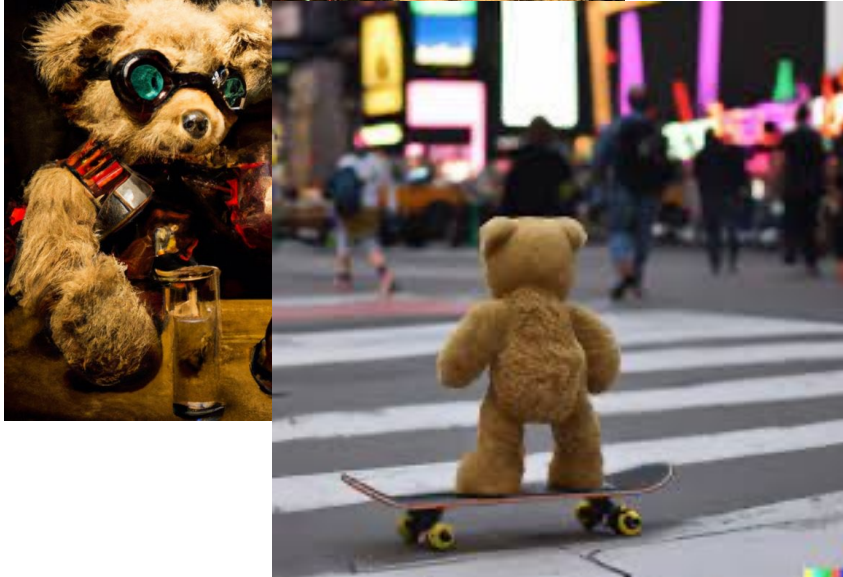# Image Editing with Optimization (part II)

Jun-Yan Zhu
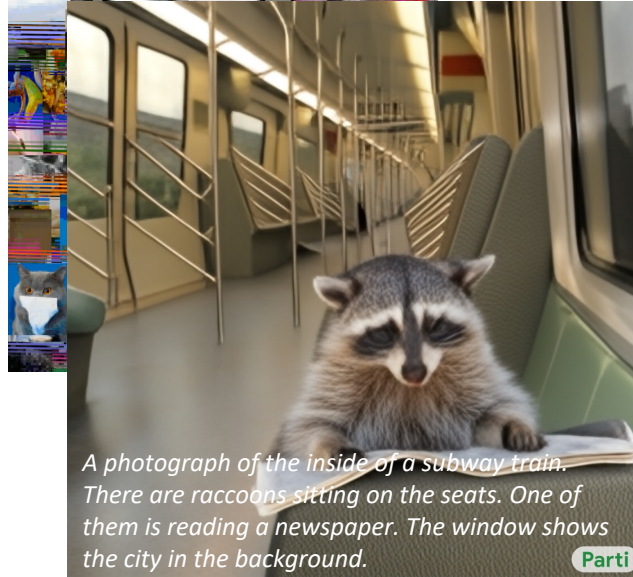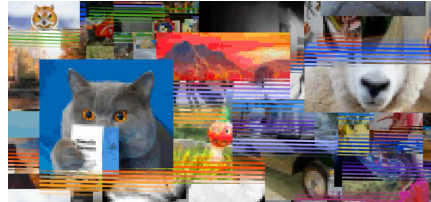
16-726 Spring 2025

# Large-scale Text-to-Image Models



Diffusion models
(DALL-E 2, Imagen, SD)

Autoregressive models
(Image GPT, Parti)

GANs, Masked GIT
(GigaGAN, MUSE)

# Limitations of Text-to-Image Models

<u>Linguistic bottleneck</u>: not everything can be described by text

<u>Data bottleneck</u>: many things are not included in the dataset:

1. Not in the public domains (e.g., personal concepts)
2. Have not been created (e.g., new concepts)

# Text-to-image isn't perfect…

Stable Diffusion



Photo of a moongate

# Text-to-image isn't perfect…



Stable Diffusion

Actual **moongate** images

Photo of a **moongate**

# Text-to-image isn't perfect…

Stable Diffusion



Actual moongate images



Photo of a moongate

# Customization



Actual **moongate** images

Stable Diffusion



Photo of a `moongate`

# Customization



Actual **moongate** images

Customized
Diffusion

Photo of a **moongate**

# Unseen contexts



Actual **moongate** images

Customized Diffusion

**Moongate** in the middle of highway

# Unseen contexts



Actual **moongate** images

Customized Diffusion

**Moongate** in snowy ice

# Unseen contexts



Actual **moongate** images

Customized Diffusion

A puppy in front of Moongate

No knowledge of personal concepts

My dog, Stark

Stable Diffusion

A dark grey color weimaraner dog

# Customization



Jun-Yan's **dog**, Stark
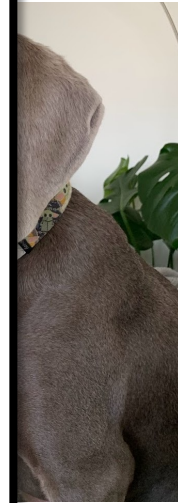
Customized
Diffusion

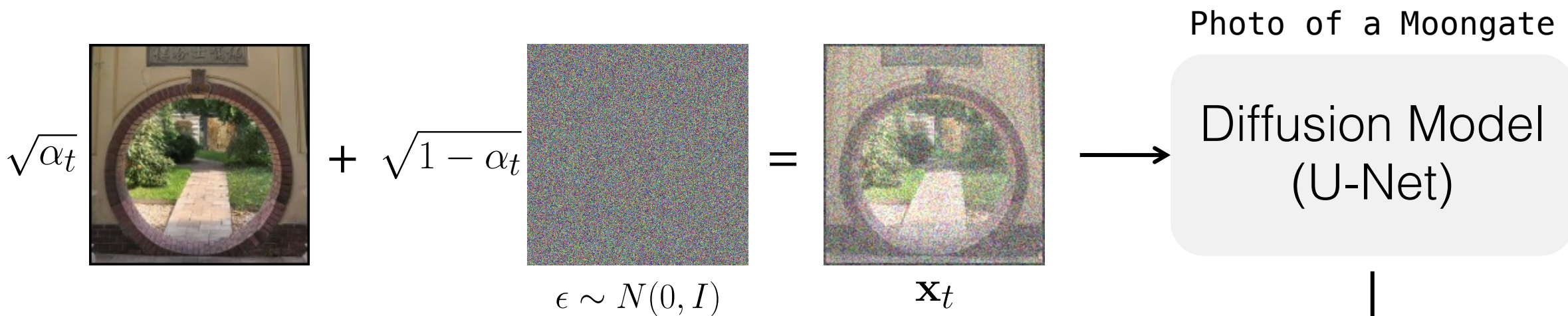V∗ **dog** wearing sunglasses

# Multiple concepts



Customized Diffusion

Actual moongate images

ark

V* **dog** wearing sunglasses in front of **moongate**

# Diffusion Model Quick Recap

# Diffusion model training

$\sqrt{\alpha_t}$ + $\sqrt{1-\alpha_t}$

$\epsilon \sim N(0, I)$  =  $\mathbf{x}_t$

Photo of a Moongate

Diffusion Model (U-Net)

$\epsilon_\theta$

$\longrightarrow$ L2 loss $\longleftarrow$

noisy input   timestep

$$\arg\min_\theta \mathbb{E}_{\epsilon,\mathbf{x},\mathbf{c},t}[||\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}, t)||_2]$$

U-Net   caption

# Which parts shall we customize?

# Textual Inversion: Optimizing Text Embedding



Input samples $\xrightarrow{invert}$ "$S_*$"      "An oil painting of $S_*$"      "App icon of $S_*$"      "Elmo sitting in the same pose as $S_*$"      "Crochet $S_*$"

[Rinon Gal et al., ICLR 2023]

# Textual Inversion: Optimizing Text Embedding



Input samples $\xrightarrow{invert}$ "$S_*$"    "An oil painting of $S_*$"    "App icon of $S_*$"    "Elmo sitting in the same pose as $S_*$"    "Crochet $S_*$"

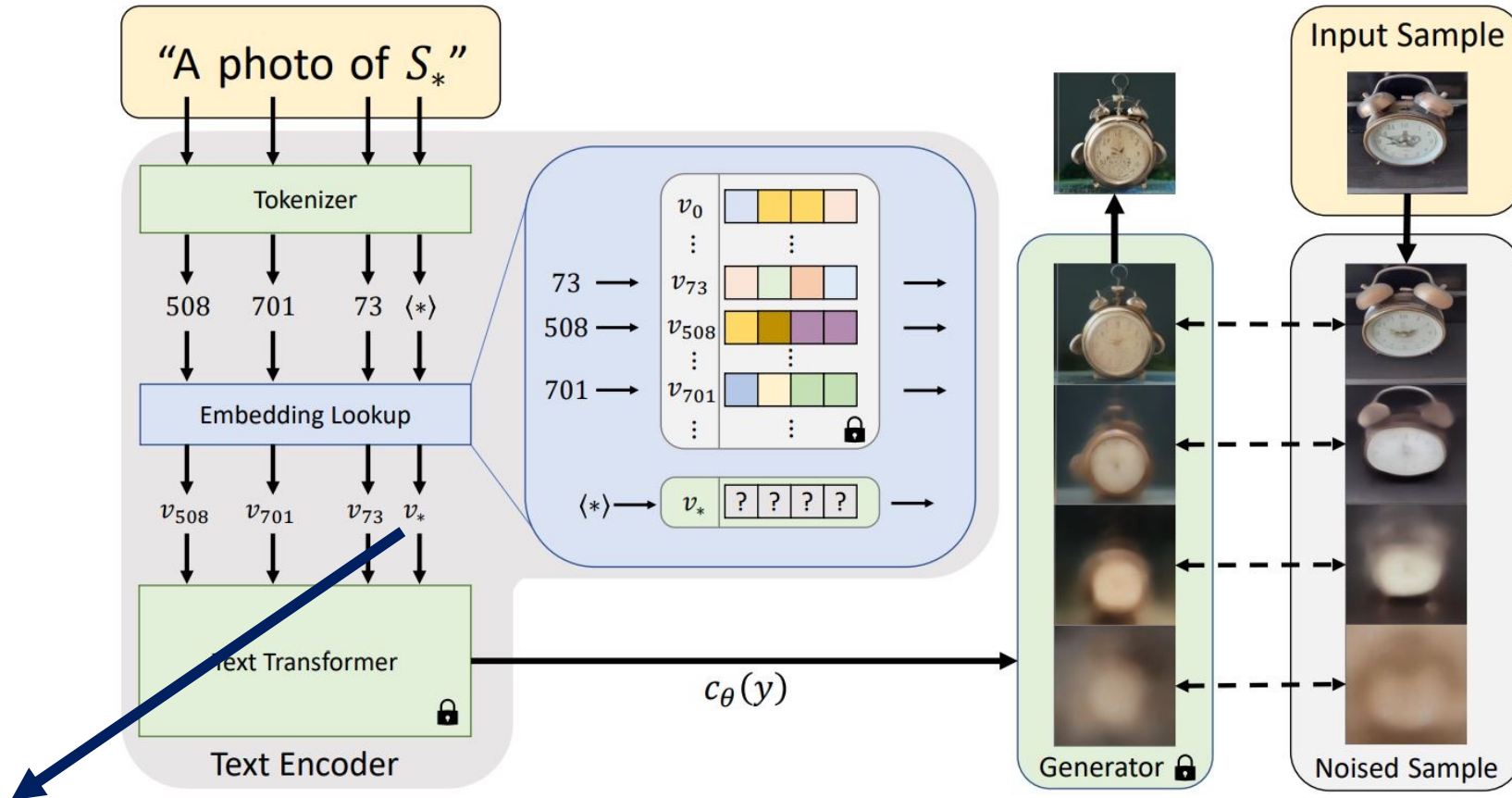Input samples $\xrightarrow{invert}$ "$S_*$"    "Painting of two $S_*$ fishing on a boat"    "A $S_*$ backpack"    "Banksy art of $S_*$"    "A $S_*$ themed lunchbox"

[Rinon Gal et al., ICLR 2023]

# Textual Inversion: Optimizing Text Embedding



$$v^* = \arg\min_{v} \mathbb{E}_{\epsilon,\mathbf{x},\mathbf{c},t}[||\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}, t)||_2]$$

GANs inversion [Zhu et al., 2016] and soft prompting [Lester et al., 2021]     [Rinon Gal et al., ICLR 2023]

# Textual Inversion Results



Input samples → "$S_*$ sports car"   "$S_*$ made of lego"   "$S_*$ onesie"   "da Vinci sketch of $S_*$"

Input samples → "Manga drawing of a steaming $S_*$"   "A $S_*$ watering can"   "$S_*$ Death Star"   "A poster for the movie 'The Teapot' starring $S_*$"

[Rinon Gal et al., ICLR 2023]

# Textual Inversion Results



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 $S_*$"

[Rinon Gal et al., ICLR 2023]

# Works well for artistic styles



Input samples

"The streets of Paris in the style of $S_*$"

"Adorable corgi in the style of $S_*$"

"Painting of a black hole in the style of $S_*$"

"Times square in the style of $S_*$"

[Rinon Gal et al., ICLR 2023]

# Cannot preserve object identity



Target images

S* cat swimming in a pool

[Rinon Gal et al., ICLR 2023]

# How to improve identity preservation?

# DreamBooth: Fine-tuning all the weights



Reconstruction Loss

"A [V] dog"

Text → Image

Shared Weights

Input images (~3-5)

Text → Image

"A dog"

Text → Image

"A dog"

Class-Specific Prior Preservation Loss

## Training Objective

$$\Delta\theta^* = \arg\min_{\Delta\theta} \mathbb{E}_{\epsilon,\mathbf{x},\mathbf{c},t}[||\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}, t)||_2]$$

where $\quad \theta = \theta_0 + \Delta\theta$

## Issues (Overfitting)
- Forget to generate subjects of the same class (e.g., dog)
- Reduce output diversity

## Regularization
- Add synthetic images of the same class.

Inspired by single-image GAN fine-tuning
GANPaint [Bau et al., 2019], PTI [Roich et al., 2021]                    [Nataniel Ruiz et al., CVPR 2023]

# DreamBooth Results



Input images

in the Acropolis

swimming

sleeping

in a doghouse

in a bucket

getting a haircut

# DreamBooth Results



Input images

A [V] backpack in the Grand Canyon

A wet [V] backpack in water

A [V] backpack in Boston

A [V] backpack with the night sky

Input images

A [V] teapot floating in milk

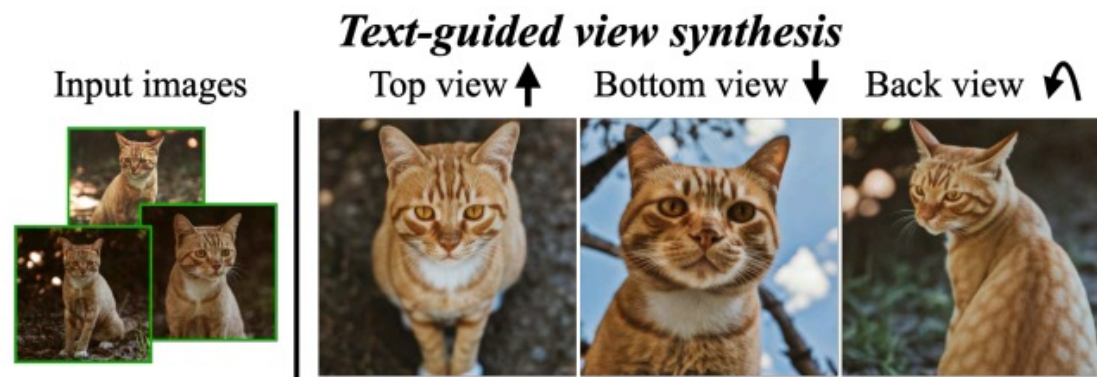A transparent [V] teapot with milk inside

A [V] teapot pouring tea

A [V] teapot floating in the sea

# DreamBooth Applications



[Nataniel Ruiz et al., CVPR 2023]

# DreamBooth vs. Textual Inversion



Input Images

DreamBooth (Imagen)

DreamBooth (Stable Diffusion)

Textual Inversion (Stable Diffusion)

[Nataniel Ruiz et al., CVPR 2023]

# Fine-tuning all model weights

Photo of a moongate

Moongate in snowy ice



**Storage requirement.** 4GB storage for each fine-tuned model.
**Compute requirement.** It requires more VRAM/training time.
**Compositionality.** Hard to combine multiple models.

# Analyze change in weights

$$\Delta_l = \frac{||\theta_l' - \theta_l||}{||\theta_l||}$$  where  $\theta_l'$ : updated weights
$\theta_l$ : pretrained weights

# Analyze change in weights

$$\Delta_l = \frac{||\theta_l' - \theta_l||}{||\theta_l||}$$

where
$\theta_l'$ : updated weights
$\theta_l$ : pretrained weights

# Analyze change in weights

$$\Delta_l = \frac{||\theta'_l - \theta_l||}{||\theta_l||} \quad \text{where} \quad \begin{array}{l} \theta'_l : \text{updated weights} \\ \theta_l : \text{pretrained weights} \end{array}$$

# Analyze change in weights

$$\Delta_l = \frac{||\theta'_l - \theta_l||}{||\theta_l||}$$

where
$\theta'_l$ : updated weights
$\theta_l$ : pretrained weights

# Text-image Cross-Attention



$$\text{Output} = \text{Softmax}\left(\frac{Q.K^T}{\sqrt{d'}}\right)V$$

# Text-image Cross-Attention



Text features only input to $W_k$ and $W_v$

Trainable   Frozen

# Only fine-tune cross-attention layers



$$\Delta W_k^*, \Delta W_v^* = \arg \min_{\Delta W_k, \Delta W_v} \mathbb{E}_{\epsilon, \mathbf{x}, \mathbf{c}, t}[||\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}, t)||_2]$$

Trainable  Frozen

[Nupur Kumari et al., CVPR 2023]

# Generated samples for target concept

Photo of a moongate



Pretrained Model

Fine-tuned Model

# Generated samples for similar concepts

## Photo of a moon



Pretrained Model



Fine-tuned Model

# How to prevent overfitting?



Photo of a {moongate}

Photo of a {moongate}

. . .

Target images

+

sky full of stars and the moon

Blood moon

. . .

Add regularization images

# Generated samples for similar concepts

## Photo of a moon



Pretrained Model

Fine-tuned Model

# Generated samples for similar concepts

## Photo of a moon



Pretrained Model

Fine-tuned Model

# Personalized concepts



Jun-Yan's dog, Stark

How to describe personalized concepts?

V* dog

Where V* is a modifier token in the text embedding space

Proposed by Textual Inversion [Rinon Gal et al.]

# Personalized concepts

Also fine-tune the modifier token **V\*** that describes the personalized concept



Trainable   Frozen

[Nupur Kumari et al., CVPR 2023]

# Single concept results



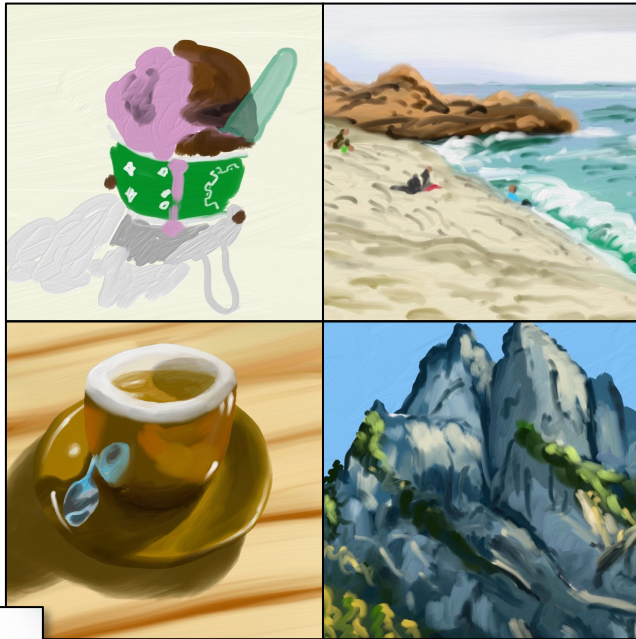V* dog wearing headphones

# Single concept results



A watercolor painting of V* tortoise plushy on a mountain

# Single concept results



V* table and an orange sofa

# Results: specific art style
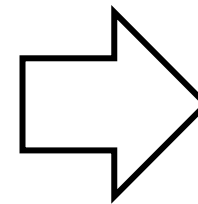


Drawings from Aaron Hertzmann

Painting of dog in the style
of V* art

# Multiple new concepts?

# Joint training

1. Combine the training dataset of multiple concepts



Target images

V∗ dog

Moongate

Regularization images
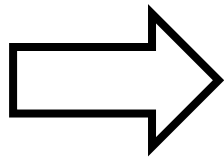
Dog          Cute dog

Wisdom moon    Gated entry

# Joint training

Requires re-training for each choice of composition

100 concepts -> 4950 combinations of **two** concepts.

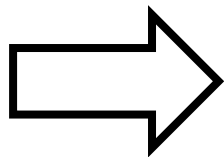100 concepts -> 161, 700 combinations of **three** concepts.

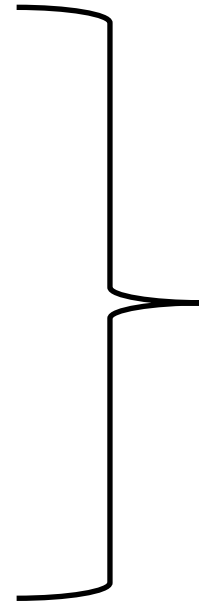# Can we merge weights of individual concepts?
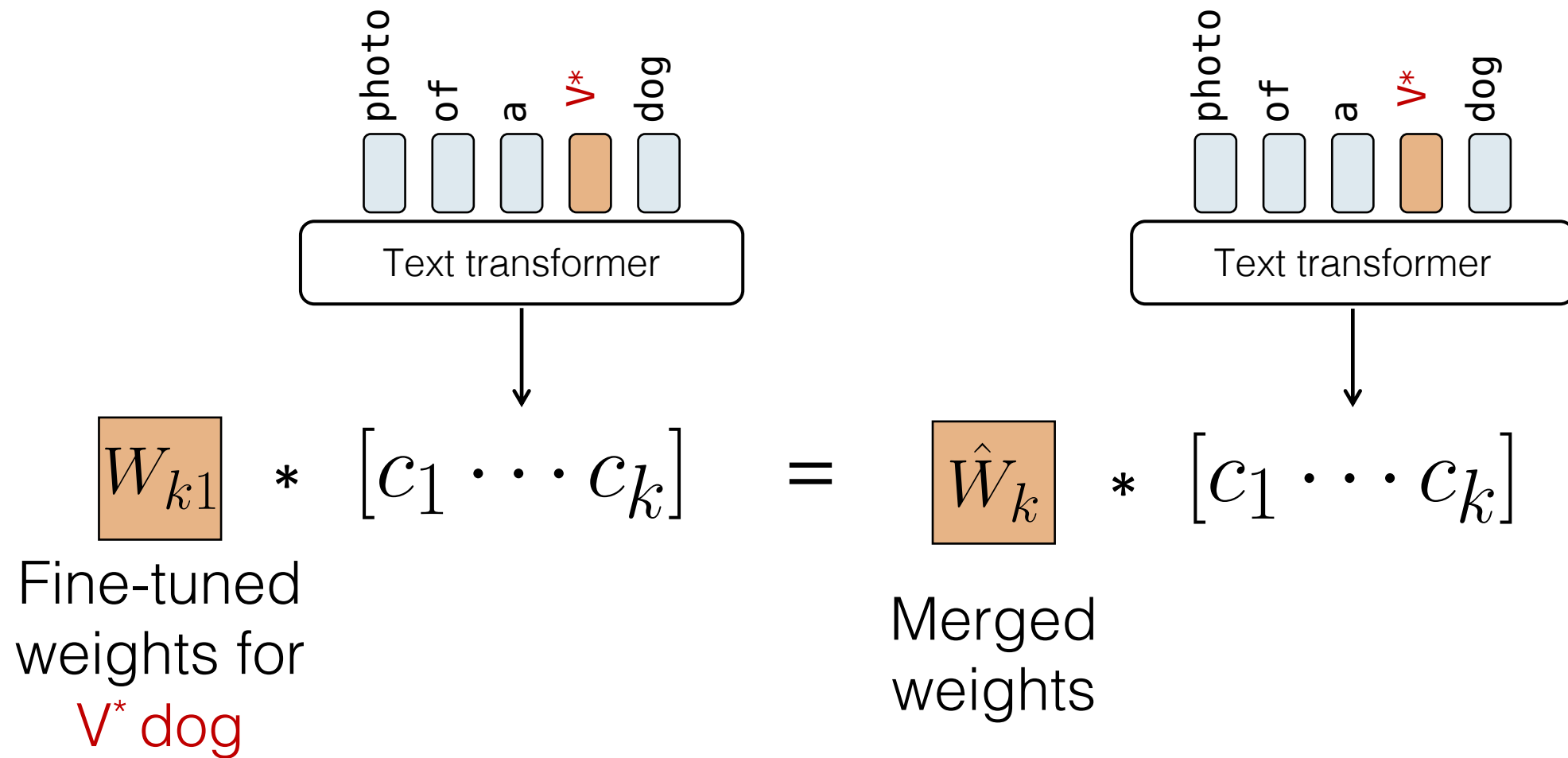


$W_{k1}$ $W_{v1}$

$+$

$W_{k2}$ $W_{v2}$

$\hat{W}_k$ $\hat{W}_v$

**V∗ dog** wearing sunglasses in front of a **moongate**

# Objective function for merging weights

photo of a V* dog

Text transformer

photo of a V* dog

Text transformer

$$W_{k1} * [c_1 \cdots c_k] = \hat{W}_k * [c_1 \cdots c_k]$$

Fine-tuned weights for V* dog

Merged weights

[Nupur Kumari et al., CVPR 2023]

# Objective function for merging weights



$$W_{k2} * [c_{k+1} \cdots c_N] = \hat{W}_k * [c_{k+1} \cdots c_N]$$

Fine-tuned weights for <span style="color:red">moongate</span>

Merged weights

[Nupur Kumari et al., CVPR 2023]

# Constrained least square problem

Stay close to pretrained weights $W_0$ for random text prompts $C_{reg}$.

$$\hat{W} = \underset{W}{\arg\min}||WC_{\text{reg}}^{\top} - W_0 C_{\text{reg}}^{\top}||_F$$

s.t. $\hat{W}[c_1 \cdots c_N] = [W_1 c_1 \cdots W_2 c_N]$

$C$: target prompts, e.g., {photo of a V* dog, photo of moongate}

[Nupur Kumari et al., CVPR 2023]

# Constrained least square problem

Constrained least square problem

$$\hat{W} = \arg \min_{W} ||WC_{\text{reg}}^{\top} - W_0 C_{\text{reg}}^{\top}||_F$$

$$\text{s.t.} \quad \hat{W}[c_1 \cdots c_N] = [W_1 c_1 \cdots W_2 c_N]$$

[Nupur Kumari et al., CVPR 2023]

# Constrained least square problem

Constrained least square problem

$$\hat{W} = \arg\min_{W} ||WC_{\text{reg}}^\top - W_0 C_{\text{reg}}^\top||_F$$

s.t. $\hat{W}[c_1 \cdots c_N] = [W_1 c_1 \cdots W_2 c_N]$

Close-form solution for solving for W and v,

$$\hat{W} = W_0 + \mathbf{v}^\top \mathbf{d}, \text{ where } \mathbf{d} = C(C_{\text{reg}}^\top C_{\text{reg}})^{-1}$$

$$\text{and } \mathbf{v}^\top = (V - W_0 C^\top)(\mathbf{d}C^\top)^{-1}$$

[Nupur Kumari et al., CVPR 2023]

# Two concept results



V₁* dog in front of moongate

# Two concept results



The $V_1$* cat is sitting inside a $V_2$* wooden pot and looking up
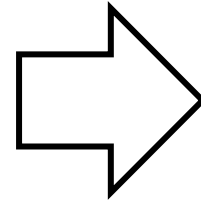
# Two concept results



$V_1$* chair with the $V_2$* cat sitting on it near a beach

# Two concept results



$V_1$* flower in the $V_2$*
wooden pot on a table

# Two concept results



Drawings from Aaron Hertzmann

$V_1$* art style painting of $V_2$* wooden pot

# Qualitative comparison (single-concept)

Target Images



V* teddybear in
Times Square??

# Qualitative comparison (single-concept)

Target Images

Custom Diffusion (Ours)

DreamBooth

Textual Inversion



V* teddybear in Times Square

# Qualitative comparison (multi-concept)

Target Images

Custom Diffusion (Ours)

DreamBooth

Textual Inversion



$V_1$* flower in the $V_2$* wooden pot on a table

# Limitations



Ours

Pretrained model

$V_1$* dog and a $V_2$* cat
playing together

dog and a cat
playing together

# Memory requirement

Each custom diffusion model: 75MB storage

Analyze the difference in pretrained and fine-tuned weights

# Compressing fine-tuned weights



| 75MB | 15MB | 0.1MB | 0.08MB |
|---|---|---|---|

Target image | Custom Diffusion | Top 20% rank | 1 Rank | 0 Rank

# Low-rank Adaptation (Lora)
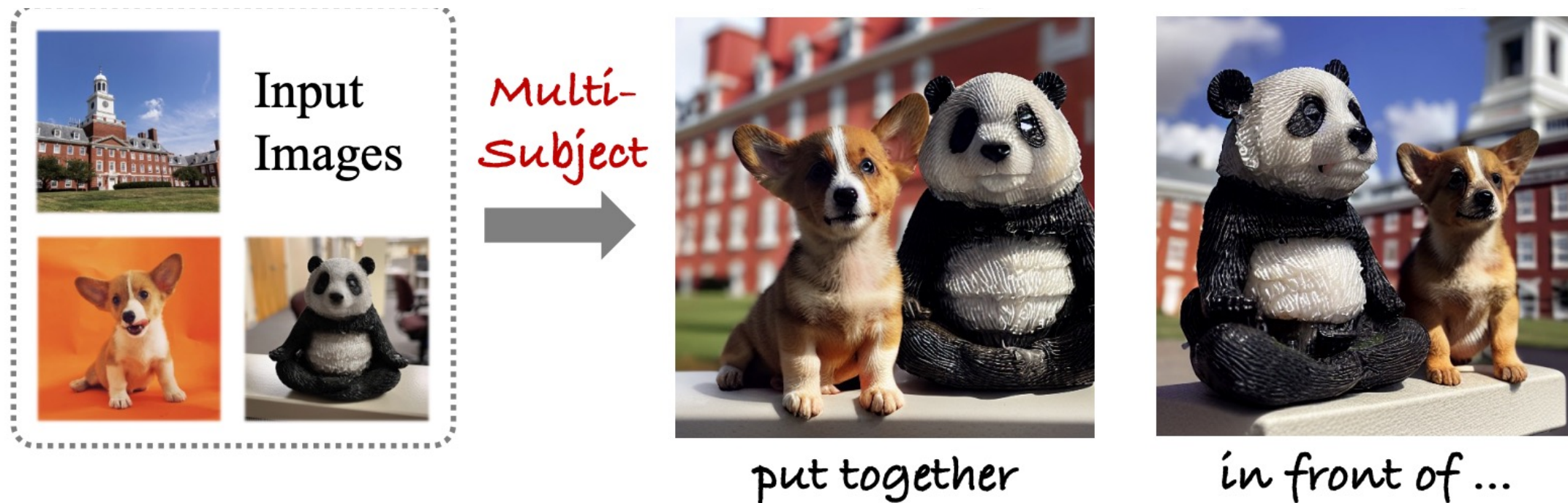


- Lora: Low-rank adaptation of large language models

Original weights

$$W = W_0 + BA$$

Low-rank difference

Lora [Edward J. Hu*, Yelong Shen*, et al., ICLR 2022]
Lora + Dreambooth (by Simo Ryu): https://github.com/cloneofsimo/lora

# Low-rank Adaptation (SVDiff)



- Composing multiple concepts

$$\Sigma_{\boldsymbol{\delta}'} = \mathrm{diag}(\mathrm{ReLU}(\boldsymbol{\sigma} + \boldsymbol{\delta}_1 + \boldsymbol{\delta}_2))$$

SVDiff [Han et all., ICLR 2022]

# Low-rank Adaptation (Rank-1)

- Rank-1 Model Editing
- Used in GAN fine-tuning [Bau et al., 2020] and
  LLM factual editing [Meng et al., 2022]

$$\hat{W} = W + \Lambda (C^{-1} \boldsymbol{i}_*)^T.$$

$$\Lambda = (\boldsymbol{o}_* - W\boldsymbol{i}_*)/[(\boldsymbol{i}_*^T (C^{-1})^T \boldsymbol{i}_*)]$$

Please see their paper for more details including key lock

Perfusion [Tewel et all., SIGGRAPH 2023]

# Optimization is too Slow!

# Encoder-based Methods

# Image Prompt Adapter (IP-Adapter)



[He Yu et al., CVPR 2024]

# Image Prompt Adapter (IP-Adapter)



[He Yu et al., CVPR 2024]

# Image Prompt Adapter (IP-Adapter)



[He Yu et al., CVPR 2024]

# Optimization + encoder (5-15 steps)

# Datasets

# DreamBooth Dataset: 30 subjects

# CustomConcept101: 101 concepts



[Nupur Kumari et al., CVPR 2023]