

Does Language Model Need Better Visual Grounding for Meaning and Understanding?

Saining Xie
Courant Institute, NYU
11/2024

Does AI need sensory grounding for meaning and understanding?

- Stevan Harnad (1990): the symbol grounding problem
- *“Symbols in AI systems must have sensory grounding (or: bodily grounding, or external grounding) to have meaning.”*



Stevan Harnad: The Symbol Grounding Problem



Do humans require sensory grounding for meaning and understanding?

- Thomas Aquinas(1200s)
- *“There’s nothing in mind that wasn’t first in the senses”*



Do humans require sensory grounding for meaning and understanding?

- Diderot, Condilac (1700s): sensim
- *“no cognition without sensasion”*



Do humans require sensory grounding for meaning and understanding?

- Avicenna (Ibn Sina) (1000s): Avicenna's Floating Man
- Avicenna: *"floating man can think about himself without ever sensing (and without sensory capabilities?)"*
- others: *"he could also think about mathematics, logic, philosophy, and could form hypotheses about external reality."*



A less philosophical question

Is sensing useful?

- Does sensing boost thinking (to a new level) in language models?
(and in intelligent creatures generally?)

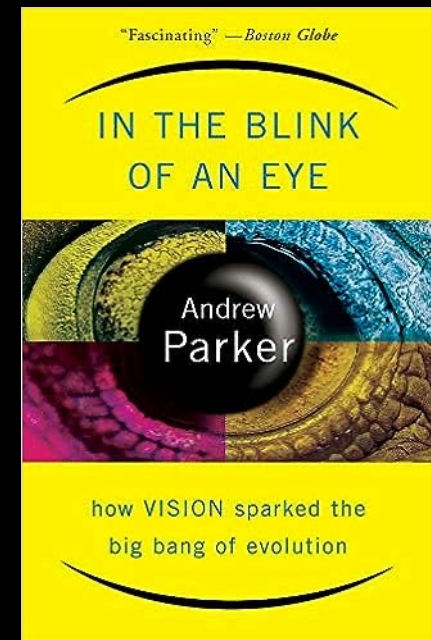
phylogeny of intelligence



538.8 million years ago
Cambrian era
“biological explosion”

“The evolution of the eye is likely to have been a catalyst for the explosion, initiating an arms race between organisms that were increasingly aware of their surroundings.”

<https://www.nhm.ac.uk/discover/eyes-on-the-prize-evolution-of-vision.html>



LLMs: knowledgeable but blindfolded



Why do we need better visual representations?

“Most of human knowledge (and almost all of animal knowledge) comes from our sensory experience of the physical world.”

- “Language is the icing on the cake. We need the cake to support the icing.”

LeCun's cake 2.0 ?



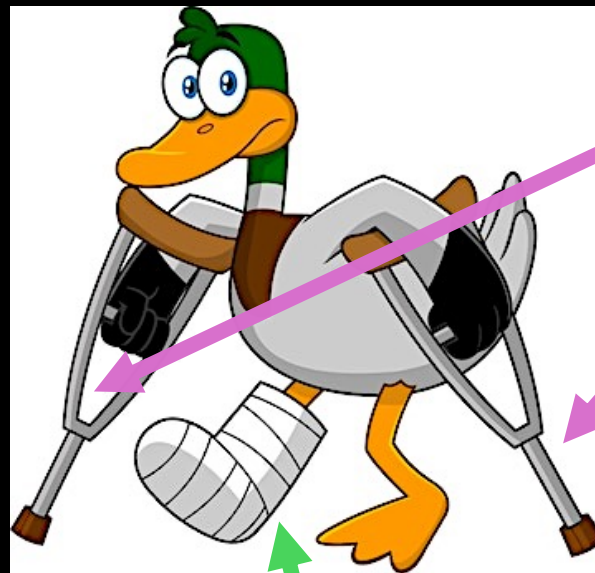
Sensory experience modeling

Language modeling



<https://lexfridman.com/yann-lecun-3-transcript/>

Relying too heavily too early on language can act as a shortcut, compensating for the deficiencies in learning effective visual representations.



Your favorite
LLMs

Our visual representations

Language vs Visual Intelligence

“Who won the game?”



[GPT-4O, OpenAI]

“Which direction leads home?”



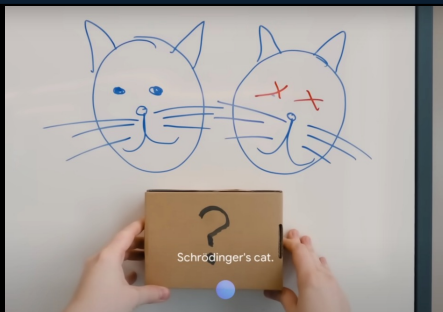
[V-IRL - ECCV 2024]

“Where can I buy this mug?”



[V* - CVPR 2024]

“what does this remind you of?”



[Project Astra, Google]

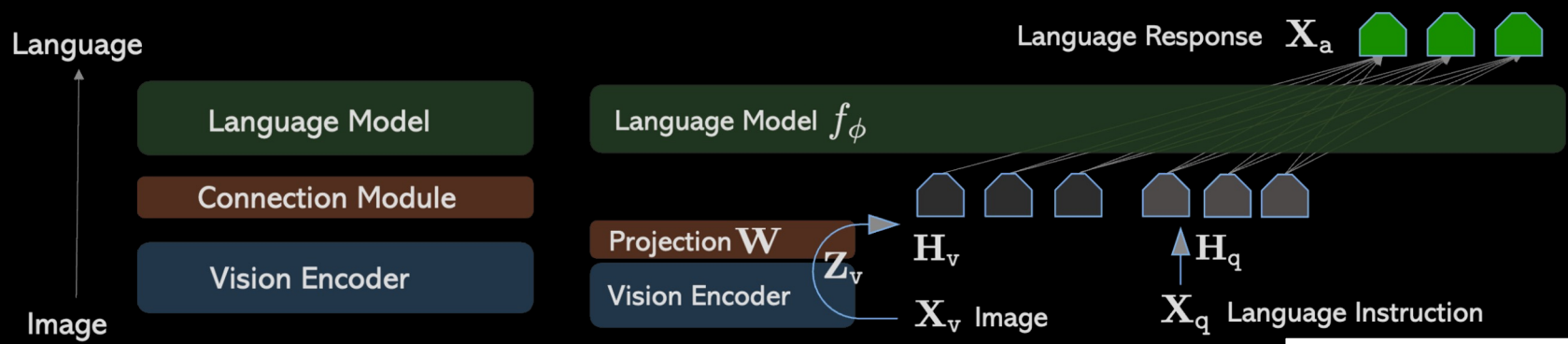
“Thinking in Space”



Tasks Requiring more
Robust Visual-Spatial Intelligence

Tasks Requiring more
Strong Language Capability

A typical MLLM pipeline (LLaVA)

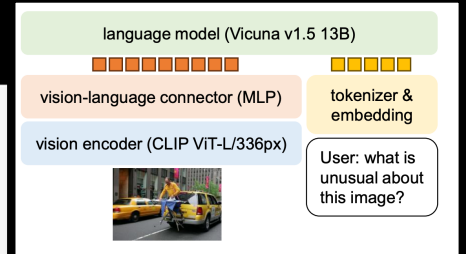


•Stage 1: Pre-training for Feature Alignment.

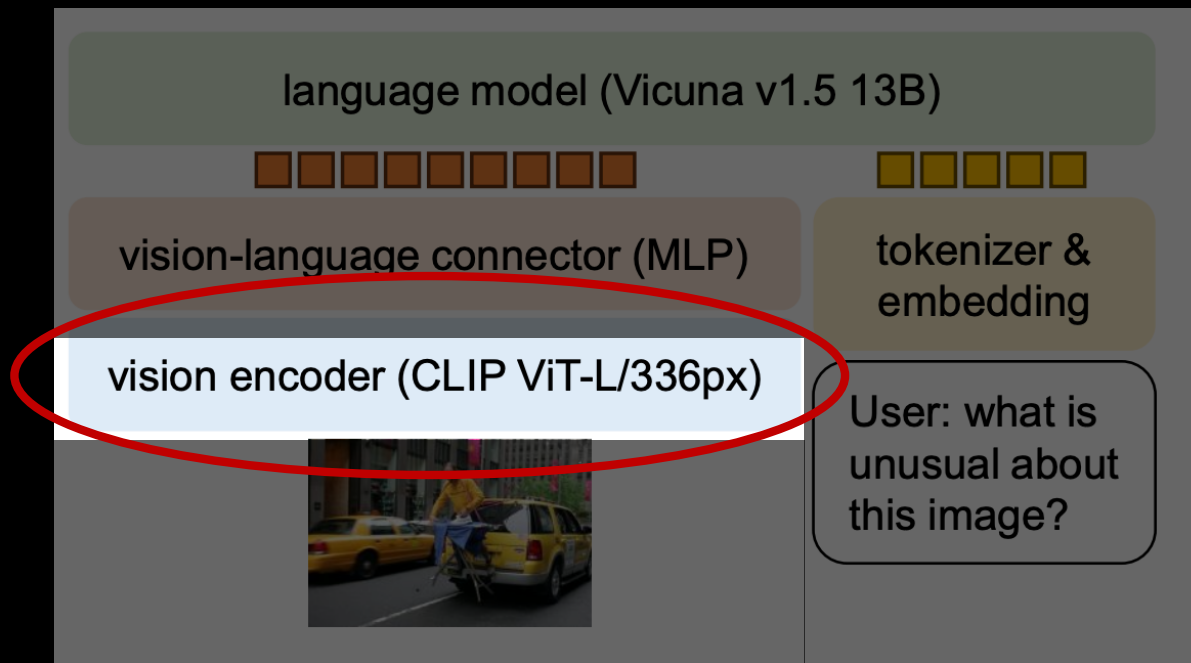
Only the projection matrix is updated, based on a subset of CC3M.

•Stage 2: Fine-tuning End-to-End. Both the projection matrix and LLM are updated

- Visual Chat:** Our generated multimodal instruction data for daily user-oriented applications.
- Science QA:** Multimodal reasoning dataset for the science domain.

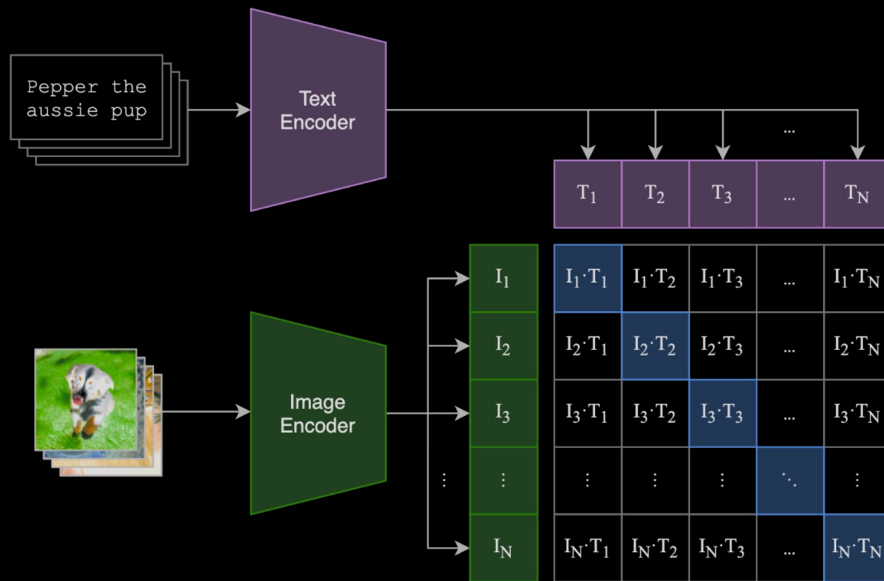


A typical MLLM pipeline (LLaVA)

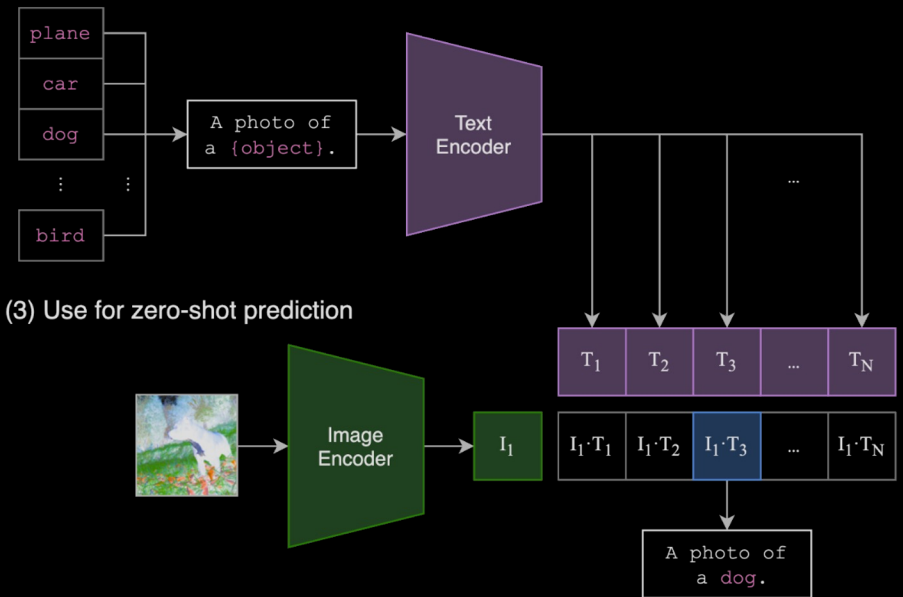


CLIP vision encoder

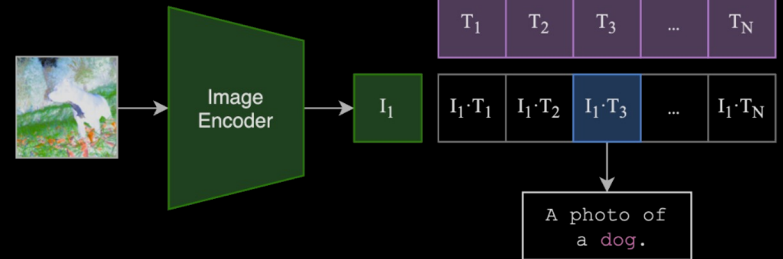
(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction



Key question #1:

Does language model need better visual grounding/representations for meaning and understanding?

Eyes Wide Shut? Exploring the Visual Shortcomings of Multimodal LLMs

Peter Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, Saining Xie
CVPR 2024 (Oral)



Q: Is the dog facing left or right from the camera's perspective?

The dog is **facing to the right** from the camera's perspective. You can tell by the position of its tail and the orientation of its body.



Q: Do you see any window in this image?

No, there are **no windows visible** in this image. The image shows a traditional red pillar box, which is a type of mailbox found in the United Kingdom.



Q: Do the individuals in the picture face the front or the back?

The individuals in the picture are **facing away from the camera**, so we see them from the back.



Q: In this image, how many eyes can you see on the animal?

The image shows **one eye** of the animal. It's a close-up of a bald eagle facing slightly to its right, which presents only one side of its face to the viewer.



Q: How many wheels can you see in the image?

In the image provided, I can see **two wheels** on the visible side of the car.



Q: Does the keyboard have a backlight?

Based on the image, it appears that the keyboard **does have a backlight**. This is indicated by the illumination seen around the keys, which suggests that there is a light source beneath them to enhance visibility in low-light conditions.



Q: Is the person in the picture on the grass or on the gravel path?

The person in the image appears to be running **on a gravel path**. You can discern this from the texture and color contrast between the light-colored path and the green grass surrounding it.



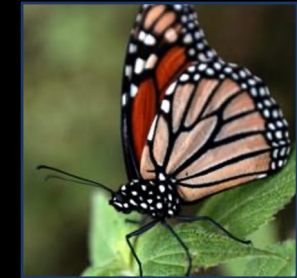
Q: Is the edge of the heart in the picture white or dark-colored?

The image you've provided is small and of low resolution, but it appears that the hearts have a **dark-colored edge** or outline. This outline varies in visibility, likely due to the quality of the image.



Q: Is the school bus parked facing the camera or away from the camera?

The school bus is parked **facing away from the camera**. You can tell because the rear lights and the back of the exterior rearview mirror are visible. The text "SCHOOL BUS" is also reversed, which is typically seen on the front of a bus, designed so that it reads correctly in a rearview mirror.



Q: Are the butterfly's feet visible?

The butterfly's feet are **not visible** in this image. The wings and part of the body of the butterfly are in view, but the feet are either hidden beneath the body or out of the frame of the picture. Butterflies have delicate, slender feet which are often not the most noticeable feature when they are perched, especially from the angle shown in this image.

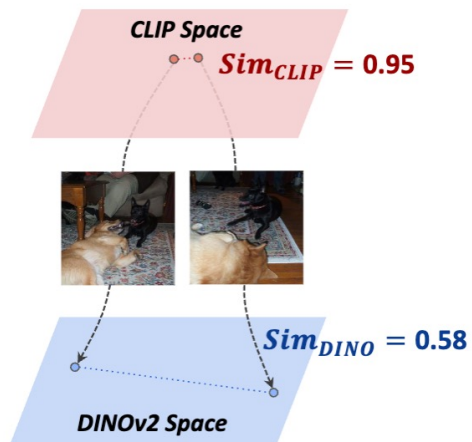
We can systematically identify instances where the visual question answering (VQA) capabilities of GPT-4V fall short.

Constructing MMVP Benchmark via CLIP-blind Pairs.

Step 1

Finding CLIP-blind $\not\sim$ pairs.

Discover image pairs that are proximate in CLIP feature space but distant in DINOv2 feature space.

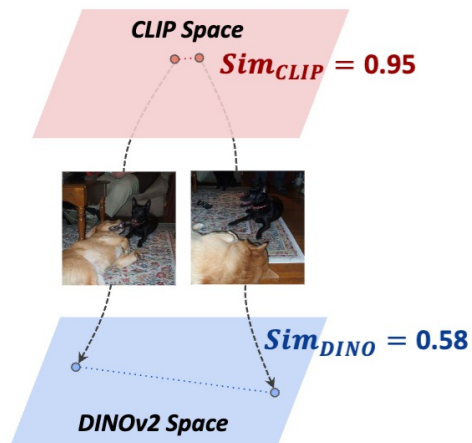


Constructing MMVP Benchmark via CLIP-blind Pairs

Step 1

Finding CLIP-blind pairs.

Discover image pairs that are proximate in CLIP feature space but distant in DINOv2 feature space.



Step 2

Spotting the difference between two images.

For a CLIP-blind pair, a human annotator attempts to spot the visual differences and formulates questions.



"The dog's head in the left image is resting on the carpet, while the dog's head in the right image is lying on the floor."

Formulating questions and options for both images.

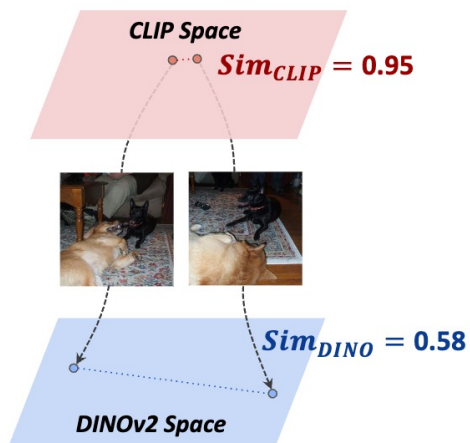
Where is the yellow animal's head lying in this image?
(a) Floor (b) Carpet

Constructing MMVP Benchmark via CLIP-blind Pairs

Step 1

Finding CLIP-blind pairs.

Discover image pairs that are proximate in CLIP feature space but distant in DINOv2 feature space.



Step 2

Spotting the difference between two images.

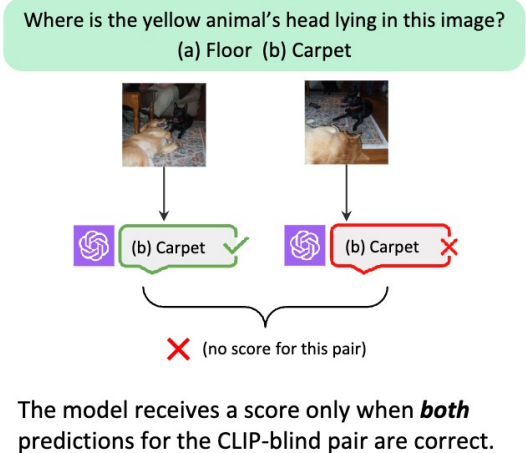
For a CLIP-blind pair, a human annotator attempts to spot the visual differences and formulates questions.



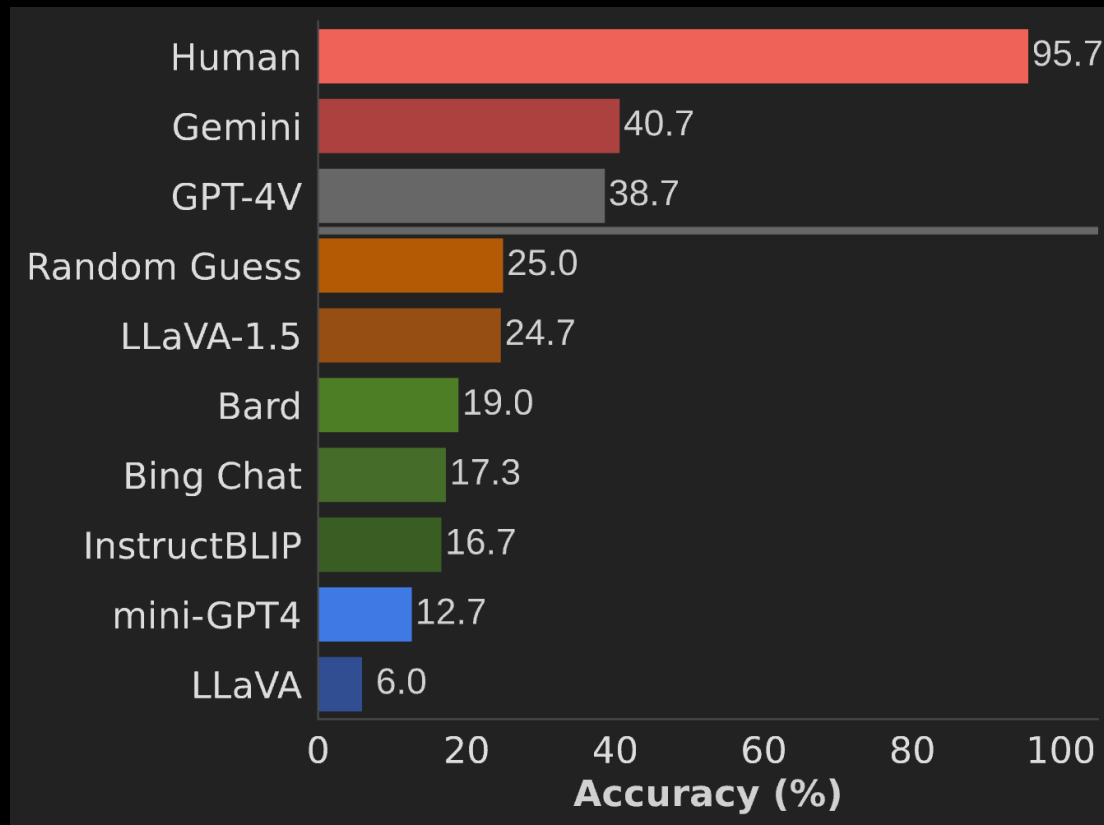
Step 3

Benchmarking multimodal LLMs.

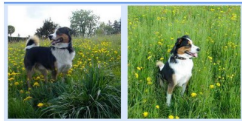
Evaluate multimodal LLMs using a CLIP-blind image pair and its associated question.



Results Of Current MLLM Models (And Humans)



Is the dog facing left or right from the camera's perspective?



(a) Left (b) Right

	(b)	(b)	✗
	(a)	(a)	✗
	(b)	(b)	✗
	(a)	(a)	✗

Is the needle pointing up or down?



(a) Up (b) Down

	(b)	(b)	✗
	(a)	(b)	✓
	(a)	(a)	✗
	(a)	(a)	✗

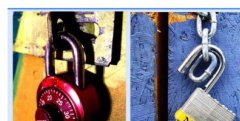
Is the cup placed on a surface or being held by hand?



(a) Placed on a surface (b) Held by hand

	(a)	(a)	✗
	(a)	(b)	✓
	(a)	(a)	✗
	(a)	(b)	✓

Is the lock locked or unlocked?



(a) Locked (b) Unlocked

	(a)	(b)	✓
	(a)	(b)	✓
	(a)	(a)	✗
	(a)	(a)	✗

Is the snail in the picture facing the camera or away from the camera?



(a) Away from the camera (b) Facing the Camera

	(b)	(b)	✗
	(a)	(a)	✗
	(b)	(b)	✗
	(a)	(a)	✗

Are the ears of the dog erect or drooping?



(a) Erect (b) Drooping

	(b)	(b)	✗
	(a)	(a)	✗
	(b)	(b)	✗
	(a)	(a)	✗

In this image, how many eyes can you see on the animal?



(a) 1 (b) 2

	(a)	(a)	✗
	(b)	(b)	✗
	(b)	(b)	✗
	(b)	(b)	✗

Is this a hammerhead shark?



(a) Yes (b) No

	(b)	(b)	✗
	(a)	(b)	✓
	(b)	(b)	✗
	(a)	(a)	✗

Are there cookies stacked on top of other cookies?



(a) Yes (b) No

	(b)	(b)	✗
	(a)	(b)	✓
	(a)	(a)	✗
	(b)	(a)	✗

Is there a hand using the mouse in this image?



(a) Yes (b) No

	(b)	(b)	✗
	(a)	(b)	✓
	(b)	(b)	✗
	(a)	(b)	✓

Are there any clouds?



(a) Yes (b) No

	(b)	(b)	✗
	(a)	(b)	✓
	(a)	(b)	✓
	(a)	(b)	✓

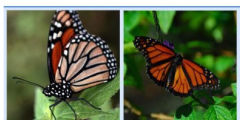
Do you see any window in this image?



(a) Yes (b) No

	(b)	(b)	✗
	(a)	(b)	✓
	(b)	(b)	✗
	(b)	(a)	✗

Are the butterfly's feet visible?



(a) Yes (b) No

	(b)	(a)	✗
	(a)	(b)	✓
	(a)	(b)	✓
	(a)	(a)	✗

Is the following statement correct: There are different colors of grapes in this image



(a) Correct (b) Incorrect

	(a)	(a)	✗
	(a)	(b)	✓
	(a)	(a)	✗
	(a)	(a)	✗

Is the following statement correct: There is no letter D on this image?



(a) Correct (b) Incorrect




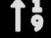





	(b)	(b)	✗
	(b)	(b)	✗
	(b)	(b)	✗
	(a)	(a)	✗

So, what is going on?

User

I am analyzing an image embedding model. Can you go through the questions and options, trying to figure out some general patterns that the embedding model struggles with? Please focus on the visual features and generalize patterns that are important to vision models
[MMVP Questions and Options]

We identify 9 visual patterns:

-  Orientation and Direction
-  Presence of Specific Features
-  State and Condition
-  Quantity and Count
-  Positional and Relational Context
-  Color and Appearance
-  Structural and Physical Characteristics
-  Text
-  Viewpoint and Perspective

Visual patterns in CLIP-blind pairs

Systematic Failures in CLIP

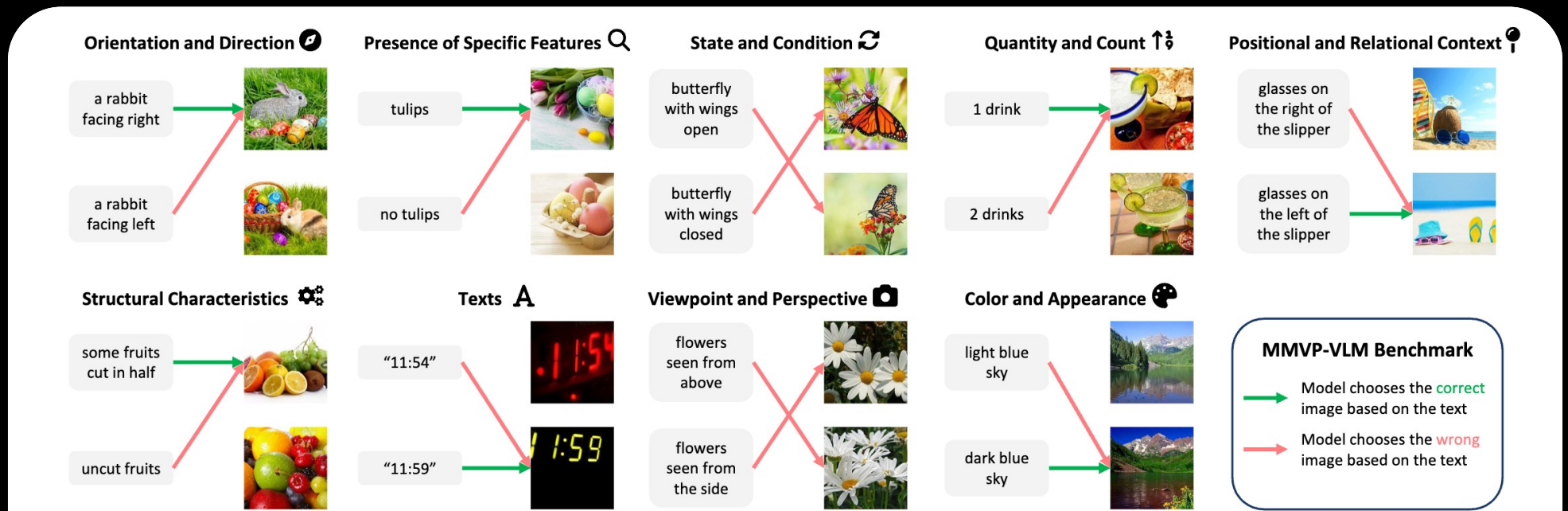









Figure 5. Examples from MMVP-VLM. MMVP-VLM consists of image pairs across nine visual patterns. The examples in the figure are from EVA01 ViT-g-14 model [54], one of the largest CLIP models that also fails to choose the right image given the text description.

from EVA01 ViT-g-14 model [24], one of the largest CLIP models that also fails to choose the right image given the text description.

Systematic Failures in CLIP-like models

	Image Size	Params (M)	IN-1k ZeroShot								A		MMVP Average
OpenAI ViT-L-14 [43]	224 ²	427.6	75.5	13.3	13.3	20.0	20.0	13.3	53.3	20.0	6.7	13.3	19.3
OpenAI ViT-L-14 [43]	336 ²	427.9	76.6	0.0	20.0	40.0	20.0	6.7	20.0	33.3	6.7	33.3	20.0
SigLIP ViT-SO-14 [66]	224 ²	877.4	82.0	26.7	20.0	53.3	40.0	20.0	66.7	40.0	20.0	53.3	37.8
SigLIP ViT-SO-14 [66]	384 ²	878.0	83.1	20.0	26.7	60.0	33.3	13.3	66.7	33.3	26.7	53.3	37.0
DFN ViT-H-14 [10]	224 ²	986.1	83.4	20.0	26.7	73.3	26.7	26.7	66.7	46.7	13.3	53.3	39.3
DFN ViT-H-14 [10]	378 ²	986.7	84.4	13.3	20.0	53.3	33.3	26.7	66.7	40.0	20.0	40.0	34.8
MetaCLIP ViT-L-14 [62]	224 ²	427.6	79.2	13.3	6.7	66.7	6.7	33.3	46.7	20.0	6.7	13.3	23.7
MetaCLIP ViT-H-14 [62]	224 ²	986.1	80.6	6.7	13.3	60.0	13.3	6.7	53.3	26.7	13.3	33.3	25.2
EVA01 ViT-g-14 [54]	224 ²	1136.4	78.5	6.7	26.7	40.0	6.7	13.3	66.7	13.3	13.3	20.0	23.0
EVA02 ViT-bigE-14+ [54]	224 ²	5044.9	82.0	13.3	20.0	66.7	26.7	26.7	66.7	26.7	20.0	33.3	33.3

CLIP Failures and MLLM Failures are Correlated

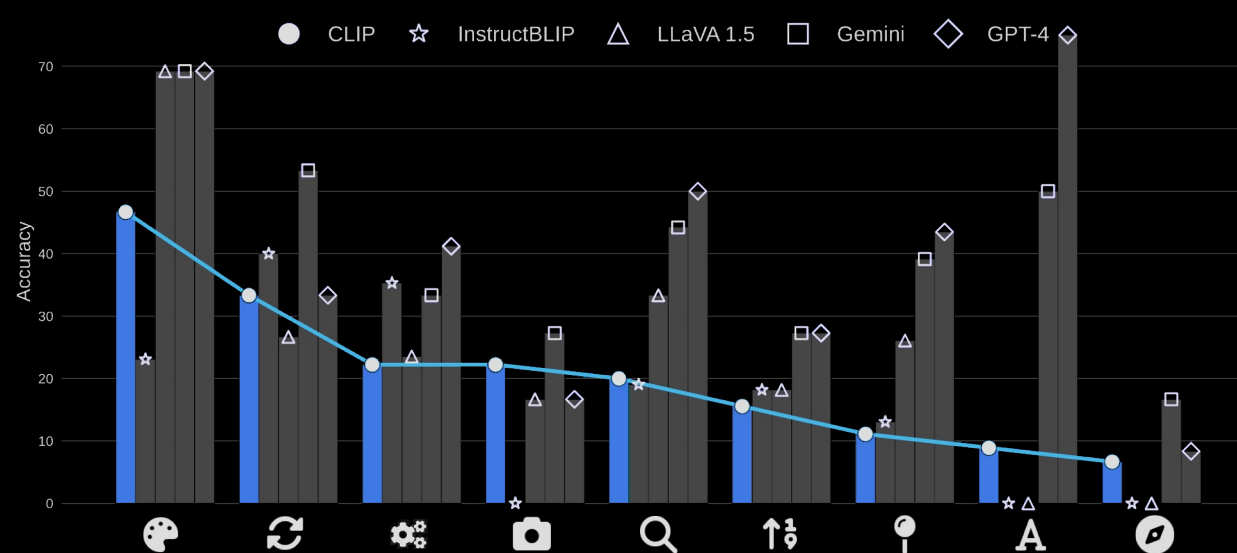

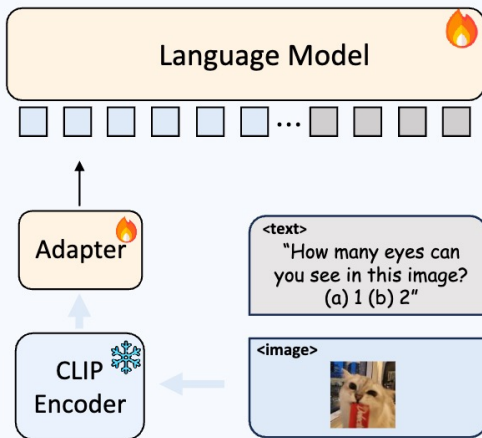


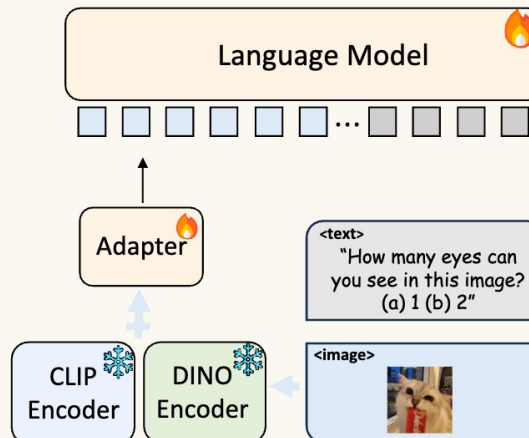
Figure 6. **CLIP and MLLM’s performance on visual patterns.** If CLIP performs poorly on a visual pattern such as “ orientation”, MLLMs also underperform on the visual pattern.

Mixture-of-Features (MoF) for MLLM

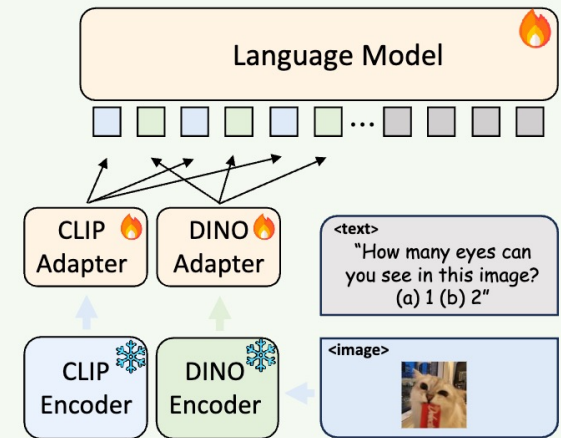
Standard MLLM



Additive-MoF MLLM



Interleaved-MoF MLLM



Additive MoF

method	SSL ratio	MMVP	LLaVA
LLaVA	0.0	5.5	81.8
	0.25	7.9 (+2.4)	79.4 (-2.4)
	0.5	12.0 (+6.5)	78.6 (-3.2)
LLaVA	0.625	15.0 (+9.5)	76.4 (-5.4)
+ A-MoF	0.75	18.7 (+13.2)	75.8 (-6.0)
	0.875	16.5 (+11.0)	69.3 (-12.5)
	1.0	13.4 (+7.9)	68.5 (-13.3)

Table 2. **Empirical Results of Additive MoF.** We use DINOv2 as the image SSL model in our work. With more DINOv2 features added, there is an improvement in visual grounding, while a decline in instruction following ability.

Interleaved MoF

method	res	#tokens	MMVP	LLaVA	POPE
LLaVA	224 ²	256	5.5	81.8	50.0
LLaVA	336 ²	576	6.0	81.4	50.1
LLaVA + I-MoF	224 ²	512	16.7 (+10.7)	82.8	51.0
LLaVA ^{1.5}	336 ²	576	24.7	84.7	85.9
LLaVA ^{1.5} + I-MoF	224 ²	512	28.0 (+3.3)	82.7	86.3

Table 3. Empirical Results of Interleaved MoF. Interleaved MoF improves visual grounding while maintaining same level of instruction following ability.

Other SSL backbones can work too

method	SSL Model	res	#tokens	MMVP	POPE
LLaVA ^{1.5}	None	336 ²	576	24.7	85.9
LLaVA ^{1.5} + I-MoF	MoCov3	224 ²	512	26.7 (+2.0)	86.1
LLaVA ^{1.5} + I-MoF	MAE	224 ²	512	27.3 (+2.6)	86.1
LLaVA ^{1.5} + I-MoF	DINOv2	224 ²	512	28.0 (+3.3)	86.3

Table 6. Results of Interleaved MoF with different vision-only SSL model

Takeaways



- Visual grounding is important for language understanding and meaning.
- CLIP's been lounging around for too long! (opportunities!)
- Vision SSL FTW!
(but we need fundamentally different ways to pursue the problem.)

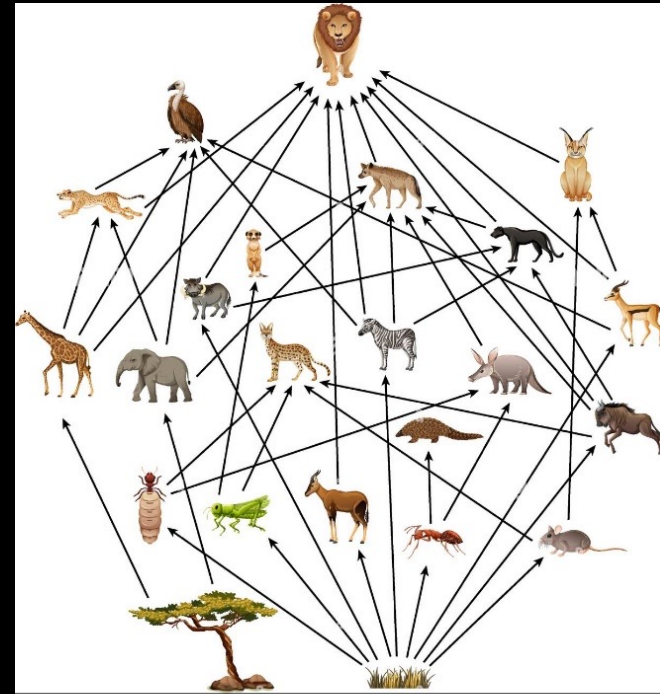
Key Question #2:

Better visual representations: beyond just static, global image feature extractors?

V^* : Guided Visual Search as a Core Mechanism in Multimodal LLMs

Penghao Wu, Saining Xie
CVPR 2024

“Deliberate” visual processing



A Concrete Example



plastic
straw



plastic
straw



plastic
straw



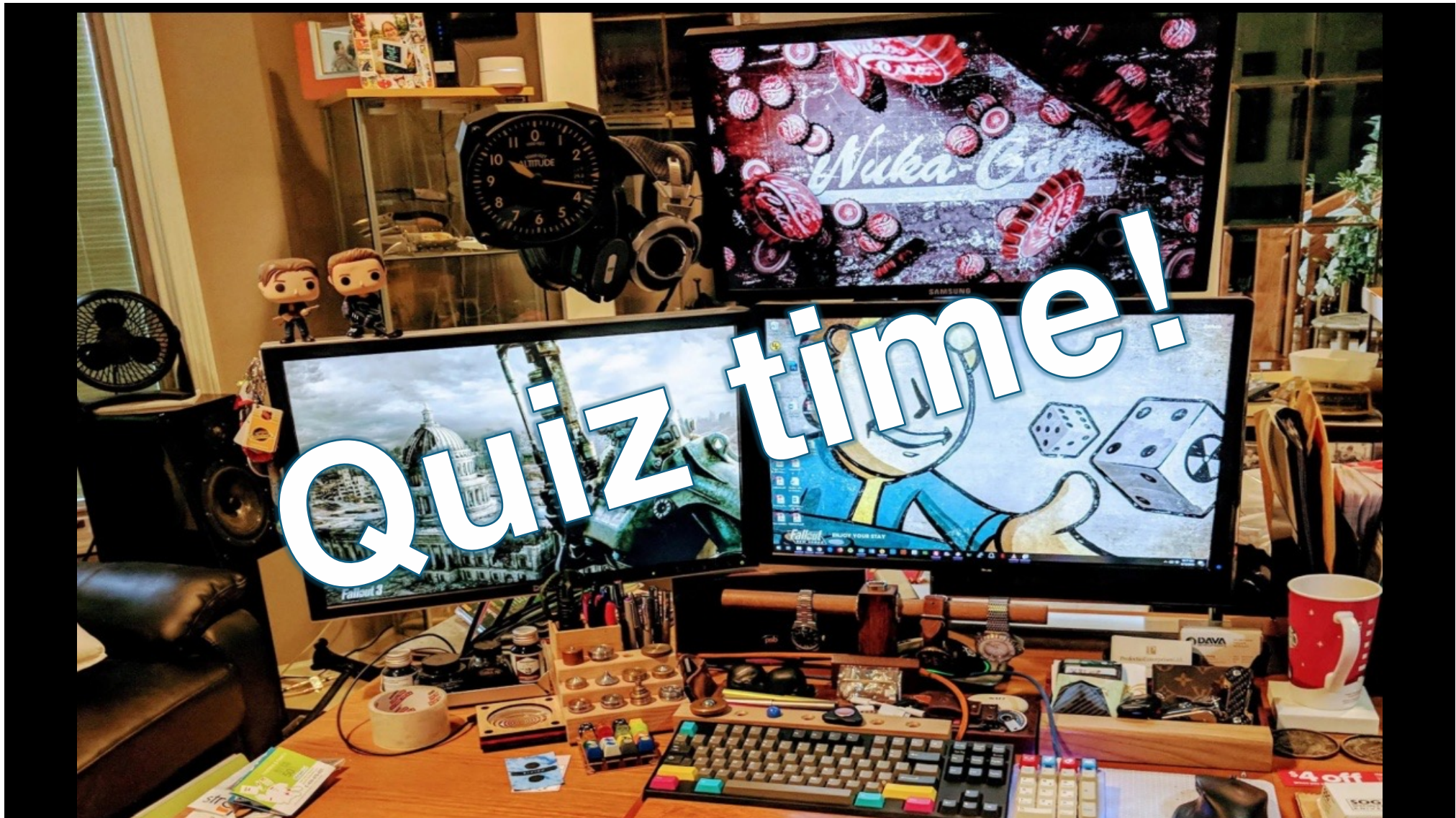
plastic
straw



plastic
straw



plastic
straw





 **GPT-4V:**

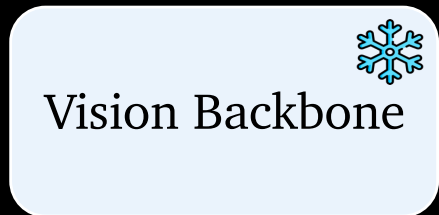
The mug in the image has a logo reminiscent of "**Nuka-Cola**" from the "Fallout" video game series. You can typically find merchandise related to "Fallout", including Nuka-Cola mugs, at online retailers such as **Amazon, eBay, or specialized gaming merchandise stores.**

 **SEAL:**

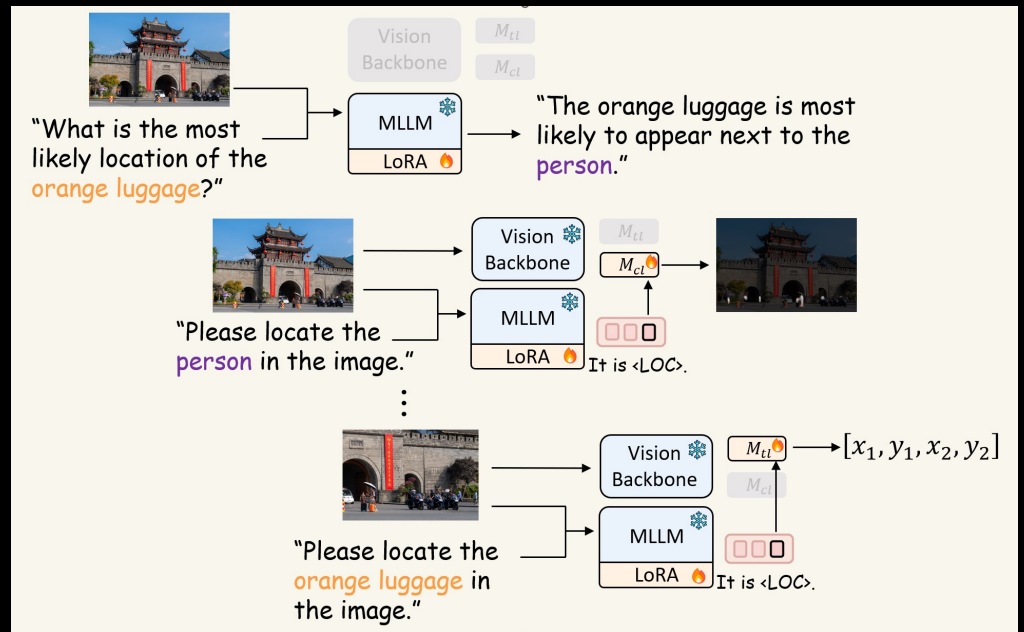
You can buy a mug like this based on its logo at a **Starbucks store** or online at the **Starbucks website.**



Visual Search Model for MLLM



Is the orange luggage on the left side of the black luggage?



Visual Search Examples



Search for the guitar

Visual Search Examples



Contextual cue

The guitar is most likely to appear on the **stage**.

Visual Search Examples



Target-specific cue

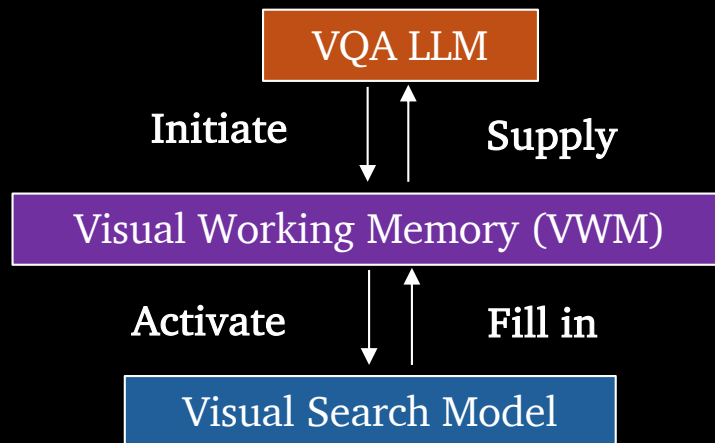
Visual Search Examples



Target Spotted!

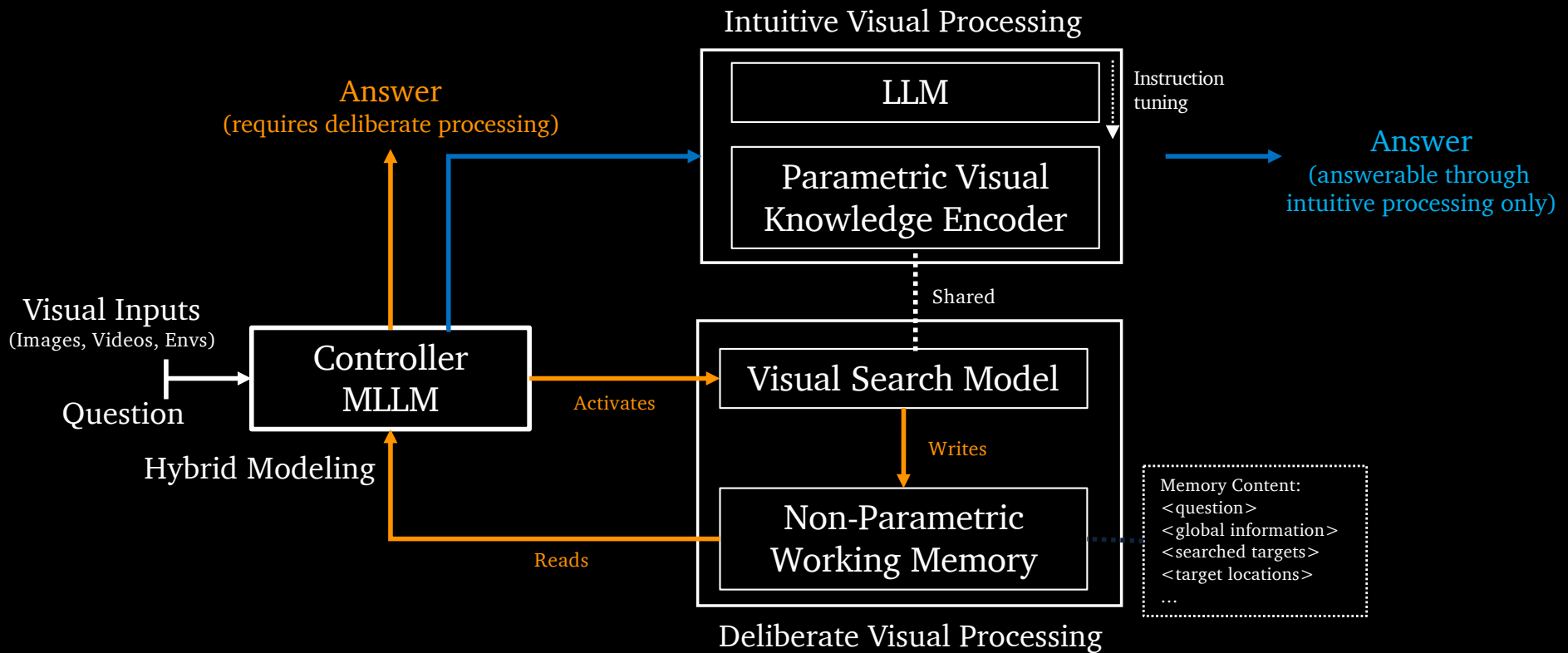
SEAL is a *Meta* Architecture for Multimodal LLMs

Show, sEArch, and Tell



VWM:
<question>
<global image>
<searched targets>
<target locations>

An architecture for hybrid visual processing



V*Bench – a Vision-centric VQA Benchmark

- 191 high-resolution images from SA-1B dataset
- Attribute recognition task (115 samples) & spatial relationship reasoning task (76 samples)
- Average image resolution 2246×1582, average target resolution 57×62
- Multiple-choice questions



Question:

What is the color of the clock?

Options:

- The color of the clock is green.
- The color of the clock is black.
- The color of the clock is red.
- The color of the clock is yellow.



Question:

What is the material of the stool?

Options:

- The material of the stool is plastic.
- The material of the stool is wood.
- The material of the stool is steel.
- The material of the stool is bamboo.



Question:

Is the red balloon above of white balloon?

Options:

- The red balloon is below the white balloon.
- The red balloon is above the white balloon.



Question:

Is the broom on the left or right side of the folded chair?

Options:

- The broom is on the left side of the folded chair.
- The broom is on the right side of the folded chair.

V*Bench – Evaluation Results of Multimodal Systems

	Attribute (%)	Spatial (%)	Overall (%)
Human	98.26	100.00	98.95
Random Guess	26.73	50.00	35.99
<i>Open-source end-to-end MLLMs</i>			
BLIP2 [23]	26.95	53.94	37.69
MiniGPT-4 [63]	30.43	50.00	38.22
LLaVA [28]	23.47	53.94	35.59
InstructBLIP [8]	25.21	47.36	34.02
Otter [22]	26.95	56.57	38.74
LLaVA-1.5 [27]	43.47	56.57	48.68
<i>LLM tool-using pipelines</i>			
MM-React [53]	34.78	51.31	41.36
VisualChatGPT [54]	30.43	48.68	37.69
Visprog [12]	31.30	56.57	41.36
<i>Commercial chatbot systems</i>			
Bard [11]	31.30	46.05	37.17
Gemini Pro [9]	40.86	59.21	48.16
GPT-4V [35]	51.30	60.52	54.97
SEAL (Ours)	74.78	76.31	75.39

} Random Guess Level

Takeaways:

A good MLLM should be able to:

1. Acknowledge that initial visual information is NOT enough, and they cannot see.
2. Explicitly list additional visual information needed.
3. Understand and integrate the search results after the visual search process.
4. Allocate more FLOPS for more complex tasks during test time.

*(*After ChatGPT-o1 🍓, I think people will be more open to this.)*

Note: these are not engineering hacks!

Key Question #3:

What are the necessary components towards building better, vision-centric MLLMs?



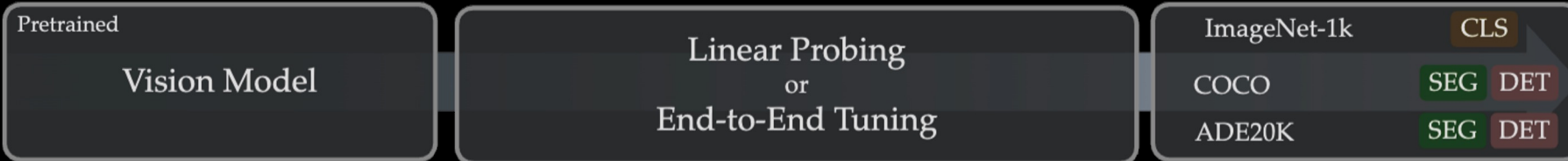
Cambrian-1

A Fully Open, *Vision-Centric* Exploration of Multimodal LLMs

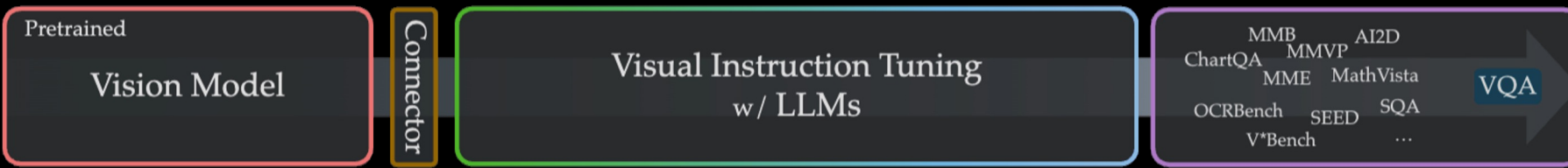
Shengbang Tong*, Ellis Brown*, Penghao Wu*,
Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang,
Shusheng Yang, Adithya Iyer, Xichen Pan, Austin Wang,
Rob Fergus, Yann LeCun, Saining Xie

NeurIPS 2024 (Oral)

Visual Representation Learning - Evaluation Protocols



Multimodal Large Language Models



Overview: Core Components in MLLMs

LLM

Vision Backbone

Vision-Language Connector

Instruction Tuning Data

Instruction Tuning Recipe

Evaluation Protocol

Overview: Core Components in MLLMs

— ~~LLM~~

Vision Backbone

Vision-Language Connector

Instruction Tuning Data

Instruction Tuning Recipe

Evaluation Protocol



Visual
Representations



Connector Design



Instruction Tuning
Data

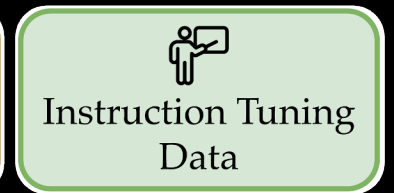
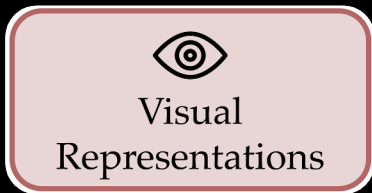
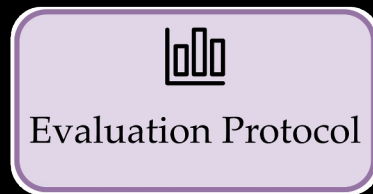


Instruction Tuning
Recipe



Evaluation Protocol

Overview: Core Components in MLLMs



Overview



Evaluation Protocol



Instruction Tuning
Recipe



Visual
Representations



Connector Design



Instruction Tuning
Data

Visual Representations for MLLMs

How to evaluate
visual reprs.?

What visual
reprs. to use?



Cowboy Hat
Sorrel
Cowboy Boot
Barrel
Revolver

Class Label Supervised
ImageNet-1K [105]

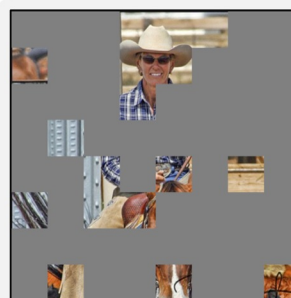


A cowboy rides a
horse at a rodeo.

Language Supervised
CLIP [102]



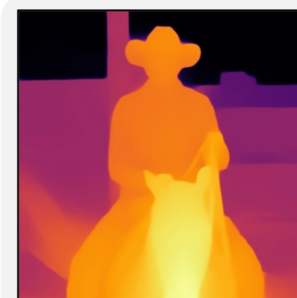
SSL-Contrastive
DINOv2 [96]



SSL-Masking
MAE [45]



Diffusion
Stable Diffusion [104]



Depth Supervised
MiDaS [13]



Segmentation Supervised
SAM [61]

Visual Representations for MLLMs

Supervision Type	Method	Architecture	Patch Size	Res.	# Tok.	Hidden Size
Language-Supervised						
Language	OpenAI CLIP	ViT-L	14	336	576	768
	DFN-CLIP	ViT-L	14	224	256	1024
	DFN-CLIP	ViT-H	14	378	729	1280
	EVA-CLIP-02	ViT-L	14	336	576	1024
	SigLIP	ViT-L	16	384	576	1024
	SigLIP	ViT-SO400M	14	384	729	1152
	OpenCLIP	ConvNeXT-L	-	512	¹ 576	1536
	OpenCLIP	ConvNeXT-L	-	1024	¹ 576	1536
	OpenCLIP	ConvNeXT-XXL	-	1024	¹ 576	3072
Self-Supervised						
Contrastive	DINOv2	ViT-L	14	336	576	1024
	DINOv2	ViT-L	14	518	¹ 576	1024
	MoCo v3	ViT-B	16	224	196	768
	MoCo v3	ViT-L	16	224	196	1024
Masked	MAE	ViT-L	16	224	196	1024
	MAE	ViT-H	14	224	256	1280
JEPA	I-JEPA	ViT-H	14	224	256	1280
Other						
Segmentation	SAM	ViT-L	16	1024	¹ 576	1024
	SAM	ViT-L	16	1024	¹ 576	1280
Depth	MiDaS 3.0	ViT-L	16	384	576	1024
	MiDaS 3.1	ViT-L	16	518	1024	1024
Diffusion	Stable Diffusion 2.1	VAE+UNet	16	512	1024	3520
Class Labels	SupViT	ViT-L	16	224	196	1024
	SupViT	ViT-H	14	224	256	1280

23 models!

Table 9 | Catalog of all vision backbones tested. ¹ denotes that the visual tokens have been interpolated down to the specified length.

Overview



Evaluation Protocol



Instruction Tuning
Recipe



Visual
Representations



Connector Design



Instruction Tuning
Data

Overview



Evaluation Protocol



Instruction Tuning
Recipe



Visual
Representations



Connector Design



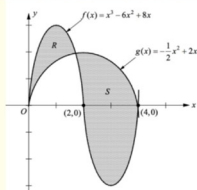
Instruction Tuning
Data

Evaluation Protocol

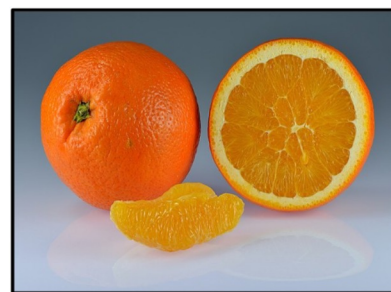
Question: <image 1> The region bounded by the graph as shown above. Choose an integral expression that can be used to find the area of R.

Options:

- (A) $\int_0^{1.5} [f(x) - g(x)] dx$
- (B) $\int_0^{1.5} [g(x) - f(x)] dx$
- (C) $\int_0^2 [f(x) - g(x)] dx$
- (D) $\int_0^2 [g(x) - x(x)] dx$



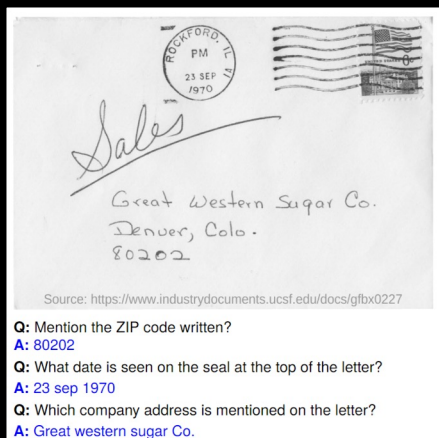
MMMU [Yue, et al. 2024]



Q: what is the color of this object?

- A. Purple
 - B. Pink
 - C. Gray
 - D. Orange
- GT: D

MM-Bench [Liu, et al. 2024]



Source: <https://www.industrydocuments.ucsf.edu/docs/gfbx0227>

Q: Mention the ZIP code written?
A: 80202

Q: What date is seen on the seal at the top of the letter?
A: 23 sep 1970

Q: Which company address is mentioned on the letter?
A: Great western sugar Co.

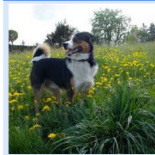
DocVQA [Mathew, et al. 2020]




Where can we go from the current lane? **A. Turn left.** B. Go straight. C. Turn left and go straight. D. Turn right.

RealWorldQA [Grok, et al. 2024]



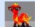

Is the dog facing left or right from the camera's perspective?



(a) Left



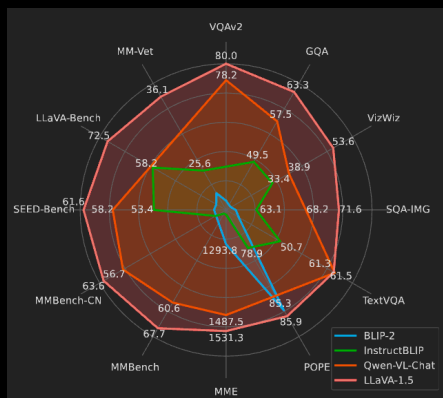
(b) Right

	(b)	(b)	✗
	(a)	(a)	✗
	(b)	(b)	✗
	(a)	(a)	✗

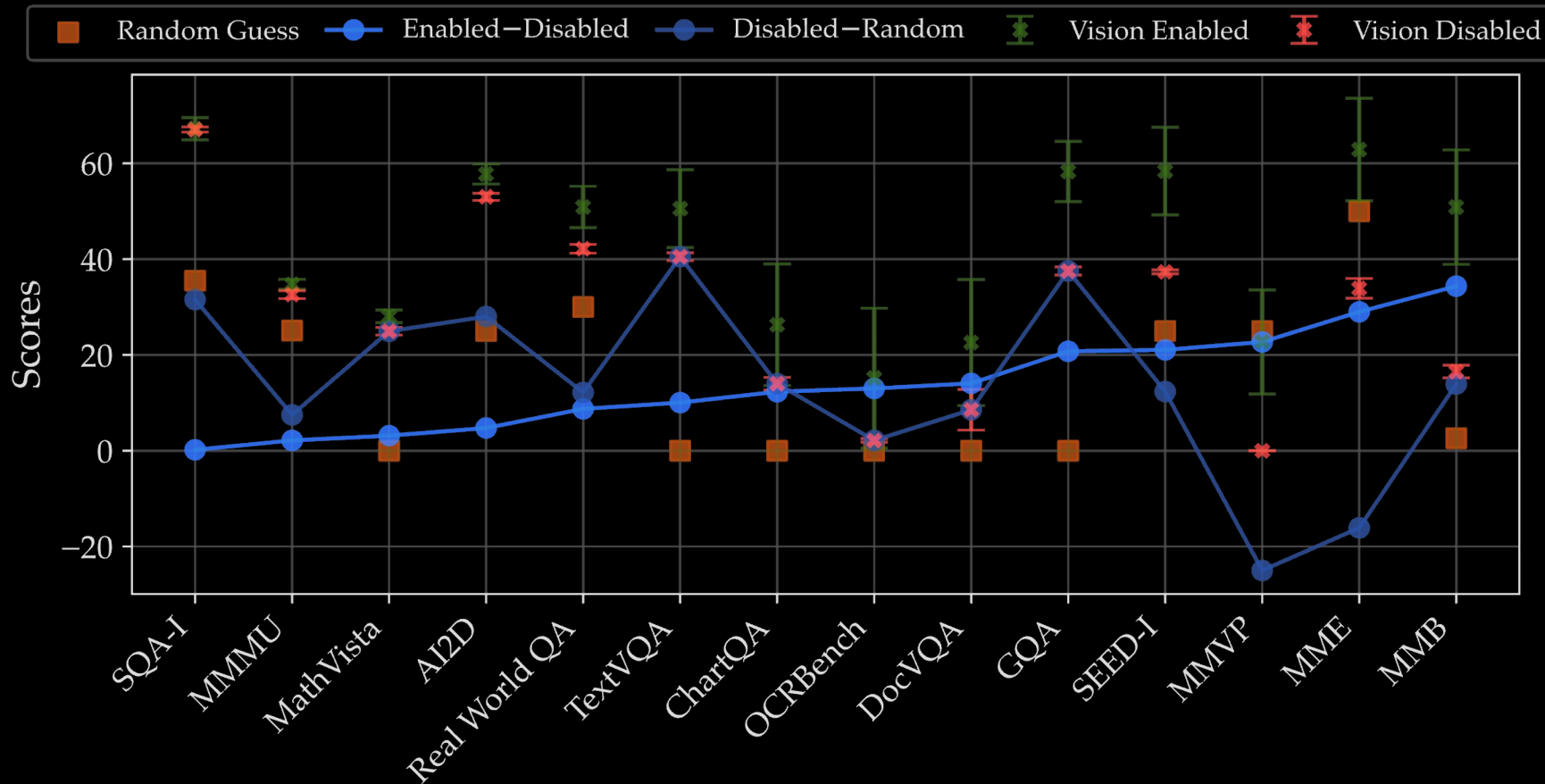
MMVP [Tong, et al. 2024]

and a lot more...

How should we systematically evaluate an MLLM and interpret the evaluation results?

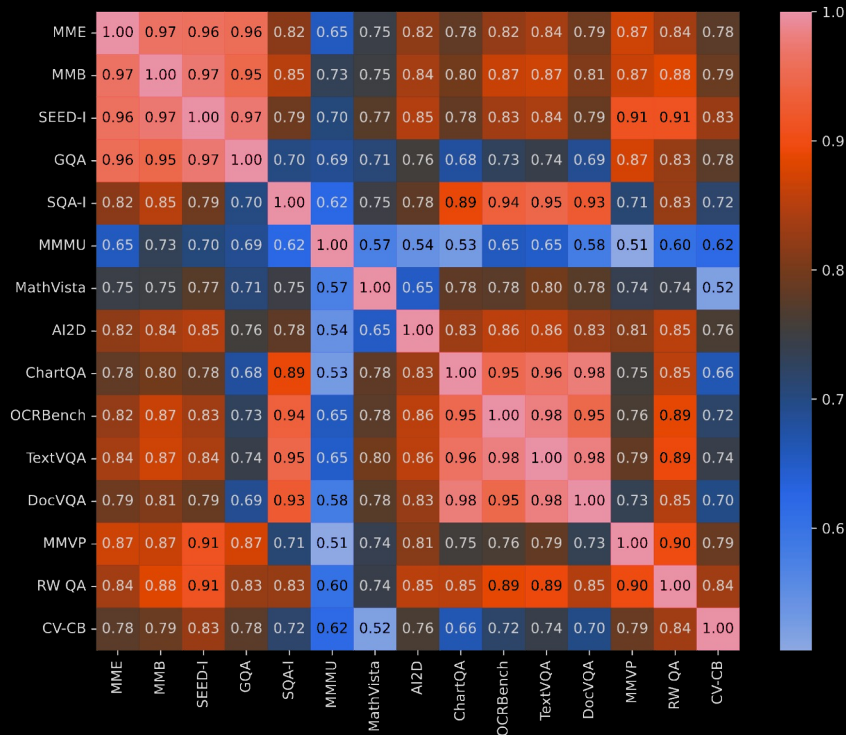


Who's answering the question: the LLM or MLLM?

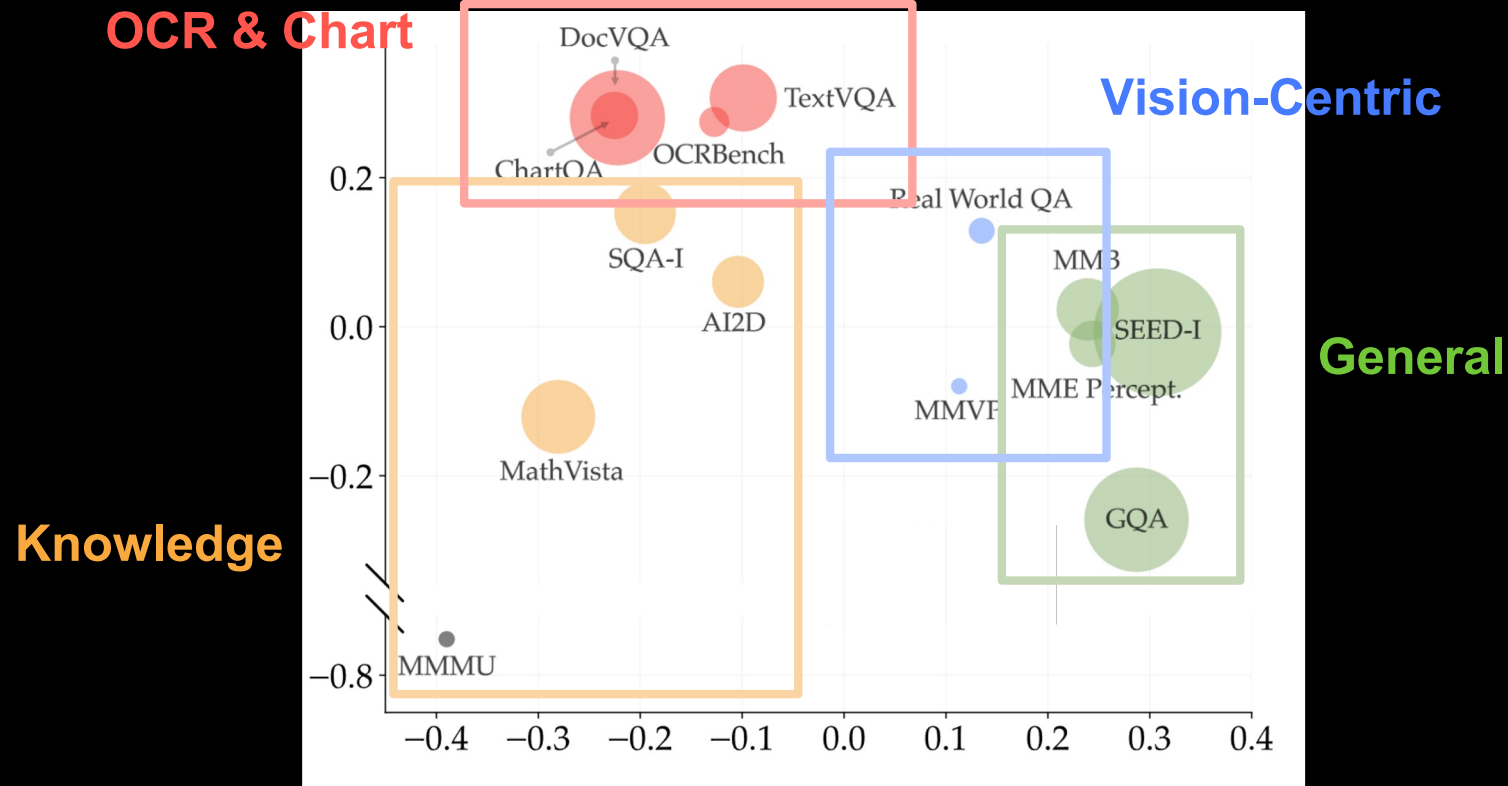


Group Benchmarks by Correlation

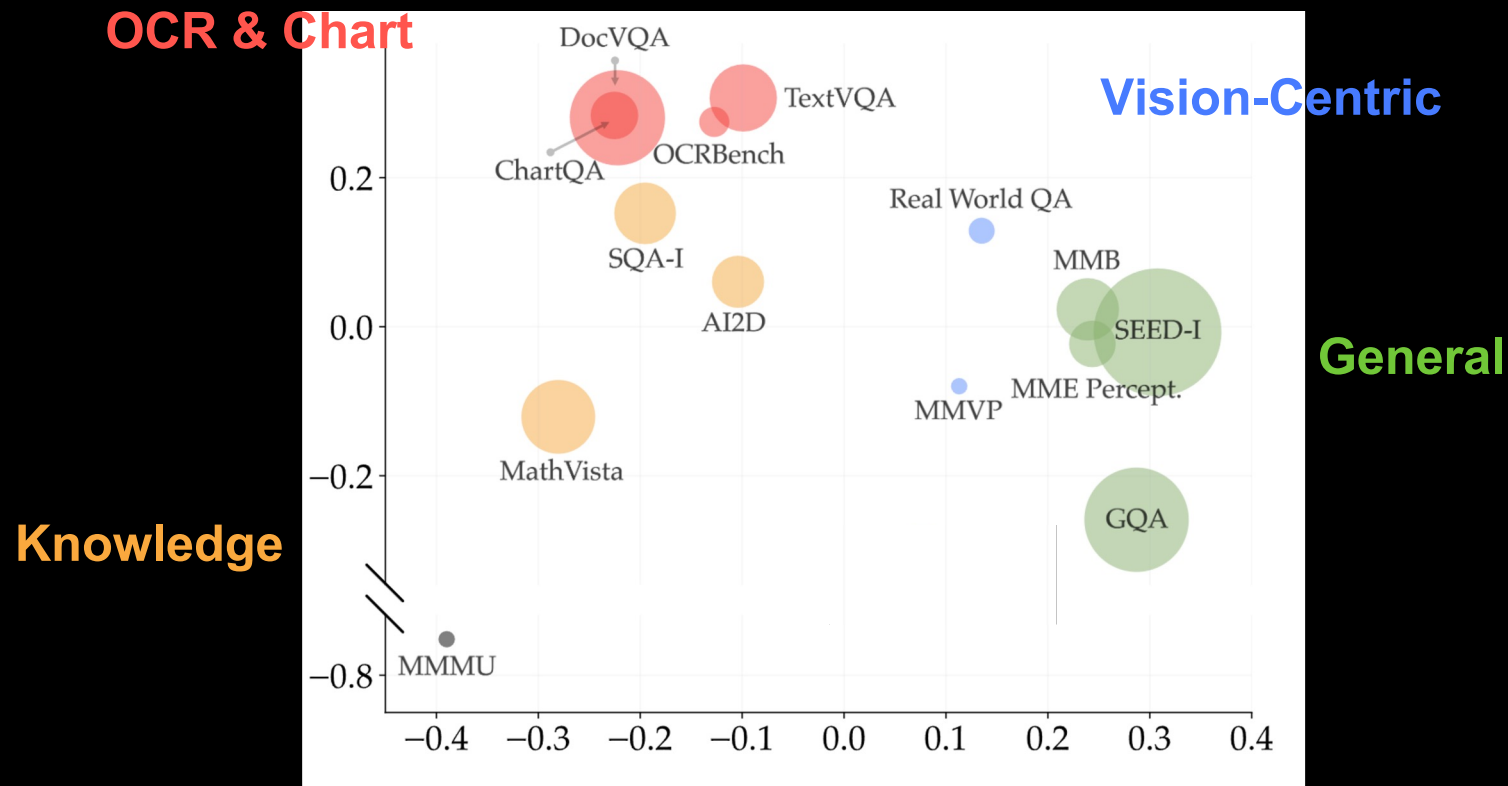
If two benchmarks evaluate on similar domains, they should have a strong correlation



Group Benchmarks by Correlation

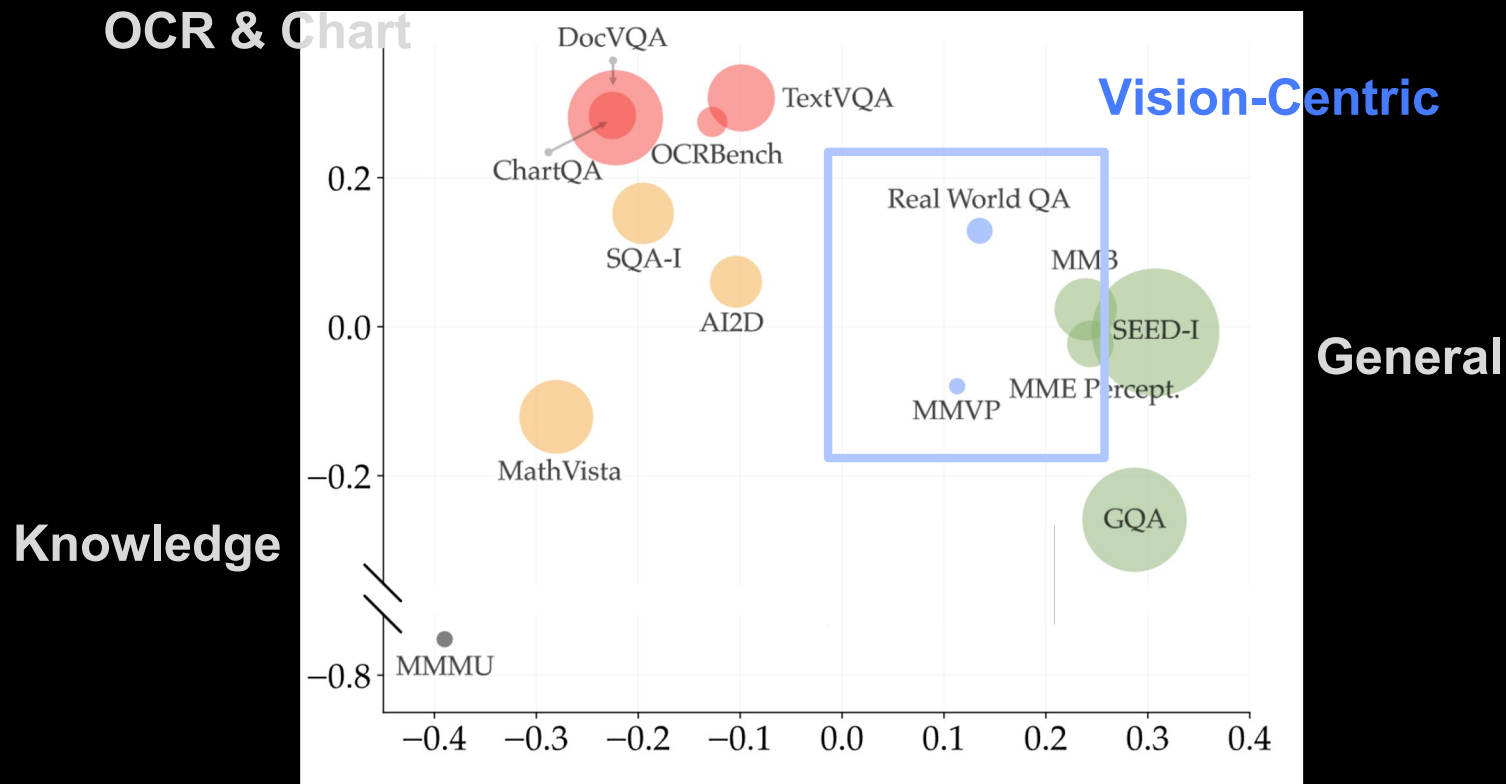


Group Benchmarks by Correlation



Group Benchmarks by Correlation

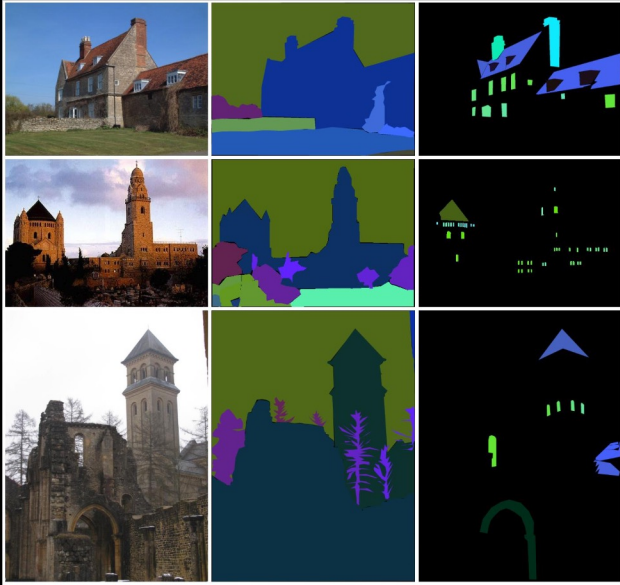
Tiny compared to others!



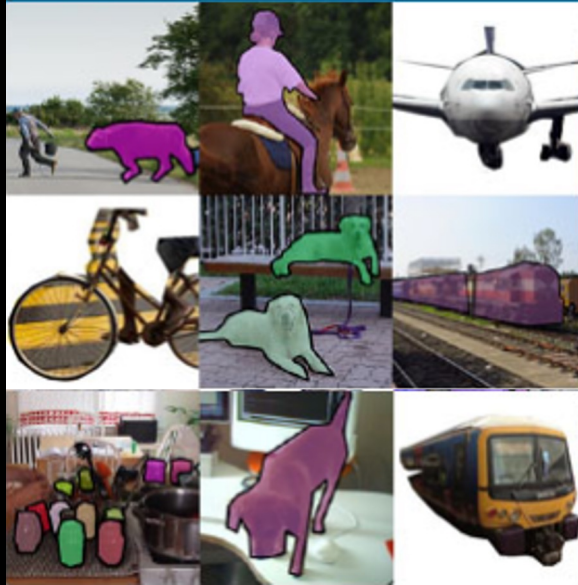
Q: How can we scalably generate *vision-centric* MLLM evaluations?

 Repurpose existing vision benchmarks!

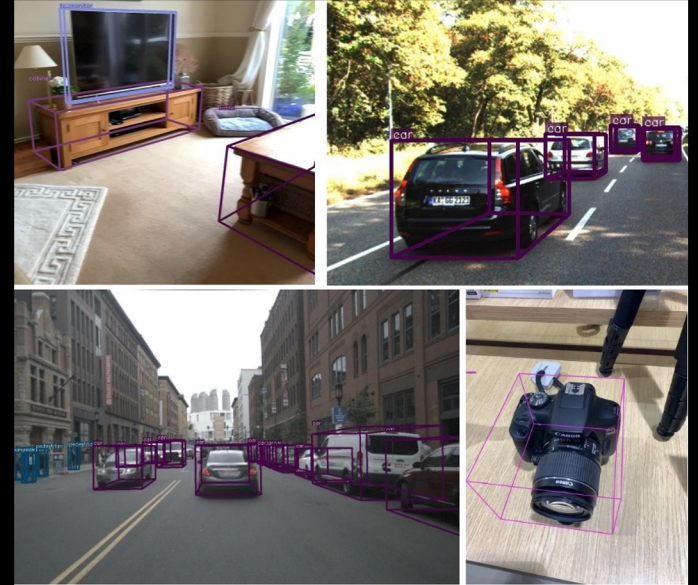
CV-Bench



ADE20K



MSCOCO



Omni3D

CV-Bench

2D

3D

Spatial Relationship



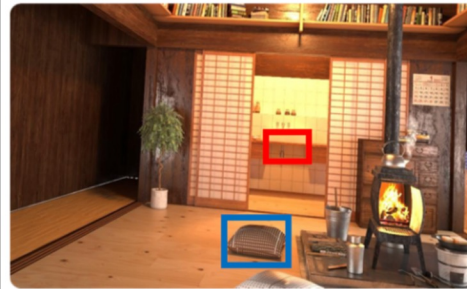
Where is the cave located with respect to the trees?

Object Count



How many cars are in the image?

Depth Order



Which is closer to the camera, **sink** or **pillow**?

Relative Distance



Which is closer to the **chair**, **refrigerator** or **door**?

Source benchmark: ADE20K [145] and COCO [72]

Source benchmark: Omini3D [16]

CV-Bench

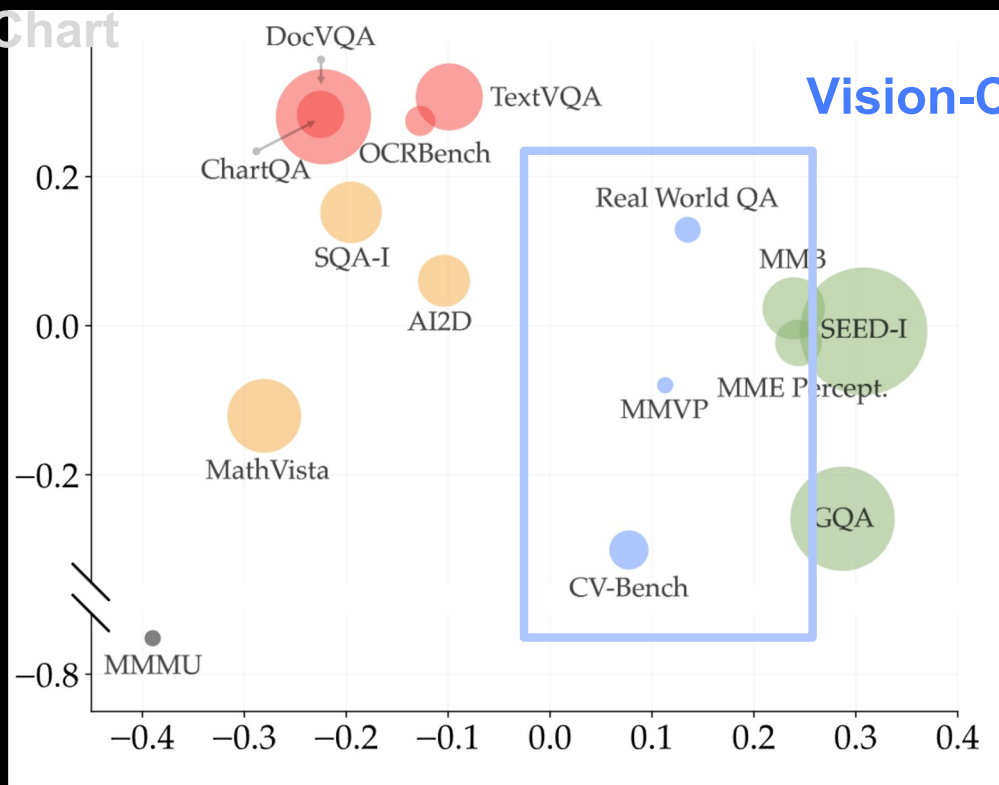
2,638 manually-inspected examples

Type	Task	Description	Sources	# Samples
2D	Spatial Relationship	Determine the relative position of an object w.r.t. the anchor object. Consider left-right or top-bottom relationship.	ADE20K COCO	650
	Object Count	Determine the number of instances present in the image.	ADE20K COCO	788
3D	Depth Order	Determine which of the two distinct objects is closer to the camera.	Omni3D	600
	Relative Distance	Determine which of the two distinct objects is closer to the anchor object.	Omni3D	600

Table 1 | Breakdown of the 2D and 3D tasks evaluated in the Cambrian Vision-Centric Benchmark (CV-Bench). The examples are sourced from ADE20K [145], COCO [72], and Omni3D [16].

3.5x more vision-centric examples!

OCR & Chart

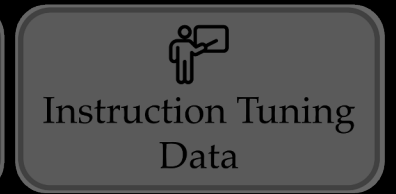
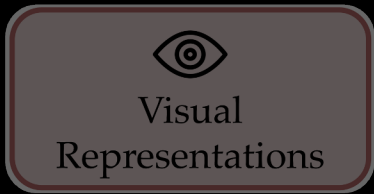
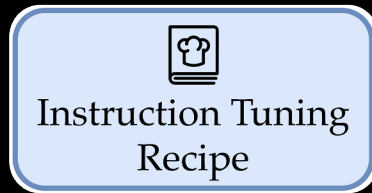


Vision-Centric

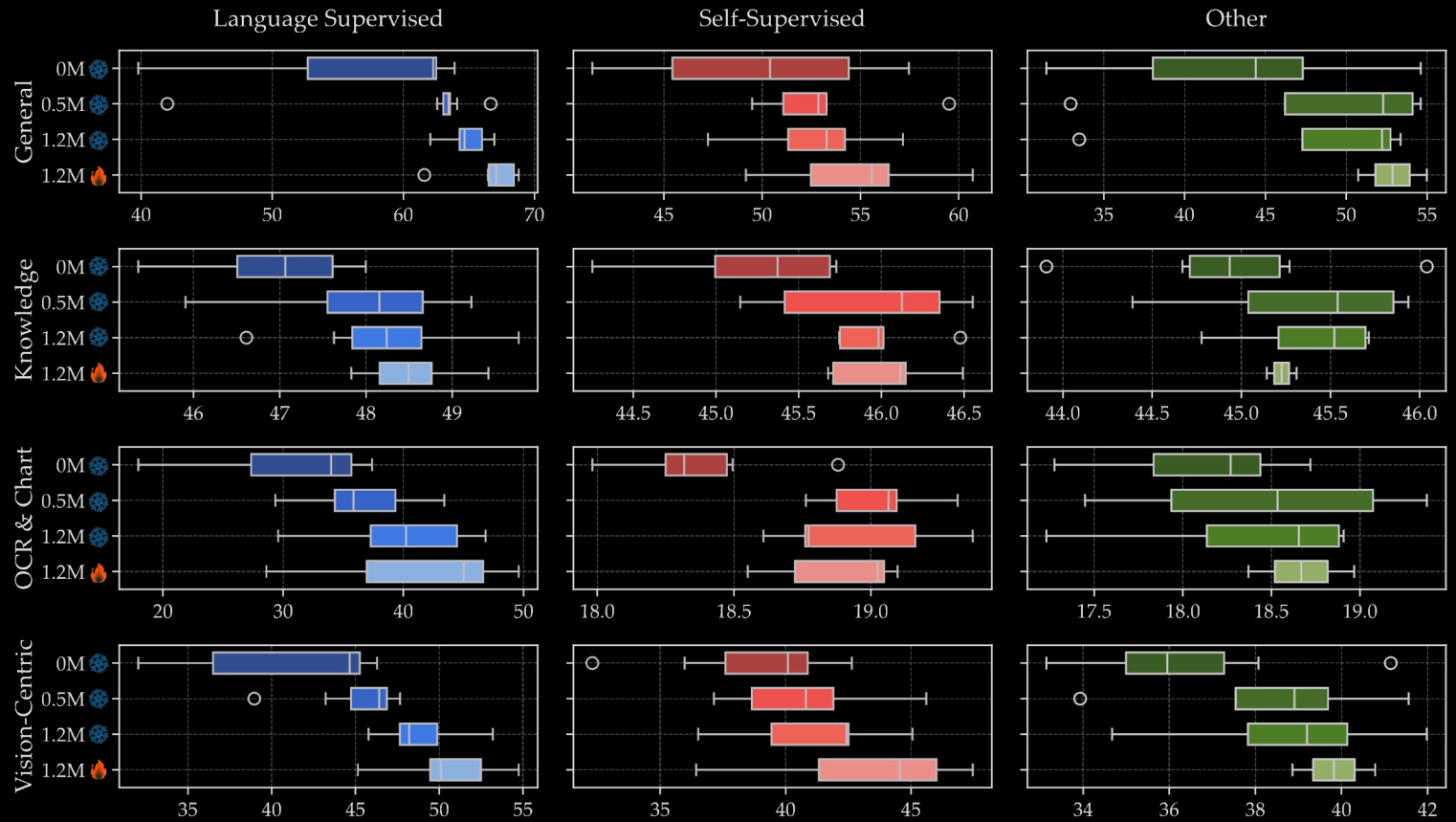
General

Knowledge

Overview

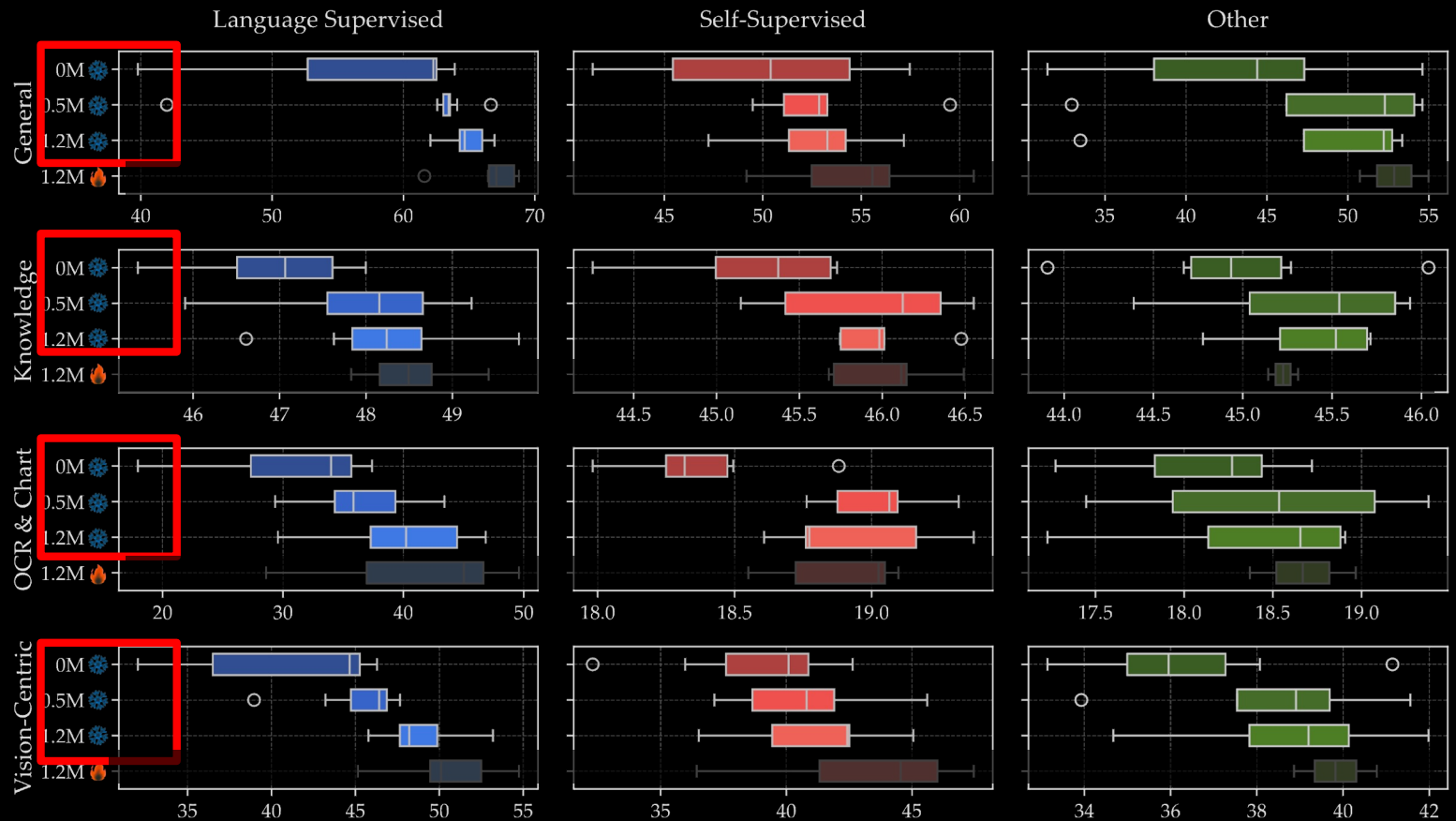


Instruction Tuning Recipe



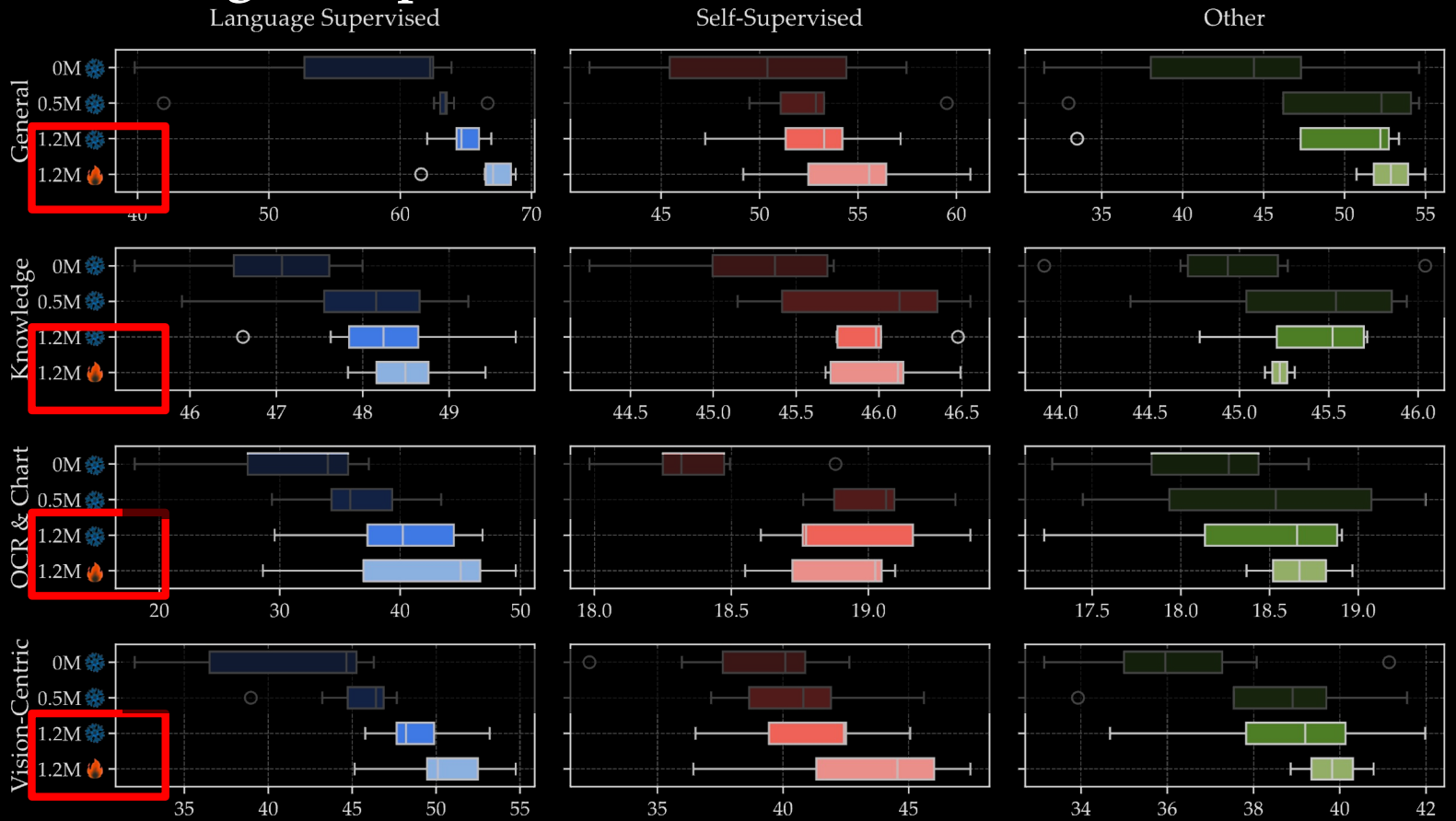
Instruction Tuning Recipe

More
Alignment
Data helps!

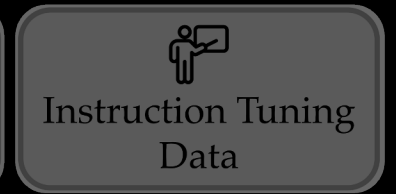
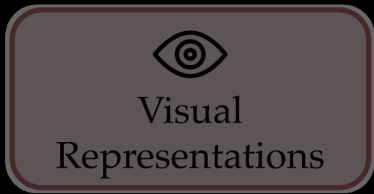
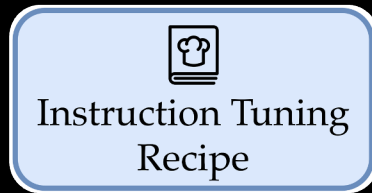


Instruction Tuning Recipe

Unfreezing
Vision
Encoder
Helps 🔥



Overview



Overview



Evaluation Protocol



Instruction Tuning
Recipe



Visual
Representations

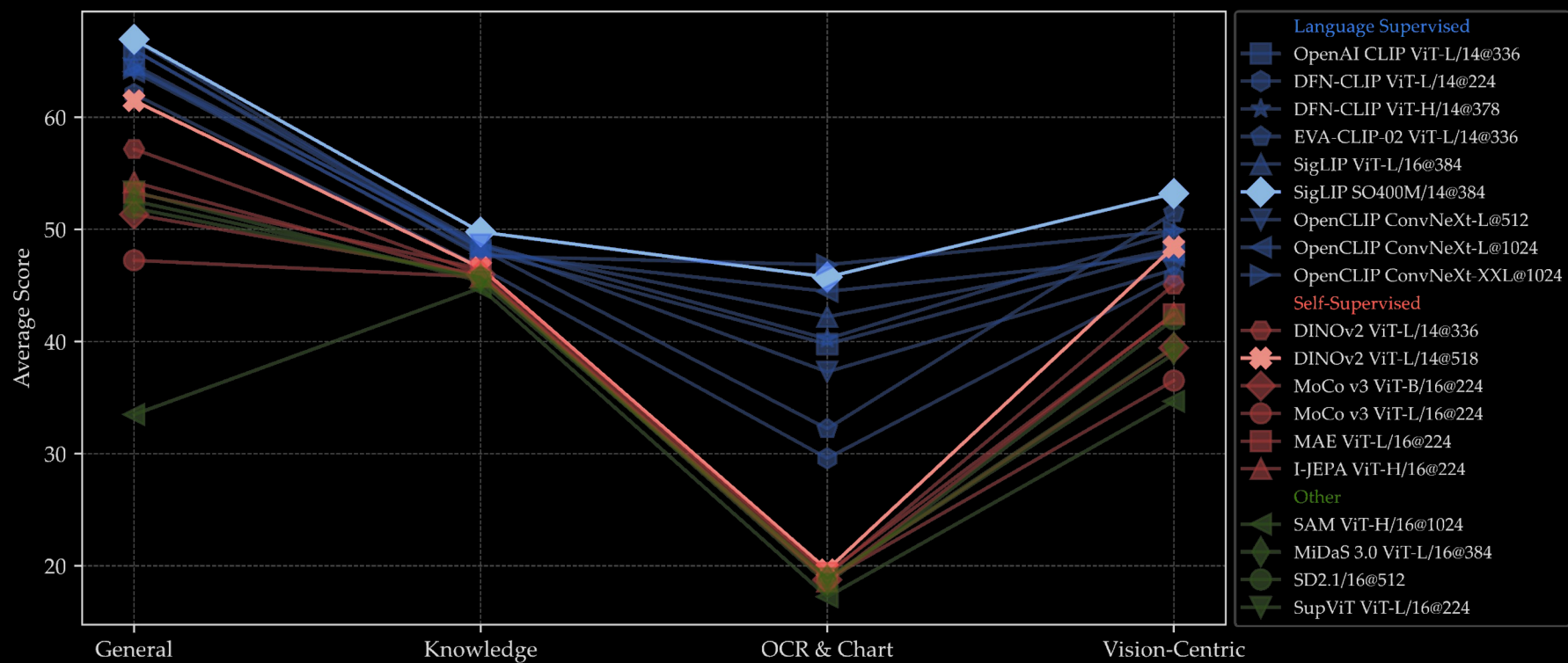


Connector Design



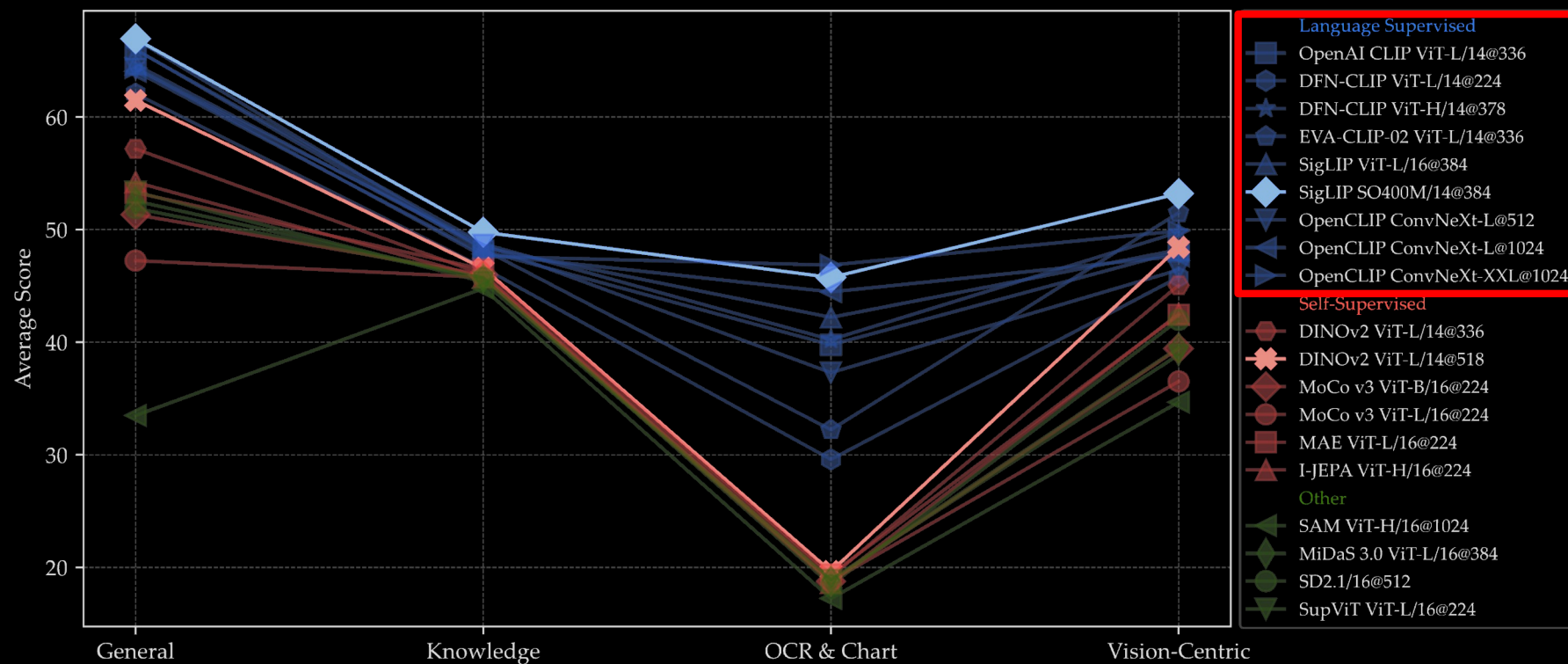
Instruction Tuning
Data

Visual Representation



Visual Representation

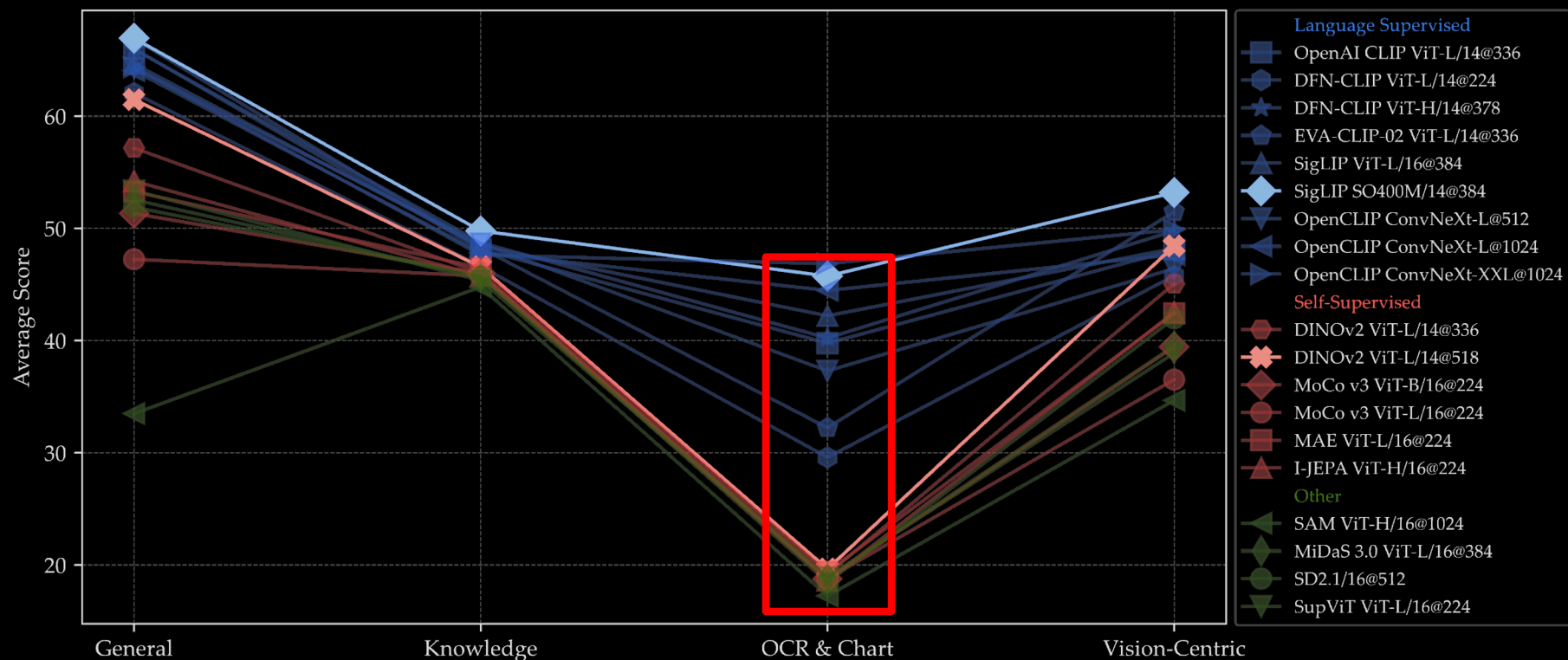
#1 Language Supervised Models are better



Visual Representation

#1 Language Supervised Models are better

#2 Gap is largest in OCR & Chart



Visual Representation

A ConvNet for the 2020s

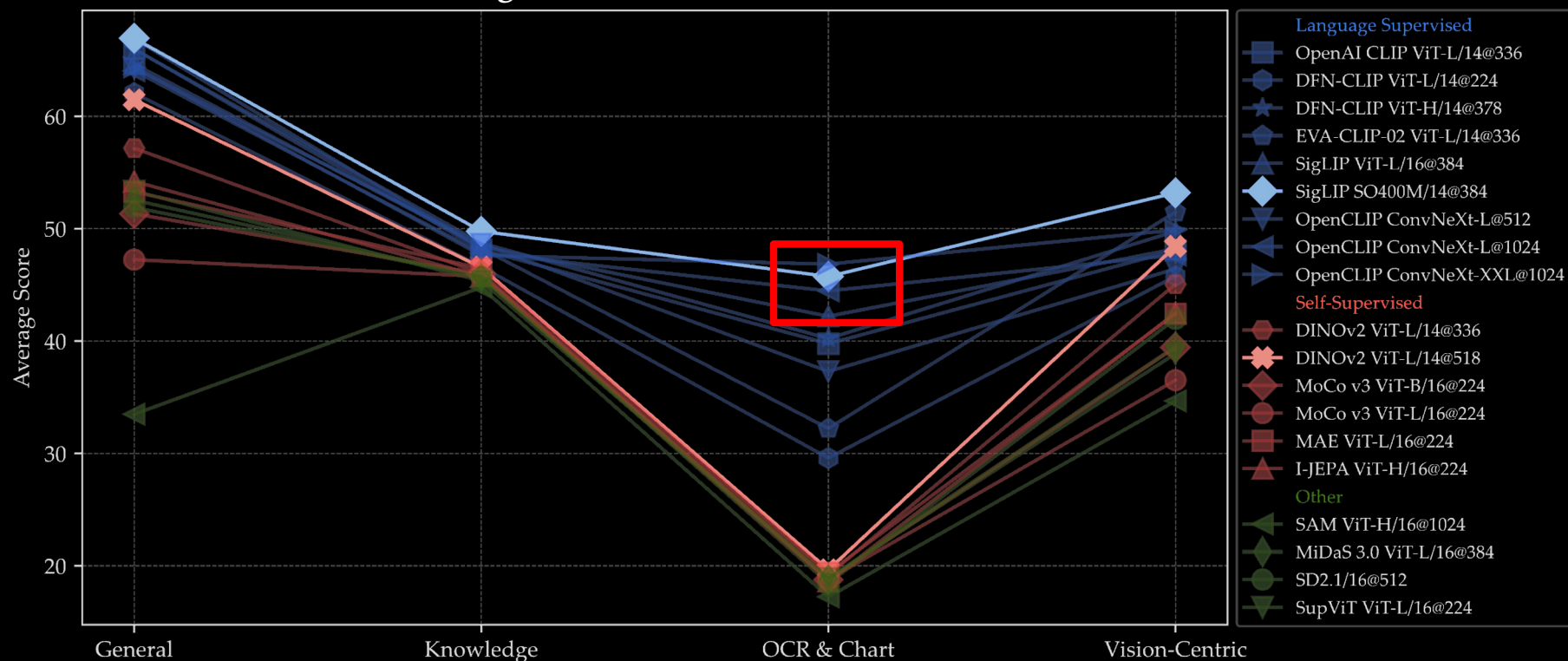
Zhuang Liu^{1,2*} Hanzi Mao¹ Chao-Yuan Wu¹ Christoph Feichtenhofer¹ Trevor Darrell² Saining Xie^{1†}

¹Facebook AI Research (FAIR) ²UC Berkeley

Code: <https://github.com/facebookresearch/ConvNeXt>

#1 Language Supervised Models are better
#3 ConvNets (ConvNeXt) are good at OCR

#2 Gap is largest in OCR & Chart



Visual Representation

#1 Language Supervised Models are better

#2 Gap is largest in OCR & Chart

#3 ConvNets are good at OCR

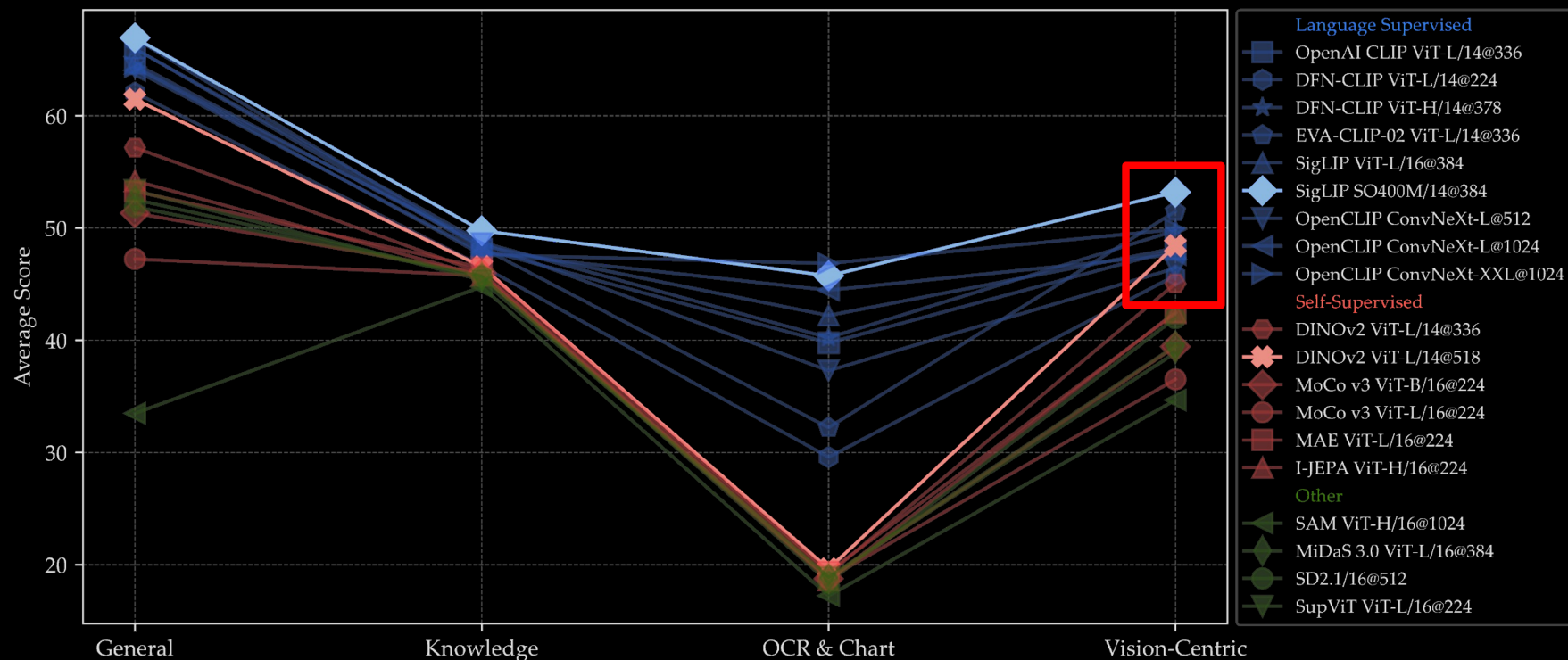
Language Supervised

Model	Architecture	All	G	K	O	V
SigLIP	ViT-SO400M/14@384	1	1	1	2	1
OpenCLIP	ConvNeXt-XXL@1024	2	6	8	1	3
DFN-CLIP	ViT-H/14@378	3	4	2	5	4
OpenCLIP	ConvNeXt-L@1024	4	8	7	3	8
SigLIP	ViT-L/16@384	5	5	4	4	6
OpenAI CLIP	ViT-L/14@336	6	3	6	6	7
EVA-CLIP-02	ViT-L/14@336	7	2	5	8	2
OpenCLIP	ConvNeXt-L@512	8	7	3	7	9
DFN-CLIP	ViT-L/14@224	9	9	9	9	10
DINOv2*	ViT-L/14@518	10	10	10	10	5

Visual Representation

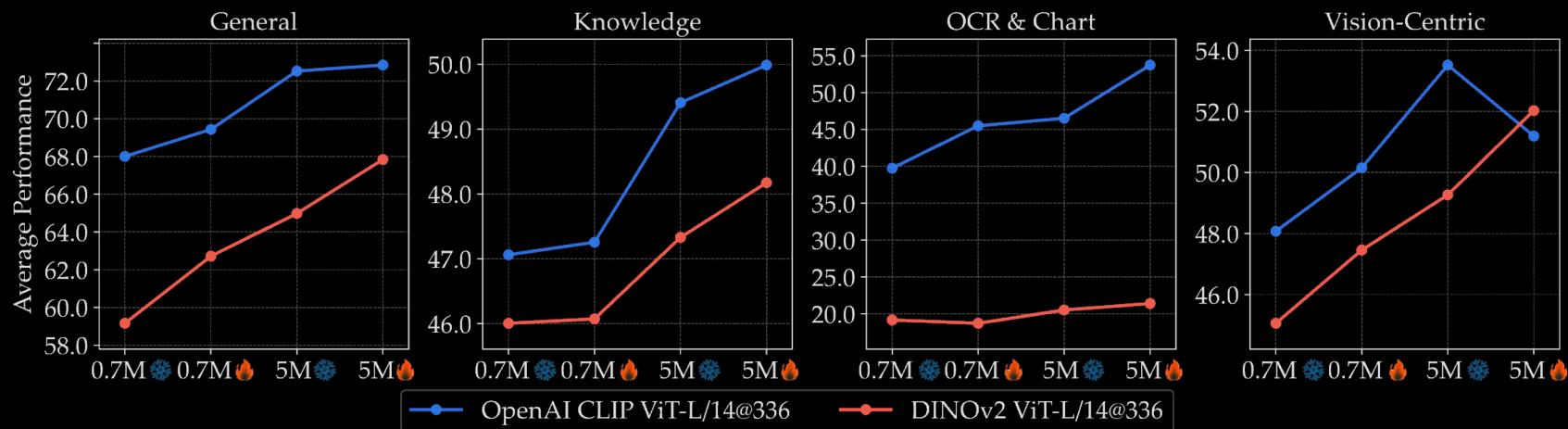
#1 Language Supervised Models are better
#3 ConvNets are good at OCR

#2 Gap is largest in OCR & Chart
#4 Best SSL model is good at Vision-Centric tasks



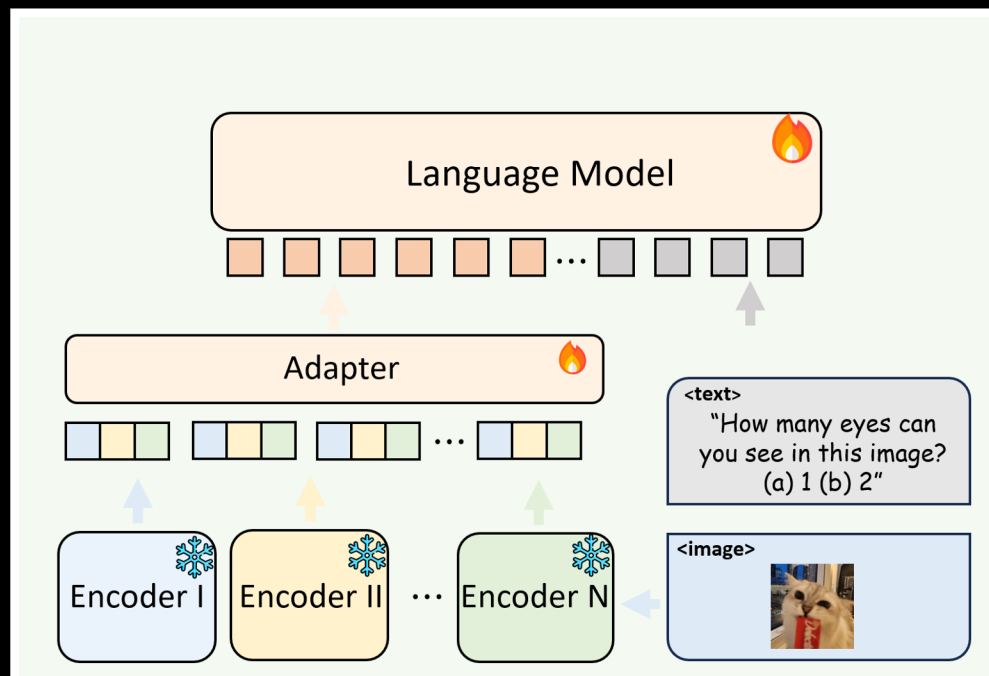
Visual Representation

Training with more data narrows the gap between Language-Supervised and Self-Supervised Models



Visual Representation

Combining different models improves performance

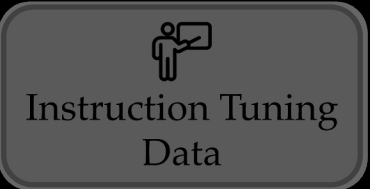
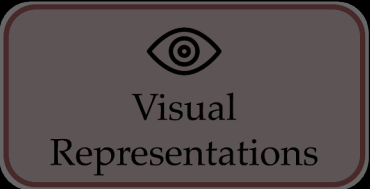


Visual Representation

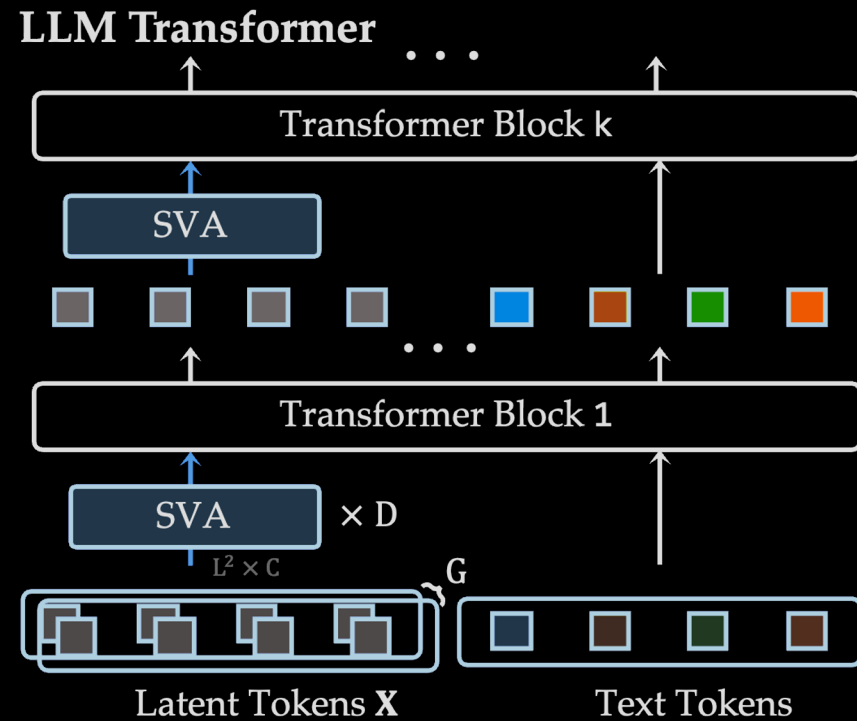
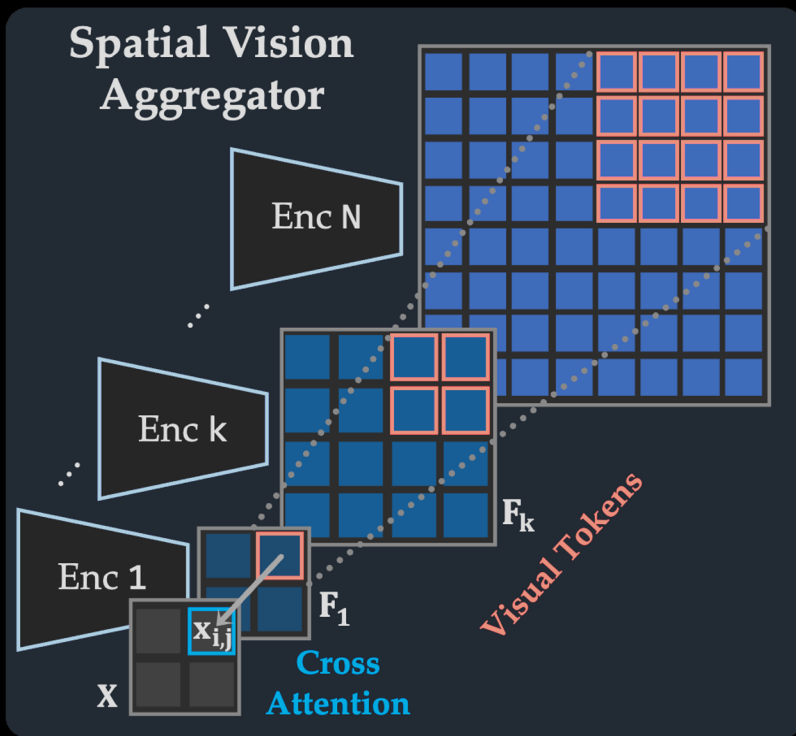
Combining different models improves performance (using a simple linear connector)

Vision Backbone	Average	General				Knowledge				OCR & Chart				Vision-Centric			
		MME ^P	MMB	SEED ^I	GQA	SQA ^I	MMMU ^V	MathVista ^M	AI2D	ChartQA	OCRBench	TextVQA	DocVQA	MMVP	RealWorldQA	CV-Bench ^{2D}	CV-Bench ^{3D}
Method																	
SigLIP+DINOv2	51.61	1,432.02	61.28	65.99	63.30	68.82	35.69	29.40	60.01	43.00	35.70	60.40	37.54	30.00	53.99	55.52	53.58
SigLIP+DINOv2+ConvNext	54.52	1,503.51	63.83	67.97	63.95	70.40	35.99	29.30	60.69	48.20	36.90	64.97	45.53	34.67	58.69	55.74	60.33
SigLIP+DINOv2+ConvNext+CLIP	54.74	1,479.46	63.32	67.63	64.04	71.39	35.49	29.10	59.88	50.24	39.60	64.55	46.12	32.67	58.95	58.54	60.42
SigLIP+ConvNext	54.53	1,494.97	64.60	67.98	63.58	71.05	34.90	29.80	60.85	50.64	38.00	64.53	46.52	32.00	57.91	58.83	56.58
CLIP+ConvNext	54.45	1,511.08	63.83	67.41	63.63	70.80	35.09	30.40	59.91	51.32	35.00	64.45	47.88	33.33	57.25	56.32	59.08
SigLIP+DINOv2+ConvNext	53.78	1,450.64	63.57	67.79	63.63	71.34	34.80	30.20	61.04	49.32	37.70	64.05	45.83	30.00	56.21	58.08	54.33
SigLIP+CLIP+ConvNext	54.53	1,507.28	63.23	68.64	63.63	71.10	35.89	30.90	59.97	52.36	38.50	65.40	47.92	28.67	57.25	57.66	55.92

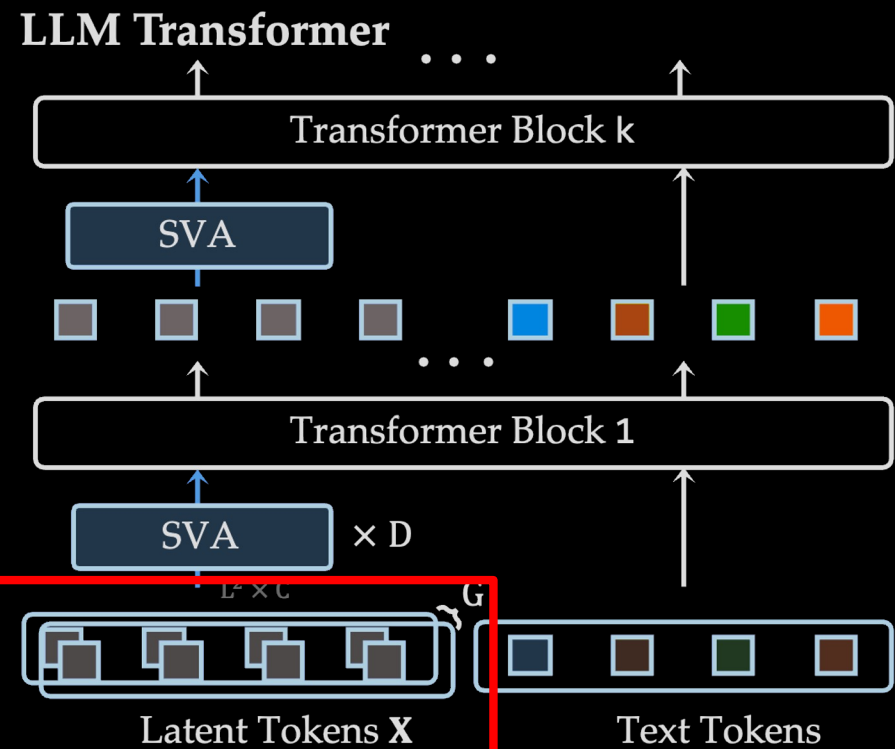
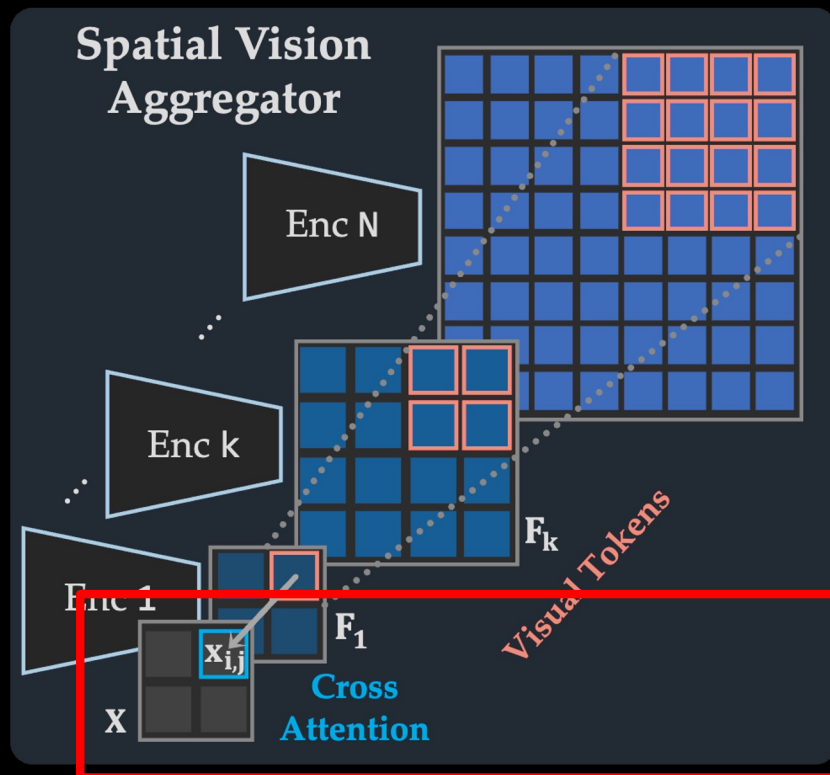
Overview



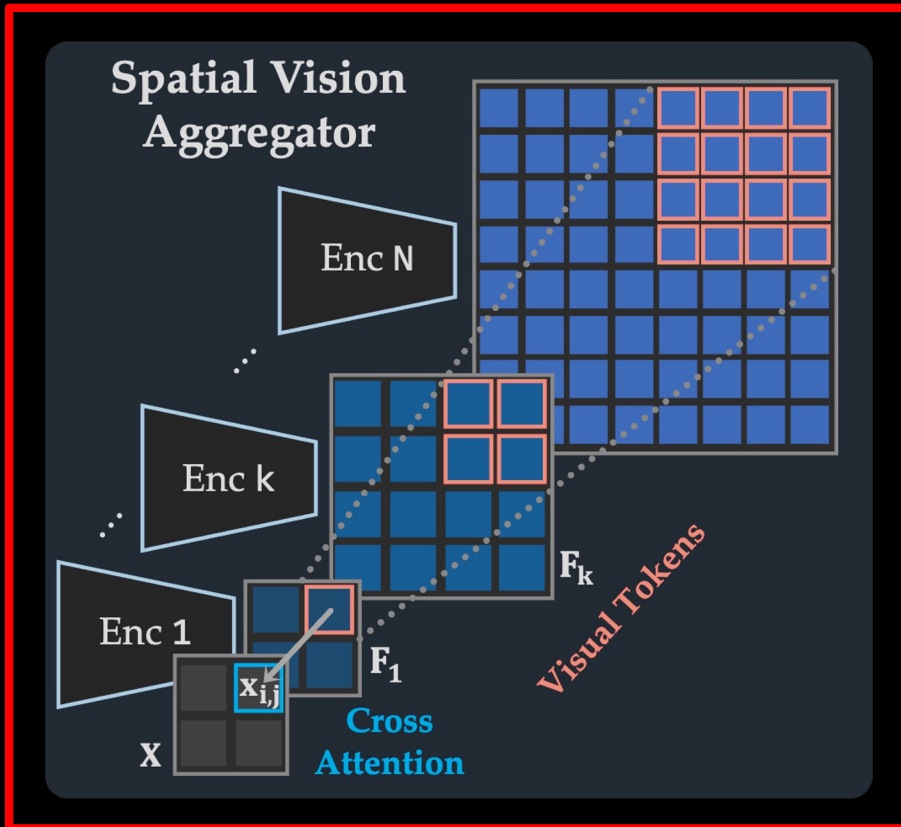
Connector Design - SVA



Connector Design - SVA

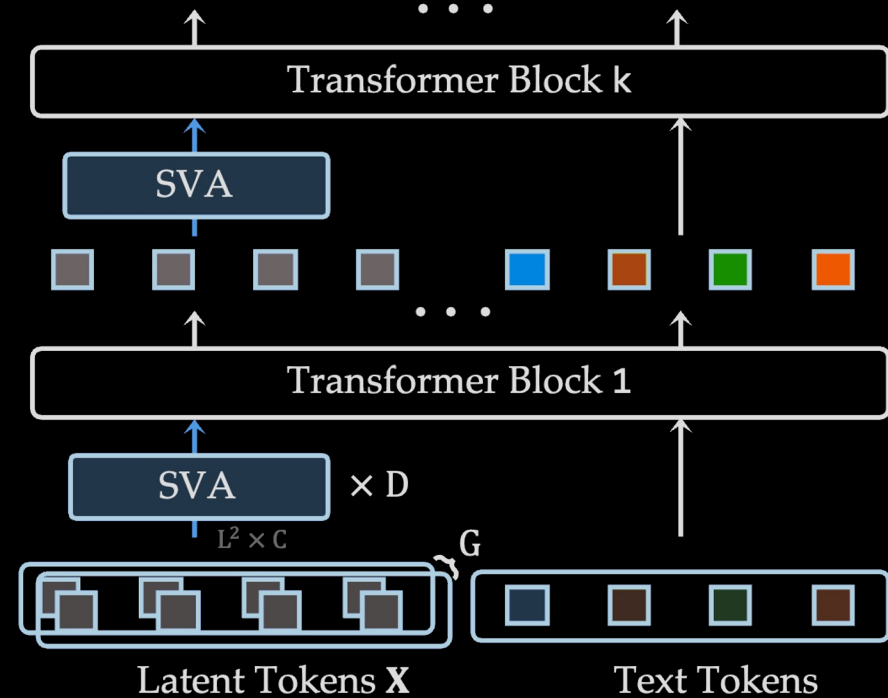


Connector Design - SVA

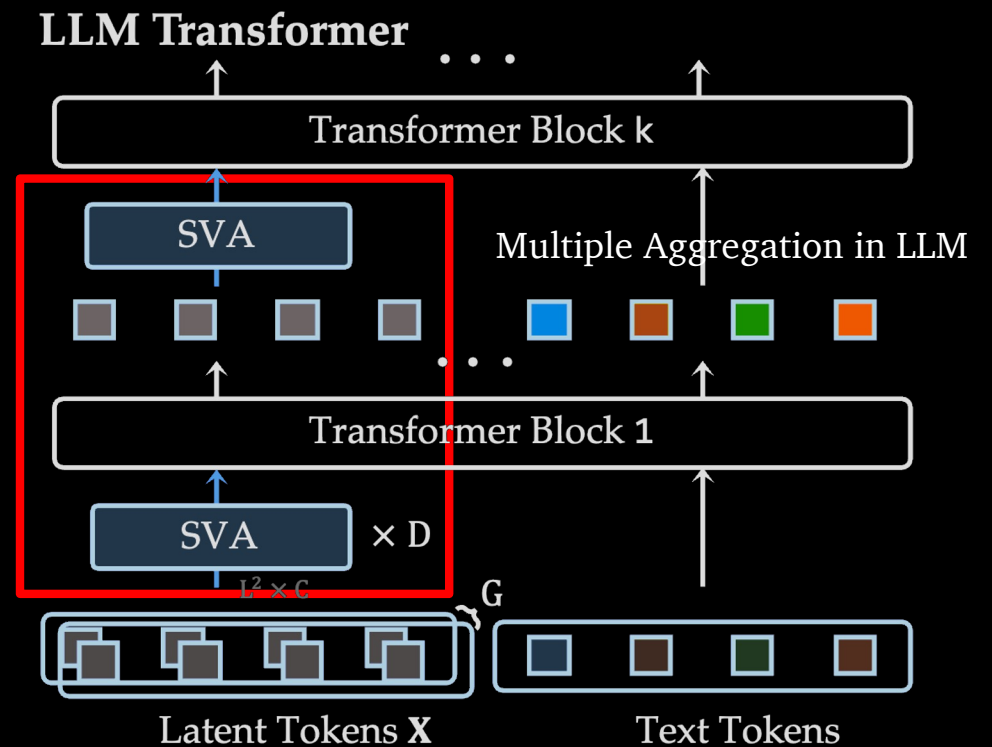
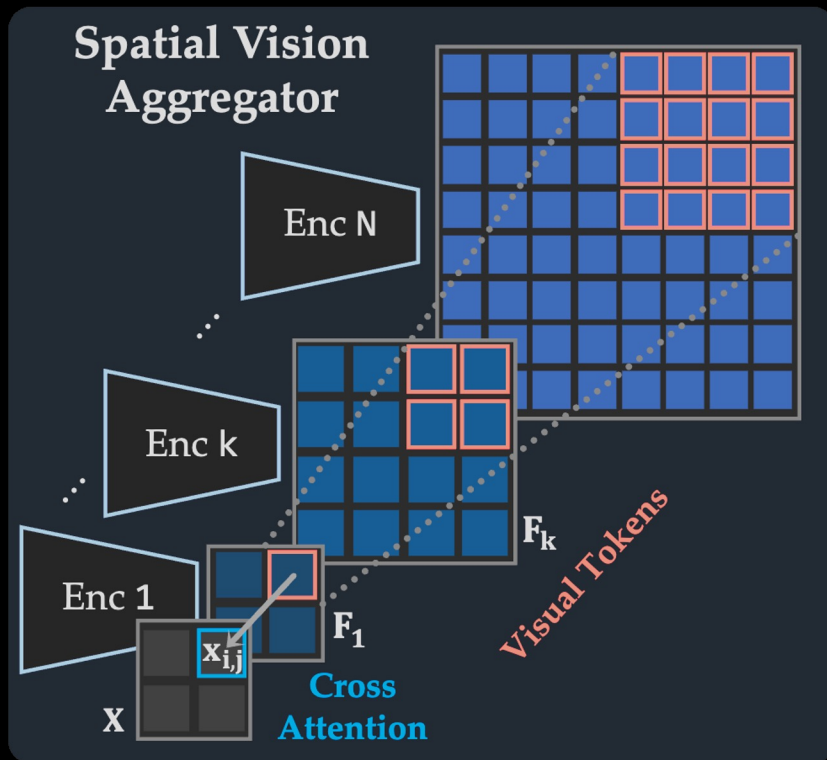


Spatial Inductive Bias

LLM Transformer



Connector Design - SVA



Connector Design - SVA

Spatial Inductive Bias is important especially for OCR&Chart and Vision-Centric Domains

Connector	General	Knowledge	OCR & Chart	Vision-Centric
Concat. [117]	67.2	48.9	50.1	52.6
Resampler [51]	63.1	46.5	27.1	42.6
SVA-no-multi-agg	68.0	49.5	55.2	52.6
SVA	68.5	49.7	55.5	53.2

Connector Design - SVA

Spatial Inductive Bias is important especially for OCR&Chart and Vision-Centric Domains

Multiple Aggregation in LLM further improves performance

Connector	General	Knowledge	OCR & Chart	Vision-Centric
Concat. [117]	67.2	48.9	50.1	52.6
Resampler [51]	63.1	46.5	27.1	42.6
SVA-no-multi-agg	68.0	49.5	55.2	52.6
SVA	68.5	49.7	55.5	53.2


Overview



Evaluation Protocol




Instruction Tuning Recipe



Visual Representations



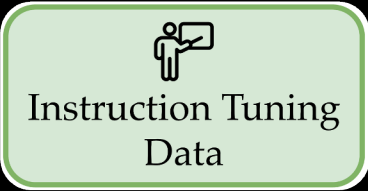
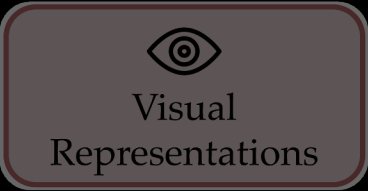
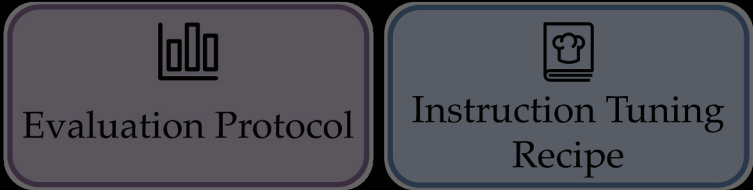
Connector Design



Instruction Tuning Data

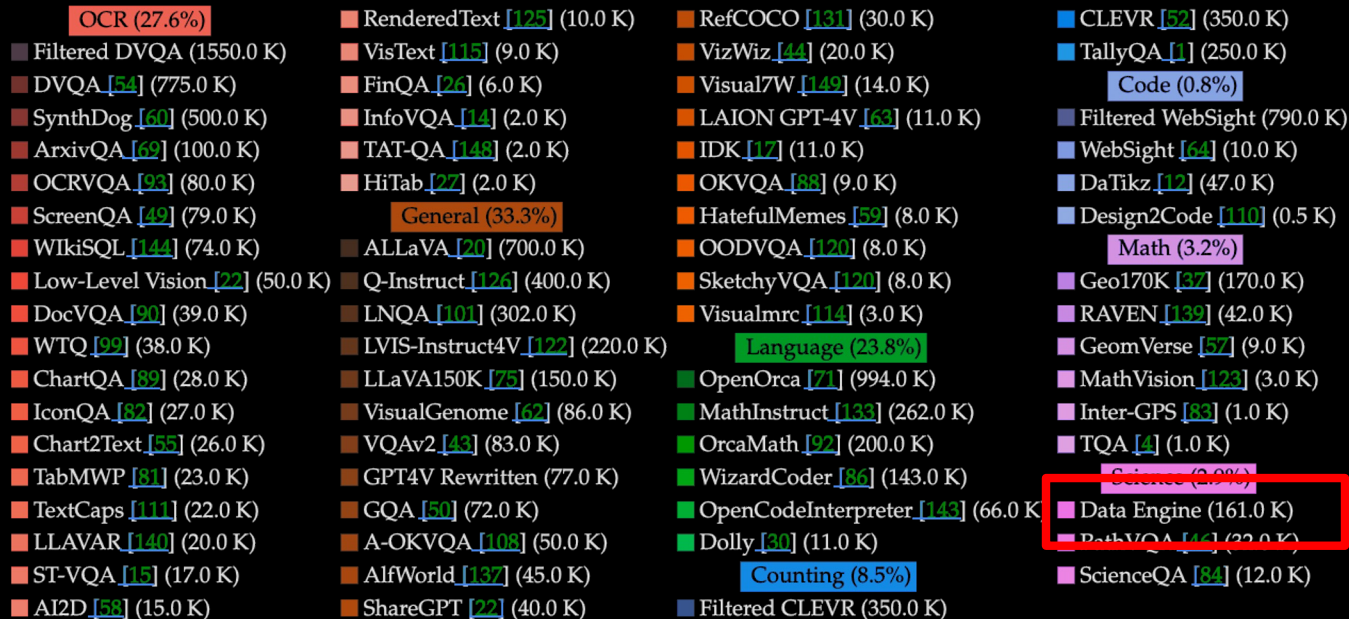


Overview

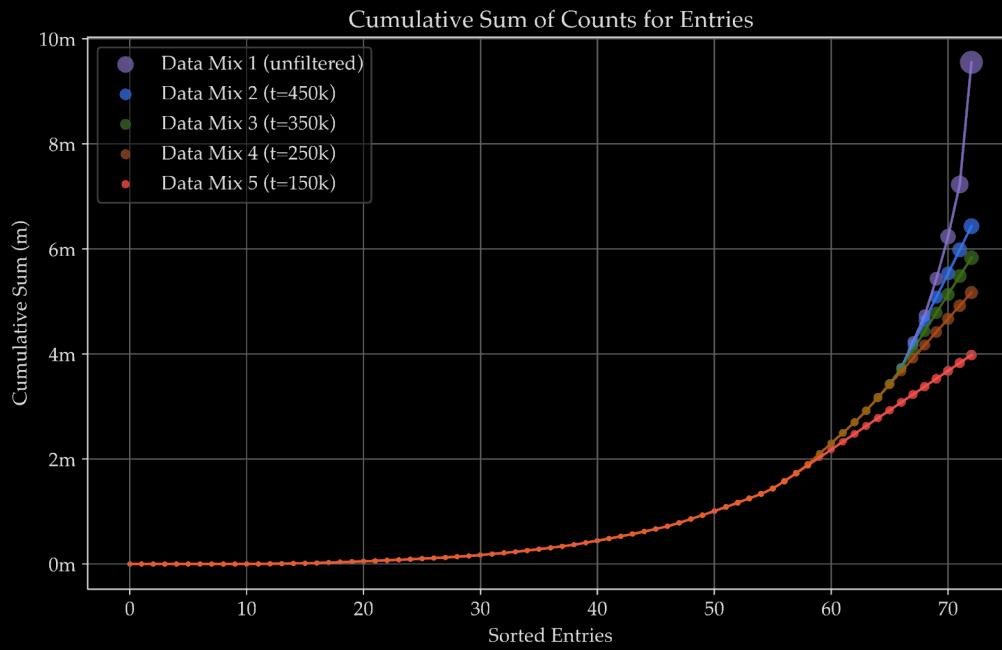


Instruction Tuning Data

Collect all-potential Instruction Tuning Data



Data Balancing (Filtering)

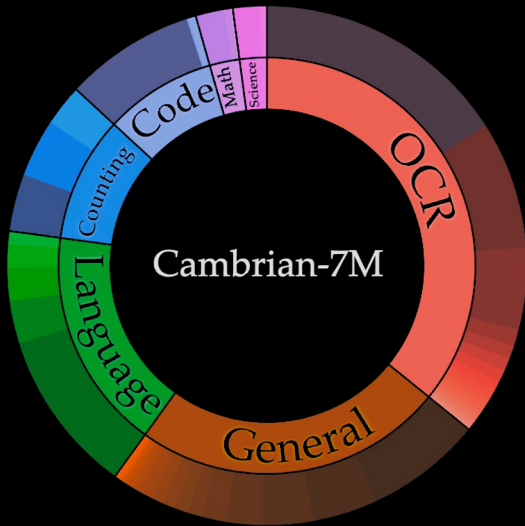


	Average	General	Knowledge	OCR & Chart	Vision-Centric
150k	53.7	68.0	51.3	45.2	50.5
250k	54.3	68.1	51.5	45.3	52.2
350k	54.3	67.4	51.4	46.0	52.3
450k	54.2	68.0	52.2	45.5	50.7

Data Mixing Ratio



Effect of Data Curation



Data quality matters

	Average	General	Knowledge	OCR & Chart	Vision-Centric
LLaVA-665K	40.7	64.7	45.2	20.8	32.0
Cambrian-10M	54.8	68.7	51.6	47.3	51.4
Cambrian-7M	55.9	69.6	52.6	47.3	54.1

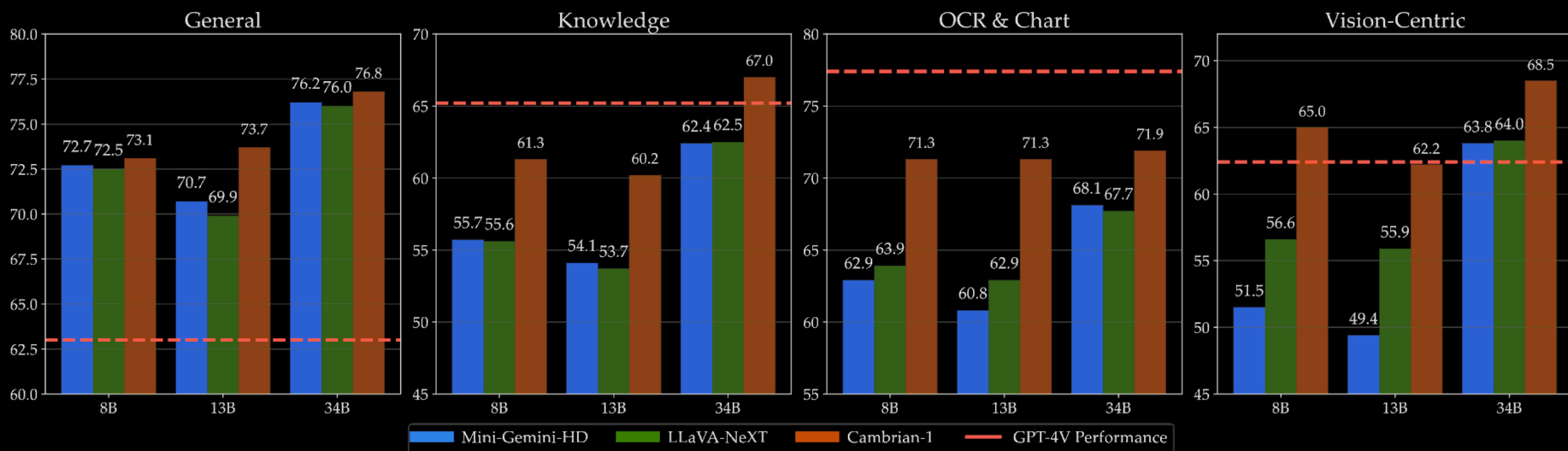


Cambrian-1 Models

“SOTA” Performance

Model		General					Knowledge					OCR & Chart					Vision-Centric				
Method	# Vis Tok.	Avg	MME ^P	MMB	SEED ^I	GQA	Avg	SQA ^I	MMMU ^V	Math Vista ^M	A12D	Avg	ChartQA	OCRBench	TextVQA	DocVQA	Avg	MMVP	RealworldQA	CV-Bench ^{2D}	CV-Bench ^{3D}
GPT-4V	UNK.	63.0	1409.4	75.8	69.1	36.8	65.2	75.7	56.8	49.9	78.2	77.4	78.5	64.5	78.0	88.4	62.4	50.0	61.4	64.3	73.8
Gemini-1.0 Pro	UNK.	-	1496.6	73.6	70.7	-	-	79.5	47.9	45.2	-	-	-	65.9	-	-	-	-	-	-	-
Gemini-1.5 Pro	UNK.	-	-	-	-	-	-	-	58.5	52.1	80.3	-	81.3	-	73.5	86.5	-	-	67.5	-	-
Grok-1.5	UNK.	-	-	-	-	-	-	-	53.6	52.8	88.3	-	76.1	-	78.1	85.6	-	-	68.7	-	-
MM-1-8B	144	-	1529.3	72.3	69.9	-	-	72.6	37.0	35.9	-	-	-	-	-	-	-	-	-	-	-
MM-1-30B	144	-	1637.6	75.1	72.1	-	-	81.0	44.7	39.4	-	-	-	-	-	-	-	-	-	-	-
<i>Base LLM: Llama-3-Ins-8B</i>																					
Mini-Gemini-HD-8B	2880	72.7	1606.0	72.7	73.2	64.5	55.7	75.1	37.3	37.0	73.5	62.9	59.1	47.7	70.2	74.6	51.5	18.7	62.1	62.2	63.0
LLaVA-NeXT-8B	2880	72.5	1603.7	72.1	72.7	65.2	55.6	72.8	41.7	36.3	71.6	63.9	69.5	49.0	64.6	72.6	56.6	38.7	60.1	62.2	65.3
Cambrian-1-8B	576	73.1	1,547.1	75.9	74.7	64.6	61.3	80.4	42.7	49.0	73.0	71.3	73.3	62.4	71.7	77.8	65.0	51.3	64.2	72.3	72.0
<i>Base LLM: Vicuna-1.5-13B</i>																					
Mini-Gemini-HD-13B	2880	70.7	1597.0	68.6	70.6	63.7	54.1	71.9	37.3	37.0	70.1	60.8	56.6	46.6	70.2	69.8	49.4	19.3	57.5	53.6	67.3
LLaVA-NeXT-13B	2880	69.9	1575.0	70.0	65.6	65.4	53.7	73.5	36.2	35.1	70.0	62.9	62.2	51.4	67.1	70.9	55.9	36.0	59.1	62.7	65.7
Cambrian-1-13B	576	73.7	1,610.4	75.7	74.4	64.3	60.2	79.3	40.0	48.0	73.6	71.3	73.8	61.9	72.8	76.8	62.2	41.3	63.0	72.5	71.8
<i>Base LLM: Hermes2-Yi-34B</i>																					
Mini-Gemini-HD-34B	2880	76.2	1659.0	80.6	75.3	65.8	62.4	77.7	48.0	43.4	80.5	68.1	67.6	51.8	74.1	78.9	63.8	37.3	67.2	71.5	79.2
LLaVA-NeXT-34B	2880	76.0	1633.2	79.3	75.9	67.1	62.5	81.8	46.7	46.5	74.9	67.7	68.7	54.5	69.5	78.1	64.0	47.3	61.0	73.0	74.8
Cambrian-1-34B	576	76.8	1689.3	81.4	75.3	65.8	67.0	85.6	49.7	53.2	79.7	71.9	75.6	60.0	76.7	75.5	68.5	52.7	67.8	74.0	79.7

“SOTA” Performance



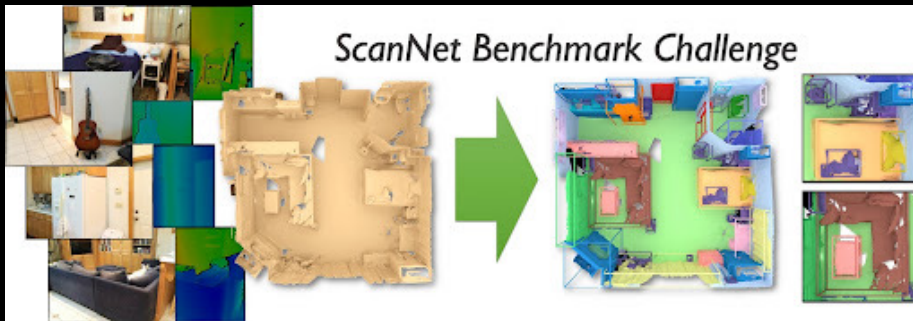
Key Question #4:

Can multimodal LLMs *think in space*?

In computer vision...

We study space, but not thinking...

We study thinking, but not in space



Video-MME

On what date did the individual in the video leave a place that Simon thought was very important to him?
A. May 31, 2022. B. June 9, 2021. C. May 9, 2021. D. June 31, 2021.

The date of Day 1 is May 31, 2021. [In Frames] Simon is the camera man. [In Frames] Yosemite National Park did mean a lot more to Simon. [In Subs/Audio] Depart Yosemite on Day 10. [In Frames]

01:10 02:22 Full Video Link: youtu.be/VPm0BRCp9A 04:12 27:52 31:16

Video-MME

How did the man wearing a bandage and holding an envelope, who appeared in the latter part of this video, sustain his injury?
A. One of his hands was hit by a firework while he was setting it off.
B. His arms got injured while he was attempting to put out the fire at a burning house.
C. His hands were injured from falling down to the ground while he was chasing Wayne's motorcycle.
D. One of his arms was dragged down by a dog lured with food by Wayne, while he was insulting Wayne's father.

Dragged down by a dog. [Option D] The man wearing a bandage and holding an envelope. Chasing Wayne's motorcycle. [Option C] A burning house. [Option B] Hit by a firework. [Option A]

03:35 Full Video Link: youtu.be/p84Og1Ap_IM 27:30 27:58 28:10 30:35

Thinking in Space: How Multimodal LLMs See, Remember and Recall Spaces

See a video of an apartment



See a laboratory

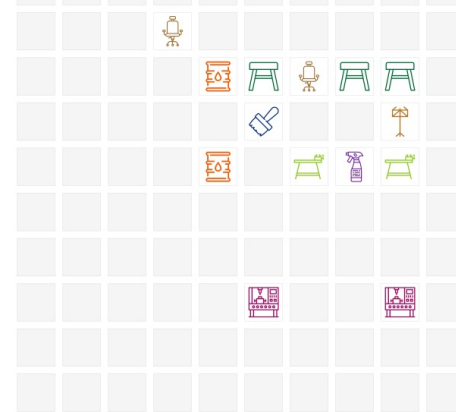
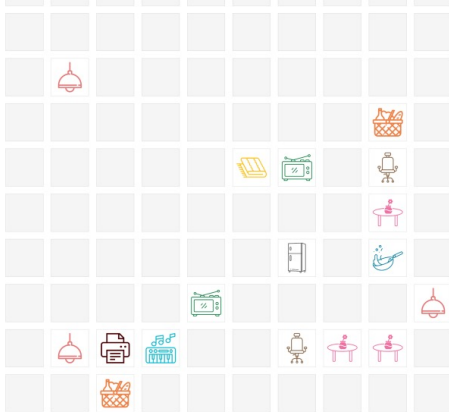


See a factory



Remember?

Multimodal LLM's "cognitive map" of the space



Recall?

What is the distance between the **keyboard** and the **TV**, in meters?

How many **cabinet(s)** are in this room?

What is the height of the **stool**, in cm?

With Jihan Yang, Shusheng Yang, Anjali Gupta, Rilyn Han, and Fei-Fei Li

Apartment



Lab





Object Count

How many chairs are there in this room?

Answer: 4

Relative Distance

Measuring from the closest point of each object, which of these objects (refrigerator, sofa, ceiling light, cutting board) is the closest to the printer?

A. refrigerator B. sofa C. ceiling light D. cutting board

Appearance Order

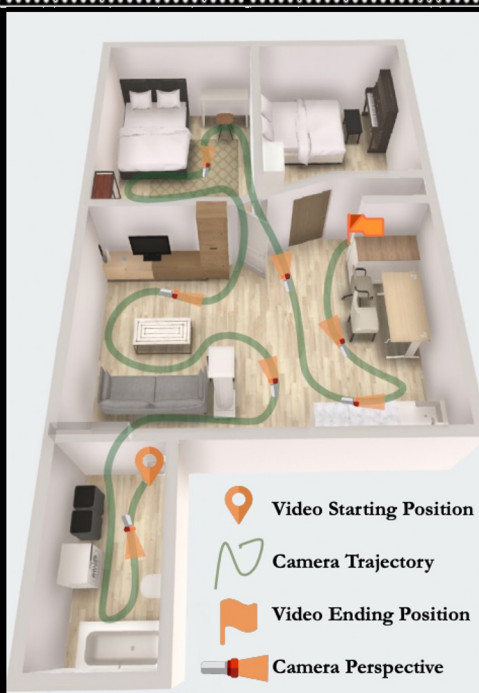
What will be the first-time appearance order of the following categories in the video: basket, printer, refrigerator, kettle?

A. kettle, basket, printer, refrigerator
 B. refrigerator, printer, basket, kettle
 C. basket, printer, refrigerator, kettle
 D. basket, refrigerator, kettle, printer

Relative Direction

If I am standing by the refrigerator and facing the sofa, is the kettle to my left, right, or back?

A. left B. right C. back



Object Size

What is the length of the longest dimension (length, width, or height) of the refrigerator in centimeters?

Answer: 119

Absolute Distance

Measuring from the closest point of each object, what is the distance between the bed and the sofa in meters?

Answer: 3.2

Room Size

What is the size of this room (in square meters)? If multiple rooms are shown, estimate the size of the combined space.

Answer: 57.6

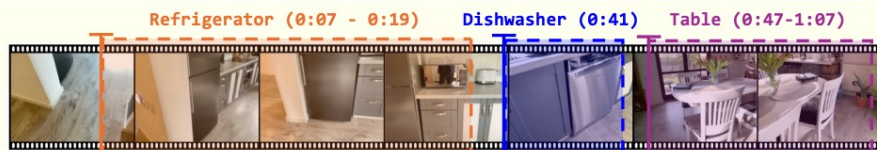
Route Plan

You are a robot beginning at the toilet and facing the washer. Navigate to the pan. Fill in this route: 1. Go forward until the washing machine 2. [?] 3. Go forward until the sofa 4. [?] 5. Go forward until the pan.

A. Turn Left, Turn Left B. Turn Left, Turn Right
 C. Turn Back, Turn Right D. Turn Right, Turn Right

Annotation *facilitated* by the 3D reconstruction datasets

Error Analysis through LLM self-explanation



Relative Direction

👍 MLLM: **Front-left** GT: **Front-left**

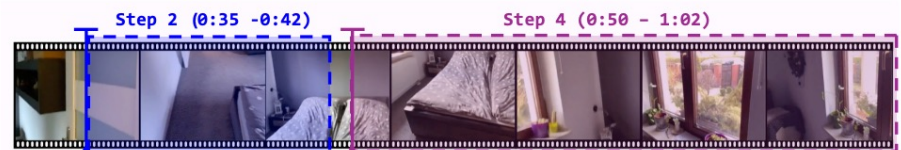
Q: If I am standing by the refrigerator and facing the table, is the dishwasher to my front-left, front-right, back-left, or back-right? The directions refer to the quadrants of a Cartesian plane (if I am standing at the origin and facing also

MLLM: Her you the at at 00:

Spatial reasoning is the primary bottleneck for MLLM performance on VSI-Bench.

at [00:00:41]. 3. Visualize the quadrants: Since you're facing the table, the area in front of you is the positive y-axis, and the area behind you is the negative y-axis. Your right is the positive x-axis, and your left is the negative x-axis. 4. Determine the dishwasher's quadrant: The dishwasher is located in the quadrant that is in front of you and to your left....

👍 Human-like Reasoning 👍 Video Understanding 👍 Internal World Model



Route Plan

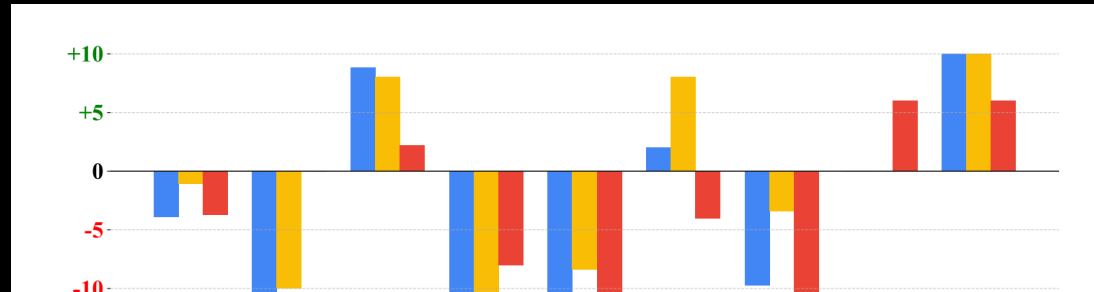
👎 MLLM: **Turn Right, Turn Left** GT: **Turn Left, Turn Right**

Q: You are beginning at the door and facing the floor. You want to navigate to the window. You will perform the following actions : 1. Go forward until the edge of bed 2. [please fill in] 3. Go forward until the wall 4. [please fill in] 5. Go

at the door of the bed window is left would back would send you back towards the door. 3. Go forward until the wall: You proceed forward until you reach the wall (0:42-0:50). 4. Turn Left: The window is to your left (0:50-1:02). Turning right would lead you along the wall. 5. Go forward until the window: You move forward and reach the window (1:02). Therefore, the correct sequence of turns is **Turn Right, Turn Left...**

👎 Egocentric-Allocentric Transformation 👎 Relational Reasoning

Language Prompting Techniques

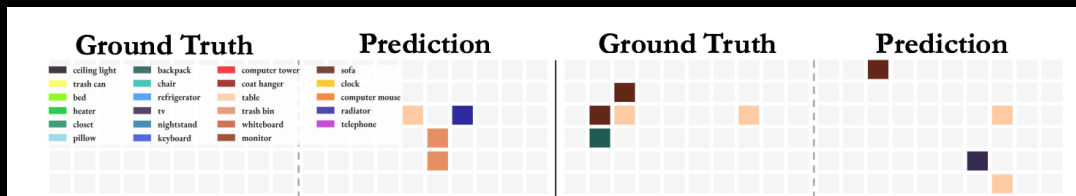


Linguistic prompting techniques, although effective in language reasoning tasks, are primarily harmful for spatial reasoning.

Tree-of-Thought compared to the baseline. All three prevailing prompting techniques fail on average on our benchmark, and, in some cases, task performance becomes *much worse* after applying them. This implies that VSI-Bench cannot be solved by solely improving linguistic capabilities.

Note: works on general video understanding tasks

A *Local* “World Model”



When remembering spaces, a MLLM forms a series of local world models in its mind from a given video, rather than a unified global model.

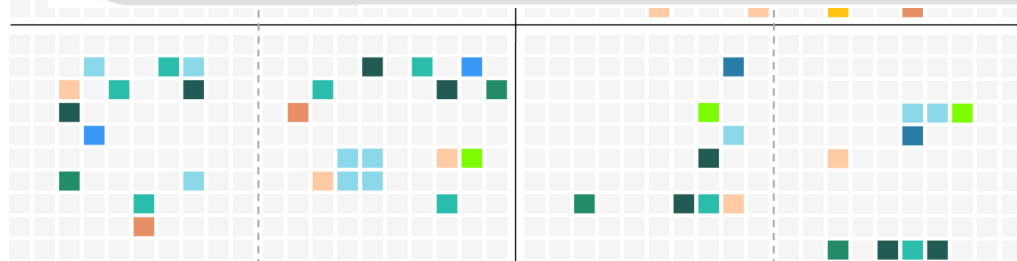
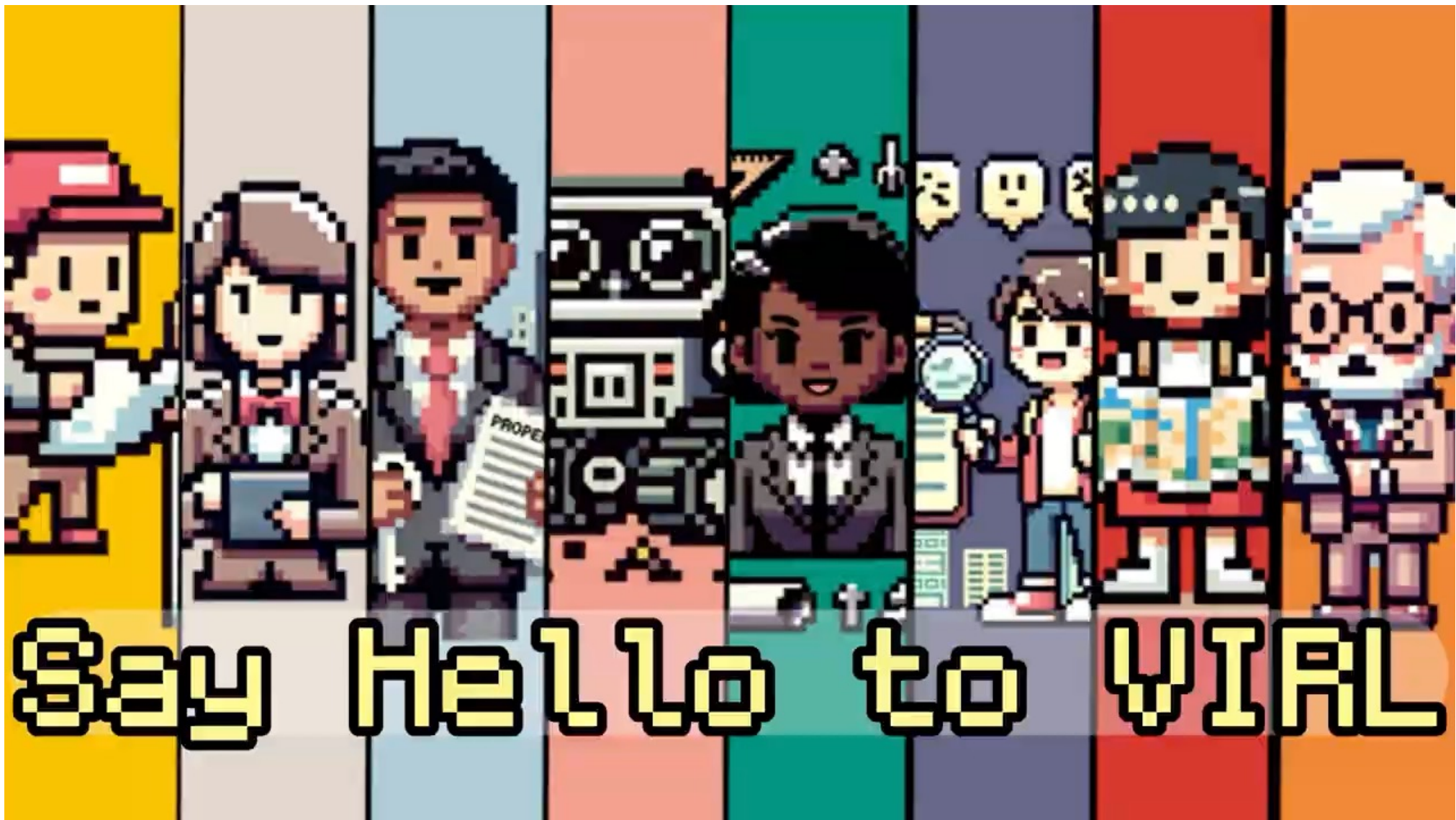


Figure 10. Visualization of cognitive maps from MLLM and GT.

Close Far
 Figure 11. **Locality of the MLLM's predicted cognitive maps.** The MLLM's map-distance accuracy decreases dramatically with increasing object distance.

Key Question #5:

How to ground multimodal agents in real life?



Say Hello to VIAL



Thank You

*Improved vision is not just about seeing farther,
but about understanding more deeply.*