

# **Conditional Generation of Diffusion Models**

Conditional generation

**User input:**

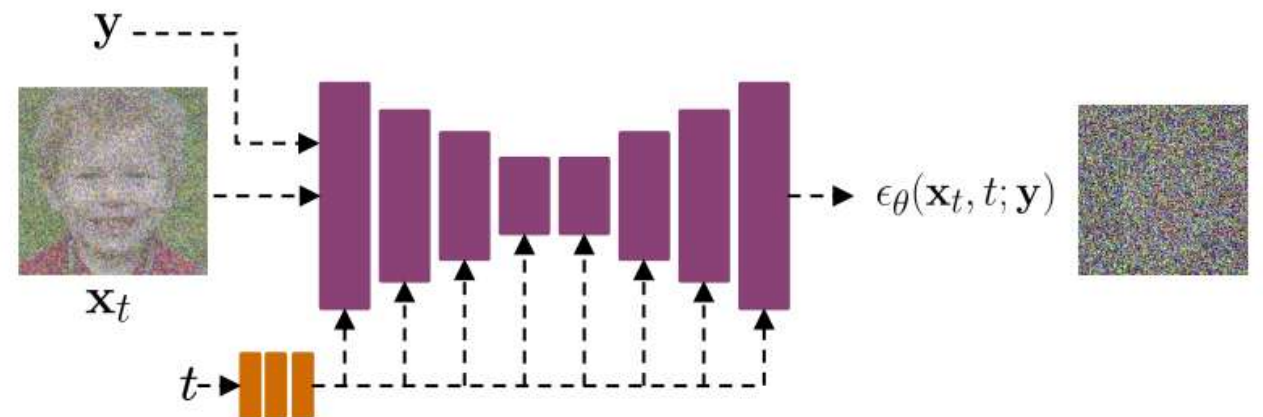
An astronaut riding a horse



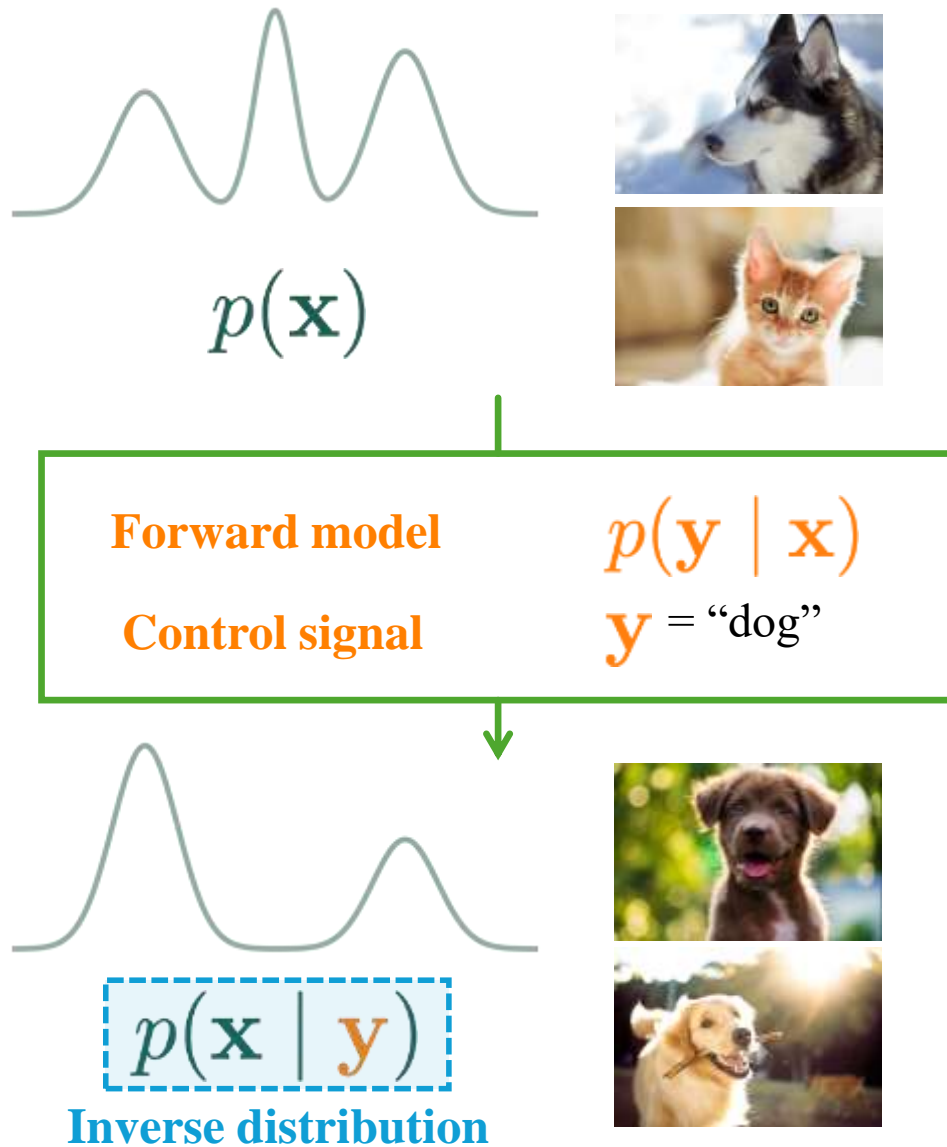
# Conditional generation

- Let  $(x,y)$  denote (image,caption) pairs
- Training a conditional generative model involves learning  $p(x | y)$
- Train score model for the image  $x$  conditional on caption  $y$

$$\mathbb{E}_{(x,y) \sim p_{\text{data}}(x,y)} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{E}_{t \sim \mathcal{U}[0,T]} \|\epsilon_{\theta}(\mathbf{x}_t, t; \mathbf{y}) - \epsilon\|_2^2$$



# Control the generation process



Bayes' rule:

$$p(\mathbf{x} | \mathbf{y}) = \frac{p(\mathbf{x}) p(\mathbf{y} | \mathbf{x})}{p(\mathbf{y})}$$

Bayes' rule for score functions:

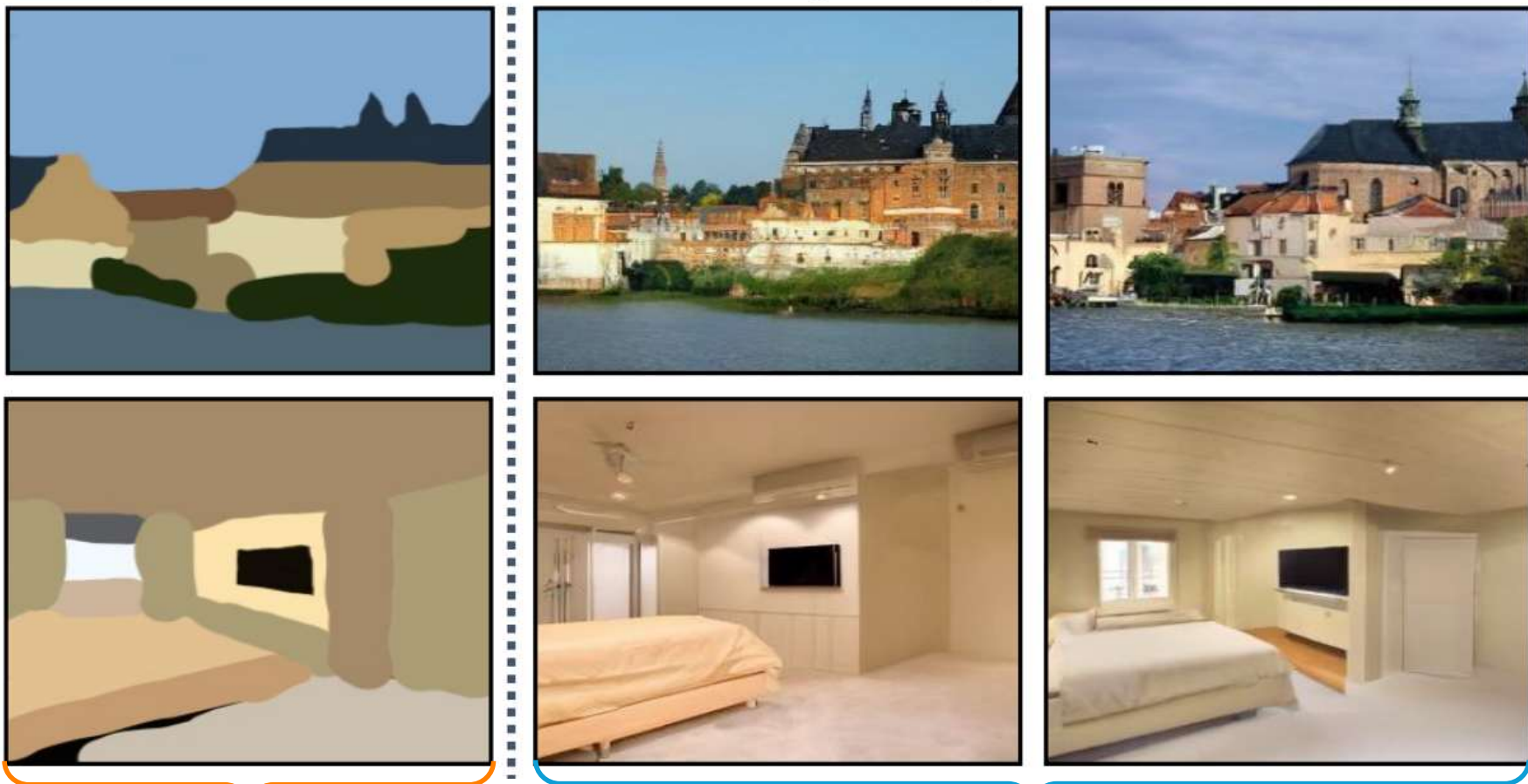
$$\begin{aligned} \nabla_{\mathbf{x}} \log p(\mathbf{x} | \mathbf{y}) &= \nabla_{\mathbf{x}} \log p(\mathbf{x}) \\ &\quad + \nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}) \\ &\quad - \nabla_{\mathbf{x}} \log p(\mathbf{y}) \quad \mathbf{0} \\ &= \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}) \end{aligned}$$



Plug in different forward models for the same score model

# Stroke to image synthesis

## Stroke Painting to Image



Forward model  
 $p(\mathbf{y} \mid \mathbf{x})$   
can be specified.

Stroke paintings

$\mathbf{y}$

Sampled images

$\mathbf{x} \mid \mathbf{y}$

# Language-guided image generation

$y$        $x \mid y$

**(Prompt)**

Treehouse in the style of  
Studio Ghibli animation



Forward model

$$p(y \mid x)$$

is an image captioning  
neural network.



[ Work by @danielrussruss ]

# Classifier Guidance

Bayes' rule:

$$p(\mathbf{x} | \mathbf{y}) = \frac{p(\mathbf{x}) p(\mathbf{y} | \mathbf{x})}{p(\mathbf{y})}$$

The equation is annotated with a blue checkmark above  $p(\mathbf{x})$ , a blue checkmark above  $p(\mathbf{y} | \mathbf{x})$ , and a grey 'X' below  $p(\mathbf{y})$ . Dashed blue boxes enclose the numerator terms.

Bayes' rule for score functions:

$$\begin{aligned} \nabla_{\mathbf{x}} \log p(\mathbf{x} | \mathbf{y}) &= \nabla_{\mathbf{x}} \log p(\mathbf{x}) \\ &\quad + \nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}) \\ &\quad - \nabla_{\mathbf{x}} \log p(\mathbf{y}) \quad \mathbf{0} \\ &= \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}) \end{aligned}$$

Annotations: A blue arrow points from the first term to "Conditional score". A blue arrow points from the second term to "Classifier obtained as the difference". A blue arrow points from the third term to "Unconditional score".

# Classifier-Guided Diffusion

To explicitly incorporate class information into the diffusion process, trained a classifier  $f_\phi(y|\mathbf{x}_t, t)$  on noisy image the use gradients to guide diffusion sampling process toward the condition

$$\begin{aligned}\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t, y) &= \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log q(y|\mathbf{x}_t) \\ &\approx -\frac{1}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) + \nabla_{\mathbf{x}_t} \log f_\phi(y|\mathbf{x}_t) \\ &= -\frac{1}{\sqrt{1-\bar{\alpha}_t}} (\epsilon_\theta(\mathbf{x}_t, t) - \sqrt{1-\bar{\alpha}_t} \nabla_{\mathbf{x}_t} \log f_\phi(y|\mathbf{x}_t))\end{aligned}$$

A new classifier-guided predictor  $\bar{\epsilon}_\theta$  would take the form as following,

$$\bar{\epsilon}_\theta(\mathbf{x}_t, t) = \epsilon_\theta(\mathbf{x}_t, t) - \sqrt{1-\bar{\alpha}_t} \nabla_{\mathbf{x}_t} \log f_\phi(y|\mathbf{x}_t)$$

To control the strength of the classifier guidance, we can add a weight  $w$  to the delta part,

$$\bar{\epsilon}_\theta(\mathbf{x}_t, t) = \epsilon_\theta(\mathbf{x}_t, t) - \sqrt{1-\bar{\alpha}_t} w \nabla_{\mathbf{x}_t} \log f_\phi(y|\mathbf{x}_t)$$

# Classifier-Guided Diffusion

---

**Algorithm 1** Classifier guided diffusion sampling, given a diffusion model  $(\mu_\theta(x_t), \Sigma_\theta(x_t))$ , classifier  $f_\phi(y|x_t)$ , and gradient scale  $s$ .

---

Input: class label  $y$ , gradient scale  $s$

$x_T \leftarrow$  sample from  $\mathcal{N}(0, \mathbf{I})$

**for all**  $t$  from  $T$  to 1 **do**

$\mu, \Sigma \leftarrow \mu_\theta(x_t), \Sigma_\theta(x_t)$

$x_{t-1} \leftarrow$  sample from  $\mathcal{N}(\mu + s\Sigma \nabla_{x_t} \log f_\phi(y|x_t), \Sigma)$

**end for**

**return**  $x_0$

---

---

**Algorithm 2** Classifier guided DDIM sampling, given a diffusion model  $\epsilon_\theta(x_t)$ , classifier  $f_\phi(y|x_t)$ , and gradient scale  $s$ .

---

Input: class label  $y$ , gradient scale  $s$

$x_T \leftarrow$  sample from  $\mathcal{N}(0, \mathbf{I})$

**for all**  $t$  from  $T$  to 1 **do**

$\hat{\epsilon} \leftarrow \epsilon_\theta(x_t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{x_t} \log f_\phi(y|x_t)$

$x_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \left( \frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}$

**end for**

**return**  $x_0$

---

# Classifier-Guided Diffusion



# Classifier-Free Guidance

- Train both a conditional and an unconditional score model (by randomly dropping the caption during training)
- Combine the two models as follows

$$\begin{aligned}(1 + w)\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) &= (1 + w)(\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \\ &= (1 + w)\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) - w\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)\end{aligned}$$

- $w$  is the classifier-guidance strength

# Effect of classifier guidance



Increased classifier guidance strength ( $w$ )

# Text-to-Image

<https://laion.ai/blog/laion-5b/>

A cat in  
the snow



Text-to-image  
Generator



ImageNet

1M



LAION

5.85B

Backend url:

<https://knn5.laion>

Index:

laion\_5B

french cat

[Clip retrieval](#) works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

Display captions

Display full captions

Display similarities

Safe mode

Hide duplicate urls

Hide (near) duplicate images

Search over

image

Search with multilingual clip



french cat



french cat



How to tell if your feline is french. He wears a b...



Winter cat



網友挑戰「加幾筆畫出最創意貓咪圖片」，第10位名之後我出



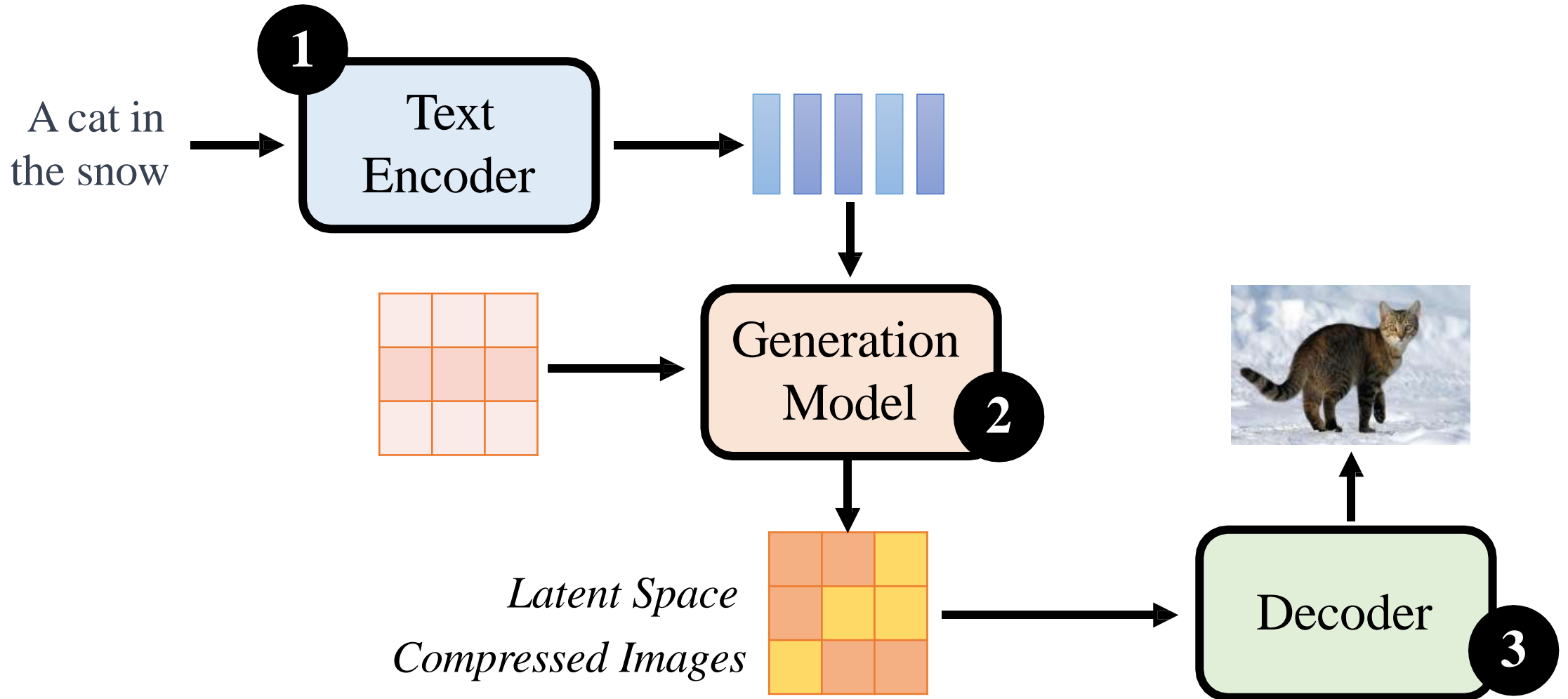
cat in a suit Georgian sells tomatoes

# Framework

A cat in  
the snow

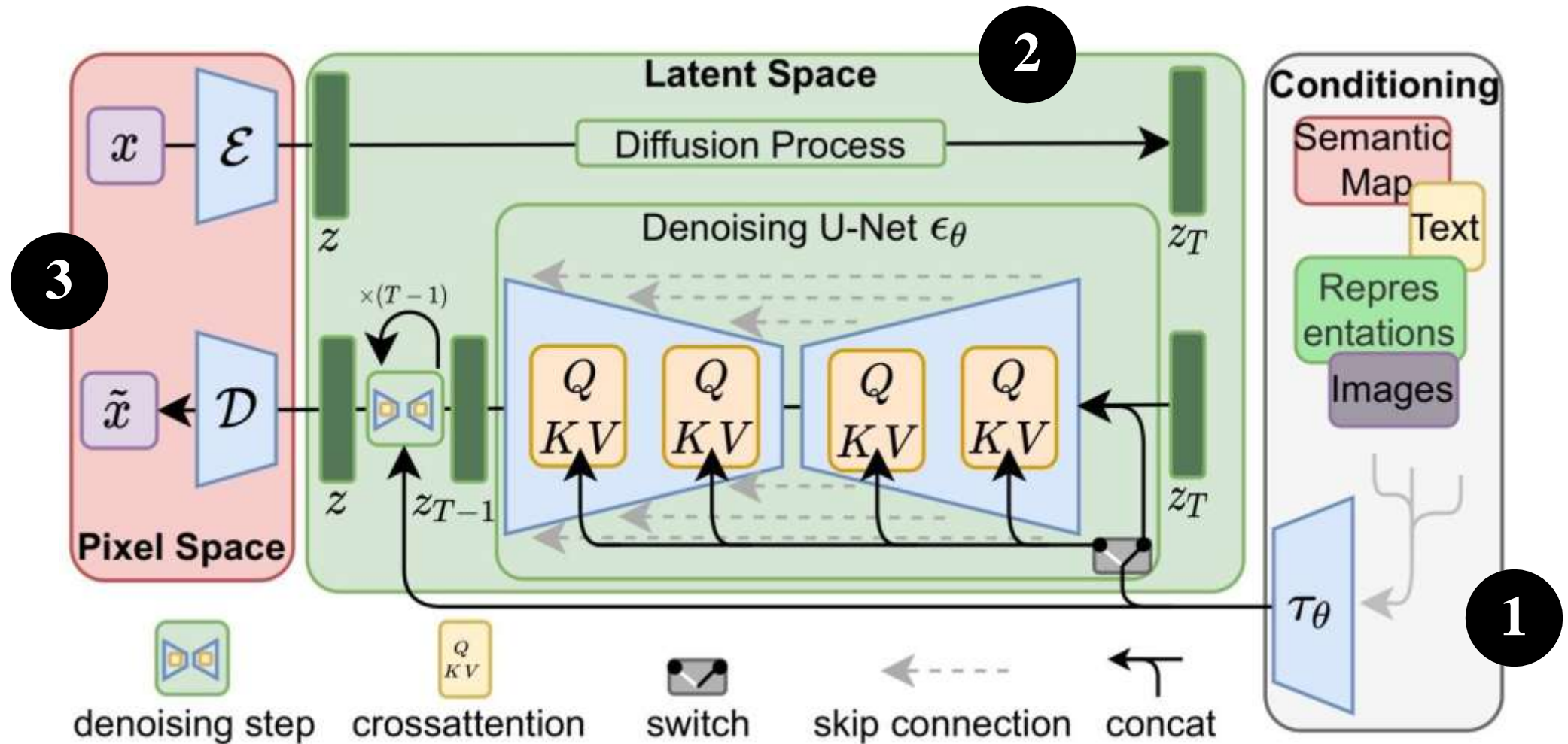


Text-to-image  
Generator



# Stable Diffusion

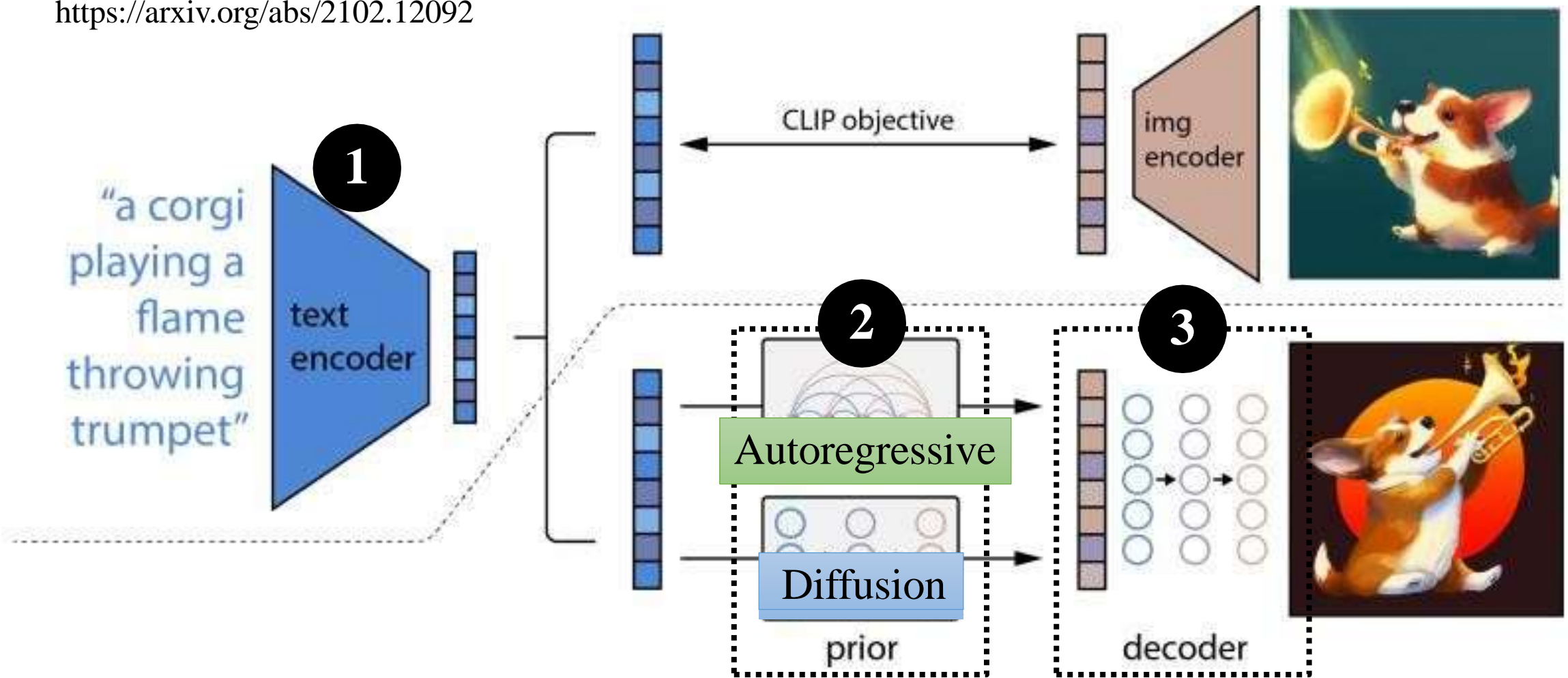
<https://arxiv.org/abs/2112.10752>



# DALL-E series

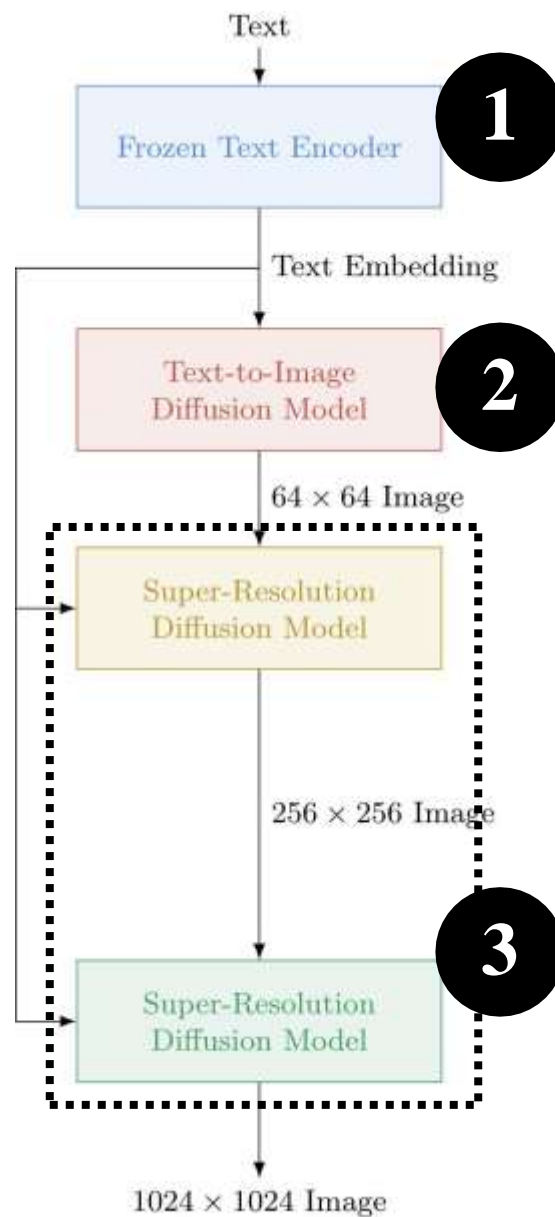
<https://arxiv.org/abs/2204.06125>

<https://arxiv.org/abs/2102.12092>

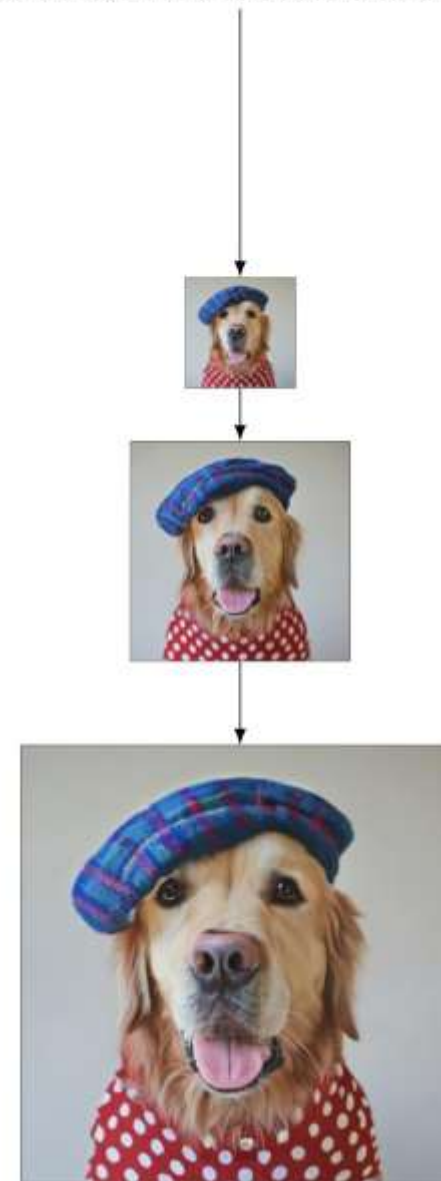


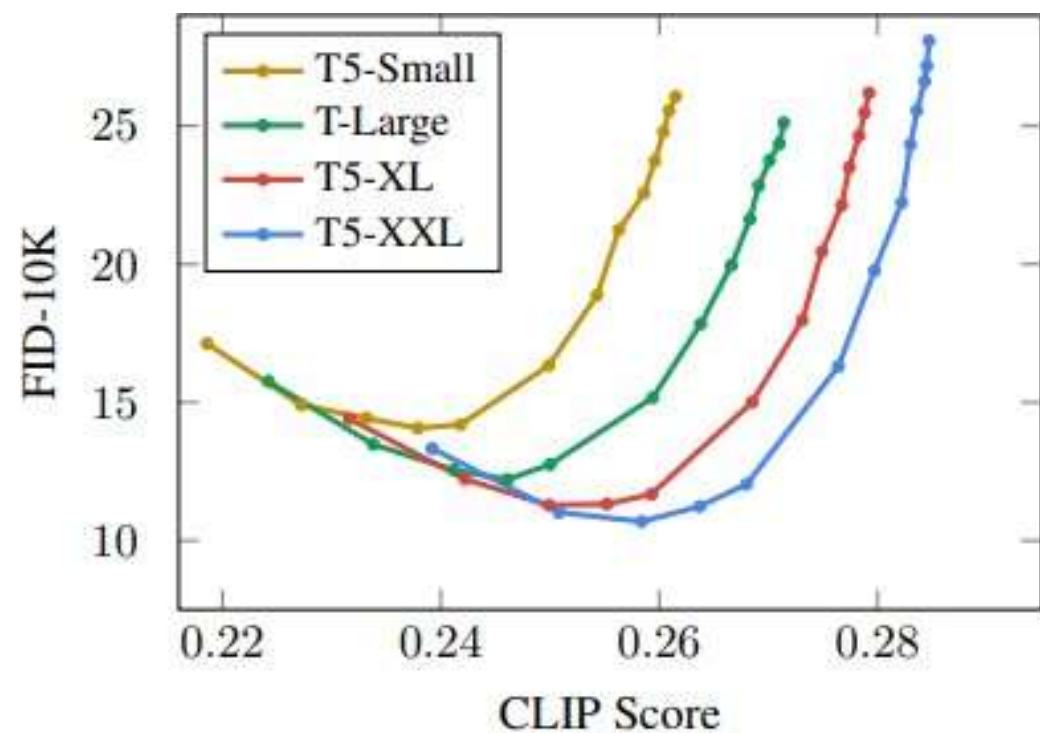
# Imagen

<https://imagen.research.google/>  
<https://arxiv.org/abs/2205.11487>

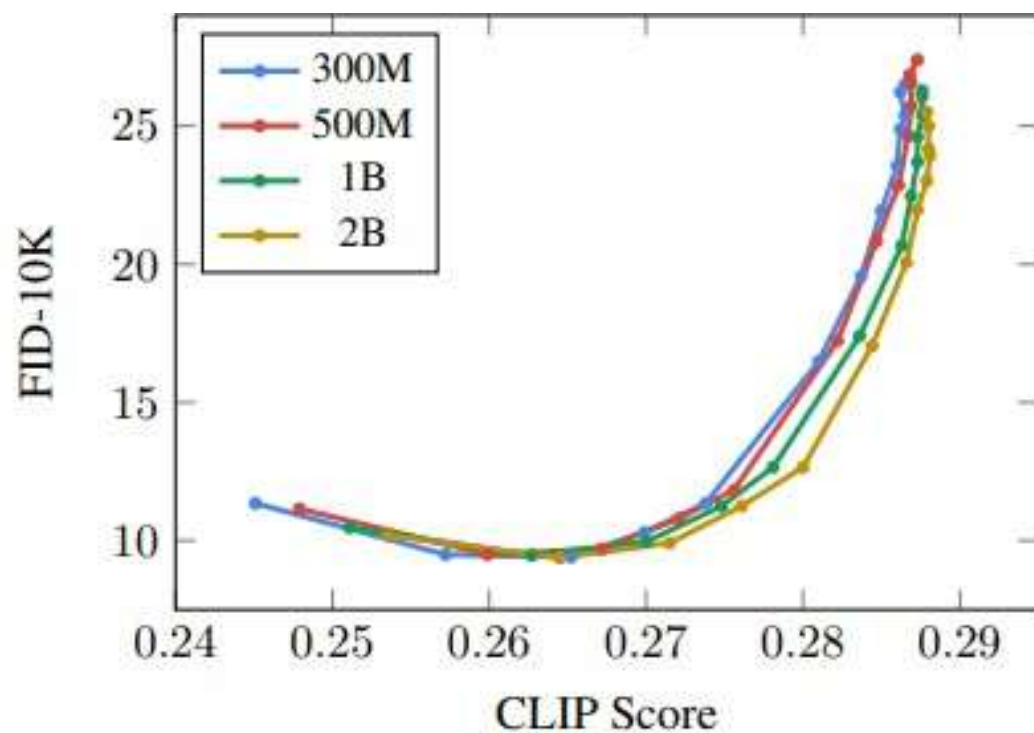


“A Golden Retriever dog wearing a blue checkered beret and red dotted turtleneck.”





(a) Impact of encoder size.



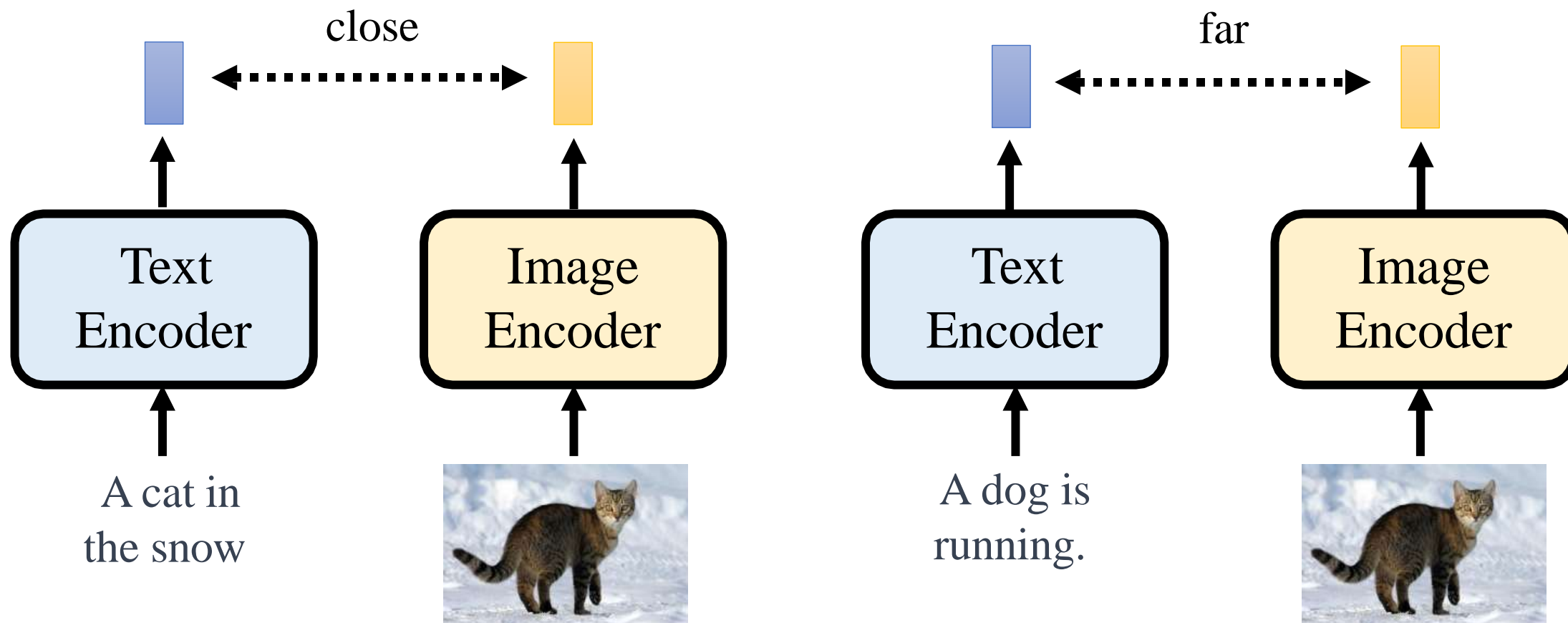
(b) Impact of U-Net size.

<https://arxiv.org/abs/2205.11487>

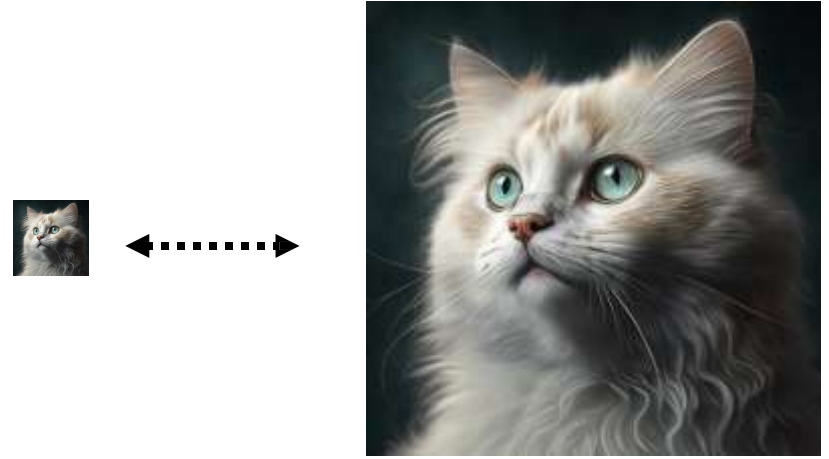
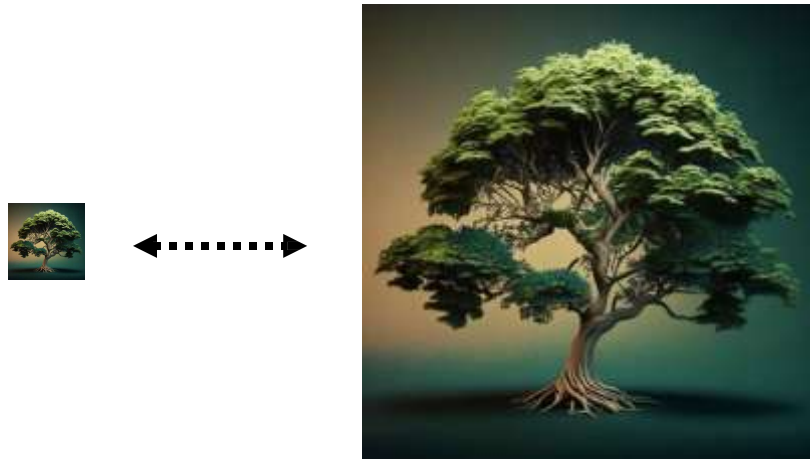
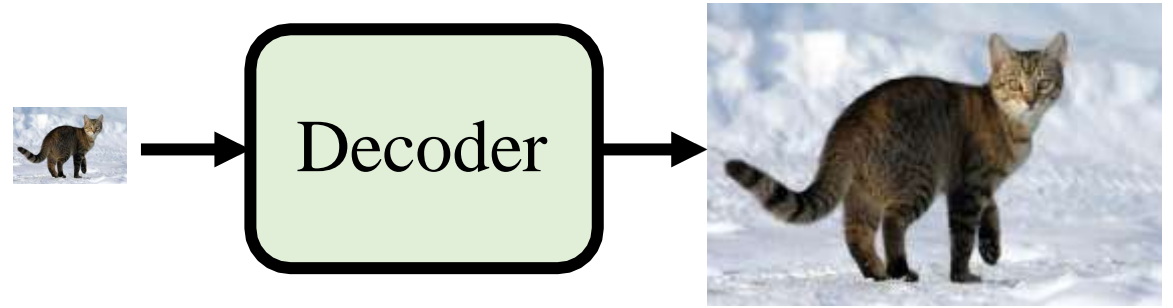
# Contrastive Language-Image Pre-Training (CLIP)

<https://arxiv.org/abs/2103.00020>

**400 million image-text  
pairs**



# *Progressive/Multi-Scale*

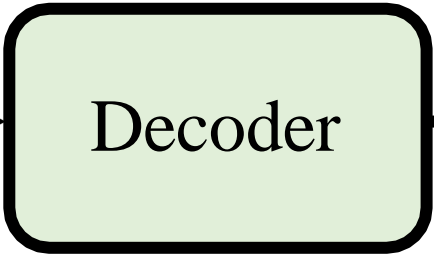
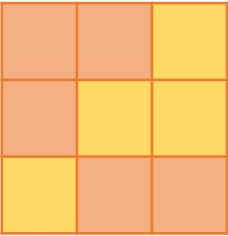


(Images are generated by Midjourney)

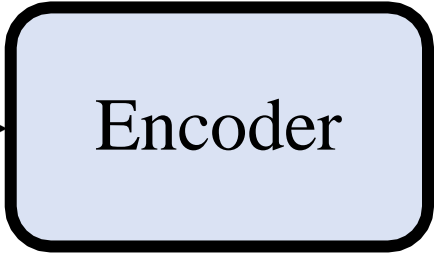
# Latent Diffusion

Auto-encoder

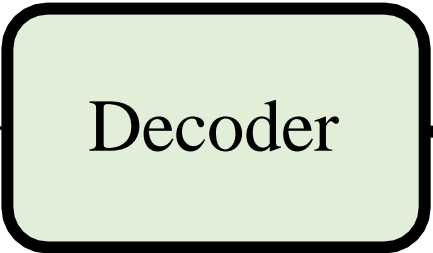
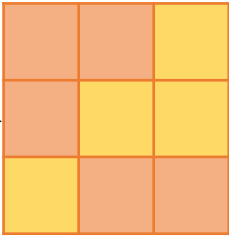
Latent Representation



$H \times W \times 3$

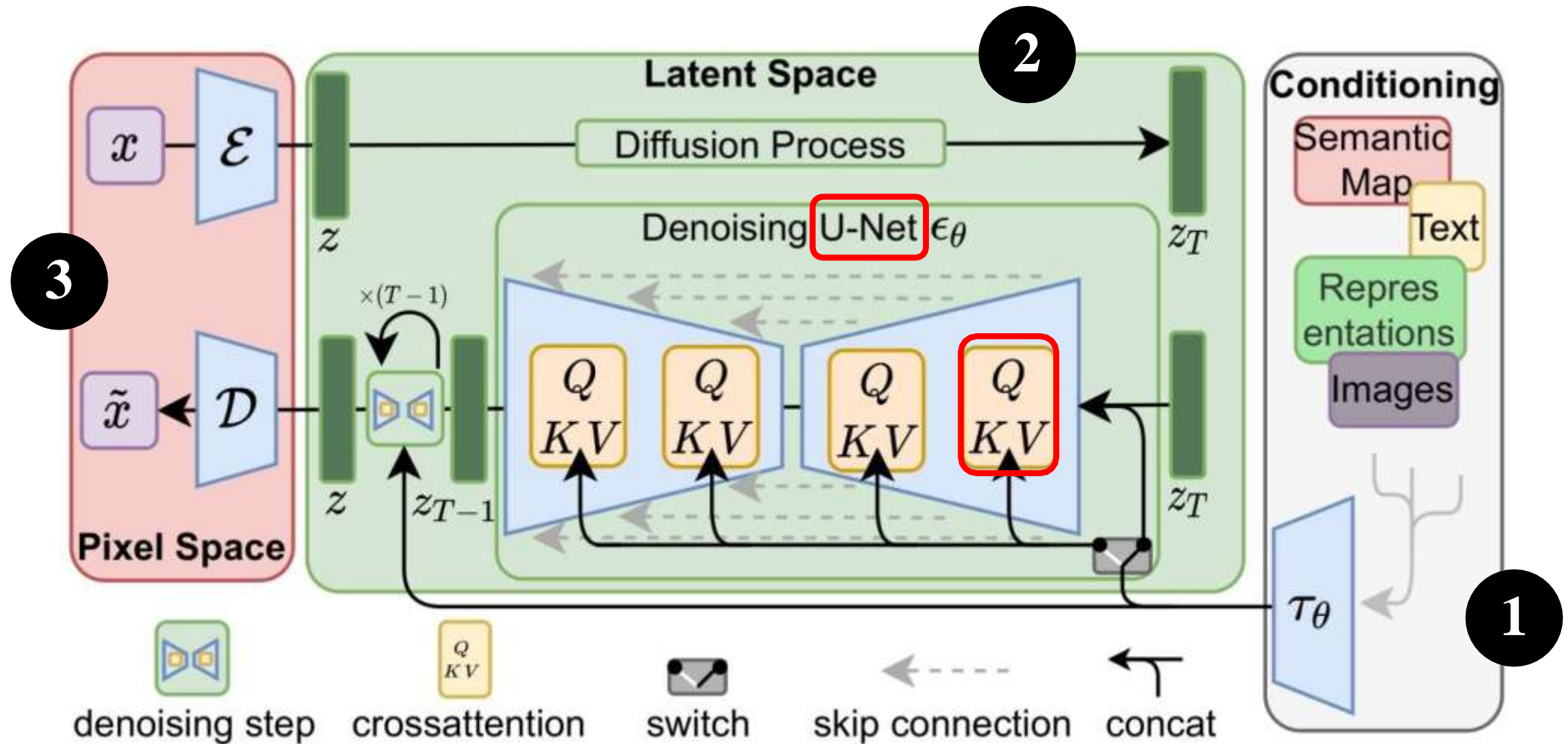


$h \times w \times c$

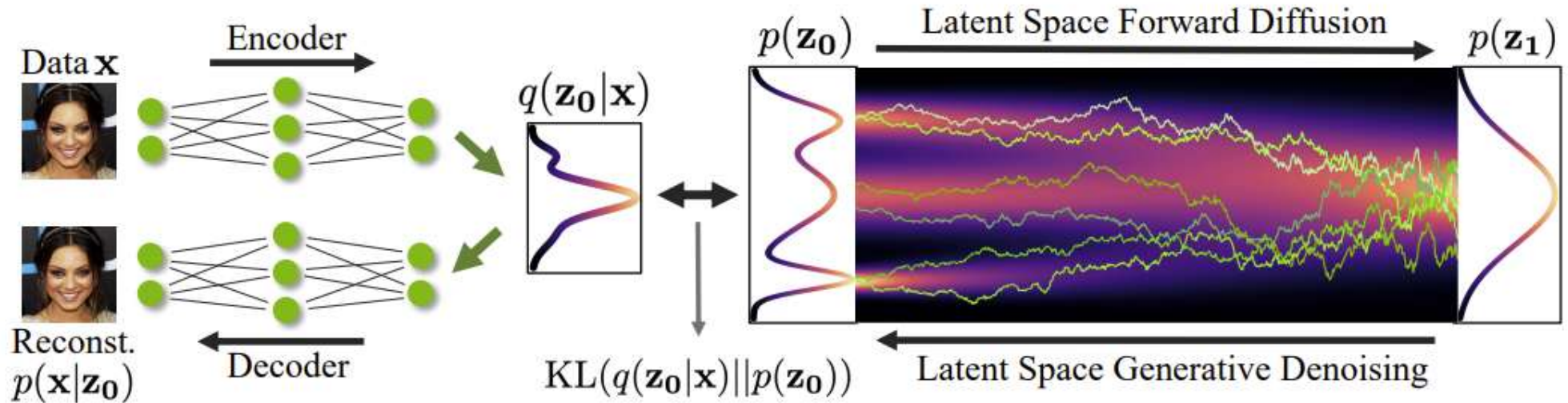


# Stable Diffusion

<https://arxiv.org/abs/2112.10752>



# Latent diffusion model

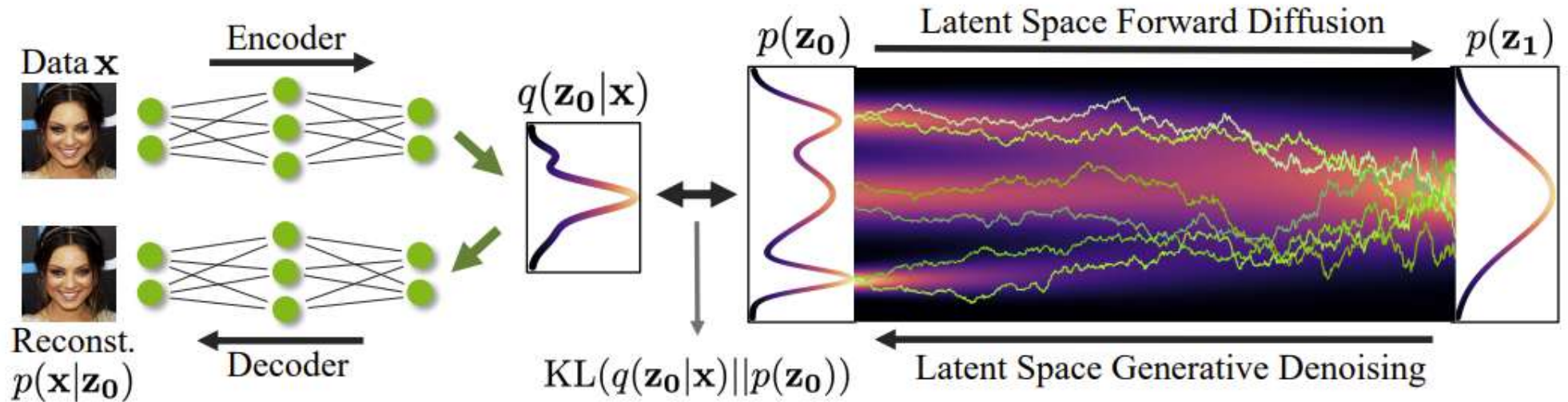


VAE mapping data to lower dimensional space

1. **Faster**
2. **Can be applied to any data type (e.g. discrete)**

Diffusion model prior over the latent space of the autoencoder

# Stable diffusion text2image model



VAE mapping data to lower dimensional space

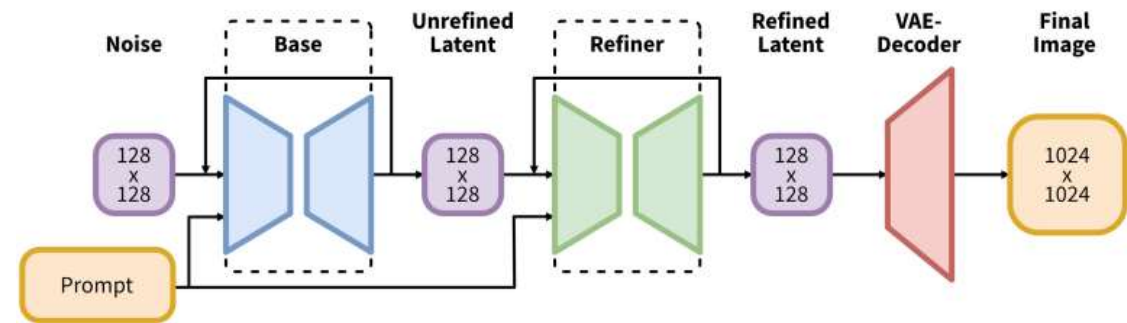
1. Pre-trained, focus on reconstruction (autoencoder)

Diffusion model prior over the latent space of the autoencoder

2. Trained in the second stage, keeping initial autoencoder fixed

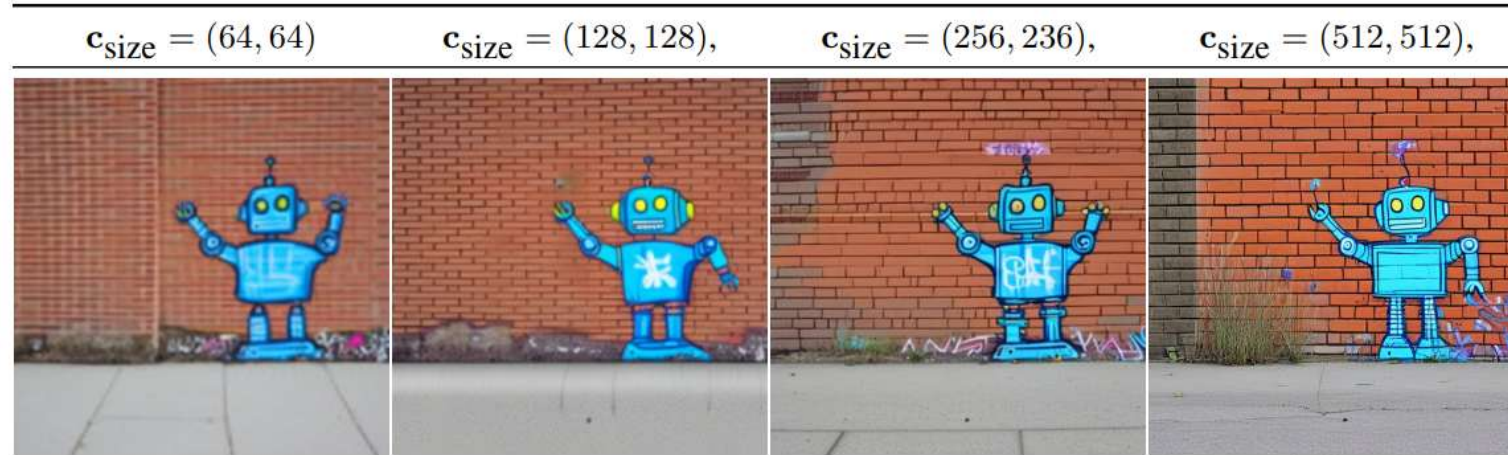
Large scale, open source model, widely adopted

# SDXL: Improving Latent Diffusion Models for High-Resolution Image Synthesis



<https://arxiv.org/pdf/2307.01952>

# SDXL: Improving Latent Diffusion Models for High-Resolution Image Synthesis

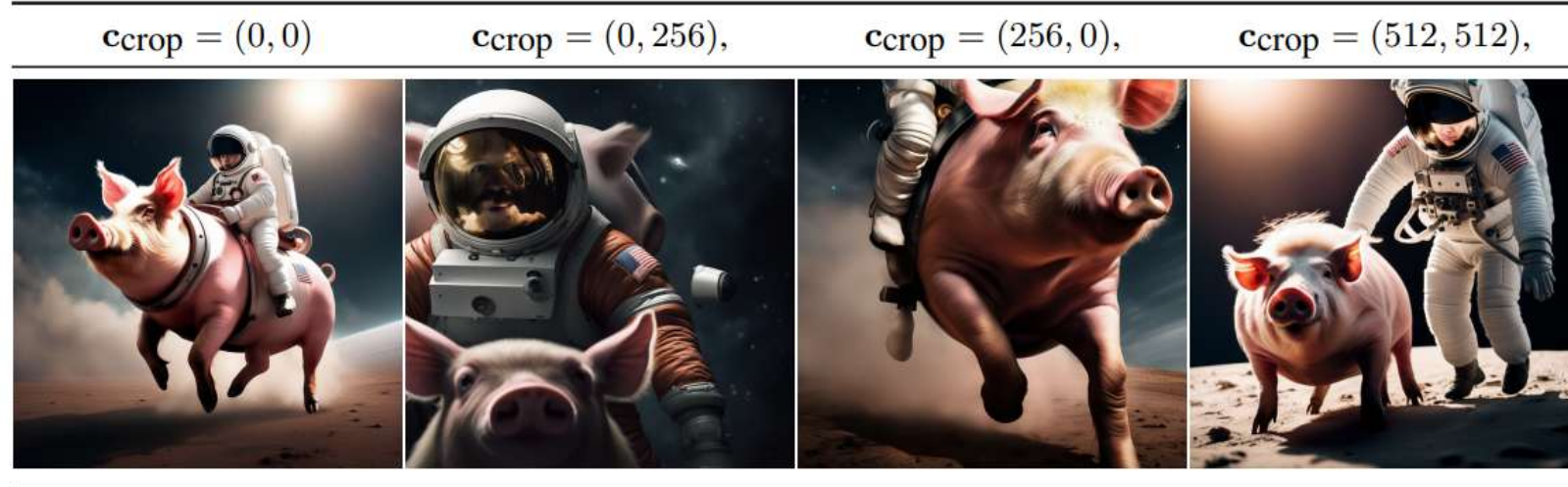


*'A robot painted as graffiti on a brick wall. a sidewalk is in front of the wall, and grass is growing out of cracks in the concrete.'*



*'Panda mad scientist mixing sparkling chemicals, artstation.'*

# SDXL: Improving Latent Diffusion Models for High-Resolution Image Synthesis



*'An astronaut riding a pig, highly realistic dslr photo, cinematic shot.'*



*'A capybara made of lego sitting in a realistic, natural field.'*

# SDXL: Improving Latent Diffusion Models for High-Resolution Image Synthesis

*'A propaganda poster depicting a cat dressed as french emperor napoleon holding a piece of cheese.'*

*'a close-up of a fire spitting dragon, cinematic shot.'*

SD 1-5

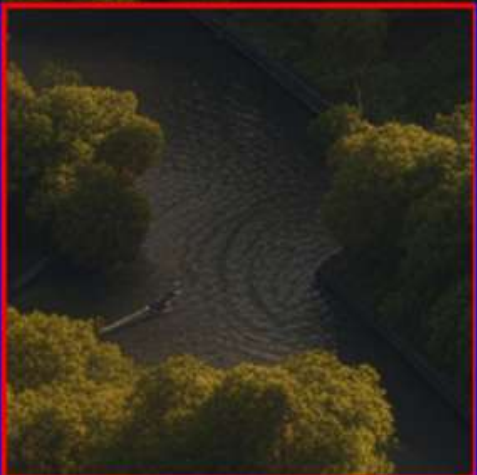


SD 2-1

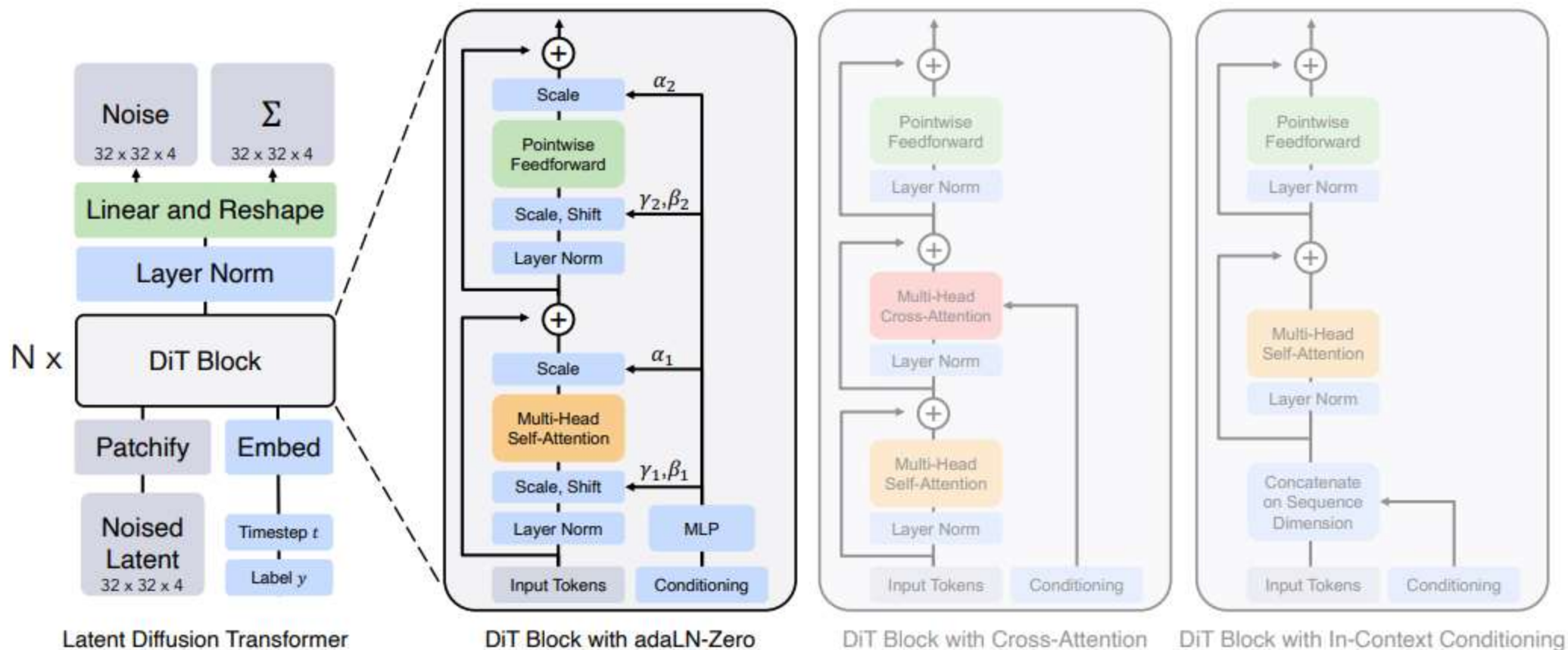


SDXL

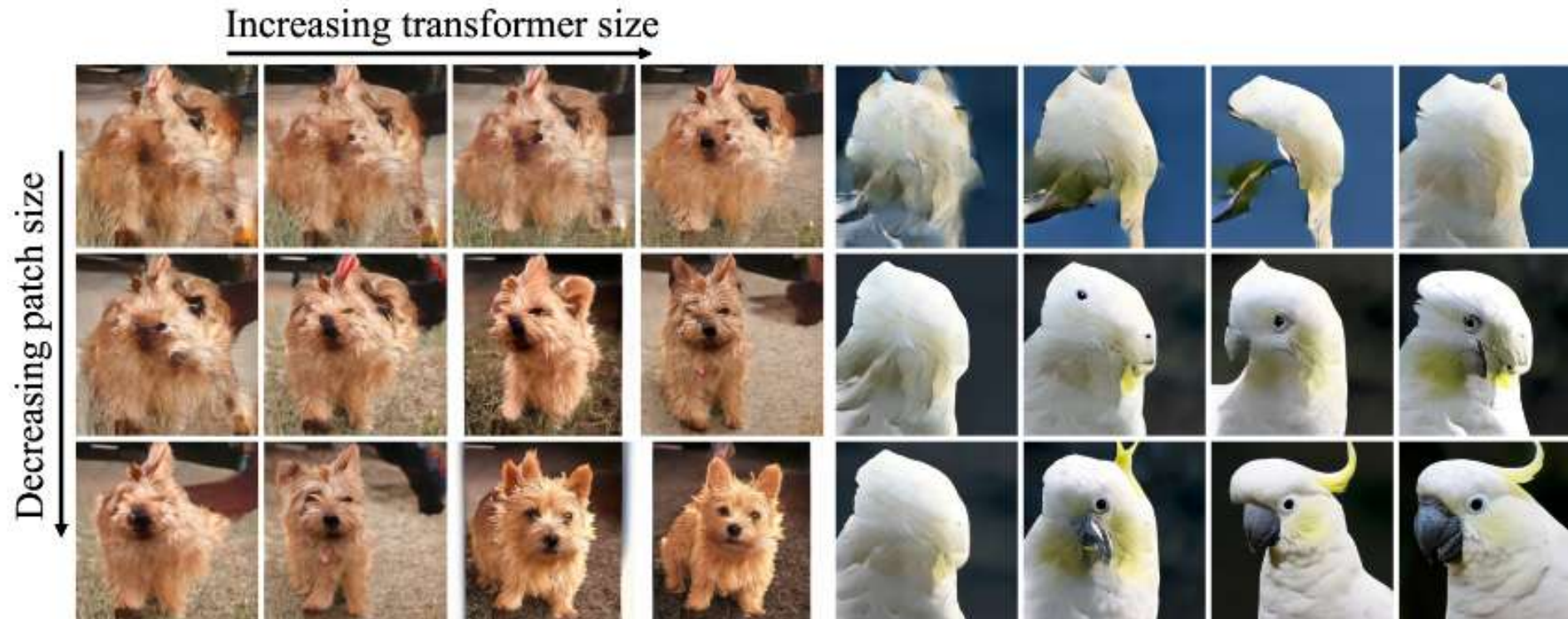




# DiT: Scalable Diffusion Models with Transformers



# DiT: Scalable Diffusion Models with Transformers



# DiT: Scalable Diffusion Models with Transformers



DiT "Dog-Ball"



BigGAN "Dog-Ball"

# DiT: Scalable Diffusion Models with Transformers



<https://arxiv.org/pdf/2212.09748>

# CLIP is good but...

It cannot tell these apart!



(a) some plants surrounding a lightbulb



(b) a lightbulb surrounding some plants

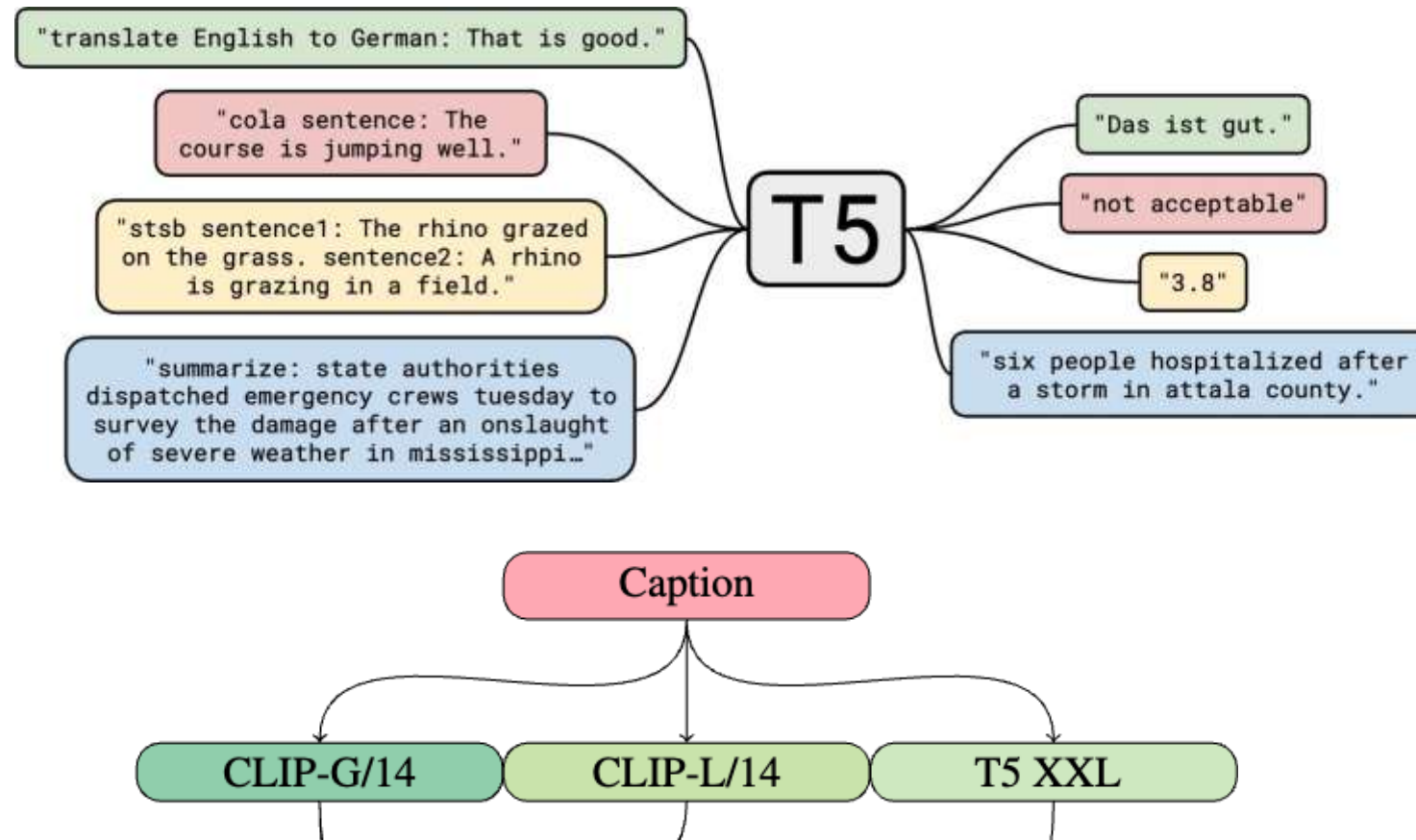
# CLIP is good but...

- It cannot handle spatial relationships well
- It cannot handle negations well
- It cannot handle counts well
- The length limit is 77 tokens
- Sometimes ignores some details

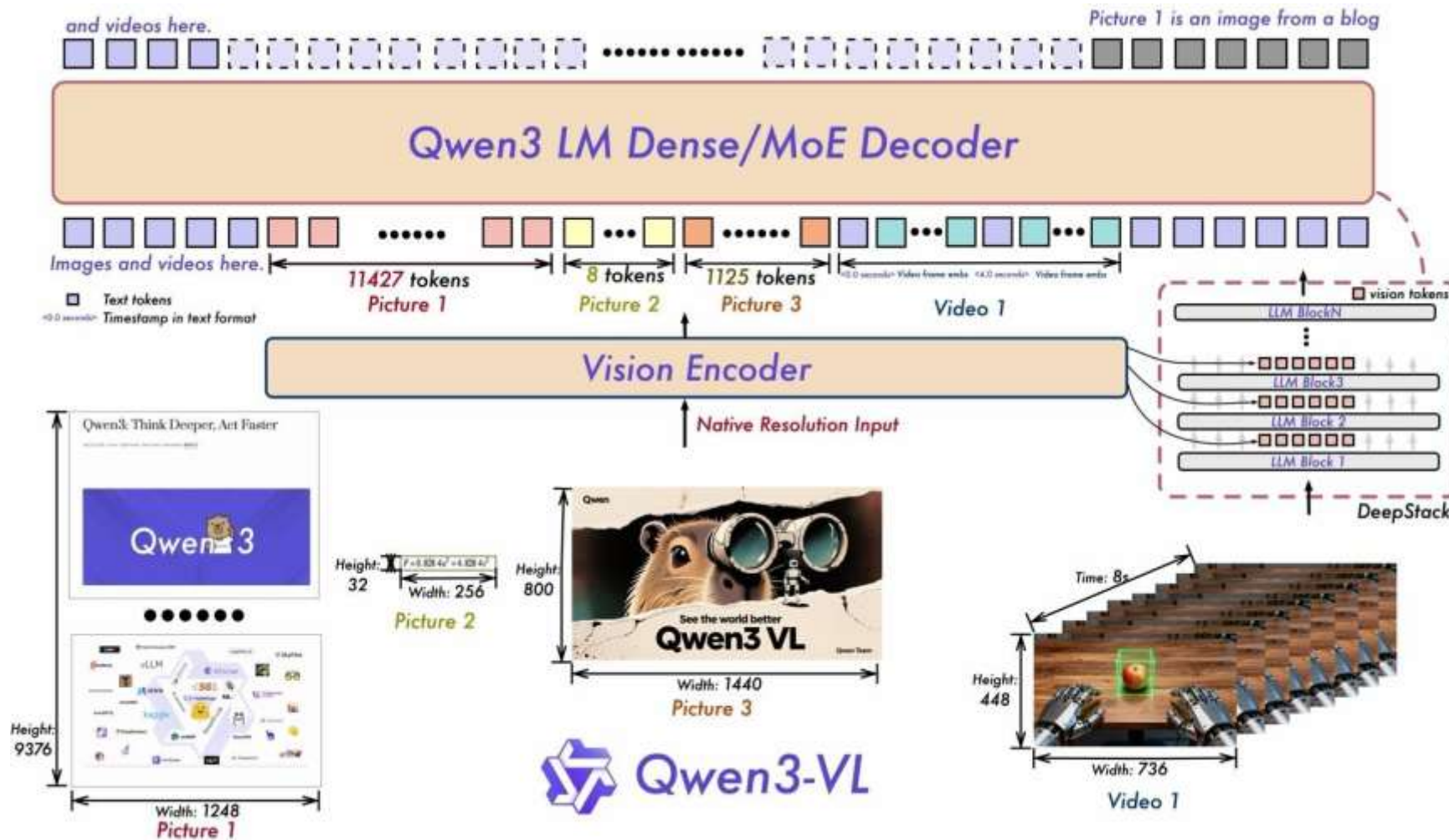
# CLIP is not a very good text encoder!

- What we want:
- Can capture semantics well ✘
- Can encode long sentences ✘
- Can attend to details ✘
- Can differentiate between different spatial relationships described in text ✘

# Attempt : Add another text encoder on top of CLIP for better language understanding



# Attempt : Why not just use an LLM/VLM/MLLM



# How to input text into an image generative model?

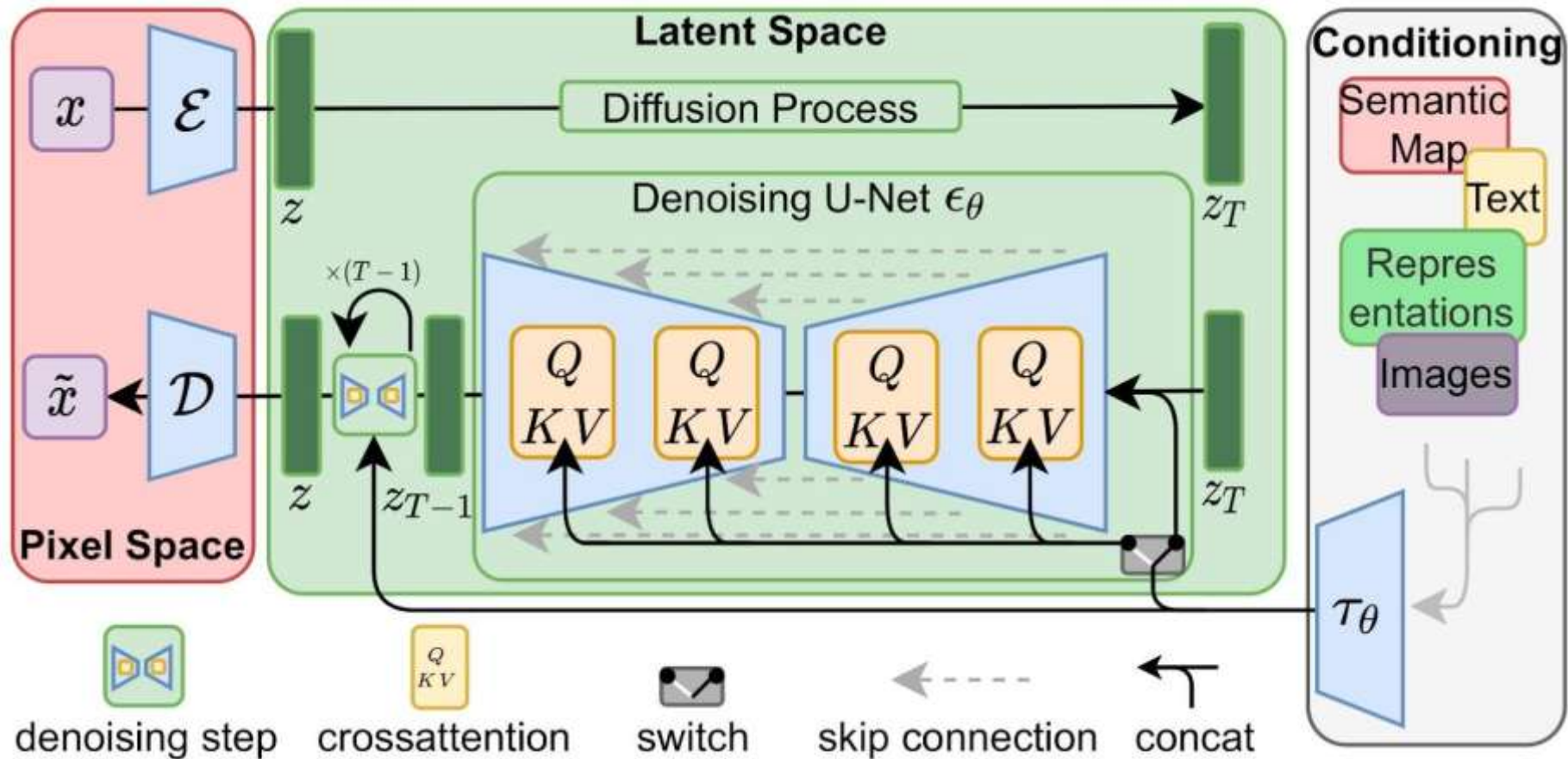
1. Somehow encode text into some feature vectors



2. Somehow squish the encoded text features into the diffusion model

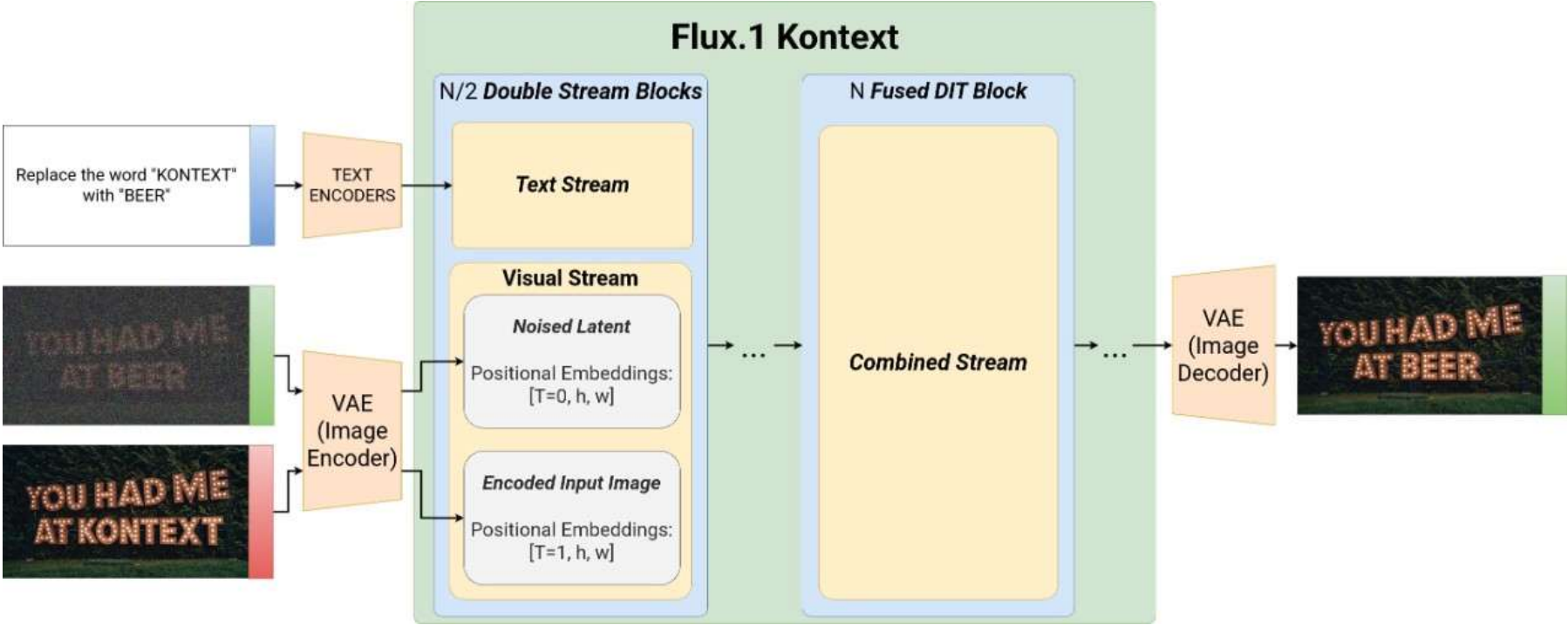


# Attempt 1: Cross Attention

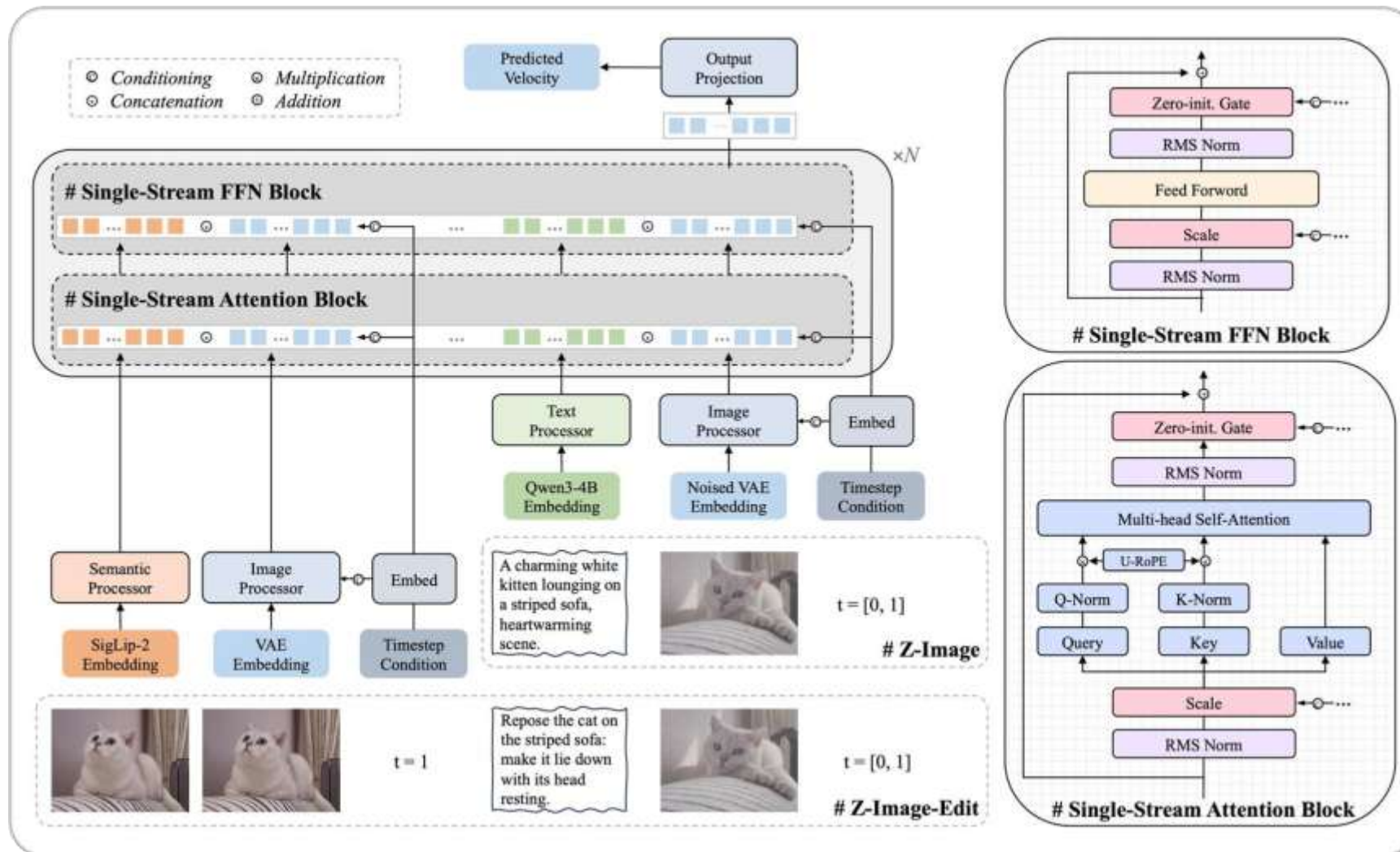




# Attempt 2: Double stream -> merged stream MM-DiT



# Attempt 3: Single stream MM-DiT



# Attempt 3.5: “Native multimodal model”

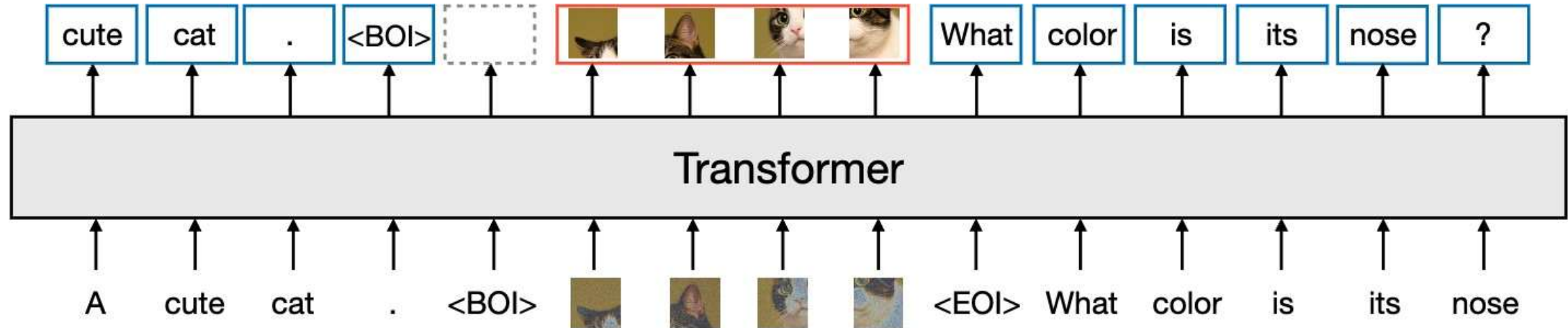
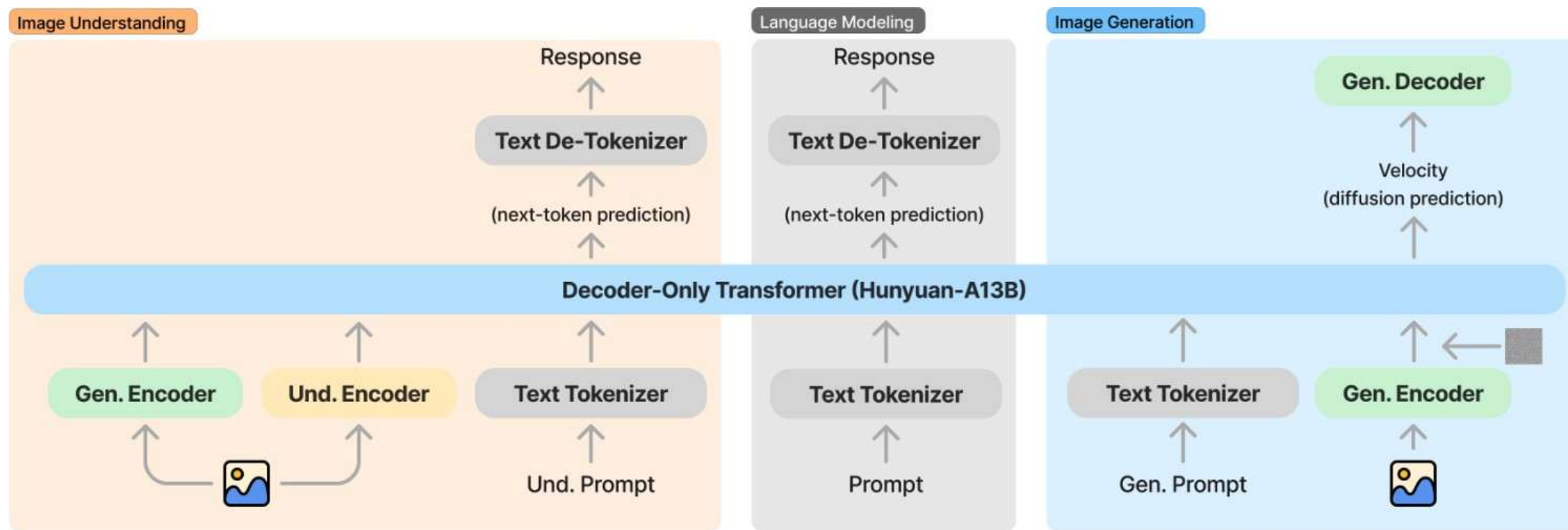
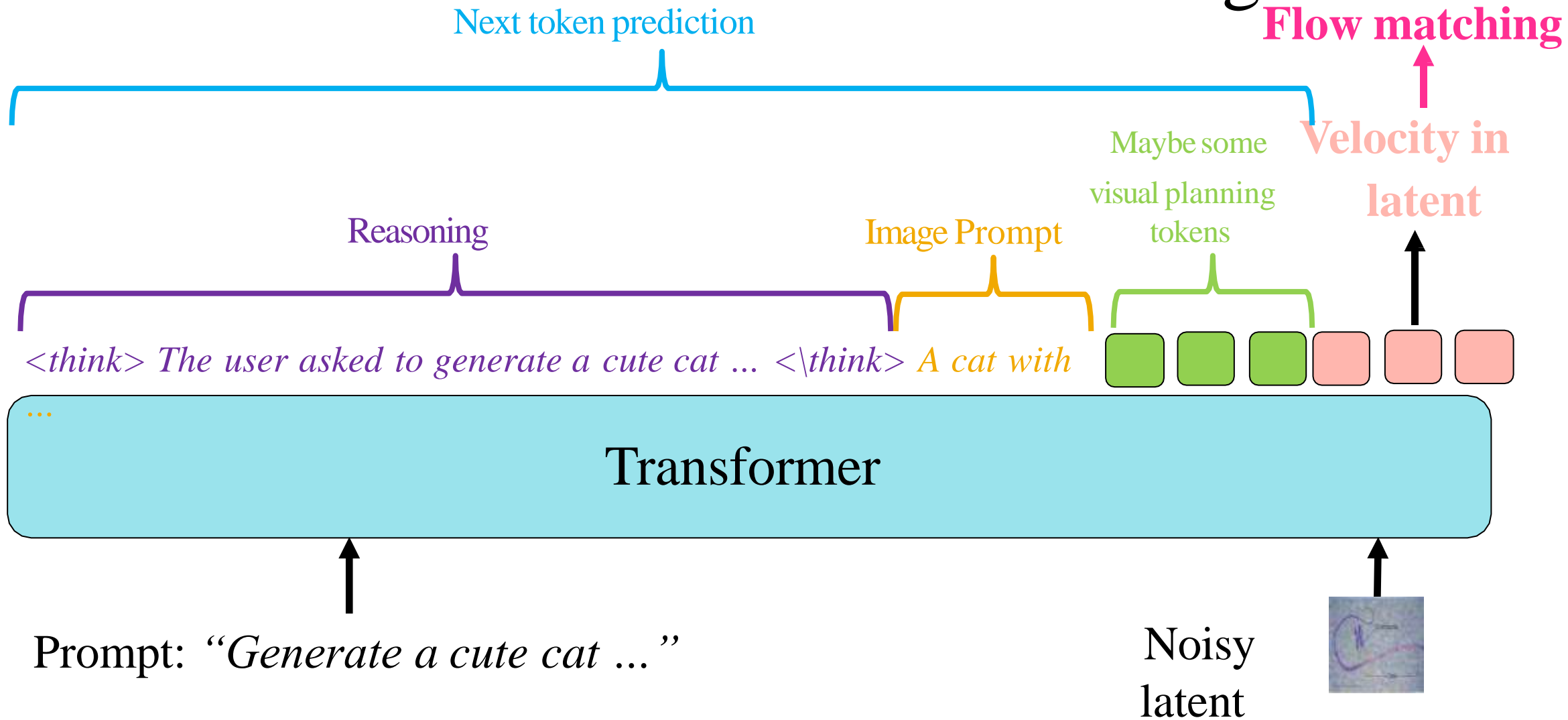


Figure 1: A high-level illustration of Transfusion. A single transformer perceives, processes, and produces data of every modality. Discrete (text) tokens are processed autoregressively and trained on the **next token prediction** objective. Continuous (image) vectors are processed together in parallel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

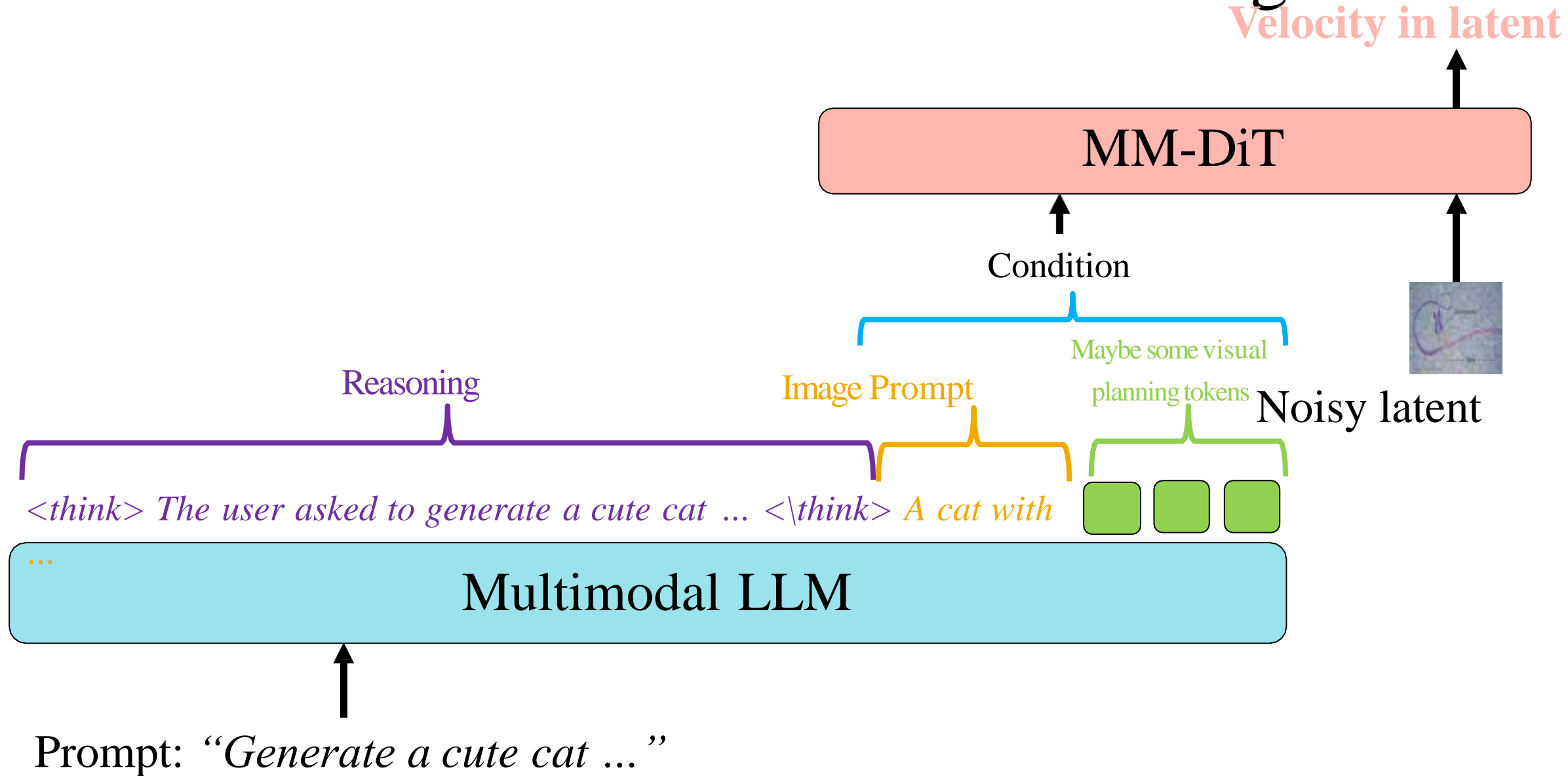
# Attempt 3.5: “Native multimodal model”



# 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)

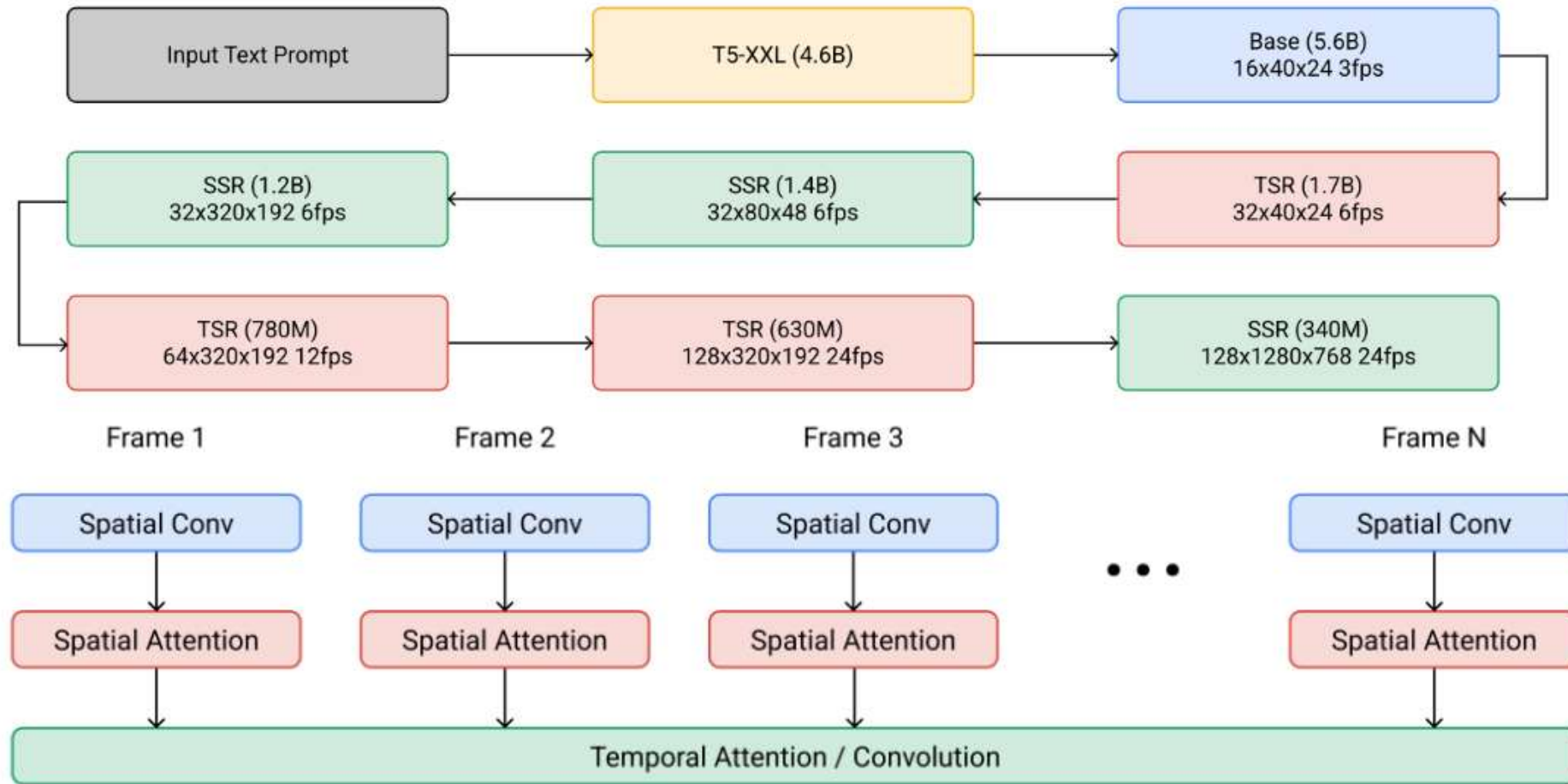
# **Video Generation**

# Sora





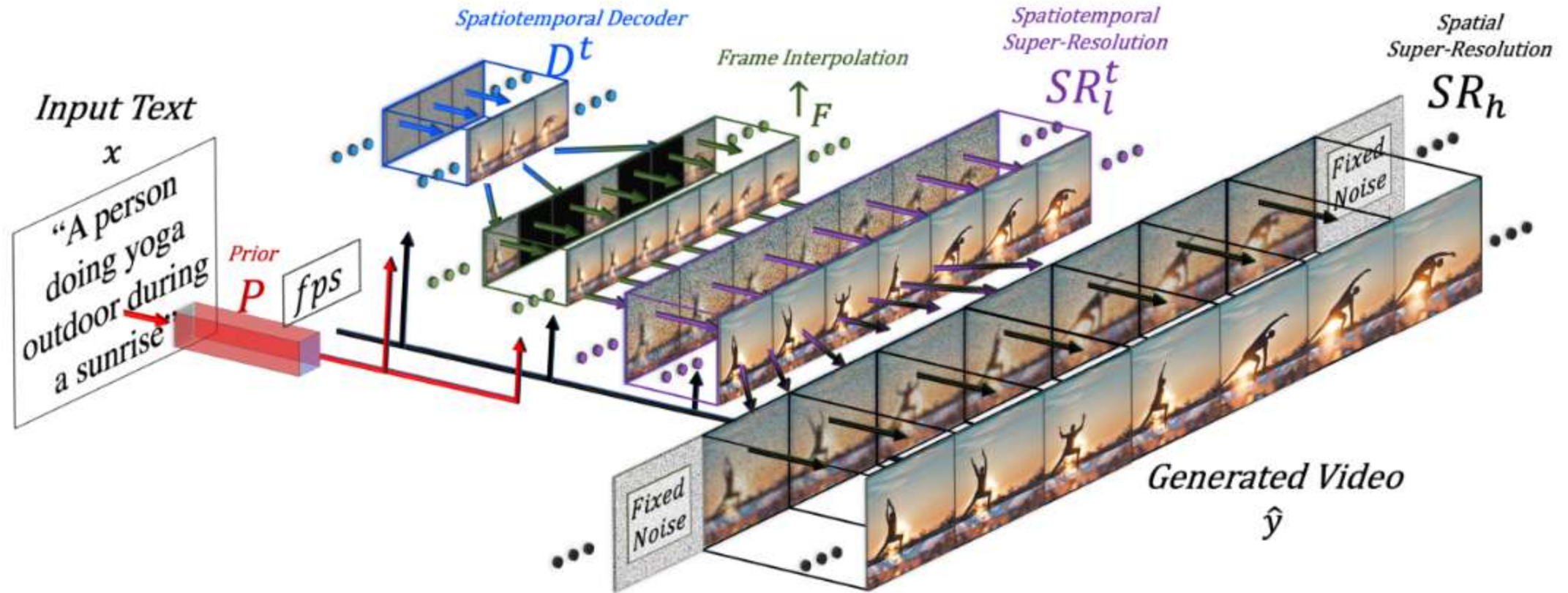
# Imagen Video



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.

# Adapting Image Models to Generate Videos



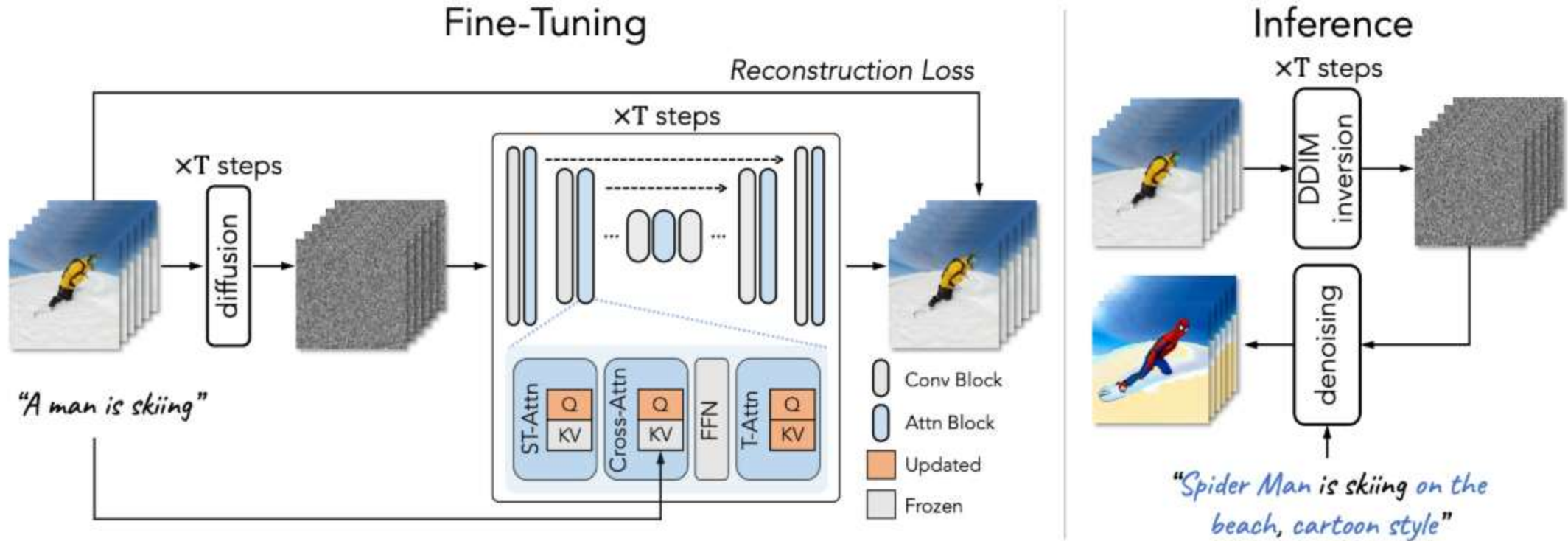
$$\hat{y}_t = SR_h \circ SR_l^t \circ \uparrow_F \circ D^t \circ P \circ (\hat{x}, CLIP_{\text{text}}(x))$$

# Adapting Image Models to Generate Videos



<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

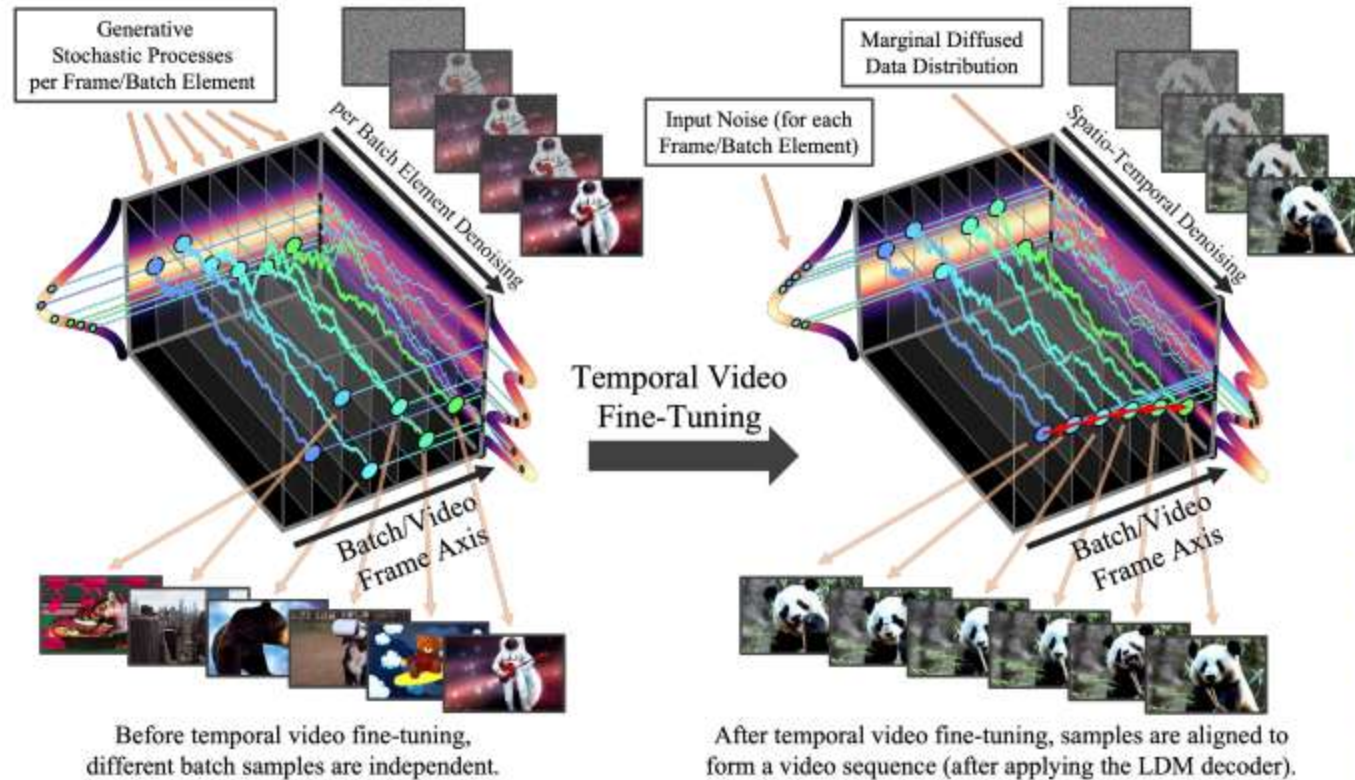


Figure 2. **Temporal Video Fine-Tuning.**

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

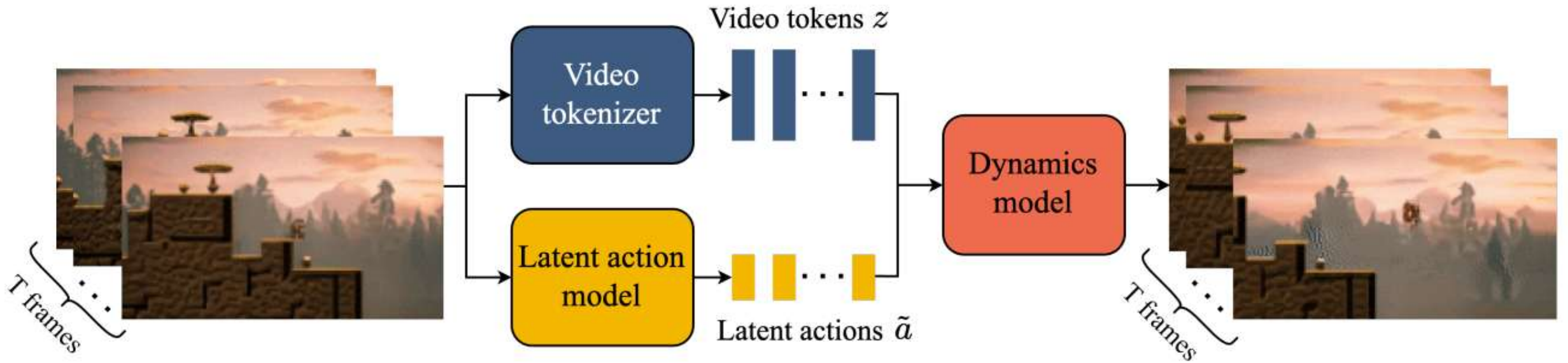


# Genie: Generative Interactive Environments



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



“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

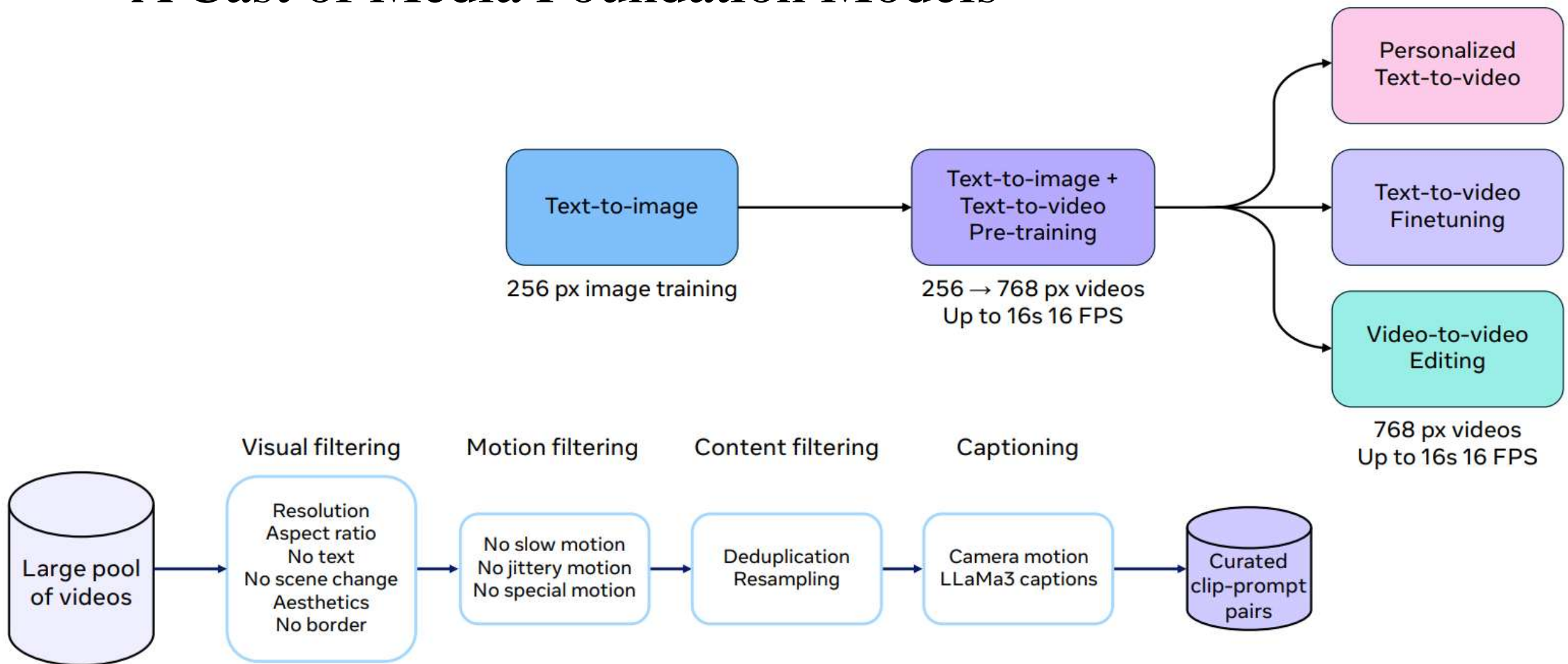
“Trajectories with the same latent action sequence typically display similar behaviors.”

# Movie Gen: A Cast of Media Foundation Models

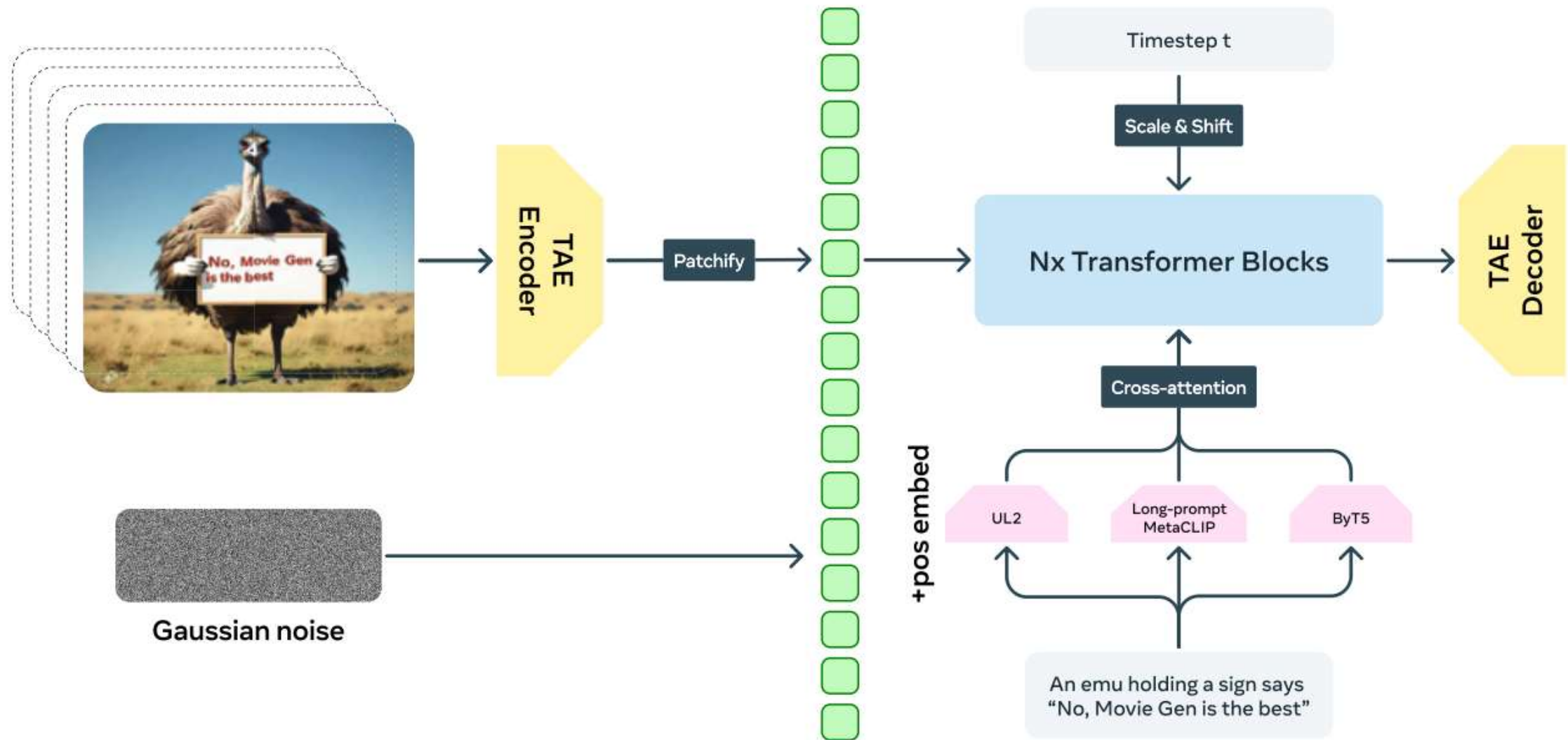


<https://ai.meta.com/blog/movie-gen-media-foundation-models-generative-ai-video/>

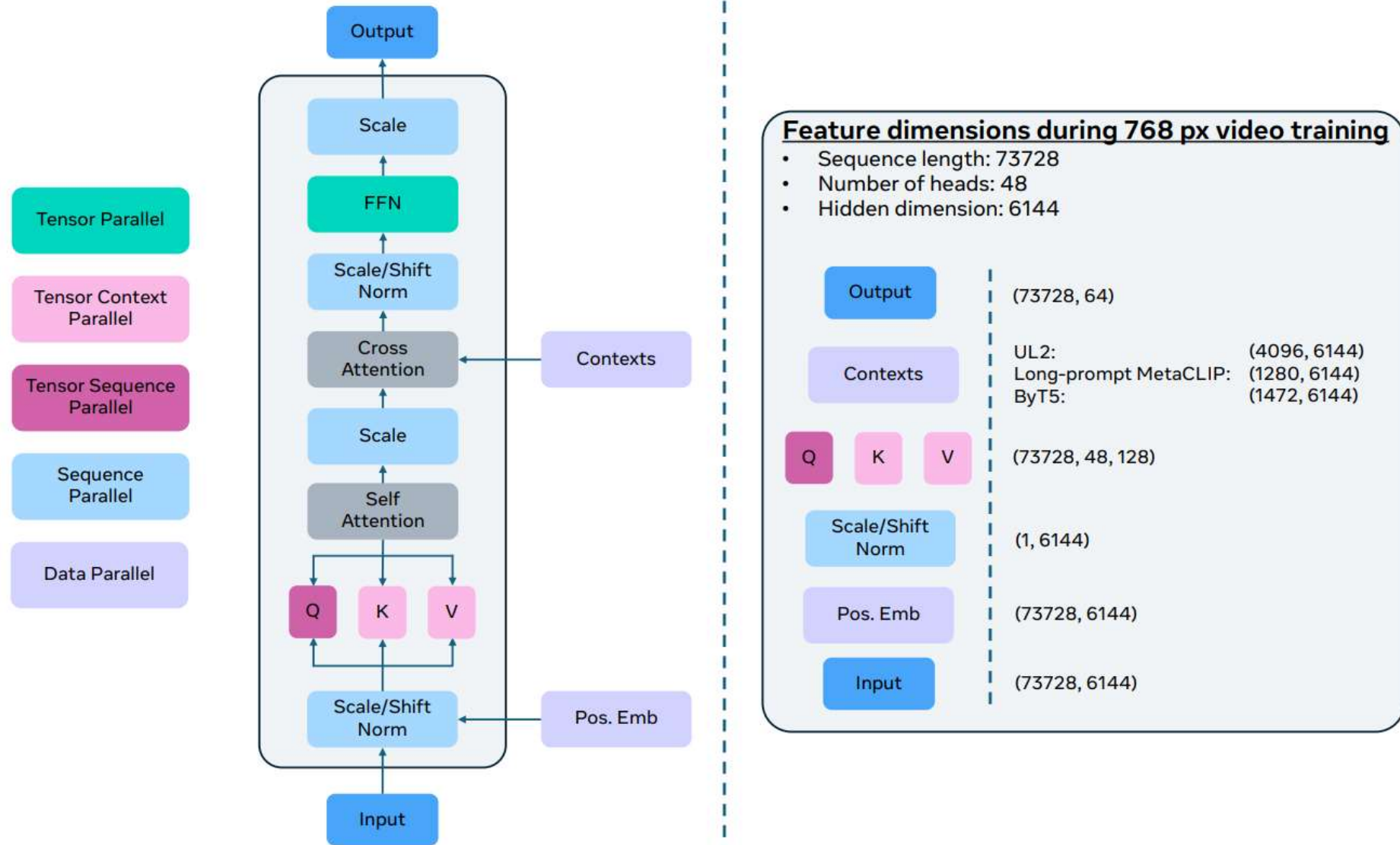
# Movie Gen: A Cast of Media Foundation Models



# Movie Gen: A Cast of Media Foundation Models



# Movie Gen: A Cast of Media Foundation Models



# One Model to Workflow

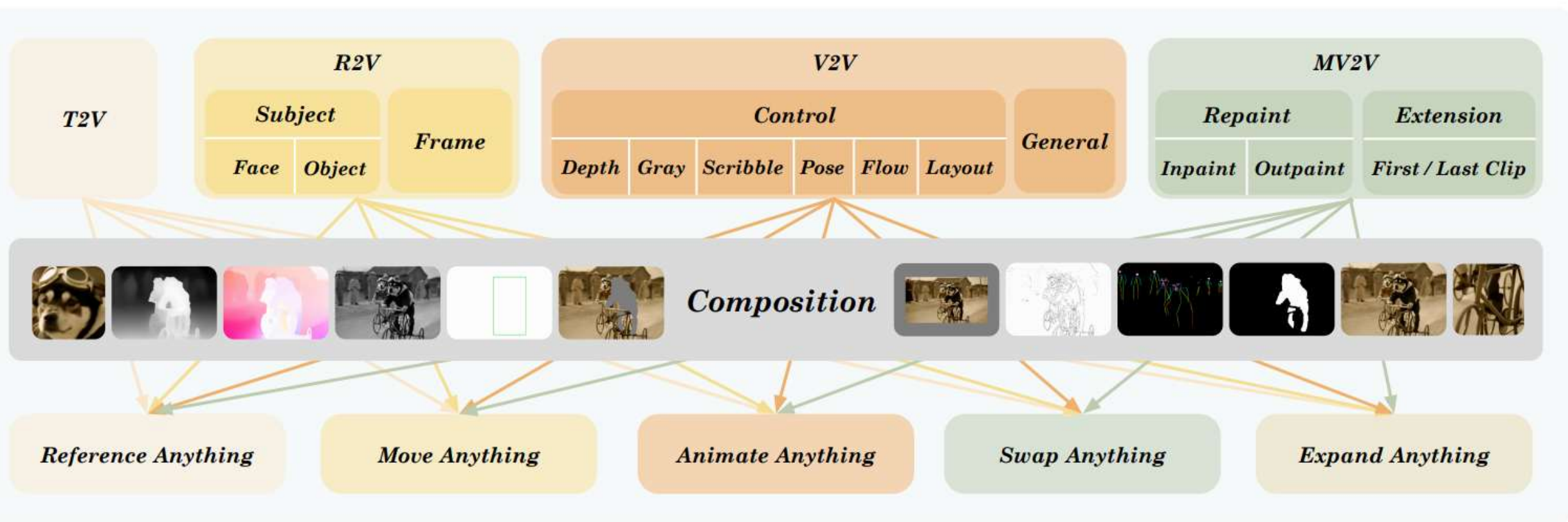


\*image taken from ChatGPT 4o

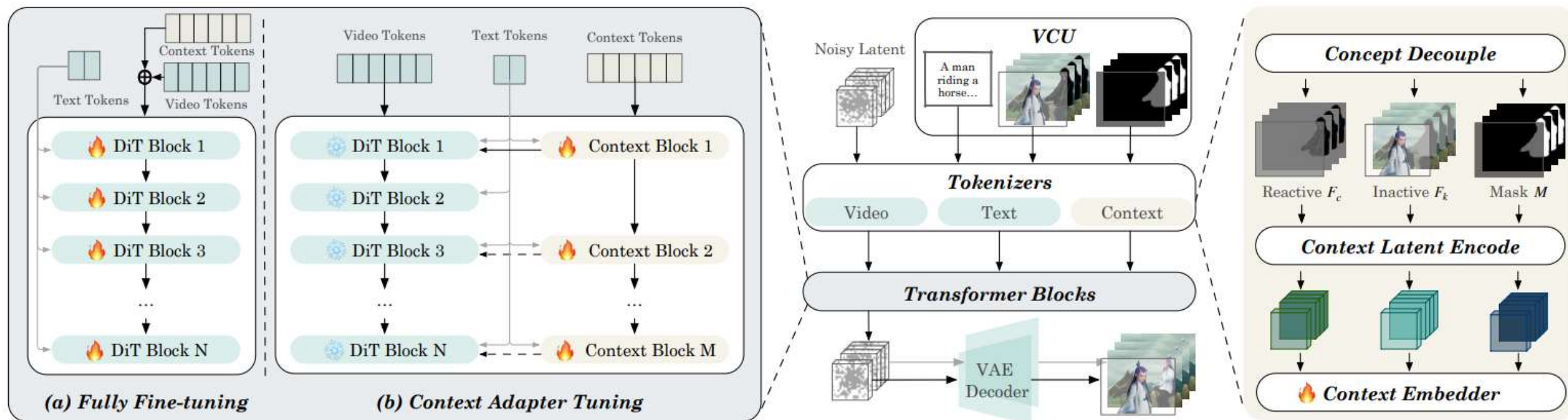
# VACE/Wan



# 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.

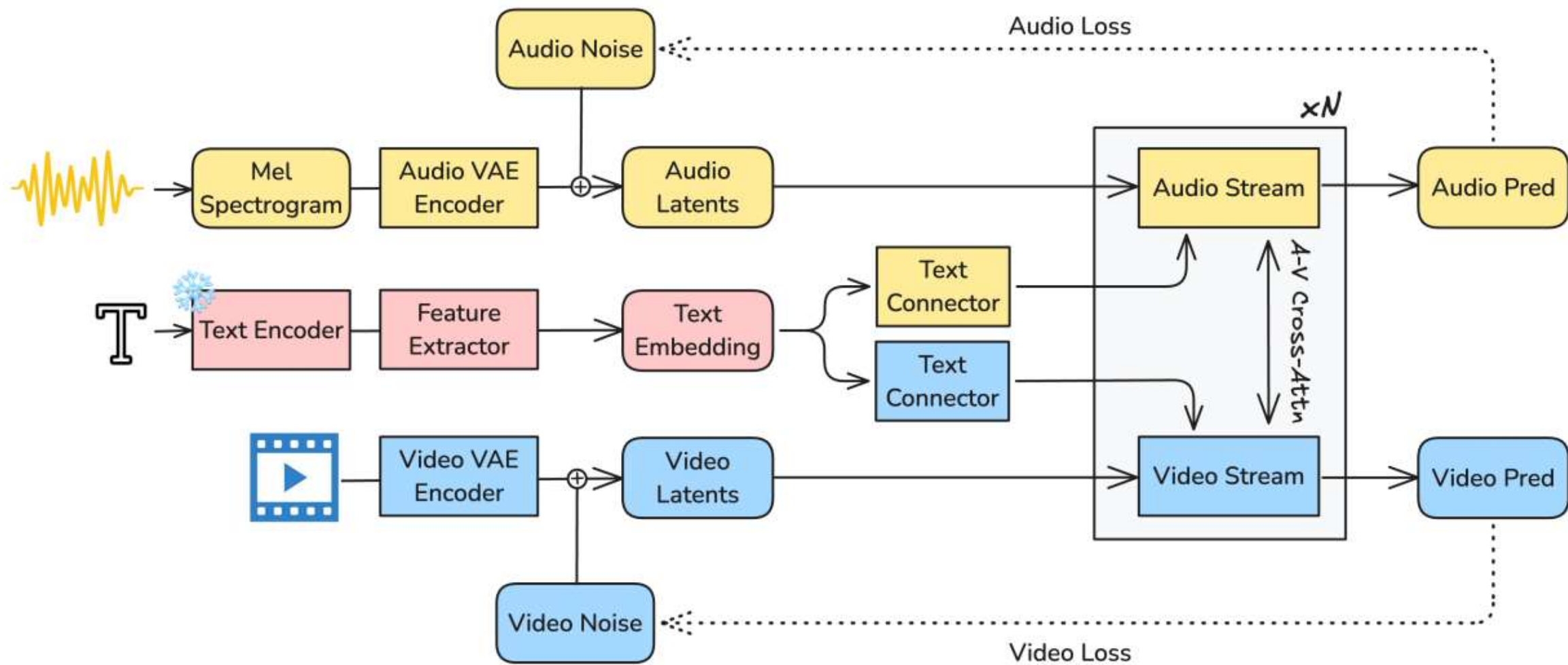
- Sora 2 improves physics, realism, synchronized audio, and steerability.
- Prompts now combine subject, camera, pacing, and audio intent.
- Storyboard, remix, stitch, extend, and characters define the workflow.

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

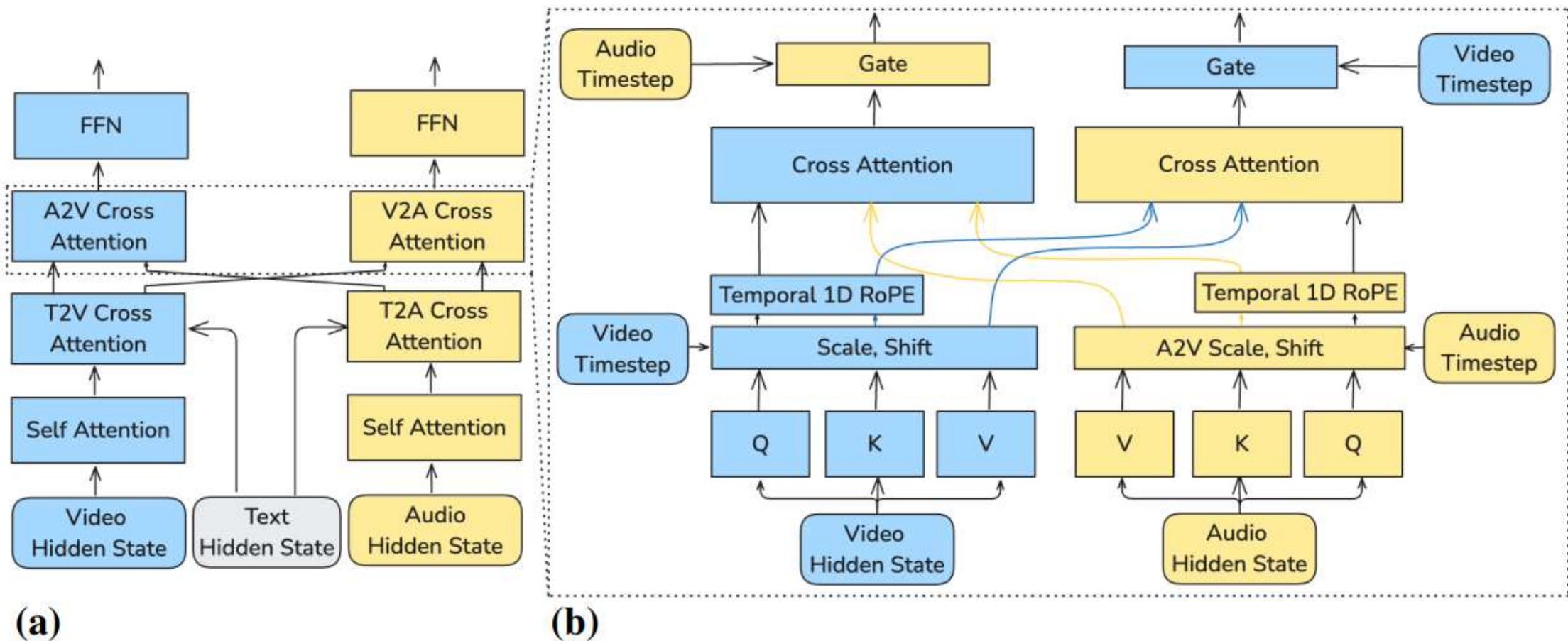
- Veo 3.1 adds richer native audio and stronger prompt adherence.
- Reference images, first/last frames, and extension expose creative control.
- A strong Veo prompt reads like a director's brief.

[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

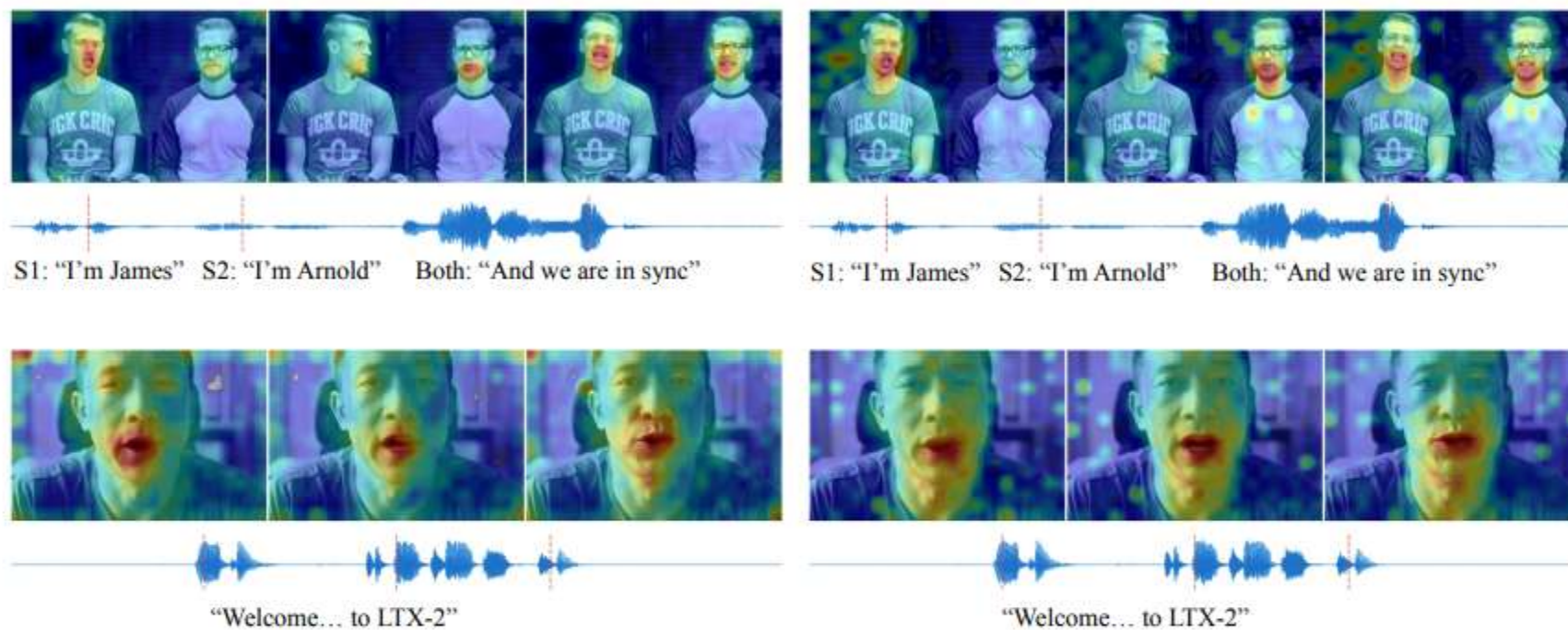


Figure 3: **Visualization of AV cross-attention maps.** The maps are averaged across attention heads and the model layers; V2A and A2V maps correspond to the first and last 1/3 of inference steps, respectively. Red vertical lines on the audio waveform mark the timestamps of the displayed frames. The 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