

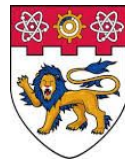


# Attention Interpolation of text-to-image Diffusion

NeurIPS 2024

Qiyuan He<sup>1</sup>, Jinghao Wang<sup>2</sup>, Ziwei Liu<sup>2</sup>, Angela Yao<sup>1</sup>

<sup>1</sup>National University of Singapore, <sup>2</sup>Nanyang Technological University



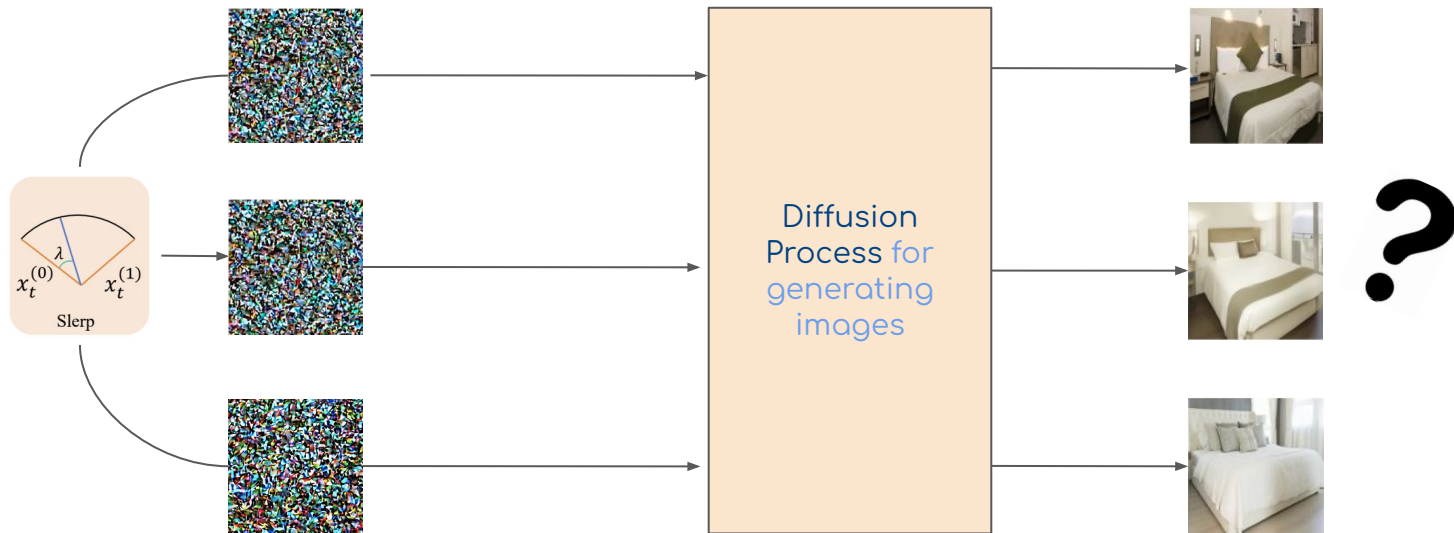
**NANYANG  
TECHNOLOGICAL  
UNIVERSITY**  
SINGAPORE



**NUS**  
National University  
of Singapore

# Background on unconditional deep interpolation

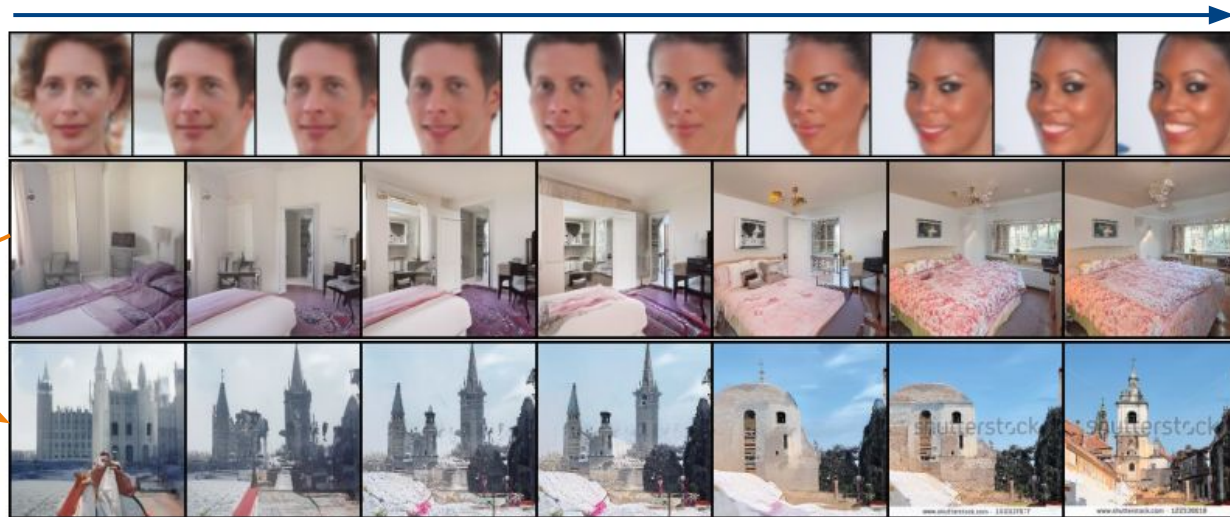
- Target at generating image sequence given two generated images
- Mainly focus on interpolation along the latent space of initial seeds



# 01 Conditional deep interpolation

- What if we want to interpolate between different conditions?

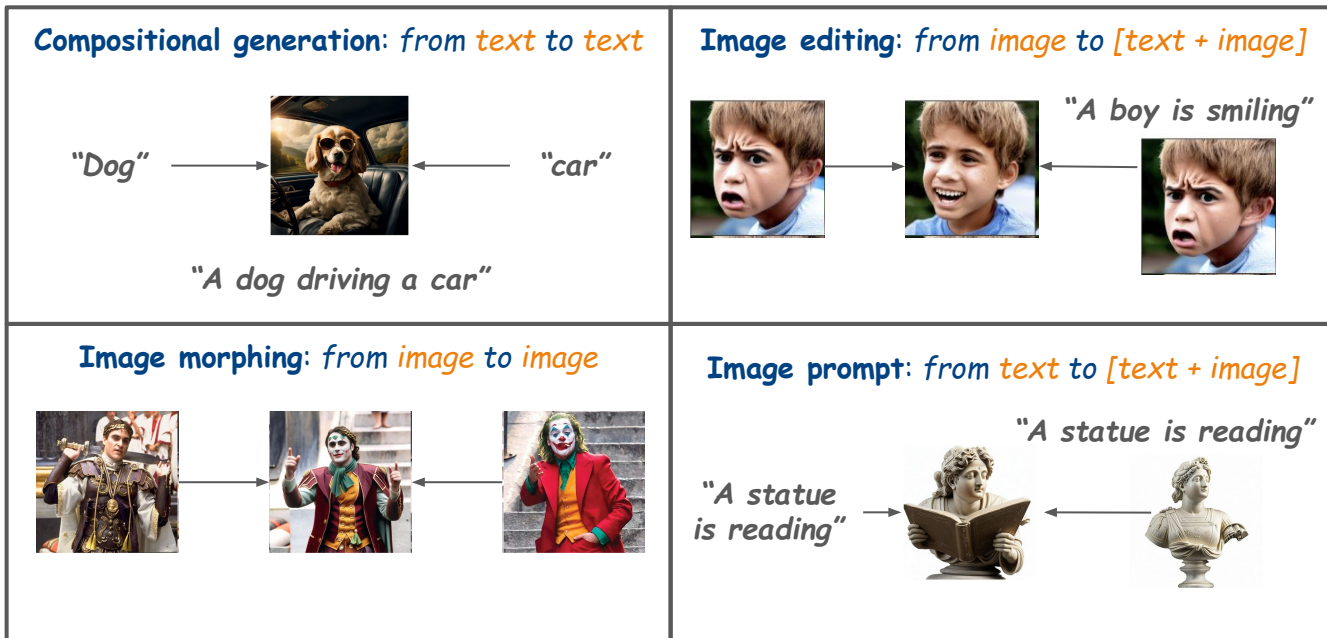
Interpolation between different initial seeds in latent space



How do we interpolate between different conditions?

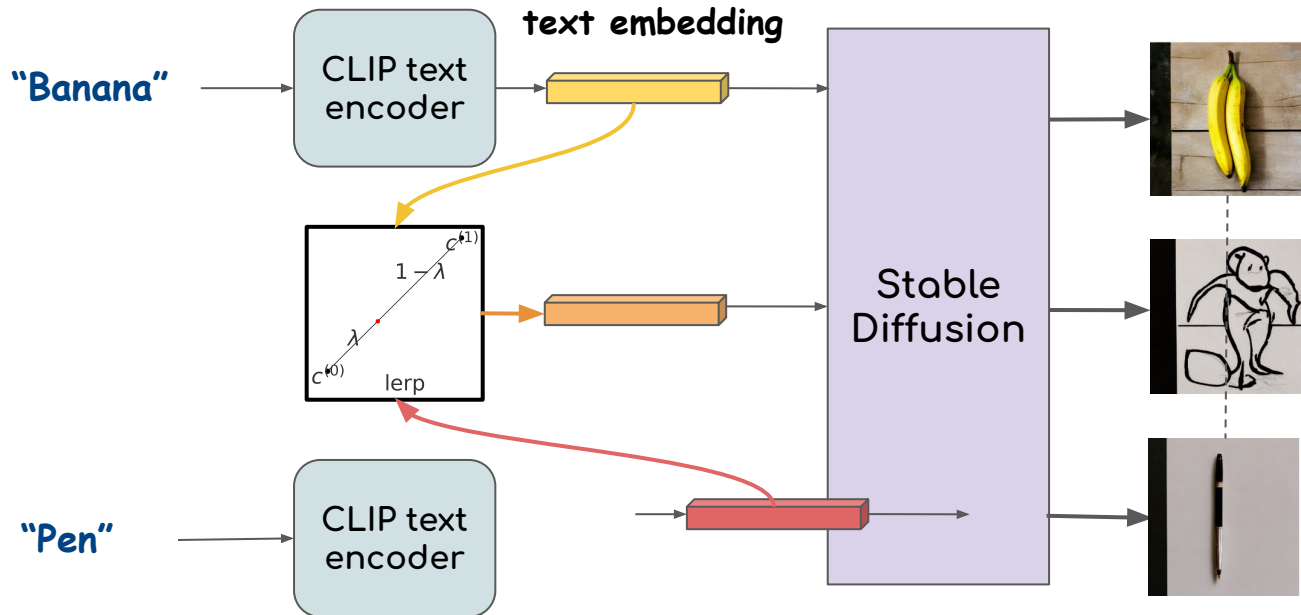
# Motivation for conditional deep interpolation

- Deep interpolation different two conditions is ***important*** in many applications
- Rarely explored as independent research question



## Failure on text embedding interpolation

- Text embedding interpolation is implicitly used in many methods for application
- but **not roundly evaluated**
- The interpolated image does not look *“nice”*?



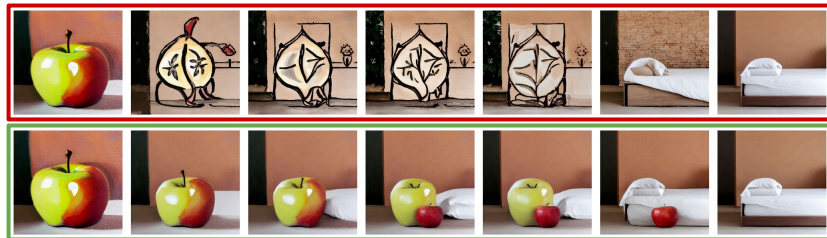
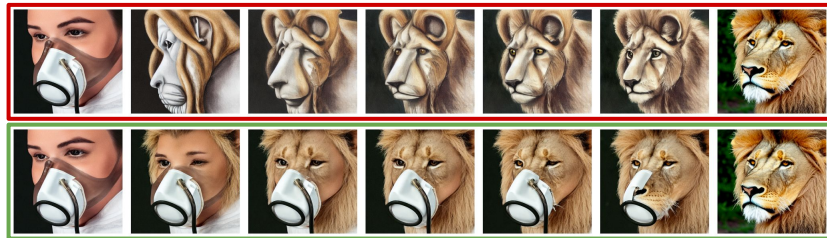
## Evaluation on the quality of deep interpolation

- This is not a good interpolation intuitively, but what is exactly a “good” or “bad” interpolation?



# Evaluation on the quality of deep interpolation

- Perceptual consistency
  - How direct the transition is?
  - Learned Perceptual Image Patch Similarity (LPIPS) model [1] for perceptual distance
- Perceptual smoothness
  - How smooth the transition is?
  - LPIPS model [1] for perceptual distance
- Fidelity
  - What about the quality of interpolated image?
  - Fréchet Inception Distance (FID) [2] as the quality proxy



03

# How to improve this?



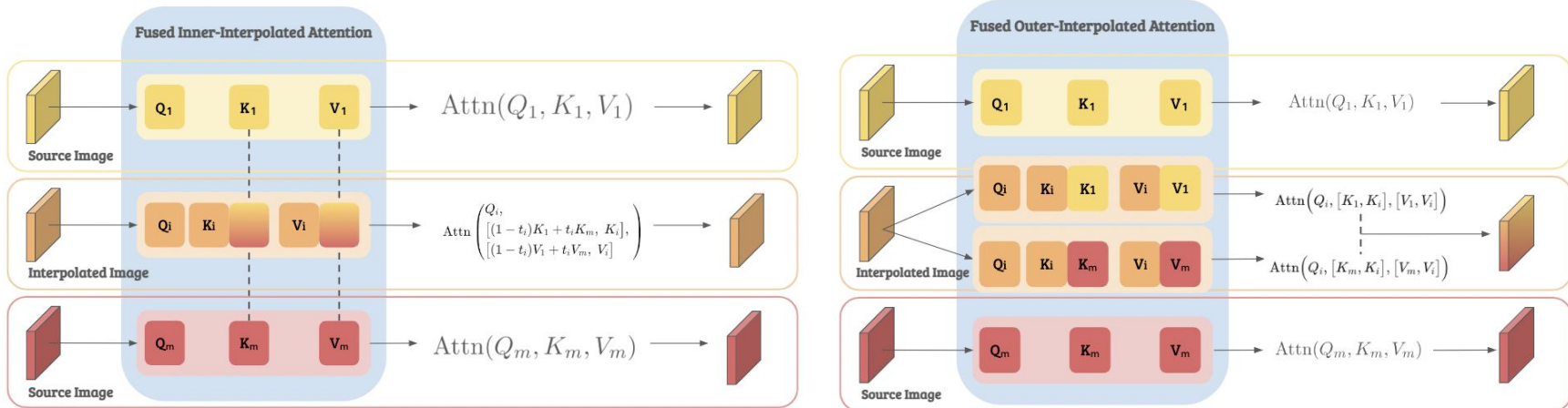


## PAID: Prompt-guided Attention Interpolation of Diffusion

- Fully training-free
- Fused interpolated attention
  - “Interpolated”: boost consistency
  - “Fused”: boost quality
- Beta prior selection
  - Boost smoothness
- Prompt guidance
  - Enable selecting interpolation path via the third condition

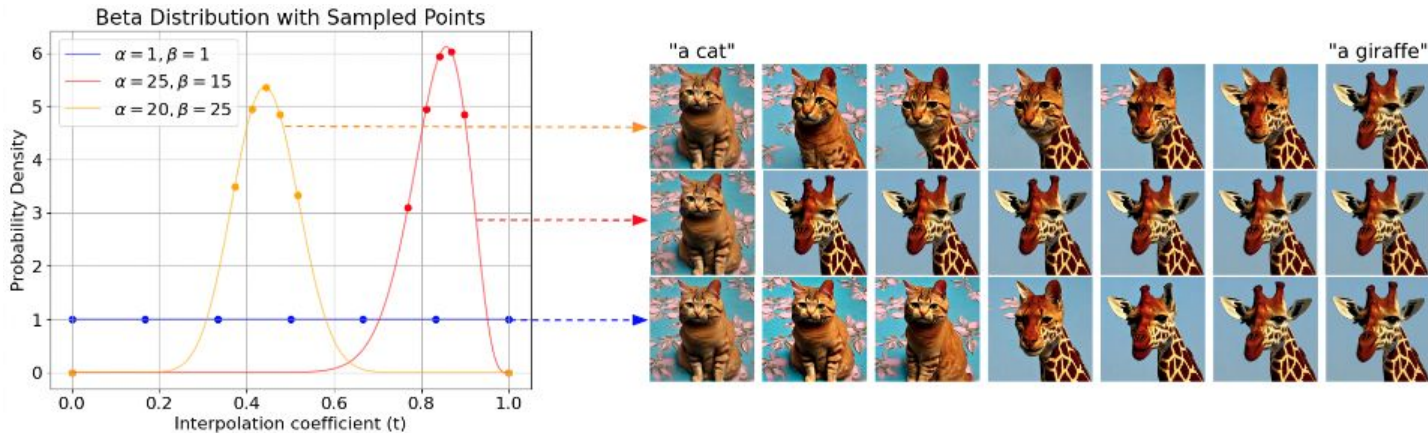
## PAID Workflow: Fused interpolated attention

- Expand interpolation to **both cross-attention & self-attention**
- Concatenate with original self-attention to enhance fidelity
- Inner / outer interpolation



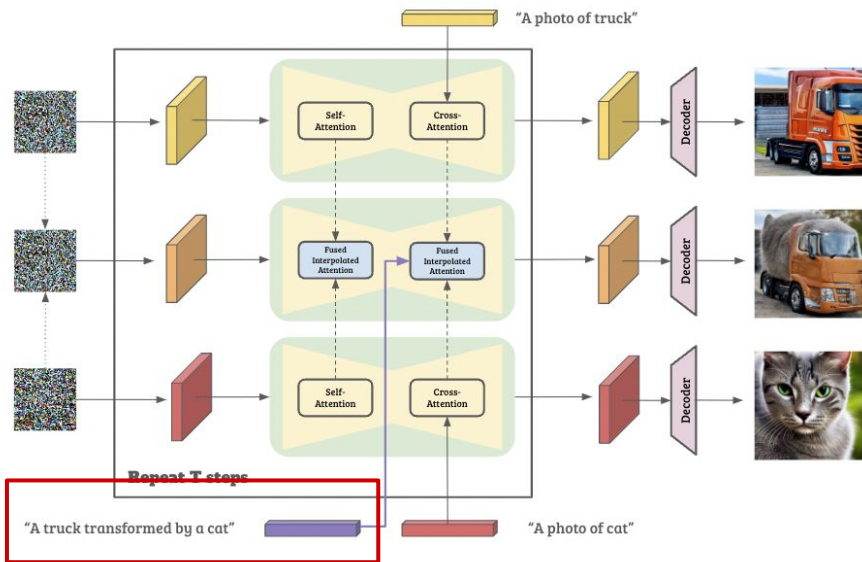
# 04 PAID Workflow: Beta prior selection

- Beta prior is good for three reasons:
  - Bell shape: encourage more selection in a small range
  - Uniform distributed point as the lower bound
  - Can make up for the bias towards one side
- We use dynamic selection to auto-fit the alpha and beta of the prior



# 04 PAID Workflow: Prompt guidance

- Since self-attention already poses spatial constraints
  - Loose semantic constraints in cross-attention
- Inject the guidance prompt in cross-attention



(f) “Photo of a dog” to “Photo of a car”, guided with “A dog driving car” (top), “A car with dog furry texture” (middle), and “A toy named dog-car” (bottom).

# 05 Experiment: Ablation Study

attention interpolation



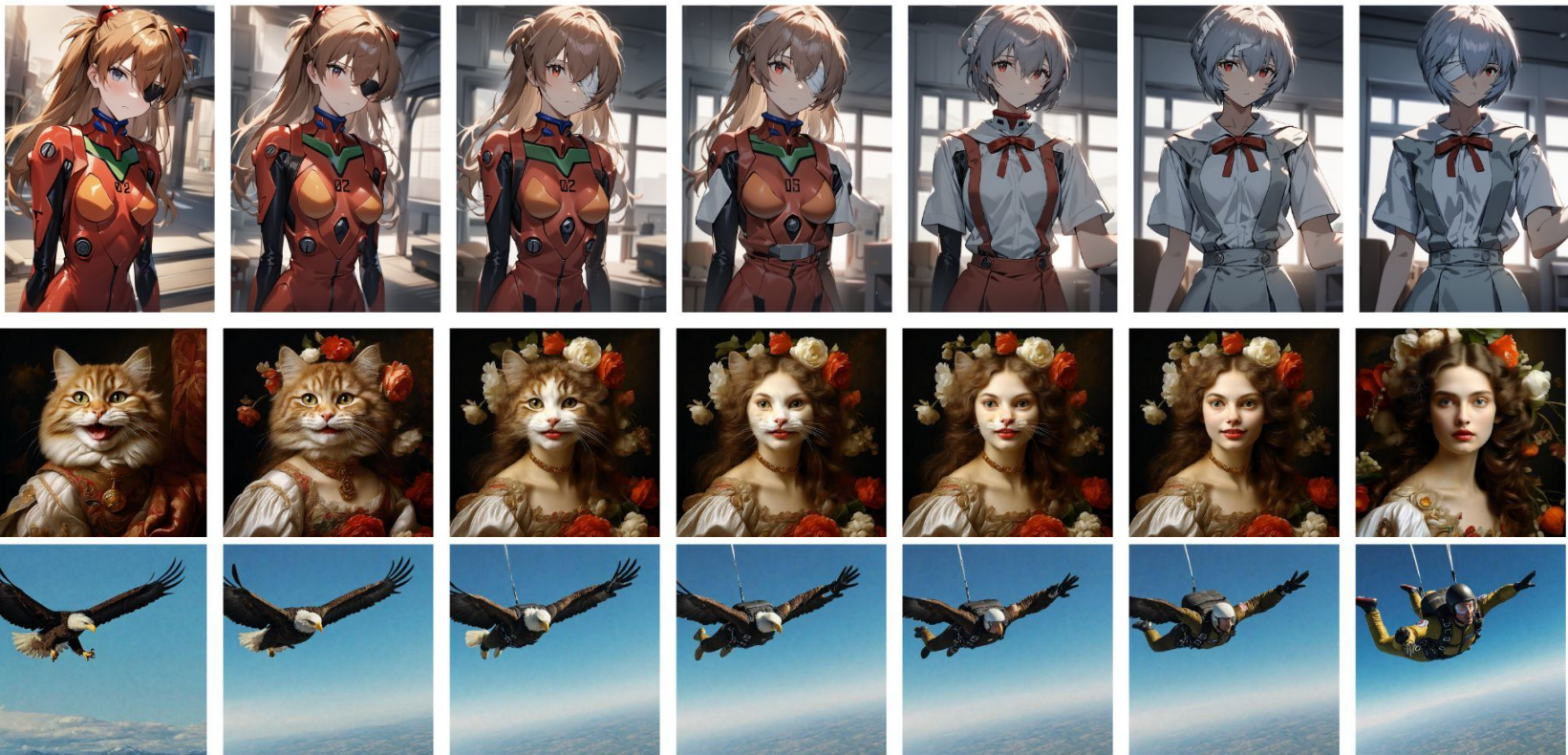
Fused interpolated attention



Fused interpolated attention  
+ Beta prior selection



# Experiment: Qualitative results



## Experiment: Quantitative Results

- Boost three evaluation metrics by a large margin
  - Especially consistency & fidelity

| Dataset          | Method | Smoothness ( $\uparrow$ ) | Consistency ( $\downarrow$ ) | Fidelity ( $\downarrow$ ) |
|------------------|--------|---------------------------|------------------------------|---------------------------|
| CIFAR-10         | TEI    | 0.7531                    | 0.3645                       | 118.05                    |
|                  | DI     | 0.7564                    | 0.4295                       | 87.13                     |
|                  | AID-O  | 0.7831                    | <b>0.2905*</b>               | <b>51.43*</b>             |
|                  | AID-I  | <b>0.7861*</b>            | 0.3271                       | 101.13                    |
| LAION-Aesthetics | TEI    | 0.7424                    | 0.3867                       | 142.38                    |
|                  | DI     | 0.7511                    | 0.4365                       | 101.31                    |
|                  | AID-O  | 0.7643                    | <b>0.2944*</b>               | <b>82.01*</b>             |
|                  | AID-I  | <b>0.8152*</b>            | 0.3787                       | 129.41                    |

## 05 Experiment: Human study

- Dominantly preferred by human evaluation

| Interpolation method | Near Object   | Far Object | Scene      | Object+Scene  |
|----------------------|---------------|------------|------------|---------------|
| TEI                  | 8.75%         | 1.16%      | 0%         | 1.26%         |
| AID-I                | <b>53.75%</b> | <b>50%</b> | 45.2%      | 45.57%        |
| AID-O                | 36.25%        | 46.5%      | <b>50%</b> | <b>51.90%</b> |
| Hard to determine    | 1.25%         | 2.32%      | 4.76%      | 1.26%         |



# 05 Applications

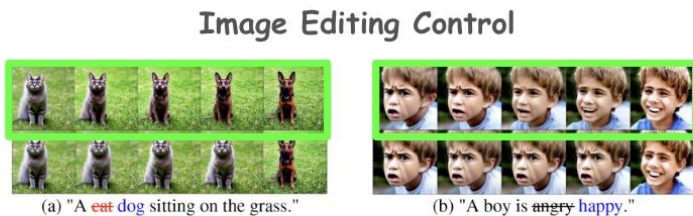
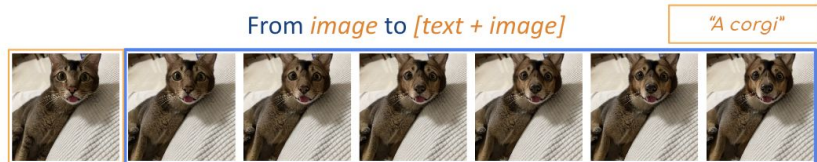
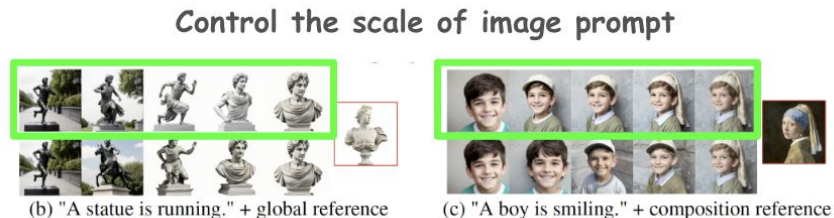


Figure 5: Results of image editing control. Our method boosts the controlling ability over editing. The first row of (a) and (b) is generated by P2P + AID while the second row is P2P + TEI.





- We formulate a problem called conditional deep interpolation and propose corresponding evaluation metric
- We analyze the behaviour of attention interpolation and tackle the formulated problem by a training-free method PAID
  - Interpolated attention → consistency
  - Fused interpolated attention → fidelity
  - Beta prior → smoothness
  - Prompt guidance → one-to-many answers
- This approach provides stronger control ability on various application
  - Text-to-Text: compositional generation
  - Image-to-[Text+Image]: image editing
  - Image-to-Image: image morphing
  - Text-to-[Text+Image]: image-control generation