

# Scaling Transformers: Architectures, Longer Contexts, Better Data

Juan Carlos Niebles



# Scaling Transformers

**Train new  
Architecture  
Designs**



**Inference for  
Long/Streaming  
Video**



**Dataset &  
Benchmarking  
Visual Generation**



# Scaling Transformers

**Train new  
Architecture  
Designs**



**Inference for  
Long/Streaming  
Video**



**Dataset &  
Benchmarking  
Visual Generation**



# Pretraining is expensive!

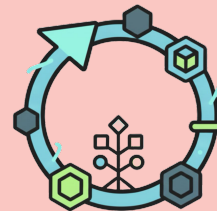
Pretraining

Handcrafted Designs



AlexNet, ResNet,  
Transformers

Automated Methods



NASNet, EfficientNet,  
EVOLUTION, STAR

Cost

# Pretraining is expensive!

Pretraining

Synthetics/ Small-scale



Zoology, MAD

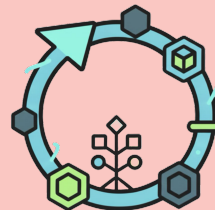
Assumptions don't  
always hold

Handcrafted Designs



AlexNet, ResNet,  
Transformers

Automated Methods



NASNet, EfficientNet,  
EVOLUTION, STAR

Cost

# Pretraining is expensive!

Pretraining

Synthetics/ Small-scale



Assumptions don't  
always hold

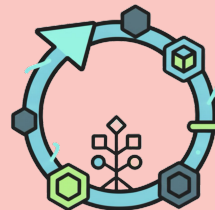
Zoology, MAD

Handcrafted Designs



AlexNet, ResNet,  
Transformers

Automated Methods



NASNet, EfficientNet,  
EVOLUTION, STAR

**WHAT DO WE WANT?**

**New (large-scale) model  
architecture designs**

**Small compute budgets**

**Cost**

# ~~Pretraining is expensive!~~ New model architectures Post-training?

Pretraining

Synthetics/ Small-scale



Assumptions don't  
always hold

Zoology, MAD

Handcrafted Designs



AlexNet, ResNet,  
Transformers

Automated Methods



NASNet, EfficientNet,  
EVOLUTION, STAR

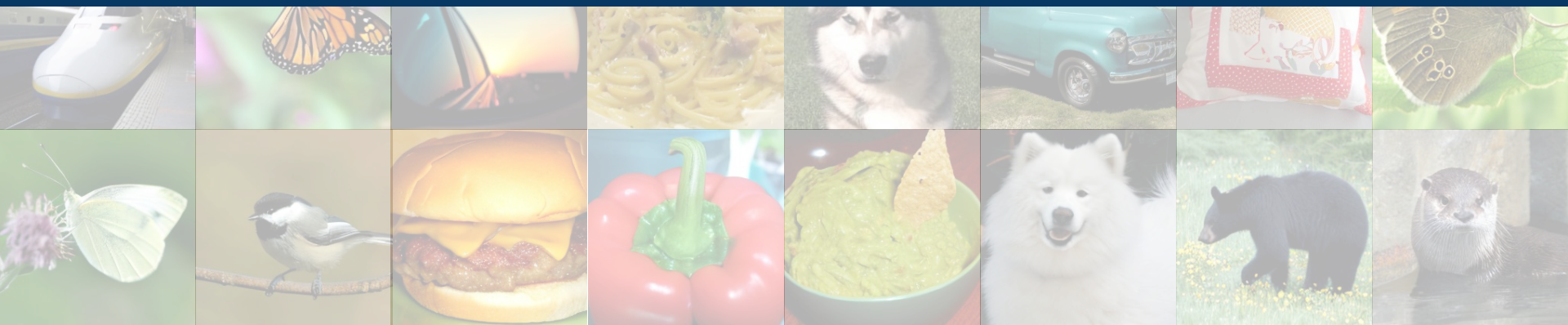
Post-training

Can we derive  
new model architectures  
using pretrained models  
under small compute budgets?

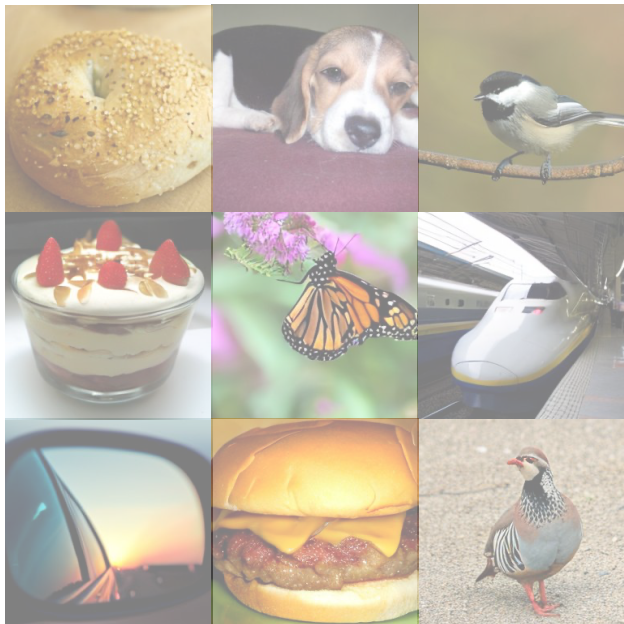
Cost



# Grafting = Model Architecture Editing



# Diffusion x Transformers



Image



Video

Boat fry bear  
rock tree ABCs  
thing sleep  
boat tennis

Text

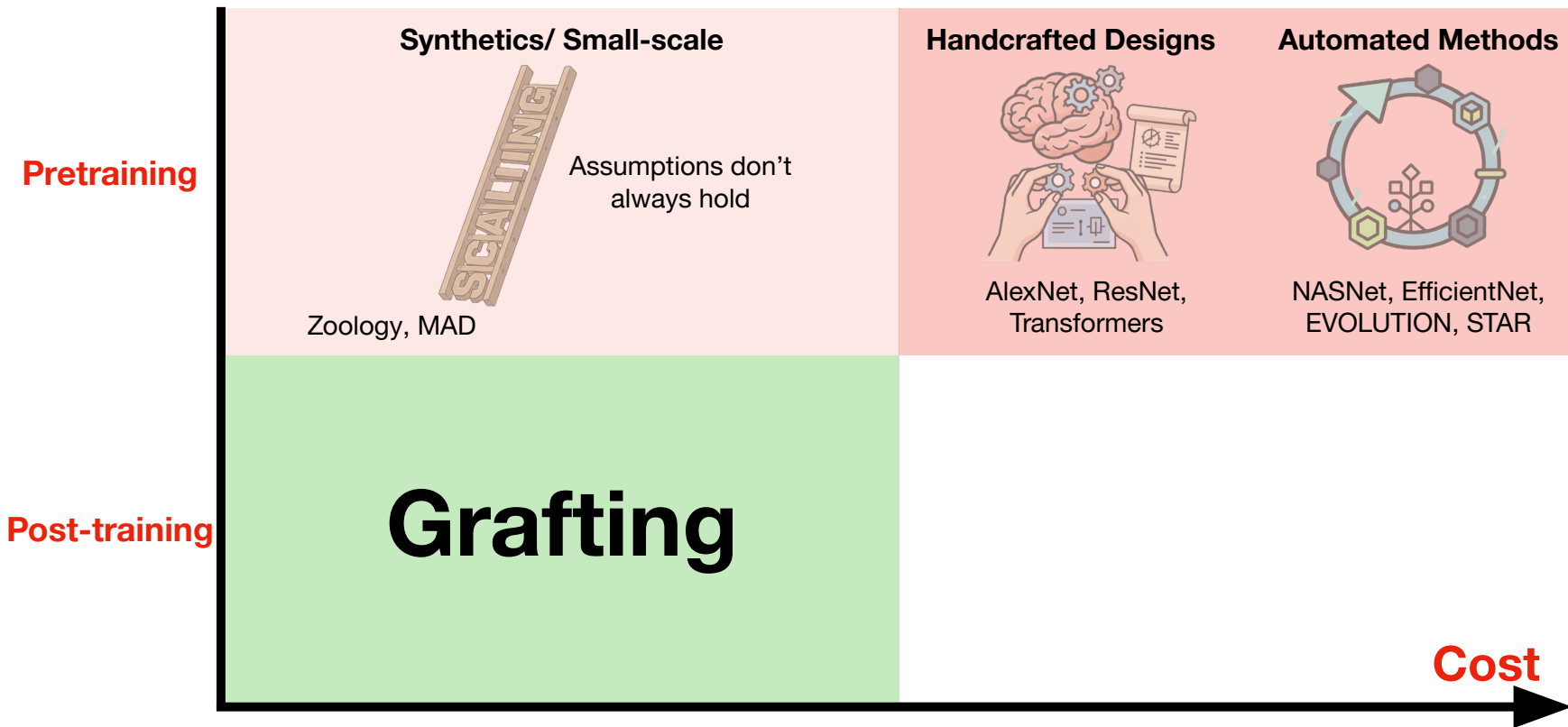
Peebles, William, and Saining Xie. "Scalable diffusion models with transformers." ICCV 2023

Brooks, Tim, et al. "Video Generation Models as World Simulators." *OpenAI*

Lou, Aaron, Chenlin Meng, and Stefano Ermon. "Discrete Diffusion Modeling by Estimating the Ratios of the Data Distribution." *ICML 2024*

[K. Chandrasegaran et al. Exploring Diffusion Transformer Designs via Grafting. NeurIPS 2025]

# ~~Pretraining is expensive!~~ New model architectures Post-training?



# Grafting: How to edit a computational graph?

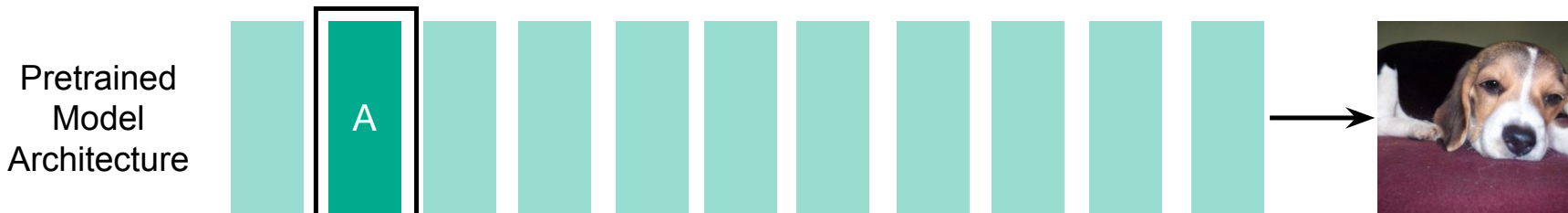
Pretrained  
Model  
Architecture



A specific  
design you care  
about!



# Core problem: How to edit a computational graph?



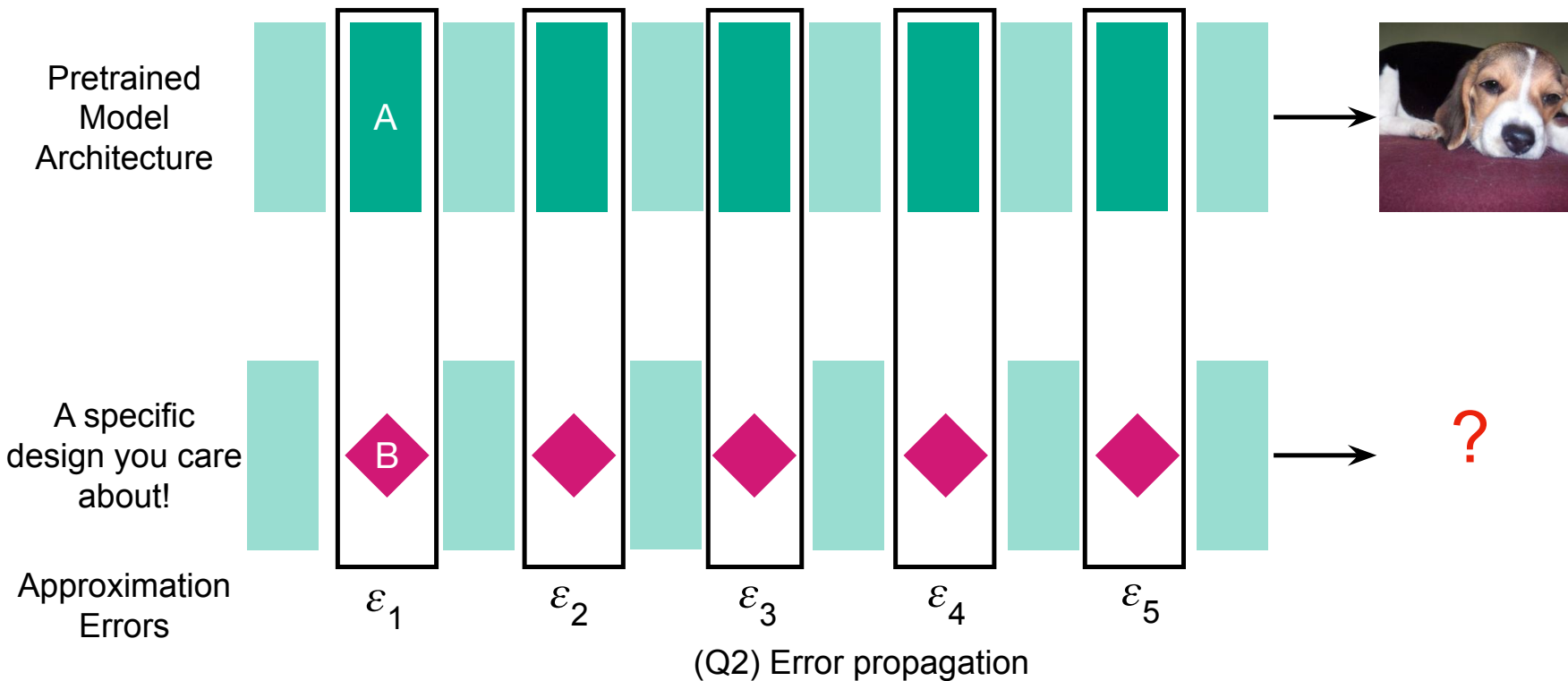
A specific design you care about!



(Q1) Operator initialization

How to initialize  $\diamond B$  such that  $\diamond B \cong \text{A}$

# Core problem: How to edit a computational graph?



How to mitigate  $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ ?

# Grafting is a simple two-stage approach to architecture editing

(Q1) Operator initialization

How to initialize  $\diamond B$  such that  $\diamond B \cong \square A$

(Q2) Error propagation

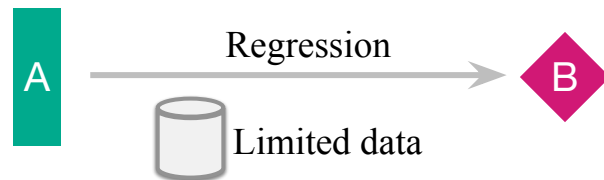
How to mitigate  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ ?

# Grafting is a simple two-stage approach to architecture editing

(Q1) Operator initialization

How to initialize  such that   $\cong$  

**(Stage 1) Activation Distillation**



(Q2) Error propagation

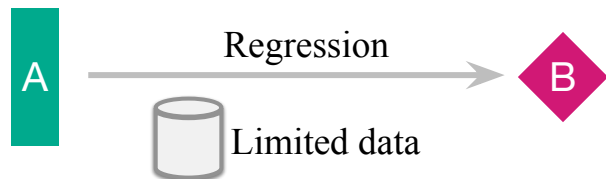
How to mitigate  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ ?

# Grafting is a simple two-stage approach to architecture editing

(Q1) Operator initialization

How to initialize  such that   $\cong$  

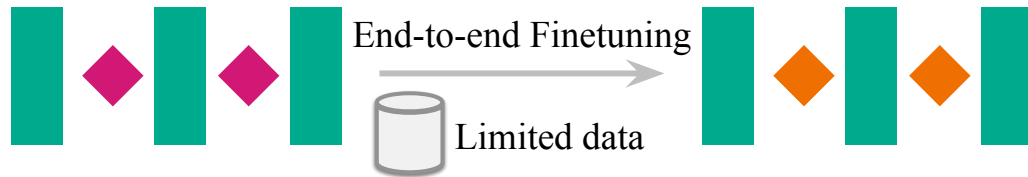
**(Stage 1) Activation Distillation**



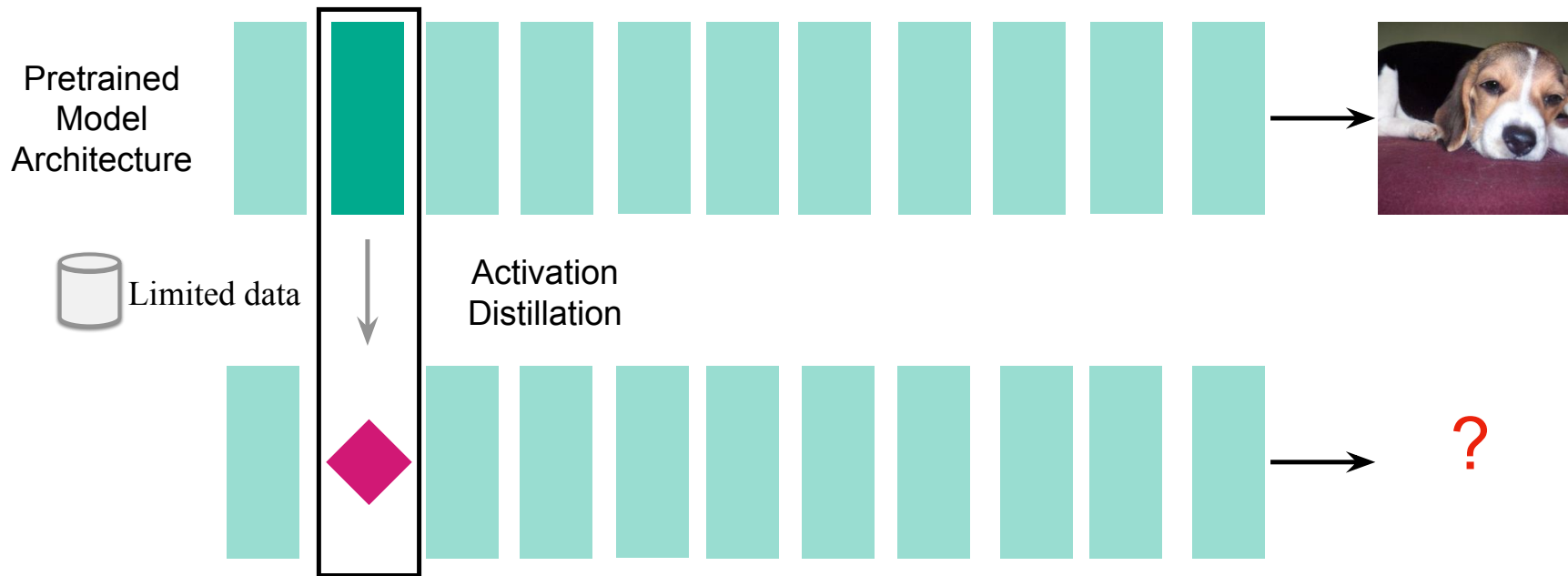
(Q2) Error propagation

How to mitigate  $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ ?

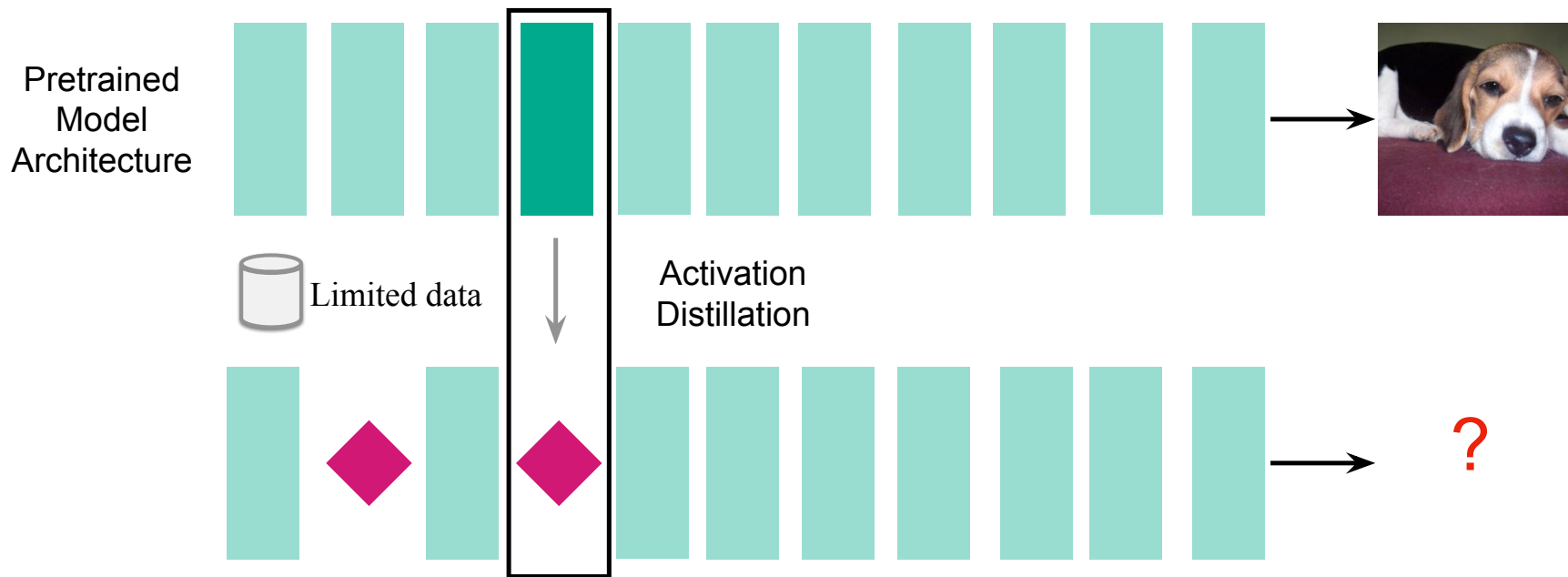
**(Stage 2) Lightweight Finetuning**



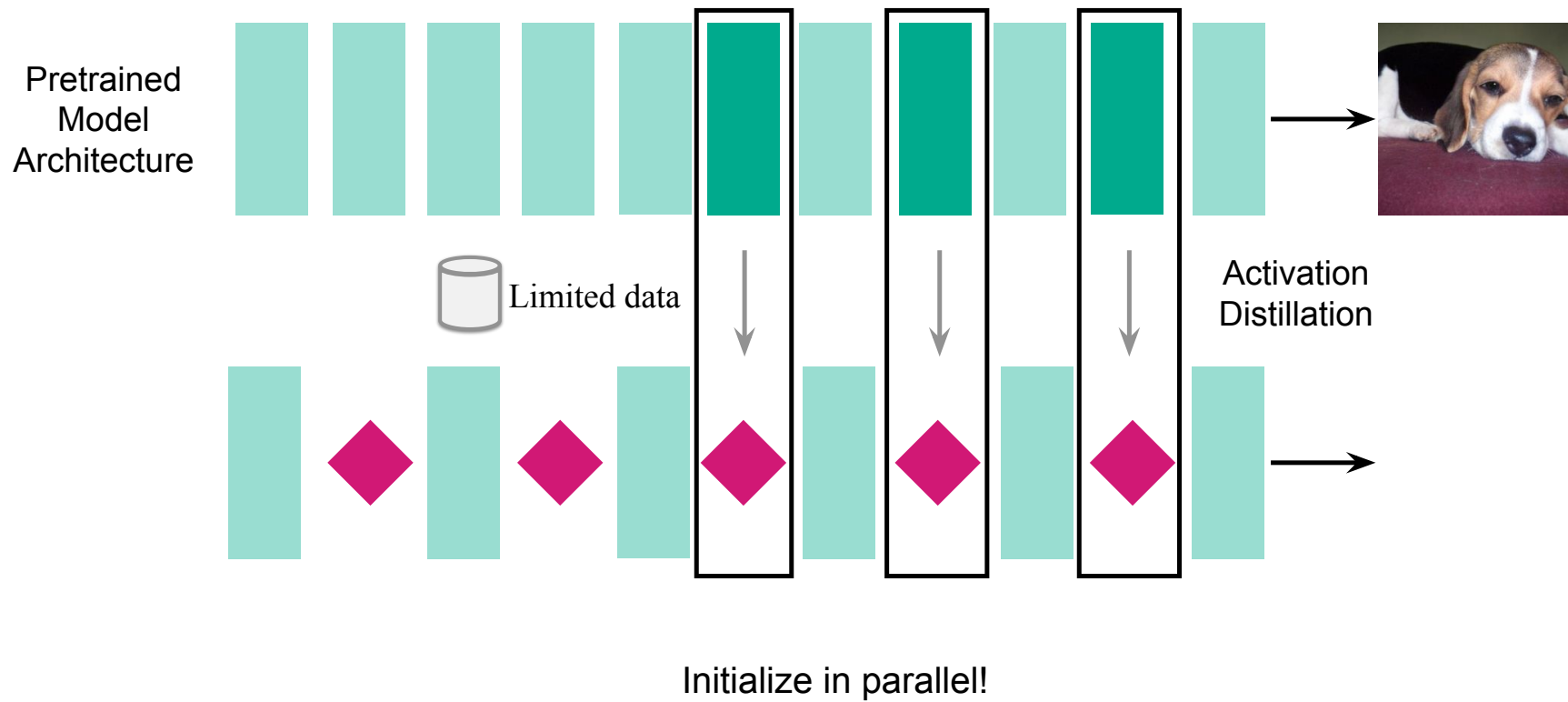
# Grafting Stage 1 - Activation Distillation



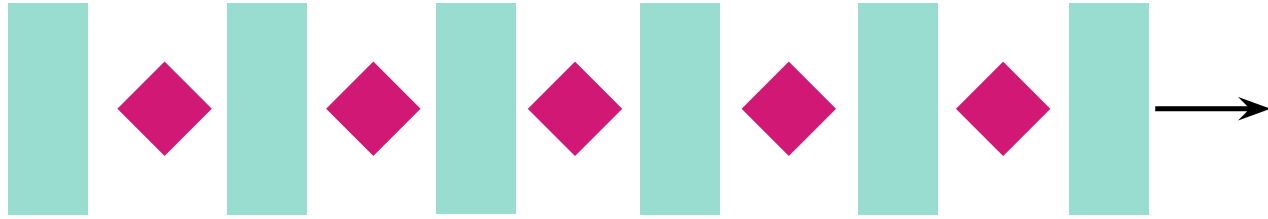
# Grafting Stage 1 - Activation Distillation



# Grafting Stage 1 - Activation Distillation



# Grafting Stage 2 - Lightweight finetuning



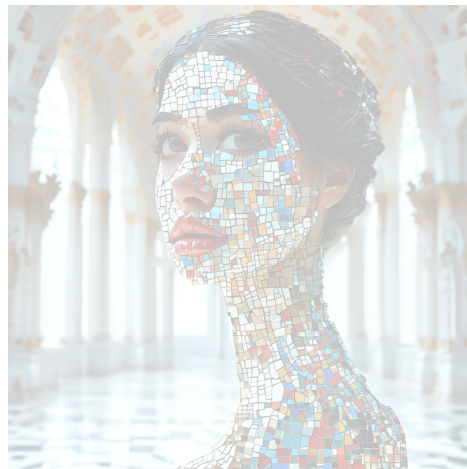
# Grafting - Stage 2



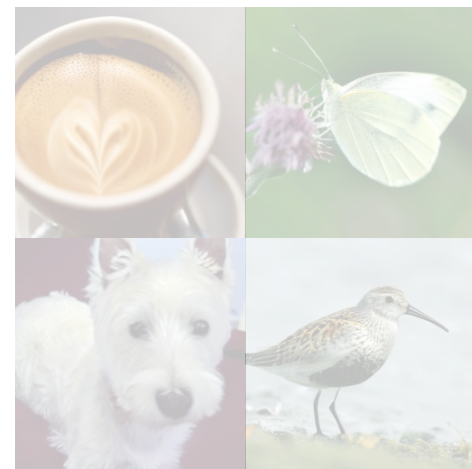
# Results



**Efficient Hybrid Architectures  
for Class-conditional  
Image Generation**



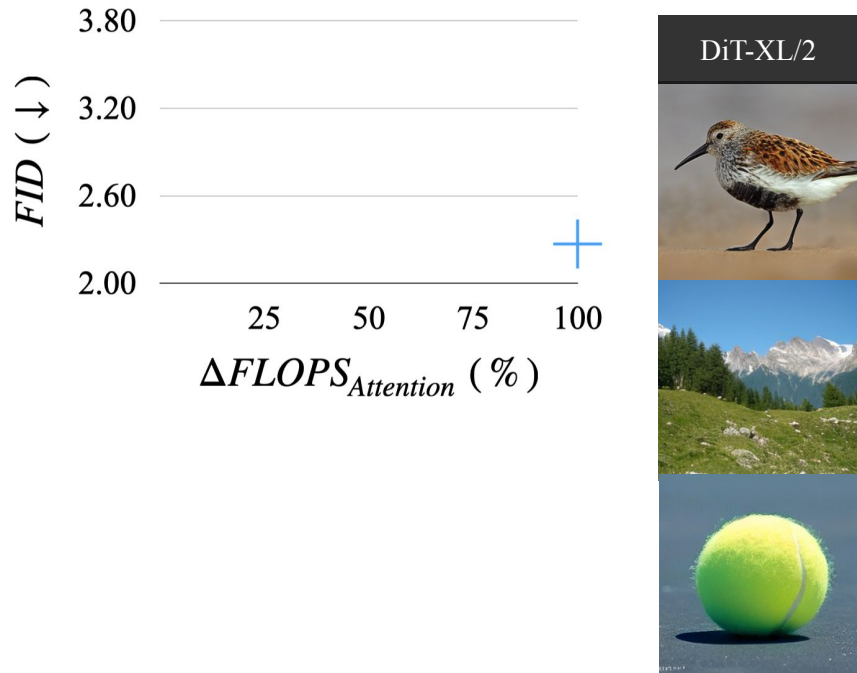
**Faster Text-to-Image  
Generation**



**Restructure Model  
Architectures (e.g., Convert  
Model Depth to Width)**

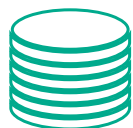
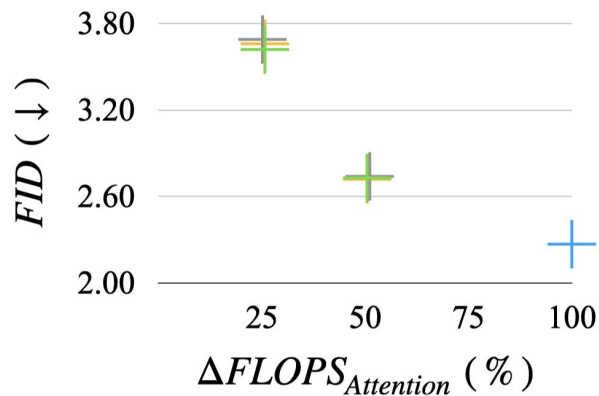
# [MHA] Experiments I: Image Generation with Hybrid Architectures obtained via Grafting

+ DiT-XL/2



# [MHA] Experiments I: Image Generation with Hybrid Architectures obtained via Grafting

+ DiT-XL/2 + Hyena-SE + Hyena-X  
+ Hyena-Y



10%



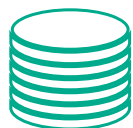
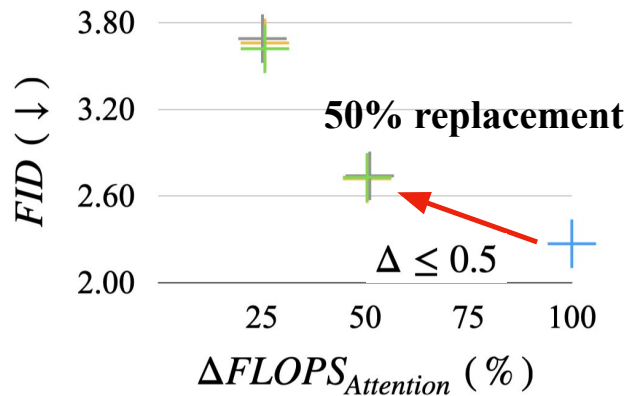
<2%



<24 hours

# [MHA] Experiments I: Image Generation with Hybrid Architectures obtained via Grafting

+ DiT-XL/2 + Hyena-SE + Hyena-X  
+ Hyena-Y



10%



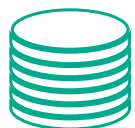
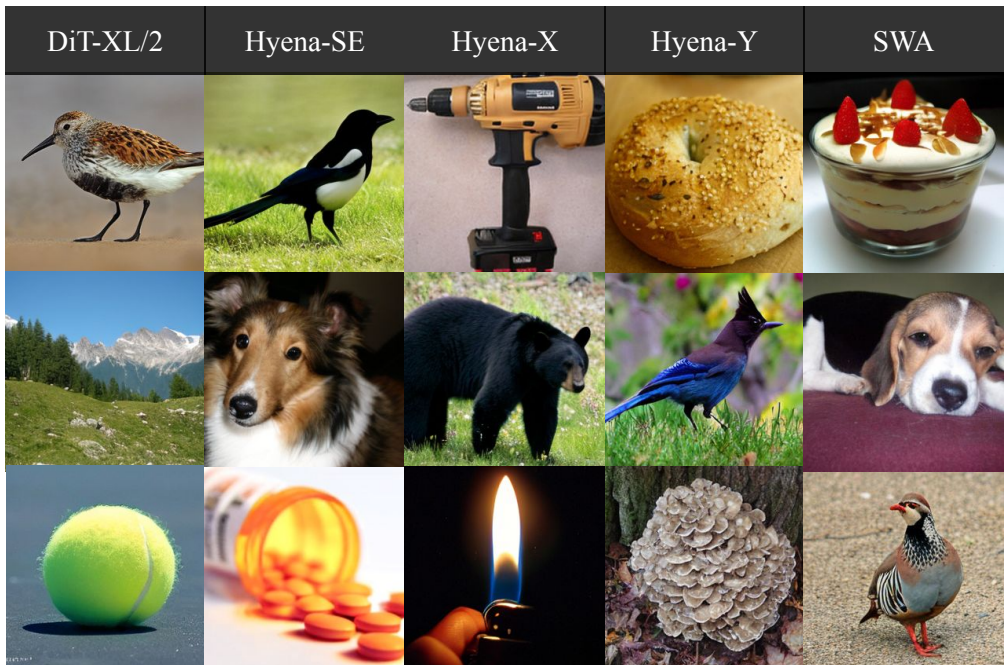
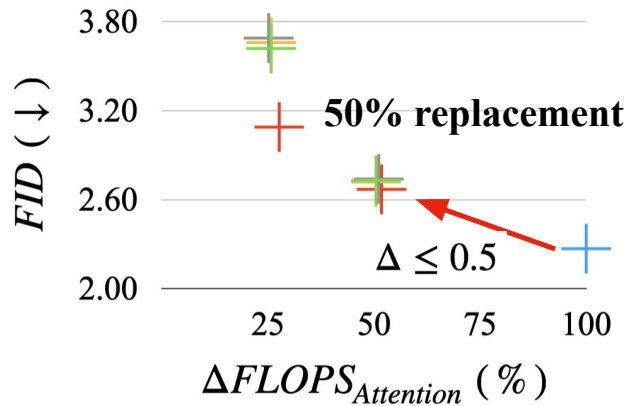
<2%



<24 hours

# [MHA] Experiments I: Image Generation with Hybrid Architectures obtained via Grafting

+ DiT-XL/2    + Hyena-SE    + Hyena-X  
+ Hyena-Y    + SWA



10%



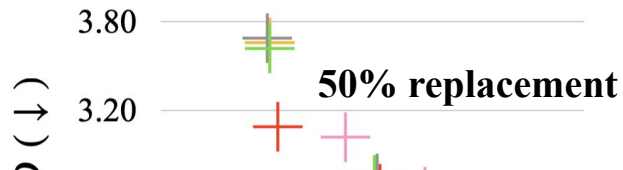
<2%



<24 hours

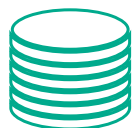
# [MHA] Experiments I: Image Generation with Hybrid Architectures obtained via Grafting

+ DiT-XL/2    + Hyena-SE    + Hyena-X  
+ Hyena-Y    + SWA    + Mamba-2



Takeaway 1: Grafting is effective for constructing efficient hybrid architectures with good generative quality under small compute budgets.

$\Delta FLOPS_{Attention}$  (%)



10%



<2%



<24 hours

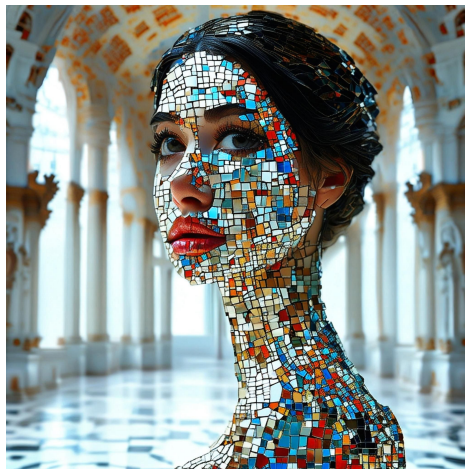


In practice, Mamba-2 is expensive due to projections (Sequence length=256 in this setup)

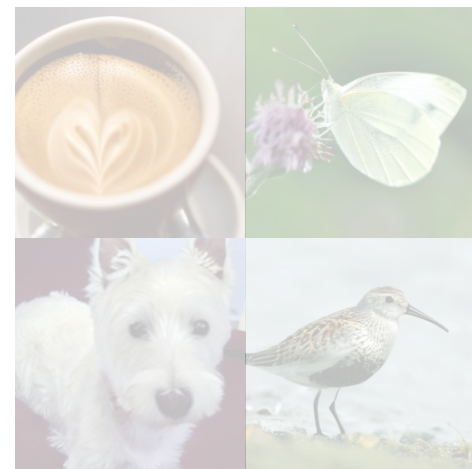
# Experiments



Efficient Hybrid Architectures  
for Class-conditional  
Image Generation



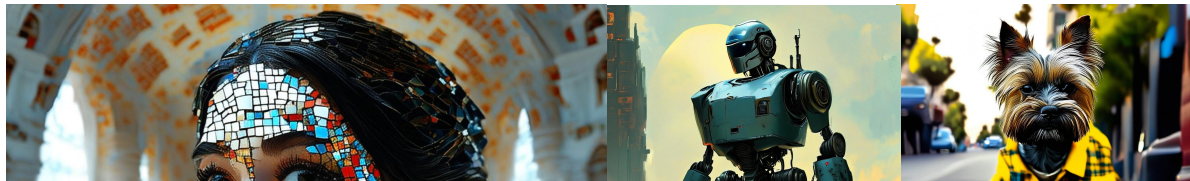
**Faster Text-to-Image  
Generation**



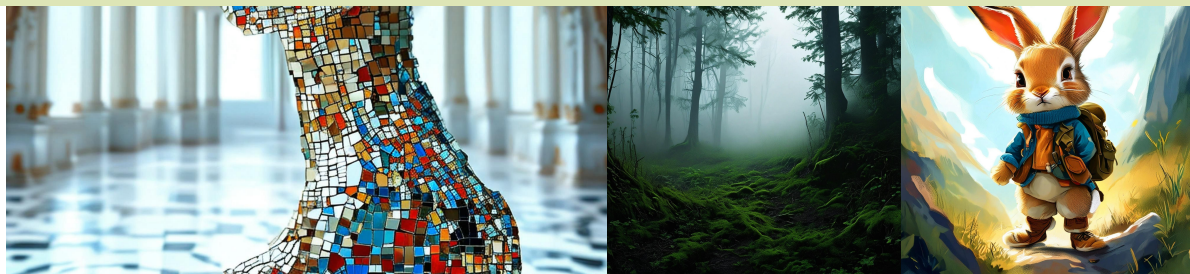
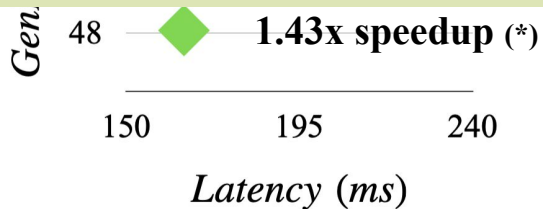
Restructure Model  
Architectures (e.g., Convert  
Model Depth to Width)

# Experiments II: Grafting Text-to-Image Diffusion Transformers

- ◆ Baseline
- ◆ Grafting (Hyena-X)



Takeaway 2: We graft high-resolution text-to-image DiTs, constructing hybrid architectures with meaningful speedups and minimal quality drop.



(\*) Speedup measured @ Batch size=2 / H100

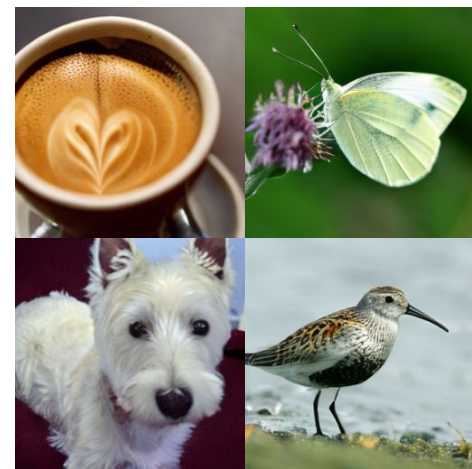
# Experiments



Efficient Hybrid Architectures  
for Class-conditional  
Image Generation



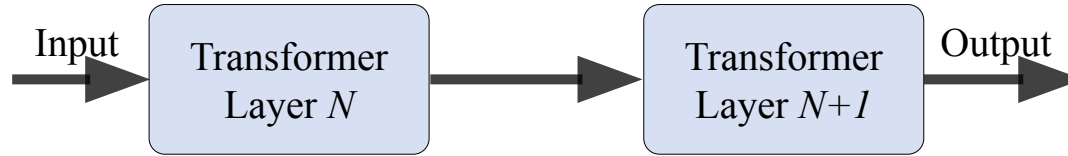
Faster Text-to-Image  
Generation



**Restructure Model  
Architectures (e.g., Convert  
Model Depth to Width)**

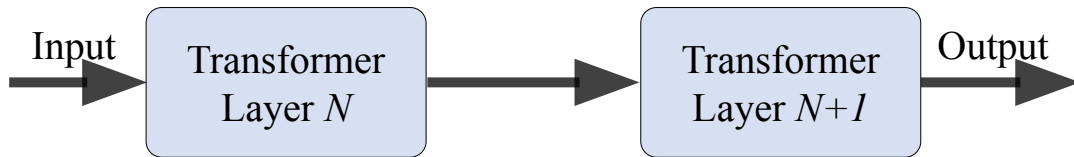
# Case study: Converting Model Depth to Width via Grafting

**(a) Sequential Transformer Blocks**

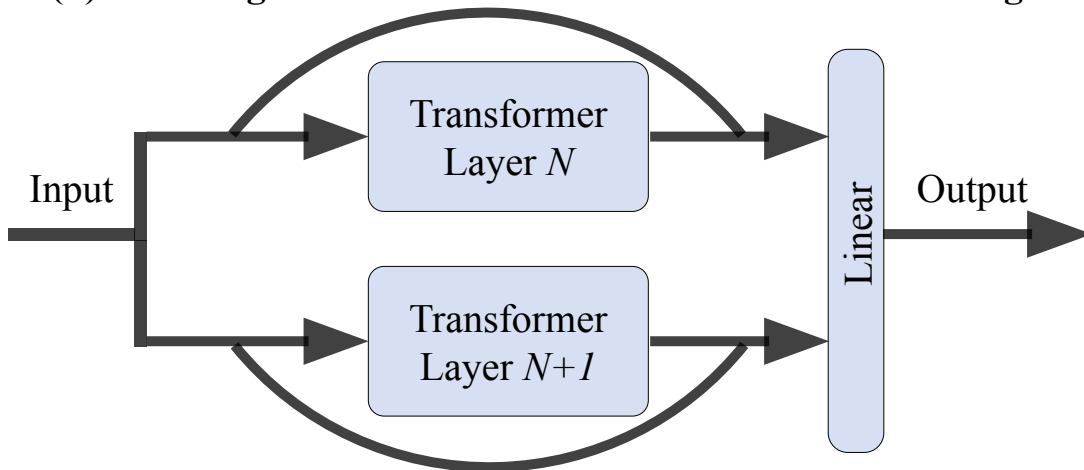


# Case study: Converting Model Depth to Width via Grafting

**(a) Sequential Transformer Blocks**

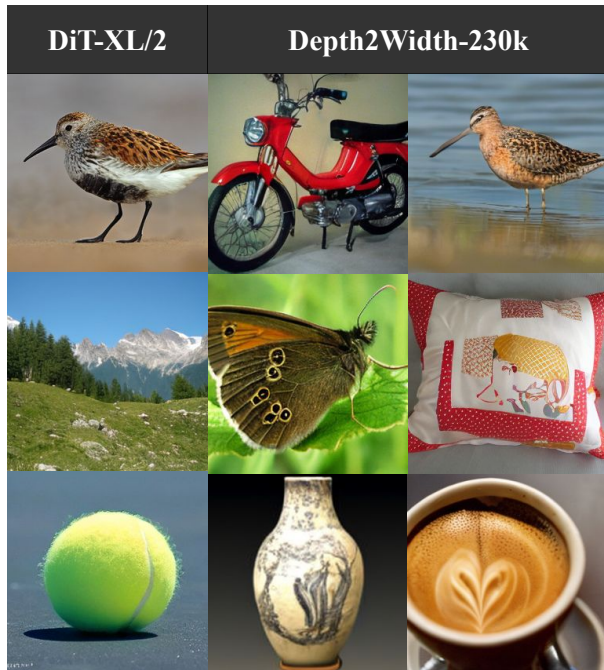
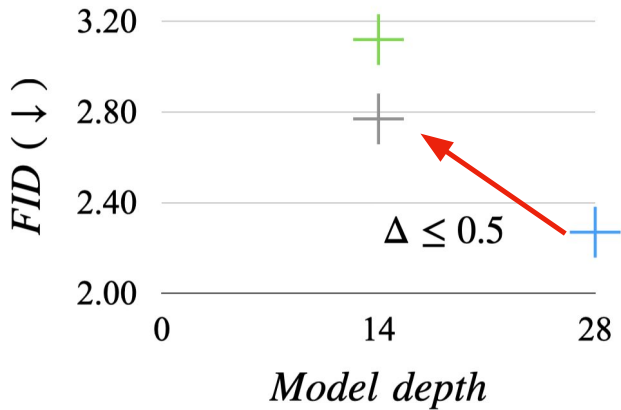


**(b) Rewiring Transformer Blocks in Parallel via Grafting**



# Case study: Converting Model Depth to Width via Grafting

- + DiT-XL/2
- + Depth2Width-230k
- + Depth2Width-100k



25%



<5%



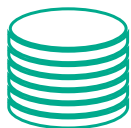
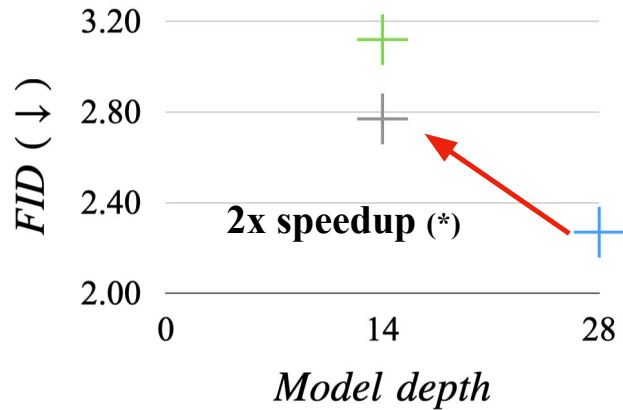
<48 hours

# Case study: Converting Model Depth to Width via Grafting

+ DiT-XL/2

+ Depth2Width-100k

+ Depth2Width-230k



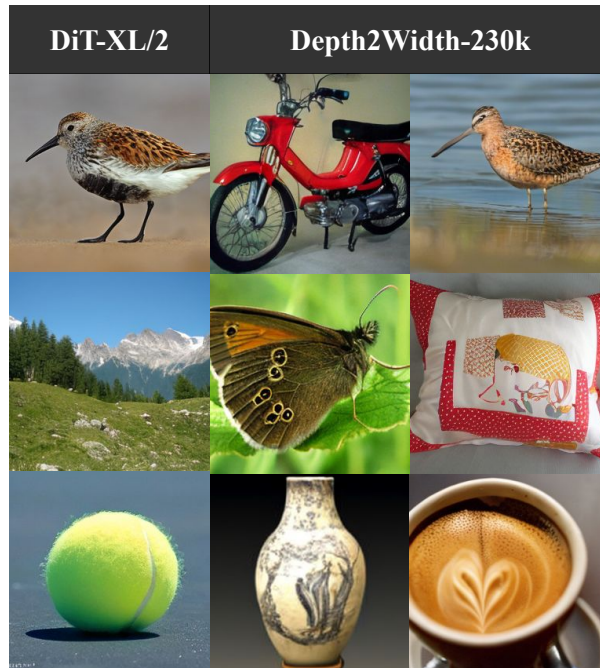
25%



<5%



<48 hours



(\*) Speedup measured @ Batch size=2 / H100

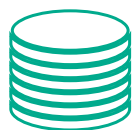
# Case study: Converting Model Depth to Width via Grafting

- + DiT-XL/2
- + Depth2Width-100k
- + Depth2Width-230k



Takeaway 3: Grafting enables architectural restructuring at the transformer block level, allowing model depth to be traded for width.

0      14      28  
*Model depth*



25%



<5%



<48 hours



(\*) Speedup measured @ Batch size=2 / H100

# Grafting - Summary

- Grafting is a simple approach to post-training architecture editing.
- Results on:
  - Hybrid Class-conditioned image generation
  - Efficient Text-to-image Generation
  - Model restructuring

Models/ code/ demo/ samples  
[grafting.stanford.edu](https://grafting.stanford.edu)



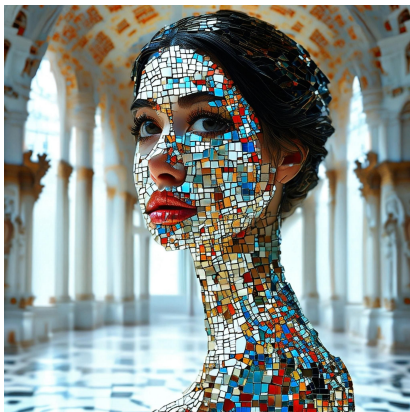
# Scaling Transformers

**Train new  
Architecture  
Designs**

**Inference for  
Long/Streaming  
Video**

**Dataset &  
Benchmarking  
Visual Generation**

[Grafting, NeurIPS 25]



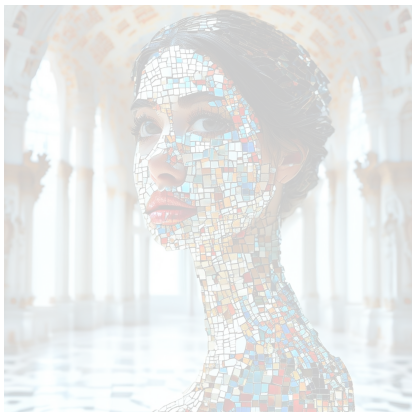
# Scaling Transformers

Train new  
Architecture  
Designs

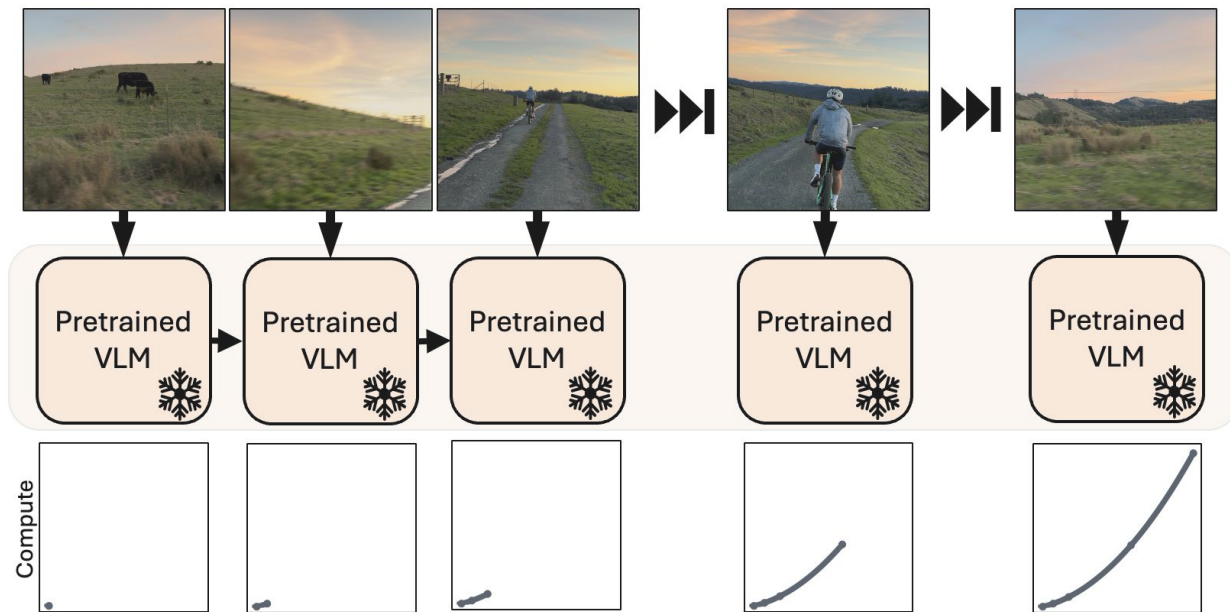
Inference for  
Long/Streaming  
Video

Dataset &  
Benchmarking  
Visual Generation

[Grafting, NeurIPS 25]



# VLMs for Video Understanding

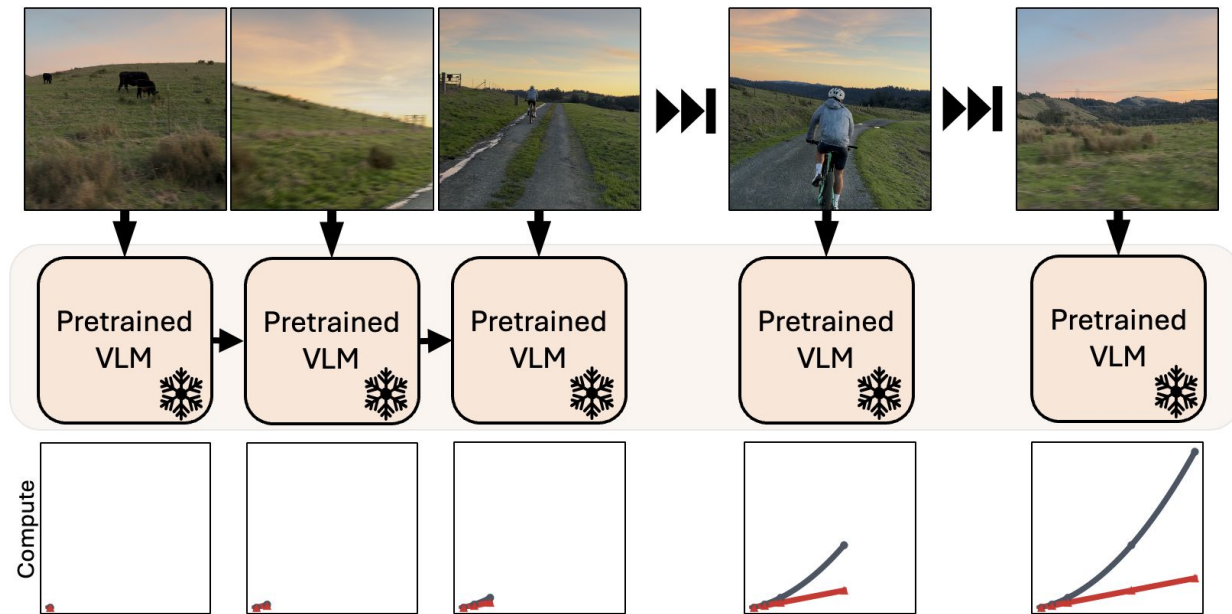


What we have

**$O(N^2)$**

*dense self-attention over video tokens*

# StateKV: linear scaling of pretrained VLMs



What we have

**$O(N^2)$**

*dense self-attention over video tokens*

What we need

**$O(N)$**

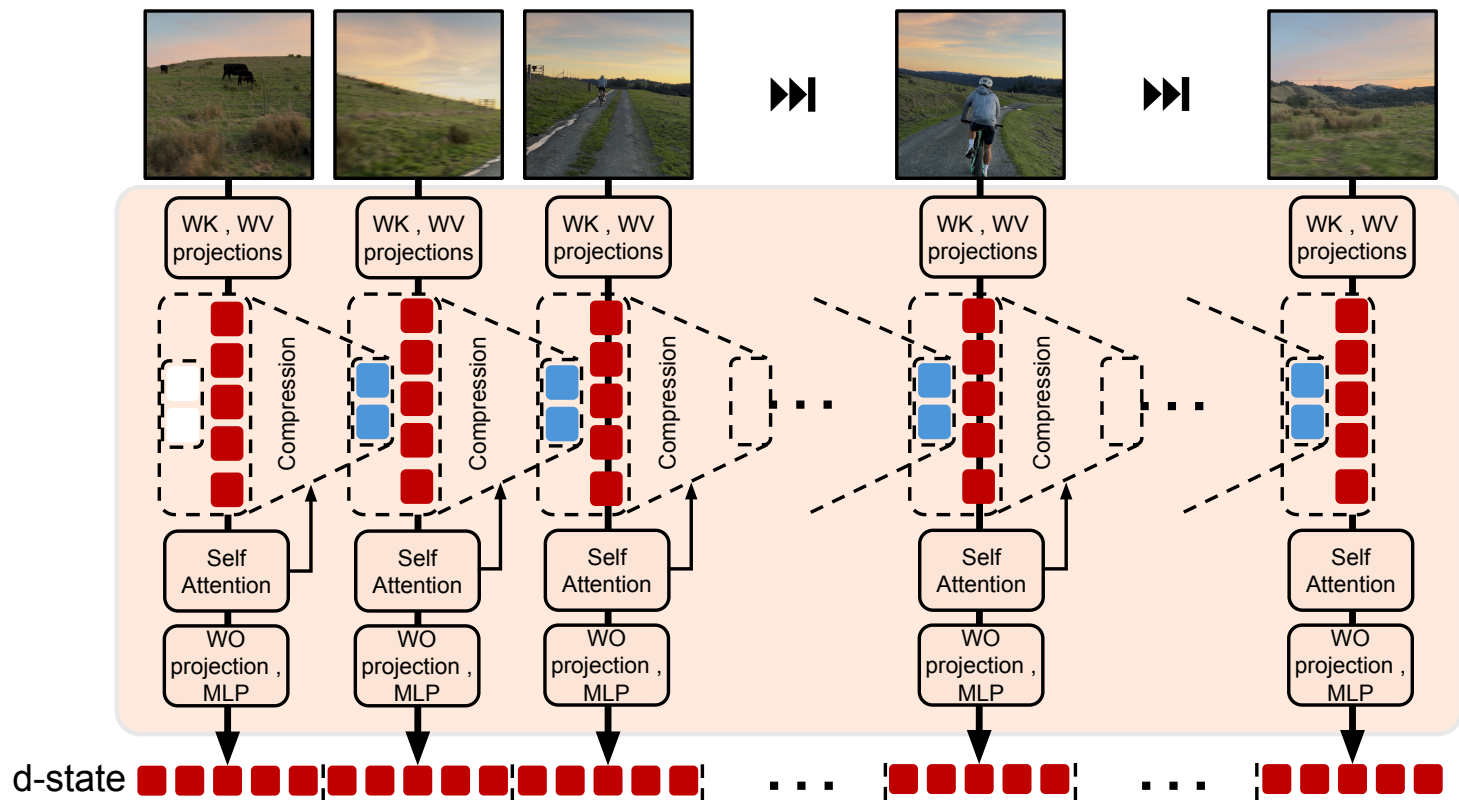
*constant cost per frame · linear total*

# StateKV: linear scaling of pretrained VLMs

Two key assumptions that we verify empirically:

- 1. Most of the past does not matter**
  - a. most cross-frame attention is concentrated on a surprisingly small subset of tokens
- 2. The important tokens change slowly**
  - a. instead of recomputing an optimal memory from scratch, we can maintain and update a compact state as the video progresses

# StateKV algorithm: compressed memory - cstate



# StateKV algorithm: update state recurrently

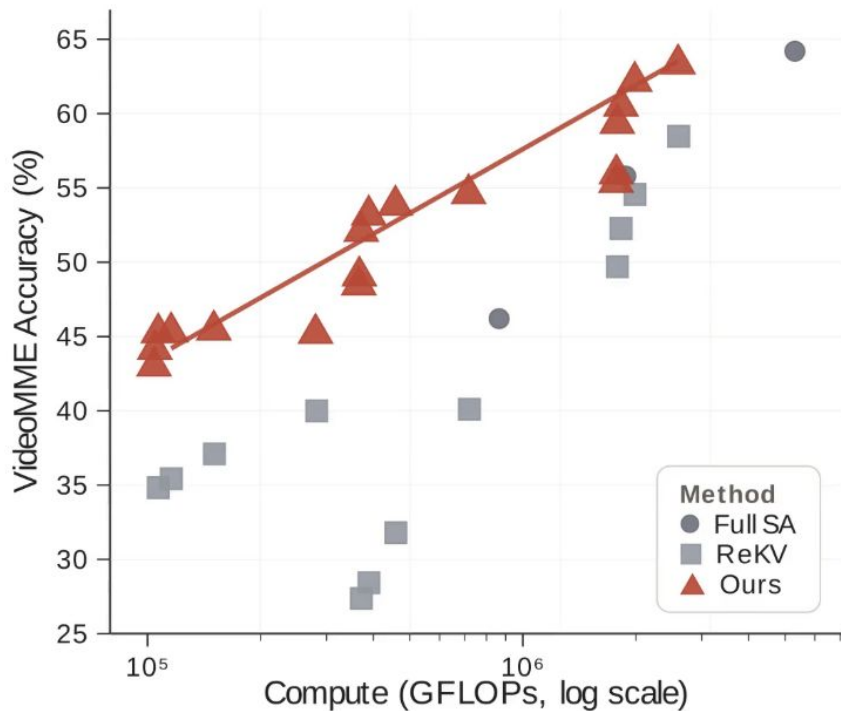


# StateKV algorithm: update state recurrently



# Accuracy–compute tradeoff

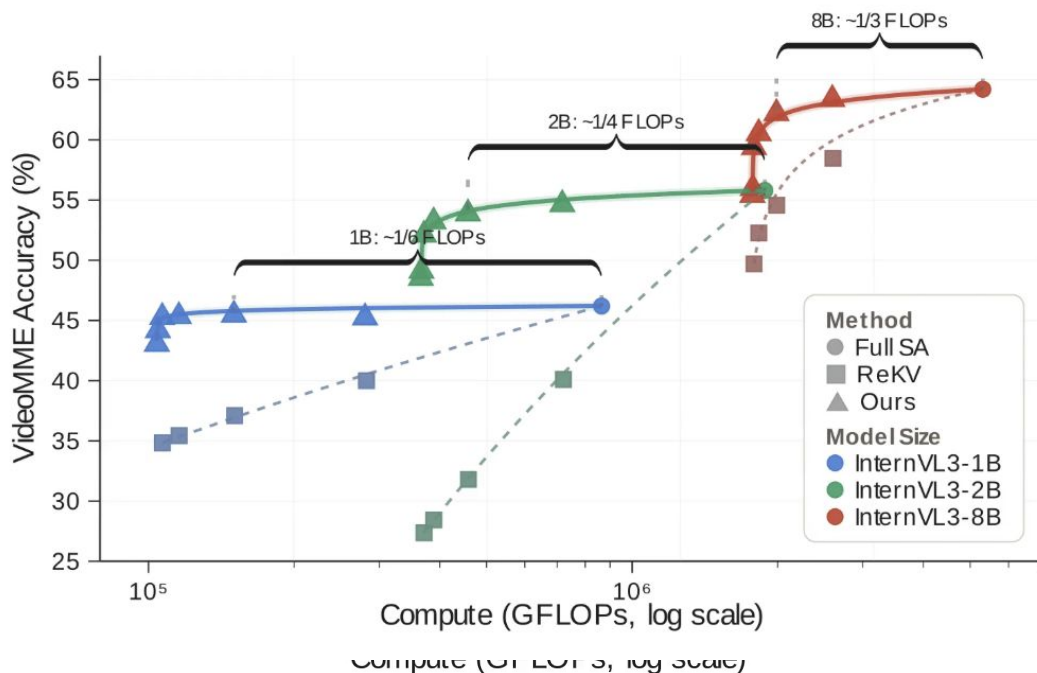
## StateKV: Compute–Accuracy Pareto Frontier



[C. Eyzaguirre et al. Linear Scaling Video VLMs for Long Video Understanding. arxiv. 2026.]

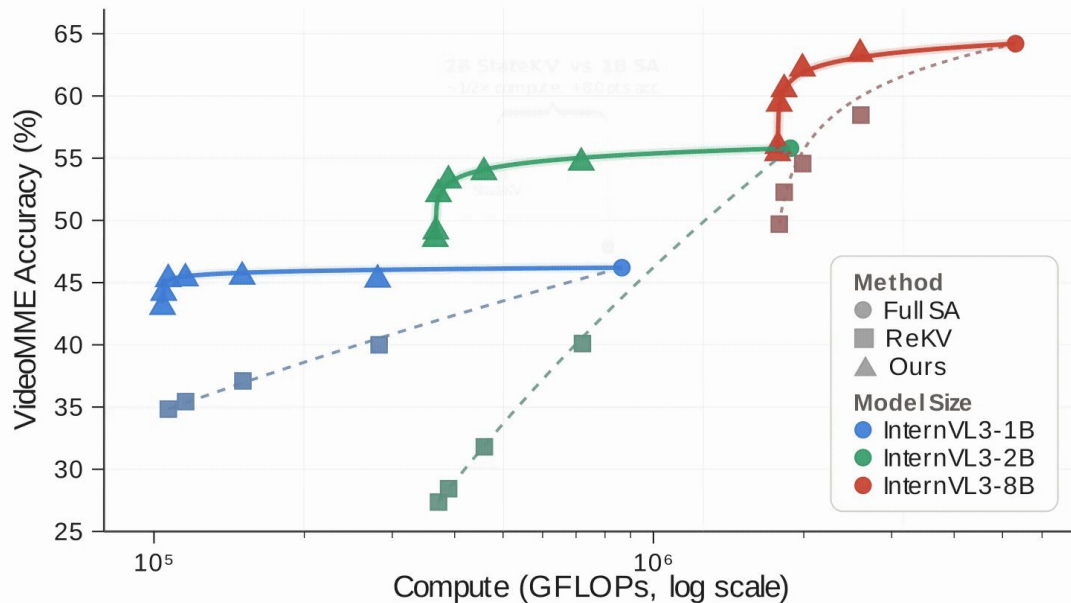
# Accuracy–compute tradeoff

## StateKV: Compute–Accuracy Pareto Frontier



# Accuracy—larger models with same compute

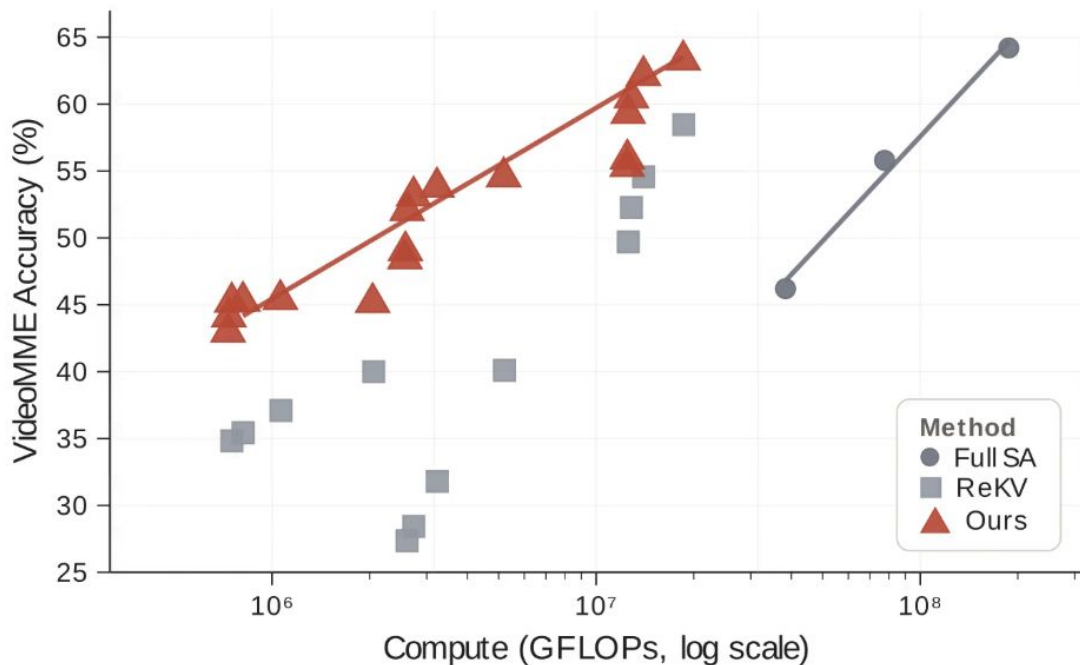
## StateKV: Compute–Accuracy Pareto Frontier



# The long-video regime

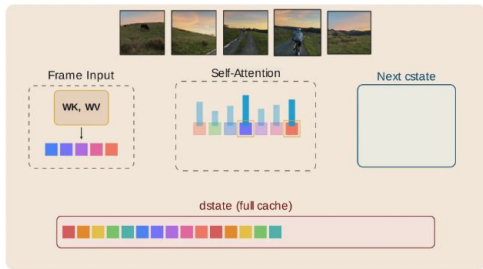
## The compute advantage grows with video length

3,600 frames (1 hr at 1 fps)



# StateKV - Summary

1. StateKV enables long-video VLM inference to scale linearly
2. It preserves most of full self-attention performance
3. Enables running larger models at a the cost of smaller models



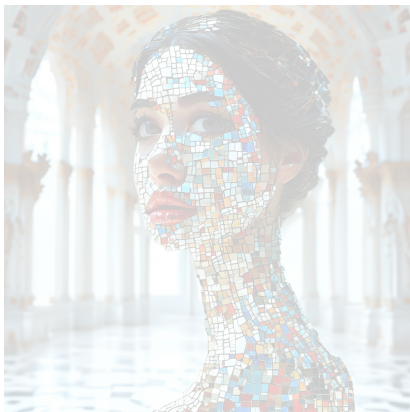
Blog / paper/ code  
[ceyzaguirre4.github.io/StateKV](https://ceyzaguirre4.github.io/StateKV)



# Scaling Transformers

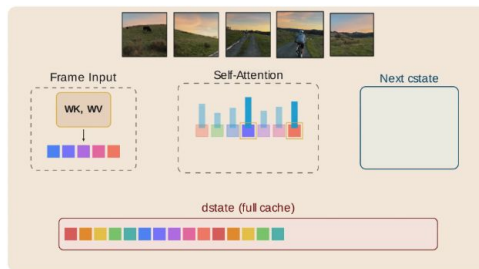
Train new  
Architecture  
Designs

[Grafting, NeurIPS 25]



Inference for  
Long/Streaming  
Video

[StateKV, arxiv 26]



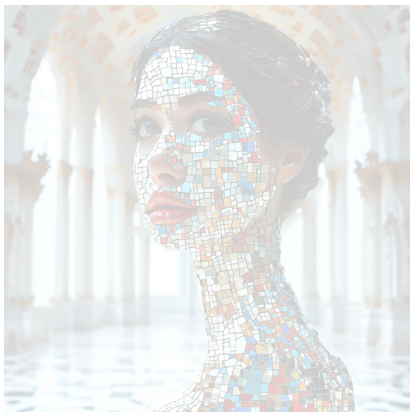
Dataset &  
Benchmarking  
Visual Generation



# Scaling Transformers

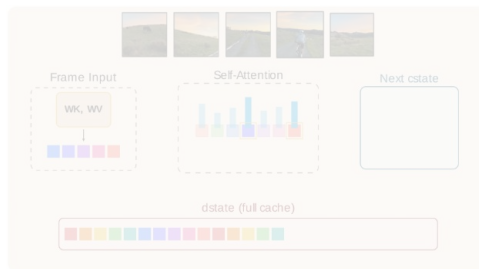
Train new  
Architecture  
Designs

[Grafting, NeurIPS 25]



Inference for  
Long/Streaming  
Video

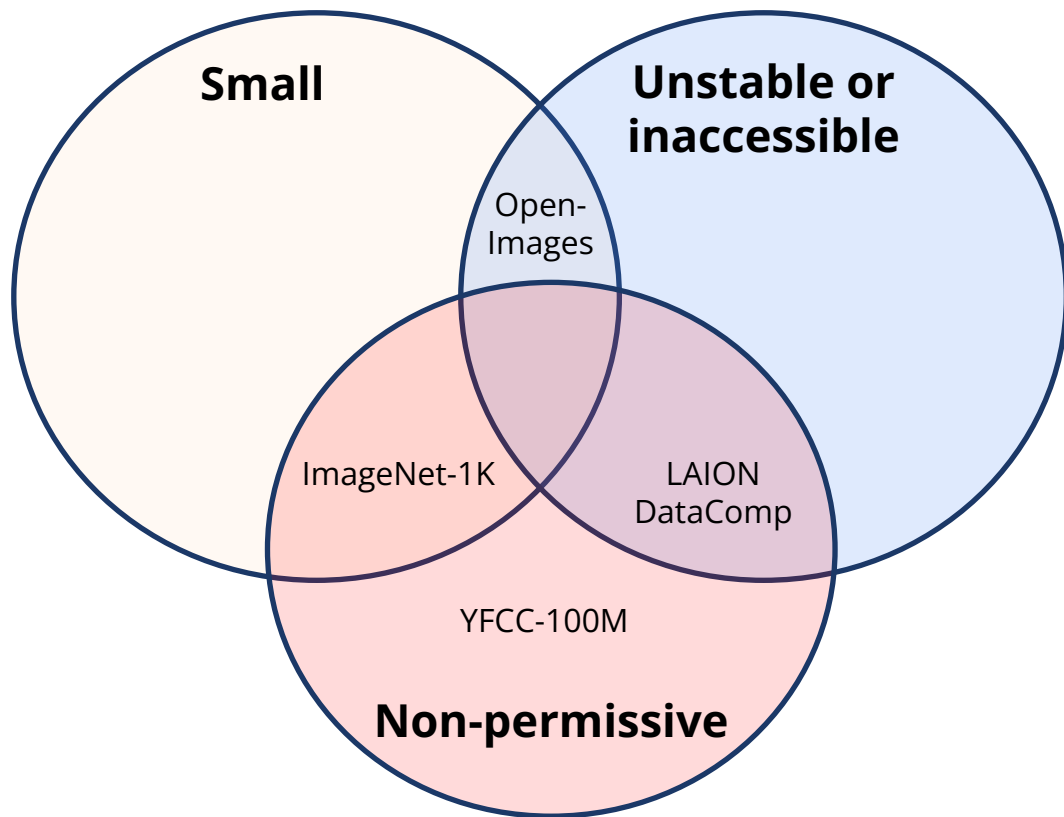
[StateKV, arxiv 26]



Dataset &  
Benchmarking  
Visual Generation



# Current datasets for visual generation



# What properties should a new dataset have?

**Small**



100M images

**Unstable or  
inaccessible**



Host the raw data in sharded  
tarfiles in a stable repo

**Non-permissive**



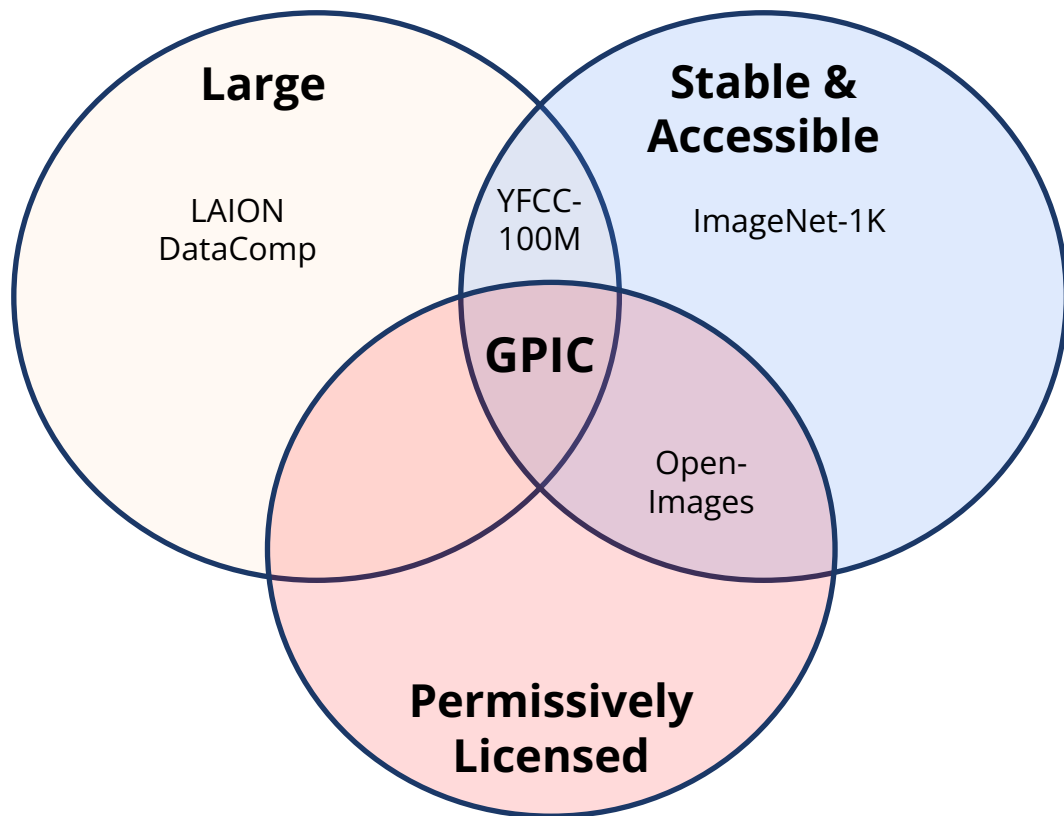
All CC, CC-BY, or Public  
domain licensed.

# *GPIC: A Giant Permissive Image Corpus for Visual Generation*

ImageNet-1K



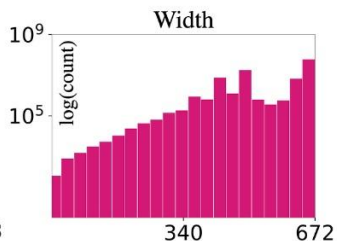
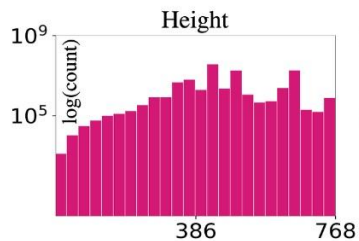
# Current datasets for visual generation



# GPIC: A Giant Permissive Image Corpus for Visual Generation

## Dataset Statistics

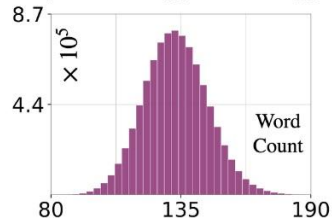
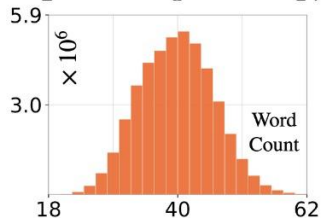
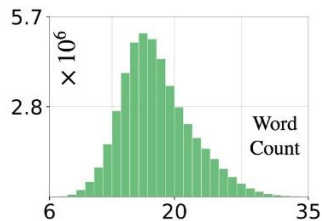
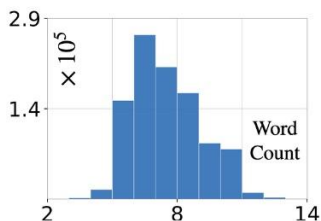
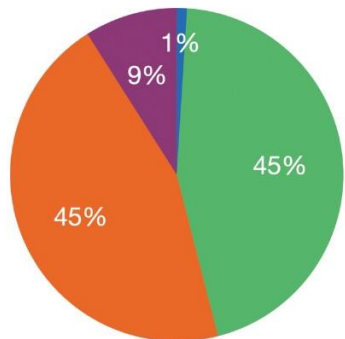
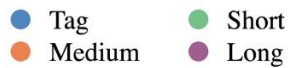
### Image Resolution



### Image License

CC BY	80.92%
Public Domain	13.85%
CC0	4.42%
No restrictions	0.81%

## Caption Statistics



## GPIC Dataset Release



Centrally Hosted



12.9TB / 8000 shards



MIT License

## GPIC Dataset Splits

Train

**100M**  
Image-Text  
Pairs

Test

**1M**  
Image-Text  
Pairs

Validation

**200K**  
Image-Text  
Pairs

## GPIC Benchmark Scales

GPIC-Full

**100M**  
Image-Text  
Pairs

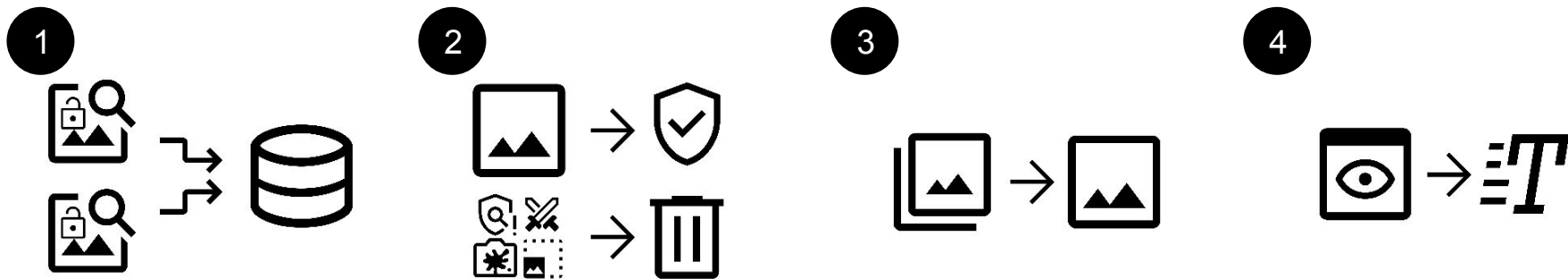
GPIC-Lite

**10M**  
Image-Text  
Pairs

GPIC-Nano

**1M**  
Image-Text  
Pairs

# GPIC Construction Pipeline



**Source Pool/  
Licensing:** 86.4 million Flickr, 15.6 million Wikimedia. Licenses in {CC, CC-BY, Public Domain} only (commercial use ok, research use ok).

**Image Filtering:** Basic image quality, extreme aspect ratios/sizes and Safety + NSFW.

**Deduplication:** SSCD copy-detection to remove near-duplicates, and verified with SHA-256 hashing

**Caption Generation:** Qwen3-VL-4B-Inst across four formats: tag, short, medium, and long.

# Image-Text Examples from GPIC



Medium

A motocross rider in full gear rides an orange dirt bike through muddy terrain. Another rider is visible in the background on a similar trail, with trees and open field surrounding the track.

# Image-Text Examples from GPIC



Medium

The back of a red and white FDNY ambulance is parked on a street. It has “KEEP BACK” written in bold letters and emergency lights along the top, with trees and a building visible in the background.

# Image-Text Examples from GPIC



Short

A gray fighter jet with "LN" and "301" markings sits on an airfield runway.

# Image-Text Examples from GPIC



Short

A bird perches on a bare tree branch against a cloudy sky. Distant mountains are visible in the background.

# Image-Text Examples from GPIC



Short

A person blows a large flame into the dark night while others watch from lounge chairs.

# Image-Text Examples from GPIC



Long

A two-story house with a beige exterior and a dark green metal roof stands at the center of the scene, featuring white trim around its windows and porch. A wide front porch with white railings extends across the lower level, and a set of stairs leads up to the entrance. To the right of the house, a bright pink pop-up canopy is set up on a grassy lawn, with a few people standing nearby under it. The house is surrounded by lush green trees and manicured bushes, with some flowering plants in red and pink near the walkway. A paved path curves from the left side toward the house's front steps, and a trash can sits near the bushes. In the background, more trees and a faint string of lights are visible, suggesting an outdoor gathering or event.

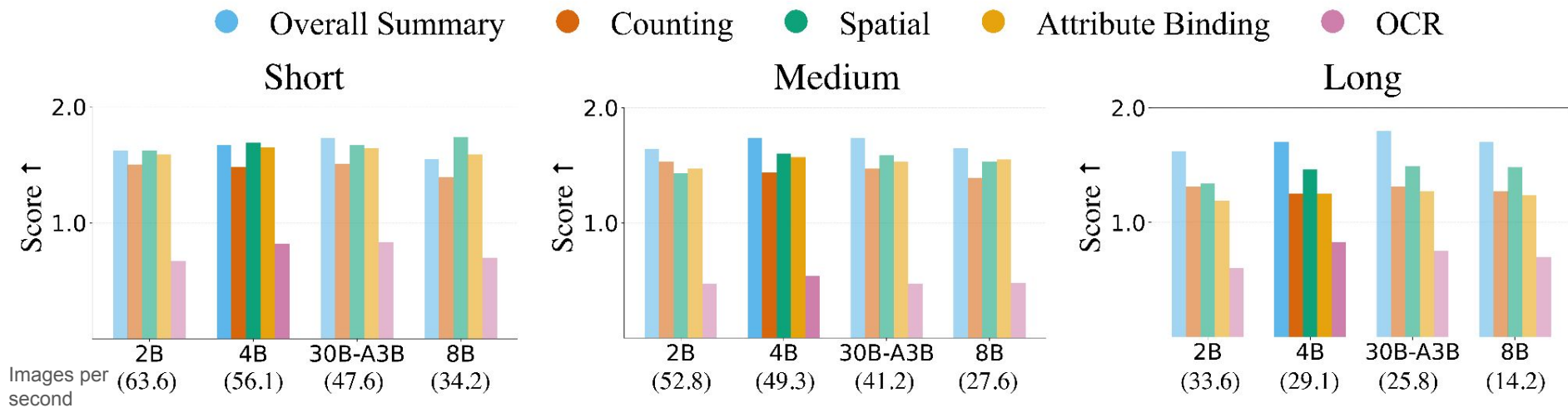
# Image-Text Examples from GPIC



Tag

blue car, yellow car, showroom, people, vintage

# Caption Quality



Qwen3-VL at 2B, 4B, 8B, and 30B scales.

Throughput (images per second, 1xH100) is shown in parentheses for each model.

# GPIC Benchmarking Protocol for Image Generation

## ImageNet-1K Benchmark

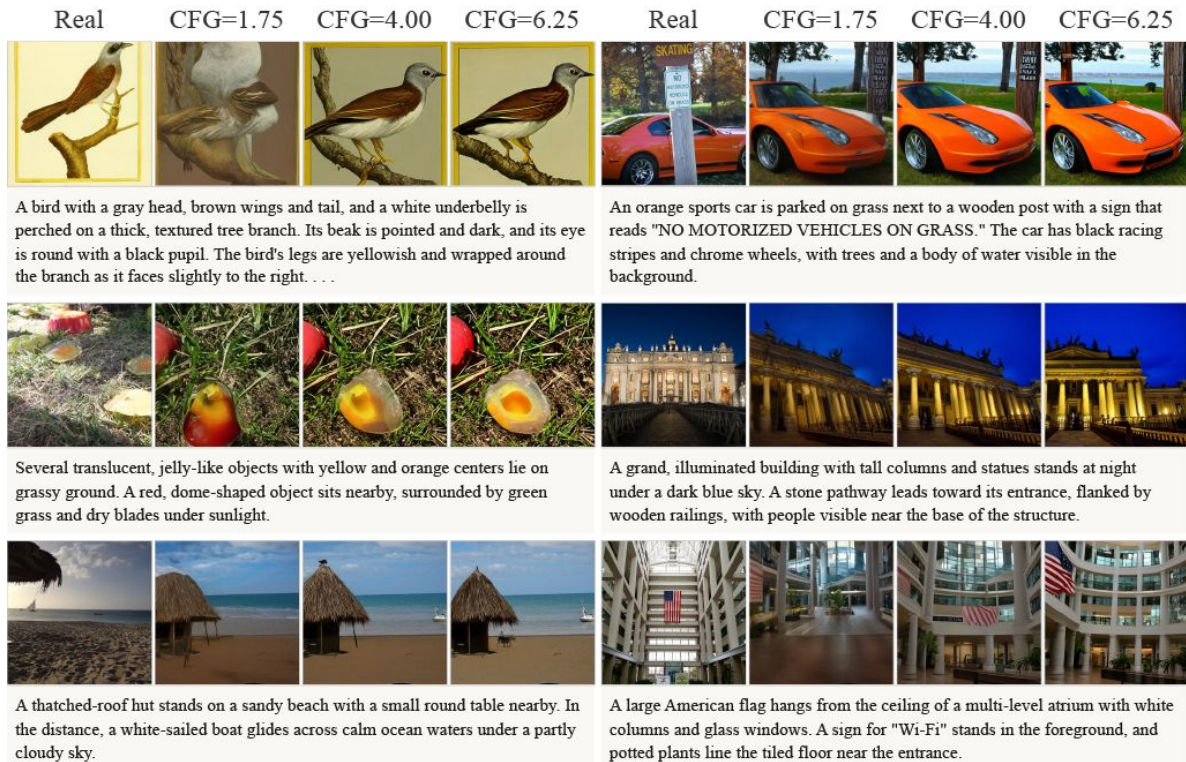
- Uses FID which is saturated.
- Many recent models now obtain lower FID than held-out real images.
- Reference statistics are computed over the ImageNet-1K training set.

## GPIC Benchmark

- Uses FD-DINOv2, which is unsaturated.
- FD-DINOv2 correlates better with human judgments than FID (Stein et. al).
- Reference statistics are computed over the 1M GPIC test set.

# Can we train a model on GPIC? Yes:

JiT-T2I (1.1B), pixel-space flow matching, trained for one epoch on GPIC-Full.



# Can we train a model on GPIC? Yes:

JiT-T2I (1.1B) [PixelGen-XXL/16, Ma et al arxiv 2026]

pixel-space flow matching, trained for one epoch on GPIC-Full.

<b>CFG</b>	<b>FD</b> ↓	<b>Precision</b> ↑	<b>Recall</b> ↑	<b>Density</b> ↑	<b>Coverage</b> ↑
1.75	204.01	0.917	0.530	1.034	0.806
4.00	87.80	0.933	0.765	1.012	0.906
6.25	76.25	0.942	0.792	1.014	0.908

Table 3: **JiT-T2I baseline results on GPIC-Full after training for one epoch.** We report FD, Precision, Recall, Density, and Coverage for three classifier-free guidance scales. We use 50-step Euler sampling for all generations. All metrics are computed against the 1M GPIC test set.



# *GPIC: A Giant Permissive Image Corpus for Visual Generation*

Dataset, Benchmark, Models

GPIC Dataset Release

**Start  
pretraining  
today!**

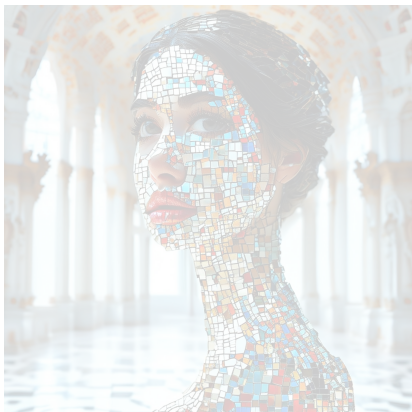
ation  
Image-Text  
Pairs



# Scaling Transformers

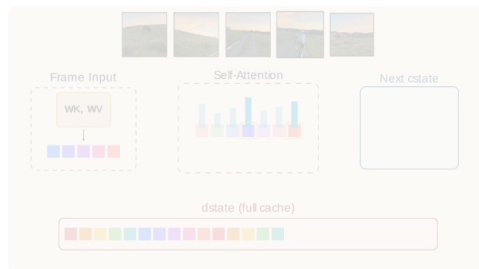
Train new  
Architecture  
Designs

[Grafting, NeurIPS 25]



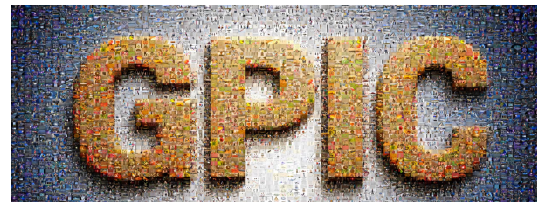
Inference for  
Long/Streaming  
Video

[StateKV, arxiv 26]



Dataset &  
Benchmarking  
Visual Generation

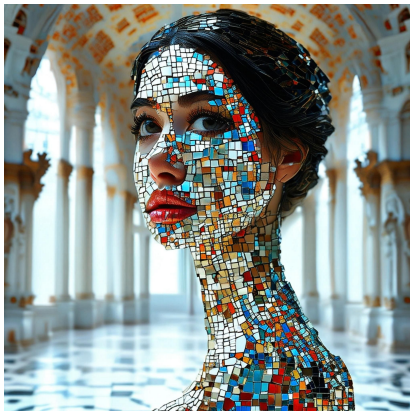
[GPIC, arxiv 26]



# Scaling Transformers

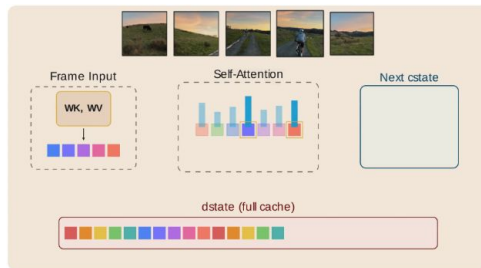
**Train new  
Architecture  
Designs**

[Grafting, NeurIPS 25]



**Inference for  
Long/Streaming  
Video**

[StateKV, arxiv 26]



**Dataset &  
Benchmarking  
Visual Generation**

[GPIC, arxiv 26]



# Thank you!



- K. Chandrasegaran, M. Poli, D. Y Fu, D. Kim, L. M. Hadzic, M. Li, A. Gupta, S. Massaroli, A. Mirhoseini, J. C. Niebles, S. Ermon, L. Fei-Fei. **“Exploring Diffusion Transformer Designs via Grafting”**. NeurIPS 2025
- C. Eyzaguirre, J. Wu, and J. C. Niebles. **Linear Scaling Video VLMs for Long Video Understanding**. arxiv. May 2026
- K. Chandrasegaran, K. Sargent, S. Agarwal, M. Jang, M. Poli, J. C. Niebles, J. Johnson, J. Wu, and L. Fei-Fei. **GPIC: A Giant Permissive Image Corpus for Visual Generation**. arxiv. May 2026

 [slides/code/data/models/papers](#)

