

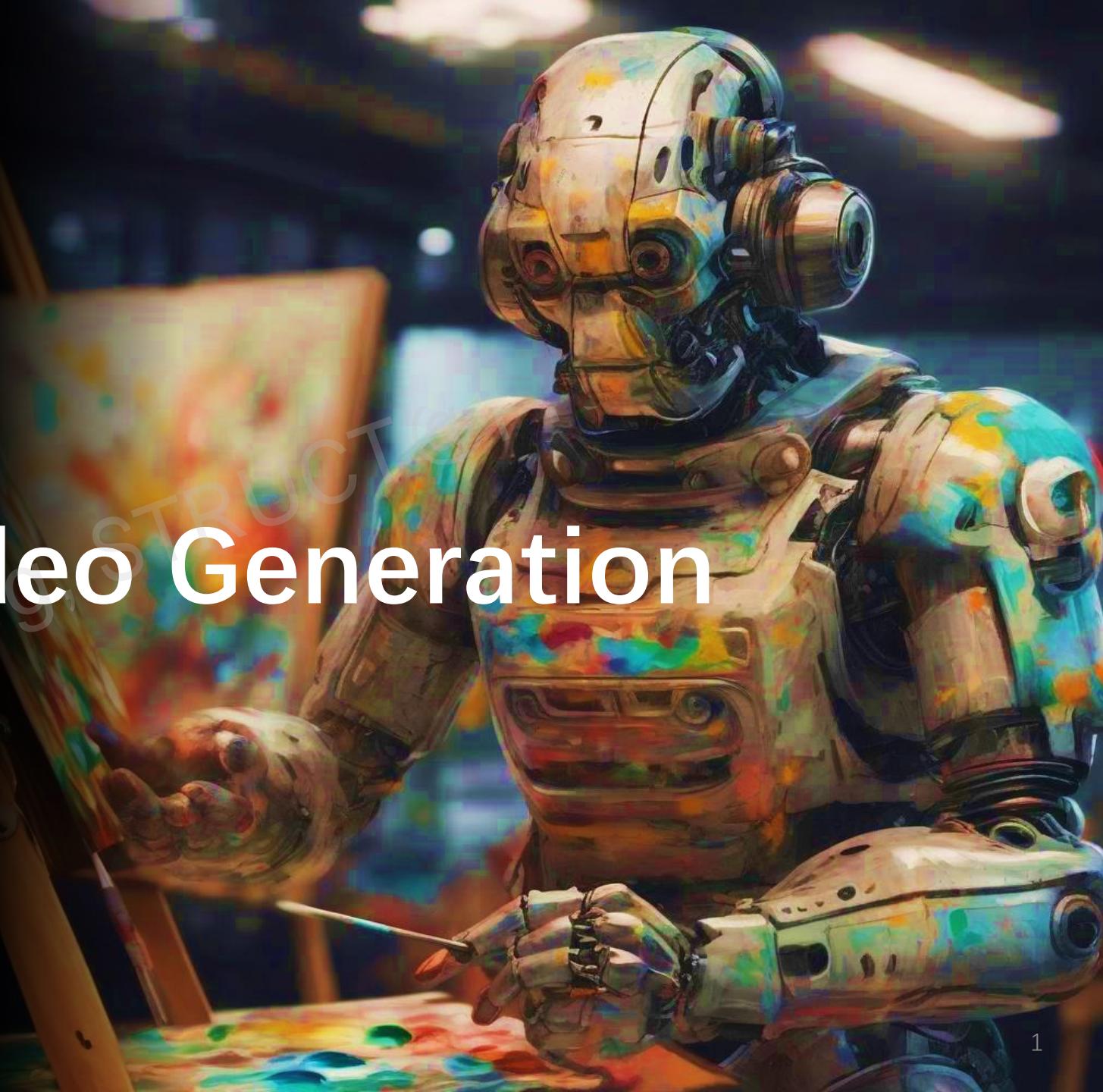
ISCAS 2024

AIGC for Image and Video Generation

Presenter: Shuai Yang



北京大学
PEKING UNIVERSITY



AI for image generation

Development of generative model

- Higher resolution
- Richer content
- More controllable



VAE
(2013)



GAN
(2014)



DCGAN
(2016)



ProGAN
(2018)



StyleGAN2
(2020)



StableDiffusion
(2022)

Traditional Era



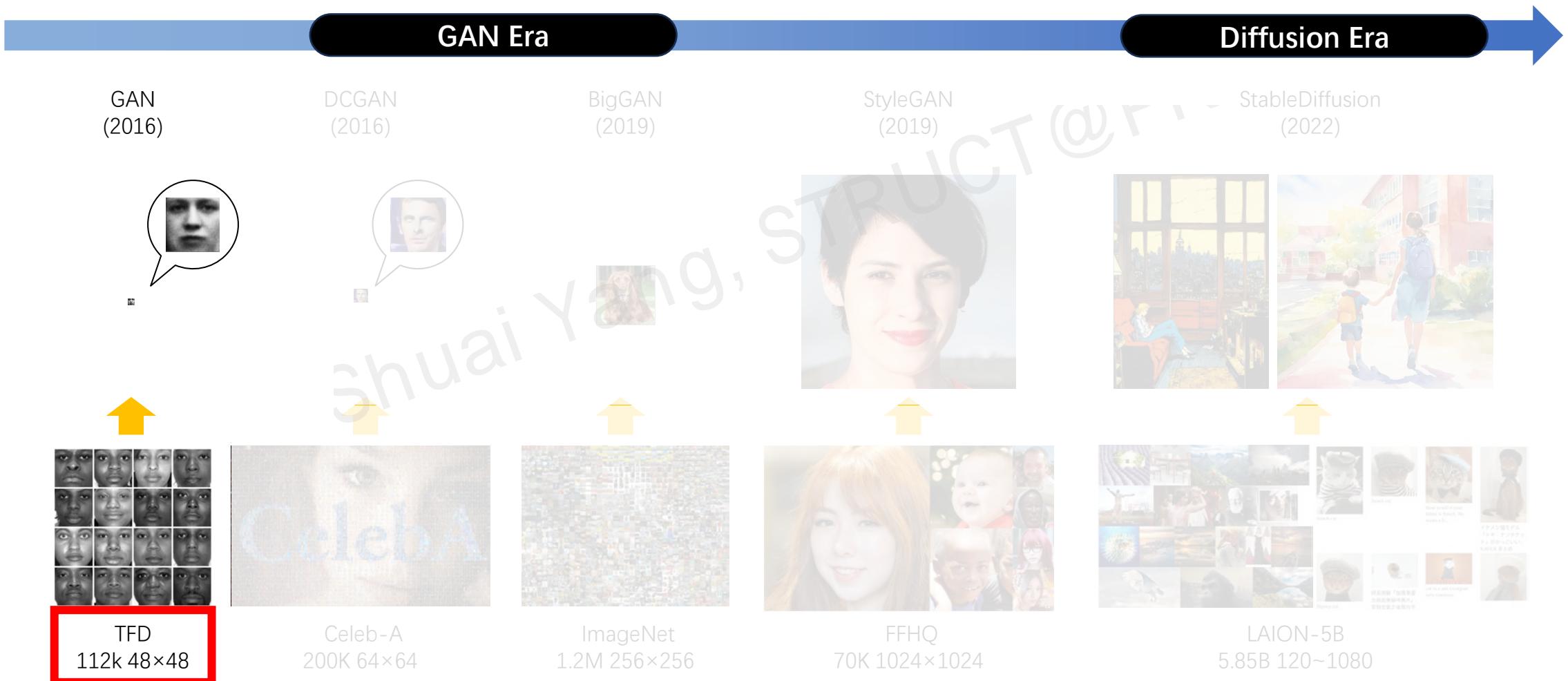
GAN Era



Diffusion Era

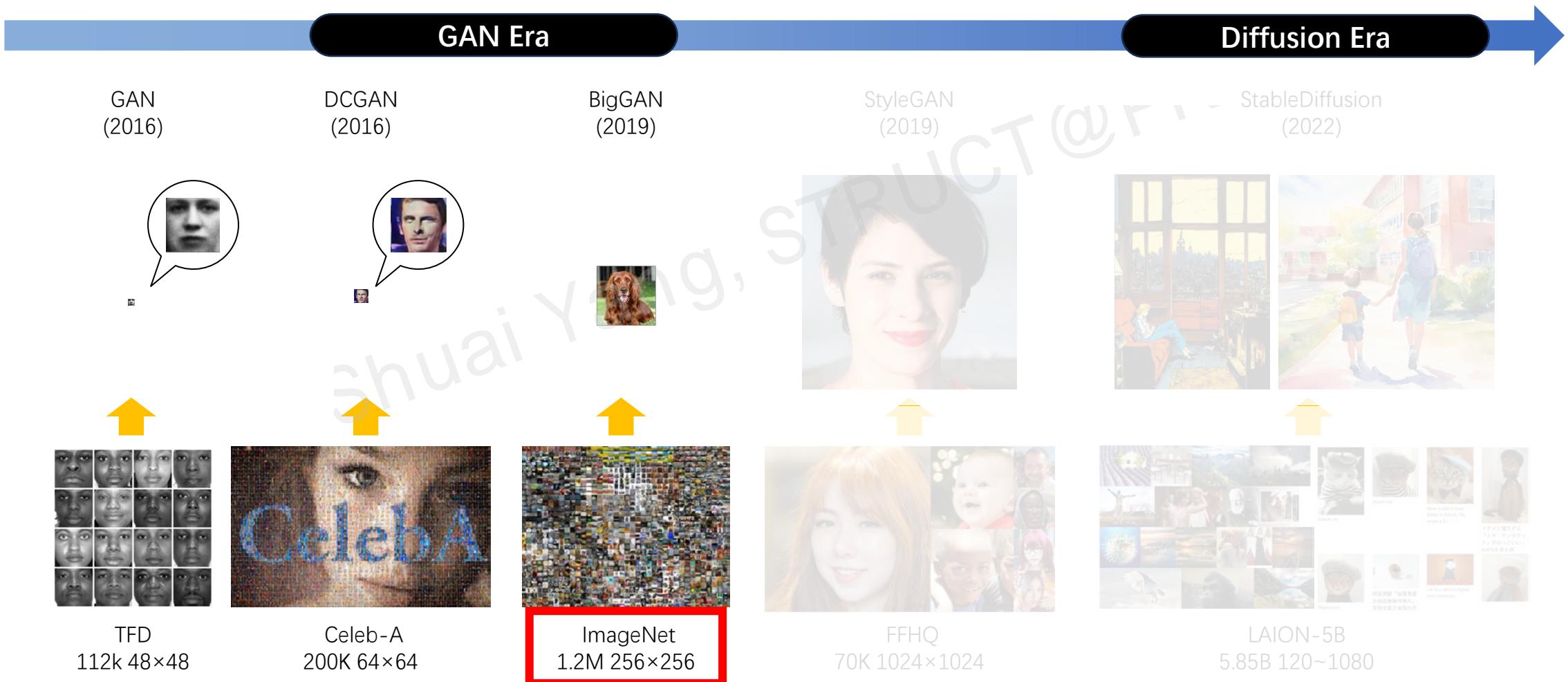
AI for image generation

As the data scale increases, the generation quality has been greatly improved.



AI for image generation

As the data scale increases, the generation quality has been greatly improved.



AI for image generation

As the data scale increases, the generation quality has been greatly improved.

GAN Era

GAN
(2016)



DCGAN
(2016)



BigGAN
(2019)



StyleGAN
(2019)

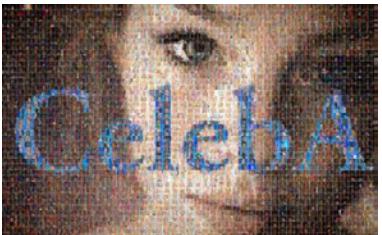


Diffusion Era

StableDiffusion
(2022)



TFD
112k 48×48



Celeb-A
200K 64×64



ImageNet
1.2M 256×256



FFHQ
70K 1024×1024



LAION-5B
5.85B 120~1080

AI for image generation

The enhancement of model representation capabilities benefits from the development of computing power and model scale.

GAN Era

DCGAN
(~70MB)



StyleGAN
(~300MB)



BigGAN
(~200MB)



Diffusion Era

StableDiffusion
(~4GB)



Single theme per model

Human face, animal face, indoor, vehicle

1000 themes per model

Arbitrary themes per model

AI for image generation

Opportunities in Diffusion Era

- Any theme: draw what I want
- **Diverse style support**

Diversity of Theme

greater expressive power than ever



AI for image generation

Opportunities in Diffusion Era

- Any theme: draw what I want
- **Diverse style support**

- What you say is what you get
- **Interactive generation**

Interactive and Controllable Generation

combined with LLMs to achieve interactive and convenient generation



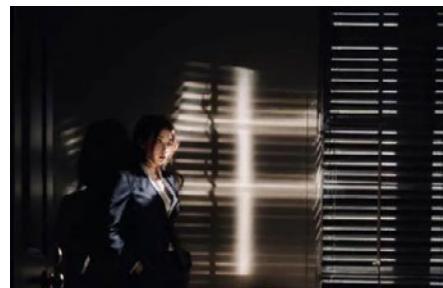
AI for image generation

Opportunities in Diffusion Era

- Any theme: draw what I want
- **Diverse style support**
- What you say is what you get
- **Interactive generation**
- Creative generation and fusion
- **Powerful generation ability**

Challenge “Only human can do creative work”

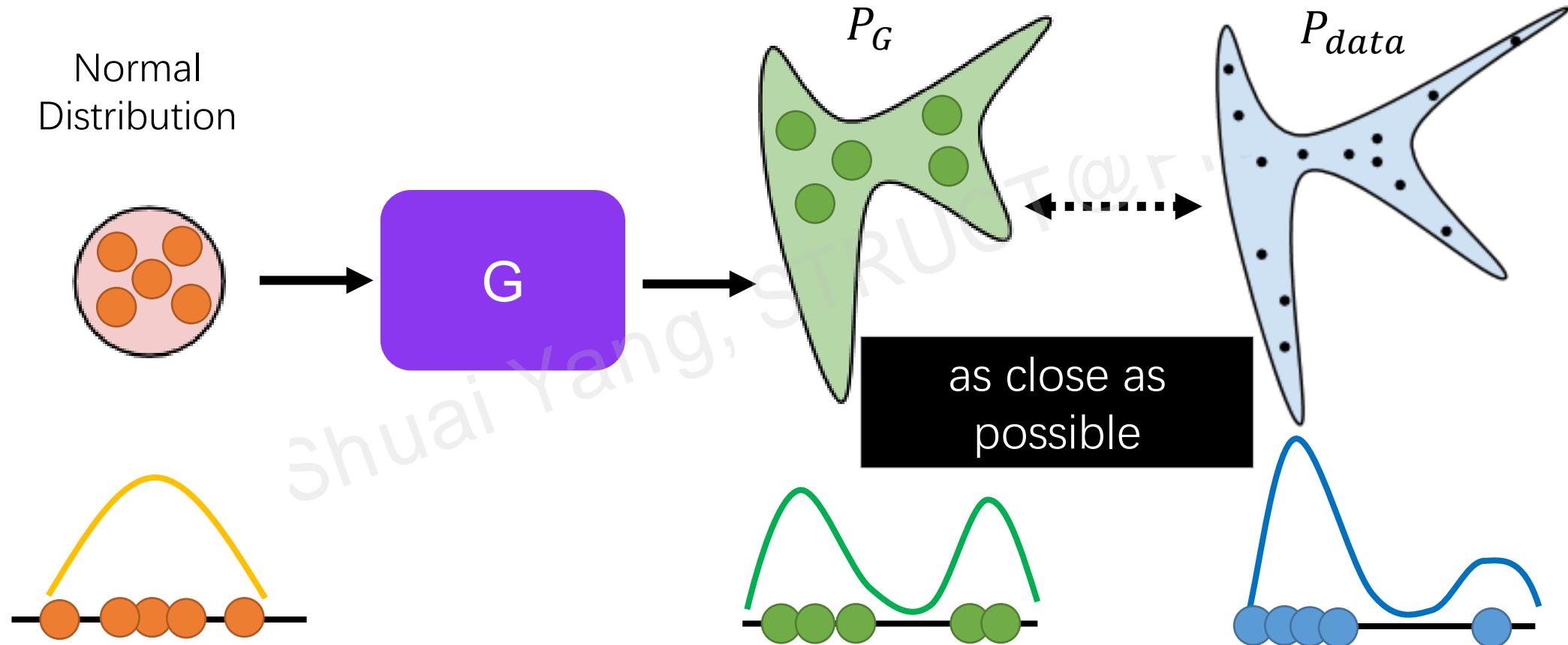
creative integration and emergence driven by large-scale multi-modal data



A robot with a metallic, multi-colored body is painting a landscape on a canvas. The landscape features a blue sky with white clouds, green trees, and a red sunset. The robot's paintbrush is visible, applying color to the scene. The background is dark and abstract.

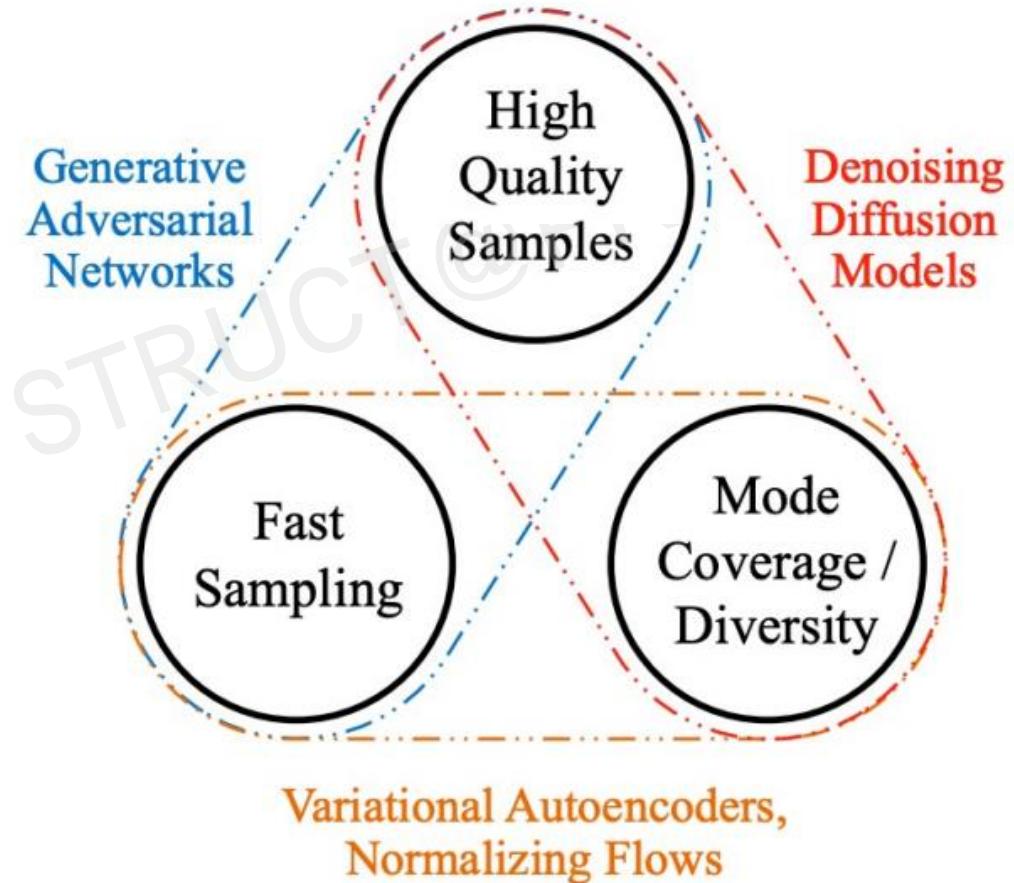
Diffusion

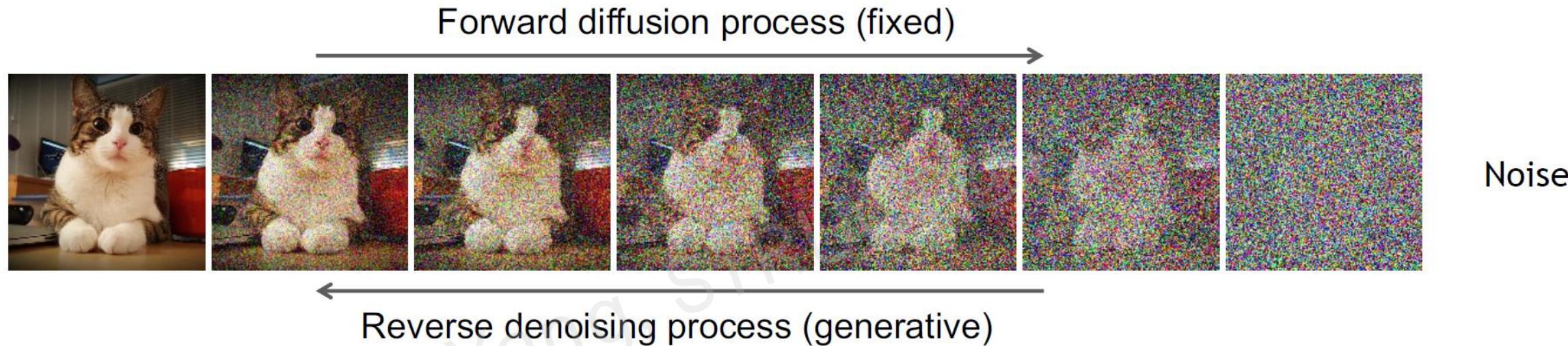
Yang, STRUCTURE



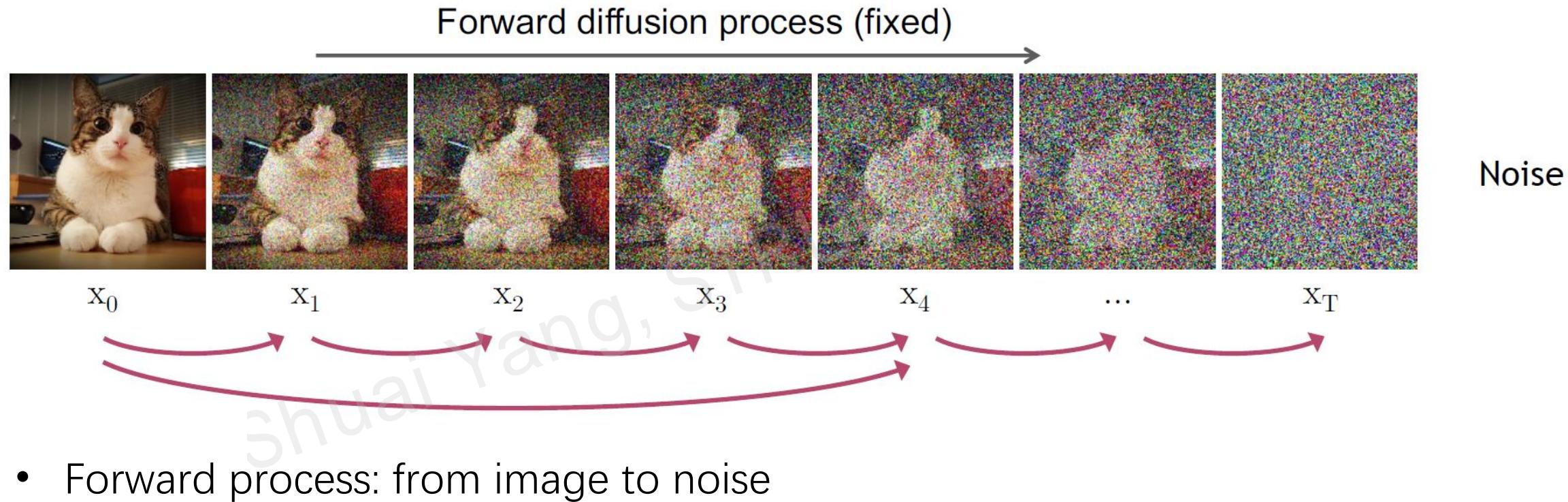
Previous generative models:

- VAEs: Lack of image quality
- GANs: Limited diversity

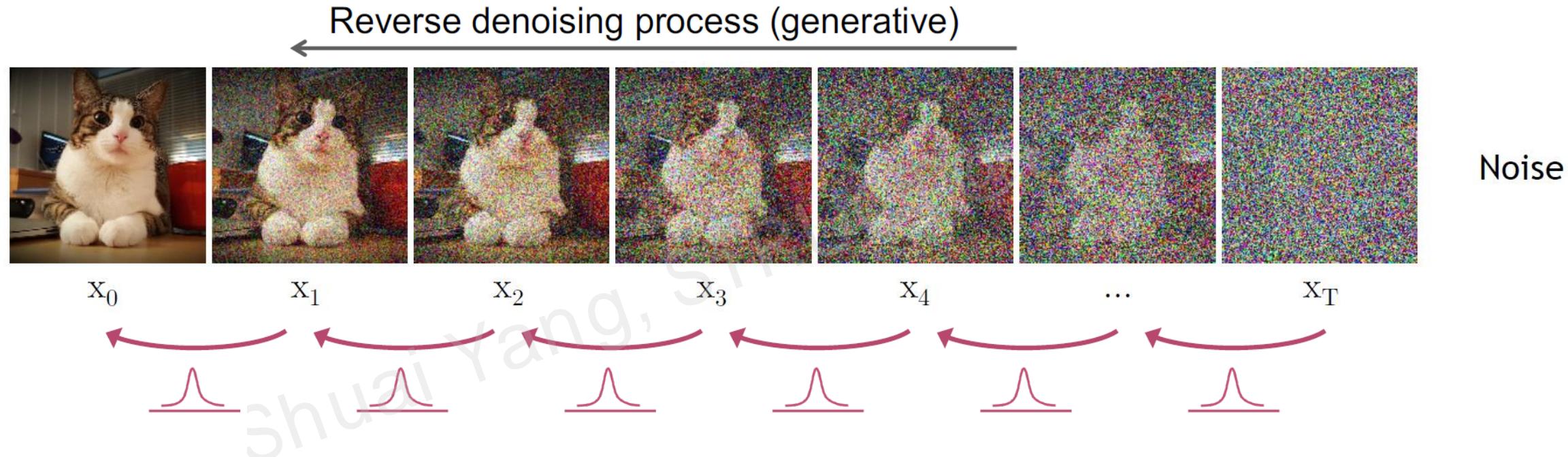




- Adding noise: Images “diffuse” in noise, like ink “diffuses” in water.
- Build a bridge between Gaussian noise and image manifold.

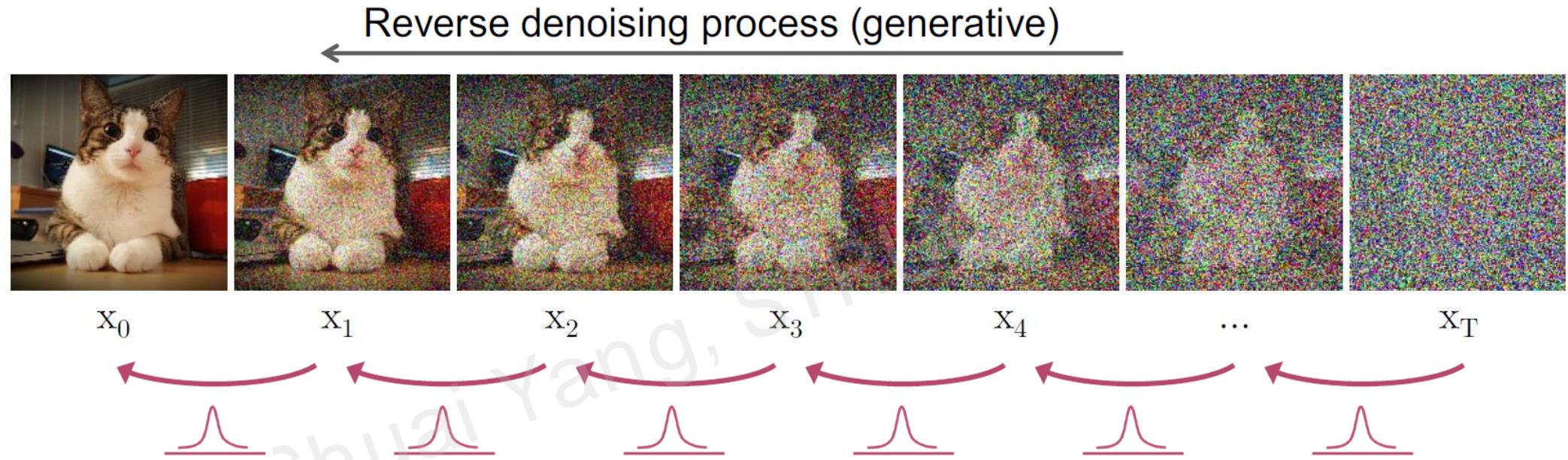


$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$



- Reverse process: gradually generate images from noise.

$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I}) \quad p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boxed{\mu_{\theta}(\mathbf{x}_t, t)}, \sigma_t^2 \mathbf{I})$$



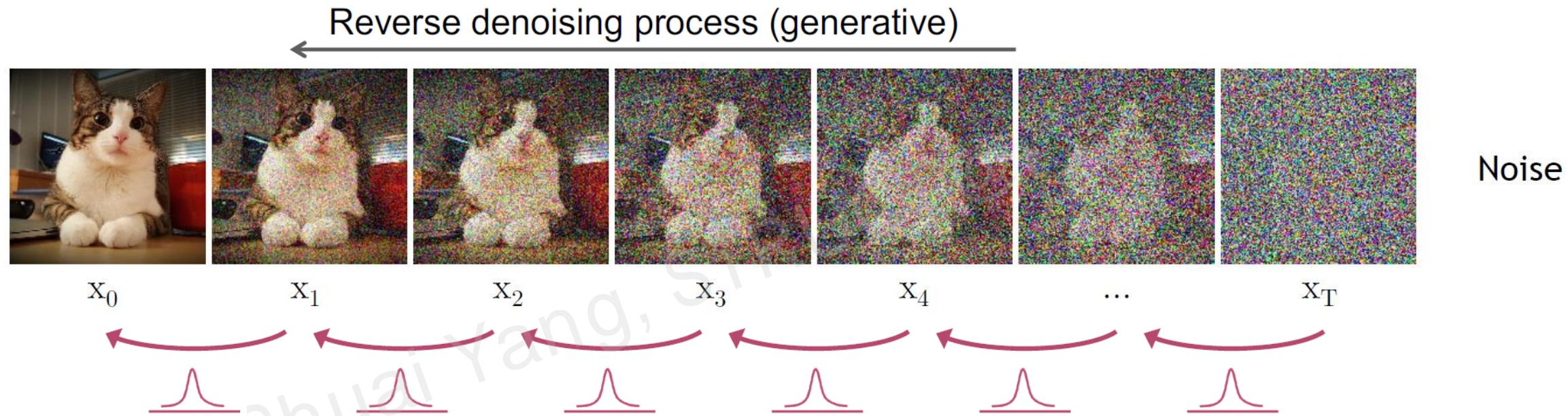
- Reverse process: given the original image \mathbf{x}_0 , we have

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I}),$$

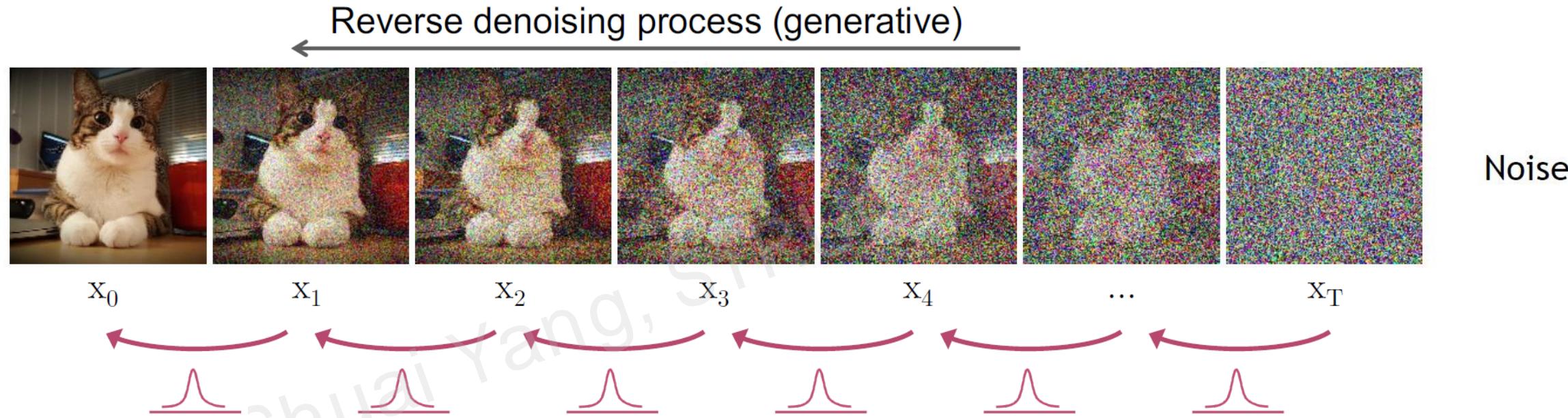
where $\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) := \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{1 - \beta_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t$ and $\tilde{\beta}_t := \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$

Preliminaries: Diffusion Models

Diffusion Models

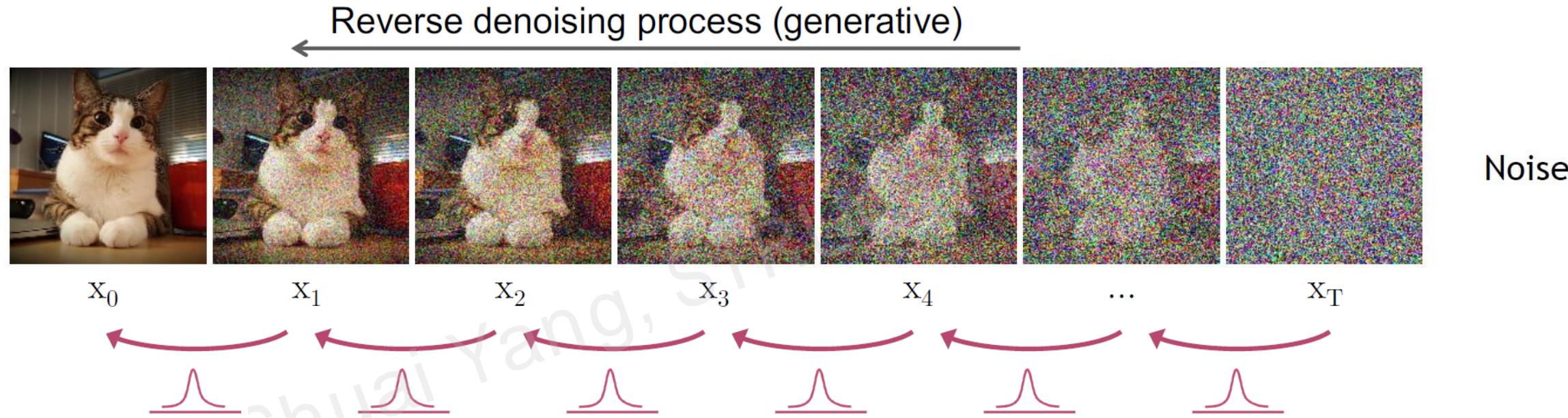


- Reverse process: replace the \mathbf{x}_0 as the predicted one:



- Reverse process: predicting the noise in \mathbf{x}_t

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right)$$



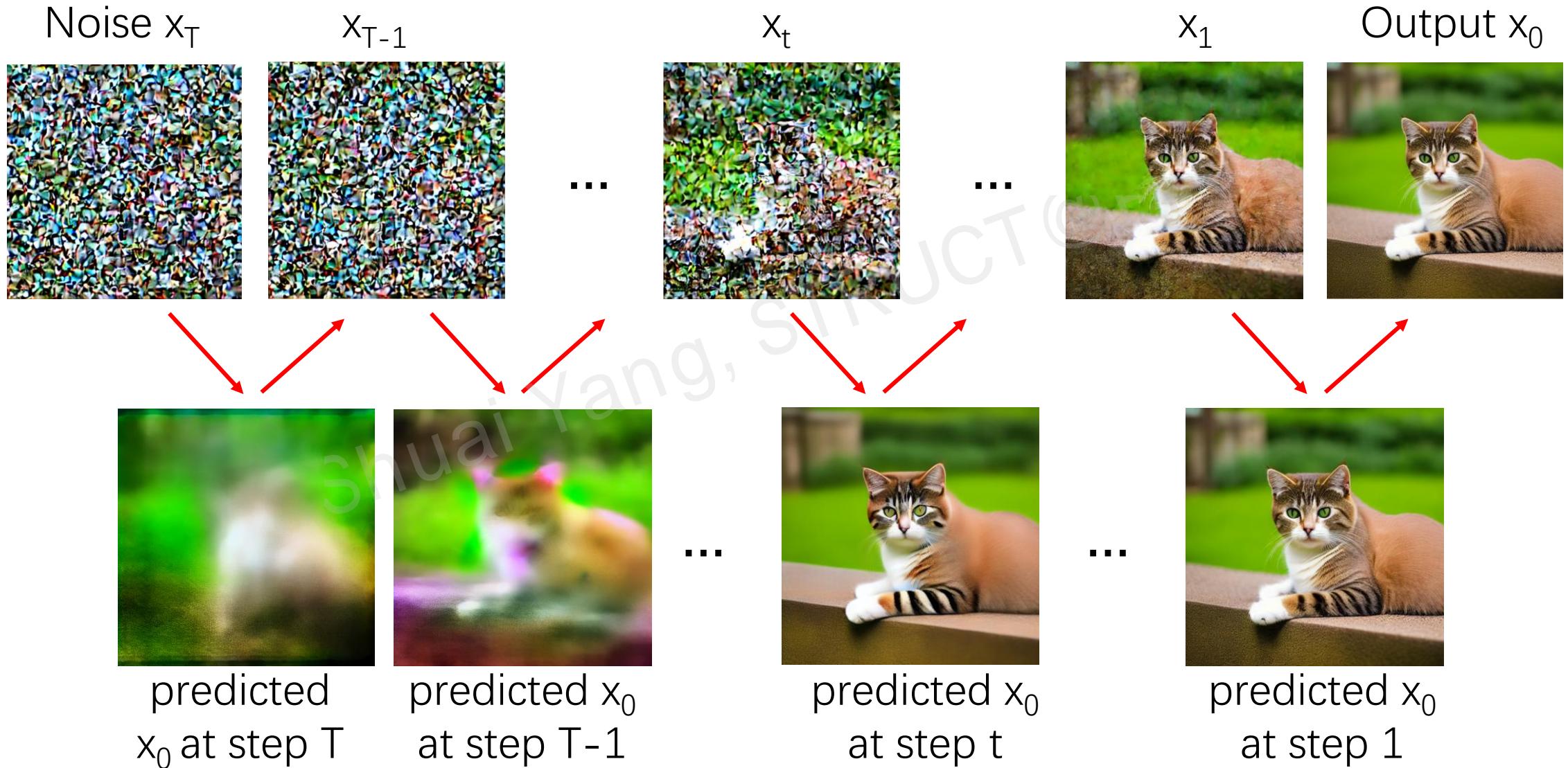
- Reverse process: training loss

$$\mathbb{E} [-\log p_\theta(\mathbf{x}_0)] =: L$$

$$L_{\text{simple}} = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t \sim \mathcal{U}(1, T)} [||\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)||^2]$$

Preliminaries: Diffusion Models

Diffusion Models



Disadvantage: Sampling speed

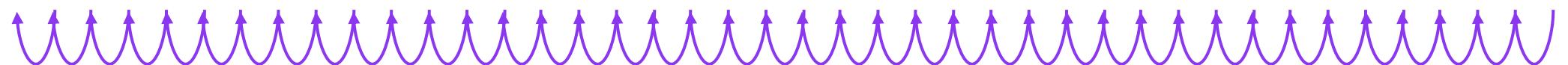
- The generation process is gradual



Data

Noise

Reverse denoising process (generative)



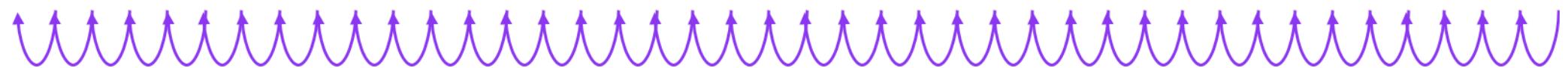
Divide the process into 1000 steps

Disadvantage: Sampling speed

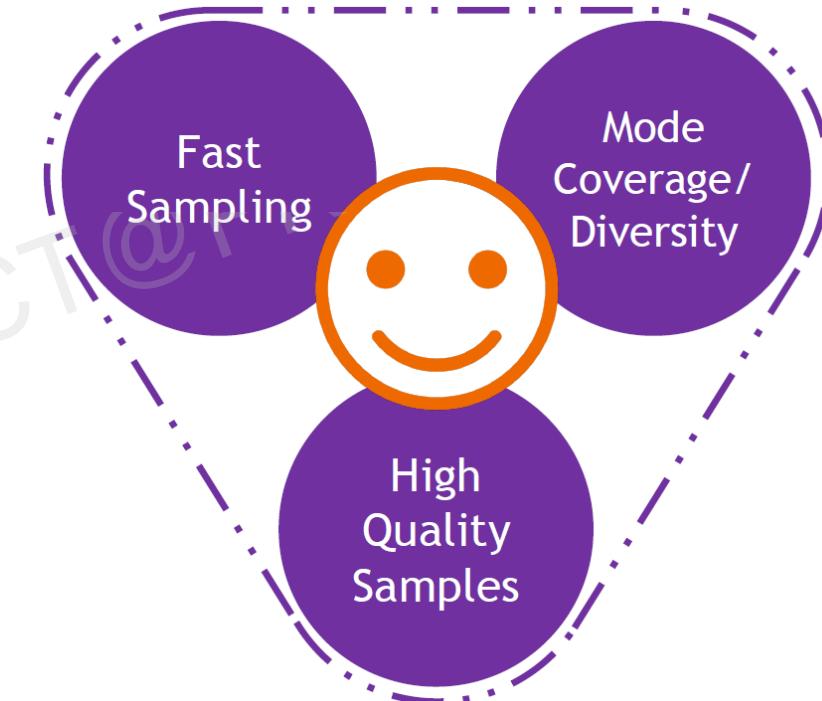
- The generation process is gradual

Accelerating:

- Determined sampling process
- Larger denoising step



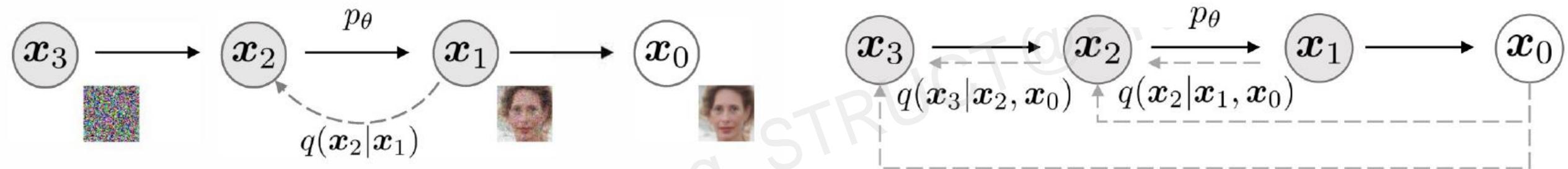
Divide the process into 50 or fewer steps



Preliminaries: Diffusion Models

DDIM

Non-Markovian process



$$q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}}\mathbf{x}_0 + \sqrt{1 - \alpha_{t-1} - \boxed{\sigma_t^2}} \cdot \frac{\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_0}{\sqrt{1 - \alpha_t}}, \boxed{\sigma_t^2}\mathbf{I}\right)$$

Determined when σ_t set to 0

Determined sampling process means –
The start noise can also be obtained by the image.

$$x_{t+1} - x_t = \sqrt{\bar{\alpha}_{t+1}} \left[\left(\sqrt{1/\bar{\alpha}_t} - \sqrt{1/\bar{\alpha}_{t+1}} \right) x_t + \left(\sqrt{1/\bar{\alpha}_{t+1} - 1} - \sqrt{1/\bar{\alpha}_t - 1} \right) \epsilon_\theta(x_t) \right]$$

The noise can be utilized in many ways (e.g., interpolating).



Problem: How to guide the diffusion models?

- A simple way: feed the noise-prediction network with the condition

$$\epsilon_{\theta}(x_t, t) \longrightarrow \epsilon_{\theta}(x_t, \boxed{y}, t)$$

Problem: How to guide the sampling process?

- Classifier guidance: employing an external classifier:

$$\mu \xrightarrow{\quad} \mu + s \sum \nabla_{x_t} \log p_\phi(y|x_t)$$

External classifier

- The classifier “pulls” the noisy image to the direction of the condition.

Problem: How to guide the sampling process?

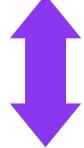
- Classifier guidance: employing an external classifier:

$$\mu \xrightarrow{\text{blue arrow}} \mu + s \sum \nabla_{x_t} \log p_\phi(y|x_t)$$

Hard to train!

Problem: How to guide the sampling process?

- Classifier-free guidance: implicit classifier

$$\tilde{\epsilon}_\theta(\mathbf{z}_\lambda, \mathbf{c}) = (1 + w)\epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c}) - w\epsilon_\theta(\mathbf{z}_\lambda)$$

$$\epsilon_\theta(\mathbf{z}_\lambda, \emptyset)$$

- Drop the condition randomly
- Apply the distance of conditioned and unconditioned

$$\tilde{\epsilon}_\theta(\mathbf{z}_\lambda, \mathbf{c}) = (1 + w)\epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c}) - w\epsilon_\theta(\mathbf{z}_\lambda)$$

a cute cat in the garden



W=-1

W=0

W=1

W=3

W=6.5

W=9

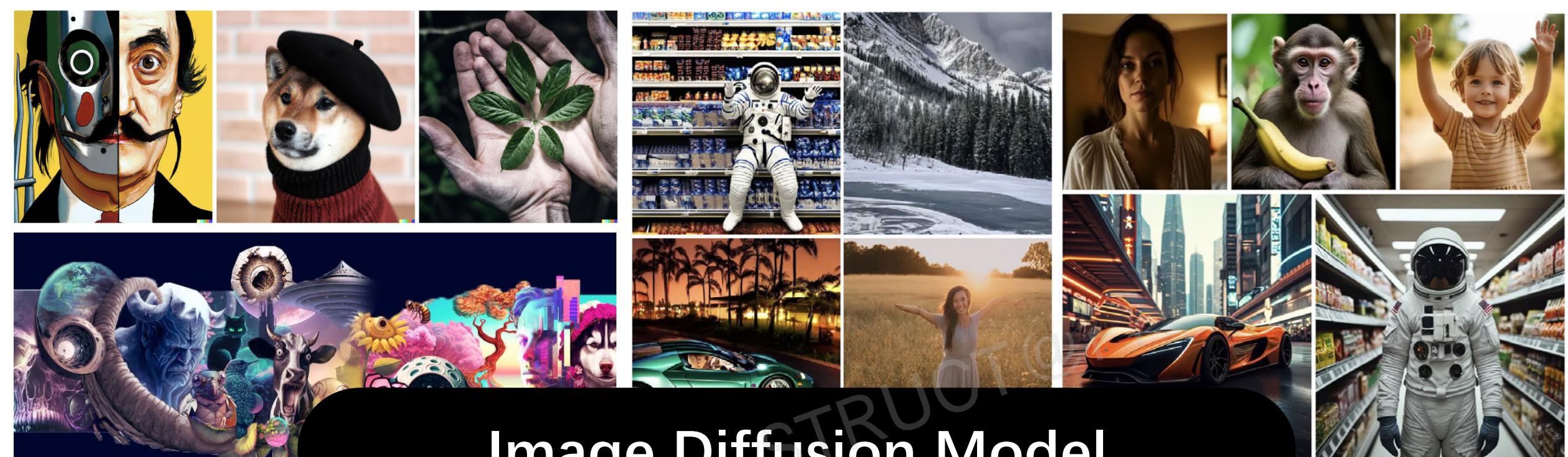


Image Diffusion Model

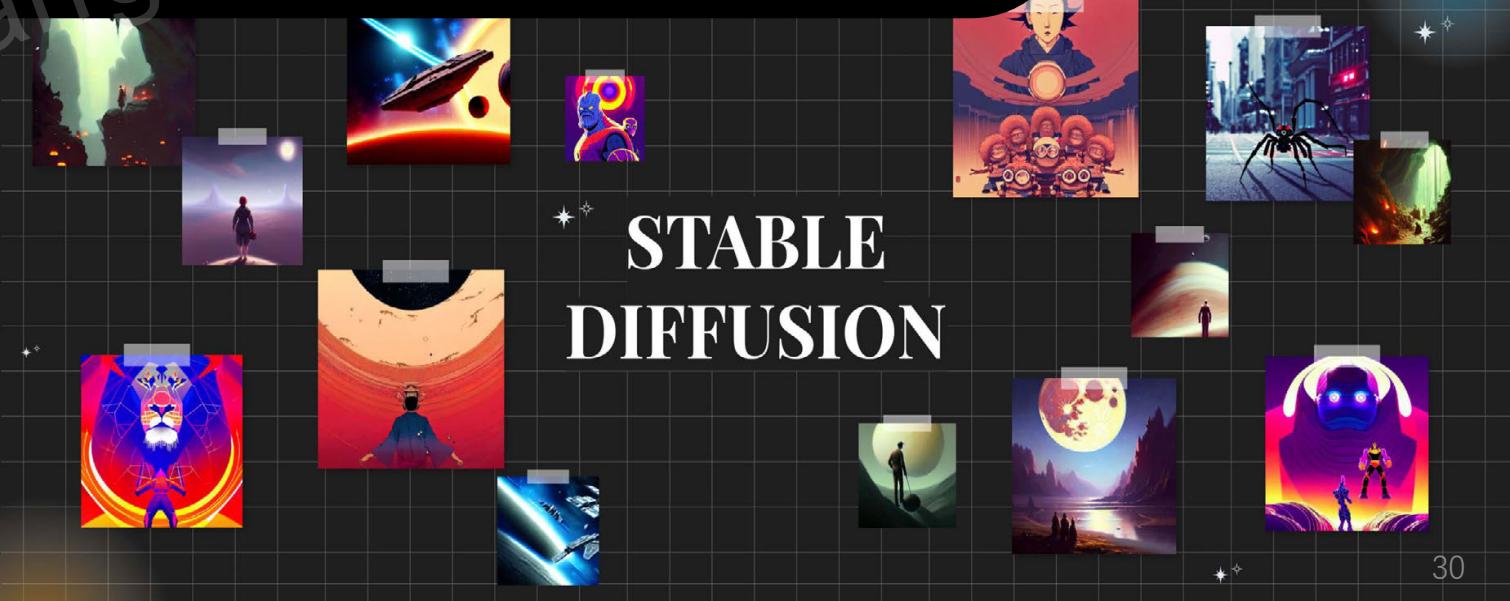


Image diffusion model

The World's Smartest Artificial Intelligence
Just Made Its First Magazine Cover

wide-angle shot from below of a female astronaut with an athletic feminine body walking with swagger toward camera on Mars in an infinite universe, synthwave digital art

digital artist — Karen X. Cheng

DALL·E 2

COSMOPOLITAN

the A.I. issue

Meet the
World's
First
Artificially
Intelligent
Magazine
Cover



Image diffusion model

DALL·E 2

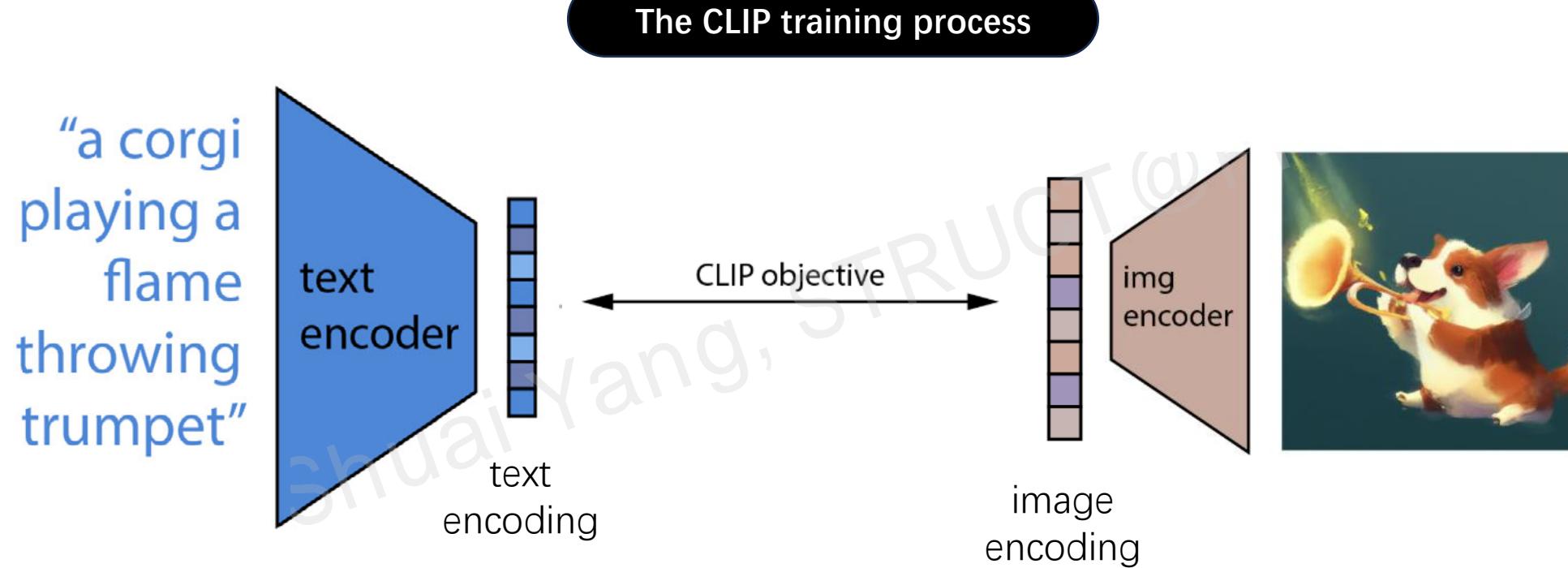


Image diffusion model

DALL·E 2

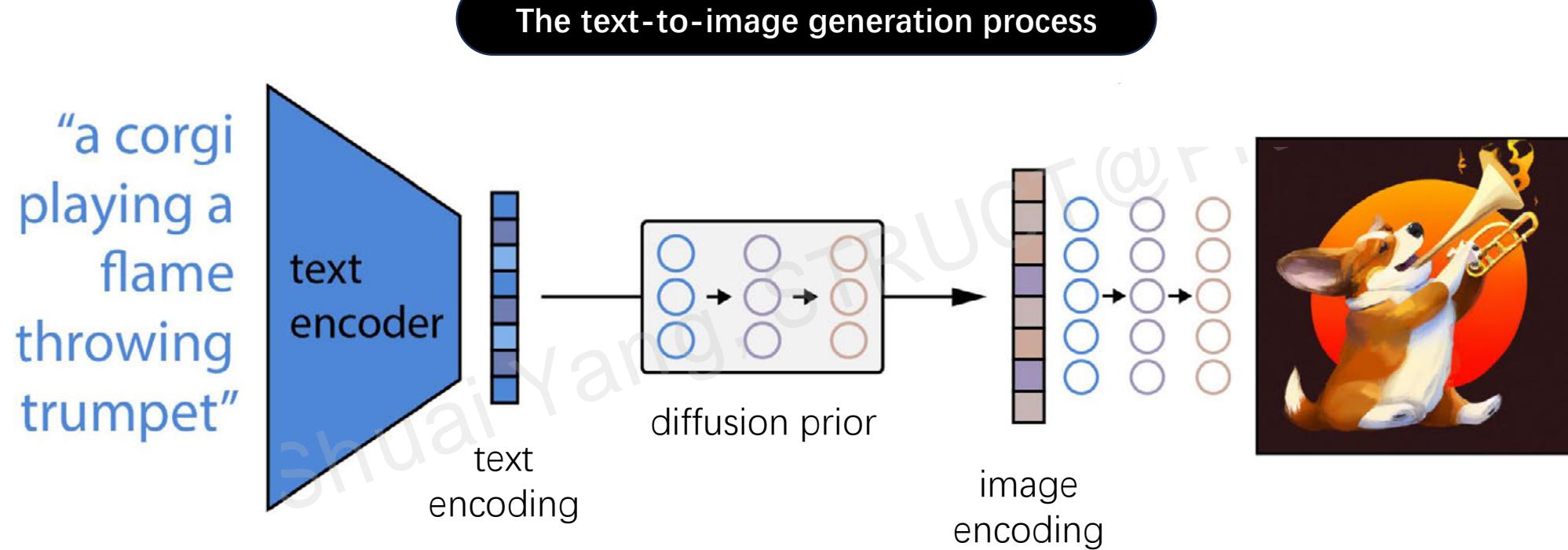
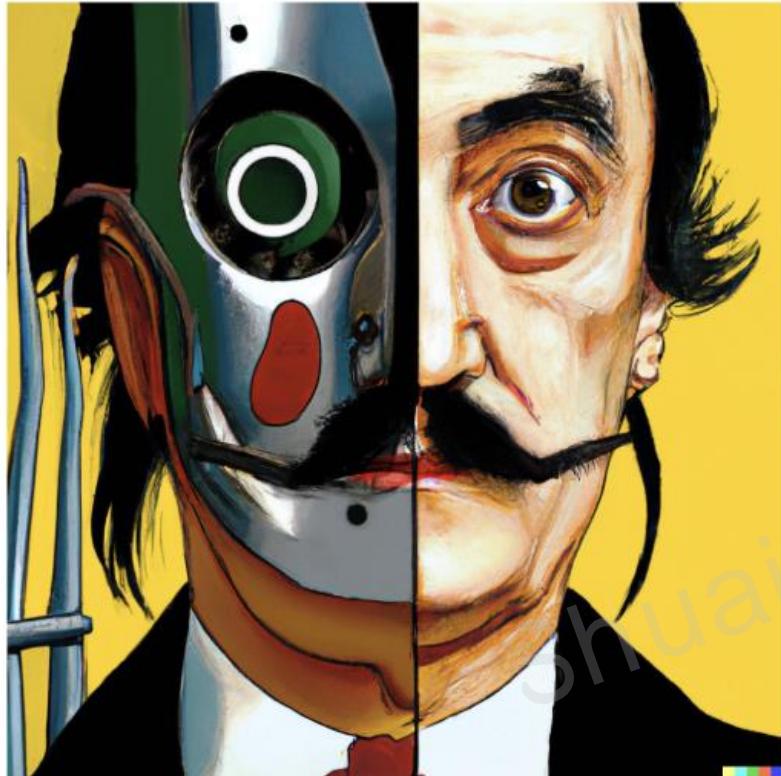


Image diffusion model

DALL·E 2



vibrant portrait painting of Salvador Dalí with a robotic half face



a shiba inu wearing a beret and black turtleneck



a close up of a handpalm with leaves growing from it

Image diffusion model

DALL·E 2

Encoding with CLIP and then decoding with a diffusion model

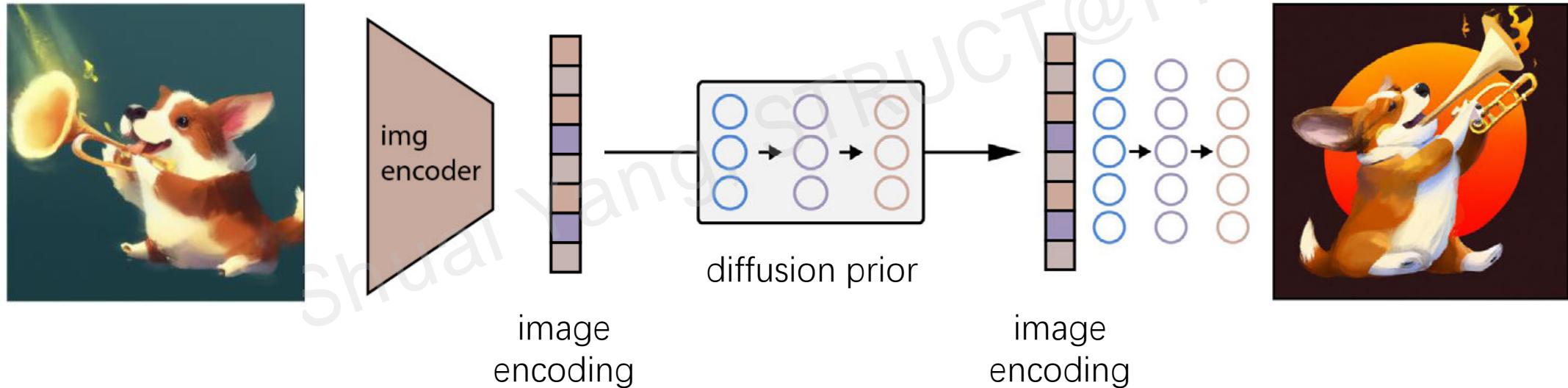


Image diffusion model

DALL·E 2

Encoding with CLIP and then decoding with a diffusion model



Image diffusion model

DALL·E 2

Interpolating CLIP image embeddings and then decoding with a diffusion model



Image diffusion model

DALL·E 2

Interpolating CLIP image embeddings and then decoding with a diffusion model



To train diffusion models on **limited computational resources** while retaining quality and flexibility

- Not directly in the pixel space
- But using **the potential space**

Image diffusion model

Latent Diffusion Models

To train diffusion models on **limited computational resources** while retaining quality and flexibility

- Not directly in the pixel space
- But using **the potential space**

Perceptual and semantic compression

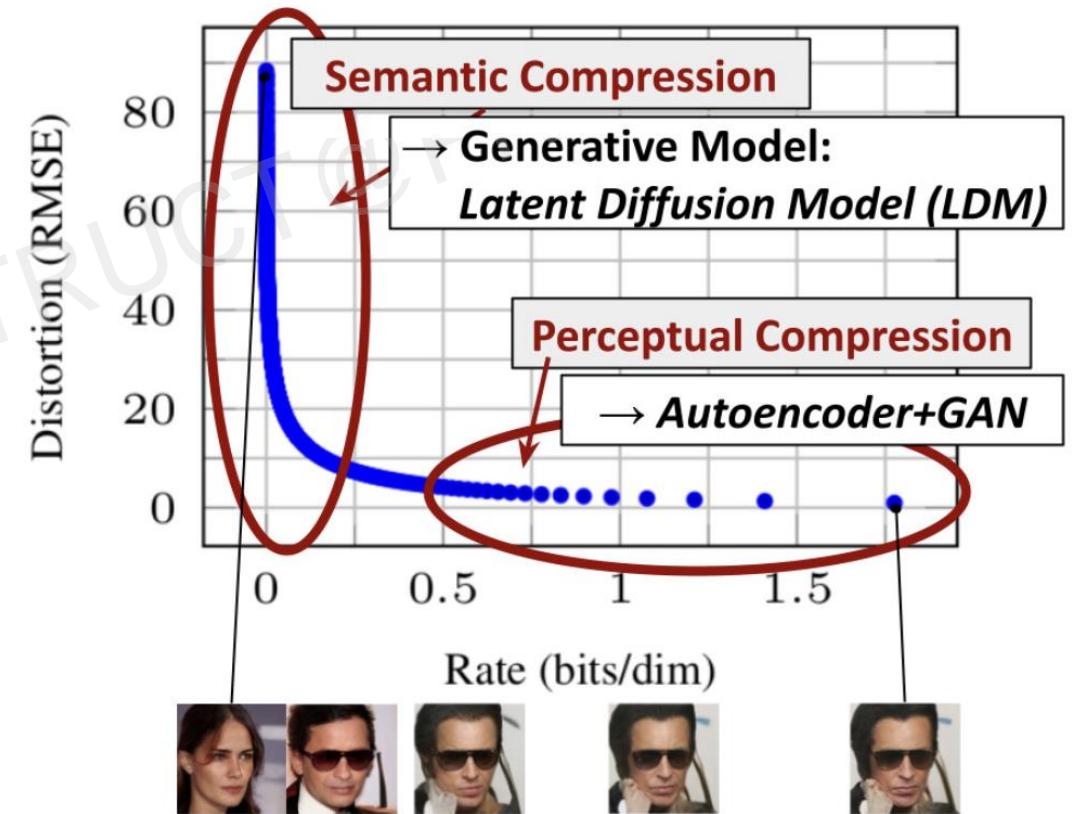


Image diffusion model

Latent Diffusion Models

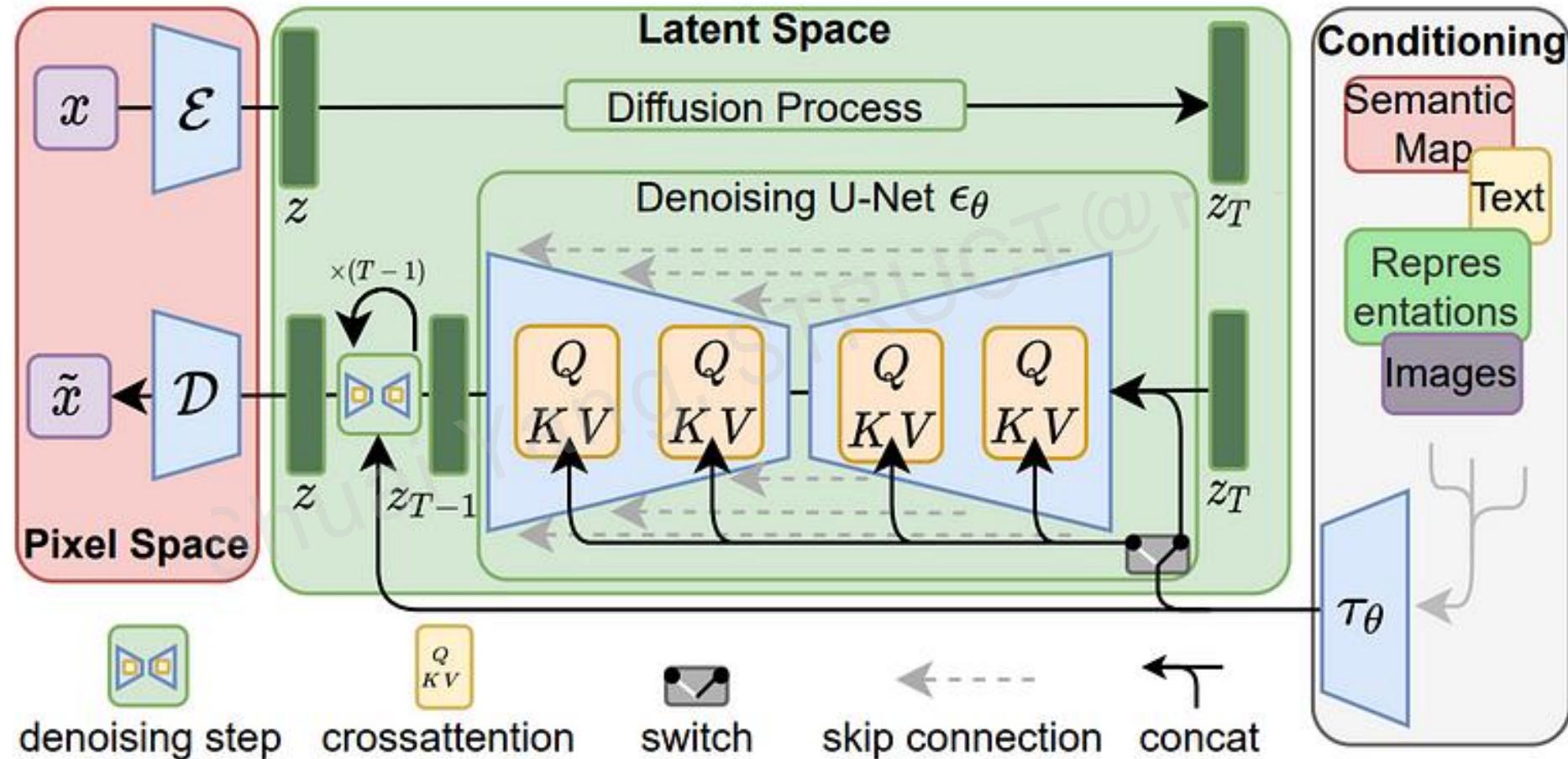


Image diffusion model

Latent Diffusion Models

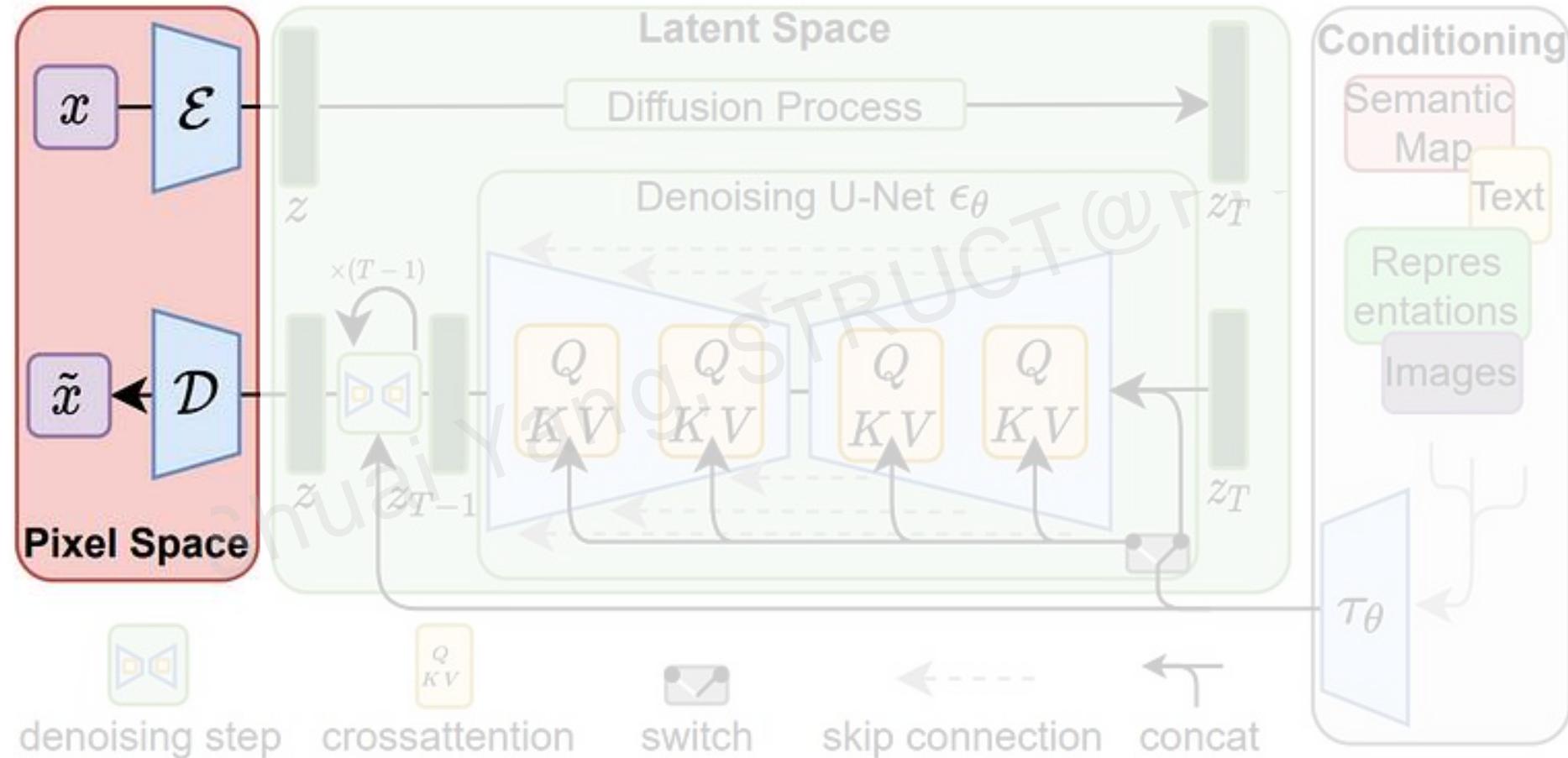


Image diffusion model

Latent Diffusion Models

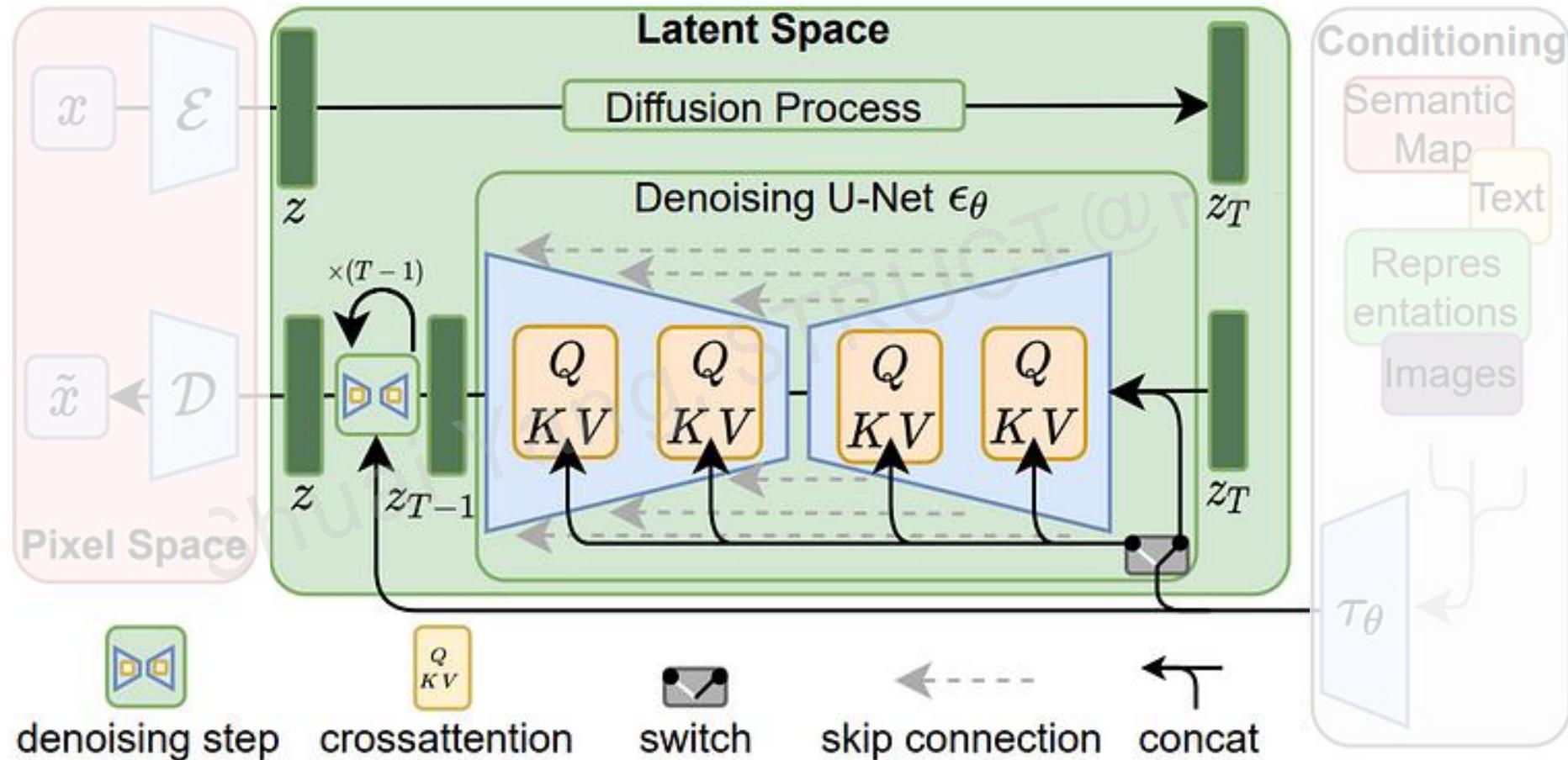
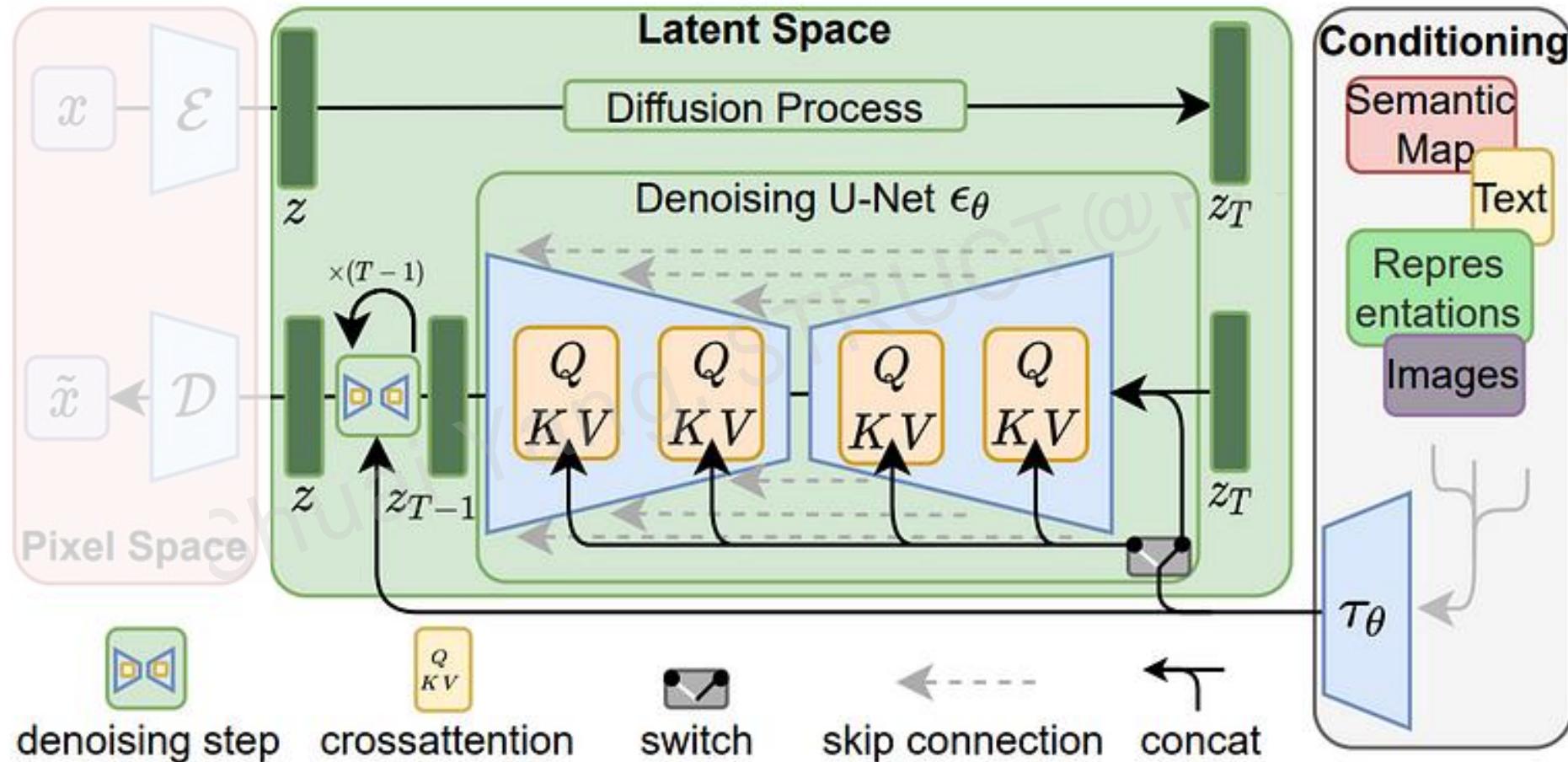


Image diffusion model

Latent Diffusion Models



STABLE DIFFUSION



Image diffusion model

Stable Diffusion series

June 2022

August 2022

October 2022

November 2022

December 2022



SD 1.1

SD 1.2

SD 1.3

SD 1.4

SD 1.5

SD 2.0

SD 2.1



SD 1.x Models: resolution of 512x512 pixels and use a ViT-L/14 CLIP model for text conditioning

Image diffusion model

Stable Diffusion series

June 2022

August 2022

October 2022

November 2022

December 2022



SD 1.1

SD 1.2

SD 1.3

SD 1.4

SD 1.5

SD 2.0

SD 2.1



SD 2.x Models: resolution of 768x768 pixels and use a different text encoder: ViT-H/14 OpenCLIP, allowing the prompts to be more expressive.

Image diffusion model

Stable Diffusion series

December 2022

July 2023

November 2023

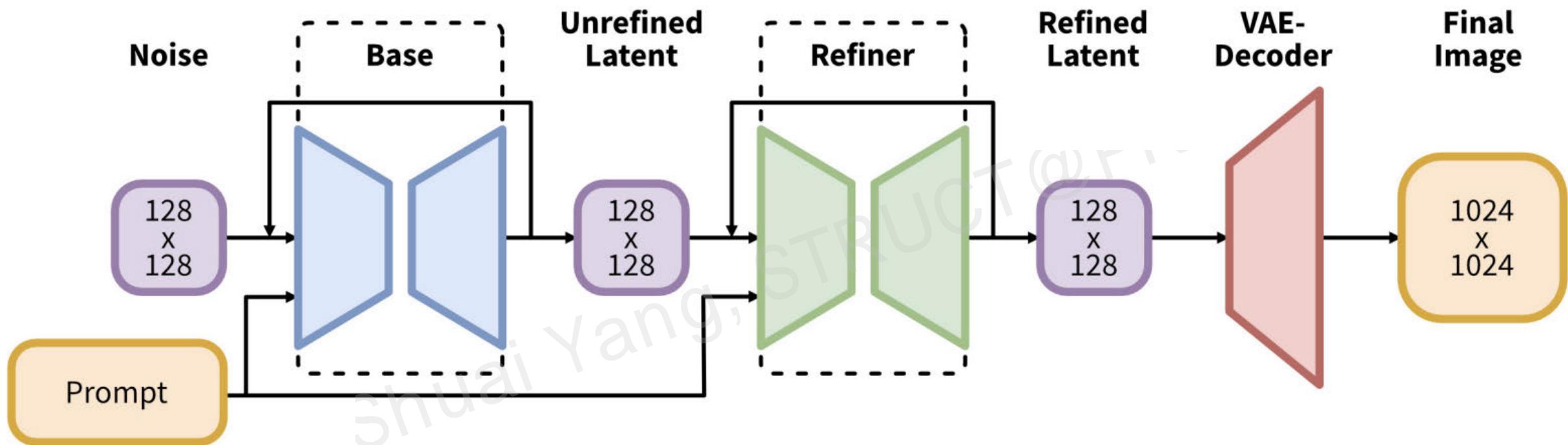
February 2024

March 2024



Image diffusion model

Stable Diffusion XL



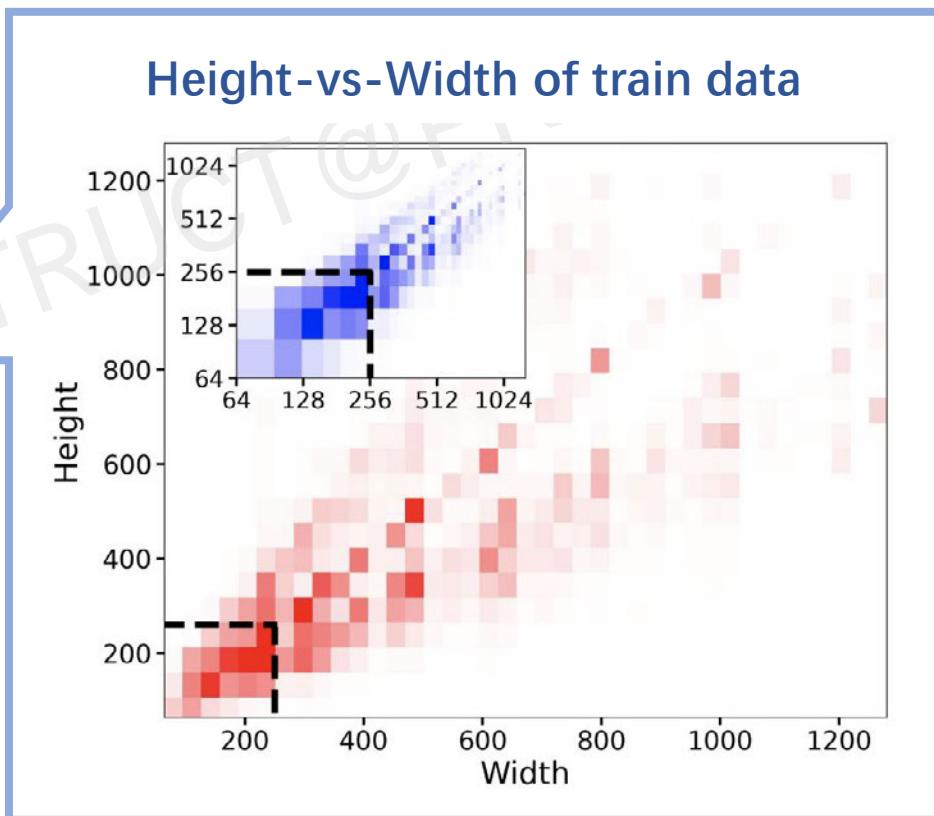
Number of parameters compared with older SD models:

- SD XL 2.6B v.s. SD 1.4/1.5 860M v.s. SD 2.0/2.1 865M

The shortcoming of the LDM: training a model requires a minimal image size

Common solutions:

- Discard all training images below a certain minimal resolution.
→ Discard a lot of training data
- Upscale images that are too small.
→ Upscaling artifacts



Micro-Conditioning: 1. Size conditioning

$\mathbf{c}_{\text{size}} = (64, 64)$



$\mathbf{c}_{\text{size}} = (128, 128),$



$\mathbf{c}_{\text{size}} = (256, 256),$



$\mathbf{c}_{\text{size}} = (512, 512),$



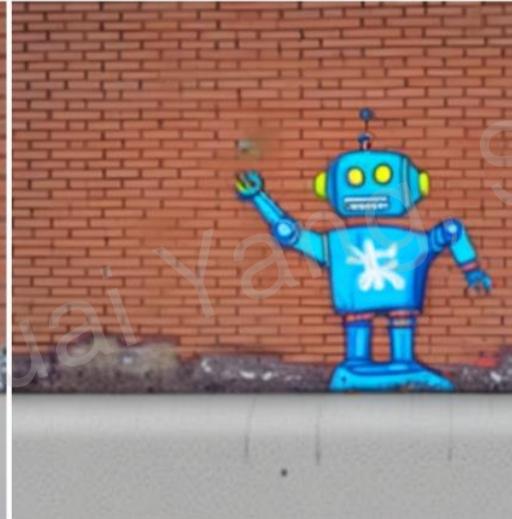
“Panda mad scientist mixing sparkling chemicals, artstation.”

Micro-Conditioning: 1. Size conditioning

$\mathbf{c}_{\text{size}} = (64, 64)$



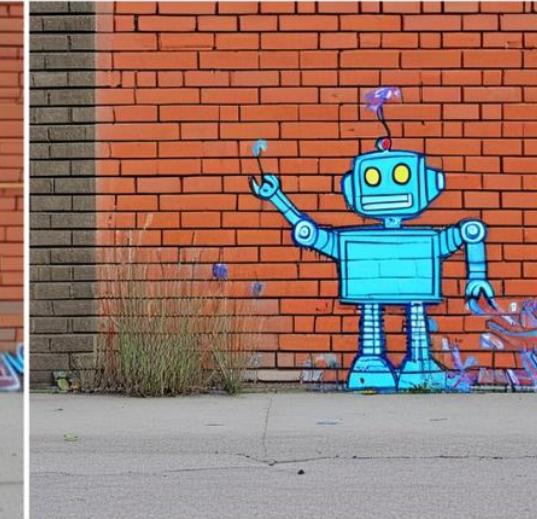
$\mathbf{c}_{\text{size}} = (128, 128),$



$\mathbf{c}_{\text{size}} = (256, 236),$



$\mathbf{c}_{\text{size}} = (512, 512),$



'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.'

Failure cases of previous SD models

"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

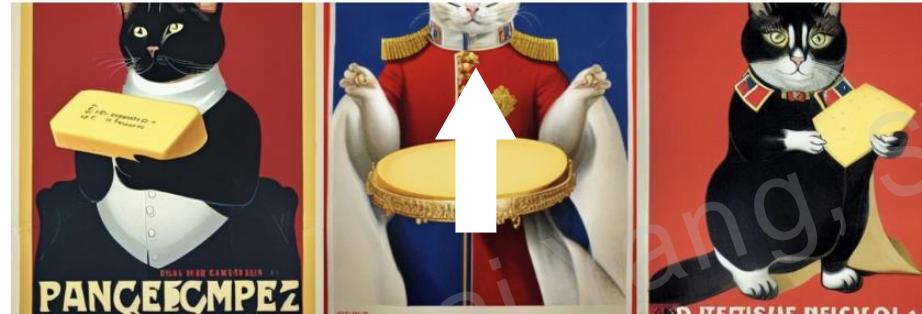


SD 2.1



Failure cases of previous SD models

"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



Micro-Conditioning: 2. Cropping parameters ($c_{\text{top}}, c_{\text{left}}$)

"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



SDXL



Micro-Conditioning: 2. Cropping parameters ($c_{\text{top}}, c_{\text{left}}$)

$\mathbf{ccrop} = (0, 0)$



$\mathbf{ccrop} = (0, 256),$



$\mathbf{ccrop} = (256, 0),$



$\mathbf{ccrop} = (512, 512),$



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

Multi-Aspect Training

- Partition the data into buckets of different aspect ratios
- Keep the pixel count close to 1024^2 pixels
- Bucket size as a conditioning

Height	Width	Aspect Ratio	Height	Width	Aspect Ratio
512	2048	0.25	1024	1024	1.0
512	1984	0.26	1024	960	1.07
512	1920	0.27	1088	960	1.13
512	1856	0.28	1088	896	1.21
576	1792	0.32	1152	896	1.29
576	1728	0.33	1152	832	1.38
576	1664	0.35	1216	832	1.46
640	1600	0.4	1280	768	1.67
640	1536	0.42	1344	768	1.75
704	1472	0.48	1408	704	2.0
704	1408	0.5	1472	704	2.09
704	1344	0.52	1536	640	2.4
768	1344	0.57	1600	640	2.5
768	1280	0.6	1664	576	2.89
832	1216	0.68	1728	576	3.0
832	1152	0.72	1792	576	3.11
896	1152	0.78	1856	512	3.62
896	1088	0.82	1920	512	3.75
960	1088	0.88	1984	512	3.88
960	1024	0.94	2048	512	4.0

Image diffusion model

Stable Diffusion XL



Image diffusion model

Stable Diffusion XL

*'Monster Baba yaga house with in a forest,
dark horror style, black and white.'*

*'A young badger delicately sniffing a
yellow rose, richly textured oil painting.'*

SD 1.5



SD 2.1



SDXL



Image diffusion model

Stable Diffusion XL

'Cute adorable little goat, unreal engine, cozy interior lighting, art station, detailed' digital painting, cinematic, octane rendering.'

'A portrait photo of a kangaroo wearing an orange hoodie and blue sunglasses standing on the grass in front of the Sydney Opera House holding a sign on the chest that says "SDXL"!'

SD 1.5



SD 2.1



SDXL



Image diffusion model

Stable Diffusion series

December 2022

July 2023

November 2023

February 2024

March 2024



SD 2.1

SD XL 1.0

SD XL Turbo

SD Cascade

SD 3.0

Adversarial
Diffusion
Distillation

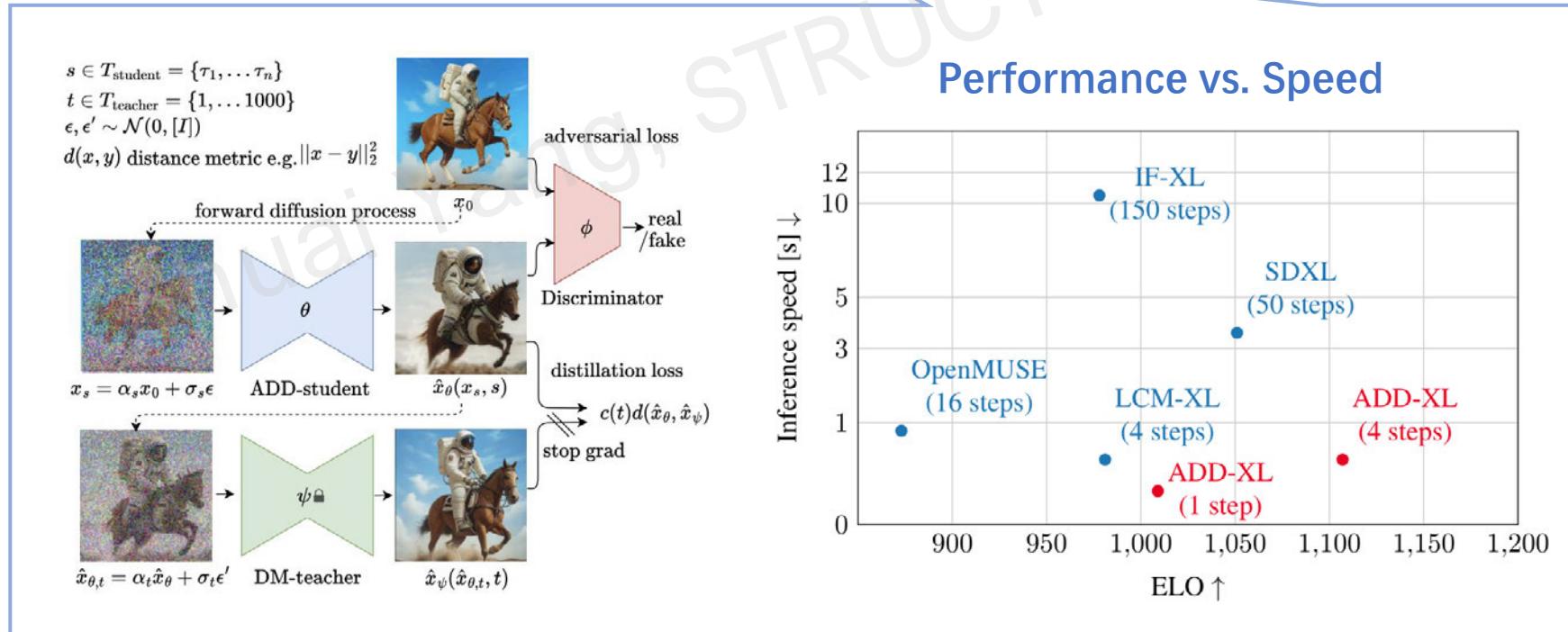


Image diffusion model

Stable Diffusion series

December 2022

July 2023

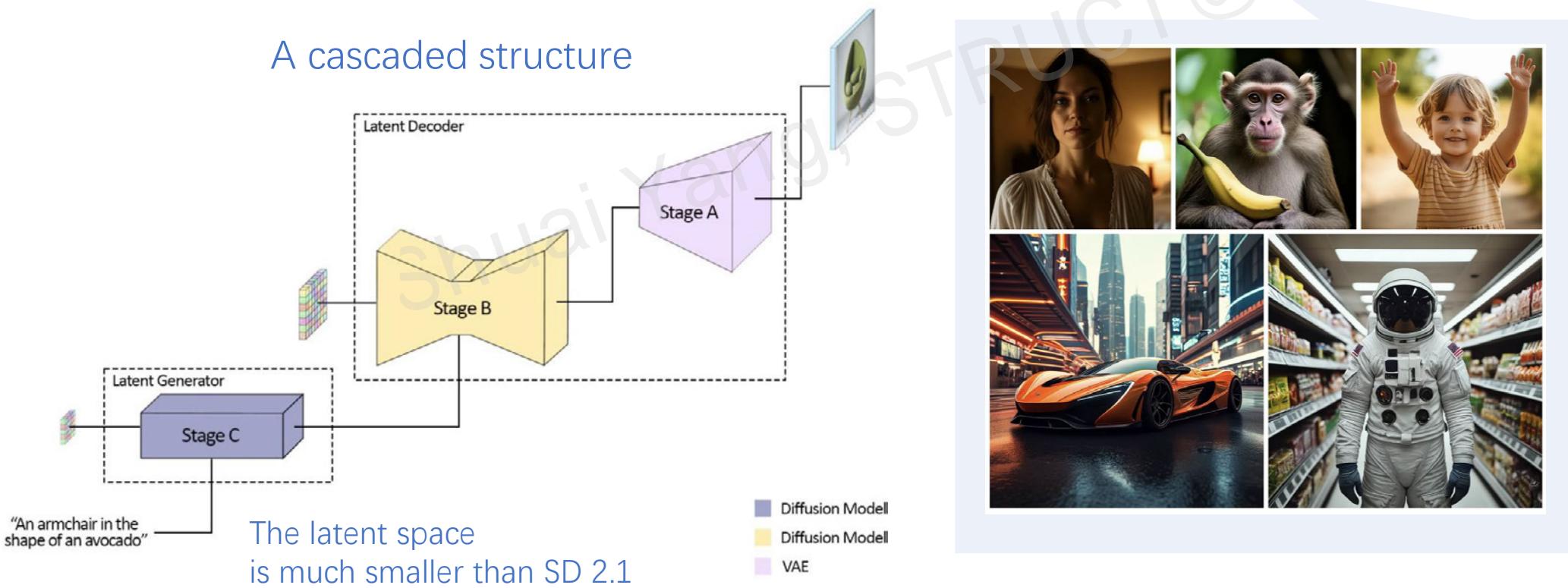
November 2023

February 2024

March 2024



A cascaded structure

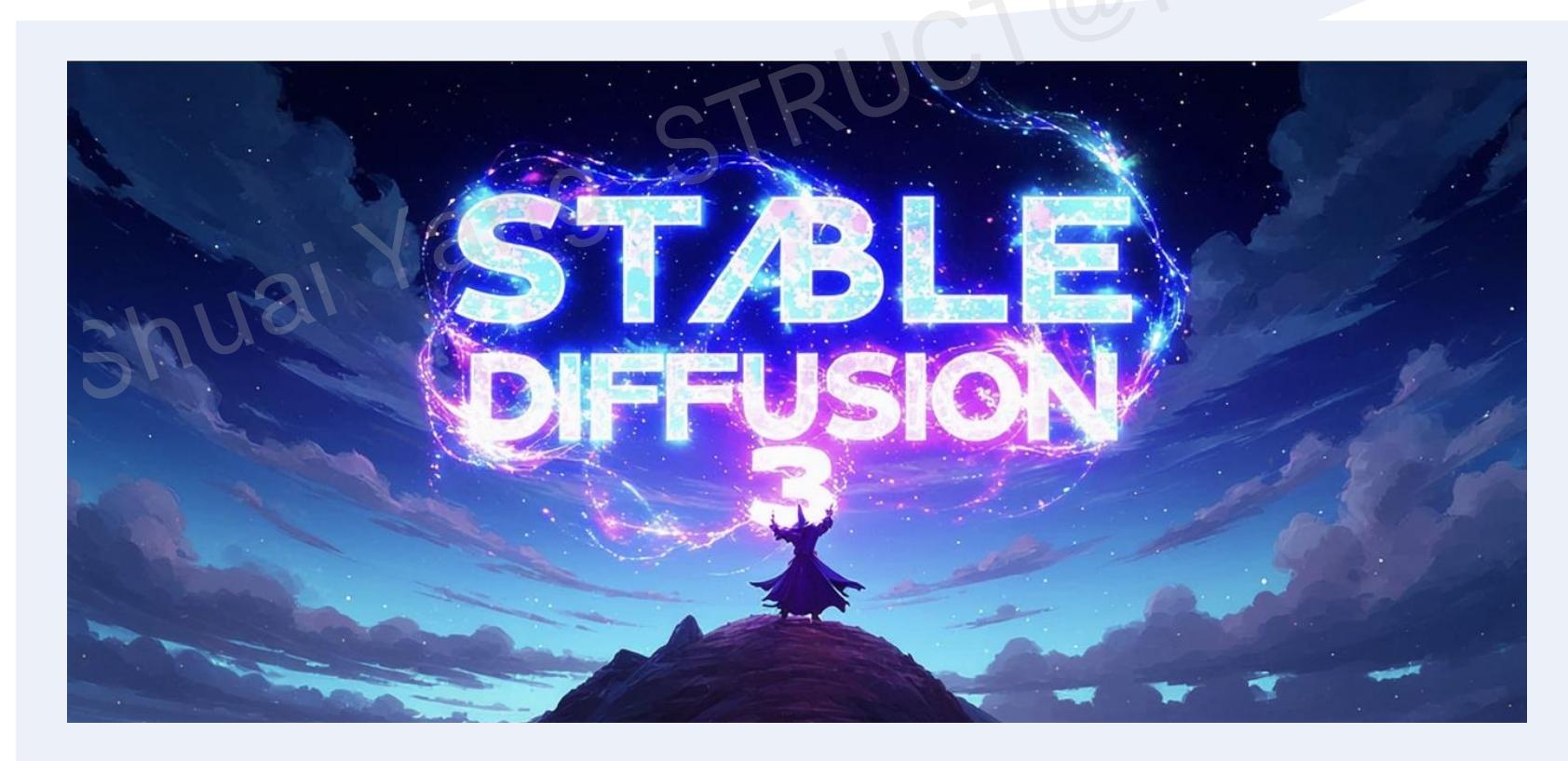
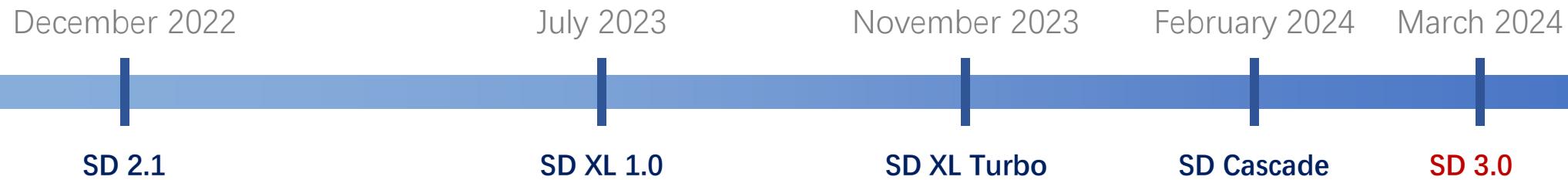


The latent space
is much smaller than SD 2.1

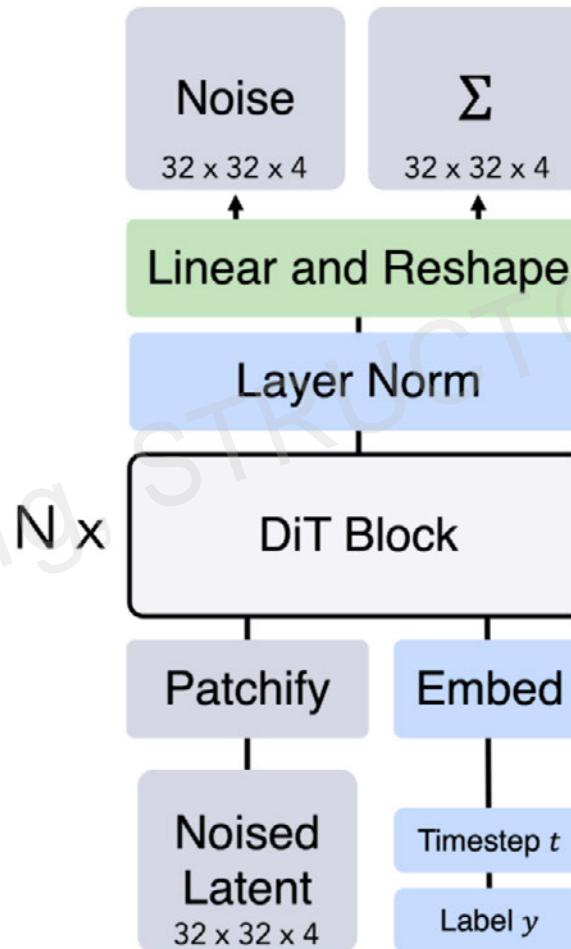


Image diffusion model

Stable Diffusion series



Latent Diffusion + Transformer



Patchify: convert the spatial input
into a sequence of T tokens

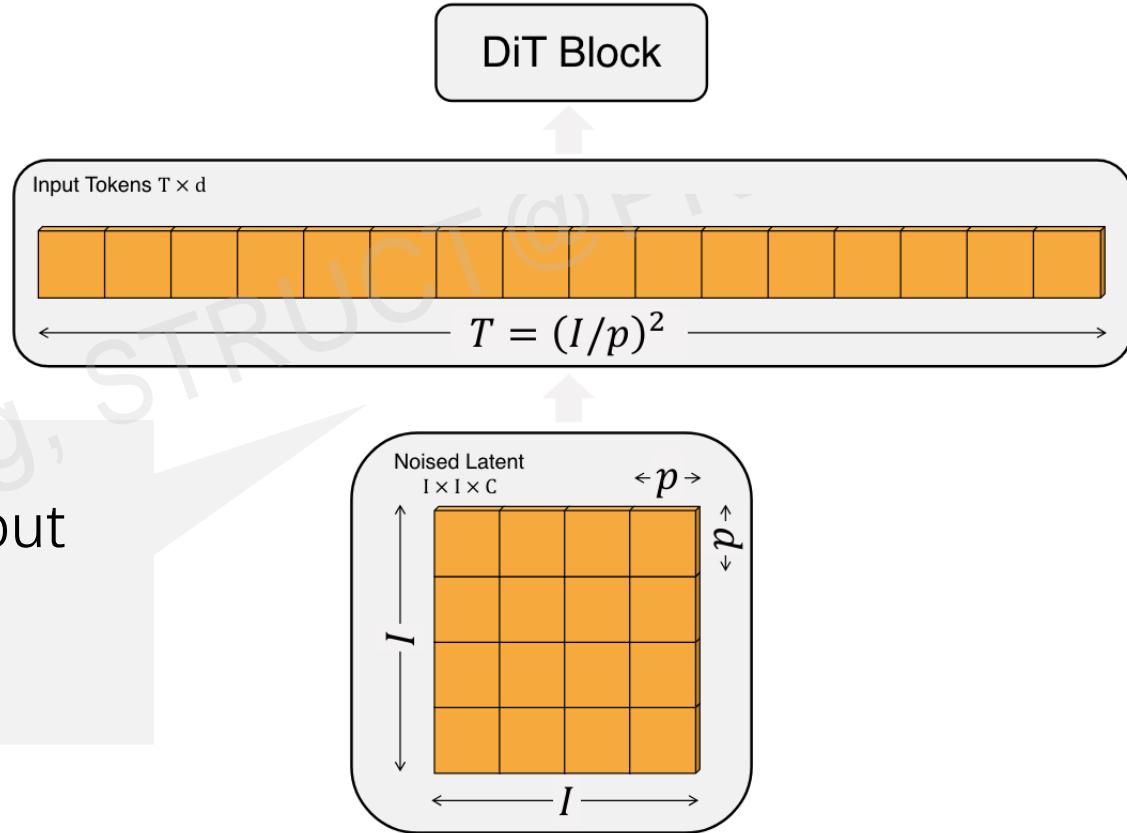
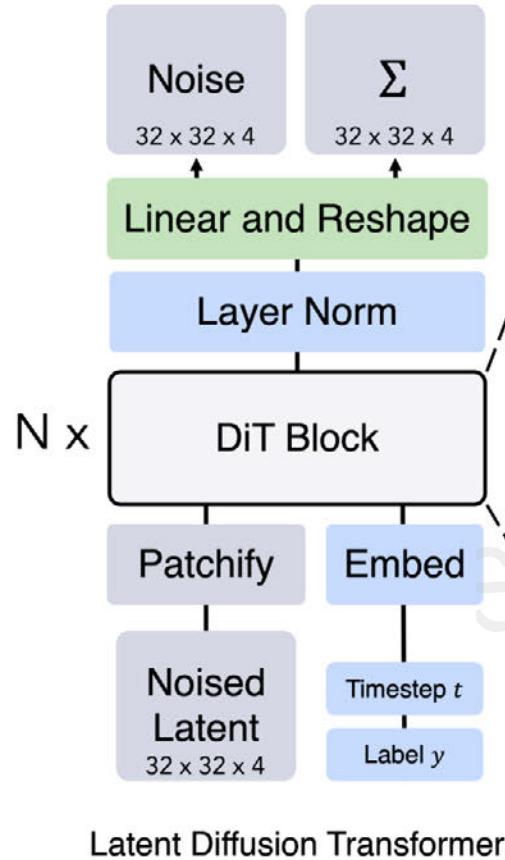
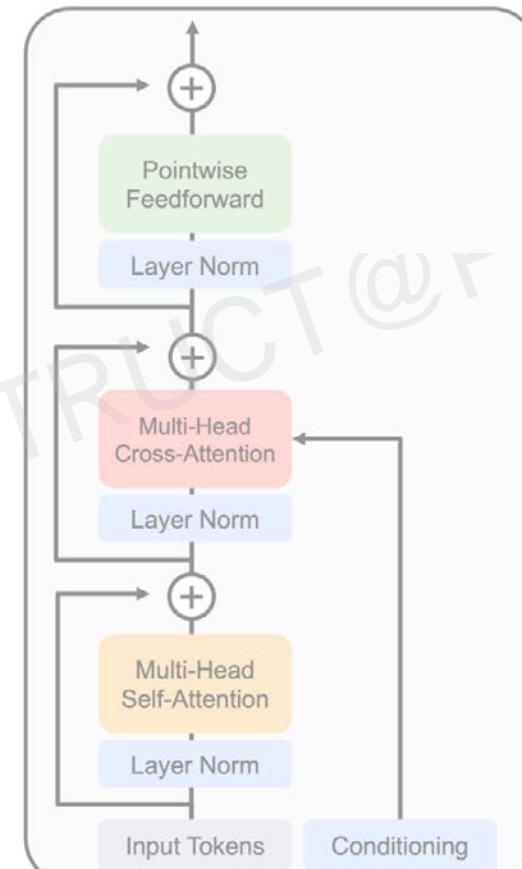


Image diffusion model

Diffusion Transformers



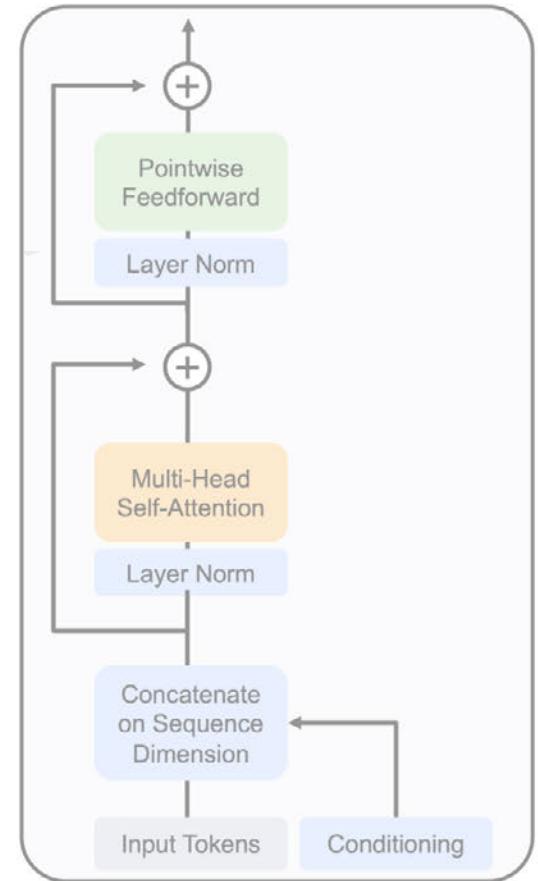
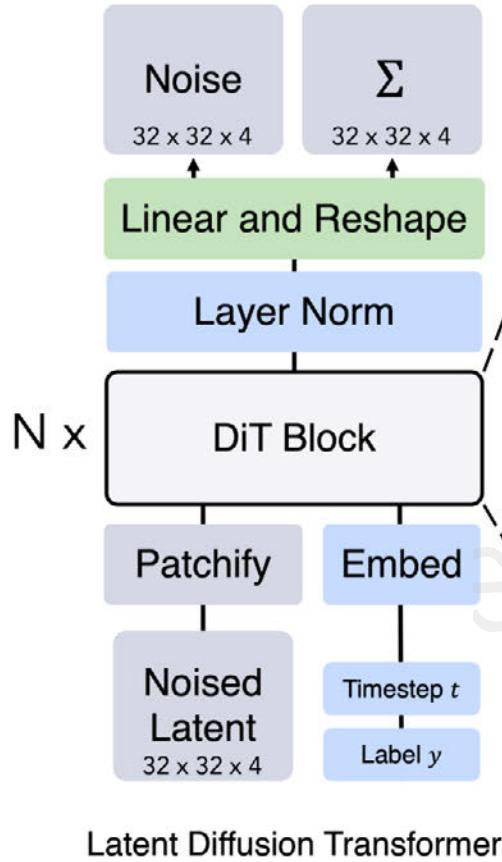
Latent Diffusion Transformer



DiT Block with Cross-Attention

Image diffusion model

Diffusion Transformers



DiT Block with In-Context Conditioning

Image diffusion model

Diffusion Transformers

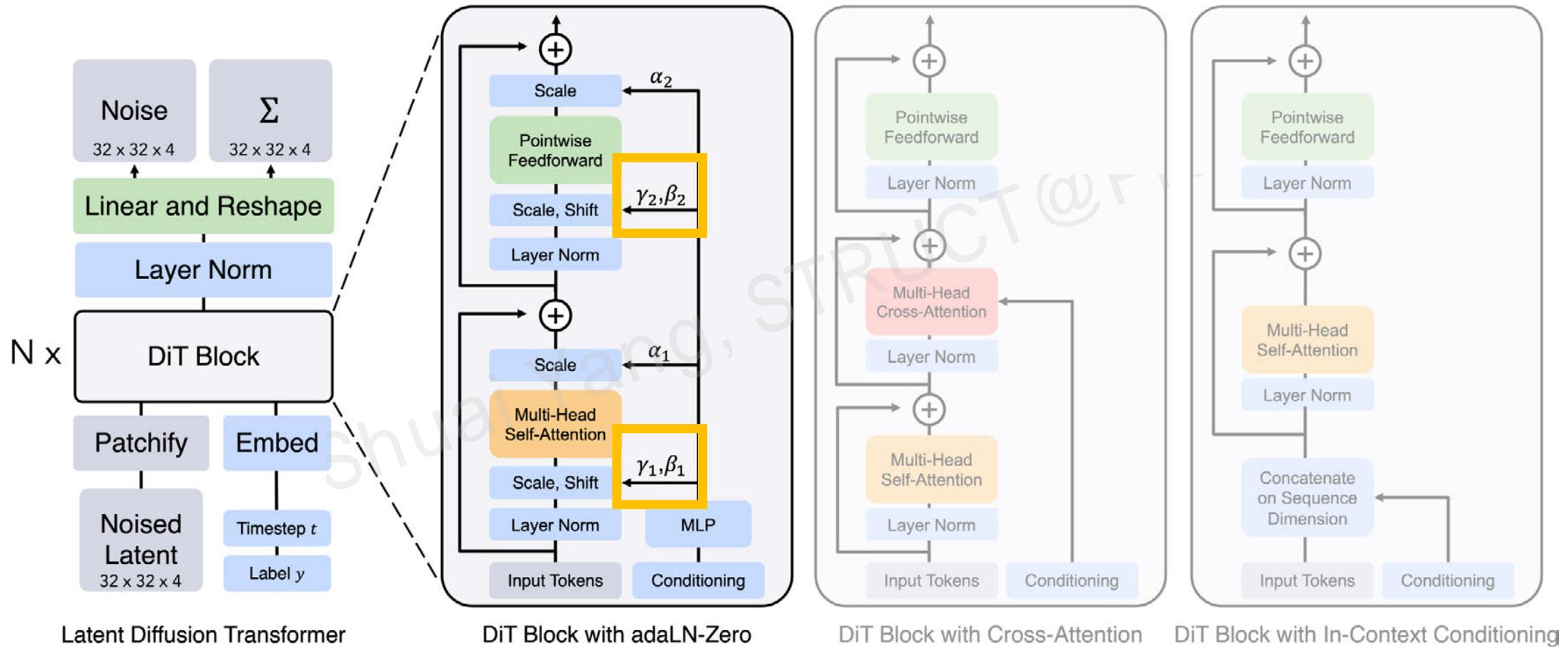


Image diffusion model

Diffusion Transformers

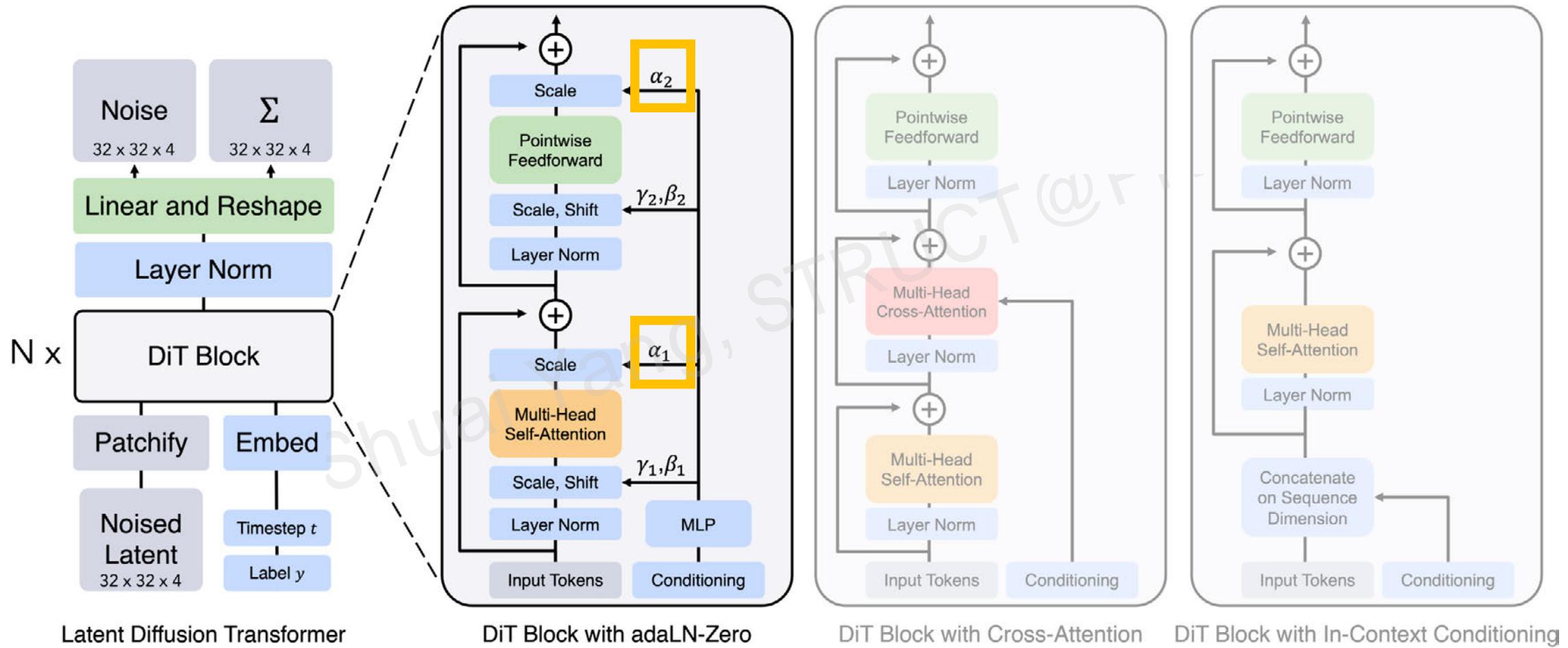


Image diffusion model

Stable Diffusion 3

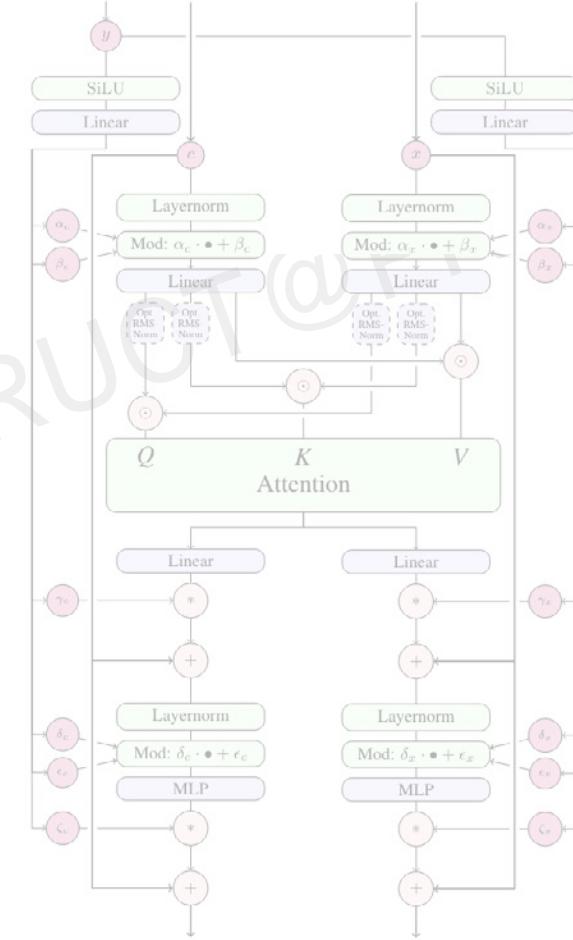
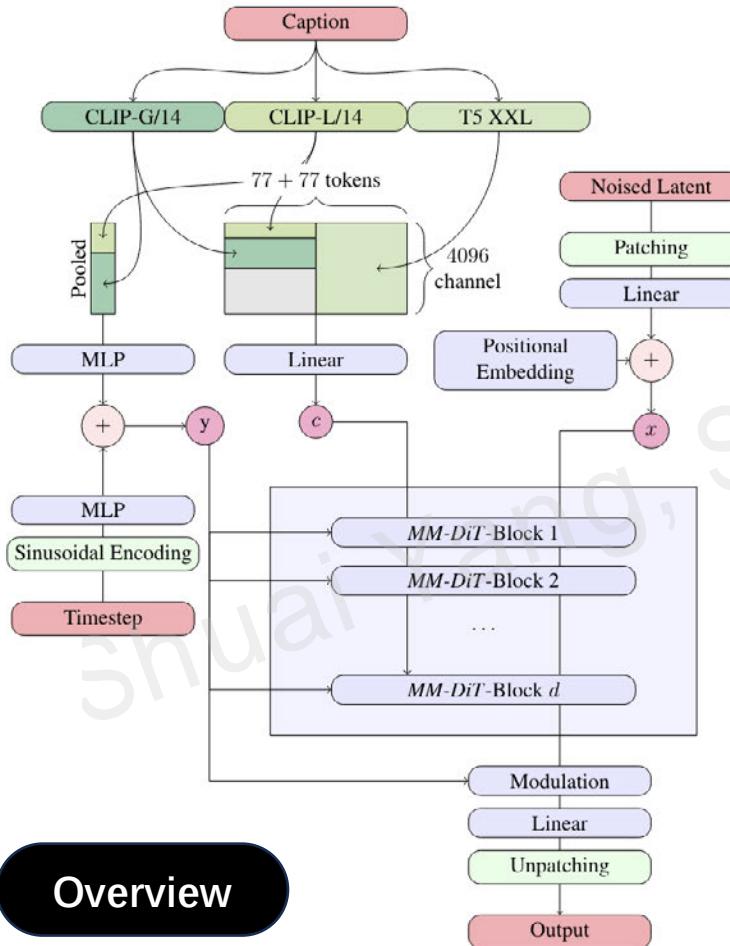
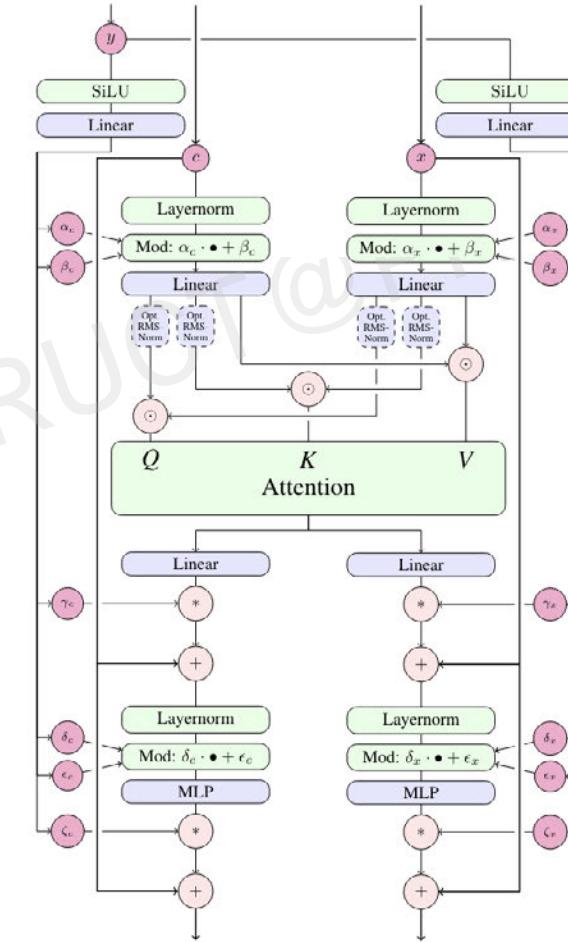
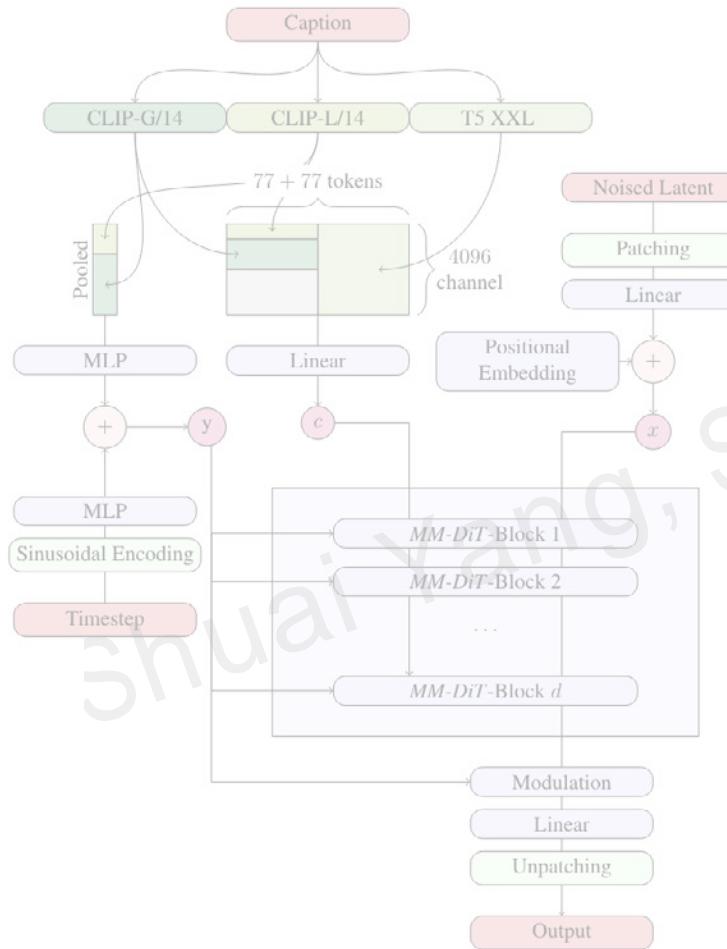


Image diffusion model

Stable Diffusion 3



Two separate sets
of weights for
image and text

Image diffusion model



an old rusted robot wearing pants and a jacket riding skis in a supermarket.

Stable Diffusion 3



smiling cartoon dog sits at a table, coffee mug on hand, as a room goes up in flames. "This is fine," the dog assures himself.

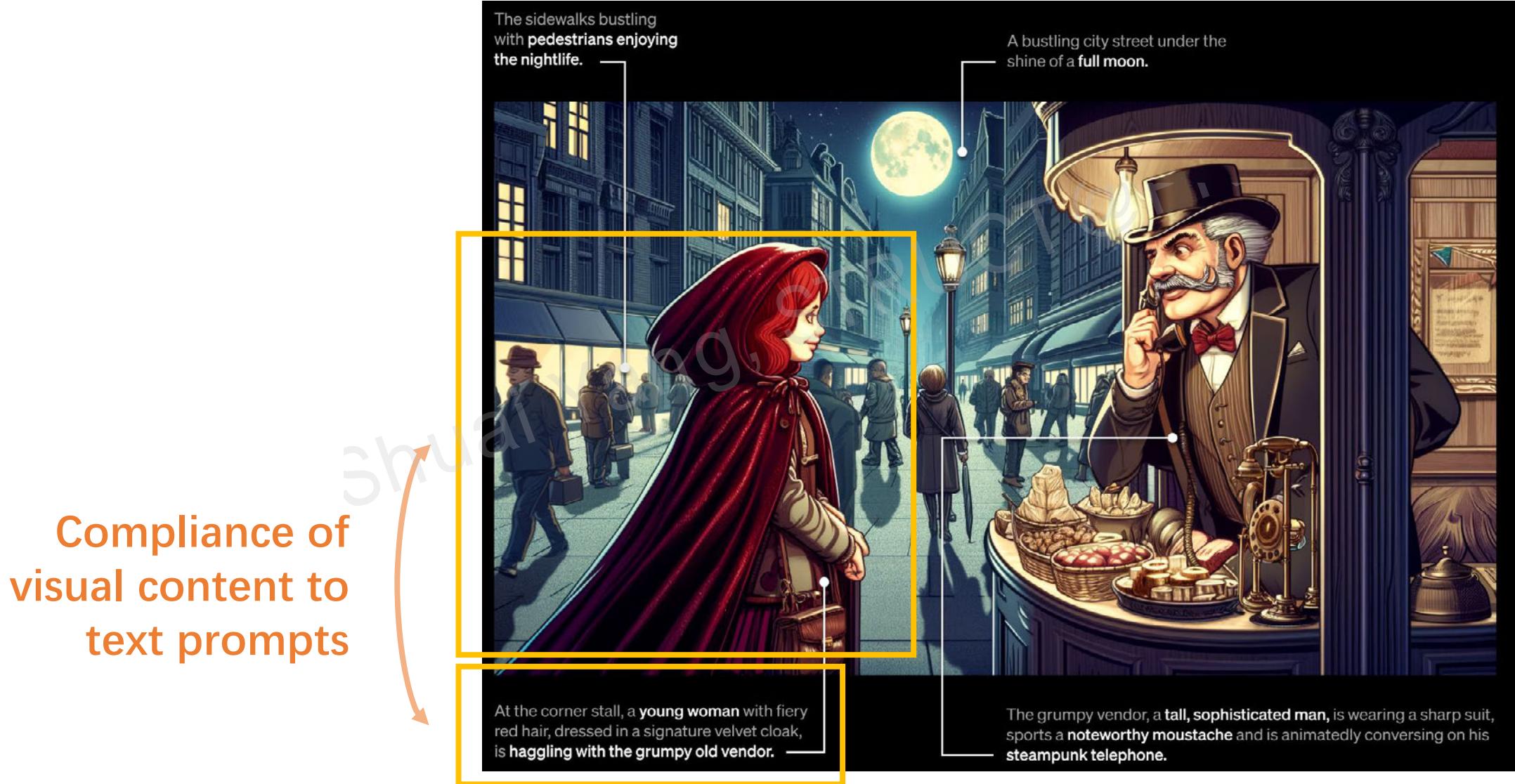
Image diffusion model

Stable Diffusion 3



Image diffusion model

DALL·E 3



Problem of data scraped from the internet



Alt Text: now at victorianplumbing.co.uk

- Problem 1: Incomplete descriptions.
- Problem 2: Irrelevant descriptions.

Dataset Recaptioning



Short Synthetic Captions:

a white modern bathtub sits on a wooden floor

Descriptive Synthetic Captions:

this luxurious bathroom features a modern freestanding bathtub in a crisp white finish. the tub sits against a wooden accent wall with glass-like panels, creating a serene and relaxing ambiance. three pendant light fixtures hang above the tub, adding a touch of sophistication. a large window with a wooden panel provides natural light, while a potted plant adds a touch of greenery. the freestanding bathtub stands out as a statement piece in this contemporary bathroom.

Image diffusion model

DALL·E 3

DALL·E 3 + ChatGPT = Interactive and more user-friendly

ChatGPT



Image diffusion model

Midjourney





Image Supporting Model



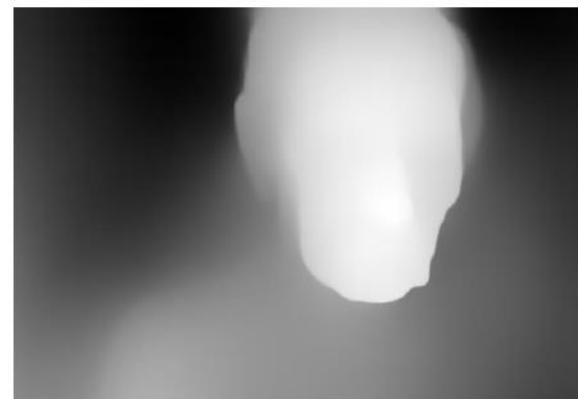
in the jungle on red fabric



at Mt. Fuji on top of snow



with Eiffel Tower



The boulevards are crowded today.

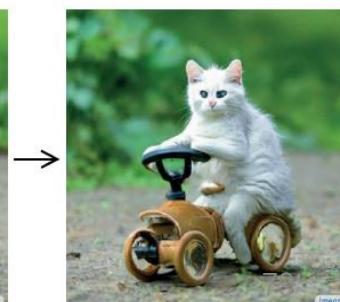


Photo of a cat riding on a bicycle.

Image supporting model

Approaches

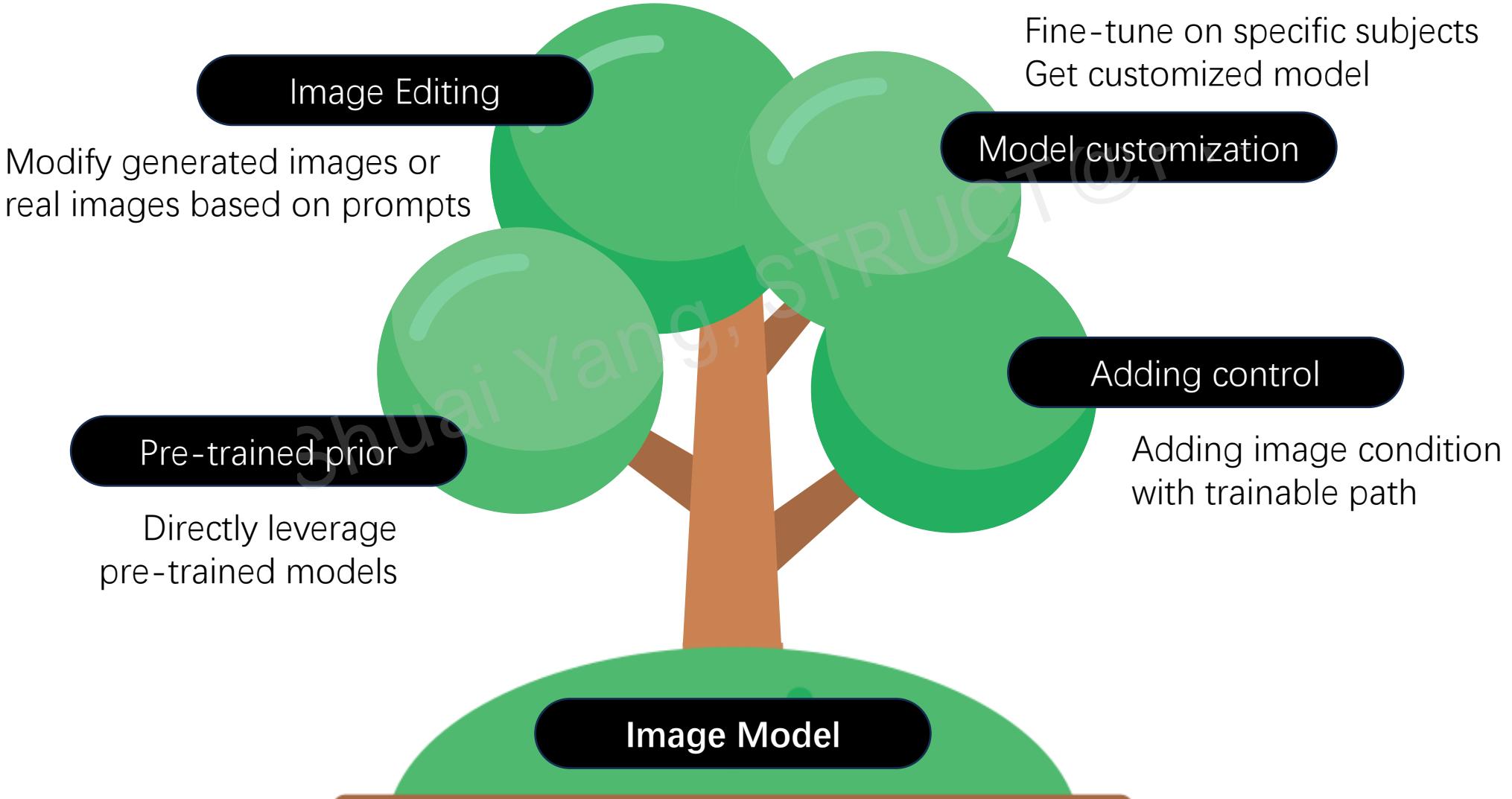
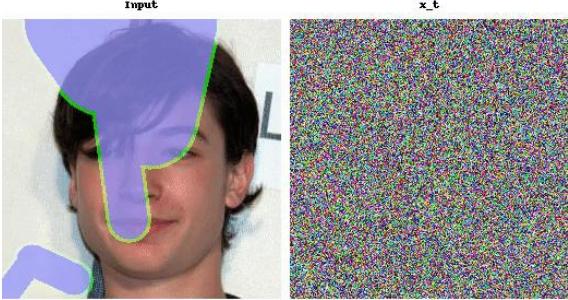


Image supporting model



Pre-trained prior

- No additional training
- Directly leveraging generative prior to specific task
- Limited customizability



Image editing

- No additional training
- Use cross-attention map between text and image feature
- Cannot introduce new modalities



Input images

in the Acropolis

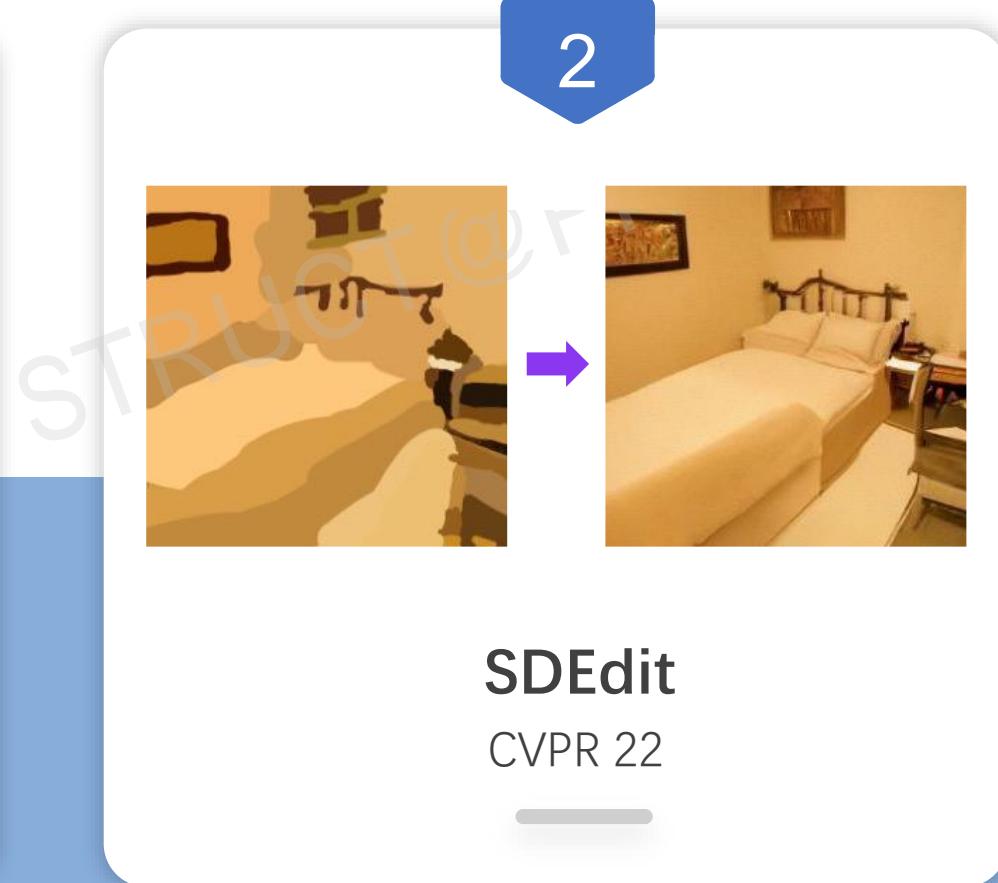


Adding control

- Subject-specific generation
- Highly customized
- Needs finetune on each case

- Diverse modalities as guidance
- Highly controllable
- Needs training on customized dataset

Pre-trained models as prior



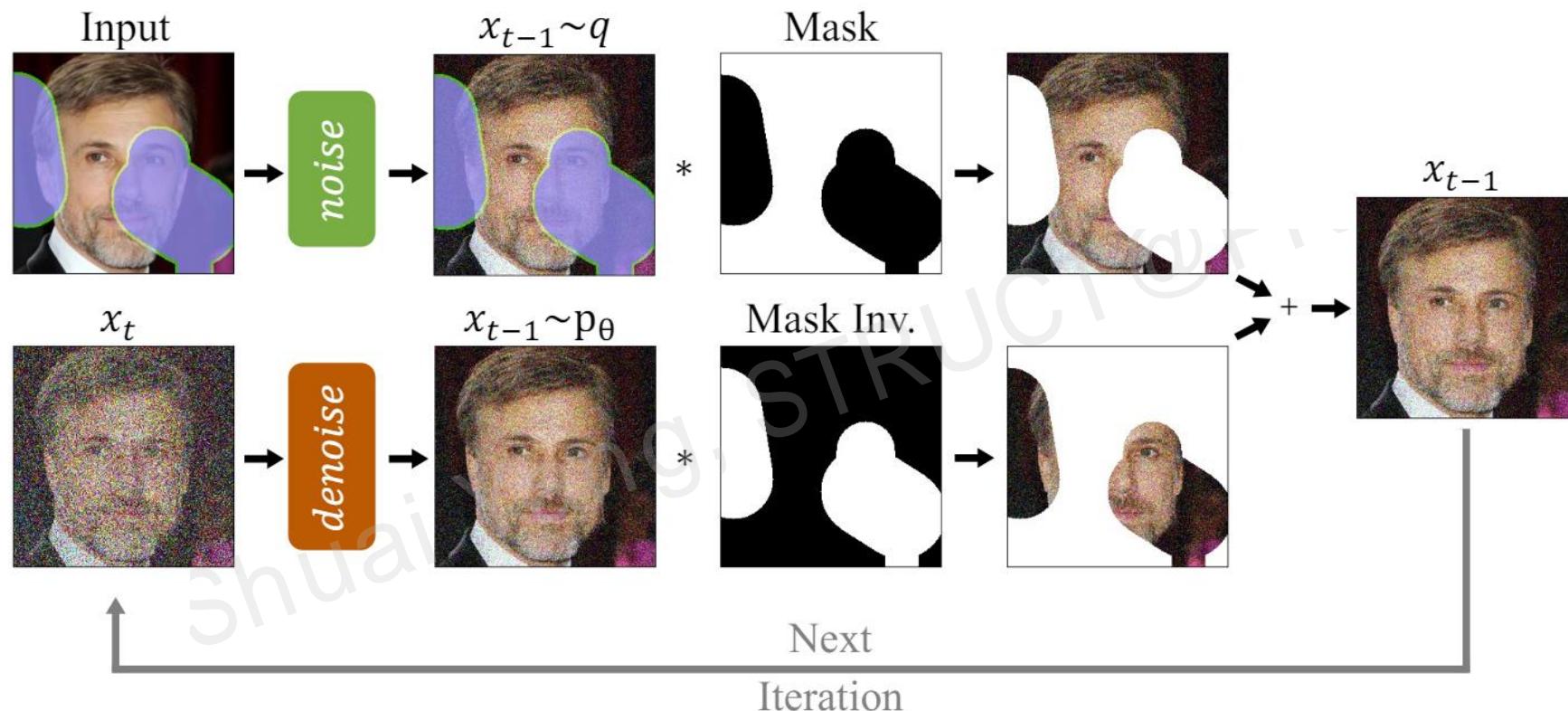
Pre-trained models as prior

1 RePaint



Pre-trained models as prior

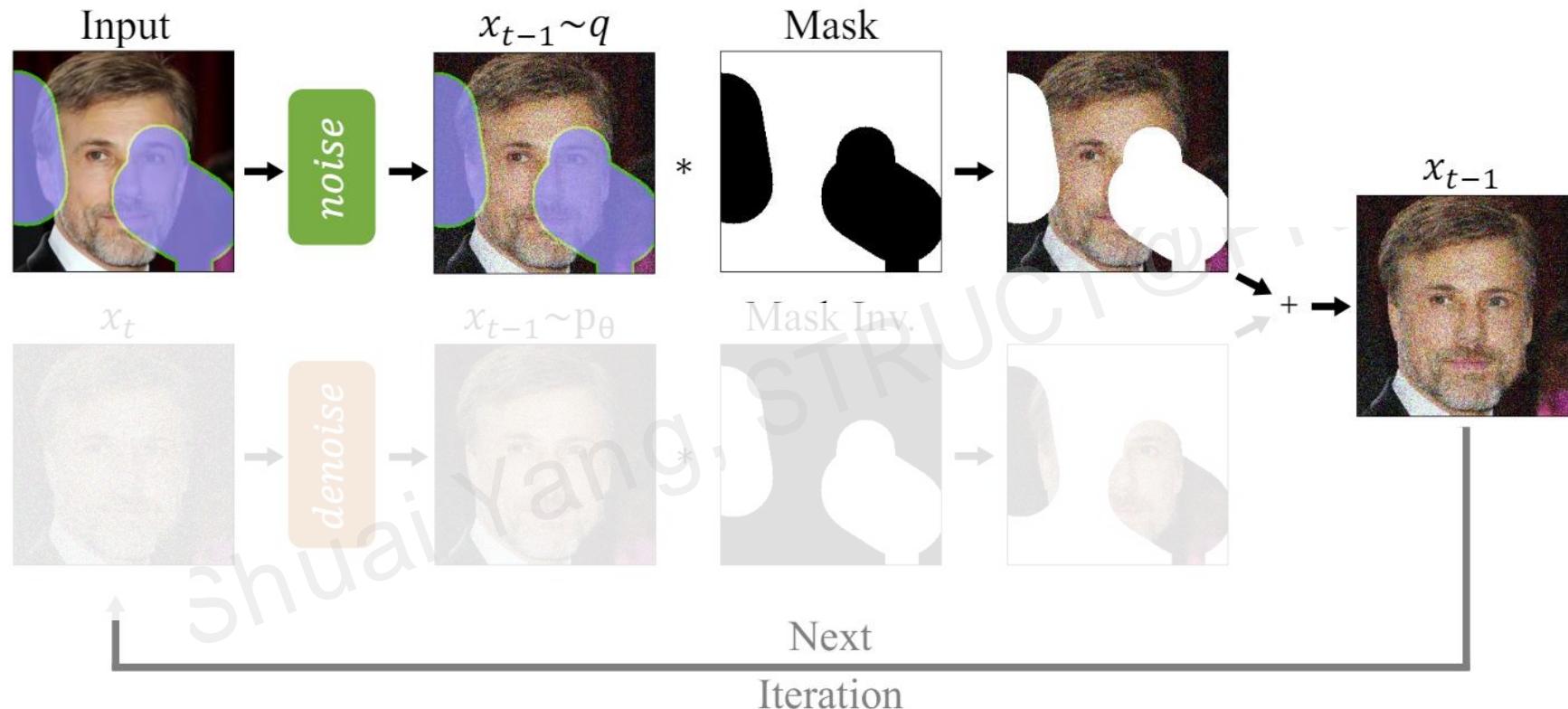
1 RePaint



Iteratively generate missing region based on visible parts

Pre-trained models as prior

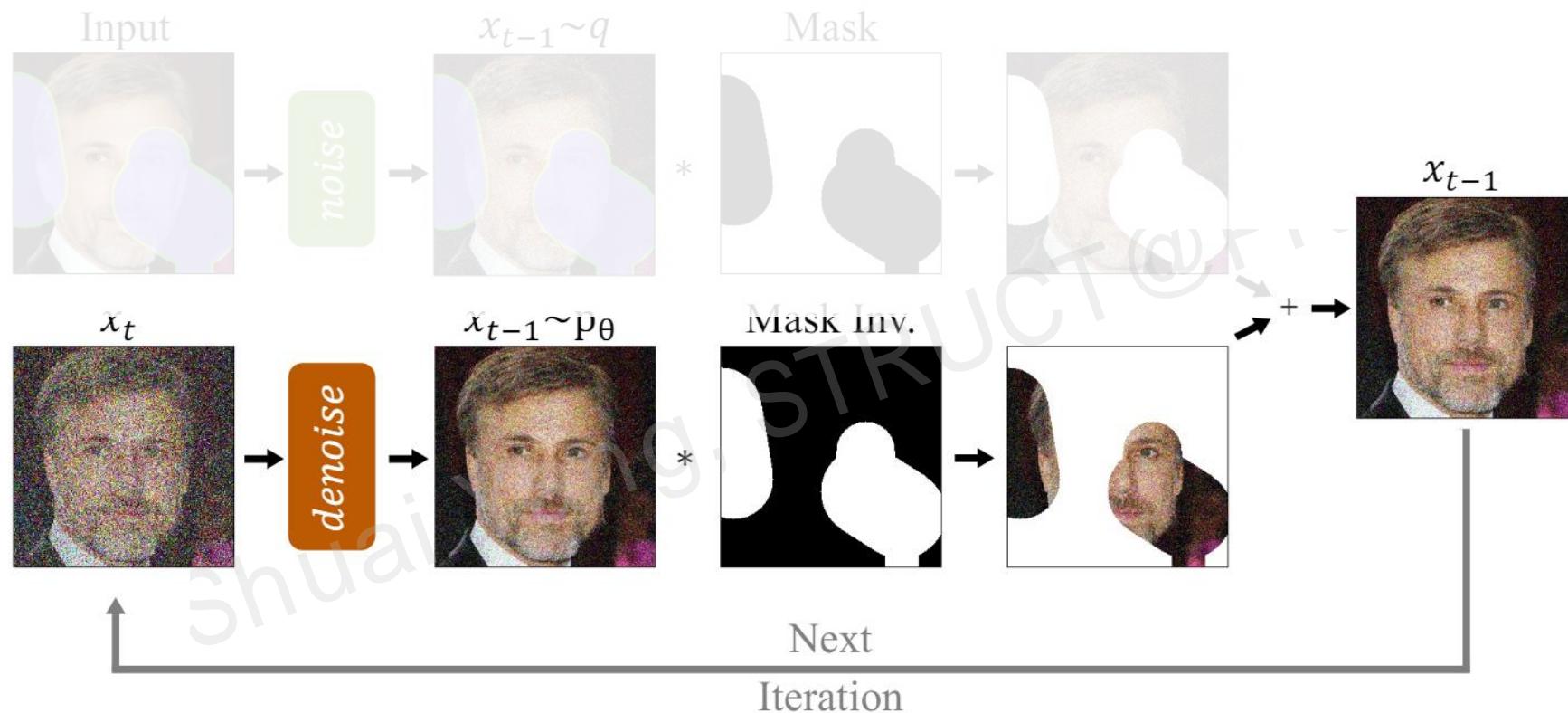
1 RePaint



Iteratively generate missing region based on visible parts

Pre-trained models as prior

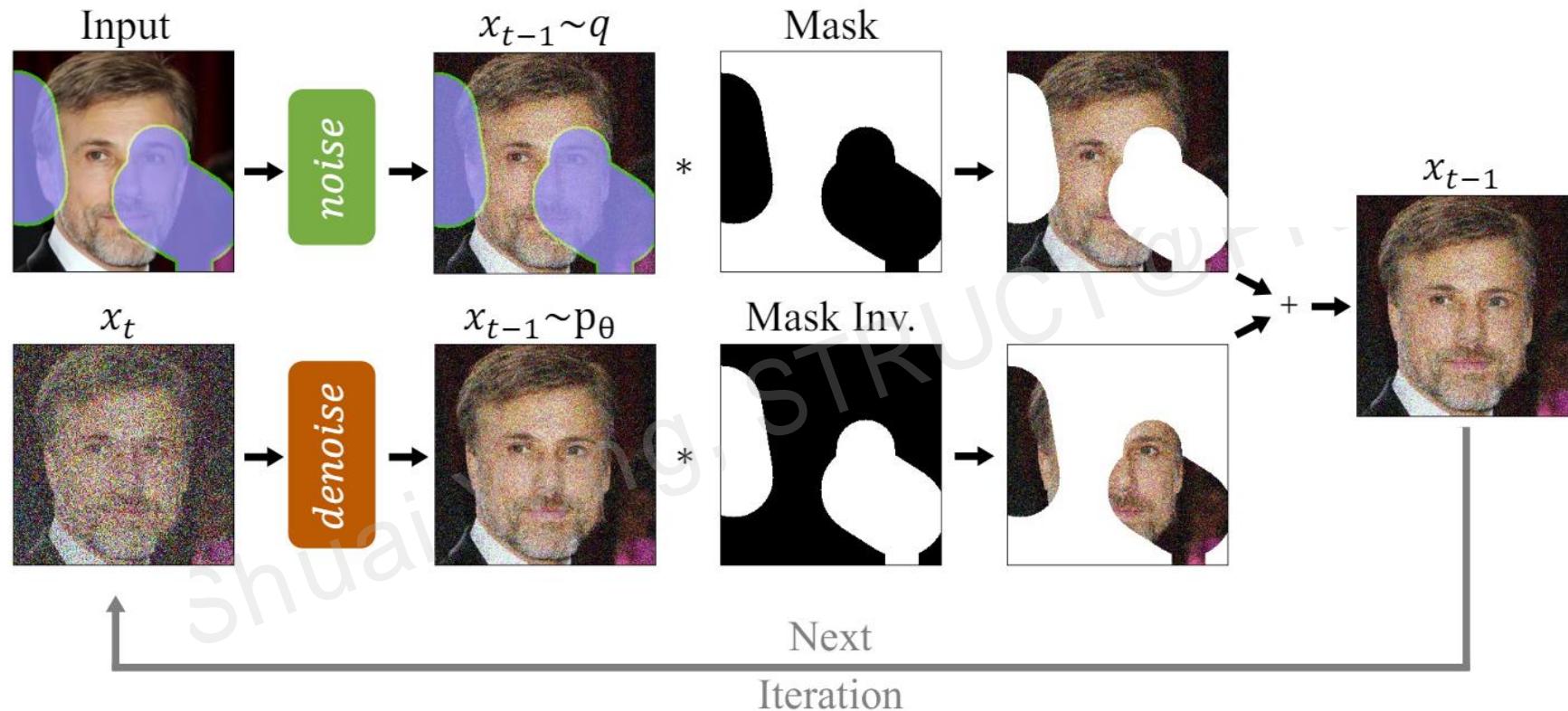
1 RePaint



Iteratively generate missing region based on visible parts

Pre-trained models as prior

1 RePaint



Iteratively generate missing region based on visible parts

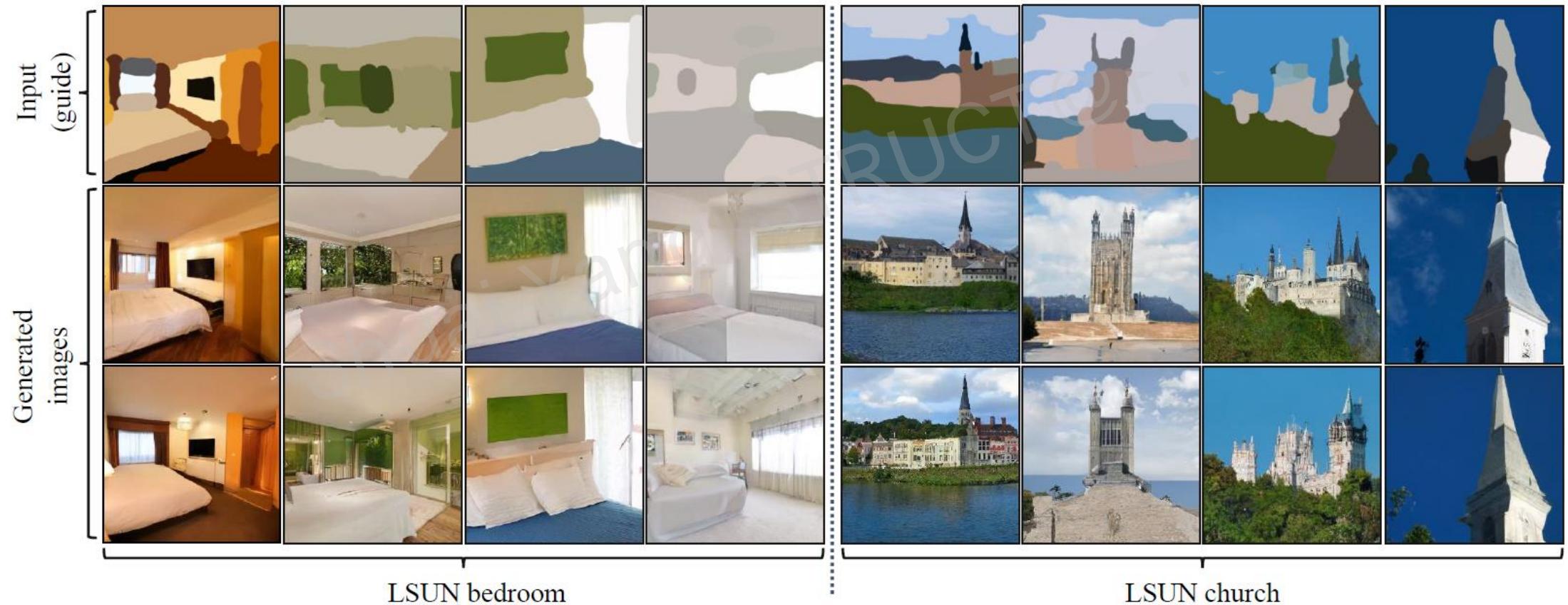
Multiple tasks as inpainting



Pre-trained models as prior

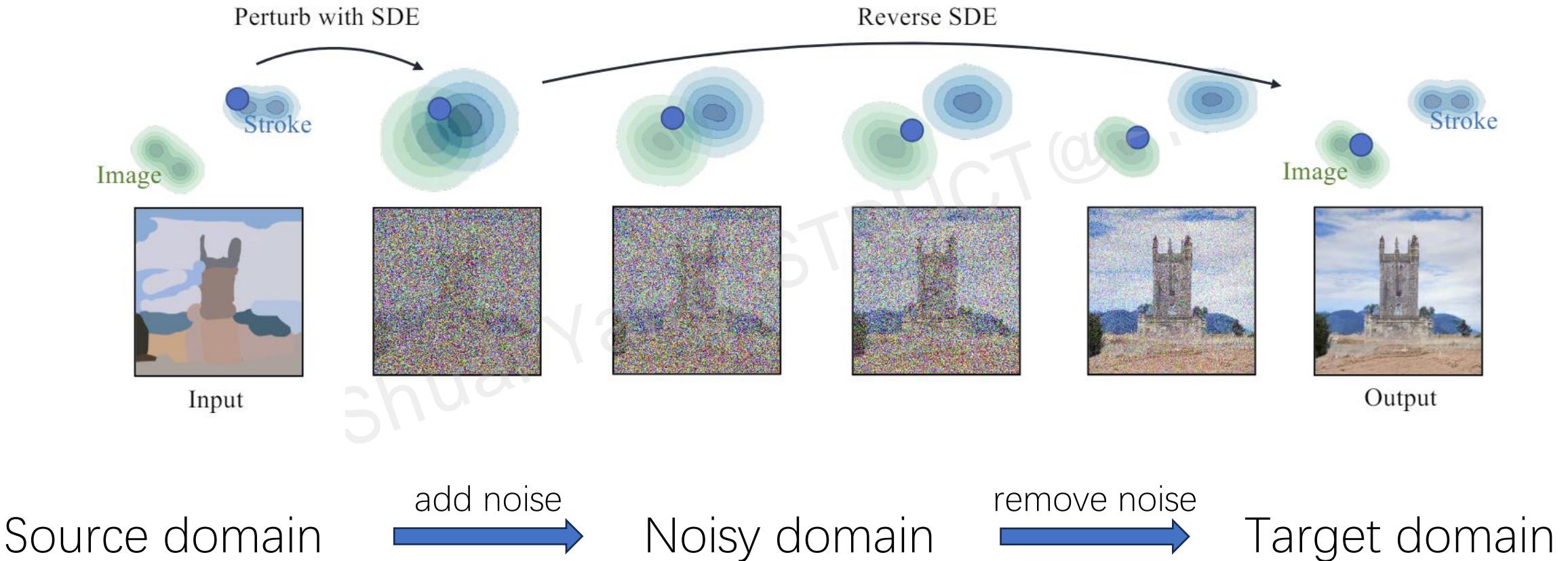
2 SDEdit

Cross-domain transformation



Pre-trained models as prior

2 SDEdit



Pre-trained models as prior

2 SDEdit

Source
Input (guide)
SDEdit (Ours)

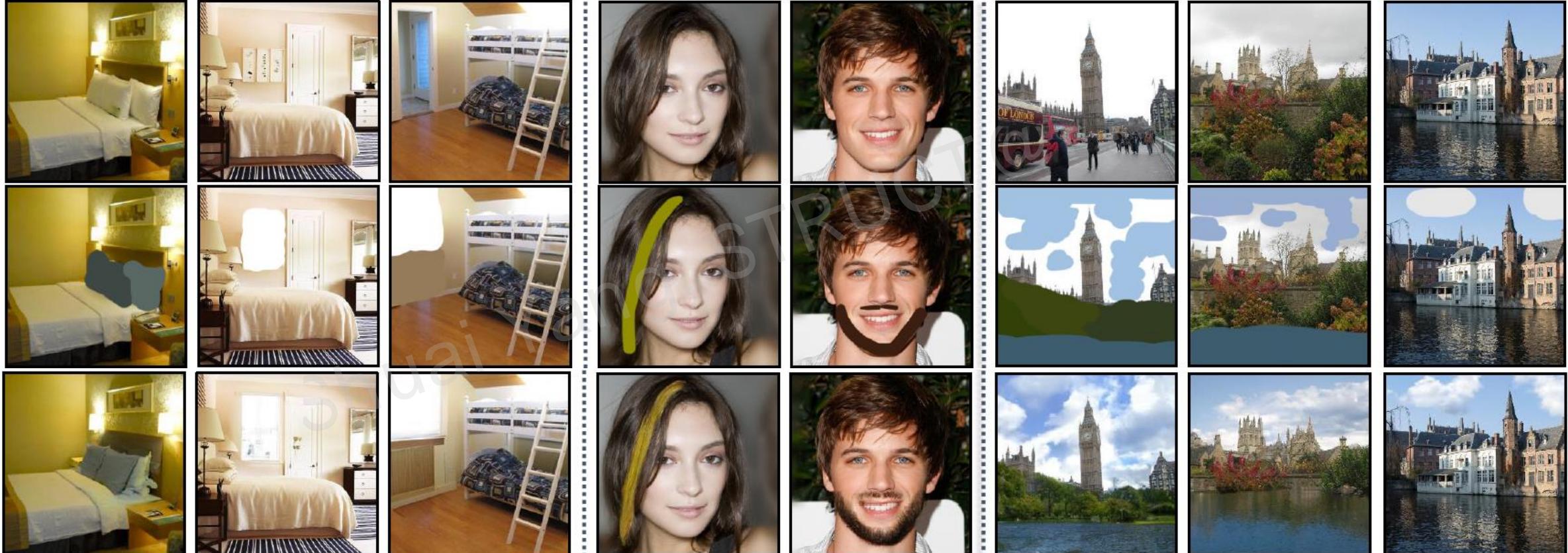


Image editing

1



"The boulevards are crowded today."

2



"a photo of a pink toy horse on the beach"

Prompt2Prompt

ICLR 23

Plug-and-Play

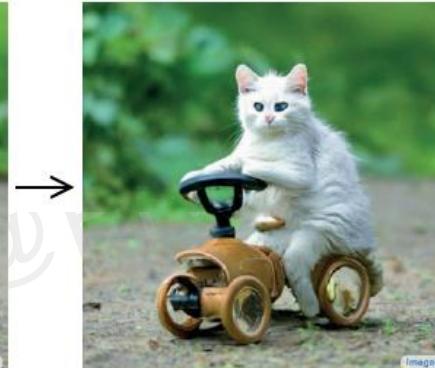
CVPR 23

Image editing

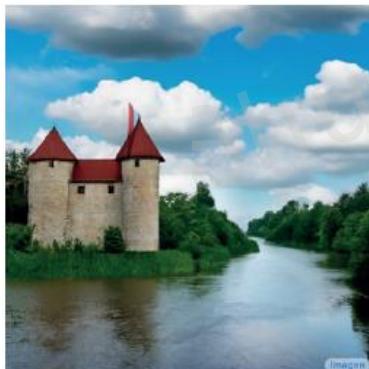
1 Prompt-to-Prompt



"The boulevards are **crowded** today."



"Photo of a cat riding on a **bicycle**."
~~car~~



"Children drawing of a castle next to a river."



"a cake with decorations."
jelly beans

Image editing

1 Prompt-to-Prompt

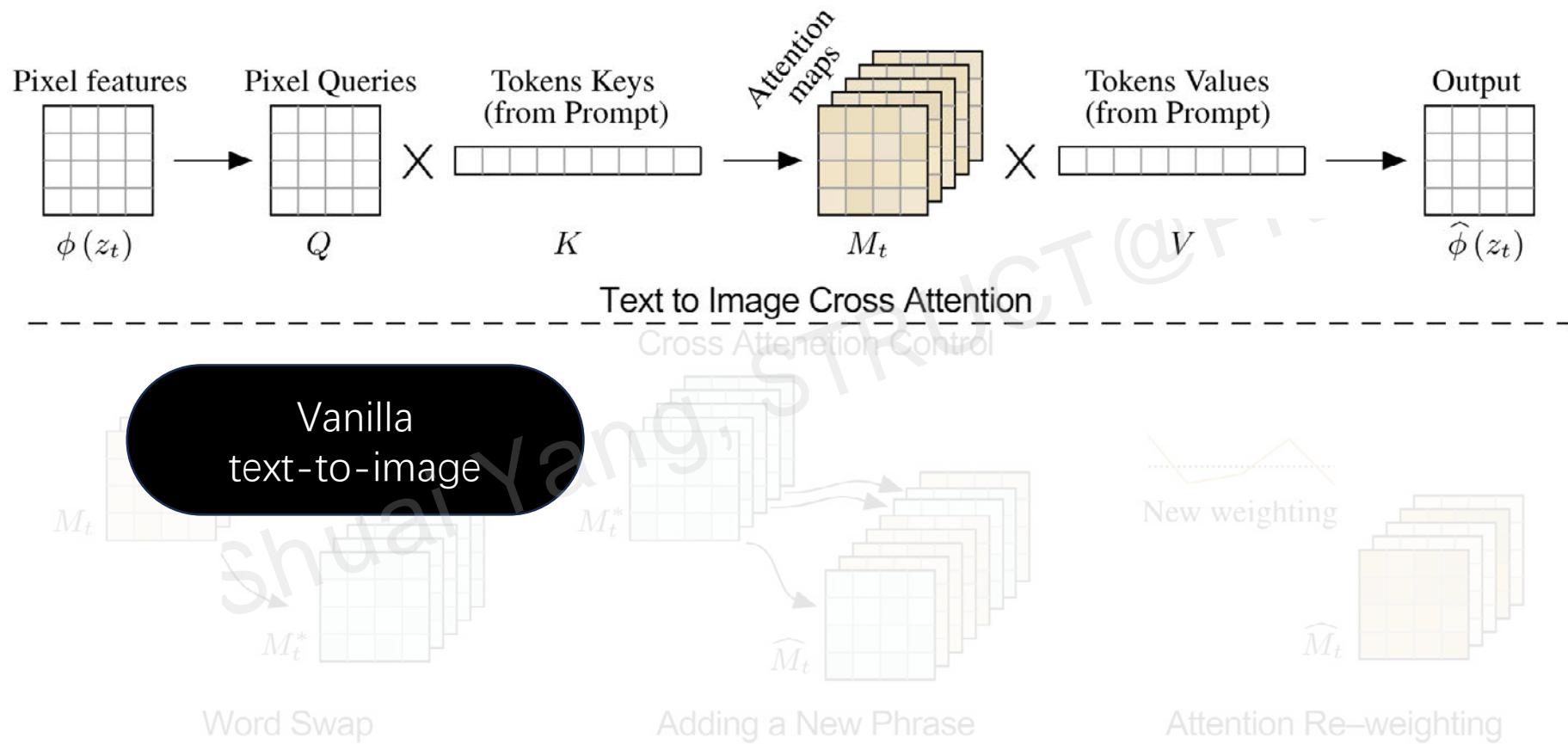


Image editing

1 Prompt-to-Prompt

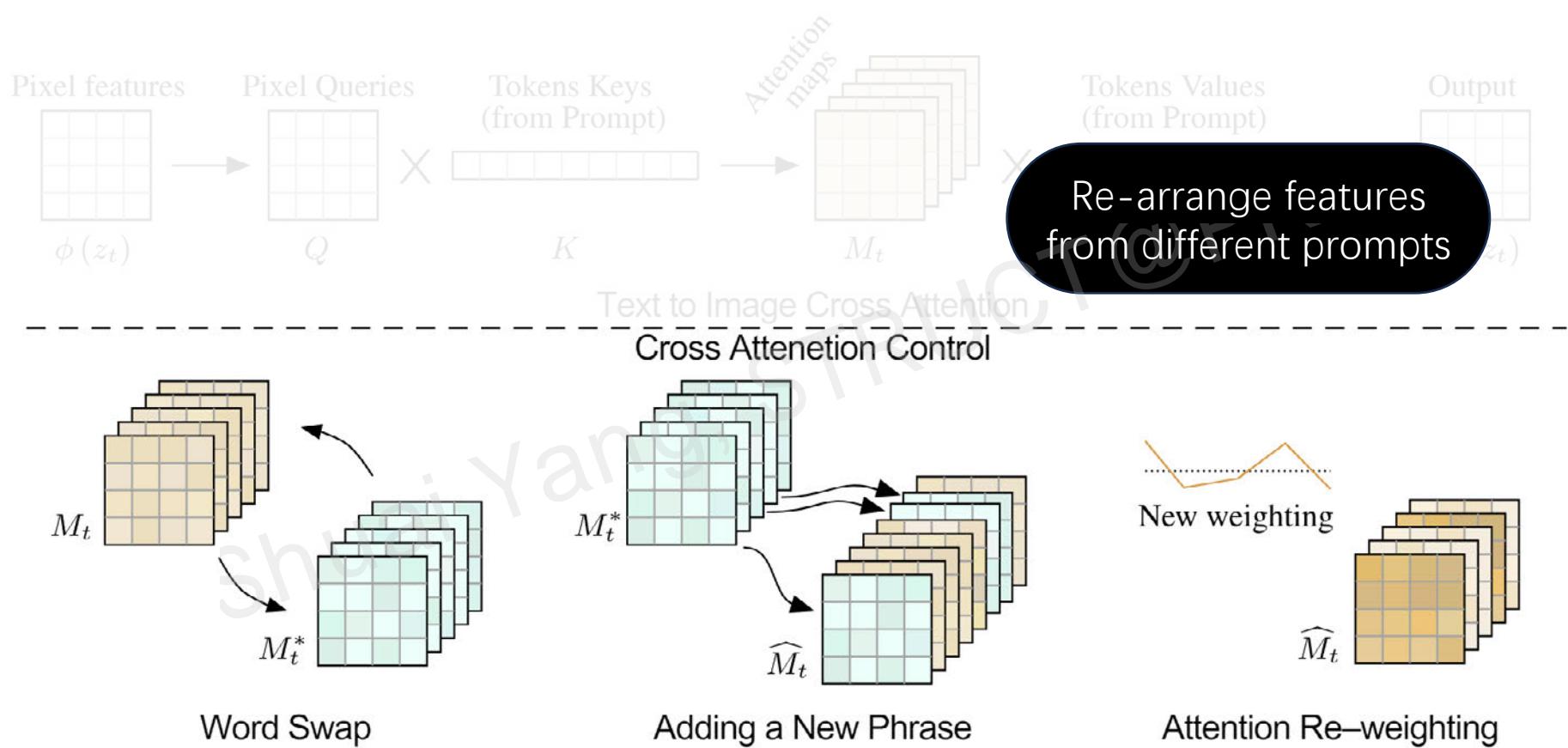


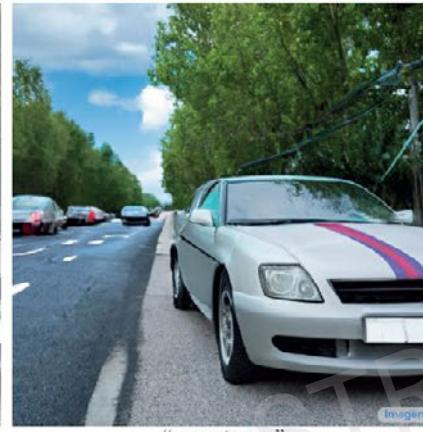
Image editing

1 Prompt-to-Prompt

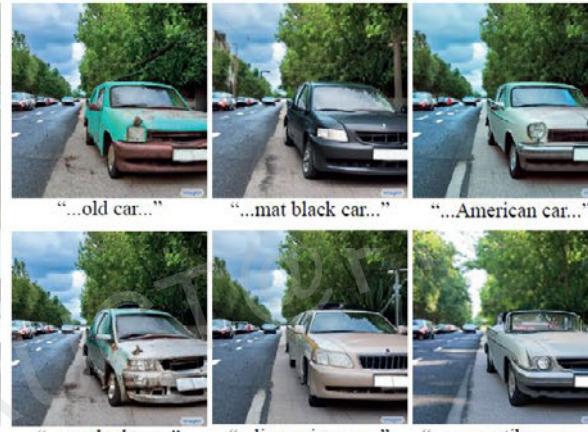
“A car on the side of the street.”



source image



“...sport car...”



“...old car...” “...mat black car...” “...American car...”



“...crushed car...” “...limousine car...” “...convertible car...”

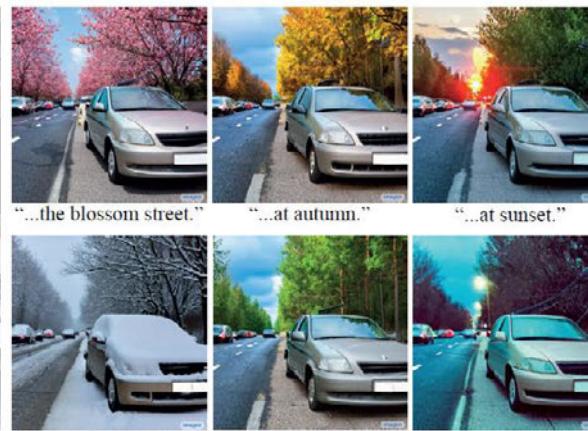
Local description
Global description



“...the flooded street.”



“...in Manhattan.”



“...the blossom street.”

“...at autumn.”

“...at sunset.”



“...in the snowy street.”

“...in the forest.”



“...at evening.”

Image editing

1 Prompt-to-Prompt

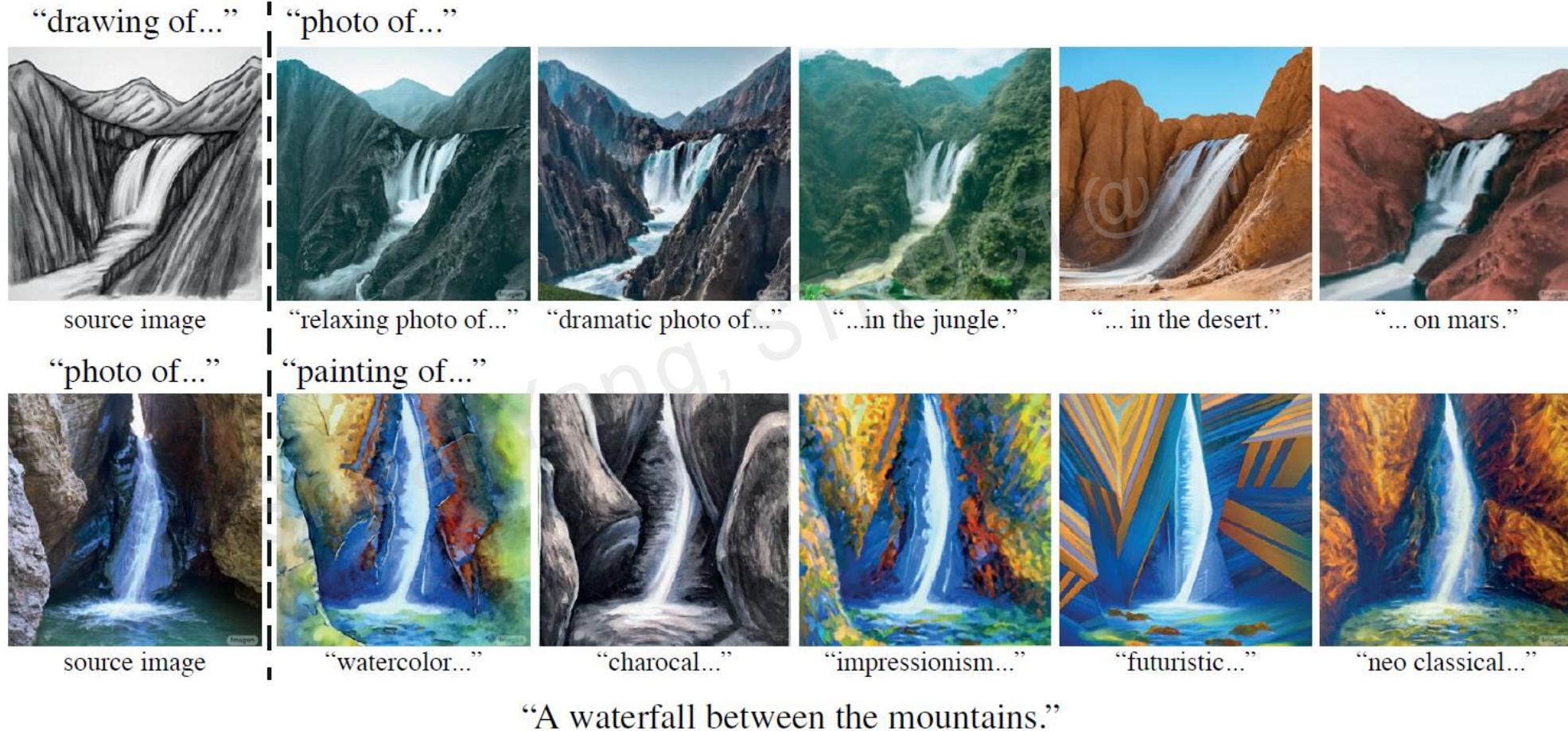


Image editing

1 Prompt-to-Prompt

“A black bear is walking in the grass.”



“Landscape image of trees in a valley...”



Real-world image editing

Image editing

2 Plug-and-Play



Input Real Image



"a photo of a bronze horse in a museum"



"A photo of a pink horse on the beach"



"A photo of a robot horse"



Input Real Image



"A wooden sculpture of a couple dancing"



"A cartoon of a couple dancing"



"a photo of robots dancing"

Image editing

2 Plug-and-Play



Image editing

2 Plug-and-Play

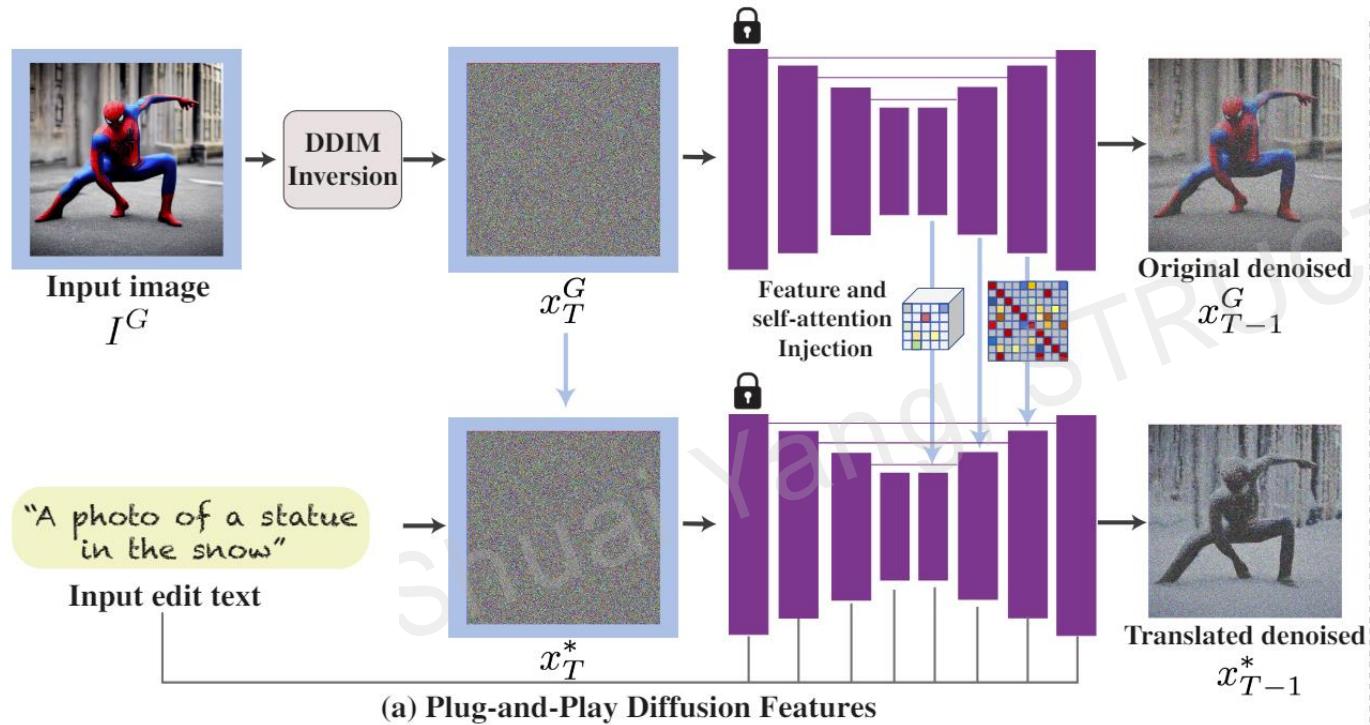


Image editing

2 Plug-and-Play

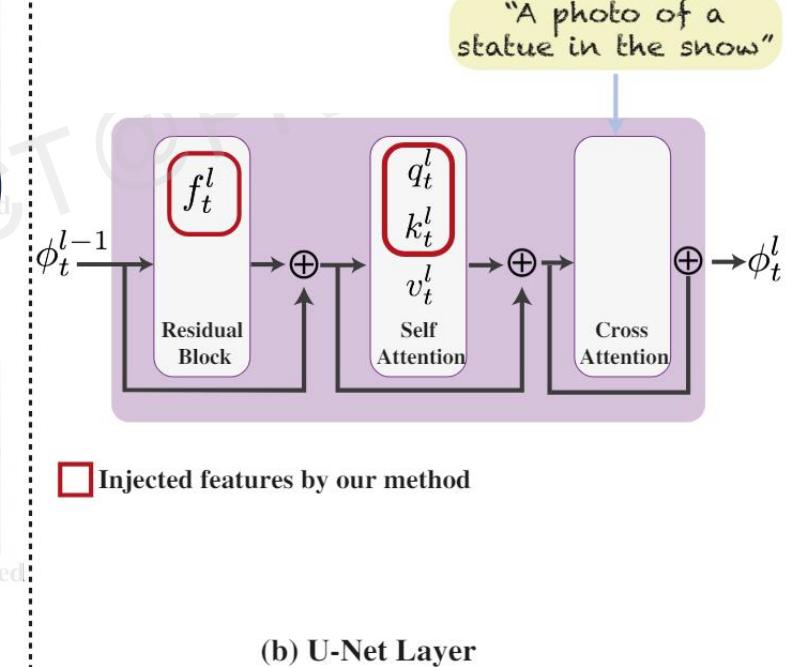
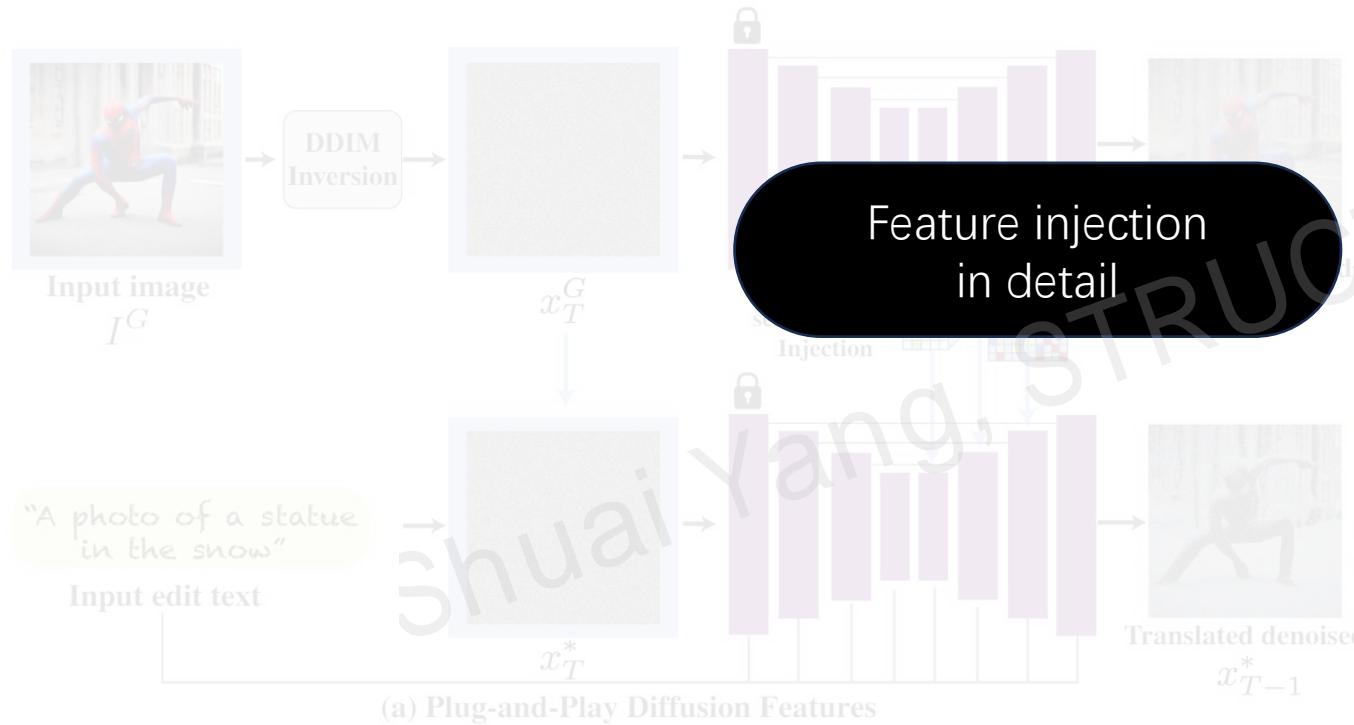


Image editing

2 Plug-and-Play

Guidance



"a photo of a golden robot"

"a photo of a wooden statue"

"a photo of a sand sculpture"

"a photo of a golden sculpture in a temple"

"a photo of a wooden statue"

"a photo of a silver robot"



Input Real Image



"A polygonal illustration of a cat and a bunny"



"A photo of bear cubs in the snow"



Input Generated Image



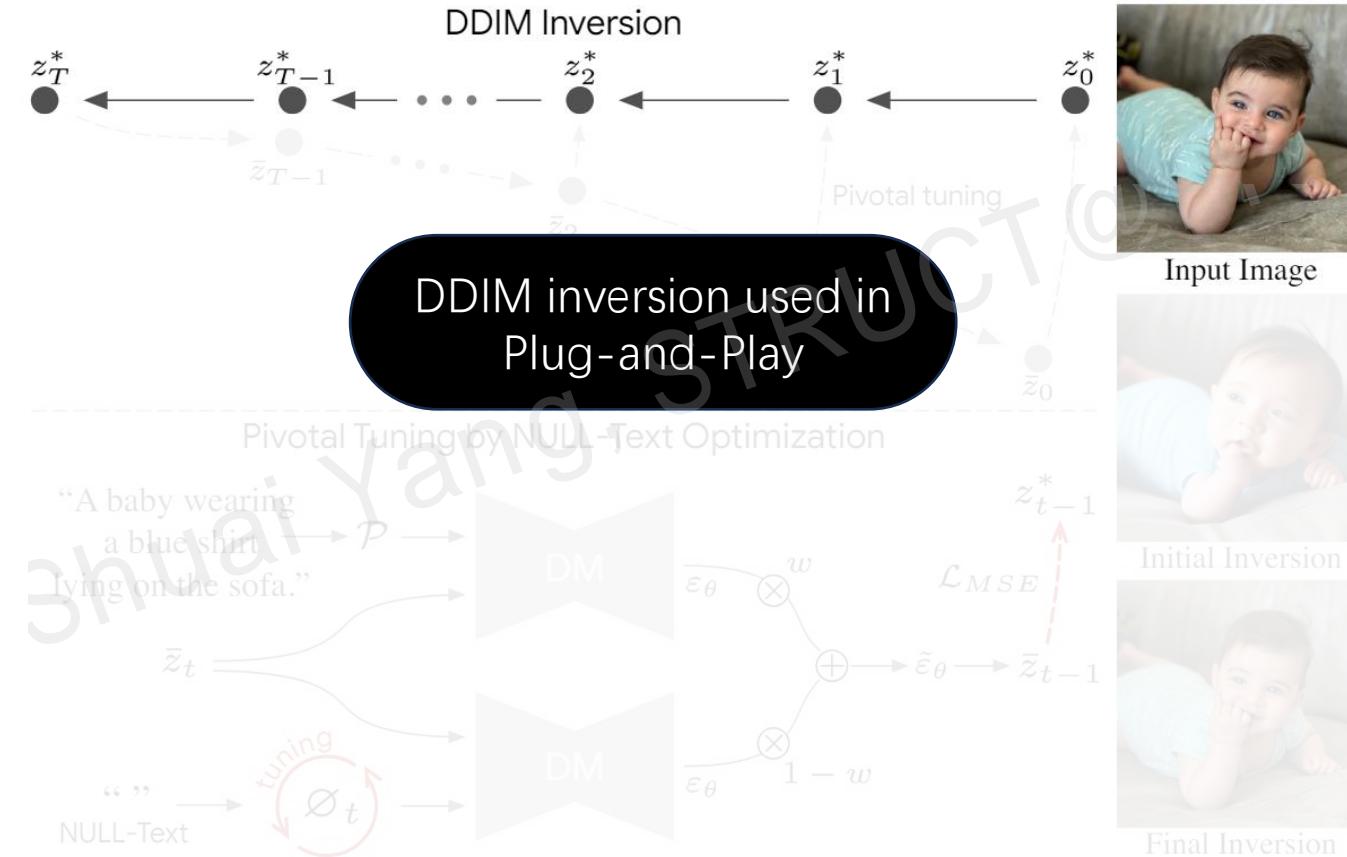
"A polygonal illustration of fish in the ocean"



"A photo of sharks in the ocean"

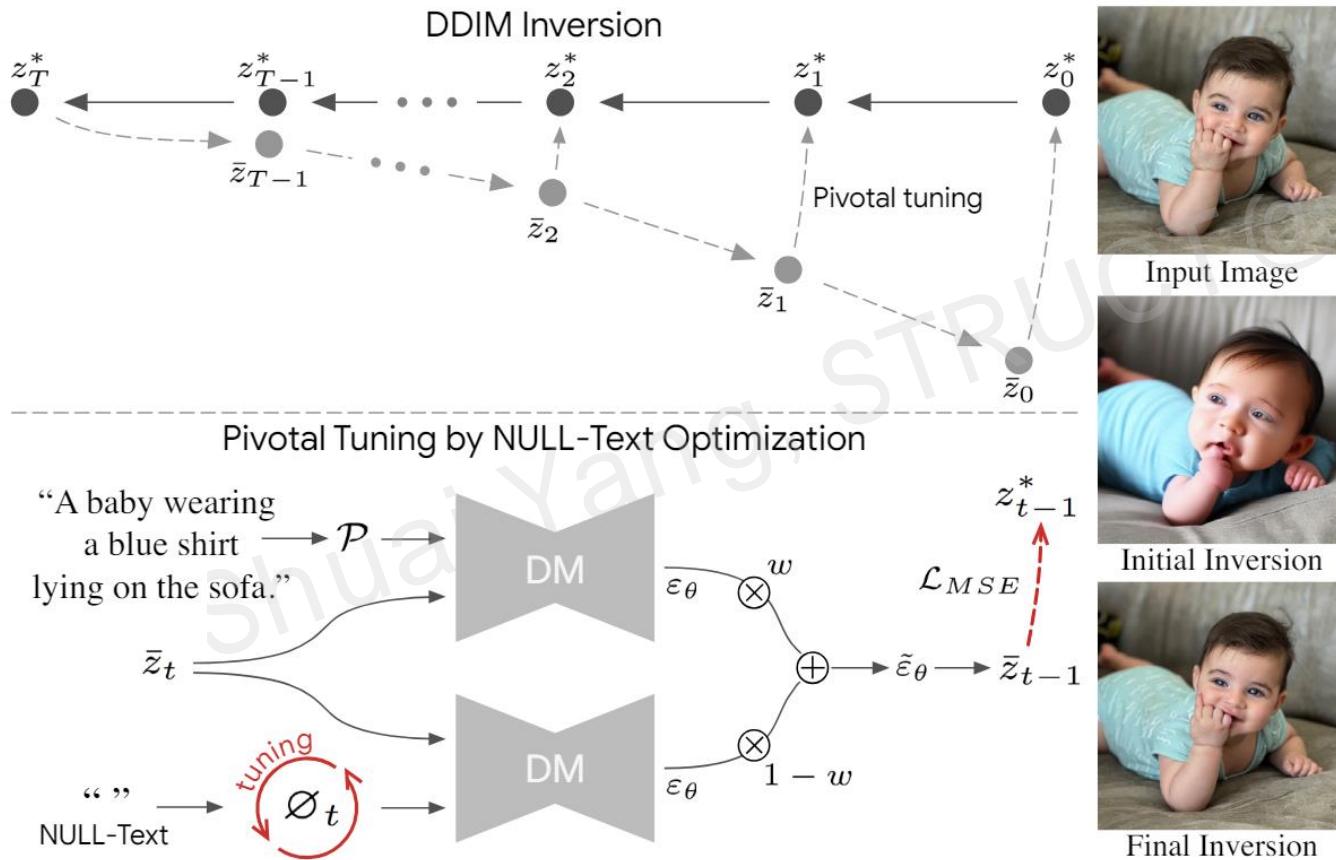
Better inversion for editing

Null-text inversion



Better inversion for editing

Null-text inversion



Finetuning “Null-Text”
for better inversion

Better inversion for editing

Null-text inversion

Input caption: “A baby wearing a blue shirt lying on the sofa.”



Input Image



“... blond baby...”



“... floral shirt...”



“... golden shirt...”



“... sleeping baby...”



“baby” → “robot”



“sofa” → “grass”



“sofa” → “ball pit”

Null-Text Inversion



Prompt to Prompt

Better inversion for editing

Null-text inversion

Input caption: “A man in glasses eating a doughnut in the park.”



Input Image



“... red-haired man...”



“glasses” → “sunglasses”



“angry man...”



“doughnut” → “pizza”



“glasses” → “Joker mask”



“...the park at sunset.”



“park” → “desert”

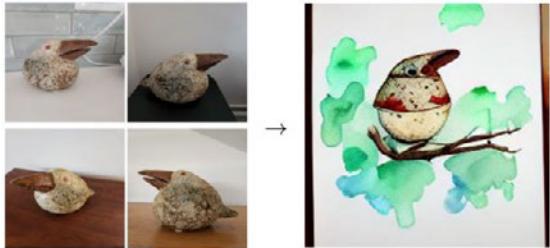
Null-Text Inversion



Prompt to Prompt

Model customization

1



"Watercolor painting of *S.* on a branch"

Textual Inversion

ICLR 23

2



DreamBooth

CVPR 23

3



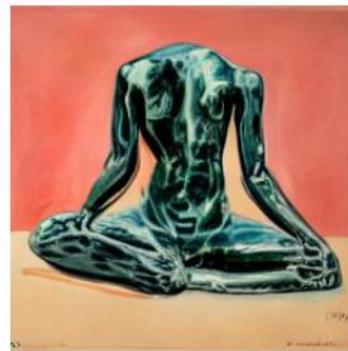
GitHub Repo

Model customization

1 Textual Inversion



→



Input samples $\xrightarrow{\text{invert}}$ “ S_* ”

“An oil painting of S_* ”

“App icon of S_* ”

“Elmo sitting in
the same pose as S_* ”

“Crochet S_* ”



→



Input samples $\xrightarrow{\text{invert}}$ “ S_* ”

“Painting of two S_*
fishing on a boat”

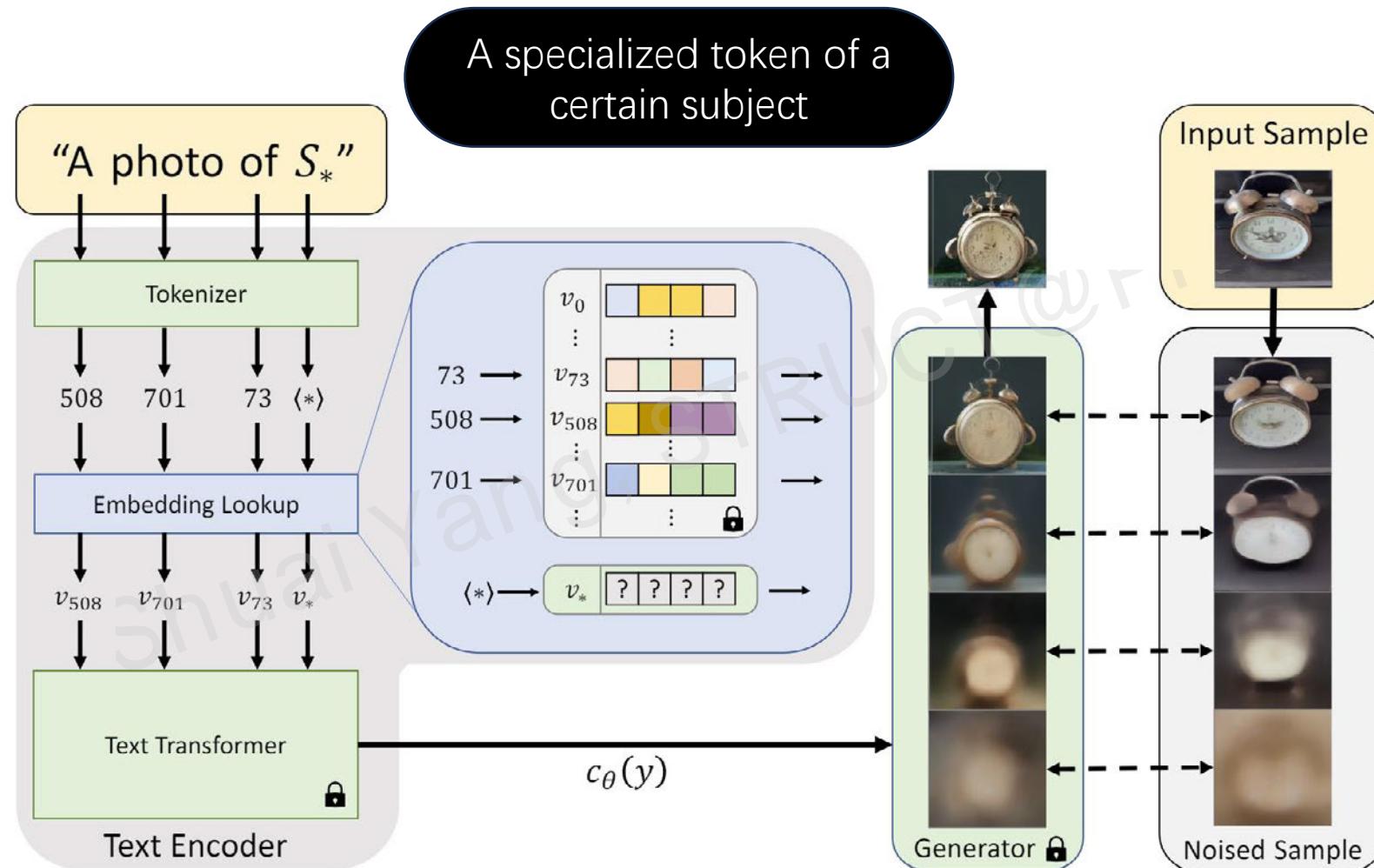
“A S_* backpack”

“Banksy art of S_* ”

“A S_* themed lunchbox”

Model customization

1 Textual Inversion



Model customization

1 Textual Inversion



Input samples



"Watercolor painting of S_* on a branch"



"A house in the style of S_* "



"Grainy photo of S_* in angry birds"



" S_* made of chocolate"



Input samples



"A mosaic depicting S_* "



"Death metal album cover featuring S_* "



"Masterful oil painting of S_* hanging on the wall"



"An artist drawing a S_* "

Model customization

1 Textual Inversion



Input samples

→



"A photo of S_* full of cashew nuts"



"A mouse using S_* as a boat"



"A photo of a S_* mask"



"Ramen soup served in S_* "



Input samples

→



"A carpet with S_* embroidery"



"A S_* stamp"



" S_* fedora"



"Felt S_* "

Model customization

2 DreamBooth



Input images



in the Acropolis



swimming



sleeping



getting a haircut



Input images



worn by a bear



in the jungle



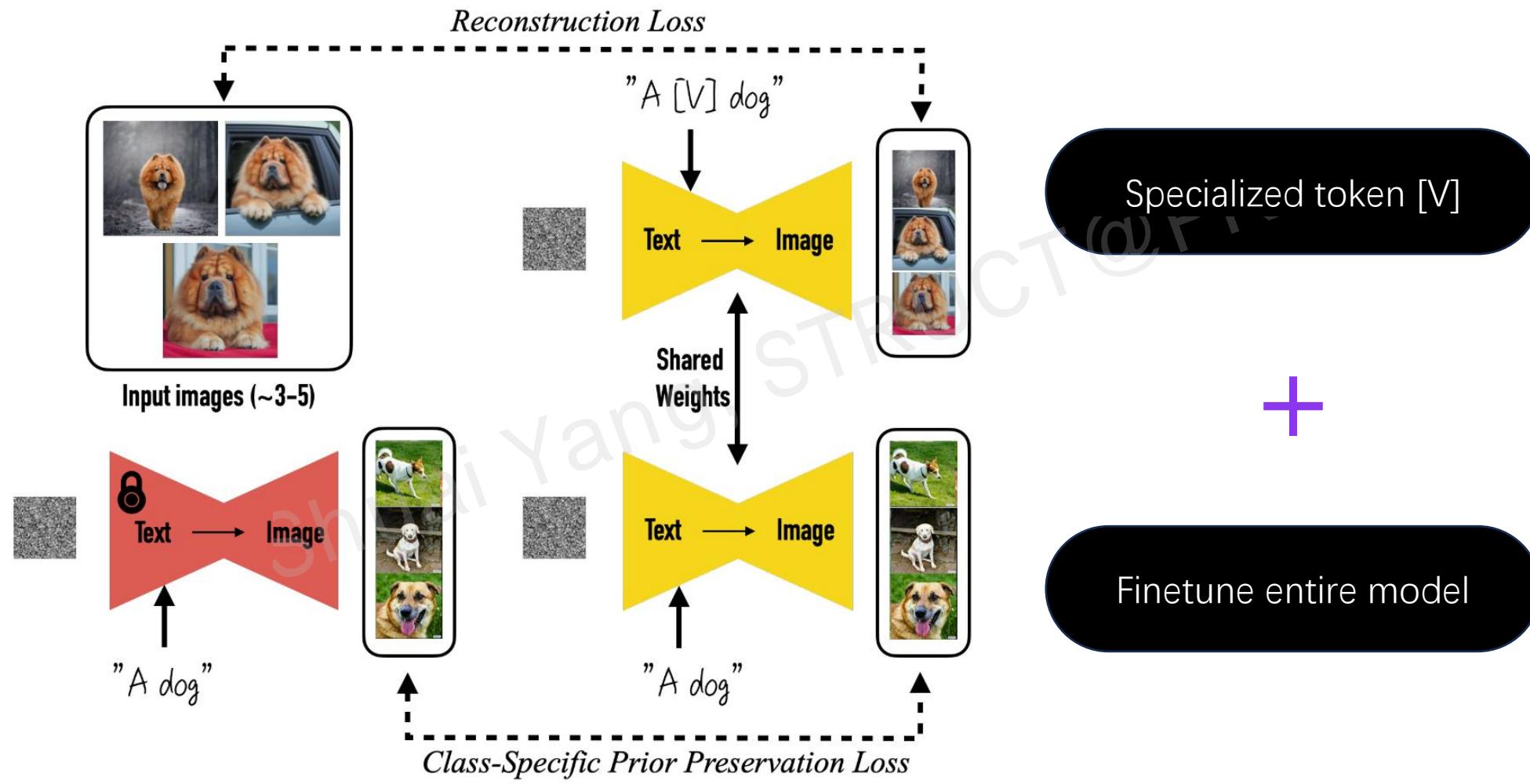
on red fabric



with Eiffel Tower

Model customization

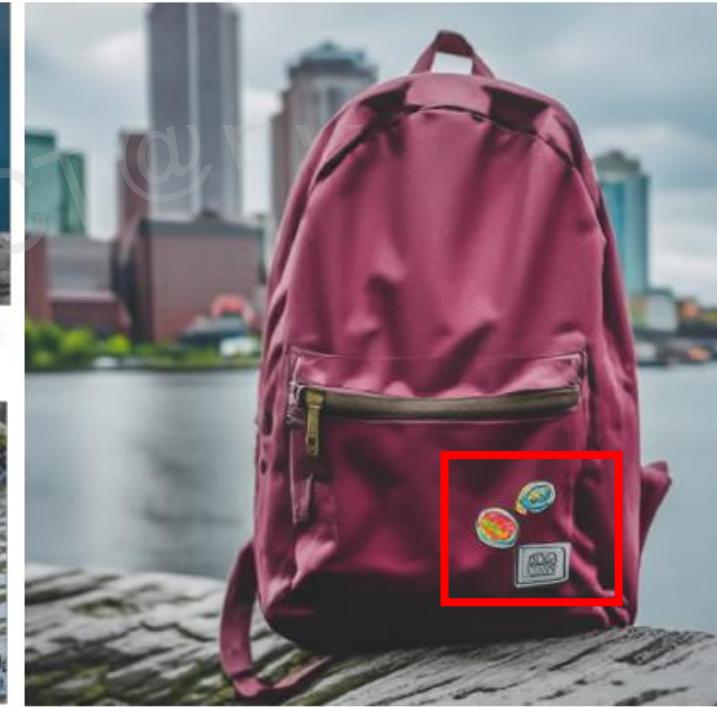
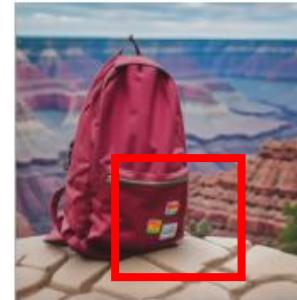
2 DreamBooth



Model customization

2 DreamBooth

Input images



Model customization

2 DreamBooth

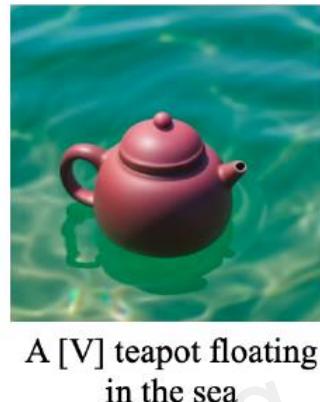
Input images



Model customization

2 DreamBooth

Input images



A [V] teapot floating
in the sea



A [V] teapot floating
in milk



A bear pouring from
a [V] teapot



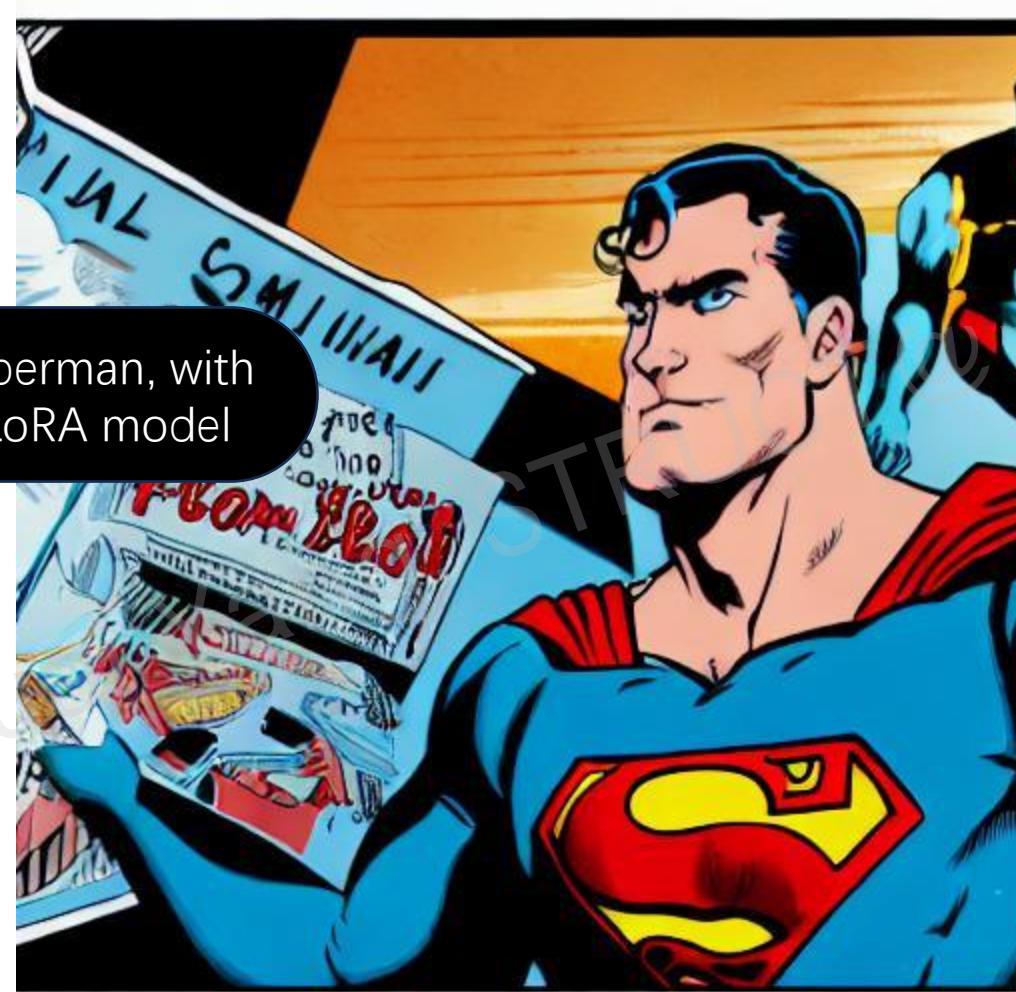
A transparent [V] teapot
with milk inside



A [V] teapot pouring tea

Model customization

3 LoRA

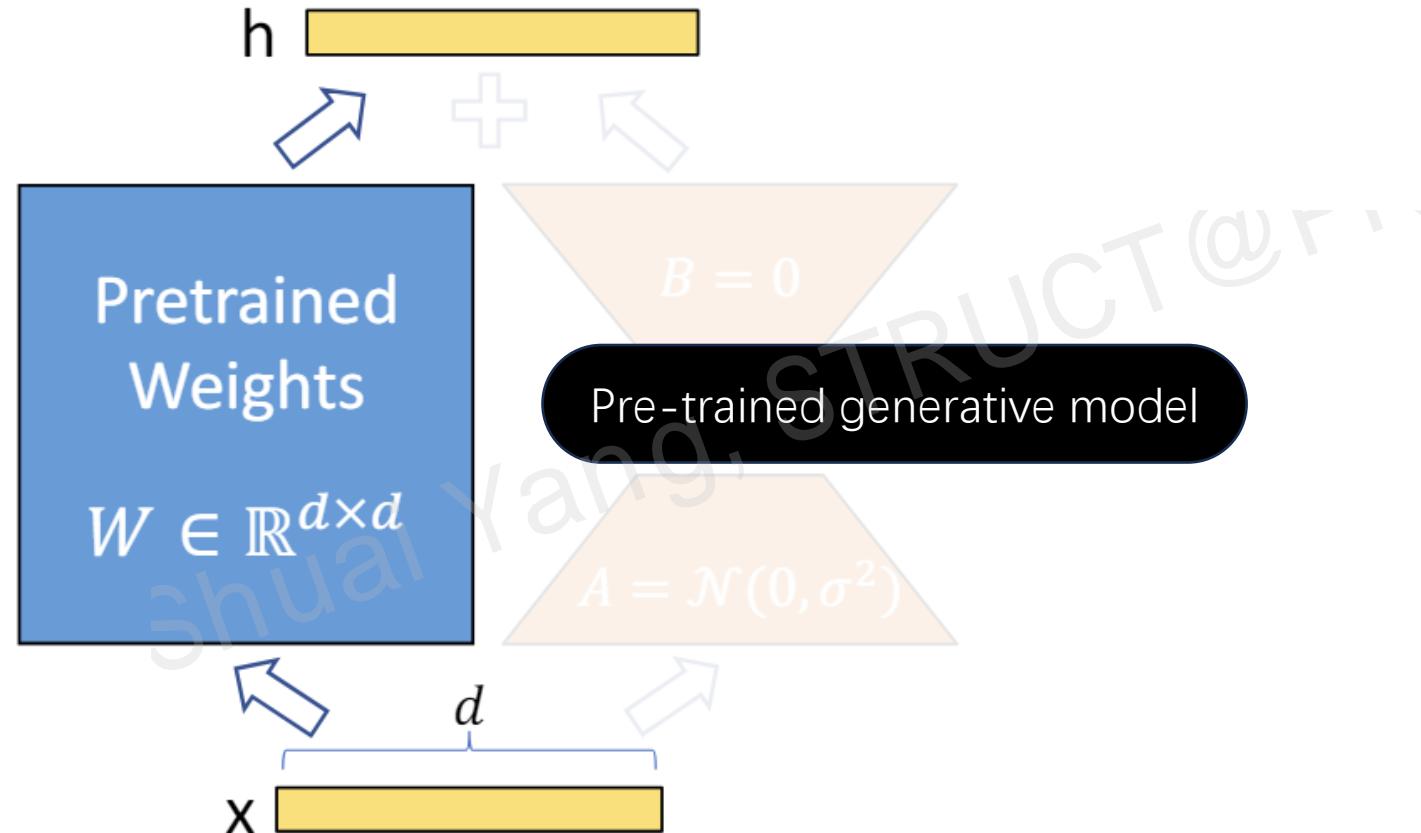


An image of superman, with
pop-art style LoRA model

Low-rank Adaptation (LoRA) for Fast Text-to-Image Diffusion Fine-tuning. GitHub repo
LoRA: Low-Rank Adaptation of Large Language Models. ICLR'22

Model customization

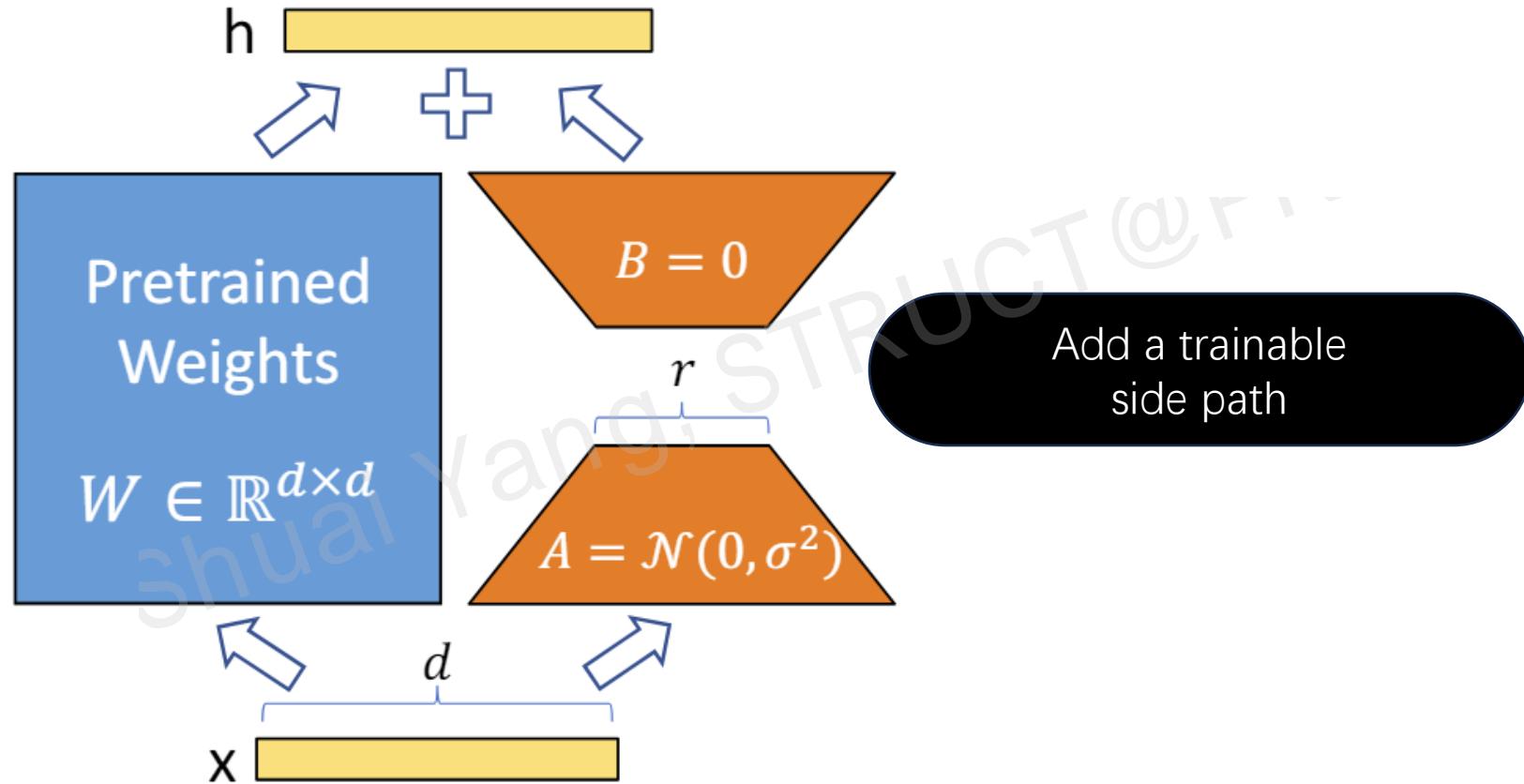
3 LoRA



Low-rank Adaptation (LoRA) for Fast Text-to-Image Diffusion Fine-tuning. GitHub repo
LoRA: Low-Rank Adaptation of Large Language Models. ICLR'22

Model customization

3 LoRA



Low-rank Adaptation (LoRA) for Fast Text-to-Image Diffusion Fine-tuning. GitHub repo
LoRA: Low-Rank Adaptation of Large Language Models. ICLR'22

Model customization

3 LoRA



An image of baby lion, with
pop-art style LoRA model

Model customization

3 LoRA



Pokémon, with
Pokémon style LoRA model

Low-rank Adaptation (LoRA) for Fast Text-to-Image Diffusion Fine-tuning. GitHub repo
LoRA: Low-Rank Adaptation of Large Language Models. ICLR'22

Model customization for editing

Imagic

Input Image



"A sitting dog"



"A jumping dog"

Edited Images



"A dog lying down"



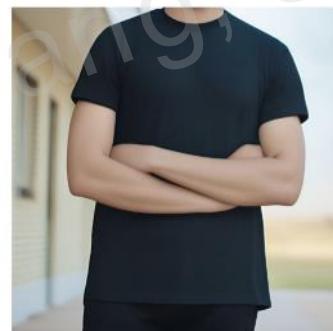
"A dog playing with a toy"



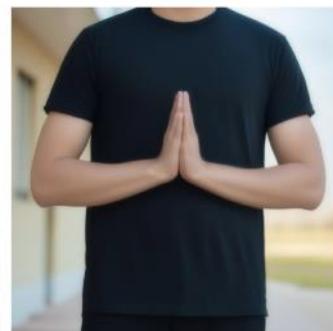
"A jumping dog holding a frisbee"



"A person giving the thumbs up"



"A person with crossed arms"



"A person in a greeting pose to Namaste hands"



"A person holding a cup"

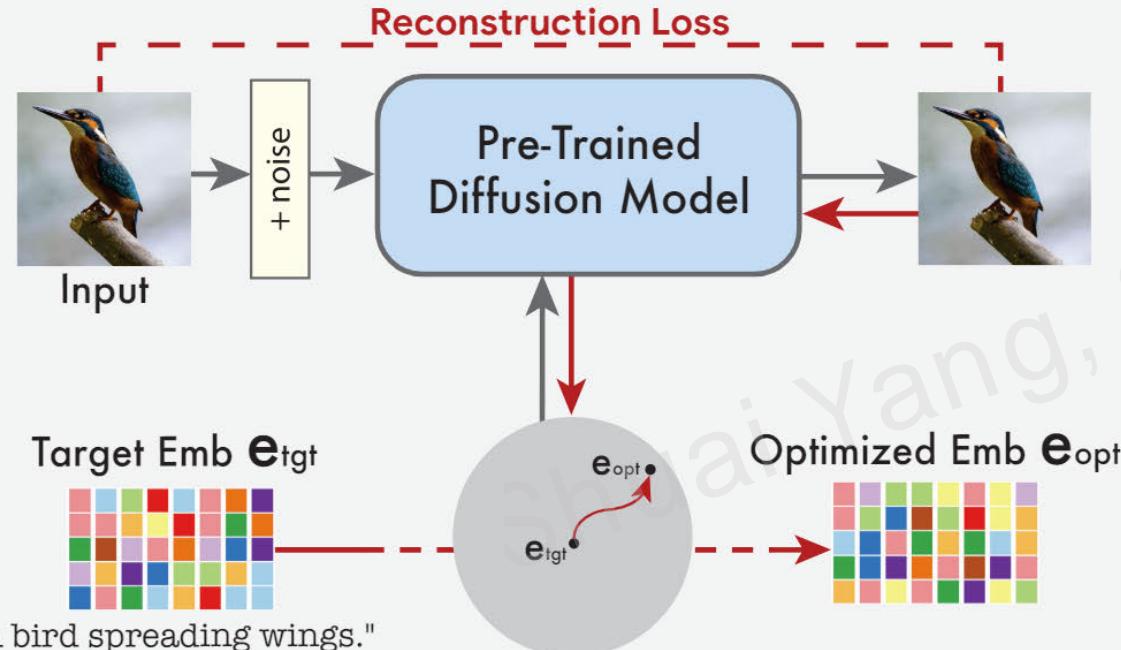


"A person making a heart sign"

Model customization for editing

Imagic

(A) Text Embedding Optimization



(B) Model Fine-Tuning



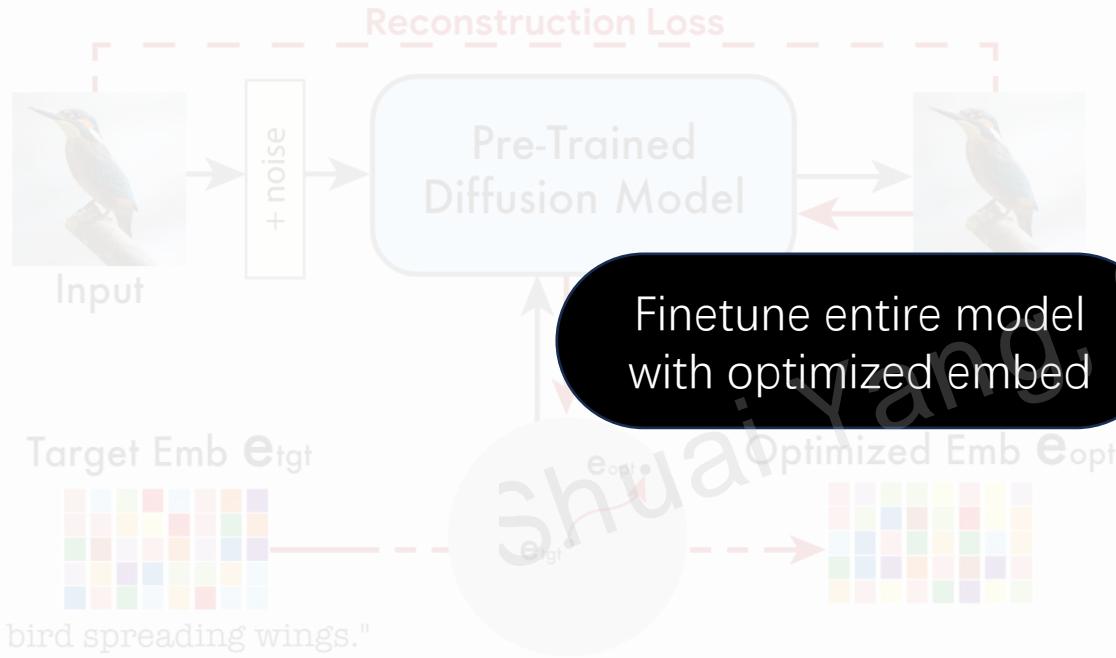
(C) Interpolation & Generation



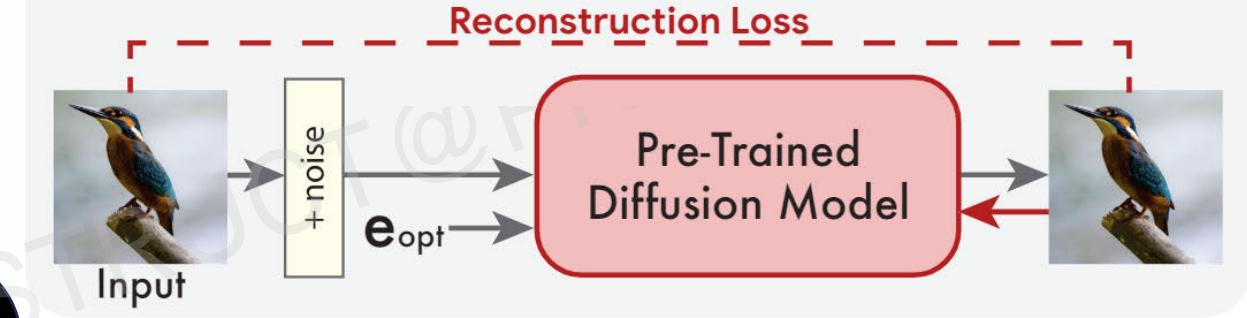
Model customization for editing

Imagic

(A) Text Embedding Optimization



(B) Model Fine-Tuning



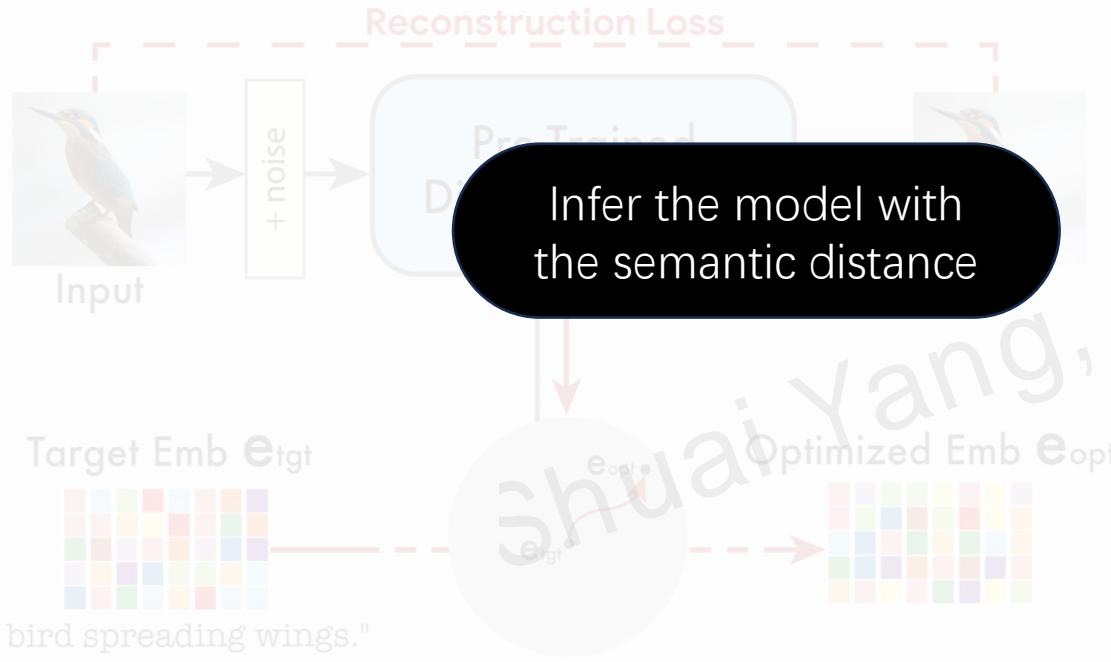
(C) Interpolation & Generation



Model customization for editing

Imagic

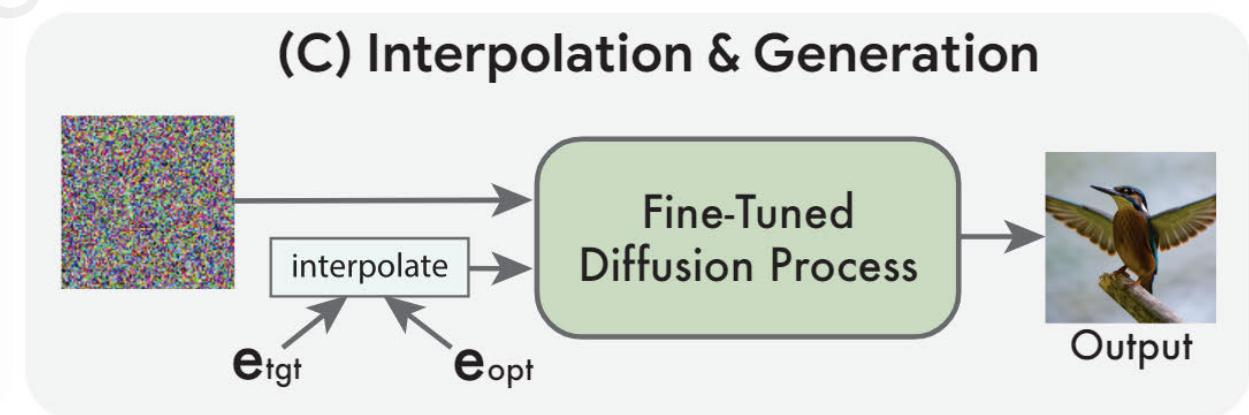
(A) Text Embedding Optimization



(B) Model Fine-Tuning



(C) Interpolation & Generation



Model customization for editing

Imagic



"A cat wearing a hat"



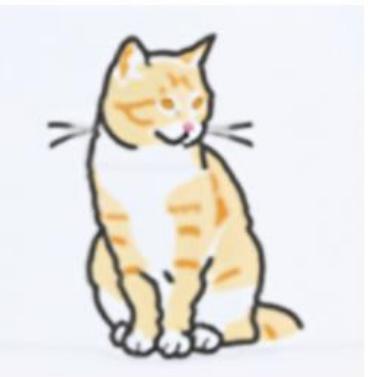
"A cat wearing an apron"



"A cat wearing a necklace"



"A cat wearing a jean jacket"



"A drawing of a cat"



"A zebra"



"A horse with a saddle"



"A horse with its head down"



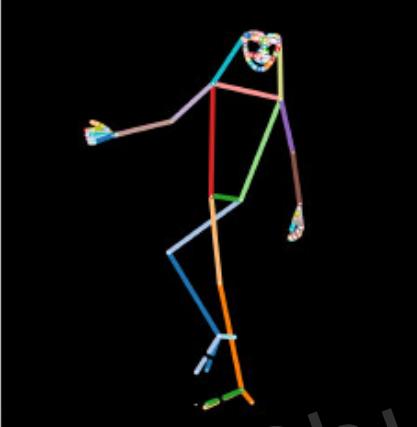
"A brown horse in a grass field"



"A cartoon of a horse"

Controlling pre-trained models

1



ControlNet

ICCV 23

2



T2I Adapter

AAAI 24

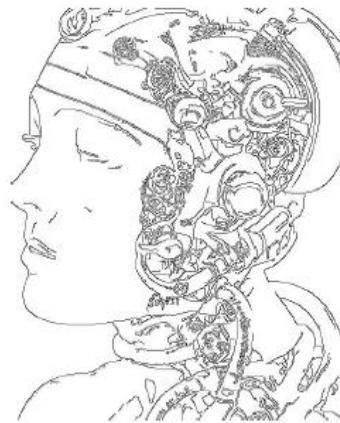
Controlling pre-trained models

1 ControlNet



"a man standing on top of a cliff"

"man on hill watching a meteor, cartoon artwork"

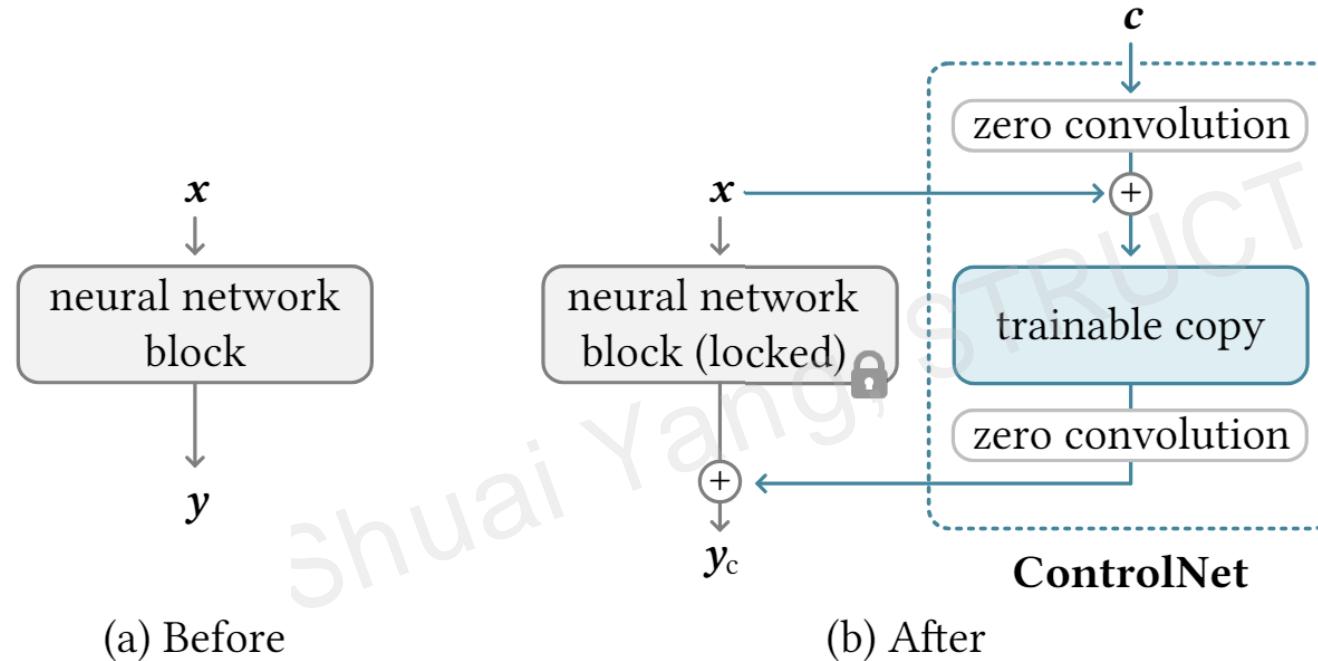


"a robot head with gears"

"robot, cybernetic, cyberpunk, science fiction"

Controlling pre-trained models

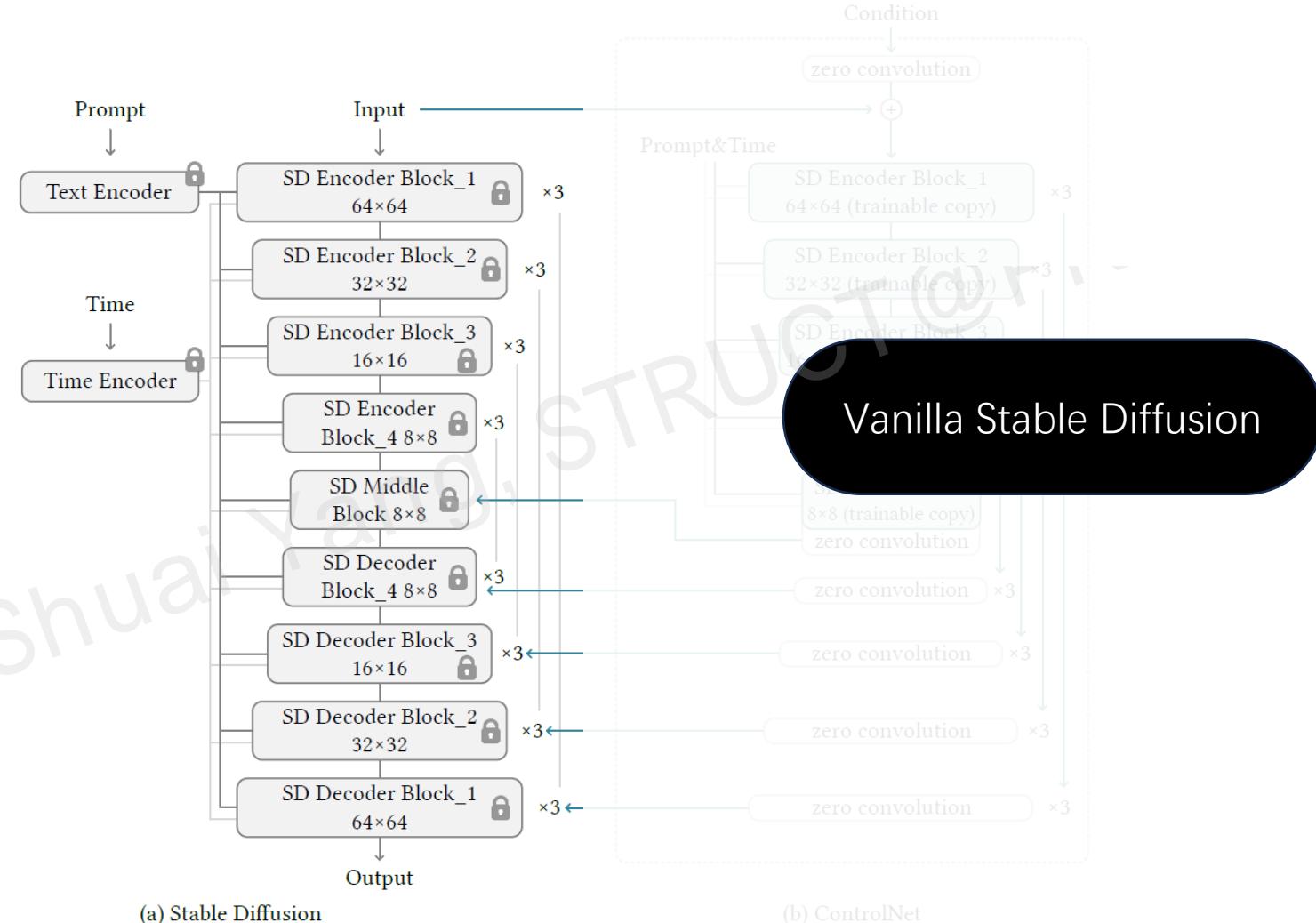
1 ControlNet



Employing zero-conv
to inject guidance

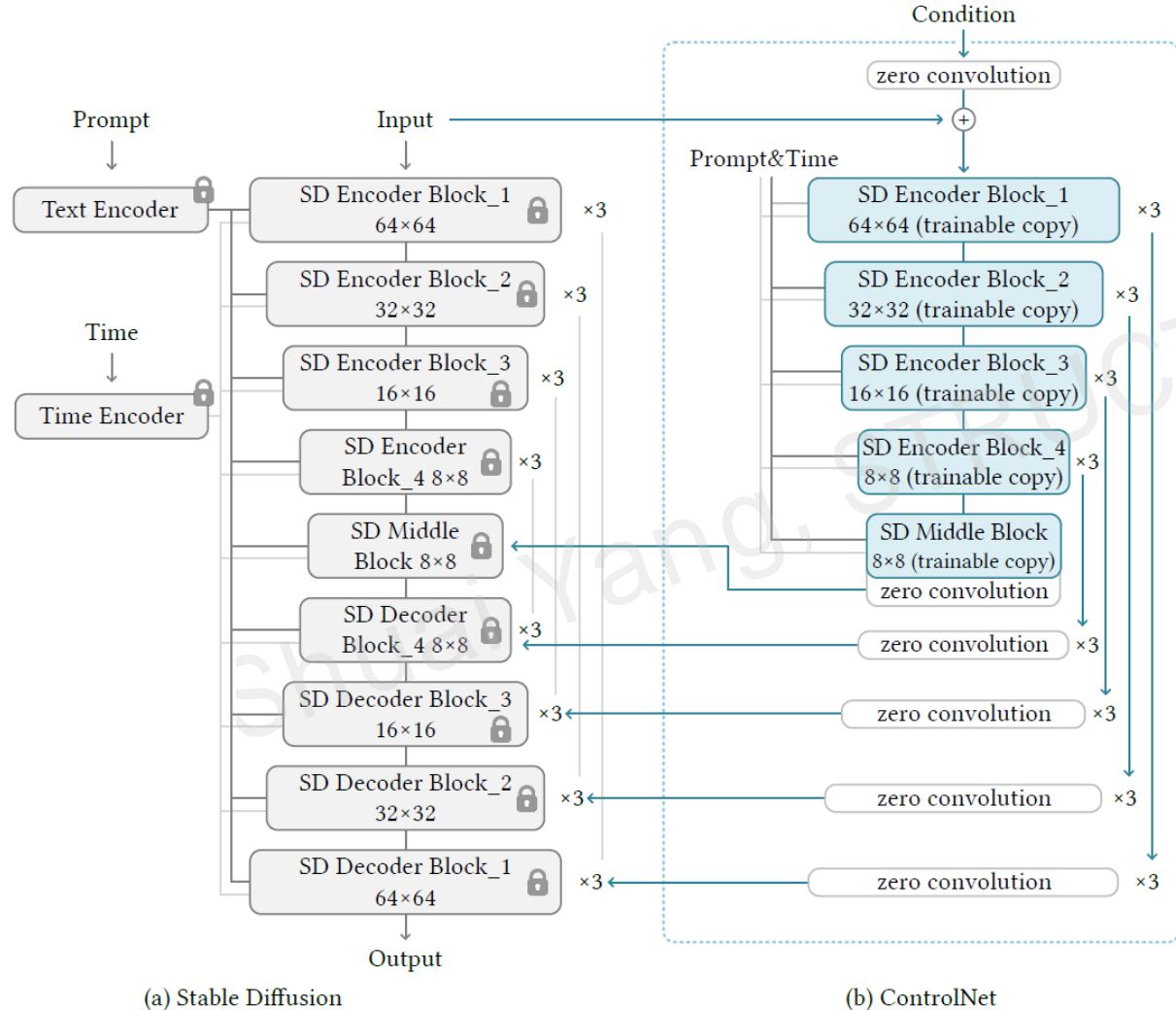
Controlling pre-trained models

1 ControlNet



Controlling pre-trained models

1 ControlNet

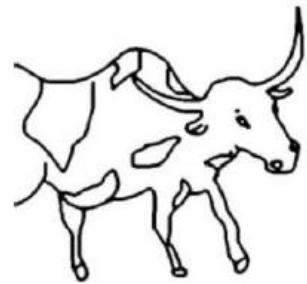
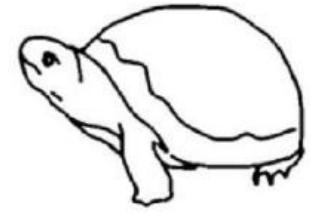


Implement ControlNet
on SD

Controlling pre-trained models

1 ControlNet

Input (User Scribble)

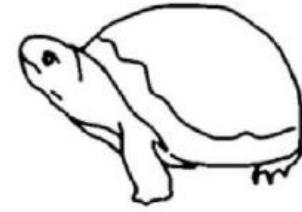


Shuai Yang, STRUCT@rui

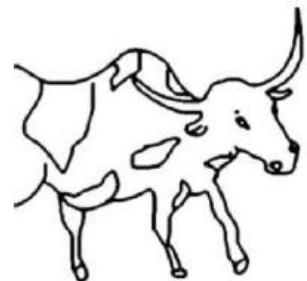
Controlling pre-trained models

1 ControlNet

Input (User Scribble)

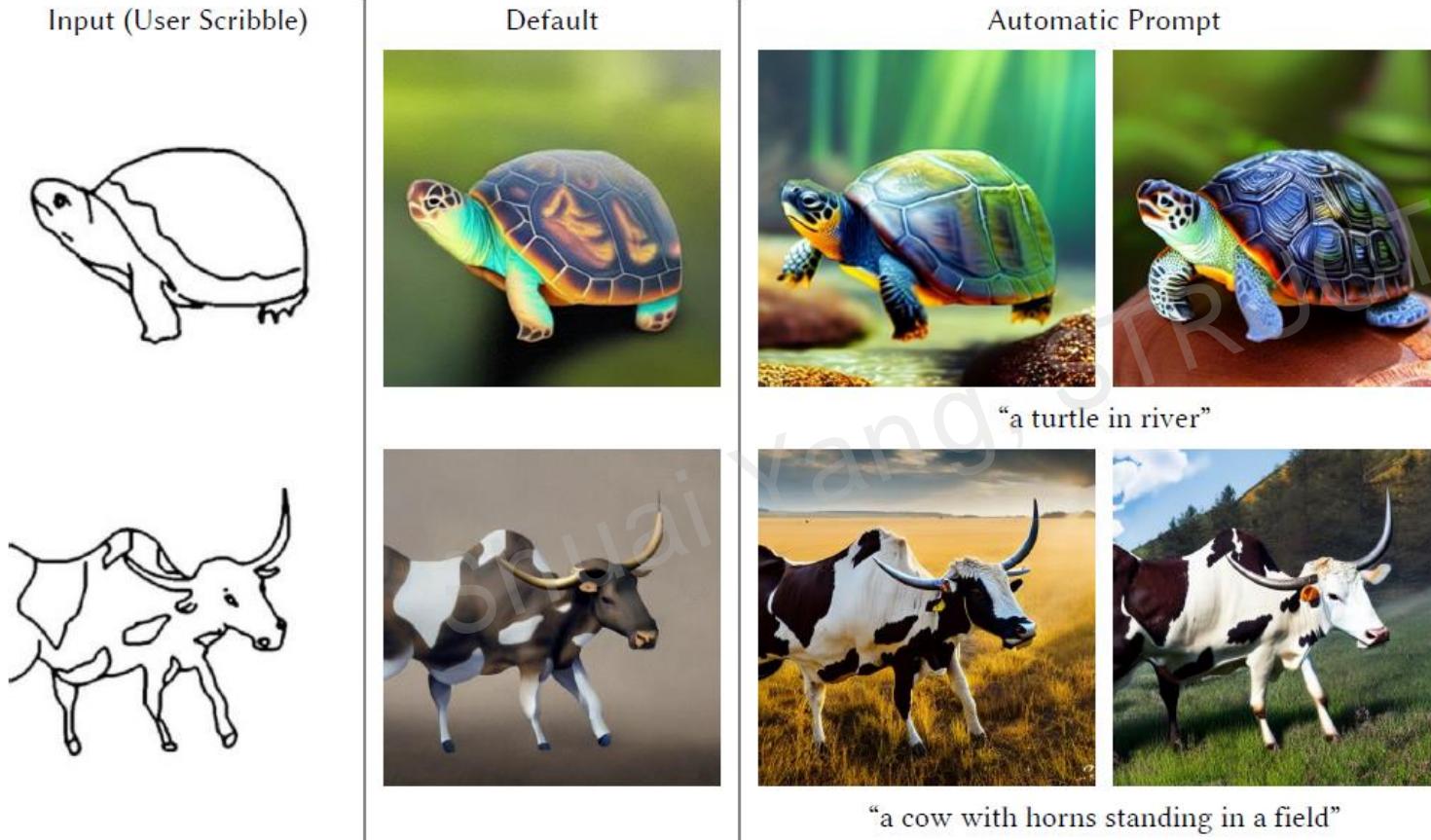


Default



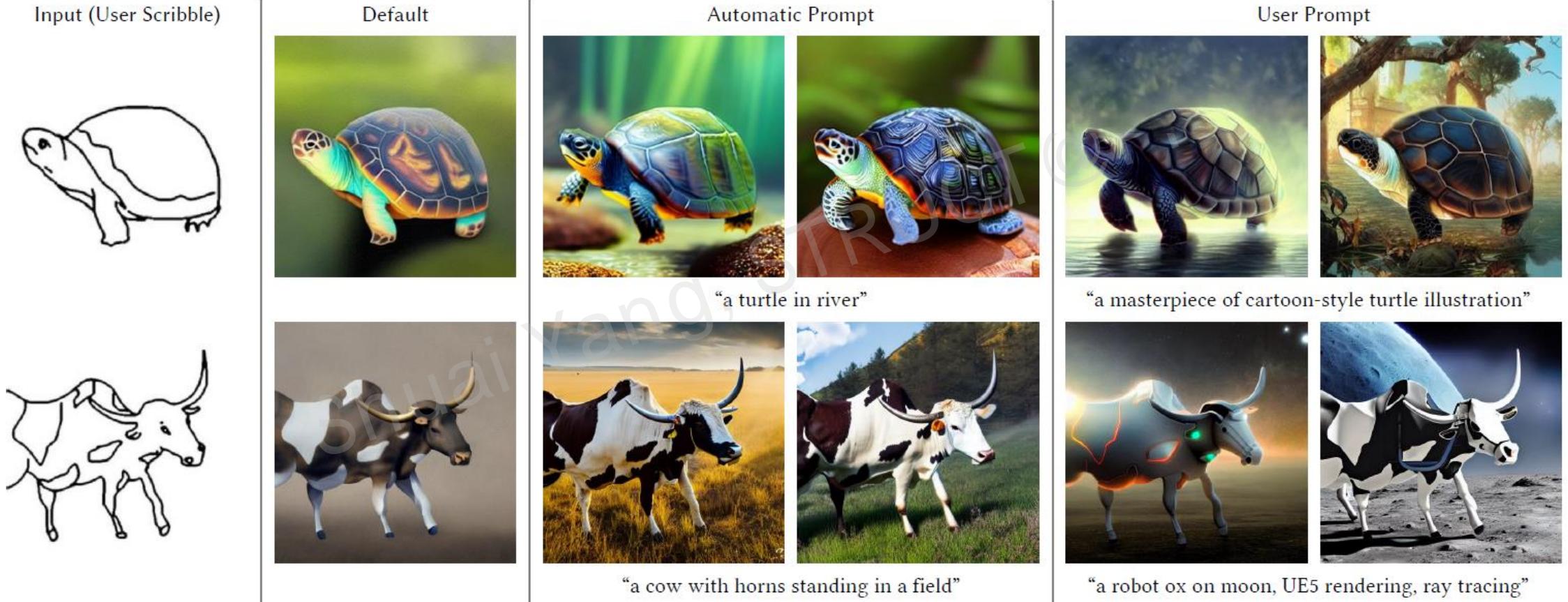
Controlling pre-trained models

1 ControlNet



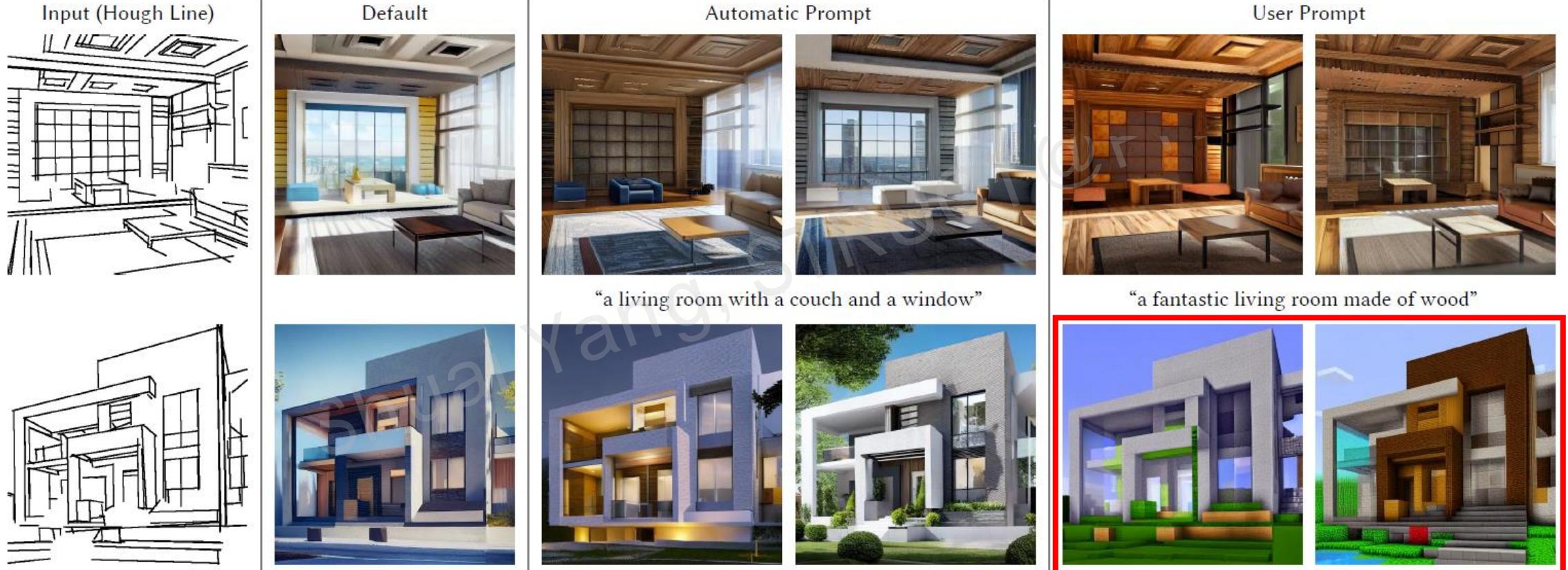
Controlling pre-trained models

1 ControlNet



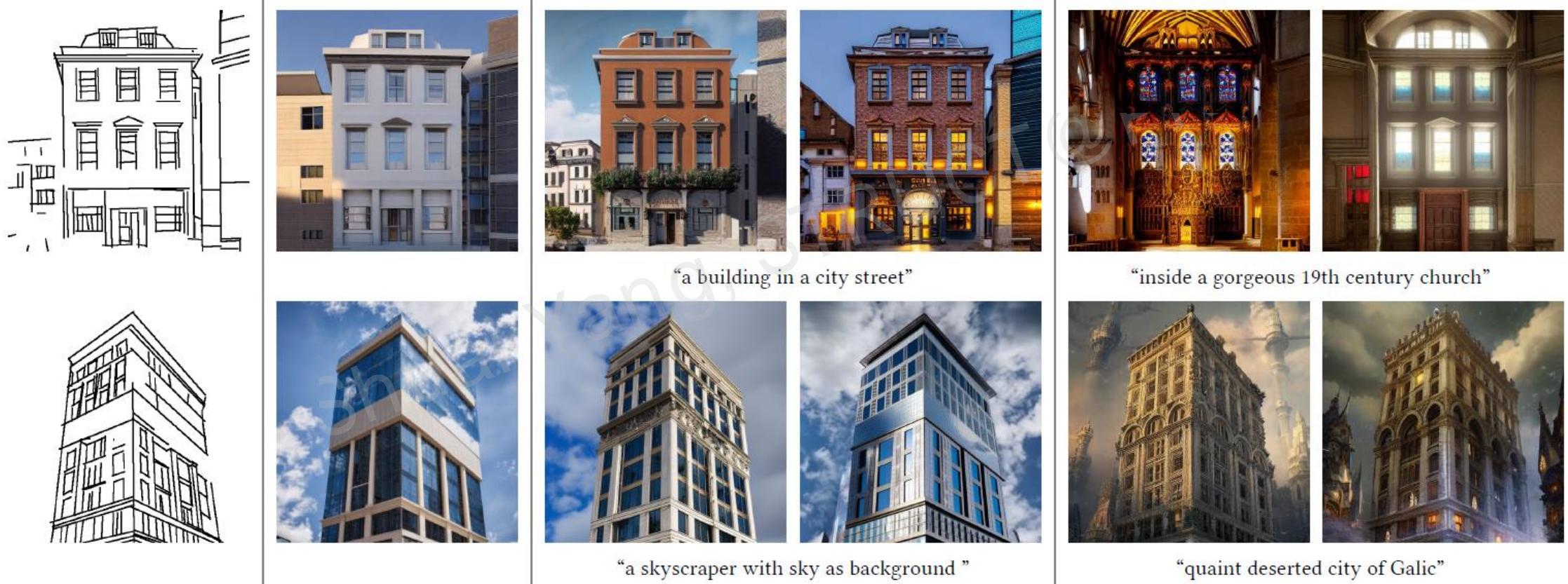
Controlling pre-trained models

1 ControlNet



Controlling pre-trained models

1 ControlNet



Controlling pre-trained models

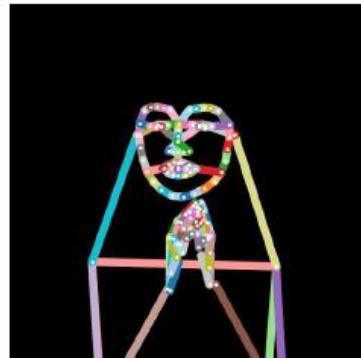
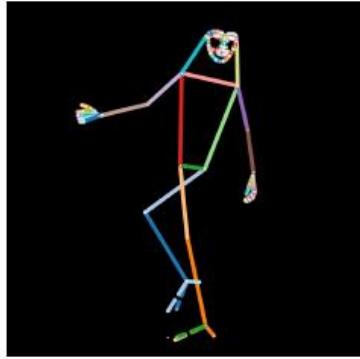
1 ControlNet



Controlling pre-trained models

1 ControlNet

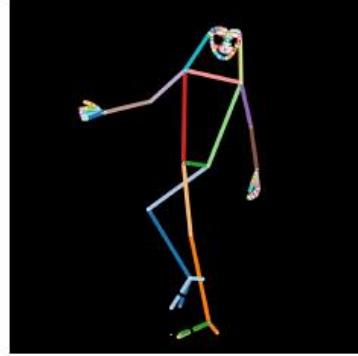
Input (openpifaf)



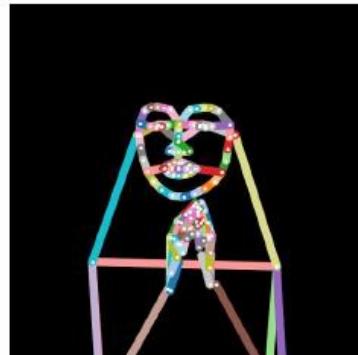
Controlling pre-trained models

1 ControlNet

Input (openpifaf)

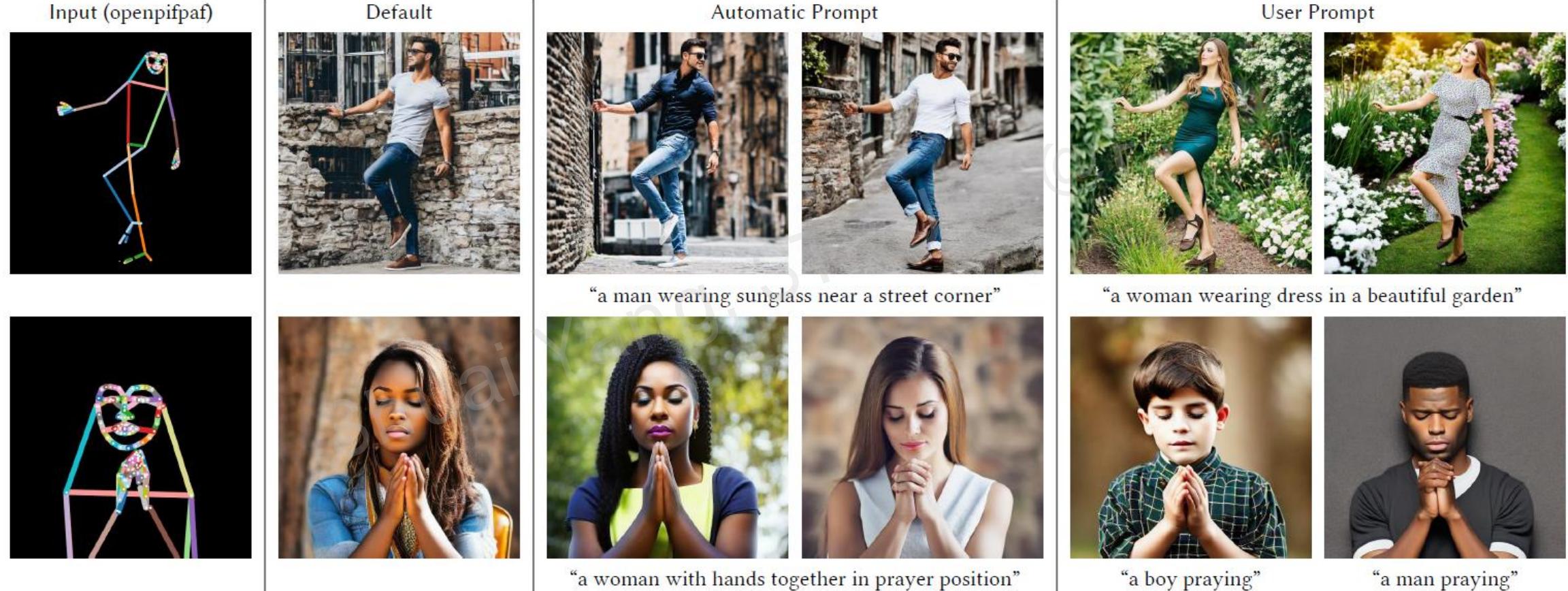


Default



Controlling pre-trained models

1 ControlNet



Controlling pre-trained models

1 ControlNet



Controlling pre-trained models

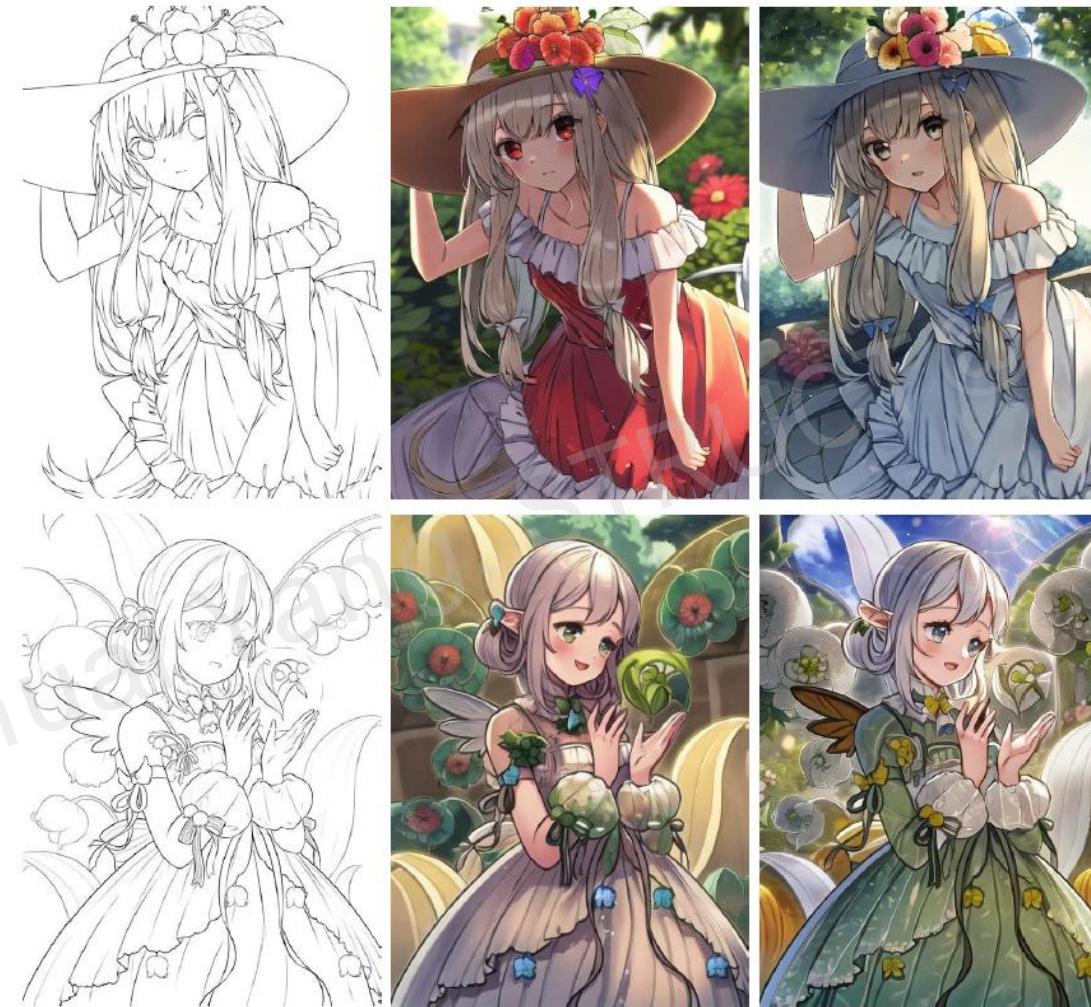
1 ControlNet



Images and Midas Depth

Controlling pre-trained models

1 ControlNet

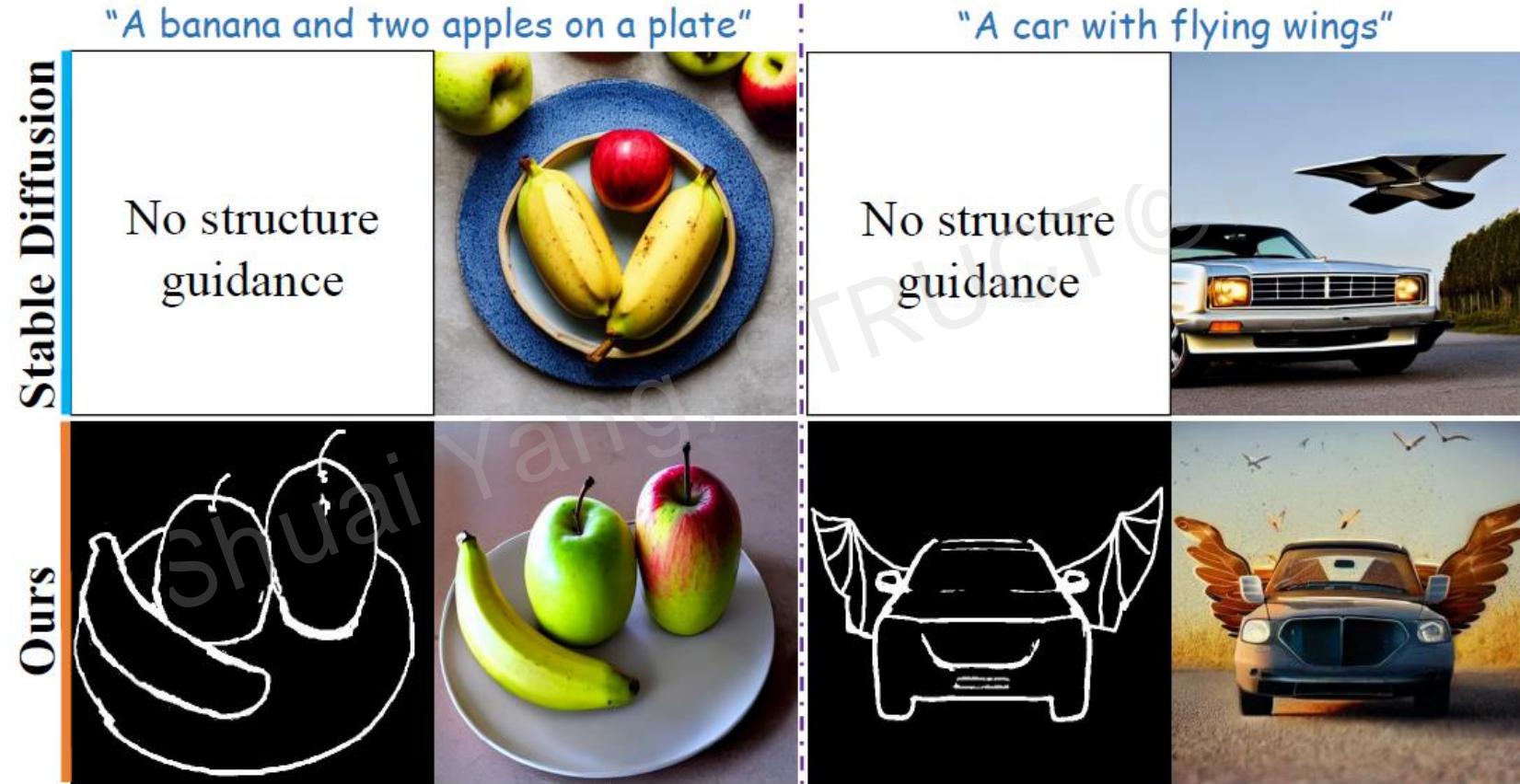


Cartoon line drawing

"1girl, masterpiece, best quality, ultra-detailed, illustration"

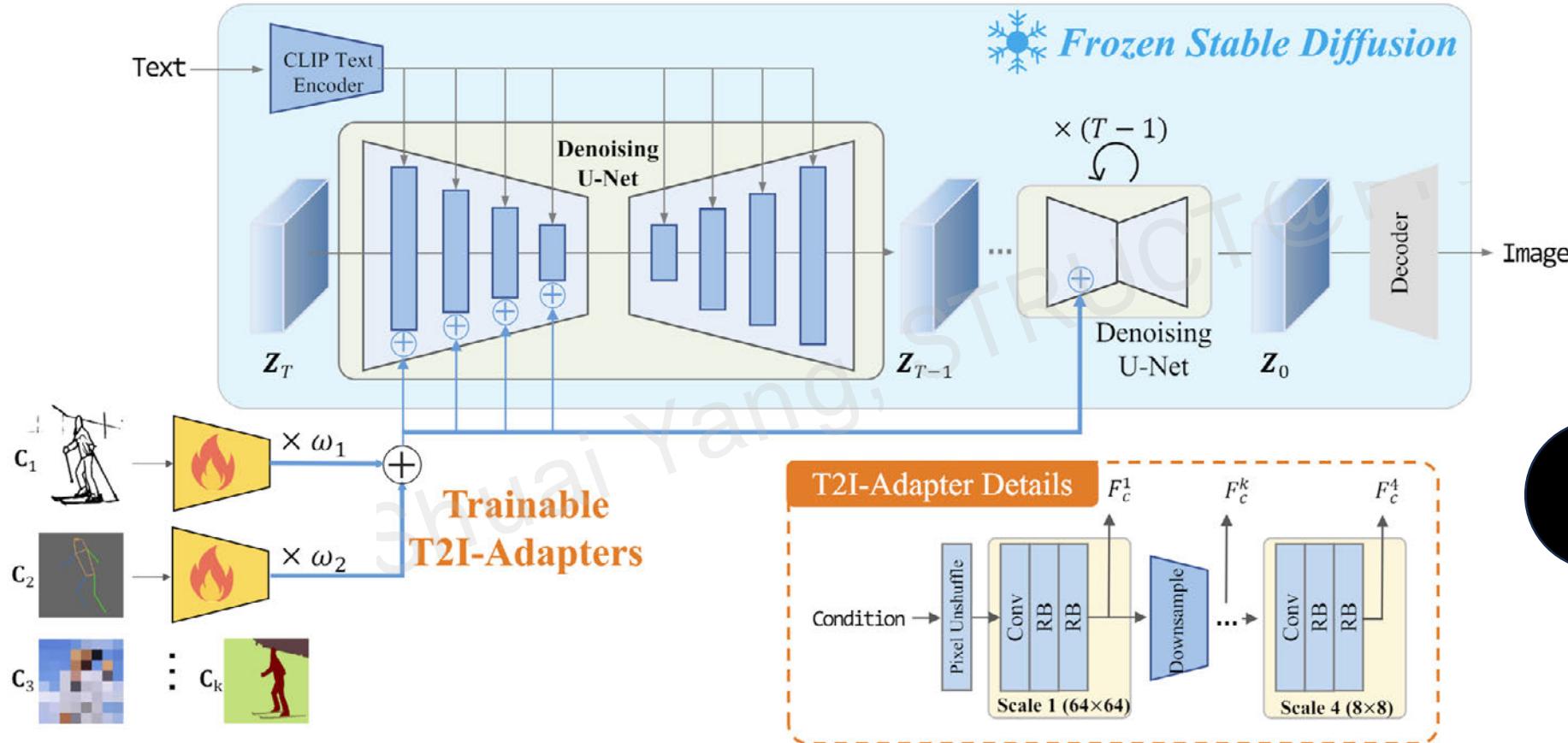
Controlling pre-trained models

2 T2I Adapter



Controlling pre-trained models

2 T2I Adapter



T2I-Adapter: Learning Adapters to Dig Out More Controllable Ability for Text-to-Image Diffusion Models. AAAI'24

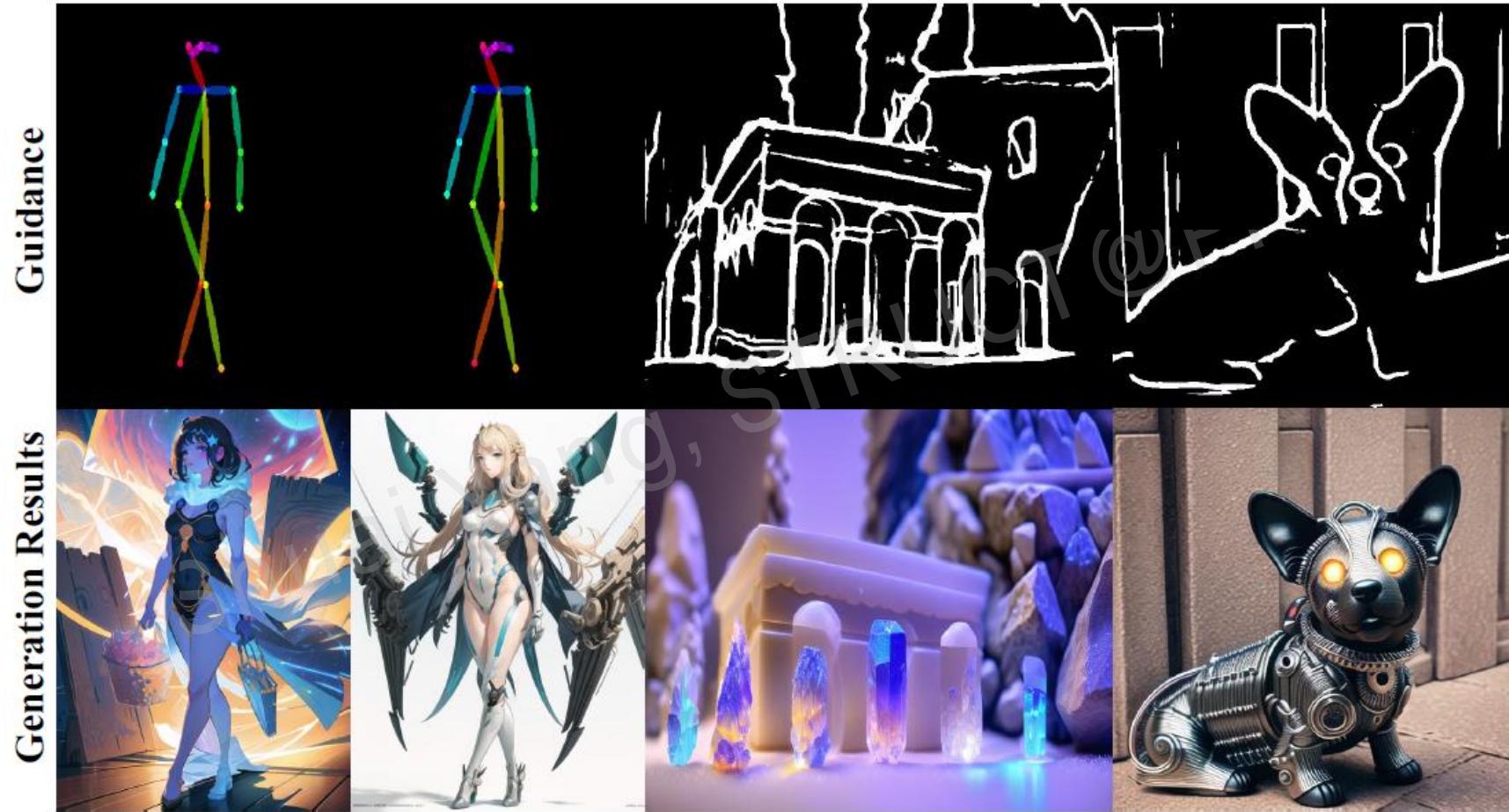
Color + Depth + Sketch + Keypoint

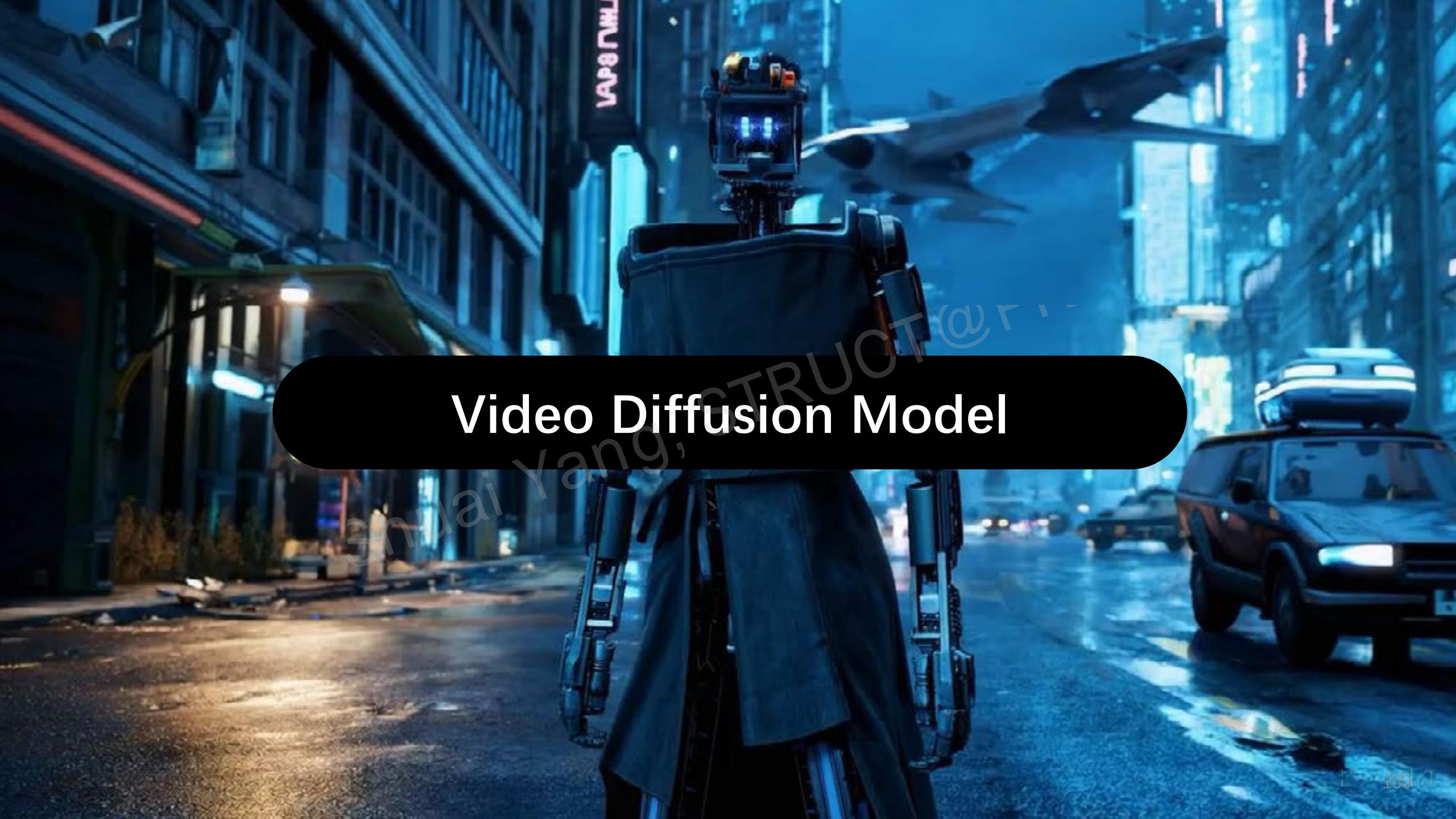


"A cool man and a Pikachu in an office, high quality"

Controlling pre-trained models

2 T2I Adapter



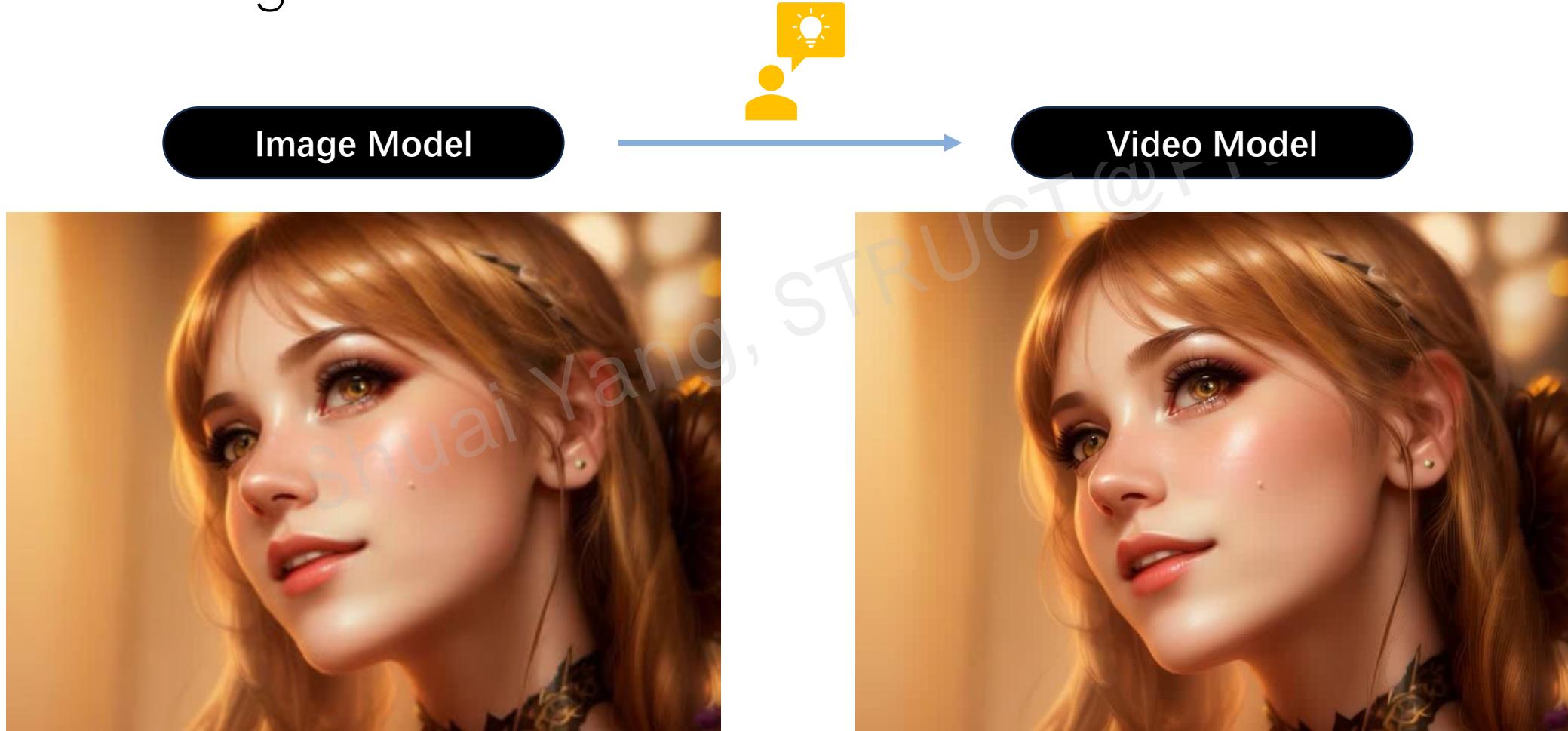
A dark, atmospheric scene set in a futuristic city at night. In the foreground, a large, metallic robot with glowing blue and orange lights on its head and chest stands on a wet, reflective street. The city buildings are tall and have various glowing signs, including one that reads "VAPORWALL". A car is parked on the right side of the street. The overall lighting is low, with most light coming from the city's own neon signs and the robot's own illumination.

Video Diffusion Model

enmai Yang, GTRUCT@UIUC

Video Diffusion Model

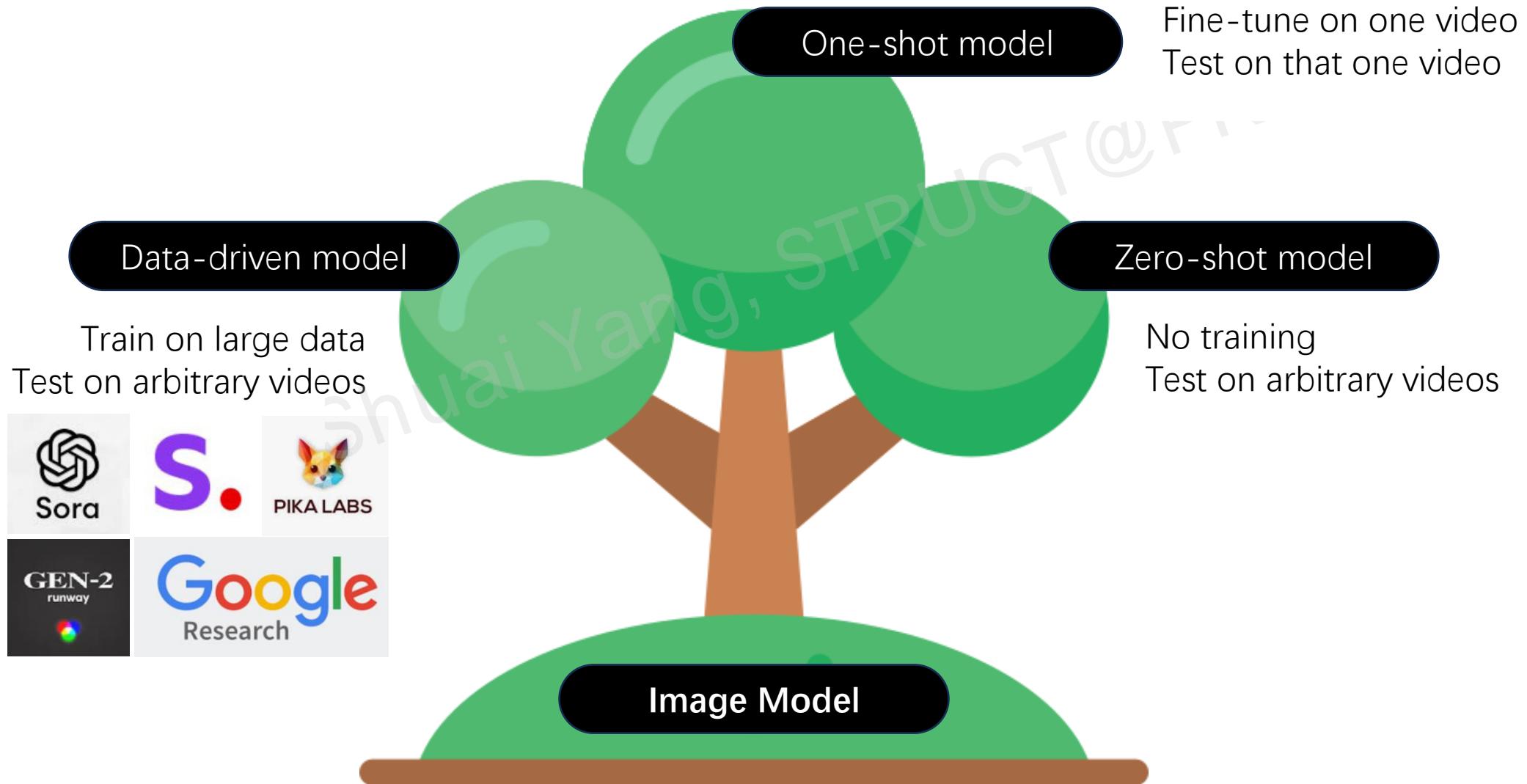
From image to video



A beautiful woman in CG style

Video Diffusion Model

From image to video



Video Diffusion Model

From image to video



Data-driven model

- Good performance
- Need large video data and high computational resources
- Good models are not open sourced



One-shot model

- Less training resource required
- Good compatibility
- Overfit to a single video
- Less efficient during testing



Zero-shot model

- No training resource required
- High compatible with image model
- Lack the knowledge of real-world motion

Video Diffusion Model

From image to video



Data-driven model

- Good performance
- Need large video data and high computational resources
- Good models are not open sourced



One-shot model

- Less training resource required
- Good compatibility
- Overfit to a single video
- Less efficient during testing



Zero-shot model

- No training resource required
- High compatible with image model
- Lack the knowledge of real-world motion

Video Diffusion Model

From image to video



Data-driven model

- Good performance
- Need large video data and high computational resources
- Good models are not open sourced



One-shot model

- Less training resource required
- Good compatibility
- Overfit to a single video
- Less efficient during testing



Zero-shot model

- No training resource required
- High compatible with image model
- Lack the knowledge of real-world motion

Video Diffusion Model

From image to video



Data-driven model

TASK

- Text-to-Video Generation
- Image-to-Video Generation
- Video Editing



One-shot model

TASK

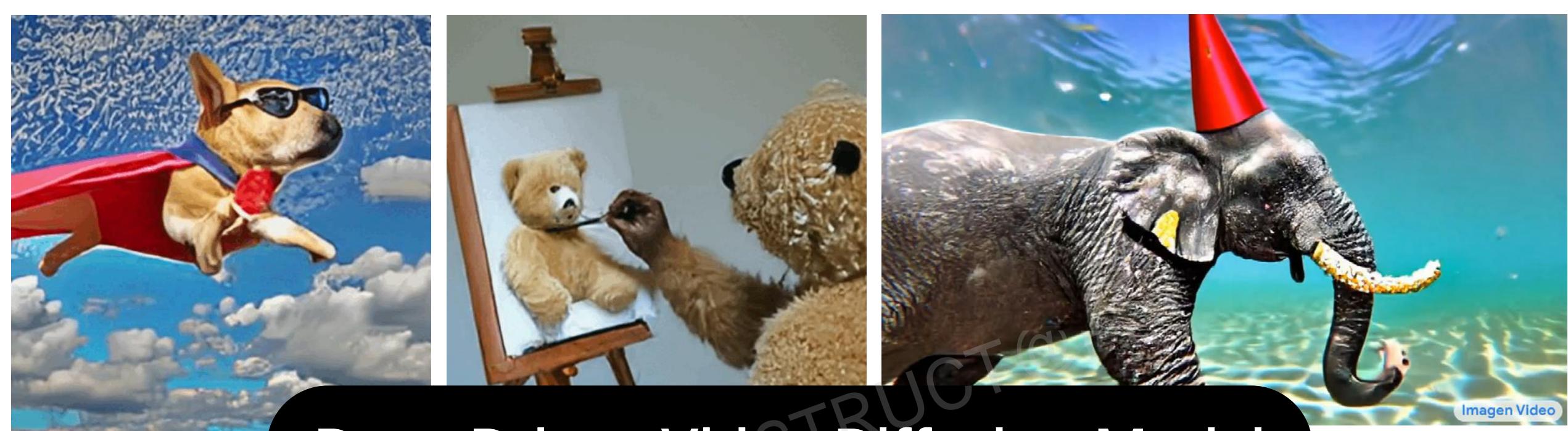
- Video Editing



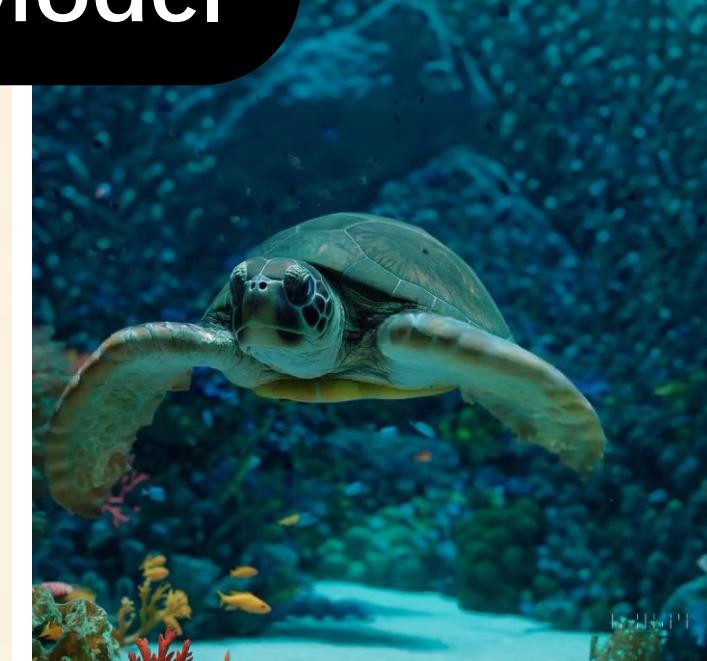
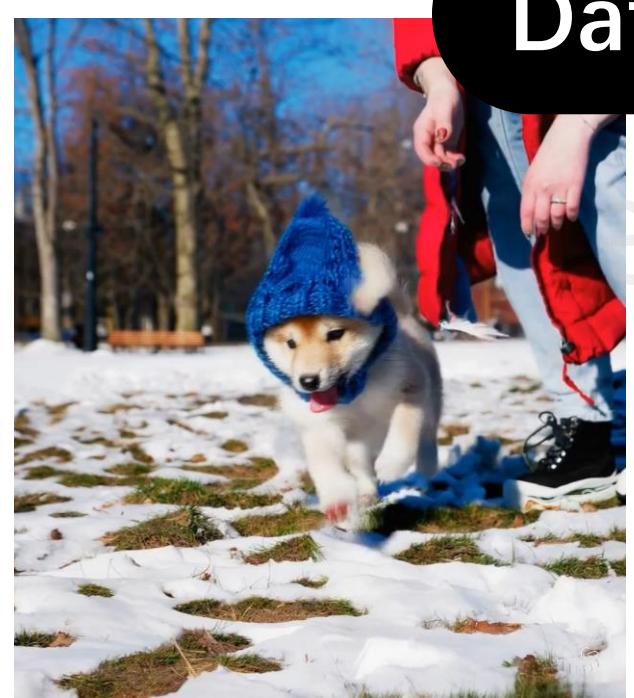
Zero-shot model

TASK

- Text-to-Video Generation
- Video Editing



Data-Driven Video Diffusion Model



Data-driven model

Non-Diffusion Era

Diffusion Era

- GAN
- Auto-regressive transformers
 - **CogVideo**



CogVideo: Large-scale Pretraining for Text-to-Video Generation via Transformers. ICLR'23

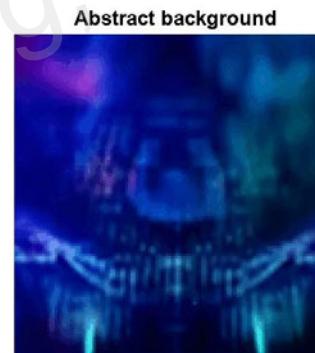
Data-driven model

1 VDM

Non-Diffusion Era

Diffusion Era

- GAN
 - Auto-regressive transformers
- Early exploration — VDM
 - Architecture: factorized space-time U-Net
 - Training: joint images and videos



Abstract background



Clouds moving



path in a tropical forest



Construction Site Activity

Data-driven model

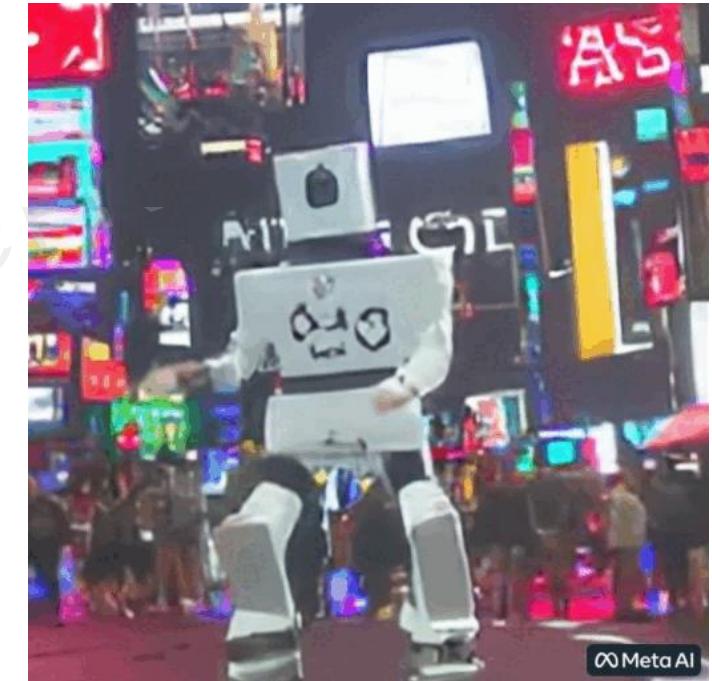
2 Make-A-Video



A dog wearing a Superhero outfit with red cape flying through the sky



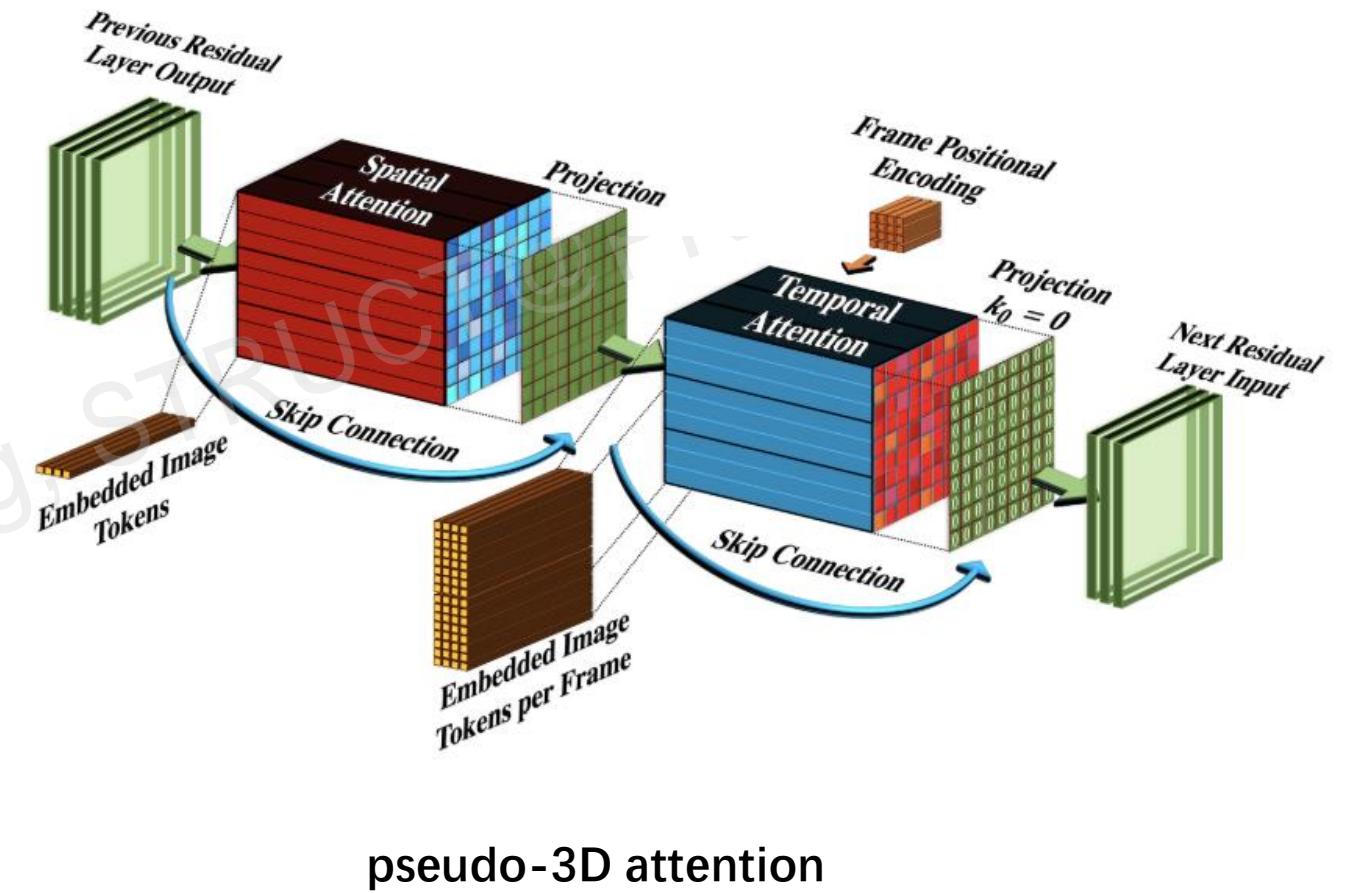
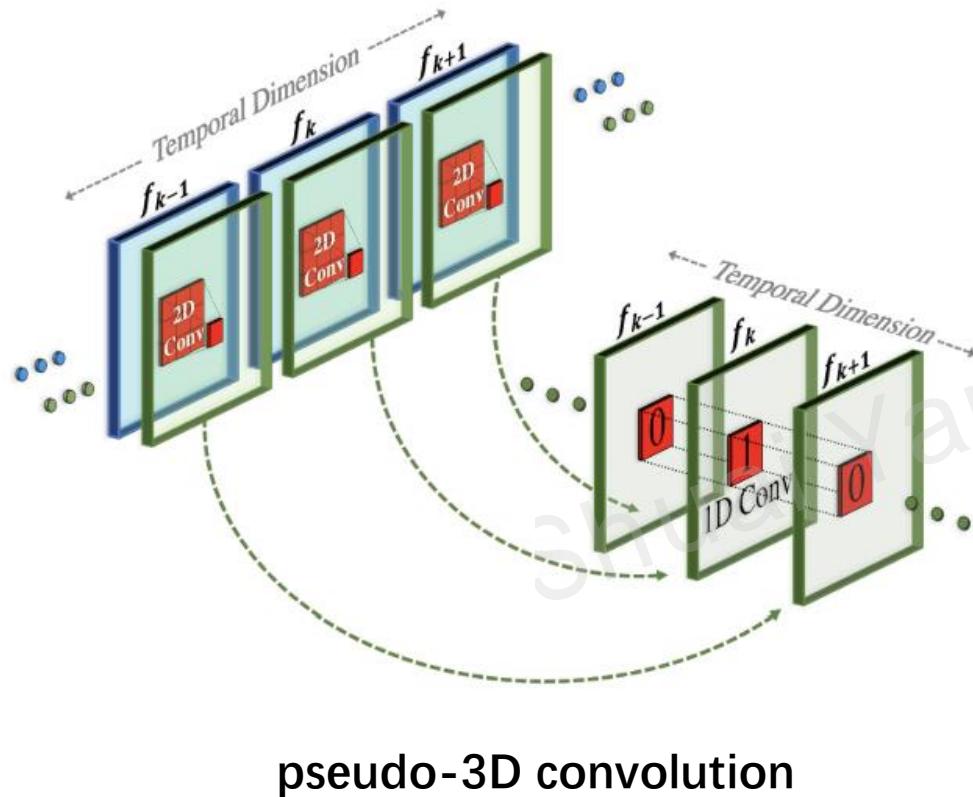
A teddy bear painting a portrait



Robot dancing in times square

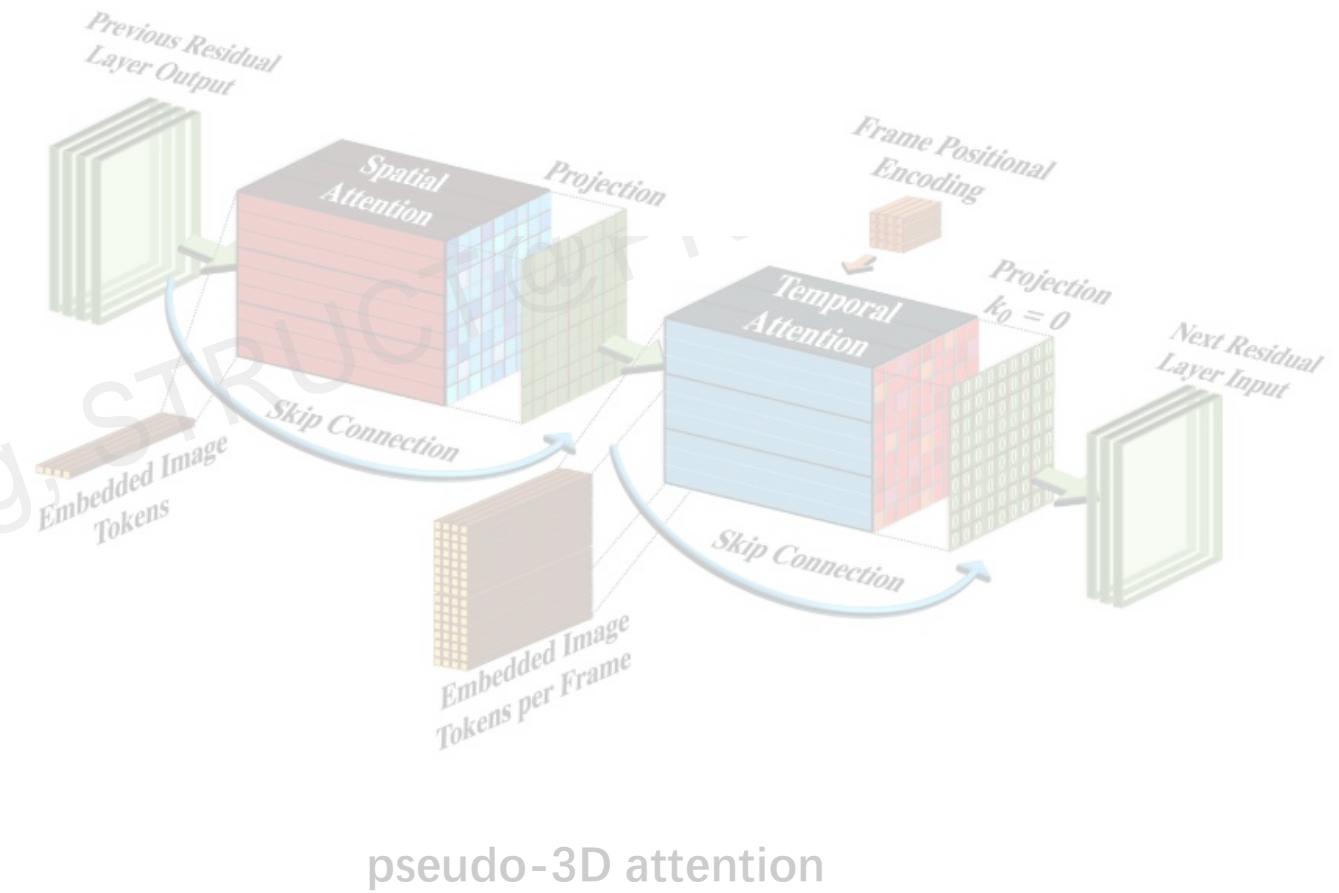
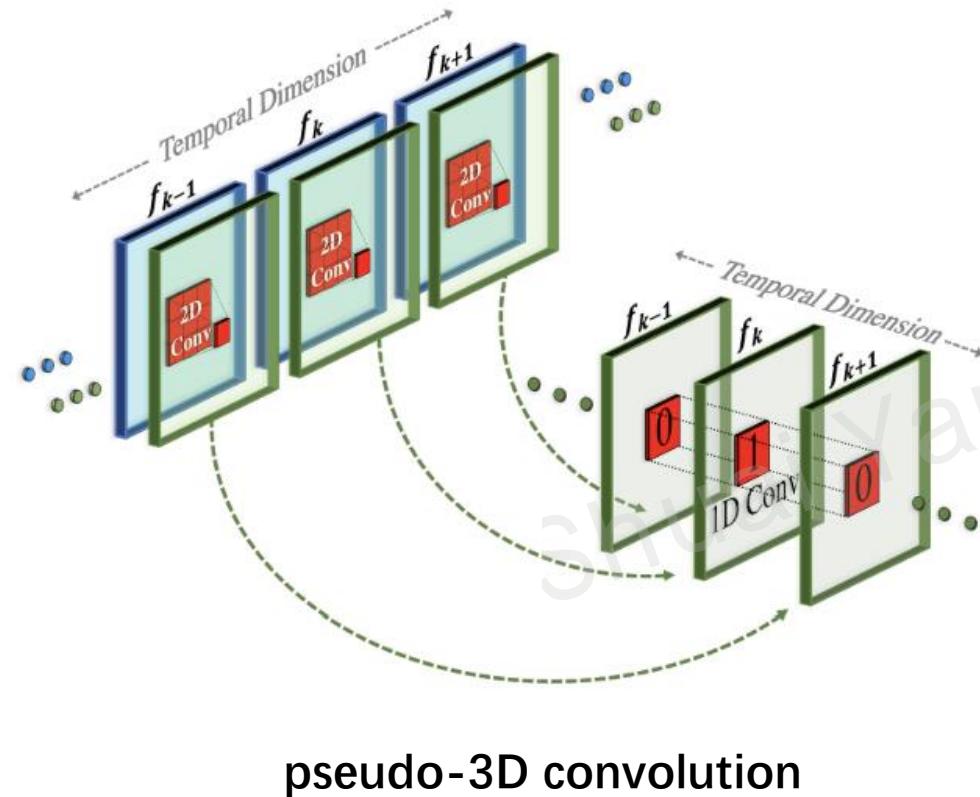
Data-driven model

2 Make-A-Video



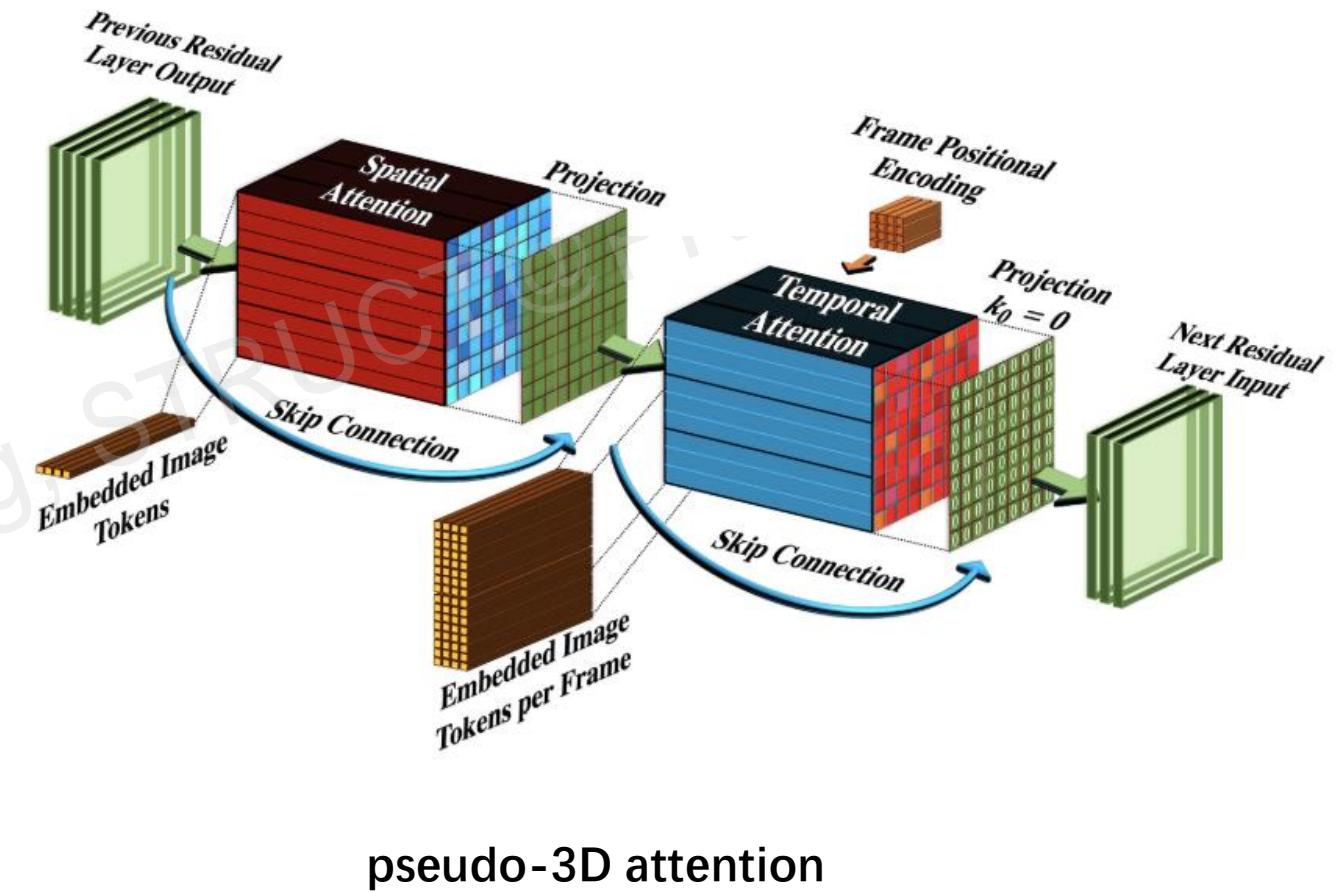
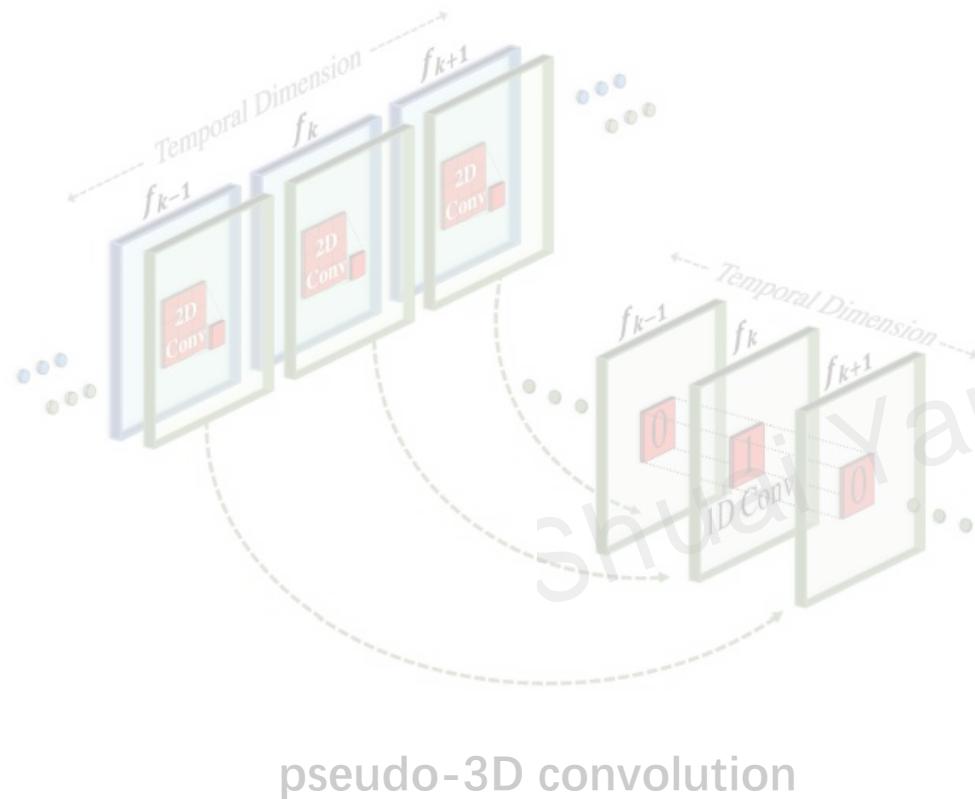
Data-driven model

2 Make-A-Video



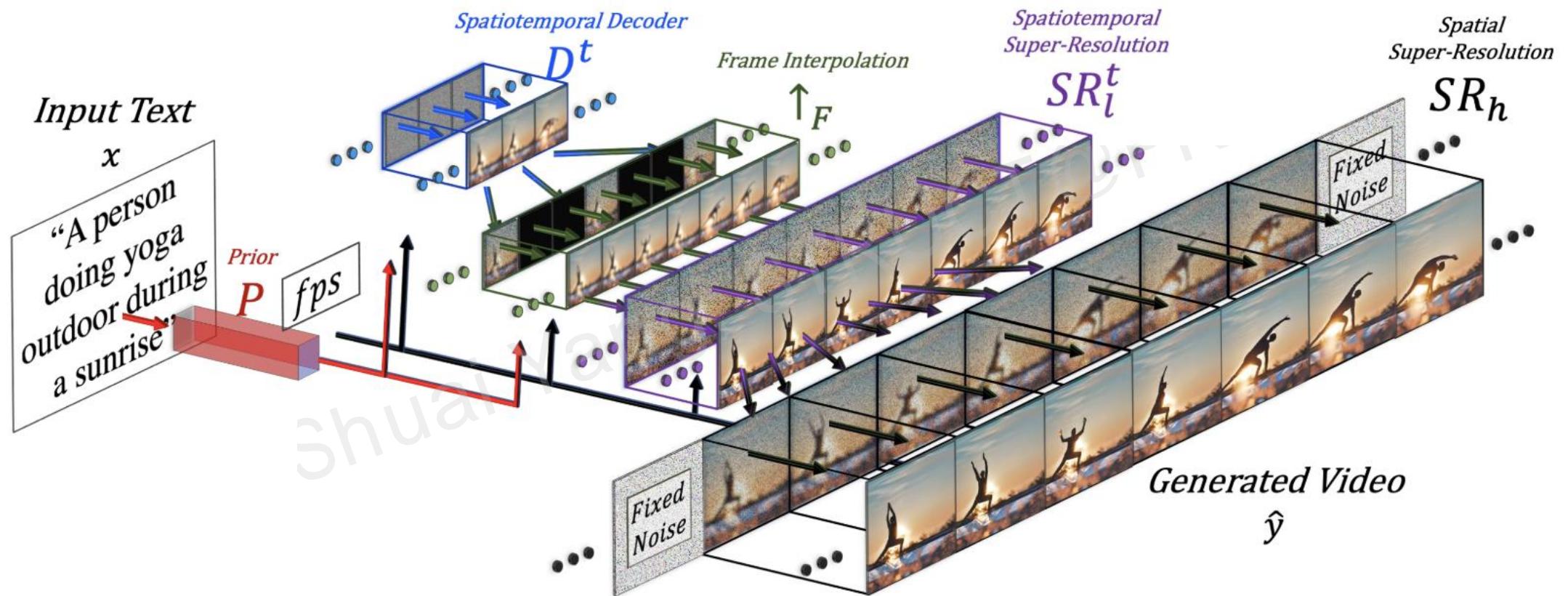
Data-driven model

2 Make-A-Video



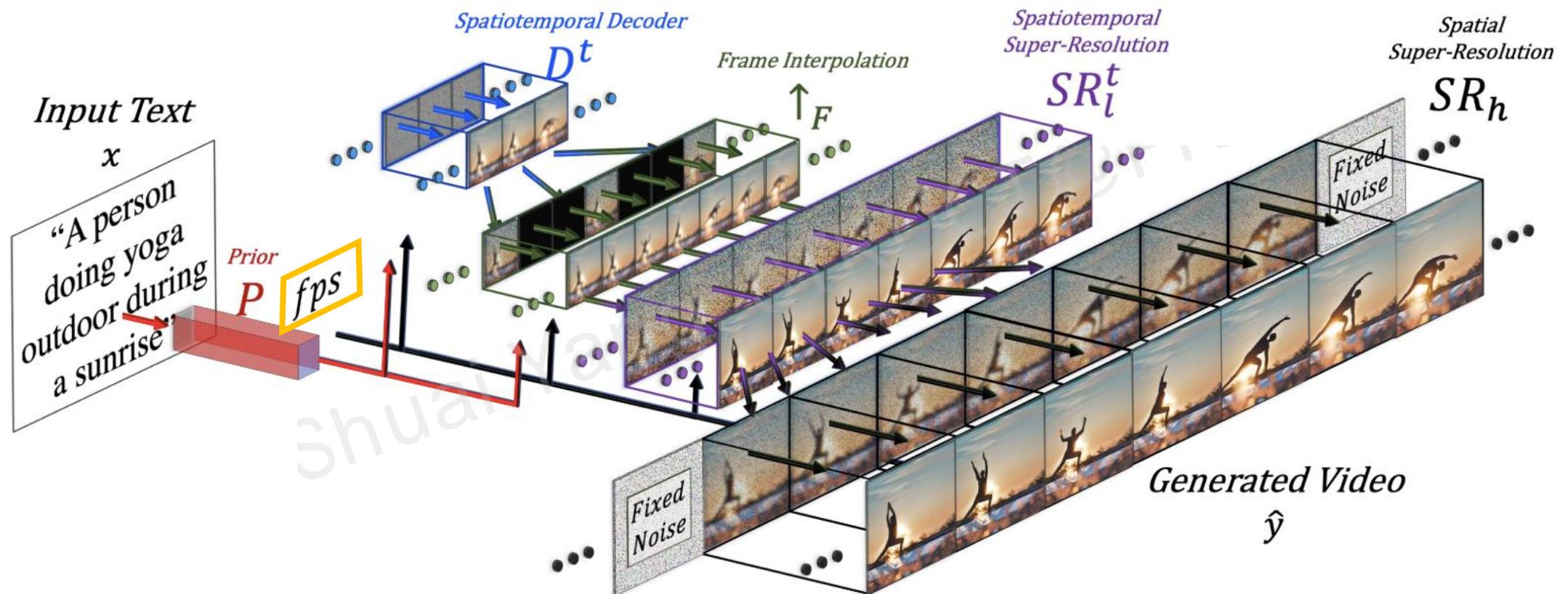
Data-driven model

2 Make-A-Video



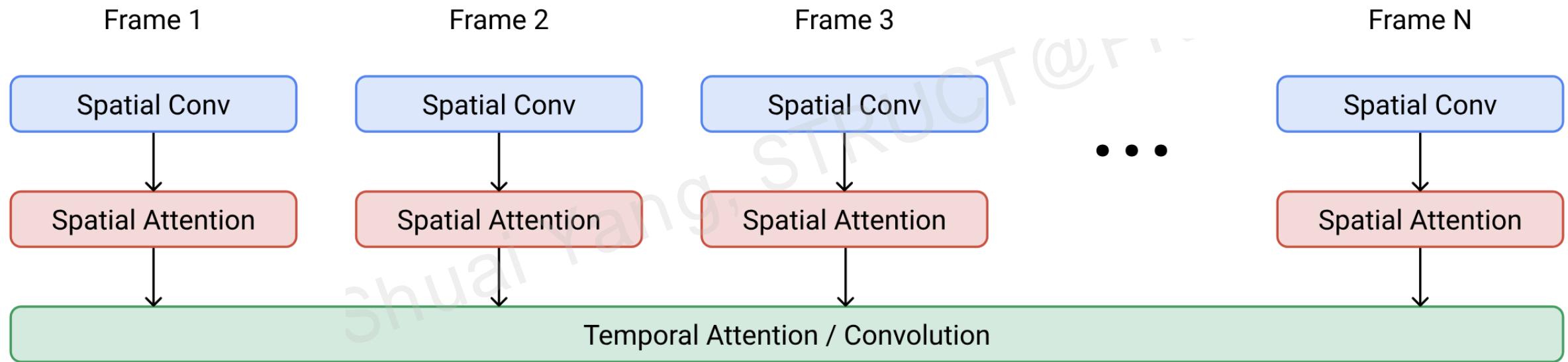
Data-driven model

2 Make-A-Video



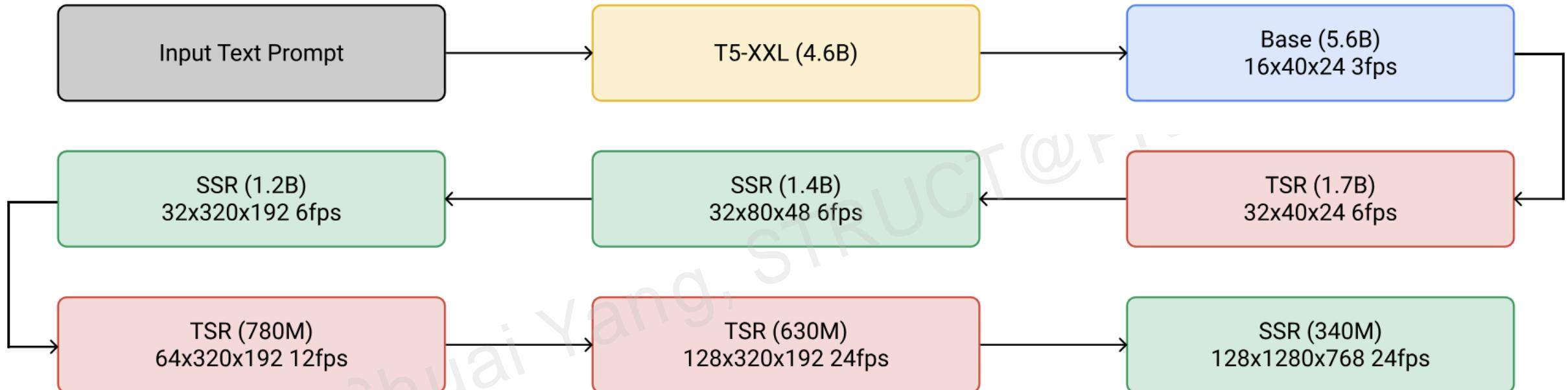
Data-driven model

3 Imagen Video



Data-driven model

3 Imagen Video



More than 11B parameters to learn

Data-driven model

3 Imagen Video



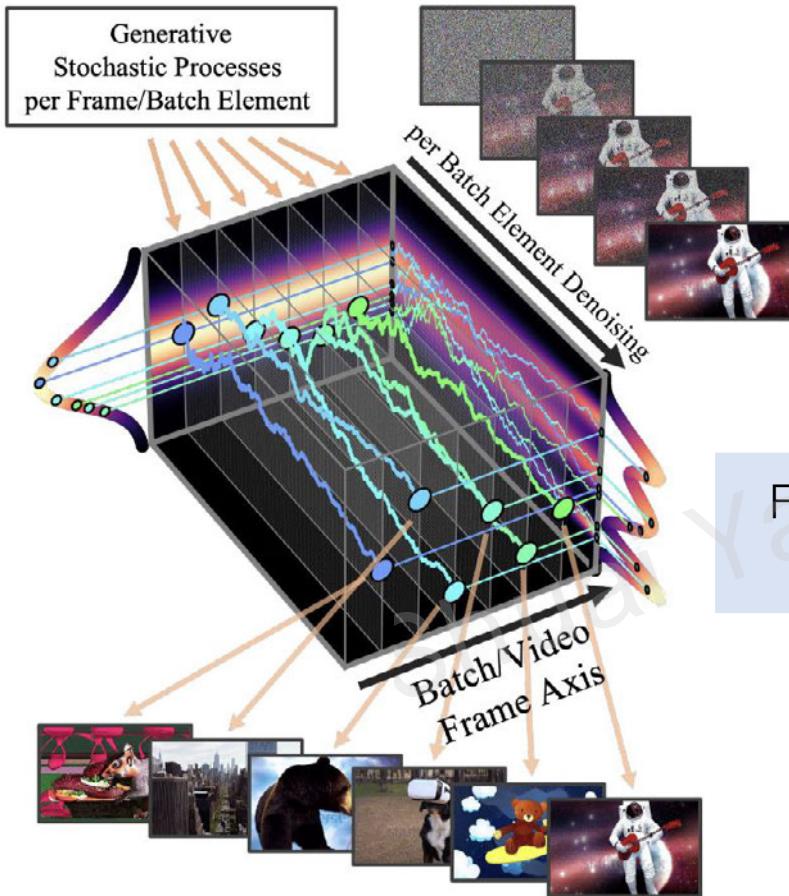
A happy elephant wearing a birthday hat walking under the sea



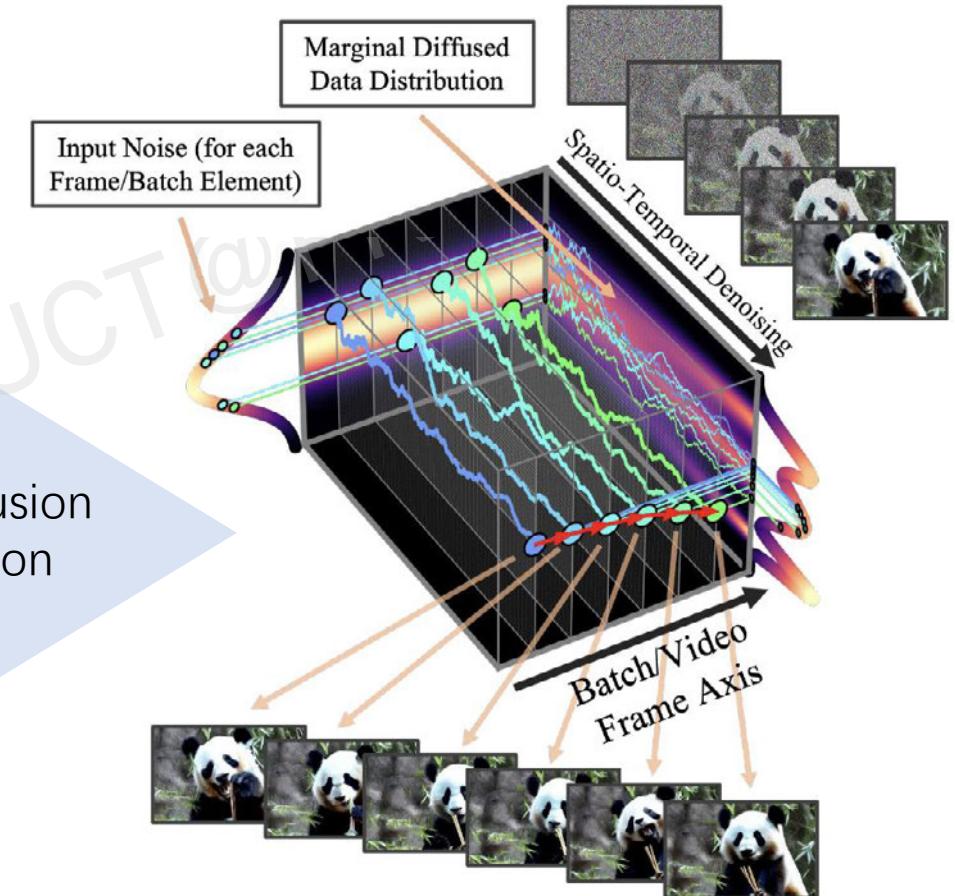
Wooden figurine surfing on a surfboard in space

Data-driven model

4 Video LDM

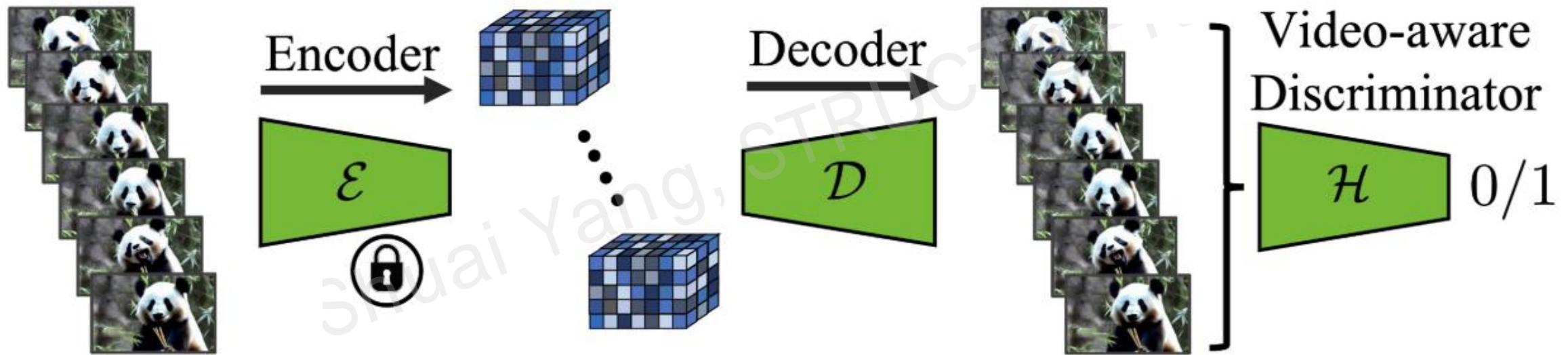


From *image* latent diffusion
to *video* latent diffusion



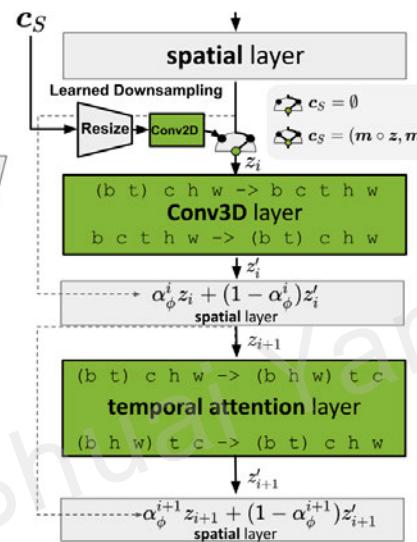
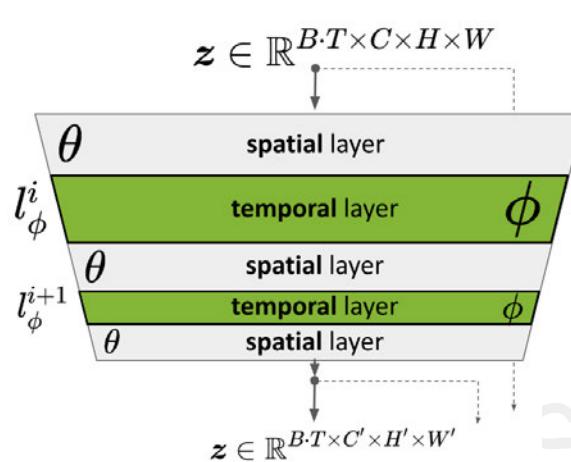
Data-driven model

4 Video LDM

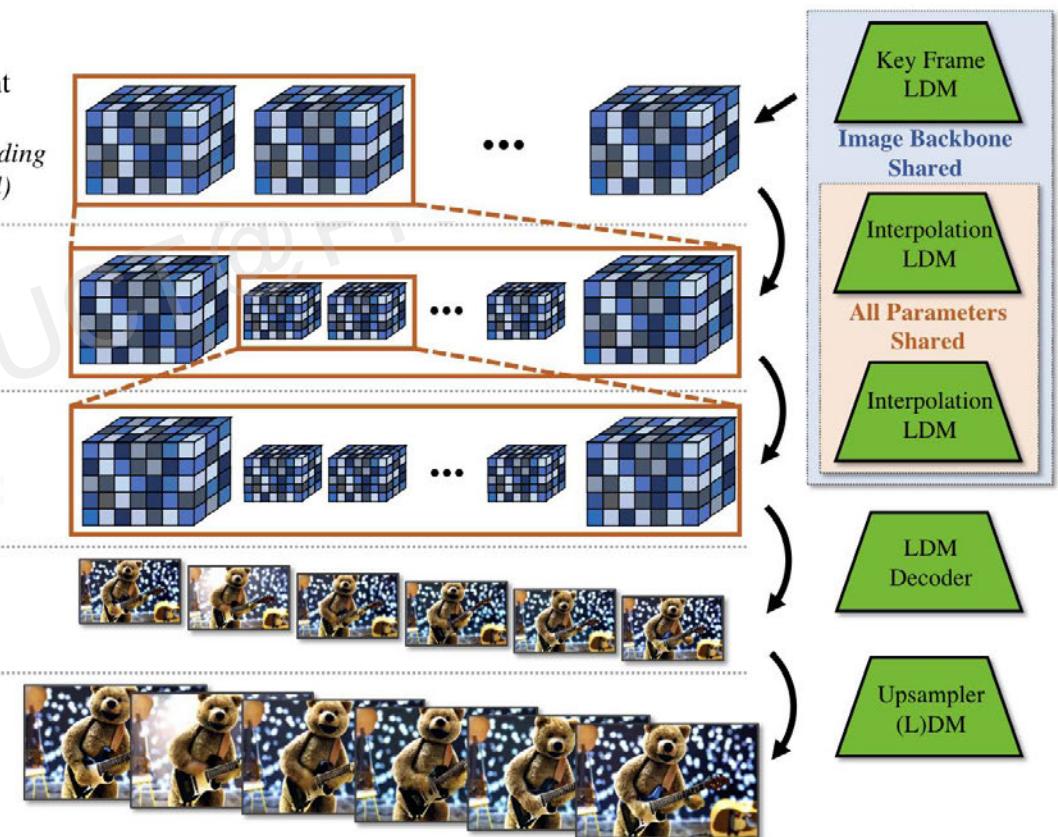


Data-driven model

4 Video LDM



1. Generate Latent Key Frames
(optionally including prediction model)
2. Latent Frame Interpolation I
3. Latent Frame Interpolation II
4. Decode to Pixel-Space
5. Apply Video Upsampler



Data-driven model

4 Video LDM

less than 3B parameters to learn



A storm trooper vacuuming the beach



A koala bear playing piano in the forest

Data-driven model

4 Video LDM



Training images for DreamBooth.



Text prompt: "A sks frog playing a guitar in a band."

Data-driven model

Problem: Lack of high-quality data

An example from the WebVid-10M dataset



Results of models trained on WebVid-10M

ModelScope T2V



LVDM



Man walking in deep snow under the branches

Data-driven model

Ideal text-to-video training data

- **Video Contents**

Clean and dynamic

- **Text Annotation**

Corresponding text description

Example:



A person is holding a long haired dachshund in their arms.

RAW internet videos

- **Video Contents**

✗ Multiple scenes

✗ Unsuitable visual contents

- **Text Annotation**

✗ No corresponding text

Example:



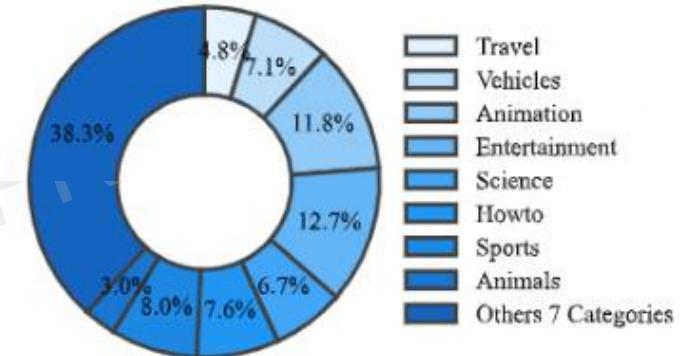
Data Collection

- Open-domain videos collected from YouTube

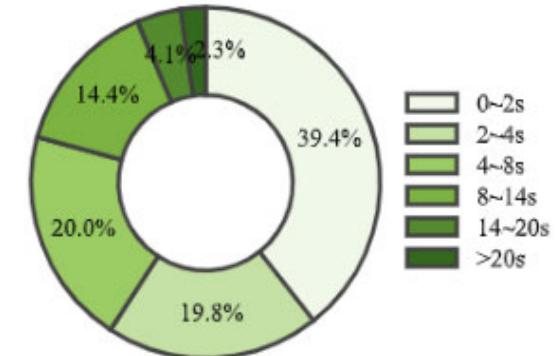
Data Processing

- Detect and split scenes → 130M video clips
- Caption video clips with BLIP-2

Video Categories

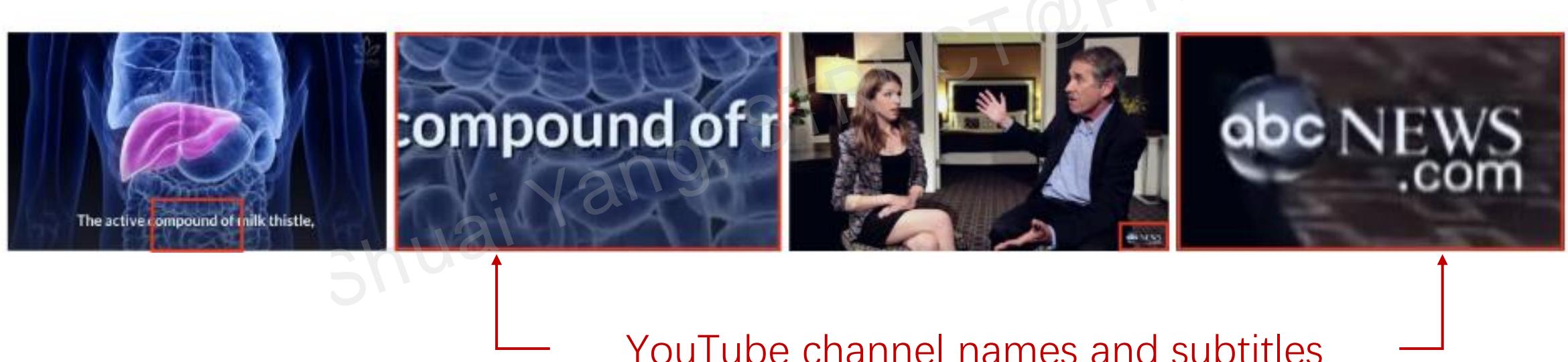


Clip Durations



Further Data Processing

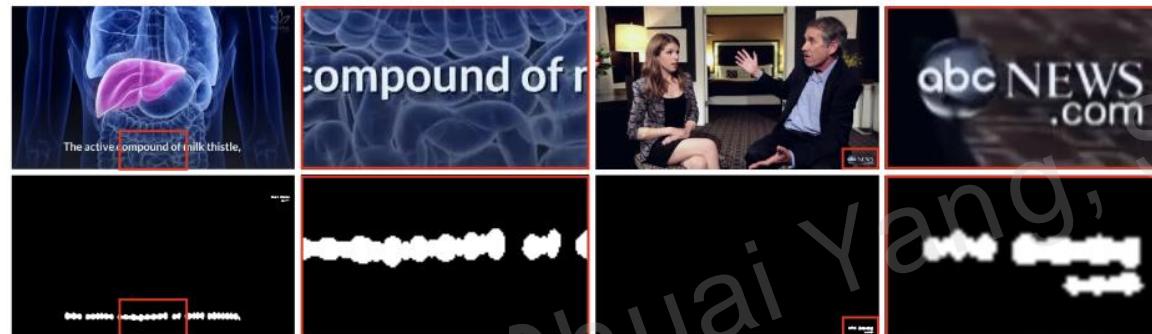
1. Text Detection



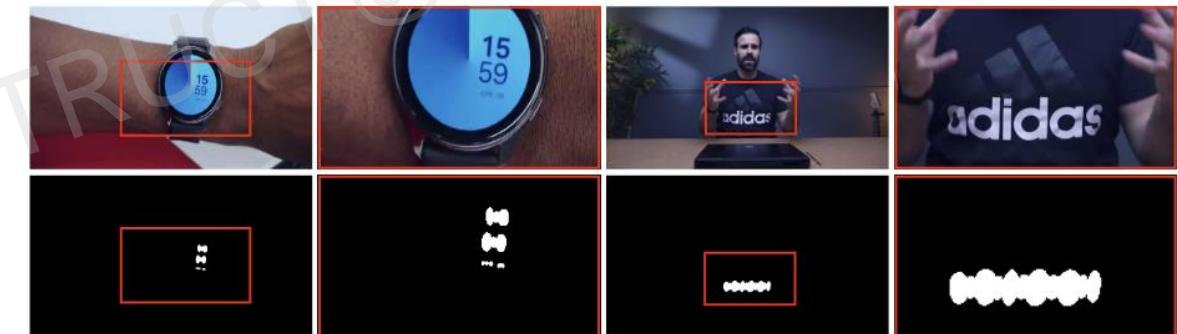
YouTube channel names and subtitles

Further Data Processing

1. Text Detection



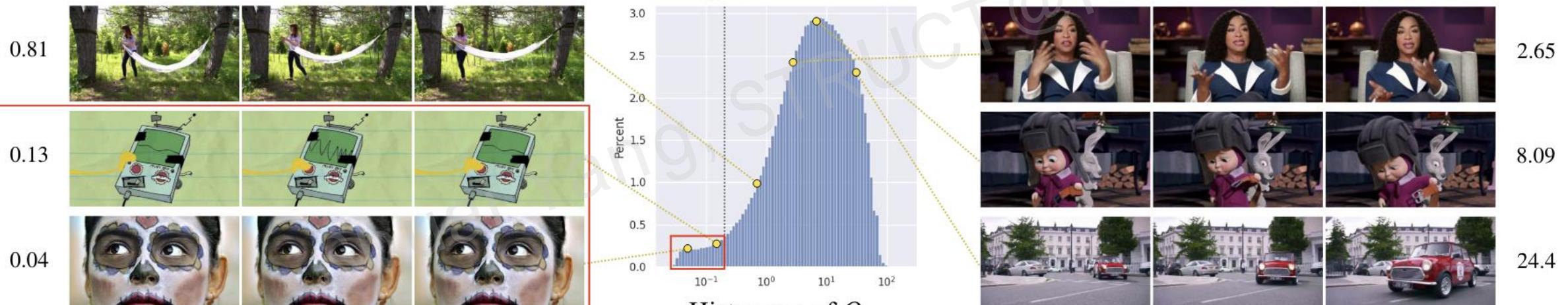
Videos filtered out due to unrelated text



Video not filtered out

Further Data Processing

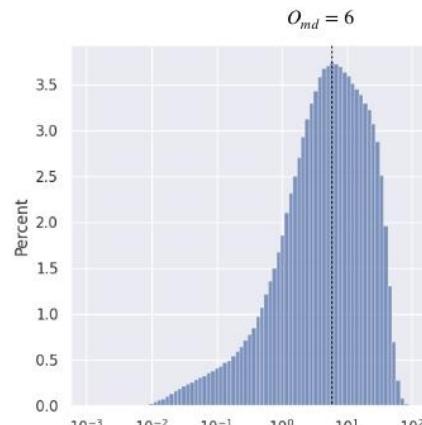
2. Motion Detection



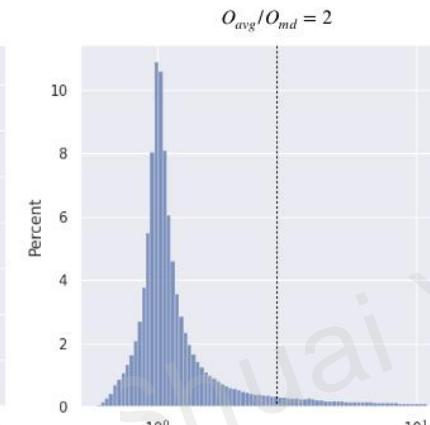
Videos lacking sufficient motion are removed

Further Data Processing

2. Motion Detection



Histogram of O_{md}



Histogram of O_{avg}/O_{md}

$$\left\{ \begin{array}{l} O_{avg} = 4.24 \\ O_{md} = 1.38 \\ \frac{O_{avg}}{O_{md}} = 3.06 \end{array} \right.$$

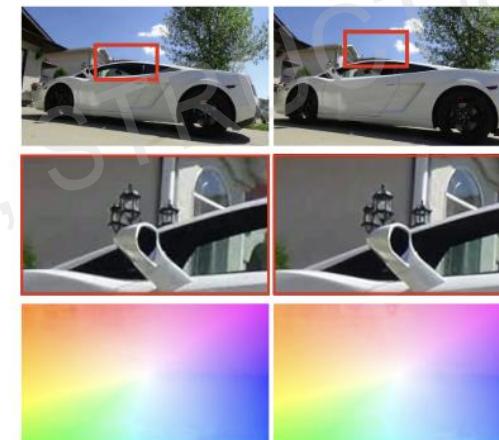


Image zooming transformation
(filtered out)

$$\left\{ \begin{array}{l} O_{avg} = 40.9 \\ O_{md} = 9.93 \\ \frac{O_{avg}}{O_{md}} = 4.11 \end{array} \right.$$



Real-world zooming
(not filtered out)

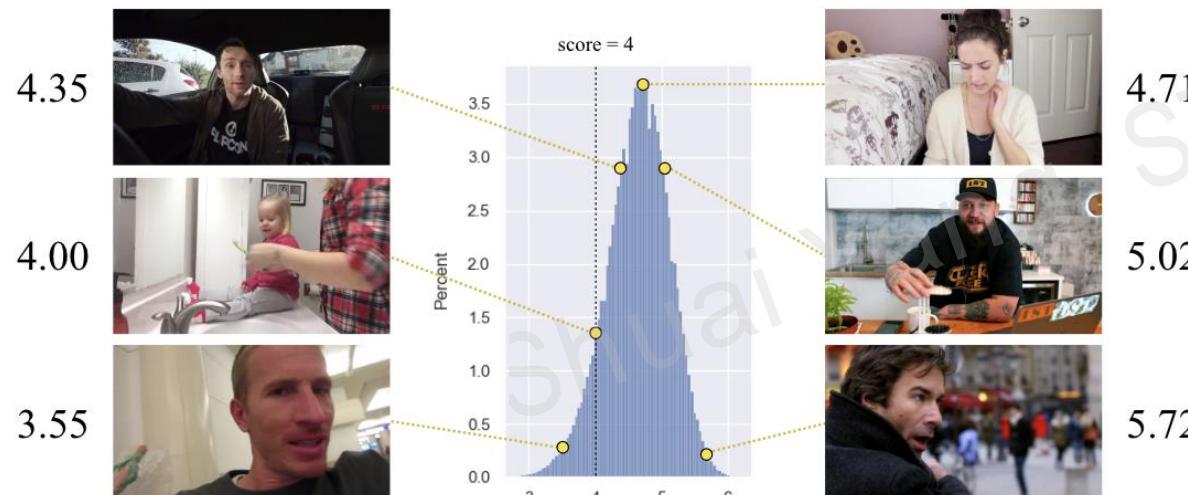
Distribution of motion features

Further Data Processing

3. Aesthetics Evaluation

Further Data Processing

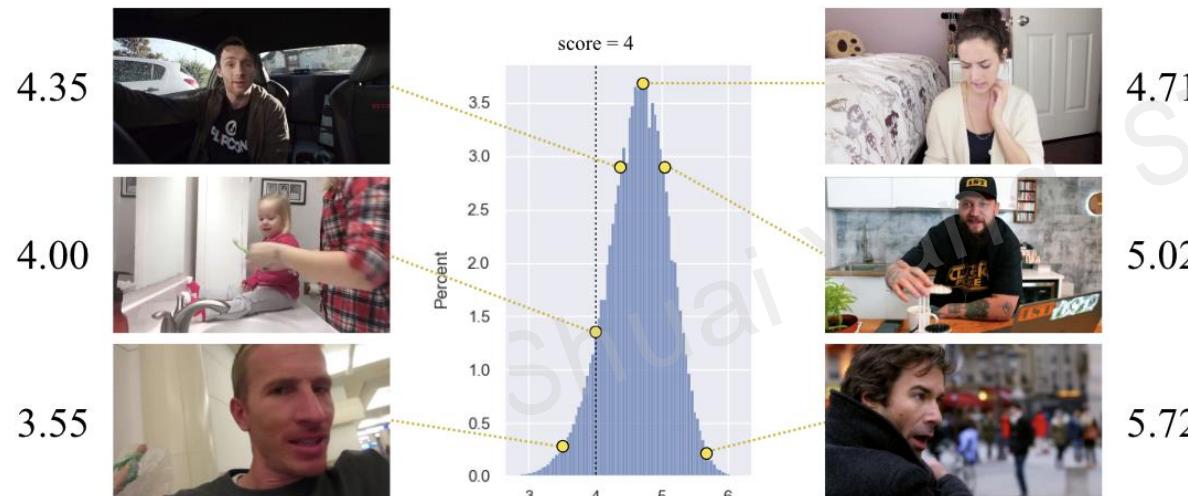
3. Aesthetics Evaluation



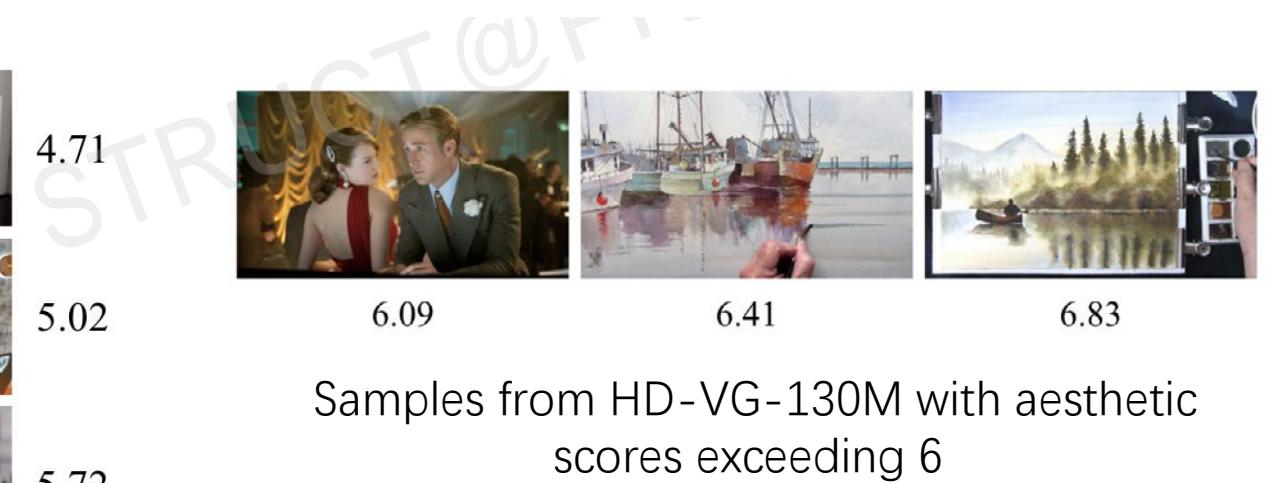
Histogram of aesthetic score

Further Data Processing

3. Aesthetics Evaluation



Histogram of aesthetic score



Samples from HD-VG-130M with aesthetic scores exceeding 6

Original Internet videos with subtitles as text labels



"He thought he was gonna get shows terrible communication on the teams part."

"Yeah, now everybody thought that we couldn't replace cat; yeah, because you're such animal lovers."

Examples from HD-VG-130M



"two men on a boat fishing on a lake"

"a group of rocks covered in seaweed"

Data Processing

Long Video



Data Processing

- **Semantics-Aware Video Splitting**

- Stage 1: shot boundary detection
- Stage 2: merging clips back according to embeddings

Data Processing

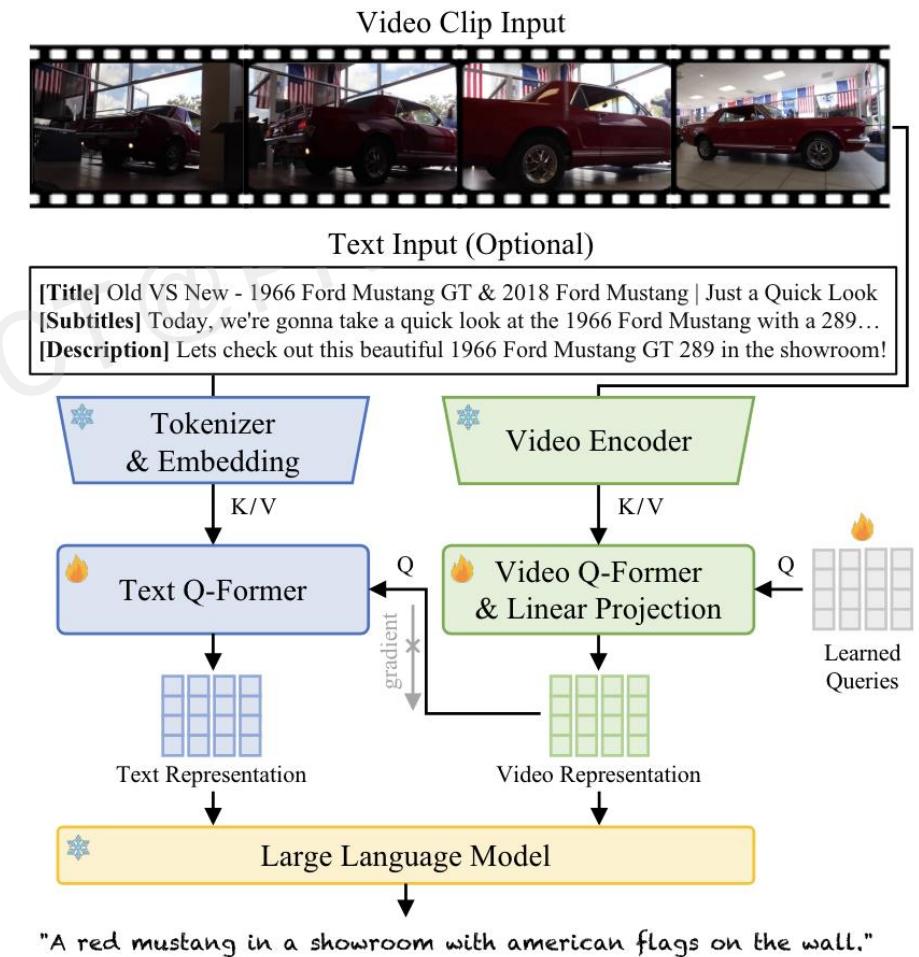
- **Semantics-Aware Video Splitting**
- **Captioning with Cross-Modality Teachers**
 - Select 8 video captioning models from 31 candidates based on user study

Data Processing

- **Semantics-Aware Video Splitting**
- **Captioning with Cross-Modality Teachers**
 - Select 8 video captioning models from 31 candidates based on user study
- **Fine-grained Video-to-Text Retrieval**
 - Collect a 100K videos with human selection

Data Processing

- Semantics-Aware Video Splitting
- Captioning with Cross-Modality Teachers
- Fine-grained Video-to-Text Retrieval
- Multimodal Student Captioning Model



Data-driven model

6 Panda-70M



A rhino and a lion are fighting in the dirt



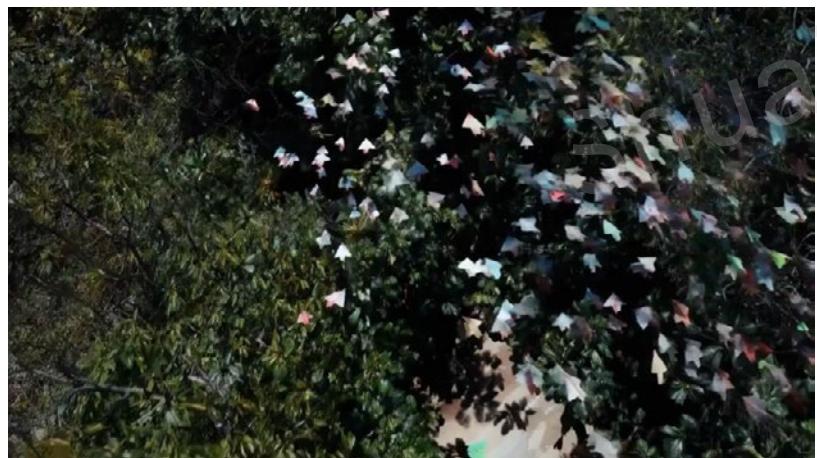
A person is making a pie crust on a table



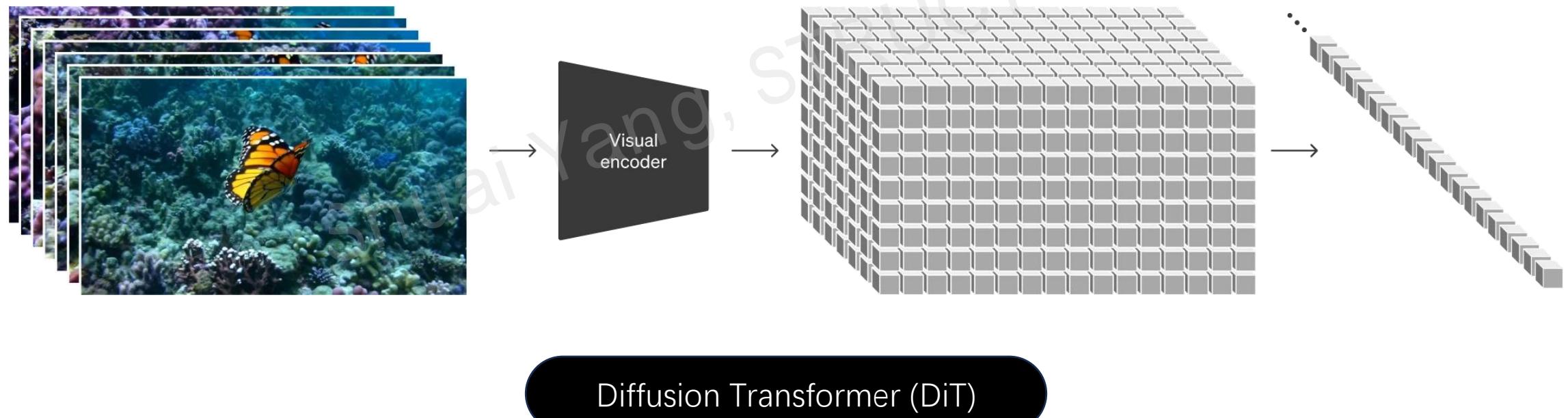
It is a rally car driving on a dirt road in the countryside, with people watching from the side of the road.

Data-driven model

7 Sora



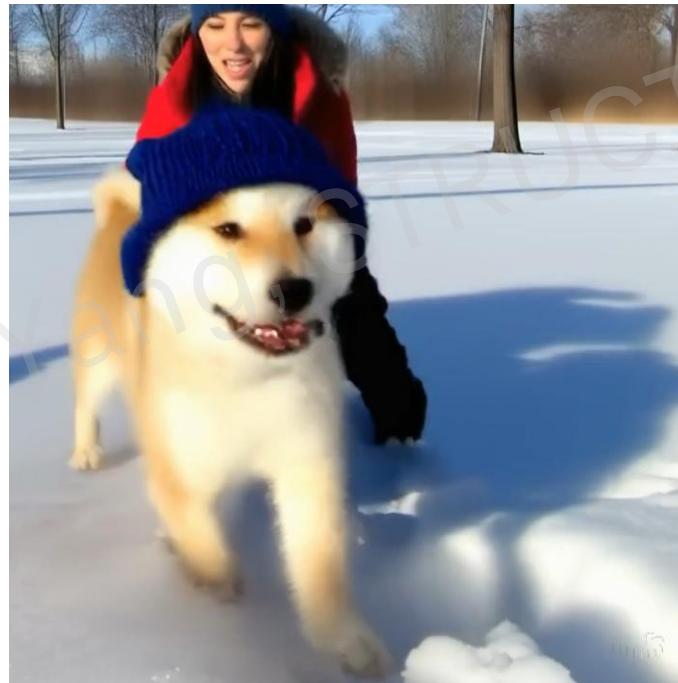
Turning visual data into patches



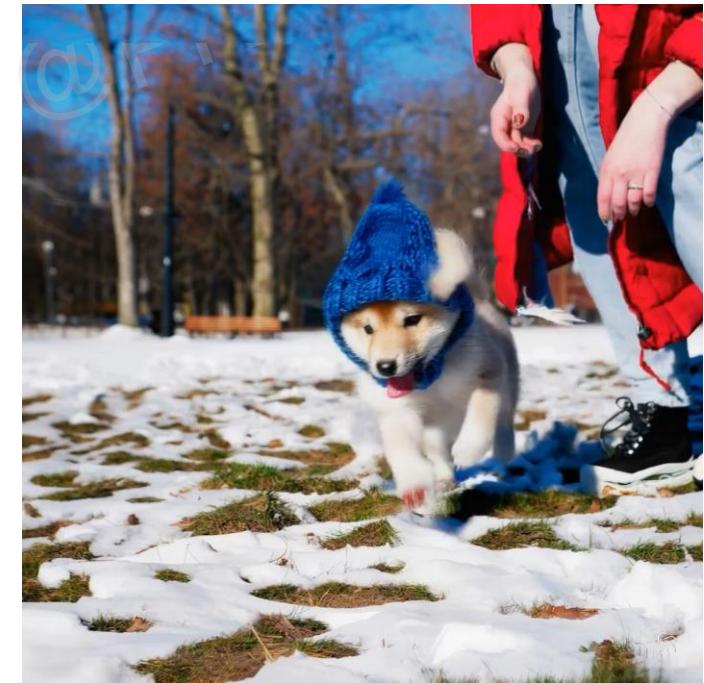
Scaling transformers for video generation



Base compute



4x compute



32x compute

Sampling flexibility



Improved framing and composition



Trained on square crops



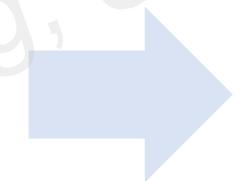
Trained at native aspect ratios

Re-captioning technique

a toy robot wearing a green dress and a sun hat taking a pleasant stroll in Johannesburg, South Africa during a winter storm.



Animating images



A Shiba Inu dog wearing a beret and black turtleneck.

Extending generated videos

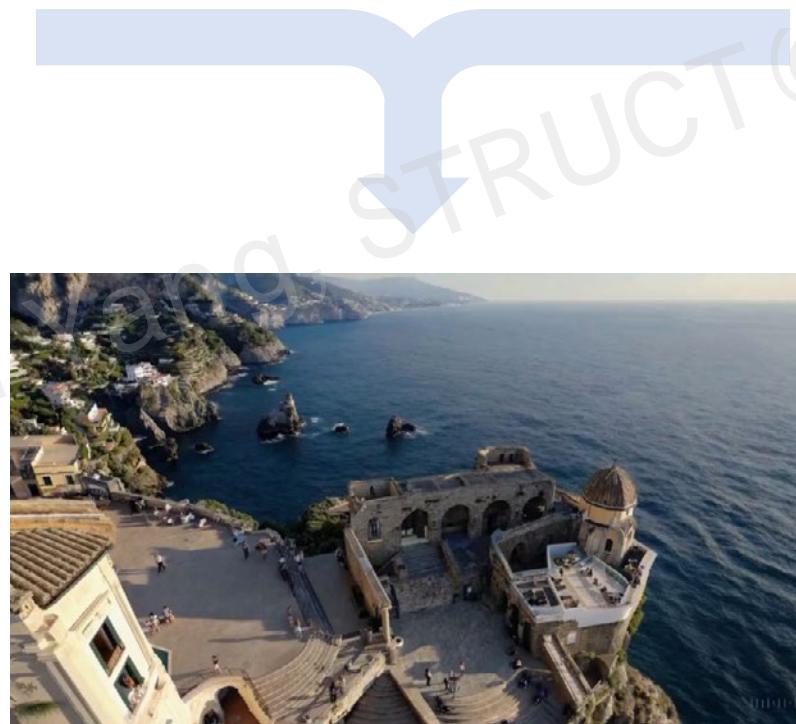


Extended backward in time starting from the same video segment

Data-driven model

7 Sora

Connecting videos

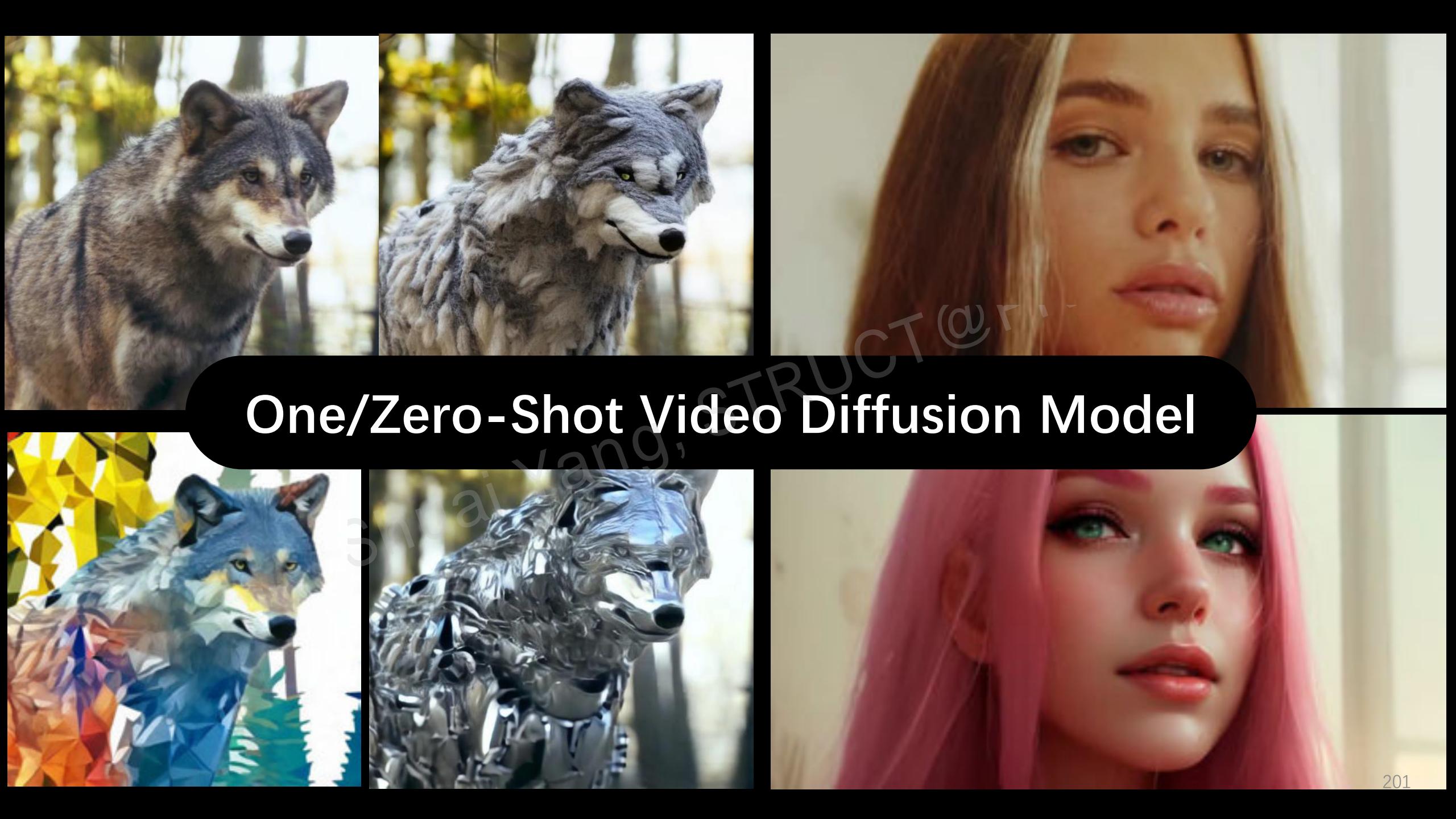


Data-driven model



Huge amount of data

Huge computational cost



One/Zero-Shot Video Diffusion Model

One-shot model

1



Tune-A-Video

ICCV 23

2



Edit-A-Video

ACML 23

3

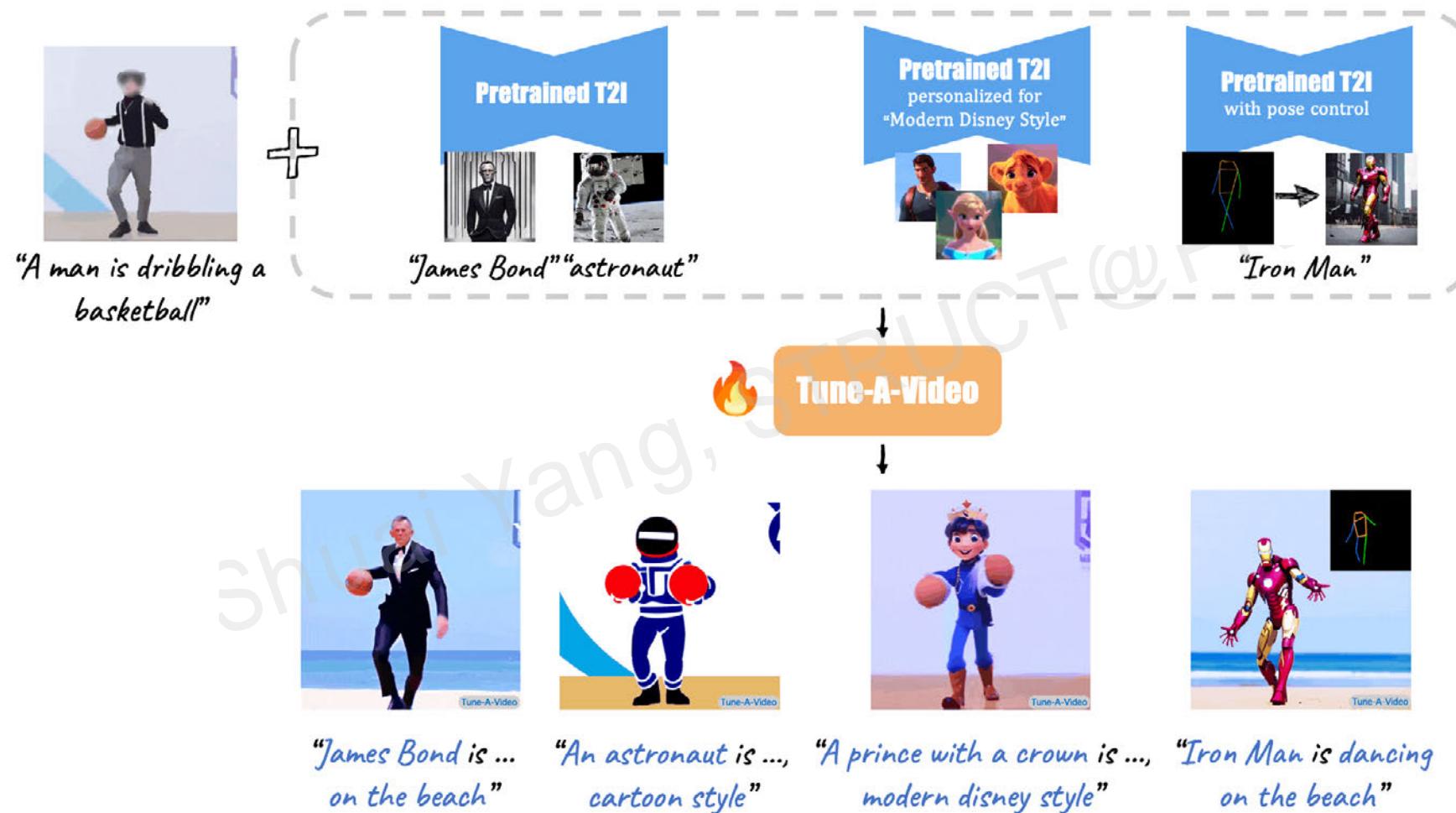


Video-P2P

CVPR 24

One-shot model

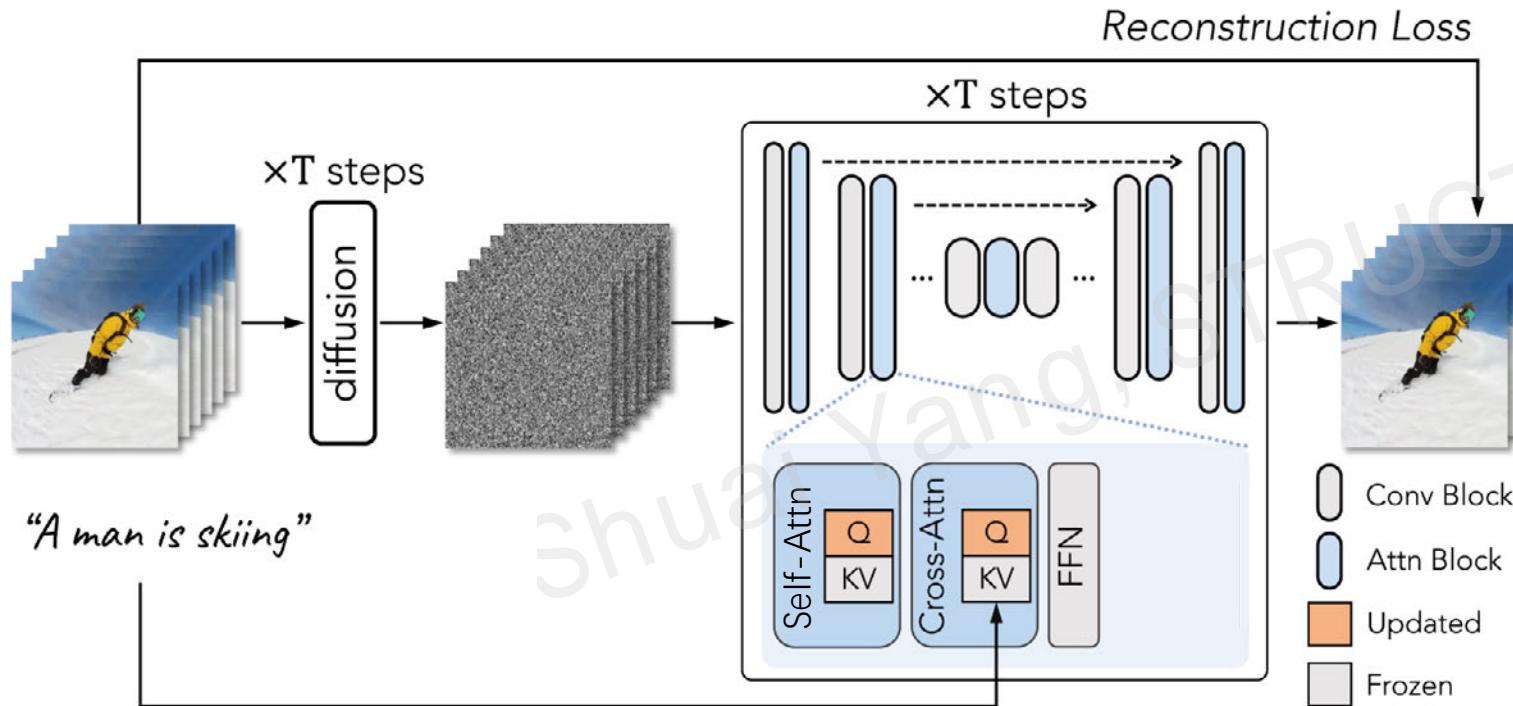
1 Tune-A-Video



One-shot model

1 Tune-A-Video

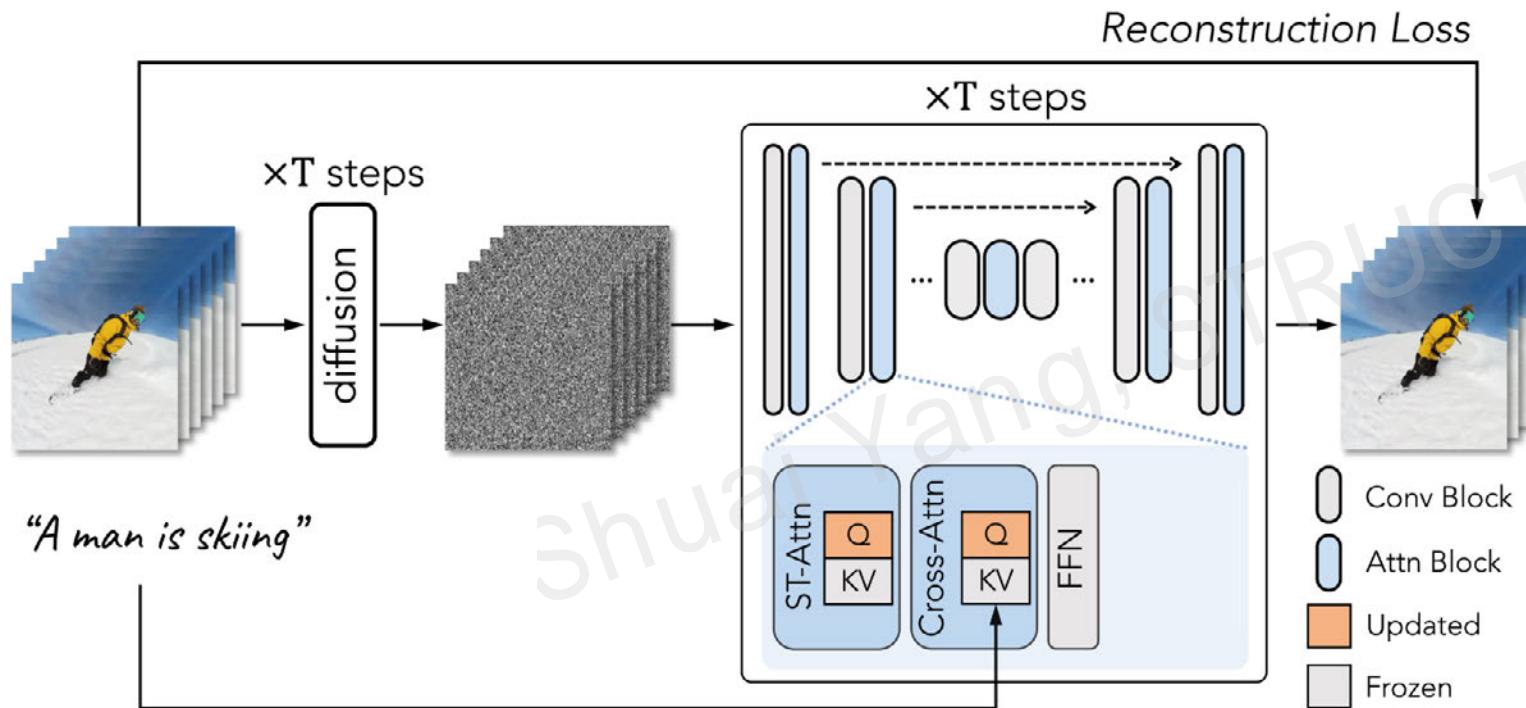
Fine-Tuning



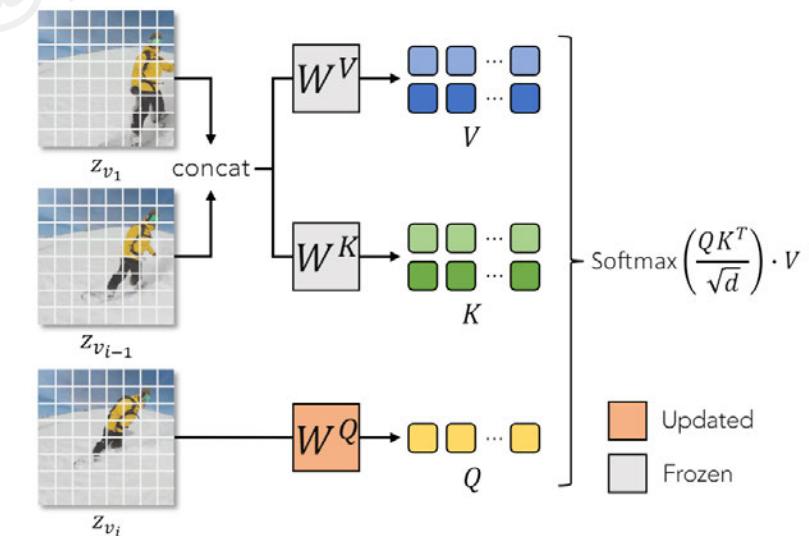
One-shot model

1 Tune-A-Video

Fine-Tuning

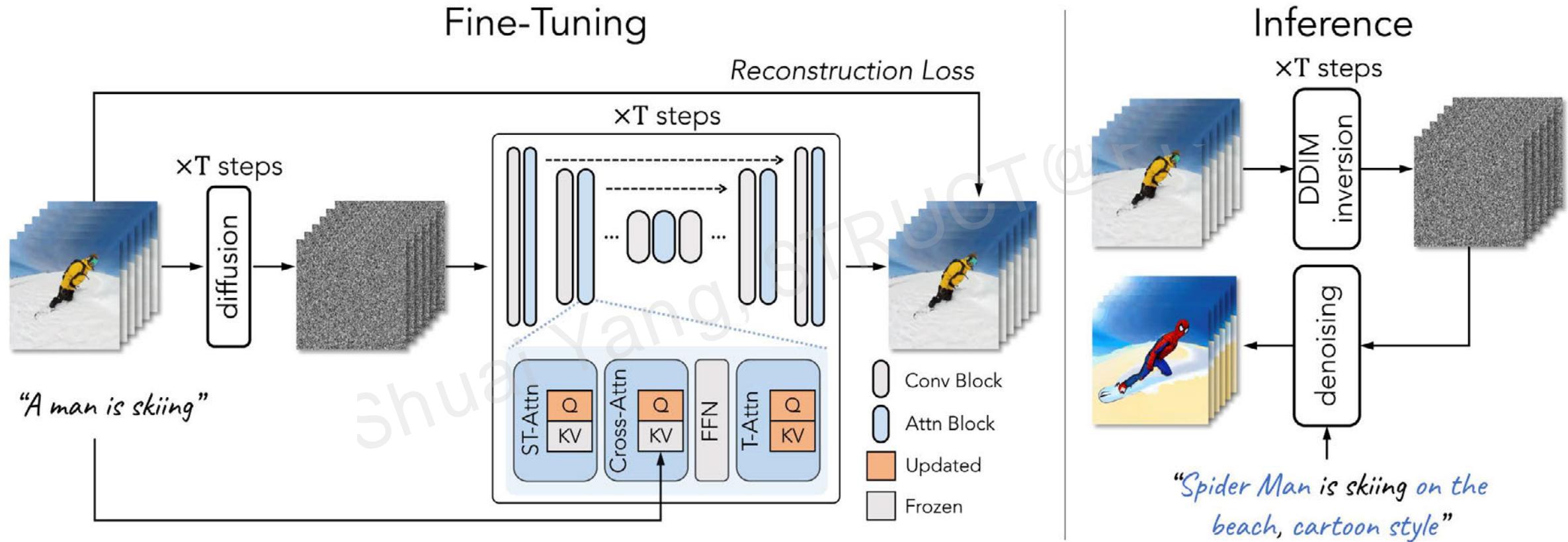


SA-Attn:
Cross-frame attention



One-shot model

1 Tune-A-Video



One-shot model

1 Tune-A-Video



A jeep car is moving on the road



*A jeep car is
moving on the road,
cartoon style*



*A jeep car is moving
on the snow*

Change too much
in the background

One-shot model

2 Edit-A-Video



*man
gorilla*



Tune-A-Video

a man is riding a motorcycle

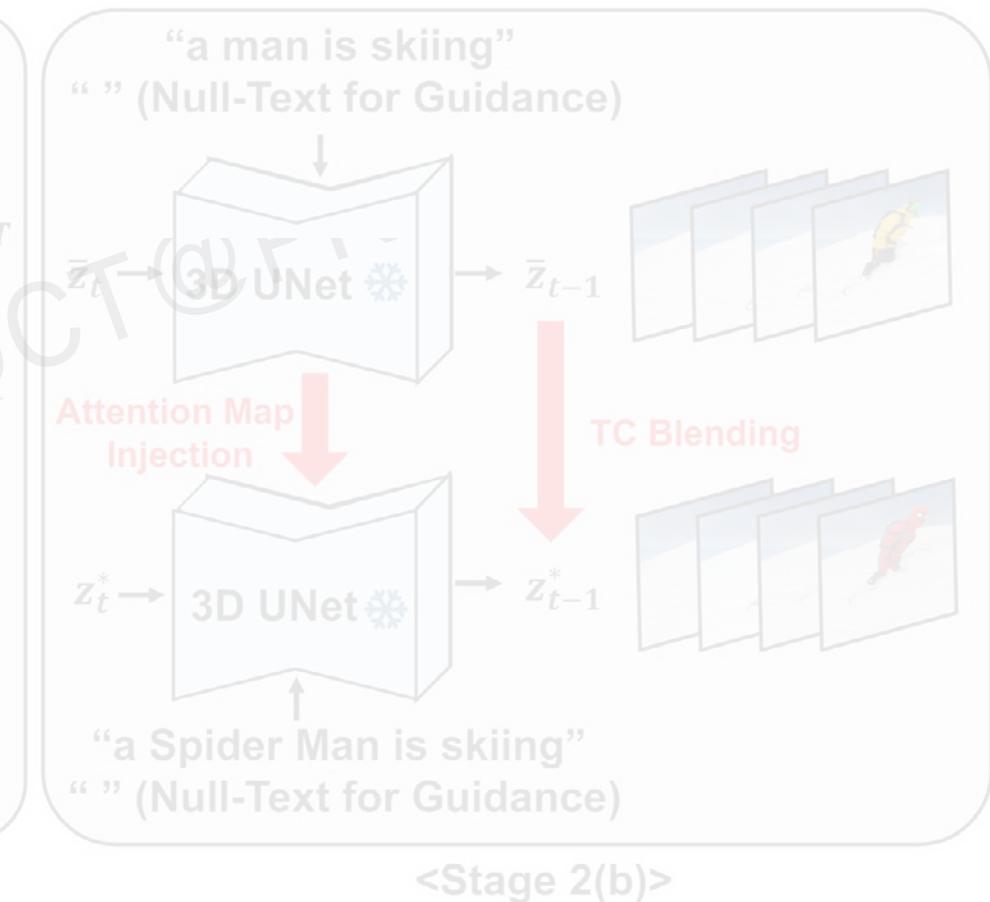
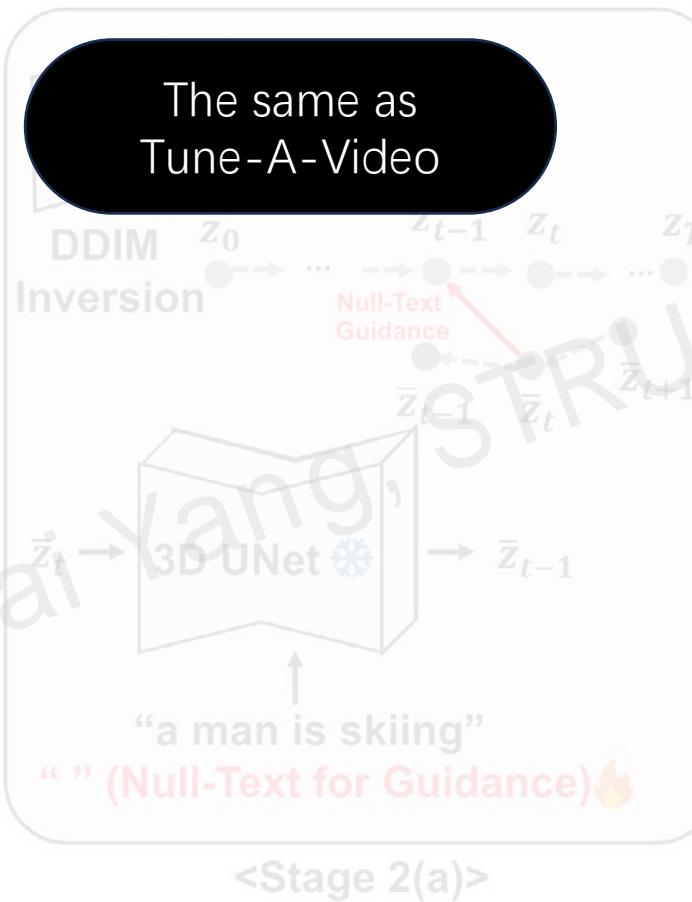
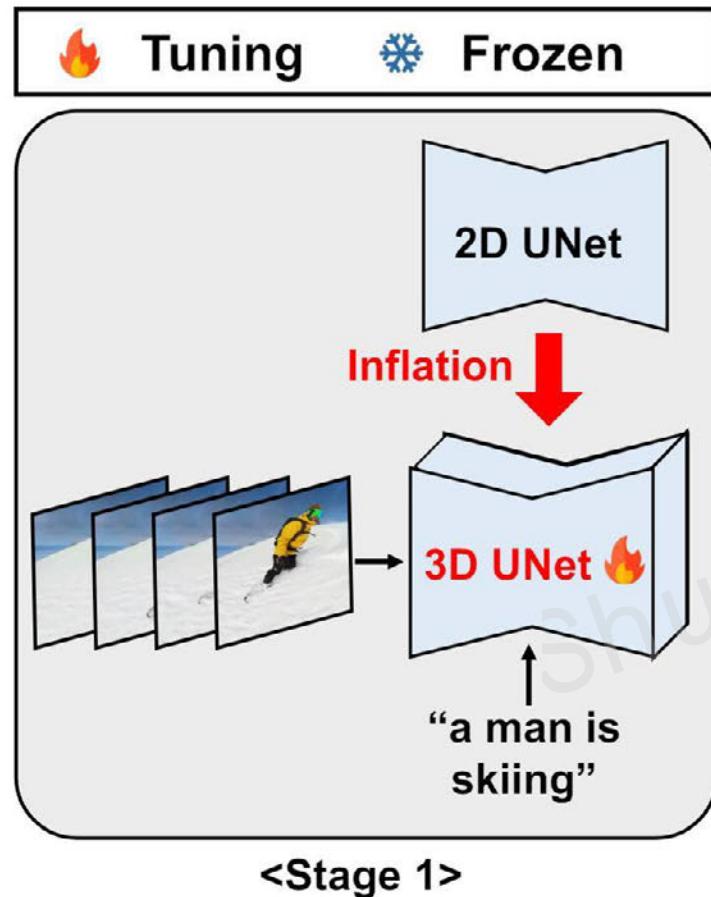


Edit-A-Video

Better background

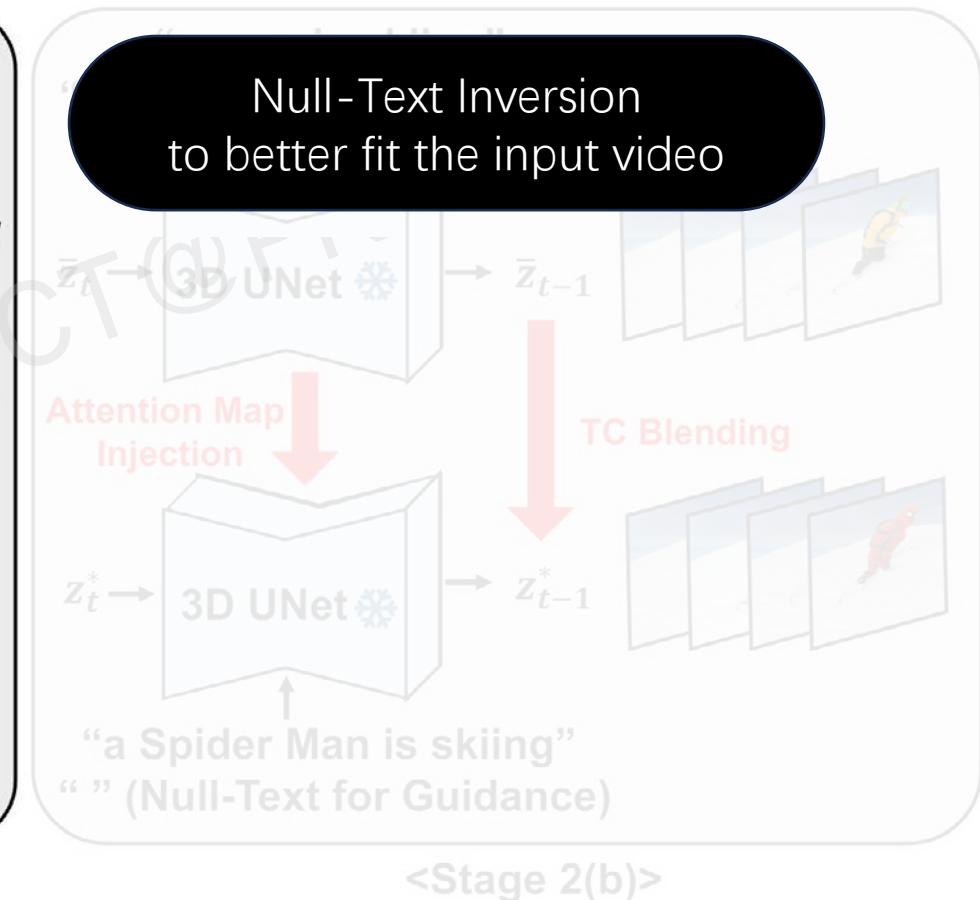
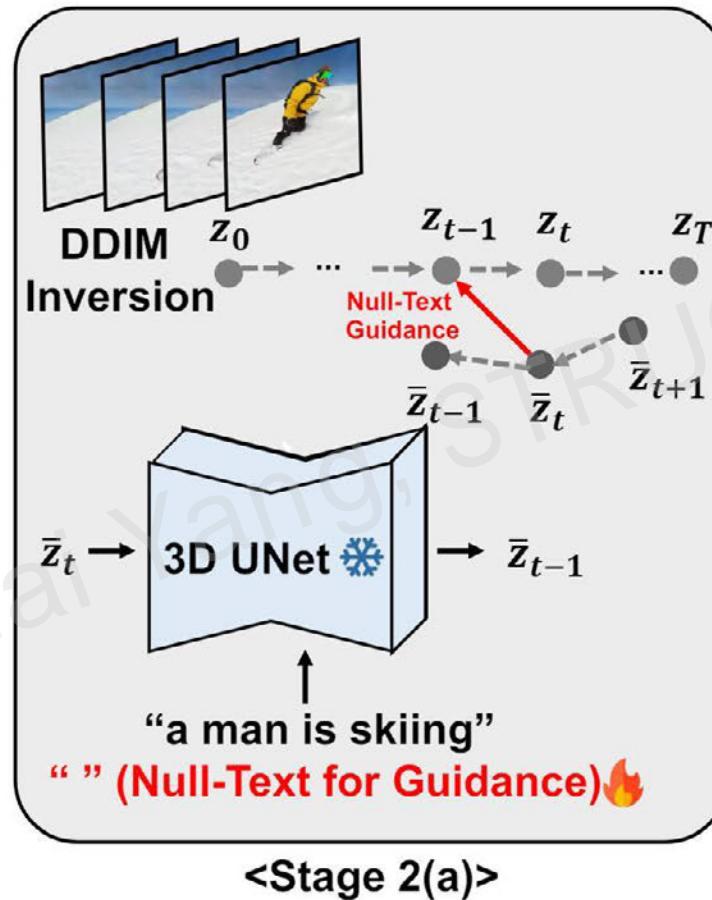
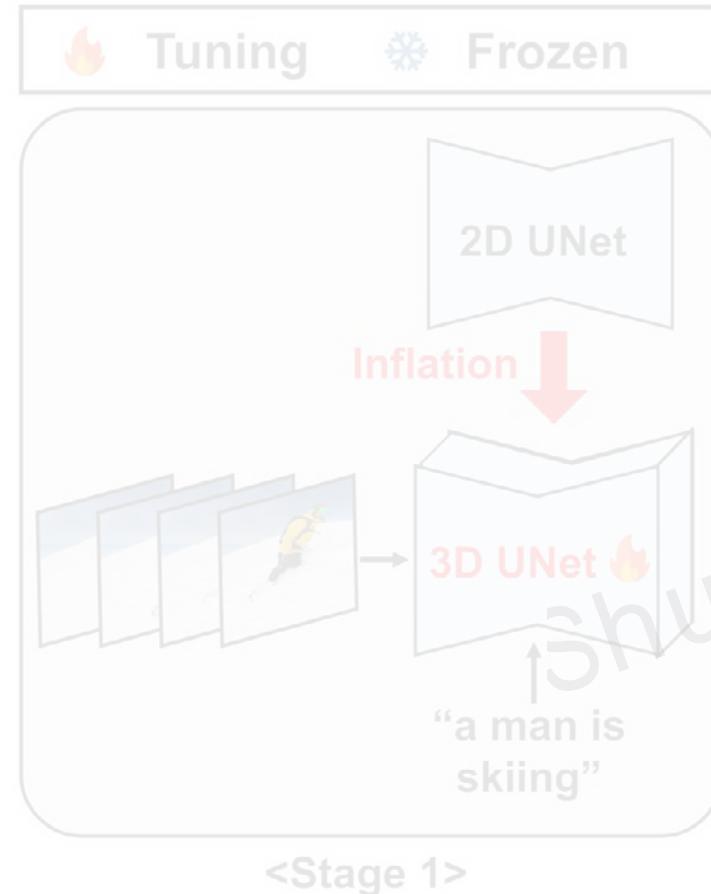
One-shot model

2 Edit-A-Video



One-shot model

2 Edit-A-Video



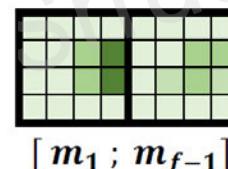
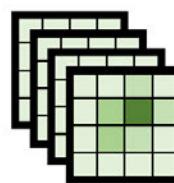
One-shot model

2 Edit-A-Video

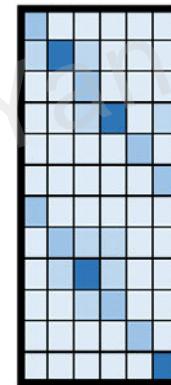
Modified Prompt-to-Prompt

Background mask should be also temporal consistent

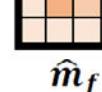
“a man is skiing”



\times

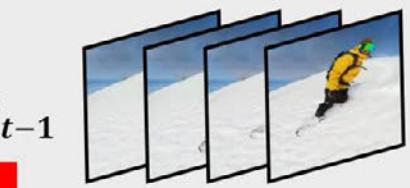
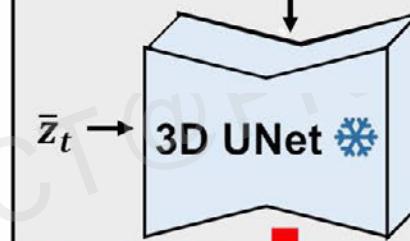


=

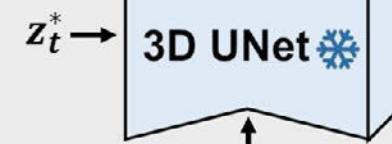


Cross Attention Maps
for “man” (m)

“a man is skiing”
“ ” (Null-Text for Guidance)



Attention Map
Injection

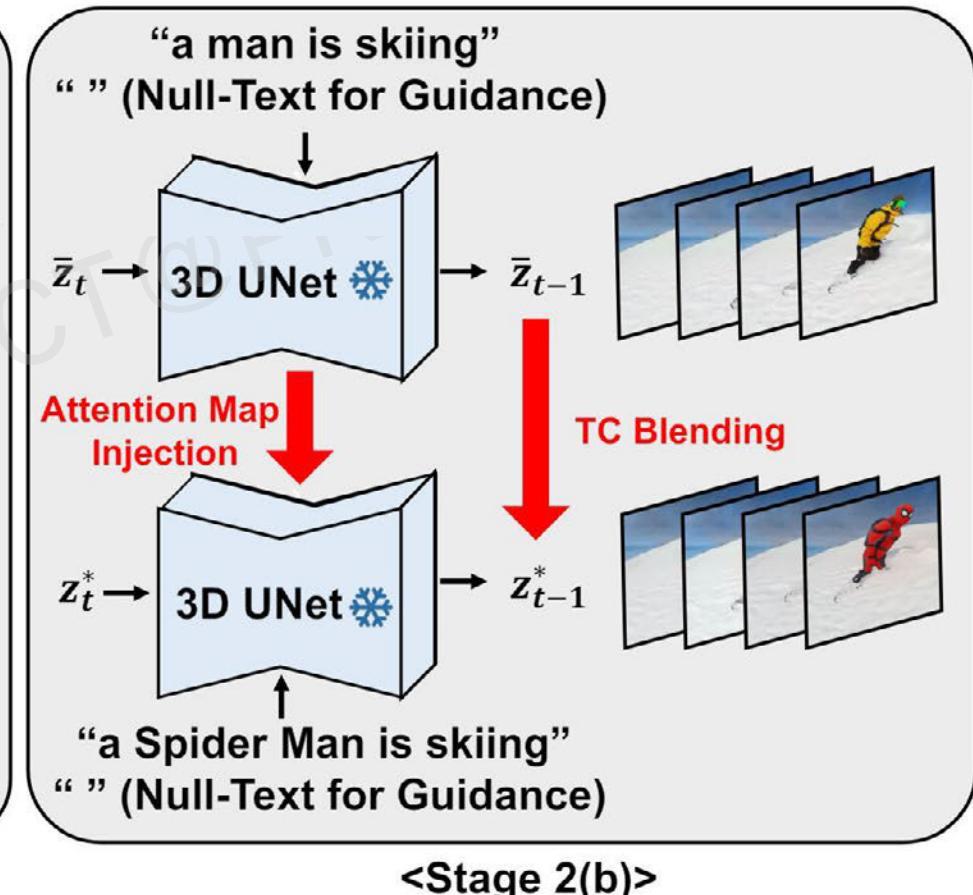
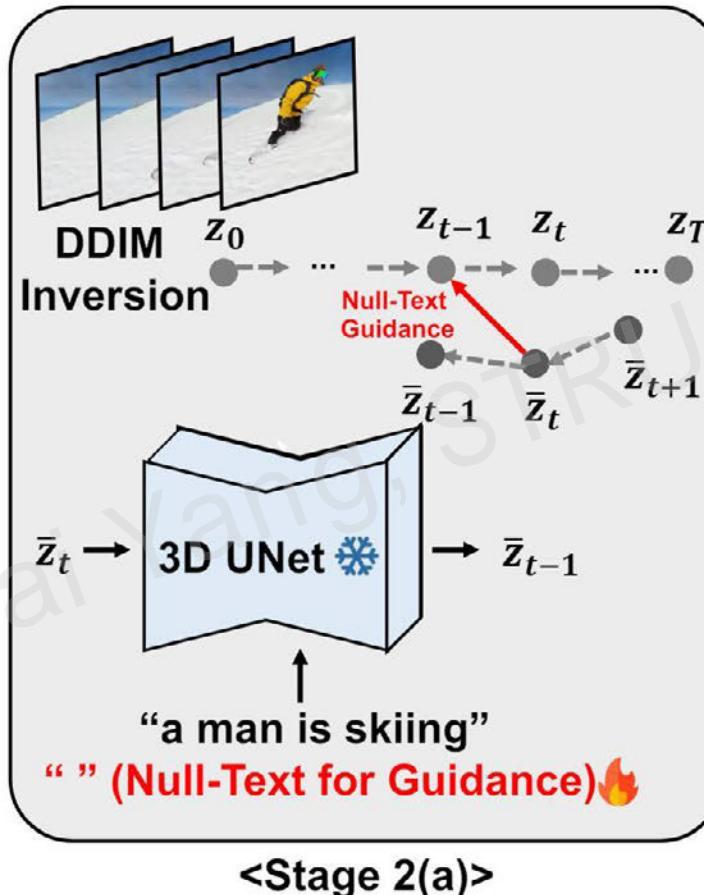
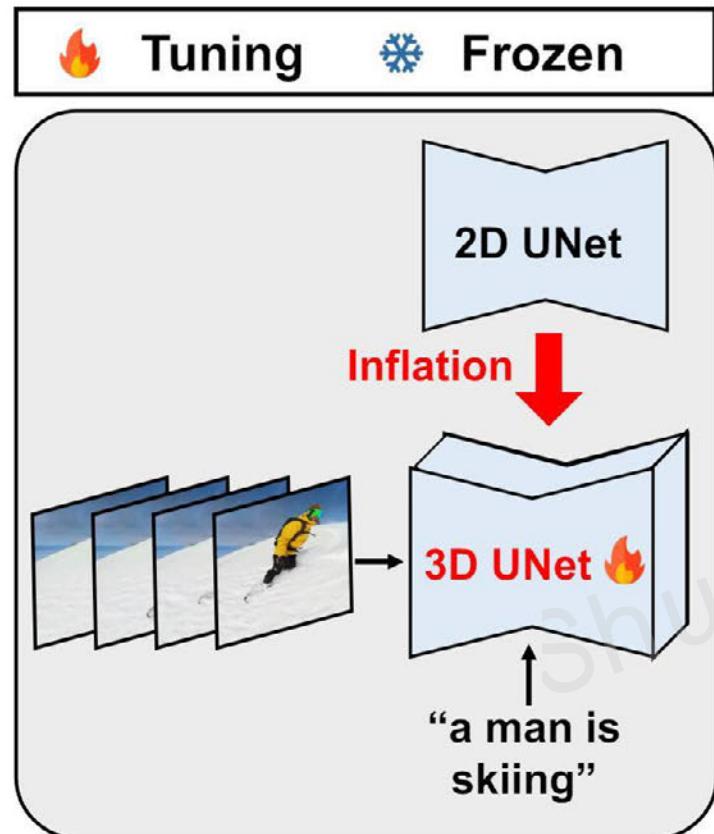


“a Spider Man is skiing”
“ ” (Null-Text for Guidance)

<Stage 2(b)>

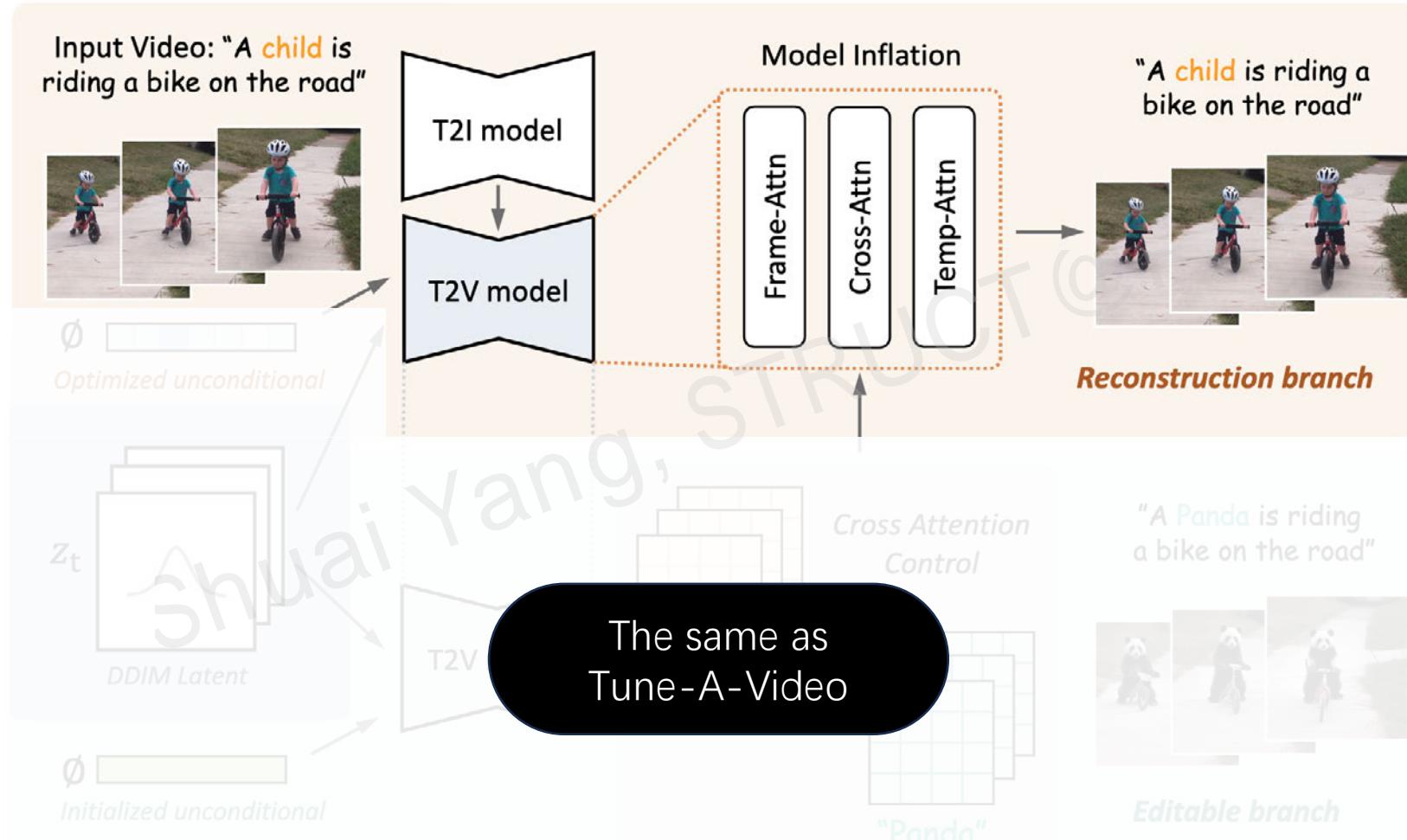
One-shot model

2 Edit-A-Video



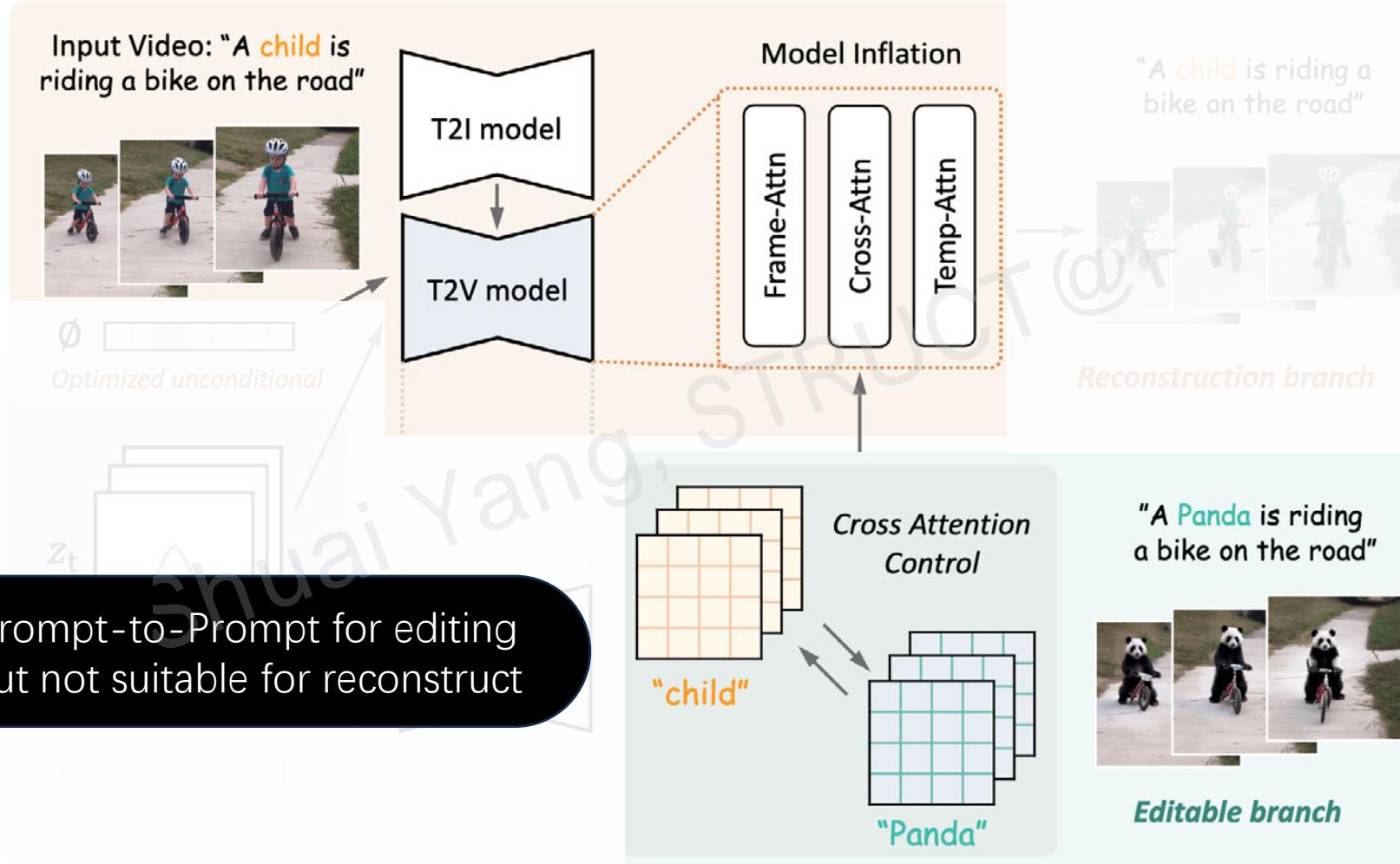
One-shot model

3 Video-P2P



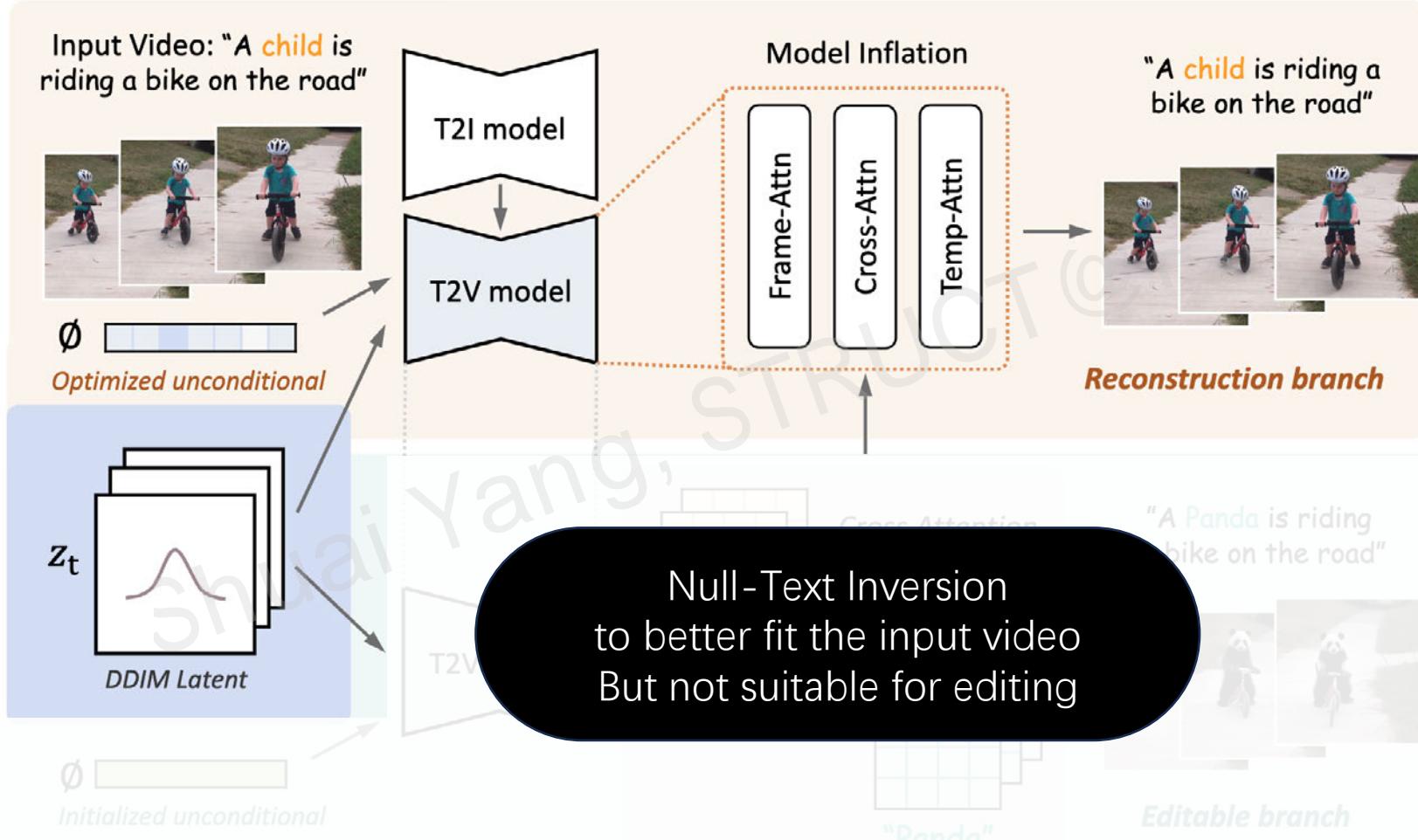
One-shot model

3 Video-P2P



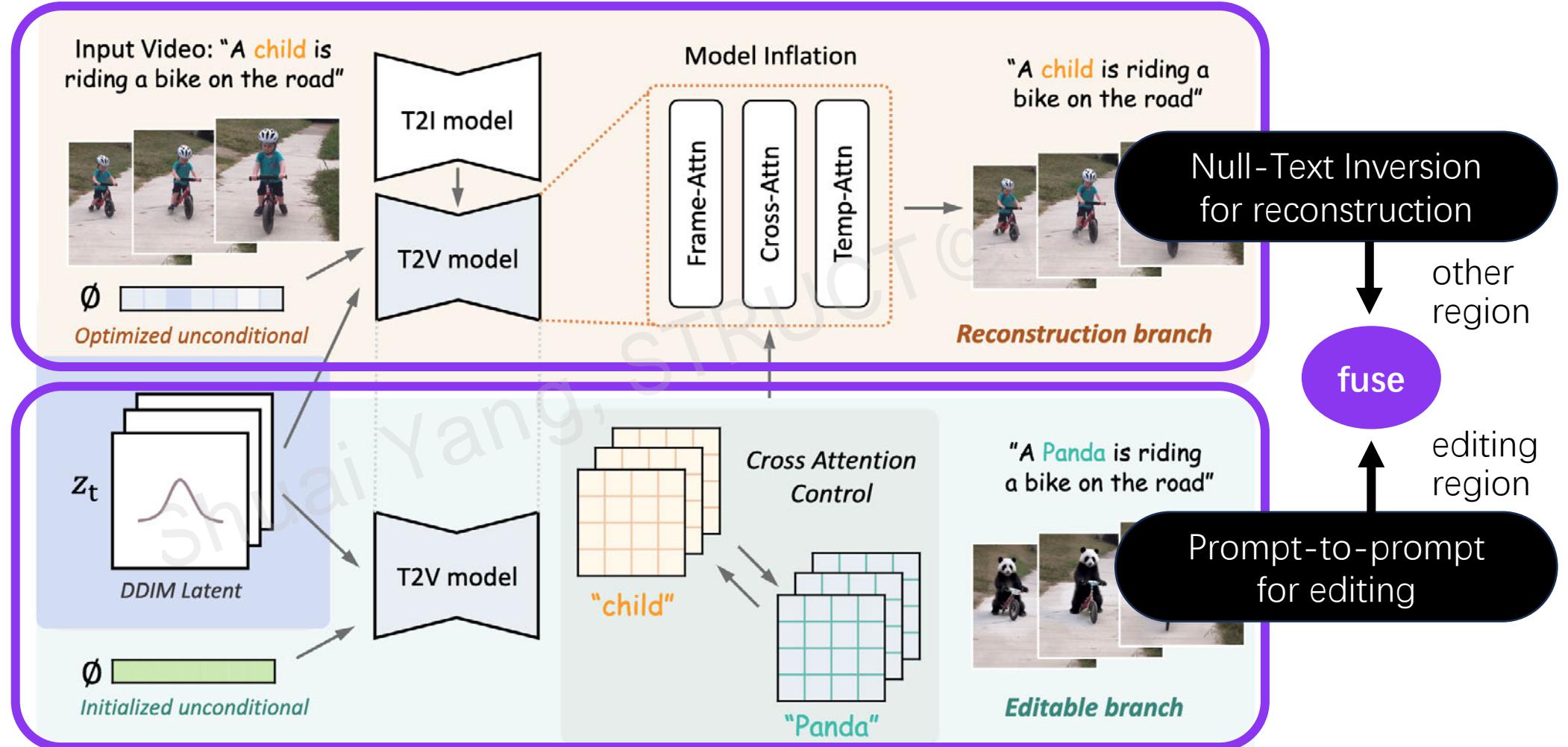
One-shot model

3 Video-P2P



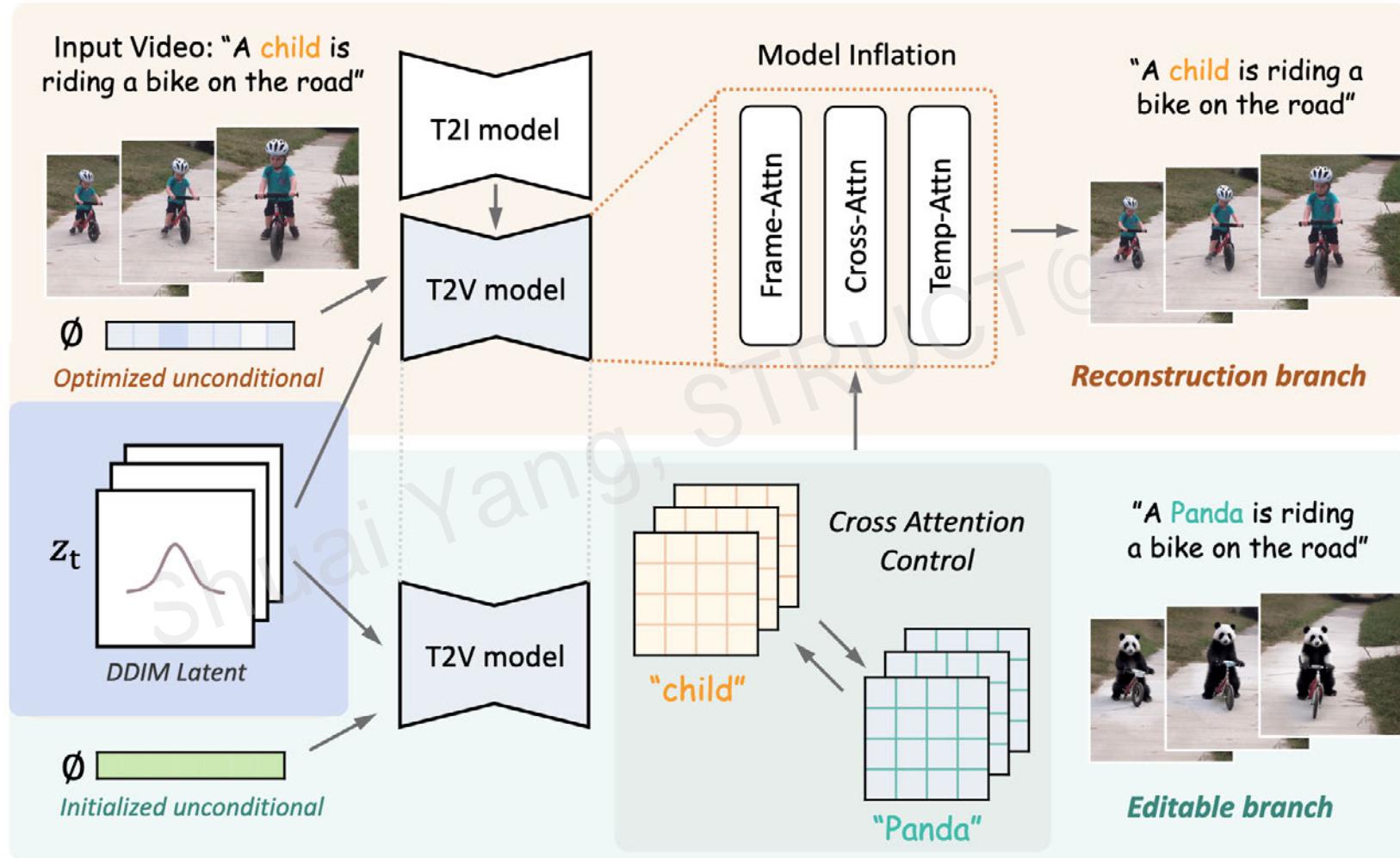
One-shot model

3 Video-P2P



One-shot model

3 Video-P2P



One-shot model

3 Video-P2P



a **crochet** penguin is running on the ice

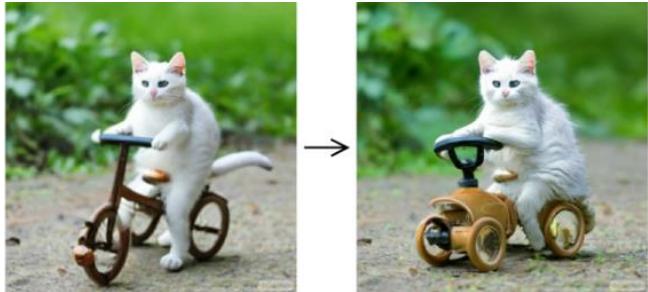


a jeep car is moving on the road **beach**



Zero-shot model

Image Editing

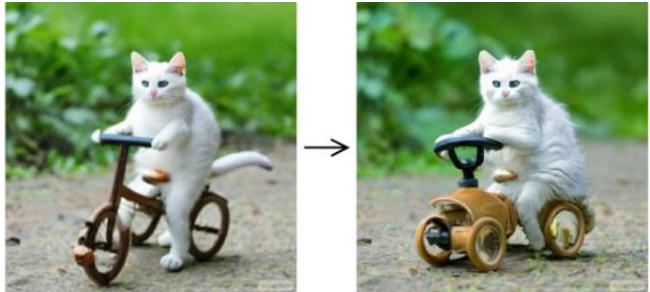


“Photo of a cat riding on a bicycle.”
car

Prompt-to-Prompt
Null-Text Inversion
Plug-and-Play

Zero-shot model

Image Editing



“Photo of a cat riding on a bicycle.”
car



Inversion-Based Model



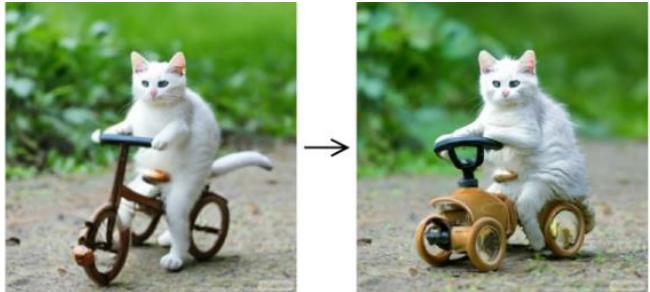
ICCV'23
ICLR'24
CVPR'24

FateZero
TokenFlow
VidToMe

Pixel2Video
FLATTEN

Zero-shot model

Image Editing



“Photo of a cat riding on a bicycle.”
car



Inversion-Based Model



ICCV'23
ICLR'24
CVPR'24

FateZero
TokenFlow

Pixel2Video
FLATTEN
VidToMe

Conditional Generation



ControlNet
SDEdit

Zero-shot model

Image Editing



“Photo of a cat riding on a bicycle.”
car



Inversion-Based Model



ICCV'23
ICLR'24
CVPR'24

FateZero
TokenFlow

Pixel2Video
FLATTEN
VidToMe

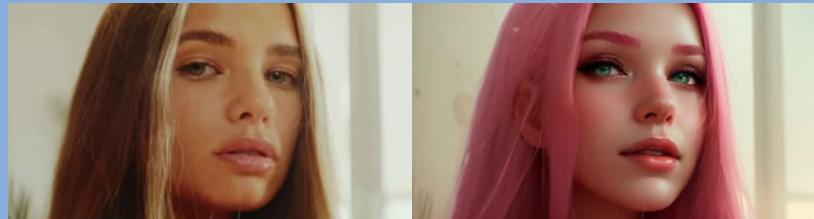
Conditional Generation



ControlNet
SDEdit



Inversion-Free Model



ICCV'23
SIGGRAPH Asia'23
ICLR'24
CVPR'24

Text2Video-Zero
Rerender-A-Video
ControlVideo
FRESCO Fairy

Zero-shot model

Image Editing



“Photo of a cat riding on a bicycle.”
car



Inversion-Based Model



ICCV'23
ICLR'24
CVPR'24

FateZero
TokenFlow

Pixel2Video
FLATTEN
VidToMe

Conditional Generation



ControlNet
SDEdit



Inversion-Free Model



ICCV'23
SIGGRAPH Asia'23
ICLR'24
CVPR'24

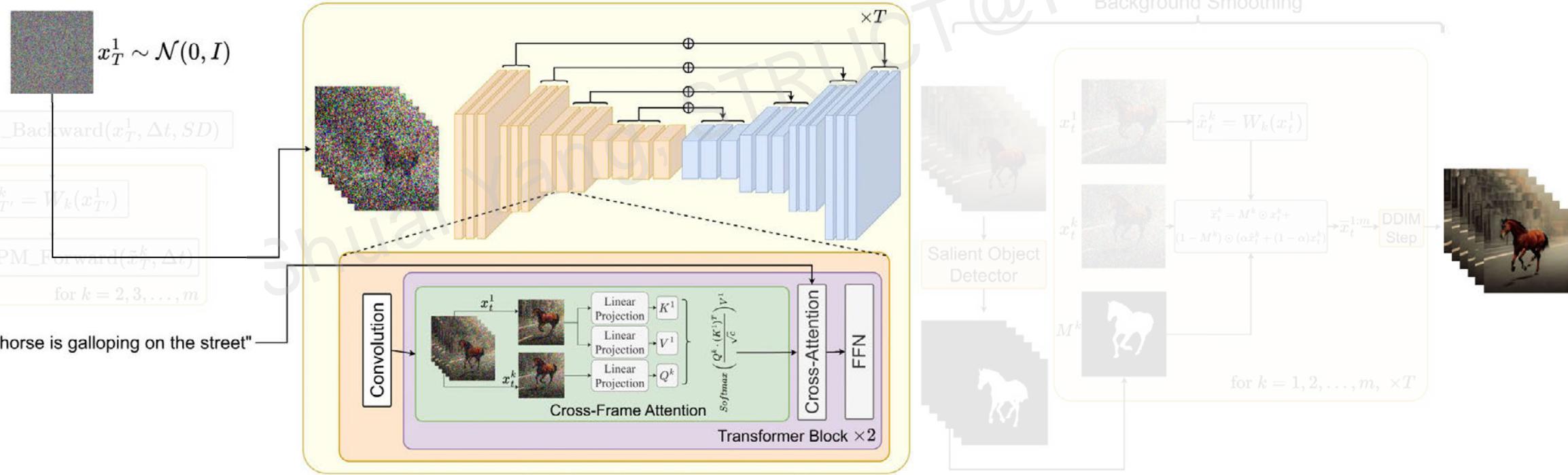
Text2Video-Zero
Rerender-A-Video
ControlVideo
FRESCO Fairy

Inversion-free zero-shot model

1 Text2Video-Zero



cross-frame attention
Not enough due to random noises

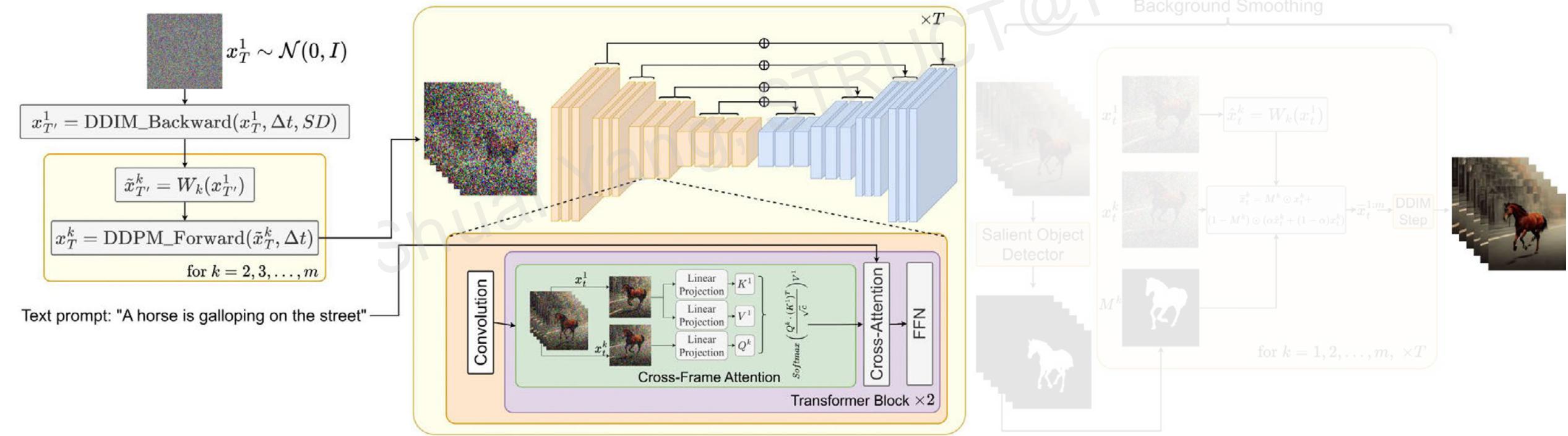


Inversion-free zero-shot model

1 Text2Video-Zero



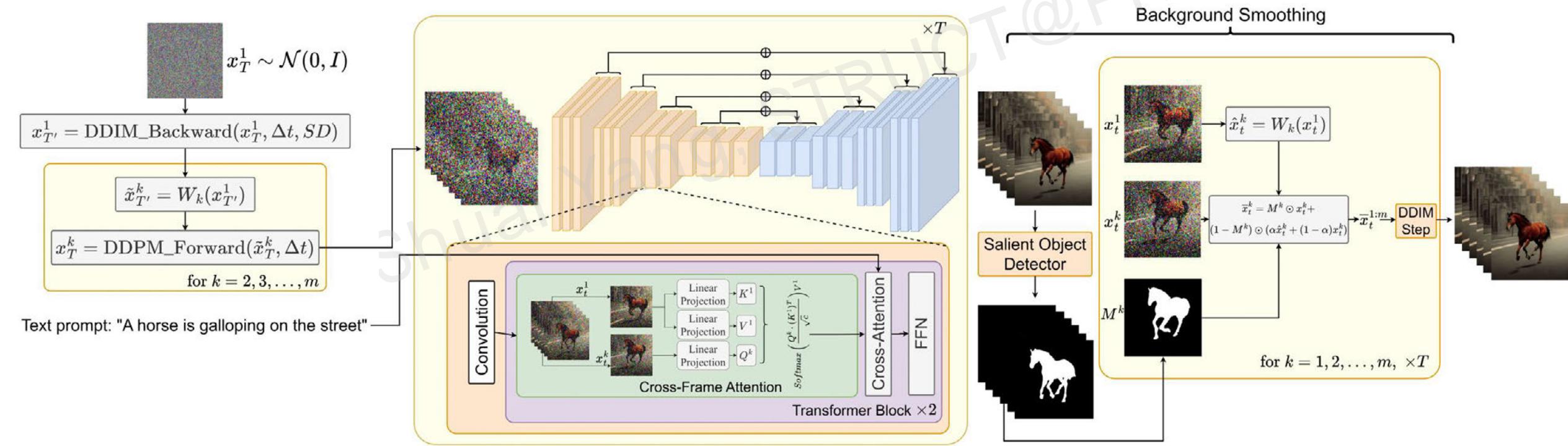
Use the same noises with translation



Inversion-free zero-shot model

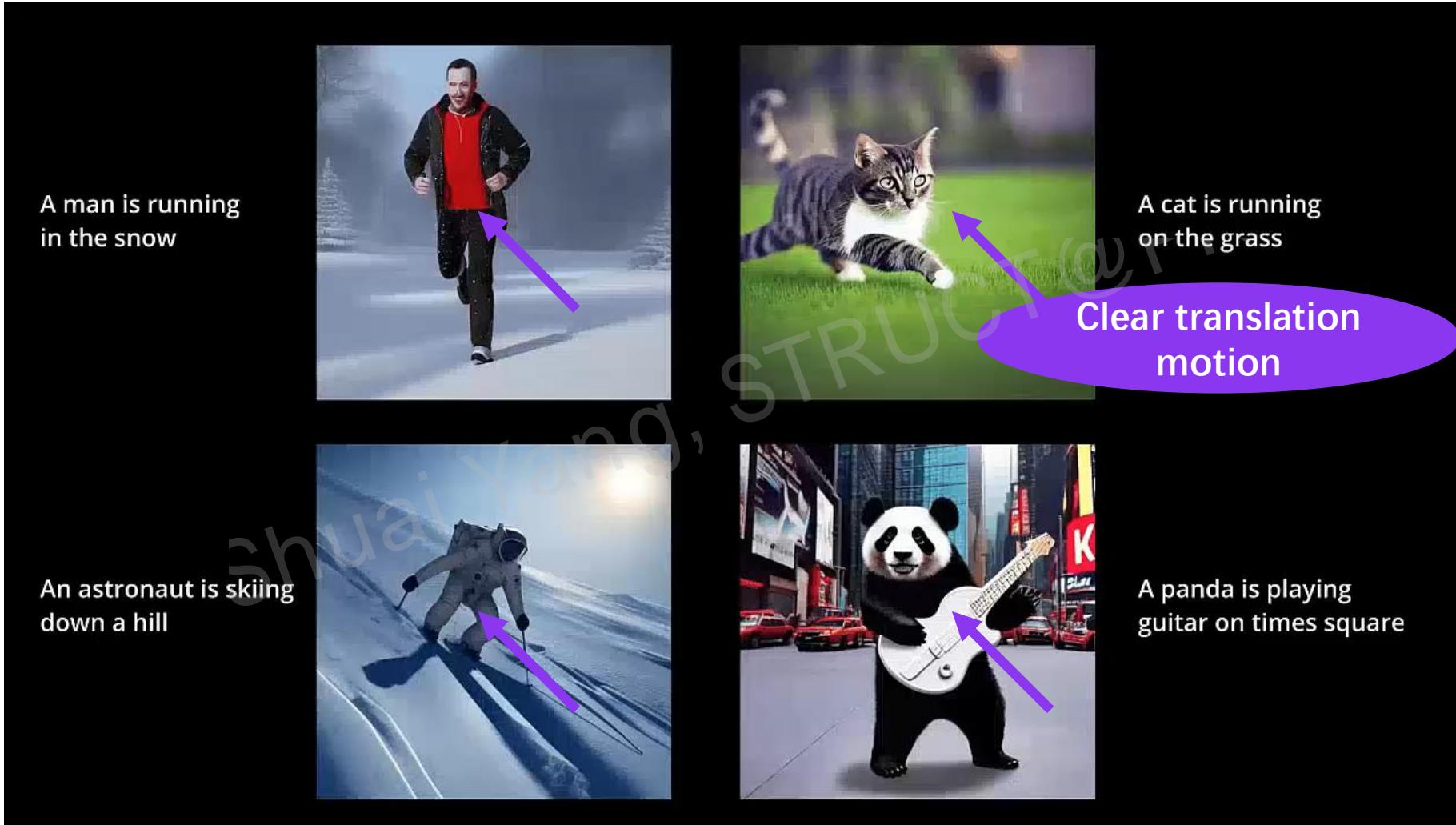
1 Text2Video-Zero

Background smooth
Keep the non-salient region the same



Inversion-free zero-shot model

1 Text2Video-Zero



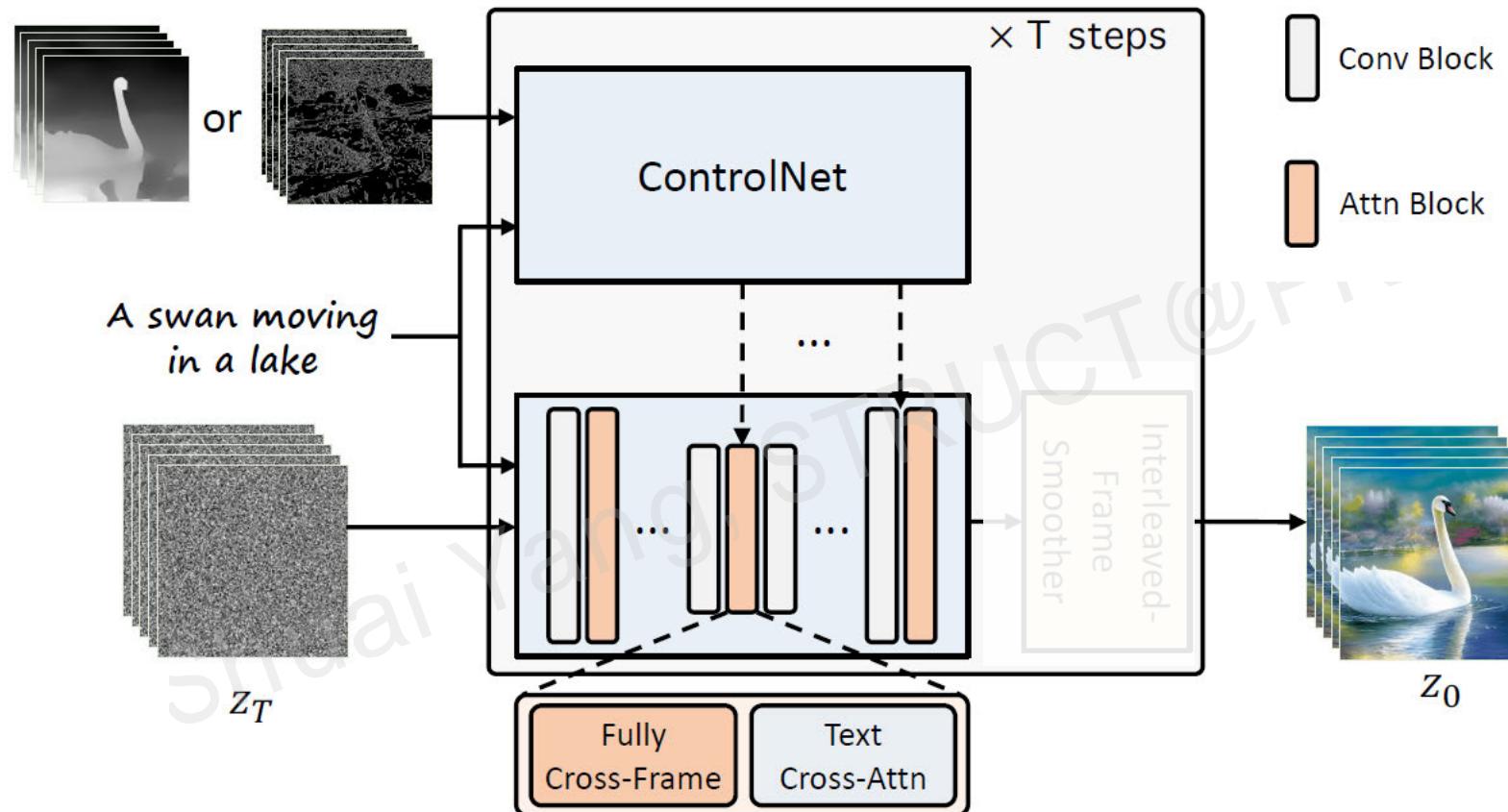
Inversion-free zero-shot model

1 Text2Video-Zero



Inversion-free zero-shot model

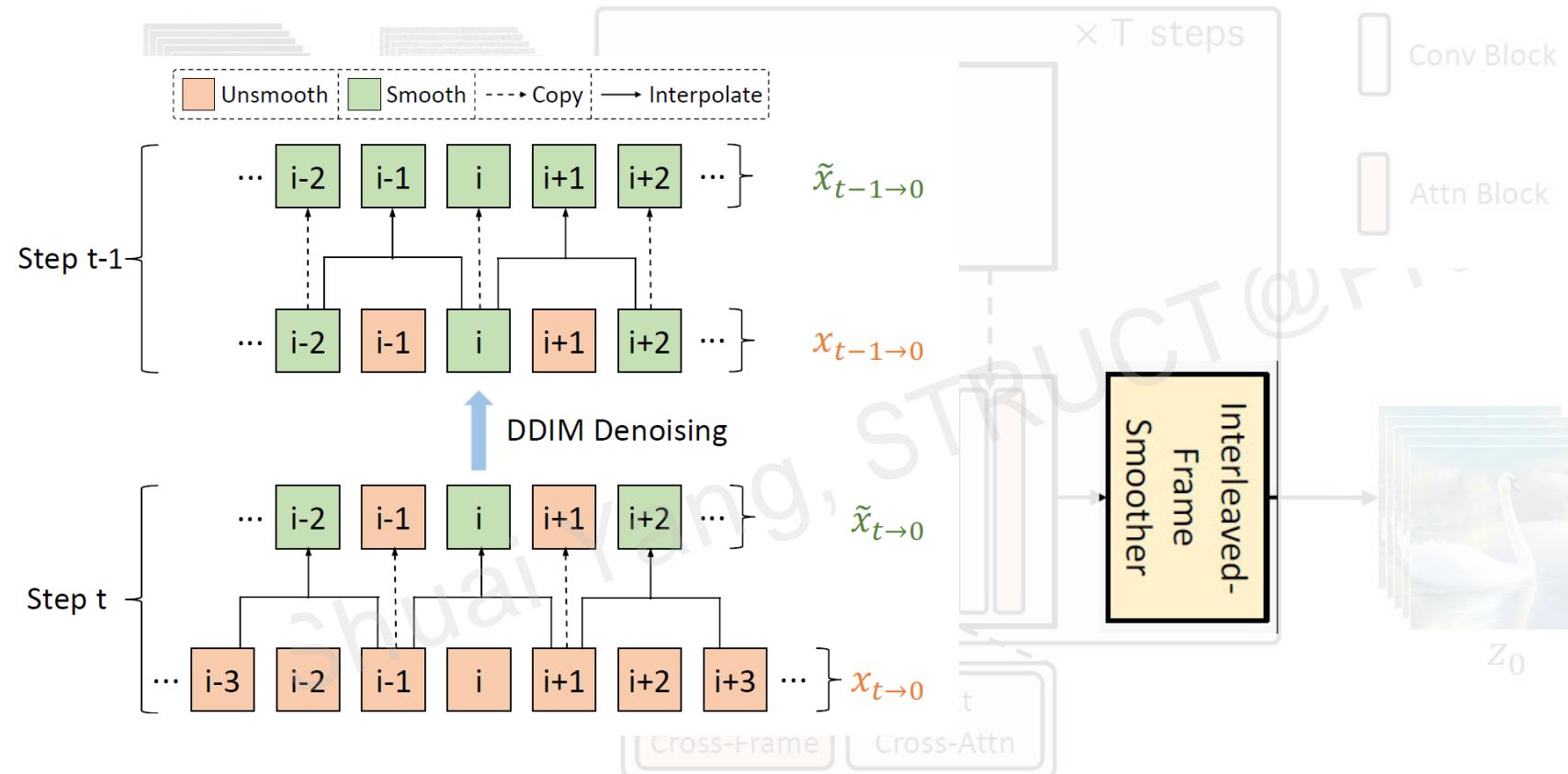
2 ControlVideo



ControlNet with cross-frame attention

Inversion-free zero-shot model

2 ControlVideo



Interleaved-frame smoother:
denoised frame interpolation

still very flickering

Can we use stronger constraints, such as ... optical flow?



Input



w/o smoother



full

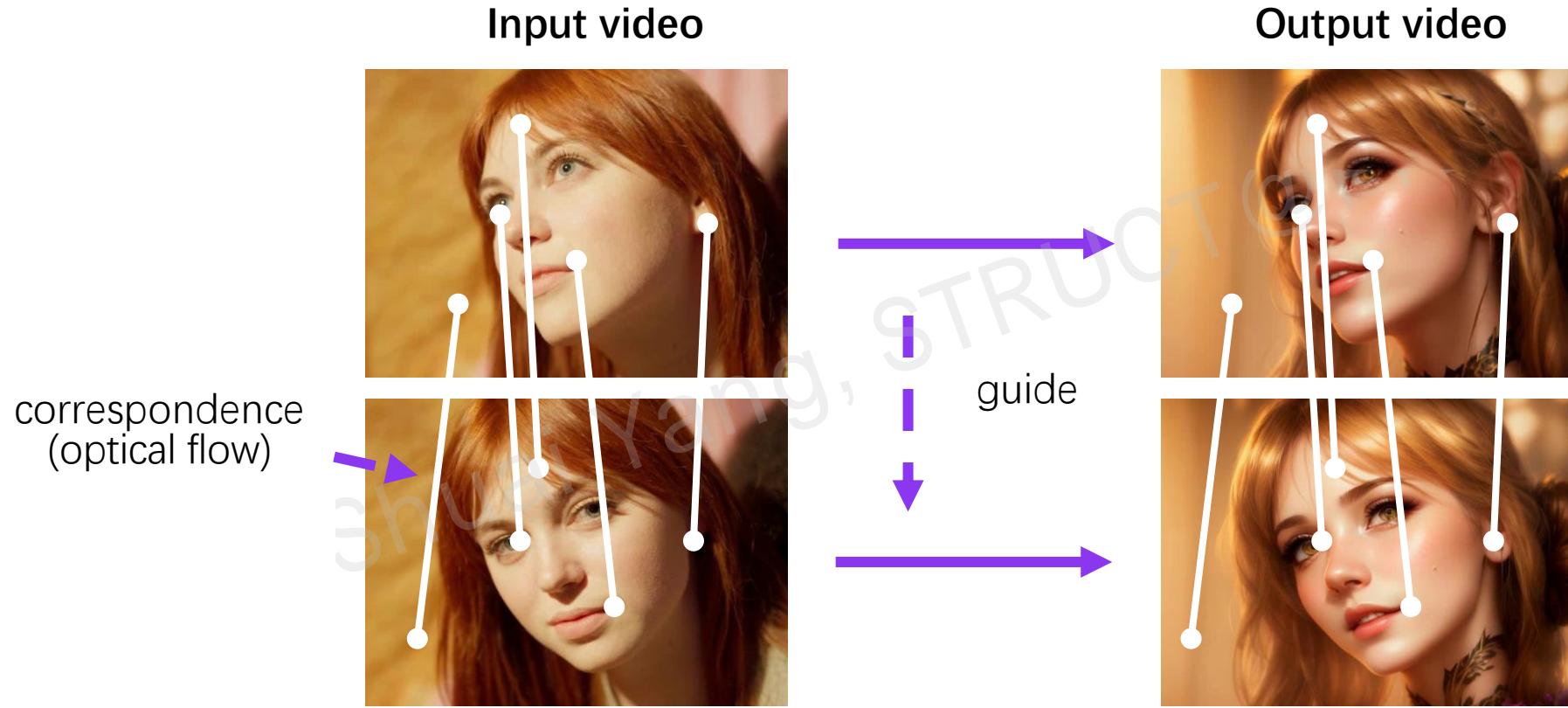
Inversion-free zero-shot model

3 Rerender-A-Video



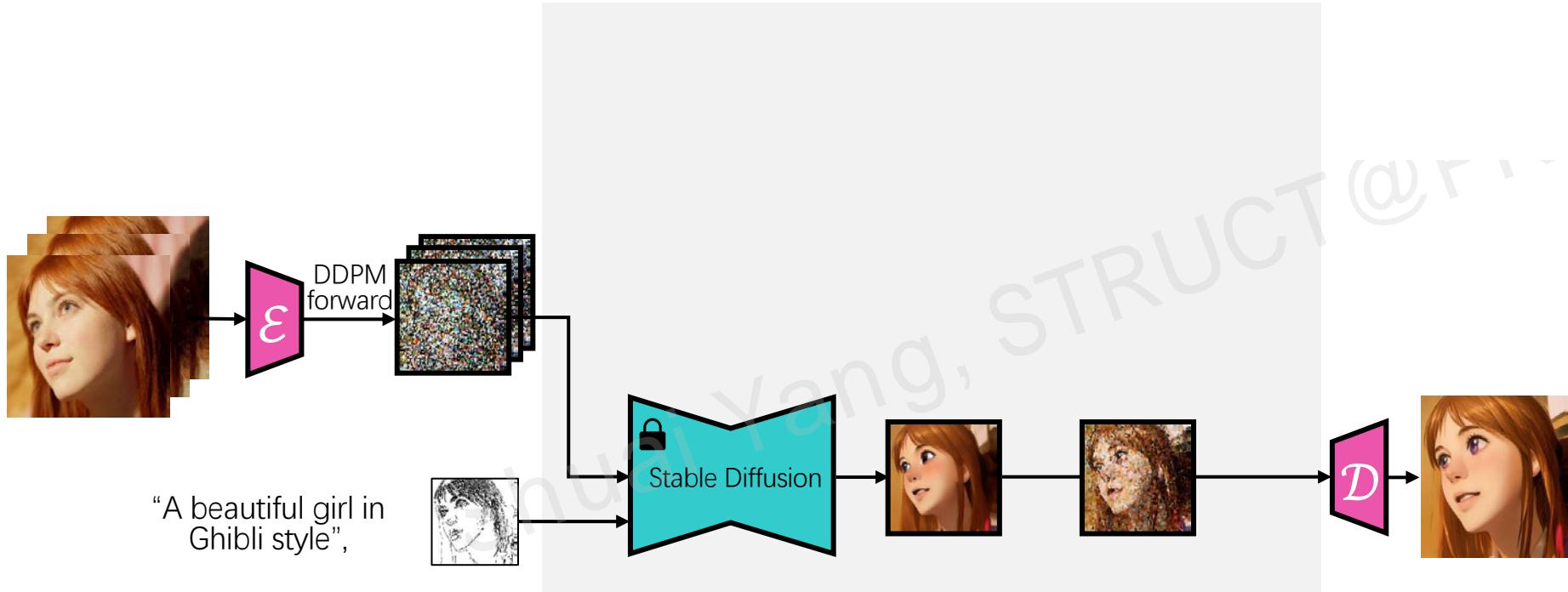
Inversion-free zero-shot model

3 Rerender-A-Video



Inversion-free zero-shot model

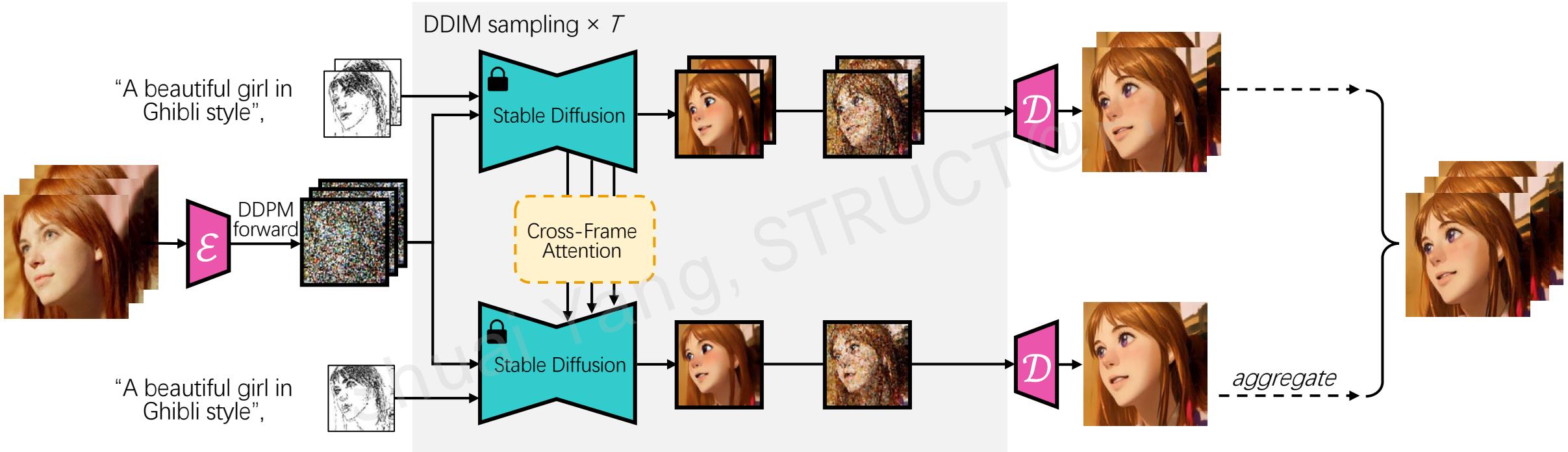
3 Rerender-A-Video



ControlNet+SDEdit image translation

Inversion-free zero-shot model

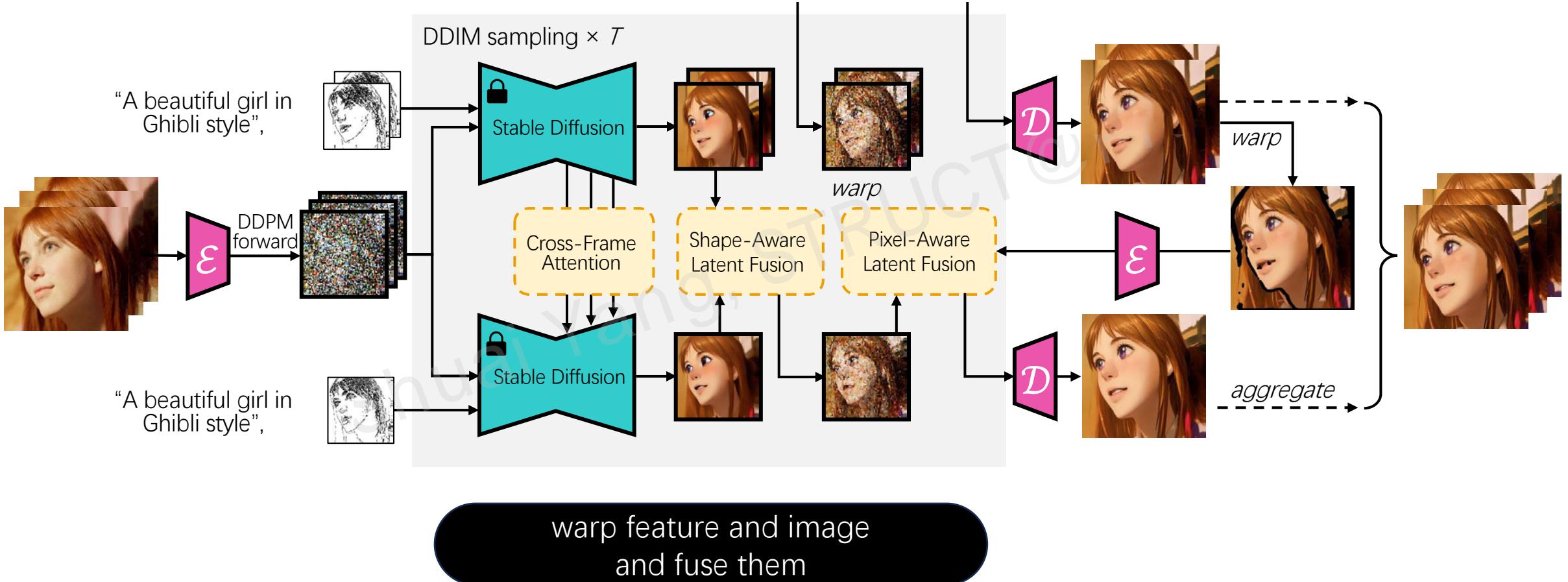
3 Rerender-A-Video



ControlNet+SDEdit image translation
with cross-frame attention

Inversion-free zero-shot model

3 Rerender-A-Video



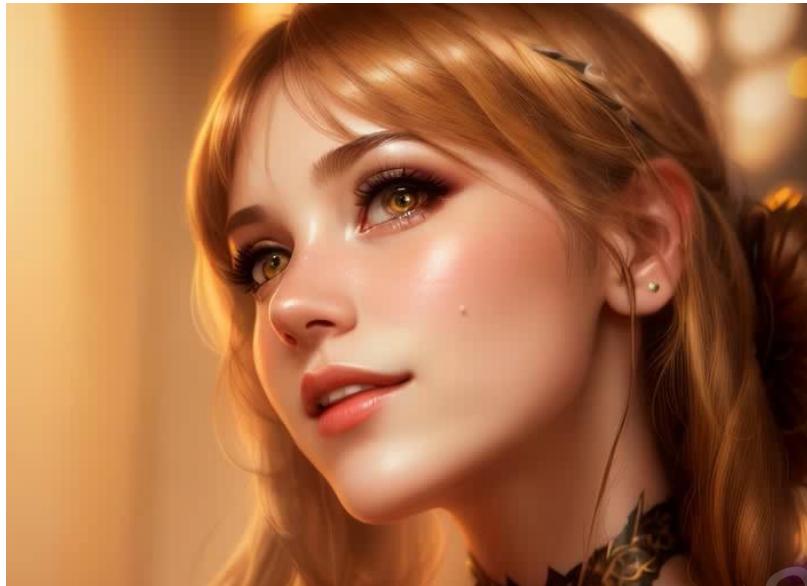
Inversion-free zero-shot model

3 Rerender-A-Video

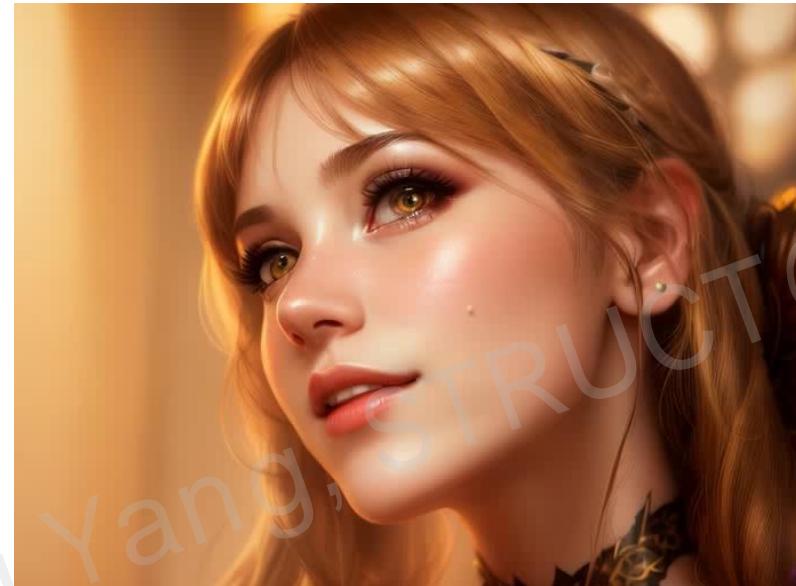


Inversion-free zero-shot model

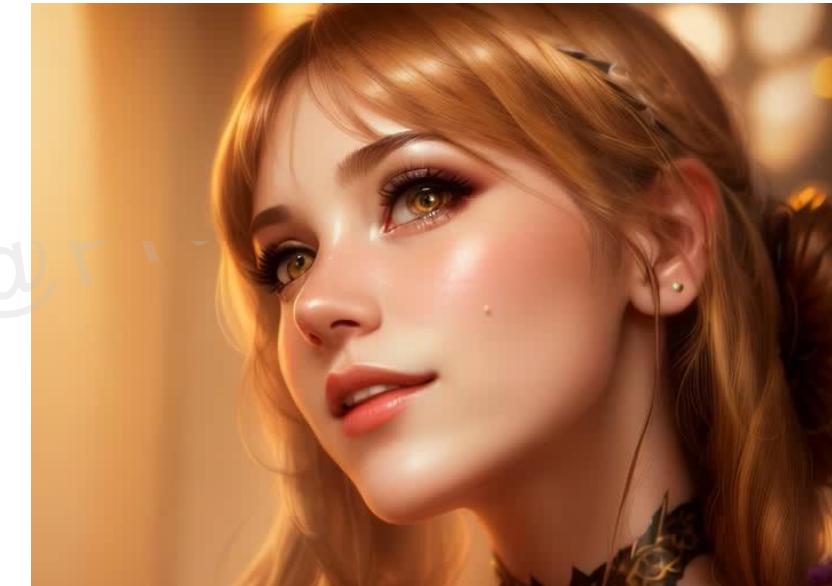
3 Rerender-A-Video



Stable Diffusion



+ Cross-frame attn

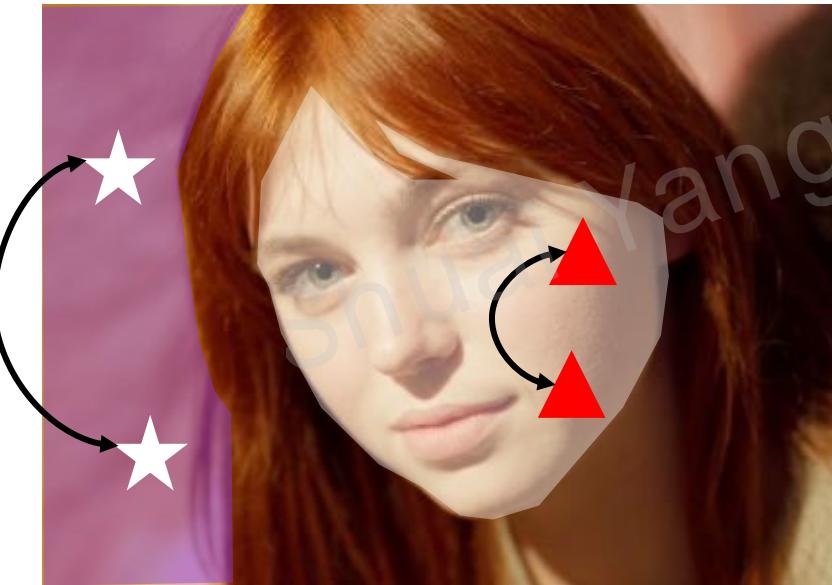


+ Optical flow

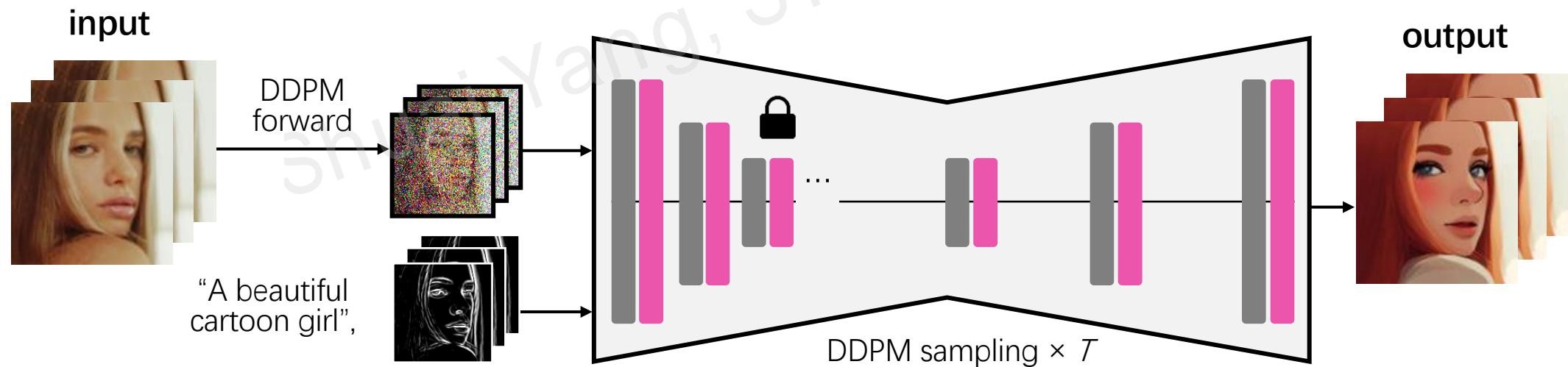
Optical flow helps a lot
but when optical flow is hard to predict, what can we do?

While previous methods primarily focus on constraining inter-frame temporal correspondence,
preserving intra-frame spatial correspondence is equally crucial

semantically similar content is manipulated cohesively,
maintaining its similarity post-translation

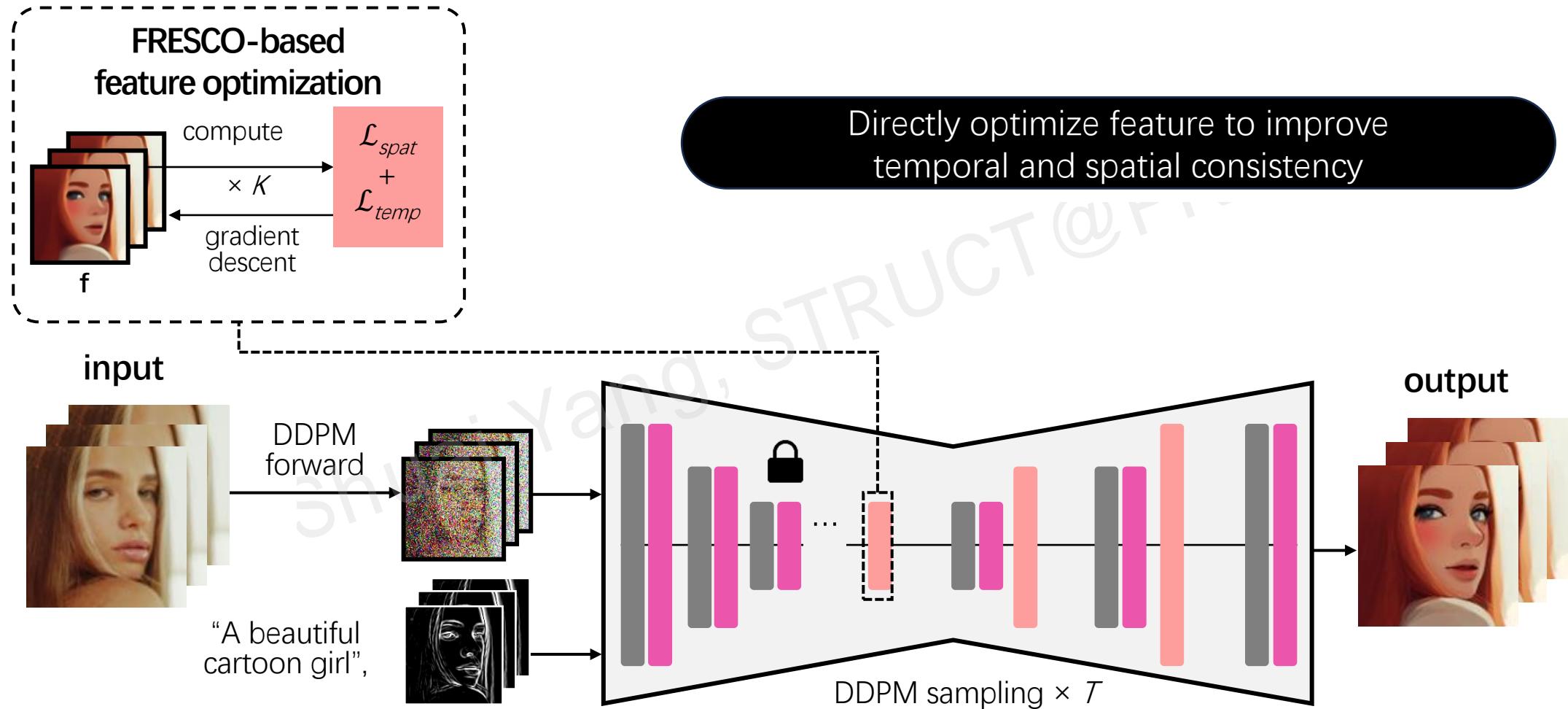


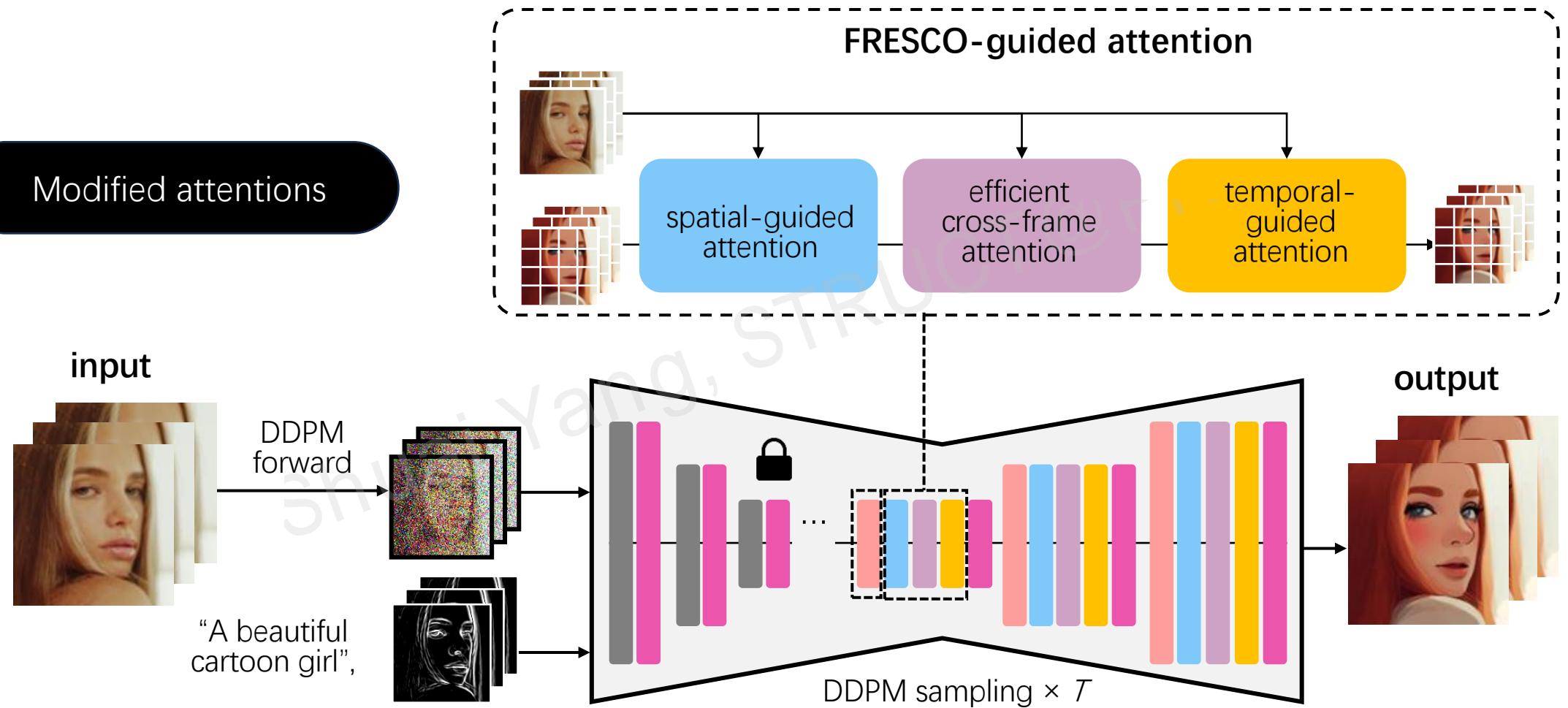
ControlNet+SDEdit image translation



Inversion-free zero-shot model

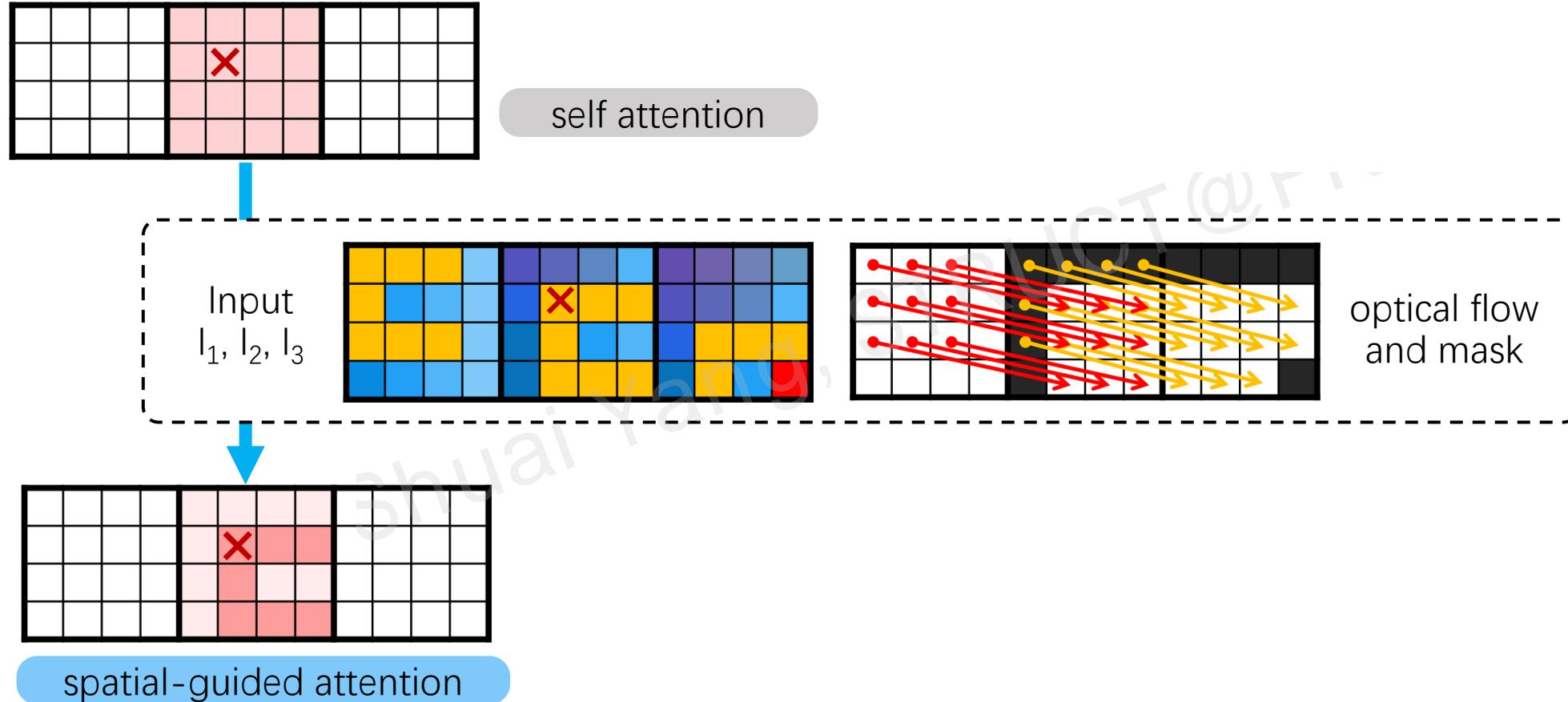
4 FRESCO





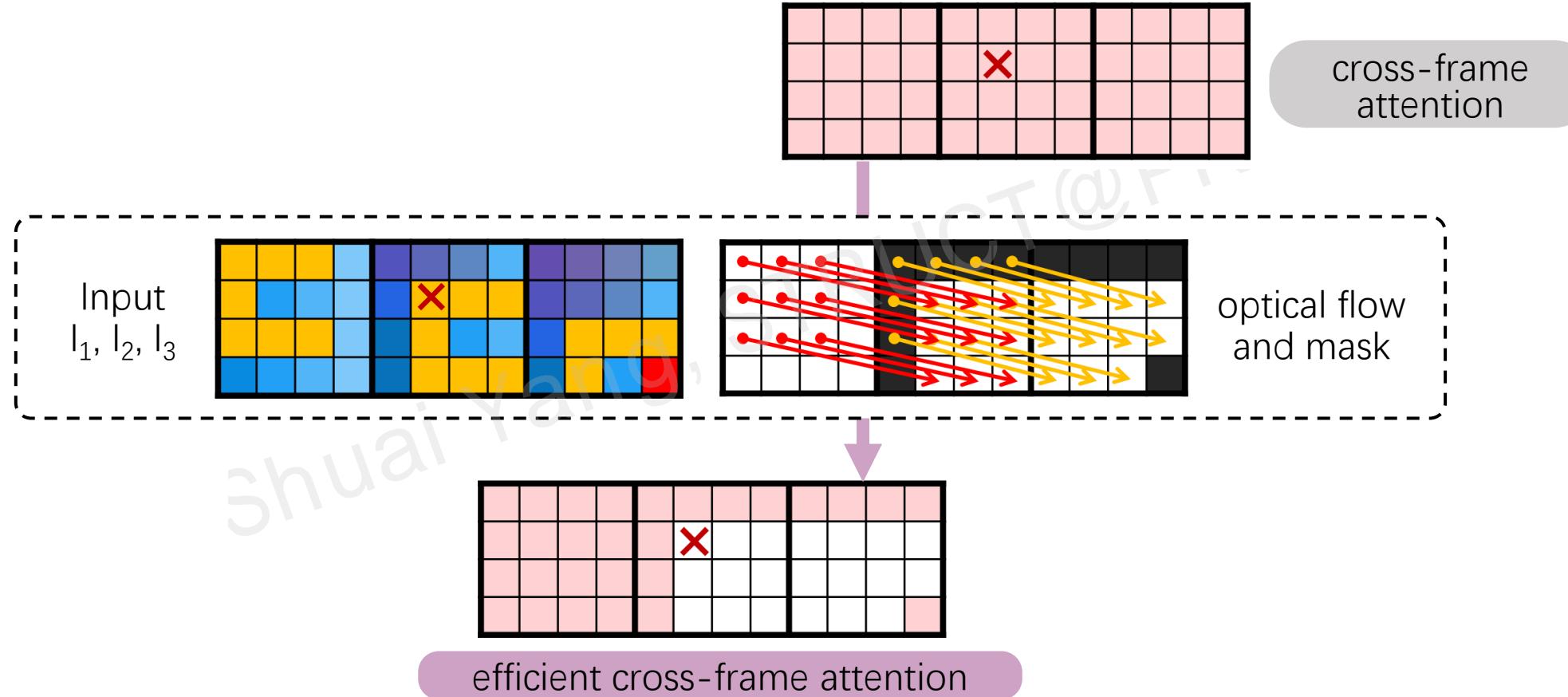
Inversion-free zero-shot model

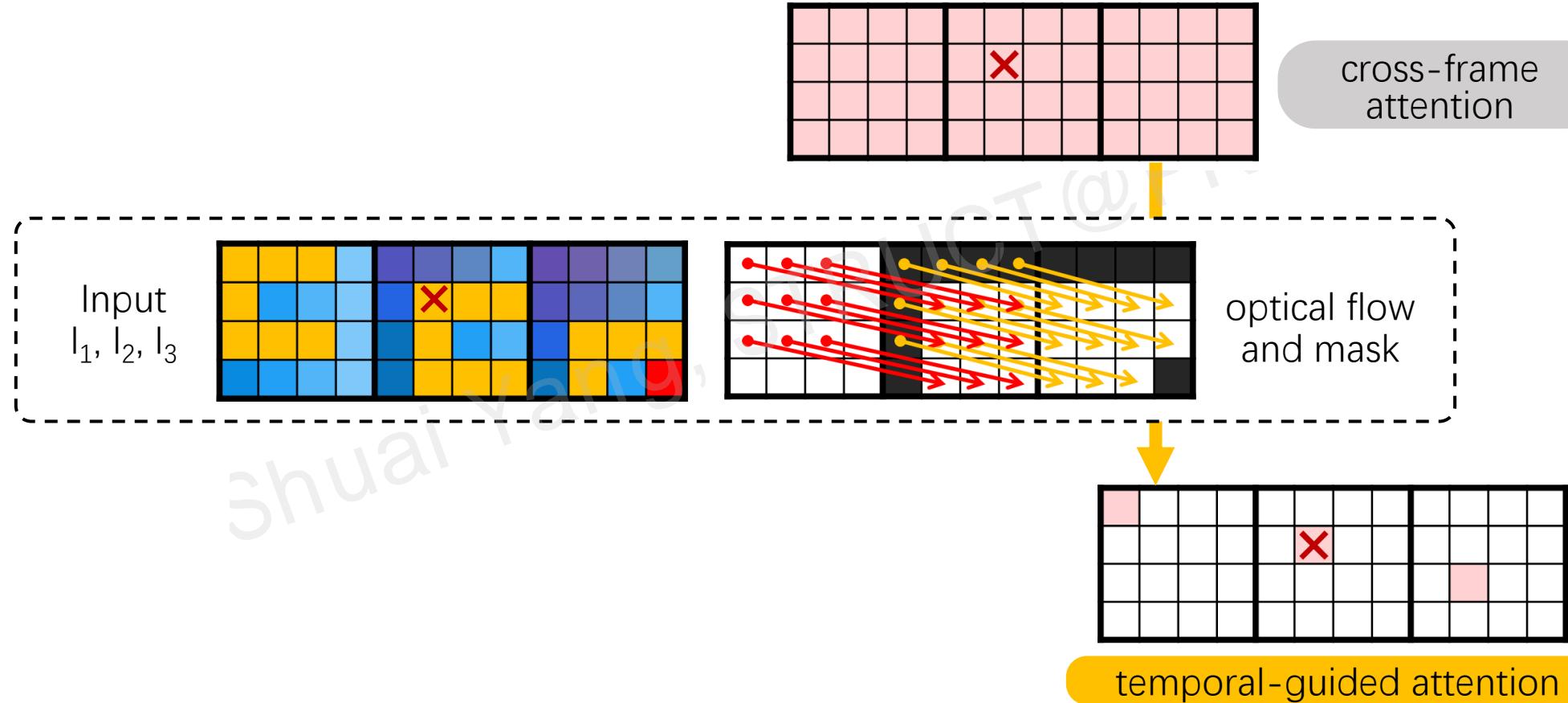
4 FRESCO



Inversion-free zero-shot model

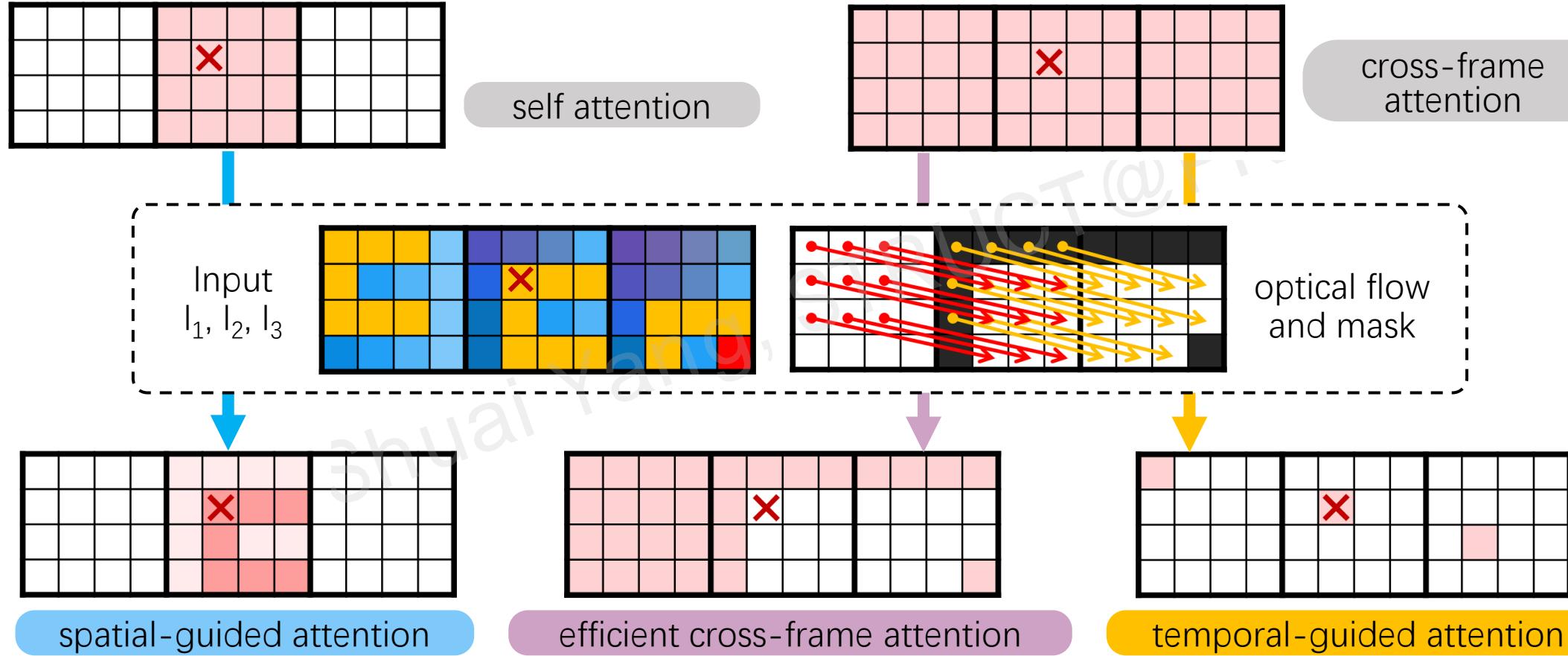
4 FRESCO





Inversion-free zero-shot model

4 FRESCO





a white cat in pink background, cartoon style



input



Text2Video-Zero



ControlVideo

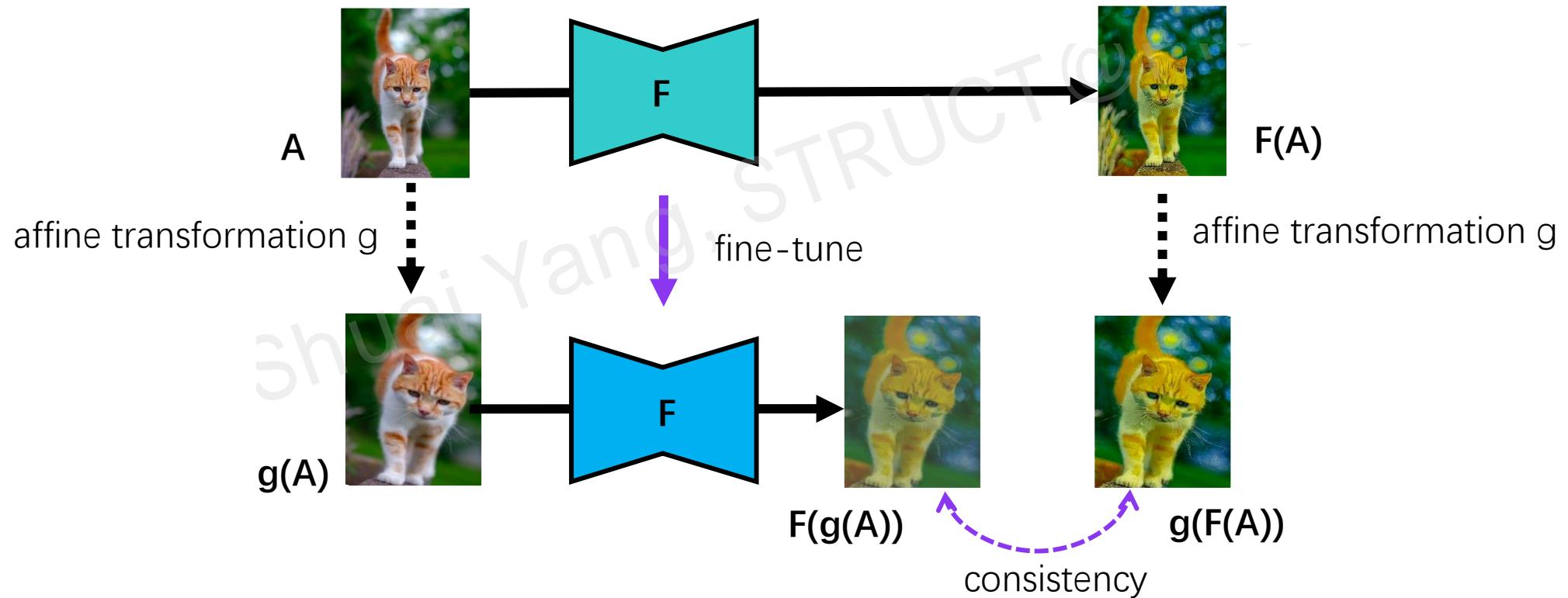


Rerender-A-Video



FRESCO

While we are always fighting against the noises,
maybe we can allow some fine-tuning to make the model less sensitive to the noises



Inversion-free zero-shot model

5 Fairy

eat fox walking



input



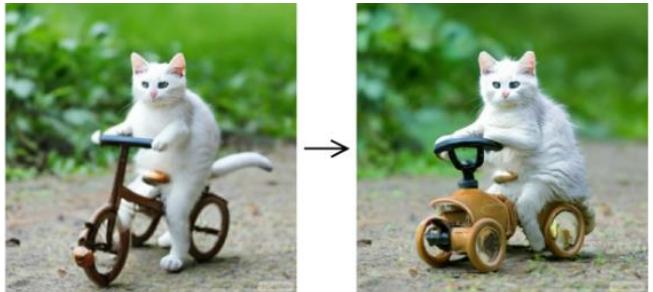
w/o fine-tune



full

Zero-shot model

Image Editing



“Photo of a cat riding on a bicycle.”
car



Inversion-Based Model



ICCV'23
ICLR'24
CVPR'24
FateZero
TokenFlow
VidToMe
Pixel2Video
FLATTEN

Conditional Generation



Inversion-Free Model



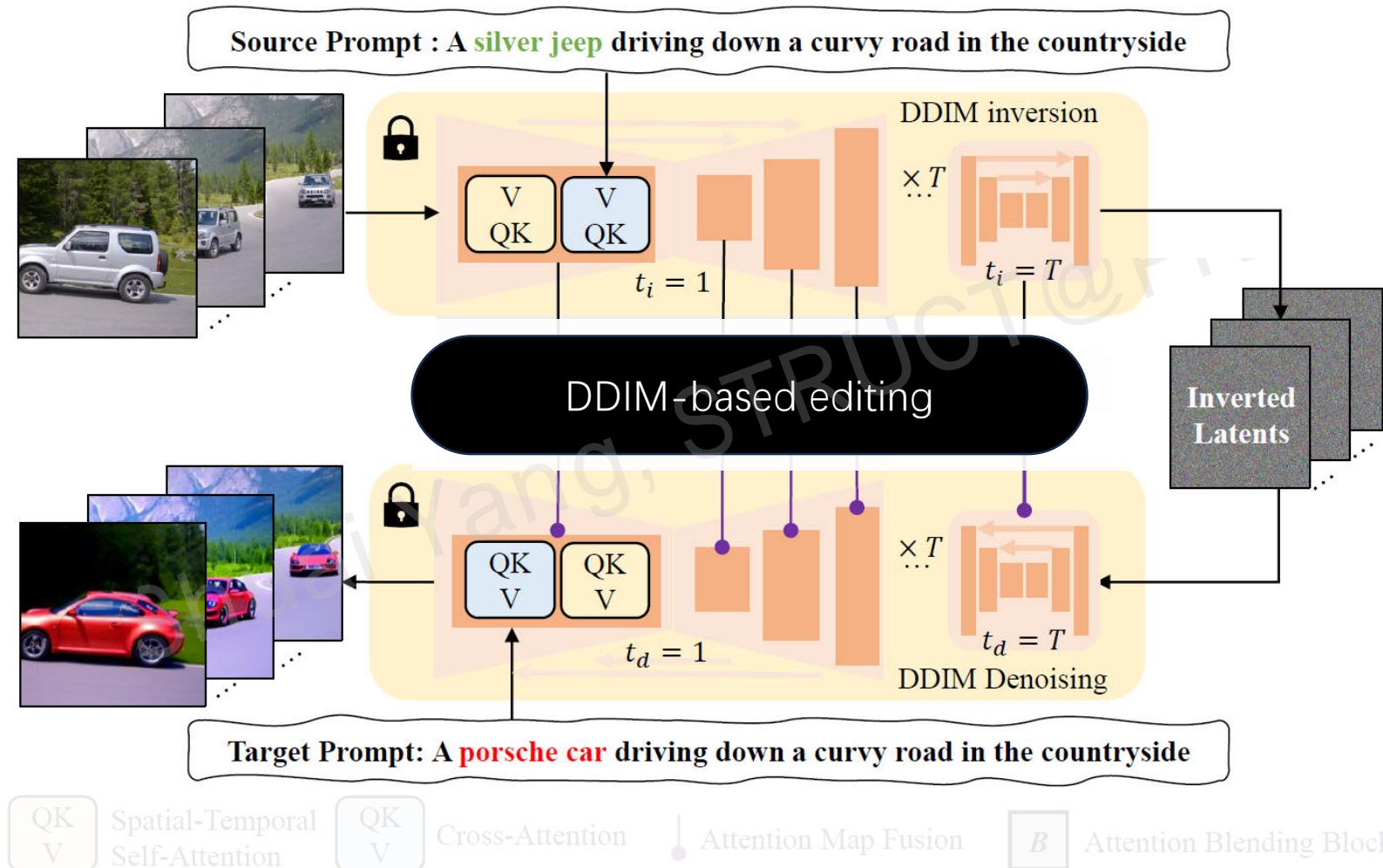
ICCV'23
SIGGRAPH Asia'23
ICLR'24
CVPR'24
Text2Video-Zero
Rerender-A-Video
ControlVideo
FRESCO Fairy

ControlNet
SDEdit



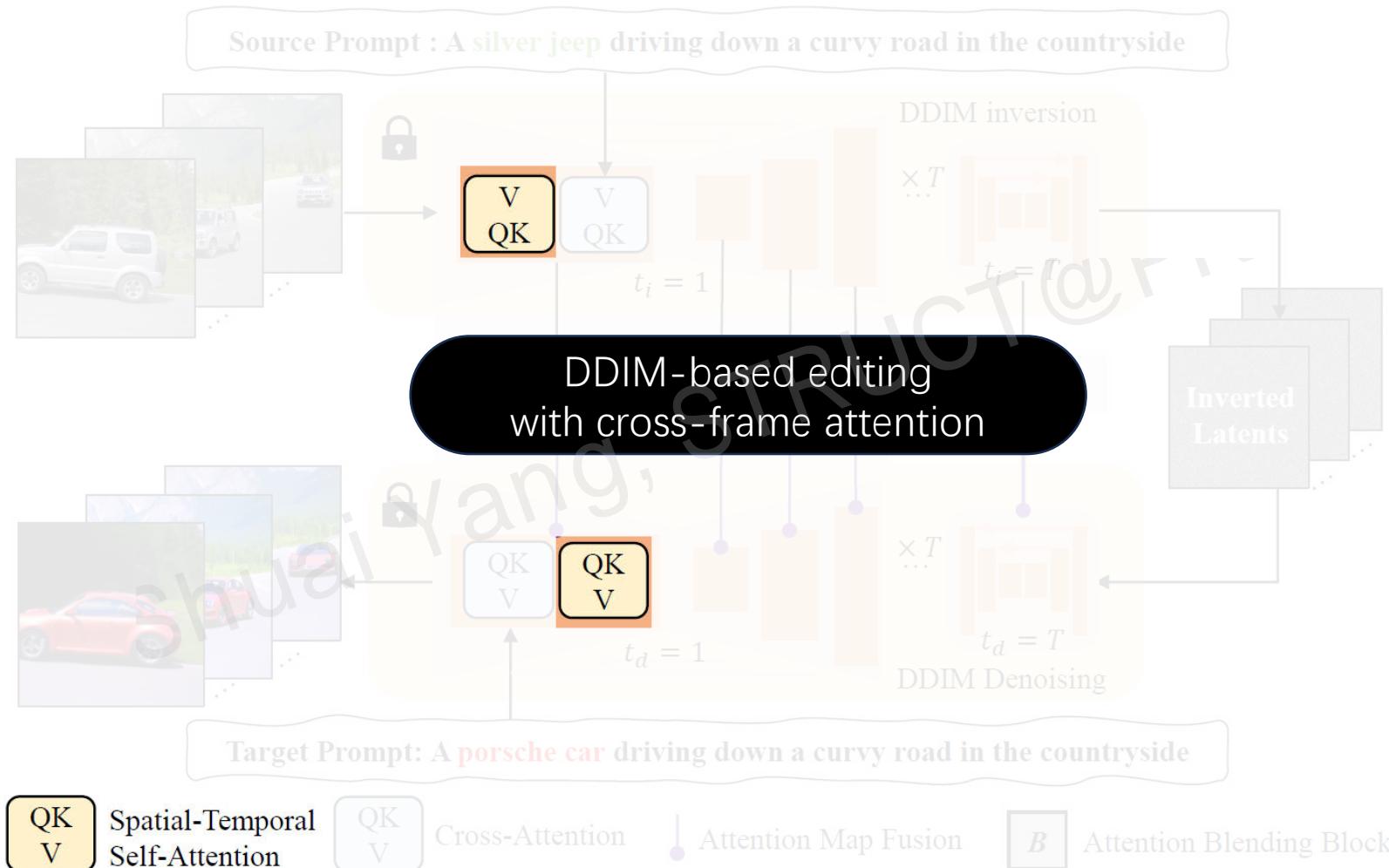
Inversion-based zero-shot model

1 FateZero



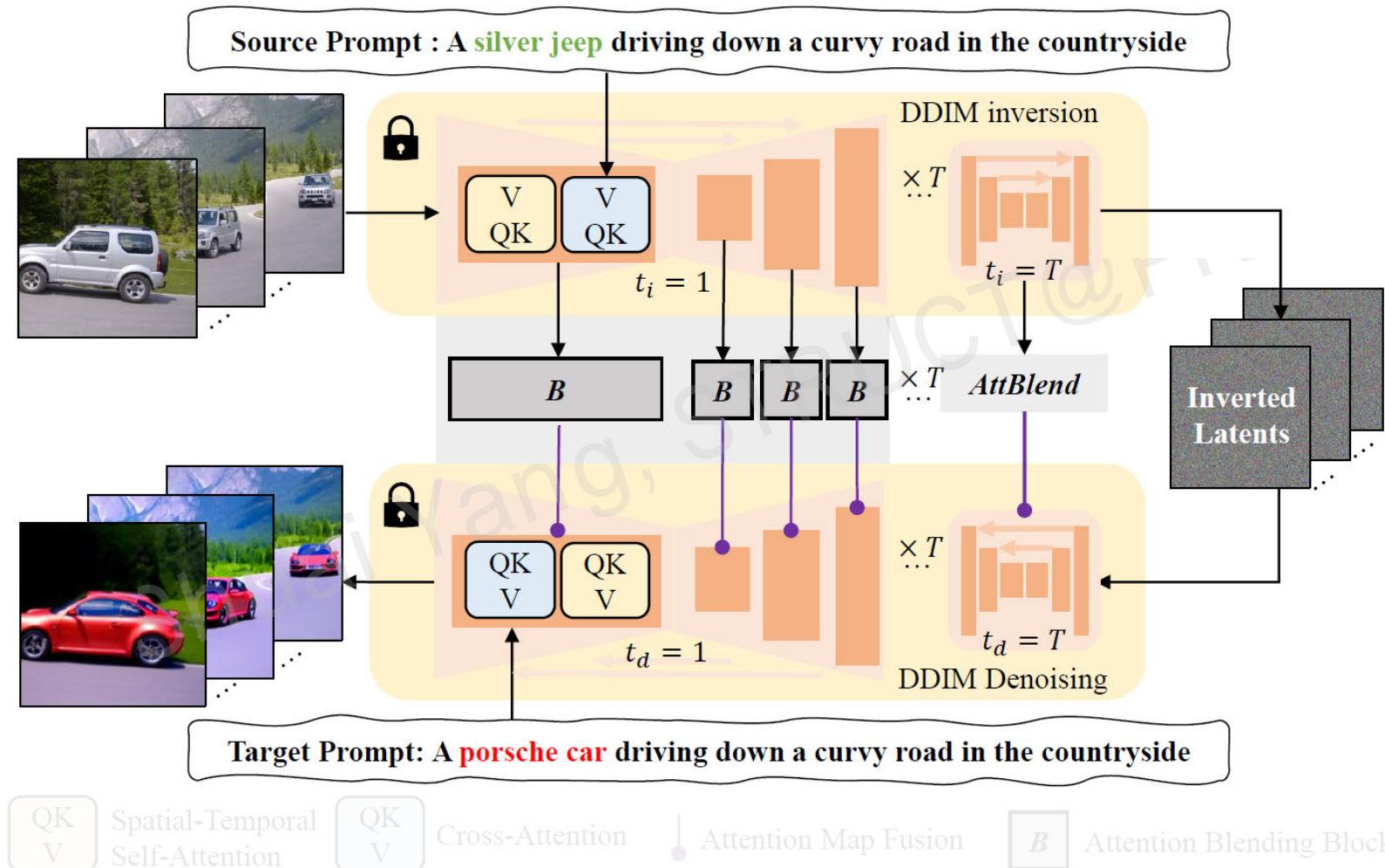
Inversion-based zero-shot model

1 FateZero



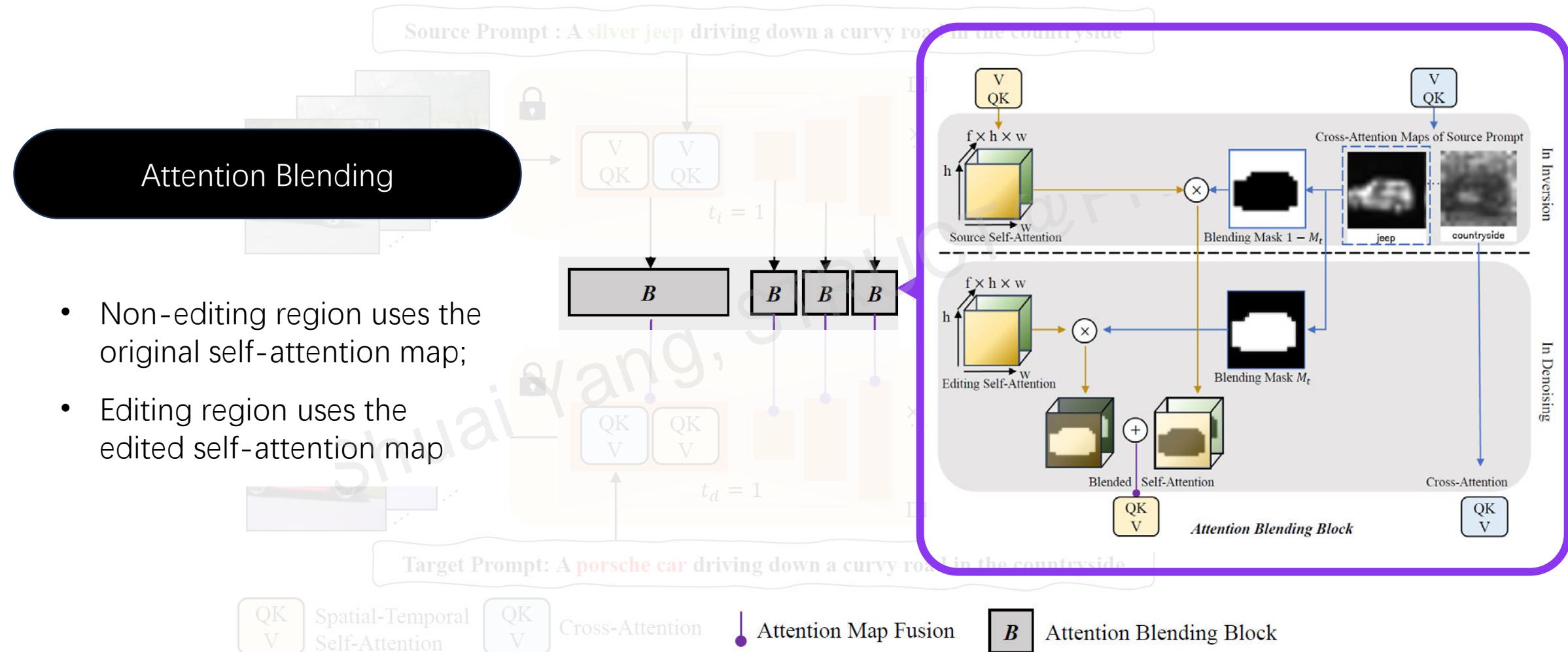
Inversion-based zero-shot model

1 FateZero



Inversion-based zero-shot model

1 FateZero



Inversion-based zero-shot model

1 FateZero

Car
Porsche Car



+
*Watercolor
Painting*



Fate-Zero

Tune-A-Video

SDEdit

Null-Text Inversion

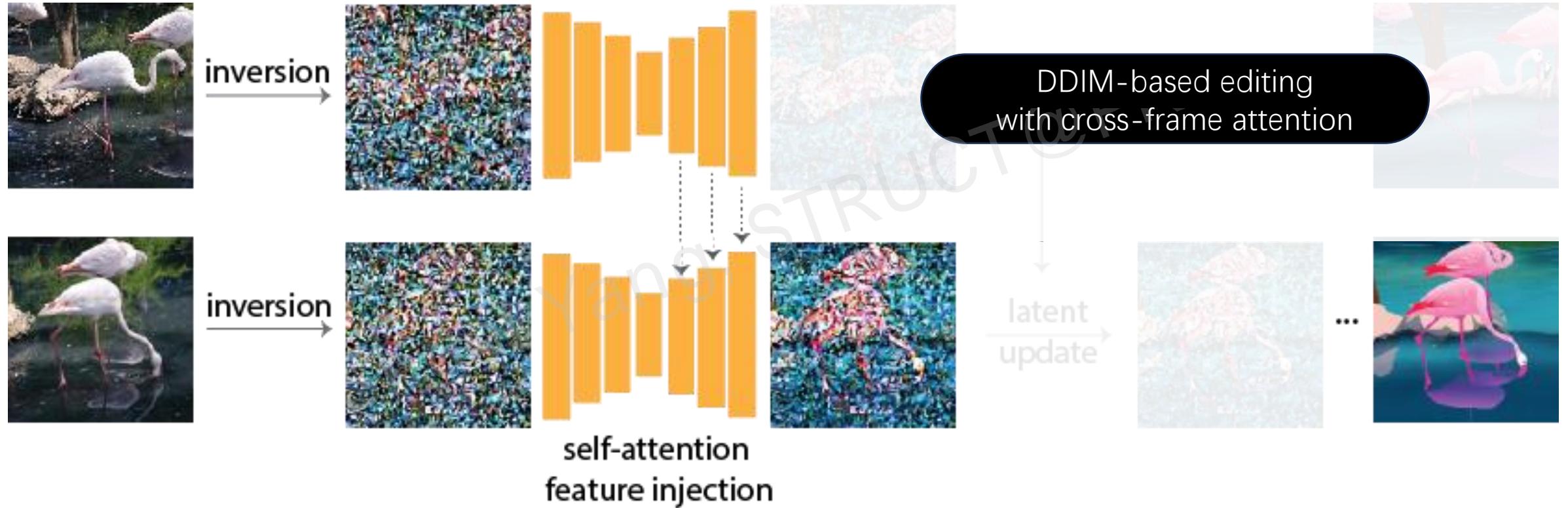
Inversion-based zero-shot model

2 Pix2Video



Inversion-based zero-shot model

2 Pix2Video



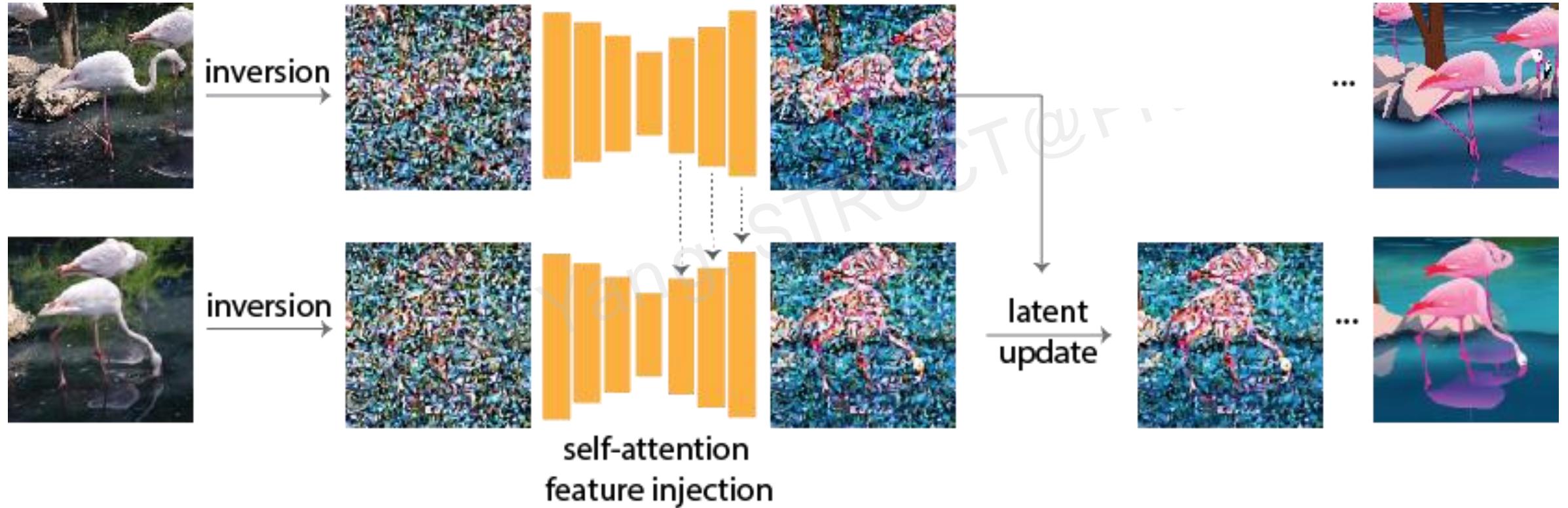
Inversion-based zero-shot model

2 Pix2Video



Inversion-based zero-shot model

2 Pix2Video



a group of pigs standing in the dirt, claymation



Input



Prompt-to-Prompt



Tune-A-Video



Pix2Video

Previous methods use inversion features globally via attention.
Can we use them more locally?



Input

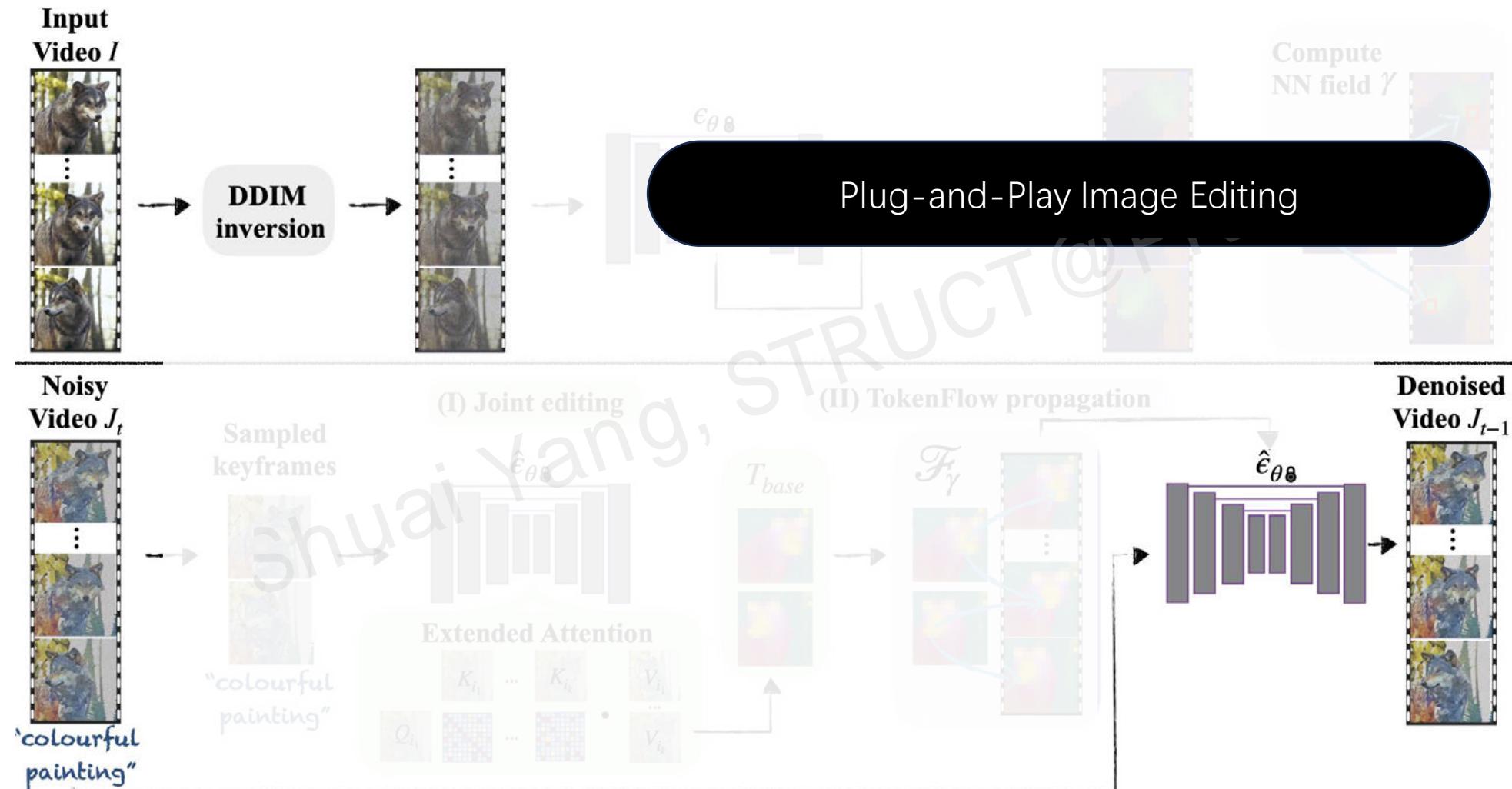


Visualized Diffusion feature

Diffusion features are semantically consistent.
Coherent features → coherent output!

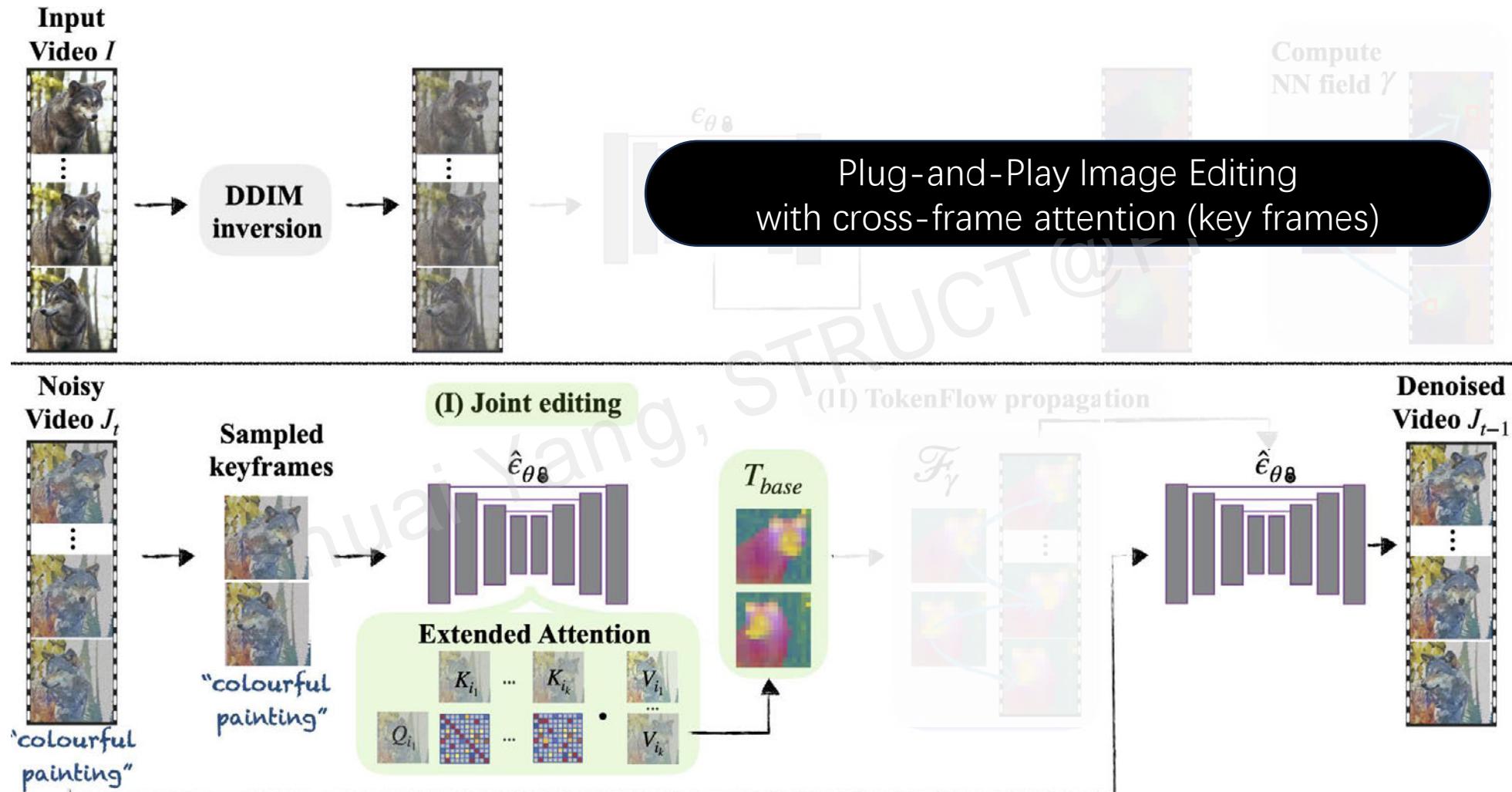
Inversion-based zero-shot model

3 TokenFlow



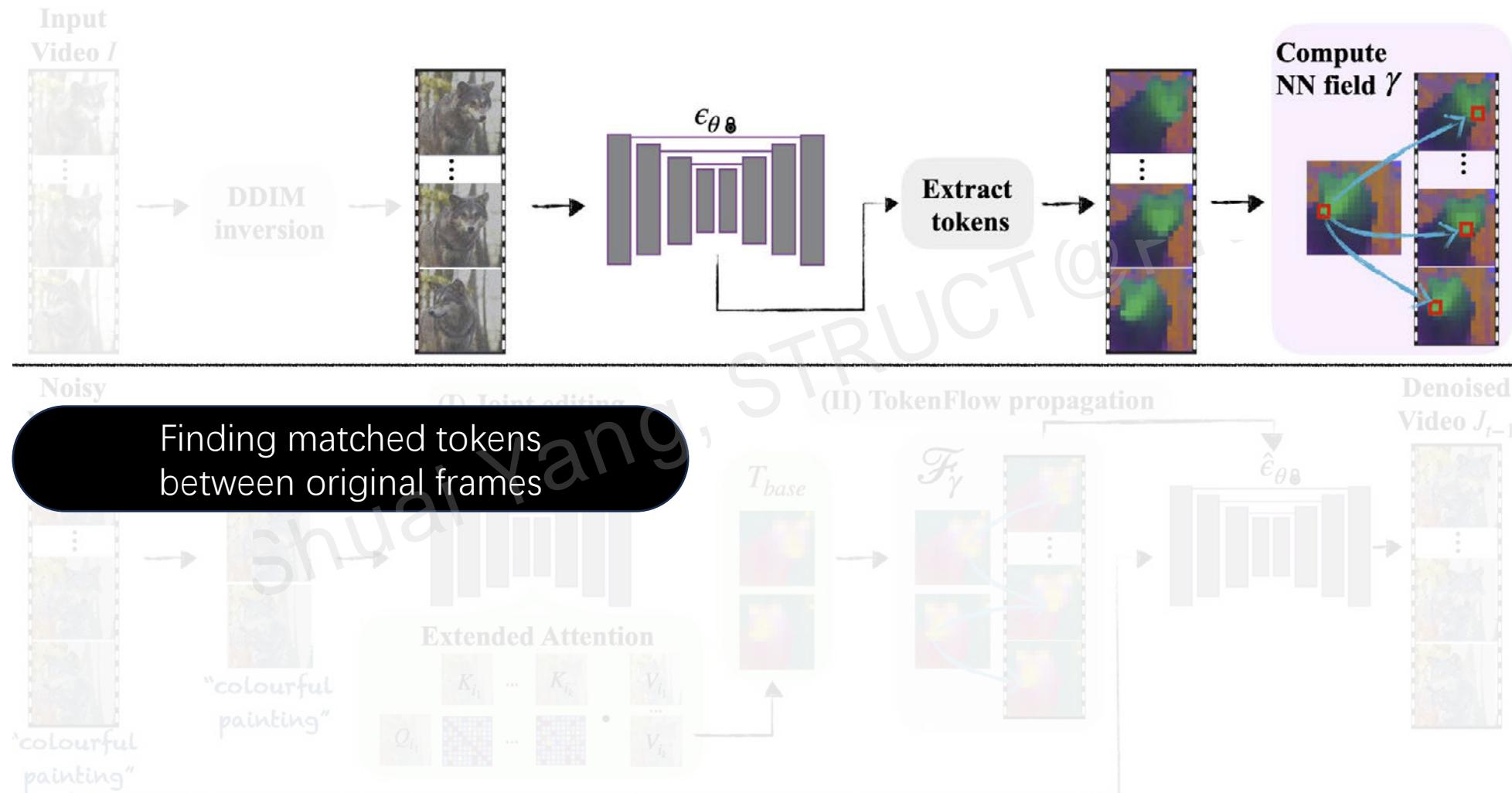
Inversion-based zero-shot model

3 TokenFlow



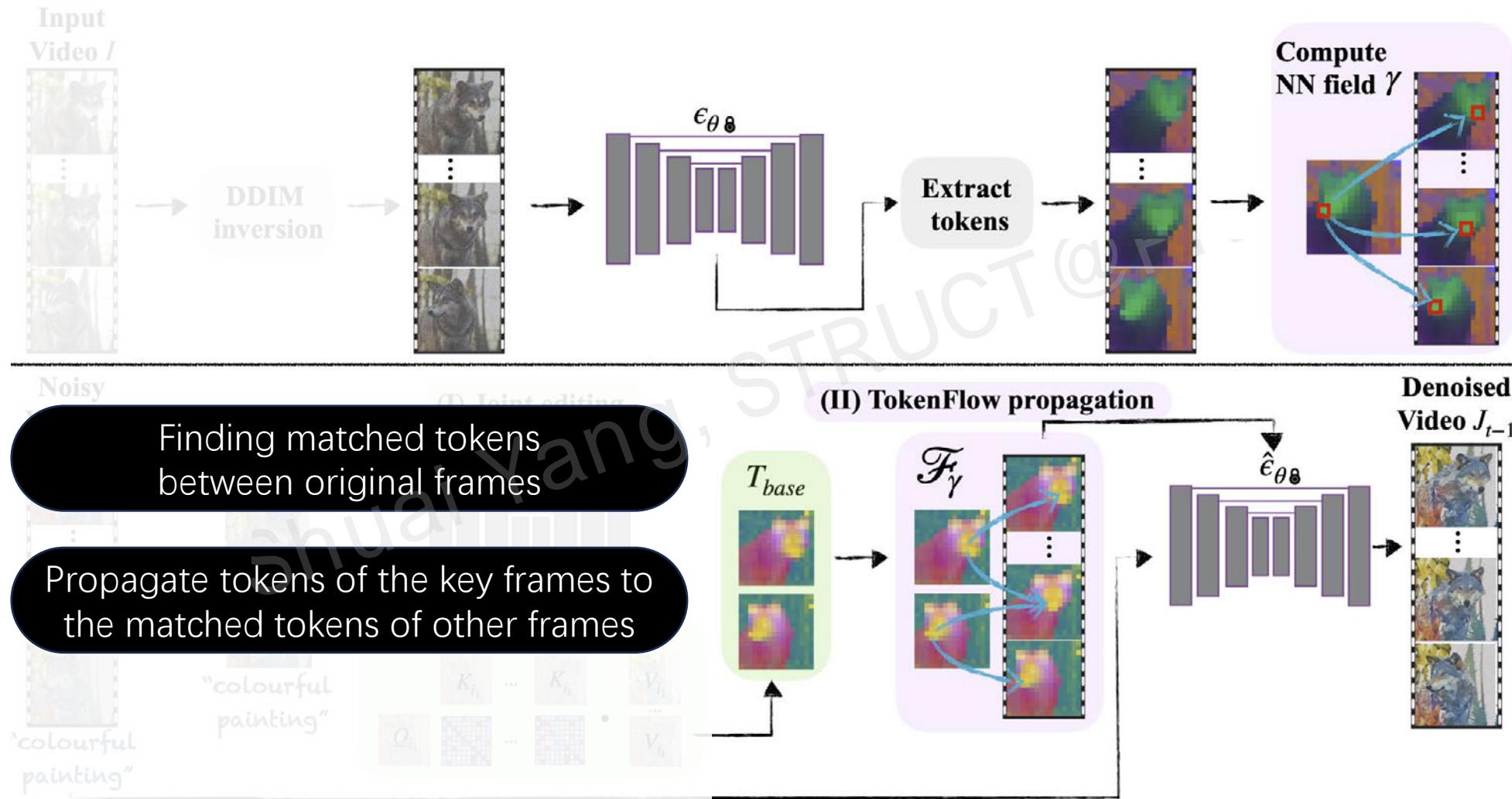
Inversion-based zero-shot model

3 TokenFlow



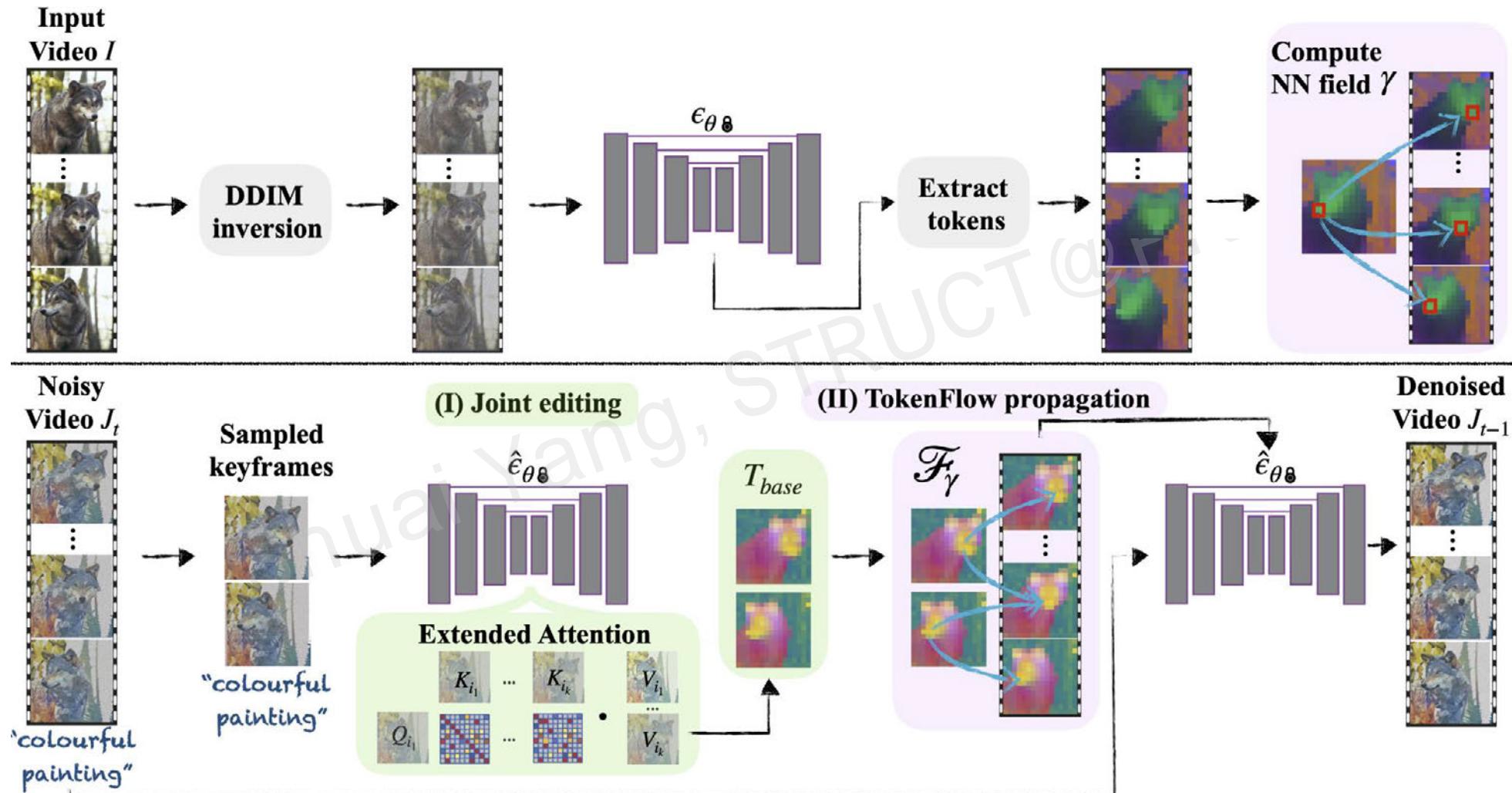
Inversion-based zero-shot model

3 TokenFlow



Inversion-based zero-shot model

3 TokenFlow



Inversion-based zero-shot model

3 TokenFlow

Input



PnP



TokenFlow



Visualized features in different layers

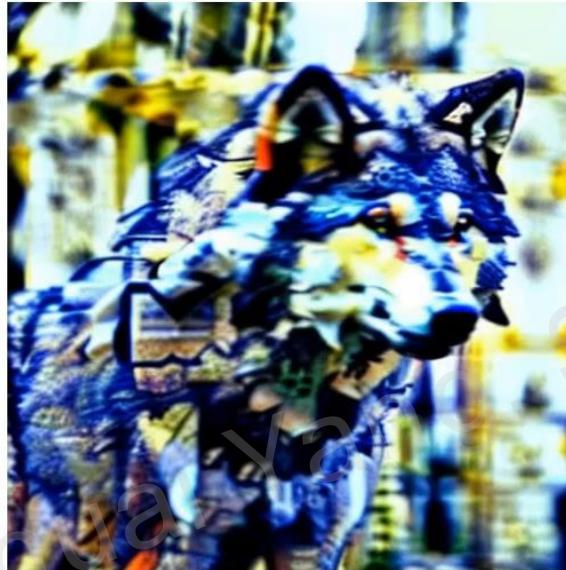
Inversion-based zero-shot model

3 TokenFlow

a fluffy wolf doll



Input



Tune-A-Video



FateZero

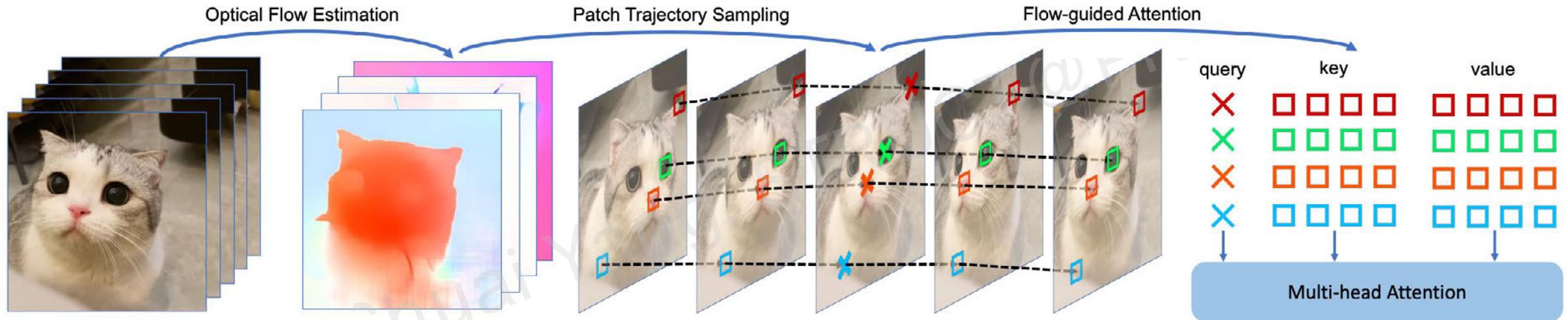


TokenFlow

Locally matched tokens can significantly improve the consistency.
What about finer-grained optical flow?

Inversion-based zero-shot model

4 FLATTEN

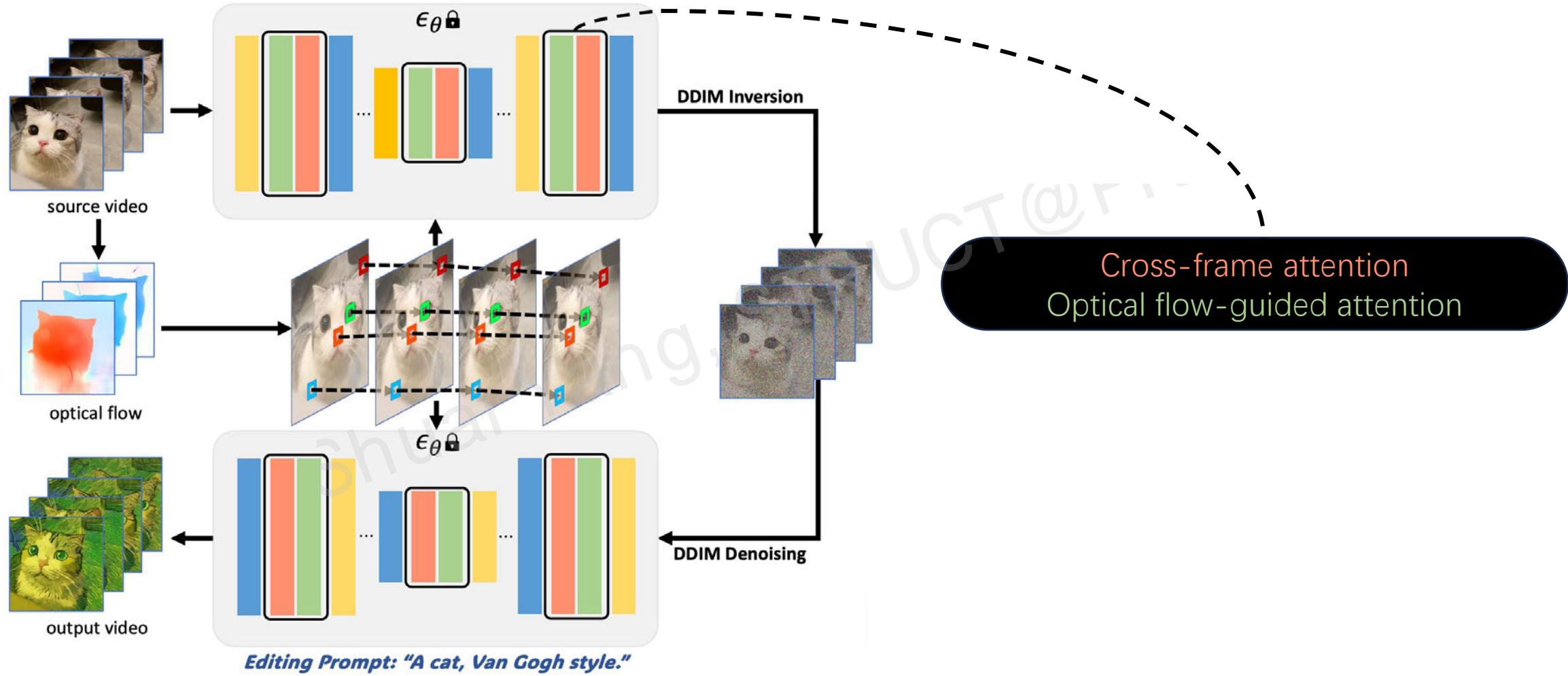


FLATTEN: optical FLow-guided ATTENtion

FLATTEN: Optical Flow-guided Attention for Consistent Text-to-Video Editing. ICLR'24

Inversion-based zero-shot model

4 FLATTEN



Wooden trucks drive on a racetrack



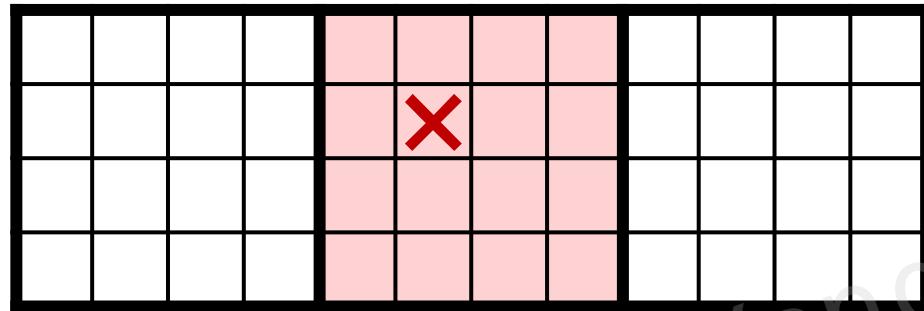
Input

FateZero

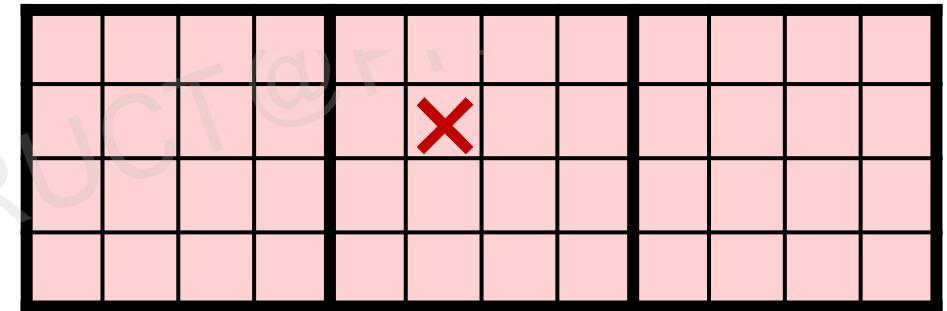
TokenFlow

FLATTEN

Attention mask or feature blending may harm editability.
Only FLATTEN changes the appearance.



Self attention

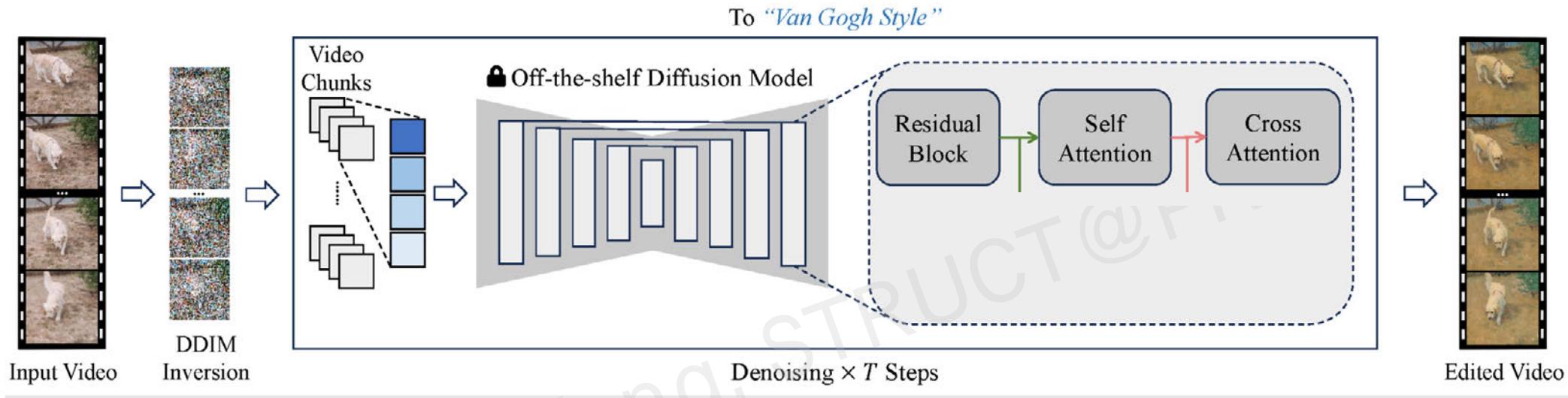


Cross-frame attention

With B frames, computation increase $B \times B$ times
Redundant tokens can be merged!

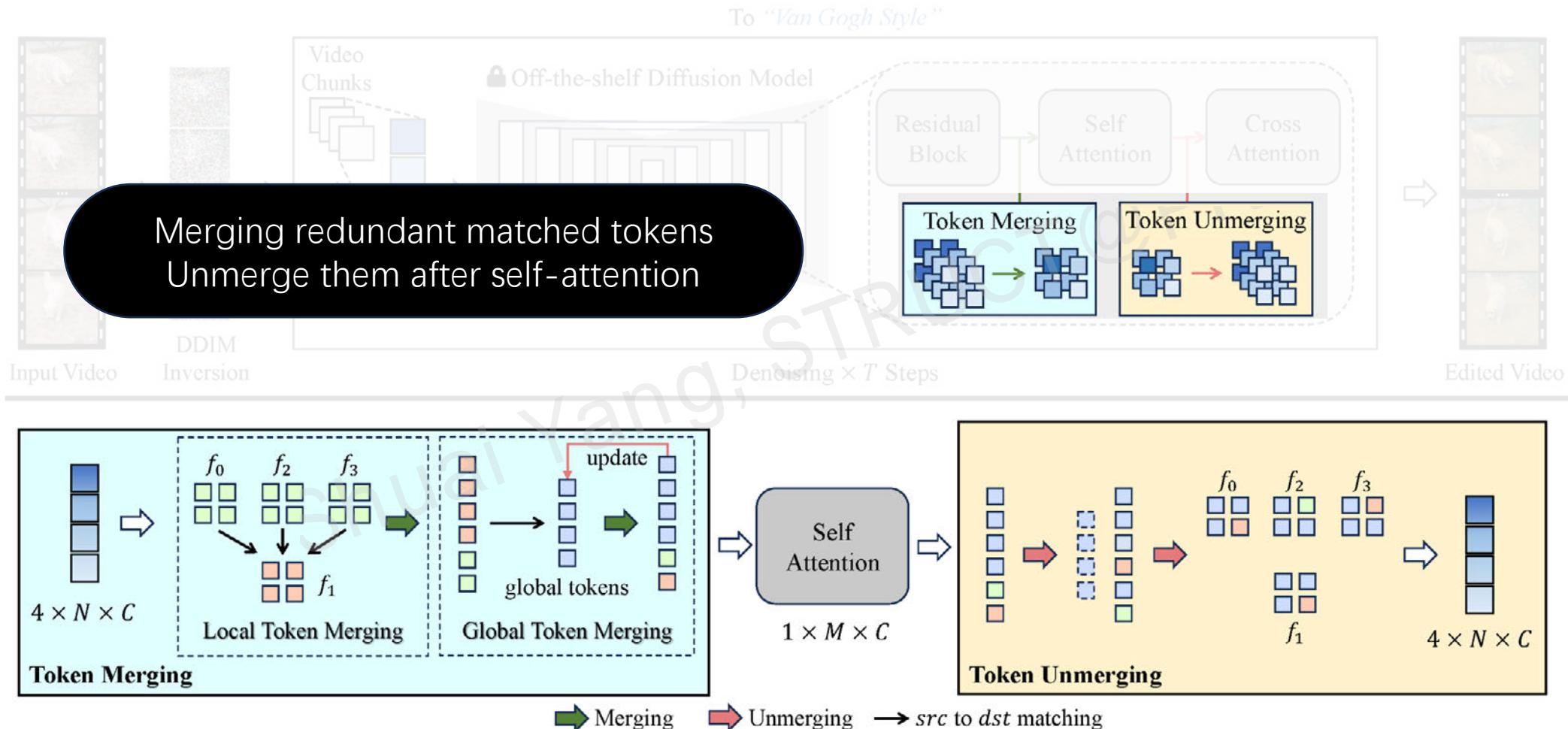
Inversion-based zero-shot model

5 VidToMe



Inversion-based zero-shot model

5 VidToMe





Zero-shot model

Image Editing



Inversion-Based Model



*a white cat in pink
background, cartoon style*

+ inversion feature fusion

Conditional Generation



Inversion-Free Model



*a white cat in pink
background, cartoon style*

+ depth condition



SUMMARY

- Big models and big data bring about key emergent capabilities
- Fundamental big models offer a lot of opportunities for creative work
- Even with limited resources, we can still use the fundamental big models for various downstream tasks
- Future: long and consistent video generation
 - make everyone a creator



A black semi-transparent rectangular overlay contains white text. The text reads "Thank you for listening!" in a large, bold, sans-serif font. Below this, in a smaller font, is the email address "williamyang@pku.edu.cn". A faint watermark reading "Shuai Yang, SJTU" is visible across the center of the image.

Thank you for listening!

williamyang@pku.edu.cn

STRUCT Group

Intelligent image computing

Wangxuan Institute of Computer Technology, Peking University

Spatial and Temporal **Restoration**,
Understanding and **Compression** Team

PI: Jiaying Liu

- Email: liujiaying@pku.edu.cn
- Website: <http://www.wict.pku.edu.cn/struct/>

