

JUNE 18-22, 2023

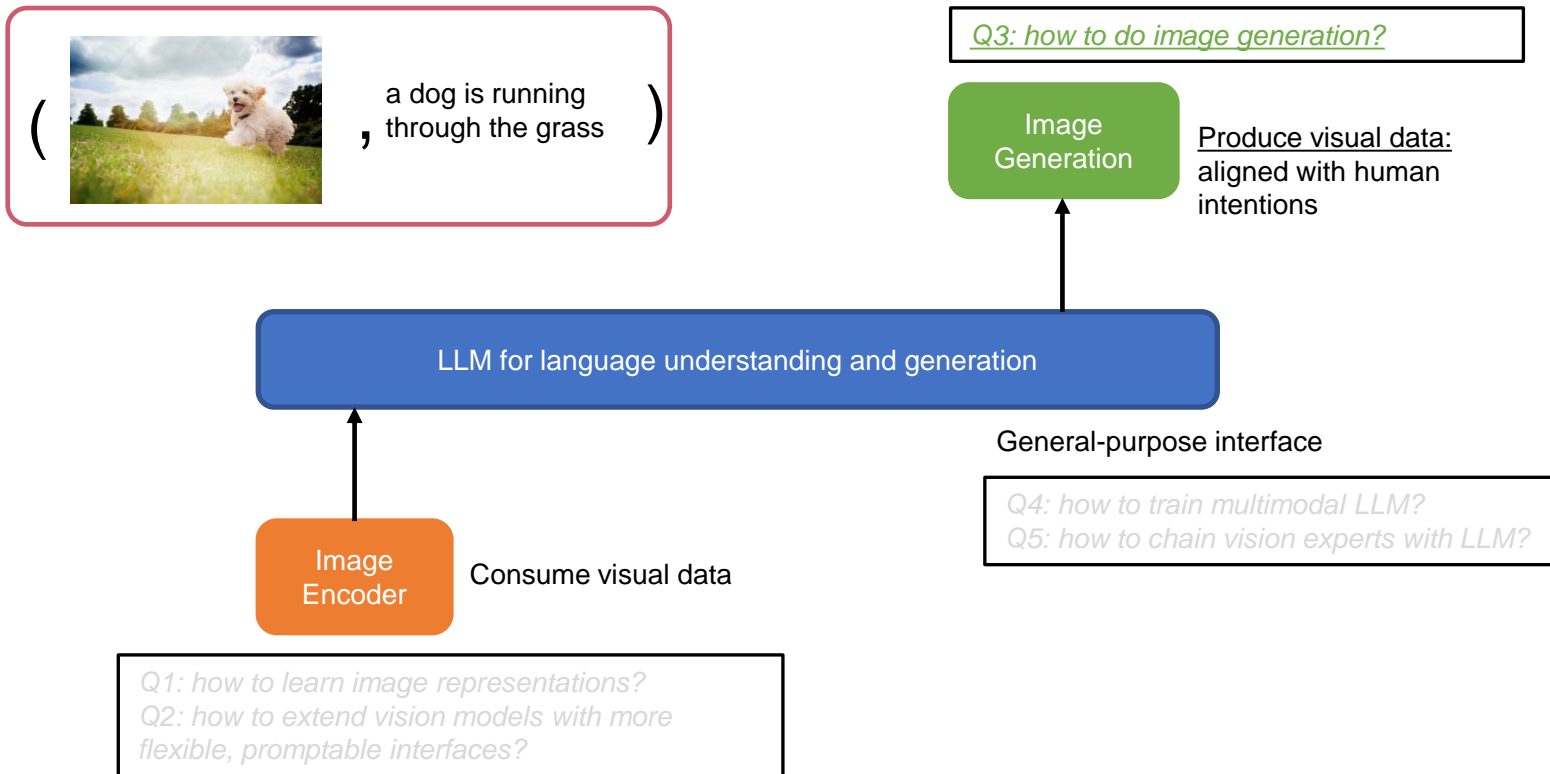


Alignments in Text-to-Image Generation

Zhengyuan Yang



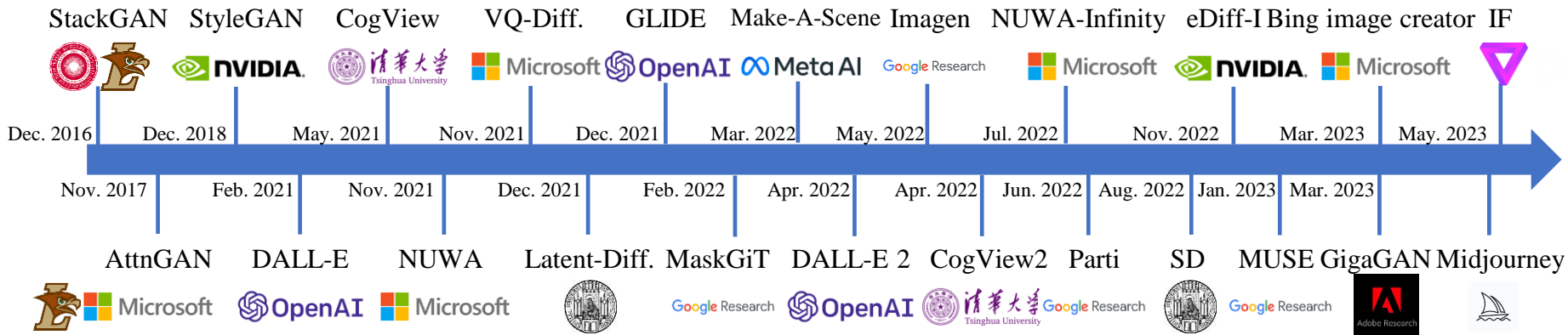
Alignments in Text-to-Image Generation



Text-to-Image Generation

- Text-to-image (T2I)
- Aligning with human intentions

Text prompt: a yellow fire hydrant with a cartoon face drawn on it.



Alignments in Text-to-Image Generation

Controllable generation

Image-level: a yellow fire hydrant with a cartoon face drawn on it.

a truck is parked next to a trash can.

a red truck is parked in a parking lot.

a yellow fire hydrant with a face on it and black eyes.



Editing

"Swap sunflowers with roses"

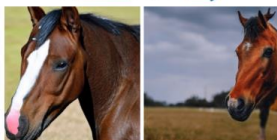


"Add fireworks to the sky"



Better following prompts

"A horse and a dog"



Stable Diffusion

"A painting of an elephant with glasses"



+Attend-and-Excite



Concept customization



Input images



in the Acropolis



swimming



sleeping



in a doghouse



in a bucket

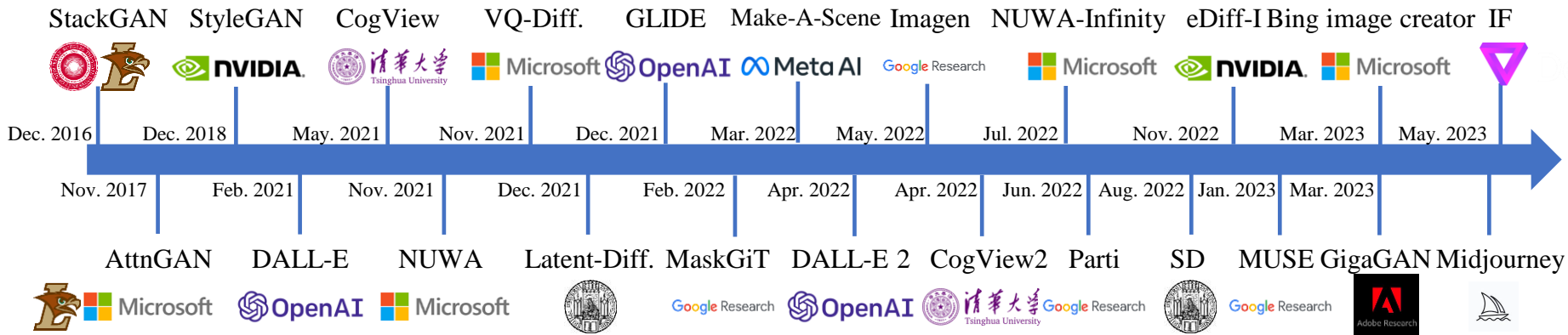
Agenda

- Text-to-image (T2I) basics
- Aligning human intentions in T2I generation
 - Controllable generation
 - Editing
 - Better following prompts
 - Concept customization
- Summary and discussion

Text-to-Image Basics

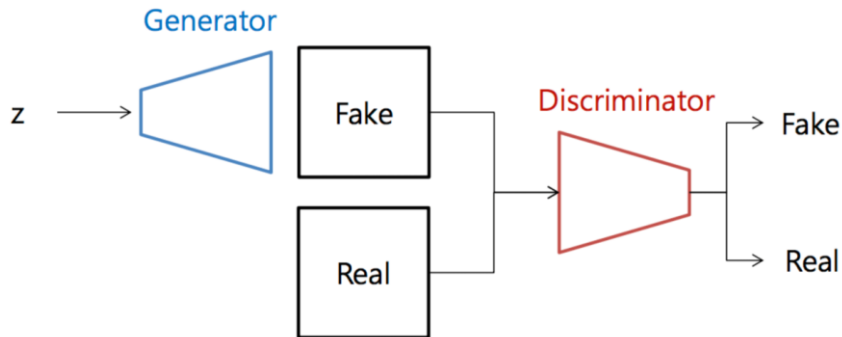
- GAN
- Auto-regressive
- Non-AR Transformer
- Diffusion

Text prompt: a yellow fire hydrant with a cartoon face drawn on it.

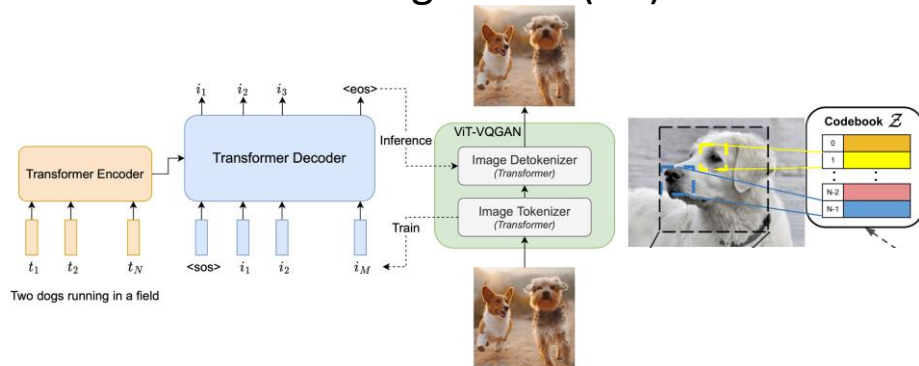


Text-to-Image Basics

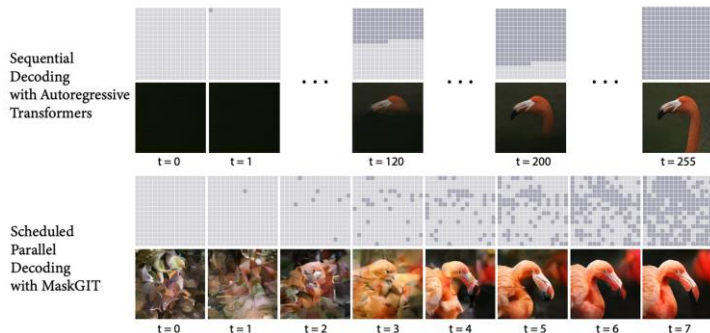
Generative Adversarial Networks (GAN)



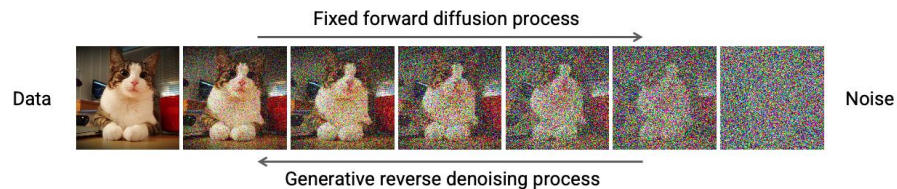
Auto-regressive (AR)



Non-AR Transformer

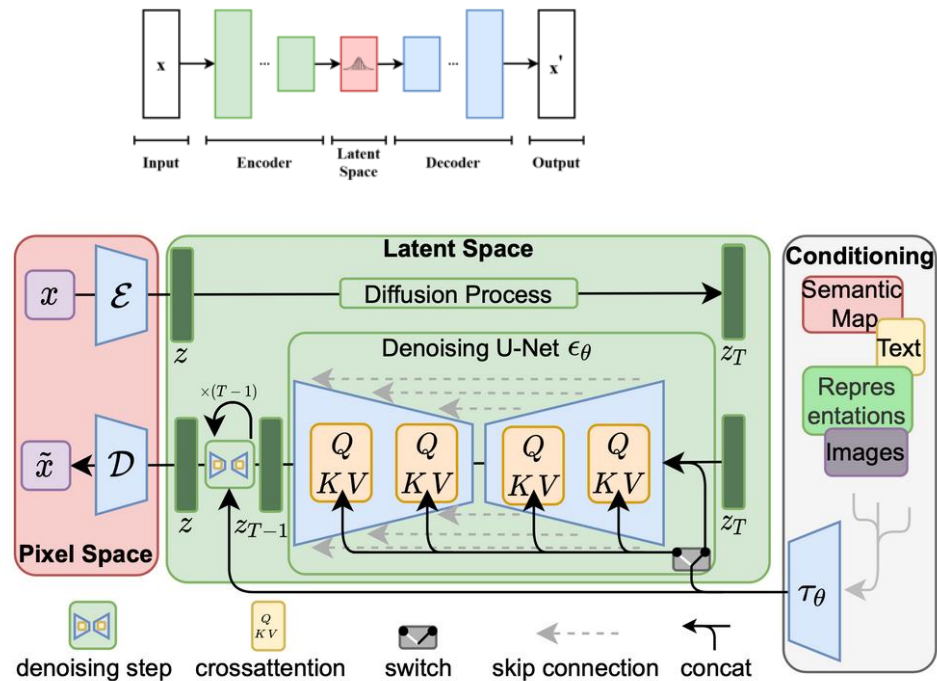


Diffusion



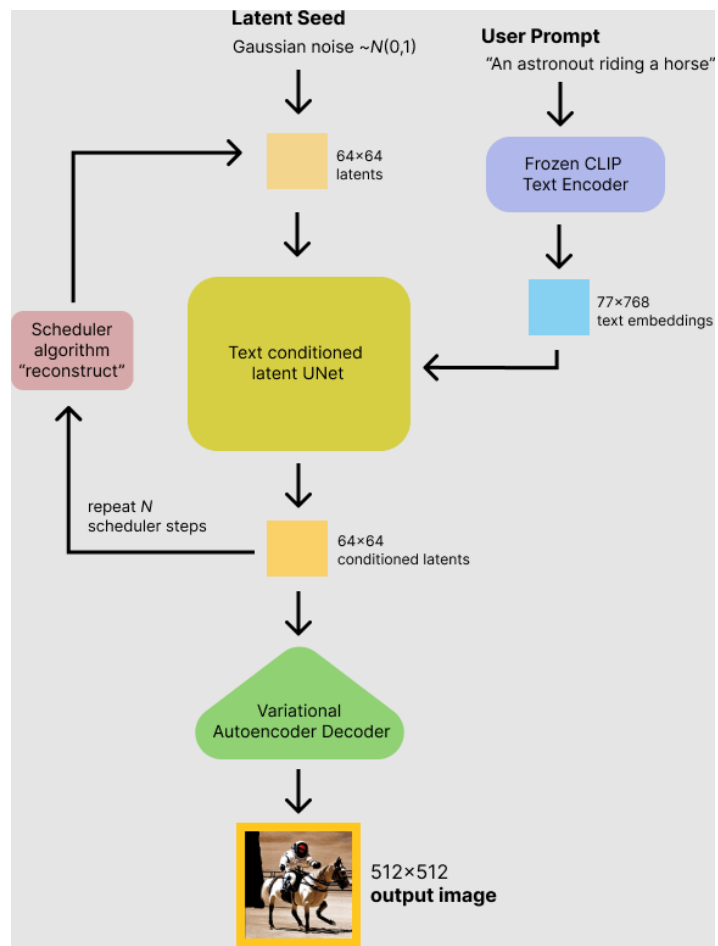
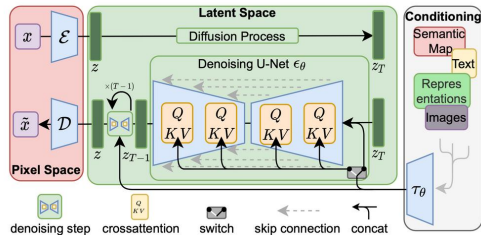
Stable Diffusion (SD) Basics

- SD overview
 - Variational autoencoder (VAE)
 - Condition encoder
 - Conditional denoising U-Net



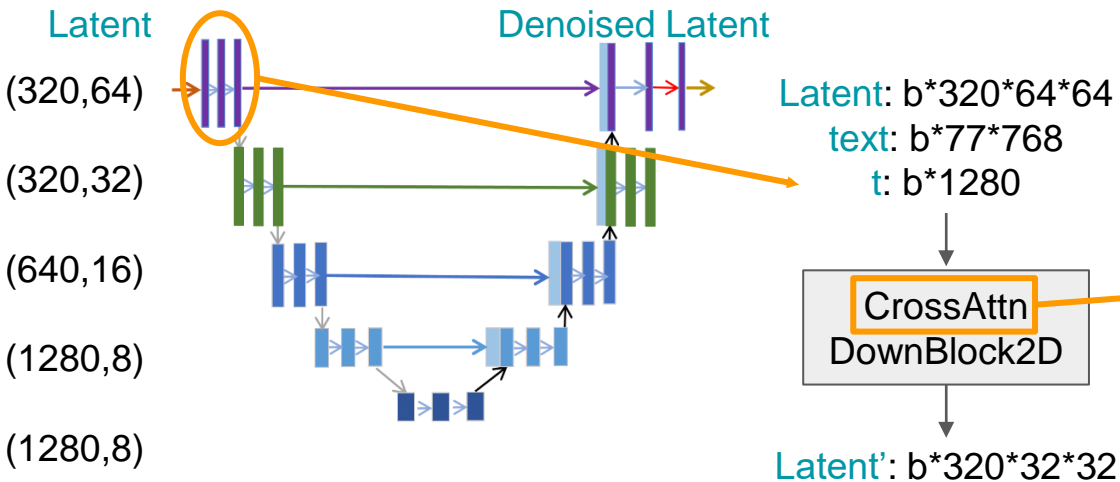
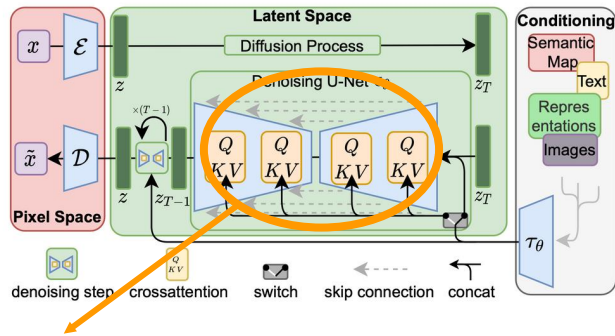
Stable Diffusion (SD) Basics

- Inference flow
 - Variational autoencoder (VAE)
 - Condition encoder
 - Conditional denoising U-Net



Stable Diffusion (SD) Basics

- Zooming into conditional U-Net: How text condition operates on image?
 - Image-text cross attention



$$Q: \text{latent} + \text{duplicate}(\text{linear}(t)) \Rightarrow b * 4096 * 320$$

$$K, V: \text{text} \Rightarrow b * 77 * 768$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) \cdot V, \text{ with}$$

$$Q = W_Q^{(i)} \cdot \varphi_i(z_t), K = W_K^{(i)} \cdot \tau_\theta(y), V = W_V^{(i)} \cdot \tau_\theta(y).$$

Image-text attention map of HW*77

Agenda

- Text-to-image (T2I) basics
- Aligning human intentions in T2I generation
 - Controllable generation
 - Editing
 - Better following prompts
 - Concept customization
- Summary and discussion

Controllable Generation

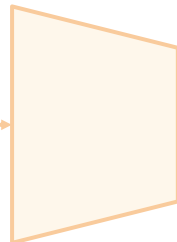
- Text+layout/box: localized description control
- Text+dense control (e.g., mask, edge, scribble, etc.)
- Inference-time guidance

Image-level: a yellow fire hydrant with a cartoon face drawn on it.

a truck is parked next to a trash can.

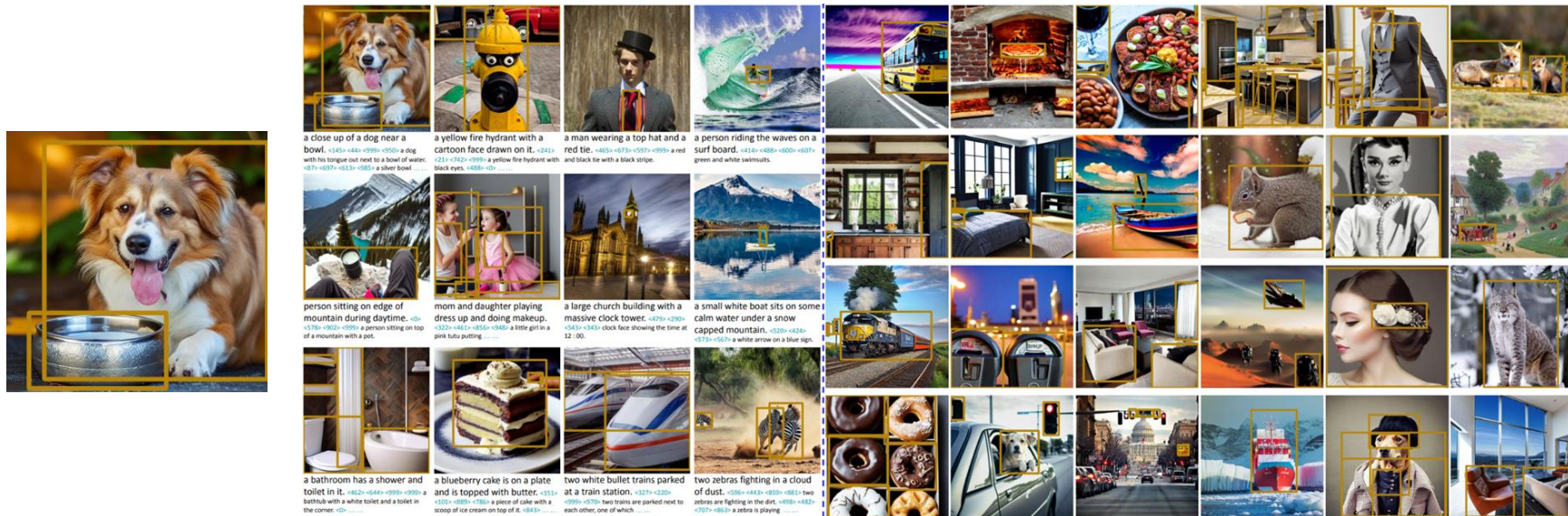
a red truck is parked in a parking lot.

a yellow fire hydrant with a face on it and black eyes.



- [1] [GLIGEN: Open-Set Grounded Text-to-Image Generation](#)
- [2] [ReCo: Region-Controlled Text-to-Image Generation](#)
- [3] [Diagnostic Benchmark and Iterative Inpainting for Layout-Guided Image Generation](#)
- [4] [Adding Conditional Control to Text-to-Image Diffusion Models](#)
- [5] [Composer: Creative and Controllable Image Synthesis with Composable Conditions](#)
- [6] [SpaText: Spatio-Textual Representation for Controllable Image Generation](#)
- [7] [T2I-Adapter: Learning Adapters to Dig out More Controllable Ability for Text-to-Image Diffusion Models](#)
- [8] [SceneComposer: Any-Level Semantic Image Synthesis](#)
- [9] [Uni-ControlNet: All-in-One Control to Text-to-Image Diffusion Models](#)
- [10] [UniControl: A Unified Diffusion Model for Controllable Visual Generation In the Wild](#)
- [11] [Universal Guidance for Diffusion Models](#)
- [12] [Training-Free Layout Control with Cross-Attention Guidance](#)

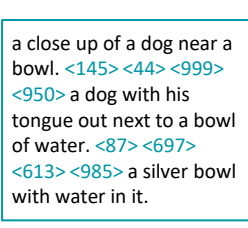
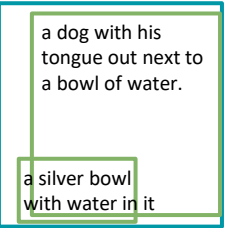
ReCo: Region-Controlled T2I Generation



Text: global image text description



Text: grounded global and regional descriptions (Grounded Region-Controlled texts)



a close up of a dog near a bowl.

a close up of a dog near a bowl.

ReCo: Region-Controlled T2I Generation

- **Input sequence expansion:** box tokens
- **Grounded:** box tokens operate on the text to follow
- Finetune T2I to understand box tokens

a person standing at the plate in mid swing of a bat
 <687> <204> <999> <833> baseball player ... jersey.
 <21> <447> <433> <840> a catcher in gray ... ball.
 <0> <323> <123> <827> a baseball player ... jersey.

Region-Controlled Input Sequence



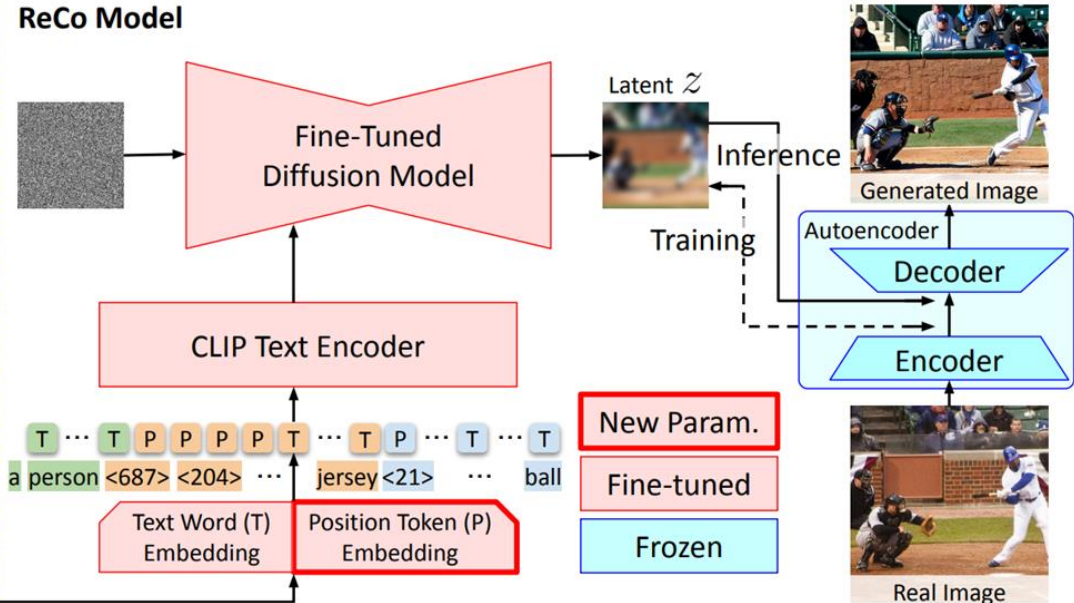
- #1 Region Description: baseball player is swinging a bat and wearing a blue and white jersey.
- #2 Region Description: a catcher in a gray and black uniform is crouching and ready to catch the ball.
- #3 Region Description: a baseball player with the number 19 on his jersey.

Image Description: a person standing at the plate in mid swing of a bat.

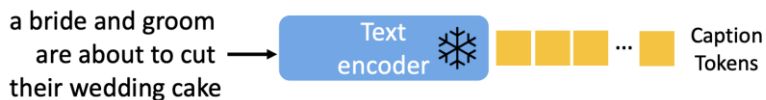
Input Sequence = Image Description + [Region-Controlled Text] * #Regions:

a person standing at the plate in mid swing of a bat
 <687> <204> <999> <833> baseball player ... jersey.
 <21> <447> <433> <840> a catcher in gray ... ball.
 <0> <323> <123> <827> a baseball player ... jersey.

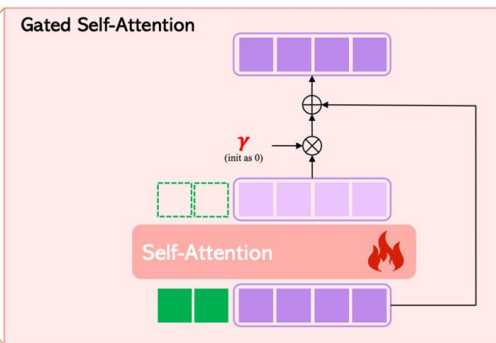
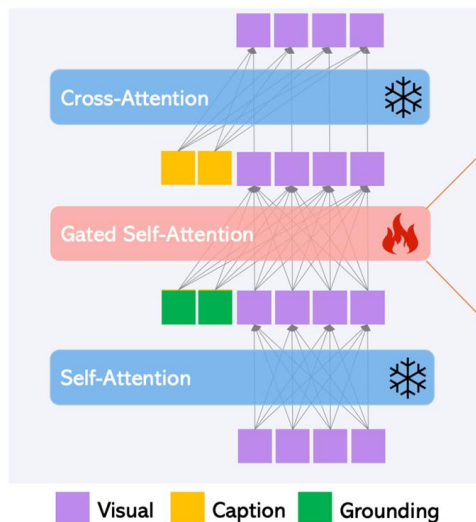
ReCo Model



GLIGEN: Open-Set Grounded T2I Generation



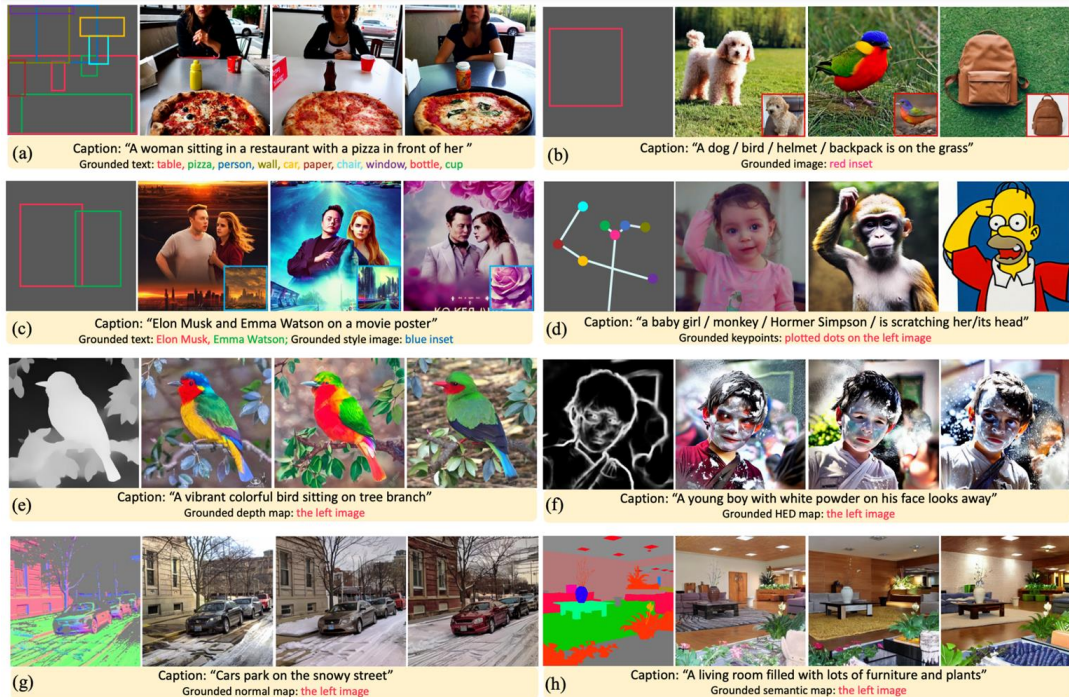
- Grounding tokens: grounded text entity + spatial location
- Gated self-attention layer with original layers frozen



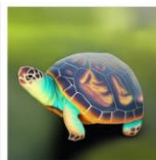
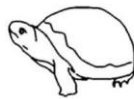
$$\mathbf{v} = \mathbf{v} + \beta \cdot \tanh(\gamma) \cdot \text{TS}(\text{SelfAttn}([\mathbf{v}, \mathbf{h}^e]))$$

GLIGEN: Open-Set Grounded T2I Generation

- Bounding box grounding
- Keypoint grounding
- Spatially-aligned dense conditions



Text+Dense Control

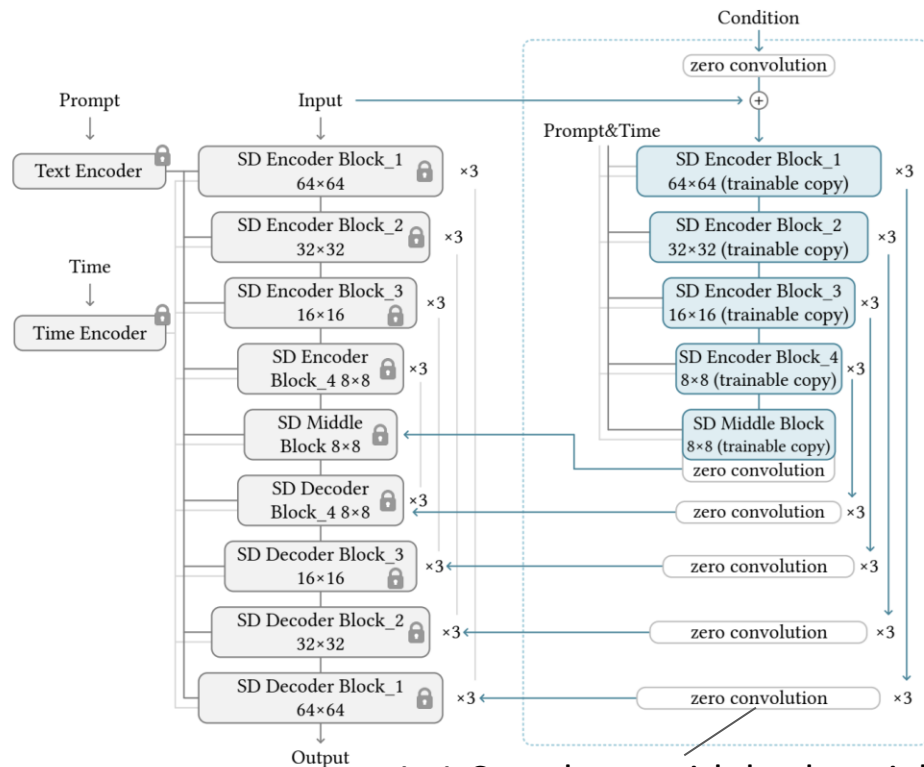


"a turtle in river"



"a masterpiece of cartoon-style turtle illustration"

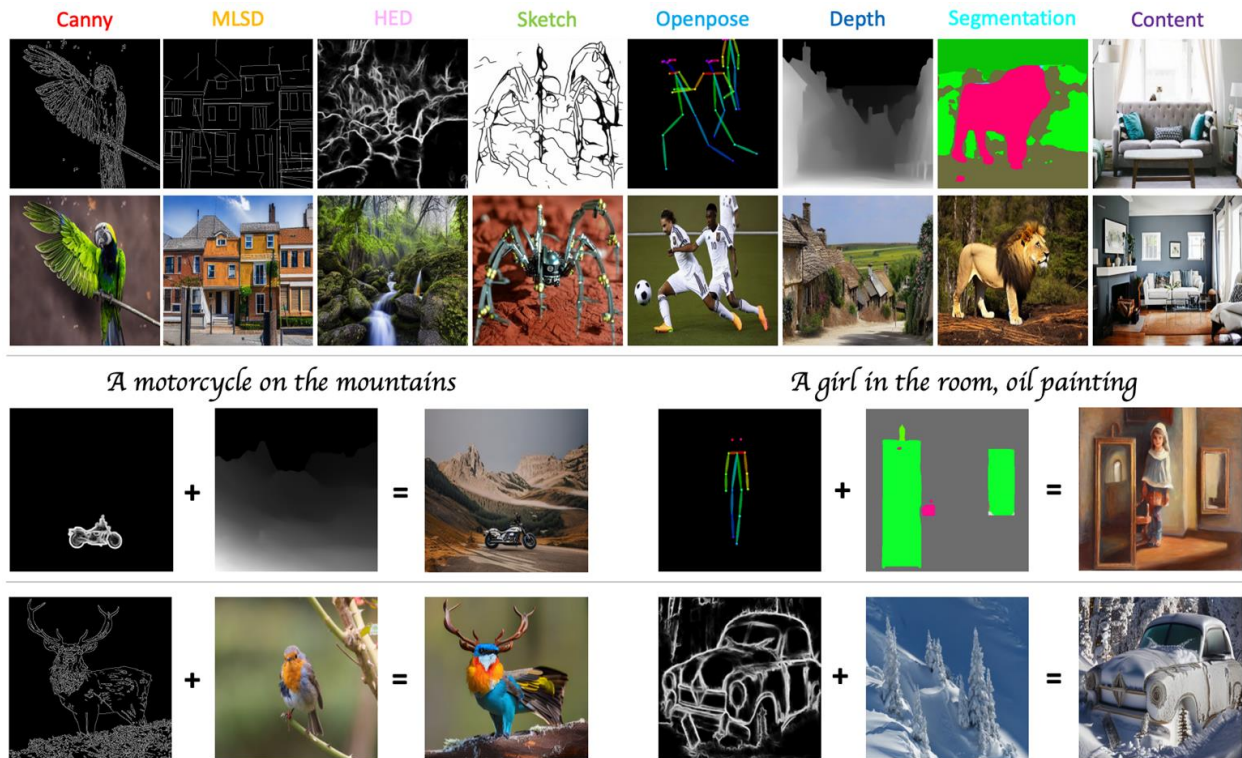
- Dense conditions:
 - Canny Edge
 - Hough Line
 - HED Boundary
 - User Sketching
 - Human Pose
 - Semantic Segmentation
 - Depth
 - Normal Maps
 - Cartoon Line Drawing



1x1 Conv layer with both weight and bias initialized with zeros 17

Uni-ControlNet, UniControl

- Unified models for different conditions
- Condition composition



Inference-time guidance

- Universal Guidance for Diffusion Models: extending classifier guidance [1] to accept any general guidance function

$$\hat{\epsilon}_\theta(z_t, t) = \epsilon_\theta(z_t, t) + s(t) \cdot \nabla_{z_t} \ell(c, f(\hat{z}_0))$$

E.g., detection:

Anchor classification, bounding box regression, and region label classification loss

Box and class labels

Faster-RCNN

Predicted “noisy” clean image



[1] Diffusion Models Beat GANs on Image Synthesis

Editing

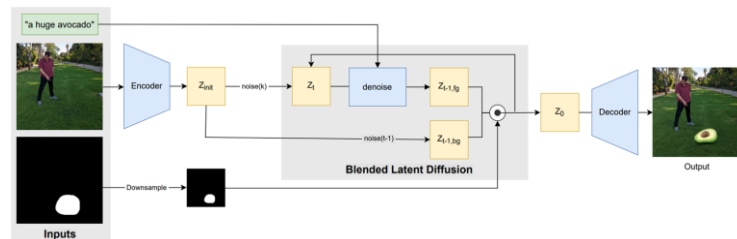
- Latents spatial blend
- Image-text attention edit
- Edit instruction
- External models



- [1] [Blended Diffusion for Text-driven Editing of Natural Images](#)
- [2] [Blended Latent Diffusion](#)
- [3] [DiffEdit: Diffusion-based semantic image editing with mask guidance](#)
- [4] [eDiff-I: Text-to-Image Diffusion Models with an Ensemble of Expert Denoisers](#)
- [5] [Region-Aware Diffusion for Zero-shot Text-driven Image Editing](#)
- [6] [Imagen Editor and EditBench: Advancing and Evaluating Text-Guided Image Inpainting](#)
- [7] [iEdit: Localised Text-guided Image Editing with Weak Supervision](#)
- [8] [EDICT: Exact Diffusion Inversion via Coupled Transformations](#)
- [9] [Prompt-to-Prompt Image Editing with Cross Attention Control](#)
- [10] [Imagic: Text-Based Real Image Editing with Diffusion Models](#)
- [11] [SINE: SINGle Image Editing With Text-to-Image Diffusion Models](#)
- [12] [InstructPix2Pix Learning to Follow Image Editing Instructions](#)
- [13] [MasaCtrl: Tuning-Free Mutual Self-Attention Control for Consistent Image Synthesis and Editing](#)
- [14] [Diffusion Self-Guidance for Controllable Image Generation](#)
- [15] [Instruct-X-Decoder](#)
- [16] [Grounded-SAM Inpainting](#)
- [17] [Inpaint Anything](#)
- [18] [Visual ChatGPT: Talking, Drawing and Editing with Visual Foundation Models](#)

Latents Spatial Blend

- Spatial editing with mask
- Image, text prompt, user input or segmented mask



$$z_t \leftarrow \underbrace{z_{fg}}_{\text{from text}} \odot m_{latent} + \underbrace{z_{bg}}_{\text{original bg image}} \odot (1 - m_{latent})$$

from text original bg image



Input image Input mask "gravestone" "toy truck" "snake"



Input image Input mask "a man with a red suit" "a man with a yellow sweater" "a muscular man with a blue shirt"



Input image Input mask a horror book named CVPR a children's book titled ECCV a romantic novel titled SIGGRAPH



Input image Input mask "beach" "big mountain" "The Great Pyramid of Giza"



Input image Input mask Prediction 1 Prediction 2 Prediction 3



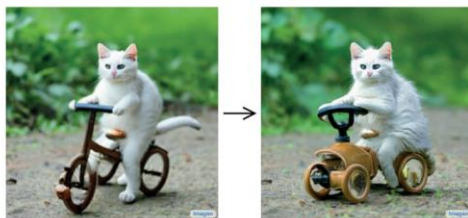
Original image Scribbled image Mask Prediction 1 Prediction 2

Image-text Attention Edit

- Edit generated images
- Manipulate image-text cross-attention map
- Word swap, adding new phrase, attention re-weighting

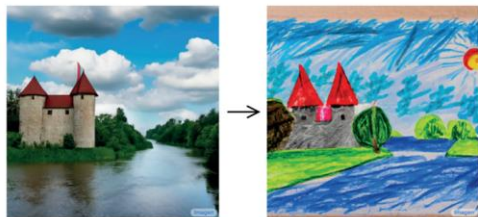


"The boulevards are crowded today."



"Photo of a cat riding on a bicycle."

Word swap



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



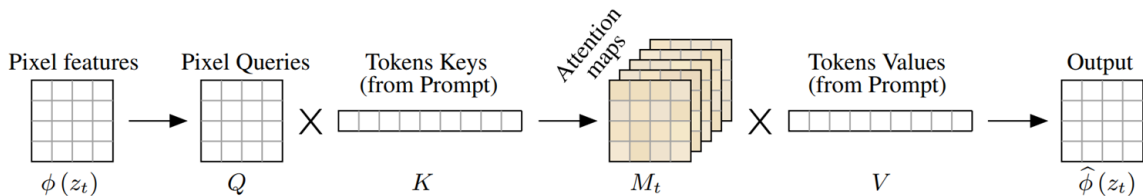
"a cake with decorations."

Adding new phrase

jelly beans

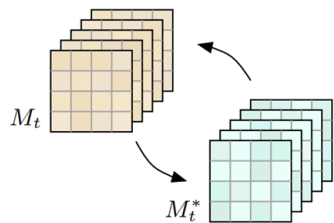
Image-text Attention Edit

- Maintaining two sets of cross-attention maps for edit:
Original prompt: M_t Edited prompt: M_t^*

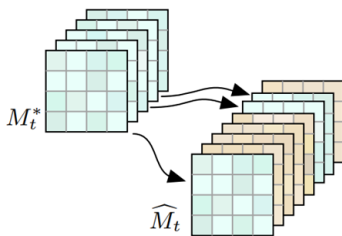


Text to Image Cross Attention

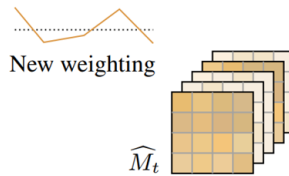
Cross Attention Control



Word Swap



Adding a New Phrase



Attention Re-weighting

$$Edit(M_t, M_t^*, t) := \begin{cases} M_t^* & \text{if } t < \tau \\ M_t & \text{otherwise.} \end{cases}$$

$$(Edit(M_t, M_t^*, t))_{i,j} := \begin{cases} (M_t^*)_{i,j} & \text{if } A(j) = None \\ (M_t)_{i,A(j)} & \text{otherwise.} \end{cases}$$

$$(Edit(M_t, M_t^*, t))_{i,j} := \begin{cases} c \cdot (M_t)_{i,j} & \text{if } j = j^* \\ (M_t)_{i,j} & \text{otherwise.} \end{cases}$$

Goal



"lemon cake."



"monster cake."



M_t

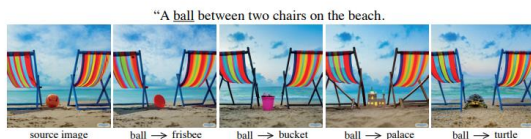
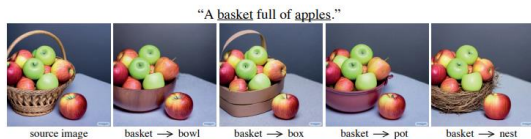
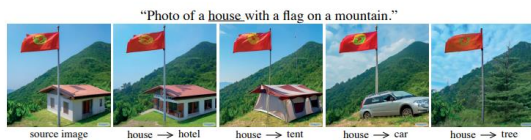
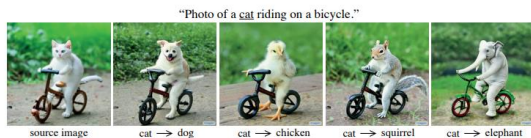
"lemon cake."



M_t^*

"monster cake."

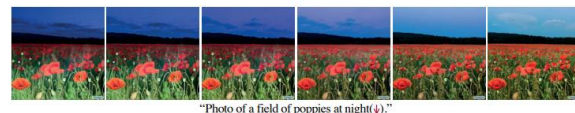
Image-text Attention Edit



Word Swap



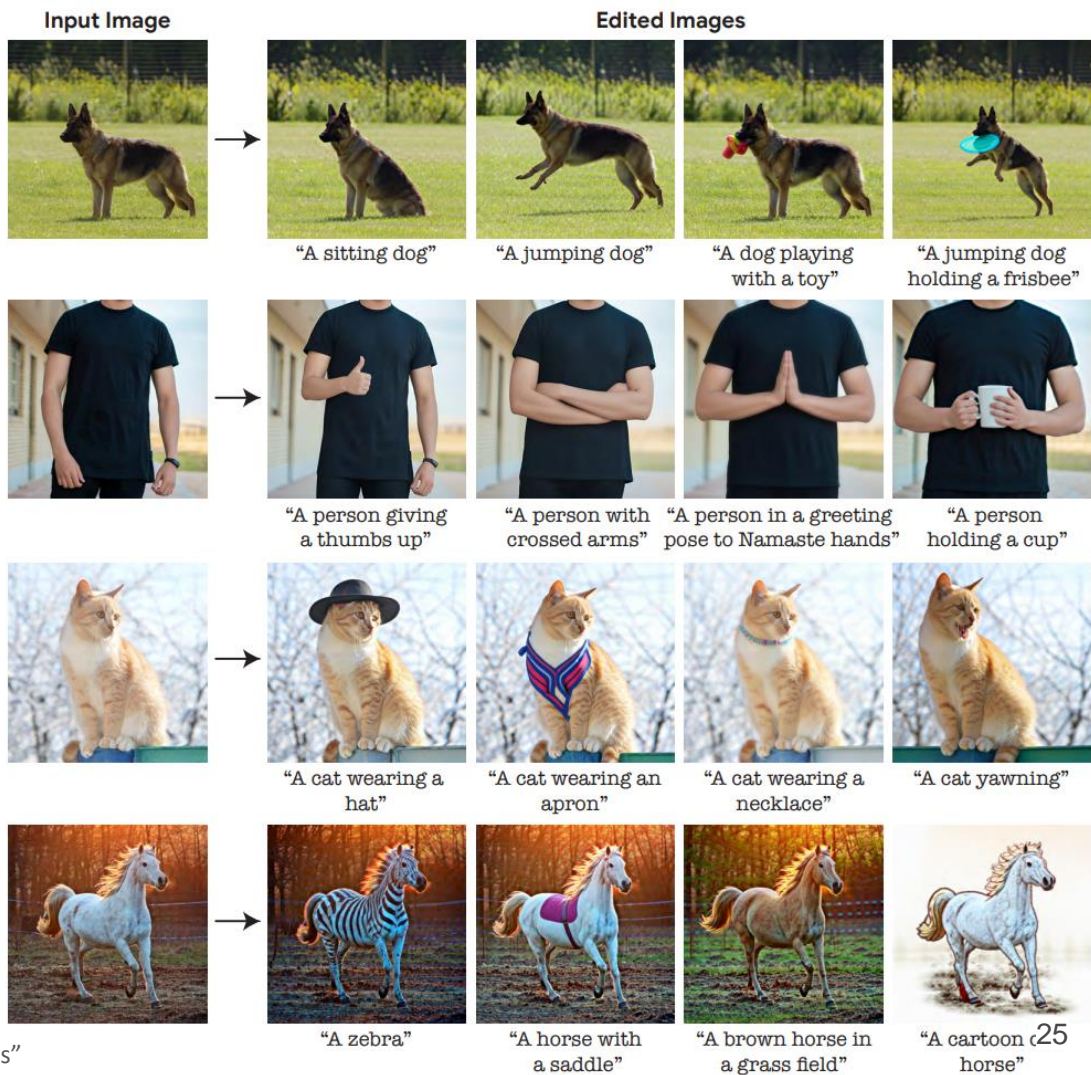
Adding a New Phrase



Attention Re-weighting

Imagic

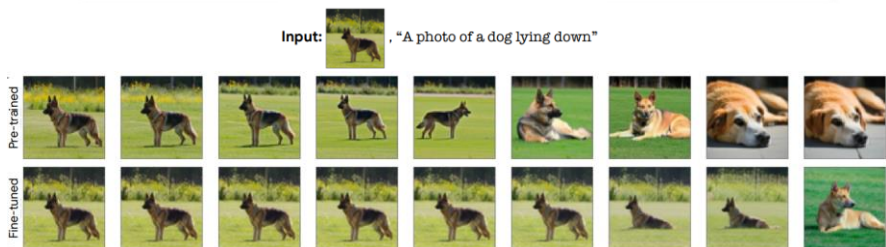
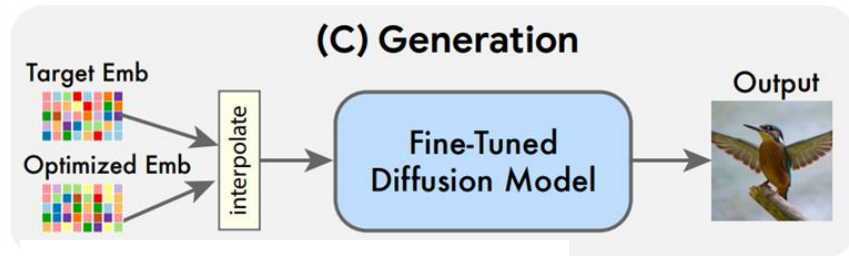
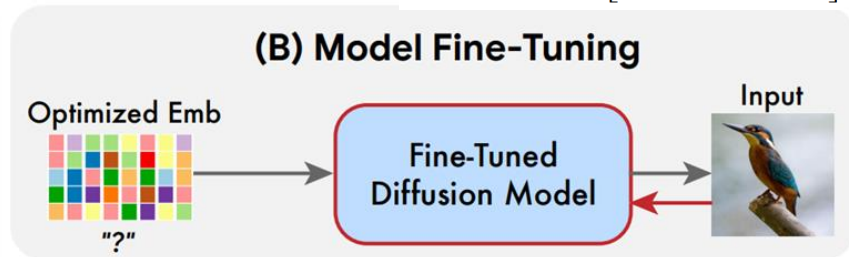
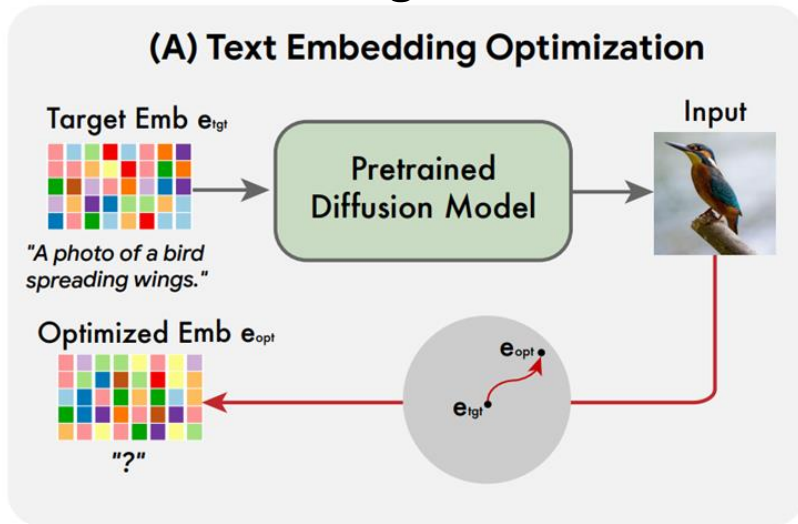
- Generated => natural image edits
- E.g., different dogs



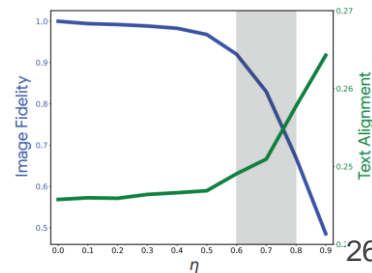
Imagic

- Obtain original text

$$\mathcal{L}(\mathbf{x}, \mathbf{e}, \theta) = \mathbb{E}_{t, \epsilon} \left[\|\epsilon - f_{\theta}(\mathbf{x}_t, t, \mathbf{e})\|_2^2 \right]$$



$$\bar{\mathbf{e}} = \eta \cdot \mathbf{e}_{tgt} + (1 - \eta) \cdot \mathbf{e}_{opt}$$



InstructPix2Pix

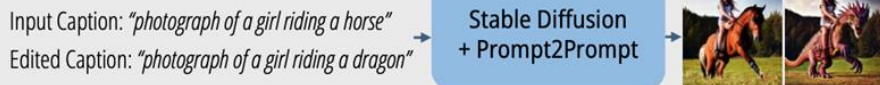
- Obtain original text => Instruction-style text
“a bird standing”, “a bird spreading wings” => “have wings spread”

Training Data Generation

(a) Generate text edits:



(b) Generate paired images:



(c) Generated training examples:



Instruction-following Diffusion Model

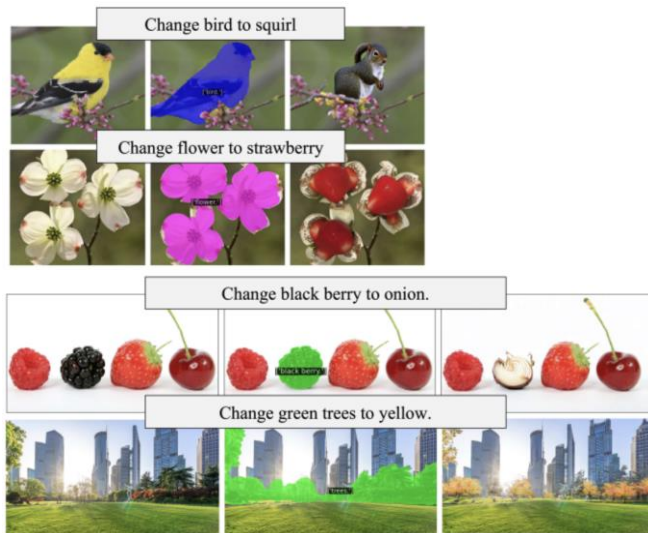
(d) Inference on real images:

“*turn her into a snake lady*”

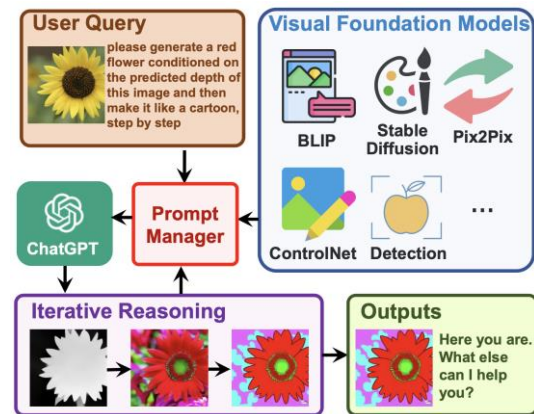


Editing Systems with External Models

- Segmentation

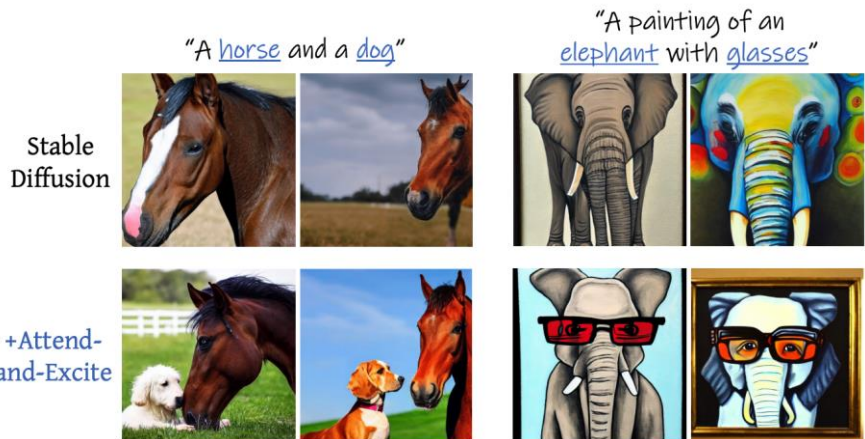


- LLM



Better Following Prompts

- Test-time latents
- Test-time attention
- Alignment tuning



- [1] [Training-Free Structured Diffusion Guidance for Compositional Text-to-Image Synthesis](#)
- [2] [Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models](#)
- [3] [Grounded Text-to-Image Synthesis with Attention Refocusing](#)
- [4] [Compositional Visual Generation with Composable Diffusion Models](#)
- [5] [Aligning Text-to-Image Models using Human Feedback](#)
- [6] [ImageReward: Learning and Evaluating Human Preferences for Text-to-Image Generation](#)
- [7] [Training Diffusion Models with Reinforcement Learning](#)
- [8] [DPOK: Reinforcement Learning for Fine-tuning Text-to-Image Diffusion Models](#)
- [9] [Better Aligning Text-to-Image Models with Human Preference](#)

StructureDiffusion

**Stable
Diffusion**



Ours

A red car and a white sheep.

Attribute leakage



*A brown bench sits in front of
an old white building*

Interchanged attributes

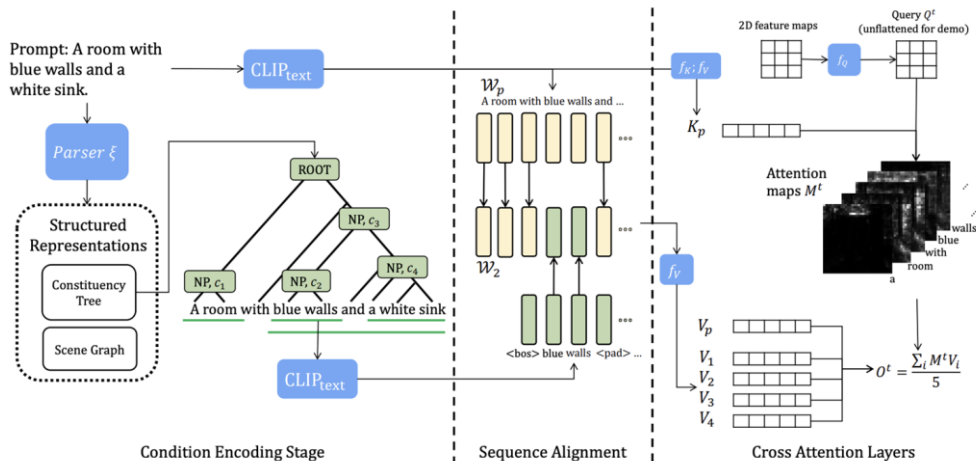


*A blue backpack and a brown
elephant*

Missing objects

StructureDiffusion

- Manipulating values in cross-attention based on linguistic parsing tree to enforce language structure
- Look at all noun phrases



Q: latent+duplicate(linear(t))
=> b*4096*320

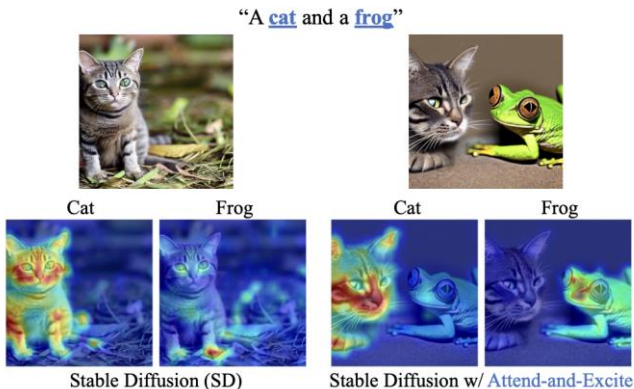
K, V: text => b*77*768

Attention(Q, K, V) = softmax($\frac{QK^T}{\sqrt{d}}$) · V, with

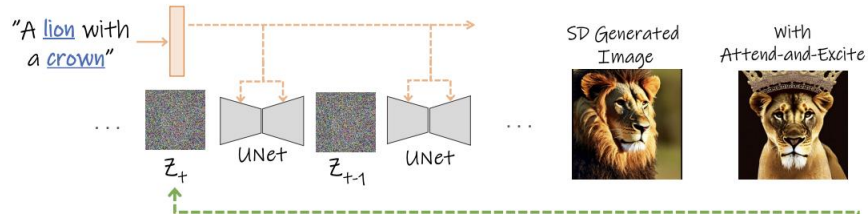
$Q = W_Q^{(i)} \cdot \varphi_i(z_t), K = W_K^{(i)} \cdot \tau_\theta(y), V = W_V^{(i)} \cdot \tau_\theta(y).$

Attend-and-Excite

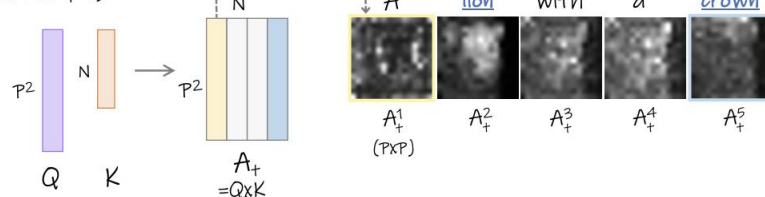
- Enhance the maximal attention for the most neglected subject token
- Updates the latent with attention loss



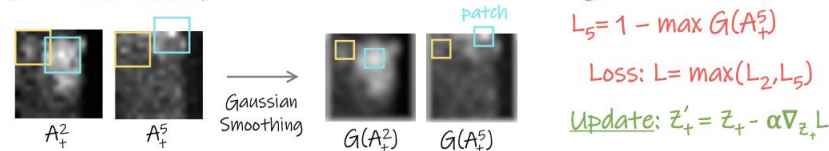
DDPM Process



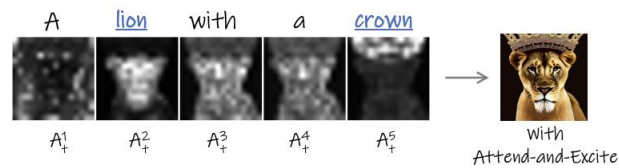
Cross Attention (timestep t)



Loss Computation (tokens “lion”, “crown”)



Final Cross-Attention Maps (timestep t=0)



Attend-and-Excite

"A horse and a dog"



Stable
Diffusion



+Attend-
and-Excite

"A painting of an elephant with glasses"

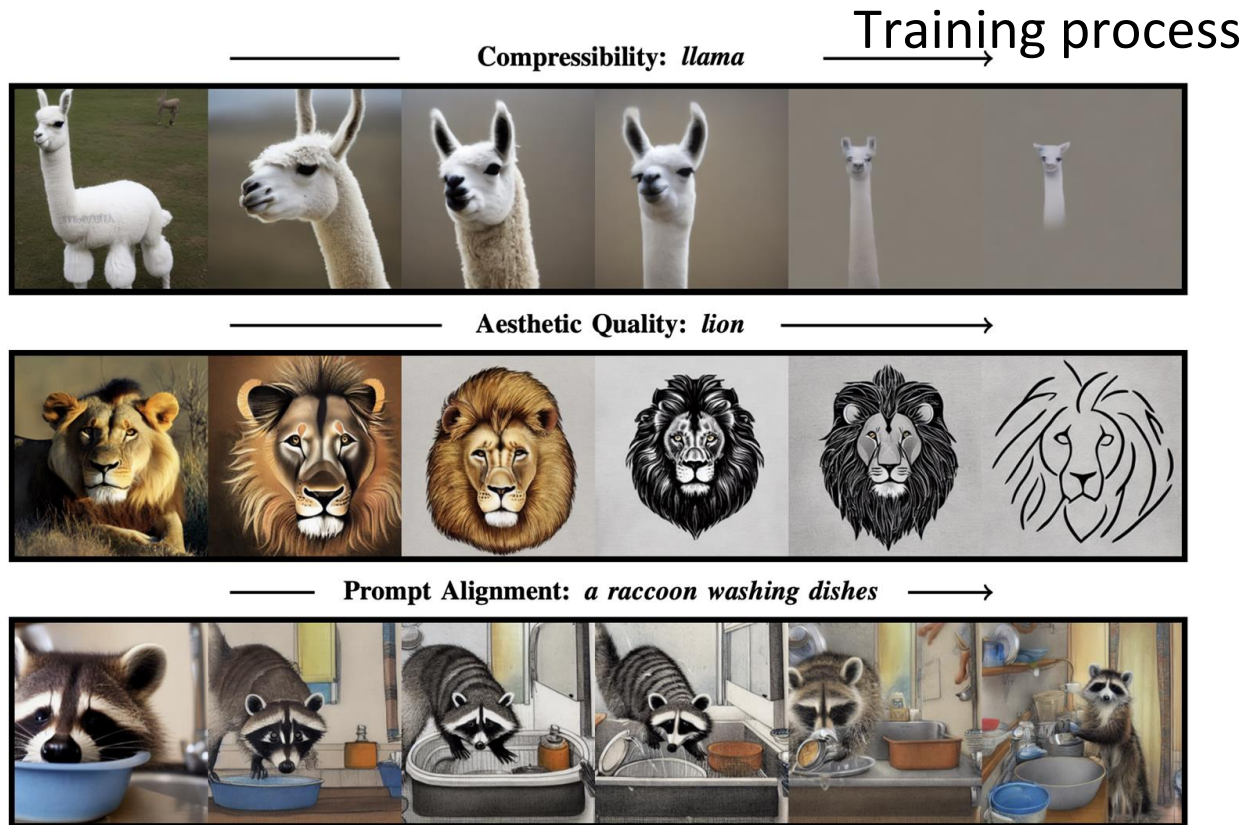


"A playful kitten chasing a butterfly in a wildflower meadow"



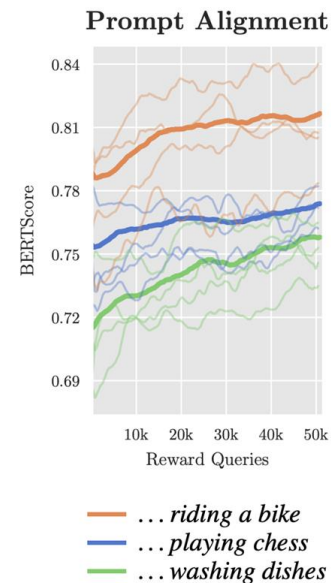
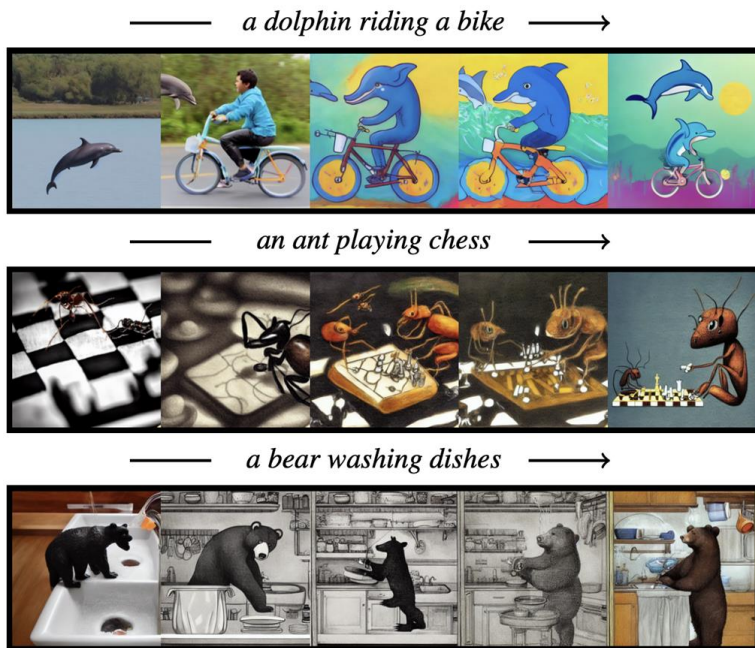
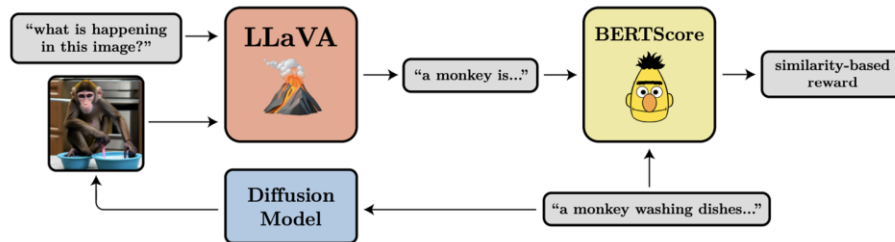
DDPO

- RL for optimizing diffusion models on different downstream objectives



DDPO

- VLM similarity reward to improve image-prompt alignment



Concept Customization

- Single-concept customization
- Multi-concept customization
- Without test-time finetuning



Input images



in the Acropolis



swimming



sleeping



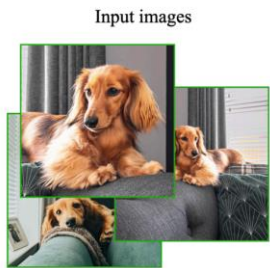
in a doghouse



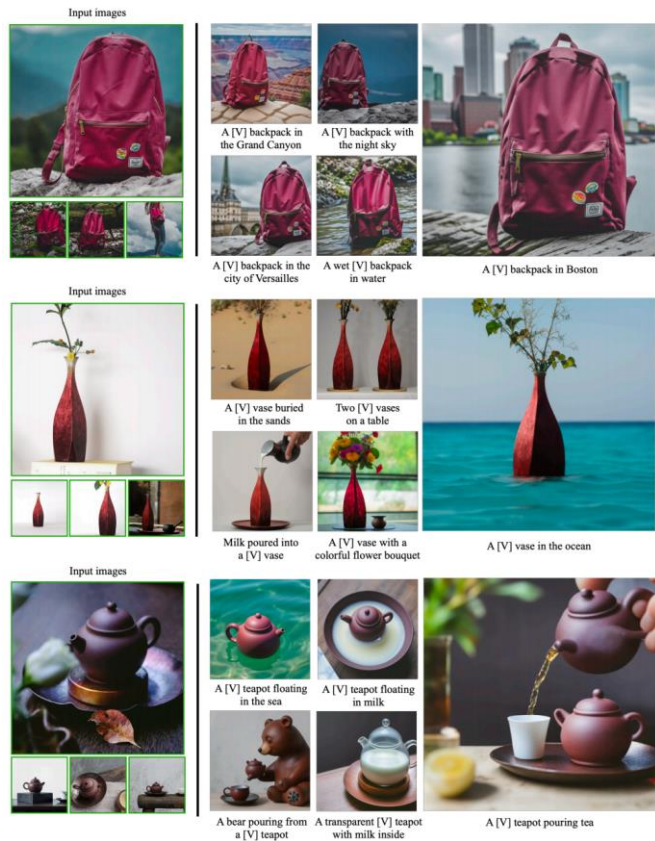
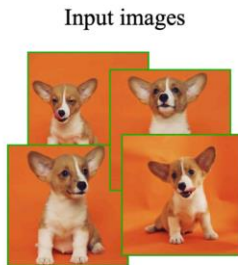
in a bucket

- [1] [An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion](#)
- [2] [DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation](#)
- [3] [Encoder-based Domain Tuning for Fast Personalization of Text-to-Image Models](#)
- [4] [ELITE: Encoding Visual Concepts into Textual Embeddings for Customized Text-to-Image Generation](#)
- [5] [Multi-Concept Customization of Text-to-Image Diffusion](#)
- [6] [Break-A-Scene: Extracting Multiple Concepts from a Single Image](#)
- [7] [Paint by Example: Exemplar-based Image Editing with Diffusion Models](#)
- [8] [BLIP-Diffusion: Pre-trained Subject Representation for Controllable Text-to-Image Generation and Editing](#)
- [9] [Face0: Instantaneously Conditioning a Text-to-Image Model on a Face](#)
- [10] [FastComposer: Tuning-Free Multi-Subject Image Generation with Localized Attention](#)
- [11] [Unified Multi-Modal Latent Diffusion for Joint Subject and Text Conditional Image Generation](#)
- [12] [Re-Imagen: Retrieval-Augmented Text-to-Image Generator](#)
- [13] [InstantBooth: Personalized Text-to-Image Generation without Test-Time Finetuning](#)
- [14] [Subject-driven text-to-image generation via apprenticeship learning](#)

Single-Concept Customization

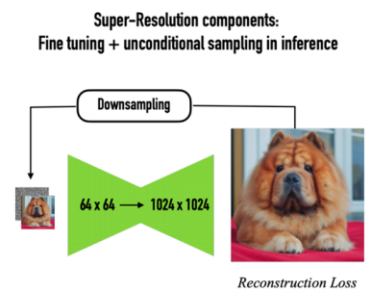
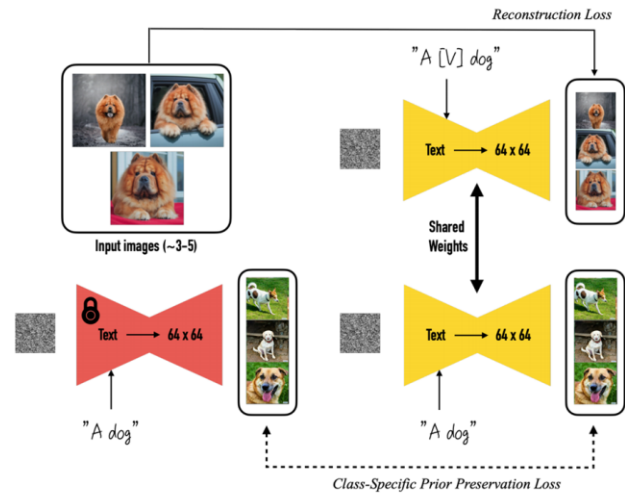
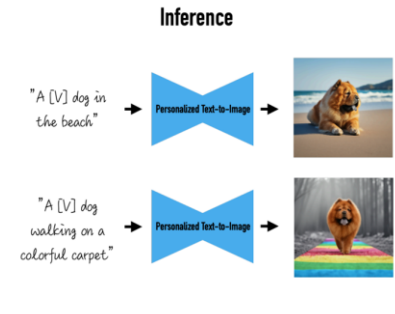
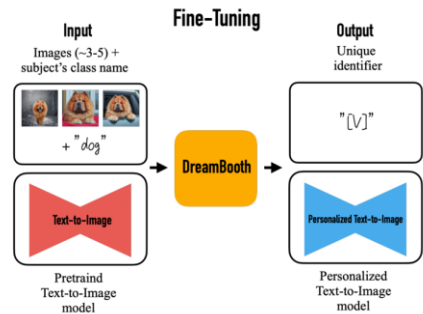


Expression modification (“A [state] [V] dog”)



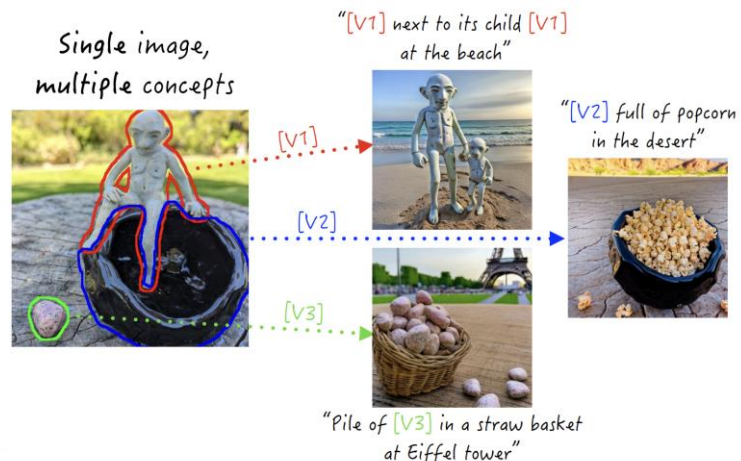
Single-Concept Customization

- Tuning unique identifier [V] for customized subject
- Originally generated samples to alleviate forgetting



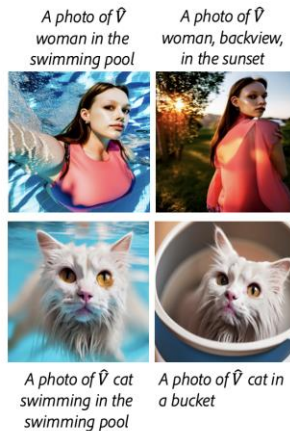
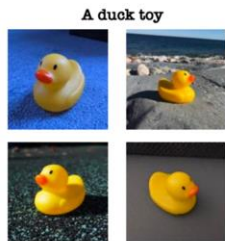
Multi-Concept Customization

- Multi-concept customization [V1], [V2], ... from single image or multiple images



Without Test-Time Finetuning

- Retrieve-augmented/ In-context generation
- Similar customization, but w/o test-time finetuning

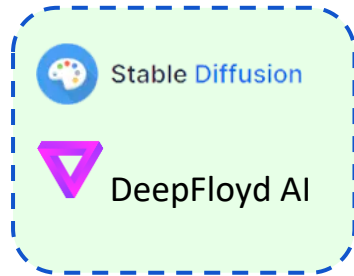


Agenda

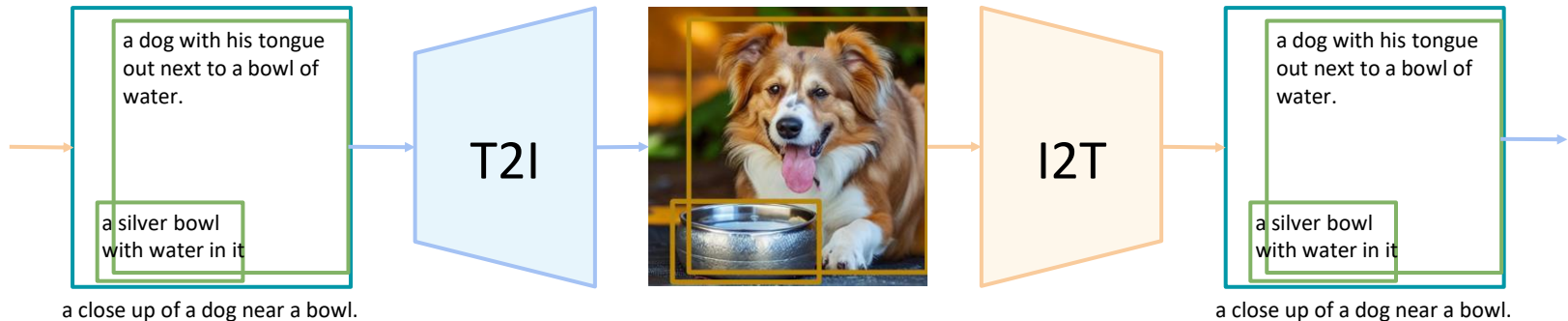
- Text-to-image (T2I) basics
- Aligning human intentions in T2I generation
 - Controllable generation
 - Editing
 - Better following prompts
 - Concept customization
- Summary and discussion

Discussion

Open-source v.s. Closed-source



Consuming and producing visual data: Understanding (I2T) and generation (T2I) loop



Thank you!