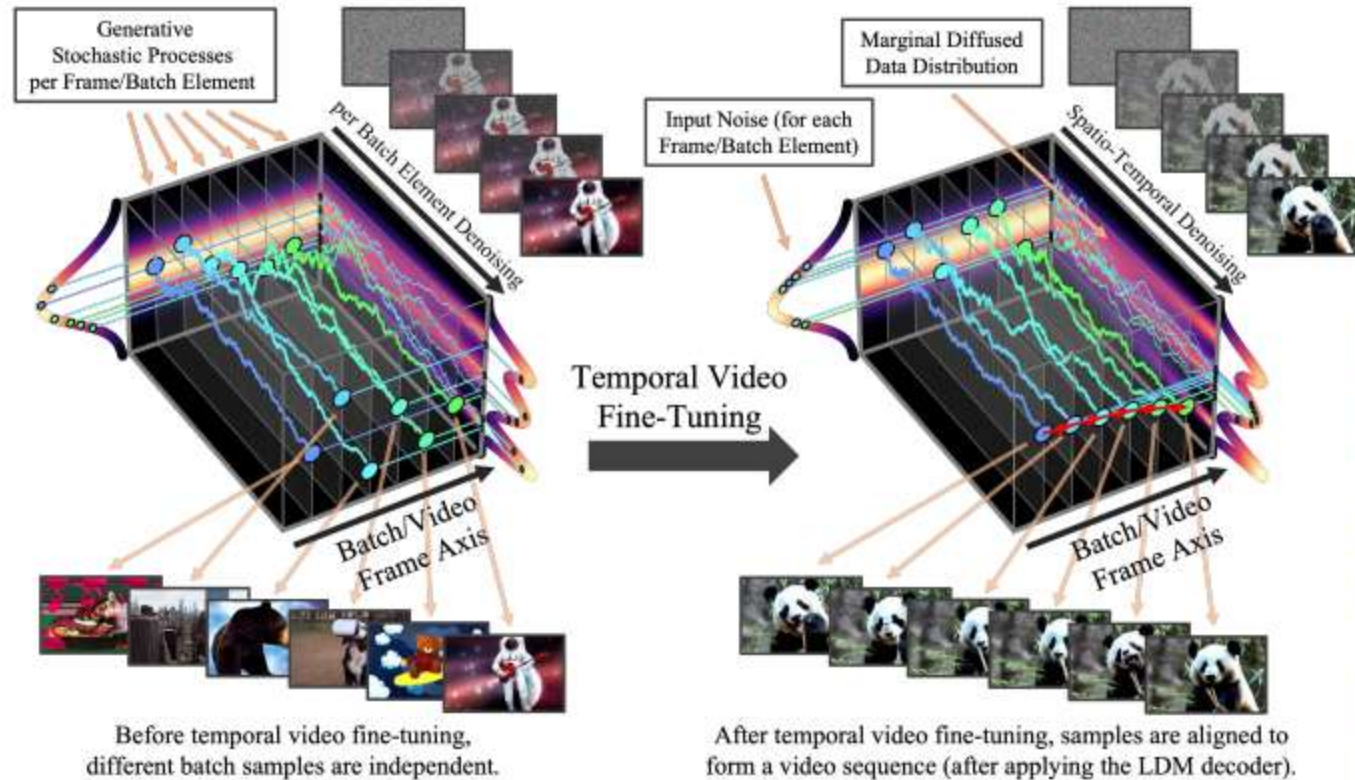


Figure 2. **Temporal Video Fine-Tuning.**

We turn pre-trained image diffusion models into temporally consistent video generators. Initially, different samples of a batch synthesized by the model are independent. After temporal video fine-tuning, the samples are temporally aligned and form coherent videos. The stochastic generation process before and after fine-tuning is visualised for a diffusion model of a one-dim. toy distribution. For clarity, the figure corresponds to alignment in pixel space. In practice, we perform alignment in LDM's latent space and obtain videos after applying LDM's decoder (see Fig. 3). We also video fine-tune diffusion model up-samplers in pixel or latent space (Sec. 3.4).

# Adapting Image Models to Generate Videos



# Adapting Image Models to Generate Videos



*"A robot dj is playing the turntables, in heavy raining futuristic tokyo, rooftop, sci-fi, fantasy"*

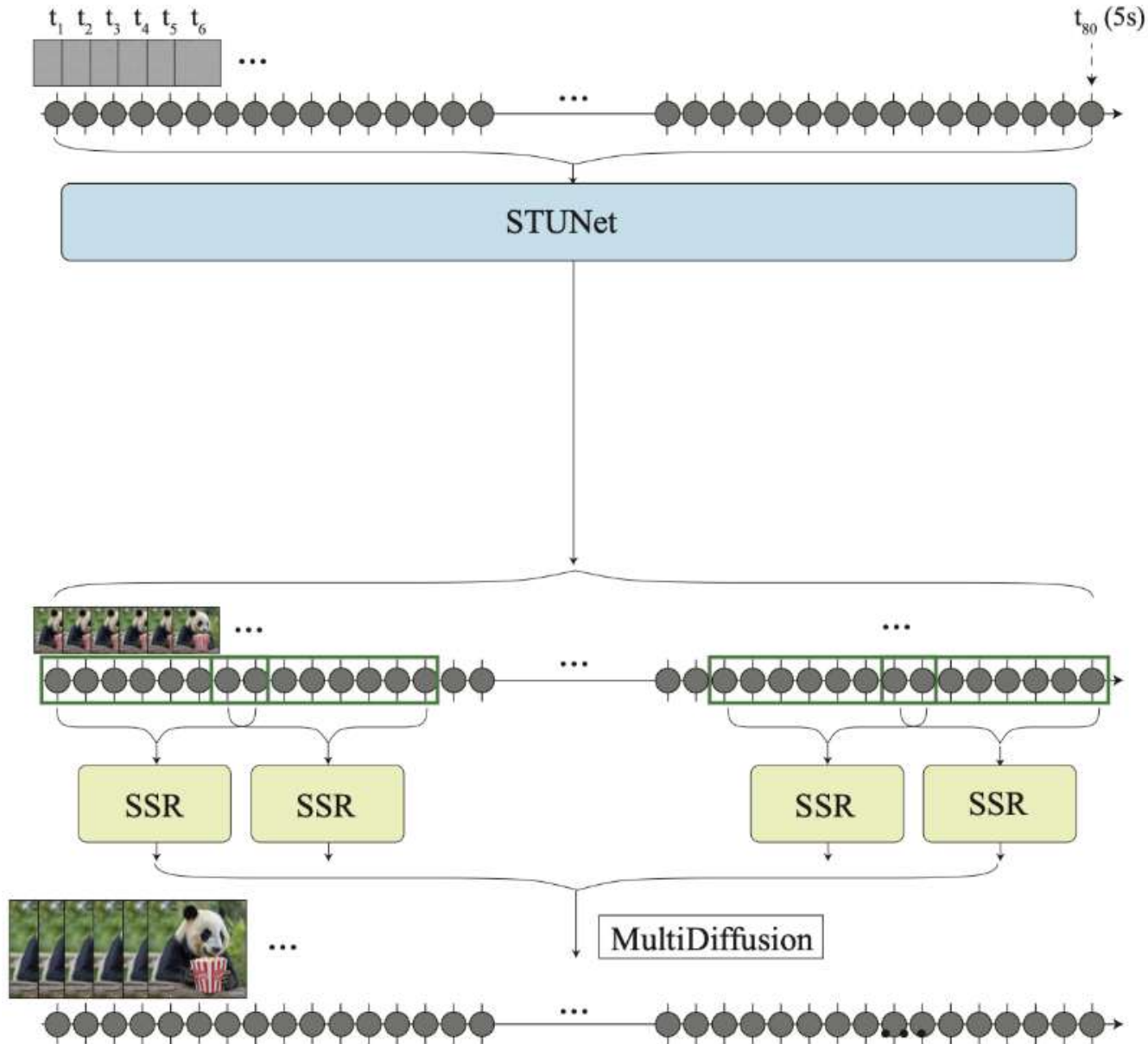


*"An exploding cheese house"*



*"A fat rabbit wearing a purple robe walking through a fantasy landscape"*

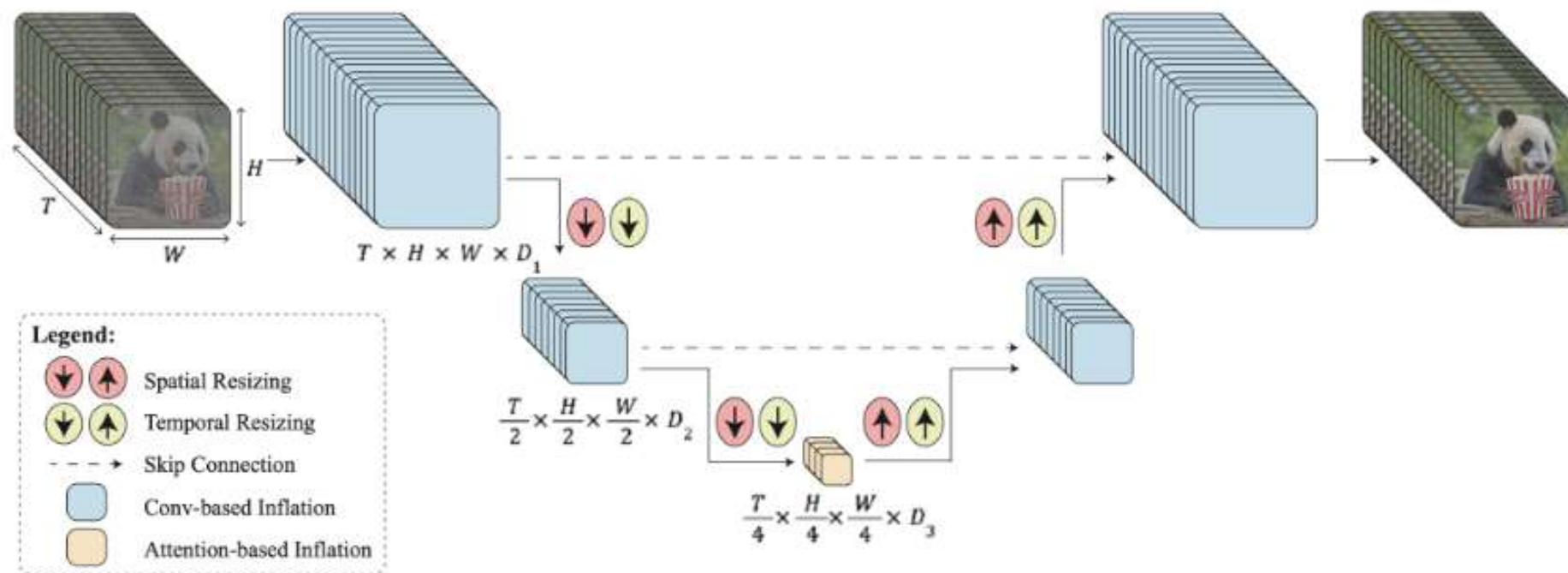
# Adapting Image Models to Generate Videos



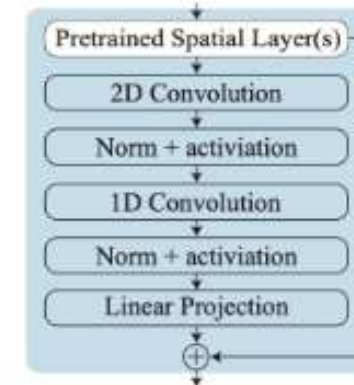
Lumiere removes TSR (temporal super-resolution) models. The inflated SSR network can operate only on short segments of the video due to memory constraints and thus SSR models operate on a set of shorter but overlapped video snippets.

# Adapting Image Models to Generate Videos

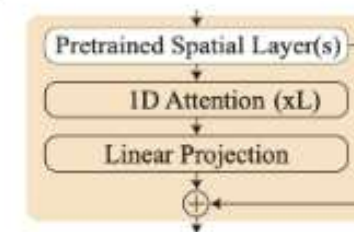
(a) Space-Time UNet (STUNet)



(b) Convolution-based Inflation Block



(c) Attention-based Inflation Block



# Adapting Image Models to Generate Videos



*wearing a gold strapless gown*

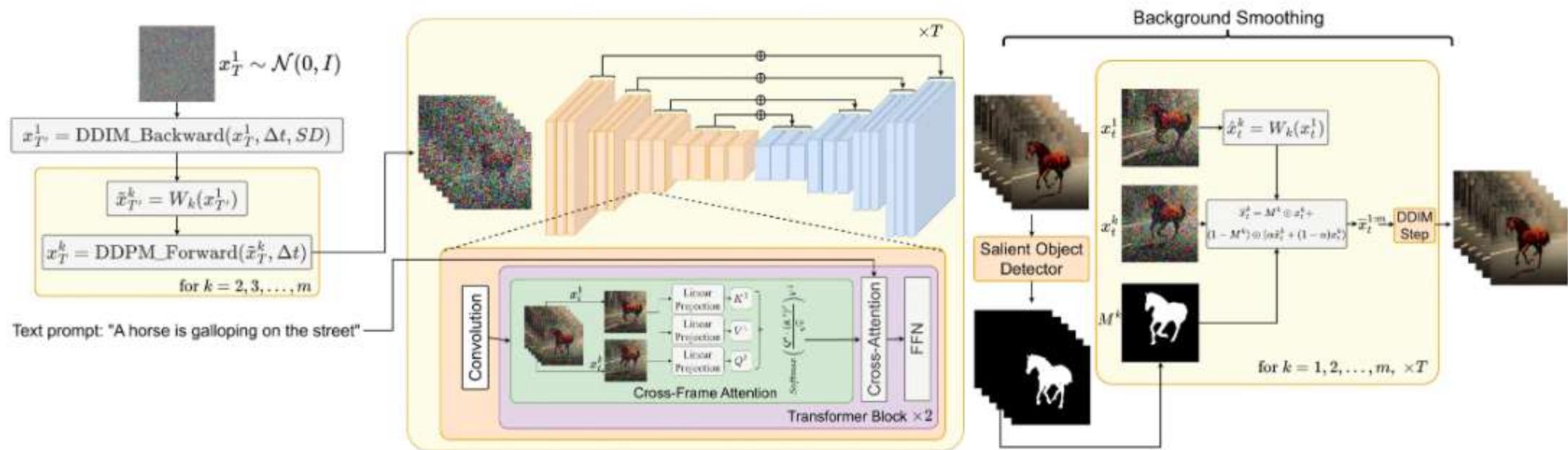


*wearing a bath robe*

# Training-Free Adaptation

**Text2Video-Zero:  
Text-to-Image Diffusion Models are  
Zero-Shot Video Generators**

# Training-Free Adaptation



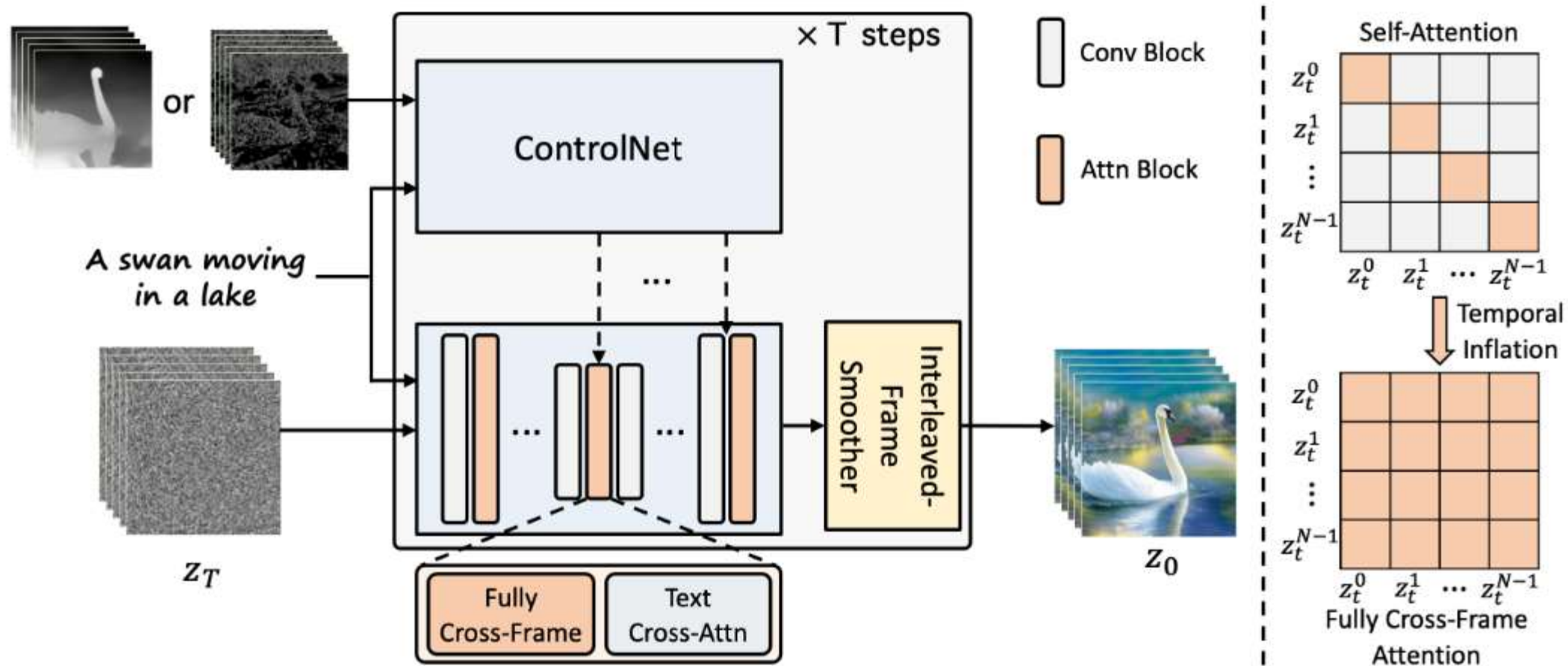
$$\mathbf{x}_{T'}^1 = \text{DDIM-backward}(\mathbf{x}_T^1, \Delta t) \text{ where } T' = T - \Delta t$$

$$W_k \leftarrow \text{a warping operation of } \delta^k = \lambda(k-1)\delta$$

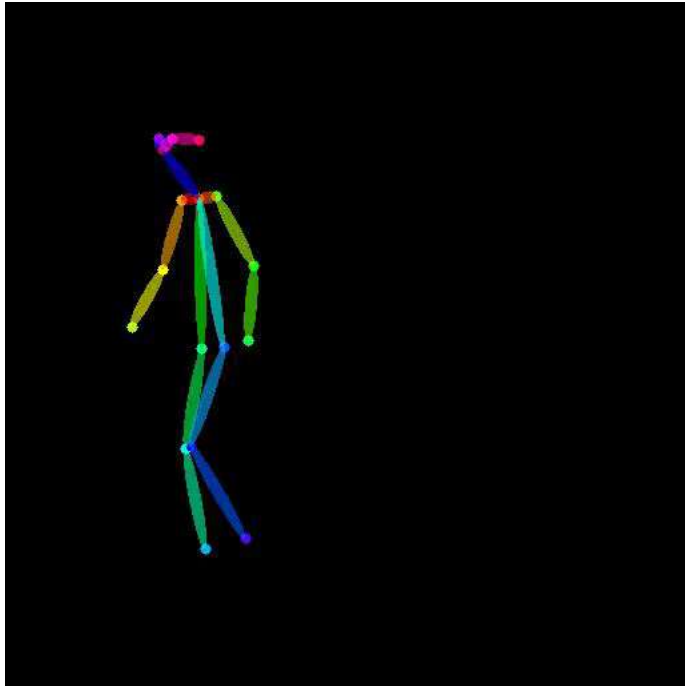
$$\tilde{\mathbf{x}}_{T'}^k = W_k(\mathbf{x}_{T'}^1)$$

$$\mathbf{x}_T^k = \text{DDIM-forward}(\tilde{\mathbf{x}}_{T'}^k, \Delta t) \text{ for } k = 2, \dots, m$$

# Training-Free Adaptation



# Training-Free Adaptation



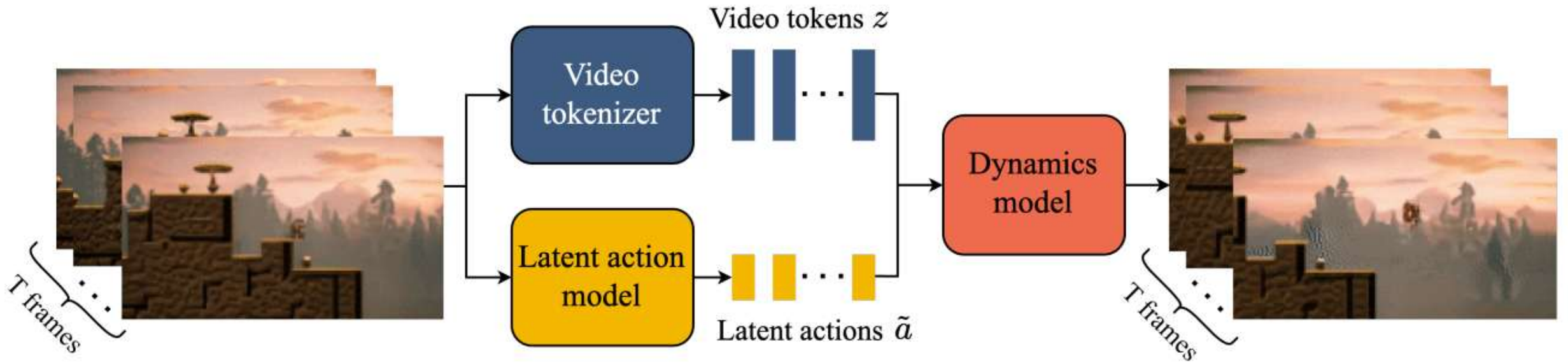
"A man, wearing pink clothes, moonwalk at sunset."

# Genie: Generative Interactive Environments



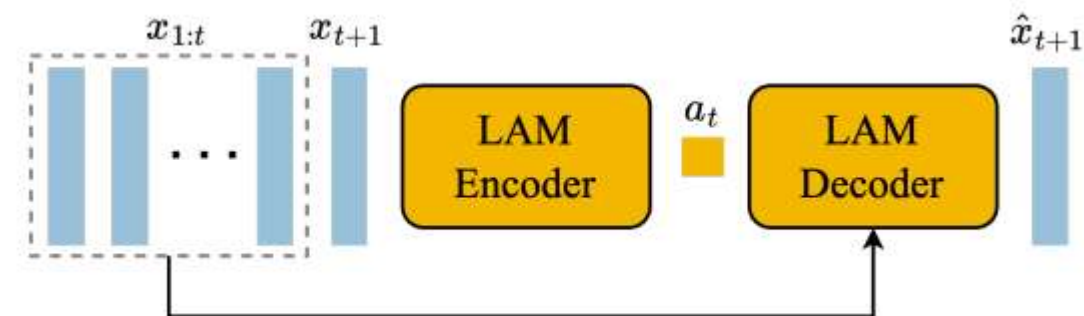
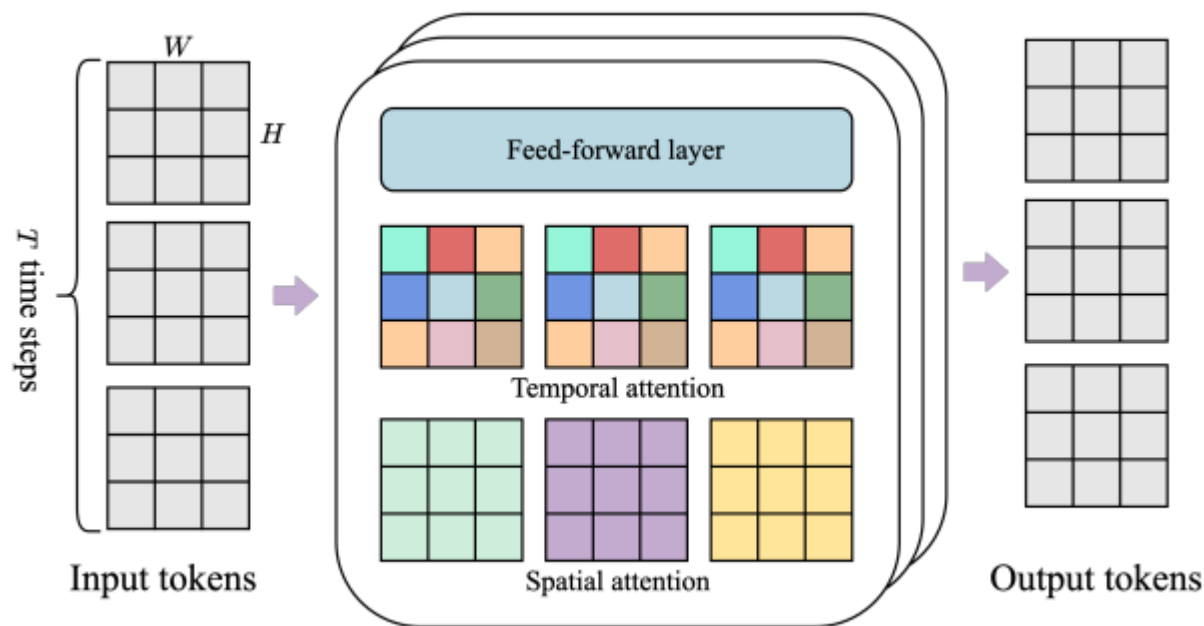
Generate a playable world  
set in a futuristic city

# Genie: Generative Interactive Environments



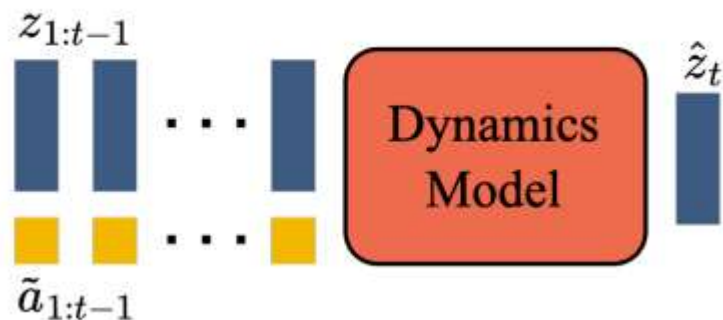
Genie takes in  $T$  frames of video as input, tokenizes them into discrete tokens  $\mathbf{z}$  via the video tokenizer, and infers the latent actions  $\tilde{\mathbf{a}}$  between each frame with the latent action model. Both are then passed to the dynamics model to generate predictions for the next frames in an iterative manner.

# Genie: Generative Interactive Environments

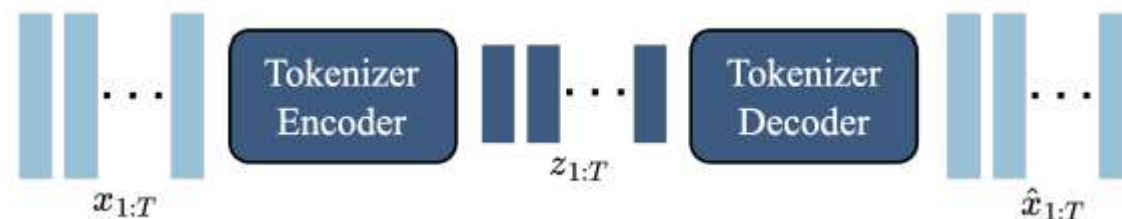


Latent action model

ST-transformer architecture



Dynamics model



Video tokenizer

# Genie: Generative Interactive Environments



“Remarkably, Genie learns not only which parts of an observation are generally controllable, but also infers diverse latent actions that are consistent across the generated environments. Note here how the same latent actions yield similar behaviors across different prompt images. ”

# Genie: Generative Interactive Environments



The future of generative virtual worlds

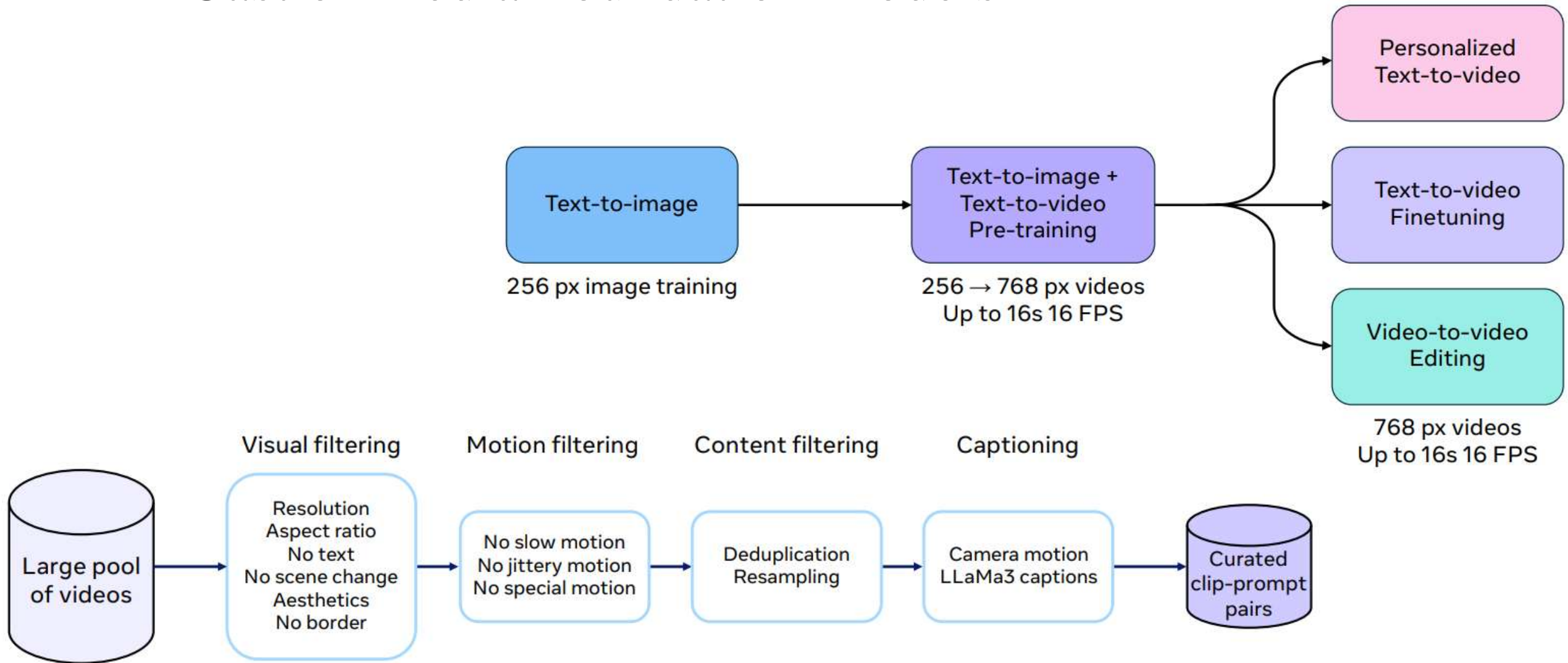
“Trajectories with the same latent action sequence typically display similar behaviors.”

# Movie Gen: A Cast of Media Foundation Models

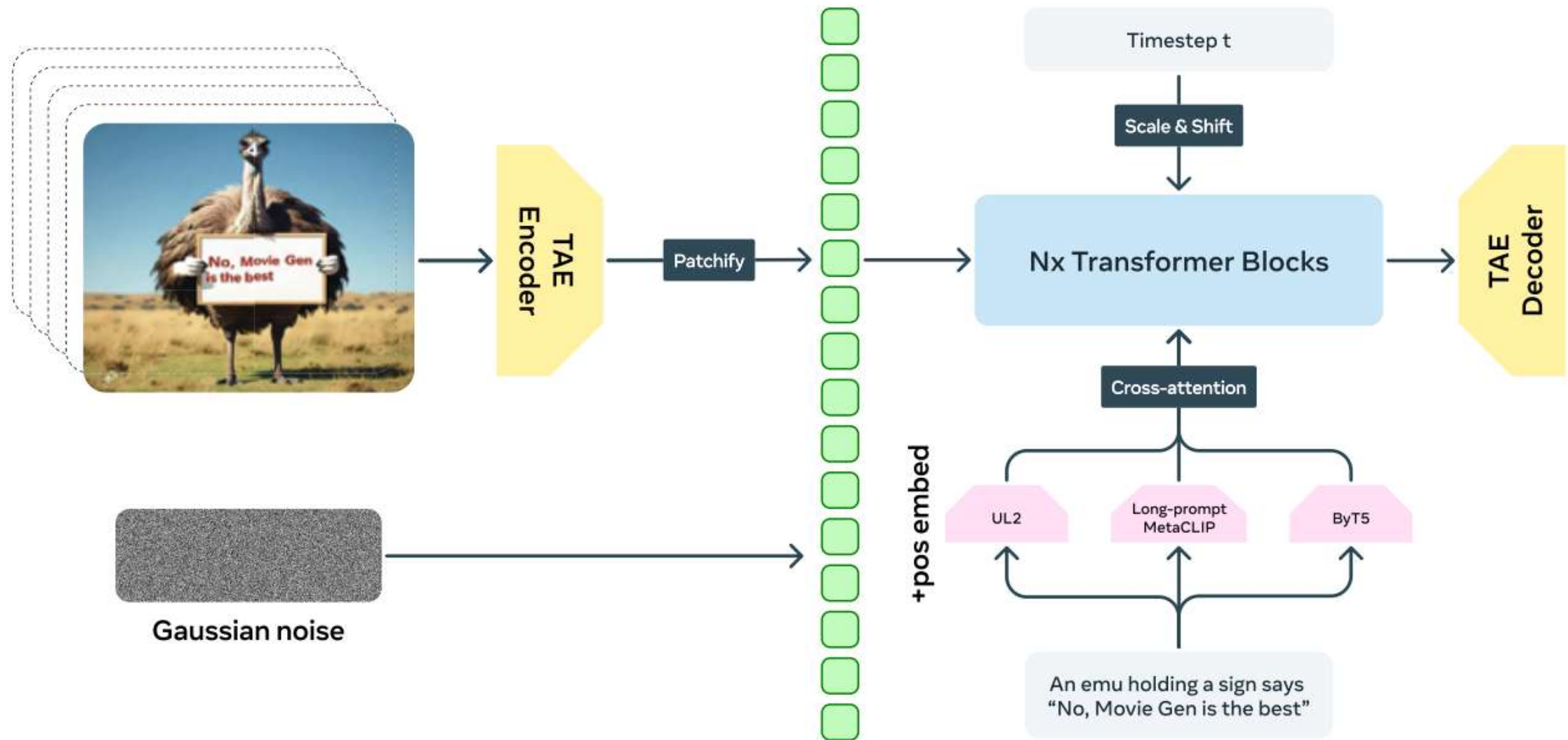


<https://ai.meta.com/blog/movie-gen-media-foundation-models-generative-ai-video/>

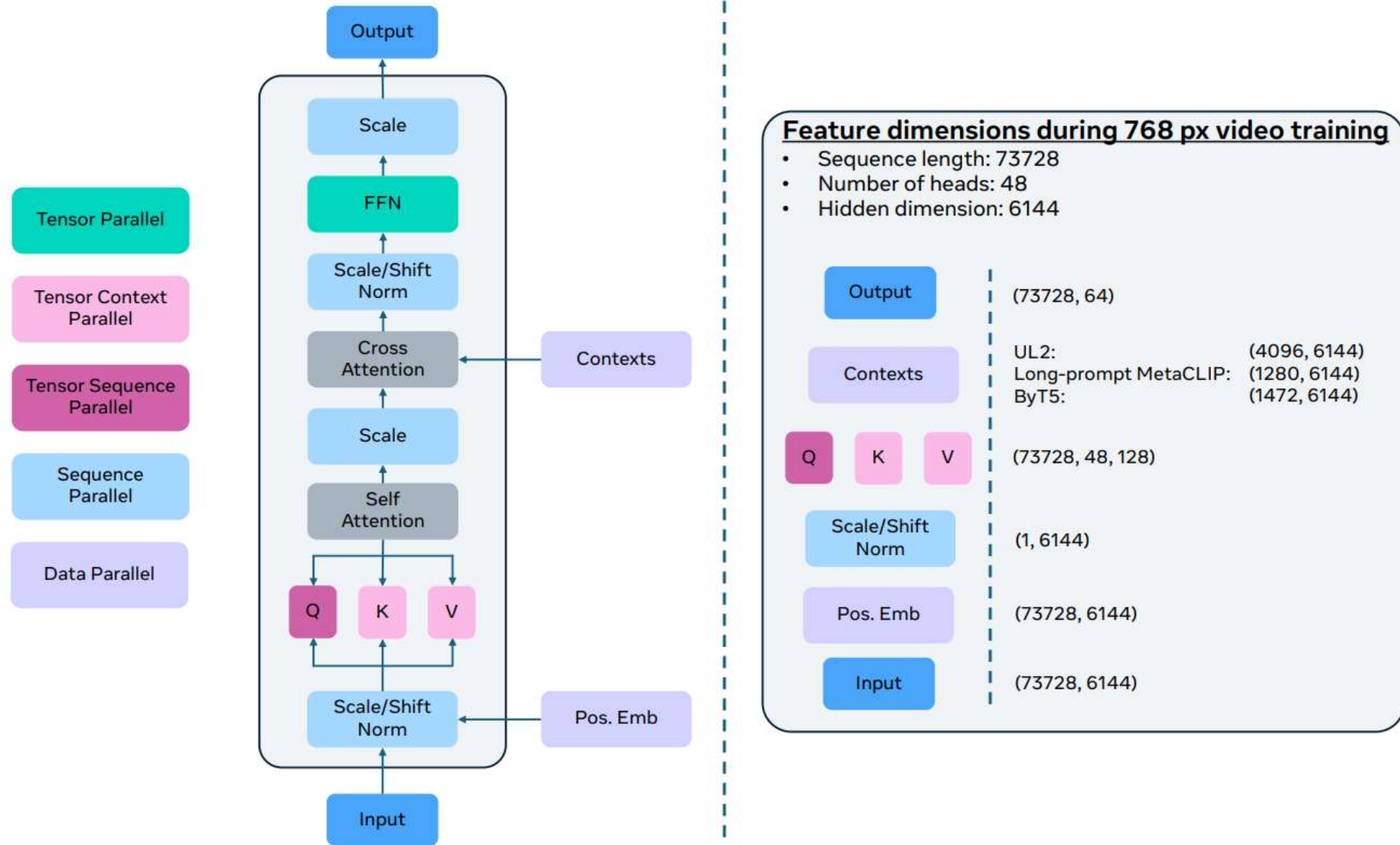
# Movie Gen: A Cast of Media Foundation Models



# Movie Gen: A Cast of Media Foundation Models



# Movie Gen: A Cast of Media Foundation Models



# One Model to Workflow

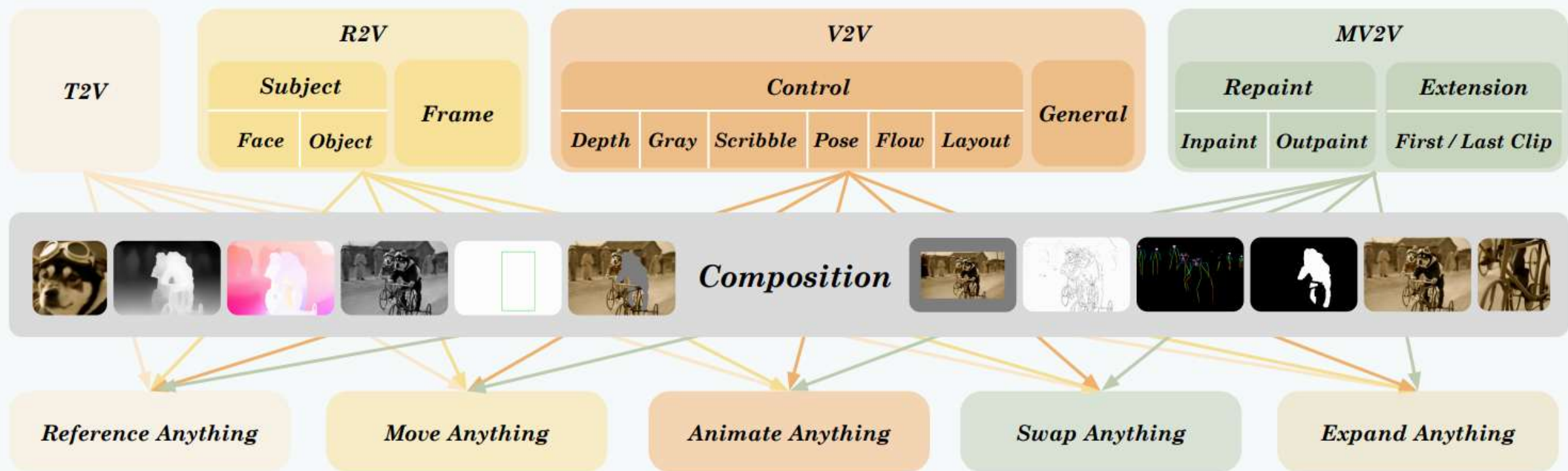


\*image taken from ChatGPT 4o

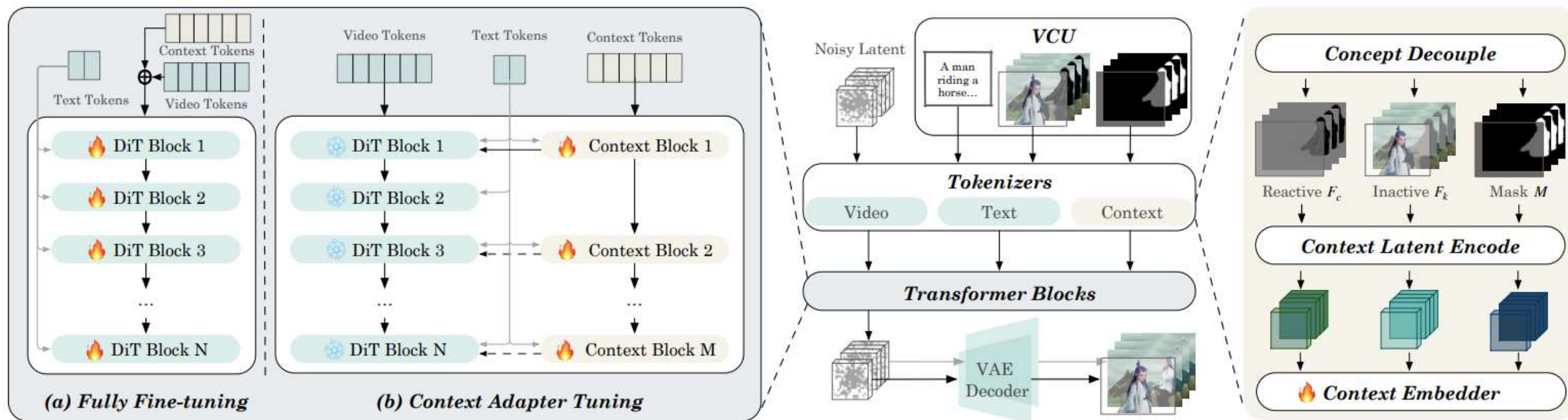
# VACE/Wan



# VACE



# VACE



# VACE



# Sora2 & Veo3.1: Omni Model / “WORLD MODEL”

## From Silent Video to Native Audiovisual Generation

- The target is synchronized audiovisual events.
- Dialogue, ambience, and sound effects must cohere with motion.
- New bottlenecks: lip-sync, event-sound causality, and cross-shot continuity.

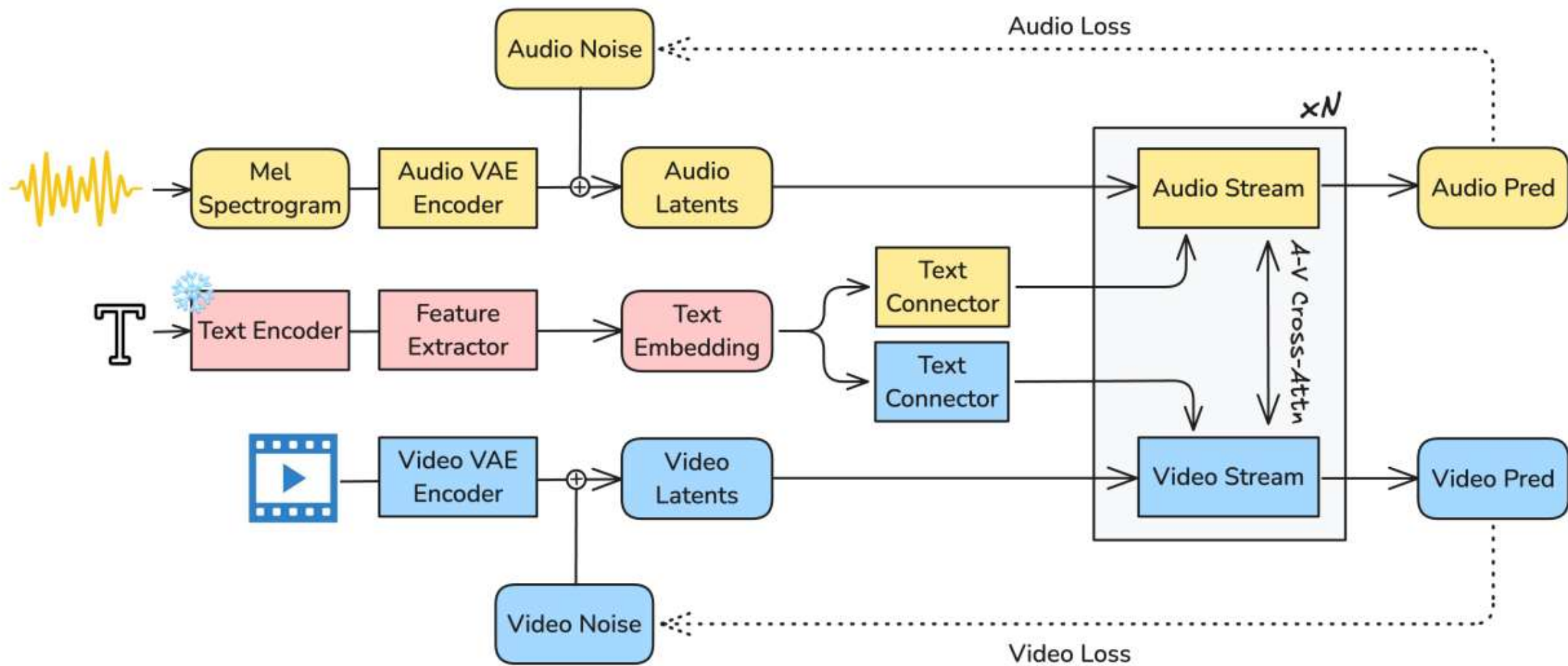
- Sora 2 improves physics, realism, synchronized audio, and steerability.
- Prompts now combine subject, camera, pacing, and audio intent.
- Storyboard, remix, stitch, extend, and characters define the workflow.

<https://sora.chatgpt.com/explore>

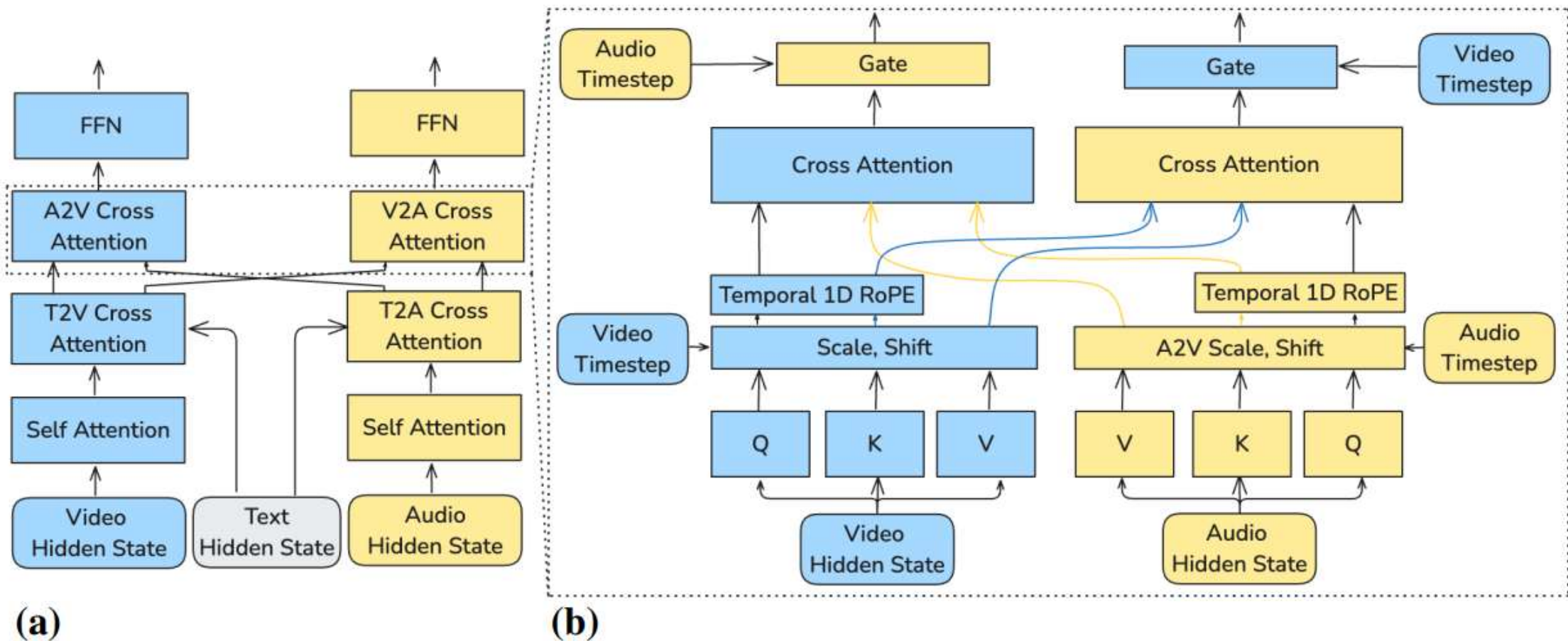
- Veo 3.1 adds richer native audio and stronger prompt adherence.
- Reference images, first/last frames, and extension expose creative control.
- A strong Veo prompt reads like a director's brief.

[https://blog.google/innovation-and-ai/products/veo-updates-flow/?utm\\_source=chatgpt.com](https://blog.google/innovation-and-ai/products/veo-updates-flow/?utm_source=chatgpt.com)

# LTX



# LTX



# LTX

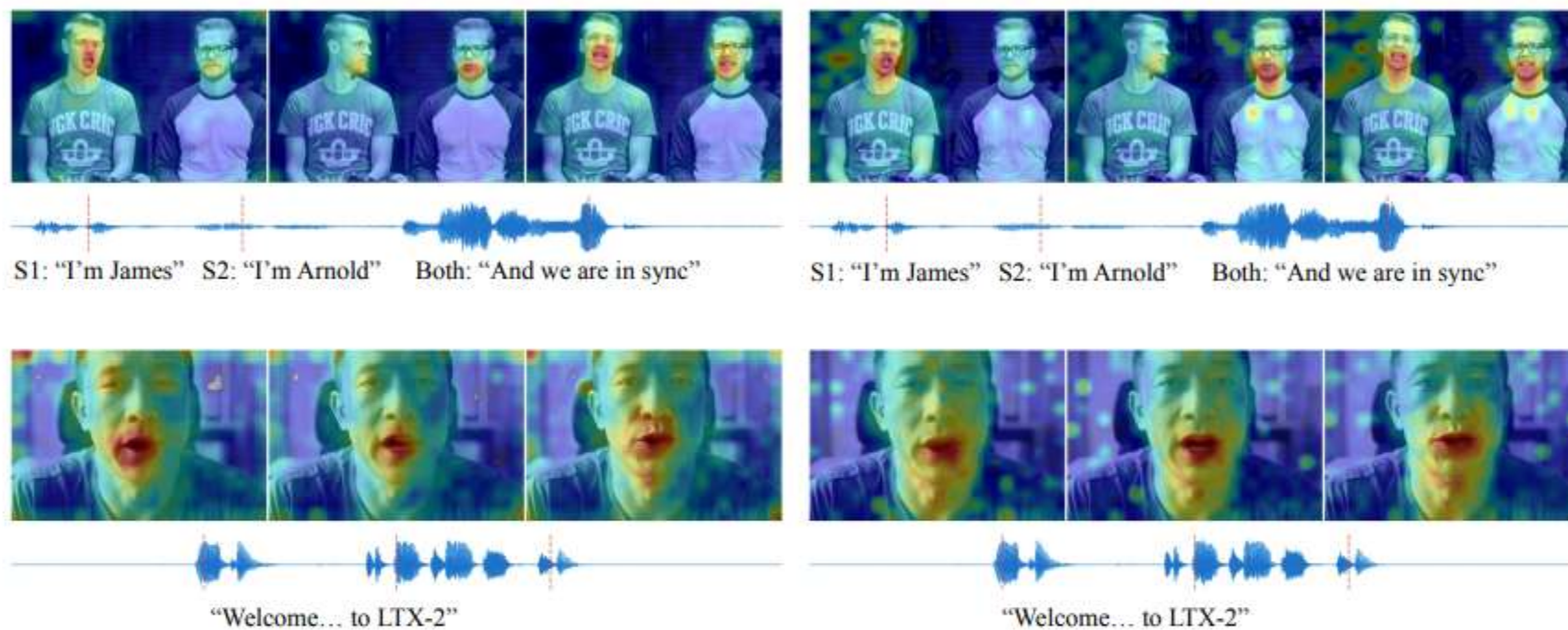


Figure 3: **Visualization of AV cross-attention maps.** The maps are averaged across attention heads and the model layers; V2A and A2V maps correspond to the first and last 1/3 of inference steps, respectively. Red vertical lines on the audio waveform mark the timestamps of the displayed frames. The visualization demonstrates the model’s ability to spatially track a moving vehicle, dynamically shift attention from one speaker to another and then to both simultaneously, and focus on the lip region during close-up speech.

# Seedance