# Case Study:
# Vision Language Models

*Samrat Mukherjee*
PhD Student

Center for Machine Intelligence and Data Science
Indian Institute of Technology, Bombay

07 March, 2024

# Outline

- **Vision Language Models (VLMs).**
  - Introduction.
  - Examples.
- **Prompt Learning.**
  - Motivation.
  - Examples.
- **VLMs for Remote-Sensing use-cases.**
  - GeoChat : Grounded Large Vision-Language Model for Remote Sensing.

# Introduction

- Deep Neural Networks (DNNs) are data hungry.
  - Rely heavily on manually **labelled data** for training. (Eg. ImageNet, MS-COCO etc)
  - Train a DNN for each single visual recognition task.
- Vision-Language Models (VLMs).
  - Combine natural language understanding with computer vision capabilities.
- VLMs learn rich vision-language correlation.
  - Uses image-text pairs available on the Internet.
  - Enables zero-shot predictions on various visual recognition tasks with a single VLM.
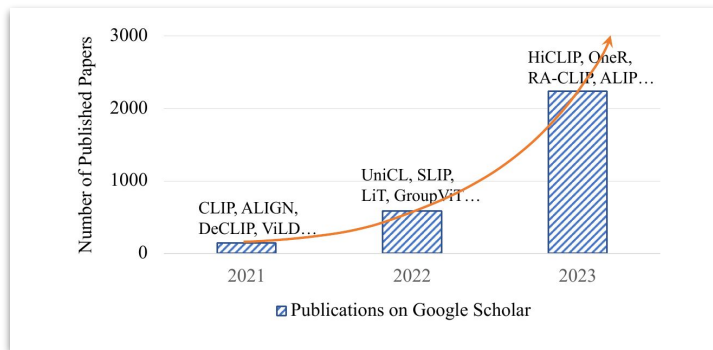


Fig 1: Trend of published work with VLMs [1]

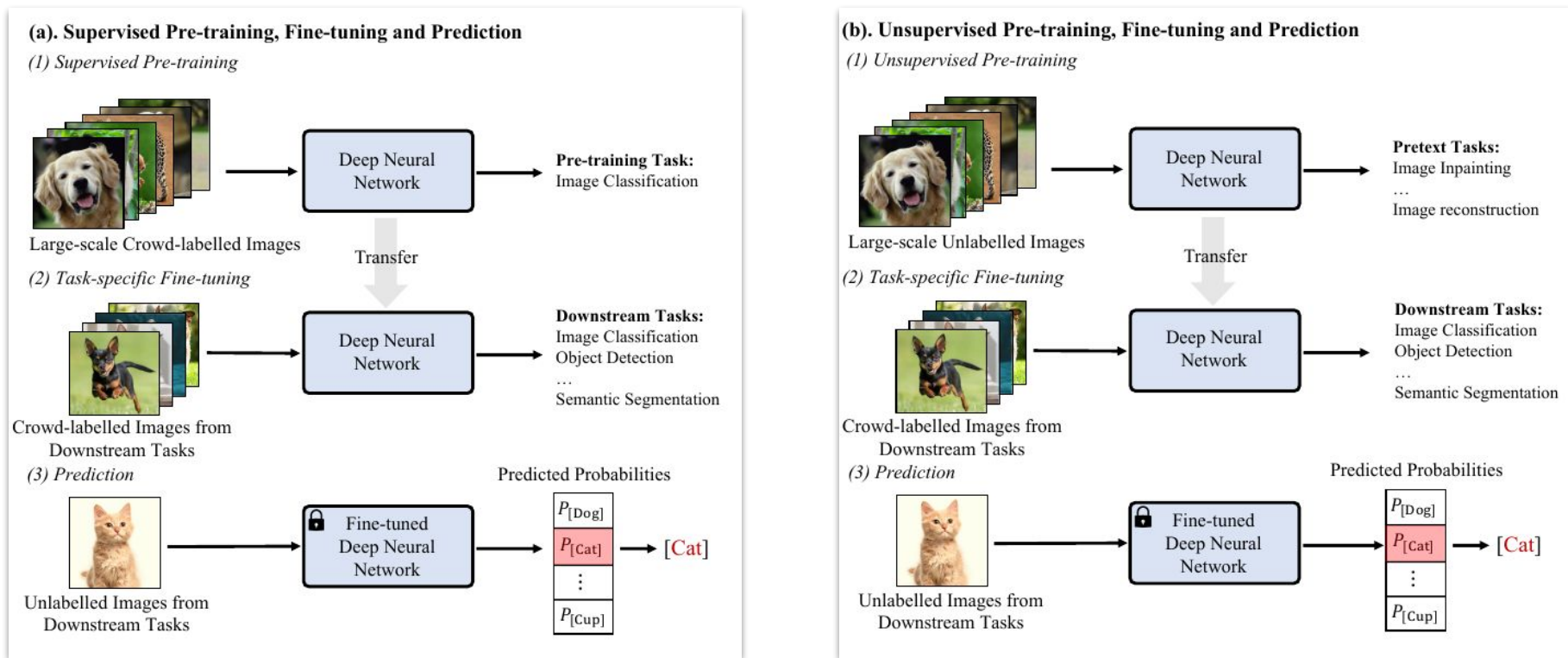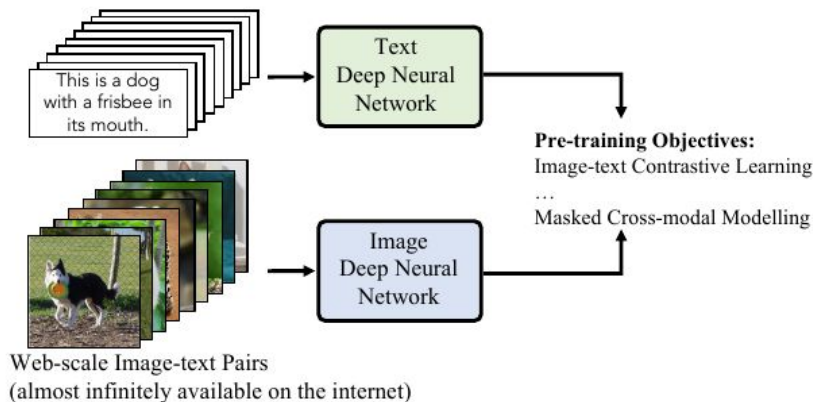# How is it different than previous approaches?



Fig.2: Deep Neural Network training paradigm (Supervised (left) and Uns .pervised (right)) [1]

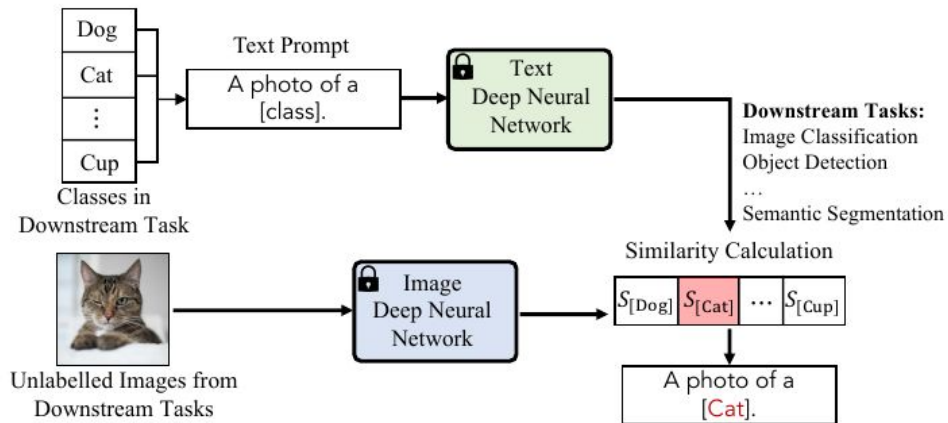# How is it different than previous approaches?



Fig. 3: DNN Training paradigm of VLMs. Observe the "absence" of Fine-tuning phase. [1]

- Pre-trained by certain vision-language objectives.
  - Captures rich vision-language knowledge.
  - Perform zero-shot predictions by matching the embeddings of images and texts.
    - Without fine-tuning on downstream tasks.

# Network architectures in VLMs

- **Architectures for learning Image features.**
  - Convolutional based architectures.
    - *ResNet*
    - VGG
    - EfficientNet
  - Transformer based architectures.
    - *Vision Transformer.*

- **Architectures for learning Text features.**
  - *Standard Transformer.*
  - GPT-2.
  - BERT.

# Pre-training objectives in VLMs

- **Contrastive objective.**
  - Image-text contrastive learning.
  - Image-text-label contrastive learning.

- **Generative objective.**
  - Masked image modelling.
  - Masked language modelling.
  - Masked cross-modal modelling.

- **Alignment objective.**
  - Image-text matching.
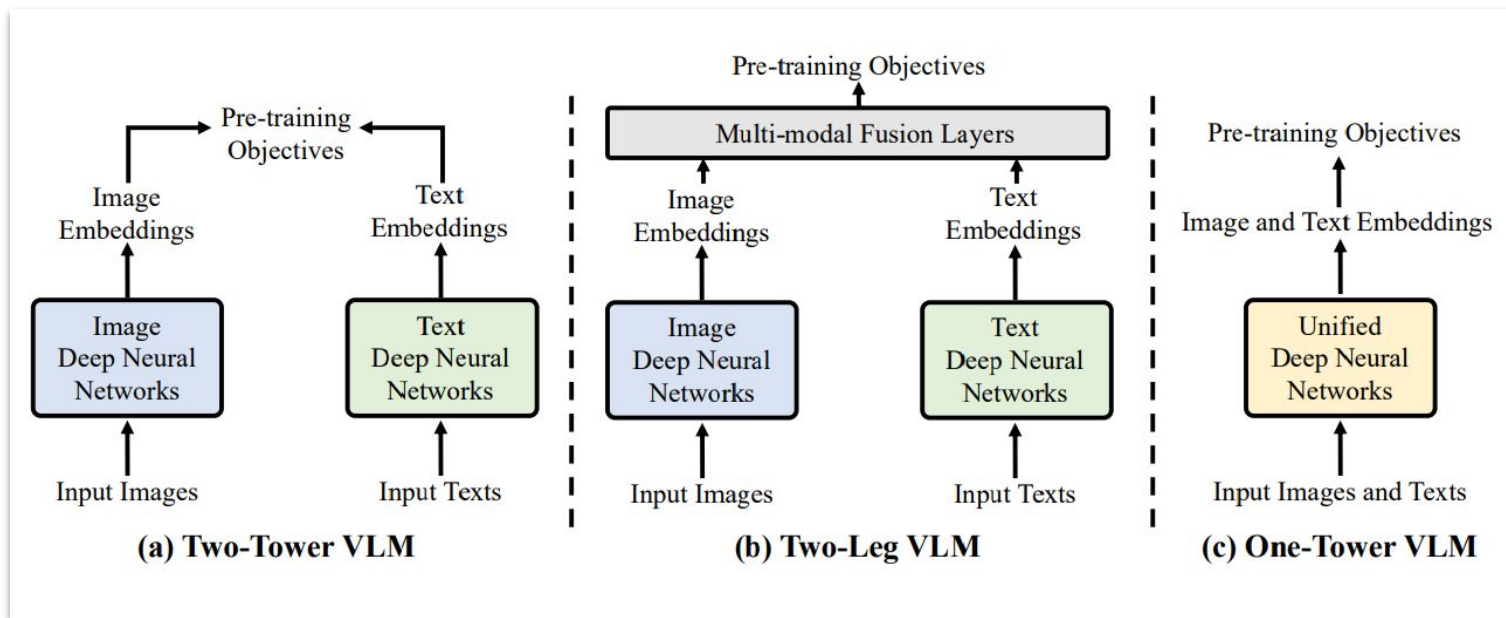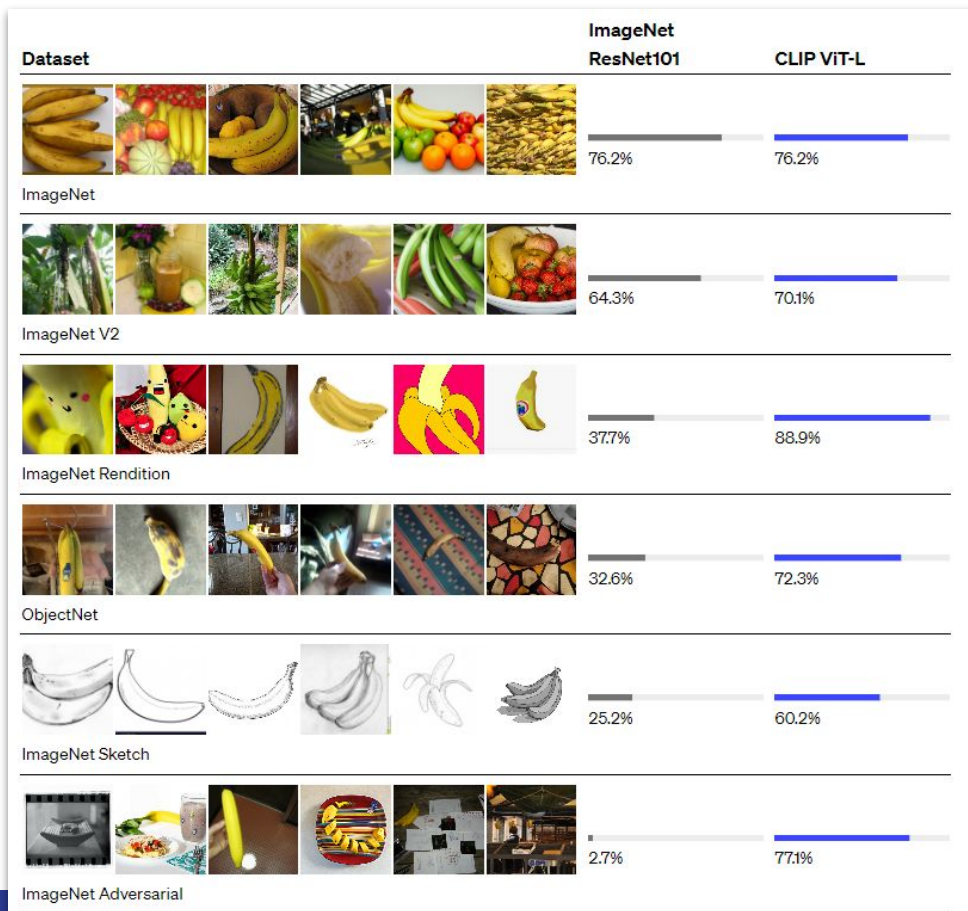  - Region-word matching.

# VLM pre-training framework



Fig. 4- Different pre-training frameworks [1]

# Contrastive Language Image Pre-training (CLIP)



Fig. 5: Motivation towards CLIP. [2]

- Introduced by OpenAI in 2021.

- Consistently performs better against ResNet101 on various dataset.

# Some more examples



Fig. 6: Some more examples and introduction to prompts [2]

# CLIP

- Scaling a simple pre-training task is sufficient to achieve competitive zero-shot performance on a great variety of image classification datasets.
- **Dataset:** the text paired with images found across the internet.
- **Proxy training task for CLIP:** Given an image, predict which out of a set of 32,768 randomly sampled text snippets, was actually paired with it in our dataset.
- To solve the task, CLIP learns broad visual concept recognition associated with their natural texts.
- As a result, CLIP models can then be applied to nearly arbitrary visual classification tasks.

- CLIP was designed to mitigate a number of major problems in the standard deep learning approach to computer vision:
  - Costly datasets
  - Narrow
  - Poor real-world performance

# Pre-training of CLIP



Fig. 7: Pre-training step for CLIP [2]

$$\mathcal{L}_{I \to T} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp\left(z_i^I \cdot z_i^T / \tau\right)}{\sum_{j=1}^{B} \exp(z_i^I \cdot z_j^T / \tau)}$$

$$\mathcal{L}_{T \to I} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp\left(z_i^T \cdot z_i^I / \tau\right)}{\sum_{j=1}^{B} \exp(z_i^T \cdot z_j^I / \tau)}$$

# Applying CLIP to new task



Fig. 8: Apply CLIP to a new task [2]

# Limitations

- CLIP struggles on more abstract or systematic tasks.
  - Counting the number of objects in an image
  - Predicting how close the nearest car is in a photo.
  - *On these two datasets, zero-shot CLIP is only slightly better than random guessing.*

- Zero-shot CLIP also struggles compared to task specific models on very fine-grained classification.
  - Such as telling the difference between car models, variants of aircraft, or flower species.

- CLIP is sensitive to wording or phrasing.
  - Require trial and error "prompt engineering" to perform well.

# Motivation towards learning the prompts

For pre-trained vision-language models, the text input, known as prompt, plays a key role in downstream datasets.

Manual prompt engineering is a non-trivial task and requires prior knowledge of the domain and the VLM.

Prompt learning research started flourishing lately in Natural Language Processing domain.

# CoOp (Context Optimization)



**Fig. 2  Overview of Context Optimization (CoOp).** The main idea is to model a prompt's context using a set of learnable vectors, which can be optimized through minimizing the classification loss. Two designs are proposed: one is unified context, which shares the same context vectors with all classes; and the other is class-specific context, which learns for each class a specific set of context vectors.

# CoOp (**Co**ntext **Op**timization)



CoOp models prompt's context words with learnable vectors, which could be initialized with either random values or pre-trained word embeddings.

CLIP prediction

$$p(y = i|\boldsymbol{x}) = \frac{\exp(\cos(\boldsymbol{w_i}, \boldsymbol{f})/\tau)}{\sum_{j=1}^{K} \exp(\cos(\boldsymbol{w_j}, \boldsymbol{f})/\tau)}$$

Using CoOp

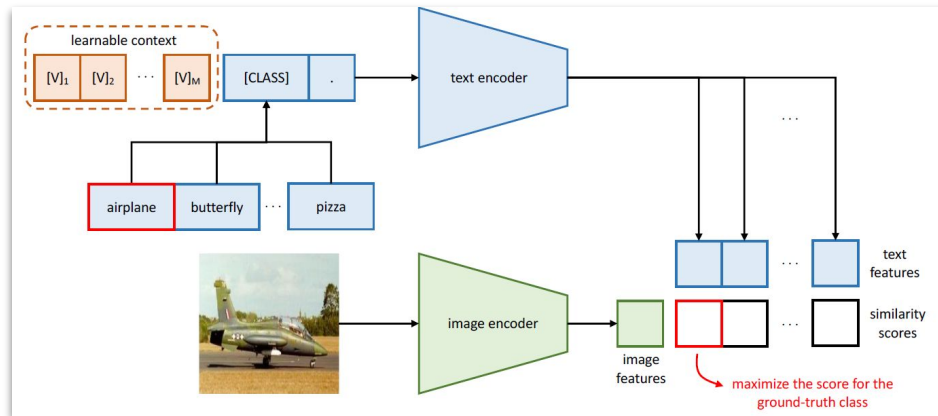$$p(y = i|\boldsymbol{x}) = \frac{\exp(\cos(g(\boldsymbol{t_i}), \boldsymbol{f})/\tau)}{\sum_{j=1}^{K} \exp(\cos(g(\boldsymbol{t_j}), \boldsymbol{f})/\tau)}$$

$$\boldsymbol{t} = [\text{V}]_1[\text{V}]_2 \dots [\text{V}]_M[\text{CLASS}],$$

$$\boldsymbol{t} = [\text{V}]_1 \dots [\text{V}]_{\frac{M}{2}}[\text{CLASS}][\text{V}]_{\frac{M}{2}+1} \dots [\text{V}]_M$$

# **Co**nditional **Co**ntext **Op**timization (CoCoOp)



(a) Both CoOp and CoCoOp work well on the base classes observed during training and beat manual prompts by a significant margin.

(b) The instance-conditional prompts learned by CoCoOp are much more generalizable than CoOp to the unseen classes.
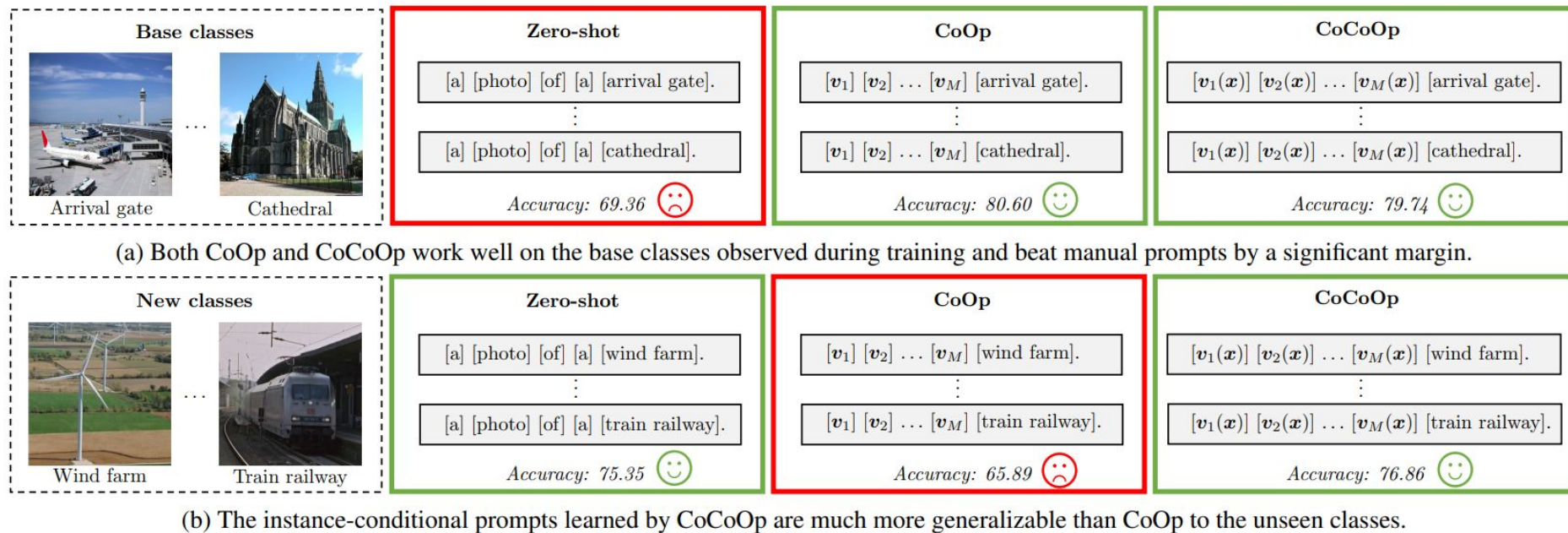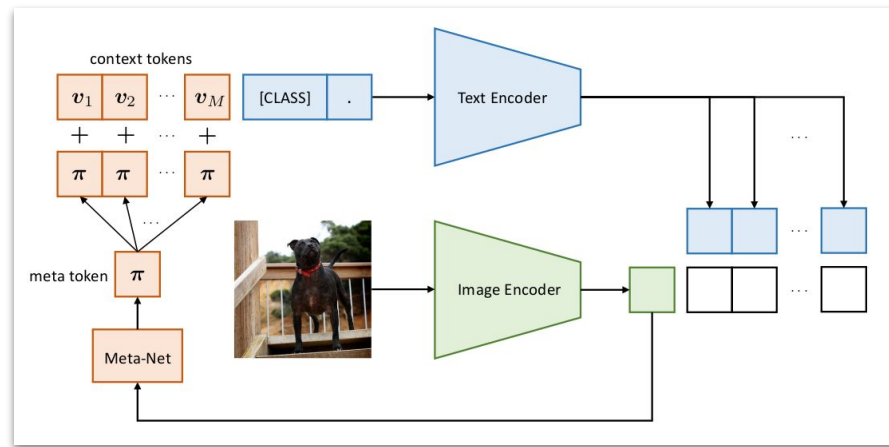
Figure 1. **Motivation of our research: to learn generalizable prompts**. The images are randomly selected from SUN397 [55], which is a widely-used scene recognition dataset.

# Conditional Context Optimization (CoCoOp)



- Instance-conditional context can generalize better because it shifts the focus away from a specific set of classes—for reducing overfitting—to each input instance. *Inspired from Image Captioning task.*

- For each input generate a conditional token from the Meta-Net which is then combined with the context vector. *(MetaNet : Linear → ReLU → Linear)*

$$p(y|\boldsymbol{x}) = \frac{\exp(\text{sim}(\boldsymbol{x}, g(\boldsymbol{t}_y(\boldsymbol{x})))/\tau)}{\sum_{i=1}^{K} \exp(\text{sim}(\boldsymbol{x}, g(\boldsymbol{t}_i(\boldsymbol{x}))/\tau)}.$$

$$\boldsymbol{t}_i(\boldsymbol{x}) = \{\boldsymbol{v}_1(\boldsymbol{x}), \boldsymbol{v}_2(\boldsymbol{x}), \dots, \boldsymbol{v}_M(\boldsymbol{x}), \boldsymbol{c}_i\}$$

$$\boldsymbol{v}_m(\boldsymbol{x}) = \boldsymbol{v}_m + \boldsymbol{\pi} \qquad \boldsymbol{\pi} = h_{\boldsymbol{\theta}}(\boldsymbol{x})$$

# VLMs in Remote Sensing

- General-domain VLMs perform poorly for Remote Sensing (RS) scenarios.
  - Fabricated information when presented with RS domain-specific queries.
- Unique challenges introduced by RS imagery.
  - High-resolution RS imagery with diverse scale changes across categories.
  - Many small objects.

- Several VLMs have been proposed for different tasks like-
  - Scene understanding
  - Region based captioning. etc.

*But not a generic VLM…..*

# Geo-Chat

- The first versatile remote sensing VLM that offers multitask conversational capabilities with high-resolution RS images.
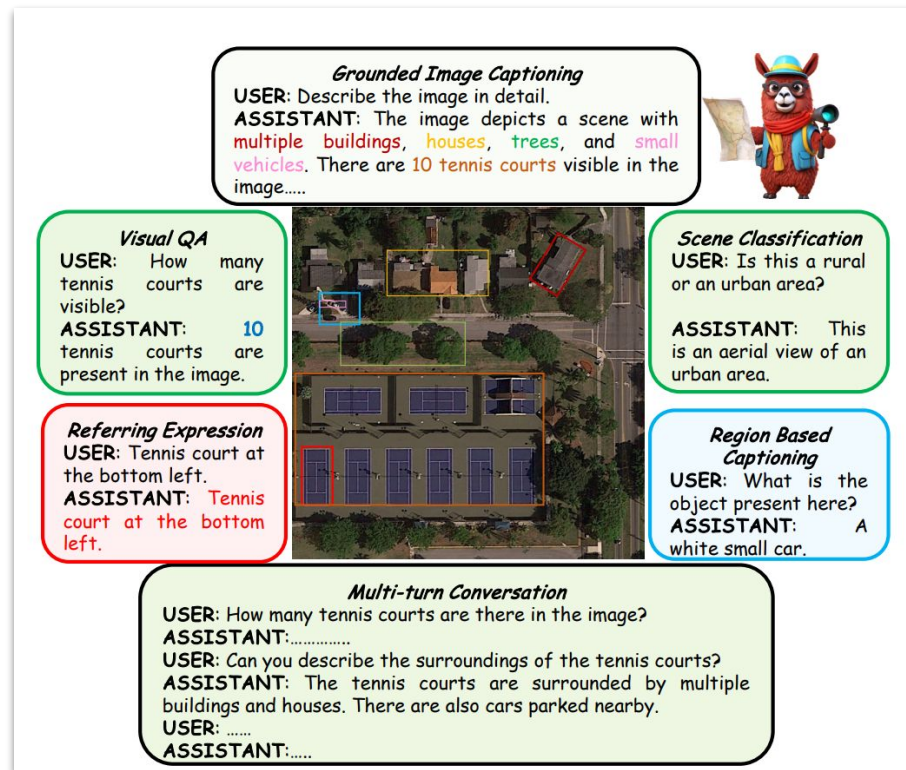- Answer image-level queries but also accepts region inputs to hold region-specific dialogue



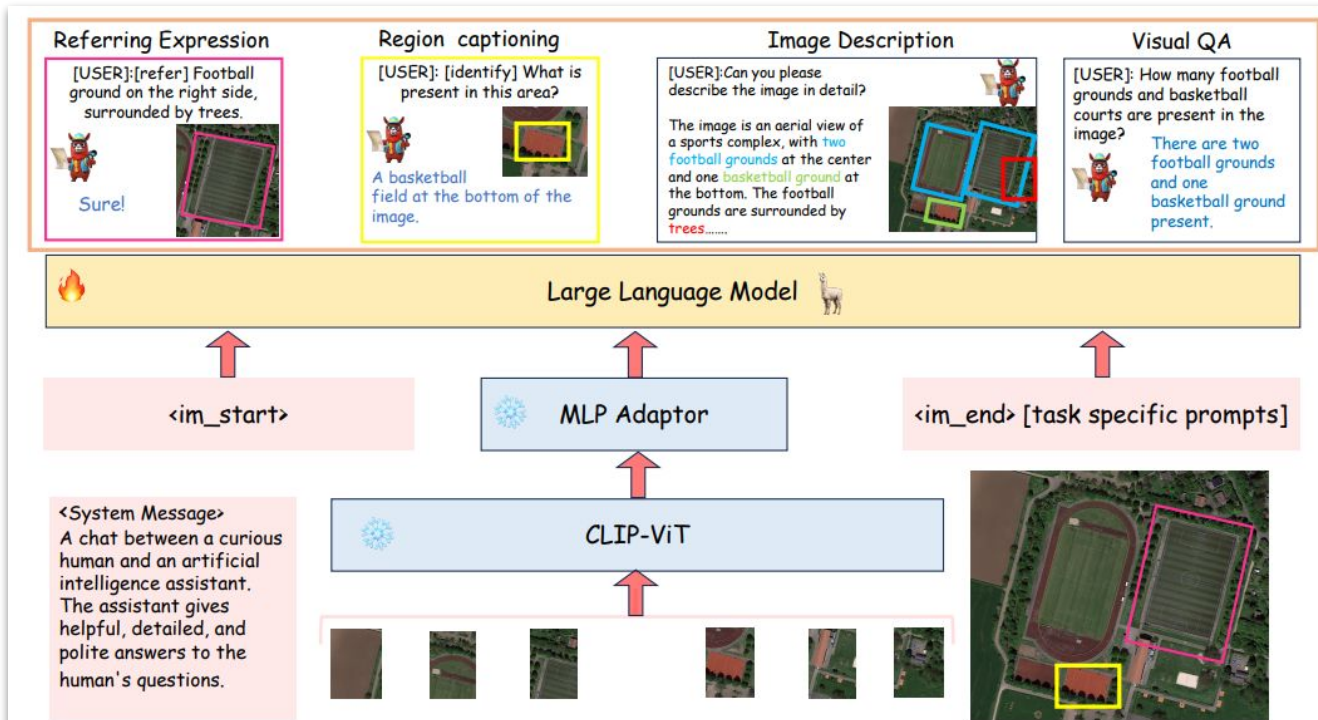Fig. 11: GeoChat and its capabilities [5]

# Geo-Chat



Figure 2. An overview of GeoChat - the first grounded large vision-language model for remote sensing. Given an image input together with a user query, a visual backbone is first used to encode patch-level tokens at a higher resolution via interpolating positional encodings. A multi-layer perceptron (MLP) is used to adapt vision-tokens to language space suitable for input to a Large Language Model (Vicuna 1.5). Besides visual inputs, region locations can also be input to the model together with task-specific prompts that specify the desired task required by the user. Given this context, the LLM can generate natural language responses interleaved with corresponding object locations. GeoChat can perform multiple tasks as shown on top e.g., scene classification, image/region captioning, VQA and grounded conversations.

# Reference

1. **Vision-Language Models for Vision Tasks: A Survey** Jingyi Zhang , Jiaxing Huang , Sheng Jin and Shijian Lu, https://arxiv.org/abs/2304.00685
2. **CLIP blog,** OpenAI, https://openai.com/research/clip
3. **Learning to Prompt for Vision-Language Models,** Kaiyang Zhou, Jingkang Yang, Chen Change Loy, Ziwei Liu, https://arxiv.org/abs/2109.01134
4. **Conditional Prompt Learning for Vision-Language Models,**Kaiyang Zhou, Jingkang Yang, Chen Change Loy, Ziwei Liu, https://arxiv.org/abs/2203.05557
5. **GeoChat : Grounded Large Vision-Language Model for Remote Sensing,** Kartik Kuckreja, Muhammad Sohail Danish, Muzammal Naseer, Abhijit Das, Salman Khan, Fahad Shahbaz Khan, https://arxiv.org/abs/2311.15826