



Multimodal Understanding & Generation with Efficiently Finetuned Foundation Models

Dr. Long Chen

Assistant Professor, CSE, HKUST

Feb. 2025

In Collaboration with: Yanghao, Wei, Lin, Zhen, Ziqi, Hongxiang (HKUST)

Yuxuan (NTU), Haoxuan (Columbia), Hongzhan (HKBU), Jiazuo (DLUT)

There are lots of pretrained foundation models...

- Large language models (LLMs)
 - ChatGPT, GPT-4, ...
- Vision language models
 - CLIP, BLIP, ...
- Visual generation models
 - Stable Diffusion, ...
- *Research Q1: How can we efficiently train or finetune foundation models.*
- *Research Q2: How can we build strong open-world multimodal understanding and generation models with these pretrained foundation models*

Efficient Finetune Foundation Models

*Parameter-
Efficient Tuning*

*Memory-
Efficient Tuning*

*Modality-
Efficient Tuning*

LLMs Can Evolve Continually on Modality for X-Modal Reasoning

**Jiazuo Yu¹, Haomiao Xiong¹, Lu Zhang^{1,*}, Haiwen Diao¹, Yunzhi Zhuge¹,
Lanqing Hong², Dong Wang¹, Huchuan Lu¹, You He³, Long Chen⁴**

¹Dalian University of Technology, ²Huawei Noah's Ark Lab

³Tsinghua University, ⁴The Hong Kong University of Science and Technology

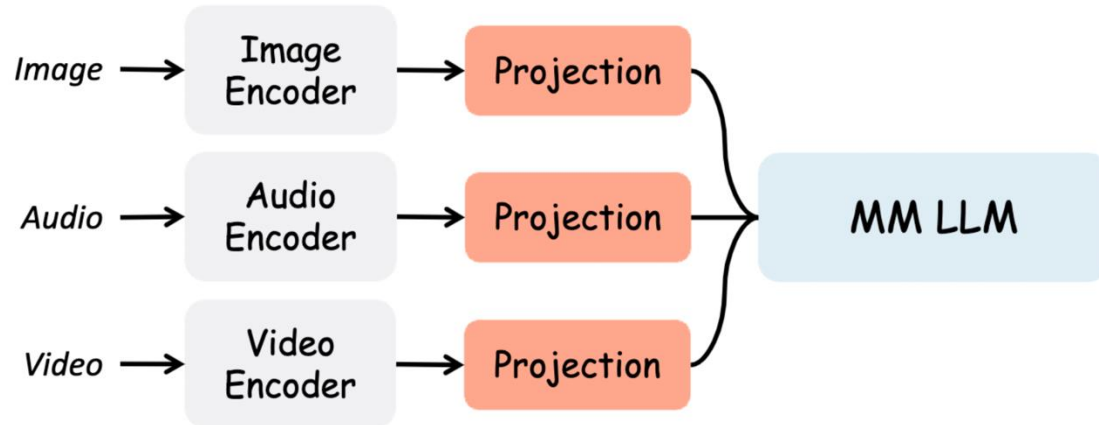
<https://arxiv.org/pdf/2410.20178>

(NeurIPS'24)

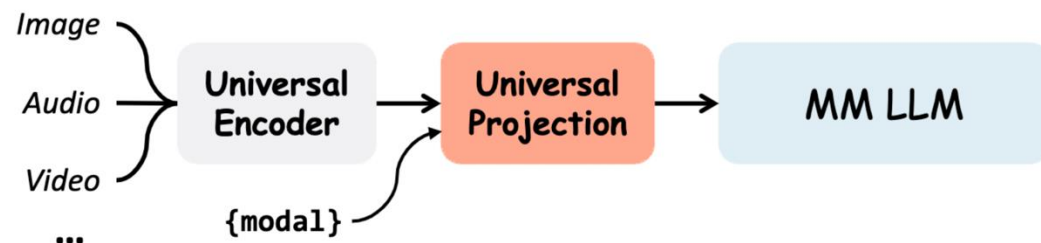
Prior Work on Multimodal LLMs



[1]LLaVA, [2]InstructBLIP, [3]InternVL...



[4]X-InstructBLIP, [5]X-LLM, [6]ChatBridge...



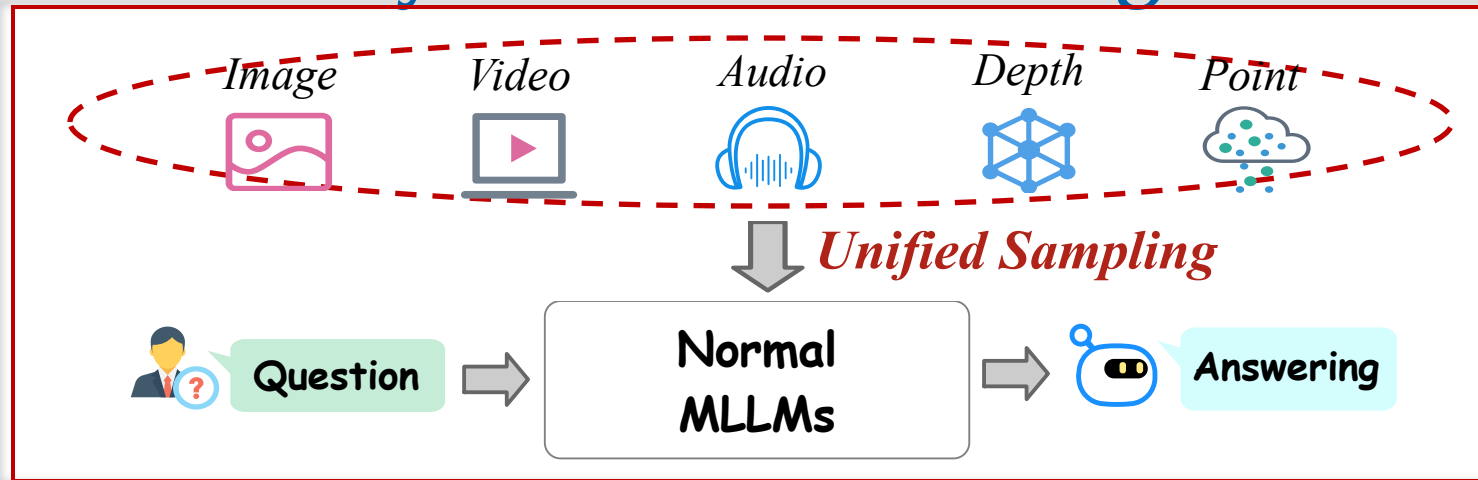
[7]OneLLM...

- Single-modality MLLM (**only image modality**)

- Multiple-modality MLLM (**different projector parameters are different**)

- Multiple-modality MLLM (**mixture multi-modal data for training**)

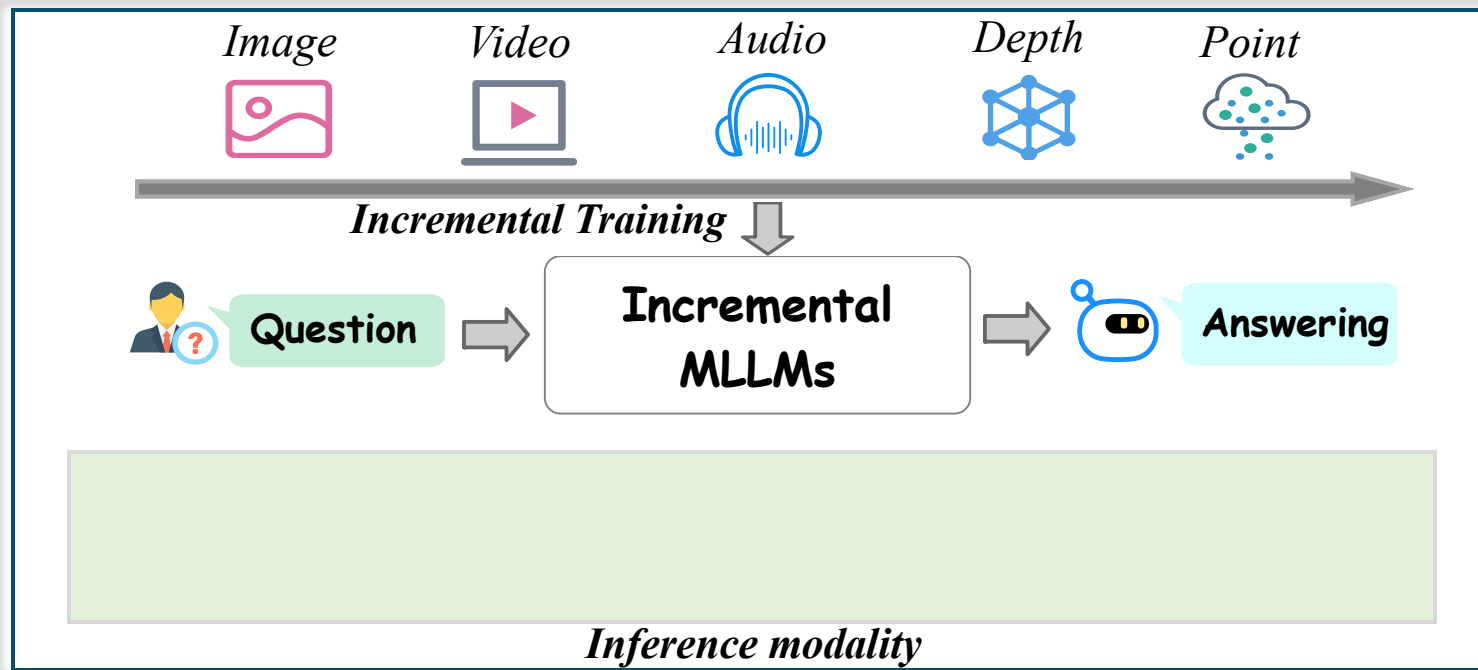
Modality-Efficient Training for X-Modal Reasoning



Existing work

Limited data

Hard to extend



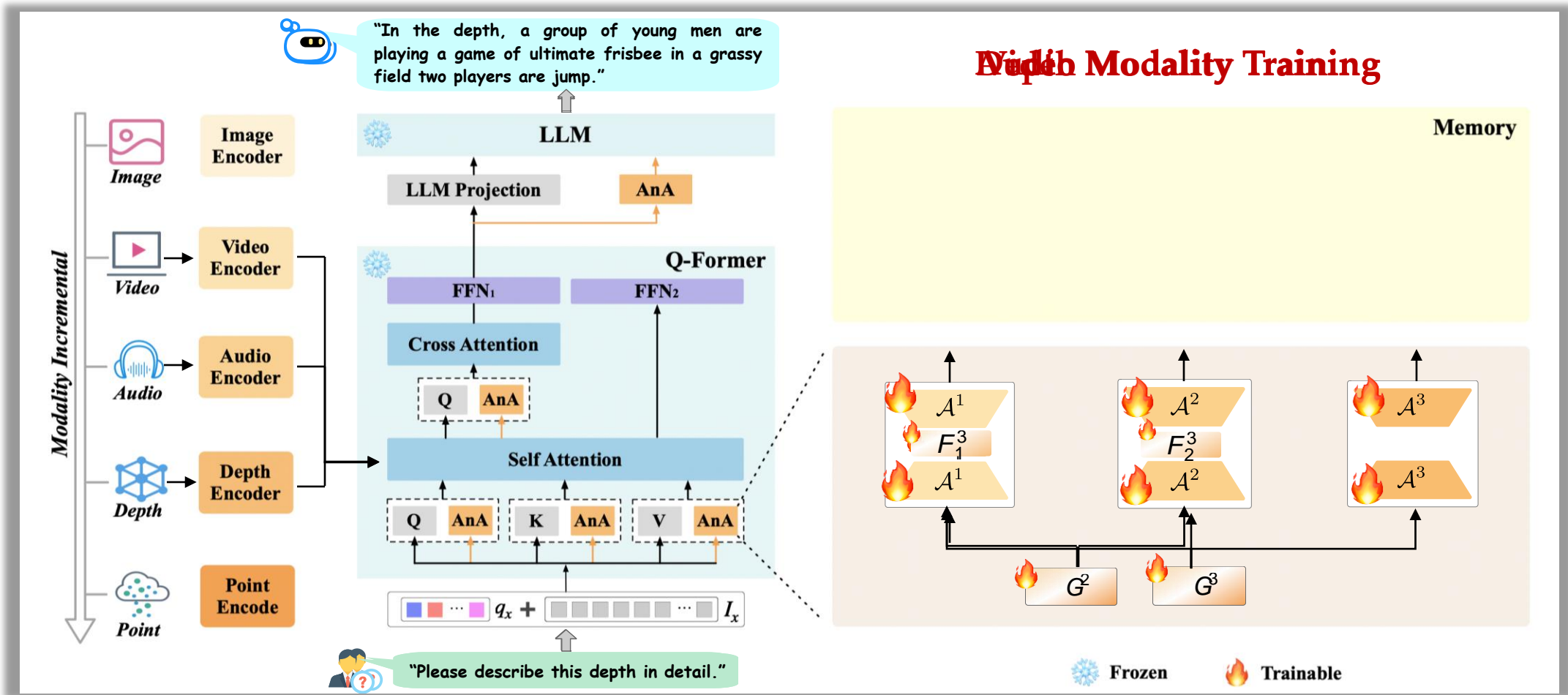
Ours

Continual training
Start from different
checkpoints



Jiazuo Yu, et al. LLMs can Evolve Continually on Modality for X-Modal Reasoning. In NeurIPS, 2024.

Modality-Efficient Training for X-Modal Reasoning



Modality-Efficient Training for X-Modal Reasoning

- **PathWeave**: It can performs well on all previous trained modalities

Diagram illustrating the PathWeave model's performance across five modalities, each with a user query and a corresponding AI response placeholder:

- Image Modality:** User query: "Please describe this image in detail." (Image: A scenic view of a lake and hills under a blue sky). AI response placeholder: [Empty light blue speech bubble]
- Image Modality:** User query: "What is the man doing?" (Image: A man wearing sunglasses walking outdoors). AI response placeholder: [Empty light blue speech bubble]
- Audio Modality:** User query: "How do you feel when you hear the audio?" (Image: A green speaker icon with sound waves). AI response placeholder: [Empty light blue speech bubble]
- Depth Modality:** User query: "What can you see in this depth?" (Image: A depth map showing green and blue areas). AI response placeholder: [Empty light blue speech bubble]
- 3D Model Modality:** User query: "What is the 3d model?" (Image: A 3D model of a mechanical part). AI response placeholder: [Empty light blue speech bubble]

Efficient Finetune Foundation Models

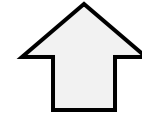
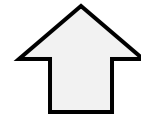
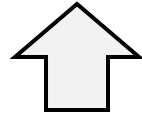
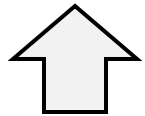
*Parameter-
Efficient Tuning*

*Memory-
Efficient Tuning*

*Modality-
Efficient Tuning*

Multimodal Understanding & Generation with Efficient Finetune Foundation Models

Open-World Perception



Parameter-Efficient Tuning *Memory-Efficient Tuning* *Modality-Efficient Tuning*
Efficiently Finetuned Foundation Models

Inversion Circle Interpolation: Diffusion-based Image Augmentation for Data-scarce Classification

Yanghao Wang, Long Chen
The Hong Kong University of Science and Technology

<https://arxiv.org/abs/2408.16266>
(Under Review)

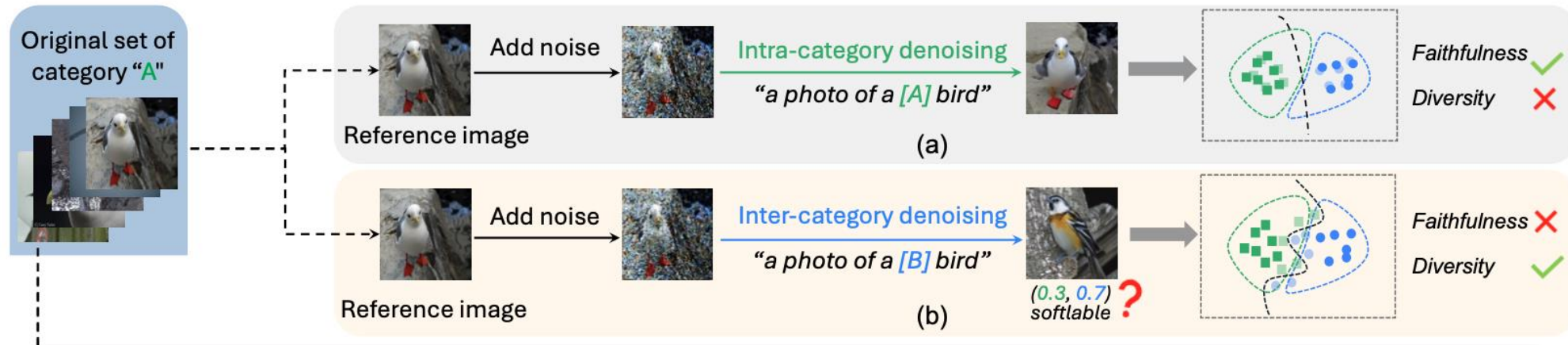
Zero-shot Visual Relation Detection via Composite Visual Cues from Large Language Models

Lin Li^{1,2}, Jun Xiao¹, Guikun Chen¹, Jian Shao¹, Yueting Zhuang¹, Long Chen^{2*}
¹Zhejiang University ²The Hong Kong University of Science and Technology

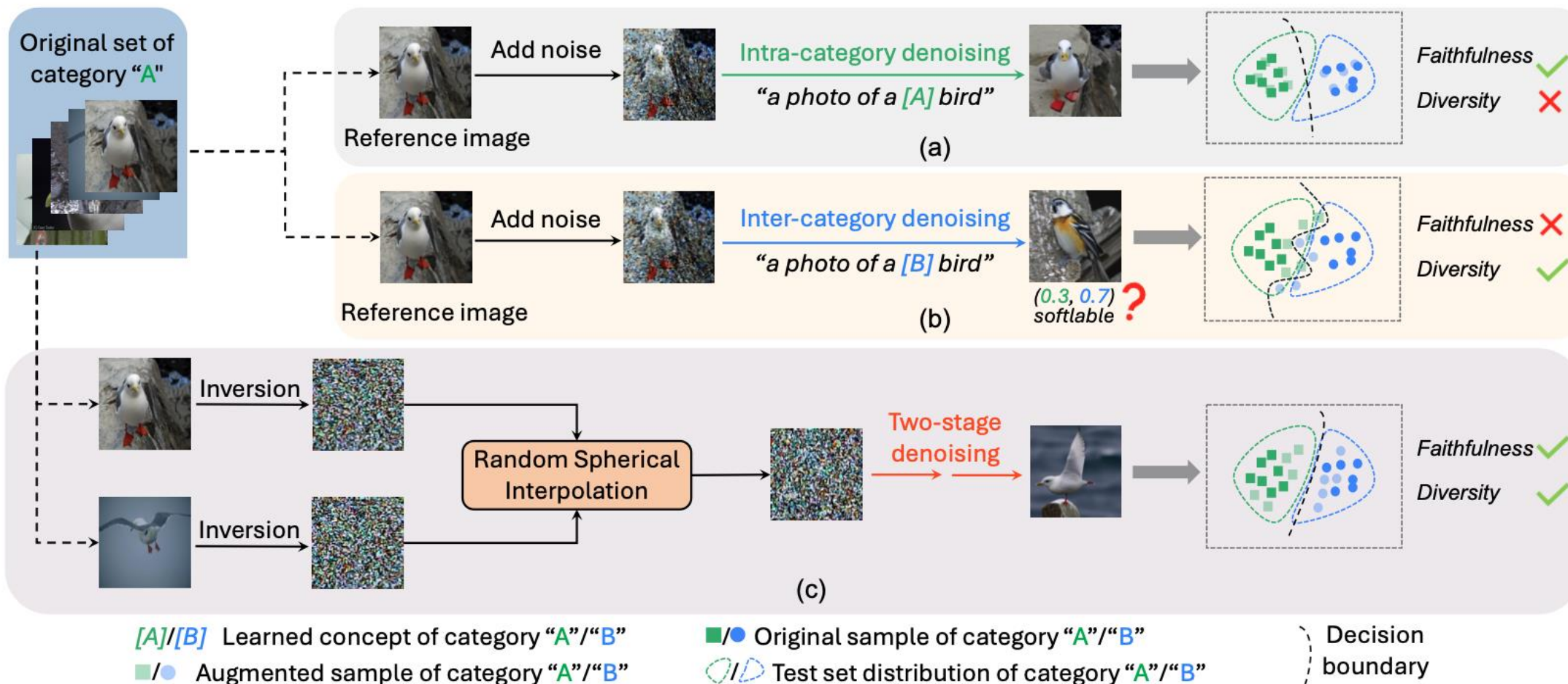
<https://arxiv.org/abs/2305.12476>
(NeurIPS'23)

Diffusion Model for Closed-set Perception

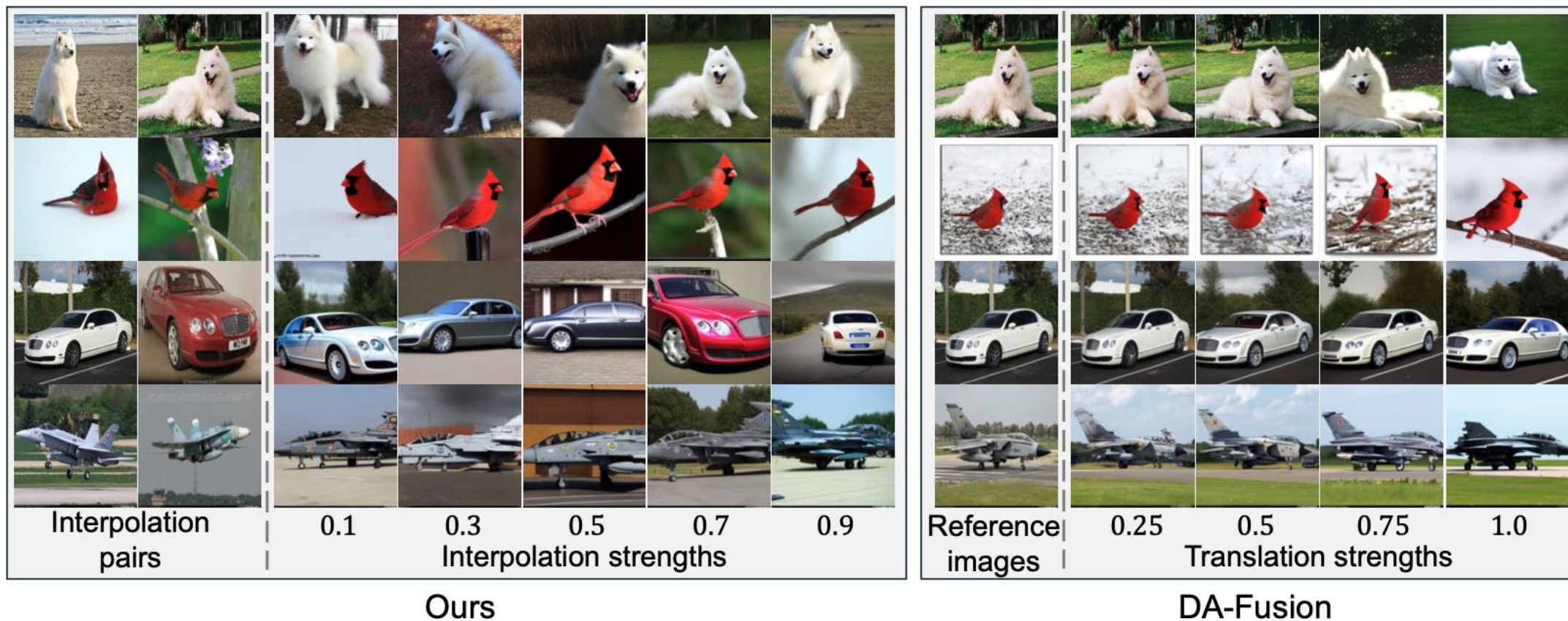
- Diffusion model is a text-to-image generation model.
- We can use diffusion model to conduct data augmentation.



Diffusion Model for Data Augmentation



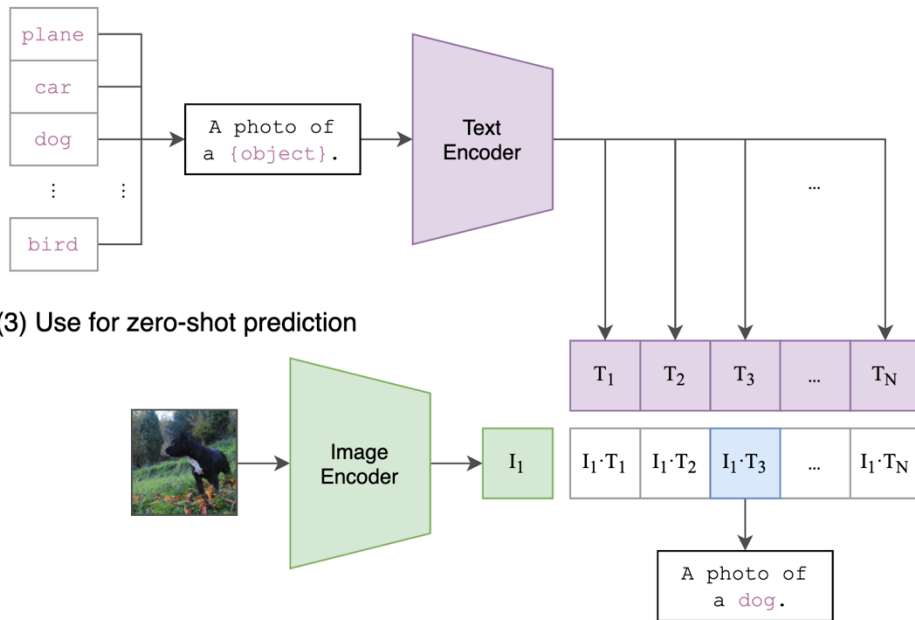
Diffusion Model for Data Augmentation



LLM + CLIP for Zero-Shot Perception

- LLMs can generate detailed descriptions to help zero-shot classification

(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

School bus

- large, yellow vehicle
- the words "school bus" written on the side
- a stop sign that deploys from the side of the bus
- flashing lights on the top of the bus
- large windows

Shoe store

- a building with a sign that says "shoe store"
- a large selection of shoes in the window
- shoes on display racks inside the store
- a cash register
- a salesperson or customer

Volcano

- a large, cone-shaped mountain
- a crater at the top of the mountain
- lava or ash flowing from the crater
- a plume of smoke or ash rising from the crater

Barber shop

- a building with a large, open storefront
- a barber pole or sign outside the shop
- barber chairs inside the shop
- mirrors on the walls
- shelves or cabinets for storing supplies
- a cash register
- a waiting area for customers

Cheeseburger

- a burger patty
- cheese
- a bun
- lettuce
- tomato
- onion
- pickles
- ketchup
- mustard

Violin

- a stringed instrument
- typically has four strings
- a wooden body
- a neck and fingerboard
- tuning pegs
- a bridge
- a soundpost
- f-holes
- a bow

Pirate ship

- a large, sailing vessel
- a flag with a skull and crossbones
- cannons on the deck
- a wooden hull
- portholes
- rigging
- a crow's nest

CLIP for zero-shot classification


LLMs can generate descriptions

Alec Radford, et al. Learning Transferable Visual Models From Natural Language Supervision. In ICML, 2021.
Sachit Menon, et al. Visual Classification via Description from Large Language Models. In ICLR, 2023.


LLM for Zero-Shot Classification

- Using LLMs to generate detailed descriptions for “challenging” tasks
 - zero-shot relation classification

(A) carrying (B) carrying



(C) holding (D) holding



Holding: a person having an object in their hands
Carrying: a person supporting an object in their hands
Which of the pictures are “holding” and “carrying”?

AC are “holding”, and BD are “carrying”

(a)

CLIP

(b)

holding carrying

a person having an object in their hands

a person supporting an object in their hands

C D A B

Describe the visual features of the predicate "sitting on" in a photo, when subject belongs to [human], object belongs to [product]:

[subject]:

- with legs.
- with hip.

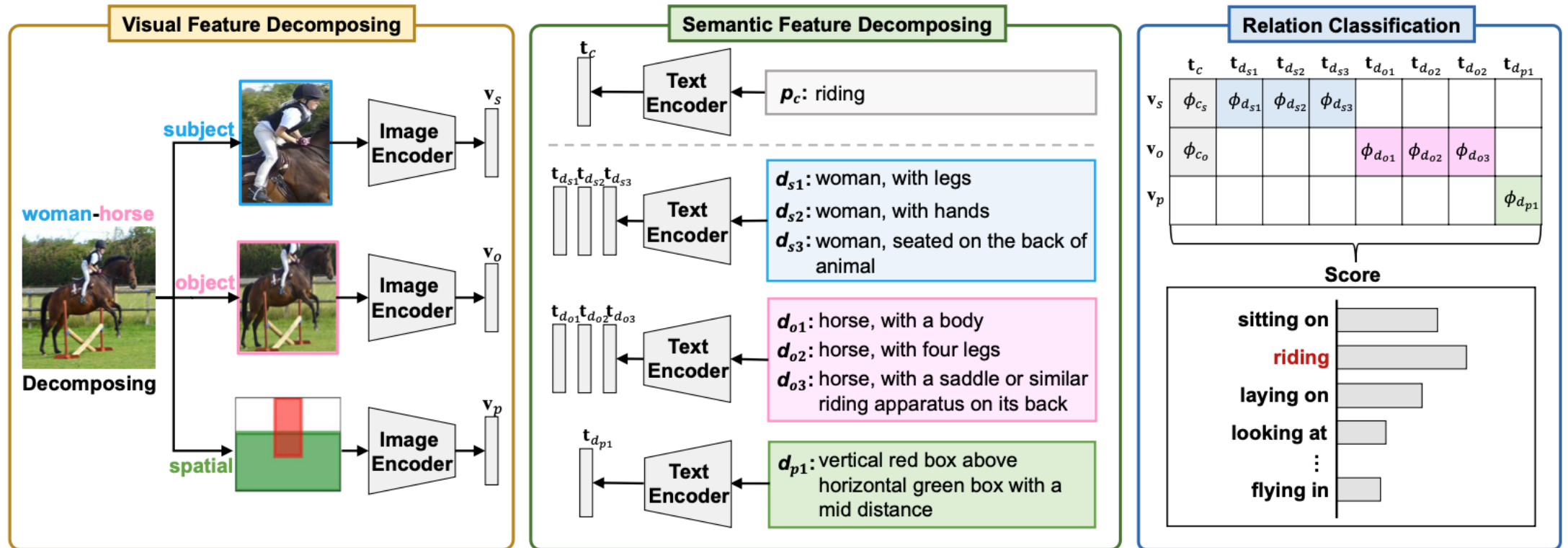
[object]:

- with flat surface.

[position]:

- square subject above horizontal object with a small distance.

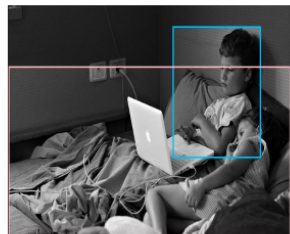
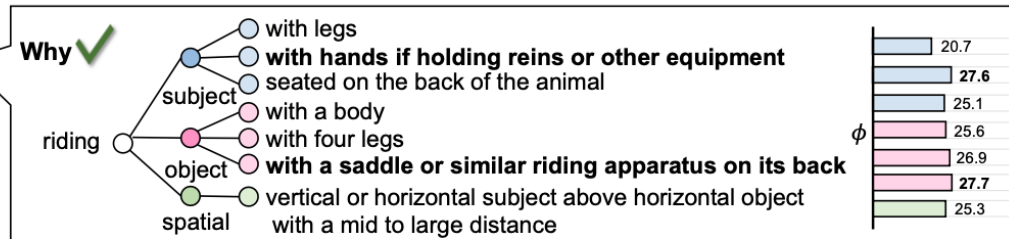
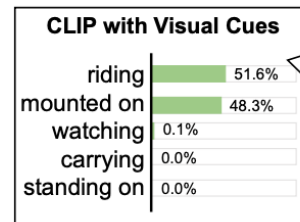
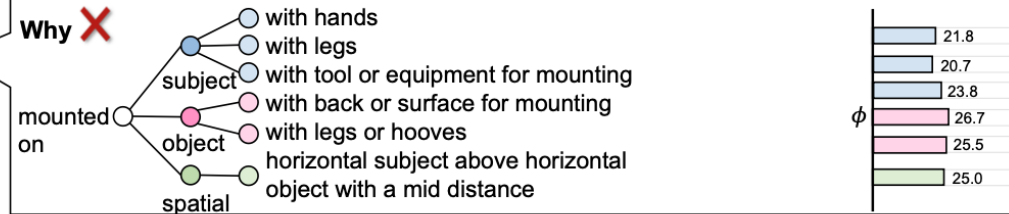
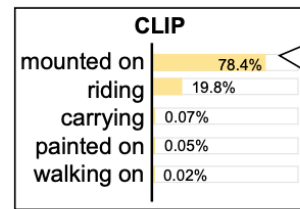
LLM + CLIP for Zero-Shot Perception



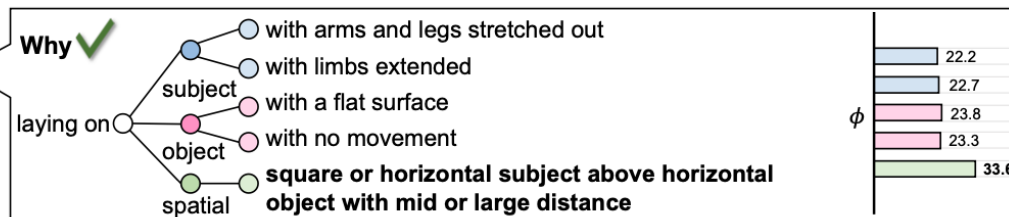
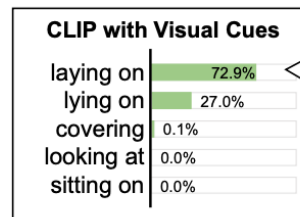
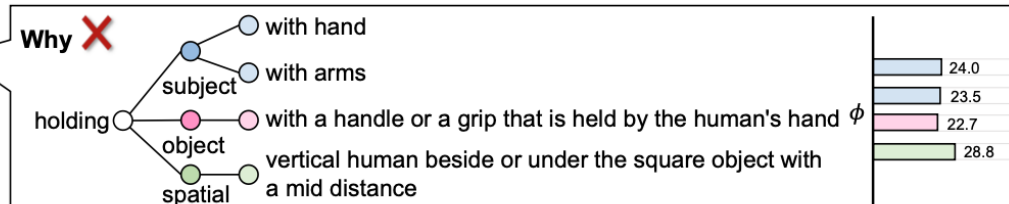
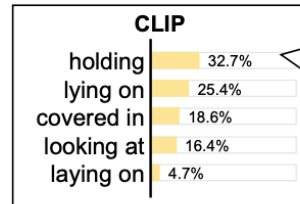
LLM + CLIP for Zero-Shot Perception



man-riding-elephant



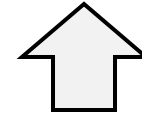
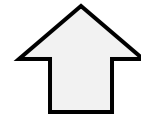
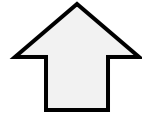
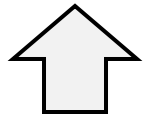
child-laying on-bed



Multimodal Understanding & Generation with Efficient Finetune Foundation Models

Open-World Perception

Multimodal Reasoning



Parameter-Efficient Tuning *Memory-Efficient Tuning* *Modality-Efficient Tuning*
Efficiently Finetuned Foundation Models

Beneath the Surface: Unveiling Harmful Memes with Multimodal Reasoning Distilled from Large Language Models

Hongzhan Lin¹, Ziyang Luo¹, Jing Ma^{1*}, Long Chen²

¹Hong Kong Baptist University

²The Hong Kong University of Science and Technology

<https://arxiv.org/abs/2312.05434>

(EMNLP'23 Findings)

IdealGPT: Iteratively Decomposing Vision and Language Reasoning via Large Language Models

Haoxuan You^{1*}, Rui Sun^{1*}, Zhecan Wang^{1*}, Long Chen², Gengyu Wang³

Hammad A. Ayyubi¹, Kai-Wei Chang⁴, Shih-Fu Chang¹

¹ Columbia University ² HKUST ³ IBM Watson ⁴ University of California, Los Angeles

<https://arxiv.org/abs/2305.14985>

(EMNLP'23 Findings)

LLMs for Harmful Memes Detection

- Harmful Memes Detection



(a) Harmful



(b) Harmful

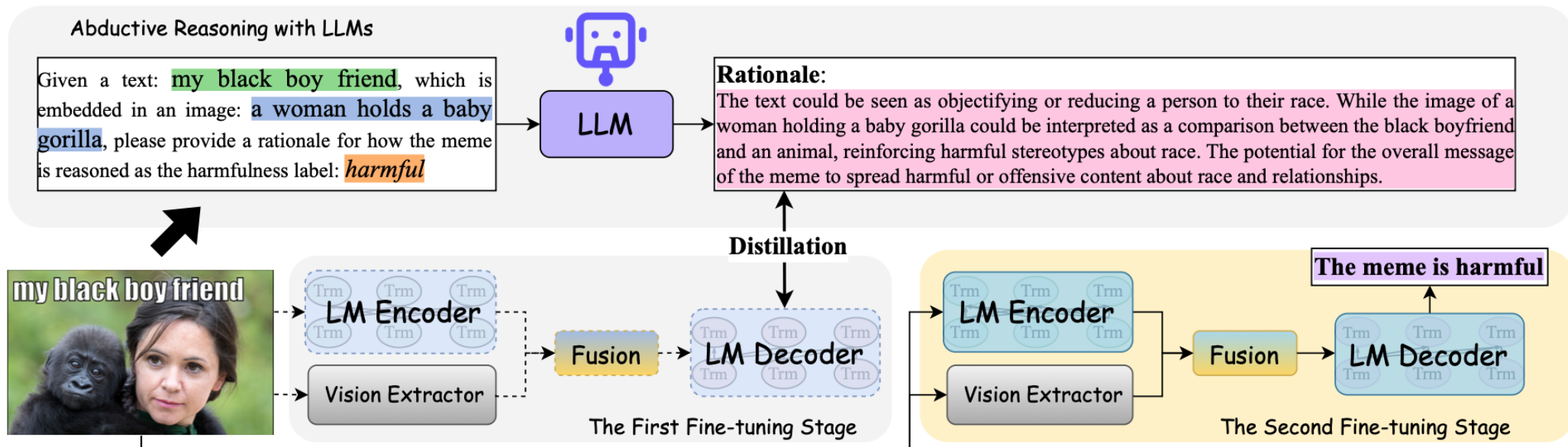


(c) Harmless

Disclaimer: This paper contains discriminatory content that may be disturbing to some readers, where meme examples and words are offensive or hateful in nature. These contents are provided for illustrative purposes only and do not represent the views and standpoints of the authors.

Hongzhan Lin, et al. Beneath the Surface: Unveiling Harmful Memes with Multimodal Reasoning Distilled from Large Language Models. In EMNLP Findings, 2023.

LLMs for Harmful Memes Detection



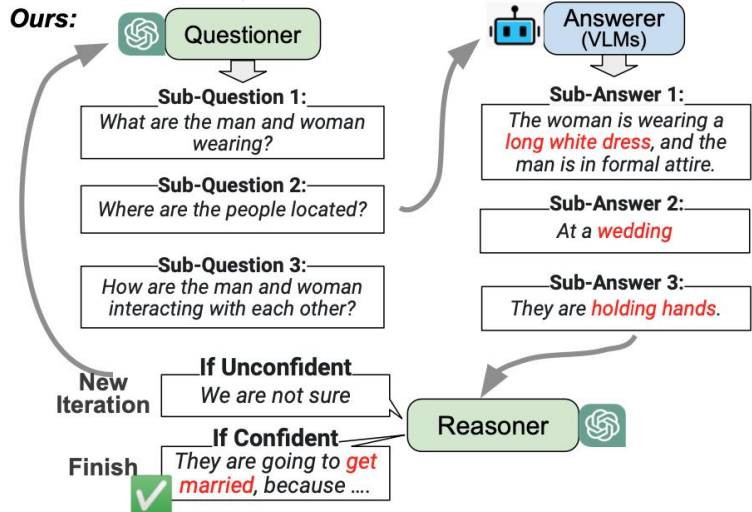
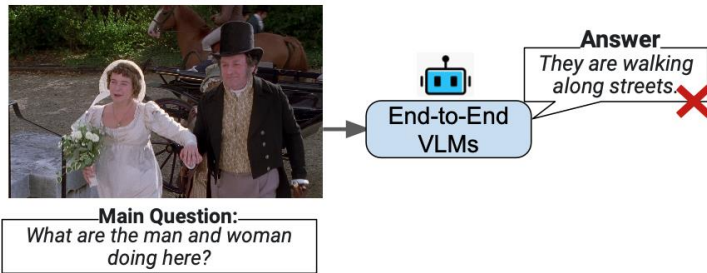
Disclaimer: This paper contains discriminatory content that may be disturbing to some readers, where meme examples and words are offensive or hateful in nature. These contents are provided for illustrative purposes only and do not represent the views and standpoints of the authors.

Hongzhan Lin, et al. Beneath the Surface: Unveiling Harmful Memes with Multimodal Reasoning Distilled from Large Language Models. In EMNLP Findings, 2023.

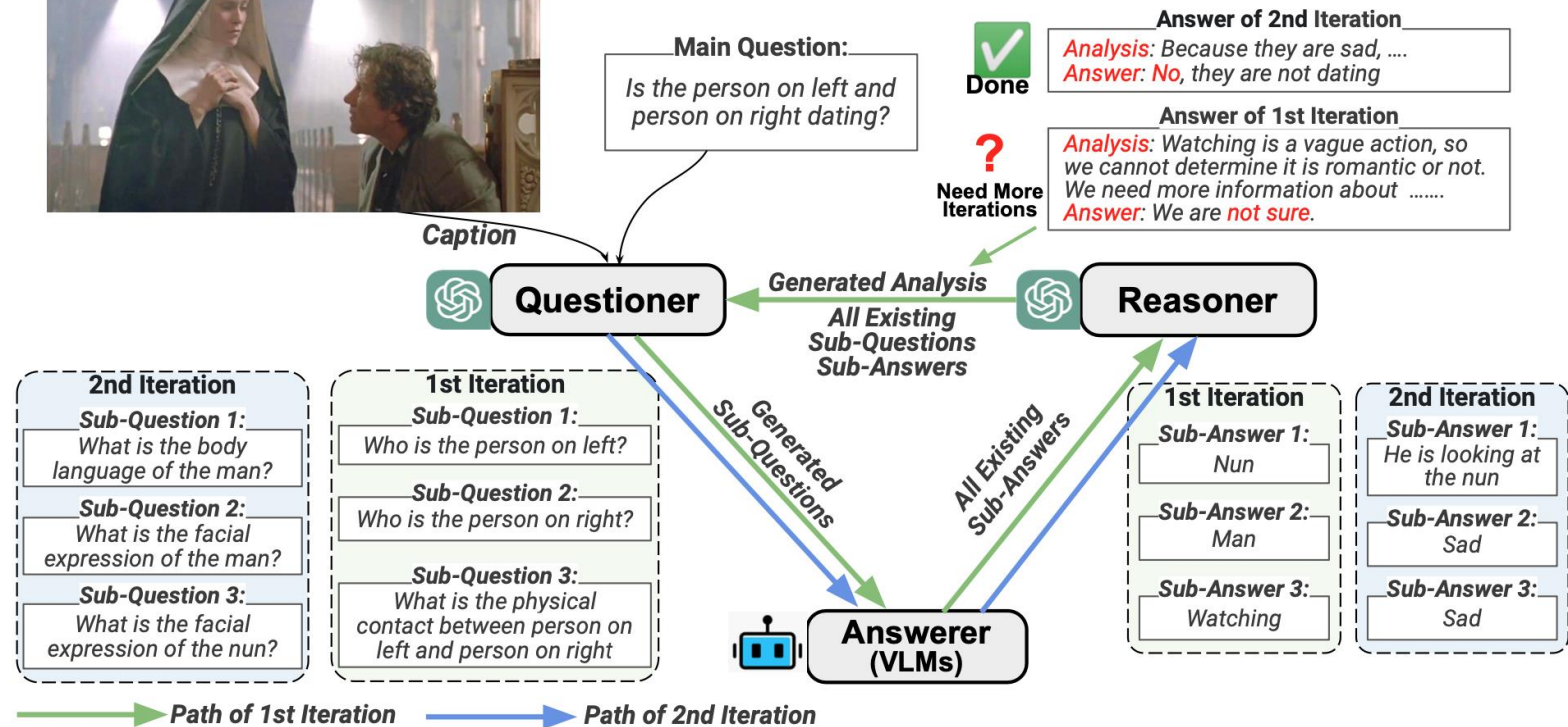
LLMs for “Complex” Visual Question Answering

- IdealGPT

End-to-End Methods:



Caption



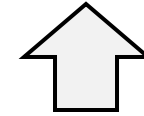
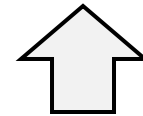
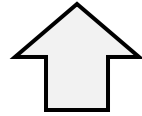
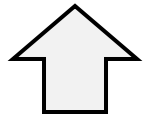
Haoxuan You, et al. IdealGPT: Iteratively Decomposing Vision and Language Reasoning via Large Language Models. In EMNLP Findings, 2023.

Multimodal Understanding & Generation with Efficient Finetune Foundation Models

Open-World Perception

Multimodal Reasoning

Visual Generation & Editing



Parameter-Efficient Tuning

Memory-Efficient Tuning

Modality-Efficient Tuning

Efficiently Finetuned Foundation Models

EVENT-CUSTOMIZED IMAGE GENERATION

Zhen Wang¹, Yilei Jiang¹, Dong Zheng¹, Jun Xiao¹, Long Chen^{2*}

¹Zhejiang University, ²The Hong Kong University of Science and Technology

CLIPDRAG: COMBINING TEXT-BASED AND DRAG-BASED INSTRUCTIONS FOR IMAGE EDITING

Ziqi Jiang, Zhen Wang, Long Chen[†]

The Hong Kong University of Science and Technology

{zjiangbl, zwangjr}@connect.ust.hk, longchen@ust.hk

DISPOSE: DISENTANGLING POSE GUIDANCE FOR CONTROLLABLE HUMAN IMAGE ANIMATION

Hongxiang Li¹, Yaowei Li¹, Yuhang Yang², Junjie Cao³, Zhihong Zhu¹, Xuxin Cheng¹, Long Chen⁴

¹Peking University ²University of Science and Technology of China

³Tsinghua University ⁴Hong Kong University of Science and Technology

View-Consistent 3D Editing with Gaussian Splatting

Yuxuan Wang^{1*}, Xuanyu Yi^{1,2*}, Zike Wu¹, Na Zhao³, Long Chen⁵, and Hanwang Zhang^{1,4}

¹Nanyang Technological University ²Institute for Infocomm Research, A*STAR

³Singapore University of Technology and Design ⁴Skywork AI

⁵Hong Kong University of Science and Technology

Nautilus: Locality-aware Autoencoder for Scalable Mesh Generation

YUXUAN WANG*, Nanyang Technological University, Singapore

XUANYU YI*, Nanyang Technological University, Singapore

HAOHAN WENG*, Tencent Hunyuan, China

QINGSHAN XU, Nanyang Technological University, Singapore

XIAOKANG WEI, The Hong Kong Polytechnic University, Hong Kong, China

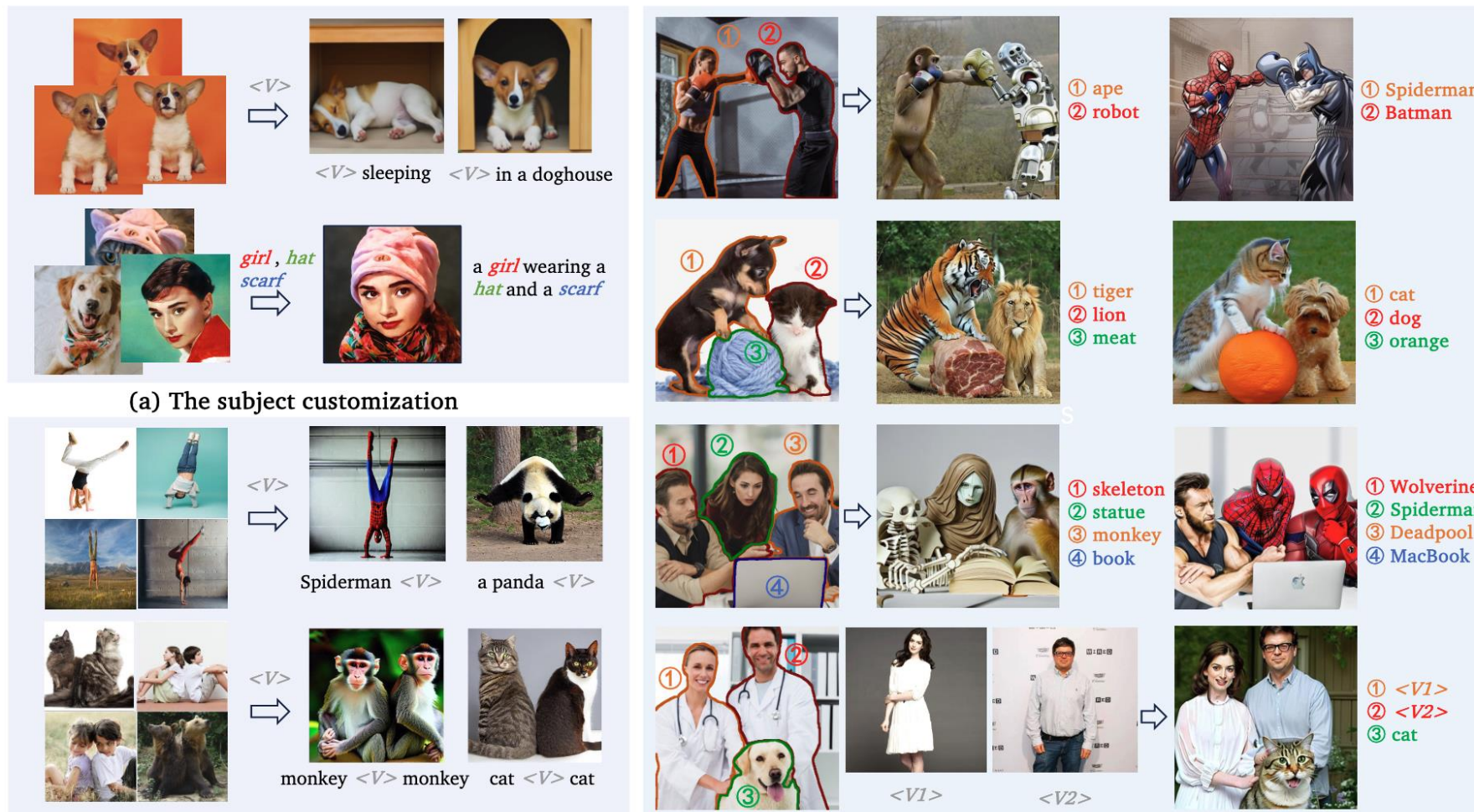
XIANGHUI YANG, Tencent Hunyuan, China

CHUNCHAO GUO, Tencent Hunyuan, China

LONG CHEN, Hong Kong University of Science and Technology, Hong Kong, China

HANWANG ZHANG, Nanyang Technological University, Singapore

Event Customization

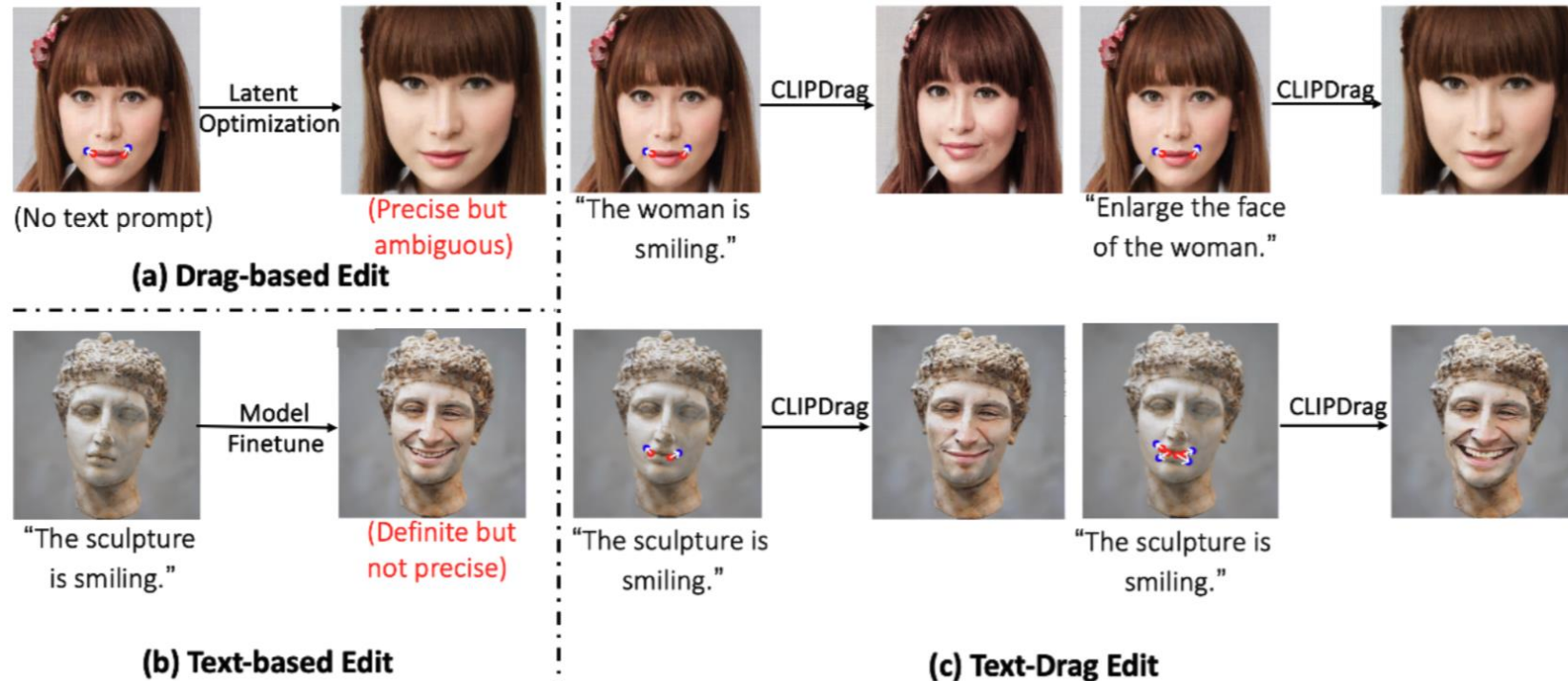


(a) The subject customization

(b) The action and interaction customization

(c) The event customization

Image Editing

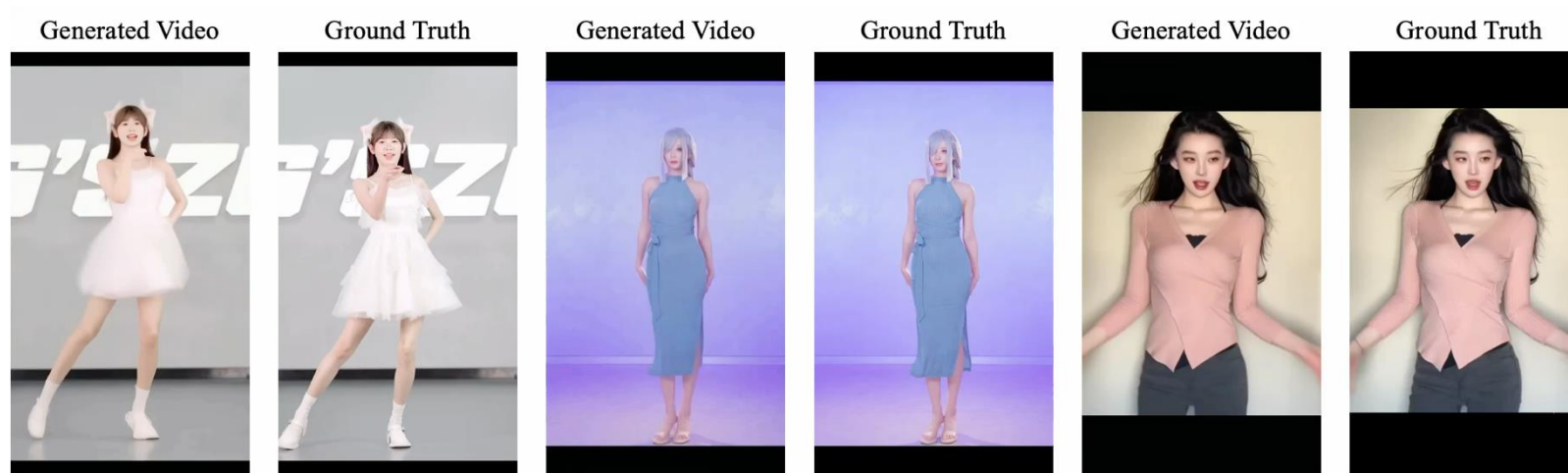


Controllable Video Generation



Github: <https://github.com/lihxxx/DisPose> (300+ stars)

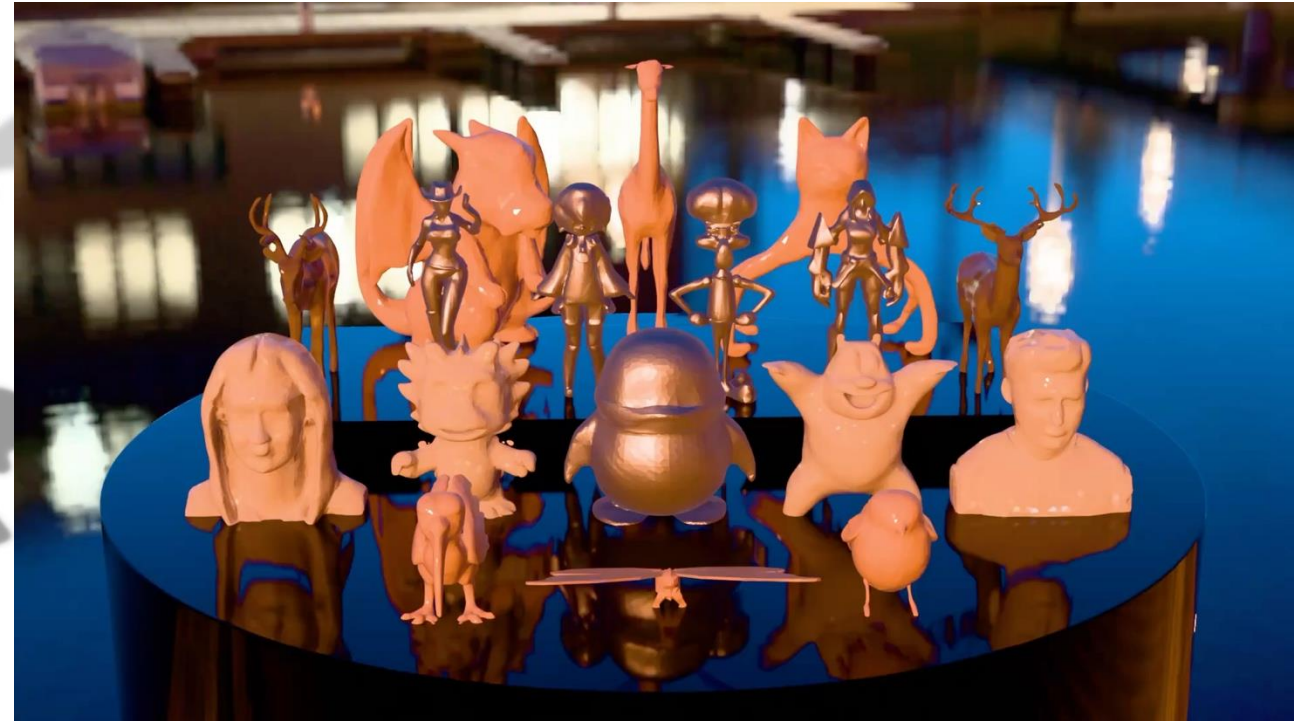
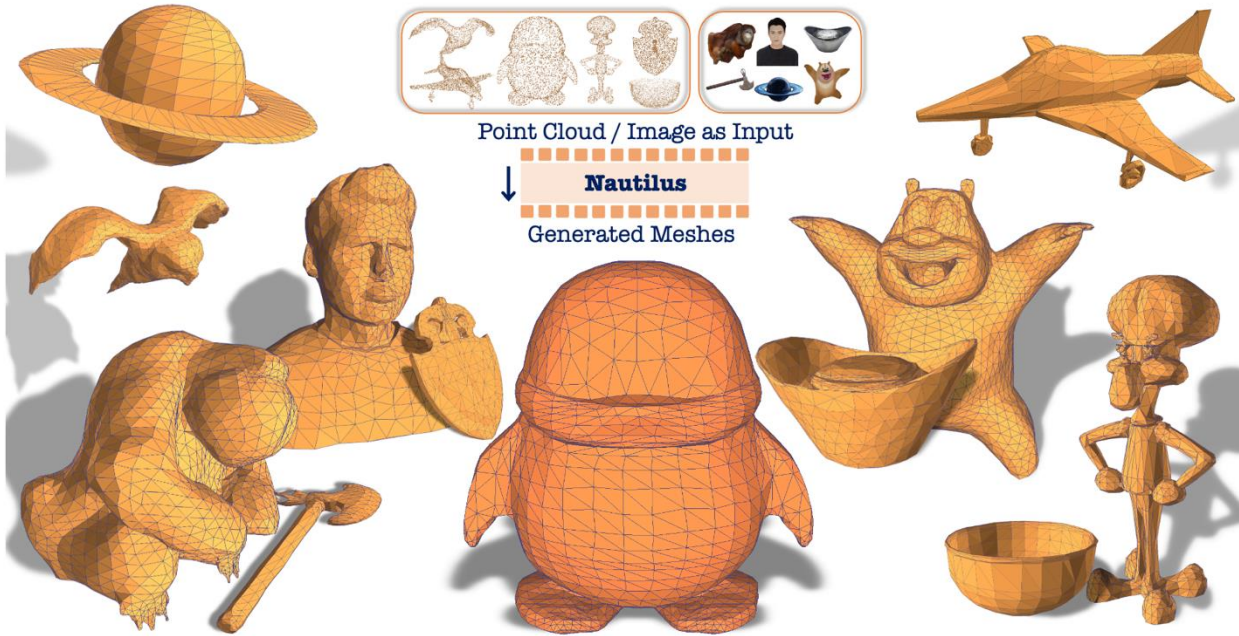
Controllable Video Generation



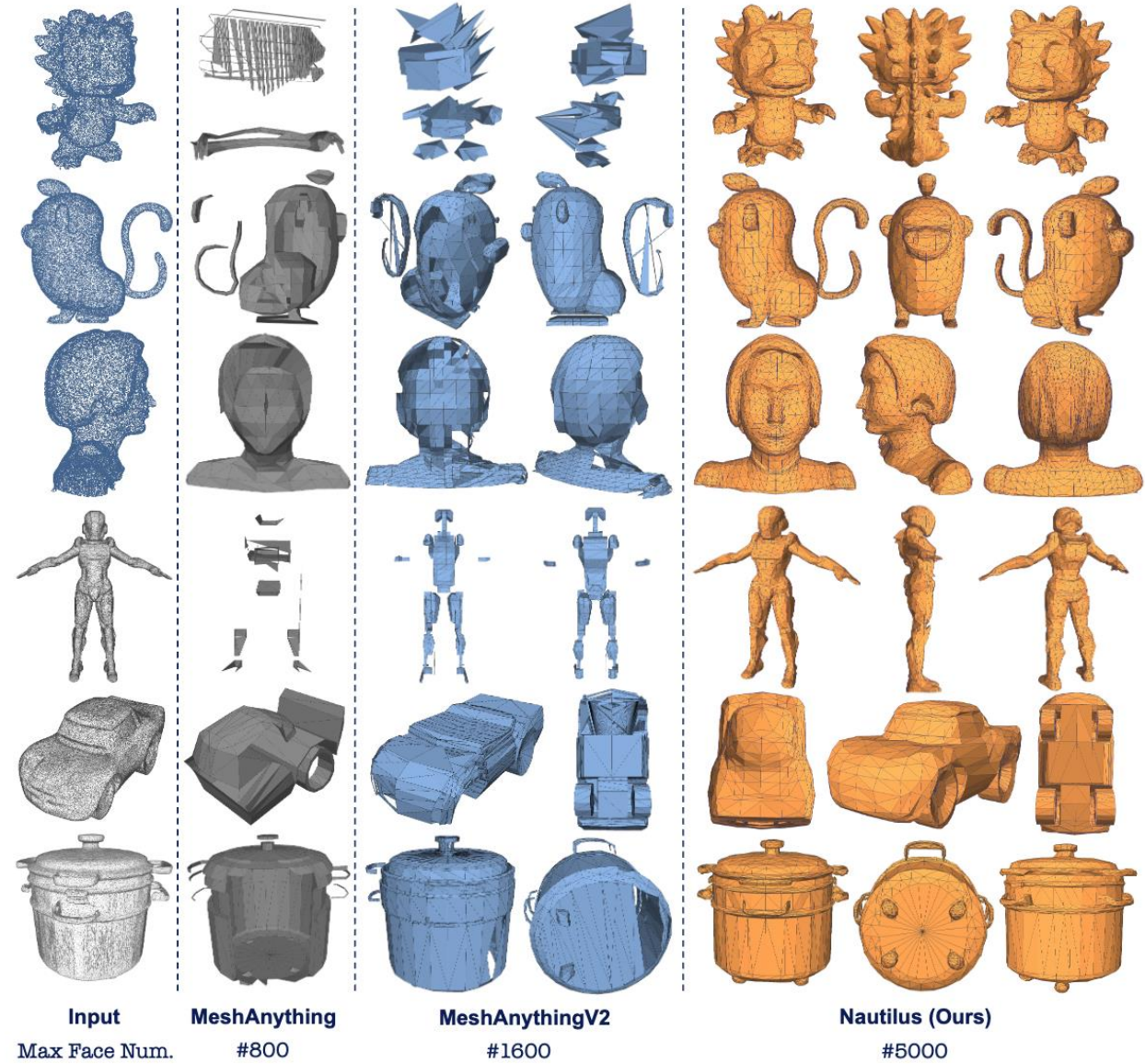
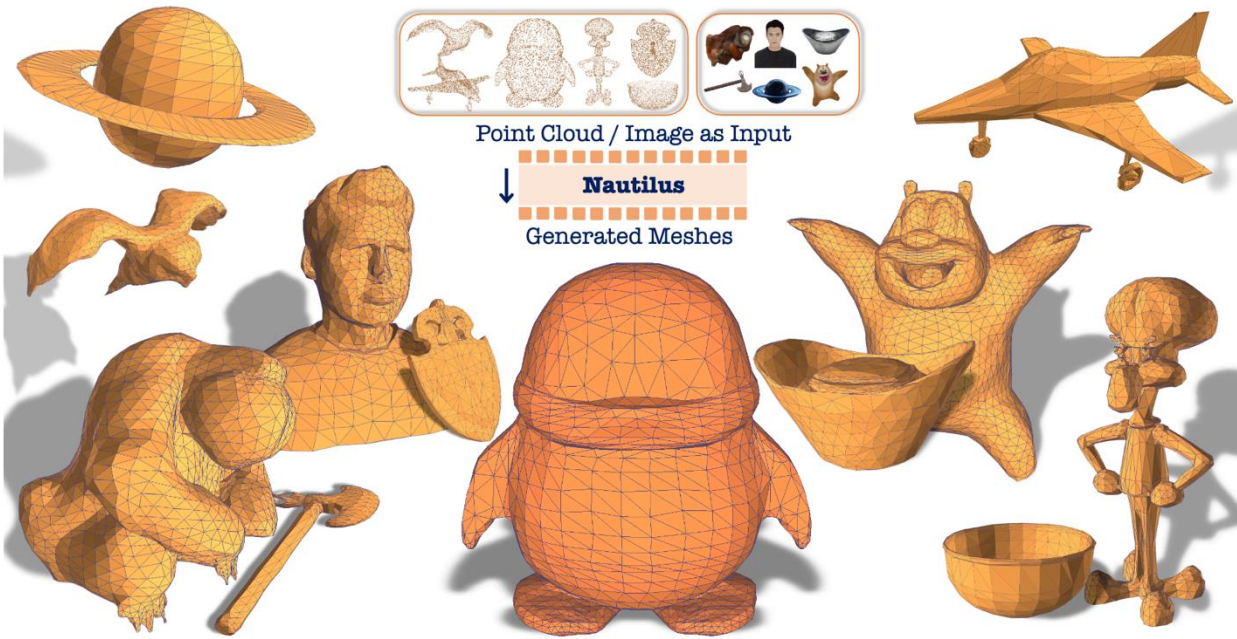
3DGS Editing



3D Mesh Generation



3D Mesh Generation



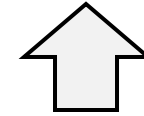
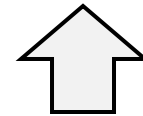
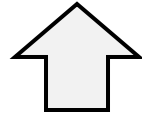
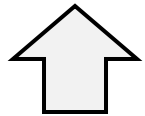
Multimodal Understanding & Generation with Efficient Finetune Foundation Models

Open-World Perception

Multimodal Reasoning

Visual Generation & Editing

Multimodal Generation



Parameter-Efficient Tuning

Memory-Efficient Tuning

Modality-Efficient Tuning

Efficiently Finetuned Foundation Models

Wei Chen^{1,2*}, Lin Li^{1*}, Yongqi Yang^{2*},
Bin Wen³, Fan Yang³, Tingting Gao³, Yu Wu^{2†}, Long Chen^{1†}
¹The Hong Kong University of Science and Technology,
²Wuhan University, ³Kuaishou Technology

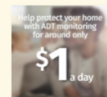
• CoMM: A New Multimodal Benchmark

(c) Sample from MMC4

When it comes to making life just a little bit smoother and safer in a wonderful place such as White Bear Lake, ADT Monitored Home Security can be your main solution for keeping your family and belongings safe and secure.



An ADT home monitoring system for your home can change your life for the better, letting you explore and enjoy all that you love in White Bear Lake. ADT home monitoring is more than just home security, it raises the bar to better living.



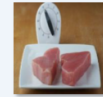
It's a chance to live your daily life without all the worry and stress. For over a century, ADT monitoring systems have paved the way in home security.



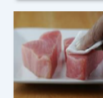
As the #1 home security provider in the country, we know all the vital facts and components that are necessary in keeping our 6+ million home and business customers safe. It's why so many people in White Bear Lake and the rest of the nation count on us for their safety.

(d) Sample from CoMM

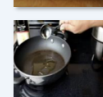
Step1. Remove the tuna steaks from the refrigerator. Allow them to sit out at room temperature for 1/2 to 1 hour before you begin cooking.



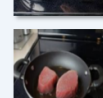
Step2. Use a paper towel to dry excess moisture from the tuna steaks. Simply pat the paper towel on both sides of the steak to remove marinade or condensation from the meat.



Step3. Add 2 tbsp. olive oil or butter to the pan over medium-high heat. Tilt the pan to spread the oil or butter across the surface and coat the pan.



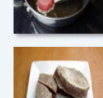
Step4. Place the tuna steaks on the hot pan. Do not crowd the tuna in the pan. Keep the steaks at least 3 to 4 inches apart to make sure they cook evenly.



Step5. Sear each side for two minutes. Do not move the meat around too much while it sears. Gently shake the pan after two minutes, and if the tuna steak releases from the surface, you can flip it on its other side. It is important to make sure that the edged are seared, but you want the inside of the tuna steak to remain rare.



Step6. Make a small cut in the center of the tuna steak with a knife. Inspect the inside to make sure that the middle of the steak stays rare and red in color.



Higher Quality Multimodal Data

- CoMM

(a) Query from CoMM



Here are step-by-step instructions with images about Perfectly Tender and Flaky Pie Crust: (1) Perfectly Tender and Flaky Pie Crust. This pastry recipe uses a combination of lard-butter-vinegar in combination with the usual suspects to produce which is all kinds of tender and flaky..... Let's get started! [Input Image 1]



(2) Tools and Materials. Ingredients: 2 1/2 cups all-purpose flour 1/2 teaspoon salt 1 to 2 tablespoons sugar 6 ounces cold butter..... A pie shell. Plastic wrap and parchment paper. [Input Image 2]



According to the above steps, can you generate the rest steps?



(2) For a specialty meal (1/2 cups) or to help me with lard, you may call 800-800-2750 to make this creamy and creamy pastry recipe. You may also call 800-800-10 cups as we are not required to have this recipe in the next couple of weeks, but that's all for us! Our freezer and veggie recipe is on sale in the United States and Canadian and we currently sell lard at other chefs stores, but for the ultimate savory and lighter dessert!



MMC4



(3) **Mix Ingredients.** Place the dough in a bowl and add some vinegar..... Press firmly.



(4) **Melt the Dough.** Dry the dough in the pan over low flame..... This will prevent the dough from sticking to the pan when you press it onto the bottom of the dish.



(5) **Cook.** Put the dough in the oven for about 45 minutes or until it looks golden.....This allows the filling to stick to the dough without being stretchy.



(6) **Finishing the Dough.** Put the dough in a bowl and set aside..... If the dough is very soft like I did it, you can use the bottom of your oven baking sheet instead.



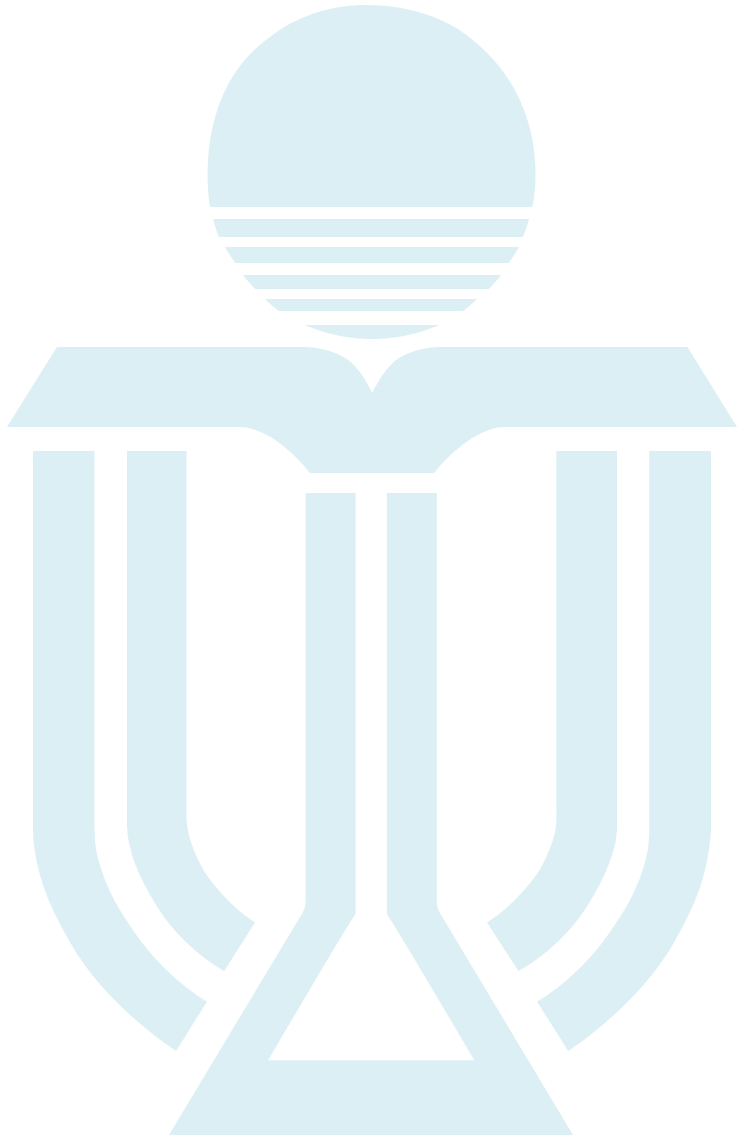
(7) **Decorating It.** With the dough still warm, you can decorate it with whatever decorating tools you want.....Make it more durable by placing a decorative piece of cake and wrapping it with a ribbon.



CoMM

Reference

- *LLMs can Evolve Continually on Modality for X-Modal Reasoning. In NeurIPS, 2024.*
- *Inversion Circle Interpolation: Diffusion-based Image Augmentation for Data-scarce Classification. In arXiv, 2024.*
- *Zero-Shot Visual Relation Detection via Composite Visual Cues from Large Language Models. In NeurIPS, 2023.*
- *Beneath the Surface: Unveiling Harmful Memes with Multimodal Reasoning Distilled from Large Language Models. In EMNLP Findings, 2023.*
- *IdealGPT: Iteratively Decomposing Vision and Language Reasoning via Large Language Models. In EMNLP Findings, 2023.*
- *Event-Customized Image Generation. In arXiv, 2024.*
- *CLIPDrag: Combining Text-based and Drag-based Instructions for Image Editing. In ICLR, 2025.*
- *DisPose: Disentangling Pose Guidance for Controllable Human Image Animation. In ICLR, 2025.*
- *Nautilus: Locality-aware Autoencoder for Scalable Mesh Generation. In arXiv, 2025.*
- *View-Consistent 3D Editing with Gaussian Splatting. In ECCV, 2024.*
- *CoMM: A Coherent Interleaved Image-Text Dataset for Multimodal Understanding and Generation. In arXiv, 2024.*



Thanks & QA!

In Collaboration with:

HKUST:

Yanghao, Wei, Lin, Zhen, Ziqi, Hongxiang

Other universities:

*Yuxuan (NTU), Haoxuan (Columbia), Hongzhan (HKBU),
Jiazu (DLUT)*