

# Large Scale Vision-Language Pretraining: Models and Applications

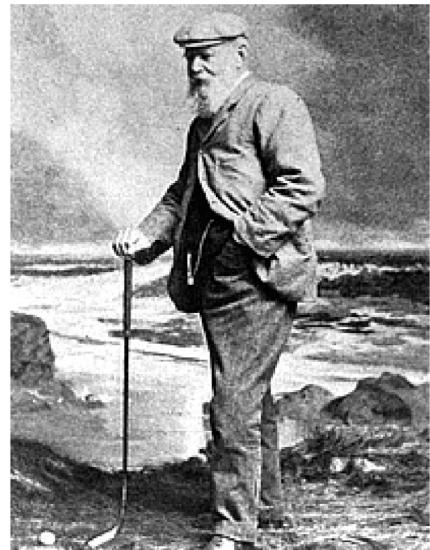
Wangmeng Zuo

Centre on Machine Learning Research  
Harbin Institute of Technology

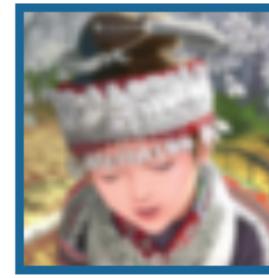
# Vision-Language Pretraining

- Vision-Language Models
- Classification: Parameter-efficient Fine-tuning
- More Challenging Tasks
- More Visual Modalities

# 典型视觉学习任务（底层视觉）



去噪

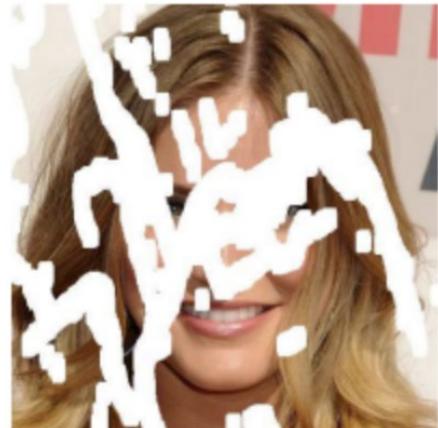


超分辨



风格迁移

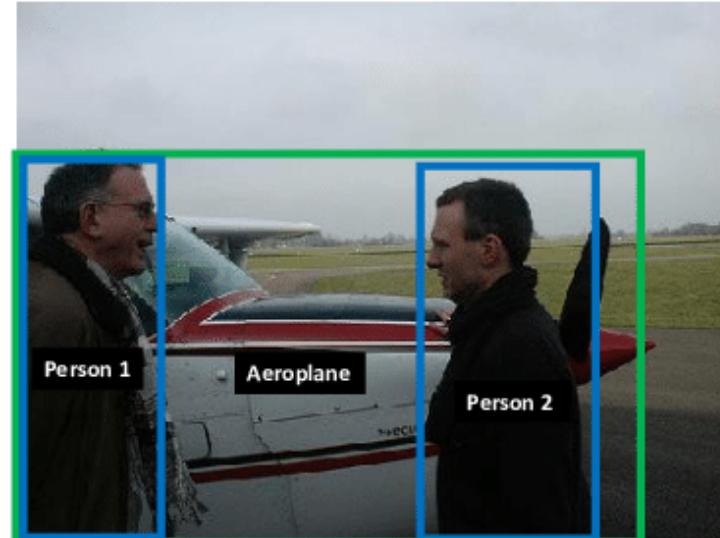
修复



# 典型视觉学习任务（视觉理解）



图像级分类



边界框级检测



物体关系预测



像素级分割

# 典型视觉任务（视觉-语言）

自然语言描述



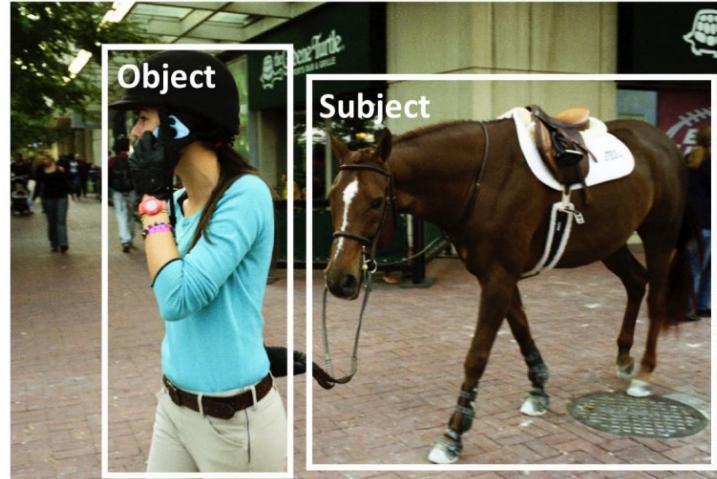
"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



指代消解

视觉问答



Q: Does this foundation have any sunscreen?  
A: yes



Q: What is this?  
A: 10 euros



Q: What color is this?  
A: green



Q: Please can you tell me what this item is?  
A: butternut squash red pepper soup



Q: Is it sunny outside?  
A: yes

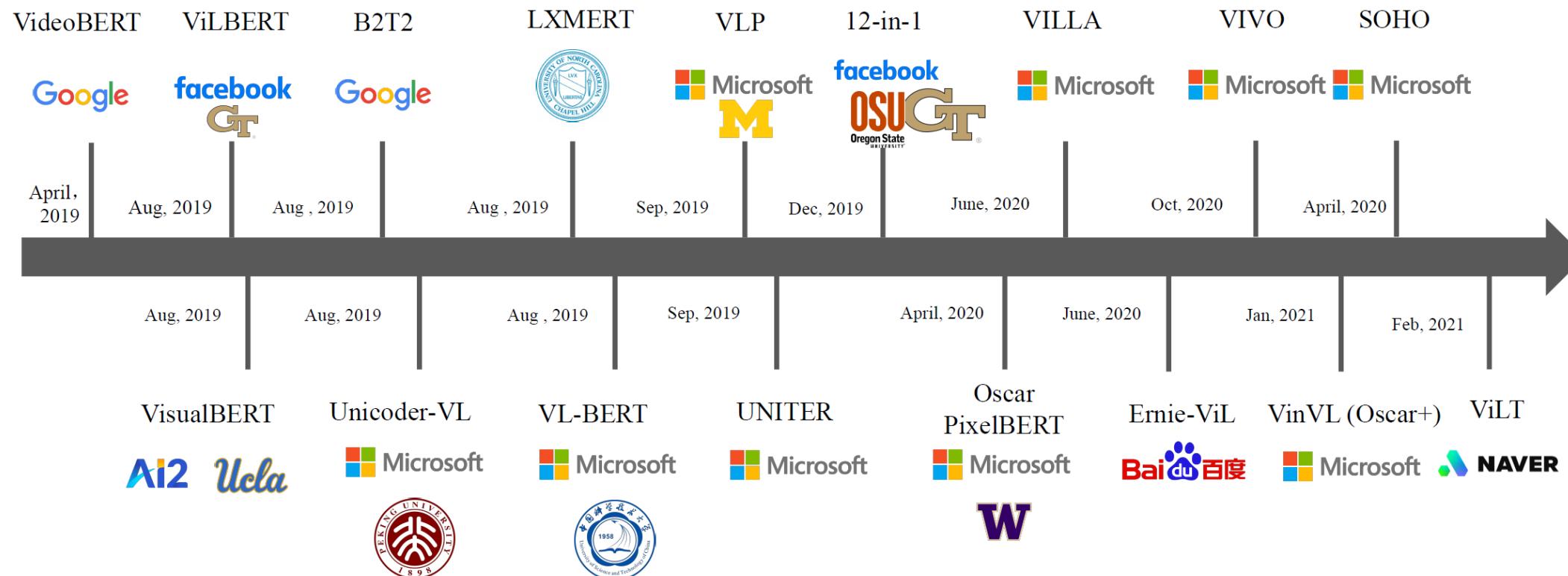


Q: Is this air conditioner on fan, dehumidifier, or air conditioning?  
A: air conditioning

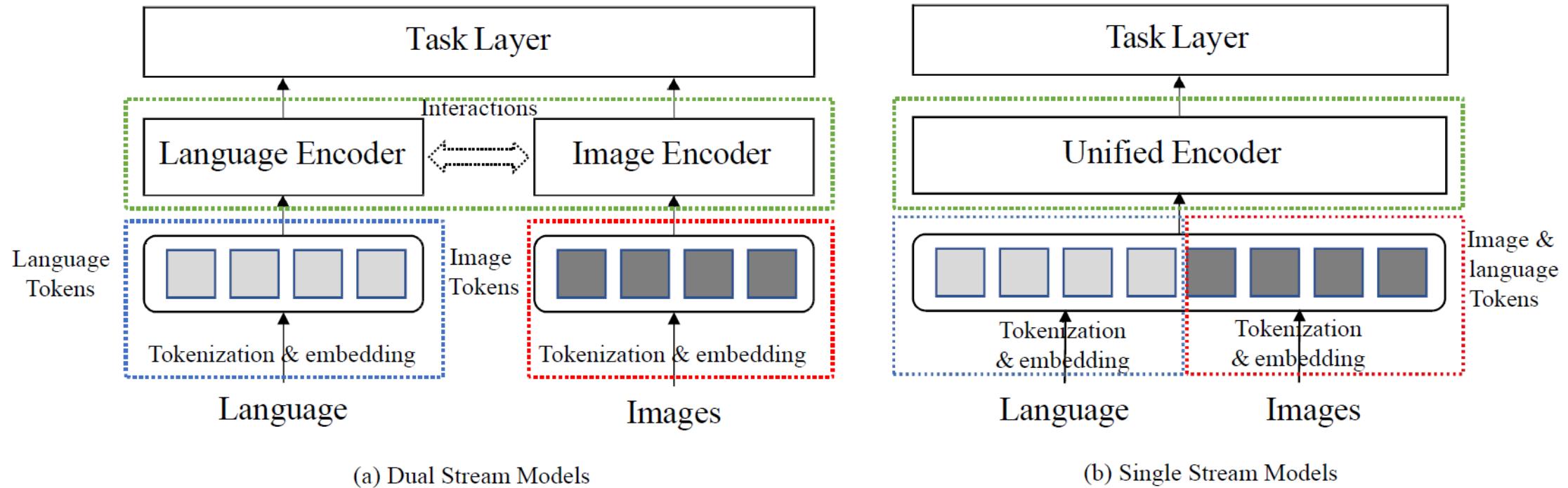
# Vision-Language Pretraining

- Early VL Pretraining
- Recent Breakthrough

# Landscape of V-L Pretraining Methods



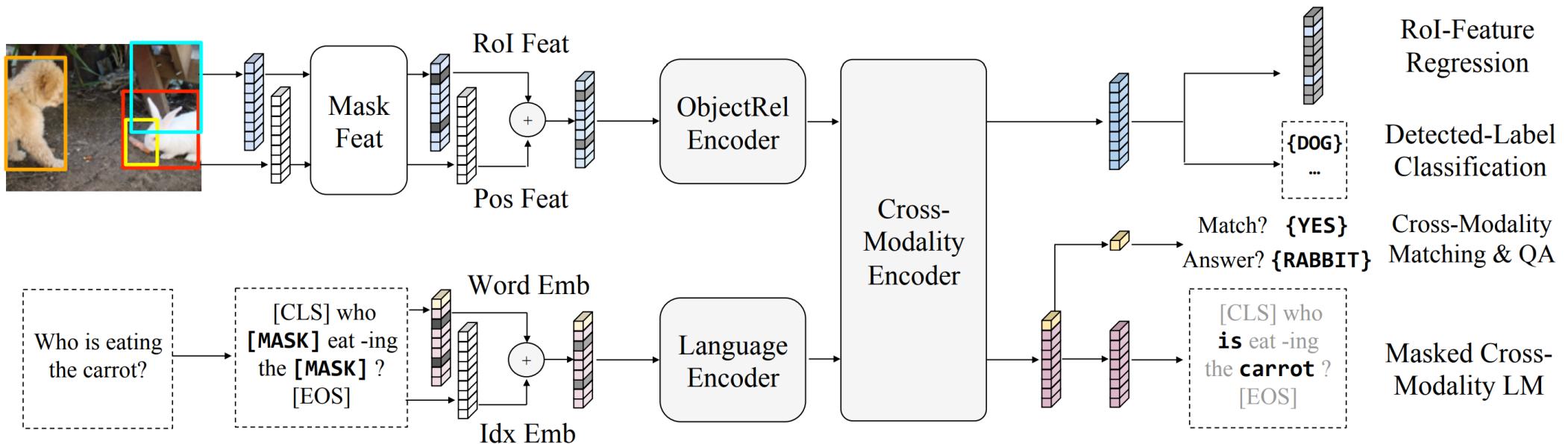
# Architectures of VLP Models



# Two-stream Architectures

|             |                  |      |   |  |                                 |
|-------------|------------------|------|---|--|---------------------------------|
| Dual-Stream | ViLBERT (2019)   | BUTD | ITM, MLM<br>MRC-kl  | CC3M   | VQA, VR, RE<br>IR, Zero-shot IR |
|             | LXMERT (2019)    | BUTD | ITM, MLM, MRP, MRC  | COCO, VG (2017a), VQA v2<br>GQA, Visual7W (2016)   | VQA, VR                         |
|             | 12-in-1 (2020)   | BUTD | ITM,MLM,MRC-kl  | VQA v2, Flickr30k<br>SNLI-VE, COCO<br>GuessWhat, VG<br>RefCOCO, RefCOCO+,RefCOCOG<br>Visual 7W, GQA, NLVR <sup>2</sup> | VQA ,IR, RE, VE, VR             |
|             | Ernie-ViL (2020) | BUTD | Object Prediction<br>Attribute Prediction<br>Relationship Prediction<br>ITM, MLM,MRC-kl | CC3M, SBU (out-of-domain)<br>COCO, VG (in-domain)  | VQA, VR, IT, TR, RE             |

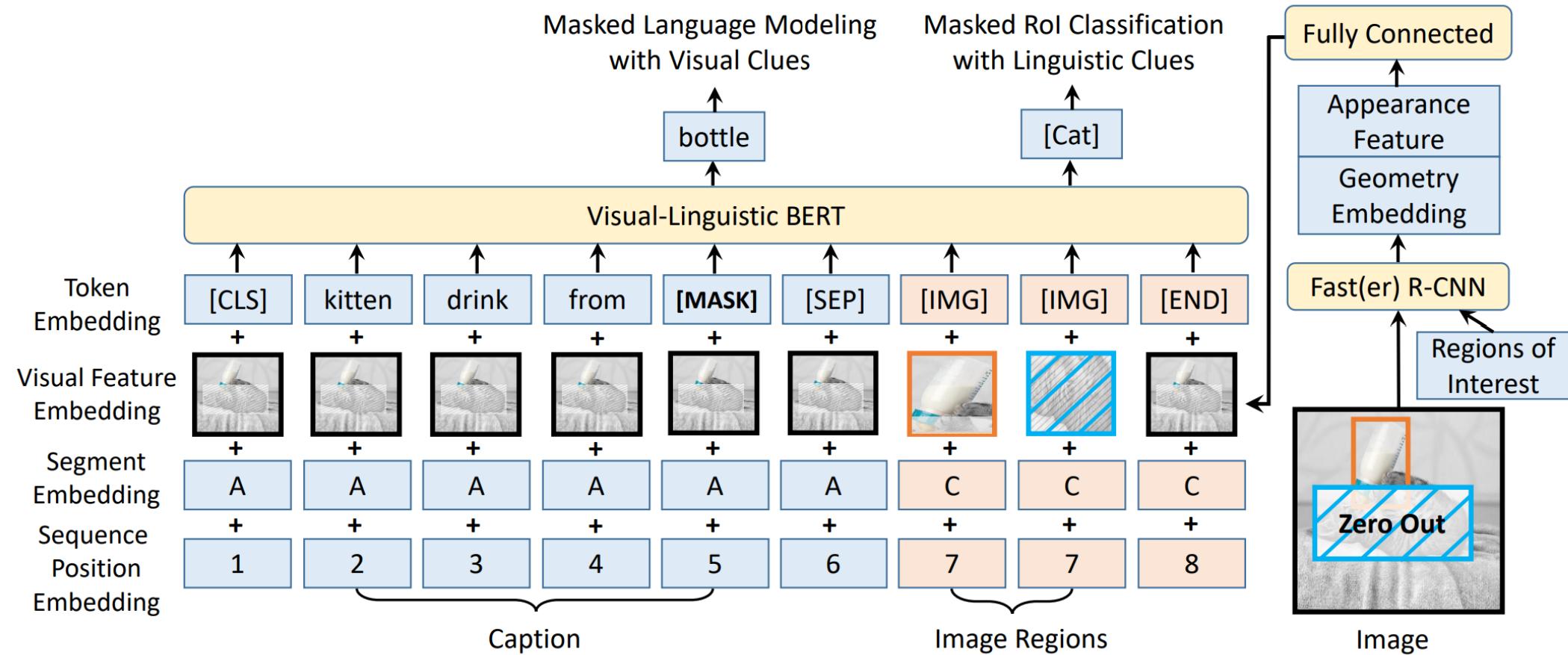
# LXMERT Training



# One-stream Architectures

|               |                     |  |  |  |   |
|---------------|---------------------|--|--|--|---|
| Single-Stream | VideoBERT (2019a)   | S3D (2017b)                                      | ITM, MLM, MVM                              | YouTube cooking videos (2017b)             | Zero-shot Action prediction (2017b)<br>Video Captioning (2017b) |
|               | VisualBERT (2019b)  | Pre-trained Fast R-CNN                           | ITM, MLM                                   | COCO                                       | VQA, VR, PG   |
|               | B2T2 (2019)         | ResNet-152 (2016b)                               | ITM, MLM                                   | CC3M                                       | VR  |
|               | Unicoder-VL (2019a) | Pre-trained Faster-RCNN (2018)                   | ITM, MLM, MRC                              | CC3M, SBU                                  | IR, TR, VR<br>Zero-shot IR/TR                                   |
|               | VL-BERT (2019)      | BUTD   | MLM, MRP                                   | CC3M, English Wikipedia BooksCorpus (2015) | VQA, VR, RE   |
|               | Unified VLP (2020)  | BUTD variant (with ResNext-101 backbone) (2017a) | MLM (Sequentially and bidirectionally)     | CC3M                                       | IC, VQA   |
|               | UNITER (2019)       | BUTD   | ITM, MLM, MRP, MRC MRC-kl                  | COCO, VG, CC3M, SBU                        | VQA, VR, VE, IR, TR, RE   |
|               | Oscar (2020b)       | BUTD+tags  | MLM (include tags)<br>ITM (pollute tags)   | COCO, VG, CC3M, SBU, flicker30k, GQA       | VQA, VR, VE, IR, TR<br>RE, IC                                   |
|               | PixelBERT (2020)    | Pixel feature embedding (2015)                   | MLM, ITM                                   | COCO, VG                                   | VQA, IR, TR, VR   |
|               | VILLA (2020)        | BUTD   | ITM, MLM, MRC-kl                           | COCO, VG, CC3M, SBU                        | VQA, VR, VE, IR, TR, RE   |
|               | VIVO (2020)         | BUTD   | Mask Tag Prediction (Hungarian match loss) | Open Images V5 (2018; 2019)                | IC  |
|               | SOHO (2021)         | VD   | ITM, MLM<br>MVMVD                          | COCO, VG                                   | IR, TR, VQA<br>VR, VE (based on VD)                             |
|               | ViLT (2021)         | Patch Projection                                 | ITM, MLM                                   | COCO, VG, SBU, CC                          | VQA, VR,<br>IR, TR,<br>Zero-shot IR&TR                          |
| Hybrid        | SemVLP (2021a)      | Pre-tained Faster-RCNN                           | MLM, MRP, VQA                              | COCO, VG, VQA v2, GQA, Visual7W            | VQA, VR, IR, TR   |

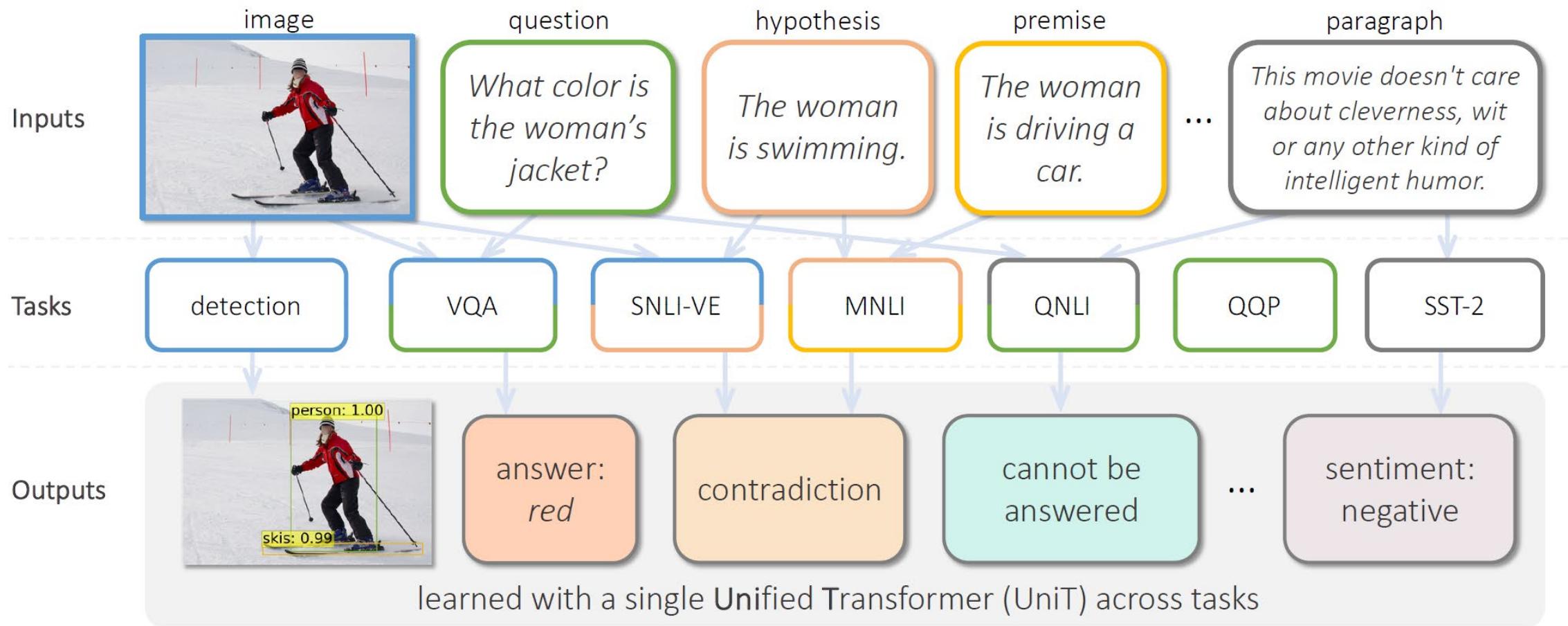
# VL-BERT



# Vision-Language Pretraining

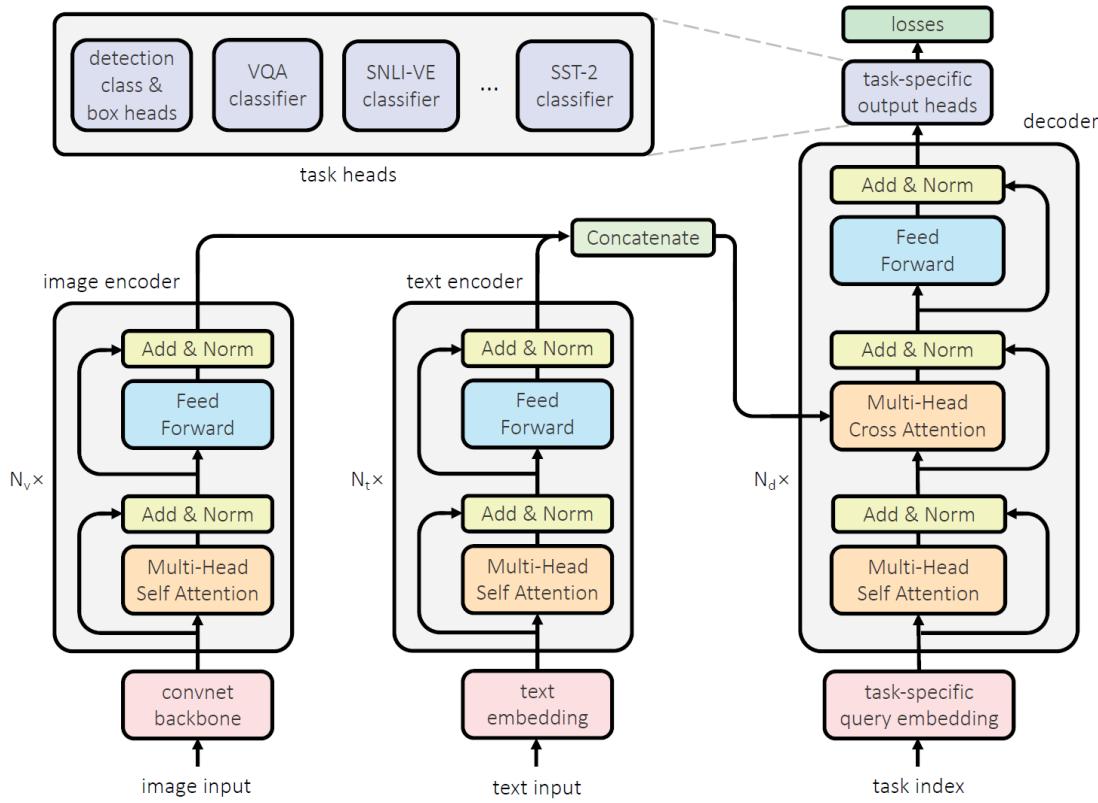
- Early VL Pretraining
- Recent Breakthrough
  - Unified Foundation Model
  - VL Pre-training

# UniT (Facebook)



Unit: Multimodal multitask learning with a unified transformer, CVPR 2021

# Unit (Facebook)



| # | decoder setup        | COCO det.<br>mAP   | VG det.<br>mAP | VQAv2<br>accuracy |
|---|----------------------|--------------------|----------------|-------------------|
| 1 | single-task training | 40.6 / –           | 3.87           | 66.38 / –         |
| 2 | separate             | <b>40.8</b> / –    | 3.91           | <b>68.84</b> / –  |
| 3 | shared               | 37.2 / –           | 4.05           | 68.79 / –         |
| 4 | shared (COCO init.)  | <b>40.8</b> / 41.1 | <b>4.53</b>    | 67.30 / 67.47     |

object detection (COCO det.)

visual question answering (VQAv2)

object detection (VG det.)

visual entailment (SNLI-VE)

QNLI

paragraph: The most important tributaries in this area are the Ill below of Strasbourg, the Neckar in Mannheim and the Main across from Mainz.  
question: What is the first major city in the stream of the Rhine?  
prediction: cannot be answered

MNLI-mm

premise: We serve a classic Tuscan meal that includes a Florentine terrine made with duck and chicken livers.  
hypothesis: We serve a meal of Florentine terrine.  
prediction: entailment

QQP

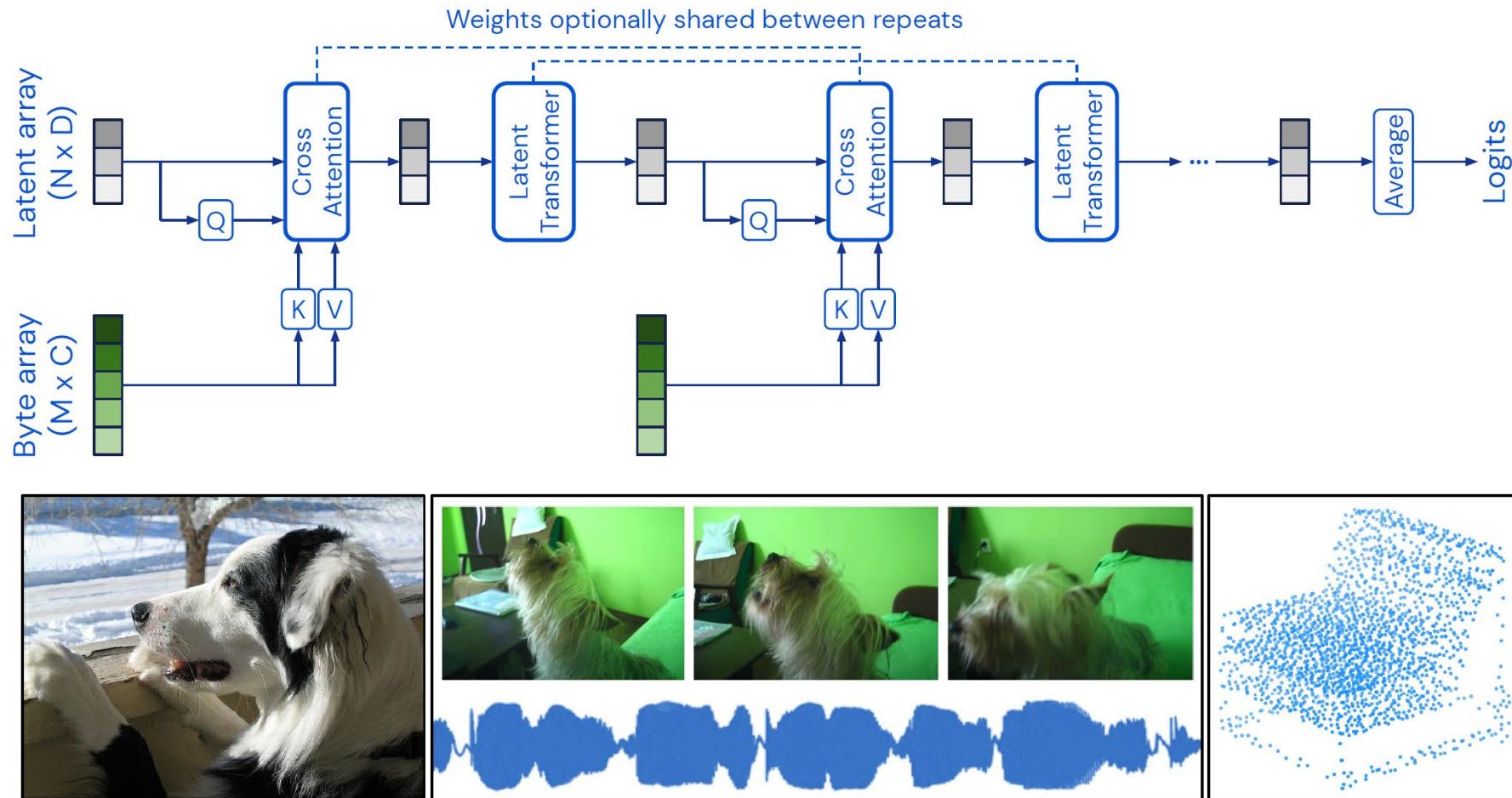
question 1: Why we do study computer fundamental in software engineering?  
question 2: Do we get to chose only one computer language when we are studying engineering?  
prediction: not equivalent

SST-2

paragraph: in exactly 89 minutes, most of which passed as slowly as if I'd been sitting naked on an igloo, formula 51 sank from quirky to jerky to utter turkey.  
sentiment: negative

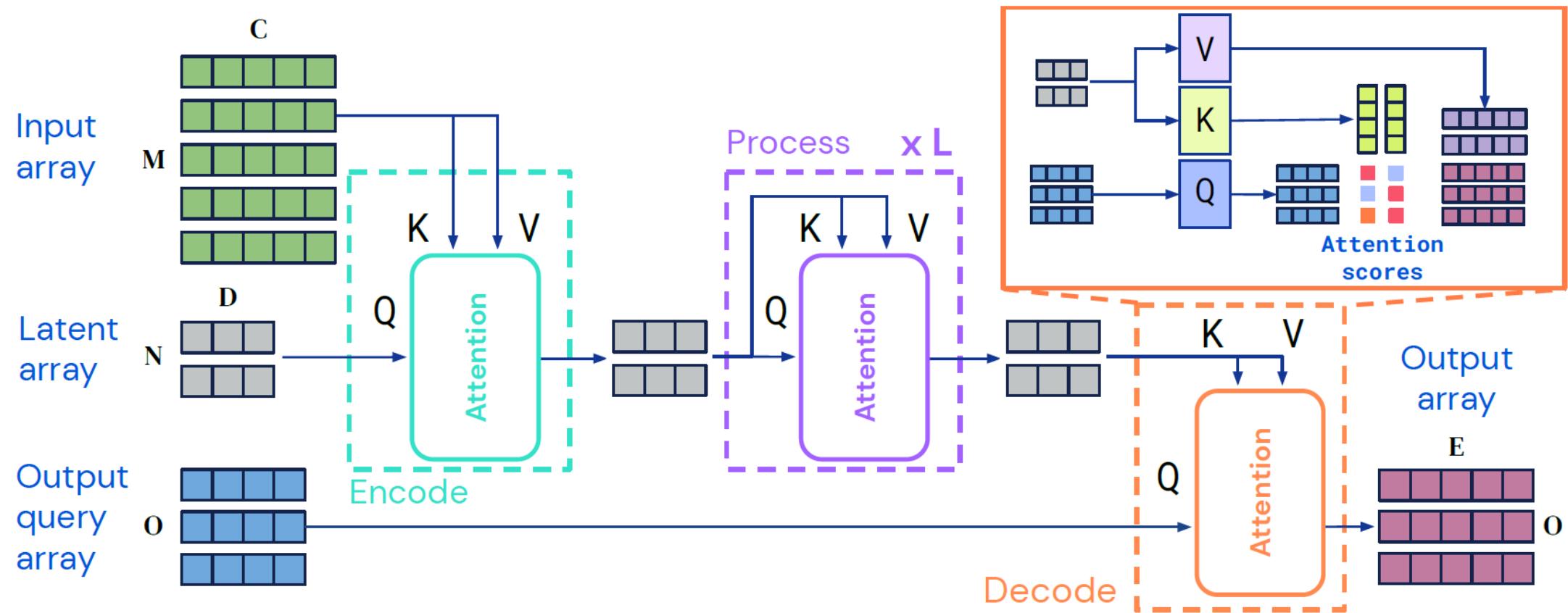
Unit: Multimodal multitask learning with a unified transformer, CVPR 2021

# Perceiver (Deepmind)



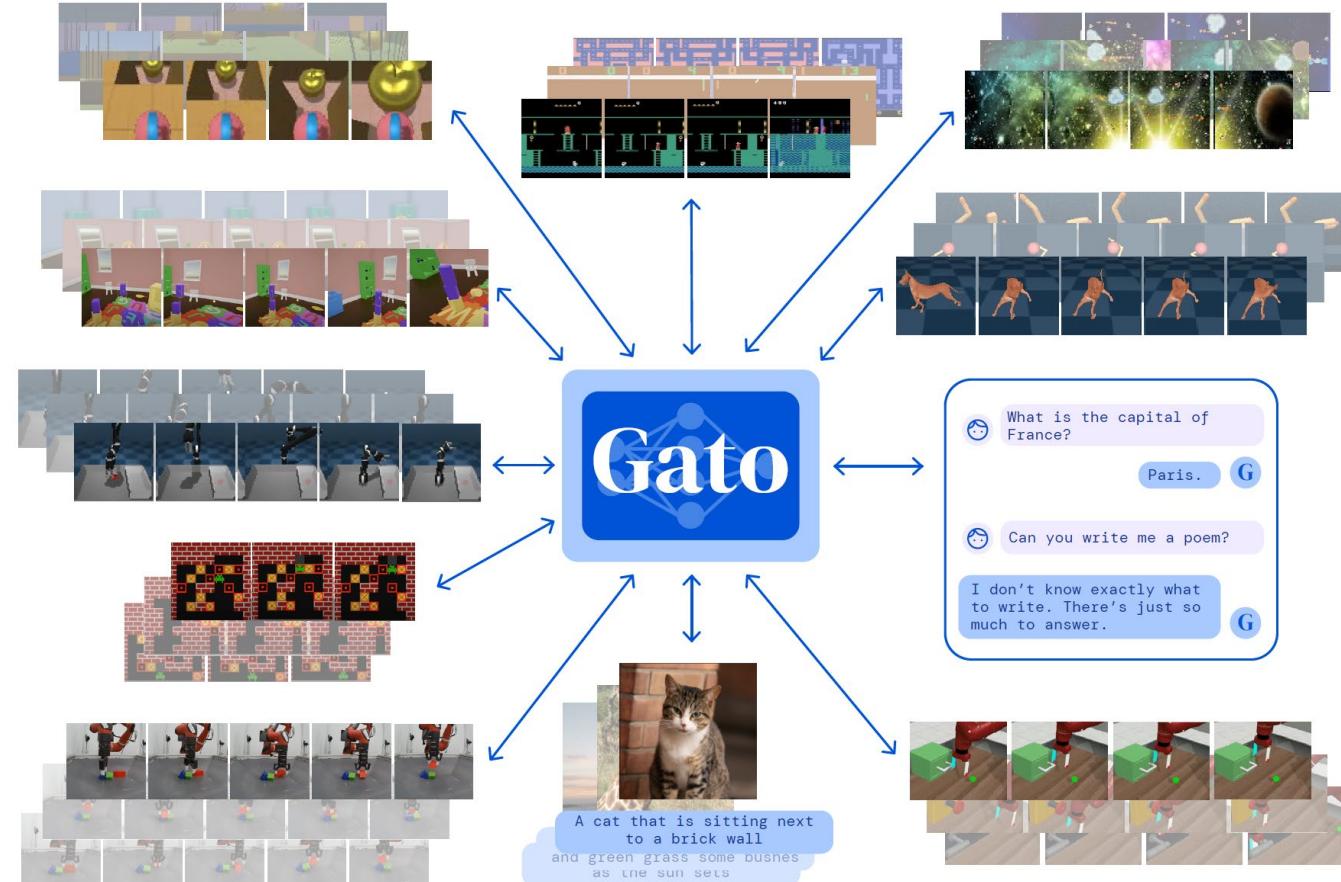
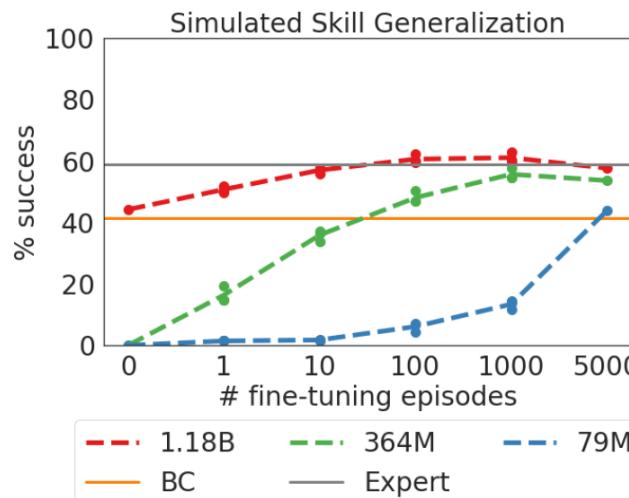
Perceiver: General Perception with Iterative Attention, Arxiv 2021

# Perceiver IO (Deepmind)



# GATO (Deepmind)

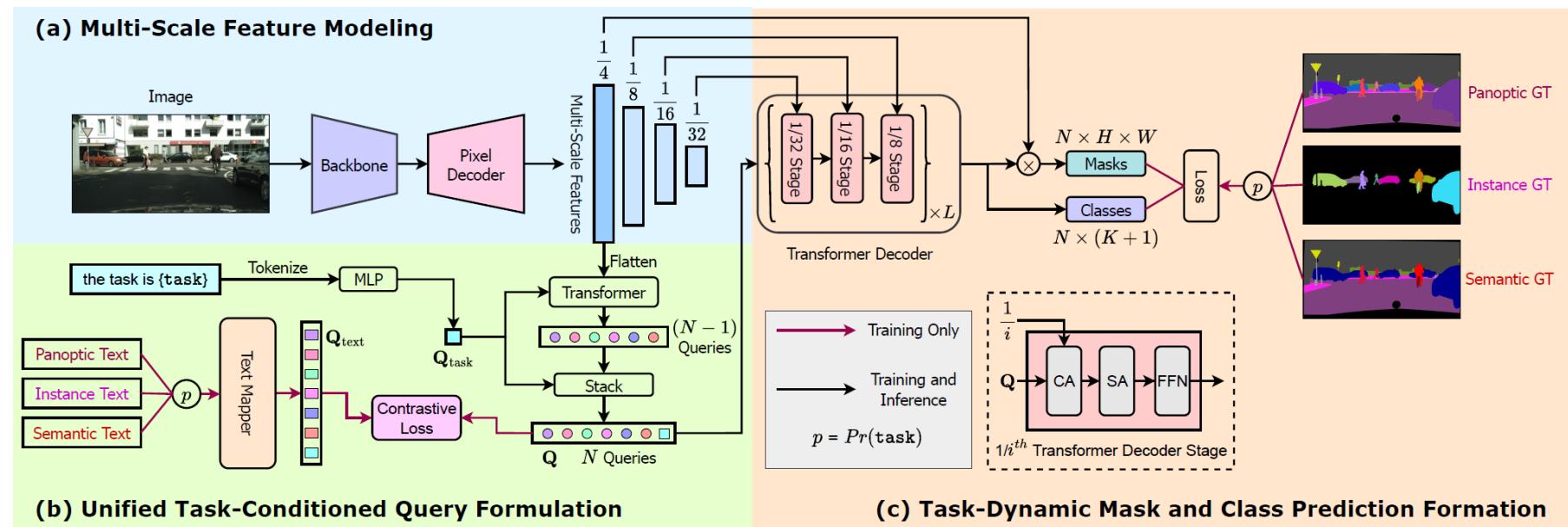
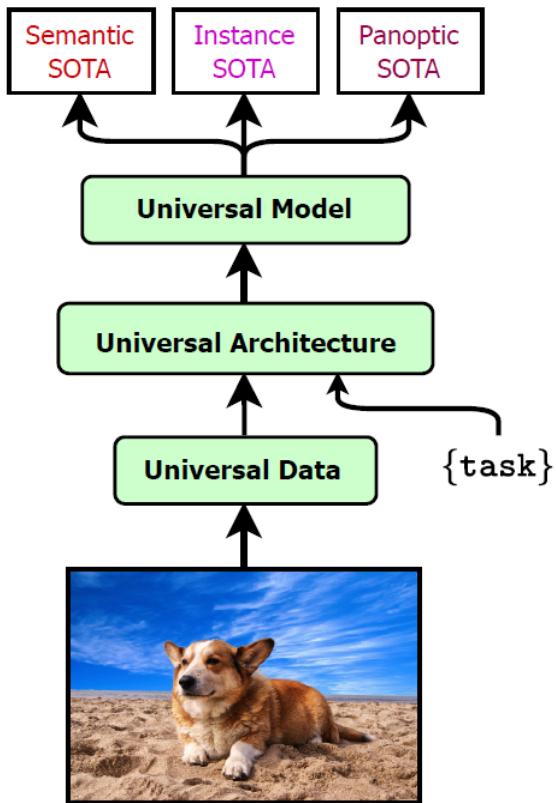
- Multi-modal
- Multi-task
- Multi-embodiment



A generalist agent, Arxiv 2022.

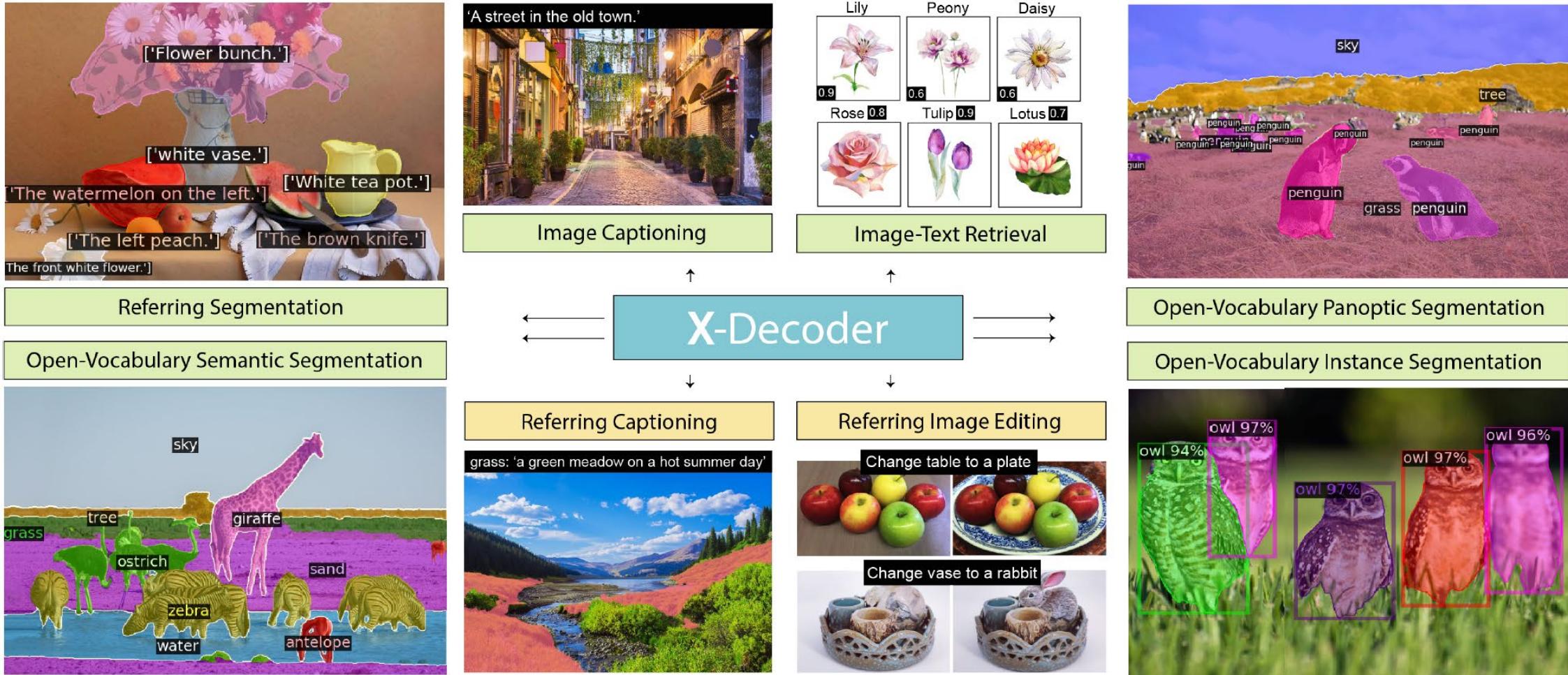
# OneFormer (UIUC)

1 architecture, 1 model & 1 dataset



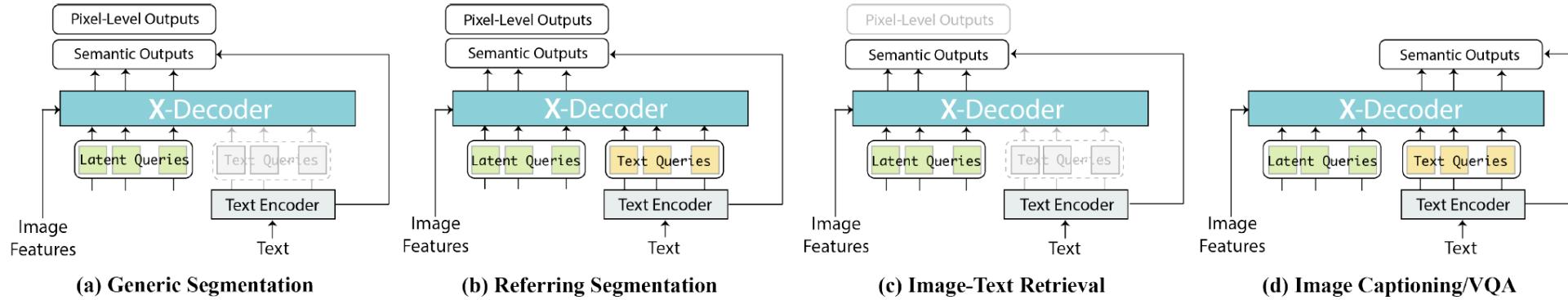
OneFormer: One Transformer to Rule Universal Image Segmentation, Arxiv 2022.

# X-Decoder (Microsoft)



Generalized Decoding for Pixel, Image, and Language, Arxiv 2022

# X-Decoder (Microsoft)



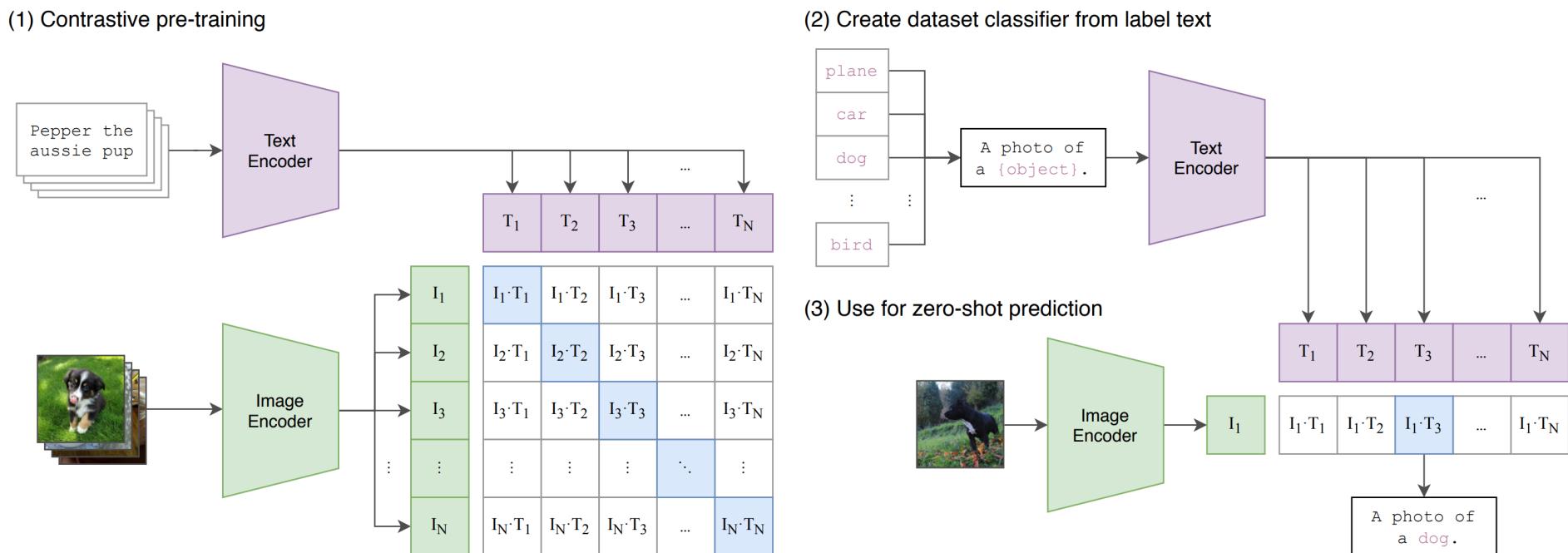
| Method                      | Type                 | Generic Segmentation |             |             |                   |             |               | Referring g-Ref | Retrieval     |      |             |             | Captioning         |       | VQA         |                |             |
|-----------------------------|----------------------|----------------------|-------------|-------------|-------------------|-------------|---------------|-----------------|---------------|------|-------------|-------------|--------------------|-------|-------------|----------------|-------------|
|                             |                      | ADE                  |             | COCO        |                   | COCO        | F30k-Karpathy |                 | COCO-Karpathy | IR@1 | TR@1        | IR@1        | TR@1               | CIDEr | BLEU        | VQAv2-test dev | std         |
|                             |                      | PQ                   | mAP         | mIoU        | PQ                | mAP         | mIoU          | g-Ref           | cIoU          | IR@1 | TR@1        | IR@1        | TR@1               | CIDEr | BLEU        | VQAv2-test dev | std         |
| Mask2Former (T) [12]        | Segmentation         | 39.7                 | 26.4        | 47.7        | 53.2              | 43.3        | 63.2          | -               | -             | -    | -           | -           | -                  | -     | -           | -              | -           |
| Mask2Former (B) [12]        |                      | *                    | *           | 53.9        | 56.4              | 46.3        | 67.1          | -               | -             | -    | -           | -           | -                  | -     | -           | -              | -           |
| Mask2Former (L) [12]        |                      | 48.1                 | 34.2        | 56.1        | 57.8              | <b>48.6</b> | 67.4          | -               | -             | -    | -           | -           | -                  | -     | -           | -              | -           |
| Pano/SegFormer (B) [47, 80] |                      | *                    | *           | 51.0        | 55.4              | *           | *             | -               | -             | -    | -           | -           | -                  | -     | -           | -              | -           |
| kMaX-DeepLab (L) [93]       |                      | 48.7                 | *           | 54.8        | <b>58.1</b>       | *           | *             | -               | -             | -    | -           | -           | -                  | -     | -           | -              | -           |
| LAVT (B) [86]               |                      | -                    | -           | -           | -                 | -           | -             | 61.2            | -             | -    | -           | -           | -                  | -     | -           | -              | -           |
| UNITER (B) [10]             | Vision Language (VL) | -                    | -           | -           | -                 | -           | -             | -               | 50.3          | 64.4 | 72.5        | 85.9        | -                  | -     | 72.7        | 72.9           | -           |
| UNITER (L) [10]             |                      | -                    | -           | -           | -                 | -           | -             | -               | 52.9          | 65.6 | 75.6        | 87.3        | -                  | -     | 73.8        | 74.0           | -           |
| VinVL (B) [96]              |                      | -                    | -           | -           | -                 | -           | -             | -               | 58.1          | 74.6 | *           | *           | 129.3              | 38.2  | 76.0        | 76.1           | -           |
| VinVL (L) [96]              |                      | -                    | -           | -           | -                 | -           | -             | -               | <b>58.8</b>   | 75.4 | *           | *           | 130.8              | 38.5  | 76.5        | 76.6           | -           |
| ALBEF-4M (B) [43]           |                      | -                    | -           | -           | -                 | -           | -             | -               | 56.8          | 73.1 | 82.8        | 94.3        | *                  | *     | 74.5        | 74.7           | -           |
| METER-Swin (B) [21]         |                      | -                    | -           | -           | -                 | -           | -             | -               | 54.9          | 73.0 | 79.0        | 92.4        | *                  | *     | 76.4        | 76.4           | -           |
| UVIM (L) [38]               | General Purpose      | *                    | *           | *           | 45.8 <sup>1</sup> | *           | *             | -               | -             | -    | -           | -           | -                  | -     | -           | -              | -           |
| UniT (T) [32]               |                      | -                    | -           | -           | -                 | -           | -             | -               | -             | -    | -           | -           | -                  | -     | 67.6        | *              | -           |
| GPV (T) [27]                |                      | -                    | -           | -           | -                 | -           | -             | -               | -             | -    | -           | -           | 102.3 <sup>2</sup> | *     | 62.5        | *              | -           |
| UniTAB (B) [85]             |                      | -                    | -           | -           | -                 | -           | -             | -               | -             | -    | -           | -           | 119.8              | 36.1  | 70.7        | 71.0           | -           |
| Pix2Seq v2 (B) [7]          |                      | -                    | *           | -           | -                 | 38.2        | -             | -               | -             | -    | -           | -           | *                  | 34.9  | -           | -              | -           |
| Unified-IO (B) [54]         |                      | -                    | *           | -           | -                 | *           | -             | -               | -             | -    | -           | -           | *                  | *     | 61.8        | *              | -           |
| Unified-IO (L) [54]         |                      | -                    | *           | -           | -                 | *           | -             | -               | -             | -    | -           | -           | *                  | *     | 67.8        | *              | -           |
| GLIPv2 (T) [95]             |                      | -                    | *           | -           | -                 | -42.0       | -             | *               | -             | -    | -           | -           | 122.1              | *     | 71.6        | 71.8           | -           |
| GLIPv2 (B) [95]             |                      | -                    | *           | -           | -                 | -45.8       | -             | *               | -             | -    | -           | -           | 128.5              | *     | 73.1        | 73.3           | -           |
| GLIPv2 (H) [95]             |                      | -                    | *           | -           | -                 | -48.9       | -             | *               | -             | -    | -           | -           | 131.0              | *     | 74.6        | 74.8           | -           |
| X-Decoder (T)               |                      | 41.6                 | 27.7        | 51.0        | 52.6              | 41.3/42.3   | 62.4          | 59.8            | 61.9          | 49.3 | 66.7        | 74.4        | 89.1               | 122.3 | 37.8        | 70.6           | 70.9        |
| X-Decoder (B)               |                      | 46.8                 | 33.5        | 54.6        | 56.2              | 45.8/45.8   | 66.0          | 62.4            | 64.5          | 54.5 | 71.2        | 80.8        | 93.2               | 129.0 | 39.6        | 74.1           | 74.2        |
| X-Decoder (L)               |                      | <b>49.6</b>          | <b>35.8</b> | <b>58.1</b> | 56.9              | 46.7/47.1   | <b>67.5</b>   | <b>64.6</b>     | <b>64.6</b>   | 58.6 | <b>76.1</b> | <b>84.4</b> | <b>94.4</b>        | 132.1 | <b>40.2</b> | <b>76.8</b>    | <b>77.0</b> |

# Vision-Language Pretraining

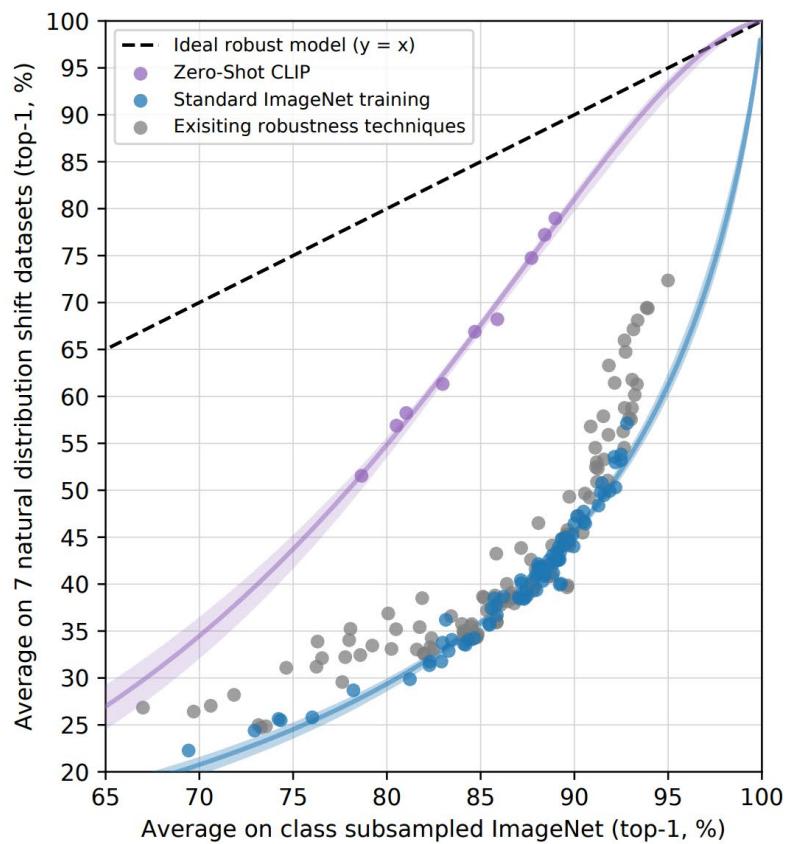
- Early VL Pretraining
- Recent Breakthrough
  - Unified Foundation Model
  - VL Pre-training

# VL Pretraining: CLIP (OpenAI)

- 400 million (image, text) pairs



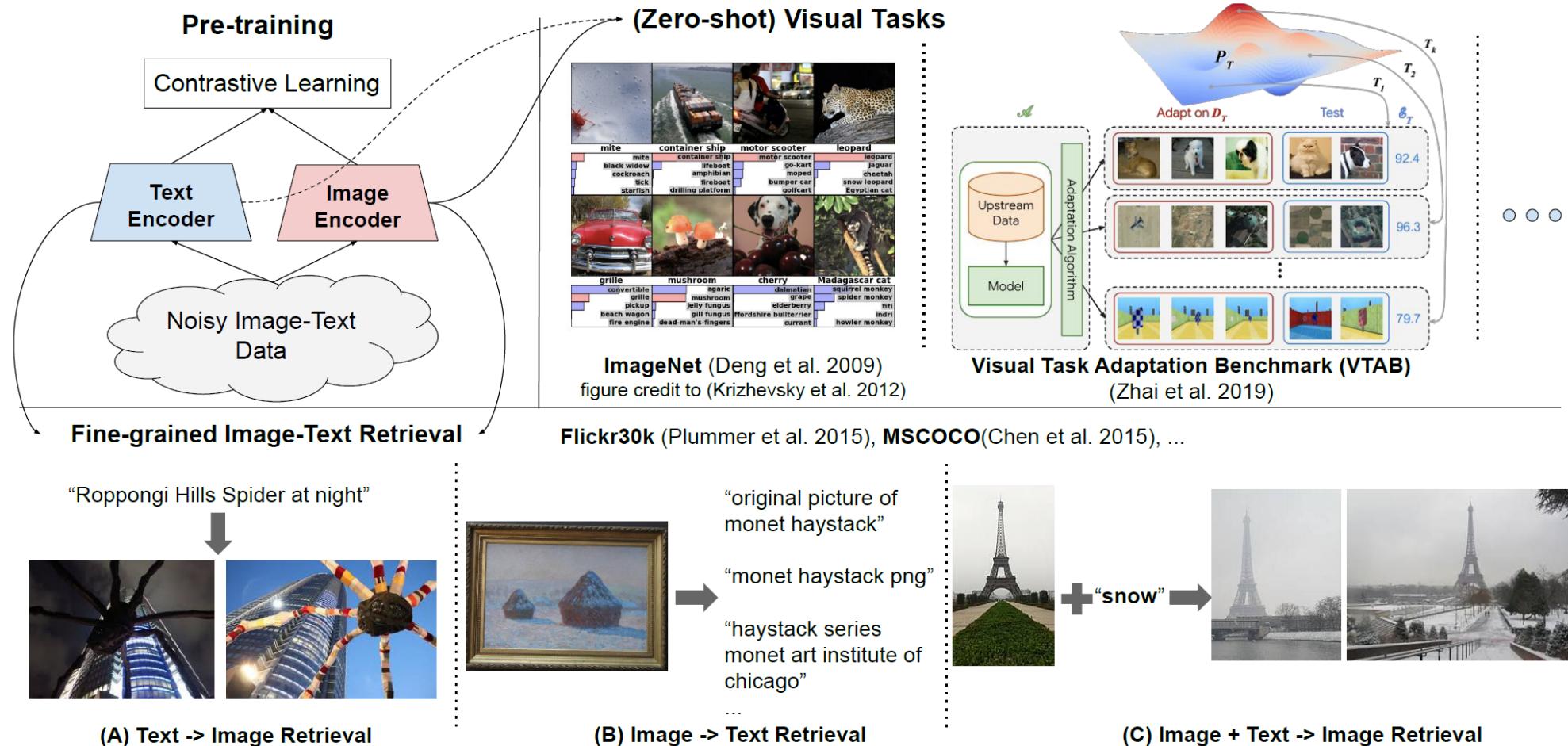
# Performance



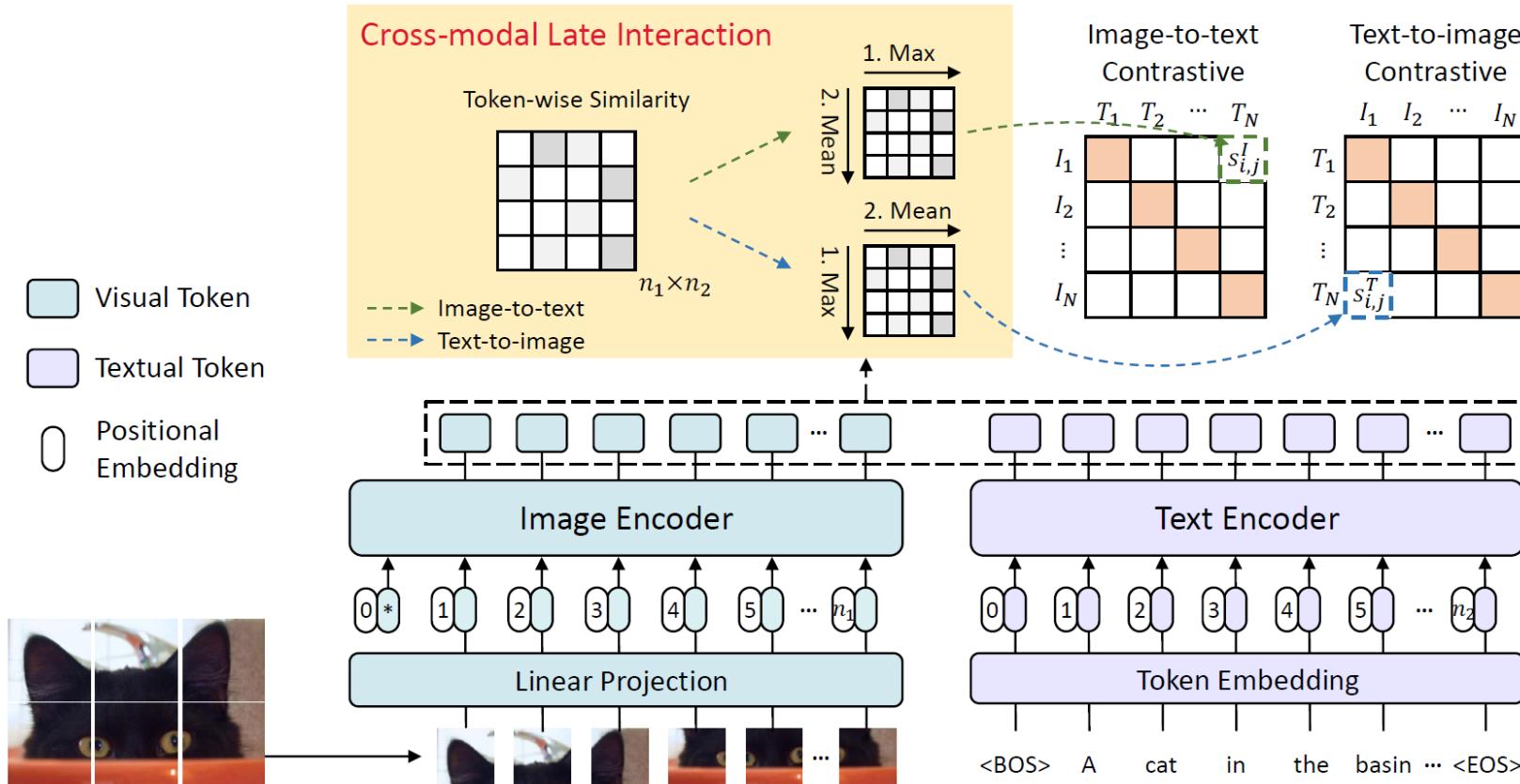
Dataset Examples

|                 | ImageNet | Zero-Shot<br>ResNet101 | CLIP        | Δ Score |
|-----------------|----------|------------------------|-------------|---------|
| ImageNet        |          | <b>76.2</b>            | <b>76.2</b> | 0%      |
| ImageNetV2      |          | 64.3                   | <b>70.1</b> | +5.8%   |
| ImageNet-R      |          | 37.7                   | <b>88.9</b> | +51.2%  |
| ObjectNet       |          | 32.6                   | <b>72.3</b> | +39.7%  |
| ImageNet Sketch |          | 25.2                   | <b>60.2</b> | +35.0%  |
| ImageNet-A      |          | 2.7                    | <b>77.1</b> | +74.4%  |

# VL Pretraining: ALIGN (Google)

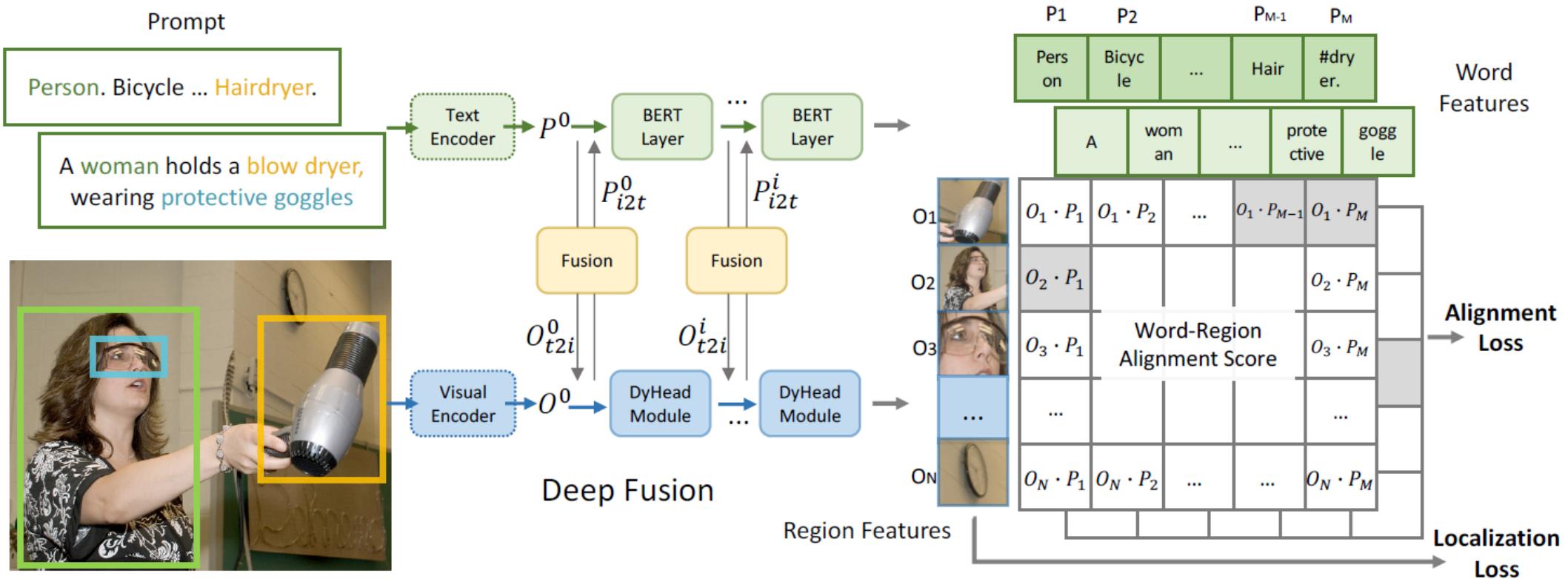


# FILIP (Huawei)



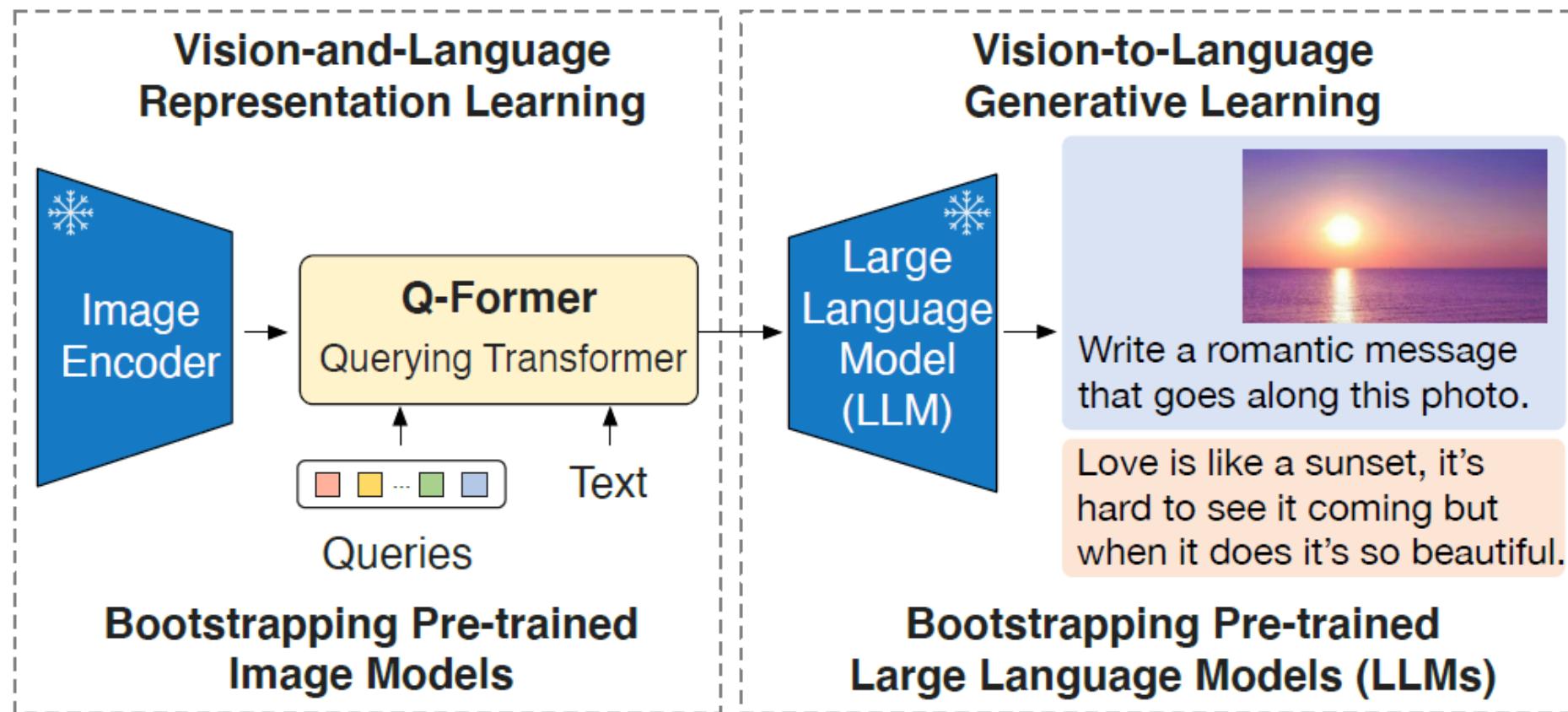
FILIP: fine-grained interactive language-image pre-training, ICLR 2022

# GLIP (Microsoft)



Grounded Language-Image Pre-training, CVPR 2022.

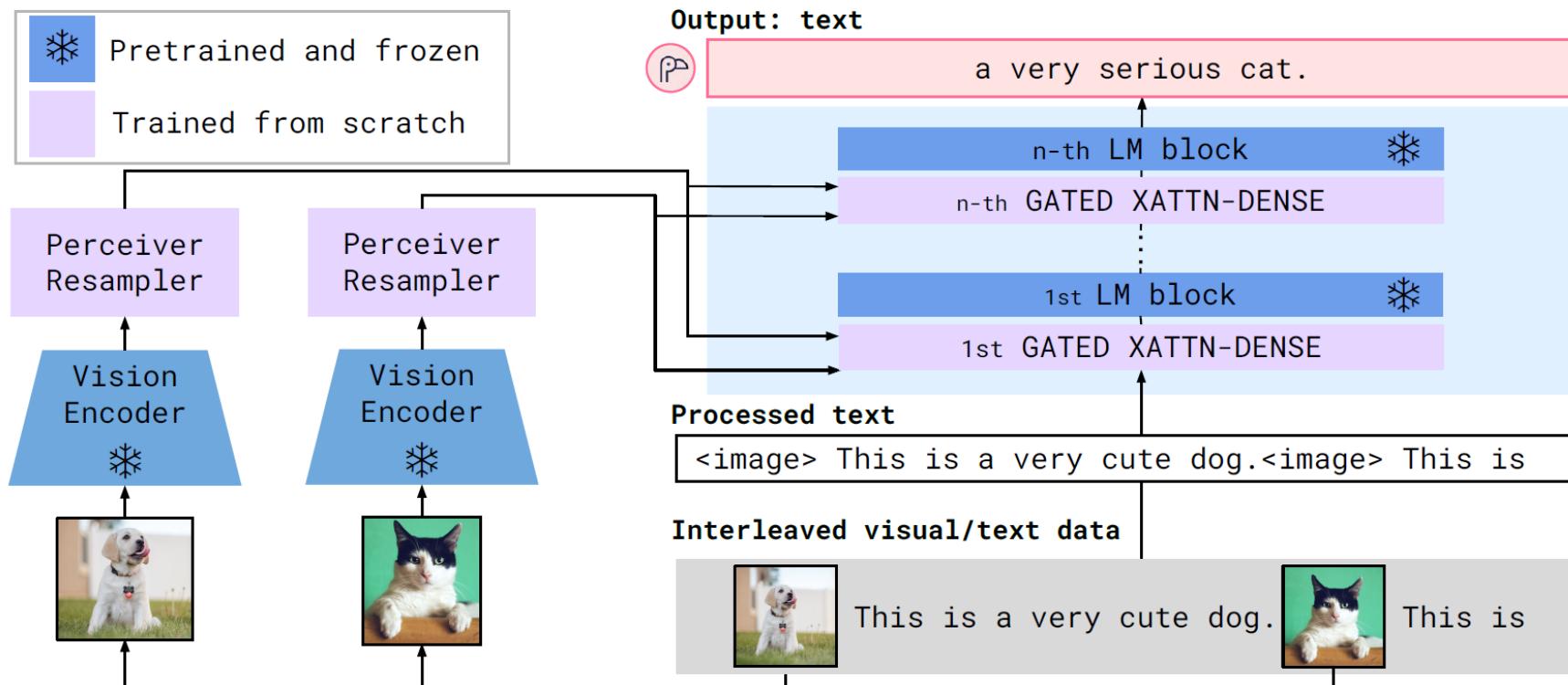
# BLIP-2 (SalesForce)



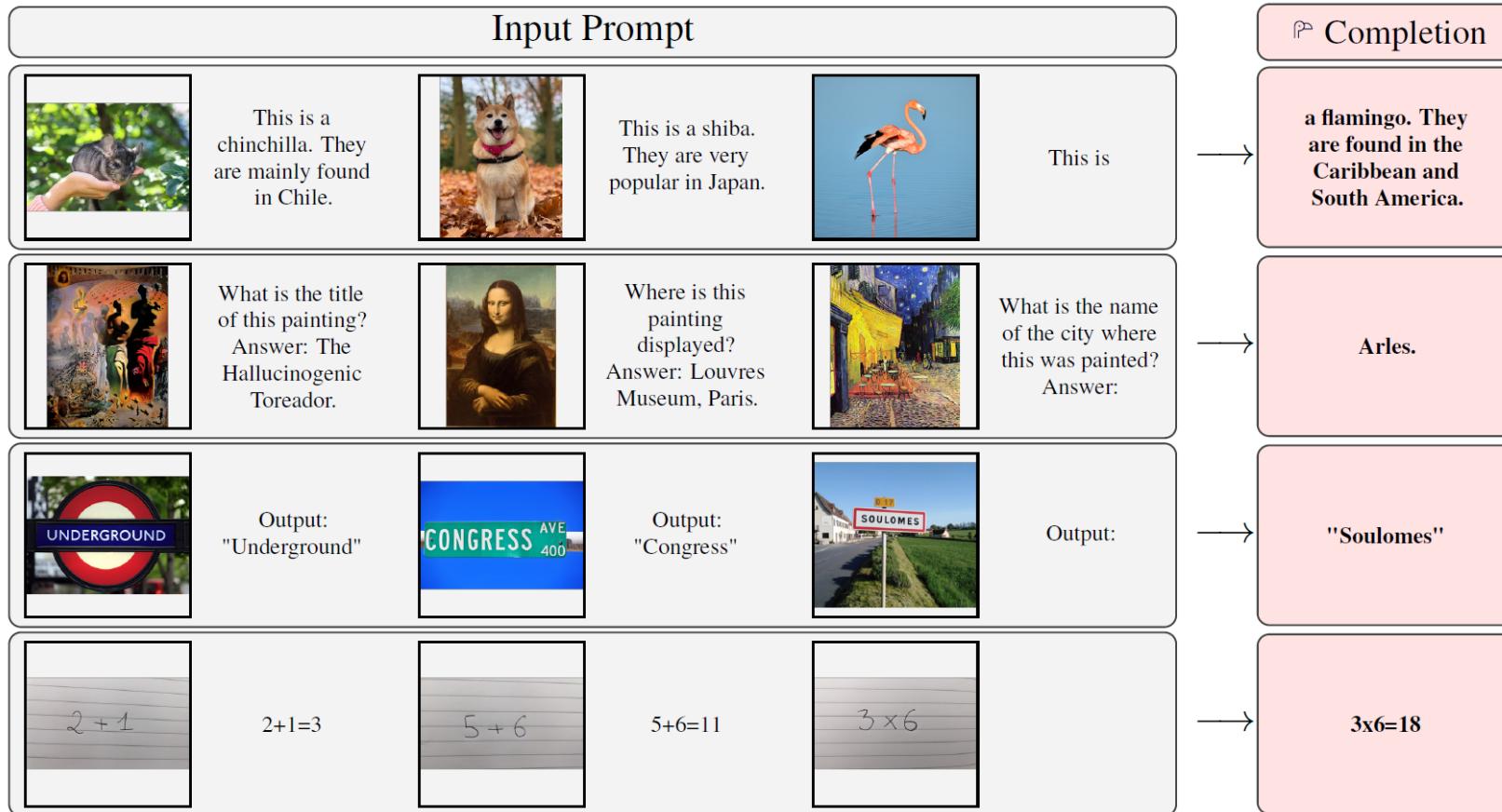
BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models,  
Arxiv 2023.

# Flamingo (Deepmind)

- In-context learning



# Flamingo (Deepmind)



# Flamingo (Deepmind)



This is a picture of two teddy bears on the moon.

What are they doing?

They are having a conversation.

What object are they using?

It looks like a computer.

Is this surprising?

Yes, it is surprising.

Why is this picture surprising to you?

I think it is surprising because teddy bears are not usually found on the moon.



What is the common thing about these three images?

They are all flamingos.

What is the difference between these three images?

The first one is a cartoon, the second one is a real flamingo, and the third one is a 3D model of a flamingo.



This is an apple with a sticker on it.

What does the sticker say?

The sticker says "iPod".

Where is the photo taken?

It looks like it's taken in a backyard.

Do you think it is printed or handwritten?

It looks like it's handwritten.

What color is the sticker?

It's white.



This is a cityscape. It looks like Chicago.

What makes you think this is Chicago?

I think it's Chicago because of the Shedd Aquarium in the background.

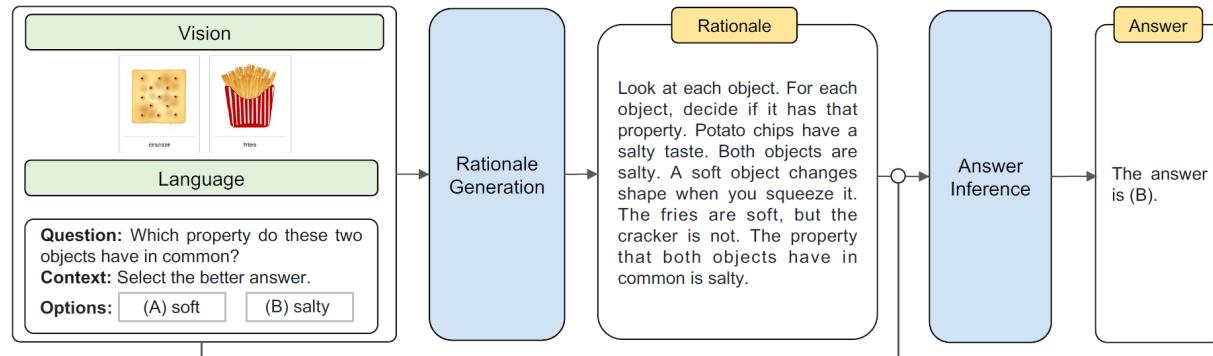


What about this one? Which city is this and what famous landmark helped you recognise the city?

This is Tokyo. I think it's Tokyo because of the Tokyo Tower.

# ChatGPT & GPT-4 (OpenAI)

- Chain of Thoughts
- RLHF
- In-context learning



Multimodal Chain-of-Thought Reasoning in Language Models, Arxiv 2023

## Example of GPT-4 visual input:

User What is funny about this image? Describe it panel by panel.



Source: <https://www.reddit.com/r/hmmm/comments/ubab5v/hmmm/>

GPT-4 The image shows a package for a "Lightning Cable" adapter with three panels.  
Panel 1: A smartphone with a VGA connector (a large, blue, 15-pin connector typically used for computer monitors) plugged into its charging port.  
Panel 2: The package for the "Lightning Cable" adapter with a picture of a VGA connector on it.  
Panel 3: A close-up of the VGA connector with a small Lightning connector (used for charging iPhones and other Apple devices) at the end.  
The humor in this image comes from the absurdity of plugging a large, outdated VGA connector into a small, modern smartphone charging port.

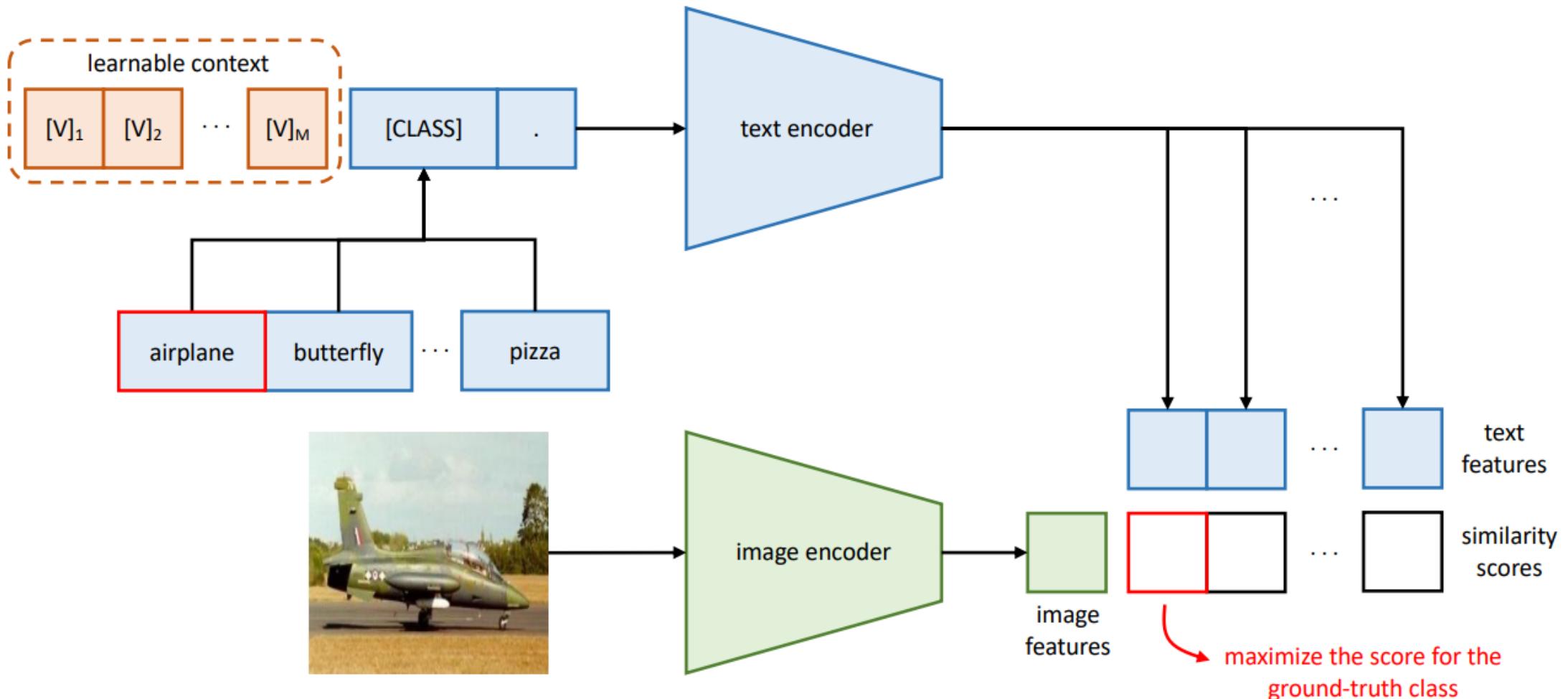
# Vision-Language Pretraining

- Vision-Language Models
- Classification: Parameter-efficient Fine-tuning
- More Challenging Tasks
- More Visual Modalities

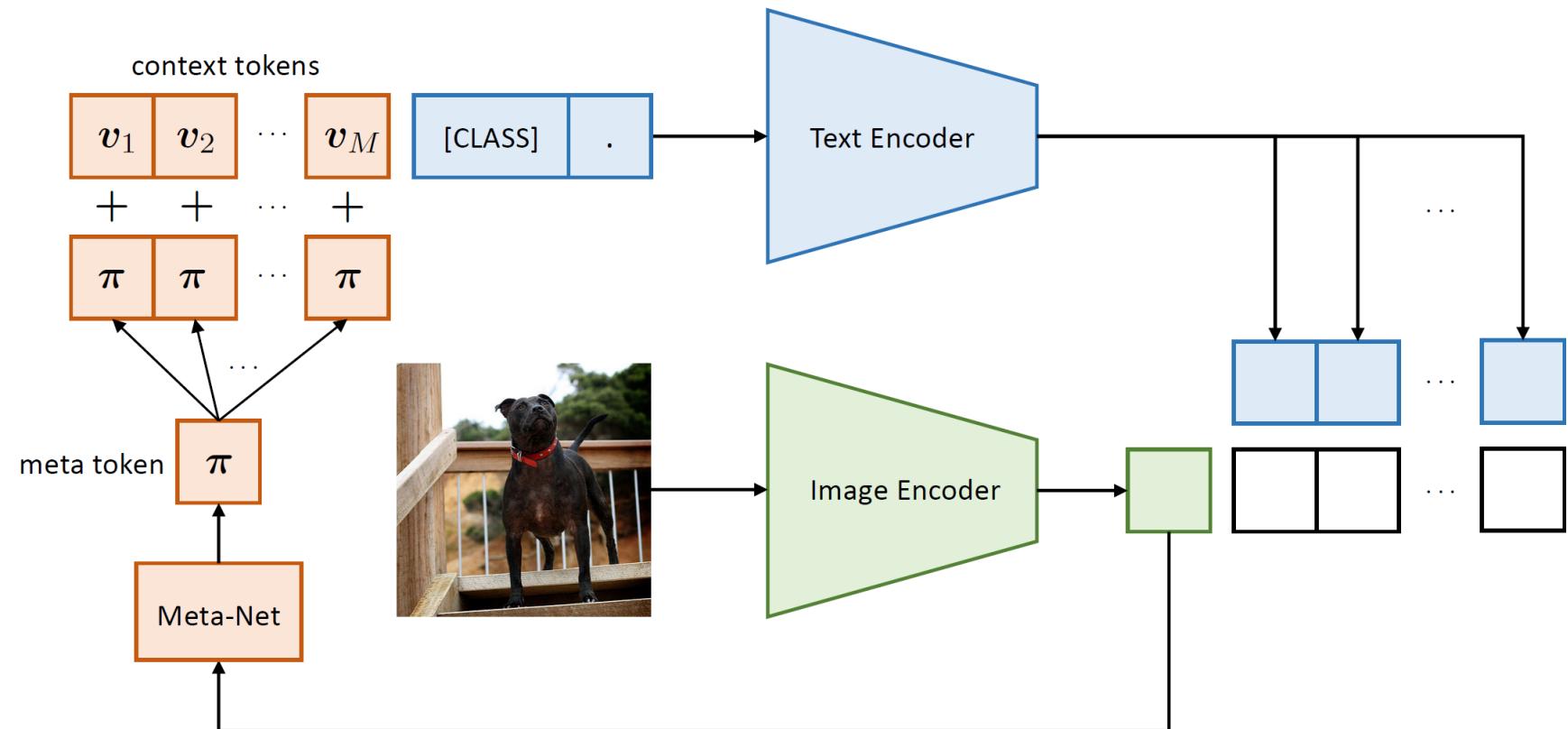
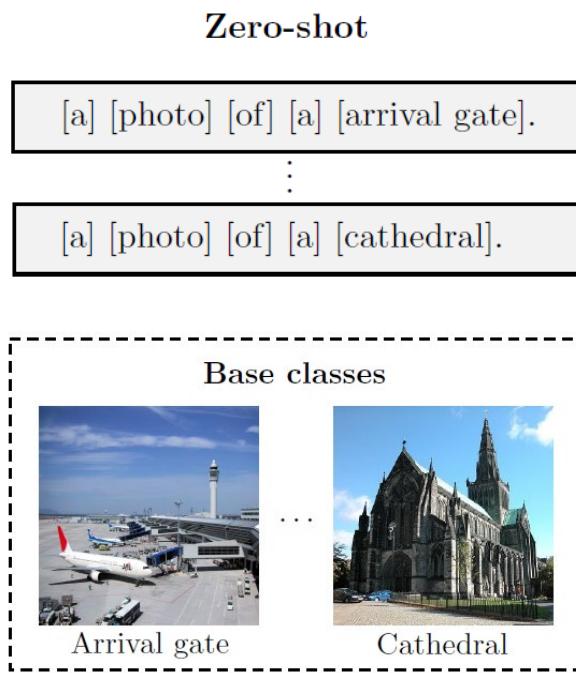
# Contents

- [Prompt](#)
- [Adapter](#)
- [Controllable FineTuning](#)

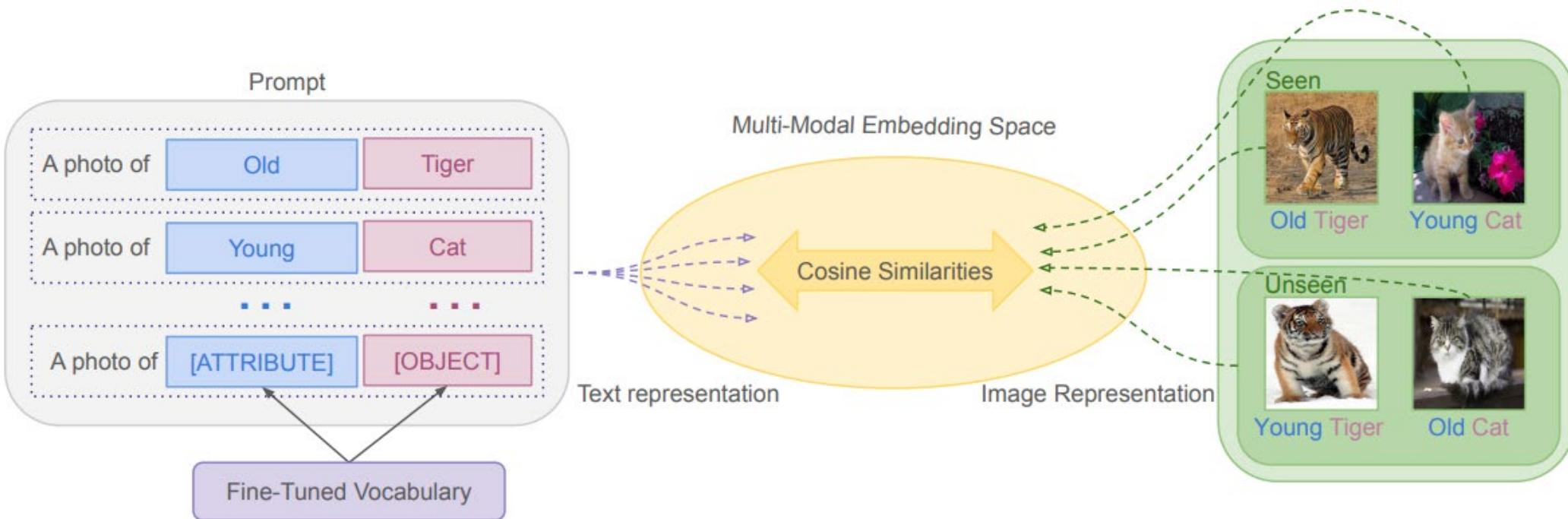
# CoOp



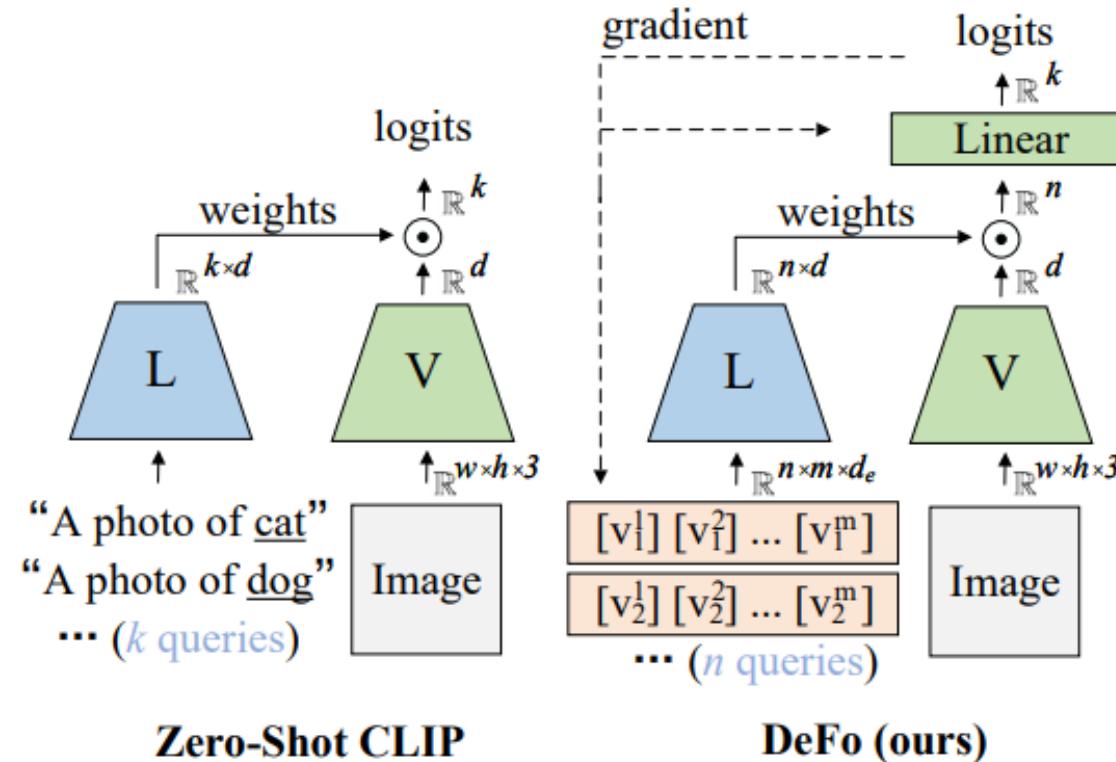
# CoCoOp: Conditional Prompt Learning



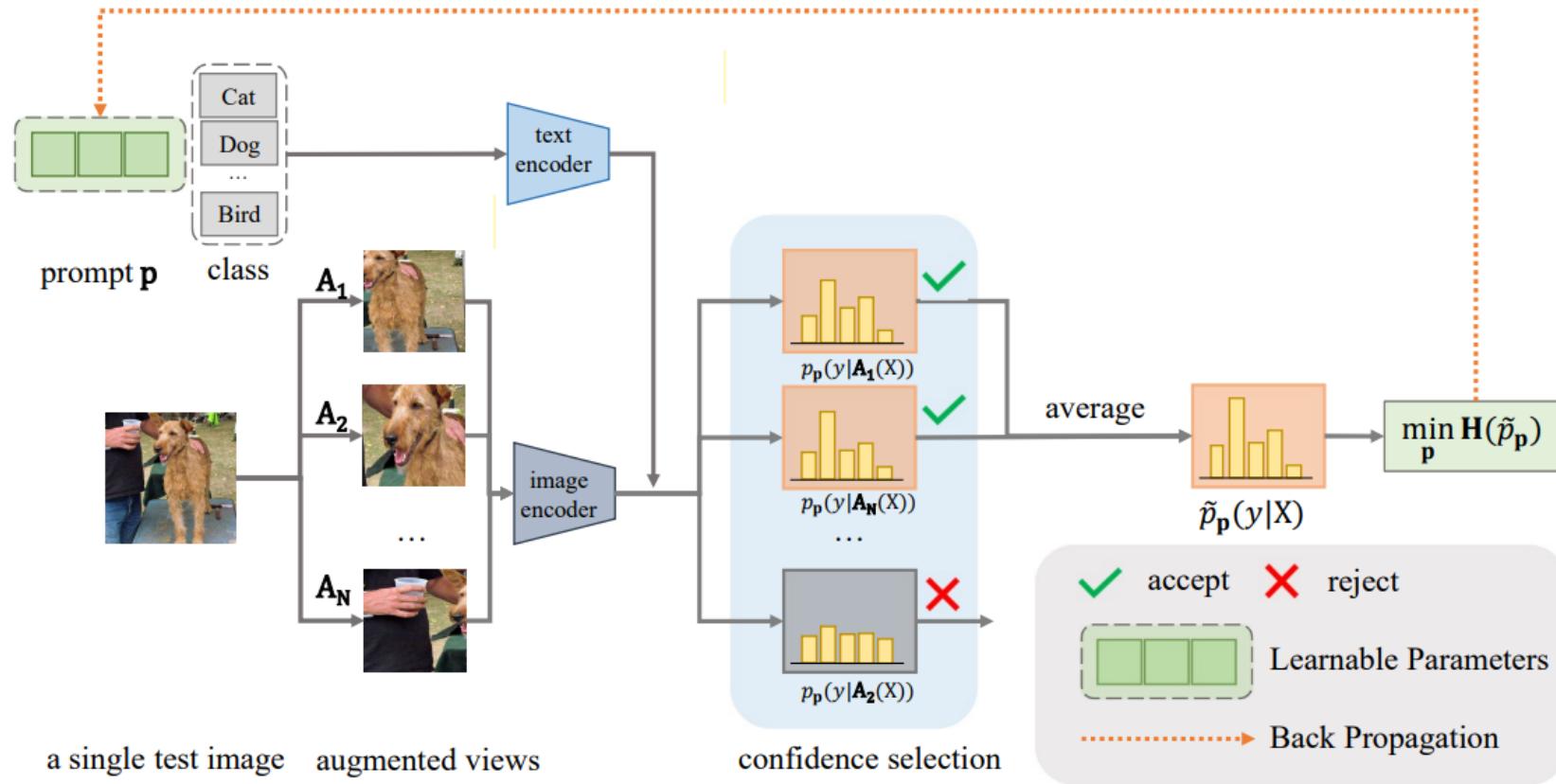
# Learning to Compose Soft Prompts for Compositional Zero-Shot Learning (ICLR2023)



# DeFo (ICLR2023)

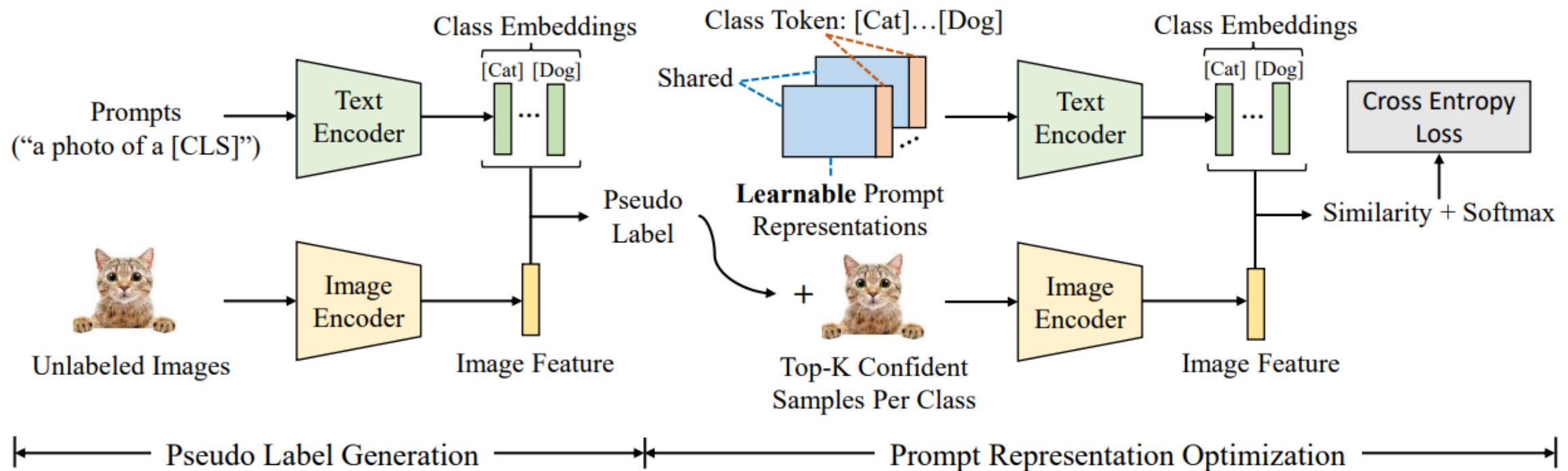


# Test-time Prompt Tuning (NIPS2022)



测试图像也是一种unsupervised setting

# Unsupervised Prompt Learning



模型确信的往往分类正确，迭代self-training有提升

# DualCoOp: Adapting CLIP to Multi-Label

- Previous single-label classification

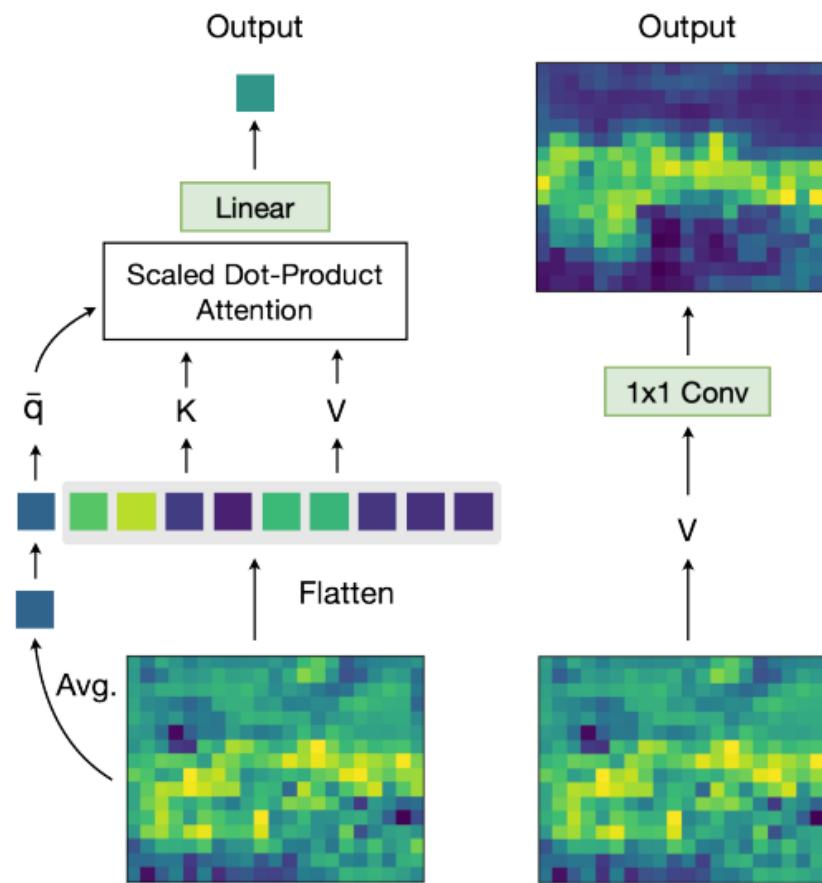
$$p(y = i | \mathbf{x}) = \frac{\exp(\langle \mathbf{w}_i, \mathbf{f} \rangle / \tau)}{\sum_{j=1}^K \exp(\langle \mathbf{w}_j, \mathbf{f} \rangle / \tau)},$$

- Multi-label binary classification

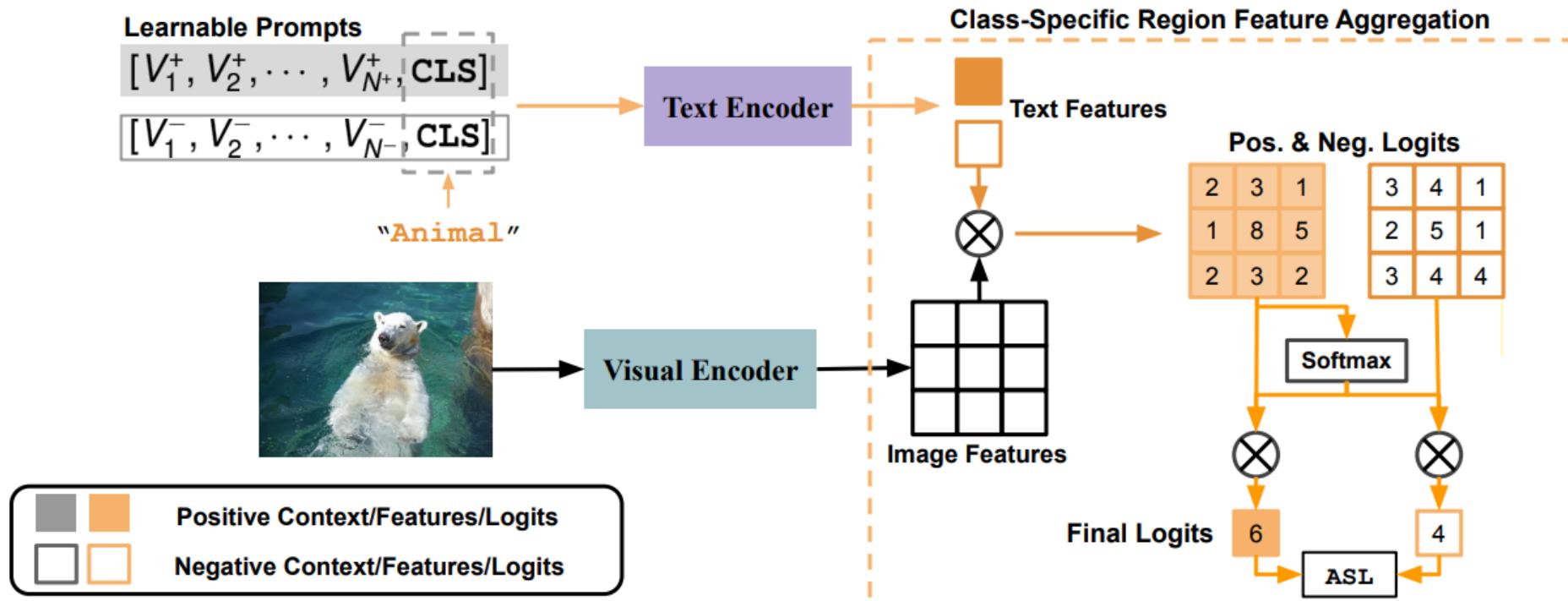
$$p(y = i | x) = \frac{\exp(\langle w_i^+, f \rangle / \tau)}{\exp(\langle w_i^+, f \rangle / \tau) + \exp(\langle w_i^-, f \rangle / \tau)}$$

# DualCoOp: Adapting CLIP to Multi-Label

- Preserve localized features



# DualCoOp: Adapting CLIP to Multi-Label



# Texts as Images



Train



$$V_1, V_2, V_3, \dots, V_M, [\text{dog}]$$



Test



A photo of a dog.  
A dog sitting on grass.  
...

Train



$$V_1, V_2, V_3, \dots, V_M, [\text{dog}]$$



Test



Woman on a horse jumping  
over a pole jump.

Train



$$V_1, V_2, V_3, \dots, V_M, [\text{person}]$$

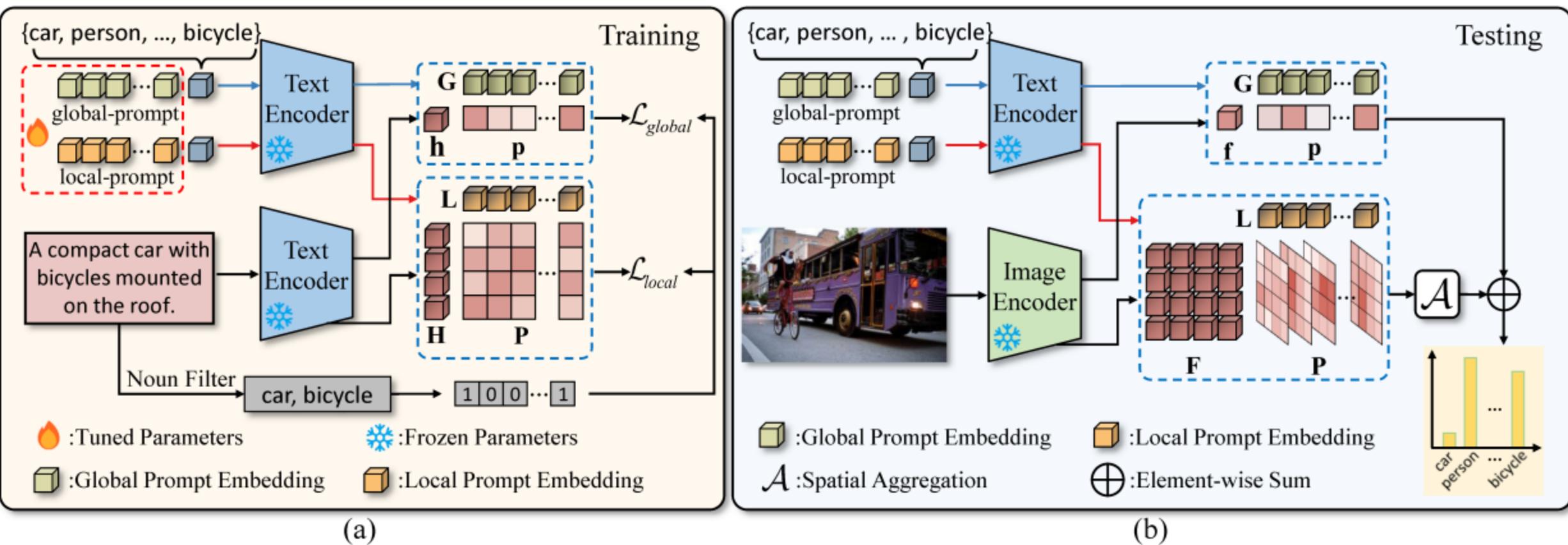
$$V_1, V_2, V_3, \dots, V_M, [\text{horse}]$$



Test



# Tal-DPT



# LOSS

- Ranking Loss

$$\mathcal{L}_{global} = \sum_{i \in \{c^+\}} \sum_{j \in \{c^-\}} \max(0, m - \mathbf{p}_i + \mathbf{p}_j),$$

$$\mathcal{L}_{local} = \sum_{i \in \{c^+\}} \sum_{j \in \{c^-\}} \max(0, m - \mathbf{p}'_i + \mathbf{p}'_j)$$

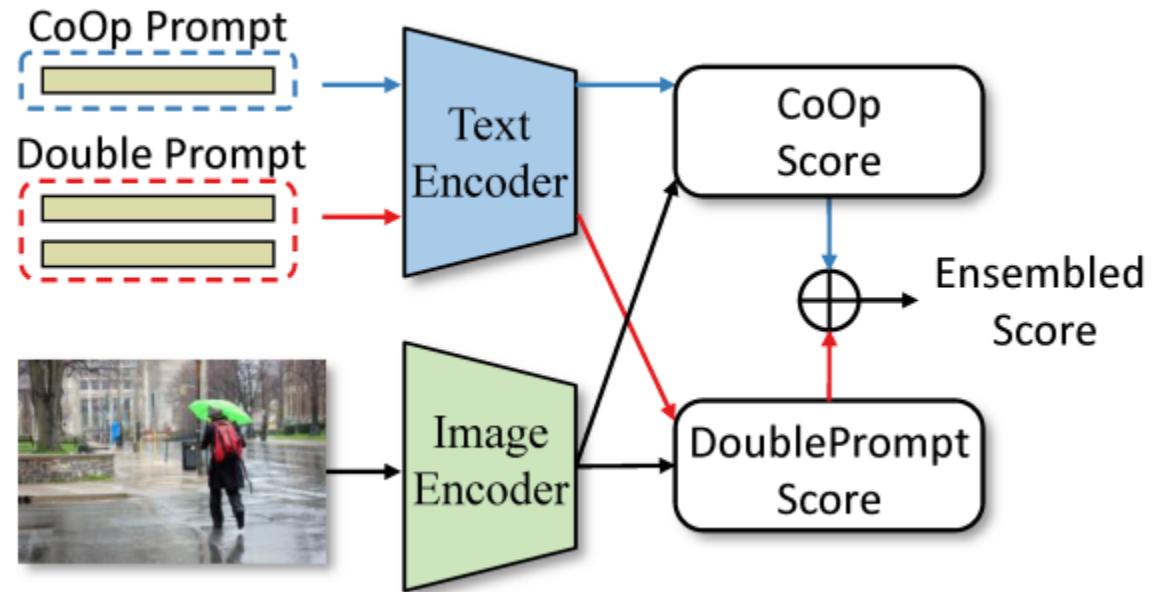
- Binary Cross Entropy Loss

$$\mathcal{L} = \text{BCE}(\mathbf{p}, \mathbf{y}) + \text{BCE}(\mathbf{p}', \mathbf{y}),$$

$$\text{BCE}(\mathbf{q}, \mathbf{y}) = -\frac{1}{C} \sum_{i=1}^C [\mathbf{y}_i \cdot \log \mathbf{q}_i + (1 - \mathbf{y}_i) \cdot \log (1 - \mathbf{q}_i)]$$

| Loss    | VOC2007     | MS-COCO     | NUSWIDE     |
|---------|-------------|-------------|-------------|
| BCE     | 84.9        | 59.0        | 40.5        |
| ASL [1] | 84.6        | 56.9        | 36.0        |
| RL [6]  | <b>88.3</b> | <b>65.1</b> | <b>46.5</b> |

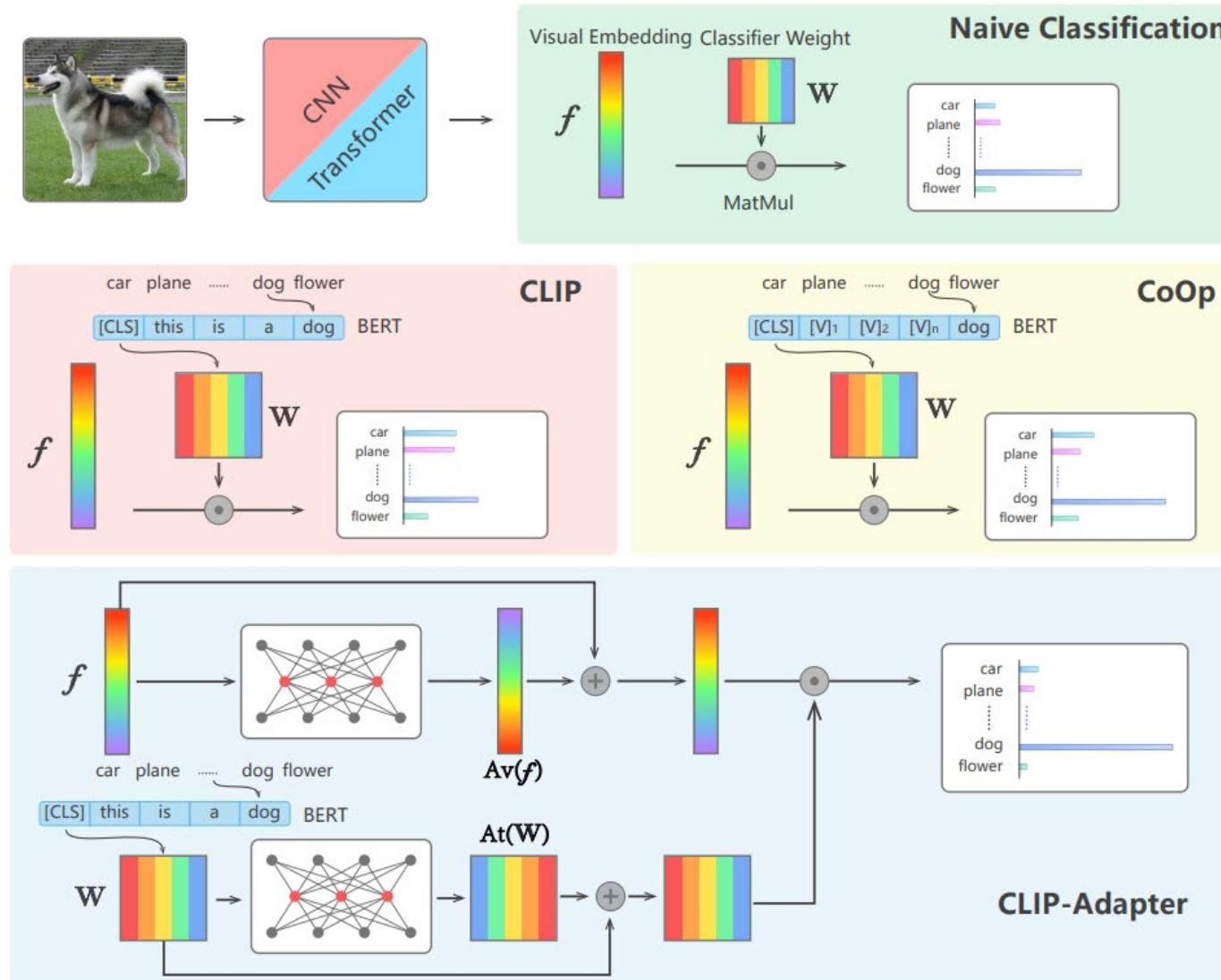
# Integration with Other Methods



# Contents

- Prompt
- Adapter
- Controllable FineTuning

# CLIP-Adapter



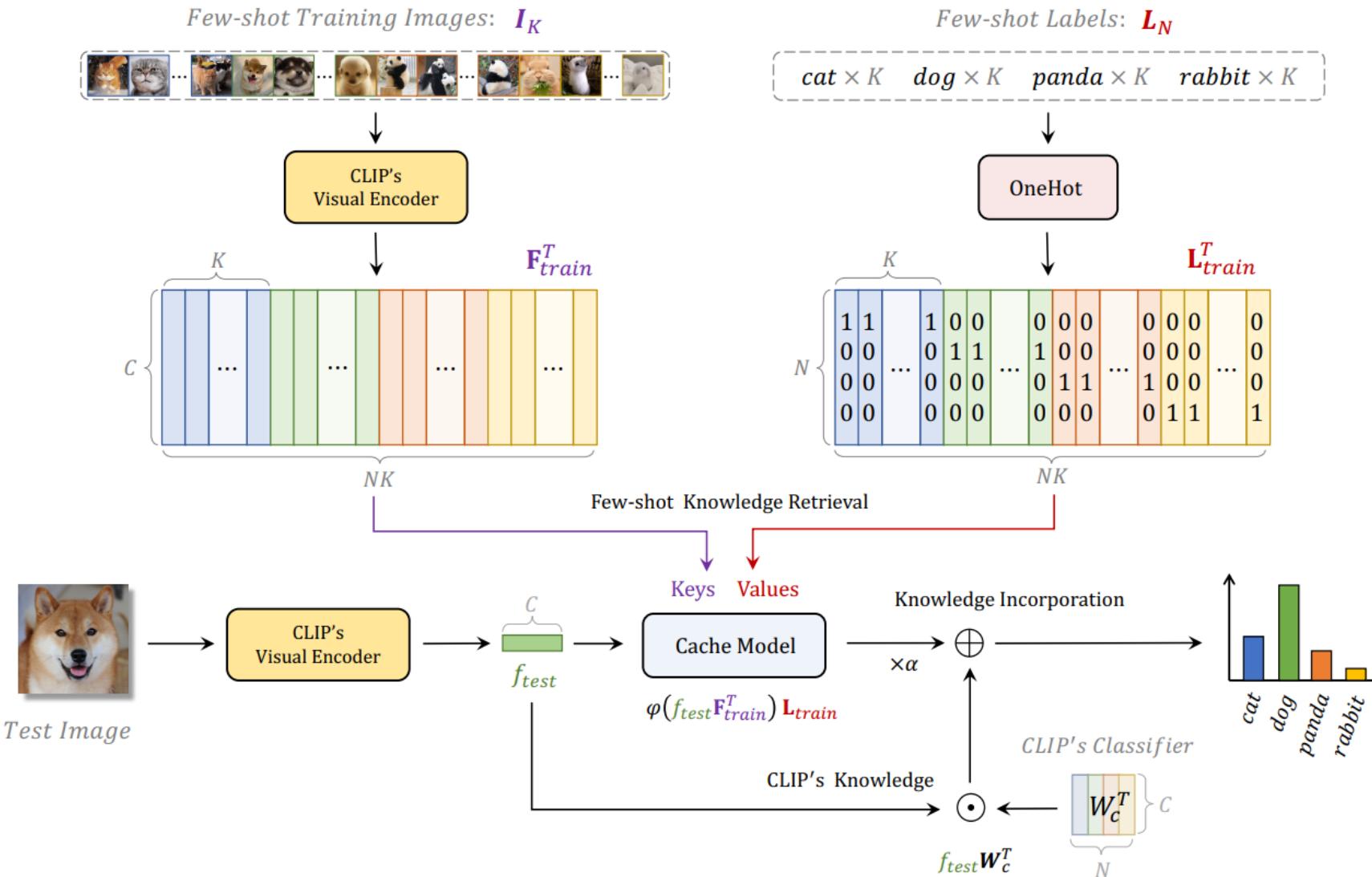
$$A_v(f) = \text{ReLU}(f^T \mathbf{W}_1^v) \mathbf{W}_2^v,$$

$$A_t(\mathbf{W}) = \text{ReLU}(\mathbf{W}^T \mathbf{W}_1^t) \mathbf{W}_2^t.$$

$$f^\star = \alpha A_v(f)^T + (1 - \alpha) f,$$

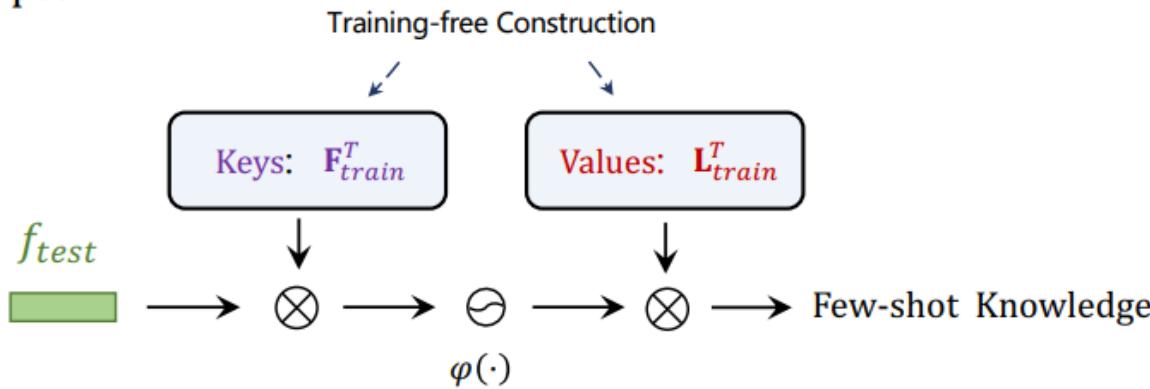
$$\mathbf{W}^\star = \beta A_t(\mathbf{W})^T + (1 - \beta) \mathbf{W}.$$

# Tip-Adapter

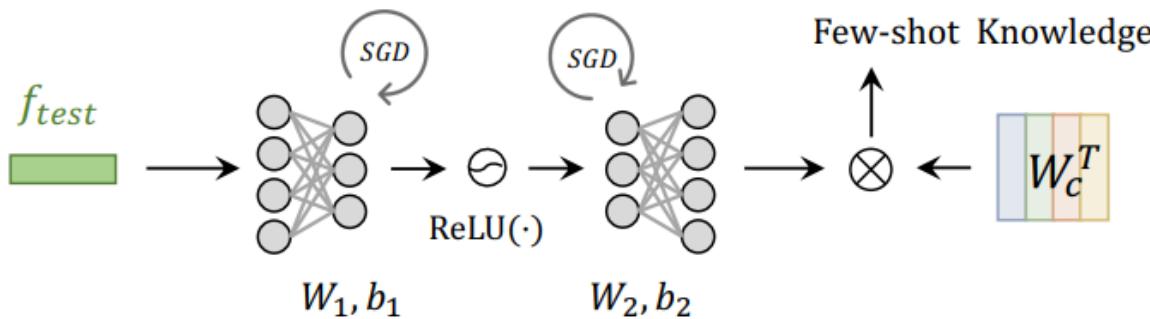


# Relations with CLIP-Adapter

Tip-Adapter:



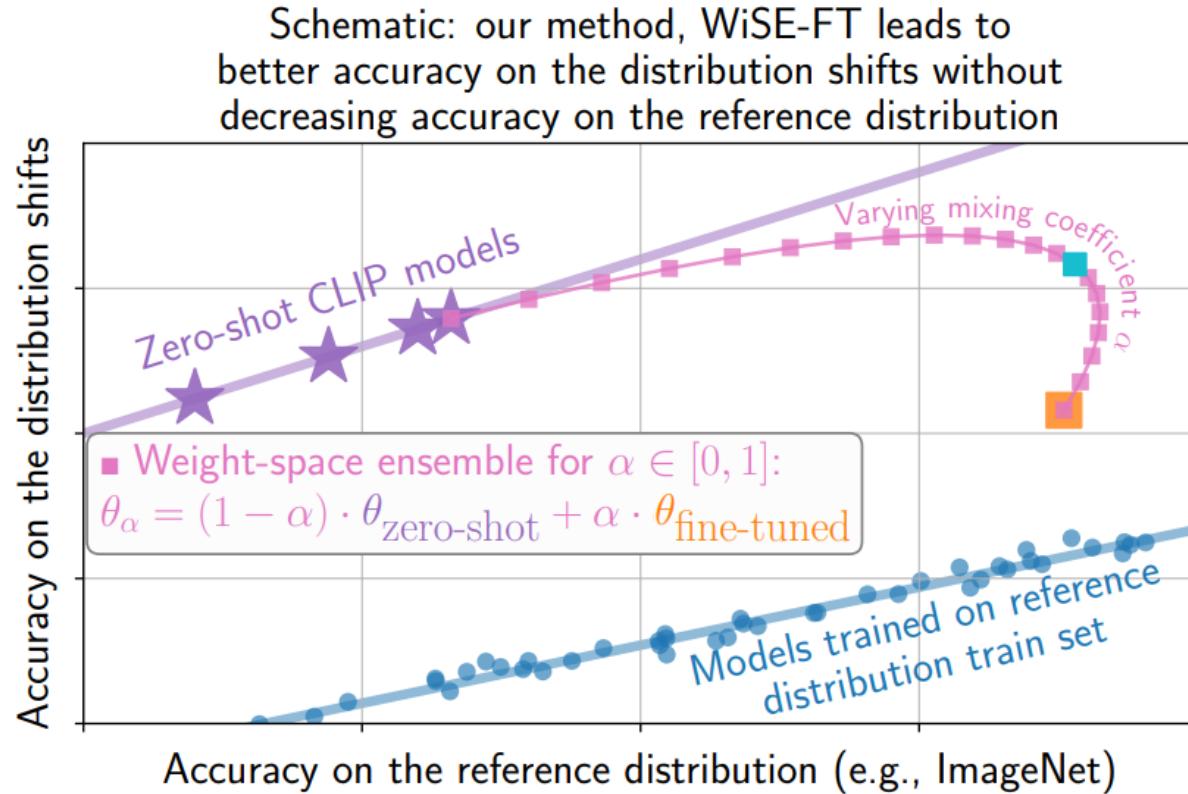
CLIP-Adapter:



# Contents

- Prompt
- Adapter
- Controllable FineTuning

# Robust fine-tuning of zero-shot models (CVPR2022)



Finetune后将模型参数插值

# Full ImageNet

| Method                                | RN-50             | RN-101            | ViT-B/32          | ViT-B/16          |
|---------------------------------------|-------------------|-------------------|-------------------|-------------------|
| Zero-Shot CLIP (Radford et al., 2021) | 58.2              | 61.5              | 62.0              | 66.9              |
| Linear-Probing CLIP                   | 72.8              | 75.5              | 76.0              | 79.5              |
| Prompt Ensembling                     | 60.4 <sup>†</sup> | 62.5 <sup>†</sup> | 63.7 <sup>†</sup> | 68.7 <sup>†</sup> |
| CoOp (Zhou et al., 2021)              | 65.6              | 67.8              | 68.0              | 72.4              |
| Target Optimization                   | 71.4              | 73.2              | 74.0              | 78.1              |
| CLIP-Adapter (Gao et al., 2021)       | 61.3 <sup>‡</sup> | -                 | -                 | -                 |
| DeFo (ours)                           | <b>73.2</b>       | <b>75.5</b>       | <b>76.2</b>       | <b>80.2</b>       |

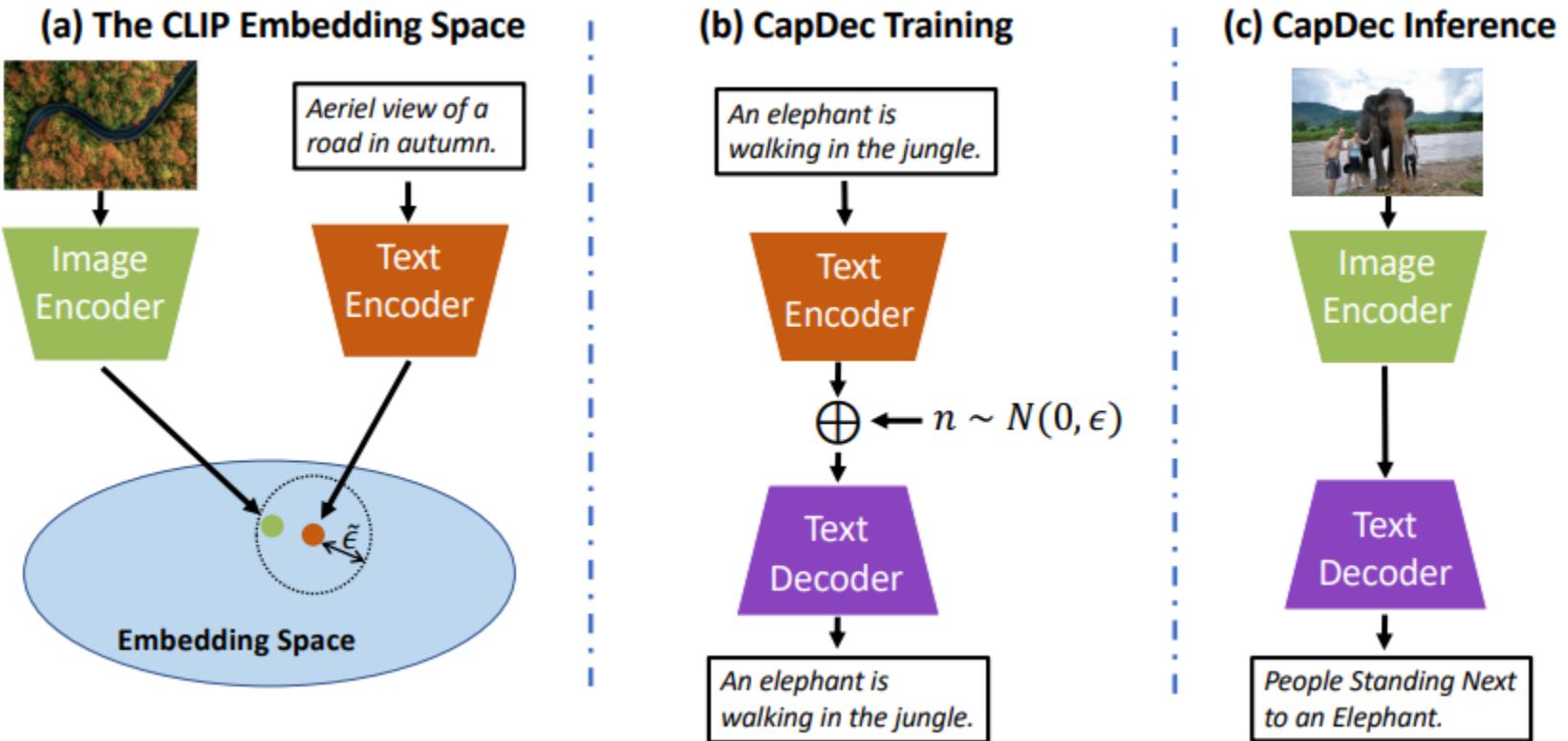
|                                 |             |
|---------------------------------|-------------|
| Wise-FT (LP, $\alpha = 0.5$ )   | <u>78.2</u> |
| Wise-FT (LP, optim. $\alpha$ )  | 80.0        |
| Wise-FT (E2E, $\alpha = 0.5$ )  | <b>82.6</b> |
| Wise-FT (E2E, optim. $\alpha$ ) | 81.7        |

数据较充足的条件下，Finetune backbone在训练数据集上效果较好

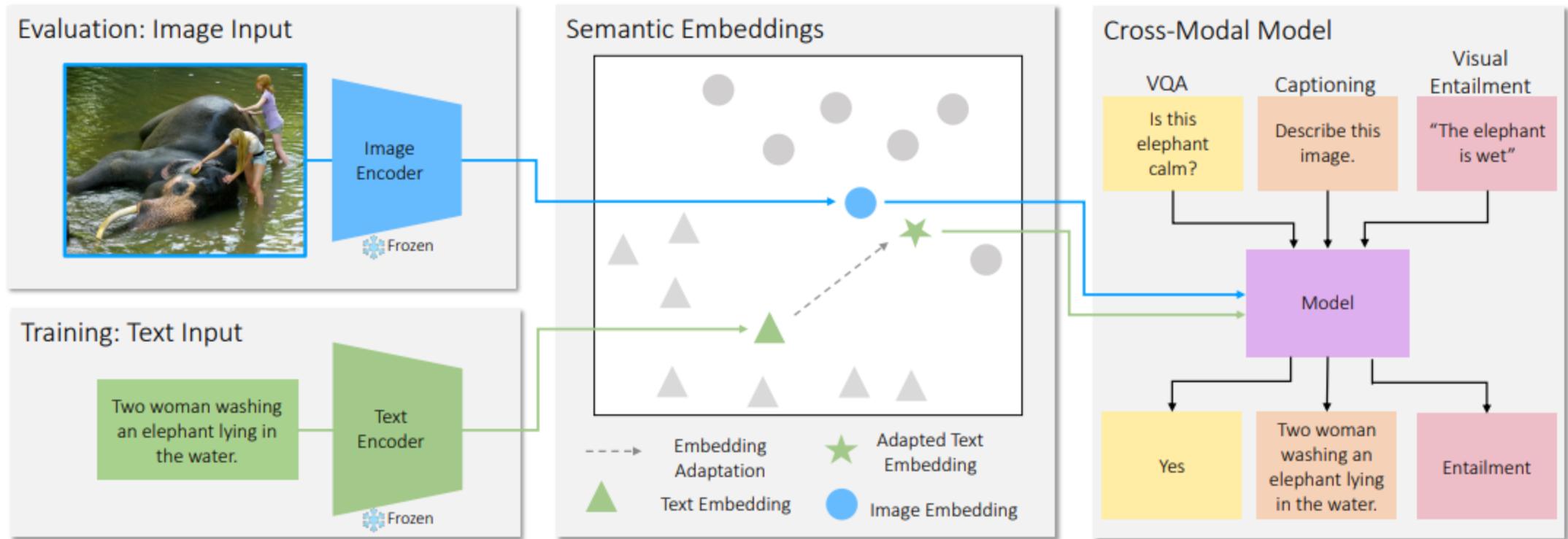
# Vision-Language Pretraining

- Vision-Language Models
- Classification: Parameter-efficient Fine-tuning
- More Challenging Tasks
- More Visual Modalities

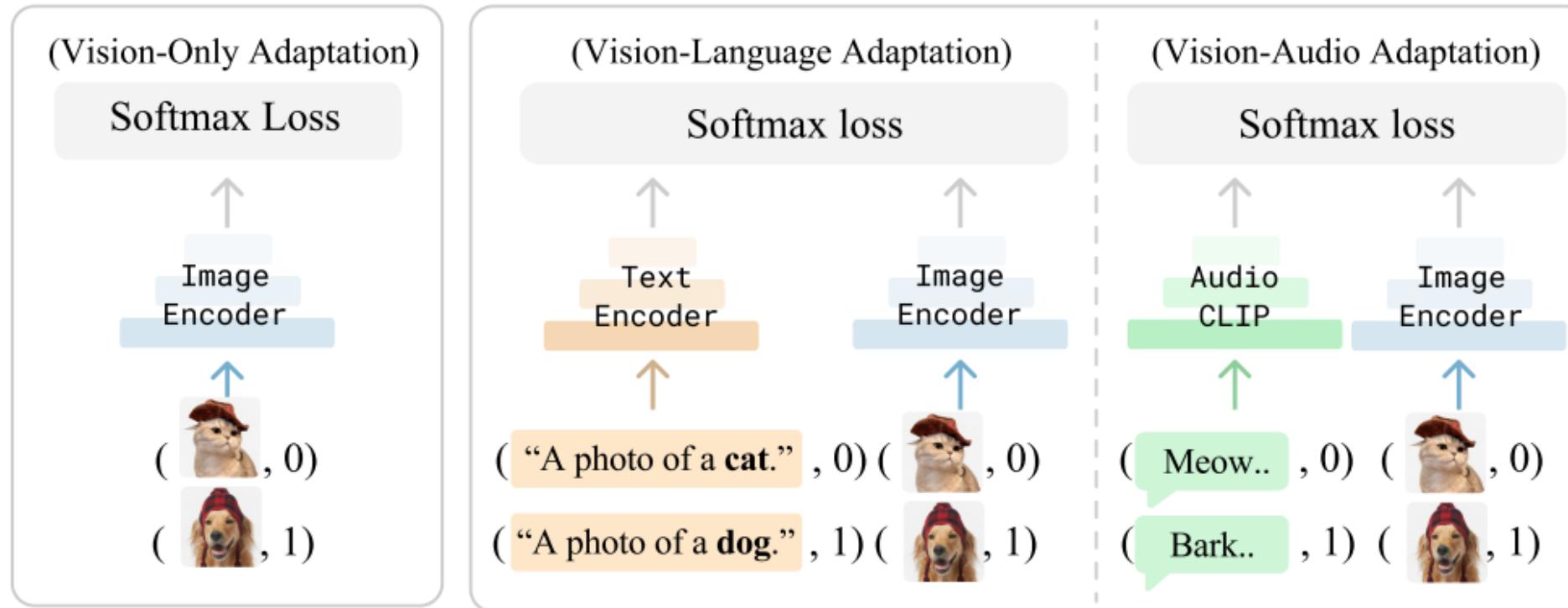
# CapDec (EMNLP2022)



# Learning Visual Tasks Using only Language Data

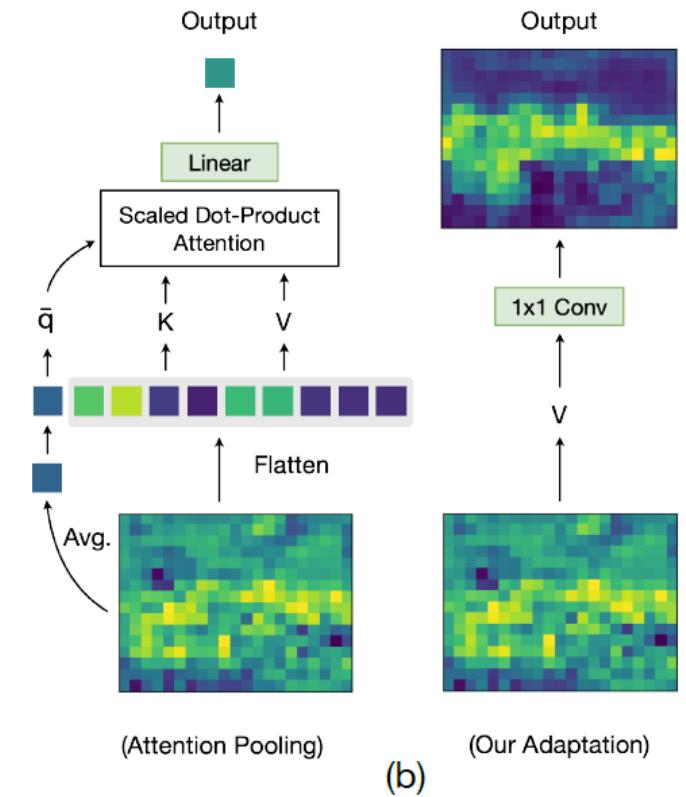
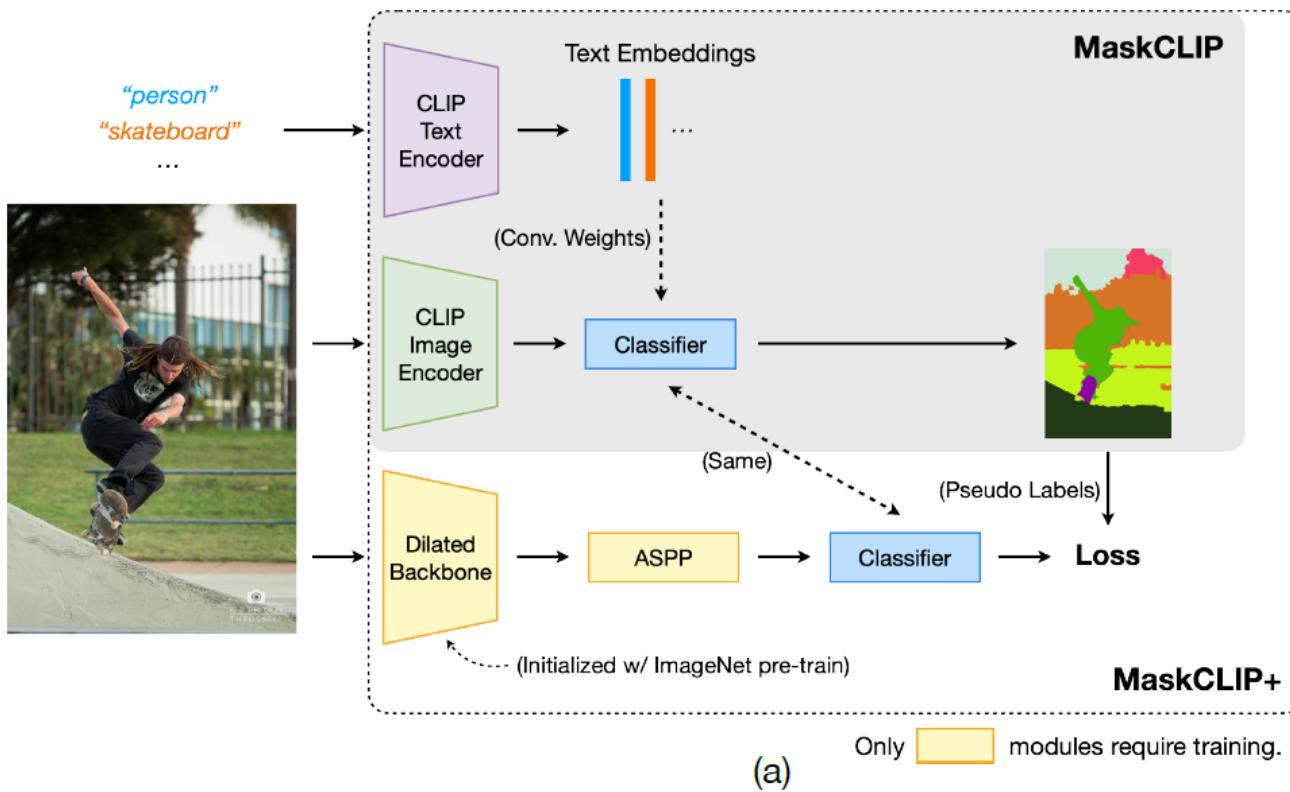


# Multimodality Helps Unimodality: Cross-Modal Few-Shot Learning with Multimodal Models

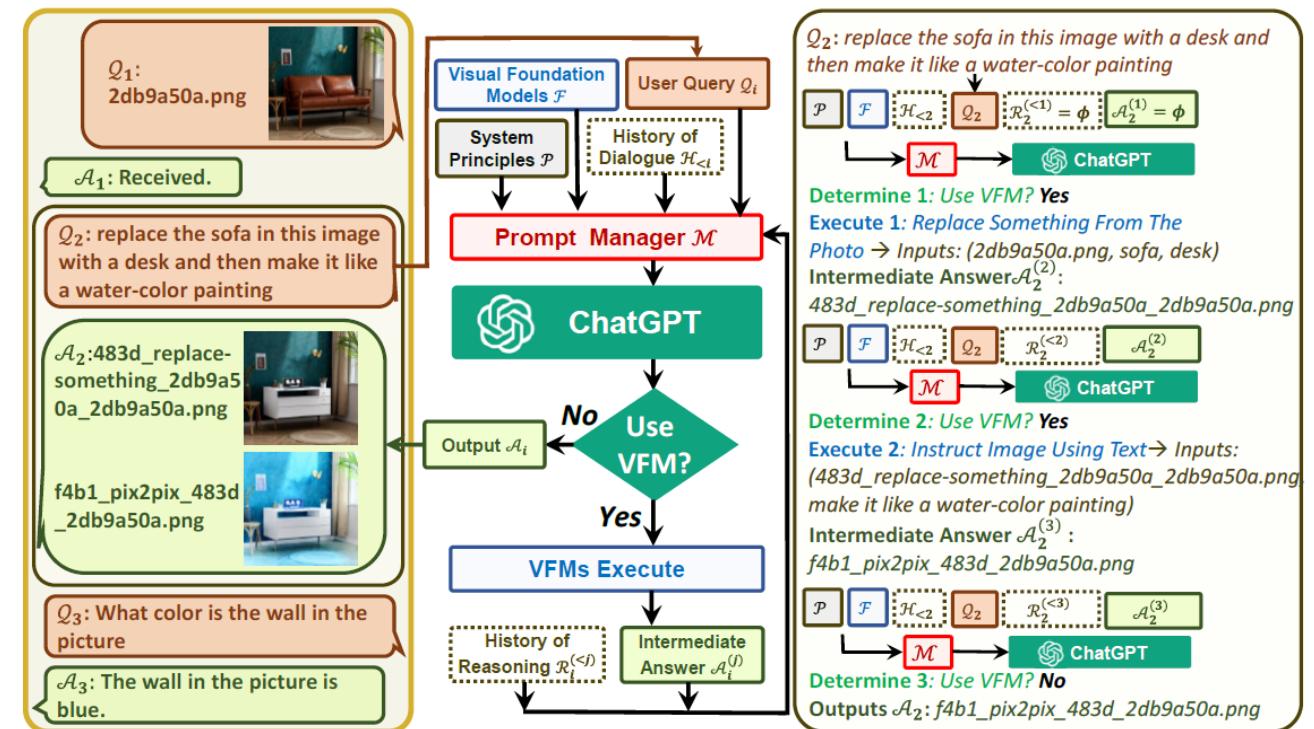
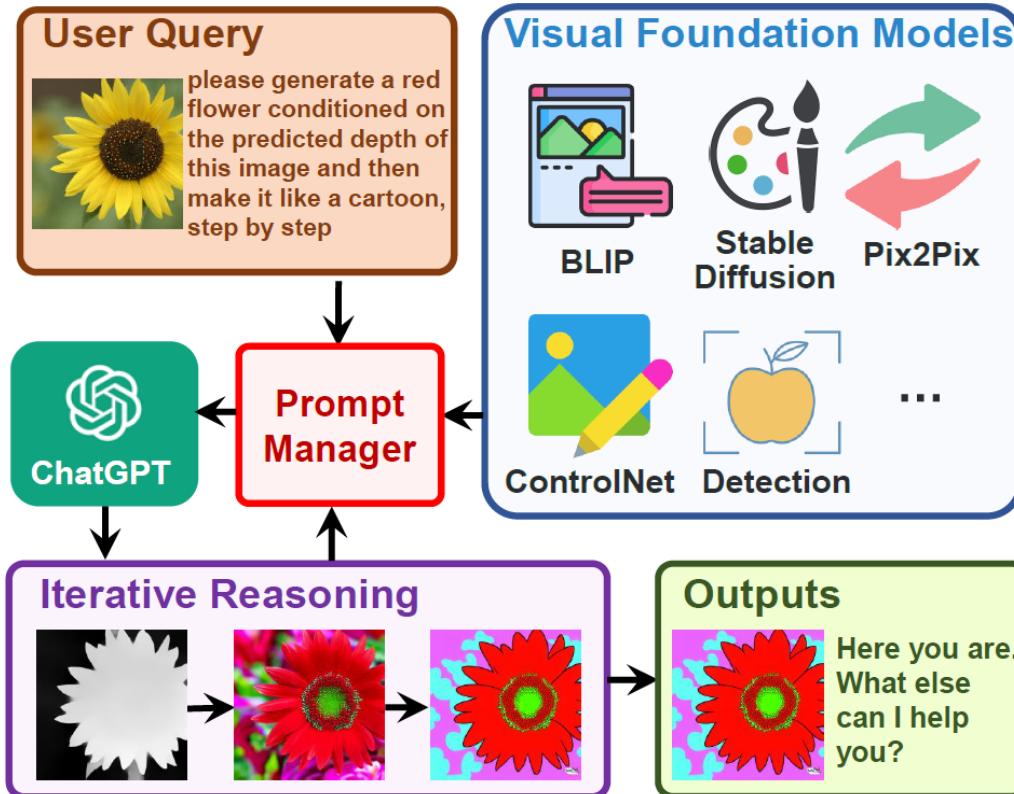


扩展了任务设定，Multi-modal training & Multi-modal testing

# MaskCLIP (NTU)



# Visual ChatGPT



# Visual ChatGPT

*Q<sub>1</sub>* : Hello, who are you?

*A<sub>1</sub>* : Hi there! I'm Visual ChatGPT. I'm here to help you with a wide range of text and visual related tasks. How can I help you?

*Q<sub>2</sub>* : I like drawing, but I'm not good at drawing, can you help me? like drawing an apple.

*A<sub>2</sub>* : Sure, I can help you draw an apple. Here is the image I generated for you:



image/a3453226.png

*Q<sub>3</sub>*: Draw an image:



*A<sub>3</sub>* : Received. What do you want to know / operate it? You can use image/dcd3af54.png to represent this image in our dialogue.

*Q<sub>4</sub>*: The image/dcd3af54.png is my sketch of an apple and a drinking glass, can you please help me to improve it?

*A<sub>4</sub>* : I generated a new image based on your sketch. The new image is saved as:



image/569b\_scribble2image\_dcd3af54\_dcd3af54.png

*Q<sub>5</sub>*: Looks good. Can you make the image/569b\_scribble2image\_dcd3af54\_dcd3af54.png into a watercolor painting?

*A<sub>5</sub>* : I have used the Instruct Image Using Text tool to generate a watercolor painting based on your input. The new image is saved as:



image/0849\_pix2pix\_569b\_dcd3af54.png

*Q<sub>6</sub>*: Wow this is very beautiful. Can you tell me what color this background is?

*A<sub>6</sub>* : The background color of the image/0849\_pix2pix\_569b\_dcd3af54.png is blue.

*Q<sub>7</sub>*: Correct answer. Can you remove this apple in this picture? and then describe the image?

*A<sub>7</sub>* : The image/2ca6\_remove-something\_0849\_dcd3af54.png now contains a drinking glass with a blue background.



image/2ca6\_remove-something\_0849\_dcd3af54.png

*Q<sub>8</sub>*: There are still some shadows on the table in the image/2ca6\_remove-something\_0849\_dcd3af54.png, can you help me to replace the table with a black table?

*A<sub>8</sub>* : I have replaced the table with a black table. The image is saved as



image/7d56\_replace-something\_2ca6\_dcd3af54.png

*Q<sub>9</sub>*: Good. Now help me remove the cup too.

*A<sub>9</sub>* : The image/6c04\_remove-something\_7d56\_dcd3af54.png now contains a black table with a blue background and no cup.



image/6c04\_remove-something\_7d56\_dcd3af54.png

*Q<sub>10</sub>*: Can you replace the background of image/6c04\_remove-something\_7d56\_dcd3af54.png with an old man sitting with a lot of vicissitudes in his face?

*A<sub>10</sub>* : The new image is saved as



image/5cd1\_replace-something\_6c04\_dcd3af54.png

*Q<sub>11</sub>*: This picture looks good. Can you describe the color of the old man's clothes?

*A<sub>11</sub>* : The old man's clothes are blue.

*Q<sub>12</sub>*: Correct. I want you to detect the depth map of the image/5cd1\_replace-something\_6c04\_dcd3af54 first, and show the depth image, then generate a new image based on this depth map, and finally describe the final image.

*A<sub>12</sub>* : The image/38a5\_depth2image\_5cd1\_dcd3af54.png, which contains 1 chair, 1 person. The person is wearing a white shirt and has a beard. The depth image of the image/5cd1\_replace-something\_6c04\_dcd3af54 is image/d149\_depth\_5cd1\_dcd3af54.png.



image/d149\_depth\_5cd1\_dcd3af54.png



image/38a5\_depth2image\_5cd1\_dcd3af54.png

*Q<sub>13</sub>*: Ok, forget apples and old people. I want you to help me generate a picture of a girl and a puppy playing on the beach.

*A<sub>13</sub>* : I have generated a picture of a girl and a puppy playing on the beach. The image is



image/ab0da279.png

*Q<sub>14</sub>*: It looks very nice. Can you replace this puppy with a kitten?

*A<sub>14</sub>* : I have replaced the puppy with a kitten in the image. The new image is



image/d1e1\_replace-something\_ab0da279\_ab0da279.png

*Q<sub>15</sub>*: Well done. Can you turn this image/d1e1\_replace-something\_ab0da279\_ab0da279.png into a pencil drawing?

*A<sub>15</sub>* : I have generated a pencil drawing of a girl and a kitten playing on the beach. The image is



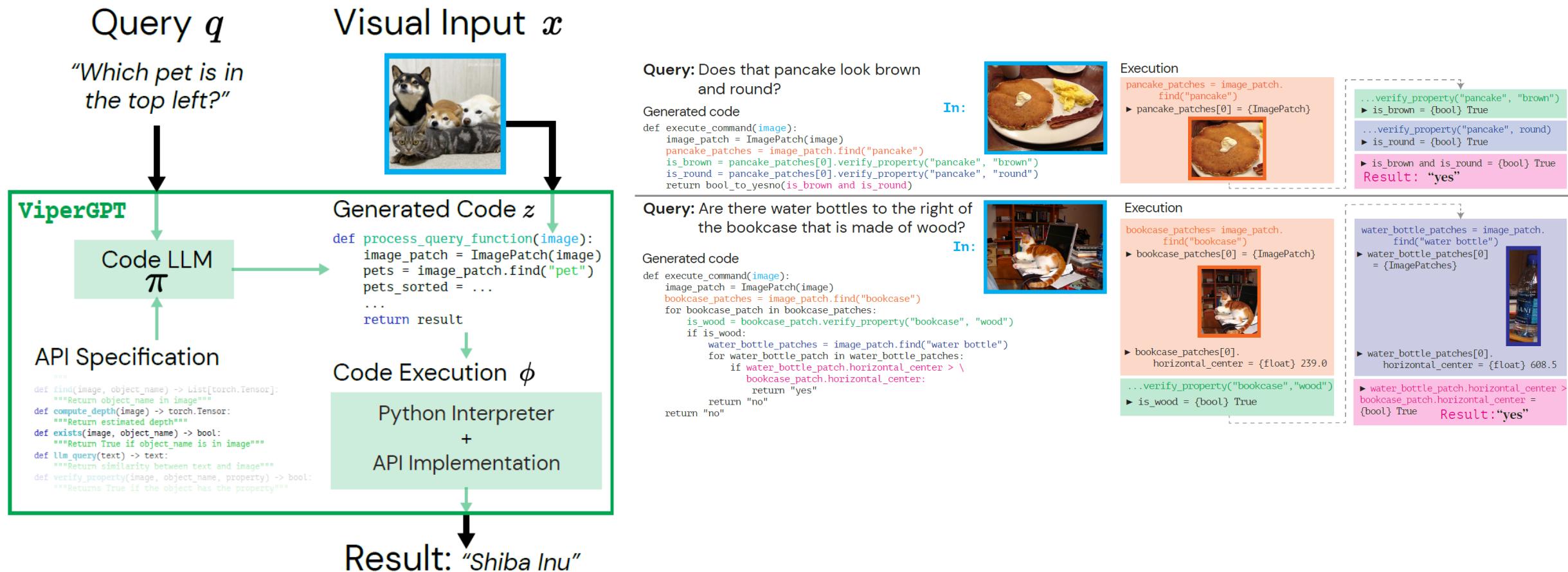
image/6c71\_pix2pix\_d1e1\_ab0da279.png

*Q<sub>16</sub>*: Thanks very much. You did a good job. Goodbye!

*A<sub>16</sub>* : You're welcome! Have a great day!

*A<sub>16</sub>* : You're welcome! Have a great day!

# ViperGPT



# Vision-Language Pretraining

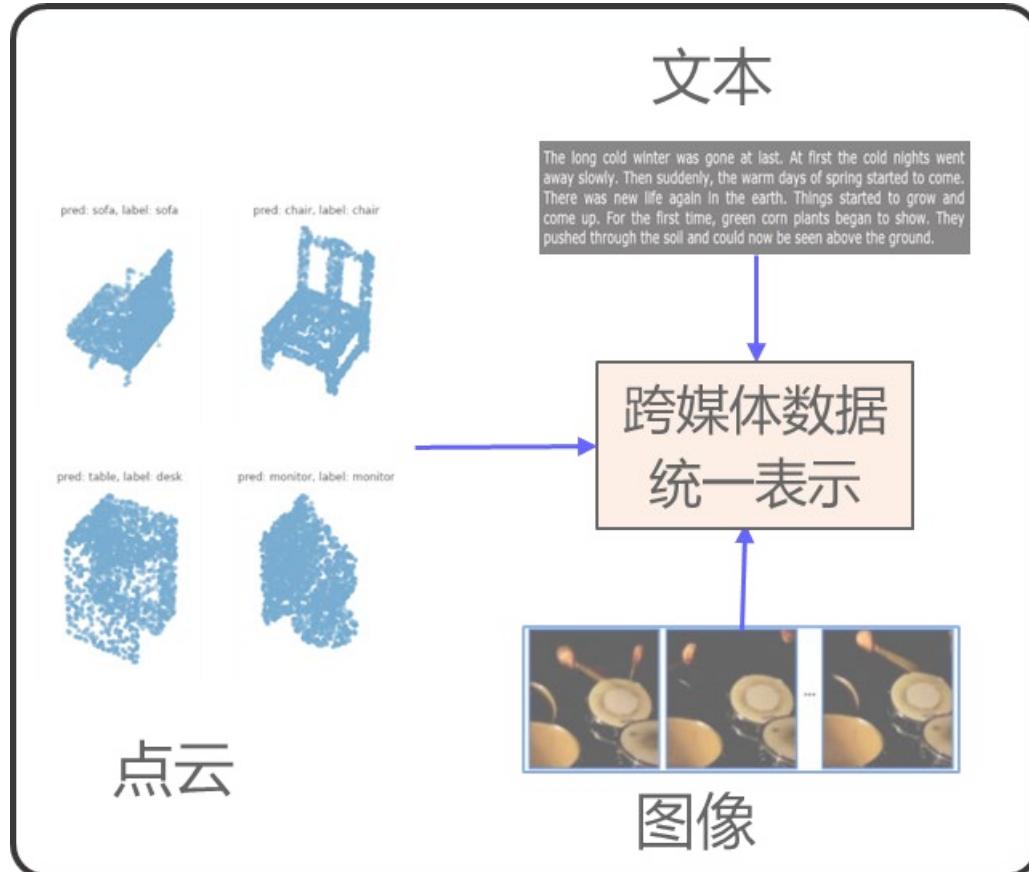
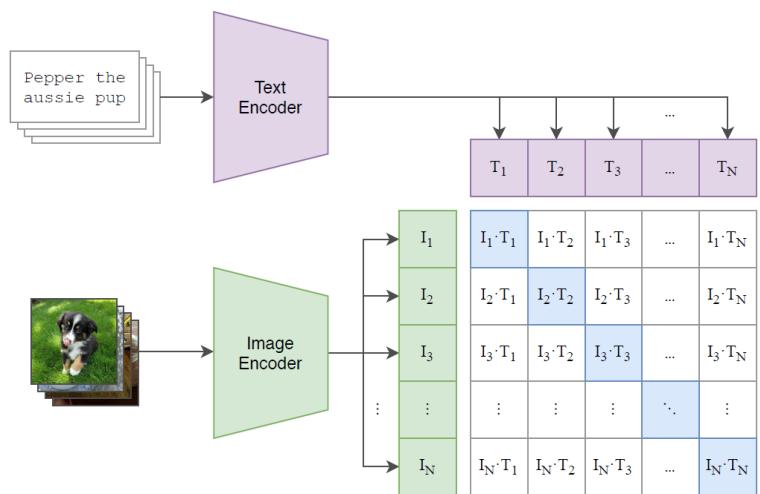
- Vision-Language Models
- Classification: Parameter-efficient Fine-tuning
- More Challenging Tasks
- More Visual Modalities

# 跨媒体数统一表示

- 图像-文本的统一表示

- CLIP
- ALIGN

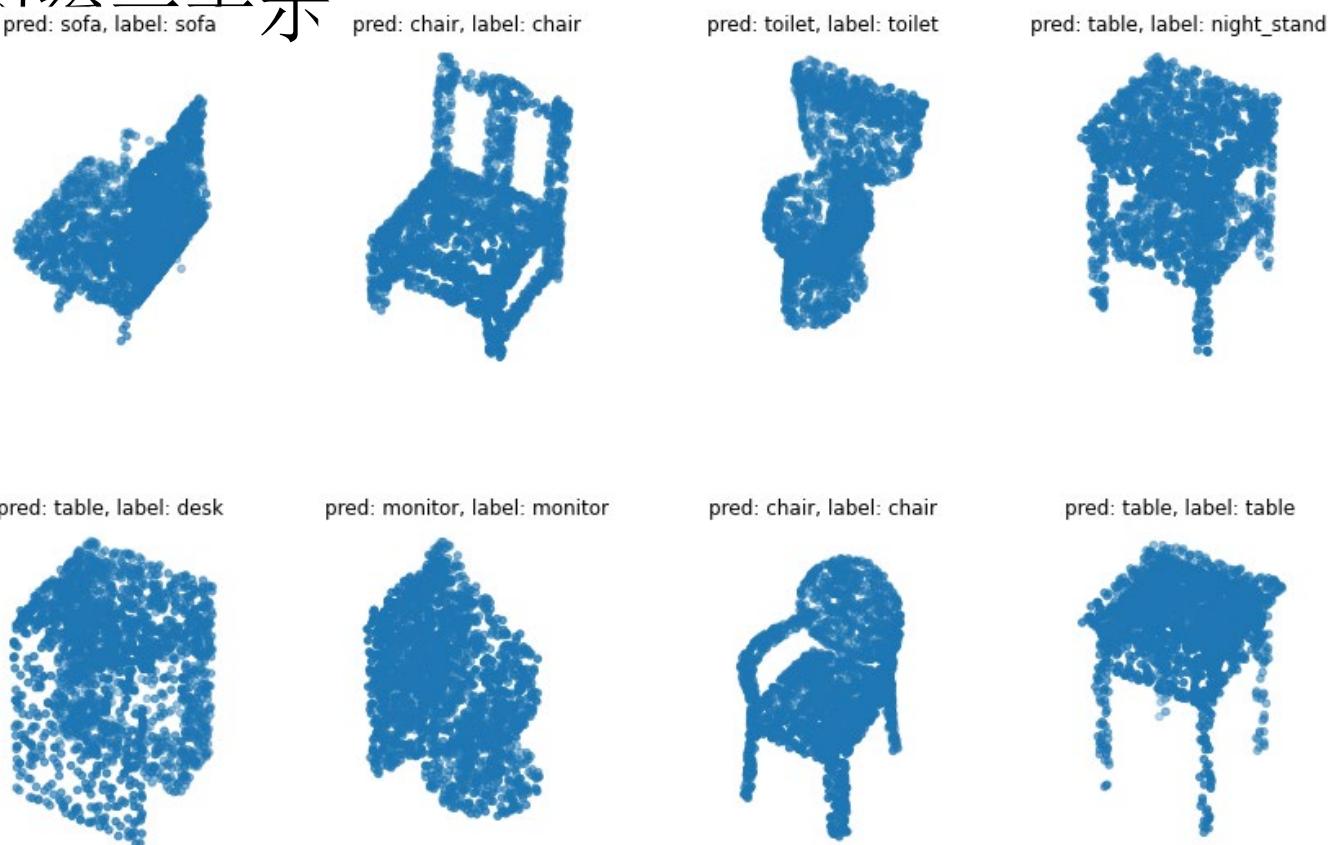
- 其他视觉模态



# 跨媒体数统一表示

- 3D点云-图像-文本的统一表示

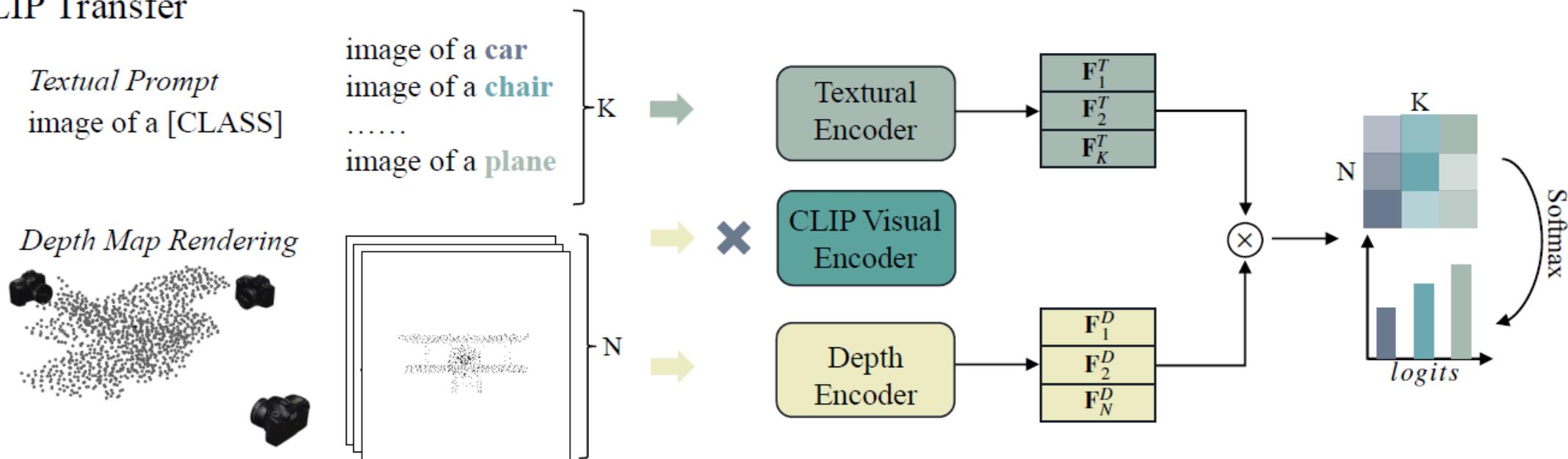
- 点云分类
- 点云检测
- 点云分割
- .....



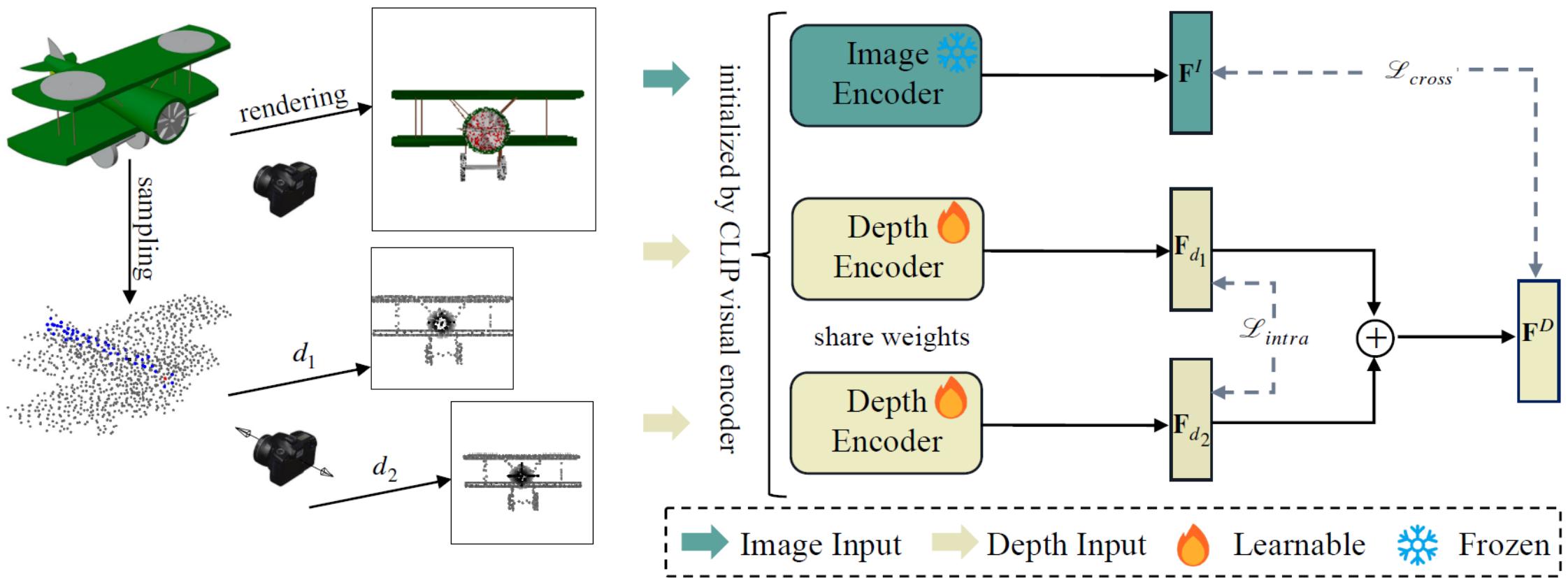
# 跨媒体数统一表示

- 直接方法存在的问题

CLIP Transfer



# CLIP2Point

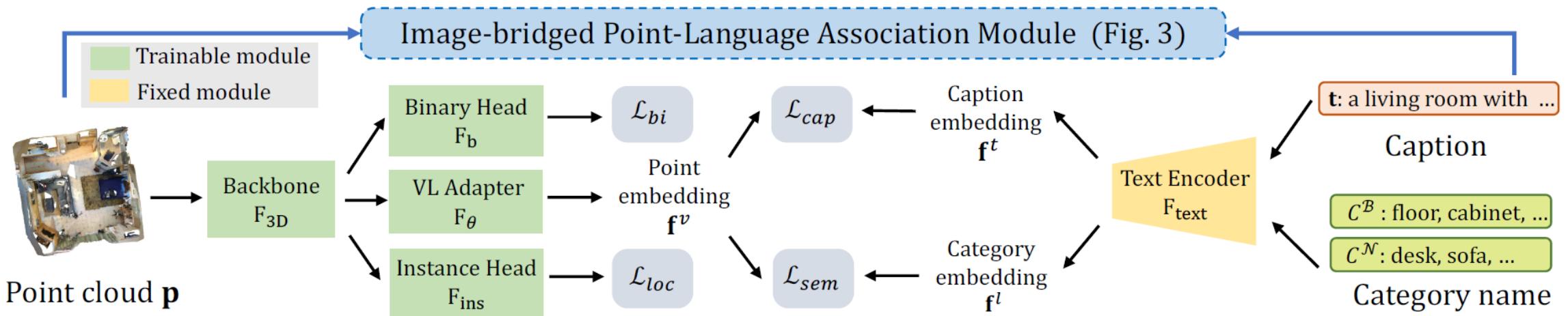


$$\mathcal{L}_{intra} = \frac{1}{2N} \sum_{i=1}^N (l_{intra}^i(d_1, d_2) + l_{intra}^i(d_2, d_1))$$

$$\mathcal{L}_{cross} = \frac{1}{2N} \sum_{i=1}^N (l_{cross}^i(D, I) + l_{cross}^i(I, D))$$

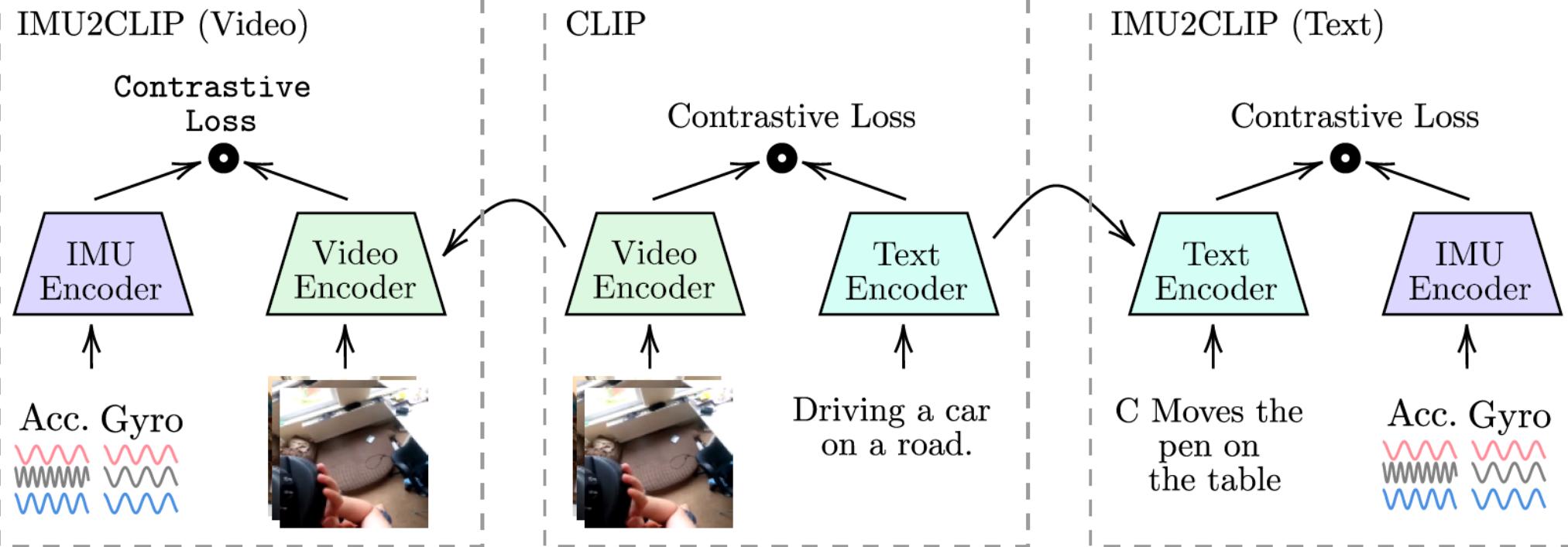
# 3D Scene Understanding

- 字节+港大

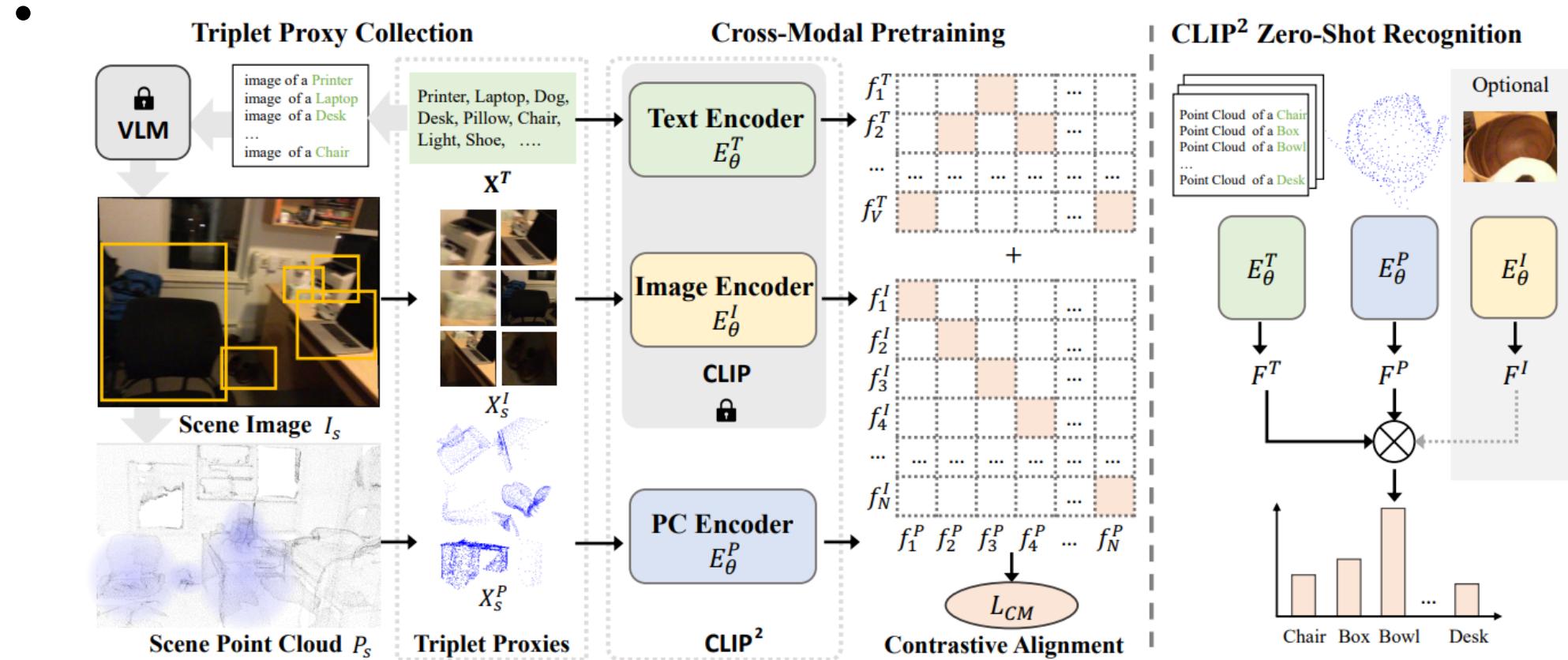


# IMU2CLIP

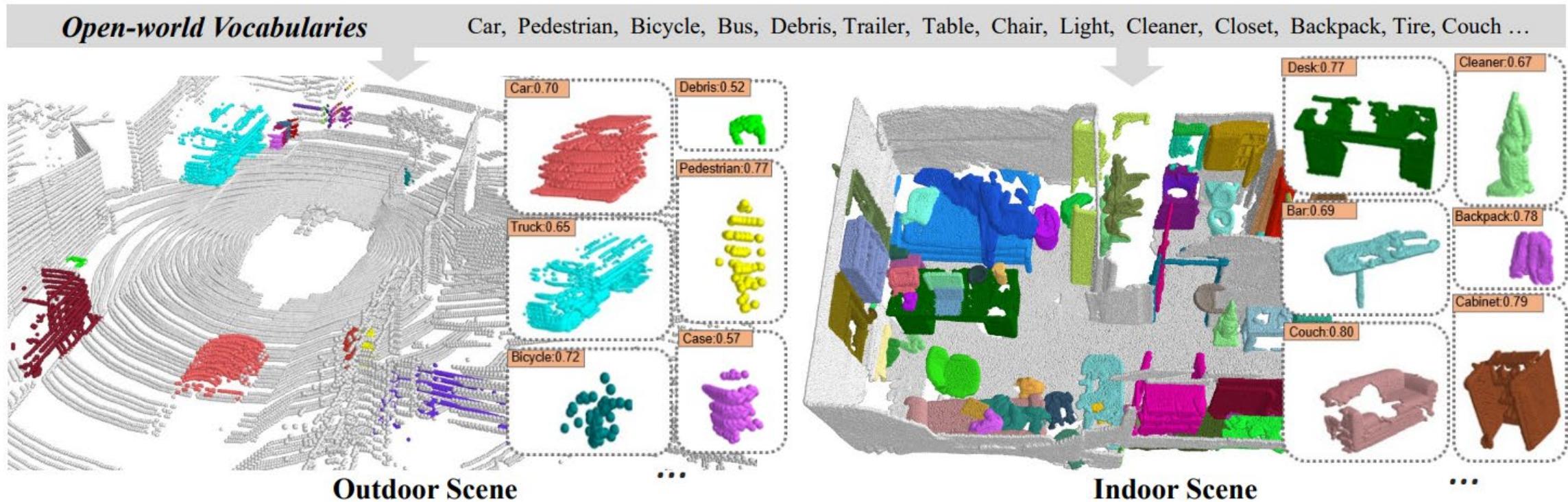
- Multimodal Learning



# CLIP<sup>2</sup>



# CLIP<sup>2</sup>



# Summary

- Rapid progress has been made in VL pre-training
- Challenges and opportunities