

# Deep Learning for Computer Vision

113-1/Fall 2024

<https://cool.ntu.edu.tw/courses/41702> (NTU COOL)

<http://vllab.ee.ntu.edu.tw/dlcv.html> (Public website)

Yu-Chiang Frank Wang 王鈺強, Professor

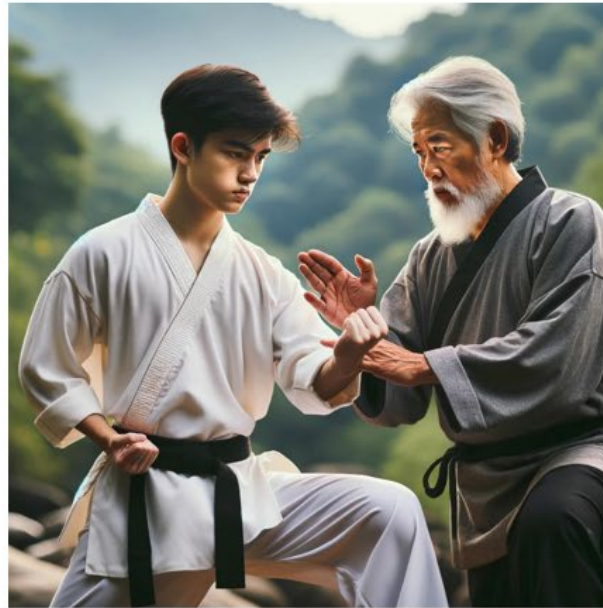
Dept. Electrical Engineering, National Taiwan University

# Pretrain & Finetune LLM/VLM/MLLM



Stage 1

Pre-training by self-supervised learning or supervised learning



Stage 2

Finetuning by downstream tasks in target domains

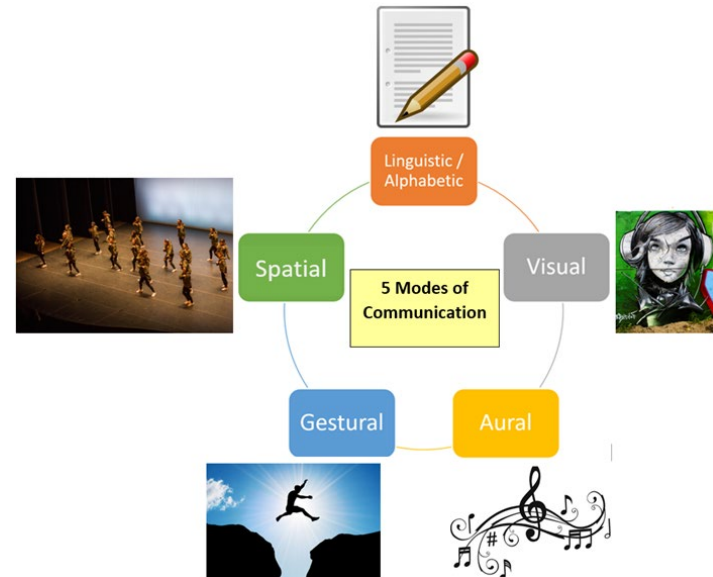
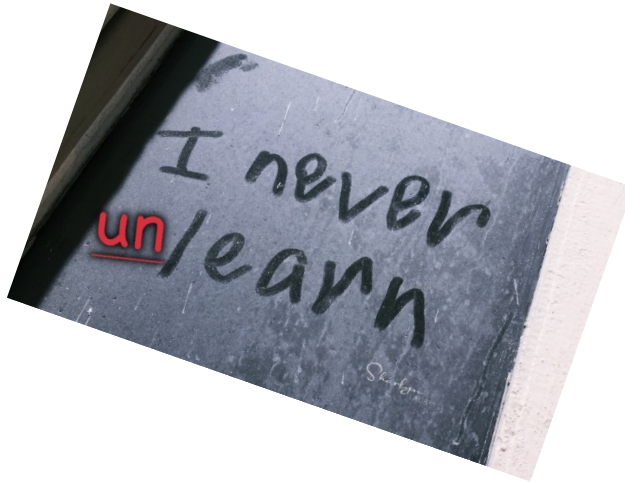
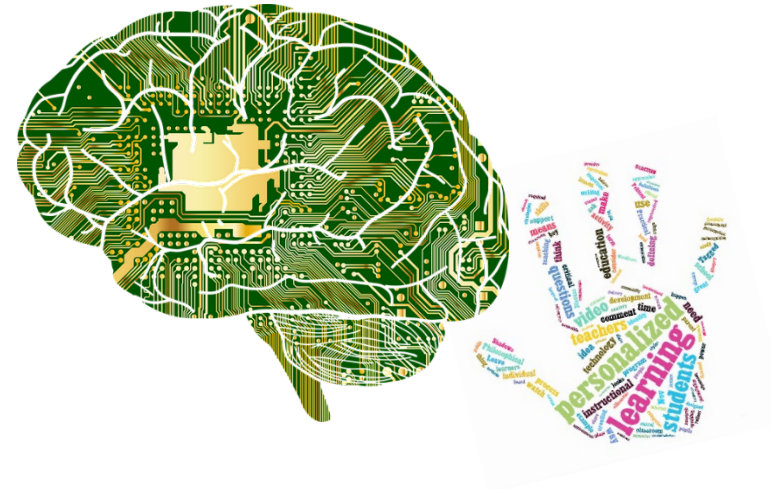


Stage 3

RLHF - Reinforcement Learning with Human Feedback  
(not covered)

# What to Be Covered Today...

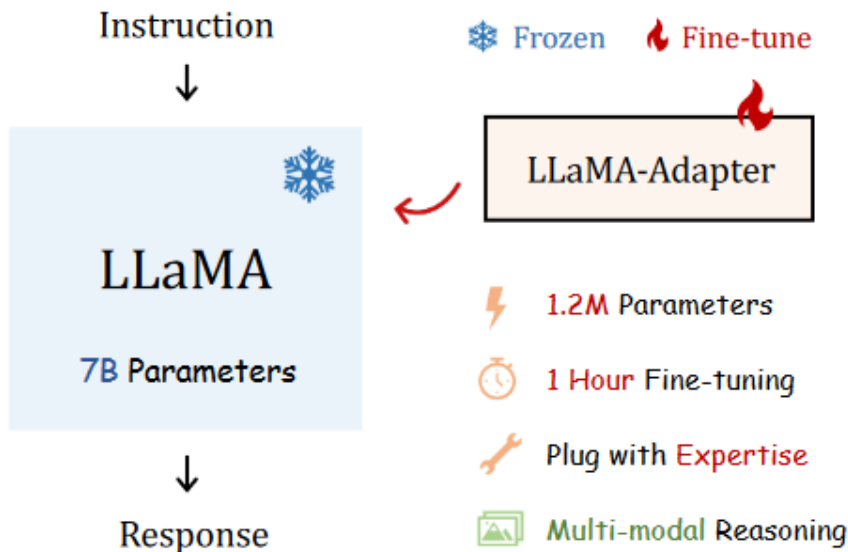
- **Multimodal LLM**
- **Advanced Topics in LLM/VLM**
  - Concept Editing
  - Concept Unlearning
  - Personalization
  - Continual Learning



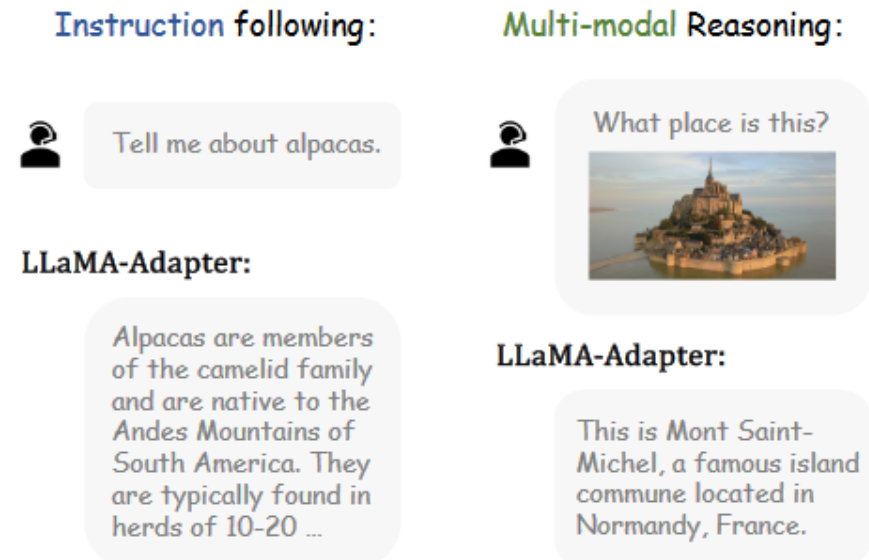
# LLaMA-Adapter (SAIL & CUHK, ICLR'24)

- A **lightweight** adaptation method to efficiently finetune LLaMA into an **multi-modal instruction-following** model
- **Result:**
  1. Equipping LLaMA the ability of understanding instruction-following data
  2. Capable of addressing multi-modal reasoning tasks

## Training



## Inference



# LLaMa-Adapter: *Method*

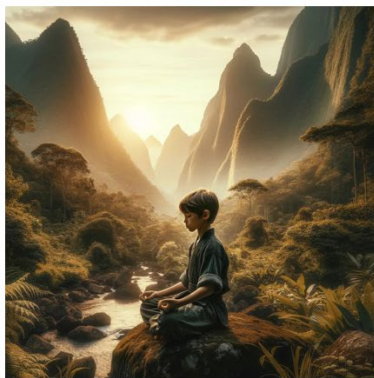
- Append K learnable prompts at each of the last L layers
- Adopt **zero-initialized attention**

for fine-tuning stability and effectiveness

- use a zero-initialized learnable gating factor  $g_l$  to control the prompts' attention scores

$$S_l = Q_l K_l^T / \sqrt{C} \in \mathbb{R}^{1 \times (K+M+1)}$$

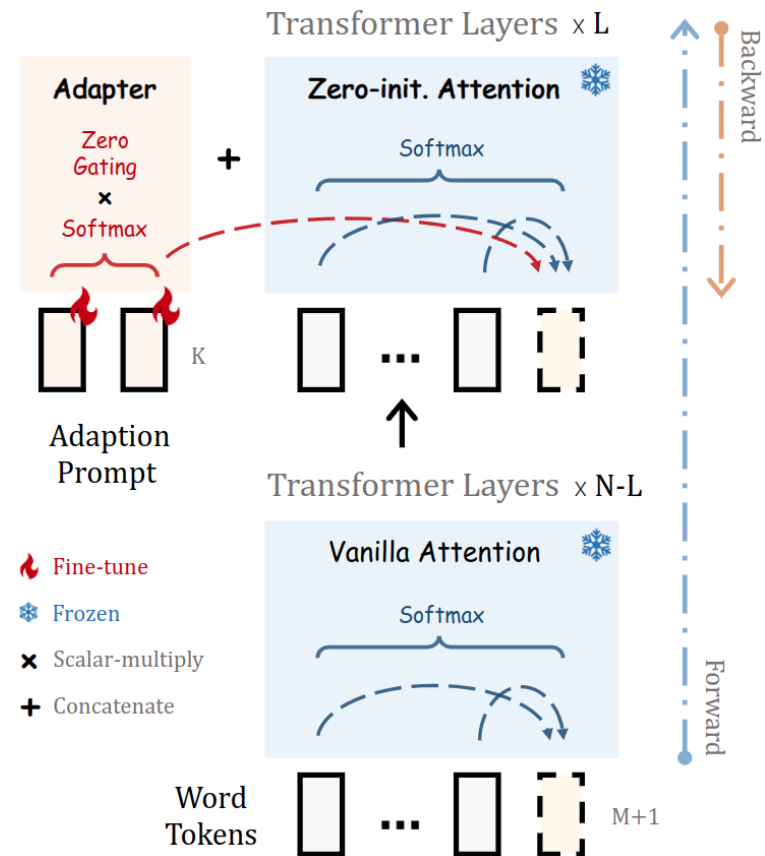
$$S_l^g = [\text{softmax}(S_l^K) \cdot g_l; \text{softmax}(S_l^{M+1})]^T$$



Stage 1

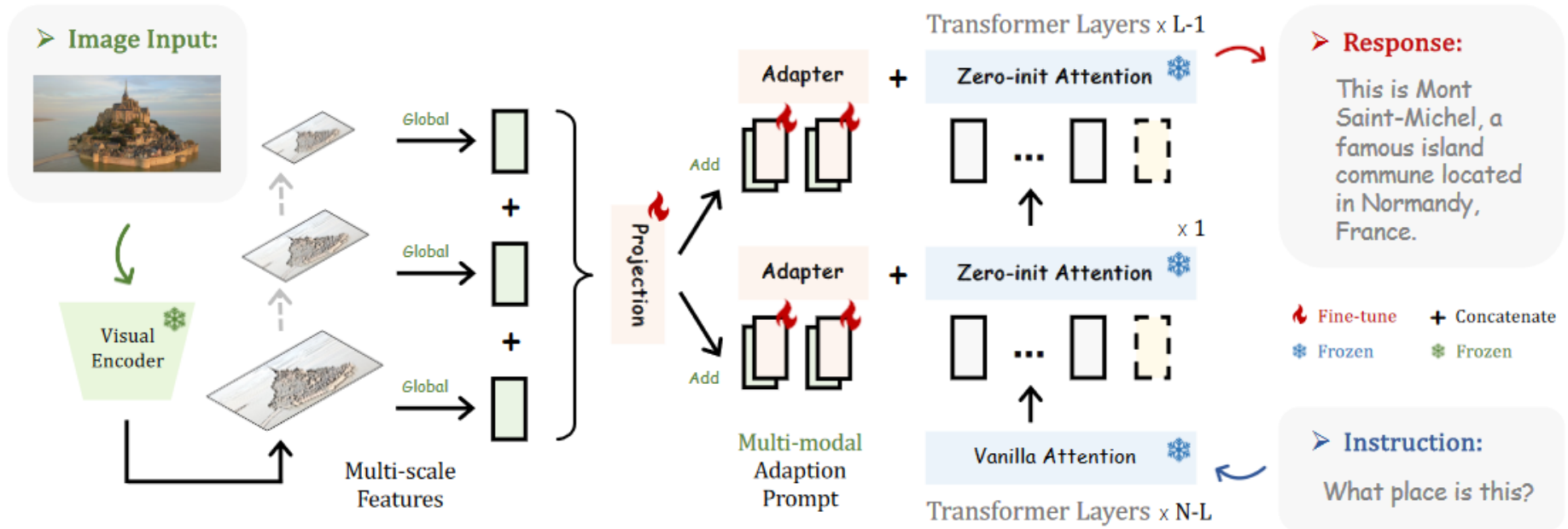


Stage 2



# Multi-Modal LLaMA-Adaptor

- When being finetuned to a multi-modal reasoning task (e.g. VQA), a pre-trained visual encoder + a learnable projection layer are additionally utilized



# Quantitative Result

- Methods comparison on the visual question answering (VQA) task (ScienceQA dataset)

Model	Tuned Params	Avg	NAT	SOC	LAN	TXT	IMG	NO	G1-6	G7-12
Random Choice [41]	-	39.83	40.28	46.13	29.25	47.45	40.08	33.66	39.35	40.67
Human [41]	-	88.40	90.23	84.97	87.48	89.60	87.50	88.10	91.59	82.42
MCAN [65]	95M	54.54	56.08	46.23	58.09	59.43	51.17	55.40	51.65	59.72
VisualBERT [33, 34]	111M	61.87	59.33	69.18	61.18	62.71	62.17	58.54	62.96	59.92
UnifiedQA [27]	223M	70.12	68.16	69.18	74.91	63.78	61.38	77.84	72.98	65.00
UnifiedQA <sub>C<sub>o</sub>T</sub>	223M	74.11	71.00	76.04	78.91	66.42	66.53	81.81	77.06	68.82
GPT-3 [4]	0M	74.04	75.04	66.59	78.00	74.24	65.74	79.58	76.36	69.87
GPT-3 <sub>C<sub>o</sub>T</sub>	0M	75.17	75.44	70.87	78.09	74.68	67.43	79.93	78.23	69.68
ChatGPT <sub>C<sub>o</sub>T</sub> [2]	0M	78.31	78.82	70.98	83.18	77.37	67.92	86.13	80.72	74.03
GPT-4 <sub>C<sub>o</sub>T</sub> [45]	0M	83.99	85.48	72.44	90.27	82.65	71.49	92.89	86.66	79.04
MM-COT <sub>T</sub> [74]	223M	70.53	71.09	70.75	69.18	71.16	65.84	71.57	71.00	69.68
MM-COT	223M	84.91	87.52	77.17	85.82	87.88	82.90	86.83	84.65	85.37
<b>LLaMA-Adapter<sub>T</sub></b>	<b>1.2M</b>	<b>78.31</b>	<b>79.00</b>	<b>73.79</b>	<b>80.55</b>	<b>78.30</b>	<b>70.35</b>	<b>83.14</b>	<b>79.77</b>	<b>75.68</b>
<b>LLaMA-Adapter</b>	<b>1.8M</b>	<b>85.19</b>	<b>84.37</b>	<b>88.30</b>	<b>84.36</b>	<b>83.72</b>	<b>80.32</b>	<b>86.90</b>	<b>85.83</b>	<b>84.05</b>

metrics:  
accuracy (%)

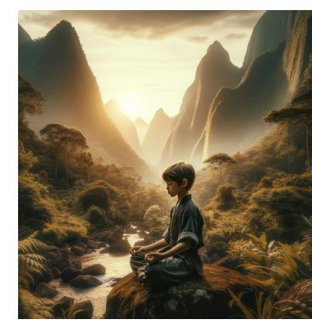
# Easy Visual Sound Localization (EZ-VSL), CMU & UWisc., ECCV 2022

- **Goal:**

A simple yet effective method to **unsupervised audio-visual sound localization**

- Task: localize visible sound sources in a video without ground-truth localization

- **Result:** SOTA performance on two popular benchmarks, Flickr SoundNet & VGG-Sound Source

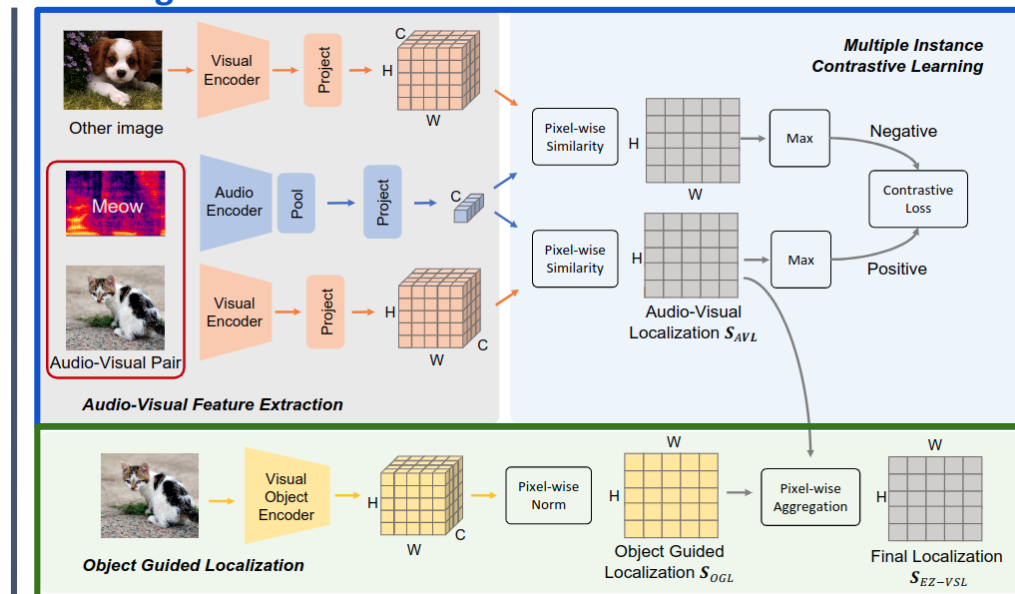


Stage 1



Stage 2

## Training



## Inference

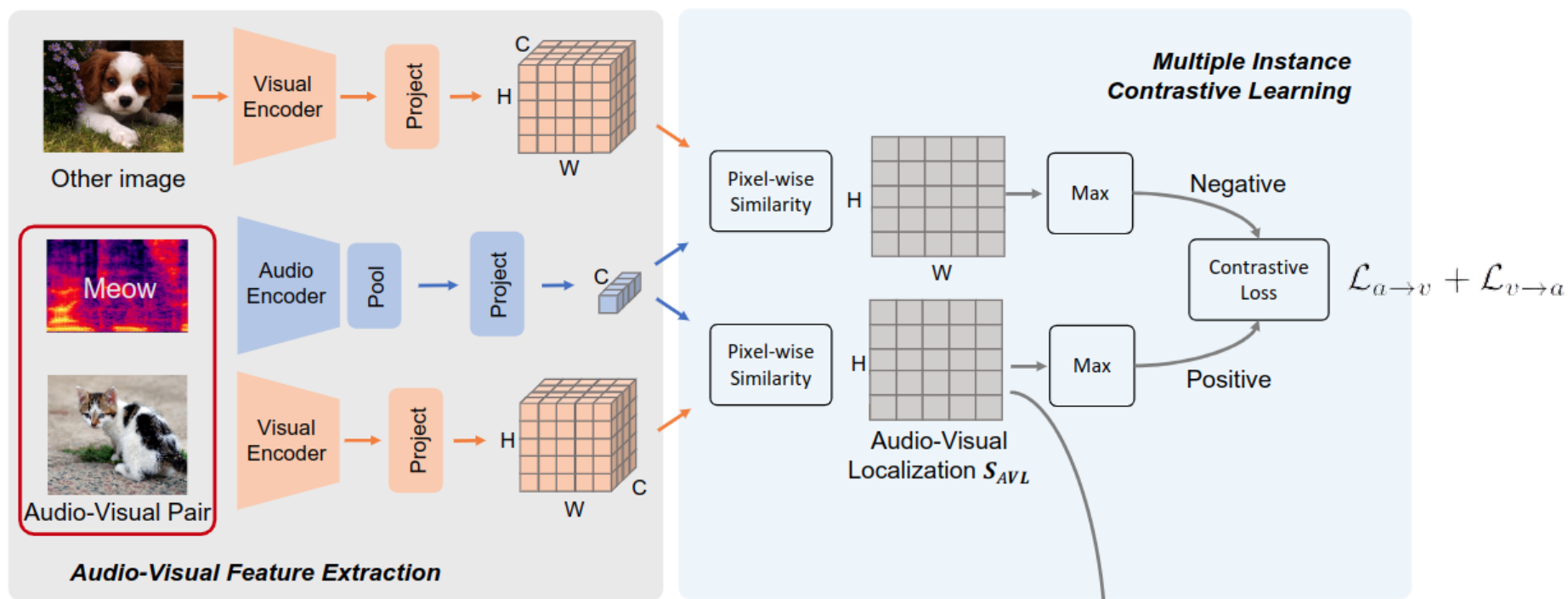
(a) audio-visual sound localization examples

(b) EZ-VSL framework



# Training Stage

- **Audio-Visual Matching by Multiple-Instance Contrastive Learning**
  - Encourage audio representation to be aligned with the associate visual representations **at least at one location** and not being associated with any locations from other images

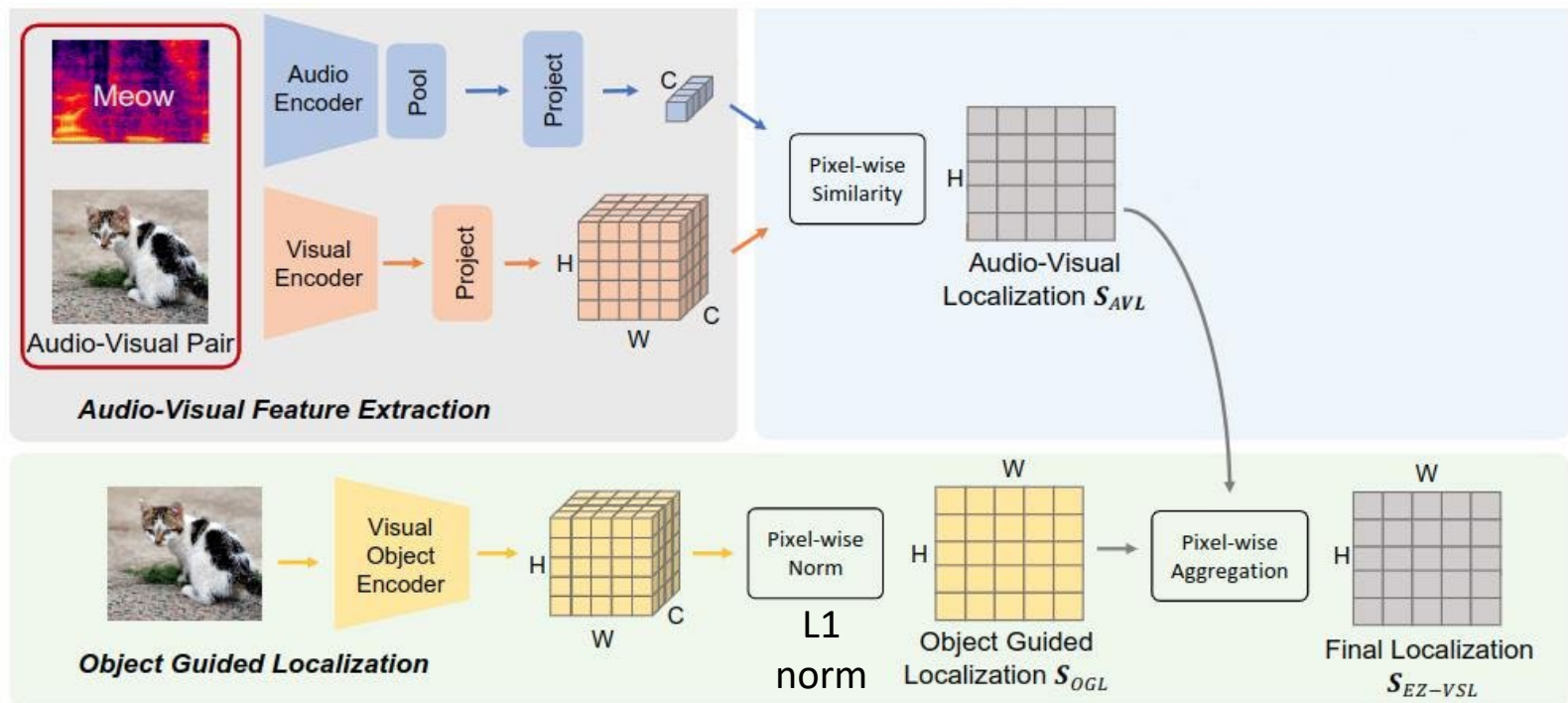


$$\mathcal{L}_{a \rightarrow v} = -\log \frac{\exp\left(\frac{1}{\tau} \max_{\hat{v} \in V_i} \text{sim}(\hat{\mathbf{a}}_i, \hat{\mathbf{v}})\right)}{\sum_k \exp\left(\frac{1}{\tau} \max_{\hat{v} \in V_k} \text{sim}(\hat{\mathbf{a}}_i, \hat{\mathbf{v}})\right)} \quad \mathcal{L}_{v \rightarrow a} = -\log \frac{\exp\left(\frac{1}{\tau} \max_{\hat{v} \in V_i} \text{sim}(\hat{\mathbf{v}}, \hat{\mathbf{a}}_i)\right)}{\sum_k \exp\left(\frac{1}{\tau} \max_{\hat{v} \in V_i} \text{sim}(\hat{\mathbf{v}}, \hat{\mathbf{a}}_k)\right)}$$

# Inference Stage

- **Object-Guided Localization**

- Improve localization precision by combining audio-visual similarity map with an **object localization map** from a pre-trained visual model



$$S_{EZ-VSL} = \alpha S_{AVL} + (1 - \alpha) S_{OGL}$$

# Quantitative Results

- Comparison results on Flickr SoundNet testset

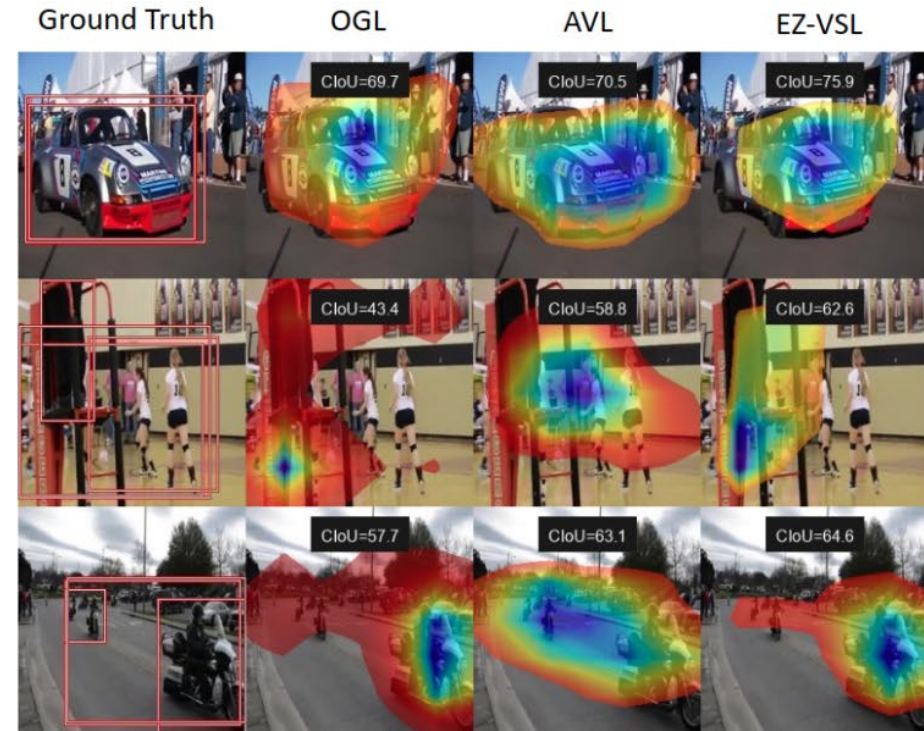
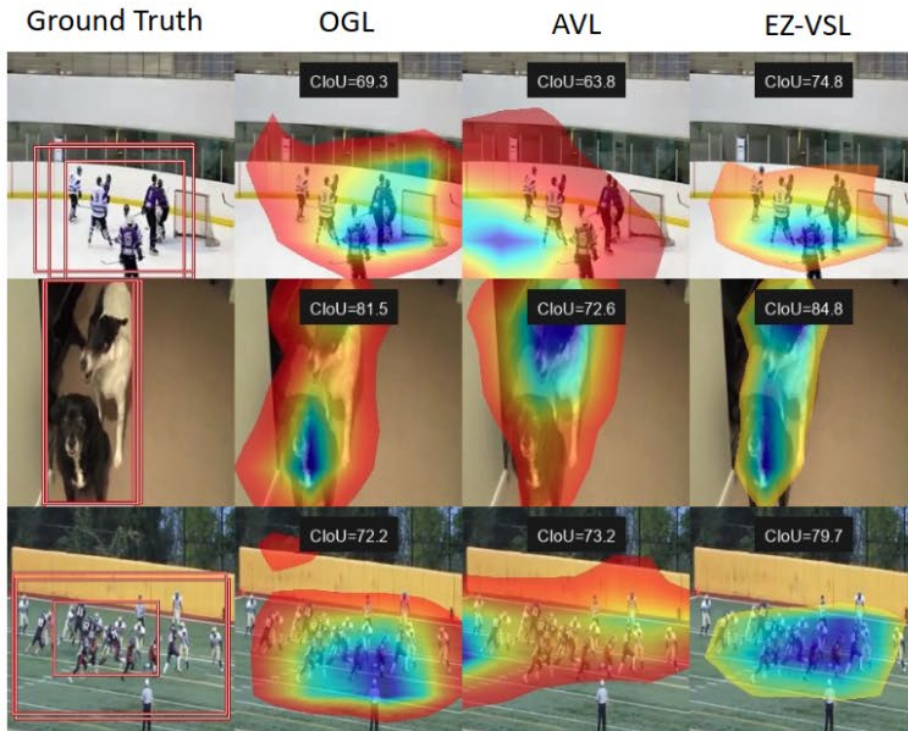
Training set	Method	CIoU(%)	AUC(%)
Flickr 10k	Attention10k [30]	43.60	44.90
	CoarsetoFine [28]	52.20	49.60
	AVObject [2]	54.60	50.40
	LVS [6]	58.20	52.50
	EZ-VSL (ours)	<b>81.93</b>	<b>62.58</b>
Flickr 144k	Attention10k [30]	66.00	55.80
	DMC [17]	67.10	56.80
	LVS [6]	69.90	57.30
	HardPos [31]	75.20	59.70
	EZ-VSL (ours)	<b>83.13</b>	<b>63.06</b>

- Methods comparison on Flickr SoundNet and VGG-SS testset

Training set	Method	Flickr-SoundNet		VGG-SS	
		CIoU(%)	AUC(%)	CIoU(%)	AUC(%)
VGG-Sound 144k	Attention10k [30]	66.00	55.80	18.50	30.20
	CoarsetoFine [28]	-	-	29.10	34.80
	AVObject [2]	-	-	29.70	35.70
	LVS [6]	73.50	59.00	34.40	38.20
	HardPos [31]	76.80	59.20	34.60	38.00
	EZ-VSL (ours)	<b>83.94</b>	<b>63.60</b>	<b>38.85</b>	<b>39.54</b>

# Visualization

- Sound source localization predictions with other methods



# AV-HuBERT, Meta, ICLR'22

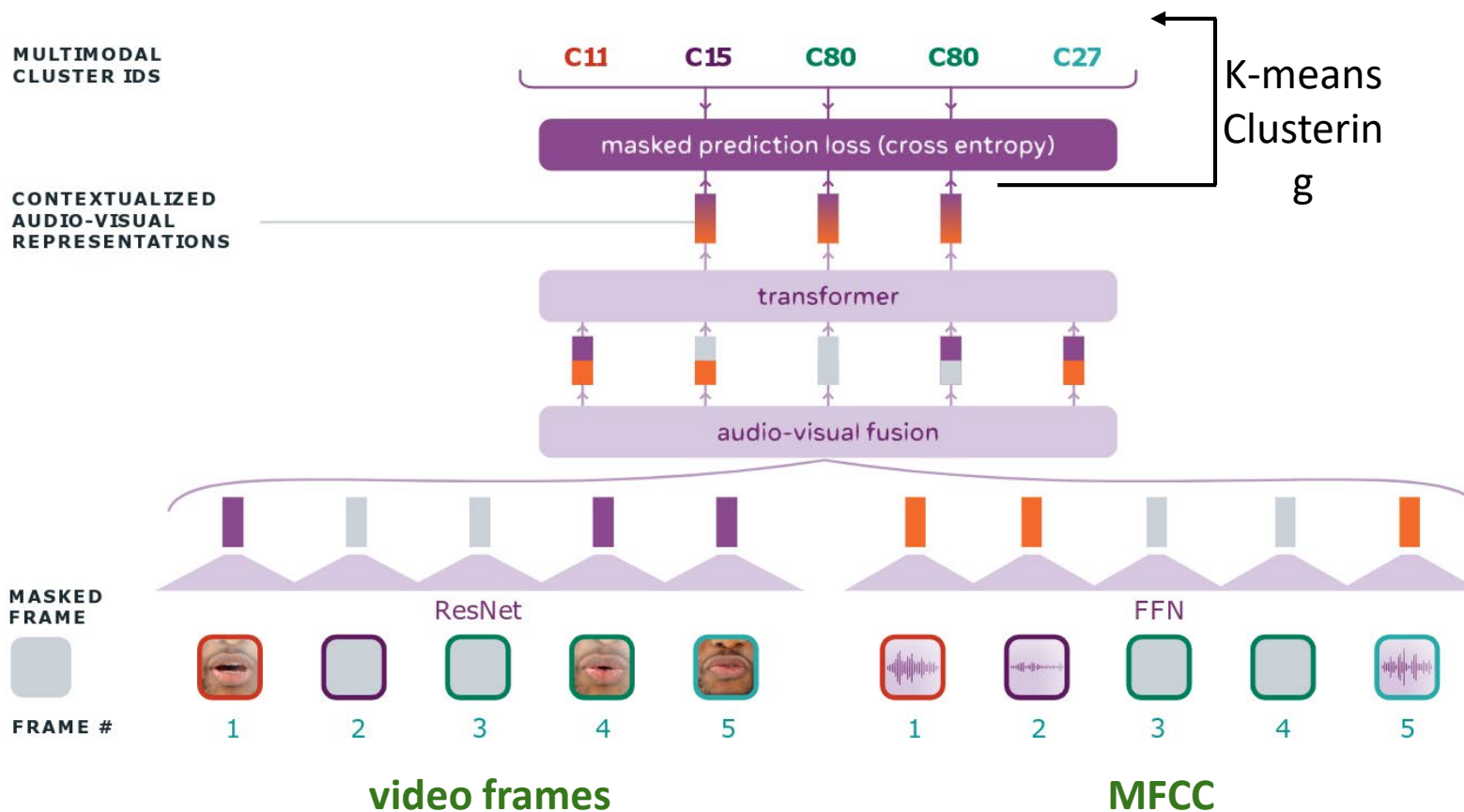
- **Goal:** A self-supervised representation learning framework for audio-visual speech recognition
- **Result:** SOTA performance on the largest downstream lip-reading benchmark



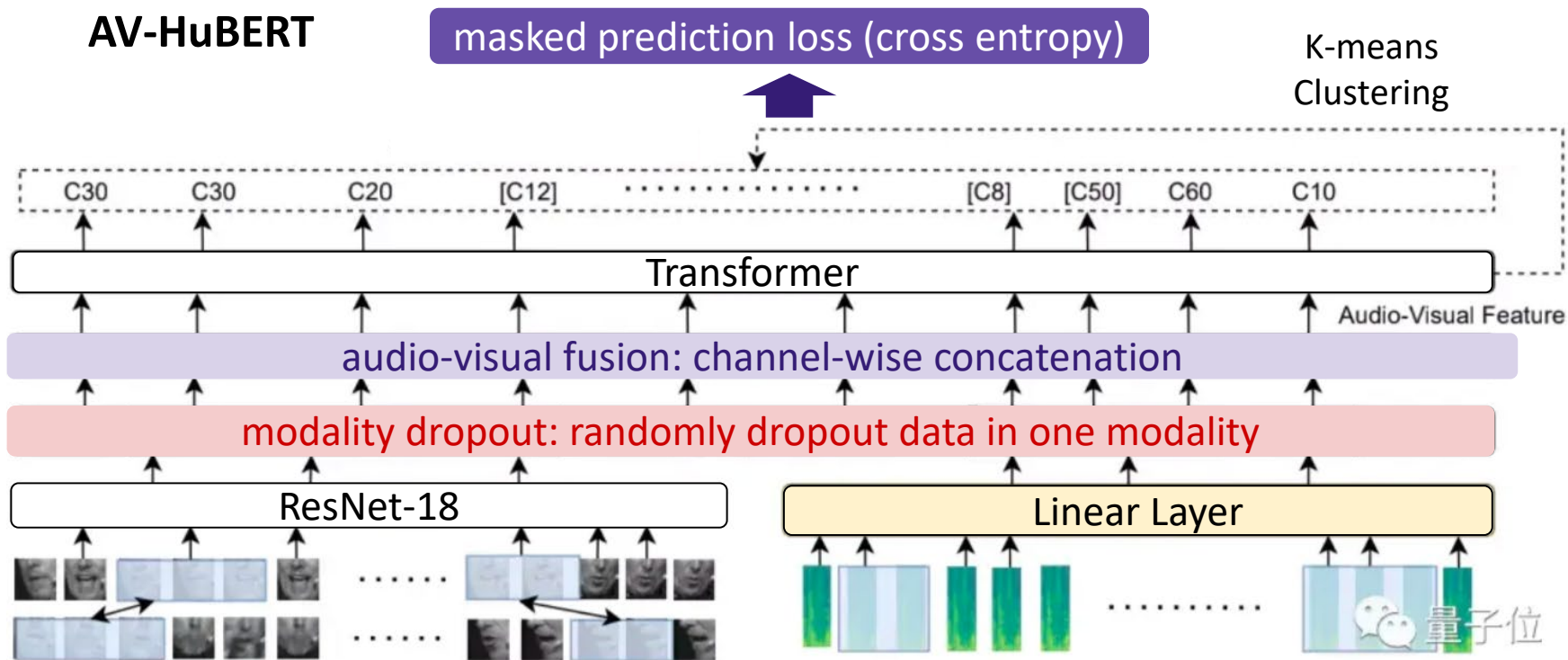
Stage 1



Stage 2



# Pre-training Framework



## Masking by Substitution

**Goal:** enforce the model to learn the temporal relationship between the frames

**Method:** replace the masked segments with random segments from the same video

<https://zhuanlan.zhihu.com/p/455426545>

# Quantitative Result

- Method comparison on lip-reading benchmark (LRS3 dataset)

Method	Backbone	Criterion	Labeled iso (hrs)	Labeled utt (hrs)	Unlabeled data (hrs)	WER (%)
<i>Supervised</i>						
Afouras et al. (2020)	CNN	CTC	157	433	-	68.8
Zhang et al. (2019b)	CNN	S2S	157	698	-	60.1
Afouras et al. (2018a)	Transformer	S2S	157	1,362	-	58.9
Xu et al. (2020)	RNN	S2S	157	433	-	57.8
Shillingford et al. (2019)	RNN	CTC	-	3,886	-	55.1
Ma et al. (2021b)	Conformer	CTC+S2S	-	433	-	46.9
Ma et al. (2021b)	Conformer	CTC+S2S	157	433	-	43.3
Makino et al. (2019)	RNN	Transducer	-	31,000	-	<b>33.6</b>
<i>Semi-Supervised &amp; Self-Supervised</i>						
Afouras et al. (2020)	CNN	CTC	157	433	334	59.8
Ma et al. (2021a)†	Transformer-BASE	S2S	-	30	433	71.9
			-	433	1,759	49.6
<i>Proposed (Self-Supervised &amp; Self-Supervised + Semi-Supervised)</i>						
			-	30	-	94.3
			-	30	433	51.8
	Transformer-BASE	S2S	-	<b>30</b>	1,759	<b>46.1</b>
			-	433	-	60.3
			-	433	433	44.0
AV-HuBERT			-	433	1,759	<b>34.8</b>
			-	30	-	92.3
			-	30	433	44.8
	Transformer-LARGE	S2S	-	30	1,759	<b>32.5</b>
			-	433	-	62.3
			-	433	433	41.6
			-	433	1,759	<b>28.6</b>
AV-HuBERT + Self-Training	Transformer-LARGE	S2S	-	30	1,759	<b>28.6</b>
			-	433	1,759	<b>26.9</b>

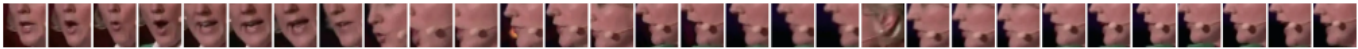

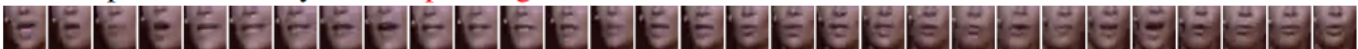
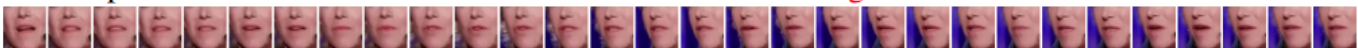
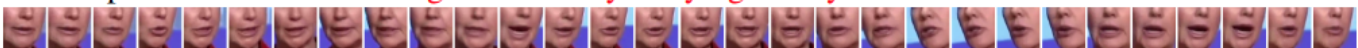
CTC: connectionist temporal classification

S2S: attention-based sequence-to-sequence cross entropy

# Qualitative Result

- Transcriptions comparison between AV-HuBERT (Proposed) and a supervised model

---

(1)	GT:	why not ask all of the states to do that instead
	Proposed:	why not ask all of <b>these things</b> to do that instead
	Supervised:	why <b>can't i actually</b> all of <b>these things</b> do <b>things and</b>
		
(2)	GT:	indeed we run the risk of making things worse
	Proposed:	indeed we <b>want</b> the risk of making things worse
	Supervised:	<b>in india</b> we <b>roughly receive money in the health world</b>
		
(3)	GT:	my desire to disappear was still very powerful
	Proposed:	my desire to disappear was still very powerful
	Supervised:	my <b>son is speaking with children about food</b>
		
(4)	GT:	the silent majority does not need to be silent
	Proposed:	the <b>same</b> majority does not need to be silent
	Supervised:	<b>this time the total disaster needs to be designed</b>
		
(5)	GT:	mortality is not going down it's going up
	Proposed:	mortality is not going down it's going up
	Supervised:	<b>we're seeing this not only carrying slowly how</b>
		

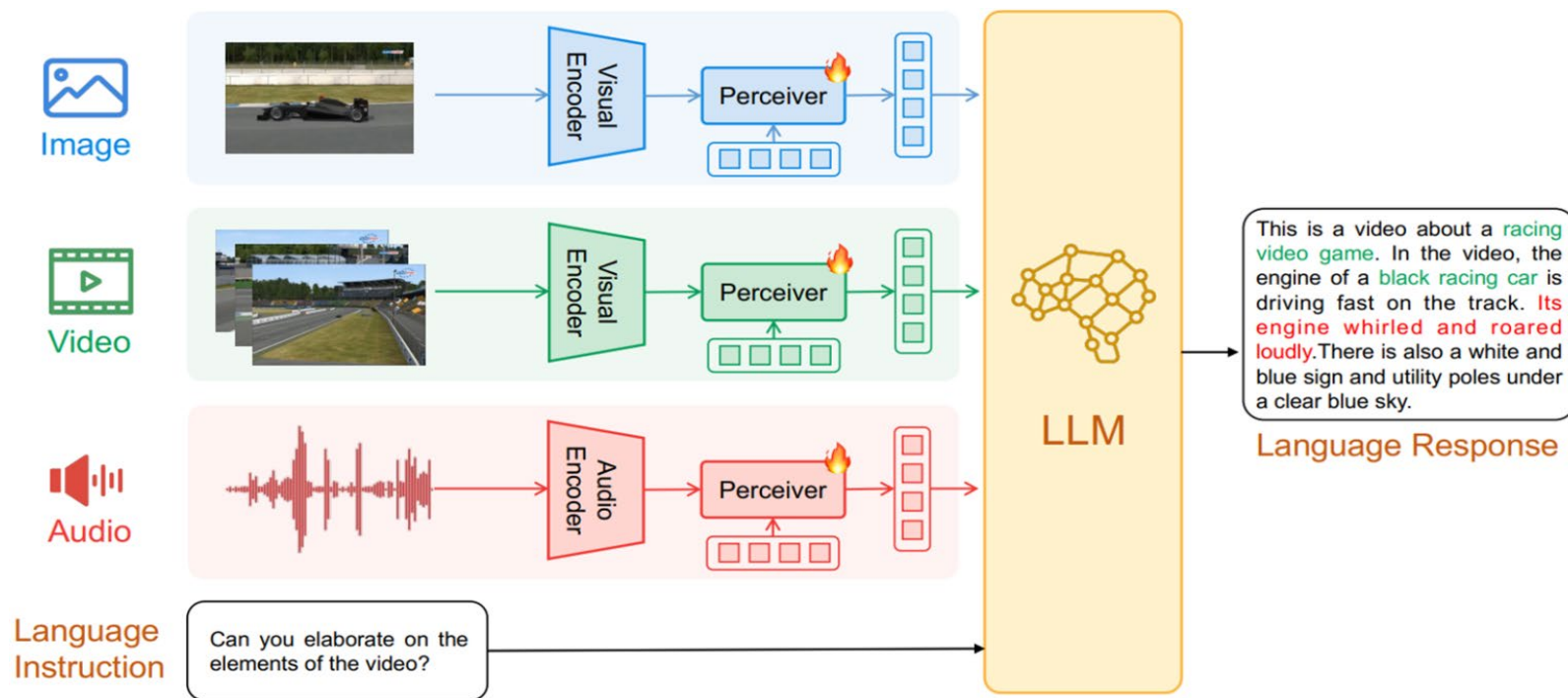
---



# ChatBridge (CAS, arxiv'23)

- **Goal:** Transform a LLM into a **multi-modal (text, image, video, audio) language model** using **language** to bridge the gap between various modalities
- **Result:** Outstanding zero-shot performance on various multi-modal tasks

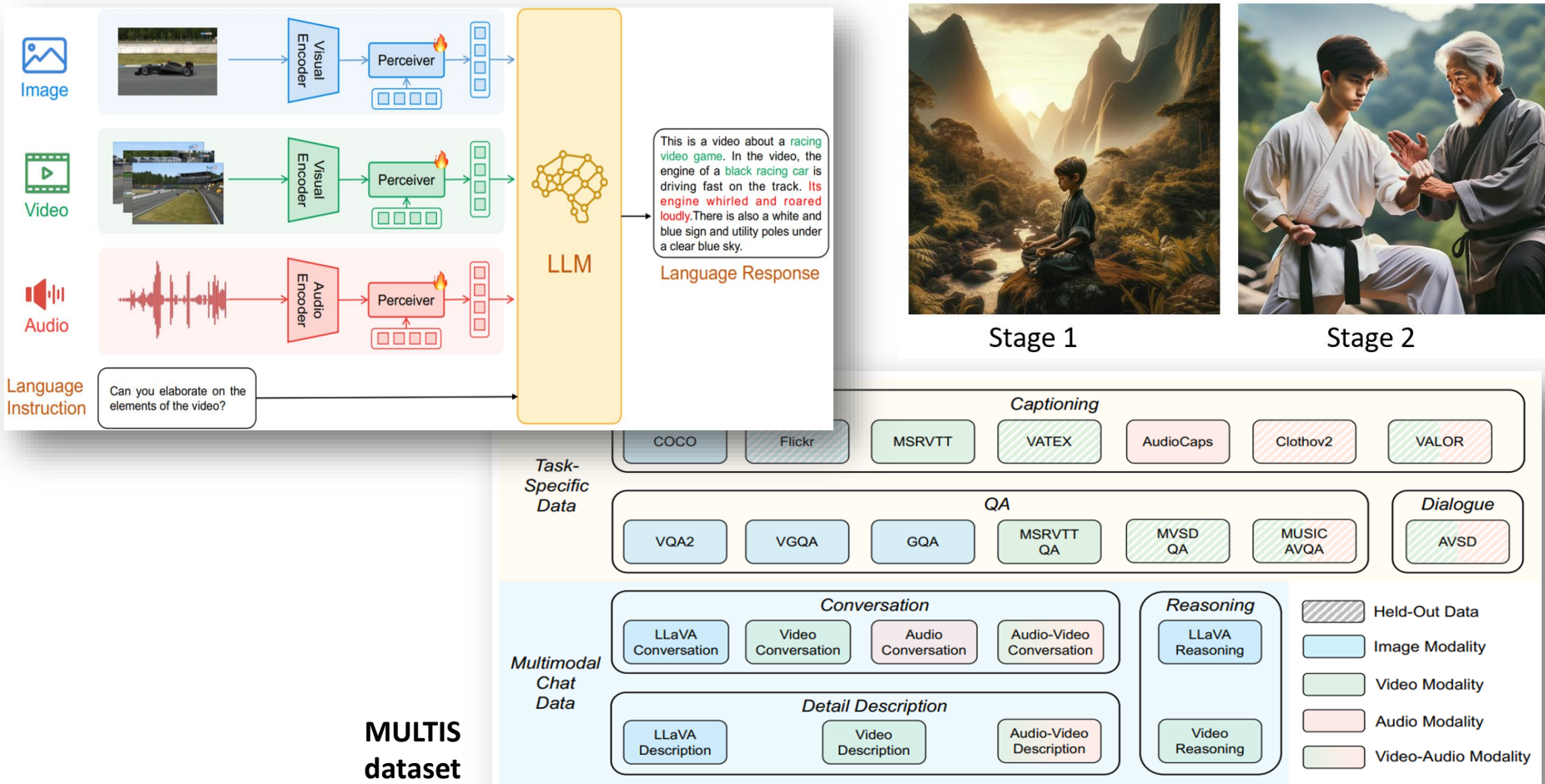
modality-shared perceiver (Q-Former)  
but modality-independent learnable queries



# Pre-training Strategy

- **Two-stage Pre-training**

- **Stage 1 multi-modal alignment:** pre-train **perceivers** with large-scale language-paired two-modality data (image-text, video-text, audio-text)
- **Stage 2 multi-modal instruction tuning:** instruction-finetune **perceivers** with the proposed multi-modal instruction-tuning dataset MULTIS



# Quantitative Result

- Zero-shot evaluation of SOTA methods on various downstream tasks


Methods	Image-Text Tasks				Video-Text Tasks		Audio-Text Tasks
	OKVQA QA	GQA QA	Flickr30k Caption	NoCaps Caption	MSVD QA	VATEX Caption	Clothov2 Caption
Finetuned SoTA	66.1 [18]	65.1 [70]	67.4 [72]	121.6 [31]	60.0 [12]	95.8 [12]	48.8 [38]
Flamingo-9B [3]	44.7	-	61.5	-	30.2	39.5	-
Flamingo-80B [3]	<b>50.6</b>	-	67.2	-	35.6	46.7	-
BLIP-2 (FlanT5-XXL) [17]	-	<b>42.4</b>	73.7	98.4	34.4	-	-
BLIP-2 (Vicuna-13B) [17]	-	32.3	71.6	103.9	20.3	-	-
ChatBridge w/o MULTIS	41.4	37.4	77.7	107.5	23.5	47.7	22.4
ChatBridge	45.2	41.8	<b>82.5</b>	<b>115.7</b>	<b>45.3</b>	<b>48.9</b>	<b>26.2</b>


metrics: accuracy for QA tasks, CIDEr score for captioning tasks


Input Modality	AVSD Dialogue		VALOR Captioning		MUSIC-AVQA
	BLEU-4	CIDEr	BLEU-4	CIDEr	Acc.
Finetuned SoTA	40.0 [44]	108.5 [44]	9.6 [12]	61.5 [12]	78.9 [12]
Video	28.3	73.1	2.8	22.3	33.1
Audio	20.2	46.2	0.3	5.2	28.9
Video+Audio	<b>29.8</b>	<b>75.4</b>	<b>4.2</b>	<b>24.7</b>	<b>43.0</b>


# Qualitative Results


- Multi-round conversation cases with audio or video inputs


Audio Content:  
(A train is coming in.) 


Describe this audio. 


ChatBridge  This audio is a recording of a **train passing by**, likely on a track. The sound is a combination of the **train's engine, wheels on the tracks, and some background noise**. The recording is fairly clear and the sound of the train is prominent.


What other sounds do trains usually make? 


ChatBridge  Trains can make a variety of sounds, including the sound of the engine, the sound of the wheels on the tracks, the sound of the brakes being applied, the sound of the whistle or horn, and the sound of the cars (boxcars, tankers, etc.) clanking or rattling. Additionally, trains can make noise from the movement of the train, such as the swaying of the cars or the rattling of the windows.


What does it mean when a train blows its whistle? 


ChatBridge  A train blowing its whistle is typically a warning to alert pedestrians, other trains, or other vehicles that the train is approaching a crossing or a station. The train's whistle is used as a safety measure to ensure that people are aware of the train's presence and take appropriate action to avoid being hit by the train.





What activities are the people involved in? 

ChatBridge  The video excerpt is from **The Godfather**.

Can you provide me with more background information about this movie? 

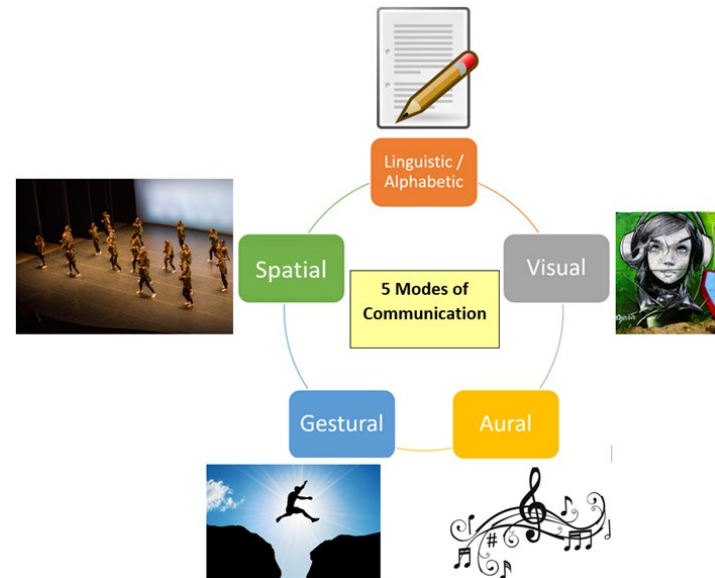
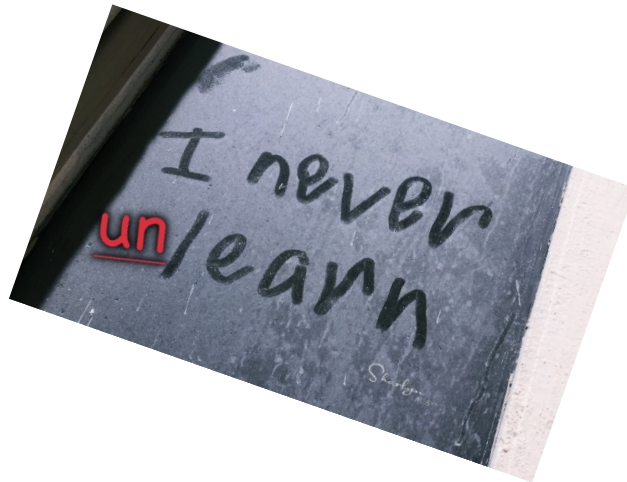
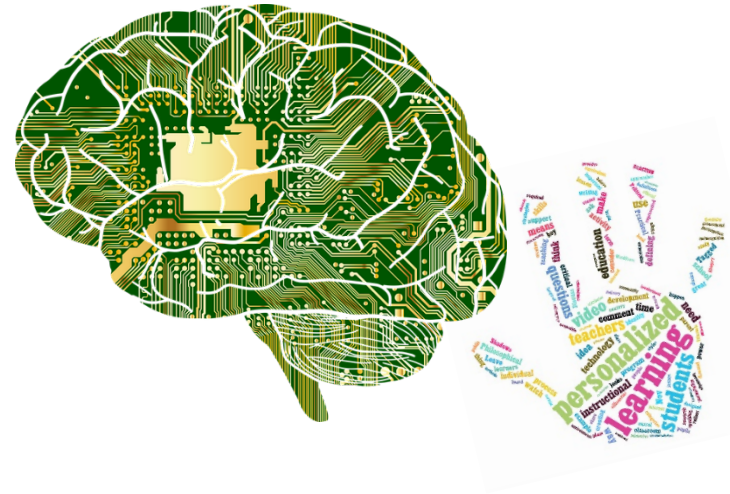
ChatBridge  The Godfather is a classic **American crime film** directed by Francis Ford Coppola and released in 1972. It is based on the novel of the same name by Mario Puzo and follows the story of the Corleone family, a powerful mafia clan in New York City. The movie is widely regarded as one of the greatest films of all time and has won numerous awards, including seven Academy Awards.

Who is the male actor appearing in the video? 

ChatBridge  The male actor appearing in the video is **Marlon Brando**, who played the character of **Don Vito Corleone** in The Godfather. Brando's portrayal of the powerful and complex patriarch of the Corleone family is widely regarded as one of the most iconic performances in film history. His performance earned him an Academy Award for Best Actor.

# What to Be Covered Today...

- Multimodal LLM
- **Advanced Topics in LLM/VLM**
  - Concept Editing
  - Concept Unlearning
  - Personalization
  - Continual Learning

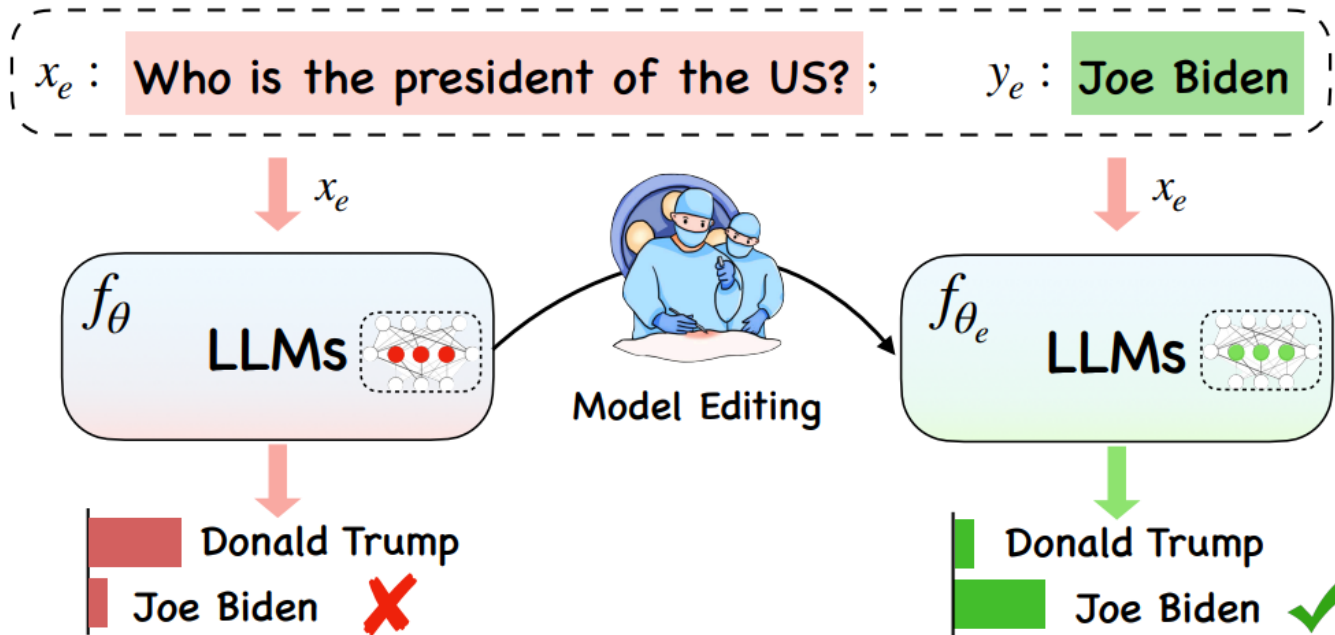


# LLM Editing (i.e., Concept/Knowledge Editing)

- **Motivation**

- The knowledge in LLM will be outdated over time.
  - E.g., The knowledge cutoff of LLaMA3-70B is Dec. 2023.

→ Need an effective and efficient way to *inject* new knowledge w/o re-training.



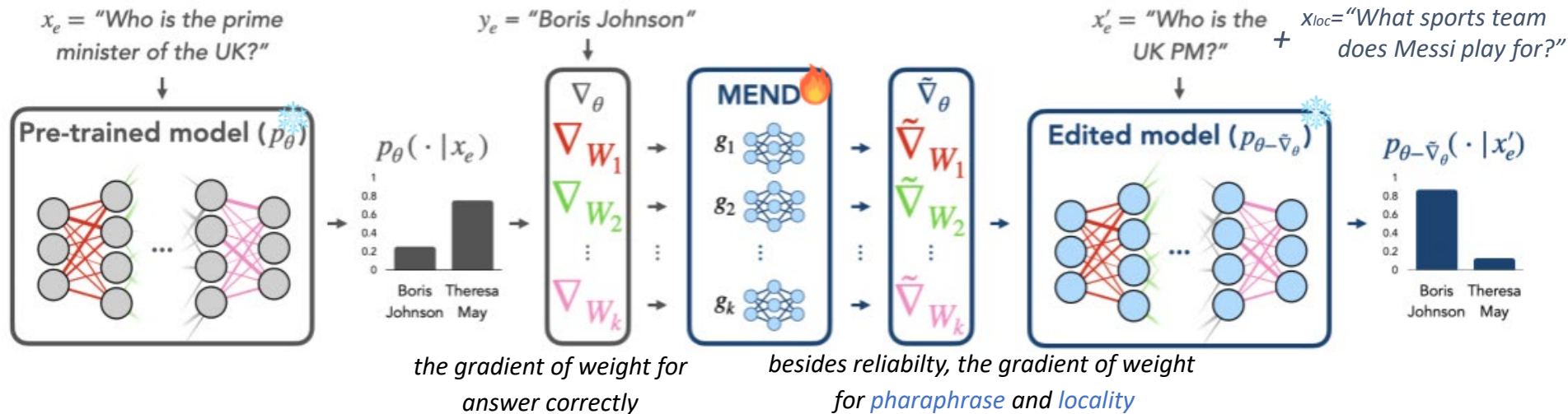
# Recent works on LLM editing

- **Editing with Hypernetwork:**
  - MEND (ICLR'22)
- **Editing with External Memory:**
  - T-Patcher (ICLR'23)
  - WISE (NeurIPS'24)

# Fast Model Editing at Scale, Stanford, ICLR'22

- **Goal:**
  - Inspired by meta-learning, update the pre-trained model by learning a hypernetwork (called Model Editor Networks with Gradient Decomposition, MEND).
- **Method:**
  - Train MEND with edited samples, supervised by equivalent (paraphrased), and unrelated (locality) samples.
- **Limitation? Batch-mode training...**

## Editing a Pre-Trained Model with MEND



**MEND losses:**  $L_e = -\log p_{\theta_{\tilde{w}}}(y'_e|x'_e), L_{loc} = \text{KL}(p_{\theta_{\tilde{w}}}(\cdot|x_{loc})||p_{\theta_{\tilde{w}}}(\cdot|x_{loc})).$  (4a,b)

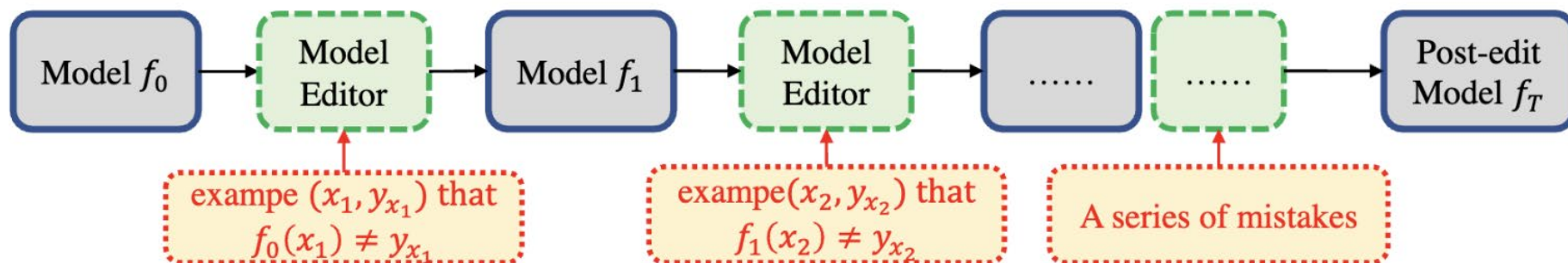


# Transformer-Patcher: One Mistake Worth One Neuron

## Mila & WeChat, ICLR'23 (1/2)

- **Idea:**

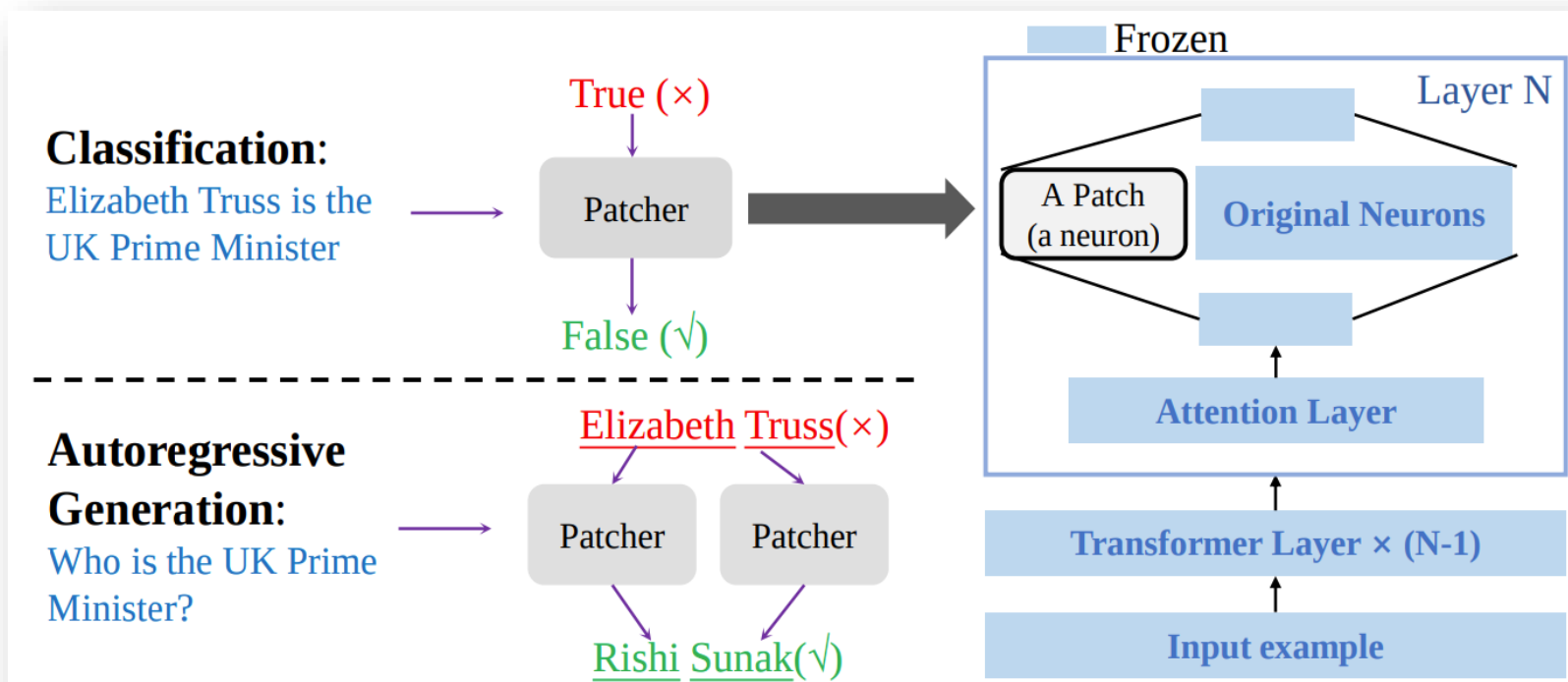
- Previous works (e.g., MEND) typically only handle “one-step” editing
- Need a method to handle “lifelong” editing
- Example: one-step vs. lifelong editing
  - We have model edited from  $x_{e1}$  to  $x_{en}$  (e.g.,  $x_{e1}$  = who is the UK prime minister?). Now, we want to edit model at (n+1)-th sample (e.g.,  $x_{e(n+1)}$  = who is the president in Taiwan?). In “one-step” editing, we need to **re-train** model **from**  $x_{e1}$  **to**  $x_{e(n+1)}$ . In “lifelong” editing, only need to **train** model **at**  $x_{e(n+1)}$  => **more efficient**.



# Transformer-Patcher: One Mistake Worth One Neuron

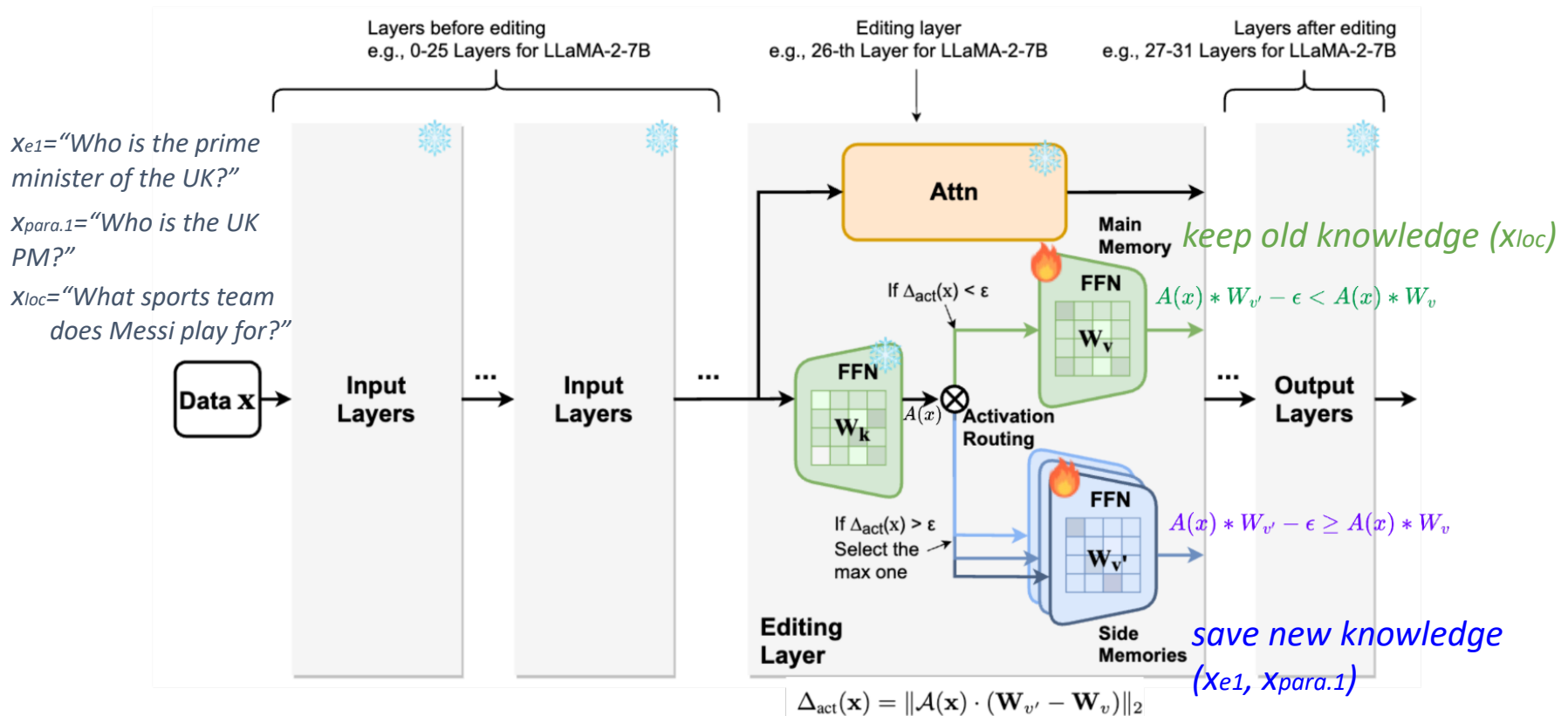
## Mila & WeChat, ICLR'23 (2/2)

- **Goal:**
  - Previous work (i.e., MEND): only handle “one-step” editing in batch modes
  - Propose T-Patcher to handle “lifelong” editing
- **Method:**
  - Patch transformers by adding **neurons** for each edited knowledge



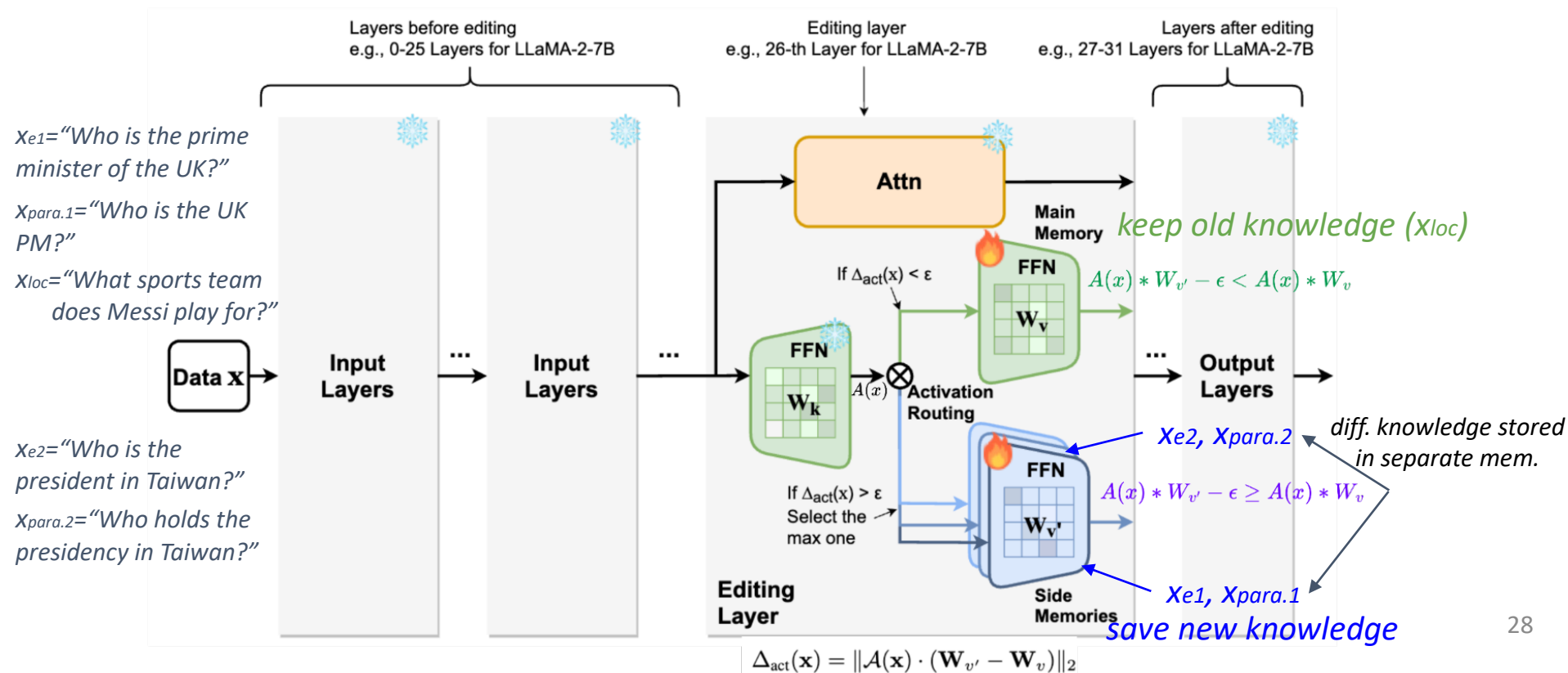
# Rethinking the Knowledge Memory for Lifelong Model Editing of Large Language Models, Alibaba, NeurIPS'24 (1/2)

- **Goal:**
  - Previous work (i.e., T-Patcher): linear growing memory complexity
- **Method:**
  - Employ *finite* side memories and design a *routing* mechanism to choose them



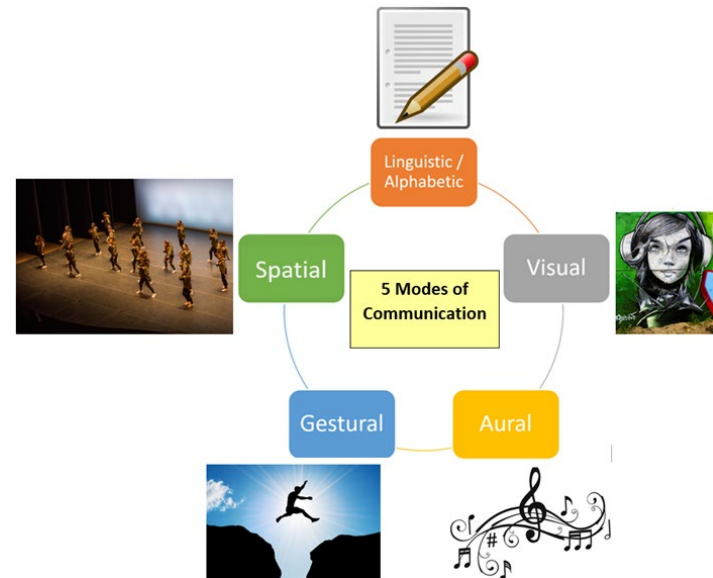
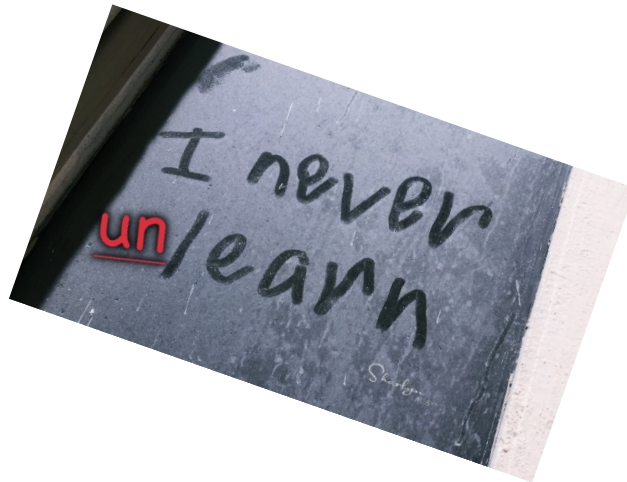
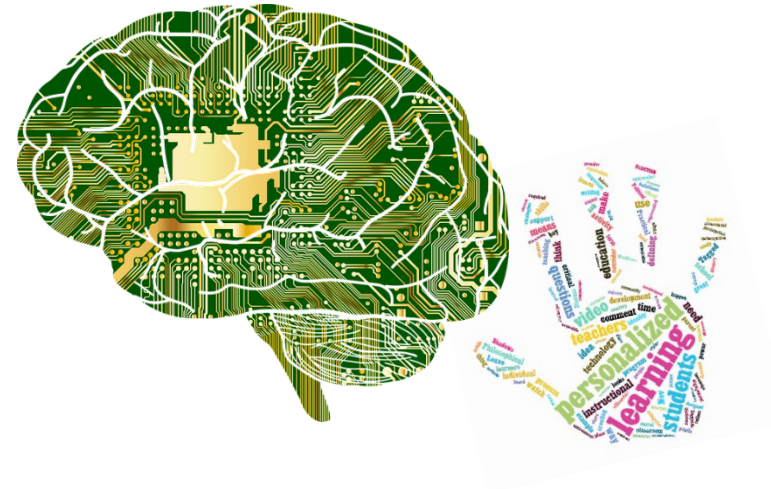
# Rethinking the Knowledge Memory for Lifelong Model Editing of Large Language Models, Alibaba, NeurIPS'24 (2/2)

- **Goal:**
  - Previous work (i.e., T-Patcher): linear growing memory complexity
- **Method:**
  - Employ *finite* side memories and design a *routing* mechanism to choose them



# What to Be Covered Today...

- Multimodal LLM
- **Advanced Topics in LLM/VLM**
  - Concept Editing
  - Concept Unlearning
  - Personalization
  - Continual Learning



# Diffusion Models Concept Unlearning

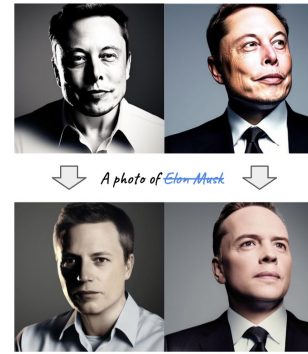
- **Motivation**

- Recently, **diffusion models** have made significant advances in GenAI.
  - E.g., DALL·E 3, Midjourney, Stable Diffusion, Sora...
- However, improper use may result in generating **harmful**, **NSFW**, or **copyrighted** content.

→ Need an effective and efficient way to **unlearn** these undesired concepts.



# Diffusion Models Concept Unlearning (cont'd)



- **Task definition:**
  - **Unlearn** the undesired concepts from pre-trained Diffusion Model, so that it no longer generates images containing that concept.
- Including high-level concept, artistic style, object, or personality.

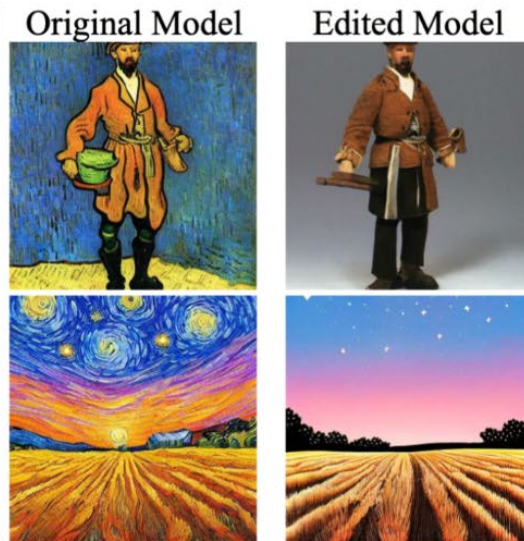
## Erasing Nudity



\* Added by authors for publication

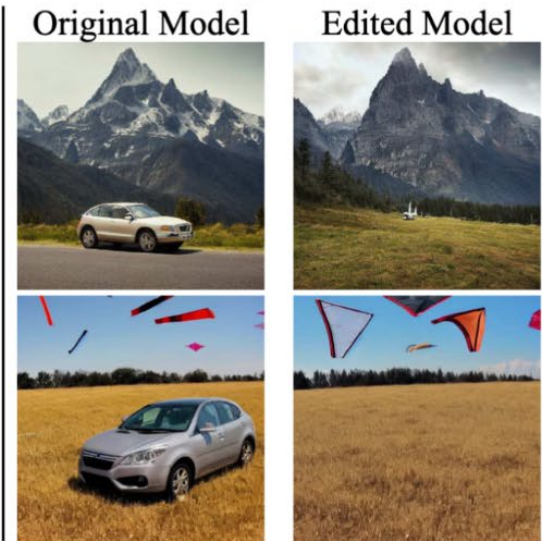
Erased from model:  
“Nudity”

## Erasing Artistic Style



Erased from model:  
“Van Gogh”

## Erasing Objects

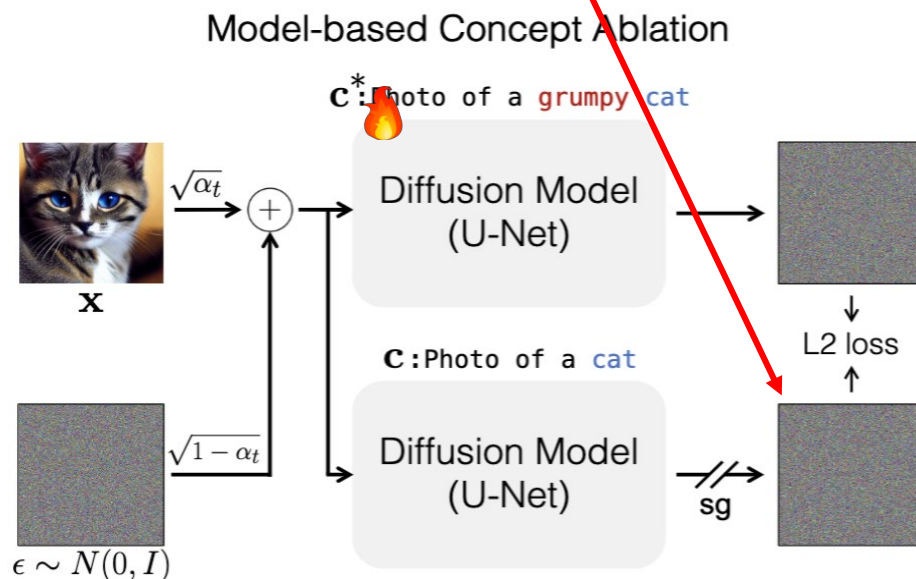


Erased from model:  
“Car”

# Ablating Concepts in Text-to-Image Diffusion Models

## CMU & Adobe Research, ICCV 2023

- **Idea:** **Replace** the output of target concept with predefined anchor concept.
- For example, after unlearning “grumpy cat”:
  - **Input:** Text of target concept (e.g., “a photo of grumpy cat”)
  - **Output:** Image of anchor concept (e.g., Image of cat)
- **Method:**
  - Simply use the predicted noise of anchor as ground truth for fine-tuning.



May be vulnerable to paraphrase!

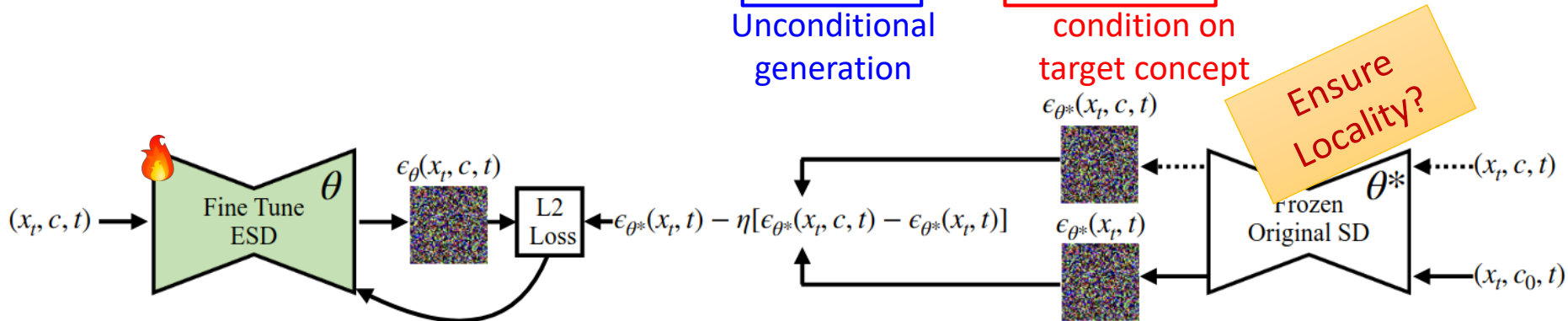


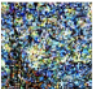
# Erasing Concepts from Diffusion Models

## Northeastern & MIT, ICCV 2023

- Idea:
  - **Reduce** the probability  $p(c|x)$  that the output image belongs to **target concept**
- Method:
  - Utilize **classifier-free guidance** in a **reverse direction** to fine-tune model.
    - Move the output distribution away from target concept.

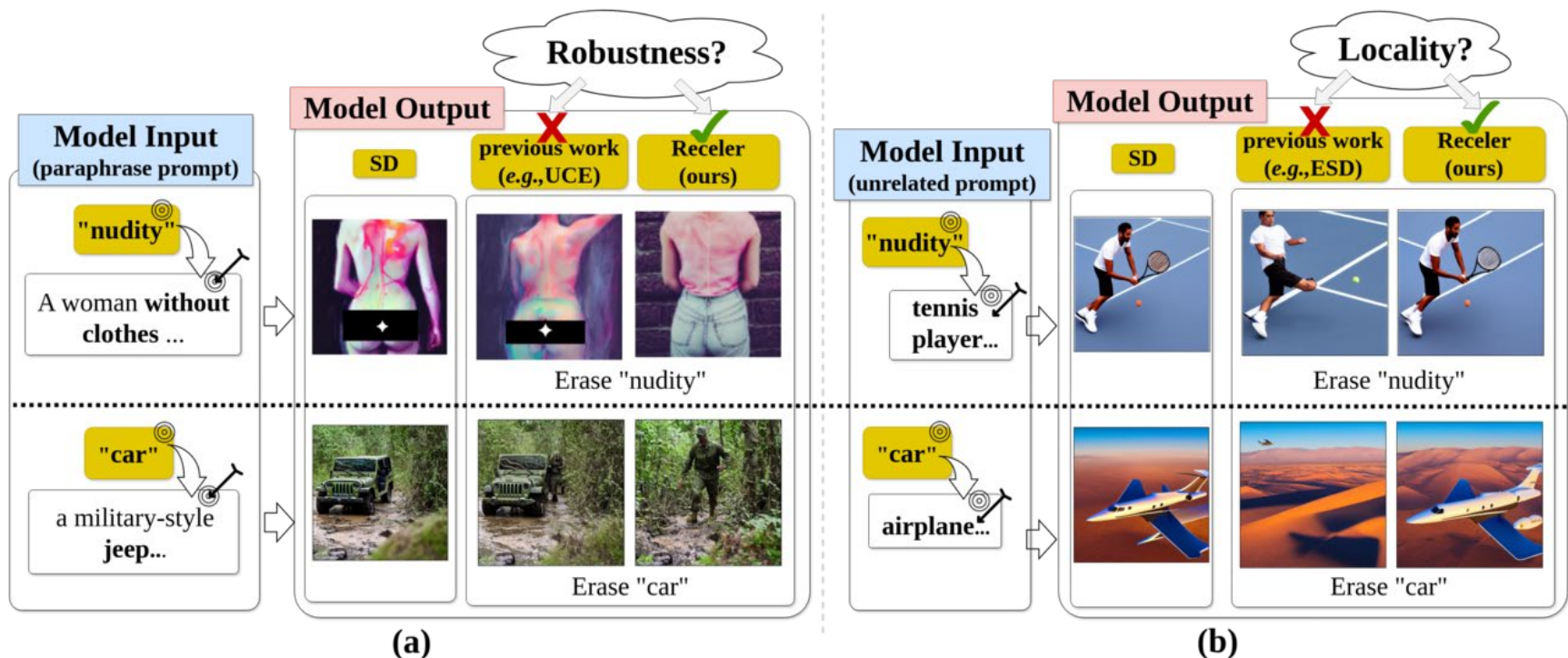
$$\epsilon_{\theta}(x_t, c, t) \leftarrow (1 + \eta) \underbrace{\epsilon_{\theta^*}(x_t, t)}_{\text{Unconditional generation}} - \eta \underbrace{\epsilon_{\theta^*}(x_t, c, t)}_{\text{condition on target concept}}$$



$x_t \rightarrow$   , generated by  $\theta$      
  $c \rightarrow$  "Van Gogh", concept to erase     
  $c_0 \rightarrow$  ""     
  $t \rightarrow$  Time step sampled uniformly

# Reliable Concept Erasing of Text-to-Image Diffusion Models via Lightweight Erasers (Receler), VL Lab (NTU), ECCV 2024

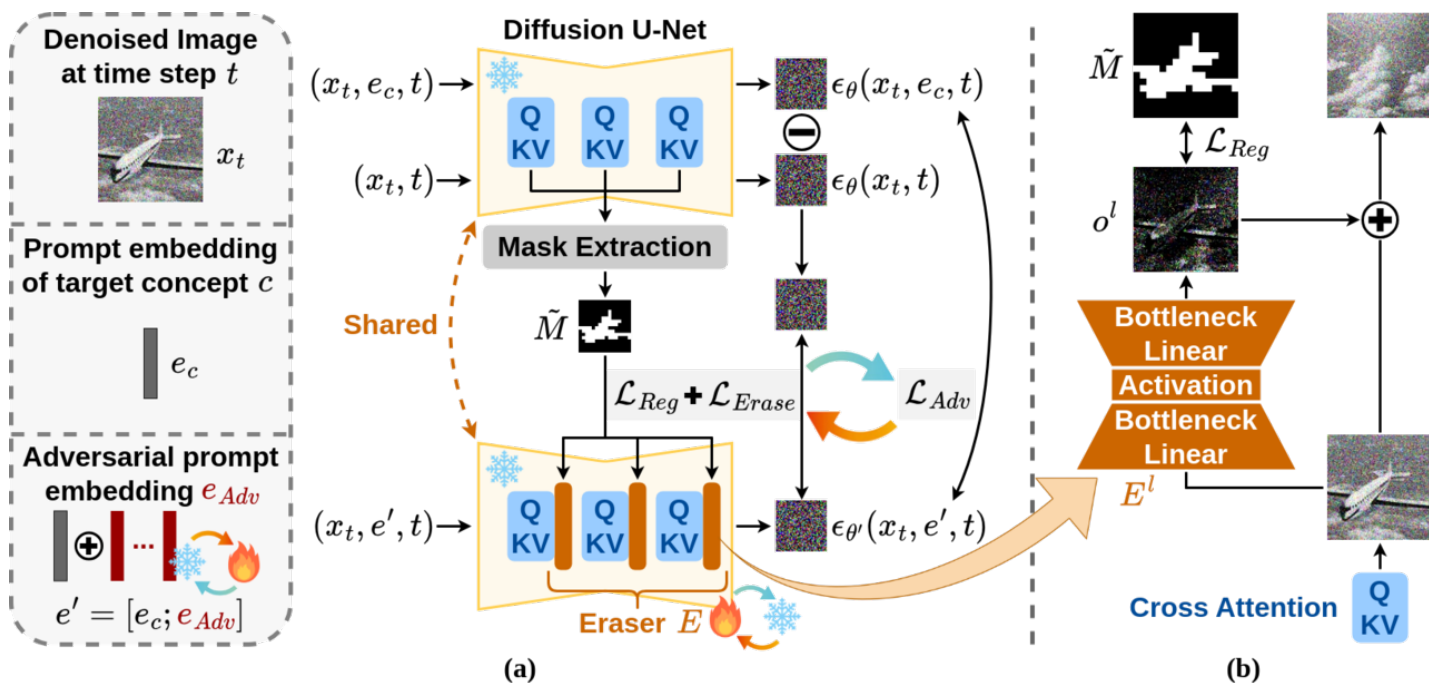
- Define **Reliable Concept Erasing** by two desirable properties:
  - **Robustness**: resistance when inputting paraphrased/adversarial attack prompts
  - **Locality**: capability on preserving the generation of non-target concepts.



# Receler, VL Lab (NTU), ECCV 2024

## • Method

- Lightweight Eraser (PEFT)
- Concept-localized Regularization
- Adversarial Prompt Learning



# Receler, VL Lab (NTU), ECCV'24 (cont'd)

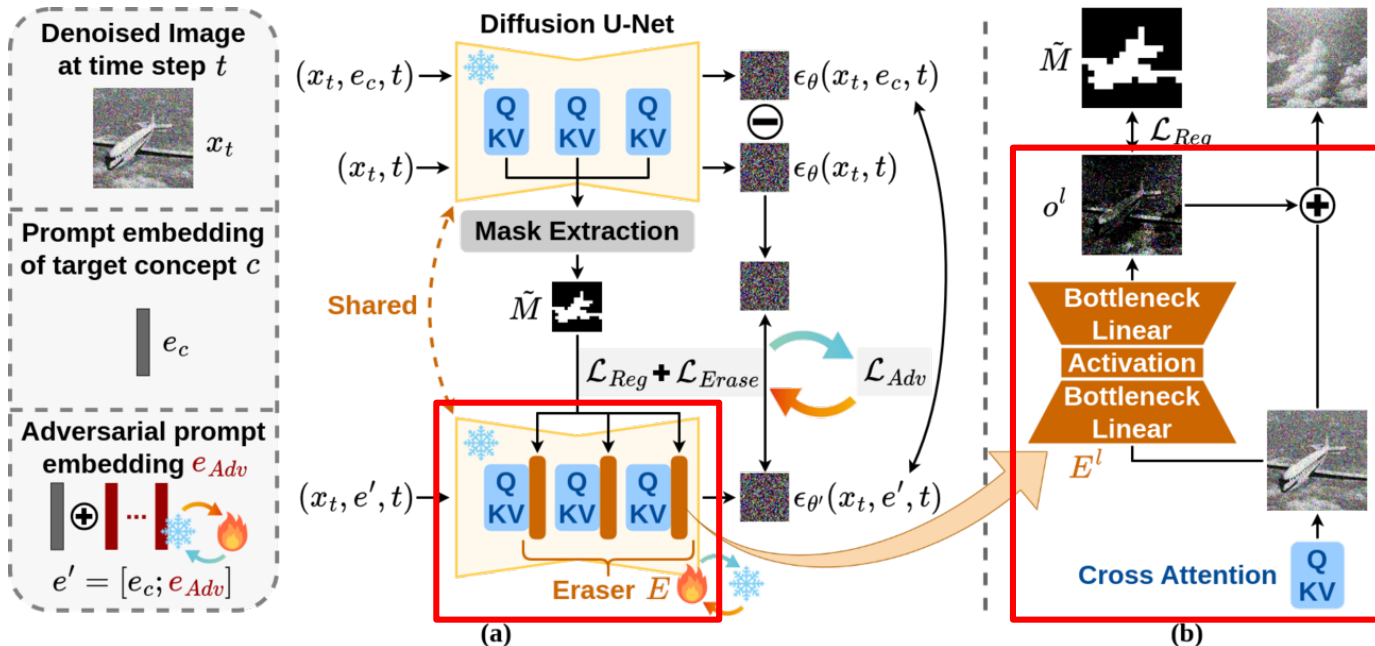
- Method (1/3)

- Lightweight Eraser (PEFT)

- Adapter-based module with only 0.37% param. of UNet
- Inserted after cross-attention layer to remove target concept features
- We use the CFG objective to fine-tune this eraser (like ESD did).

$$\mathcal{L}_{Erase} = \mathbb{E}_{x_t, t} [\|\epsilon_{\theta'}(x_t, e_c, t) - \epsilon_E\|^2],$$

where  $\epsilon_E = \epsilon_{\theta}(x_t, t) - \eta [\epsilon_{\theta}(x_t, e_c, t) - \epsilon_{\theta}(x_t, t)]$ .



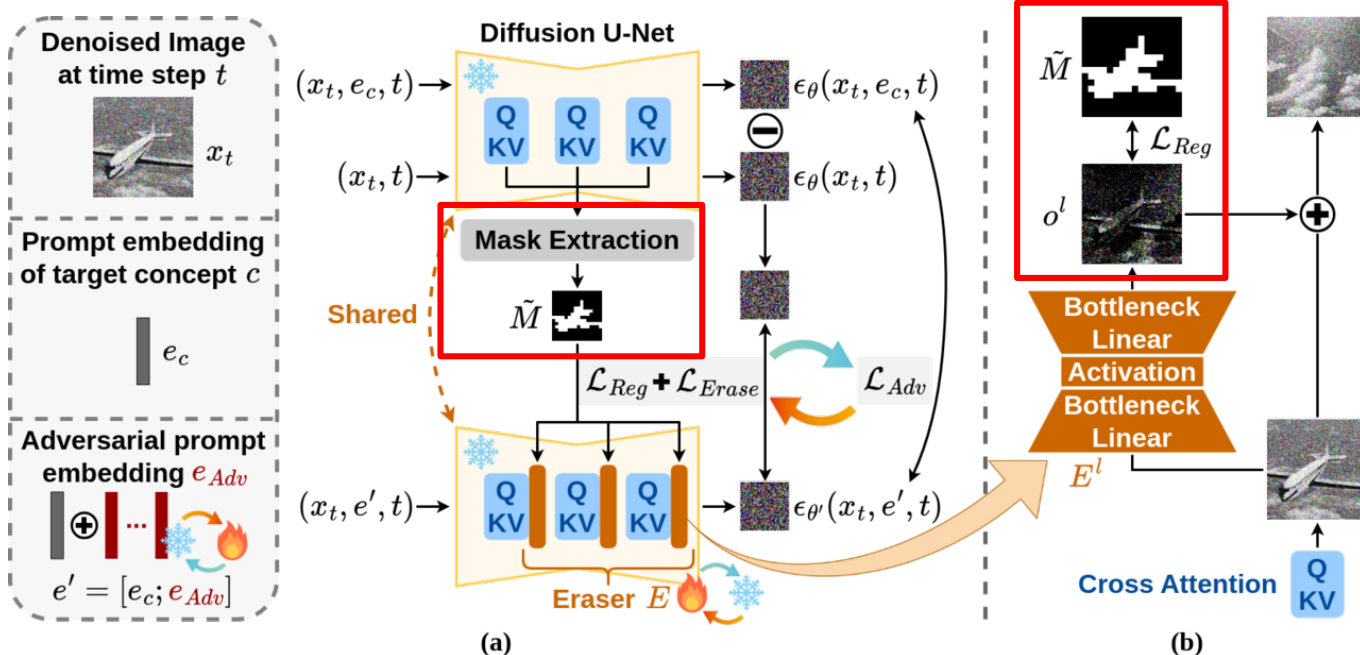
# Receler, VL Lab (NTU), ECCV'24 (cont'd)

- Method (2/3)

- Concept-localized Regularization

- To achieve **locality**, we use attention-maps extracted from UNet to regularize eraser learning
    - With this regularization, Eraser focus only on regions of target concept

$$\mathcal{L}_{Reg} = \frac{1}{L} \sum_{l=1}^L \|o^l \odot (1 - \tilde{M})\|^2$$



# Receler, VL Lab (NTU), ECCV'24 (cont'd)

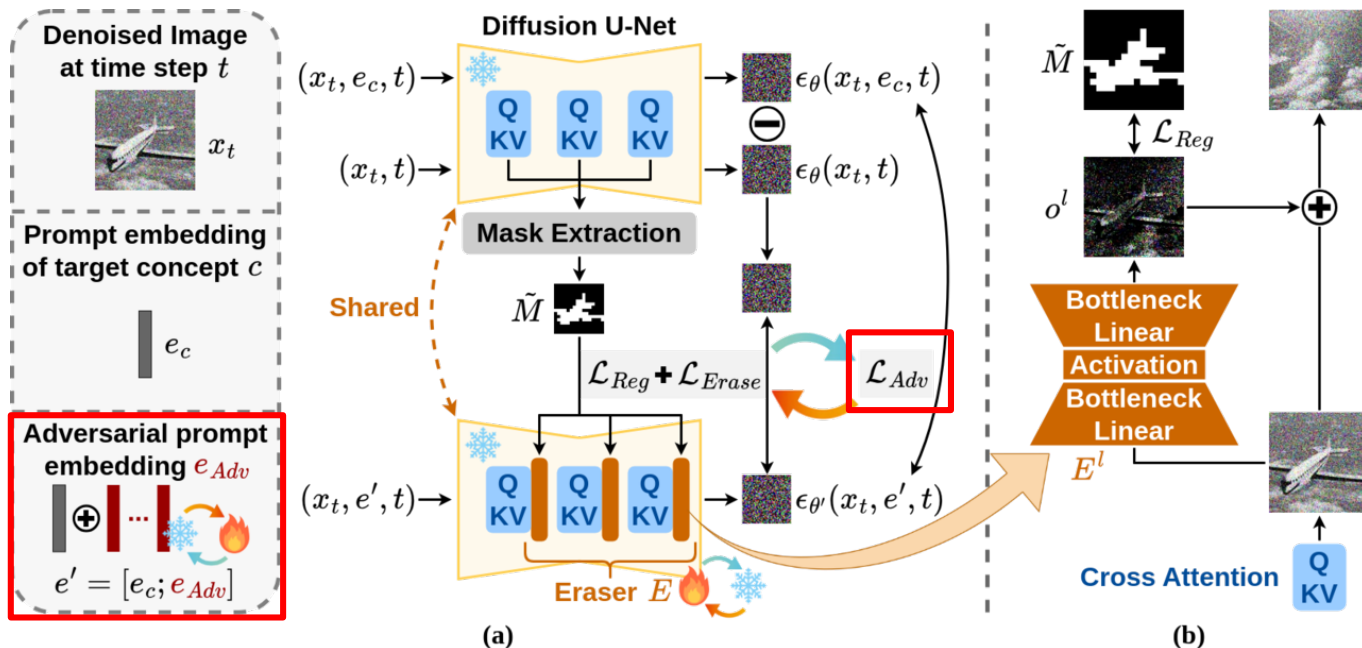
- Method (3/3)

- Adversarial Prompt Learning

- To enhance **robustness**, we train adversarial embeddings to mimic malicious attacks
- **Prompt embeddings** and **Eraser** are trained iteratively **against** each other.
  - **Prompt embeddings** ↪ try to generate target concept
  - **Eraser** ↪ try to remove target concept

Predicted noise of target concept

$$\mathcal{L}_{Adv} = \mathbb{E}_{x_t, t} [\|\epsilon_{\theta'}(x_t, e', t) - \epsilon_M\|^2]$$



# Receler, VL Lab (NTU), ECCV'24 (cont'd)

- **Quantative Results**
  - Dataset: I2P (Inappropriate Prompt dataset)
  - Any idea how to perform objective evaluation?
- Achieve **SOTA** in erasing inappropriate concepts with **robustness & locality**

Class name	Inappropriate proportion (%) (↓)					
	SD	FMN	SLD	ESD	UCE	<i>Receler</i>
Hate	44.2	37.7	22.5	26.8	36.4	28.6
Harassment	37.5	25.0	22.1	24.0	29.5	21.7
Violence	46.3	47.8	31.8	35.1	34.1	27.1
Self-harm	47.9	46.8	30.0	33.7	30.8	24.8
Sexual	60.2	59.1	52.4	35.0	25.5	29.4
Shocking	59.5	58.1	40.5	40.1	41.1	34.8
Illegal activity	40.0	37.0	22.1	26.7	29.0	21.3
Overall	48.9	47.8	33.7	32.8	31.3	<b>27.0</b>

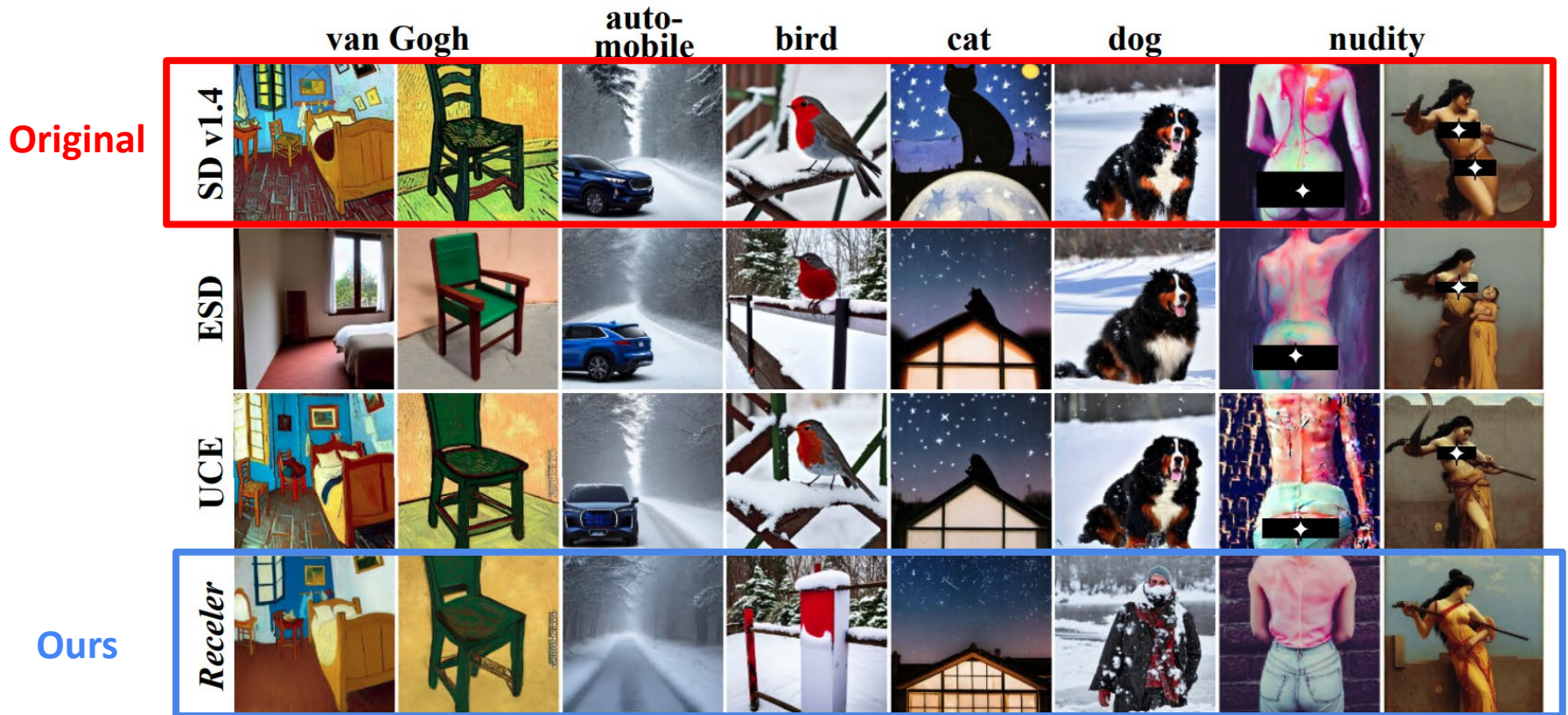
Fig 1. Erase inappropriate concepts in I2P dataset using NudeNet

Method	Robustness	Locality	
	Nudity-erased ratio(↑)	CLIP-30K(↑)	FID-30K(↓)
SD	-	31.32	14.27
FMN	44.2%	30.39	<b>13.52</b>
SLD	71.6%	30.90	16.34
ESD	81.3%	30.24	15.31
UCE	75.9%	30.85	14.07
<i>Receler</i>	<b>84.5%</b>	<b>31.02</b>	14.10

Fig 2. Robustness and locality metrics

# Receler, VL Lab (NTU), ECCV'24 (cont'd)

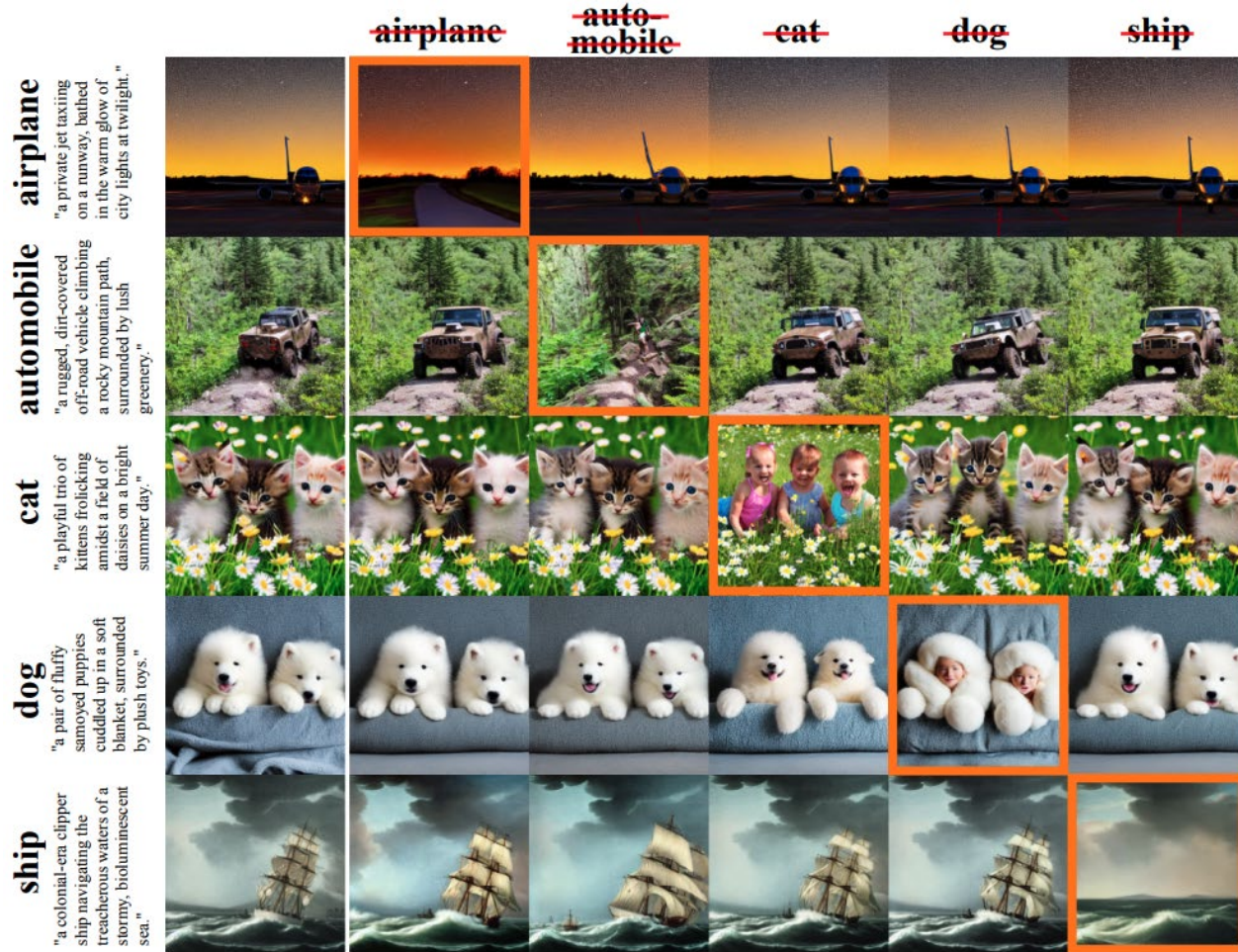
- Qualitative Results (1/2)





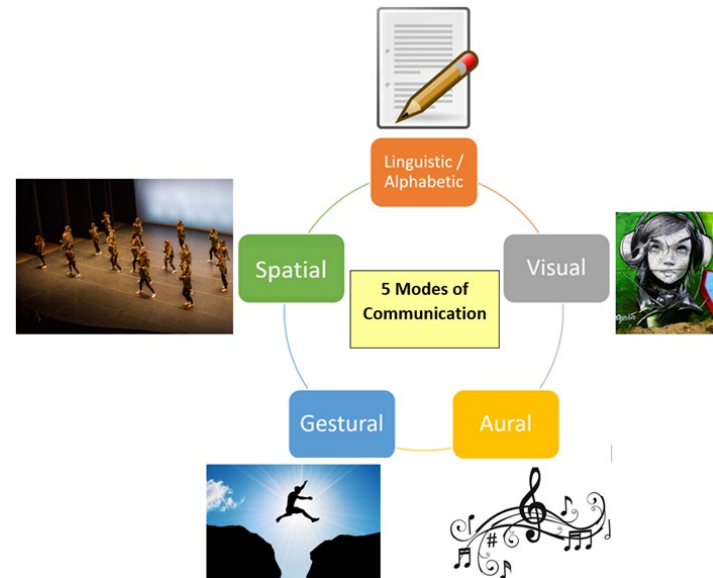
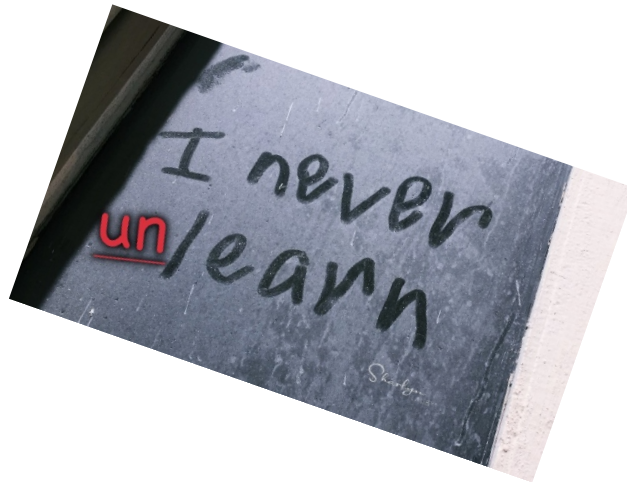
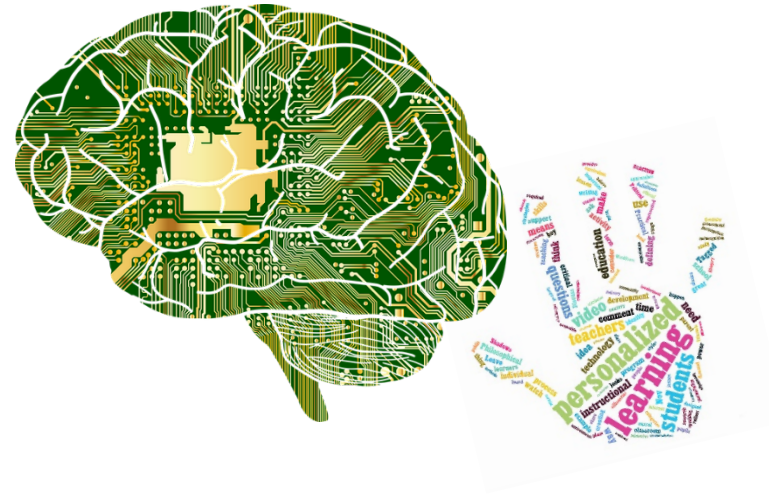
# Receler, VL Lab (NTU), ECCV'24 (cont'd)

- Qualitative Results (2/2)



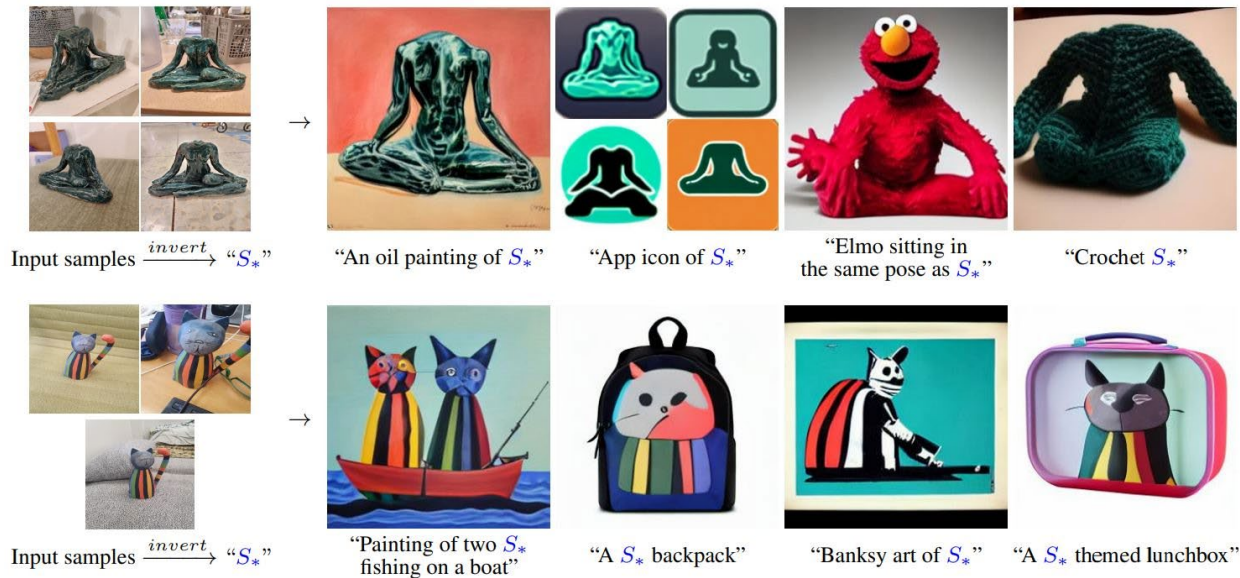
# What to Be Covered Today...

- Multimodal LLM
- **Advanced Topics in LLM/VLM**
  - Concept Editing
  - Concept Unlearning
  - Personalization
  - Continual Learning



# Diffusion Model for Personalization

- **(Recap) Single Concept**
  - Textual Inversion, ICLR'23
  - DreamBooth, CVPR'23
- **Multiple Concepts**
  - CustomDiffusion, CVPR'23
  - Mix-of-Show, NeurIPS'23
- **Beyond Image: Video Motion Customization**
  - Video Motion Customization, CVPR'24



# Single Concept Personalization

- **Definition**

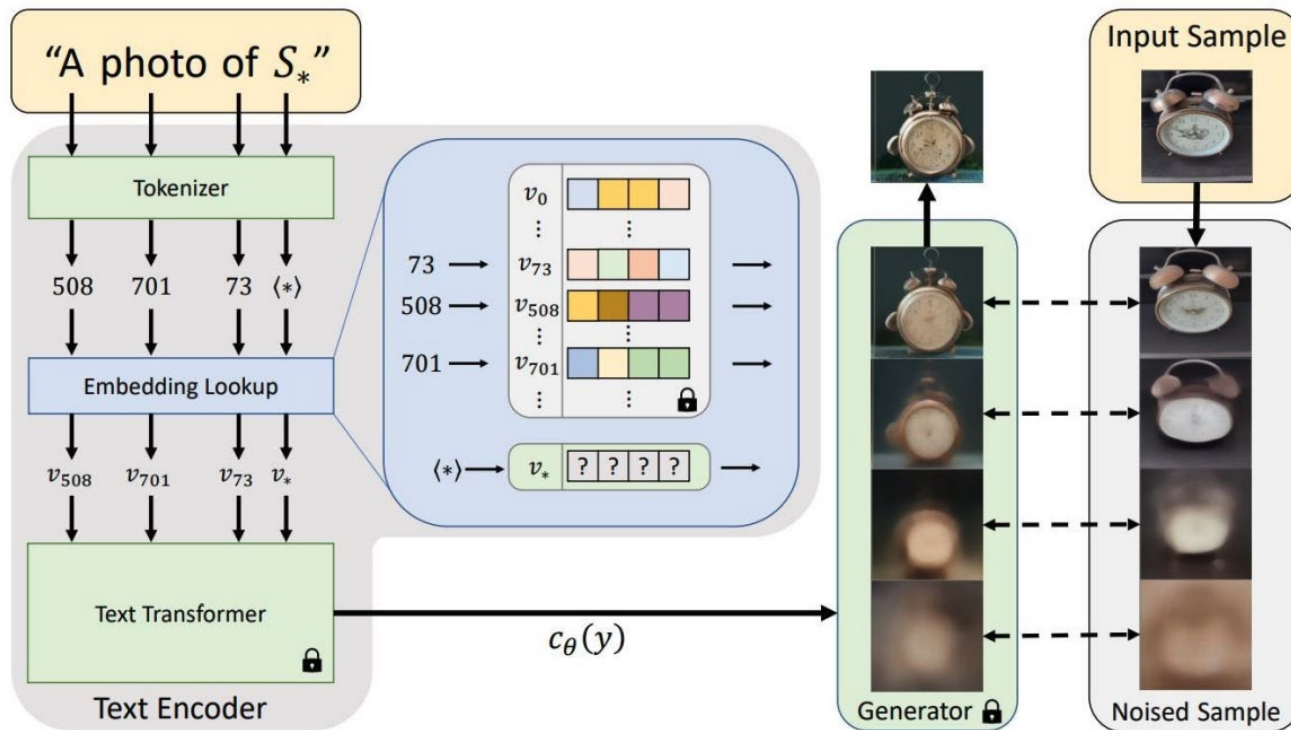
- Given a number of subject images, fine-tune a **pre-trained diffusion model** to enable the generation of that special subject.



# Recap: Textual Inversion, ICLR'23

- **Method: Learning of special token  $S^*$**

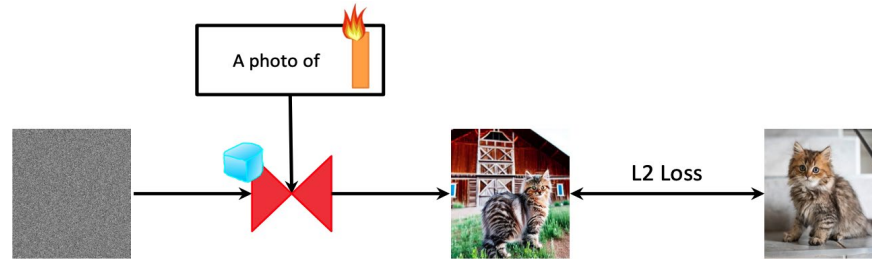
- Pre-train and fix text encoder & diffusion model (i.e., generator)
- Randomly initialize a token as the text encoder input
- Optimize this token via image reconstruction objectives



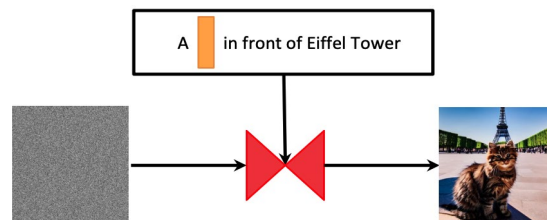
# Recap: Textual Inversion, ICLR'23

- **Method: Learning of special token  $S^*$**

- Pre-train and fix text encoder & diffusion model (i.e., generator)
- Randomly initialize a token as the text encoder input
- Optimize this token via image reconstruction objectives
- Training:



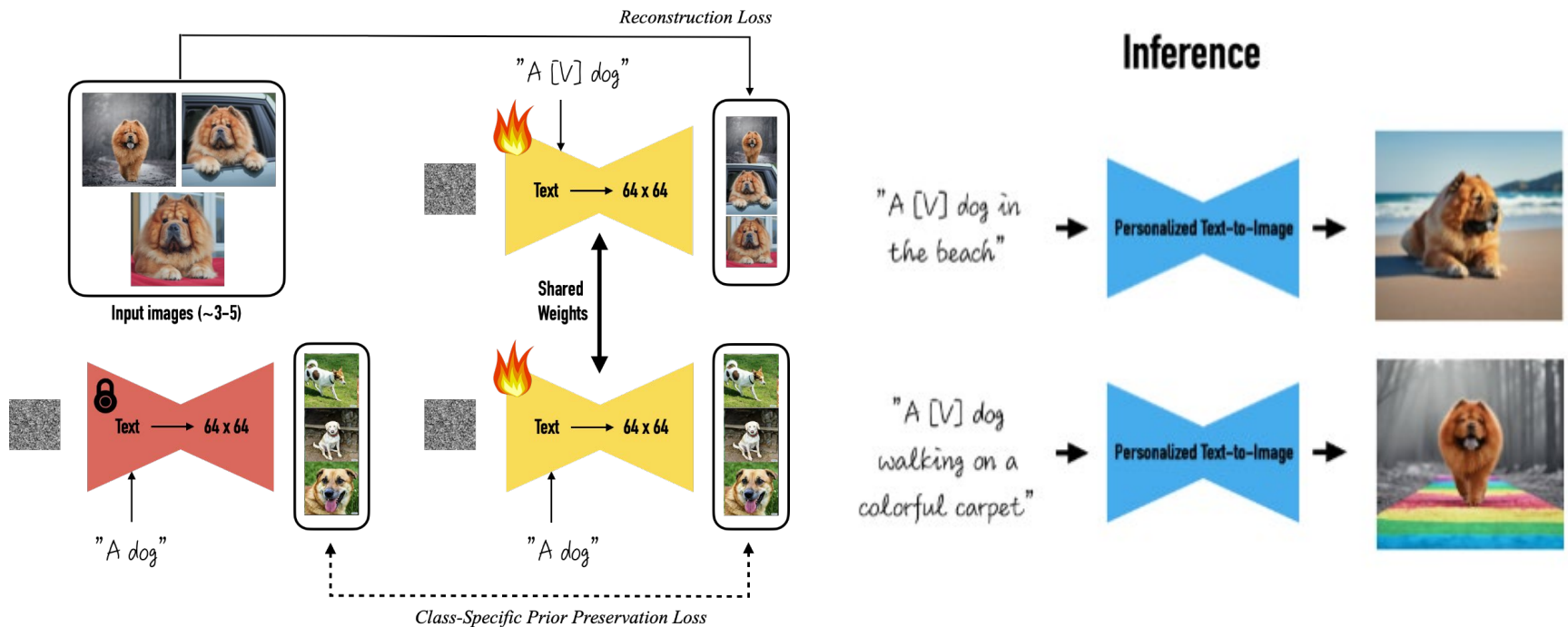
- Inferences:



- **Potential concern?**

# Recap: Single-Concept Personalization - DreamBooth

- Proposed by Google Research, CVPR 2023
- Finetune the diffusion model w/ a fixed token to represent the image concept
  - Determine and fix a rare token (e.g., [V])
  - Finetune the diffusion model for image restoration objectives
  - Enforce a class-specific prior (**why?**)



- **Any concern?**

# Diffusion Model for Personalization

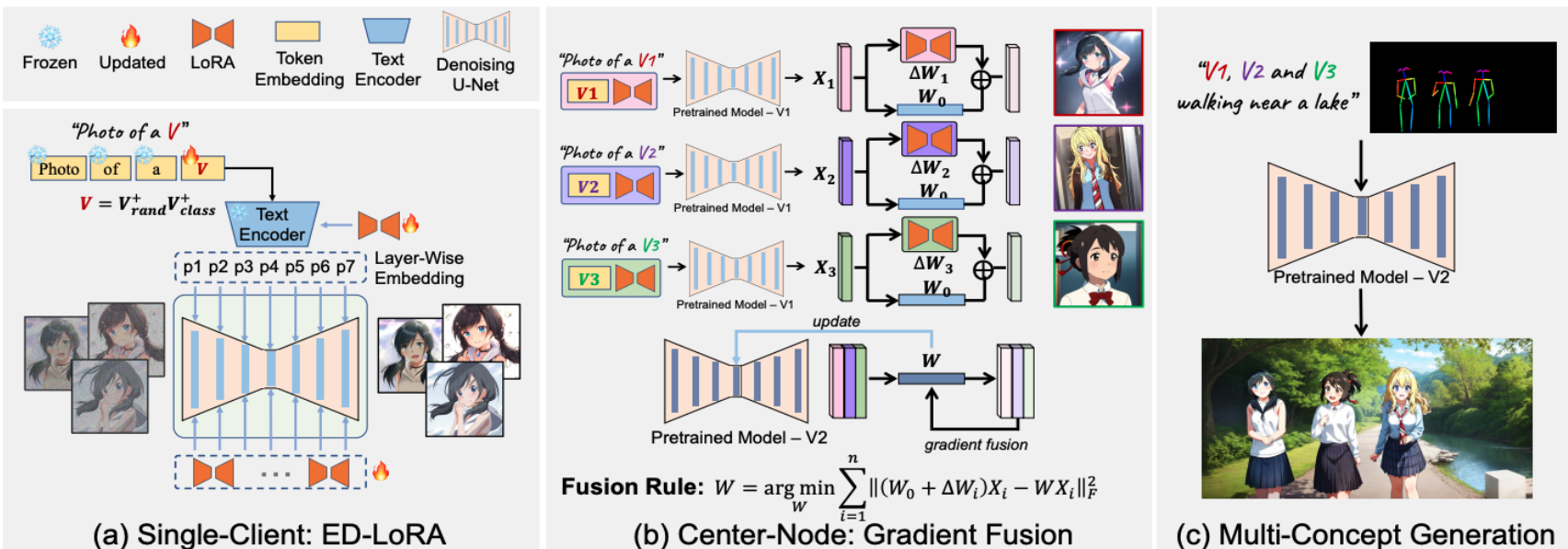
- (Recap) Single Concept
  - Textual Inversion, ICLR'23
  - DreamBooth, CVPR'23
- Multiple Concepts
  - CustomDiffusion, CVPR'23
  - Mix-of-Show, NeurIPS'23
- Beyond Image: Video Motion Customization
  - VMC, CVPR'24





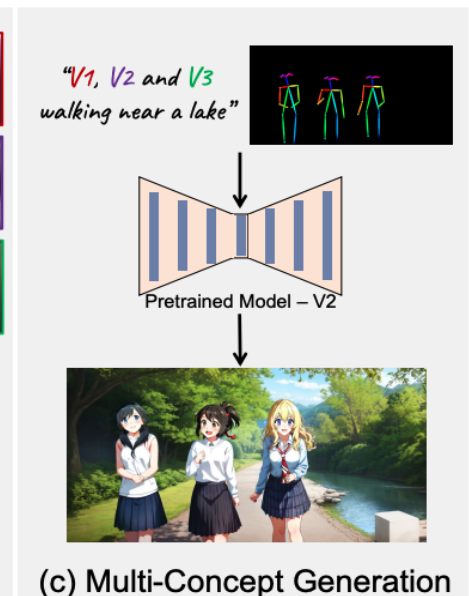
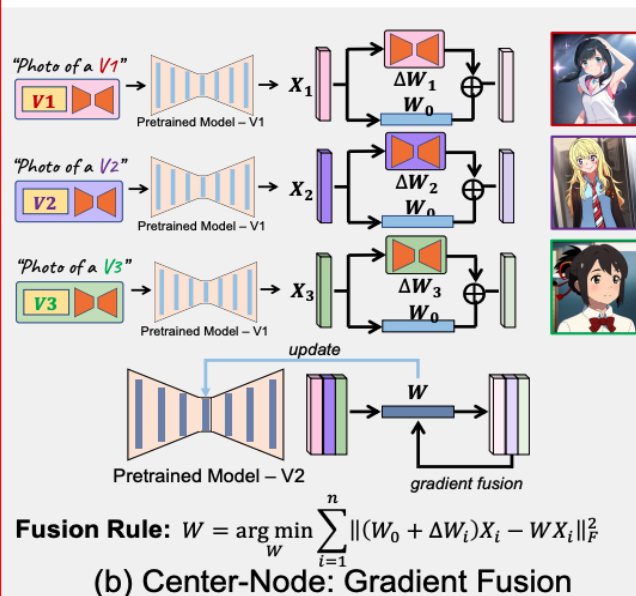
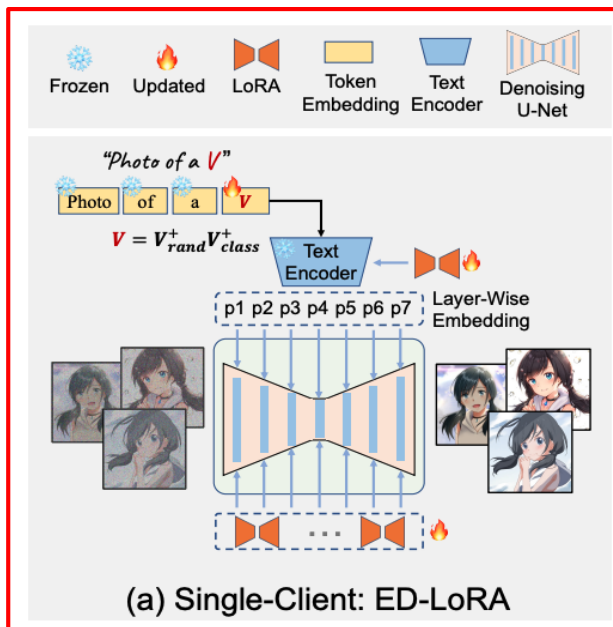
# Multi-concept Personalization: Mix-of-Show, NUS & Tecent, NeurIPS'23

- **Training Method**
  - **Single-Client Learning:** learn LoRA for each concept separately
  - **Center-Node Fusion:** LoRA fusion
- **Inference**
  - Use the fused LoRA to generate multi-concept images



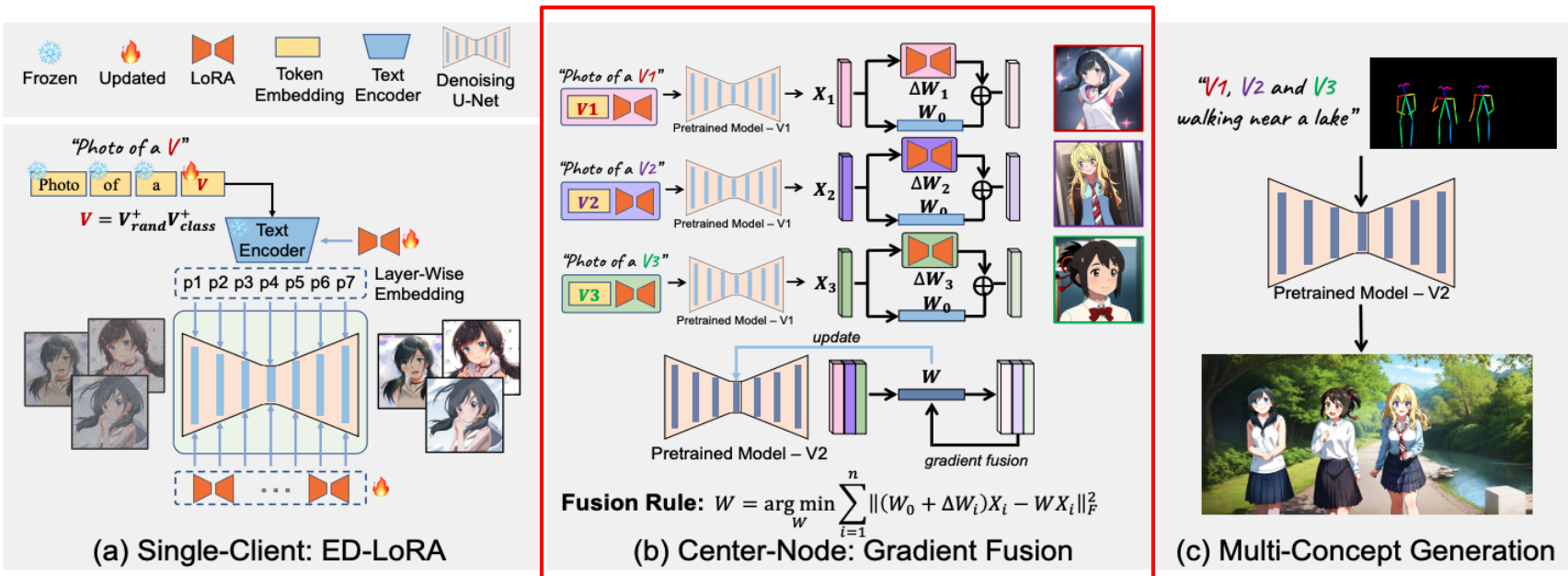
# Multi-concept Personalization: Mix-of-Show, NeurIPS'23

- Training Method
  - **Single-Client Learning:** learn LoRA for each concept separately.
  - **Center-Node Fusion:** LoRA fusion
- Inference
  - Use the fused LoRA to generate multi-concept images



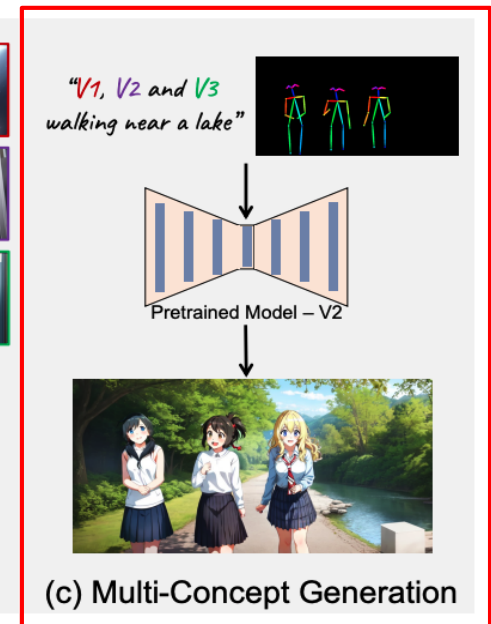
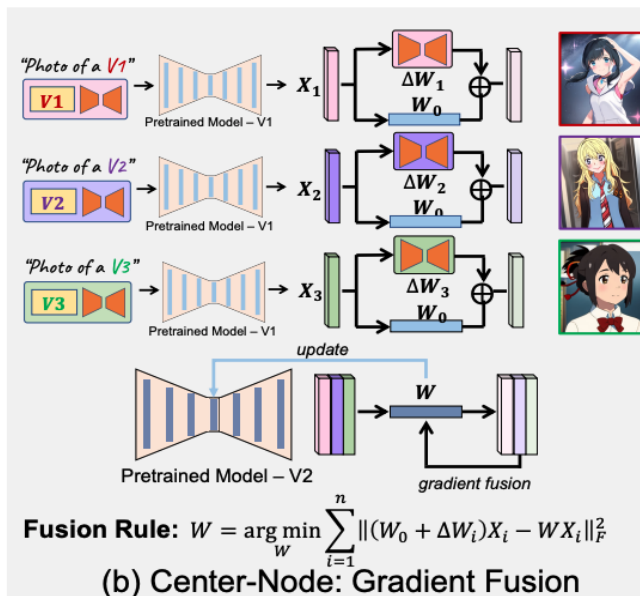
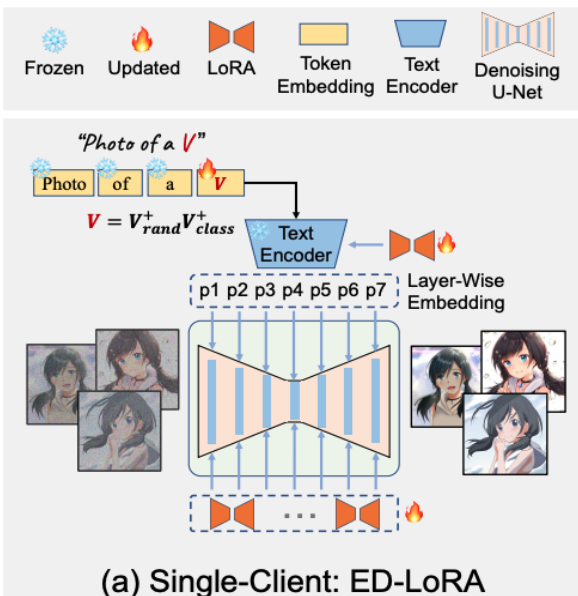
# Multi-concept Personalization: Mix-of-Show, NeurIPS'23

- **Training Method**
  - **Single-Client Learning:** learn LoRA for each concept separately.
  - **Center-Node Fusion:** fuse the above LoRAs into a single LoRA...how?
- **Inference**
  - Use the fused LoRA to generate multi-concept images



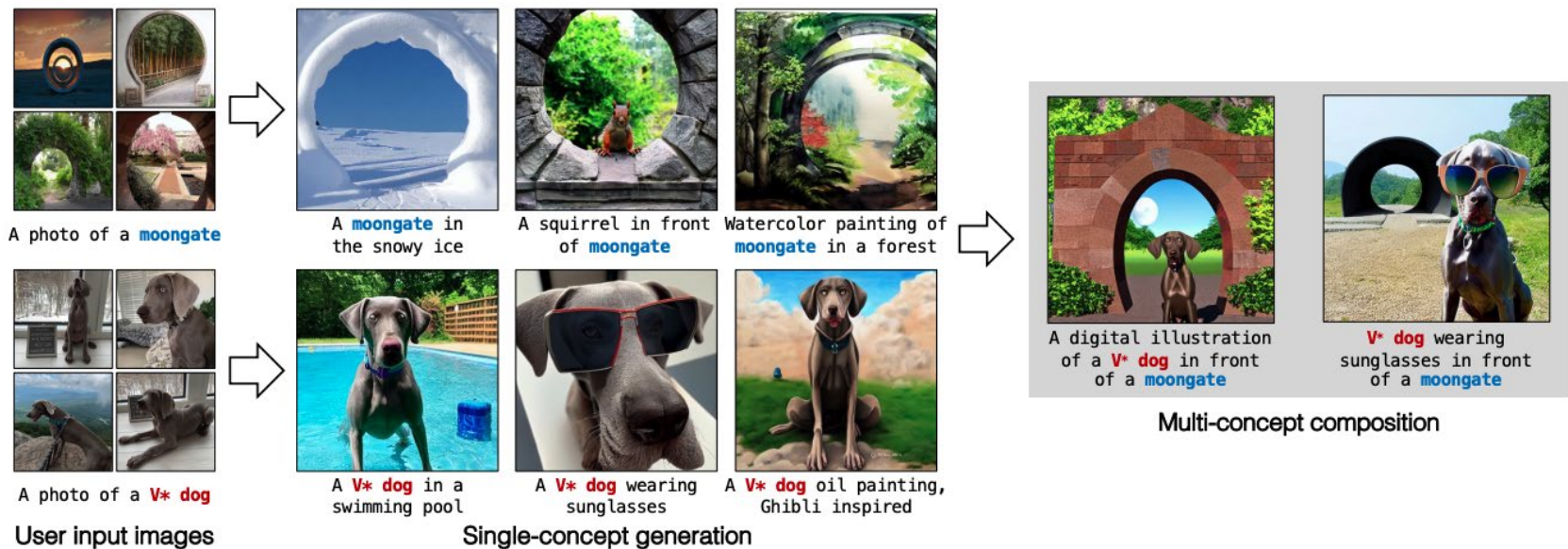
# Multi-concept Personalization: Mix-of-Show, NeurIPS'23

- **Training Method**
  - **Single-Client Learning:** learn LoRA for each concept separately.
  - **Center-Node Fusion:** fuse the above LoRAs into a single LoRA...how?
- **Inference**
  - **Use the fused LoRA to generate multi-concept images**



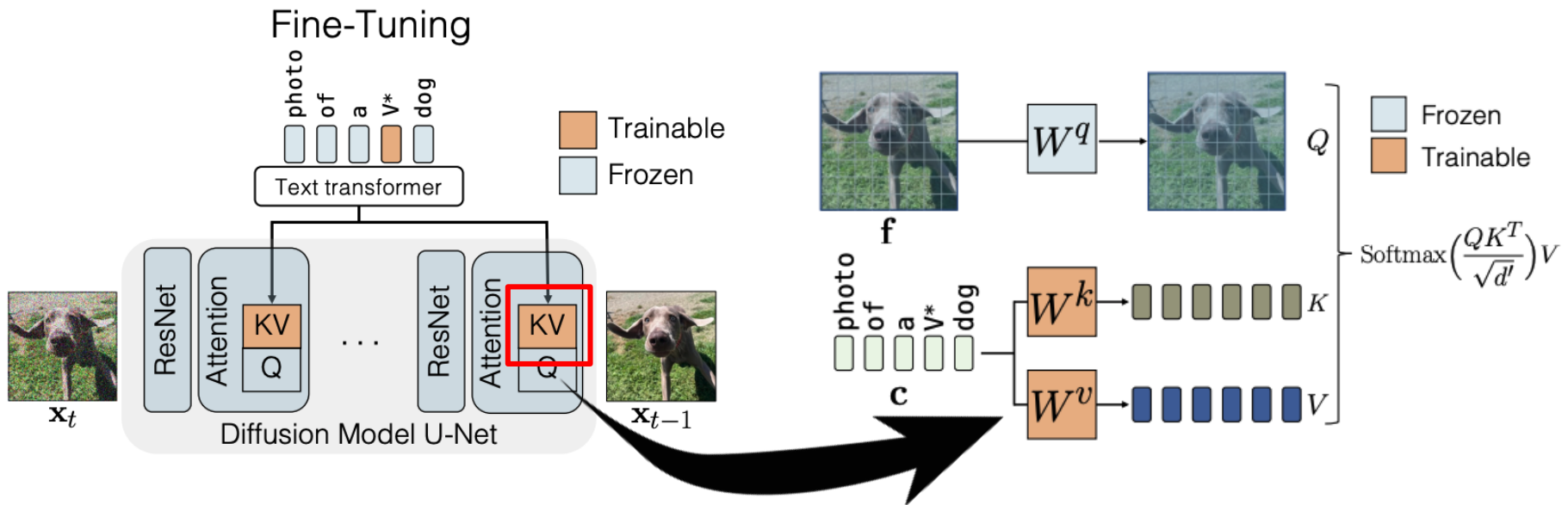
# Multi-concept Personalization: CustomDiffusion, CMU/Tsinghua/Adobe, CVPR'23

- Similar to Mix-of-Show, CustomDiffusion also has two stages:
  - Stage 1: Fine-tune for each concept separately
  - Stage 2: Merge different concepts into one



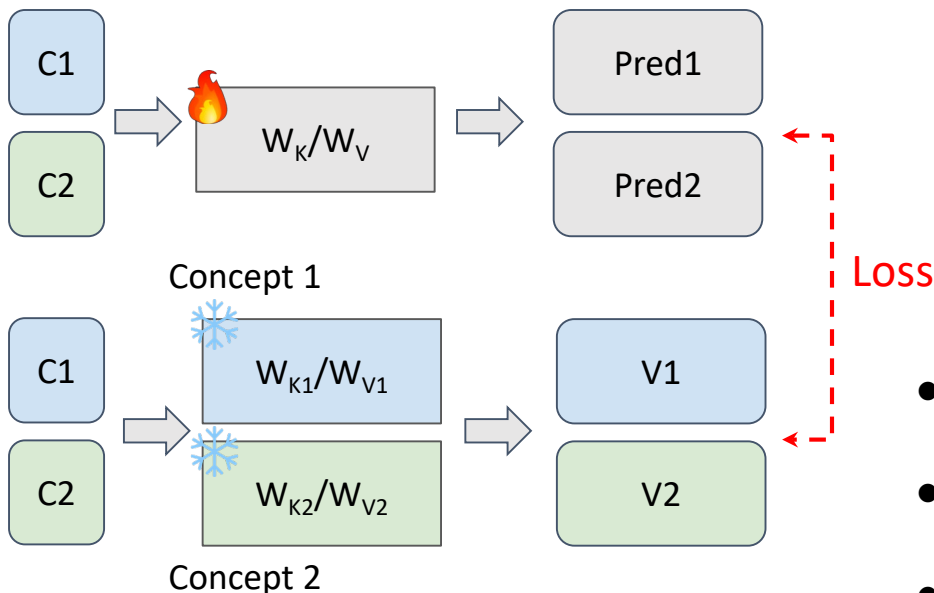
# Multi-concept Personalization: CustomDiffusion (cont'd)

- Similar to Mix-of-Show, CustomDiffusion also has two stages:
  - **Stage 1: Fine-tune  $V^*$  & cross-attention layer for each concept separately**
  - Stage 2: Merge different concepts into one



# Multi-concept Personalization: CustomDiffusion (cont'd)

- **Stage 1:** Fine-tune for each concept separately
- **Stage 2:** Merge different concepts into one
  - Fuse  $W_{K_i}$  &  $W_{V_i}$  of different concepts into a single  $W_K$  &  $W_V$  by minimizing the loss below
  - It's a closed-form solution since the  $W_K$  &  $W_V$  are linear matrices (via Lagrange multiplier)
    - ~4x faster than Mix-of-Show for merging multiple concepts



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

$$\text{s.t. } WC^T = V, \text{ where } C = [\mathbf{c}_1 \cdots \mathbf{c}_N]^T$$

$$\text{and } V = [W_1 \mathbf{c}_1^T \cdots W_N \mathbf{c}_N^T]^T.$$

- C & V are the input and output of **different concepts**
- $C_{\text{reg}}$  is text feature sampled from regularization data &  $W_0$  as pretrained model
- W is trained to fit all concepts.

# Diffusion Model for Personalization

- (Recap) Single Concept
  - Textual Inversion, ICLR'23
  - DreamBooth, CVPR'23
- Multiple Concepts
  - CustomDiffusion, CVPR'23
  - Mix-of-Show, NeurIPS'23
- Beyond Image: Video Motion Customization
  - Video Motion Customization, CVPR'24





# Video Motion Customization (VMC)

KAIST, CVPR'24

- **Task**

- **Given:** reference video + text prompt
- **Output:** video that matches  
(1) the motion of reference video & (2) the semantic of text prompt

Reference  
Video



Prompt

+  
"A duck"  
||

+  
"A robot"  
||

+  
"A dog"  
||

Result

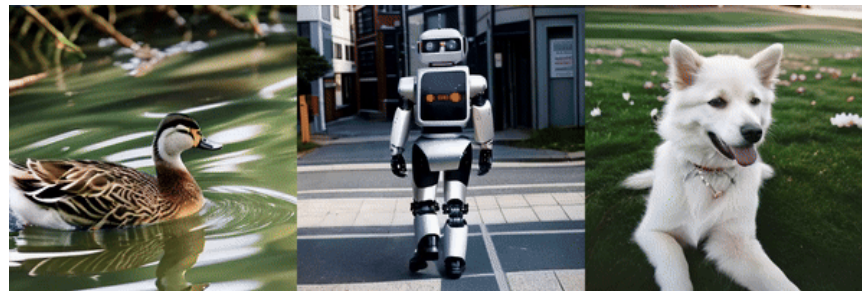
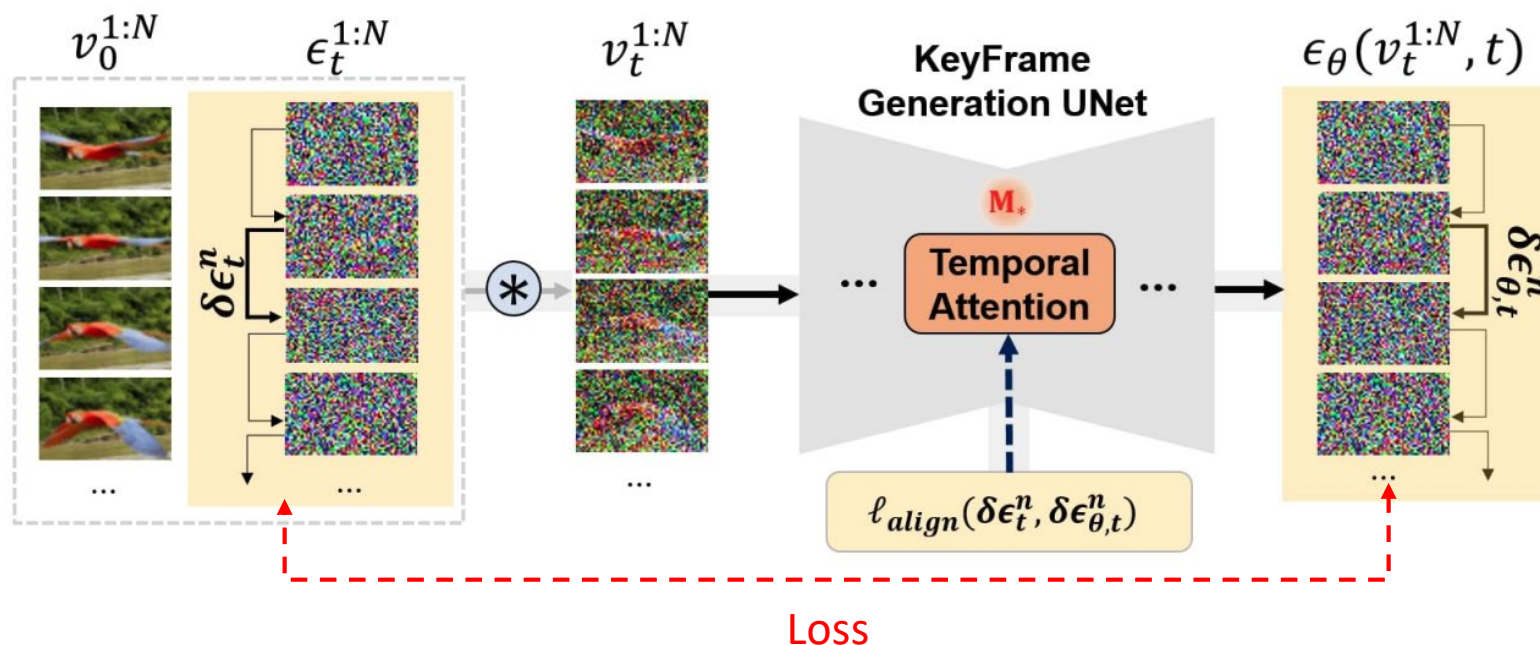


Figure Credit:  
[MotionClone](#)

# Video Motion Customization (VMC)

- **Training**

- Instead of calculating the MSE loss between noise prediction and GT, VMC calculates that btw the **noise residual** of prediction and GT
- To focus on learning temporal info (i.e., motion), only fine-tune the **temporal attention layer**

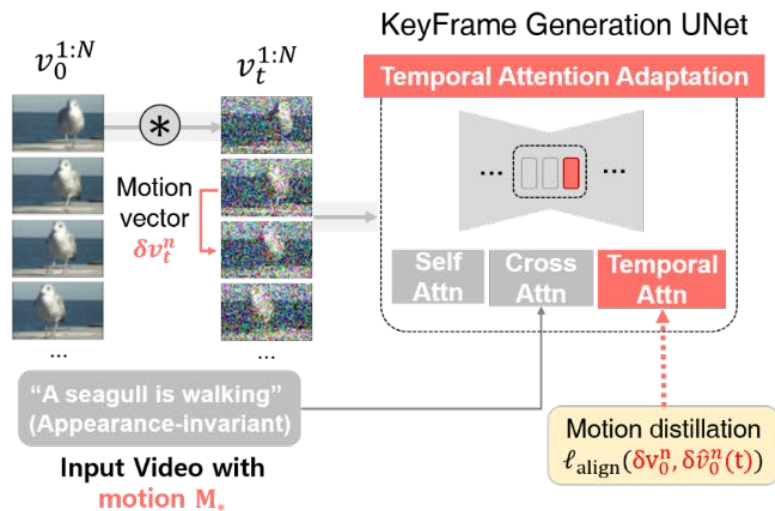


# Video Motion Customization (VMC)

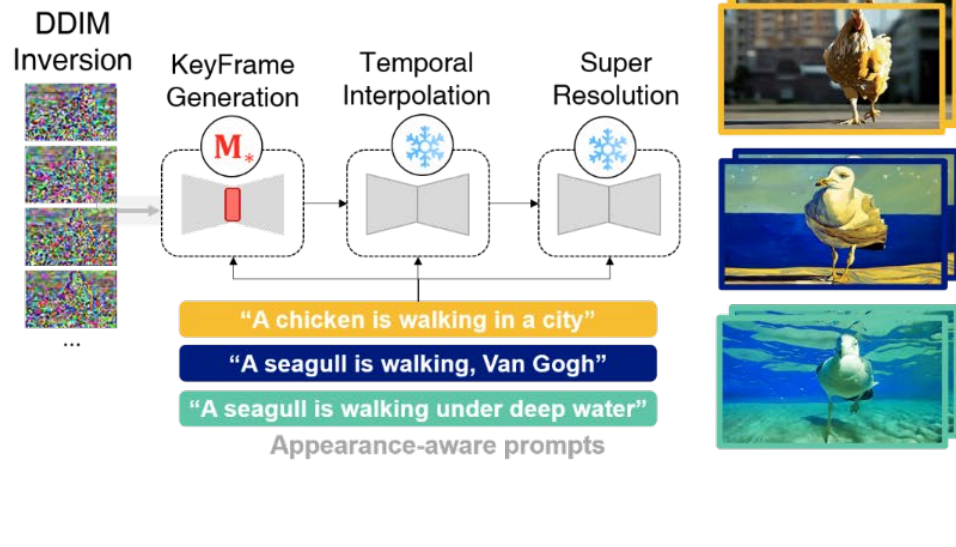
- Inference

- **Input:** DDIM inversed noise of reference video + text prompt
- **Output:** output video with desirable motion + appearance

(a) Training

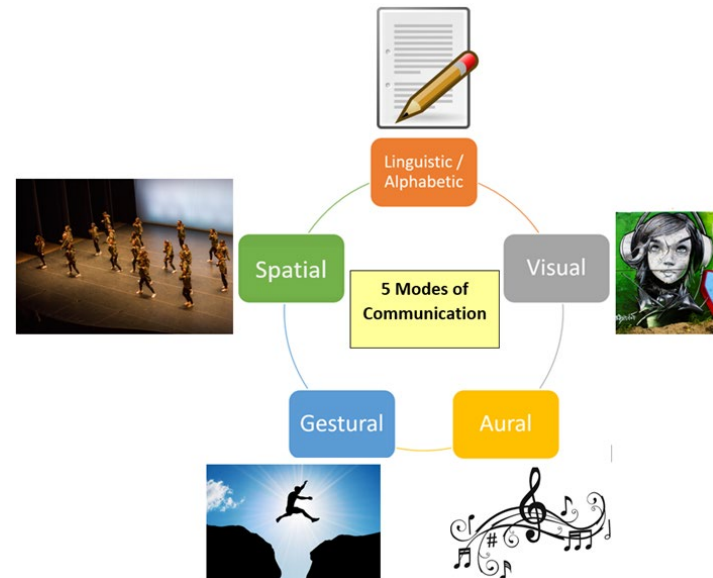
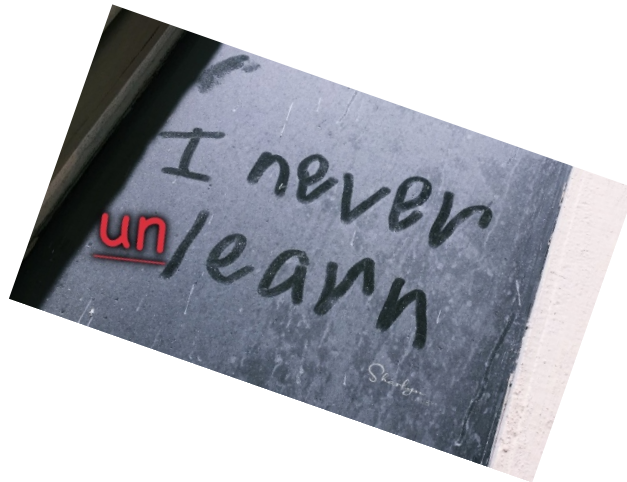
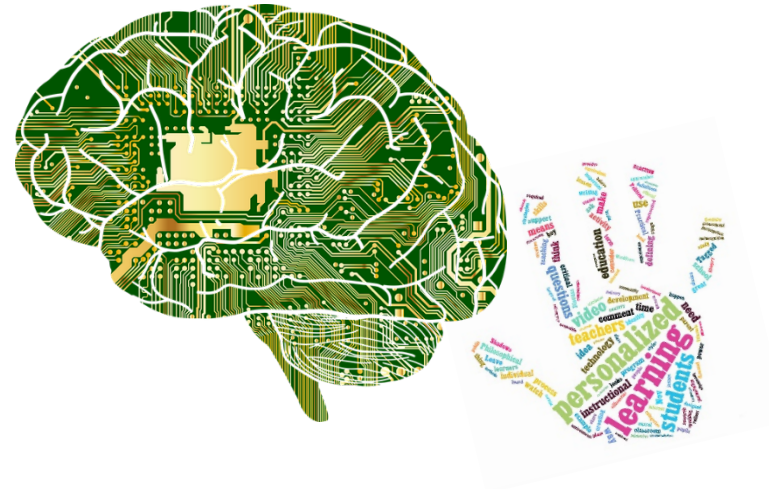


(b) Inference



# What to Be Covered Today...

- Multimodal LLM
- **Advanced Topics in LLM/VLM**
  - Concept Editing
  - Concept Unlearning
  - Personalization
  - Continual Learning

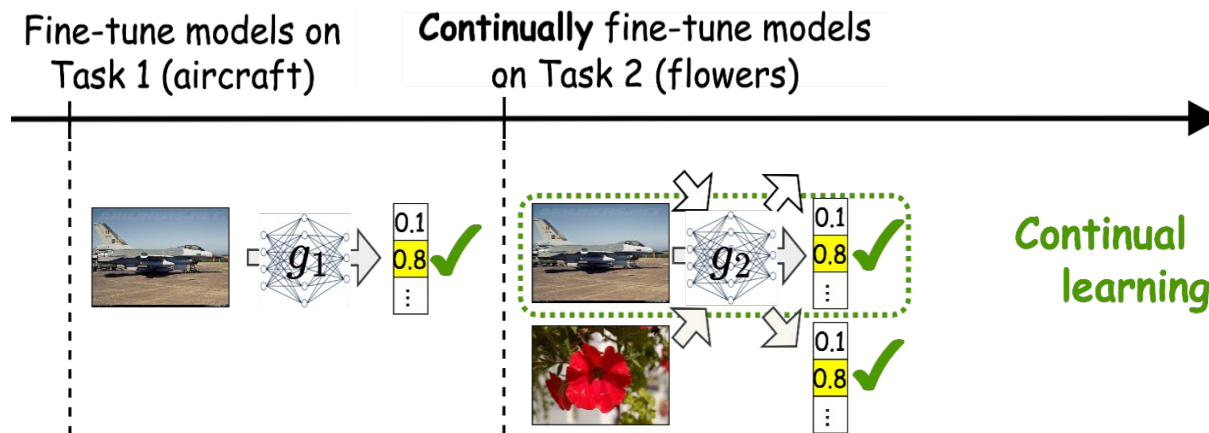


# Continual Learning (aka Incremental Learning)

- **Motivation**

- Always new dataset, knowledge, etc, to finetune the LLM/VLM
  - No practical to re-train foundation models from scratch
- It is a naive learning way, since **human is a continual learner**.

→ Goal: learn downstream tasks/datasets in a sequential (or incremental) way, while not forgetting what models have learned before.



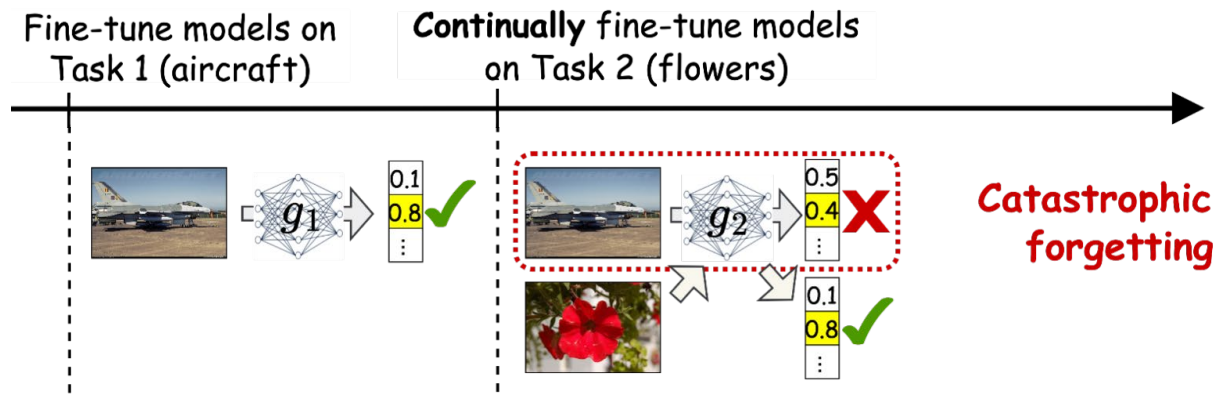
# Continual Learning (cont'd)

- **Task Definition**

- Learning a list of datasets in a sequential manner **without forgetting** previous knowledge.

- **The most straight forward strategy**

- Directly fine-tune a pre-trained model on a new dataset
- **Challenge:** Suffer from the well-known **catastrophic forgetting issue**, as the model weights can be totally distorted toward the new task only



# Previous works on Continual Learning

- **Rehearsal-based methods**
  - iCaRL (CVPR'17)
- **Regularization-based methods**
  - EWC (PNAS'17)
- **Continual Learning for open-vocab. Vision-Language Models**
  - ZSCL (ICCV'23)
  - Select and Distill (ECCV'24)

# iCaRL: Incremental Classifier and Representation Learning, Oxford, CVPR'17

- Rehearsal-based method
- Idea:
  - Maintain a subset of previous data in a class exemplar sets  $P = (P_1, \dots, P_{s-1})$  where  $\{1, 2, \dots, k-1\}$  are the learned classes
  - Joint training with the current data  $X^s, \dots, X^t$  with classes  $\{s, \dots, t\}$
- Method
  - For data in P, enforce the learned model  $\theta$  output as that of  $\theta_{\text{old}}$ .
    - Can be viewed as **Knowledge Distillation**
  - For the newer data, training with the standard cross entropy loss.

$$Y_{\text{old}} = \{f_{\theta_{\text{old}}}(x) | \forall x \in P\}$$

$$\mathcal{L}(\theta) = \sum_{(x,y) \in D} \left[ \sum_{y=s}^t \mathcal{L}(Y_{\text{new}}, \hat{Y}) + \sum_{y=1}^{s-1} \mathcal{L}(Y_{\text{old}}, \hat{Y}) \right]$$



# EWC: Overcoming catastrophic forgetting in neural networks, DeepMind, PNAS'17

- **Regularization-based method**

- **Idea:**

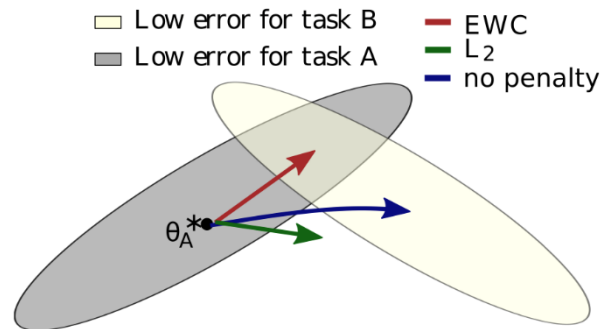
- Weight Consolidation: Restrict the updated weights **not to be too different** from the original model weights.

$$\mathcal{L}_{WC} = \sum_i (\theta_i - \bar{\theta}_i)^2$$

- Elastic Weight Consolidation: Each parameter should not be restricted with the same weights

- $i$ : the index of the model parameters.

$$\mathcal{L}_{EWC} = \sum w_i \cdot (\theta_i - \bar{\theta}_i)^2$$



# EWC, DeepMind, PNAS'17 (cont'd)

- **Method (cont'd)**

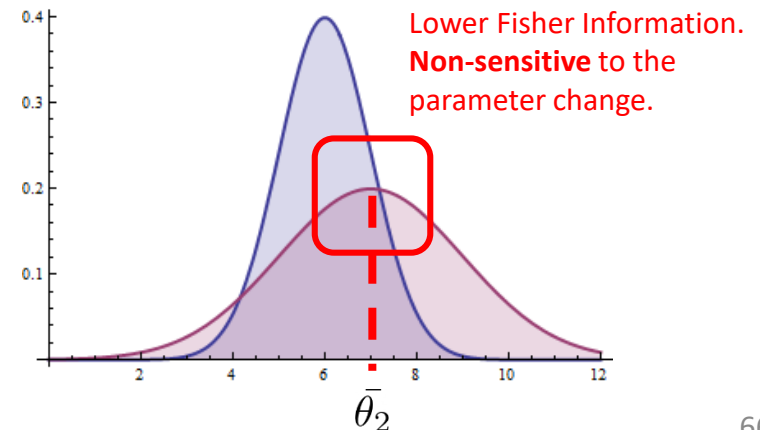
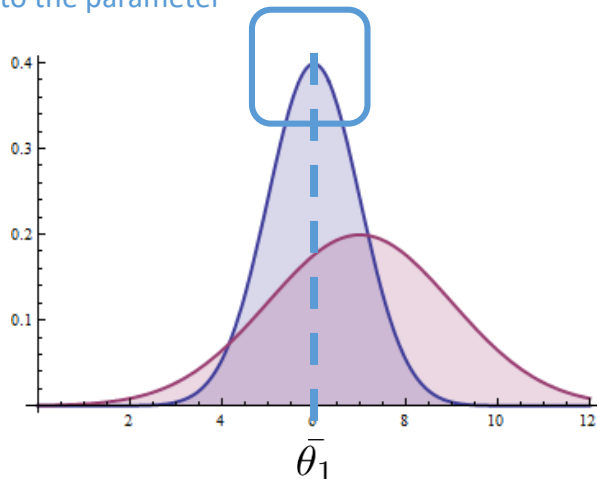
- Using Fisher Information (F) to determine the importance of a parameter to the previous task.
  - Fisher information: the expectation of second derivative of negative log-likelihood at  $\bar{\theta}$

$$\mathcal{L}_{\text{EWC}} = \sum_i \frac{\lambda}{2} F_i (\theta_i - \bar{\theta}_i)^2$$

- $\lambda$ : a hyper-parameter to determine the overall importance of previous tasks.

$$\mathcal{L}(\theta) = \mathcal{L}_B(\theta) + \sum_i \frac{\lambda}{2} F_i (\theta_i - \theta_{A,i}^*)^2$$

Higher Fisher Information.  
Sensitive to the parameter change



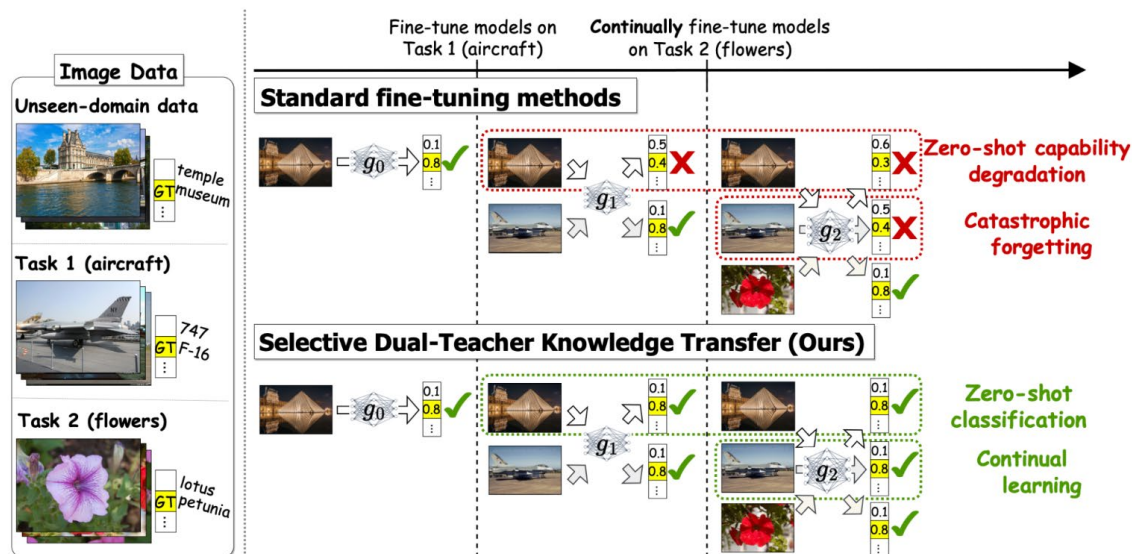
# Continual Learning for Open-Vocab Vision-Language Models

- **Motivation**

- With the prevalence of large-scale Vision-Language Models (VLMs), Continual Learning for VLMs has emerged as a potential research trends.

- **Goal**

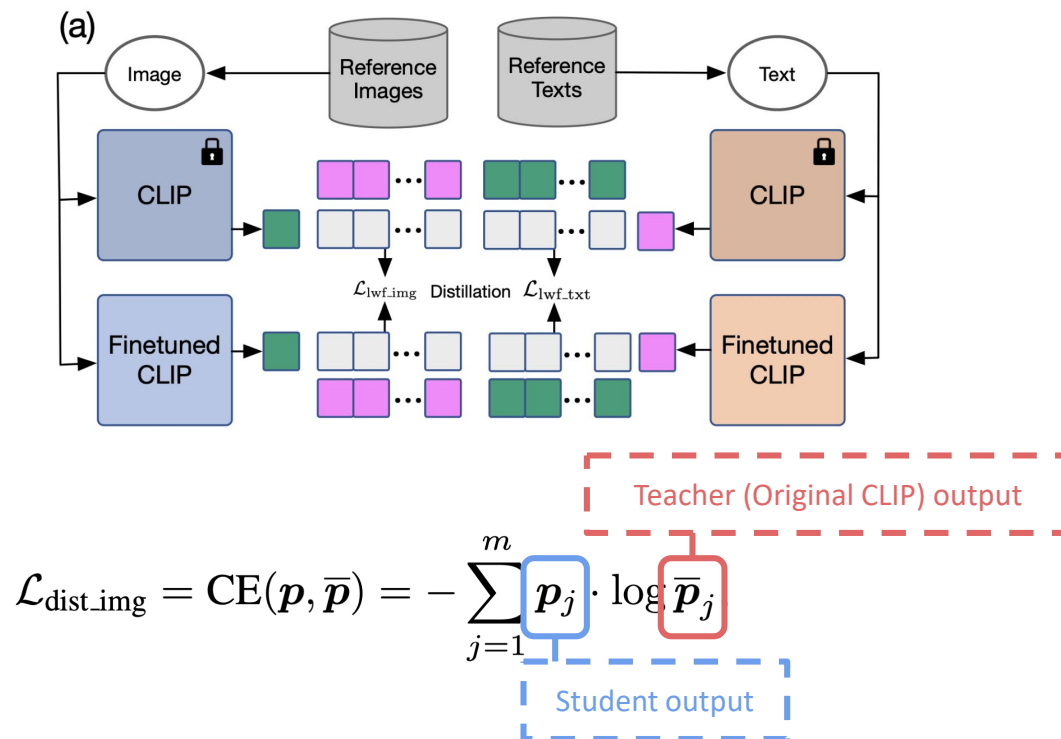
- Sequentially learning from new datasets
- Preserve the **original zero-shot ability for unseen data**
- Maintain **knowledge learned from previous stages** (as existing continual learning methods do)



# ZSCL: Preventing Zero-Shot transfer degradation in Continual Learning of vision-language models, NUS, ICCV'23

## ● Method

- Utilize an **auxiliary reference dataset** (e.g., ImageNet), and perform **Knowledge Distillation** from the **original CLIP model**.
  - (1) Distill knowledge on **both visual and textual sides**.



# ZSCL, NUS, ICCV'23 (cont'd)

- **Method (cont'd)**

- (2) WE: Weight space Ensemble to regularize the weights
  - The learned model weights would not be too different from the weights of the previous stage

$$\hat{\theta}_t = \begin{cases} \theta_0 & t = 0 \\ \frac{1}{t+1}\theta_t + \frac{t}{t+1} \cdot \hat{\theta}_{t-1} & \text{every } I \text{ iterations} \end{cases} .$$

- Same form as EMA (exponential moving average)
- Training strategy: (1) -> (2) -> (1) -> (2) -> ...

$$(1) \quad \mathcal{L} = \mathcal{L}_{ce} + \lambda \cdot (\mathcal{L}_{lwf\_img} + \mathcal{L}_{lwf\_txt})$$

$$(2) \quad \hat{\theta}_t = \begin{cases} \theta_0 & t = 0 \\ \frac{1}{t+1}\theta_t + \frac{t}{t+1} \cdot \hat{\theta}_{t-1} & \text{every } I \text{ iterations} \end{cases} .$$

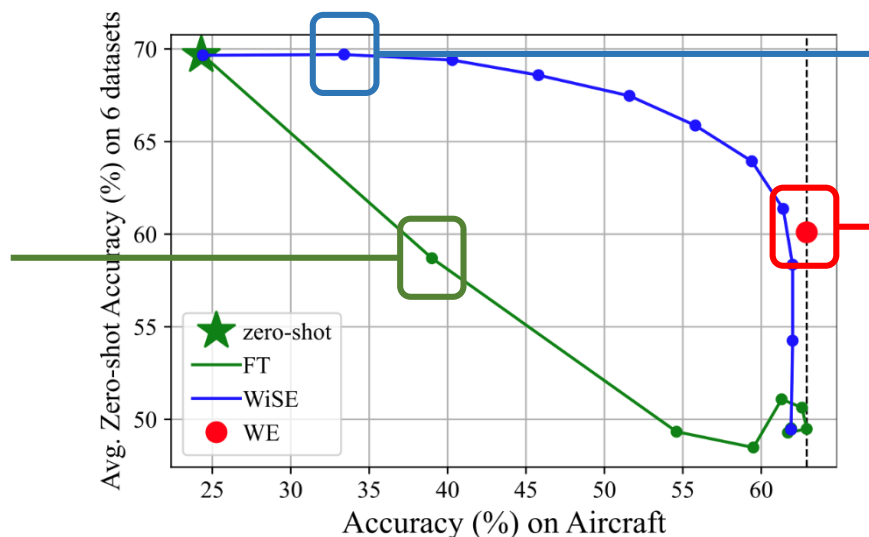
← Data from novel task + auxiliary ref dataset

# ZSCL, NUS, ICCV'23 (cont'd)

- **Comparisons**

- Zero-shot accuracy vs. accuracy on novel task

Sample every 100 iterations.  
As training progresses, the model's zero-shot capability deteriorates.



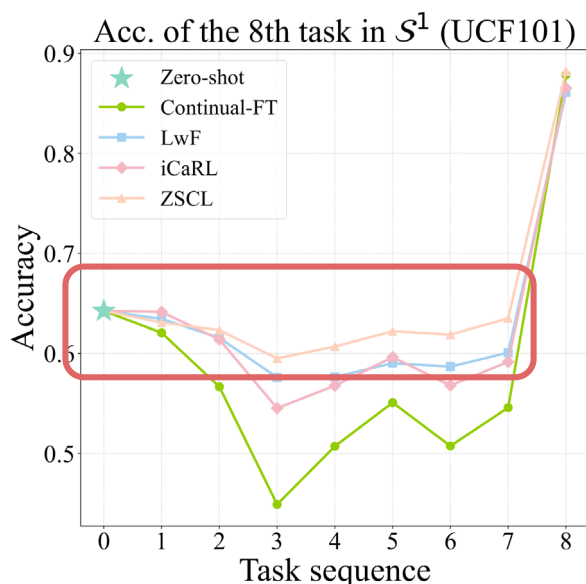
Ensemble the original and fine-tuned model weights with different ratios.  
**Sensitive to the choose of the ratio**

Iterative weight ensemble. Improved fine-tuned accuracy, with an acceptable decrease in zero-shot performance

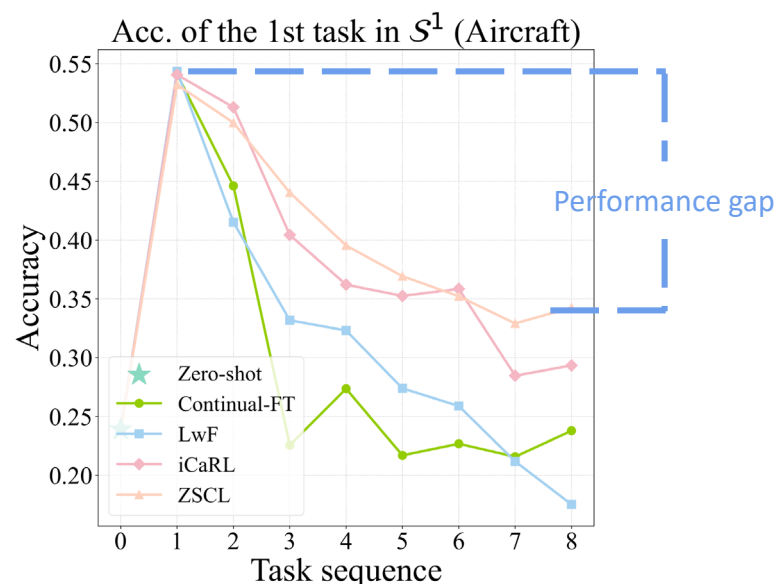
# ZSCL, NUS, ICCV'23 (cont'd)

- **Limitation**

- ZSCL still largely suffer from catastrophic forgetting for previous tasks.



ZSCL can preserve zero-shot ability for unseen data

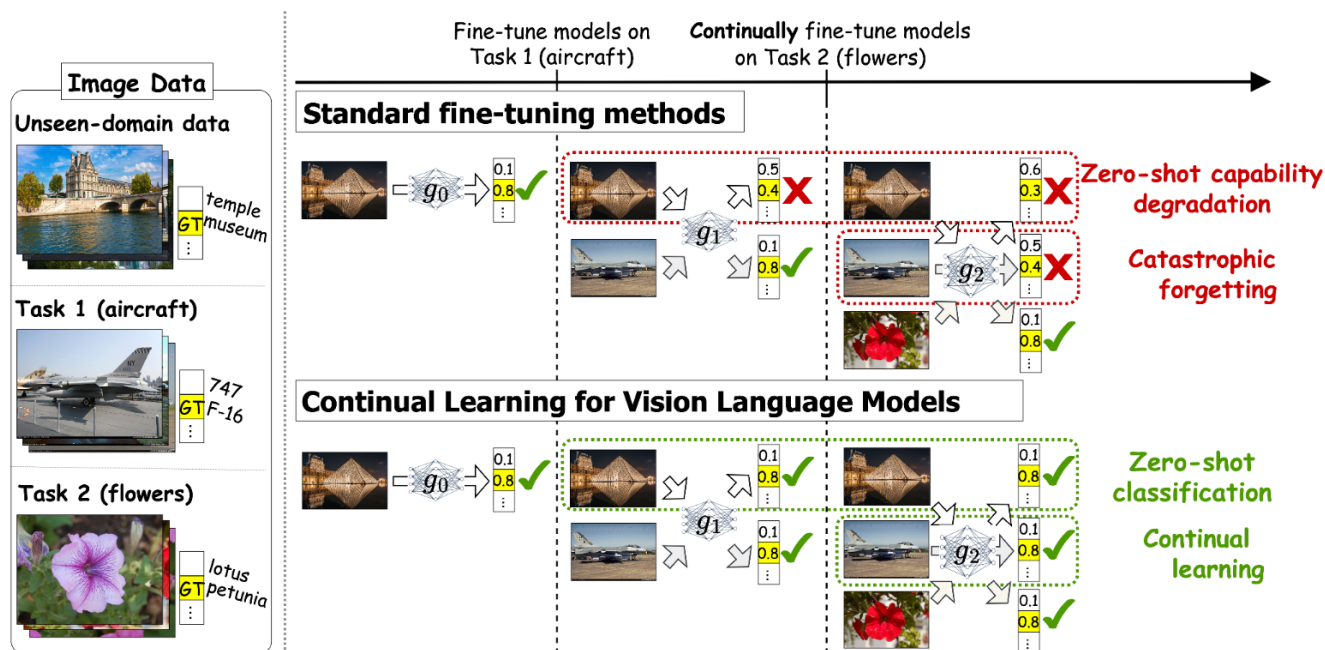


There is still a gap for previous task after training on multiple datasets

# Select and Distill: Selective Dual-Teacher Knowledge Transfer for Continual Learning on Vision-Language Models, NTU, ECCV'24

- **Goal**

- Same as ZSCK, adapt to new datasets sequentially while:
  - preserving the **original pre-trained zero-shot ability**
  - maintaining the **knowledge learned from previous stages**.





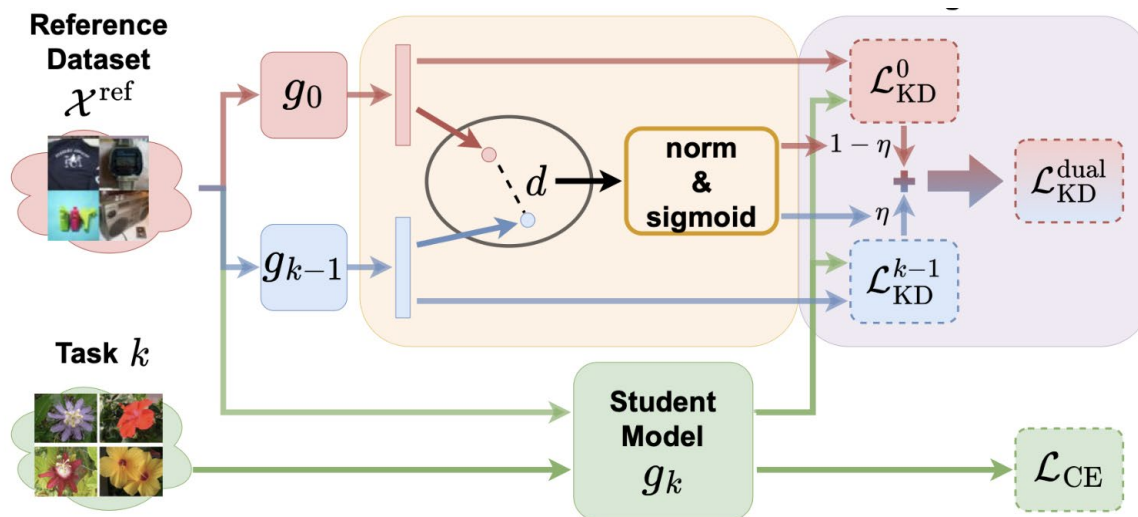
# Select and Distill, NTU, ECCV'24 (cont'd)

- Idea

- Follow ZSCL, utilize a reference dataset for knowledge distillation
- Dual-Teacher Knowledge Distillation
  - Distill from **the original pre-trained VLM** to preserve **zero-shot ability**.
  - Distill from **the most recent fine-tuned VLM** to preserve **prior knowledge**.

- Key

- For any data point in the reference dataset, we need to **select a proper model** and distill its knowledge.

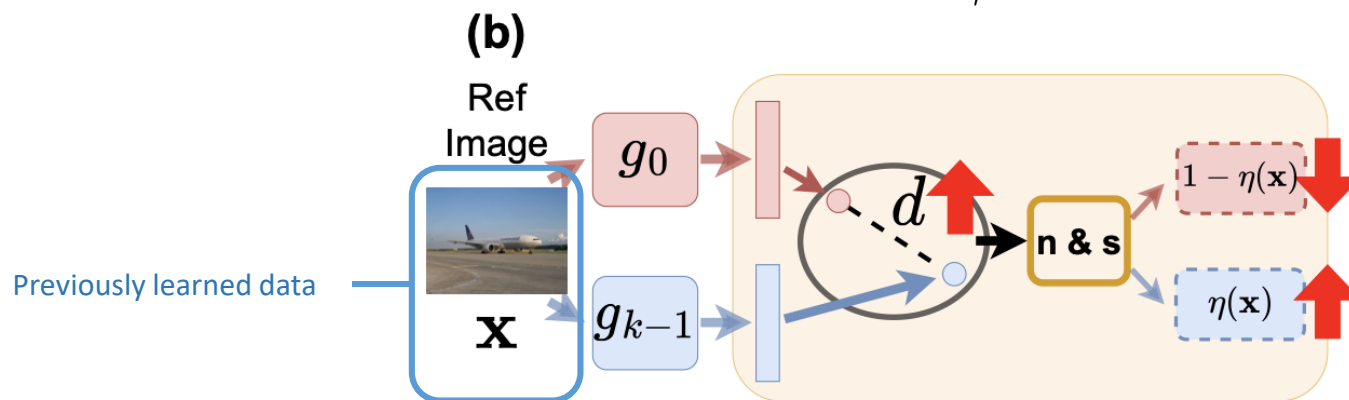


# Select and Distill, NTU, ECCV'24 (cont'd)

- **Observation**

- If a data point is a **previously learned data**.
  - It must be **seen by  $g_{k-1}$** , but **never been seen by  $g_0$**
  - The **feature distance  $d$**  between  $g_0$  and  $g_{k-1}$  can be large or small?
  - Select  $g_{k-1}$  as the teacher model to **maintain previous knowledge**
- $\eta(\mathbf{x})$  : A normalized distance between 0~1, determine how much should we distill from  $g_{k-1}$

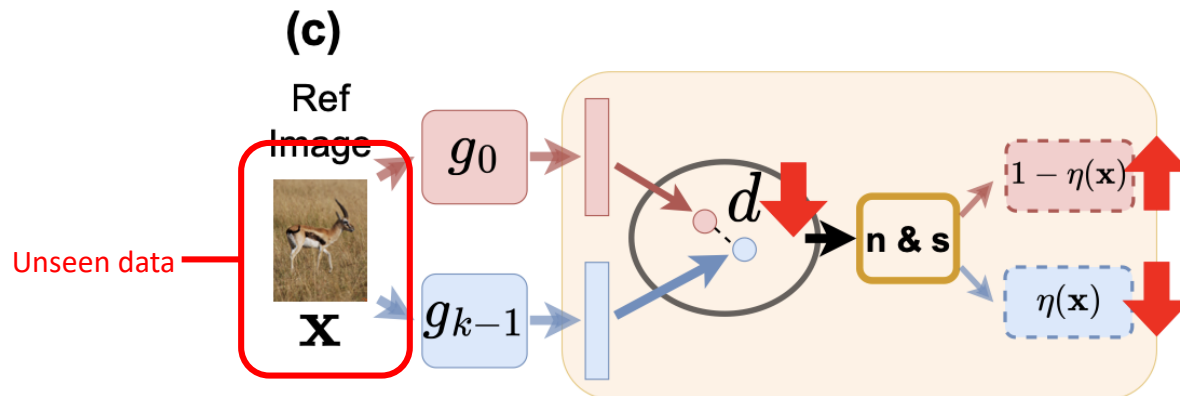
$$\eta(\mathbf{x}) = \sigma\left(\frac{d(g_{k-1}(\mathbf{x}), g_0(\mathbf{x})) - \delta}{\gamma}\right),$$



# Select and Distill, NTU, ECCV'24 (cont'd)

- **Observation**

- If a data point has never been seen by both  $g_{k-1}$  and  $g_0$  (**unseen data**)
  - The **feature distance  $d$**  between  $g_0$  and  $g_{k-1}$  can be **relatively small**
  - In this case, we should select  $g_0$  as the teacher model to **preserve the original zero-shot ability**.

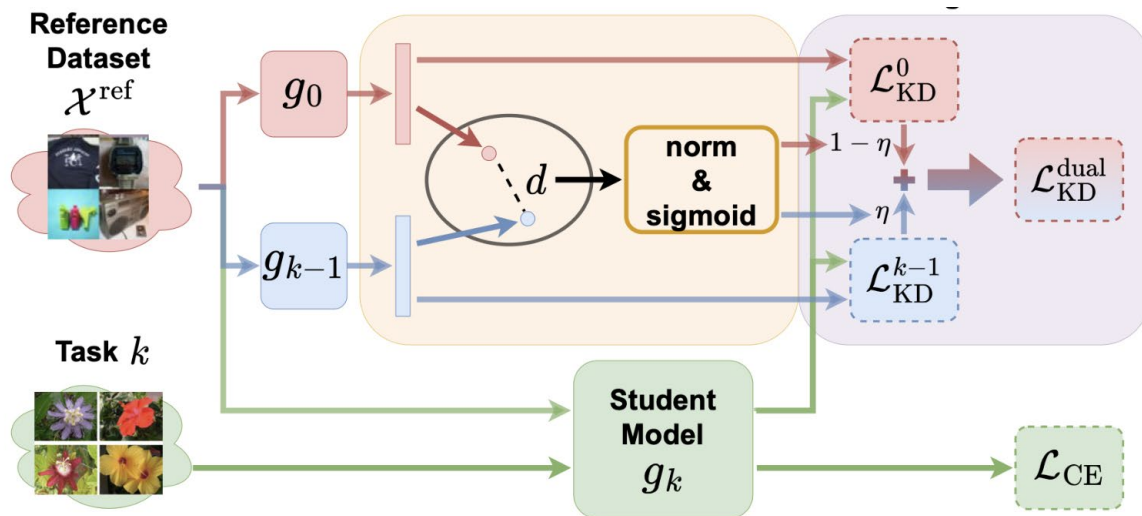


# Select and Distill, NTU, ECCV'24 (cont'd)

- Objective

$$\mathcal{L}_{\text{KD}}^{k-1} = d(g_{k-1}(\mathbf{x}), g_k(\mathbf{x})) \quad , \quad \mathcal{L}_{\text{KD}}^0 = d(g_0(\mathbf{x}), g_k(\mathbf{x}))$$

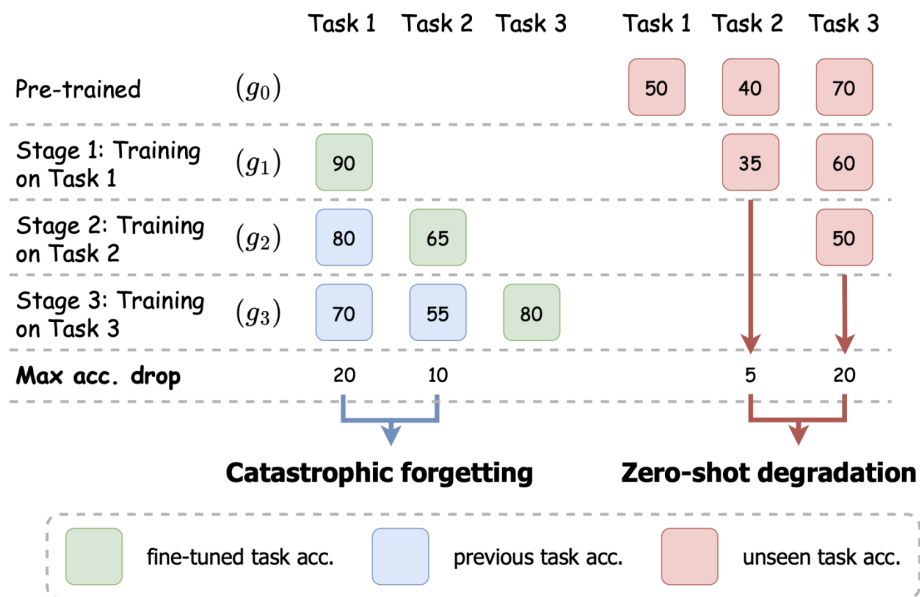
$$\mathcal{L}_{\text{KD}}^{\text{dual}} = \sum_{\mathbf{x} \sim \mathcal{X}^{\text{ref}}} \eta(\mathbf{x}) \cdot \mathcal{L}_{\text{KD}}^{k-1} + (1 - \eta(\mathbf{x})) \cdot \mathcal{L}_{\text{KD}}^0$$



# Select and Distill, NTU, ECCV'24 (cont'd)

- **Metrics**

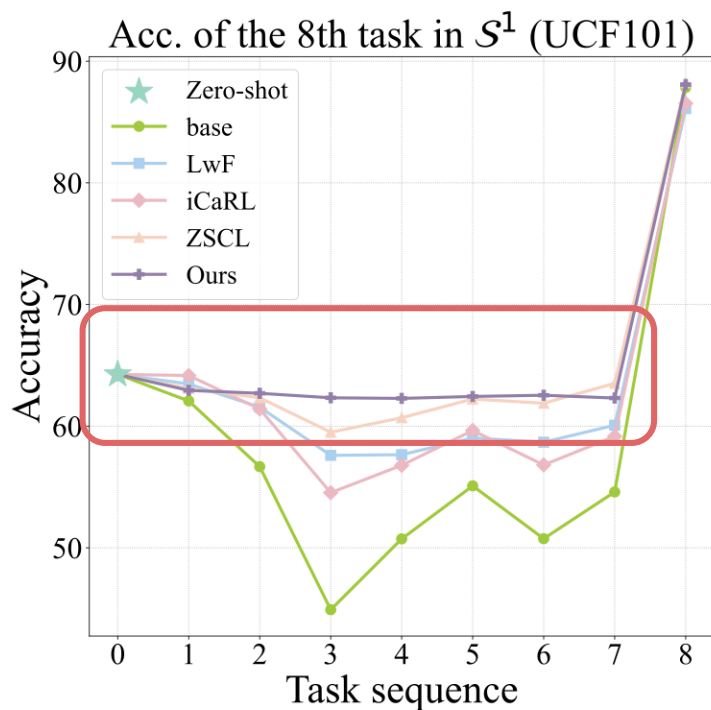
- Average Accuracy
  - Average of the last performance on each dataset
- Catastrophic forgetting
  - Max. performance gap after the task has been fine-tuned
- Zero-shot degradation
  - Max. performance gap before the task has been fine-tuned



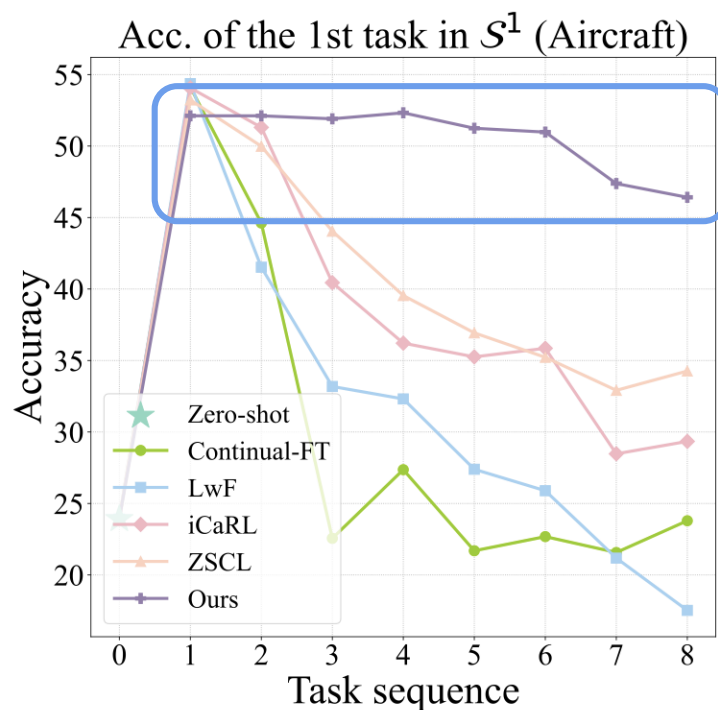
# Select and Distill, NTU, ECCV'24 (cont'd)

## ● Results

- Successfully preserve the **zero-shot ability for unseen data**
- Mitigate the **catastrophic forgetting of previously learned data**



Successfully preserve zero-shot ability for unseen data



Largely mitigate the performance gap

# Select and Distill, NTU, ECCV'24 (cont'd)

## ● Robustness

- We shuffle the training orders, producing 8 different training sequences.
- Our methods showing state-of-the-art performance on all metrics, and the results are stable across all training sequences.

Method / Sequence	$\mathcal{S}^1$	$\mathcal{S}^2$	$\mathcal{S}^3$	$\mathcal{S}^4$	$\mathcal{S}^5$	$\mathcal{S}^6$	$\mathcal{S}^7$	$\mathcal{S}^8$	Mean
<b>Catastrophic forgetting (<math>\downarrow</math>)</b>									
Continual FT	10.98	10.60	8.80	19.17	10.11	11.95	15.19	9.48	12.04
LwF [24]	10.38	6.52	6.37	10.22	7.99	7.70	10.41	8.91	8.56
iCaRL [35]	8.42	7.00	6.45	10.21	7.03	7.33	9.68	8.23	8.04
ZSCL [50]	4.67	2.35	2.13	2.97	3.15	4.28	4.89	4.70	3.64
MoE-Adapters [48]	2.74	4.71	4.28	1.15	1.50	1.60	2.94	2.77	2.71
Ours	<b>1.70</b>	<b>1.16</b>	<b>0.89</b>	<b>1.04</b>	<b>0.59</b>	<b>1.34</b>	<b>1.12</b>	<b>1.79</b>	<b>1.20</b>
<b>Zero-shot degradation (<math>\downarrow</math>)</b>									
Continual FT	24.81	23.58	19.54	16.46	22.22	19.02	19.54	24.02	21.15
LwF [24]	10.75	10.23	8.63	8.25	12.02	10.33	8.98	11.01	10.03
iCaRL [35]	13.77	12.68	11.28	12.14	13.20	13.20	13.09	14.01	12.92
ZSCL [50]	3.44	3.94	4.02	2.85	3.79	2.31	1.86	1.84	3.00
MoE-Adapters [48]	1.62	2.58	<b>1.04</b>	2.37	4.31	3.05	<b>1.77</b>	<b>0.63</b>	2.17
Ours	<b>1.55</b>	<b>2.04</b>	1.21	<b>1.92</b>	<b>2.79</b>	<b>2.18</b>	1.90	2.08	<b>1.96</b>
<b>Average accuracy (<math>\uparrow</math>)</b>									
Continual FT	76.16	76.24	78.03	68.69	76.64	75.44	72.71	77.45	75.17
LwF [24]	76.78	80.45	80.65	77.52	79.64	79.45	77.31	78.70	78.81
iCaRL [35]	77.99	79.77	79.93	76.66	79.26	79.08	77.06	78.61	78.55
ZSCL [50]	81.89	83.98	84.30	83.49	83.41	82.38	81.92	81.97	82.92
MoE-Adapters [48]	82.71	80.74	81.15	83.97	83.68	83.68	82.73	79.68	82.29
Ours	<b>84.48</b>	<b>84.92</b>	<b>84.97</b>	<b>84.89</b>	<b>85.50</b>	<b>85.07</b>	<b>85.02</b>	<b>84.52</b>	<b>84.92</b>

# What We Have Covered Today...

- **Multimodal LLM**
- **Advanced Topics in LLM/VLM**
  - Concept Editing
  - Concept Unlearning
  - Personalization
  - Continual Learning

