

# Efficient generative AI for images and video

Amirhossein Habibian

Director of Engineering  
Qualcomm Technologies Netherlands



# Agenda

A woman with dark hair pulled back, wearing a light-colored patterned blouse and beige trousers, is seated in a black leather chair. She is looking out a large window at a city skyline at night, holding a smartphone in her hands. The scene is dimly lit, with the city lights providing the primary illumination.

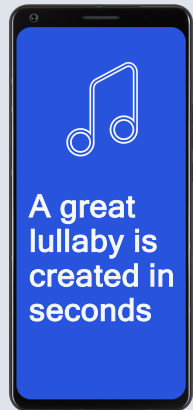
- The potential impact of efficient generative vision
- Efficient image generation
- Efficient video generation
- Efficient 3D generation
- Applications: automotive
- Q&A

## Text generation (ChatGPT, Bard, Llama, etc.)



Input prompts

“Write a lullaby about cats and dogs to help a child fall asleep, include a golden shepherd”



### Real-life application

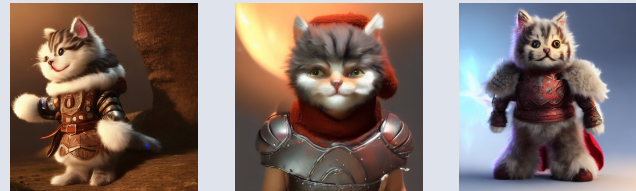
- Communications
- Journalism
- Publishing
- Creative writing
- Writing assistance

## Image generation (Stable Diffusion, MidJourney, etc.)



Input prompts

“Super cute fluffy cat warrior in armor”



### Real-life application

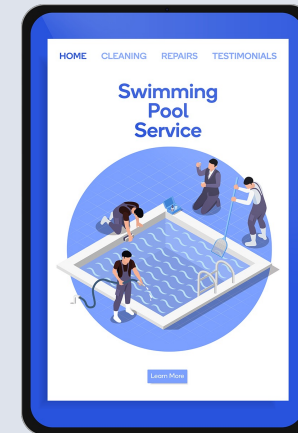
- Advertisements
- Corporate visuals
- Published illustrations
- Novel image generation

## Code generation (Codex, etc.)



Input prompts

“Create code for a pool cleaning website with tab for cleaning, repairs, and testimonials”



A beautiful website is created in seconds

### Real-life application

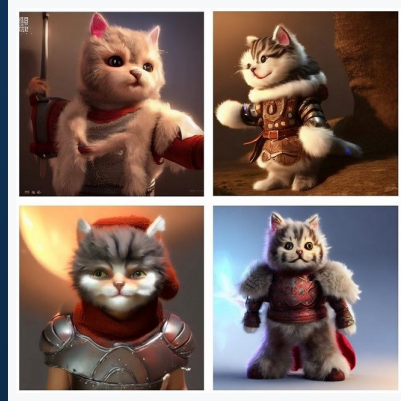
- Web design
- Software development
- Coding
- Technology

# What is generative AI?

AI models that create new and original content like text, images, video, audio, or other data

Generative AI, foundational models, and large language models are sometimes used interchangeably

# Why is generative AI for computer vision important?



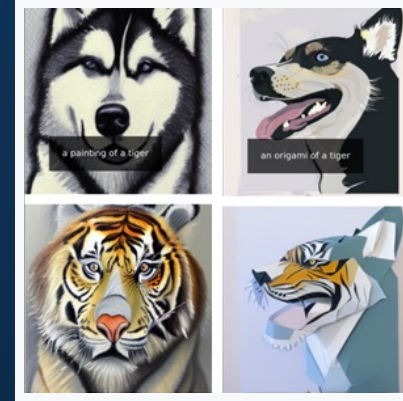
## Generating images and videos

Generative models create images and videos from scratch

Original, life-like visuals generated from textual and/or image prompts (in the case of image/text-to-image or image/text-to-video)

**Examples:**  
Stable Diffusion, ControlNet

Prompt: "Super cute fluffy cat warrior in armor, photorealistic, 4K, ultra detailed, vray rendering, unreal engine"



## Editing images and videos

Generative models change aspect of images and videos

Swap the background, change style or edit object's attribute and appearance

**Examples:**  
SDEdit, PnP, Pix2Pix



## Generating 3D content

Generative models create 3D meshes and assets

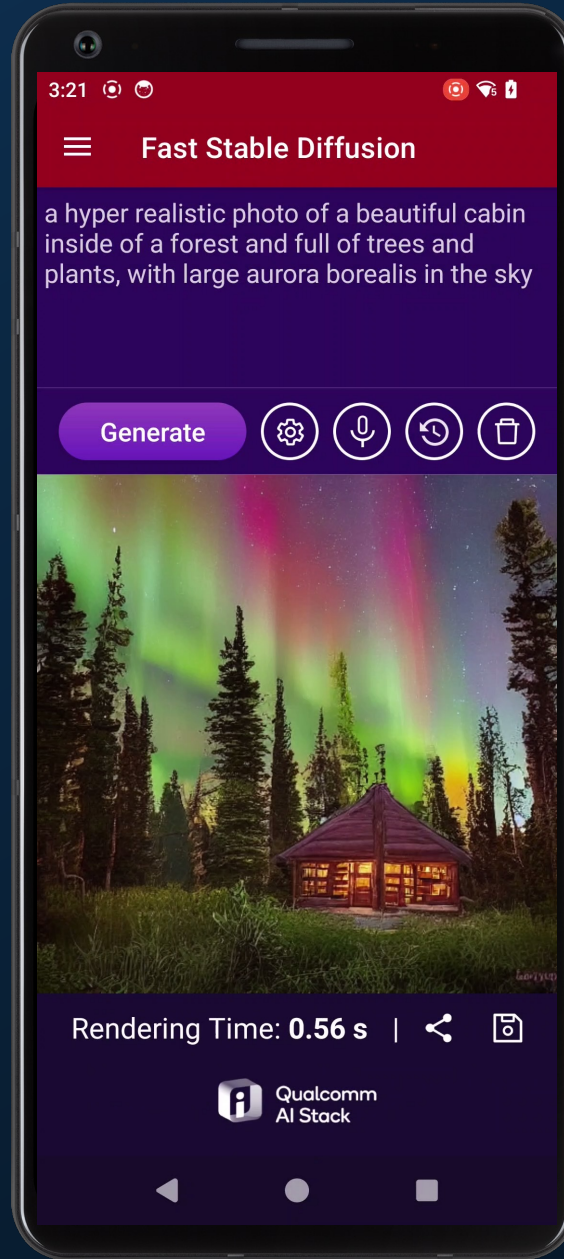
Based on textual description or a handful of images with minimal manual effort

**Examples:**  
DreamFusion, Magic3D

Prompt: "a plush dragon toy"

AT  
SNAPDRAGON  
SUMMIT  
2023

World's  
fastest AI  
text-to-image  
generative AI  
on a phone



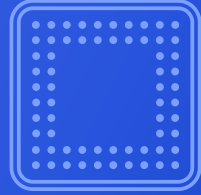
## Fast Stable Diffusion

Takes less than 0.6 seconds for generating 512x512 images from text prompts

Efficient UNet architecture, guidance conditioning, and step distillation

Full-stack AI optimization to achieve this improvement

# What are the challenges to overcome for generative AI images and videos?



## High computation and latency

Gen AI requires immense computational power and infrastructure.



## Memory costs

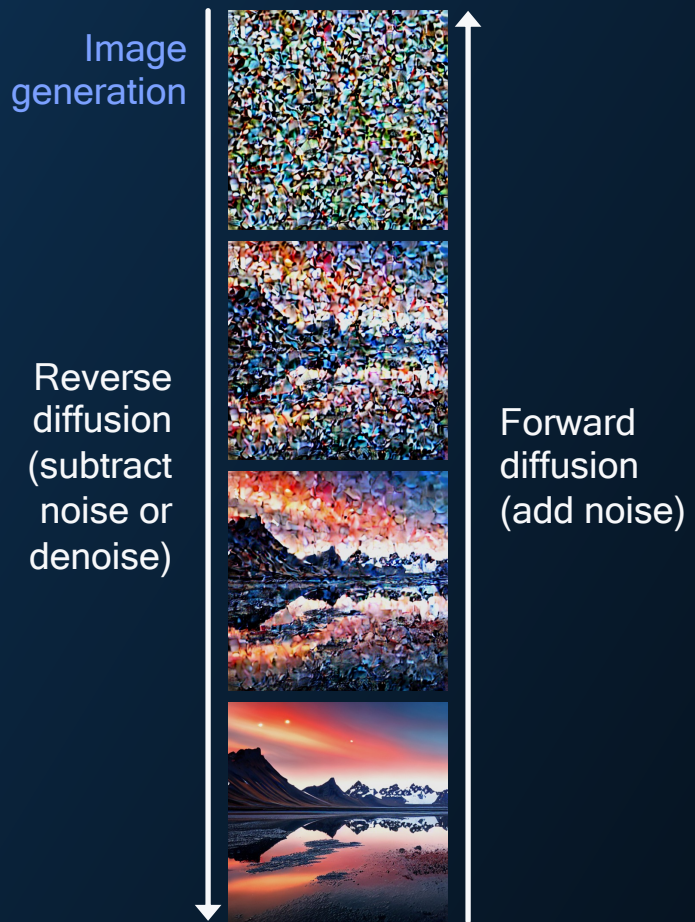
Models demand a lot of memory to perform well and sometimes need to run concurrently.



## Data inefficiency

Models require billions of training samples, that makes it hard to adapt them to new domains.

# What is diffusion?



Prompt: Panoramic view of mountains of Vestrahorn and perfect reflection in shallow water, soon after sunrise, Stokksnes, South Iceland, Polar Regions, natural lighting, cinematic wallpaper

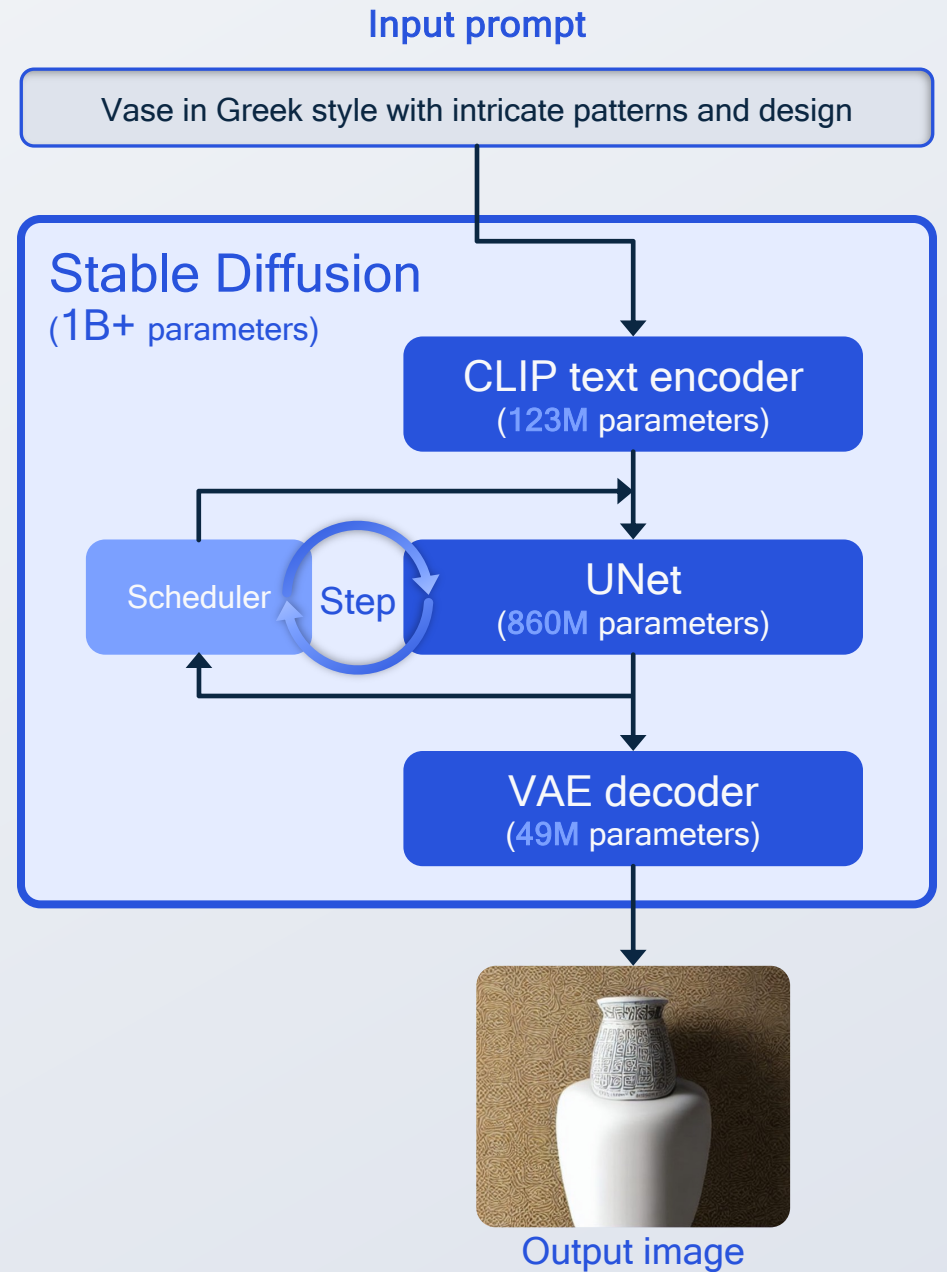
# Stable Diffusion architecture

UNet is the biggest component model of Stable Diffusion

Many steps, often 20 or more, are used for generating high-quality images

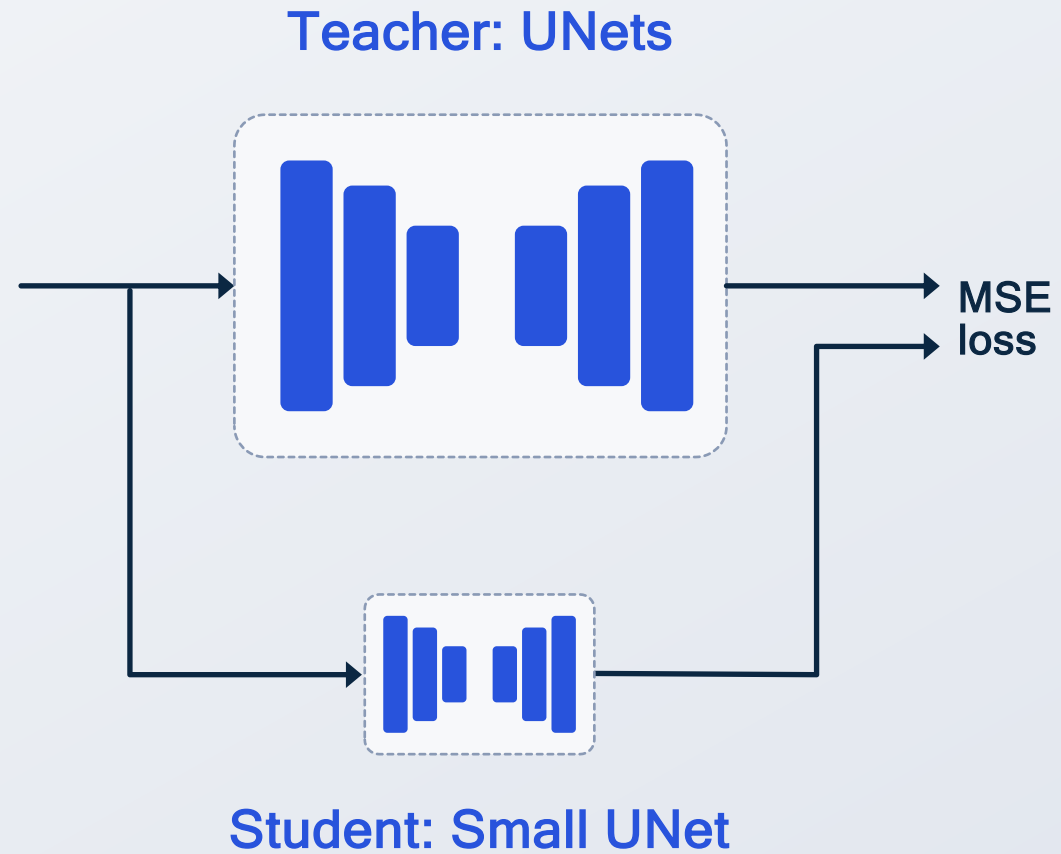
Significant compute is required

VAE: Variational Auto Encoder;  
CLIP: Contrastive Language-Image Pre-Training



Key concept:  
**Model  
distillation**

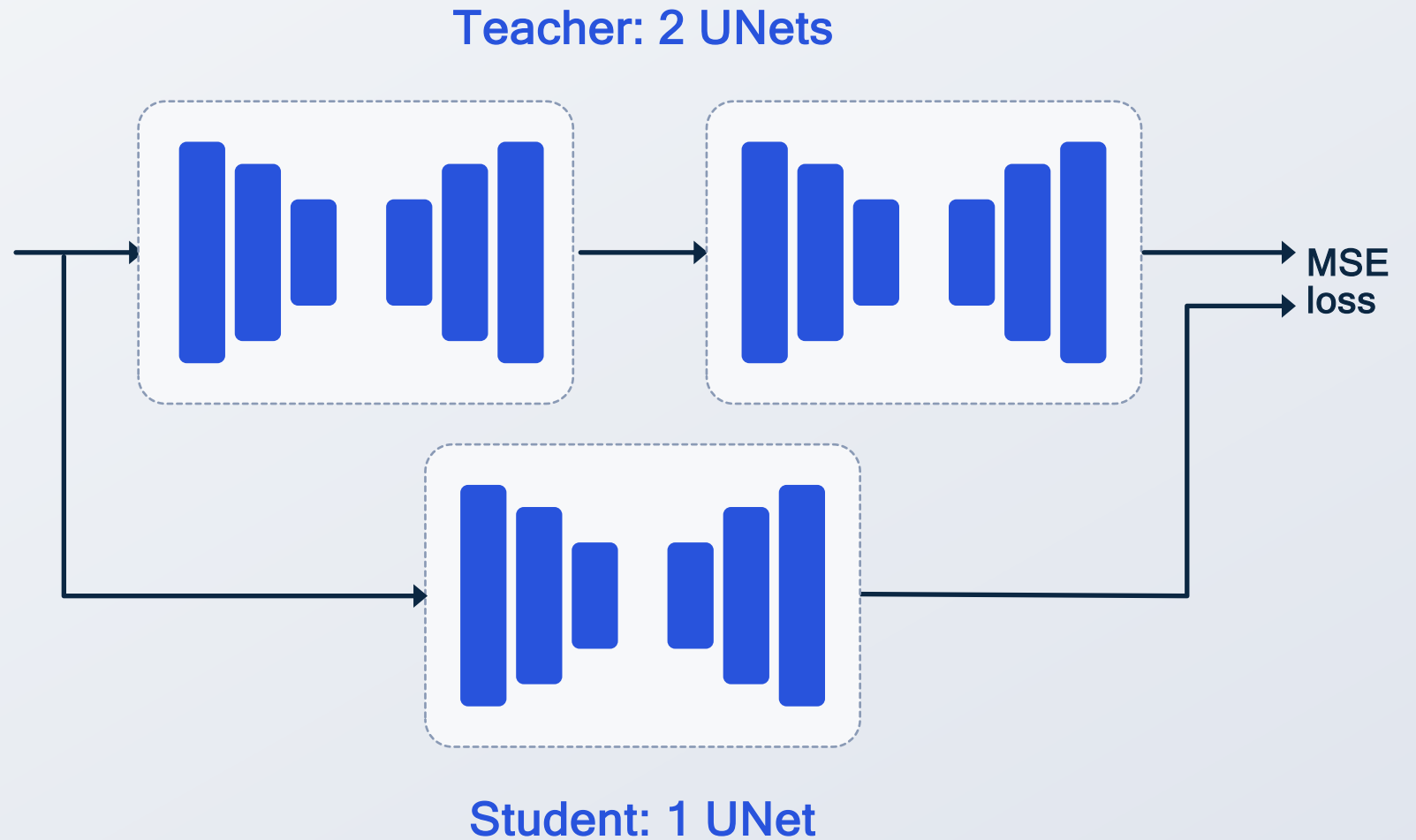
Teach the student  
model to achieve what  
the teacher achieves  
at each step





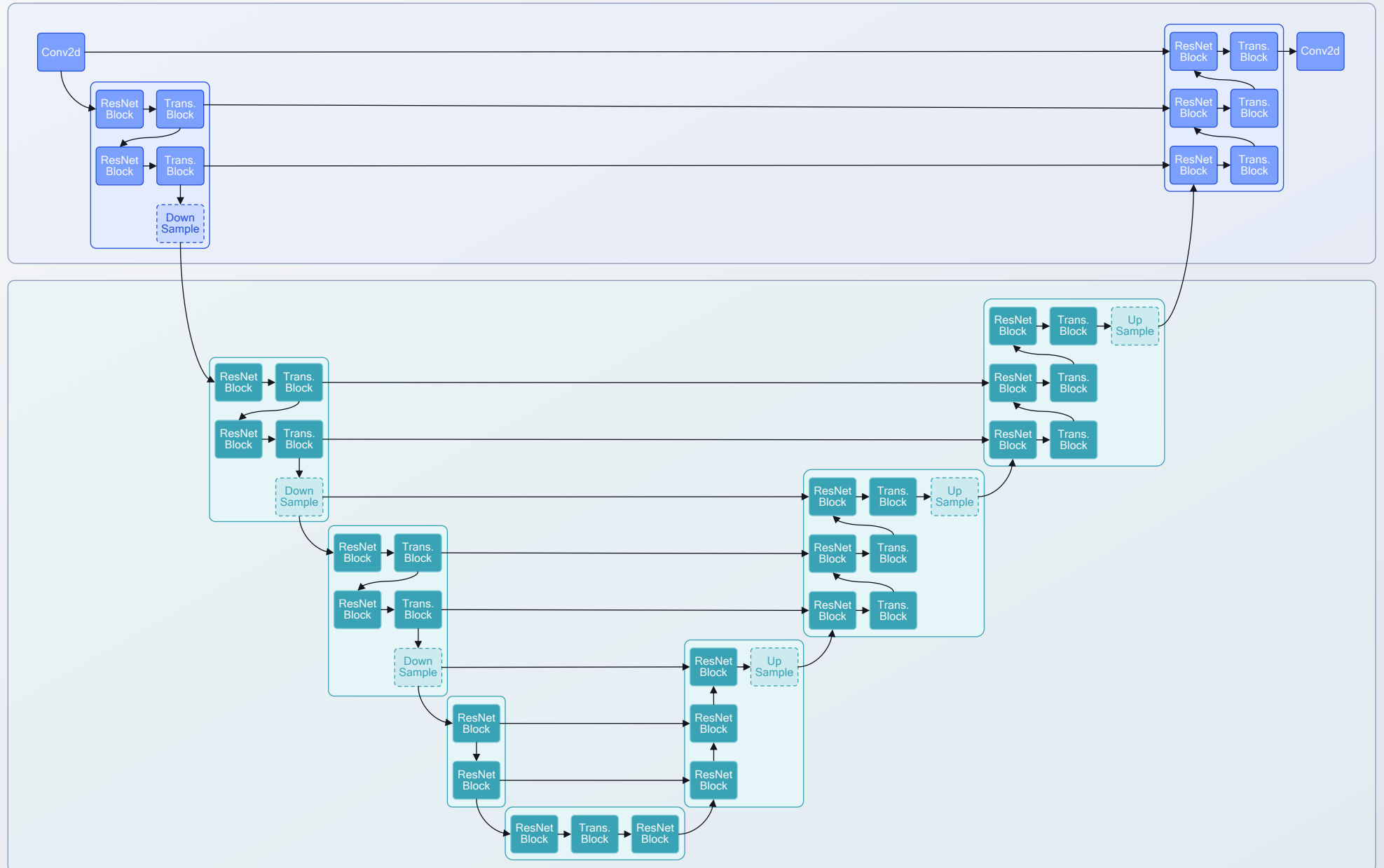
Key concept:  
**Step  
distillation**

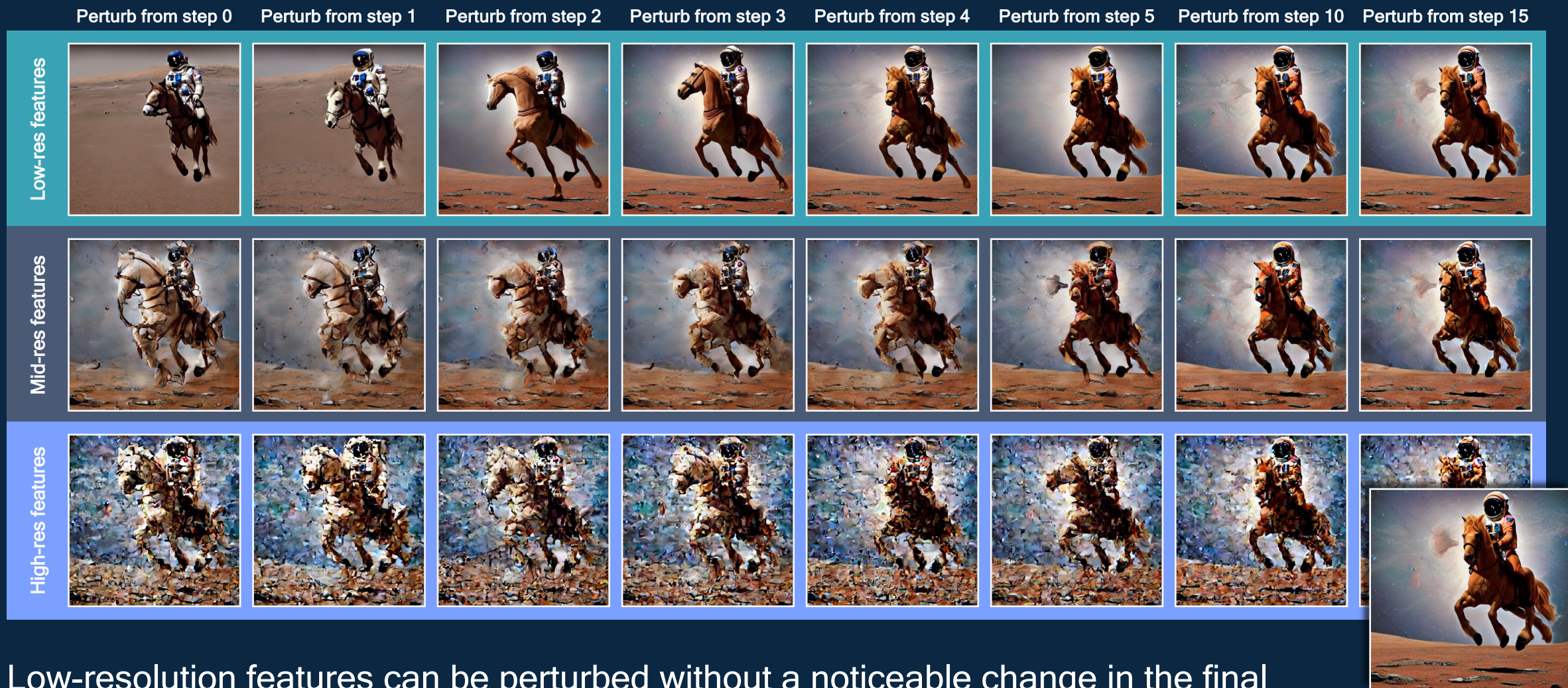
Teach the student model  
to achieve in one step  
what the teacher achieves  
in multiple steps



High-resolution representations in UNet carry high-frequency content (e.g., textures)

Low-resolution representations in UNet carry high-level structure (e.g., scene layout)

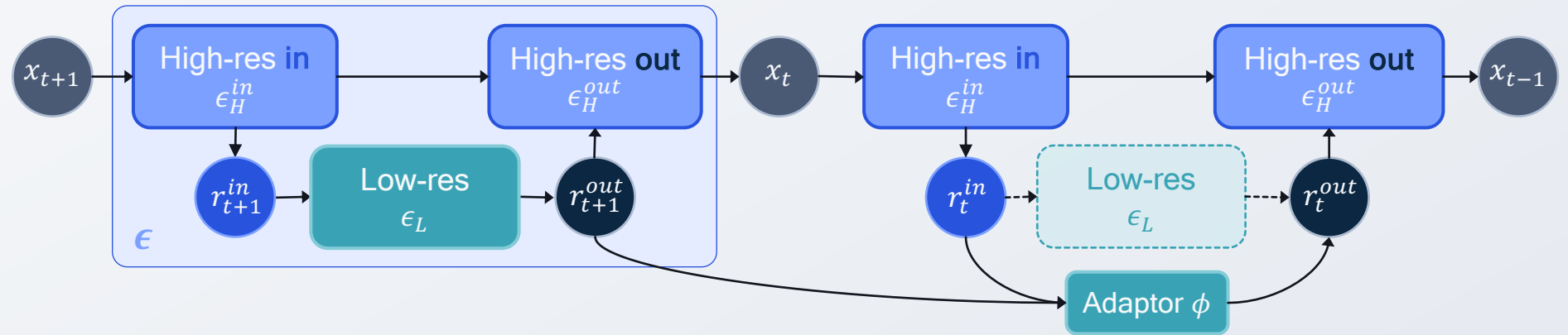




Low-resolution features can be perturbed without a noticeable change in the final output, whereas small perturbations on the high-resolution features degrade the image generation

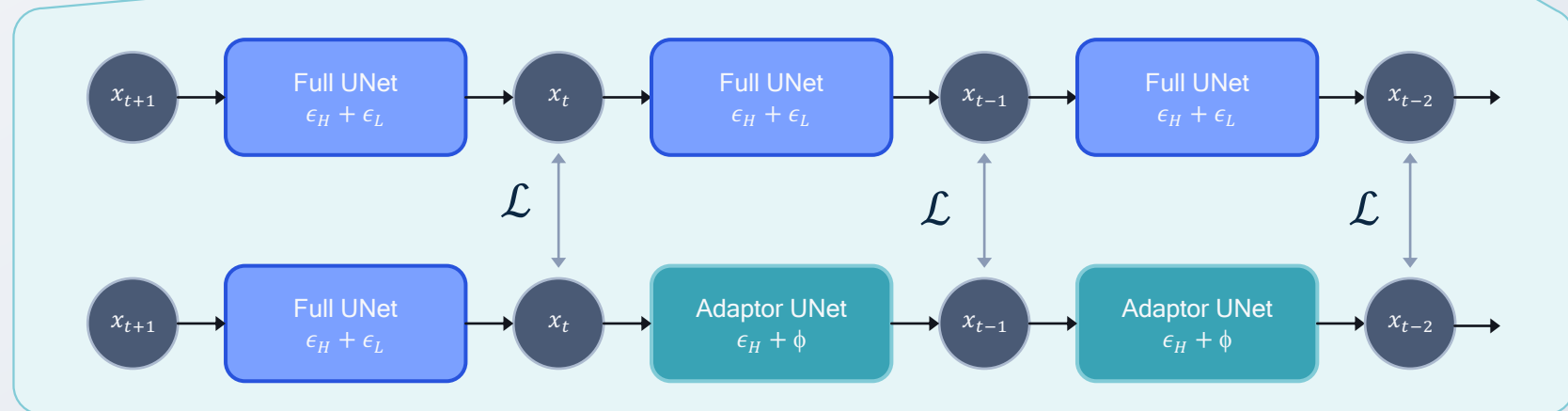
## Clockwork architecture

An efficient approximation of low-res features by adapting from previous steps



## Training the adaptor

Distillation from a full UNet over all denoising steps



How to leverage the perturbation robustness to save computations?

Clockwork  
improves any  
diffusion model

## Text-to-image generation on MS-COCO 2017-5K

Model	FID ↓	CLIP ↑	FLOPs ( $10^{12}$ ) ↓
Stable Diffusion UNet	24.64	0.300	10.8
+ Clockwork	24.11	0.295	7.3 (1.48×)
Efficient UNet	24.22	0.302	9.5
+ Clockwork	23.21	0.296	5.9 (1.61×)
Distilled Efficient UNet	25.75	0.297	4.7
+ Clockwork	24.45	0.295	2.9 (1.62×)

Clockwork  
generates high-  
quality images  
faster than state of  
the art

## Text-to-image generation on MS-COCO 2017-5K

Model	FID ↓	CLIP ↑	FLOPs ( $10^{12}$ ) ↓
InstaFlow (1 step) <sup>1</sup>	29.30	0.283	0.8
Model Distillation <sup>2</sup>	31.48	0.268	7.8
Guidance Distillation <sup>3</sup>	26.90	0.300	6.4
SnapFusion <sup>4</sup>	24.20	0.300	4.0
Clockwork	24.45	0.295	2.9

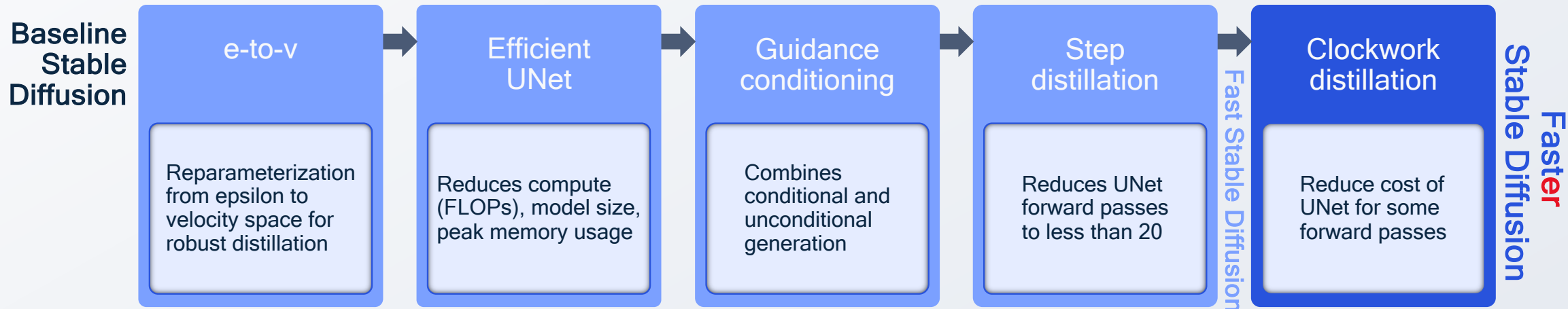
<sup>1</sup> Instaflow: One step is enough for high-quality diffusion-based text-to-image generation. arXiv'22

<sup>2</sup> On architectural compression of text-to-image diffusion models, arXiv'23

<sup>3</sup> On distillation of guided diffusion models, CVPR'23

<sup>4</sup> SnapFusion: Text-to-image diffusion model on mobile devices within two seconds, NeurIPS'23

FID = Frechet Inception Distance, CLIP = Contrastive Language-Image Pre-training



	FID↓	CLIP↑	Diffusion latency	Total latency
Fast Stable Diffusion	26.04	0.297	0.40 seconds	0.65 seconds
<b>Faster</b> Stable Diffusion	25.21	0.292	0.27 seconds	0.53 seconds

**1.2x**  
Speedup

Clockwork reduces the total latency by 1.2x while improving quality compared to Fast Stable Diffusion

# Results: state-of-the-art efficient image generation by Clockwork



Prompts: "large white bear standing near a rock", "the vegetables are cooking in the skillet on the stove.", "bright kitchen with tulips on the table and plants by the window", "red clouds as sun sets over the ocean", "a picnic table with pizza on two trays", "a couple of sandwich slices with lettuce sitting next to condiments."



# Results: Clockwork for image editing

Input



Edited by  
PnP



Edited by  
PnP +  
Clockwork



# The potential of generative video editing

Given an input video and a text prompt describing the edit, generate a new video

The edit usually changes the appearance or shape of a particular object

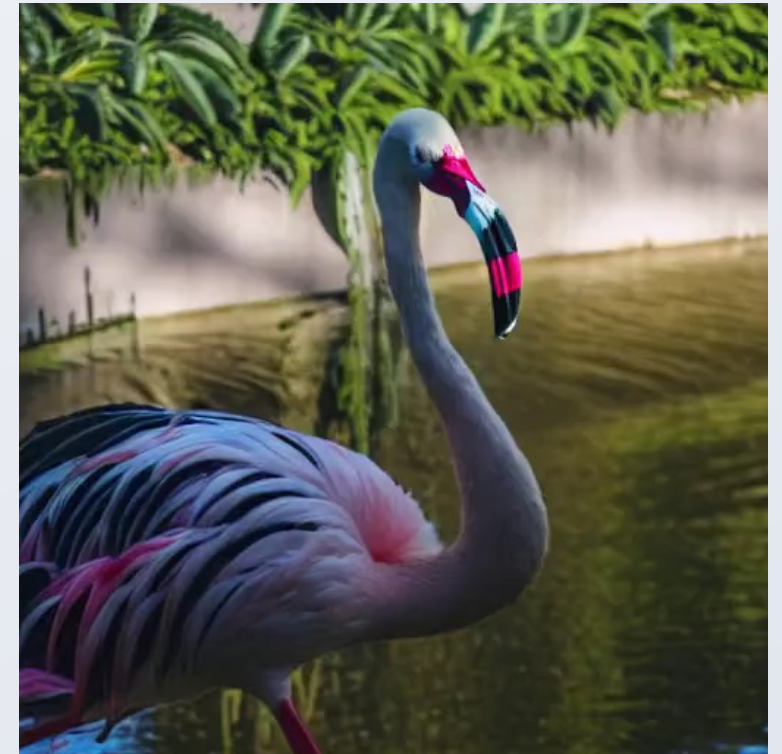
## Key challenges:

1. Temporal consistency
2. High computational cost

Input video



Edited video

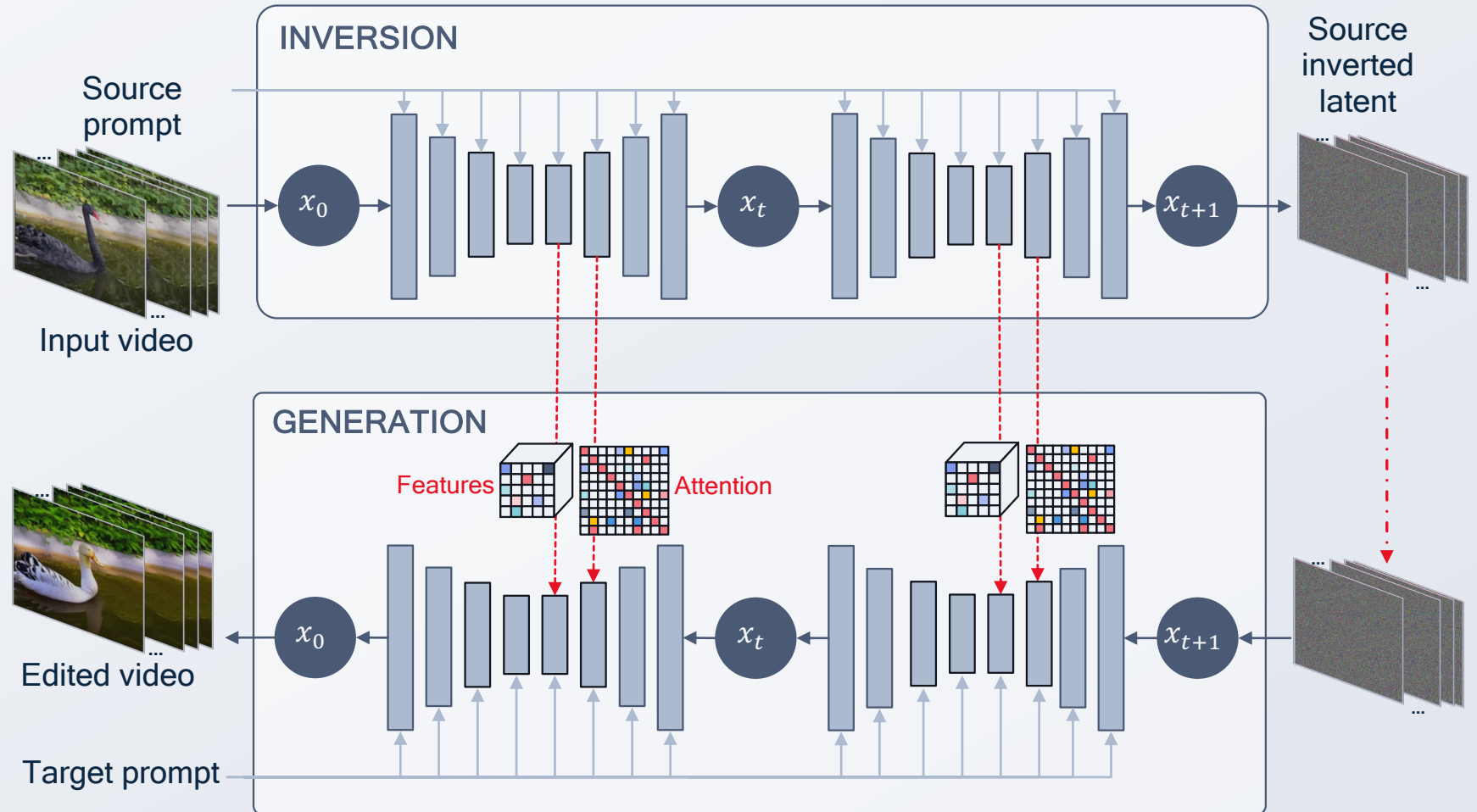


Prompt: “pink flamingo walking”

# Why is video editing so slow?

## Diffusion inversion

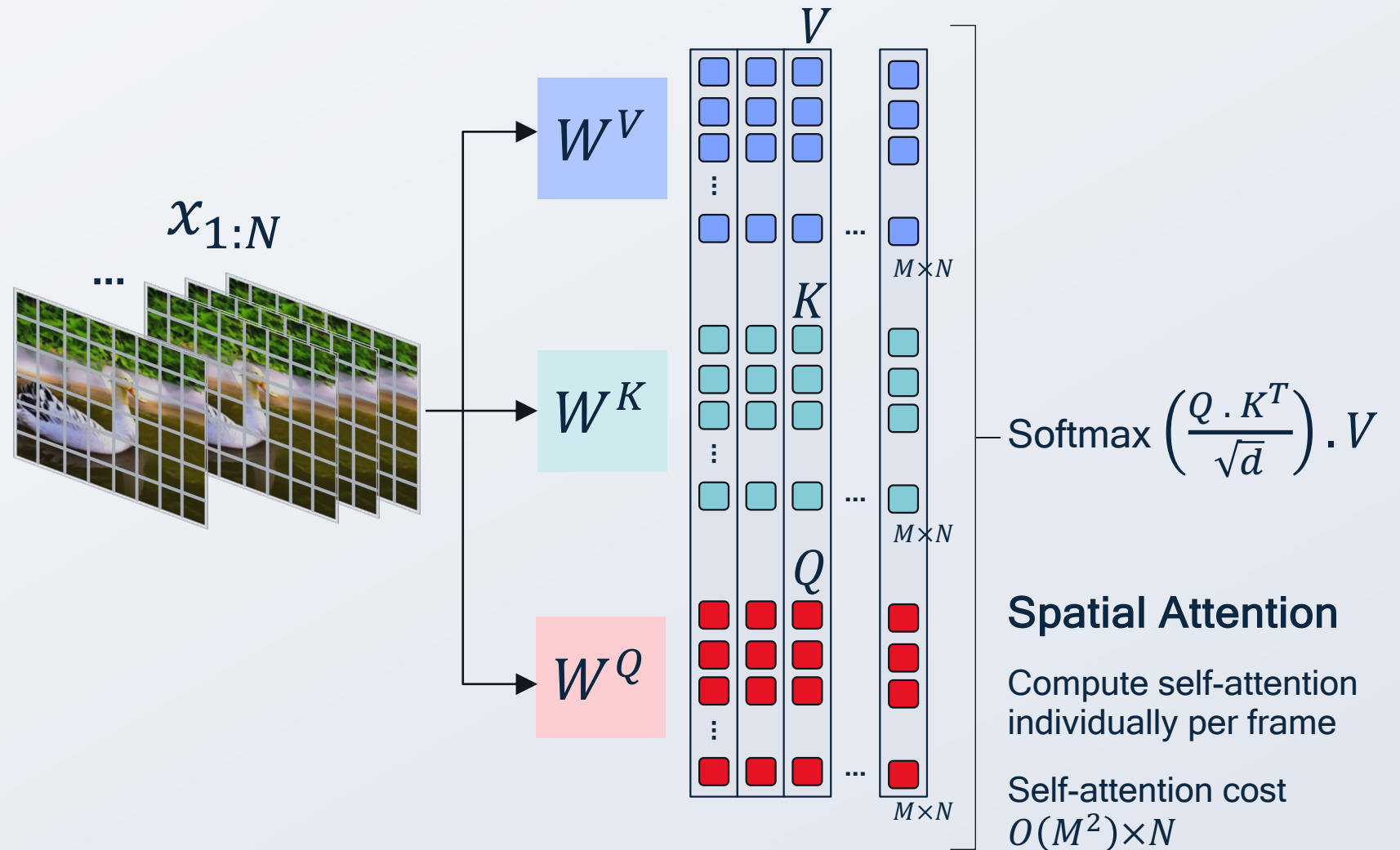
- Essential to preserve temporal consistency and details in the source video
- Comes at a high memory cost to store attention maps and feature



# Why is video editing so slow?

## Temporal attentions

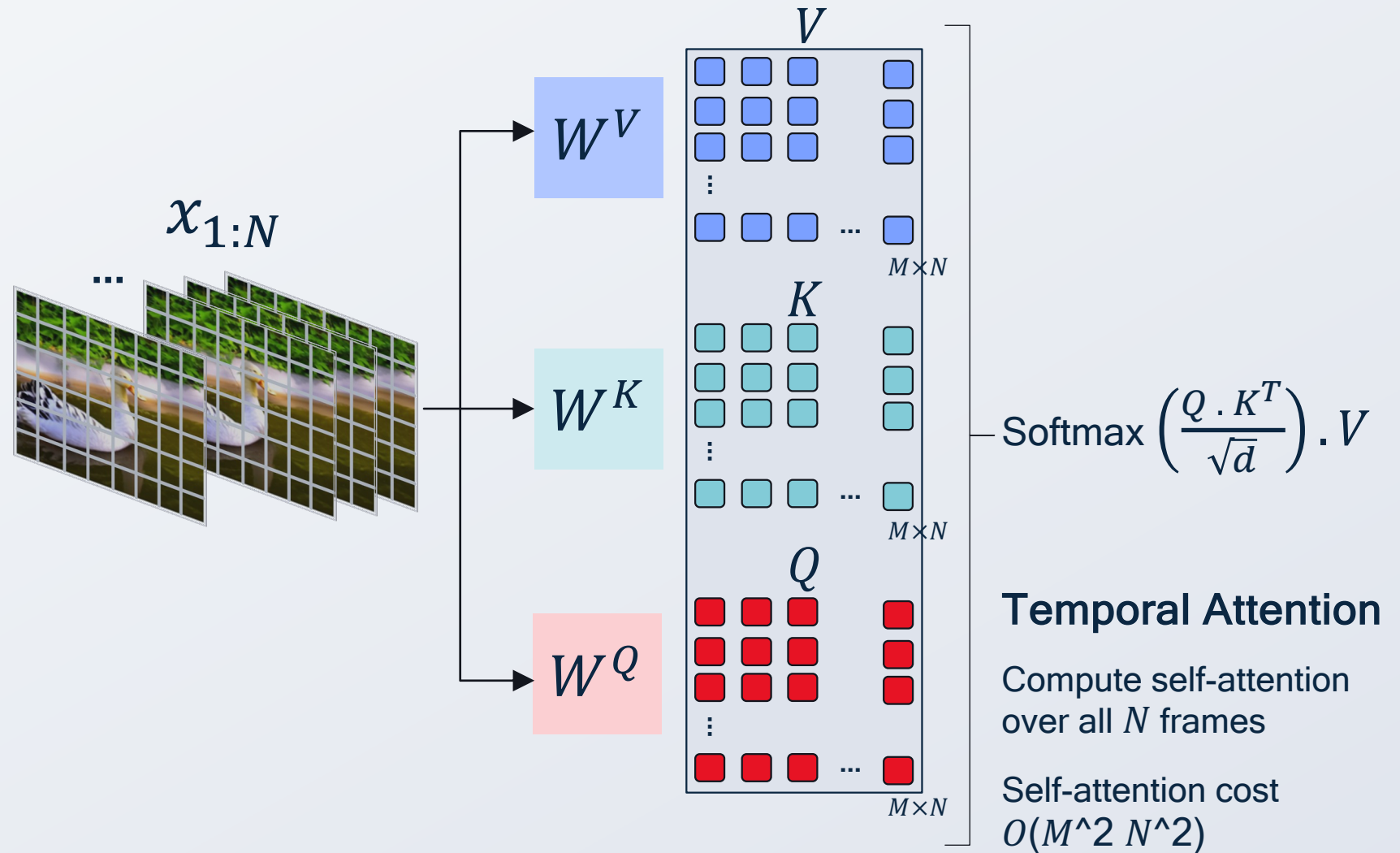
Comes at a high computational cost due to the quadratic cost with respect to video length



# Why is video editing so slow?

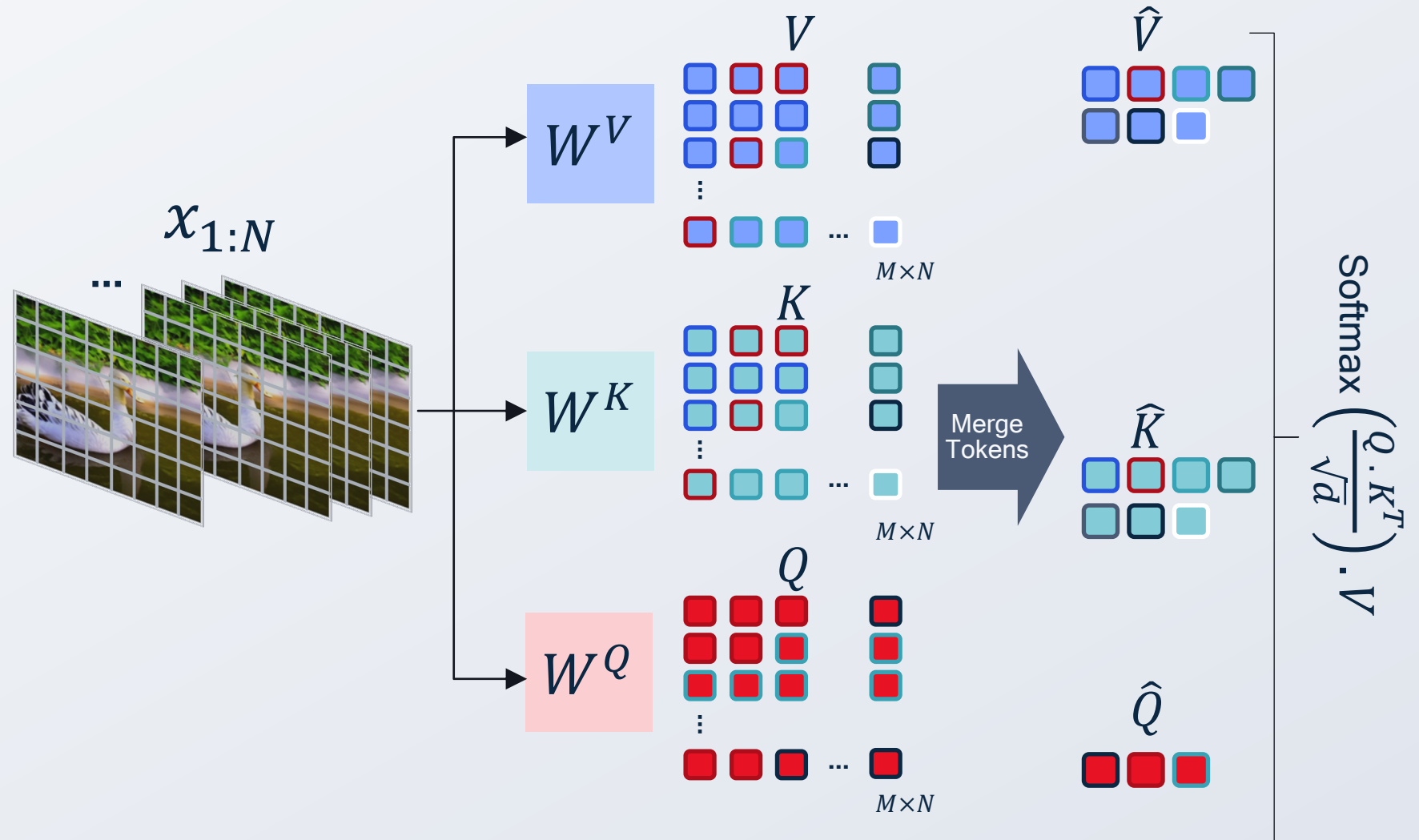
## Temporal attentions

Comes at a high computational cost due to the quadratic cost with respect to video length



# Token merging solves these two challenges

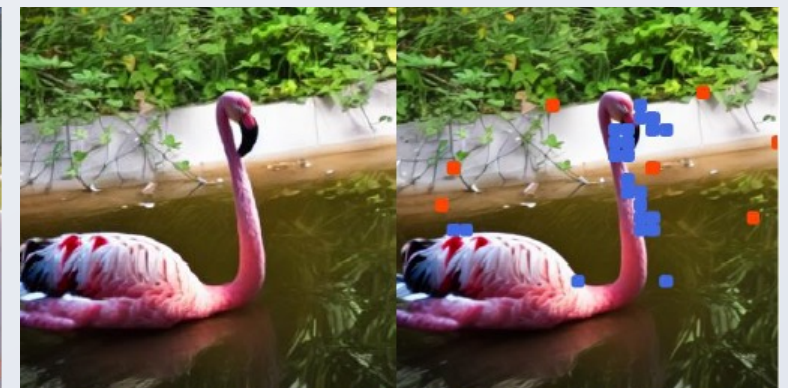
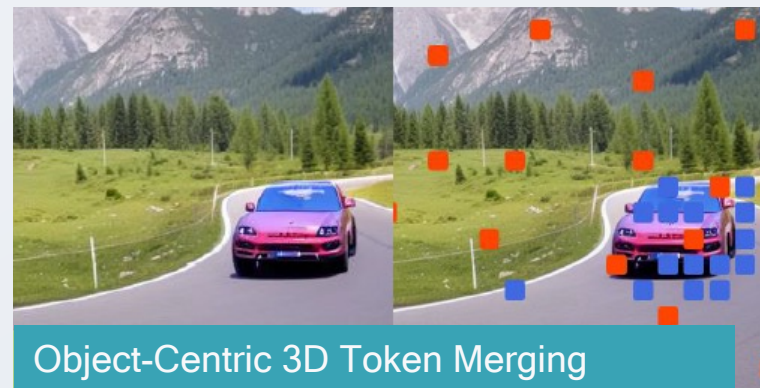
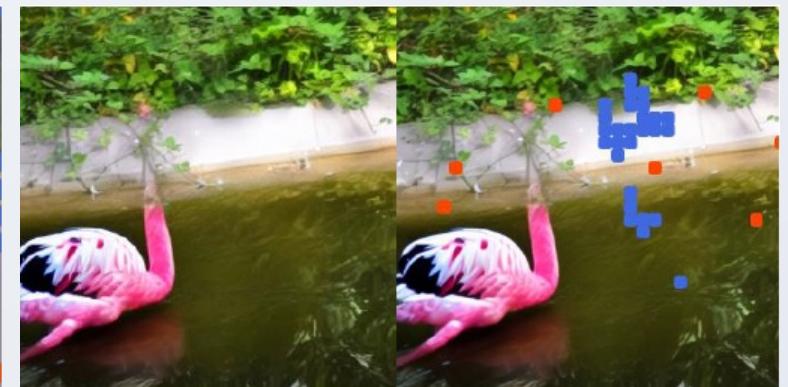
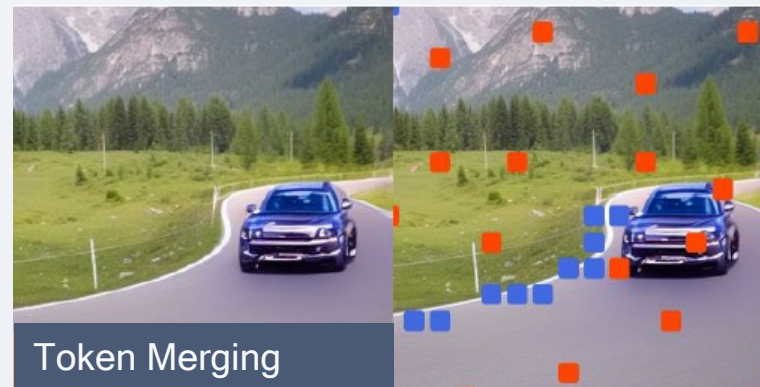
1. Merge the redundant tokens
2. Perform computation on clusters
3. Copy the output back into merged tokens



The costly attention is computed over a fraction of tokens (centroids)

Jeep → Porsche

Swan → Flamingo



# Encourages merging tokens on the background regions

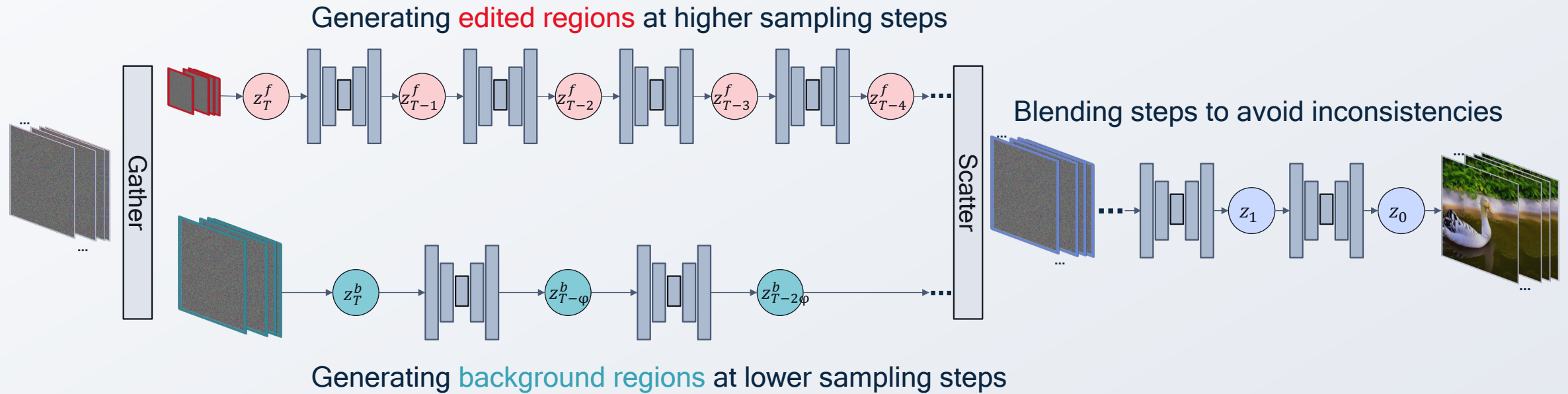
By increasing  $\eta$ , more and more foreground tokens will be left unmerged

Two tokens are merged if their similarities exceeds a threshold

Introduce different thresholds for background vs. edited regions

Using lower threshold on background regions:

- Encourages merging more tokens on background
- Leaves more unmerged tokens on foreground regions



We perform a different number of sampling steps on **edited** and **background** regions:

**Edited regions:** Are usually small, but require most synthesis (more sampling steps)

**Background regions:** Are usually large, and don't require much synthesis (less sampling steps)

Further acceleration by Object-Centric Sampling



# 6 -10x speedup with negligible drop in quality

Applied our acceleration  
on two recent video  
generation frameworks:

- FateZero
- ControlVideo

Our acceleration includes:

- Object-centric 3D token merging
- Object-centric sampling

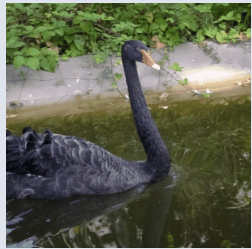
Model	Temporal Cons ↑	CLIP ↑	Latency (s) ↓		
			Inversion	Generation	Overall
FateZero	0.961	0.344	135.80	41.34	177.14
<b>+ Our acceleration</b>	0.967	0.331	8.22 (16.5×)	9.29 (4.4×)	17.51 (10×)

DAVIS benchmark: Generating 8 frames on server GPU

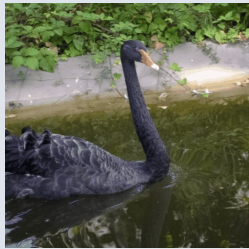
Model	Temporal Cons ↑	CLIP ↑	Latency (s) ↓
ControlVideo	0.972	0.318	152.64
<b>+ Our acceleration</b>	0.977	0.313	25.21 (6.0 ×)

CV benchmark: Generating 15 frames on server GPU

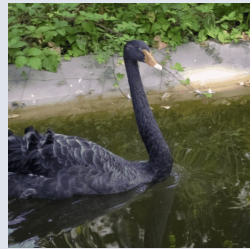
{shape, attribute, style} editing



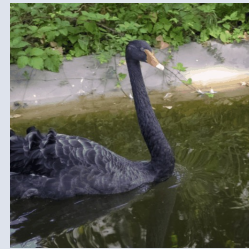
White duck



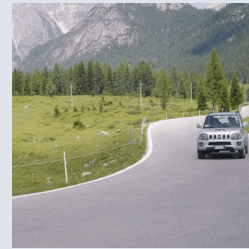
Pink flamingo walking



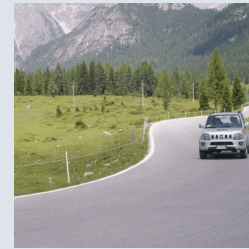
Swarovski crystal



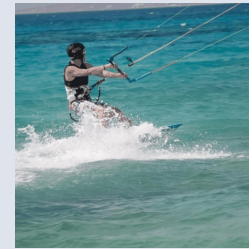
Cartoon photo



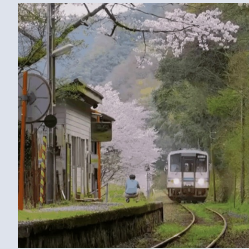
Porsche car



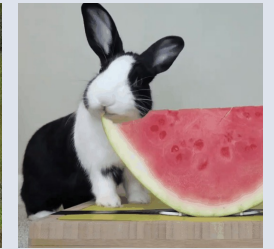
Watercolor painting



Ukiyo-e style

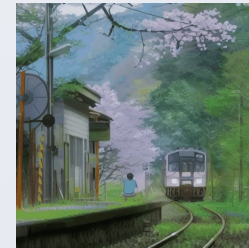
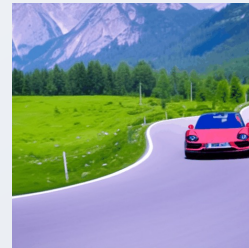
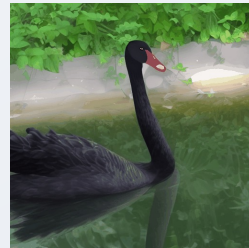
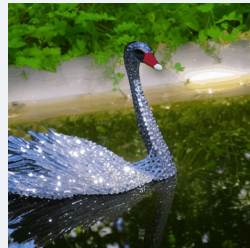
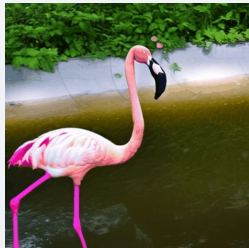
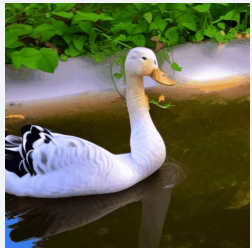


Makoto shinkai style

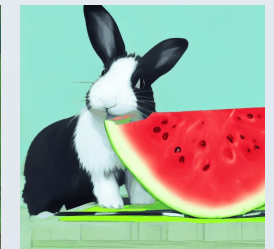
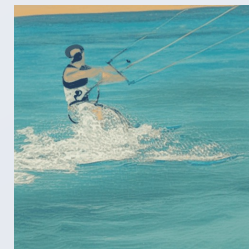
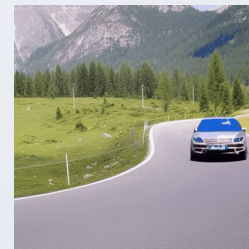
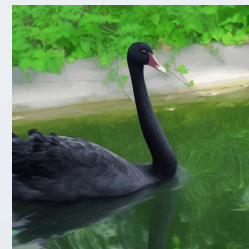
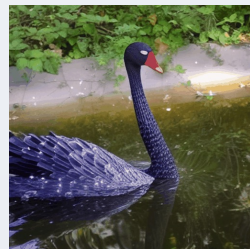
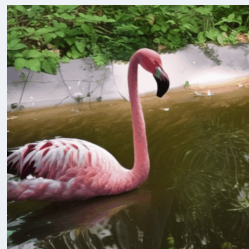
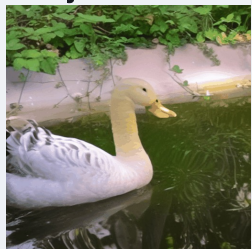


Pokemon cartoon

### FateZero



### + Object-Centric Diffusion



# 10x faster at a comparable editing quality

Research on further optimizations to enable on-target deployment of video generation models

# How does 3D generation work?

Generating 3D mesh  
from a text prompt

- Crucial for many tasks  
e.g., XR, graphics
- Manual creation of  
3D assets is costly

A plush dragon toy



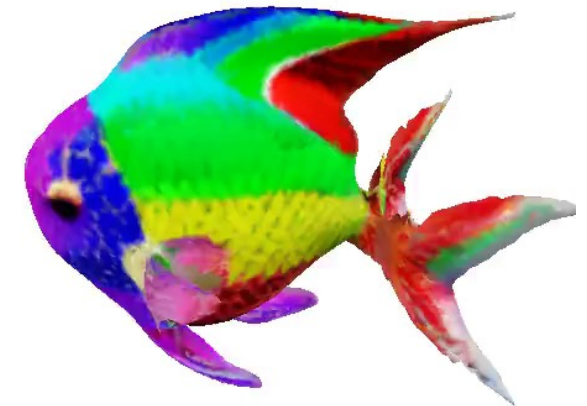
A DSLR photo of a hippo  
wearing a sweater



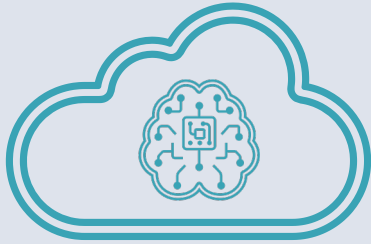
A DSLR photo of a train engine  
made out of clay



A beautiful rainbow fish



## Optimization-based approach



- Costly optimizations to fit mesh parameters for each object
- Takes **+20 min** to model a new object / scene
- Leverage a pretrained image generator to improve the optimization, i.e., score distillation sampling<sup>1</sup>

## Feed-forward approach



- Generate mesh parameters directly without any optimizations at inference
- Takes **seconds** to model a new object / scene
- Learned from scratch on the limited 3D data available

Can pretrained image generators, e.g., Stable Diffusion, improve **feed-forward** 3D generation?

---

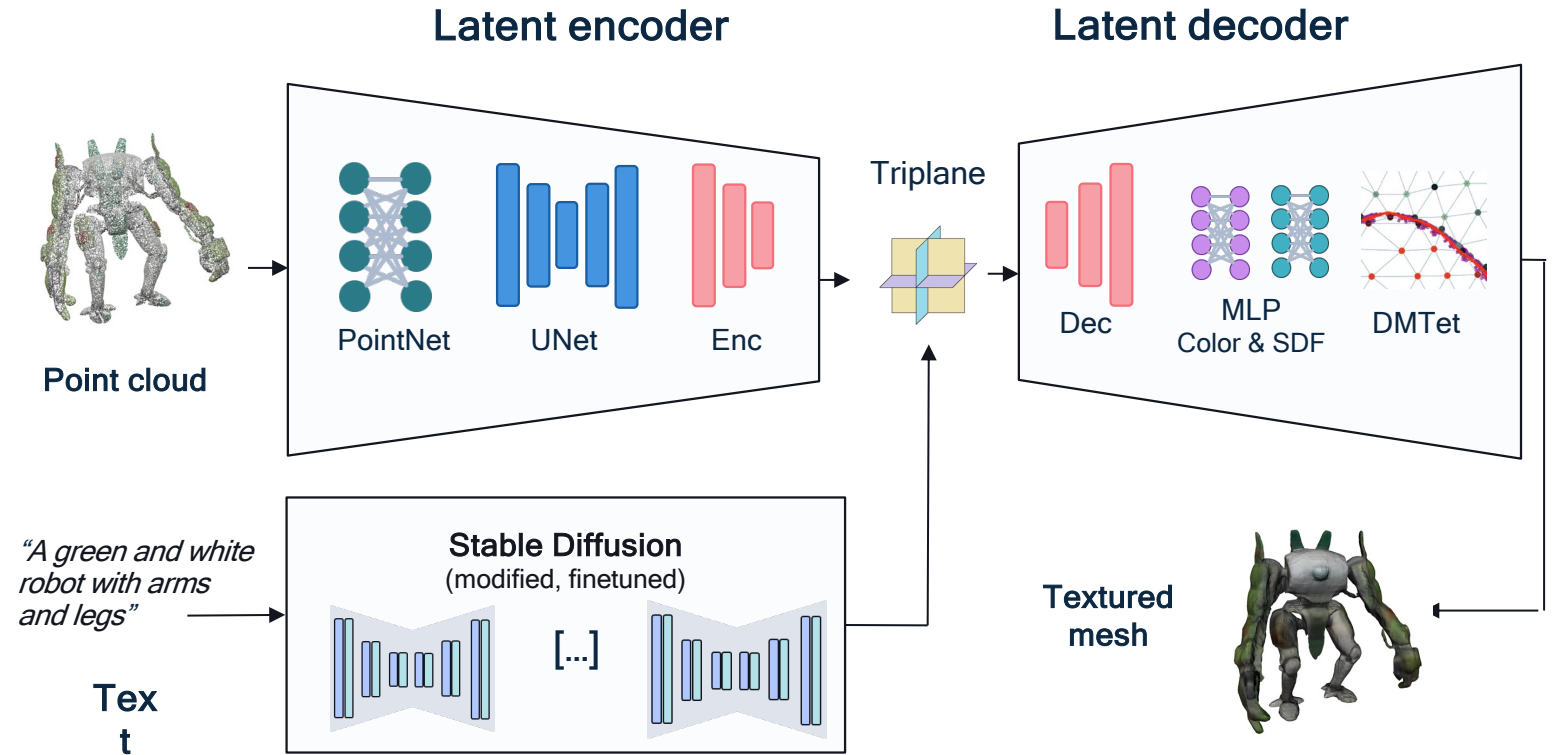
Transfer the huge diversity in 2D image datasets into 3D tasks

<sup>1</sup> DreamFusion: Text-to-3D using 2D Diffusion, arXiv'22

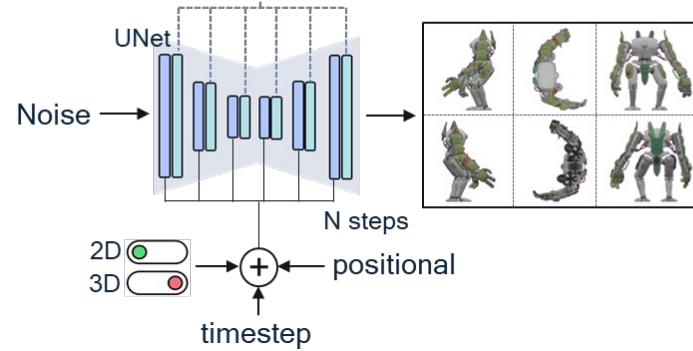
# HexaGen3D

Like other latent diffusion models, we follow a two-stage training:

1. Variational Auto-Encoder (VAE) to reconstruct meshes from point clouds:
  - Latents defined in a triplane space of the shape  $H \times 3 \cdot W \times C$
2. Conditional generation of triplanar latents from text
  - By adapting a pretrained Stable Diffusion model



*"A green and white robot with arms and legs"*



## Step 1:

### Generate Hexaview guidance

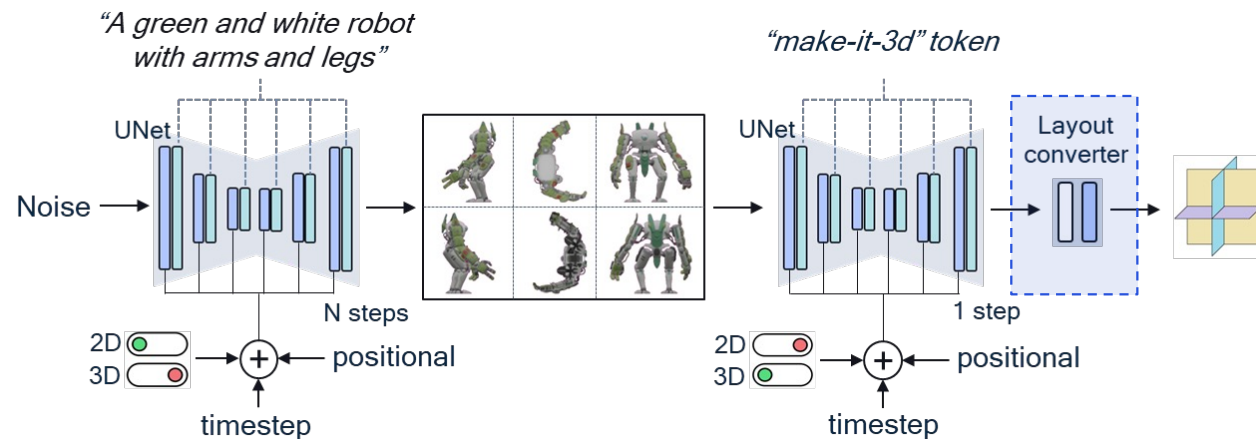
- Tile "front", "rear", "right", "left", "top" and "bottom" views into a large Hexaview image
- As an intermediate generation step, guides Stable Diffusion to generate triplanar latents

We generate triplanar latents using a pretrained Stable Diffusion Model in two steps

## Step 1:

### Generate Hexaview guidance

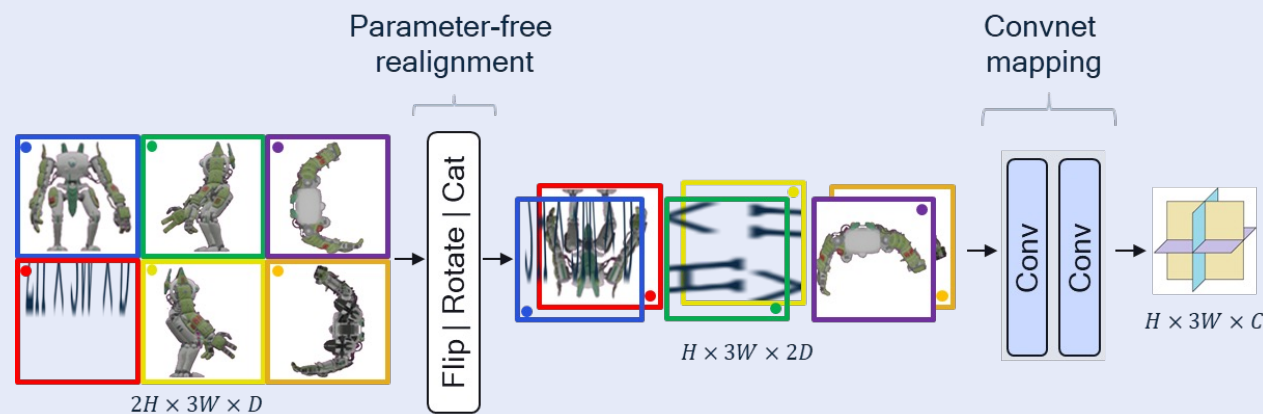
- Tile “front”, “rear”, “right”, “left”, “top” and “bottom” views into a large Hexaview image
- As an intermediate generation step, guides Stable Diffusion to generate triplanar latents



## Step 2:

### Convert Hexaviews into triplanar latents

- Split and align the views, followed by a ConvNet
- Using the same UNet (parameter efficiency) with a different prompt and 3D embedding



We generate triplanar latents using a pretrained Stable Diffusion Model in two steps

We generate high quality meshes, much faster than optimization-based methods:

7 sec vs. +22 min

Model	Latency (s) ↓	CLIP ↑	User preference ↑
MVDream	194 mins	30.35	0.97
TextMesh	23 mins	25.06	0.12
DreamFusion	22 mins	28.91	0.59
<b>HexaGen3D</b>	<b>7 secs</b>	29.58	0.73

MVDream-SDv2.1 (~194 mins)



TextMesh-SDv2.1 (~23 mins)



DreamFusion-SDv2.1 (~22 mins)



HexaGen3D-SDXL (~7 sec)



*"a bald eagle carved out of wood"*

*"A DSLR photo of a frog wearing a sweater"*

*"a brightly colored mushroom growing on a log"*

*"a DSLR photo of a mug of hot chocolate with whipped cream and marshmallows"*

*"a DSLR photo of a hippo wearing a sweater"*

*"a beautiful dress made out of fruit, on a mannequin. Studio lighting, high quality, high resolution"*

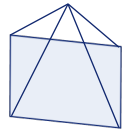
*"a blue motorcycle"*



## Novel View Synthesis (NVS) generates a novel view of an object from a target pose

- Input: an image and a camera pose

Source  
camera pose



Source  
image

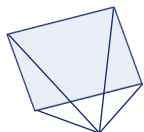


Novel View  
Synthesis

Target  
image



Target  
camera pose

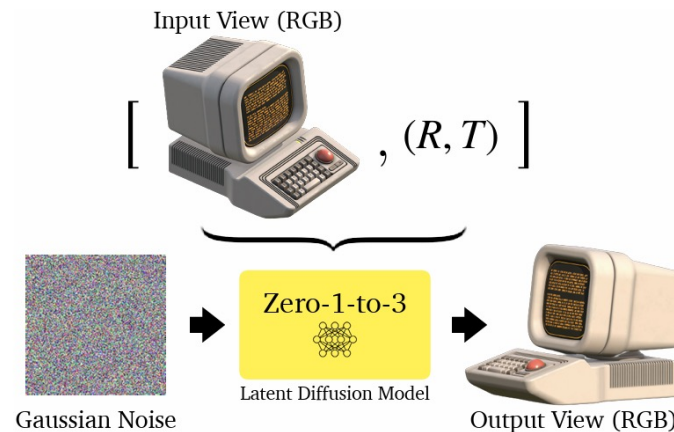


## Optimization based models i.e., NeRF

- Slow but high-quality when there are multiple views available

## Recently Stable Diffusion is adopted for Zero-shot NVS

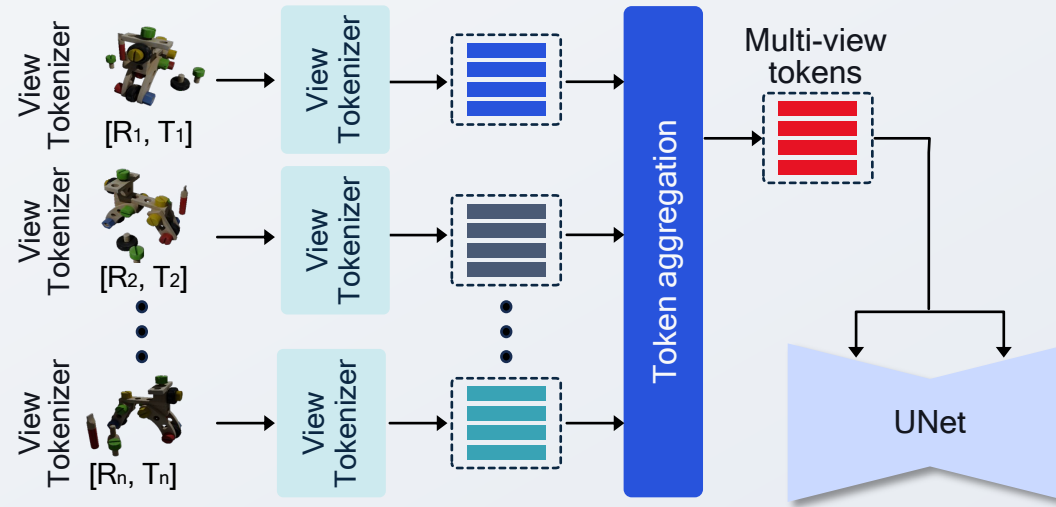
- Fast but quality is limited as consumes single view only



Zero-1-to-3: Zero-shot One Image to 3D Object, ICCV'23

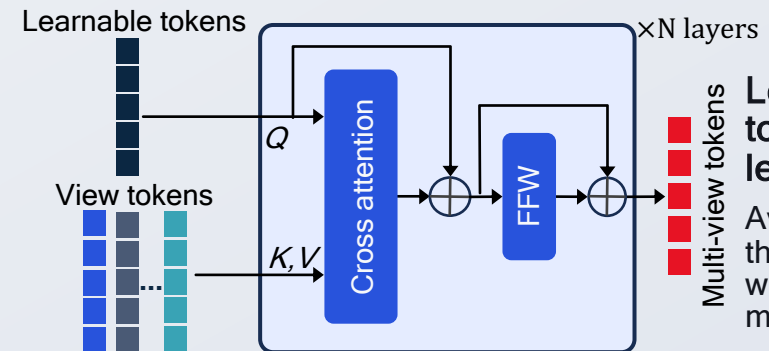
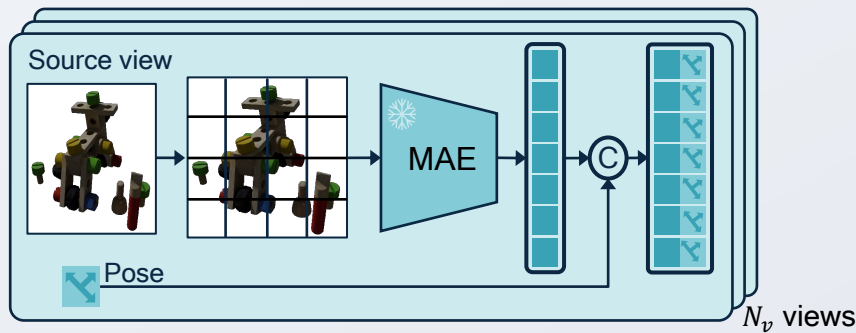
# How to enable zero-shot models to handle multiple views without increasing compute?

Tokenize each view, aggregate over views, and use the **multi-view tokens** as cross-attention conditioning



Pretrained Masked AutoEncoder (MAE) to tokenize views

Provides better representation than CLIP



**Learnable tokens to generate fixed length output**  
 Avoid increasing the computations when fusing more views

# VaLID: Variable Length Input Diffusion

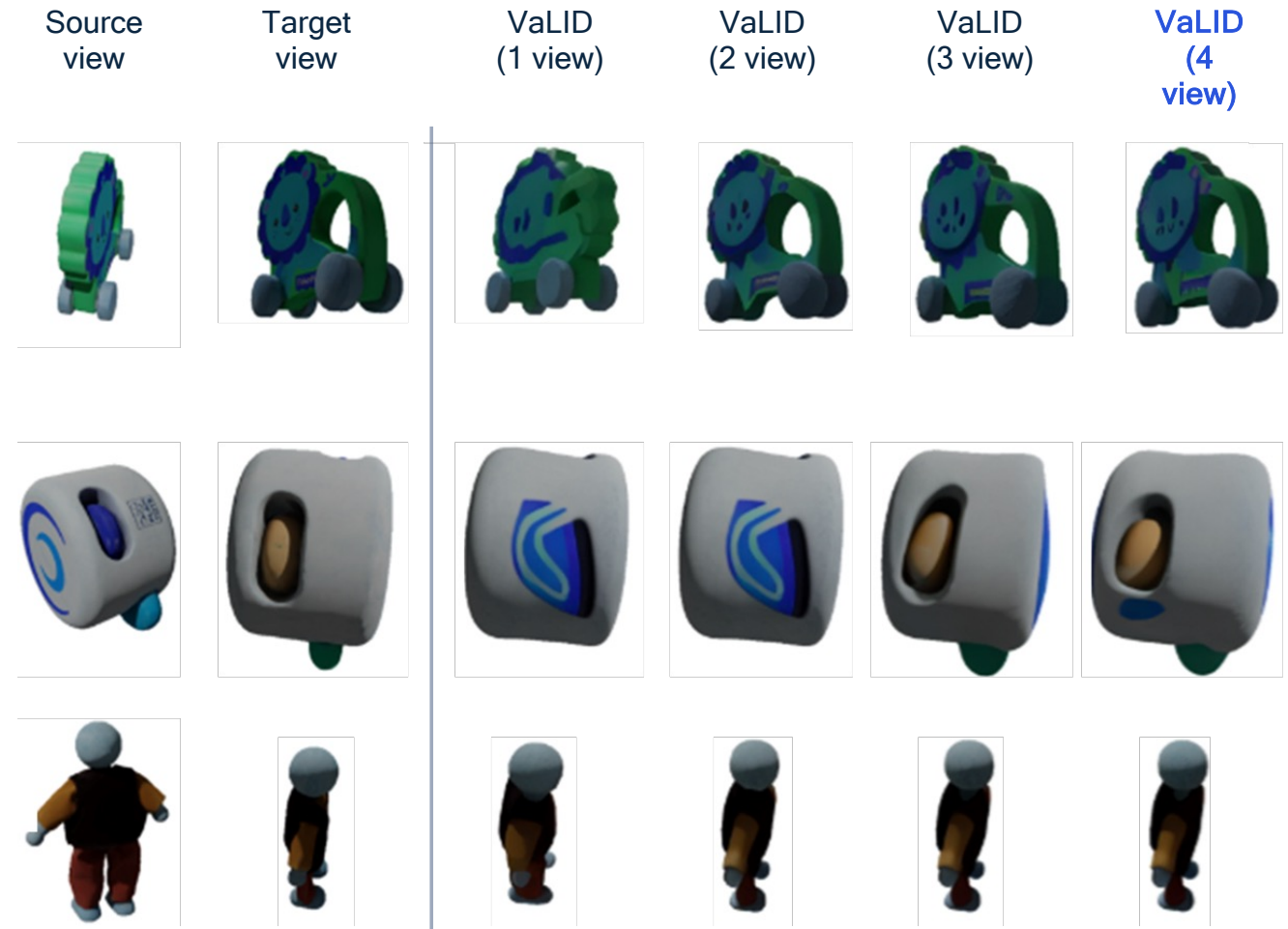
# Our method outperforms existing SOTA methods in quality

## Google Scanned Object dataset

Model	PSNR $\uparrow$	LPIPS $\downarrow$	GFLOPs $\downarrow$
DietNeRF <sup>1</sup>	8.93	0.412	High
SJC-I <sup>2</sup>	5.91	0.545	High
IV <sup>3</sup>	6.57	0.484	High
Zero123 <sup>4</sup>	19.0	0.115	Similar to ours
VaLID (1 view)	20.03	0.091	87.2
VaLID (2 view)	20.41	0.085	87.8
VaLID (3 view)	21.05	0.073	88.8
<b>VaLID (4 view)</b>	<b>21.30</b>	<b>0.069</b>	<b>91.4</b>

PSNR = Peak Signal-to-Noise Ratio

LPIPS = Learned Perceptual Image Patch Similarity



At a negligible computational cost, VaLID processes multiple views to generate more accurate views

1: Putting NeRF on a Diet: Semantically Consistent Few-Shot View Synthesis, CVPR'21

2: Stable Diffusion Image Variation, arXiv'23

3: Score Jacobian Chaining: Lifting pretrained 2d diffusion models for 3d generation, CVPR'23

4: Zero-1-to-3: Zero-shot One Image to 3D Object, ICCV'23

# Our method outperforms existing SOTA methods in quality

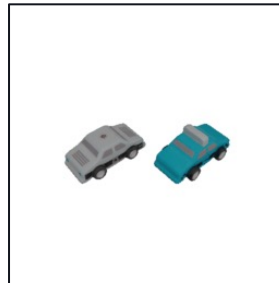
## Google Scanned Object dataset

Model	PSNR $\uparrow$	LPIPS $\downarrow$	GFLOPs $\downarrow$
DietNeRF <sup>1</sup>	8.93	0.412	High
SJC-I <sup>2</sup>	5.91	0.545	High
IV <sup>3</sup>	6.57	0.484	High
Zero123 <sup>4</sup>	19.0	0.115	Similar to ours
VaLID (1 view)	20.03	0.091	87.2
VaLID (2 view)	20.41	0.085	87.8
VaLID (3 view)	21.05	0.073	88.8
<b>VaLID (4 view)</b>	<b>21.30</b>	<b>0.069</b>	<b>91.4</b>

PSNR = Peak Signal-to-Noise Ratio

LPIPS = Learned Perceptual Image Patch Similarity

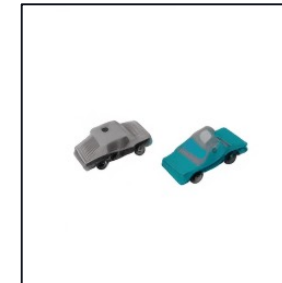
Ground Truth



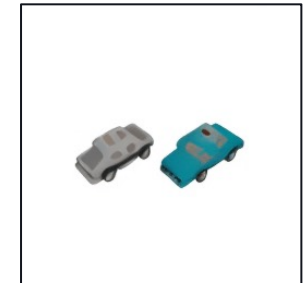
Zero123



VaLID (1 view)



VaLID (4 view)



At a negligible computational cost, VaLID processes multiple views to generate more accurate views

1: Putting NeRF on a Diet: Semantically Consistent Few-Shot View Synthesis, CVPR'21

3: Stable Diffusion Image Variation, arXiv'23

2: Score Jacobian Chaining: Lifting pretrained 2d diffusion models for 3d generation, CVPR'23

4: Zero-1-to-3: Zero-shot One Image to 3D Object, ICCV'23

# We adapt the generative model to new domain: automotive

Generative models improve graphic simulators by being:

- Realistic by being trained on real images and videos
- Scalable by sampling examples instead of manually crafting the assets/objects and scenario



Generate training data for long-tailed object classes

i.e., animals and emergency vehicles



Scale up test set by diversifying

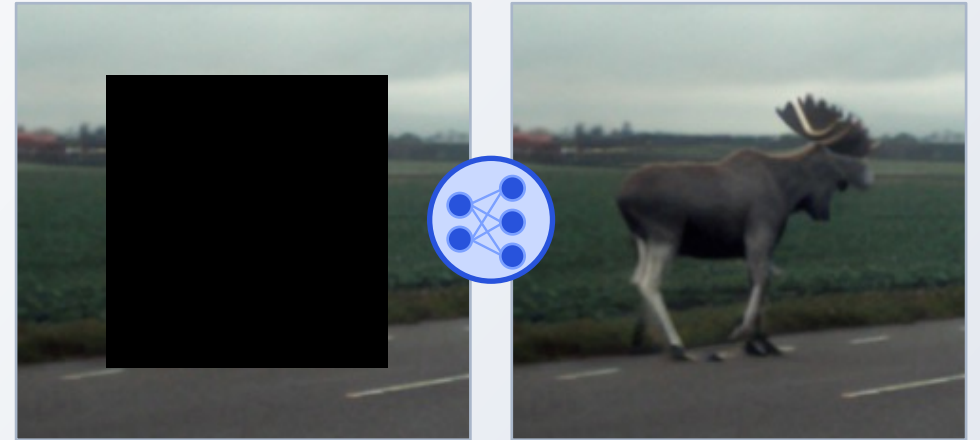
i.e., weather, object appearance and geometries



Generate safety critical test scenarios

i.e., crashes and pedestrians on road

## Animal Detection



Training data

Real images

Real + generated images

mA

P<sub>50.2</sub>

57.7




Number of generated objects

Method  Original  Ours  XPaste

Show bounding boxes  None  Generated  All

 Generate

 Gallery

## Inpainting for animal detection

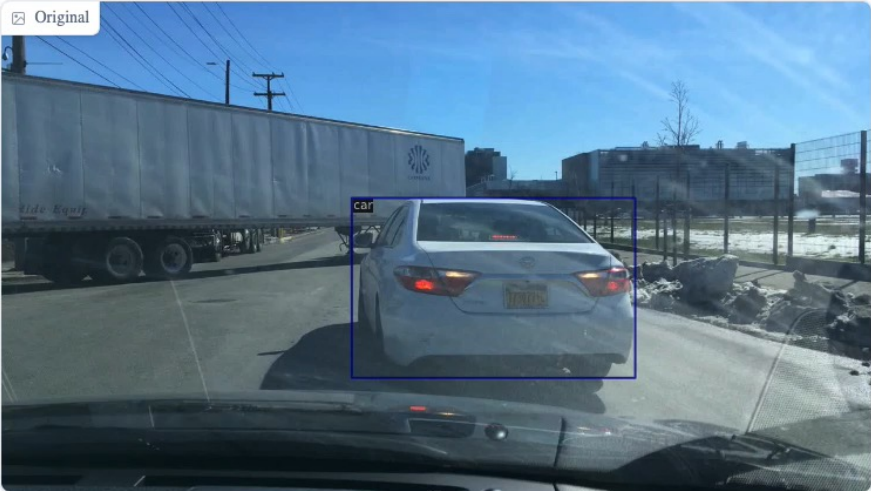
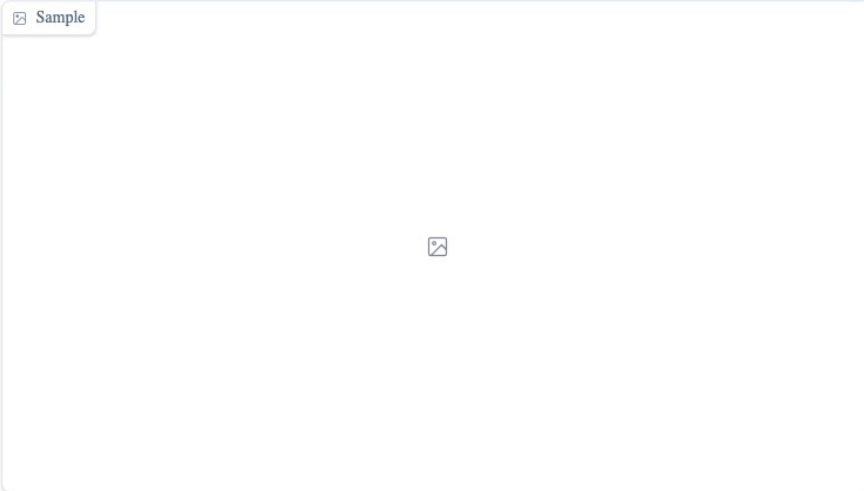
Adapting the generative model to new domains, i.e., automotive scenes

High-fidelity generation in tight bounding boxes

Putting animals at the right geometry: location and scale

Added to the training set to improve animal detector

Previous Image index 5 Next

Original   Sample 

Show which image?  
 Generated image  Generation + detections

Text prompt

Edit strength 0.7

Generate

## High-fidelity object editing

Using generative editing models to change appearance of vehicles

Diversify the test data to less common vehicle types like classics

Avoid unintended changes in appearance and geometry of vehicle and its background

# Qualcomm

Generative vision has a great potential in image and video generation across enterprise, entertainment, XR, and automotive.

Efficient generative vision is important for achieving scale, at the cloud and on device.

Qualcomm AI Research has achieved state-of-the-art results in image and video generation with novel techniques.





## Connect with us



[www.qualcomm.com/research/artificial-intelligence](http://www.qualcomm.com/research/artificial-intelligence)



[www.qualcomm.com/news/onq](http://www.qualcomm.com/news/onq)



[www.youtube.com/c/QualcommResearch](http://www.youtube.com/c/QualcommResearch)



[@QCOMResearch](https://twitter.com/QCOMResearch)



<https://assets.qualcomm.com/mobile-computing-newsletter-sign-up.html>



[www.slideshare.net/qualcommwirelessevolution](http://www.slideshare.net/qualcommwirelessevolution)

# Thank you

**Qualcomm**

Follow us on: [in](#) [twitter](#) [instagram](#) [youtube](#) [facebook](#)

For more information, visit us at:

[qualcomm.com](http://qualcomm.com) & [qualcomm.com/blog](http://qualcomm.com/blog)

Nothing in these materials is an offer to sell any of the components or devices referenced herein.

©2018-2024 Qualcomm Technologies, Inc. and/or its affiliated companies. All Rights Reserved.

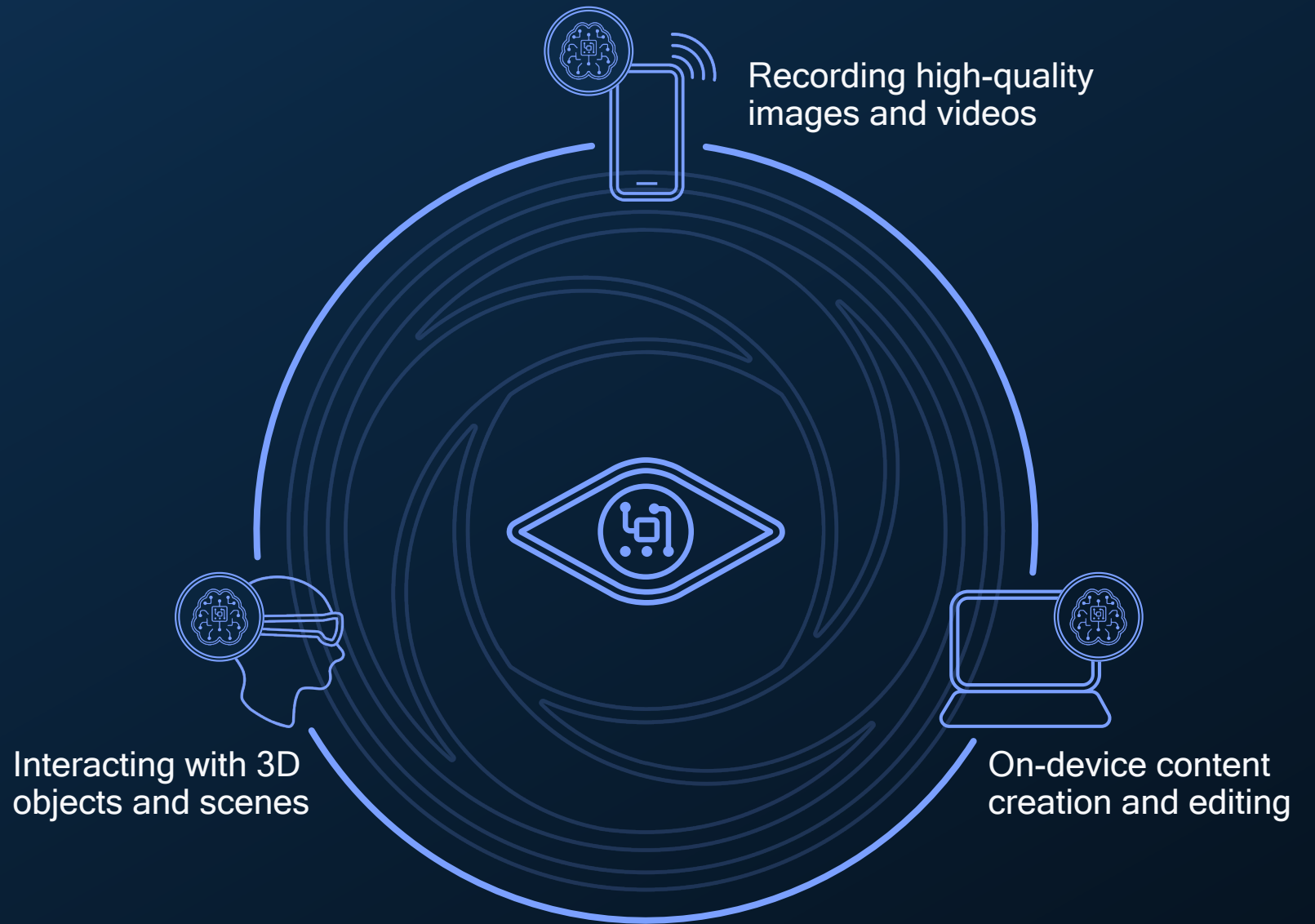
Qualcomm is a trademark or registered trademark of Qualcomm Incorporated. Other products and brand names may be trademarks or registered trademarks of their respective owners.

References in this presentation to “Qualcomm” may mean Qualcomm Incorporated, Qualcomm Technologies, Inc., and/or other subsidiaries or business units within the Qualcomm corporate structure, as applicable. Qualcomm Incorporated includes our licensing business, QTL, and the vast majority of our patent portfolio. Qualcomm Technologies, Inc., a subsidiary of Qualcomm Incorporated, operates, along with its subsidiaries, substantially all of our engineering, research and development functions, and substantially all of our products and services businesses, including our QCT semiconductor business.

Snapdragon and Qualcomm branded products are products of Qualcomm Technologies, Inc. and/or its subsidiaries. Qualcomm patented technologies are licensed by Qualcomm Incorporated.

# Efficient generative vision

Processing on the edge enables scale across devices



MVDream-SDv2.1



DreamFusion-SDv2.1



HexaGen3D-SDXL



*"a brightly colored mushroom growing on a log"*

*"a squirrel dressed like Henry VIII King of England"*

*"a hippo wearing a sweater"*

## Higher quality 3D images with HexaGen3D

Generating Hexaviews is much more effective than directly generating the triplanar latents

CLIP score ↑	<b>24.02</b> With Hexaview generation	18.47 Without Hexaview generation
--------------	--	--------------------------------------

Using the same UNet for generating and converting Hexaviews is more effective

CLIP score ↑	<b>24.02</b> With weight sharing	23.43 Without weight sharing
--------------	-------------------------------------	---------------------------------

**We generate more diverse generations (random seeds)**