

# VLM Mechanics, Issues and Variants BEV

D.A. Forsyth,

University of Illinois at Urbana Champaign

# Vision – language strategies

- Shared embedding space
  - Get
    - reliable image encoder
    - reliable language encoder
  - fine-tune them to match
    - as in CLIP
    - how? huge models, etc.
  - strategies
    - LORA
    - Prefix/Prompt tuning

# LORA finetuning

- Weights in most layers form a matrix
- Choose a (some) layer(s)
  - Fix weights
  - add low rank matrix to weight matrix
  - adjust low rank matrix
  - Very few weights
  - Significant improvements
- Curious fact:
  - Lora's are traded in the NSFW world

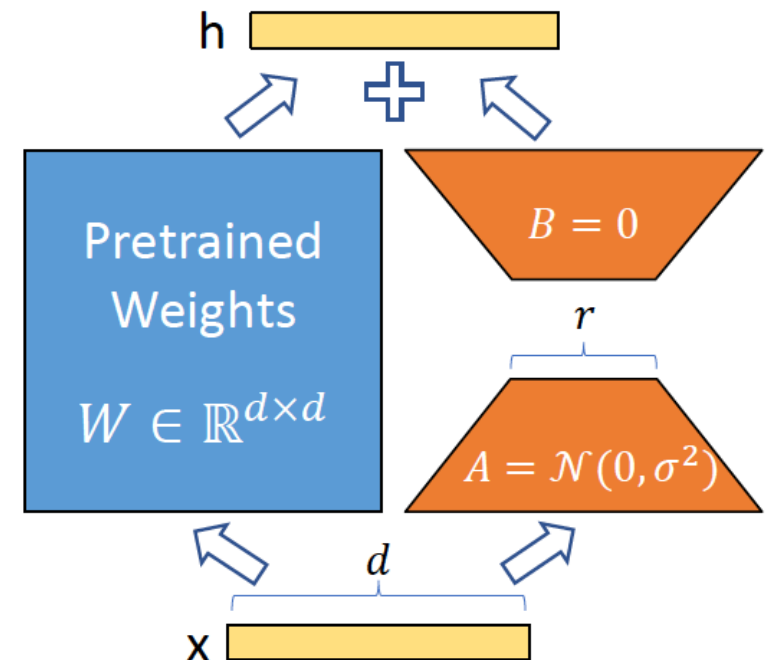


Figure 1: Our reparametrization. We only train  $A$  and  $B$ .

# LORA finetuning

Original model

$$\mathbf{x}_{\text{out}} = \text{activation}(\mathcal{W}\mathbf{x}_{\text{in}} + \mathbf{b})$$

$\mathbb{R}^{d \times d}$

LORA version

Freeze this

$$\mathbf{x}_{\text{out}} = \text{activation}([\mathcal{W} + \mathcal{A}\mathcal{B}]\mathbf{x}_{\text{in}} + \mathbf{b})$$

Learn these

$\mathbb{R}^{d \times n}$      $\mathbb{R}^{n \times d}$

# Prefix/Prompt tuning

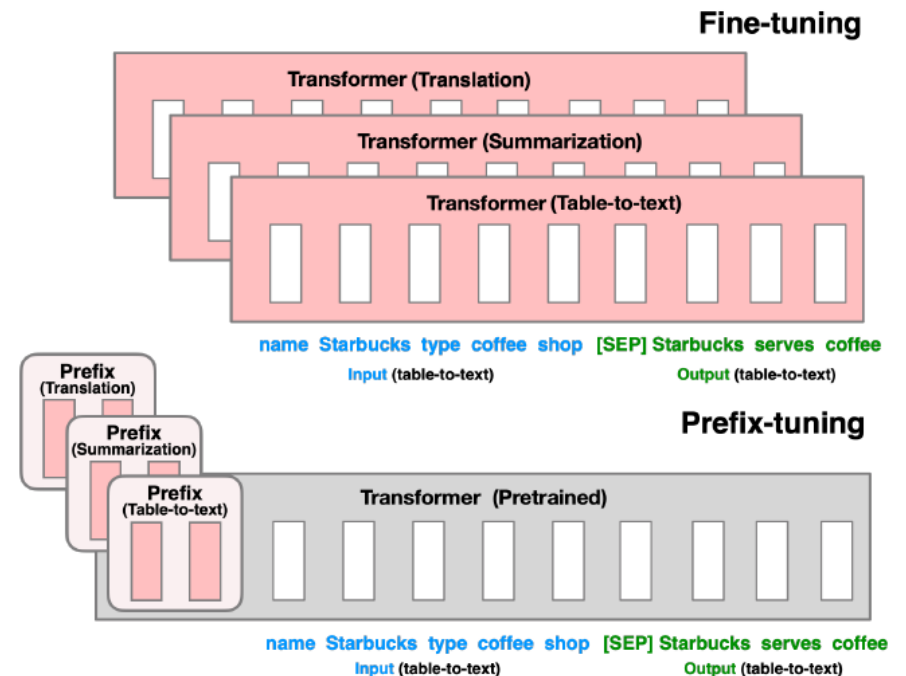
- Prefix/Prompt tuning
  - freeze model
    - attach learnable tokens to prompt
    - update these tokens to improve performance

*The Power of Scale for Parameter-Efficient Prompt Tuning*

Lester et al 2021

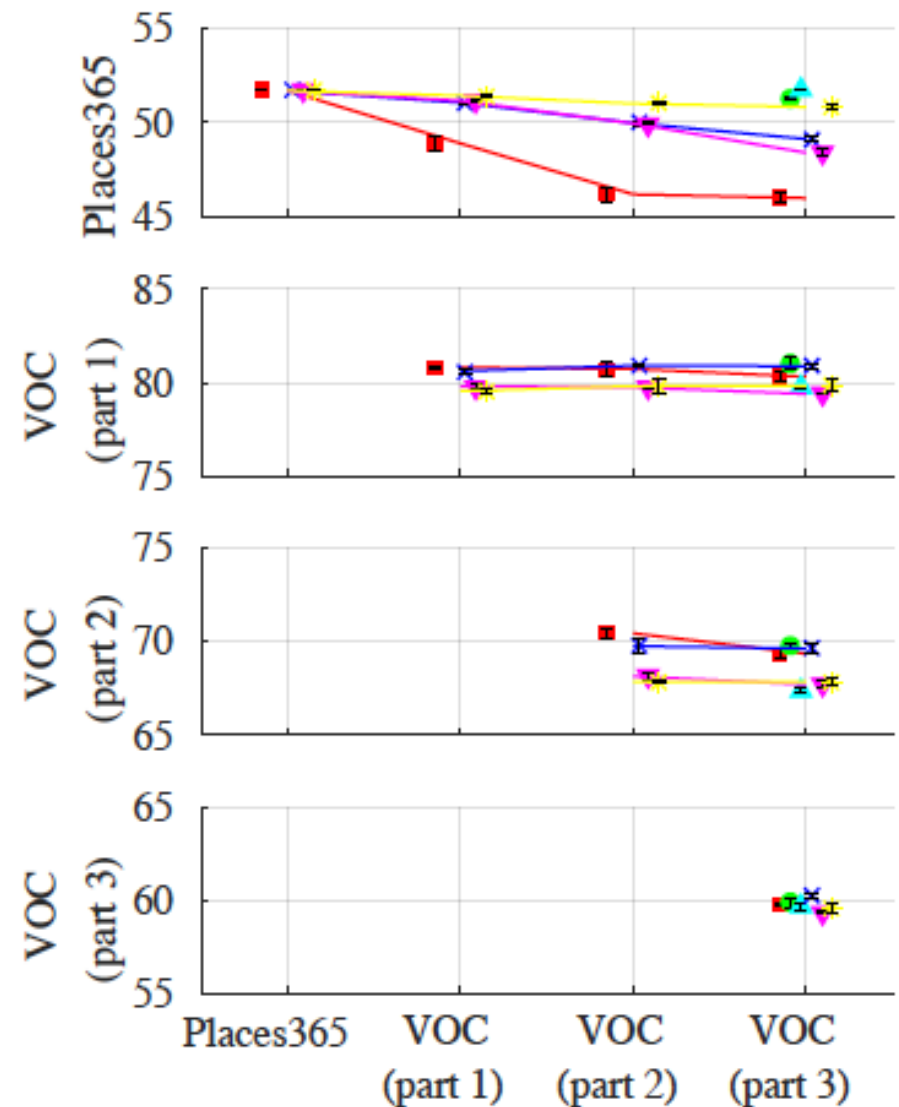
*Prefix-Tuning: Optimizing Continuous Prompts for Generation*

Li +Liang, 2021



# Nuisances: Forgetting

- Forgetting
  - Fine-tune a pretrained network on new task
  - It gets worse on previous tasks
  - Robust, widespread effect
- There are strategies



(a) Places365→VOC

# Nuisances: hallucination



**NBT:** A woman talking on a cell phone while sitting on a *bench*.  
CIDEr: **0.87**, METEOR: 0.23, SPICE: **0.22**, CHs: **1.00**, CHi: **0.33**

**TopDown:** A woman is talking on a cell phone.  
CIDEr: 0.54, METEOR: **0.26**, SPICE: 0.13, CHs: **0.00**, CHi: **0.00**



**Image Model predictions:**  
bowl, broccoli, carrot, dining table

**Language Model predictions for the last word:**  
fork, spoon, bowl

**Generated caption:** A plate of food with broccoli and a *fork*.

Figure 2: Example of image and language consistency. The hallucination error (“fork”) is more consistent with the Language Model.

# Nuisances: hallucination



**TD:** A cat is sitting on a bed in a room.

S: 12.1 M: 23.8 C: 69.7

**TD Restrict:** A bed with a blanket and a pillow on it.

S: 23.5 M: 25.4 C: 52.5



**TD:** A cat laying on the ground with a frisbee.

S: 8.0 M: 13.1 C: 37.0

**TD Restrict:** A black and white animal laying on the ground.

S: 7.7 M: 15.9 C: 17.4

Figure 5: Examples of how TopDown (TD) sentences change when we enforce that objects cannot be hallucinated: SPICE (S), Meteor (M), CIDEr (C), see Section 3.4.

- Cure strategy: force entities in captions to have corresponding detector responses

# Nuisances: localization issues

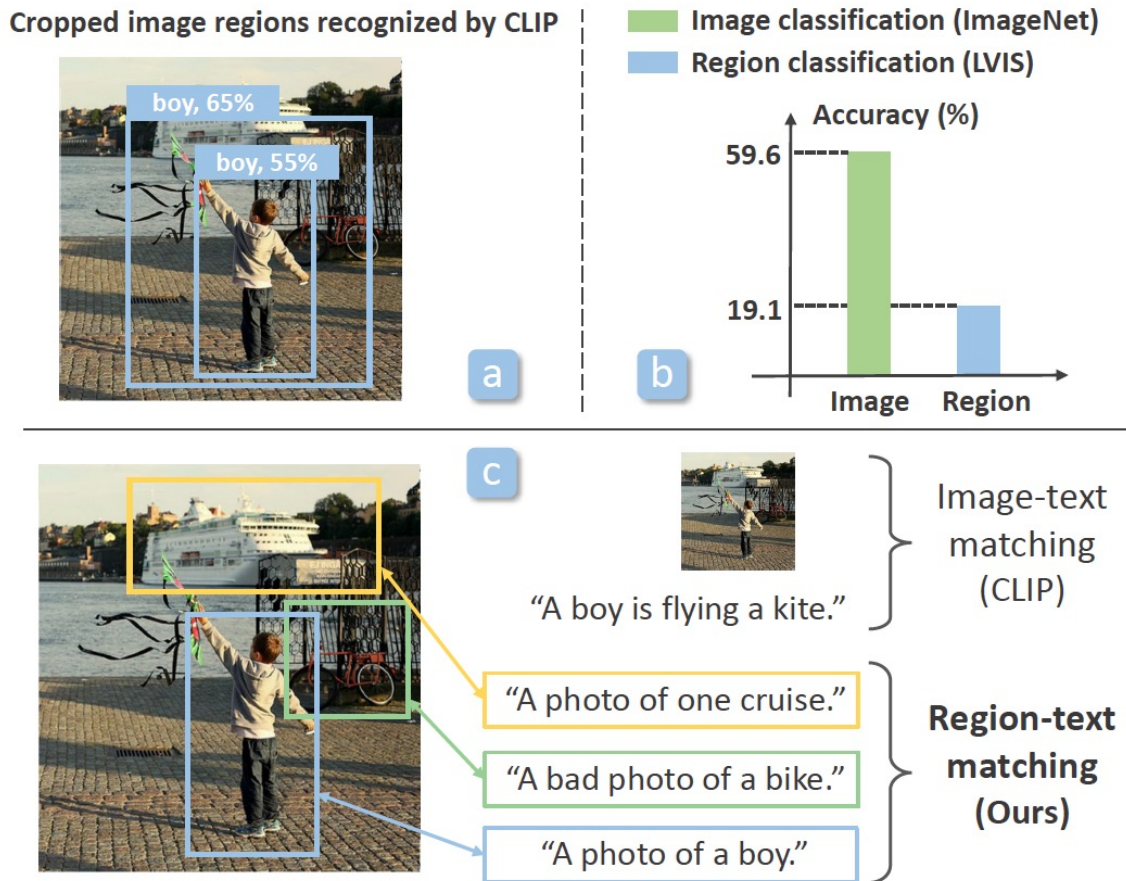


Figure 1. **(a)**. A pretrained CLIP model [37] failed to capture localization quality. **(b)**. A major drop on accuracy when using the same pretrained CLIP to classify image regions. **(c)**. Our key idea is learning to match *image regions* and their text descriptions.

# Nuisances: localization issues

- Cure:
  - use RPN to propose region boxes
  - match to word labels using visual features
  - Now have a dataset of localized regions – text
    - use that to finetune

# Improvements result

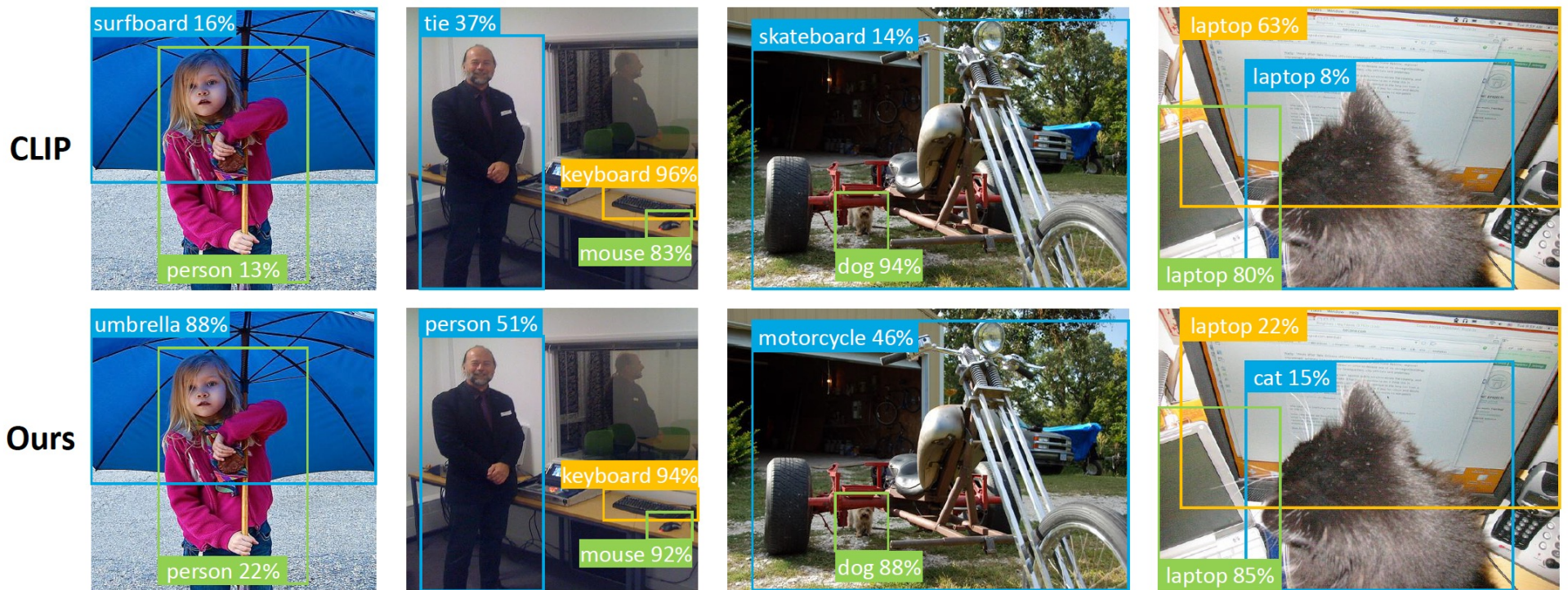


Figure 3. Visualization of zero-shot inference on COCO dataset with *ground-truth boxes*. Without finetuning, the pretrained models (top: CLIP, bottom: Ours) are directly used to recognize image regions into 65 categories in COCO. (Image IDs: 9448, 9483, 7386, 4795)

# Improvements result

## Success case:



**Ours:**  
teddy bear, 99.5%  
bear, 0.43%  
honey, 0.02%

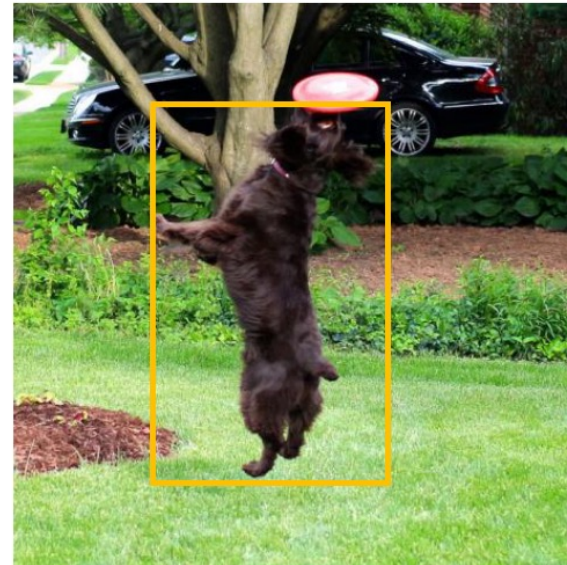
**CLIP:**  
fleece, 11.2%  
shawl, 1.9%  
turban, 1.8%



**Ours:**  
chocolate cake, 12.9%  
truffle chocolate, 12.8%  
chocolate mousse, 7.8%

**CLIP:**  
tape, 2.7%  
razorblade, 0.97%  
truffle chocolate, 0.84%

## Failure case:

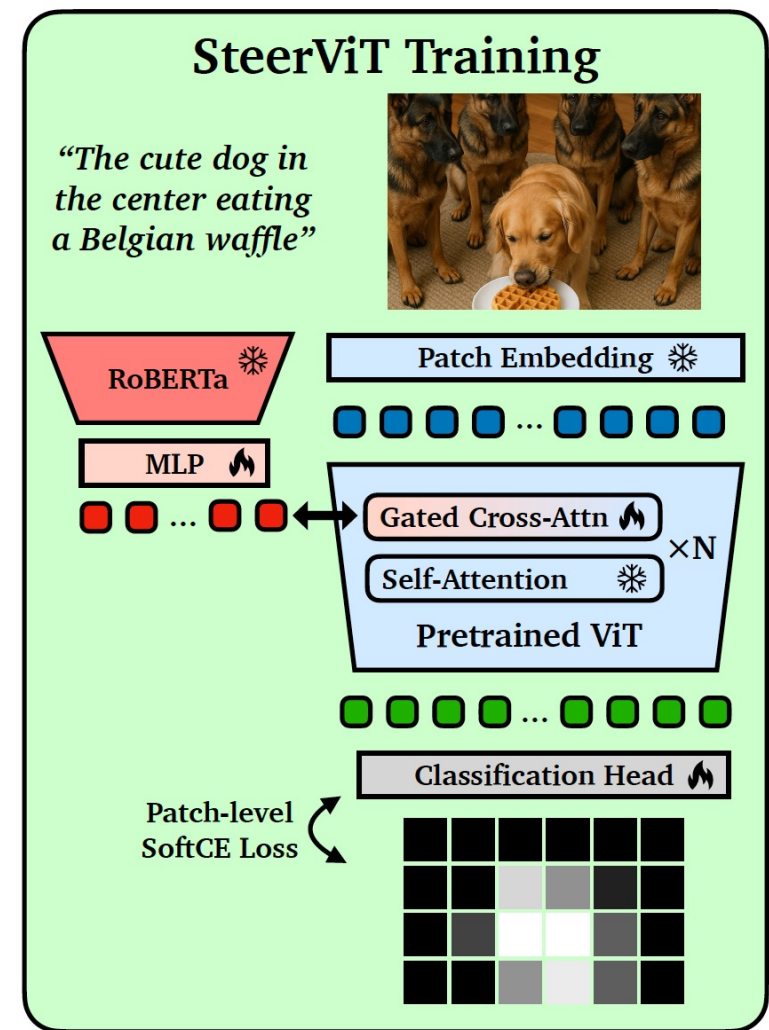


**Ours:**  
ferret, 8.8%  
cub, 8.1%  
shepherd dog, 5.4%

**CLIP:**  
grizzly, 9.3%  
cub, 8.8%  
gorilla, 8.1%

# Variants: Focus

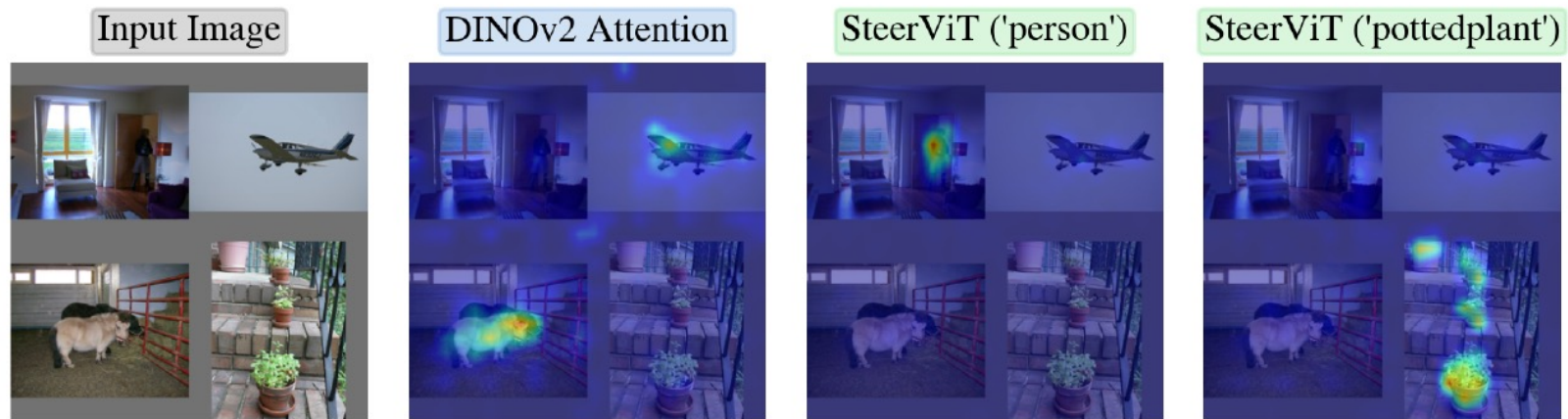
- Idea:
  - finetune vision model to accept cues from language to emphasize features
  - use cross-attention
  - train with



**Fig. 4: Steering any ViT using text conditioning.** Our method adds lightweight **vision**-to-**language** cross-attention layers within pretrained ViT blocks and applies a patch-level segmentation proxy objective to fuse prompt cues into patch tokens.

# Focus helps

- Att'n to CLS token shows steering is helpful



**Fig. 6: Text enables targeted attention.** Attention maps on a four-image mosaic demonstrate that text conditioning with **SteerViT** redirects self-attention to the queried concept whereas **DINOv2** attends to the most prominent objects. Note that the [CLS]-token of SteerViT was not directly optimized for targeted attention and remains frozen in its original state.

# Variants: Pointing

- Collect data where people point at named item
- Train model to predict these points
  - decode to points

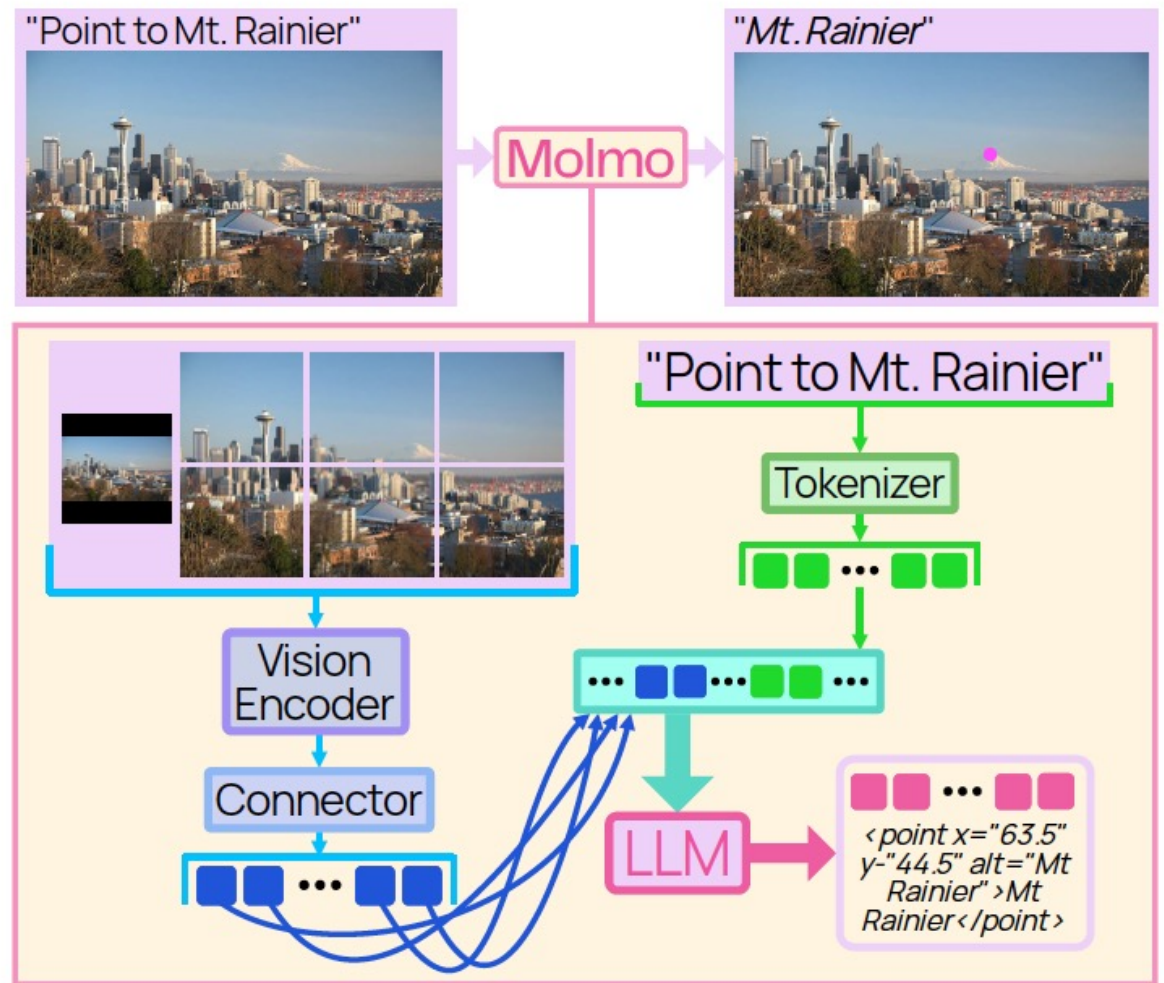


Figure 2. **Molmo** follows the simple and standard design of connecting a vision encoder and a language model.

# Pointing: advantages

- Model is discouraged from hallucination
  - as above
- Significant reduction in required training data
  - quality > quantity
- Knowing where something is is useful

# Resources

- everything I cite is on arxiv
  - search by name
- Open VLM with pretrained weights, all code, open data:
  - <https://openai.com/research/molmo>
- **Survey in:** *Vision-Language Models for Vision Tasks: A Survey, Zhang et al, 2024*
- **Detailed models in:**

*QWEN TECHNICAL REPORT, Bai et al, 2023*

*Flamingo: a Visual Language Model for Few-Shot Learning, Alayrac, 2022*

*Molmo and PixMo: OpenWeights and Open Data for State-of-the-Art Vision-Language Models, Deitke et al, 2024*

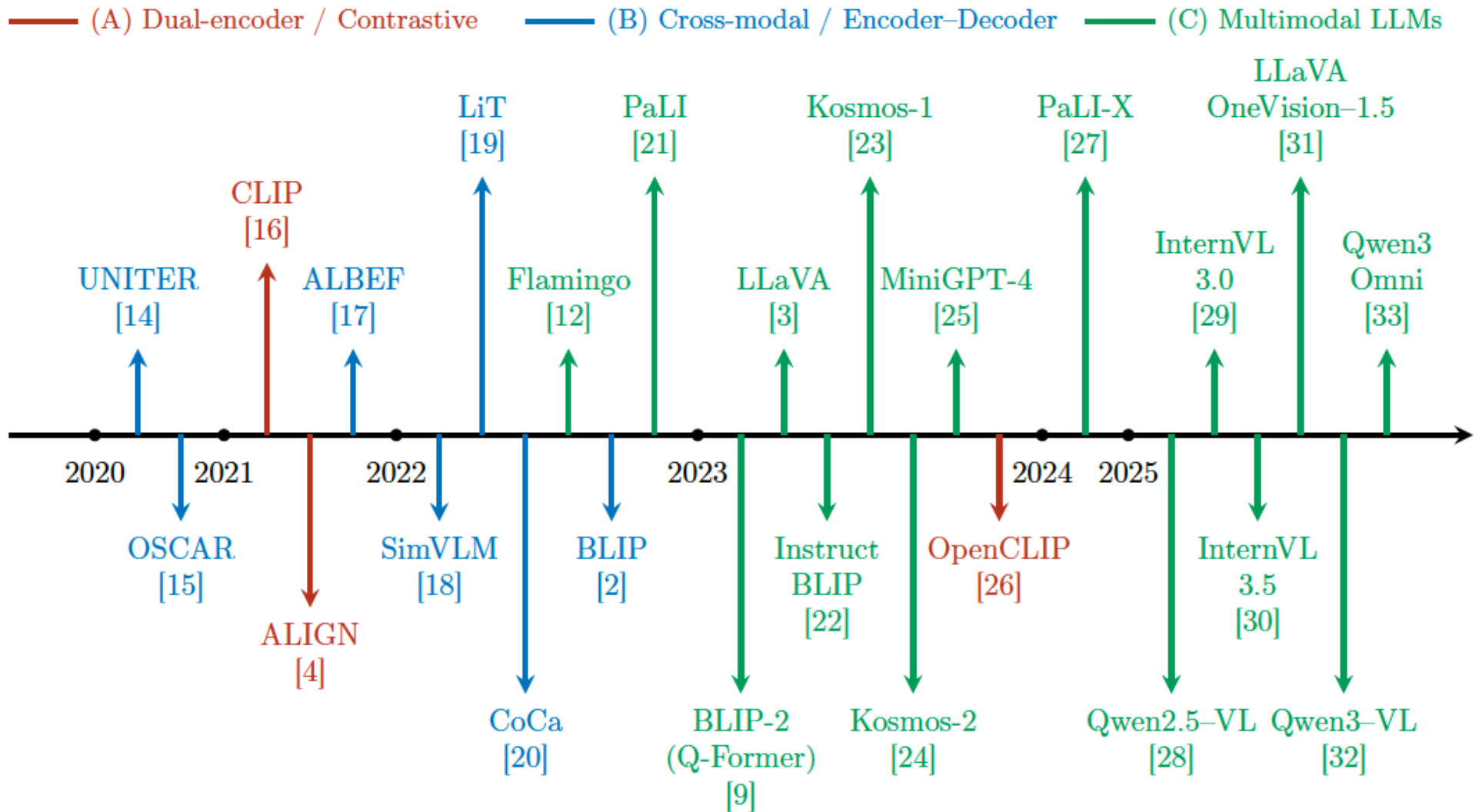


Figure 4: **Chronological overview of representative VLM / Multimodal LLM milestones (2020–2025).** Four color-coded categories: (A) Dual-encoder / Contrastive, (B) Cross-modal / Encoder-Decoder, (C) Multimodal LLMs. Each arrow is vertically offset within its year and labeled above/below with the model name and citation.

Table 2: **Named model families and design choices.** We summarize core ideas, fusion styles, typical backbones, objectives, and the nature/scale of pretraining data when explicitly reported in the cited papers. Abbreviations: ITC/ITM = image-text contrast / matching; LM = language modeling; GCA = gated cross-attention.

Model	Core Idea	Fusion	Vision & Text Backbones	Objective(s)	Pretraining Data
ALIGN [4]	Dual-encoder at web scale	None (late)	CNN/Transformer encoders	Contrastive alignment	~1B+ noisy alt-text pairs
CLIP [1]	Dual-encoder, universal embeddings	None (late)	ViT/ResNet + text Transformer	Contrastive (InfoNCE)	~400M web pairs
BLIP [2]	MED unifies understanding & generation + CapFilt	Cross-attn in encoder/decoder	ViT + BERT-style text Transformer	ITC + ITM + LM	14M–129M images (+ partial LAION)
Flamingo [12]	LVLm with GCA layers and Perceiver Resampler	Gated cross-attention	CNN/ViT features → LLM decoder	Autoregressive LM (vision-conditioned)	Interleaved image-text corpora (e.g., M3W)
LLaVA [3]	Visual instruction tuning with GPT-4 conversations	Projector → LLM tokens	CLIP-like image encoder + decoder LLM	SFT / LM on multimodal dialogs	High-quality synthetic conversations
DINO [34]	Self-distillation ViT without labels	N/A	Vision Transformer (student/teacher)	Self-distill with momentum teacher	Unlabeled images (multi-crop)
Grounding DINO [13]	Open-vocabulary detection by region-text pretraining	Cross-modal at region level	Detector + text encoder (phrase grounding)	Detection + grounding losses	Region-phrase corpora (web & annotations)
MoE [8]	Token routing to expert MLPs (sparse)	Inside decoder (sparse)	Decoder-only Transformer with experts	Autoregressive LM, sparse gating	Large text corpora (scalable)

Table 5: Training paradigms appearing in VLM papers.

Item	2023	2024	2025	Trend	Slope (pp/yr)
Pretrain + Finetune	11.6%	16.9%	16.8%	+5.2%	3.83
Prompt/Prefix	13.0%	16.4%	14.3%	+1.3%	3.43
Self-/Weak-/Semi- sup.	9.6%	2.8%	3.5%	-6.1%	-2.65
Distillation	4.2%	4.8%	4.0%	-0.4%	0.56
Instruction Tuning	1.1%	4.2%	5.0%	+3.9%	1.75
LoRA/Adapters	1.3%	4.0%	4.1%	+2.8%	1.26
Multi-task/Curriculum	2.6%	1.6%	1.9%	-0.8%	-0.08

Table 7: Curated datasets explicitly named in VLM abstracts (note under-reporting).

Item	2023	2024	2025	Trend	Slope (pp/yr)
MS-COCO	4.9%	2.1%	1.0%	-3.0%	-1.50
ImageNet	3.1%	2.4%	1.6%	-1.5%	-0.76
LAION	0.6%	0.8%	0.2%	-0.5%	0.03
RefCOCO/g/+	0.6%	0.6%	0.3%	-0.3%	-0.12
Flickr30k	0.8%	0.3%	0.2%	-0.7%	-0.28
CC3M/CC12M	0.4%	0.3%	0.3%	-0.1%	-0.19
VQA-v2/OK-VQA	0.4%	0.2%	0.3%	-0.2%	-0.09
WebVid/MSRVTT/MSVD	0.3%	0.3%	0.1%	-0.2%	-0.04
YouCook2/HowTo	0.2%	0.2%	0.1%	-0.1%	-0.16
Visual Genome	0.3%	0.1%	0.0%	-0.2%	-0.05
COCO Captions	0.1%	0.0%	0.0%	-0.0%	-0.03