

# CSCSE 638 Natural Language Processing Foundation and Techniques

## Lecture 18: Vision-Language Models

Kuan-Hao Huang

Spring 2026



# Invited Talk

- **Date:** 4/6 online @ Zoom
  - <https://tamu.zoom.us/my/khhuang?pwd=oAdWOKVOCGPAPqDbJnVtktdW2AE6nb.1>
  - **Talk:** Improving Personalization and Consistency of Large Foundation Models
  - **Speaker:** Jindong Wang, Assistant Professor, The College of William & Mary
- **Date:** 4/15 online @ Zoom
  - <https://tamu.zoom.us/my/khhuang?pwd=oAdWOKVOCGPAPqDbJnVtktdW2AE6nb.1>
  - **Talk:** Less Is More: Why Compression May Be the Missing Incentive for LLM Generalization
  - **Speaker:** Ben Zhou, Assistant Professor, Arizona State University

# Vision + Language

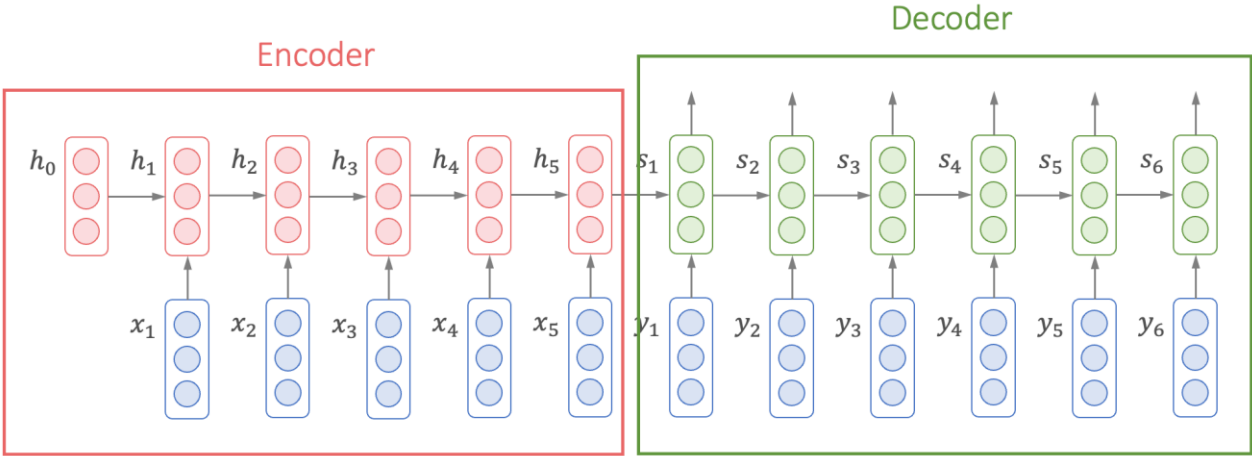
- Image captioning

<p>A young boy is playing basketball.</p> 	<p>Two dogs play in the grass.</p> 	<p>A dog swims in the water.</p> 	<p>A little girl in a pink shirt is swinging.</p> 
<p>A group of people walking down a street.</p> 	<p>A group of women dressed in formal attire.</p> 	<p>Two children play in the water.</p> 	<p>A dog jumps over a hurdle.</p> 

# Image Captioning with Encoder-Decoder Models



Replace the text encoder as an image encoder



Encoder-Decoder Model

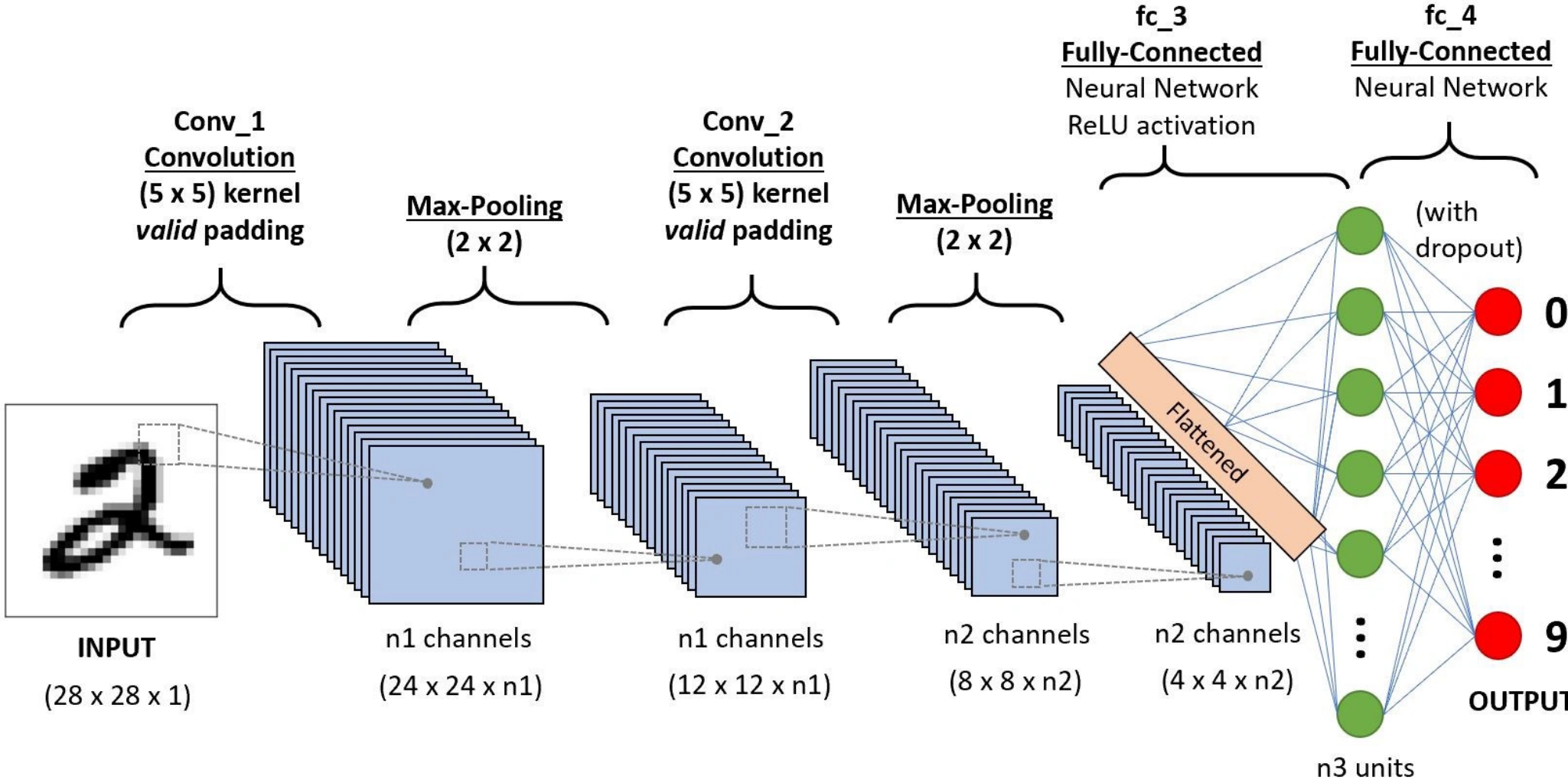
# Recap: Convolutional Neural Network (For Text)

Learnable Weight (Filter)  
Filter Size = 3

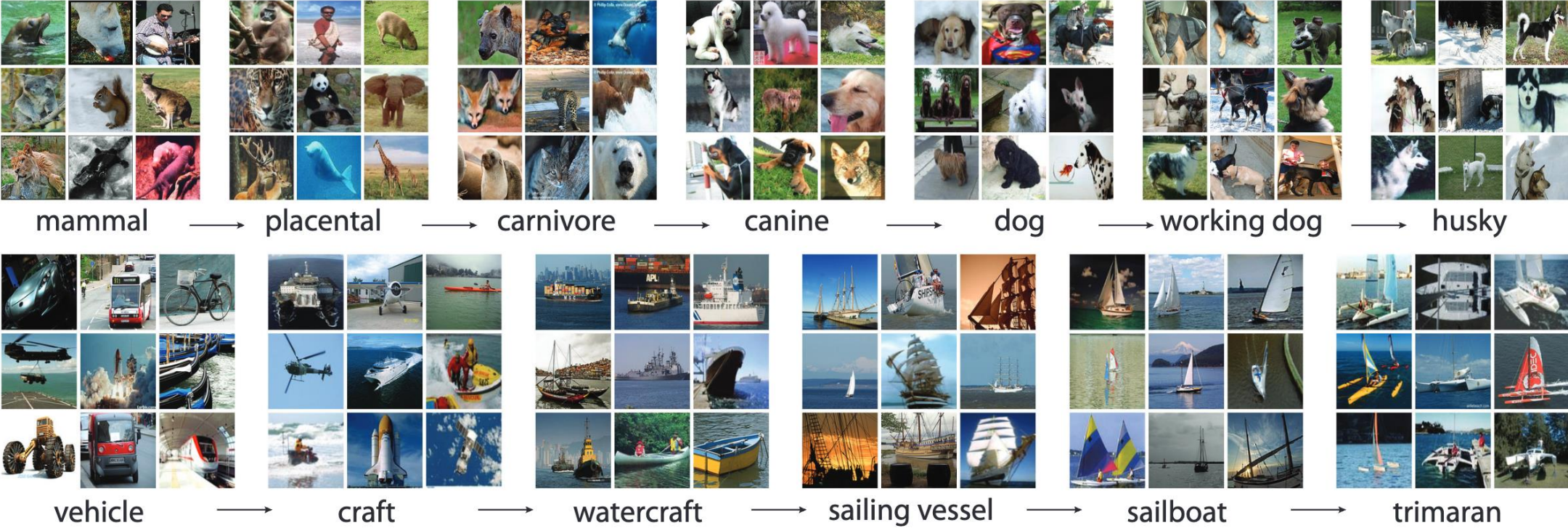
$$W = \begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \\ \dots & \dots & \dots \\ W_{4,1} & W_{4,2} & W_{4,3} \end{bmatrix}$$



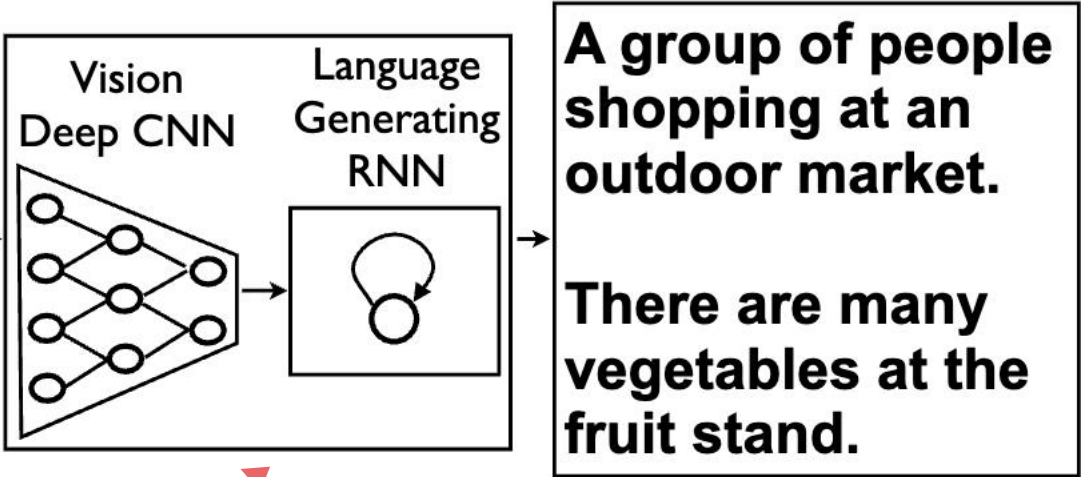
# Convolutional Neural Network (For Image)



# Pre-Trained CNN with ImageNet

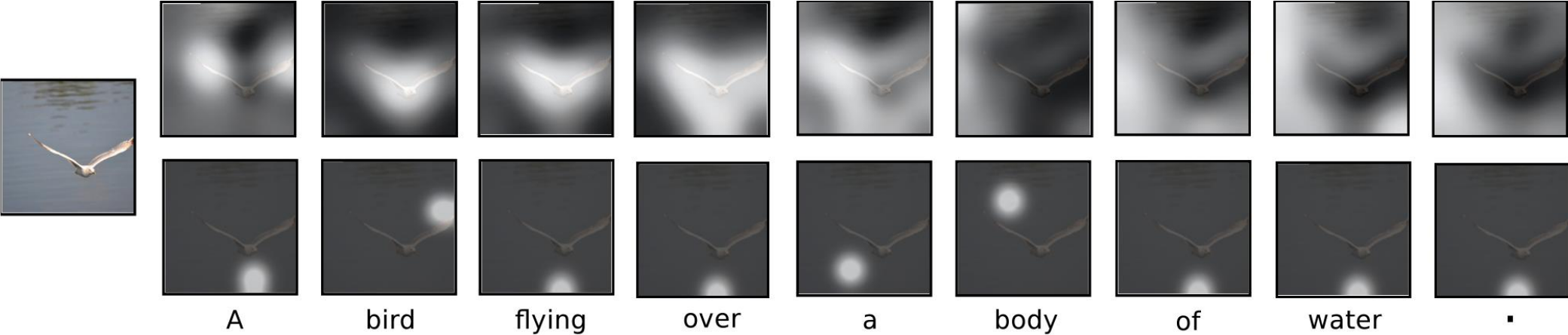


# Encoder-Decoder: CNN-RNN



Text embedding space and image embedding space can be aligned!

# CNN + Attention LSTM



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.

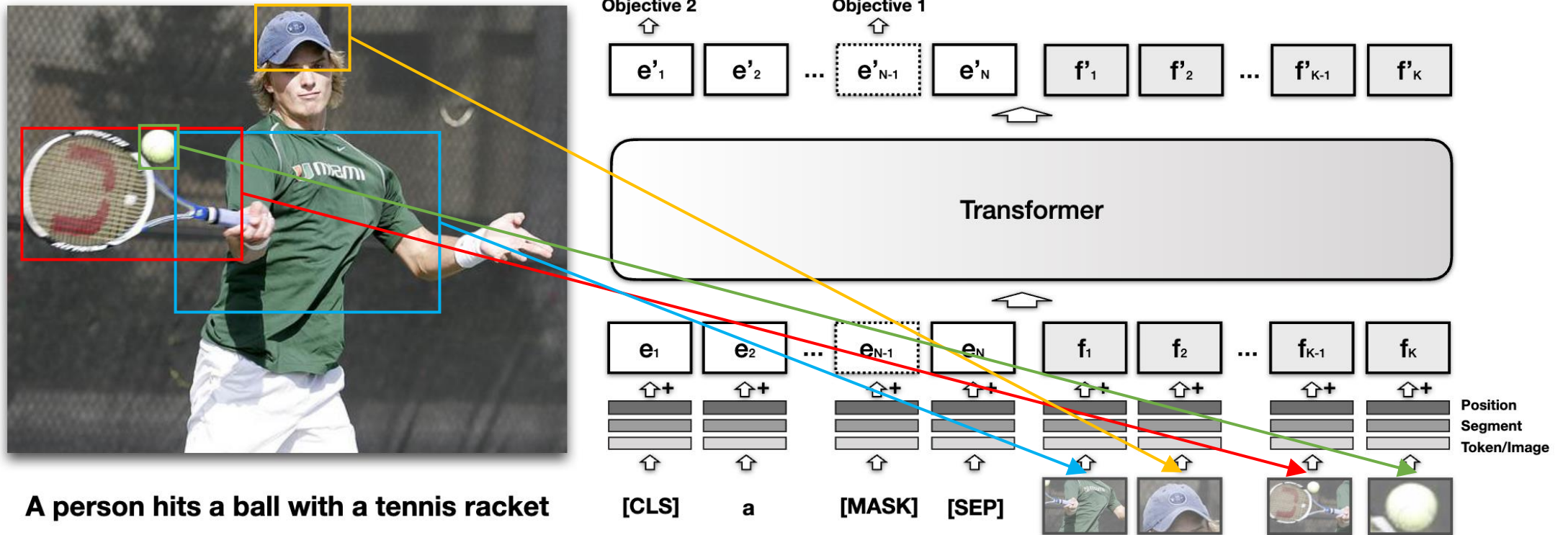


A group of people sitting on a boat in the water.



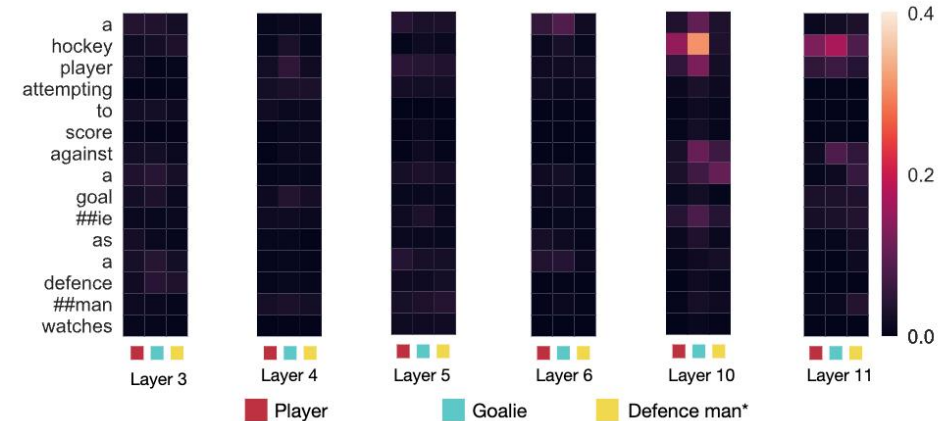
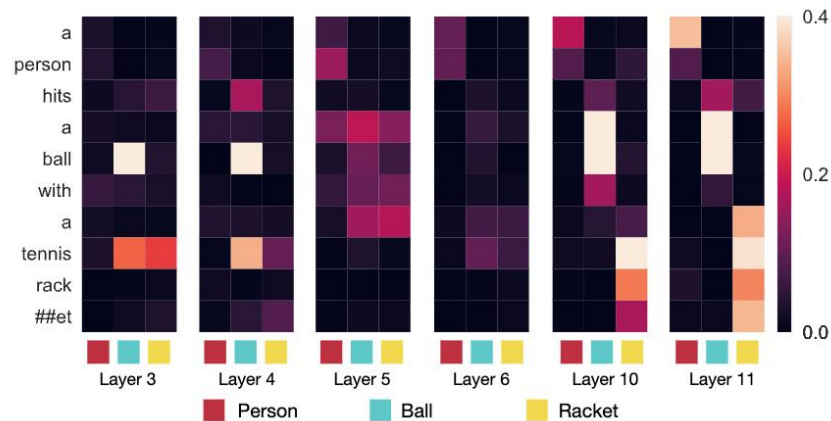
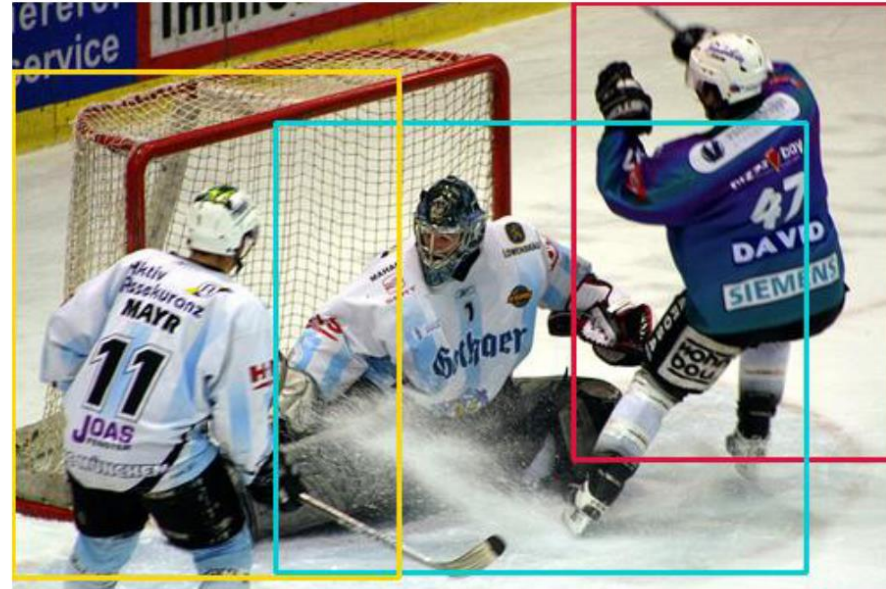
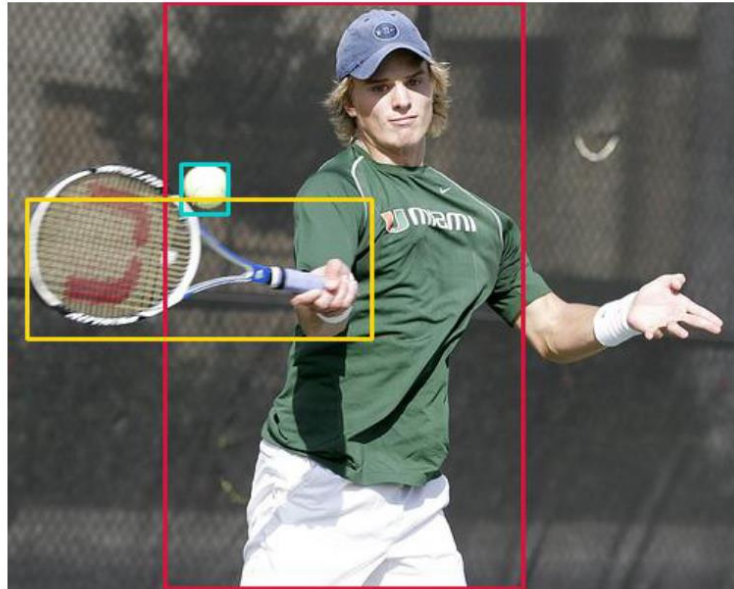
A giraffe standing in a forest with trees in the background.

# Joint Visual and Textual Embeddings: VisualBERT



Require an object detection model

# Joint Visual and Textual Embeddings: VisualBERT



# Visual Question Answering

Who is wearing glasses?

man



woman



Where is the child sitting?

fridge



arms

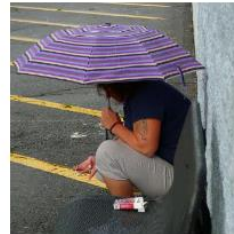


Is the umbrella upside down?

yes



no



How many children are in the bed?

2

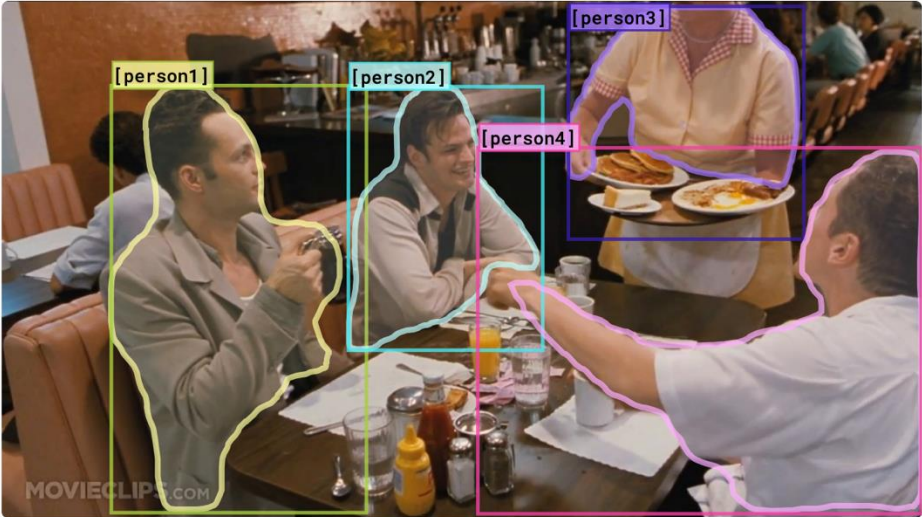


1



Model	Test-Dev	Test-Std
Pythia v0.1 (Jiang et al., 2018)	68.49	-
Pythia v0.3 (Singh et al., 2019)	68.71	-
VisualBERT w/o Early Fusion	68.18	-
VisualBERT w/o COCO Pre-training	70.18	-
VisualBERT	70.80	71.00
Pythia v0.1 + VG + Other Data Augmentation (Jiang et al., 2018)	70.01	70.24
MCAN + VG (Yu et al., 2019b)	70.63	70.90
MCAN + VG + Multiple Detectors (Yu et al., 2019b)	72.55	-
MCAN + VG + Multiple Detectors + BERT (Yu et al., 2019b)	72.80	-
MCAN + VG + Multiple Detectors + BERT + Ensemble (Yu et al., 2019b)	75.00	75.23

# Visual Commonsense Reasoning



hide all show all [person1] [person2] [person3] [person4]

more objects »

Why is [person4] pointing at [person1]?

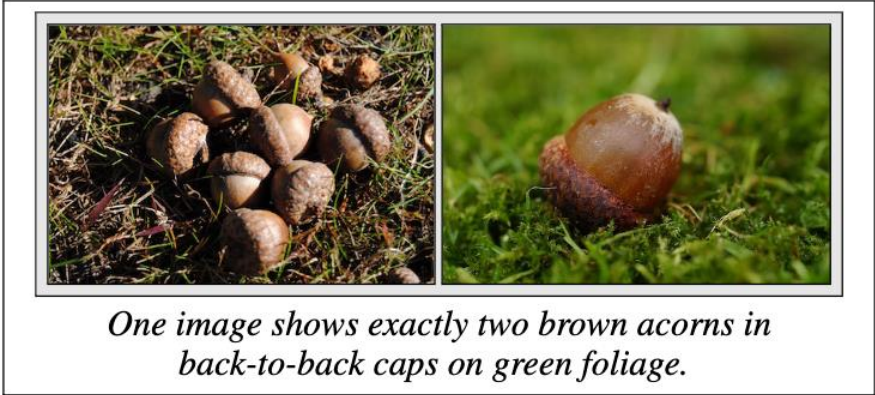
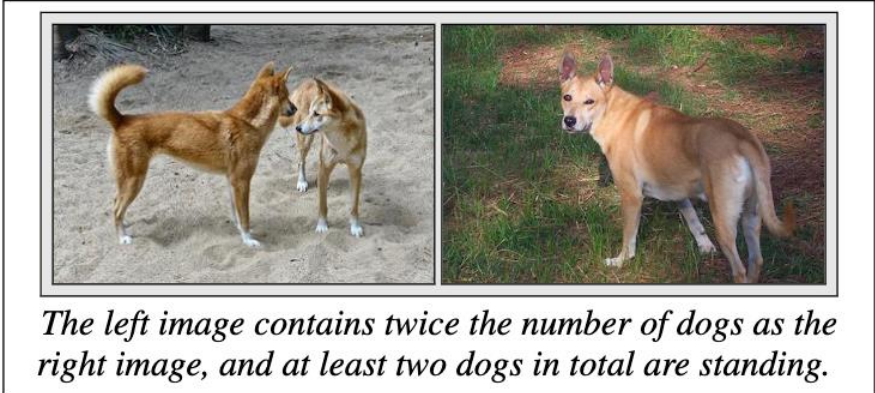
- a) He is telling [person3] that [person1] ordered the pancakes.
- b) He just told a joke.
- c) He is feeling accusatory towards [person1].
- d) He is giving [person1] directions.

Rationale: I think so because...

- a) [person1] has the pancakes in front of him.
- b) [person4] is taking everyone's order and asked for clarification.
- c) [person3] is looking at the pancakes both she and [person2] are smiling slightly.
- d) [person3] is delivering food to the table, and she might not know whose order is whose.

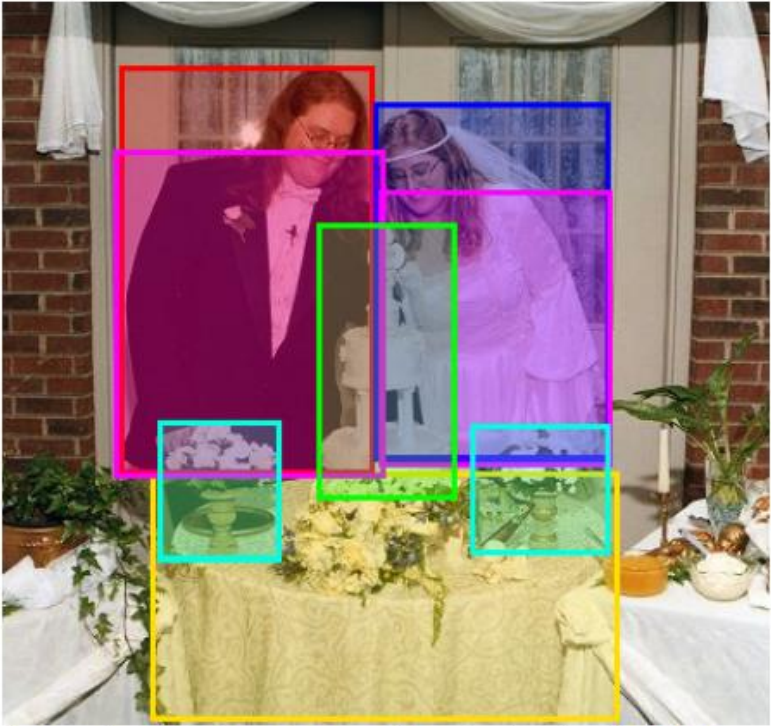
Model	Q → A		QA → R		Q → AR	
	Dev	Test	Dev	Test	Dev	Test
R2C (Zellers et al., 2019)	63.8	65.1	67.2	67.3	43.1	44.0
B2T2 (Leaderboard; Unpublished)	-	72.6	-	75.7	-	55.0
VisualBERT w/o Early Fusion	70.1	-	71.9	-	50.6	-
VisualBERT w/o COCO Pre-training	67.9	-	69.5	-	47.9	-
VisualBERT	70.8	71.6	73.2	73.2	52.2	52.4

# Natural Language Visual Reasoning



Model	Dev	Test-P	Test-U	Test-U (Cons)
MaxEnt (Suhr et al., 2019)	54.1	54.8	53.5	12.0
VisualBERT w/o Early Fusion	64.6	-	-	-
VisualBERT w/o COCO Pre-training	63.5	-	-	-
VisualBERT	67.4	67.0	67.3	26.9

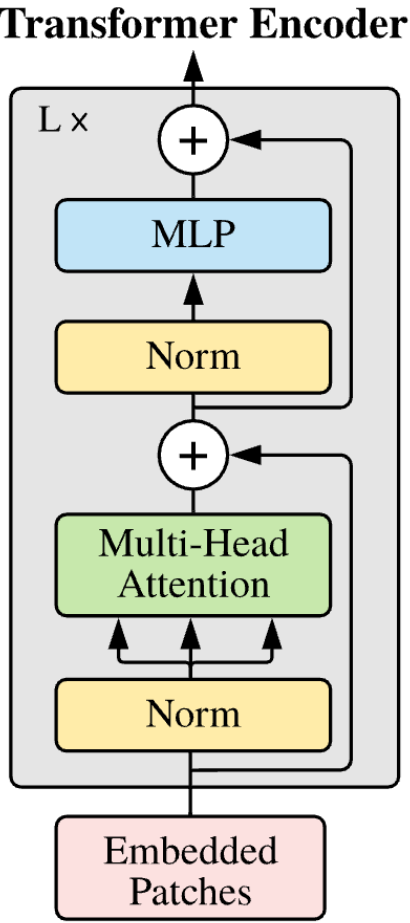
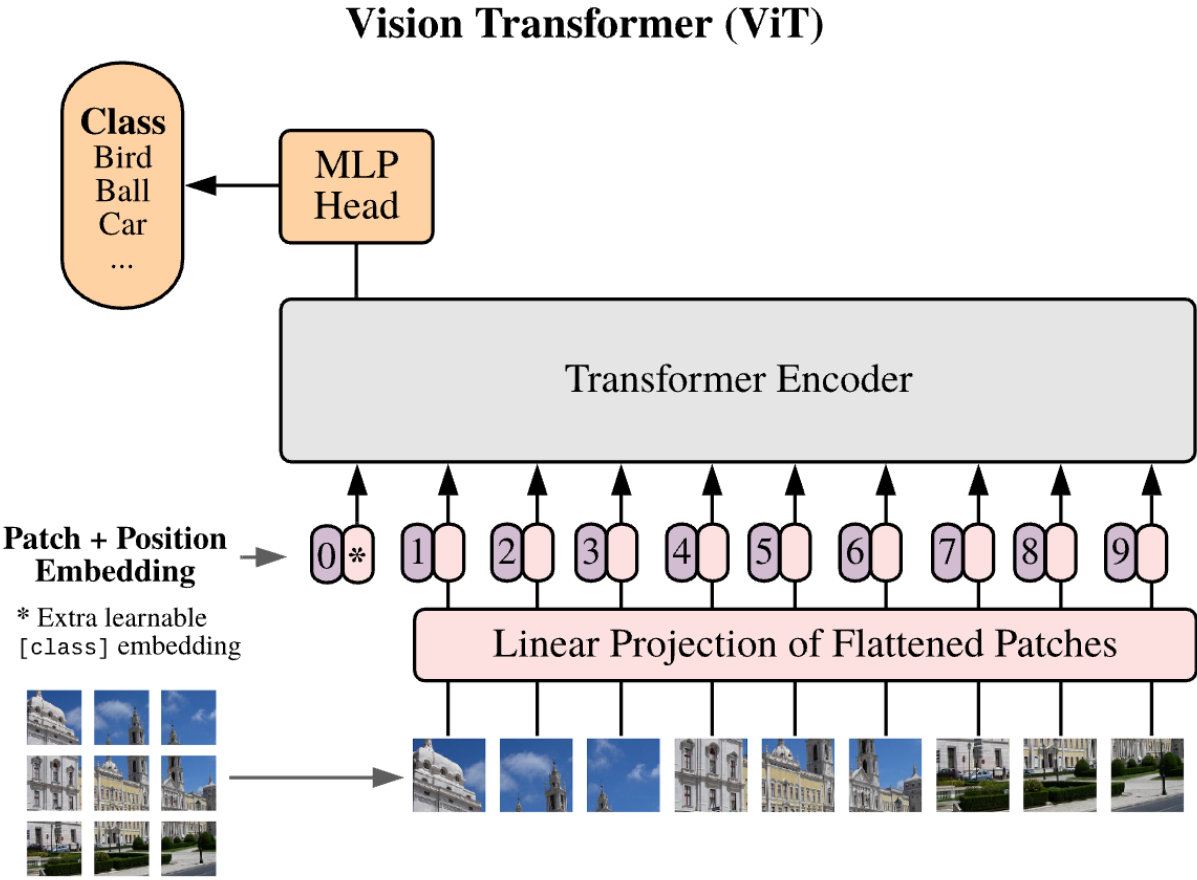
# Language Grounding



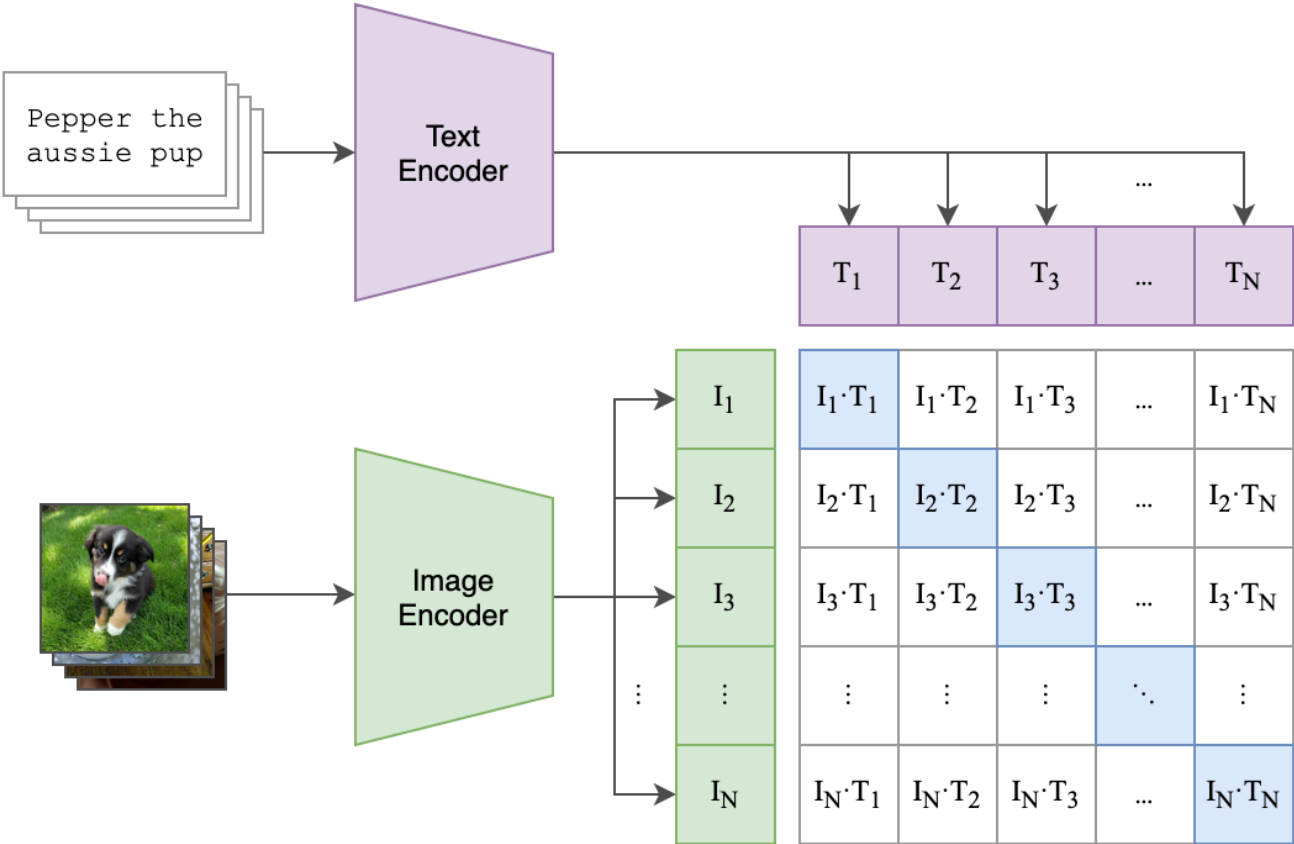
A couple in **their wedding attire** stand behind **a table** with **a wedding cake** and **flowers**.  
**A bride** and **groom** are standing in front of **their wedding cake** at their reception.  
**A bride** and **groom** smile as **they** view **their wedding cake** at a reception.  
**A couple** stands behind **their wedding cake**.  
**Man** and **woman** cutting **wedding cake**.

Model	R@1		R@5		R@10		Upper Bound	
	Dev	Test	Dev	Test	Dev	Test	Dev	Test
BAN (Kim et al., 2018)	-	69.69	-	84.22	-	86.35	86.97	87.45
VisualBERT w/o Early Fusion	70.33	-	84.53	-	86.39	-	-	-
VisualBERT w/o COCO Pre-training	68.07	-	83.98	-	86.24	-	86.97	87.45
VisualBERT	70.40	71.33	84.49	84.98	86.31	86.51	-	-

# Vision Transformer



# CLIP: Contrastive Language-Image Pre-Training

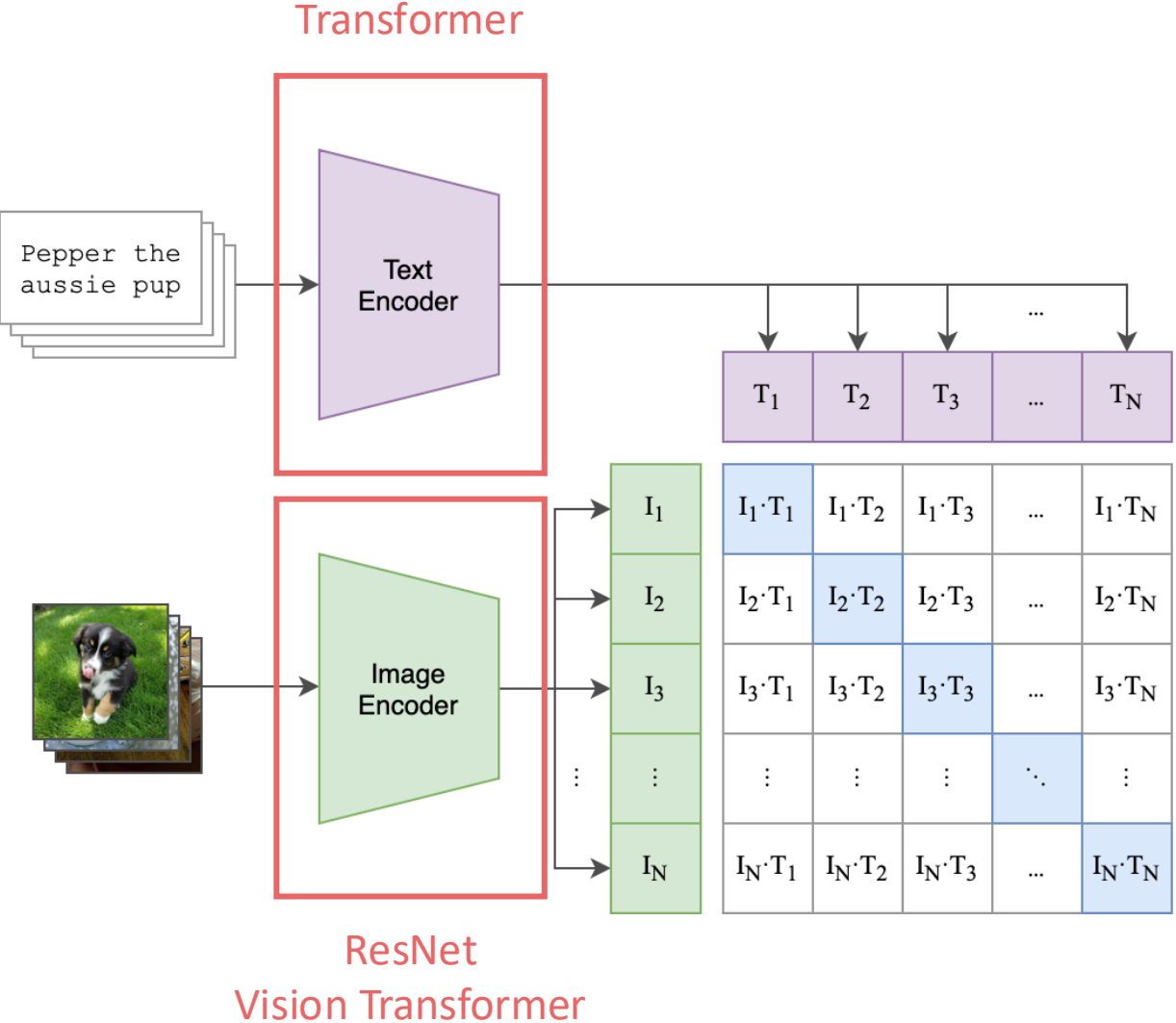


# Training with Image-Caption Pairs

Cosine similarity between text and image features



# Training Details



```
# extract feature representations of each modality
I_f = image_encoder(I) #[n, d_i]
T_f = text_encoder(T)  #[n, d_t]
```

```
# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)
```

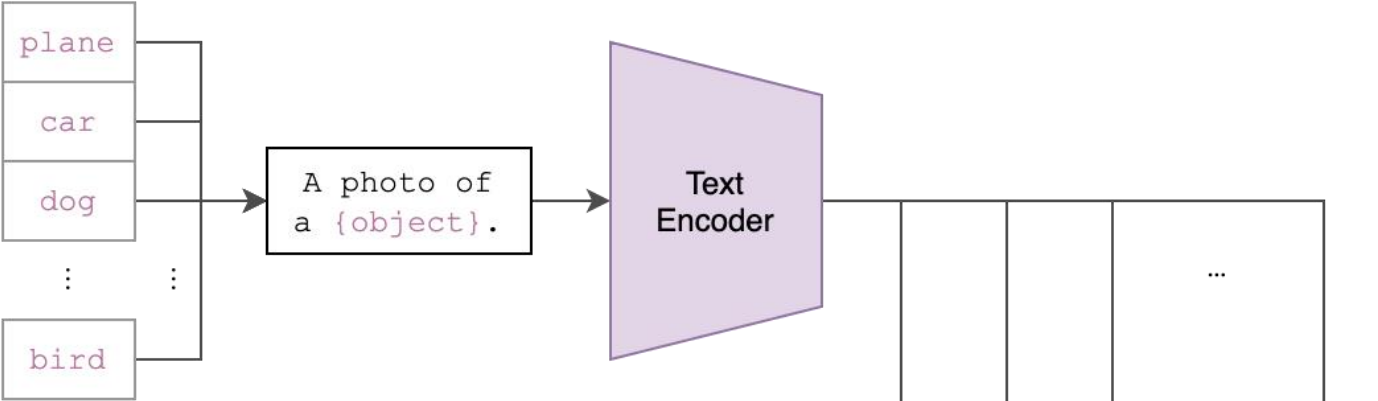
```
# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)
```

```
# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss   = (loss_i + loss_t)/2
```

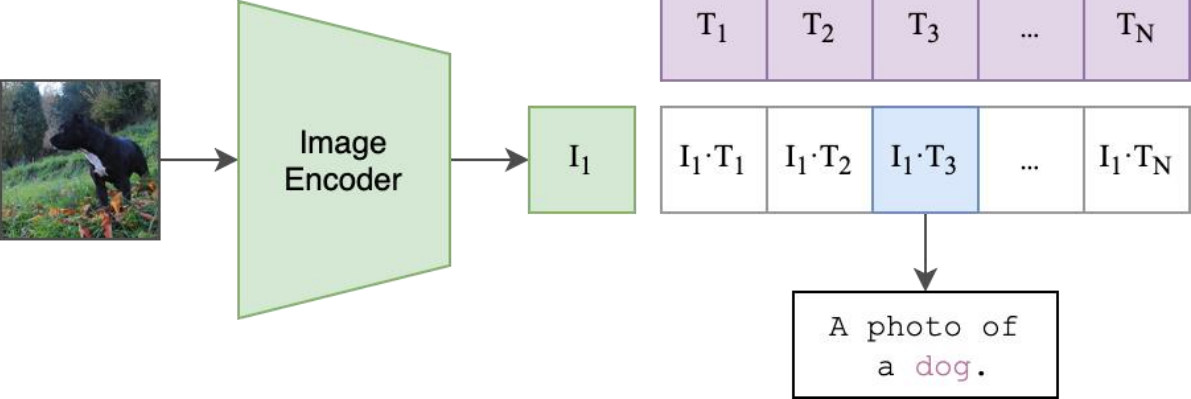
1	0	0	...	0
0	1	0	...	0
0	0	1	...	0
...	...	...	...	0
0	0	0	0	1

# Zero-Shot Prediction

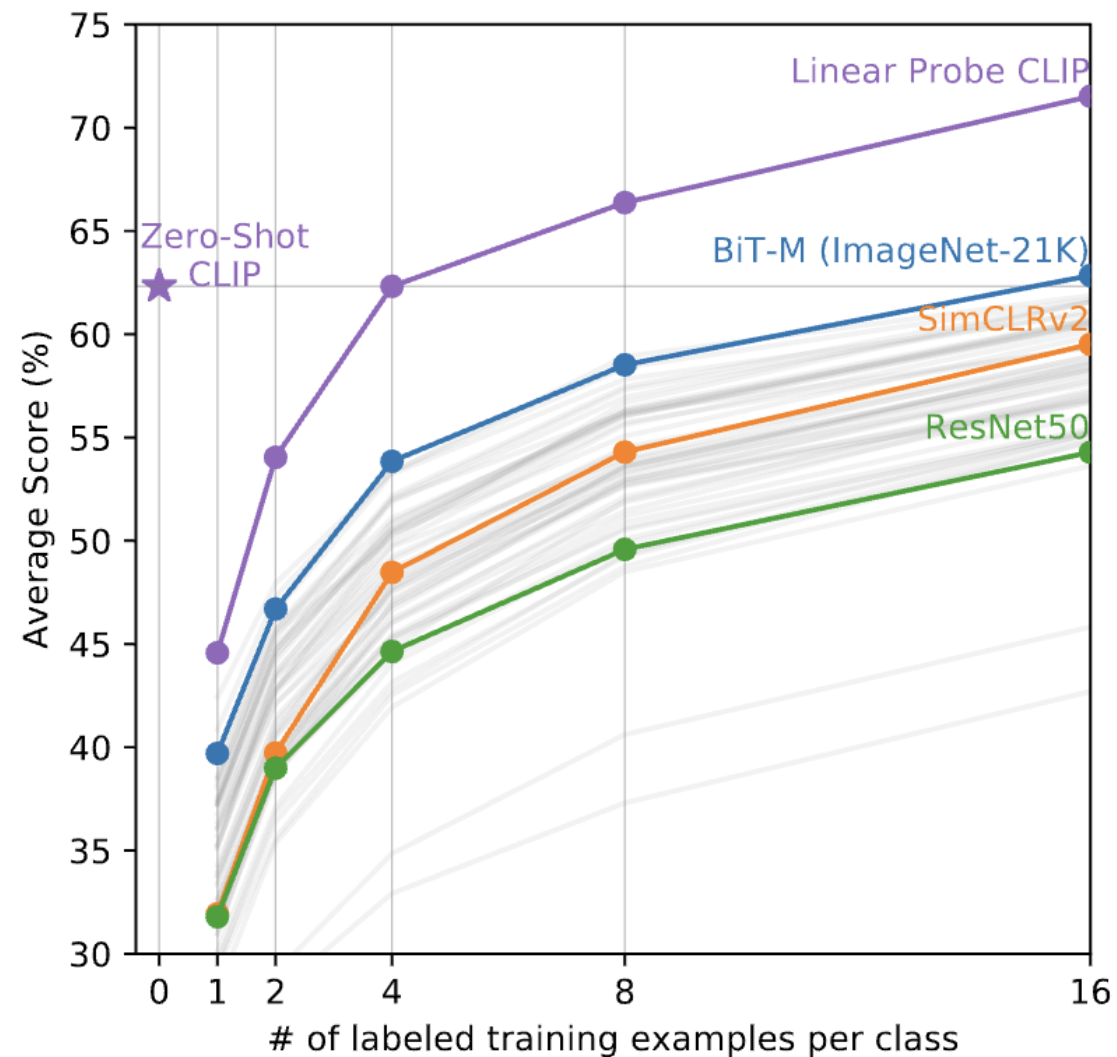
(2) Create dataset classifier from label text



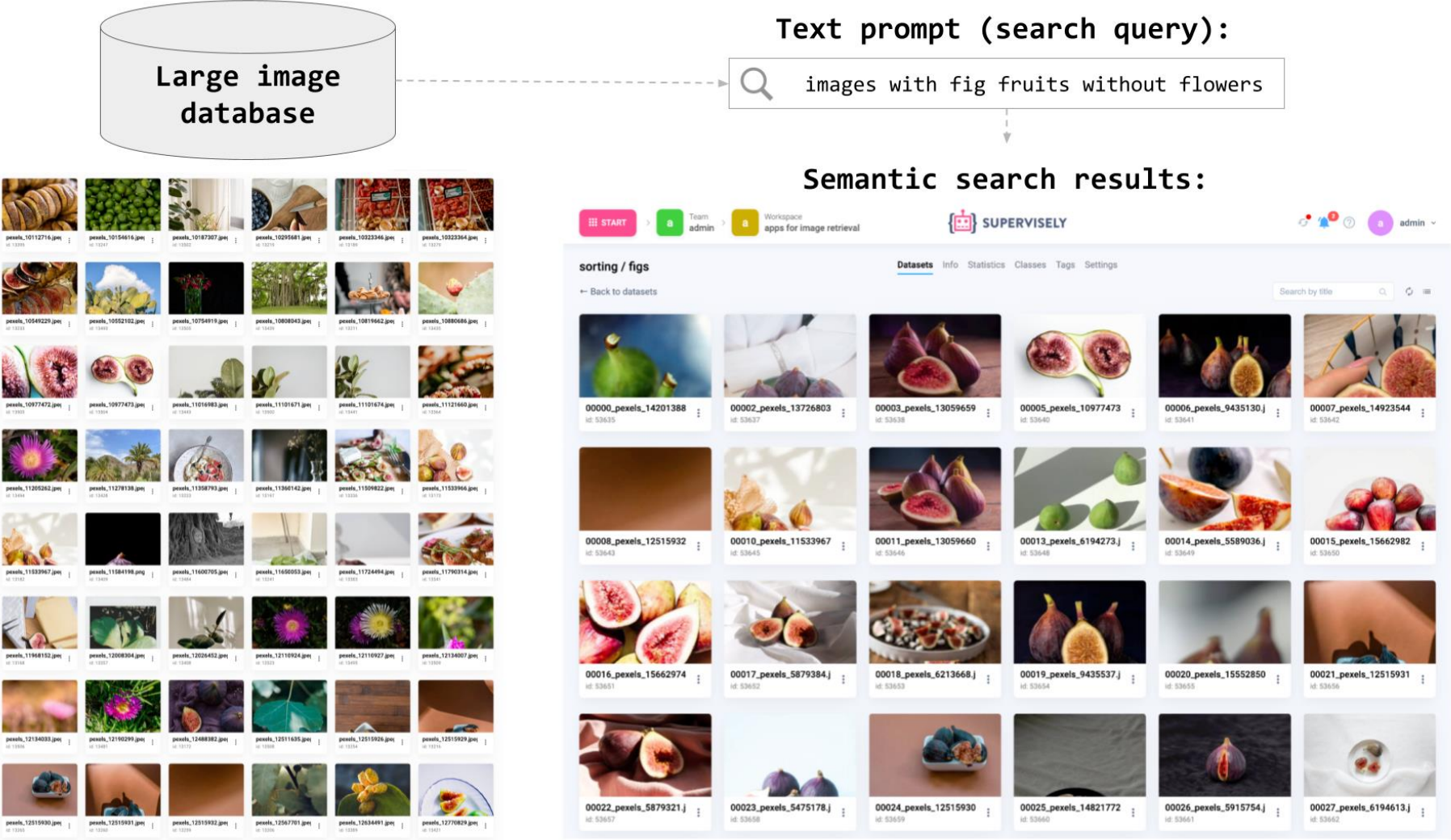
(3) Use for zero-shot prediction



# Zero-Shot CLIP vs. Few-shot Linear Probes



# Image Retrieval with Text Query



Large image database

Text prompt (search query):

images with fig fruits without flowers

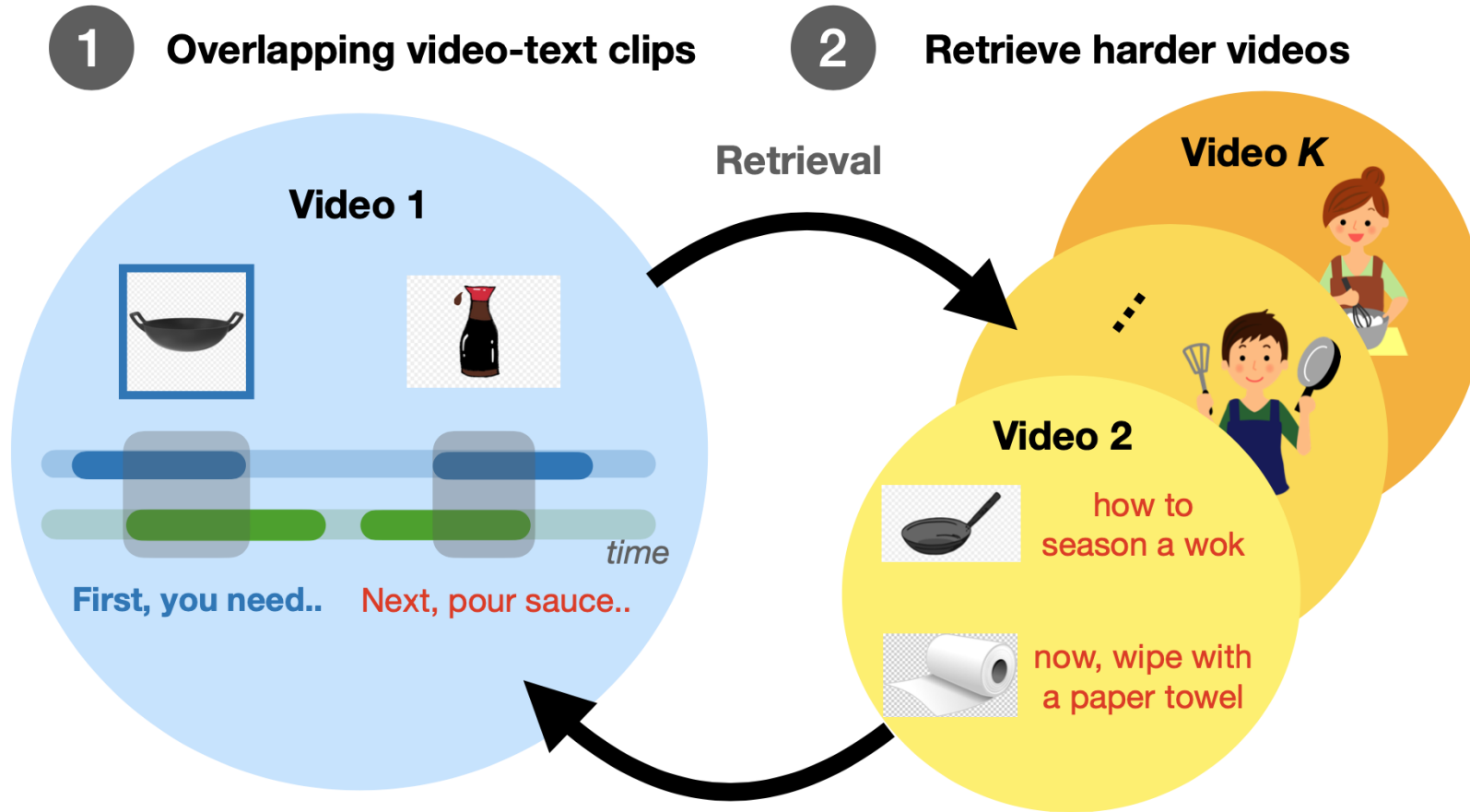
Semantic search results:

START Team admin Workspace apps for image retrieval SUPERVISELY admin

sorting / figs Datasets Info Statistics Classes Tags Settings

Back to datasets Search by title

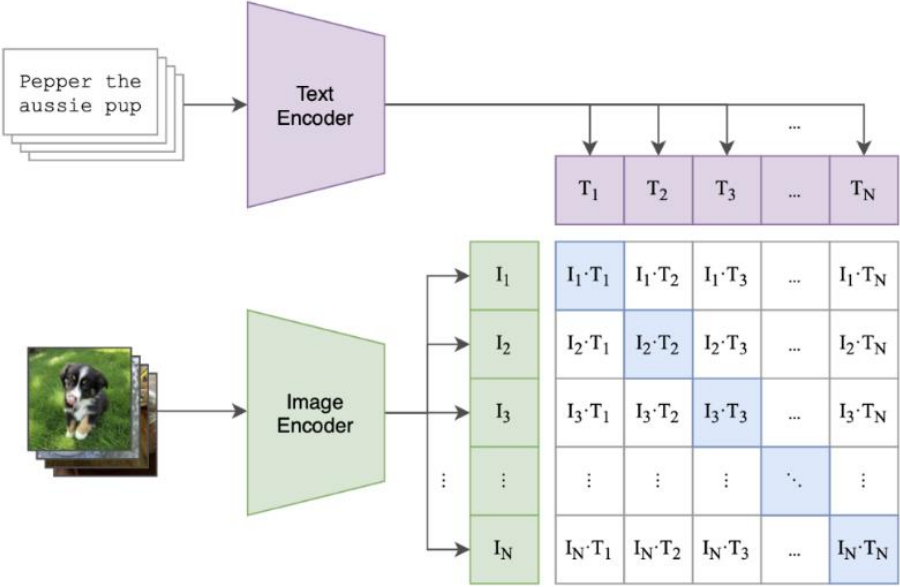

# VideoCLIP



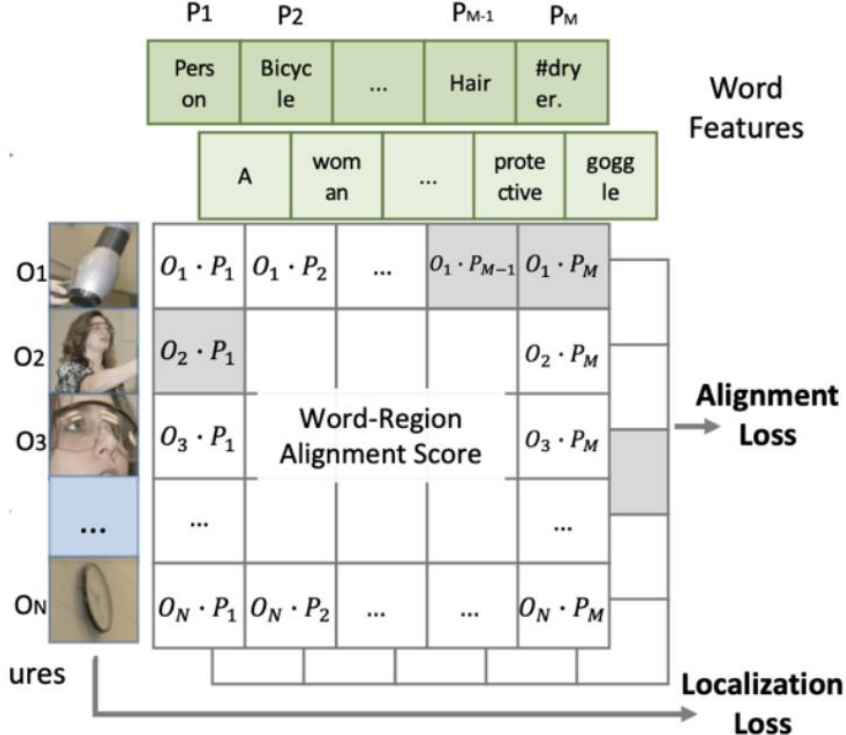
VideoCLIP: Contrastive learning with **hard-retrieved negatives** and **overlapping positives** for video-text pre-training.

# GLIP: Grounded Language-Image Pre-training

CLIP: capture information for whole image



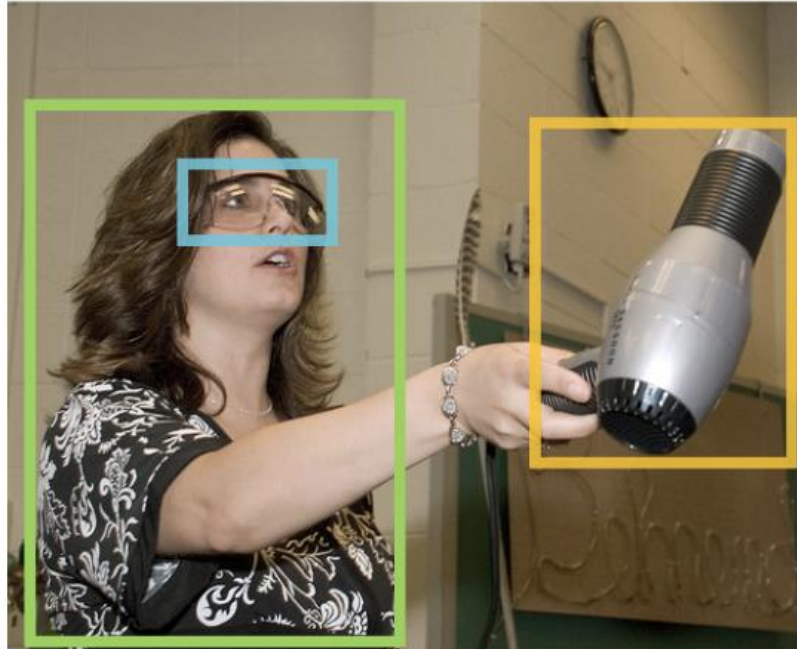
GLIP: capture information more for objects/entities



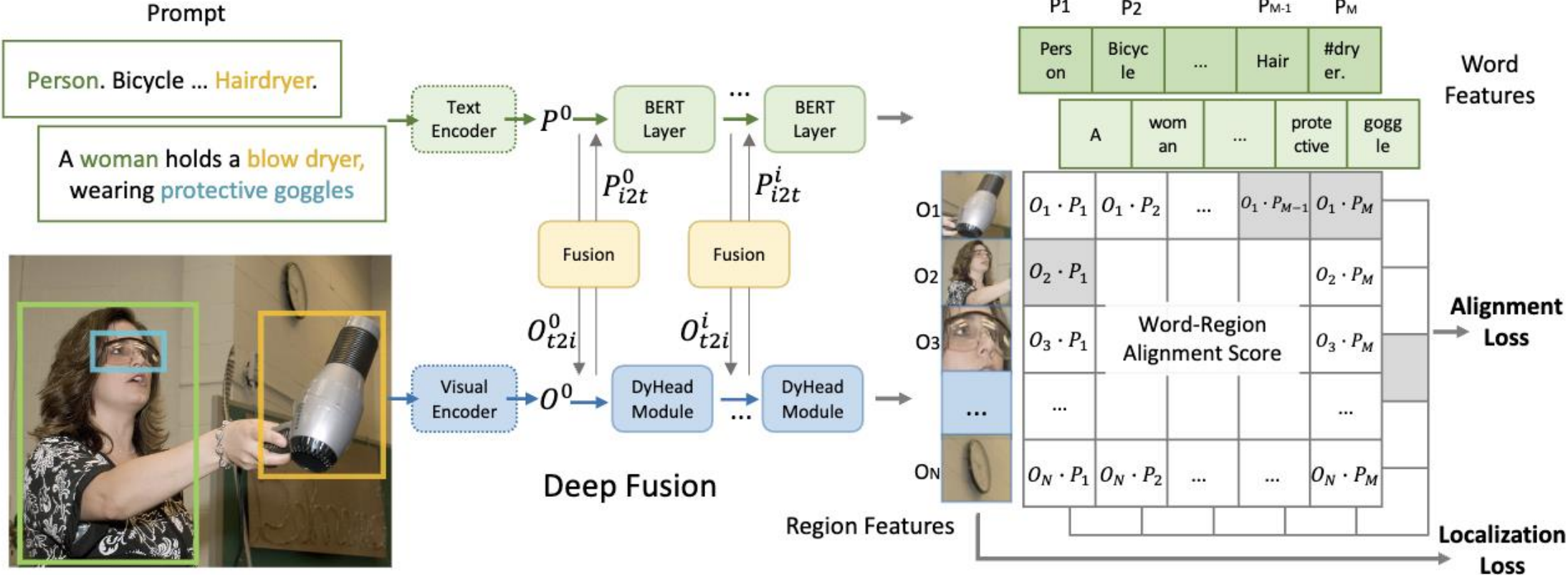
# Object Detection and Text Grounding

Person. Bicycle ... **Hairdryer.**

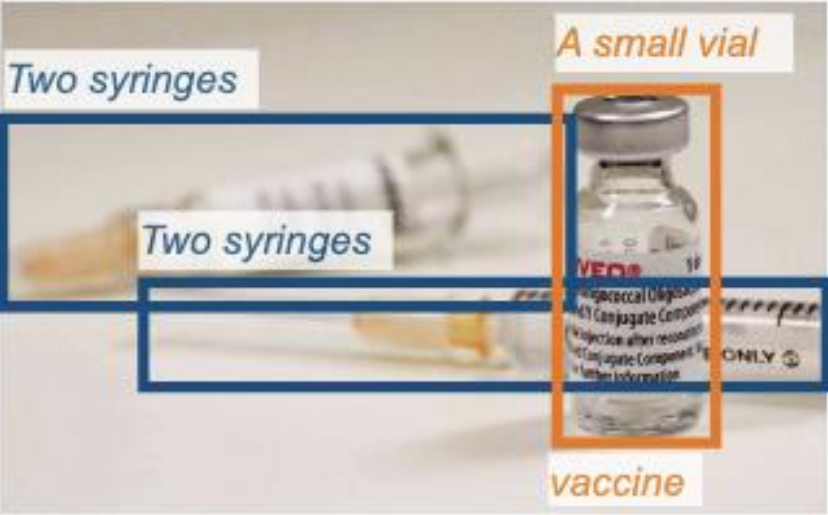
A **woman** holds a **blow dryer**,  
wearing **protective goggles**



# GLIP: Grounded Language-Image Pre-training



# Grounding Results



Two syringes and a small vial of vaccine.



playa esmeralda in holguin, cuba. the view from the top of the beach. beautiful caribbean sea turquoise

# Zero-Shot Grounding

Model	Backbone	MiniVal [23]				Val v1.0			
		APr	APc	APf	AP	APr	APc	APf	AP
MDETR [23]	RN101	20.9	24.9	24.3	24.2	-	-	-	-
MaskRCNN [23]	RN101	26.3	34.0	33.9	33.3	-	-	-	-
Supervised-RFS [15]	RN50	-	-	-	-	12.3	24.3	32.4	25.4
GLIP-T (A)	Swin-T	14.2	13.9	23.4	18.5	6.0	8.0	19.4	12.3
GLIP-T (B)	Swin-T	13.5	12.8	22.2	17.8	4.2	7.6	18.6	11.3
GLIP-T (C)	Swin-T	17.7	19.5	<b>31.0</b>	24.9	7.5	11.6	<b>26.1</b>	16.5
GLIP-T	Swin-T	<b>20.8</b>	<b>21.4</b>	<b>31.0</b>	<b>26.0</b>	<b>10.1</b>	<b>12.5</b>	25.5	<b>17.2</b>
GLIP-L	Swin-L	<b>28.2</b>	<b>34.3</b>	<b>41.5</b>	<b>37.3</b>	<b>17.1</b>	<b>23.3</b>	<b>35.4</b>	<b>26.9</b>

# DesCo: Object Recognition with Language Description

Detect with specifications for shape & subpart  
Target Object Confusable Object



Eclair, a kind of food, **long**, cylindrical pastry, filled with cream, topped with chocolate

Tart, a kind of food, **round**, could be filled with fruits, could be served with cream

Detect with specifications for relation  
Target Object Confusable Object



A clown making a balloon animal for a pretty lady

A clown kicking a soccer ball for a pretty lady

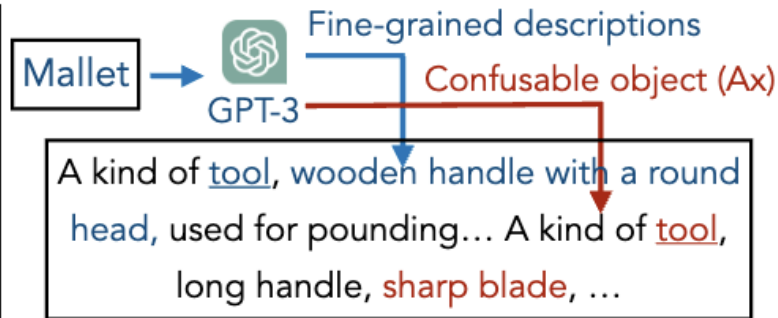
# DesCo: Description-Conditioned

Detect: Mallet.  
Bear. Cat...

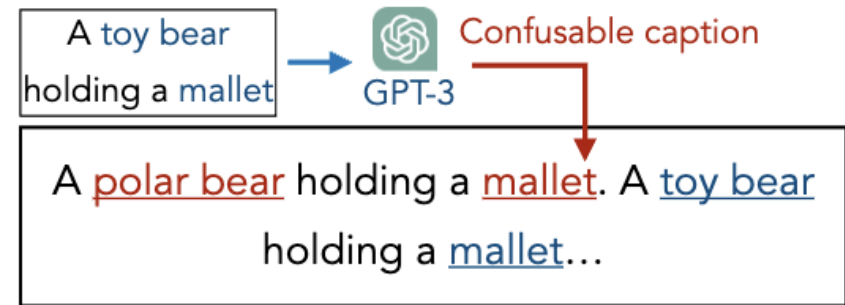
A toy bear holding  
a mallet.



Original training  
data for GLIP



	...	tool	...	tool	...
	0	1	0	0	0
	0	0	0	0	0
	0	0	0	0	0
	0	0	0	0	0



	...	polar bear	...	mallet	...	toy bear	...	mallet
	0	0	0	0	0	0	0	1
	0	0	0	0	0	0	0	0
	0	0	0	0	0	1	0	0
	0	0	0	0	0	0	0	0

Description-rich and context-sensitive data for DESCO-GLIP

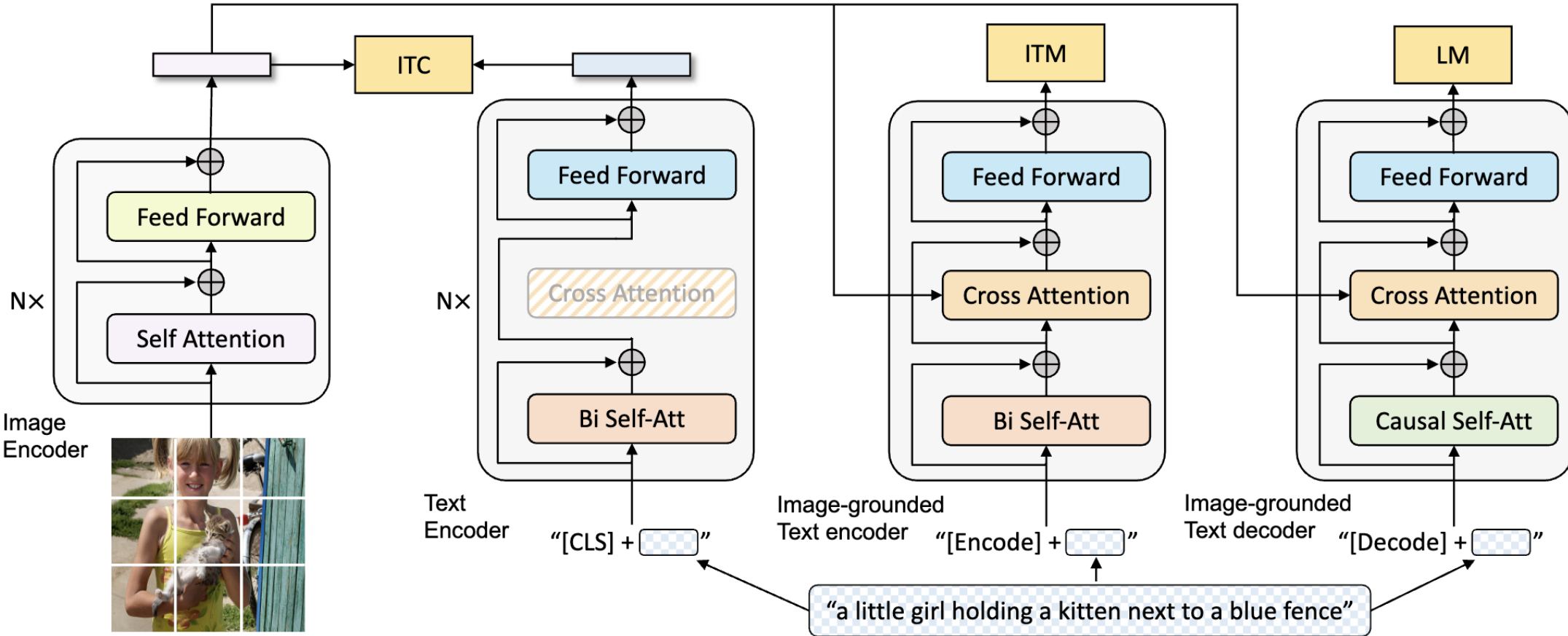
# Zero-Shot Grounding

Model	Backbone	LVIS MiniVal [16]				OmniLabel [34]			
		APr	APc	APf	AP	AP	APc	APd	APd-P
MDETR [16]	RN101	20.9	24.9	24.3	24.2	-	-	4.7	9.1
MaskRCNN [16]	RN101	26.3	34.0	33.9	33.3	-	-	-	-
RegionCLIP [50]	ResNet-50	-	-	-	-	2.7	2.7	2.6	3.2
Detic [52]	Swin-B	-	-	-	-	8.0	15.6	5.4	8.0
K-LITE [37]	Swin-T	14.8	18.6	24.8	21.3	-	-	-	-
GroundingDINO-T [25]	Swin-T	18.1	23.3	32.7	27.4	-	-	-	-
GroundingDINO-L [25]	Swin-L	22.2	30.7	38.8	33.9	-	-	-	-
GLIP-L [22]	Swin-L	28.2	34.3	41.5	37.3	25.8	32.9	21.2	33.2
GLIP-T [22]	Swin-T	20.8	21.4	31.0	26.0	19.3	23.6	16.4	25.8
DESCO-GLIP	Swin-T	<b>30.8</b>	<b>30.5</b>	<b>39.0</b>	<b>34.6</b>	<b>23.8</b>	<b>27.4</b>	<b>21.0</b>	<b>30.4</b>
FIBER-B [7]	Swin-B	25.7	29.0	39.5	33.8	25.7	30.3	22.3	34.8
DESCO-FIBER	Swin-B	<b>34.8</b>	<b>35.5</b>	<b>43.9</b>	<b>39.5</b>	<b>29.3</b>	<b>31.6</b>	<b>27.3</b>	<b>37.7</b>

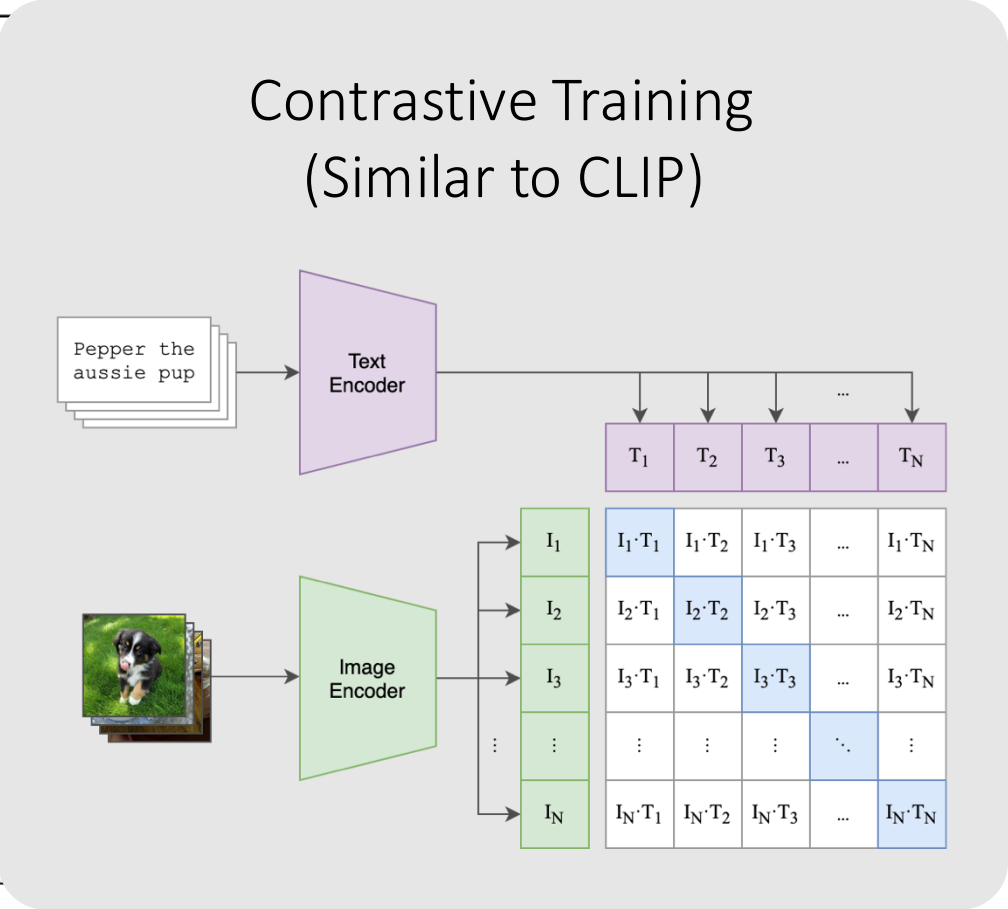
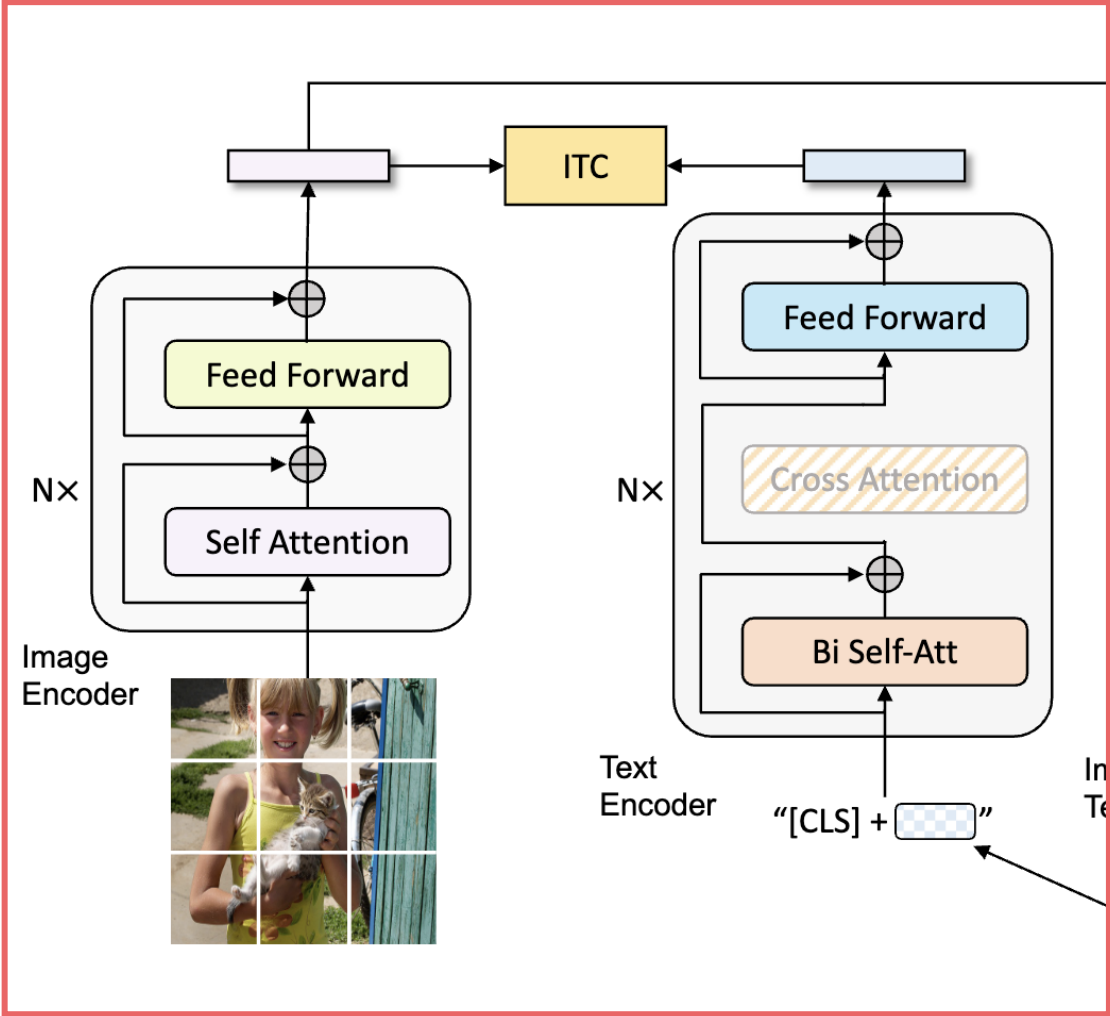
# Encoder-Only vs. Encoder-Decoder

- Encoder-only
  - CLIP, GLIP, DesCo, etc.
  - Better for image-text retrieval
- Encoder-decoder
  - Better for generation

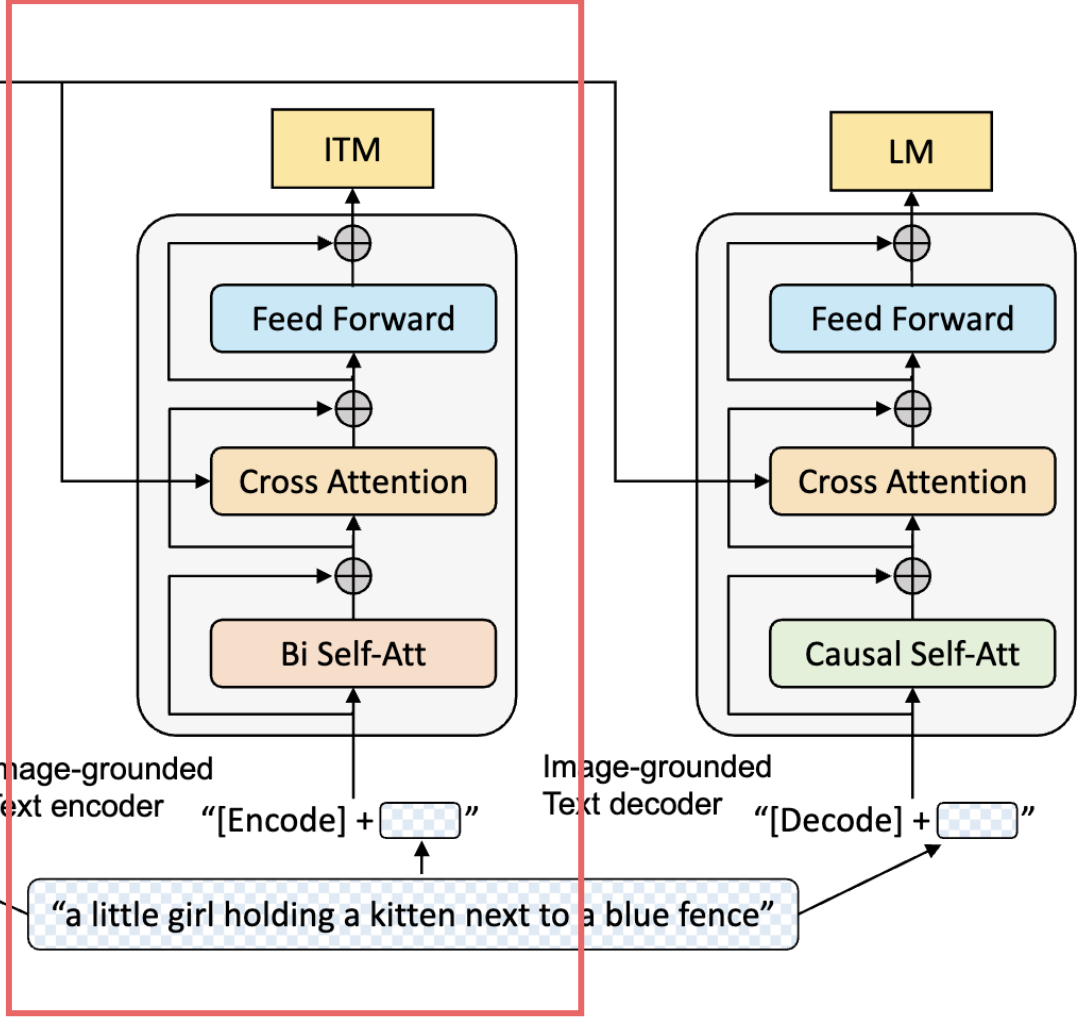
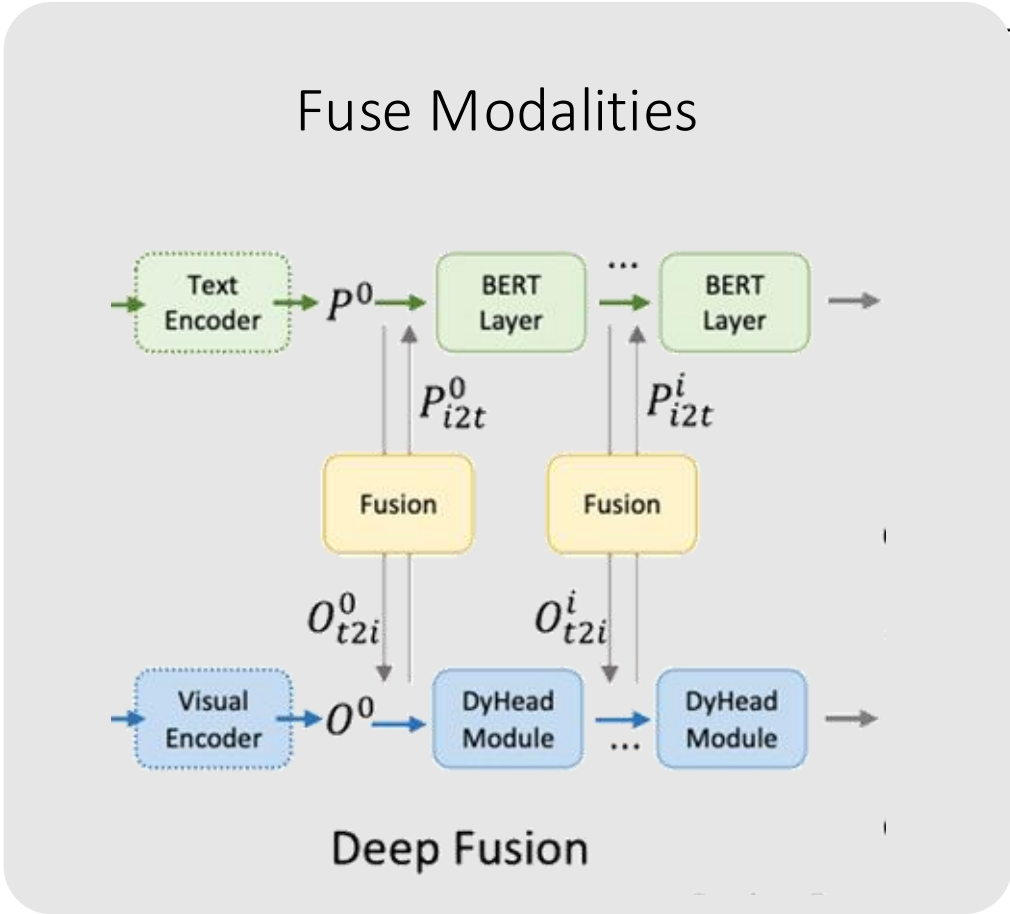
# BLIP: Bootstrapping Language-Image Pre-training



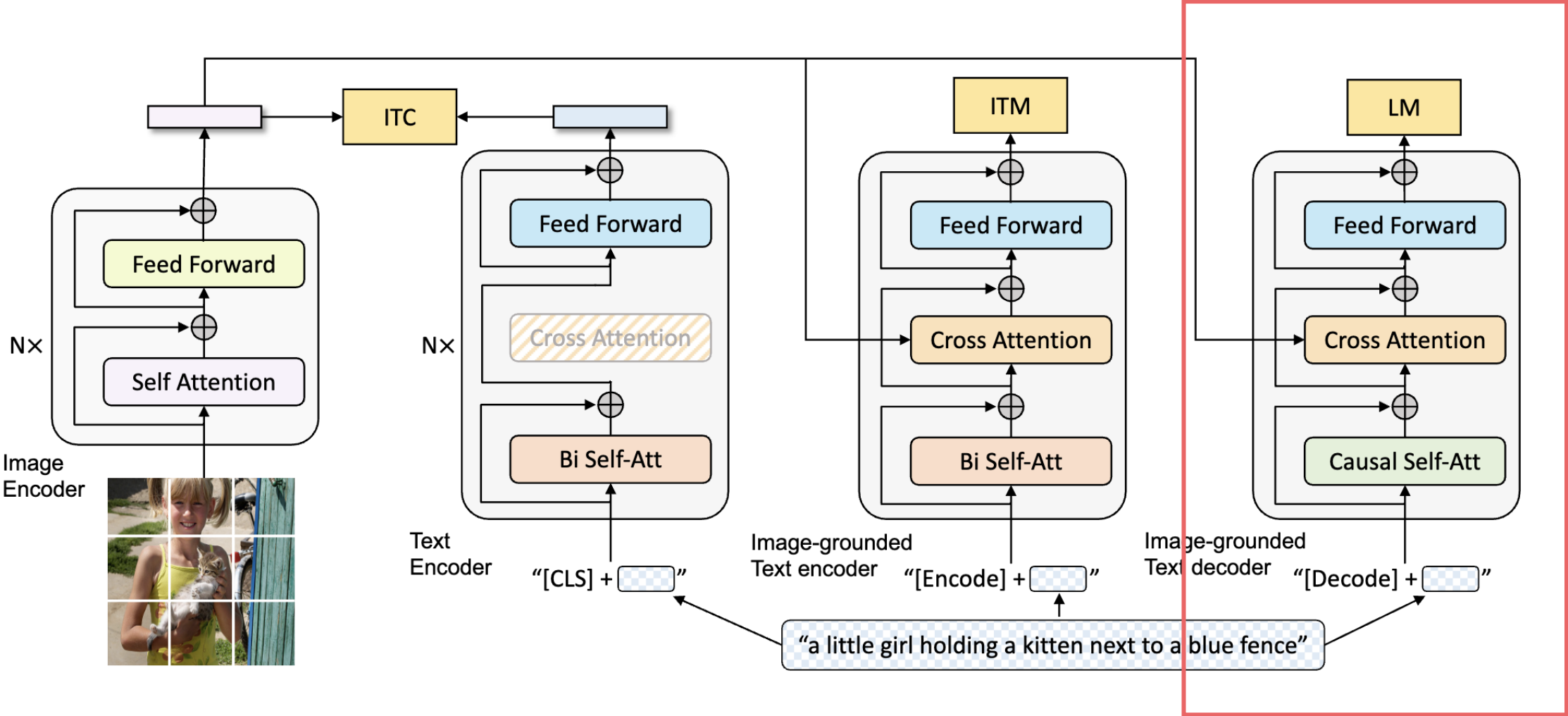
# A Unified Framework



# A Unified Framework



# A Unified Framework



# Zero-Shot Image-Text Retrieval

Method	Pre-train # Images	Flickr30K (1K test set)					
		TR			IR		
		R@1	R@5	R@10	R@1	R@5	R@10
CLIP	400M	88.0	98.7	99.4	68.7	90.6	95.2
ALIGN	1.8B	88.6	98.7	99.7	75.7	93.8	96.8
ALBEF	14M	94.1	99.5	99.7	82.8	96.3	98.1
BLIP	14M	94.8	99.7	<b>100.0</b>	84.9	96.7	98.3
BLIP	129M	<b>96.0</b>	<b>99.9</b>	<b>100.0</b>	85.0	<b>96.8</b>	98.6
BLIP <sub>CapFilt-L</sub>	129M	<b>96.0</b>	<b>99.9</b>	<b>100.0</b>	<b>85.5</b>	<b>96.8</b>	<b>98.7</b>
BLIP <sub>ViT-L</sub>	129M	96.7	100.0	100.0	86.7	97.3	98.7

# Image Captioning

Method	Pre-train #Images	NoCaps validation								COCO Caption Karpathy test	
		in-domain		near-domain		out-domain		overall		B@4	C
		C	S	C	S	C	S	C	S		
Enc-Dec (Changpinyo et al., 2021)	15M	92.6	12.5	88.3	12.1	94.5	11.9	90.2	12.1	-	110.9
VinVL† (Zhang et al., 2021)	5.7M	103.1	14.2	96.1	13.8	88.3	12.1	95.5	13.5	38.2	129.3
LEMON <sub>base</sub> † (Hu et al., 2021)	12M	104.5	14.6	100.7	14.0	96.7	12.4	100.4	13.8	-	-
LEMON <sub>base</sub> † (Hu et al., 2021)	200M	107.7	14.7	106.2	14.3	107.9	13.1	106.8	14.1	<b>40.3</b>	<b>133.3</b>
BLIP	14M	111.3	15.1	104.5	14.4	102.4	13.7	105.1	14.4	38.6	129.7
BLIP	129M	109.1	14.8	105.8	14.4	105.7	13.7	106.3	14.3	39.4	131.4
BLIP <sub>CapFilt-L</sub>	129M	<b>111.8</b>	<b>14.9</b>	<b>108.6</b>	<b>14.8</b>	<b>111.5</b>	<b>14.2</b>	<b>109.6</b>	<b>14.7</b>	39.7	<b>133.3</b>
LEMON <sub>large</sub> † (Hu et al., 2021)	200M	116.9	15.8	113.3	15.1	111.3	14.0	113.4	15.0	40.6	135.7
SimVLM <sub>huge</sub> (Wang et al., 2021)	1.8B	113.7	-	110.9	-	115.2	-	112.2	-	40.6	143.3
BLIP <sub>ViT-L</sub>	129M	114.9	15.2	112.1	14.9	115.3	14.4	113.2	14.8	40.4	136.7

# Visual Question Answering

## Visual Question Answering

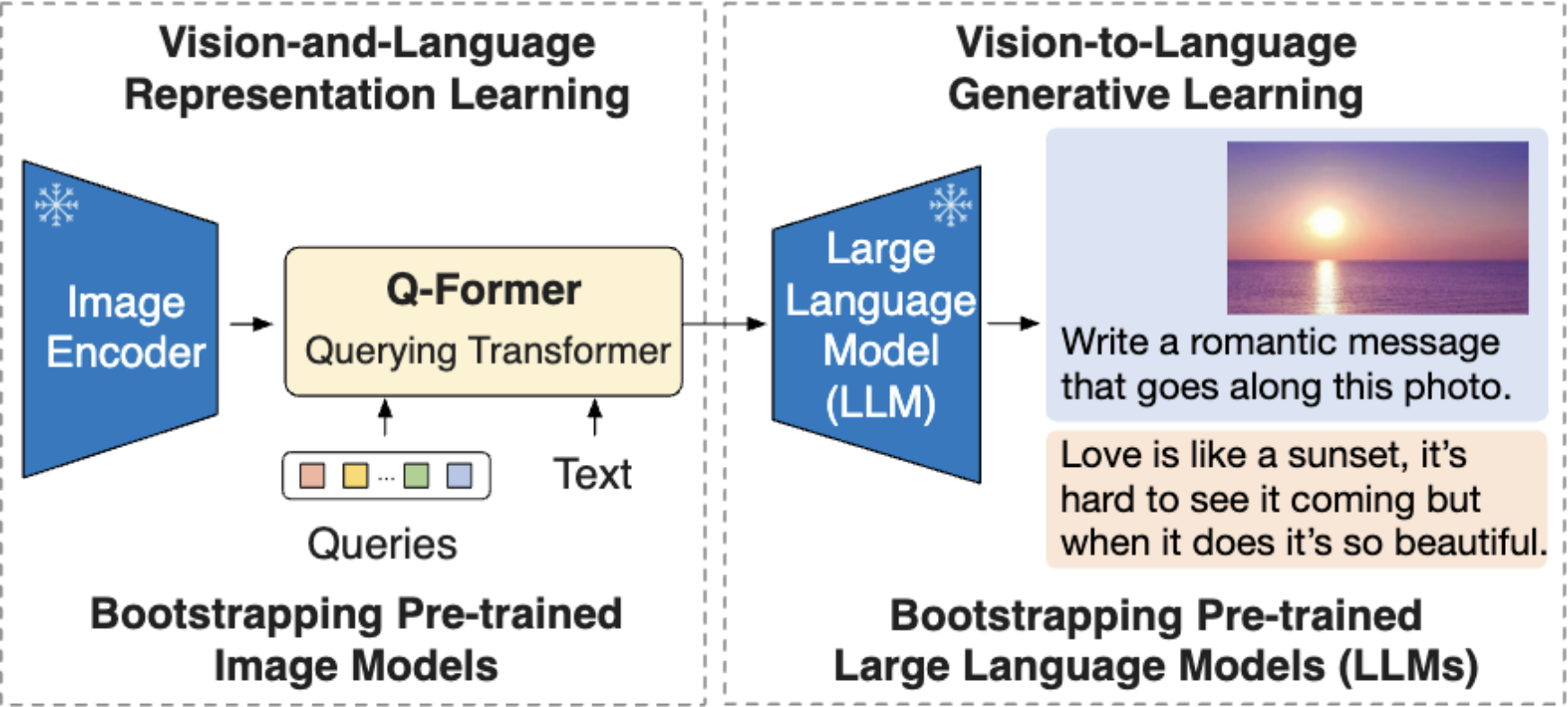


Does it appear to be rainy?

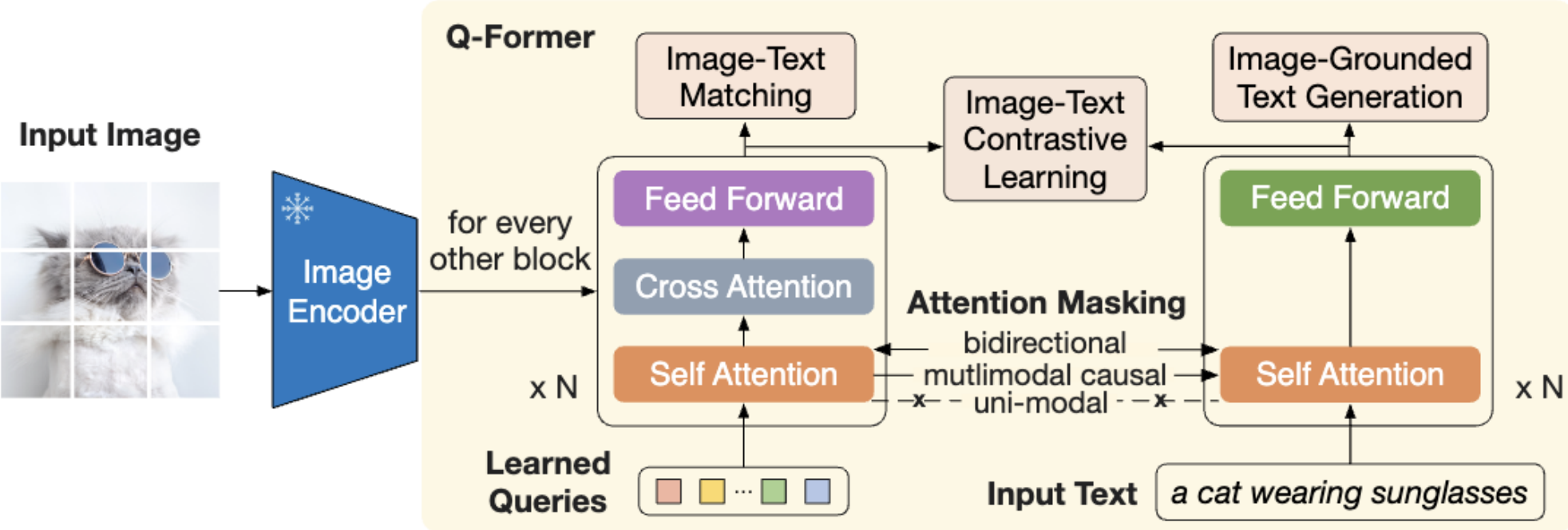
Does this person have 20/20 vision?

Method	Pre-train #Images	VQA		NLVR <sup>2</sup>	
		test-dev	test-std	dev	test-P
LXMERT	180K	72.42	72.54	74.90	74.50
UNITER	4M	72.70	72.91	77.18	77.85
VL-T5/BART	180K	-	71.3	-	73.6
OSCAR	4M	73.16	73.44	78.07	78.36
SOHO	219K	73.25	73.47	76.37	77.32
VILLA	4M	73.59	73.67	78.39	79.30
UNIMO	5.6M	75.06	75.27	-	-
ALBEF	14M	75.84	76.04	82.55	83.14
SimVLM <sub>base</sub> <sup>†</sup>	1.8B	77.87	78.14	81.72	81.77
BLIP	14M	77.54	77.62	<b>82.67</b>	82.30
BLIP	129M	78.24	78.17	82.48	<b>83.08</b>
BLIP <sub>CapFilt-L</sub>	129M	<b>78.25</b>	<b>78.32</b>	82.15	82.24

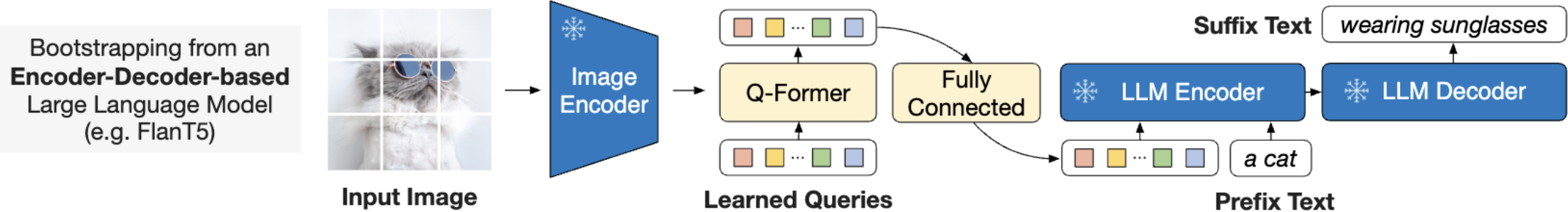
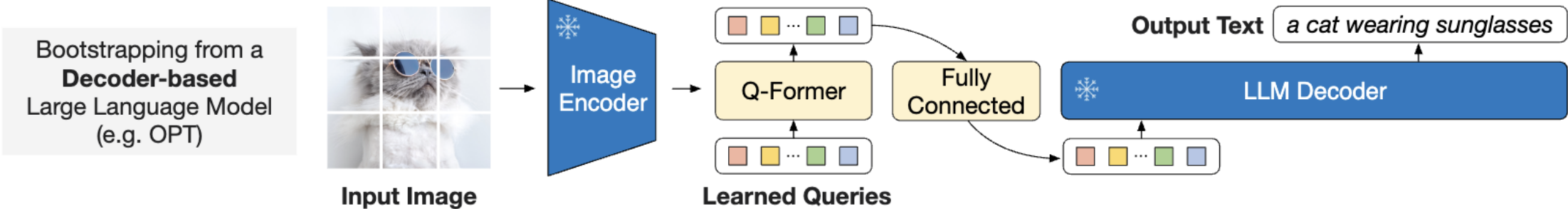
# BLIP-2: Frozen Image Encoders and Large Language Models




# Vision-Language Representation Learning



# Vision-to-Language Generative Pre-Training




# Instructed Zero-Shot Image-to-Text Generation




Explain the advantages of this product.

The audi e-tron quattro concept is a plug-in hybrid electric sports car that has a range of 310 miles and can accelerate from 0 to 60 mph in just 4 seconds.




Tell me something about the history of this place.

The great wall of china was built in 221 bc by emperor qin shi Huang to protect the emperor's capital from invaders from the north.



Write down the facts that you know about this flower.

It is a genus of flowering plants in the orchid family, orchidaceae. It is native to eastern asia and south east asia. It is a perennial flowering plant that grows up to a height of 1 m.



Is this photo unusual?

Yes, it's a house that looks like it's upside down.

How could someone get out of the house?

It has a slide on the side of the house.




What are shown in the photo?

A man and a chicken.

What does the man feel and why?

He is scared of the chicken because it is flying at him.



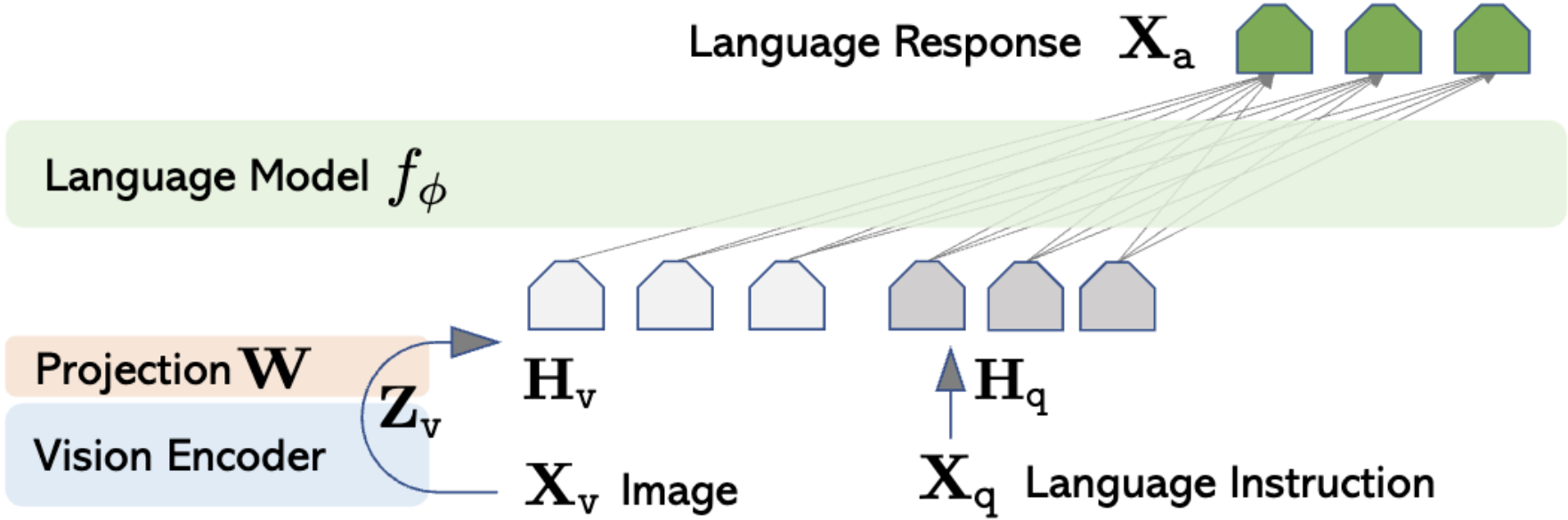
What are the ingredients I need to make this?

Pepperoni, mozzarella cheese, pizza sauce, olive oil, salt, pepper, basil.

What is the first step?

Place the pizza dough on a baking sheet, brush with olive oil, sprinkle with salt, pepper, and basil.

# LLaVA: Visual Instruction Tuning



# GPT-Assisted Visual Instruction Data Generation

## Context type 1: Captions

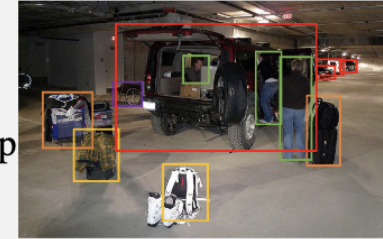
A group of people standing outside of a black vehicle with various luggage.

Luggage surrounds a vehicle in an underground parking area

People try to fit all of their luggage in an SUV.

The sport utility vehicle is parked in the public garage, being packed for a trip

Some people with luggage near a van that is transporting it.



## Context type 2: Boxes

person: [0.681, 0.242, 0.774, 0.694], backpack: [0.384, 0.696, 0.485, 0.914], suitcase: ...<omitted>

---

## Response type 1: conversation

Question: What type of vehicle is featured in the image?

Answer: The image features a black sport utility vehicle (SUV) ...<omitted>

## Response type 2: detailed description

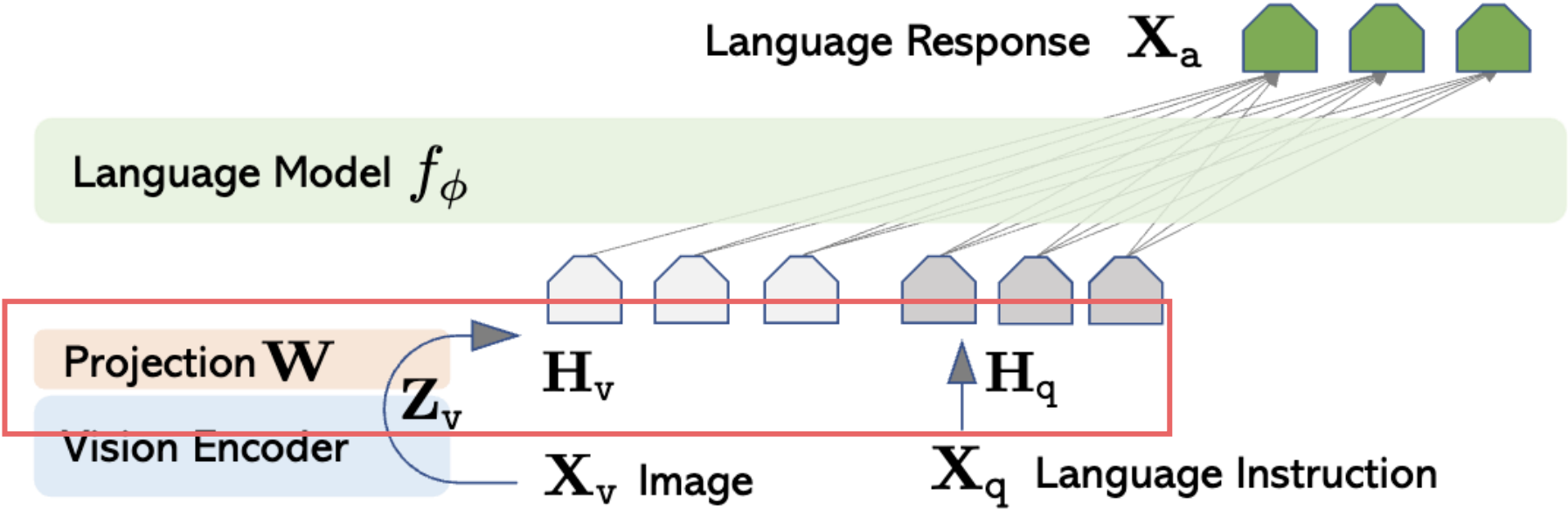
The image is an underground parking area with a black sport utility vehicle (SUV) parked. There are three people in the scene, with one person standing closer to the left side of the vehicle, another person in the middle, and the third person on the right side. They are all working together to pack their luggage into the SUV for a trip. ...<omitted>

## Response type 3: complex reasoning

Question: What challenges do these people face?

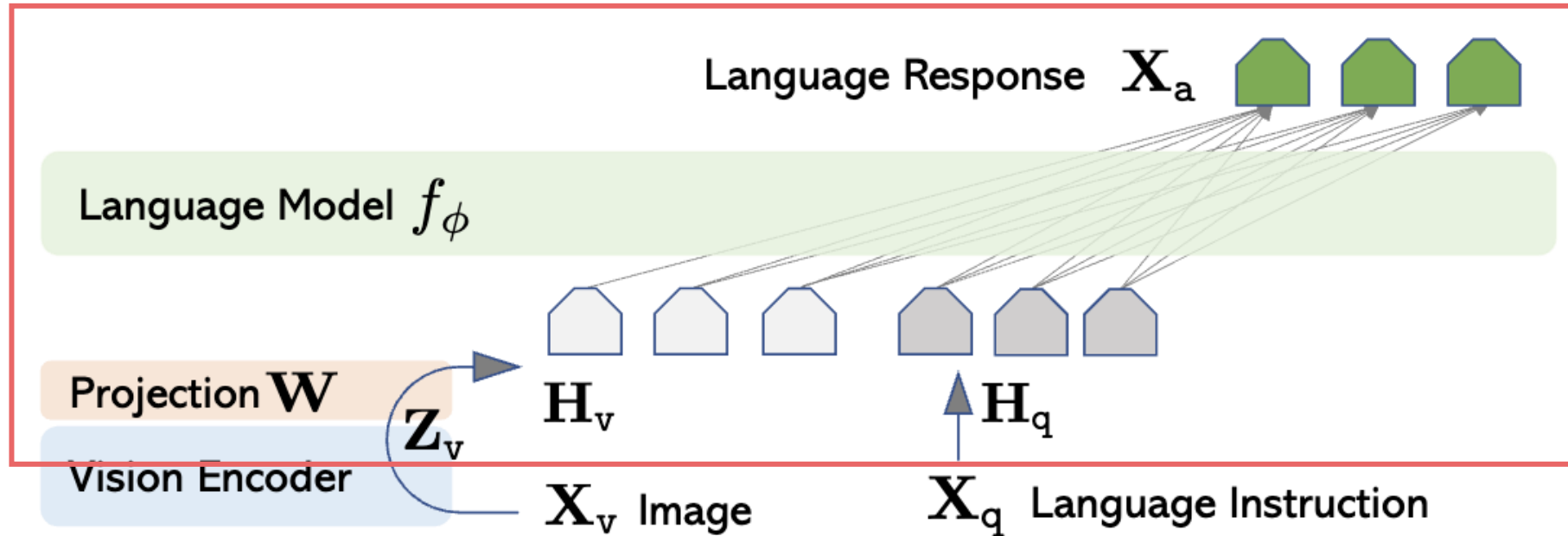
Answer: In the image, a group of people is standing outside a black SUV in a parking area, surrounded by various pieces of luggage, including suitcases and backpacks. They are facing the challenge of fitting all their luggage into the black SUV. There are multiple suitcases and backpacks to be packed, which suggests that the group has a significant amount of belongings ...<omitted>

# Stage 1: Pre-training for Feature Alignment



Train with Image-Text Pairs

# Stage 2: Fine-tuning End-to-End



Visual Chat (Visual Instruction Data) and Science QA

# Examples

---

## Visual input example, Extreme Ironing:

---



Source: <https://www.barnorama.com/wp-content/uploads/2016/12/03-Confusing-Pictures.jpg>

---

User	What is unusual about this image?
GPT-4 [36]	The unusual thing about this image is that a man is ironing clothes on an ironing board attached to the roof of a moving taxi.
User	What is unusual about this image?
BLIP-2	a man is sitting on the back of a yellow cab
User	What is unusual about this image?
OpenFlamingo	The man is drying his clothes on the hood of his car.

---

User  
LLaVA

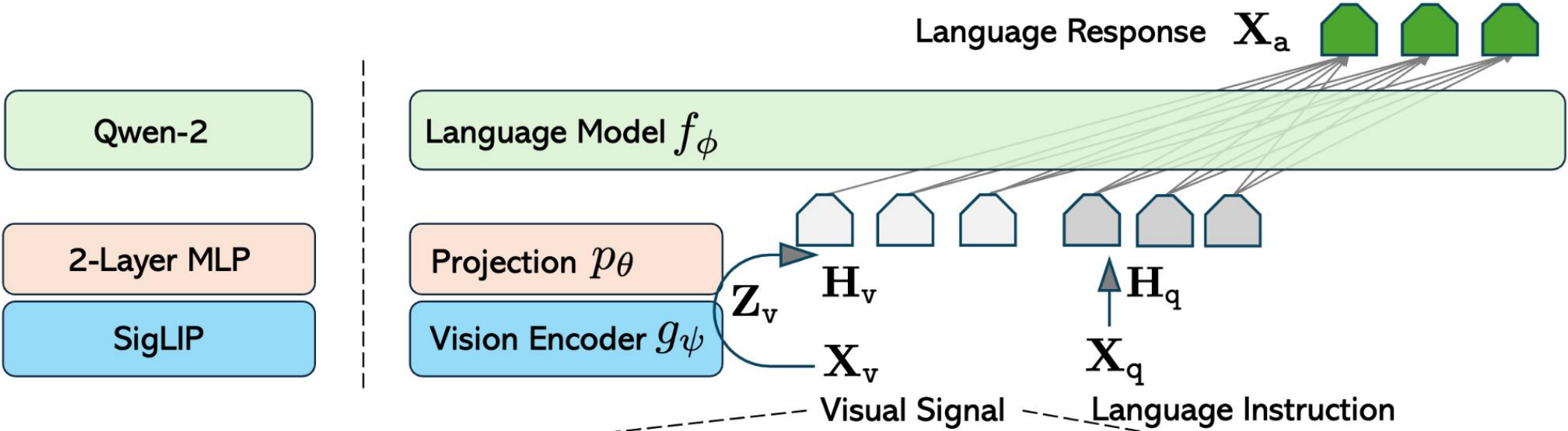
What is unusual about this image?  
The unusual aspect of this image is a man ironing clothes on the back of a minivan or van. This is not a typical place to perform this activity, as one would usually iron clothes in a more stationary and safe location, such as a home, using a regular ironing board. The scene depicted in the image is peculiar as it involves a makeshift ironing setup on a vehicle, which can be both unsafe and unconventional. Additionally, it is not clear how the man is able to maintain balance and stability while ironing clothes in such an unstable environment.

[Start a new conversation, and clear the history]

User  
LLaVA

What's happening in the scene?  
The image depicts a man standing on top of a yellow SUV in a busy city street. He is holding a portable ladder, seemingly ironing clothes while standing on the vehicle. Around the scene, there are other cars, a traffic light, and a fire hydrant. Additionally, there are two other people in the vicinity, one close to the man on the car and another further away. The man on the ladder appears to be attracting attention **due to his unconventional choice of ironing his clothes on top of a moving car**. The city street around him is bustling with activity, adding to the unique nature of the scene.

# LLaVA-OneVision



Single Image



Multi-Image



Video

# LLaVA-OneVision

 <p>Single-Image</p>	 <p>... N Crops</p> <p>729 + N * 729 Tokens</p>	<p><math>(1 + 9) * 729 = 7290</math> Tokens</p>
 <p>Multi-Image</p>	 <p>... N Images</p> <p>N * 729 Tokens</p>	<p><math>12 * 729 = 8748</math> Tokens</p>
 <p>Video</p>	 <p>... N Frames</p> <p>N * 196 Tokens</p>	<p><math>32 * 196 = 6272</math> Tokens</p>
<p>Example on Token Strategy</p>		<p>Max Tokens</p>

# LLaVA-NeXT-Interleave

## Spot the Difference

User:  
Tell the detailed differences of these two images:



LLaVA-NeXT-M3:

1. In the first image, the player on the top right is wearing a red uniform, but is wearing a blue uniform in the second image.
2. In the first image, the player on the top left has the number "9", but changes to "6" in the second image ...

## Visual Story Telling

User:  
Complete the story based on the given images and plots:



We had an exciting day at the fair!



The lights lit up the night and the rides made us all dizzy.



The arcade games had the funniest stuffed monkeys as prizes.



LLaVA-NeXT-M3:  
We threw a million darts trying to win one!

## Image Edit Instruction

User:  
Provide the image edit instruction that can transform the source image to the target image:



LLaVA-NeXT-M3:  
Replace the soccer ball with a giant sunflower that the puppy is holding in the same manner as the soccer ball.

## Social Media Post

User:  
Help me write a Twitter post providing the following images:



LLaVA-NeXT-M3:

Wandering through the enchanting streets of Milan, where the vibrant pulse of the city is as captivating as its historical architecture. From the mesmerizing ... the gentle glow of candles offers a moment of serenity in this bustling metropolis. #MilanDiaries #TravelMoments #CityLights

## Low-level Comparison

User:  
Compared to the first image, how is the sharpness of the second image, sharper, blurrier, or about the same?



LLaVA-NeXT-M3: Sharper

## Video Detailed Description

User:  
Given me a detailed description of the video:



LLaVA-NeXT-M3:

The video depicts an interior scene, specifically a hallway. The flooring is wooden with a light brown hue and a polished finish. There are two doors visible in the video, both with dark frames and glass panels. The door ...

## 3D Grounding

User:  
Where is the bathtub located?



LLaVA-NeXT-M3:  
Behind shower curtain inside bathroom

# How Reliable Vision-Language Models Are?

- ARO (Attribution, Relation and Order) Benchmark

## Visual Genome Relation

Assessing relational understanding (23,937 test cases)



- ✓ *the horse is eating the grass*
- ✗ *the grass is eating the horse*

## Visual Genome Attribution

Assessing attributive understanding (28,748 test cases)



- ✓ *the paved road and the white house*
- ✗ *the white road and the paved house*

# ARO (Attribution, Relation and Order) Benchmark

## Visual Genome Relation

Assessing relational understanding (23,937 test cases)



- ✓ *the horse is eating the grass*
- X *the grass is eating the horse*

## Visual Genome Attribution

Assessing attributive understanding (28,748 test cases)



- ✓ *the paved road and the white house*
- X *the white road and the paved house*

# ARO (Attribution, Relation and Order) Benchmark

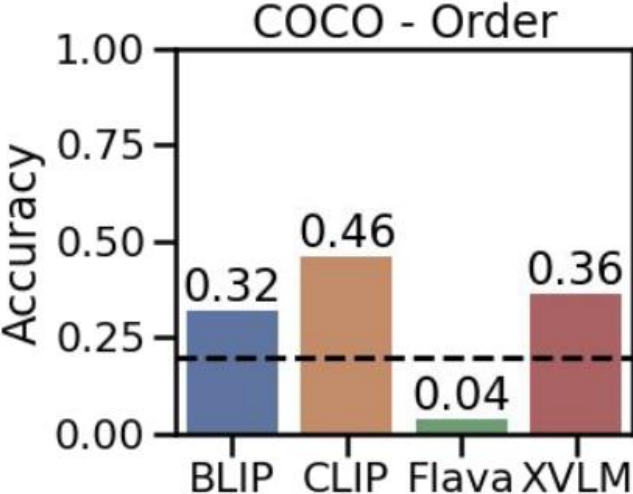
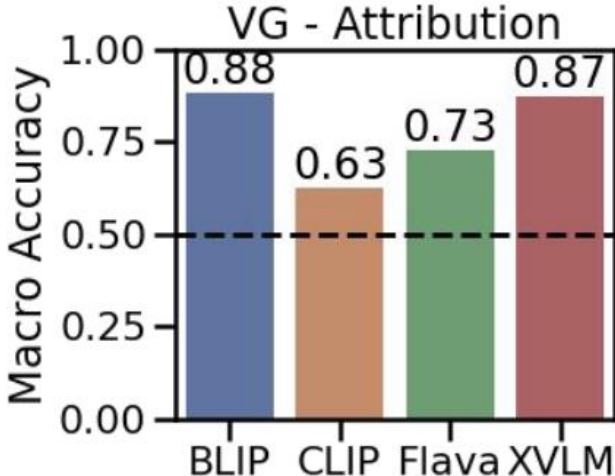
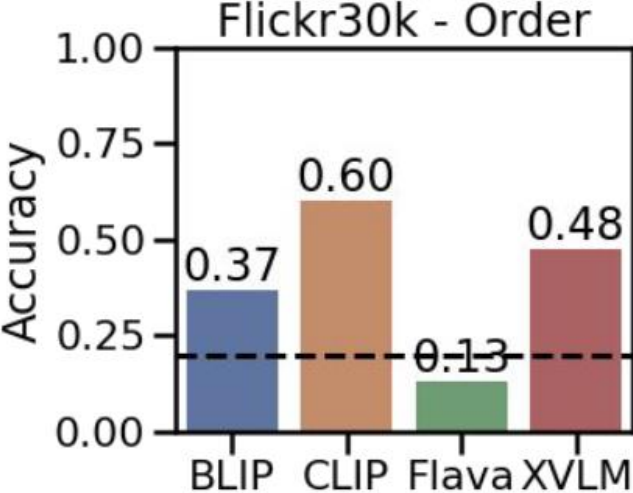
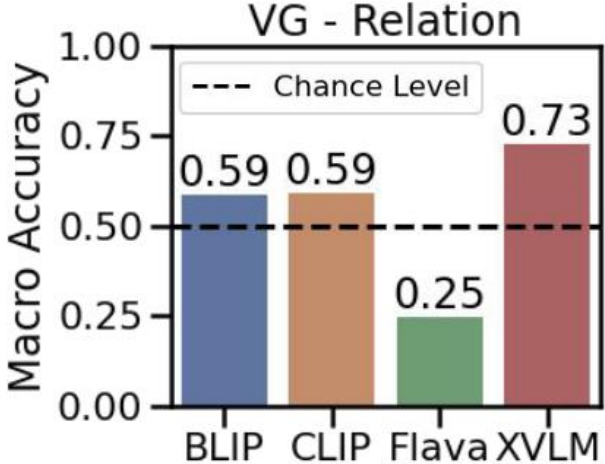
## COCO Order and Flickr Order

### Assessing sensitivity to order (6,000 test cases)



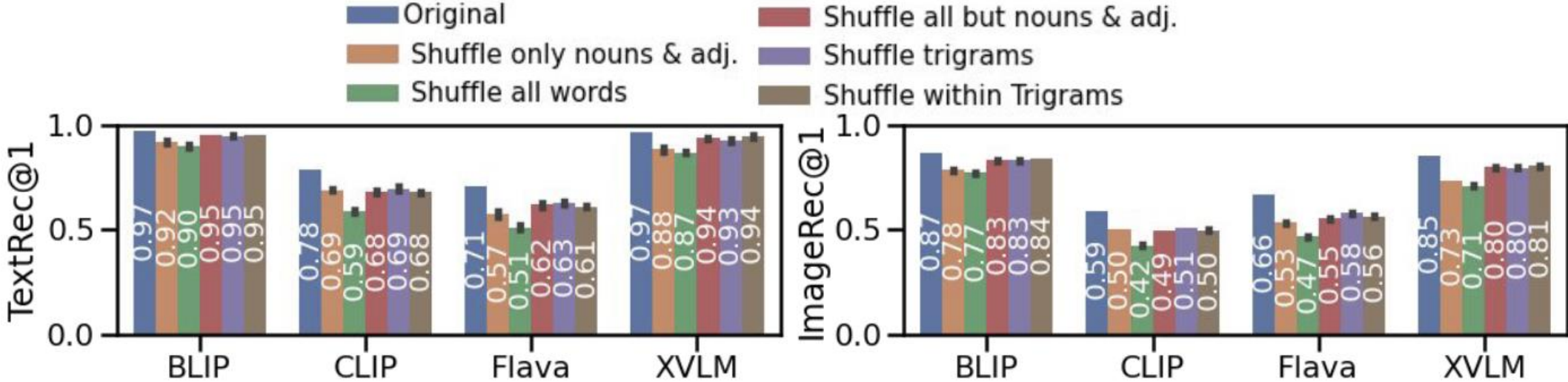
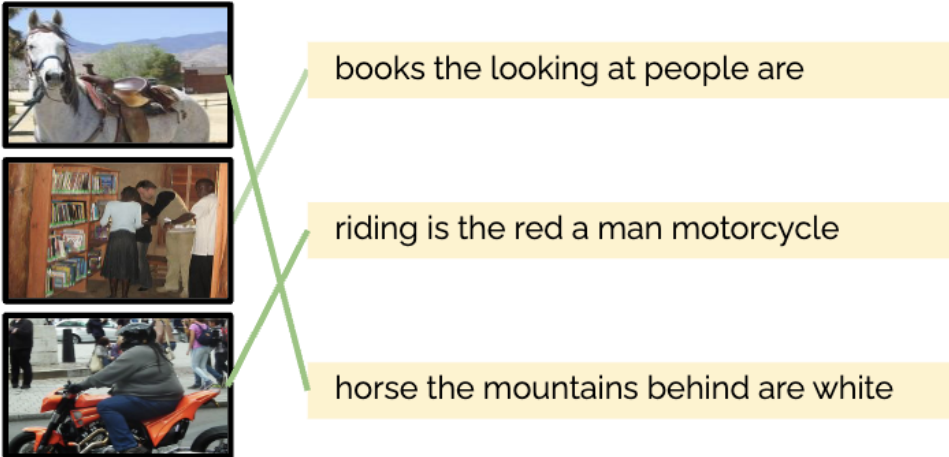
- ✓ a brown cat is looking at a gray dog and sitting in a white bathtub
- X (shuffle adjective/noun) a gray bathtub is looking at a white cat and sitting in a brown dog
- X (shuffle all but adjective/noun) at brown cat a in looking a gray dog sitting is and a white bathtub
- X (shuffle words within trigrams) cat brown a at is looking a gray dog in and sitting bathtub a white
- X (shuffle trigrams) a brown cat a white bathtub is looking at a gray dog and sitting in

# Results



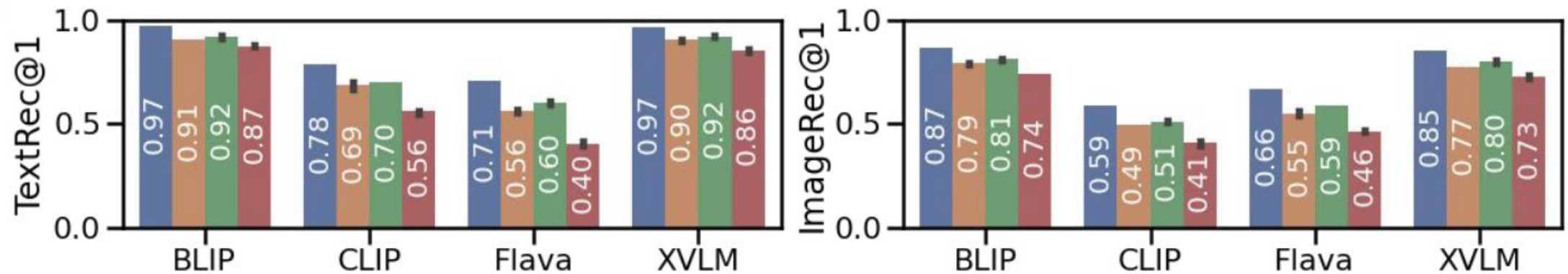
# More Analysis on Text-Image Retrieval

## Retrieval without access to word order



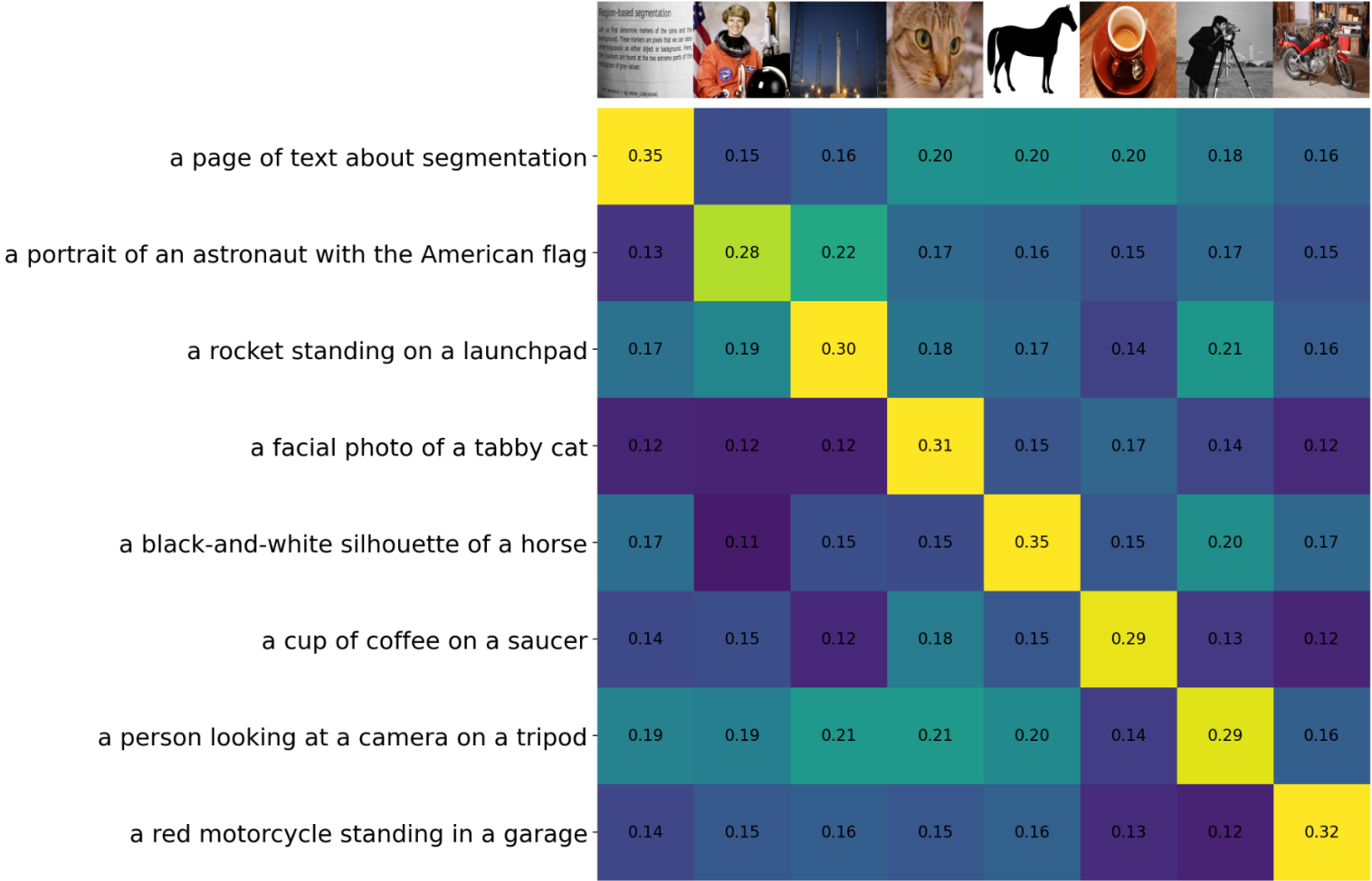
# More Analysis on Text-Image Retrieval

## Retrieval without access to visual patch order



# Contrastive Pre-Training

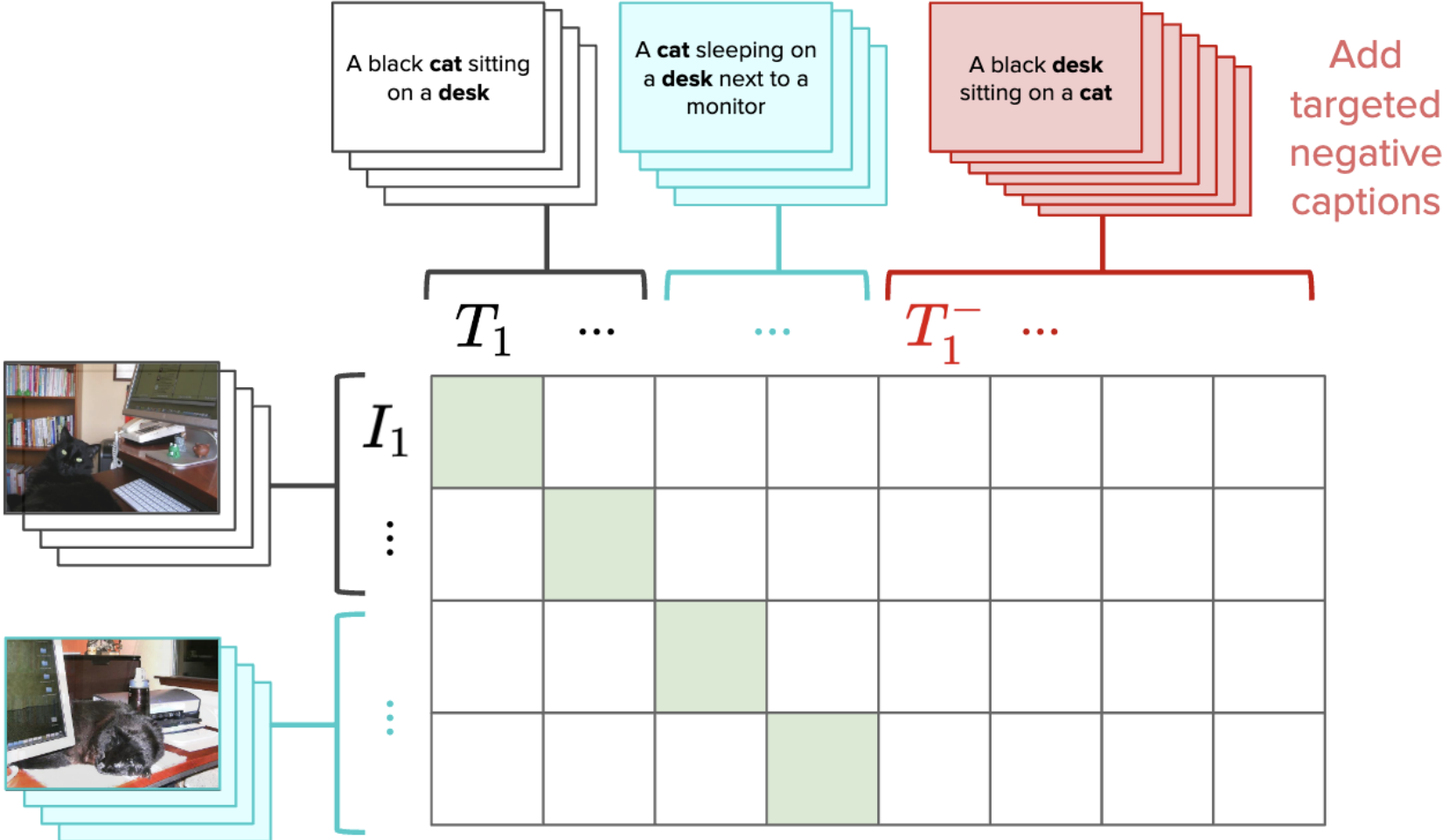
Cosine similarity between text and image features



# Solution: Composition-Aware Hard Negatives

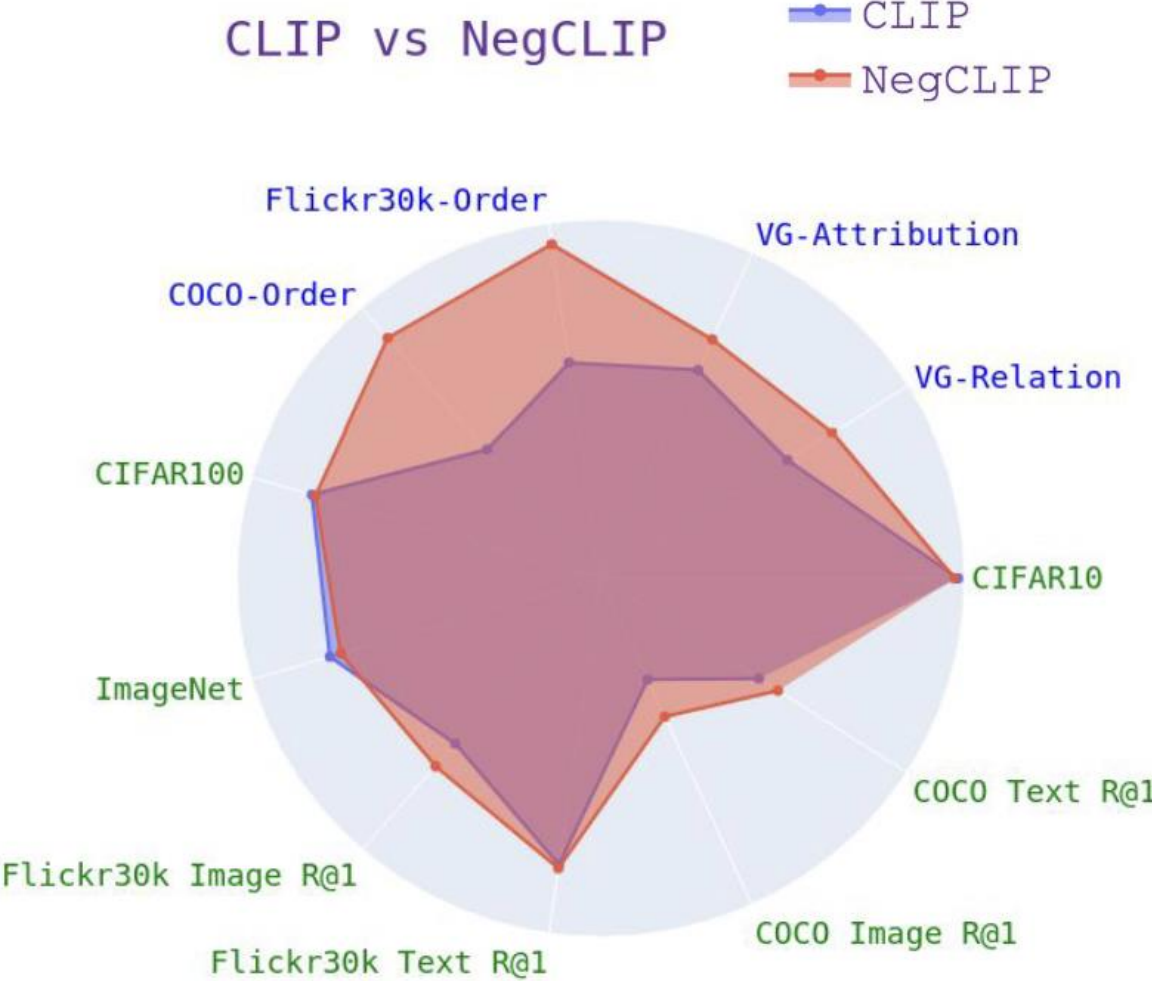
- Generation of negative captions
  - Generate a negative caption by swapping different linguistic elements: noun phrases, nouns, adjectives, adverbs, verb phrases
- Sampling strong alternative images
  - Sample one of the  $K$  nearest neighbors as the strong alternative image

# NegCLIP



Add strong alternative images

# Results



# What's Up: Benchmarks with Spatial Reasoning Questions

What'sUp (Subset A)



- A mug on a table
- A mug under a table
- A mug to the left of a table
- A mug to the right of a table



- A mug on a table
- A mug under a table
- A mug to the left of a table
- A mug to the right of a table



- A mug on a table
- A mug under a table
- A mug to the left of a table
- A mug to the right of a table



- A mug on a table
- A mug under a table
- A mug to the left of a table
- A mug to the right of a table

# What's Up: Benchmarks with Spatial Reasoning Questions

What'sUp (Subset B)



- A mug in front of a plate
- A mug behind a plate
- A mug to the left of a plate
- A mug to the right of a plate



- A mug in front of a plate
- A mug behind a plate
- A mug to the left of a plate
- A mug to the right of a plate



- A mug in front of a plate
- A mug behind a plate
- A mug to the left of a plate
- A mug to the right of a plate



- A mug in front of a plate
- A mug behind a plate
- A mug to the left of a plate
- A mug to the right of a plate

# Results

Model	Whats- Up	COCO- spatial	GQA- spatial	Avg
CLIP ViT-B/32	31.0	47.4	46.9	41.8
CLIP ViT-L/14	26.1	49.5	47.3	41.0
NegCLIP	34.4	46.9	46.0	42.4
RoBERTaCLIP	25.1	50.0	49.8	41.6
CoCa	29.4	46.7	47.1	41.0
XVLM 4M	31.5	61.7	<b>58.7</b>	50.6
XVLM 16M	<b>41.9</b>	<b>65.0</b>	58.2	<b>55.0</b>
BLIP 14M	38.5	54.0	49.8	47.5
BLIP 129M	30.4	49.3	49.0	42.9
BLIP2-ITM	37.6	53.0	49.8	46.8
BLIP2-ITC	29.0	53.7	51.0	44.6
FLAVA	30.5	52.6	51.7	44.9
CoCa-Caption	24.1	48.6	49.5	40.8
XVLM-Flickr30K	44.3	65.2	61.4	56.9
XVLM-COCO	42.1	<b>71.0</b>	<b>68.1</b>	<b>60.4</b>
BLIP-Flickr30K	33.8	54.2	48.9	45.6
BLIP-COCO	32.8	51.4	51.4	45.2
BLIP-VQA	<b>47.8</b>	62.0	58.4	56.0
Random / Text-only	25.0	50.0	50.0	41.7
Human Estimate	100.0	97.3	99.0	98.8

# Visual Analogies



$I(\text{red mug}) - I(\text{yellow mug}) + I(\text{yellow bowl})$

$I(\text{red bowl}) ?$

61%

# Visual Analogies

$$I(\text{mug on table}) - I(\text{mug under table}) + I(\text{bowl under table}) \\ I(\text{bowl on table}) ?$$

9%

# Paxion: Action Robustness for Video-Language Models

## Probing Task: Action Antonym (AA)



Original Video

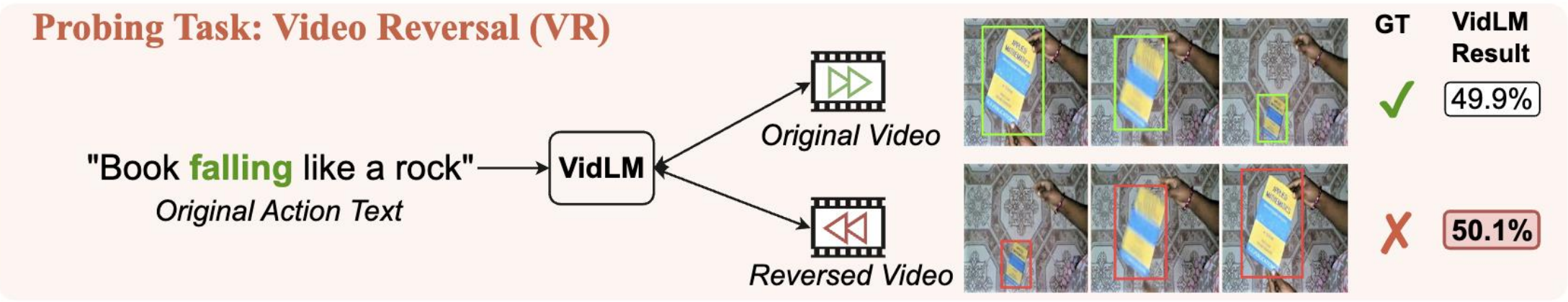
VidLM

"Book **falling** like a rock"  
*Original Action Text*

"Book **rising** like a rock"  
*Action Antonym Text*

GT	VidLM Result
✓	23.2%
✗	76.8%

# Paxion: Action Robustness for Video-Language Models



# Paxion: Action Robustness for Video-Language Models

## Baseline Task: Object Replacement (OR)



Original Video

VidLM

"**Book** falling like a rock"  
*Original Action Text*

"**Cellphone** falling like a rock"  
*Object Replaced Text*

GT

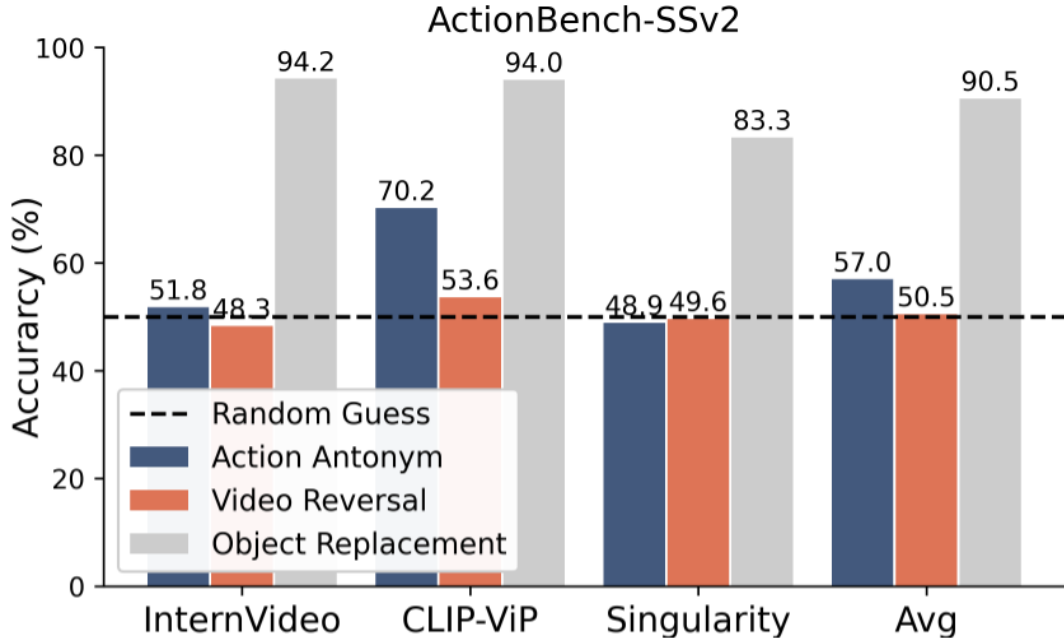
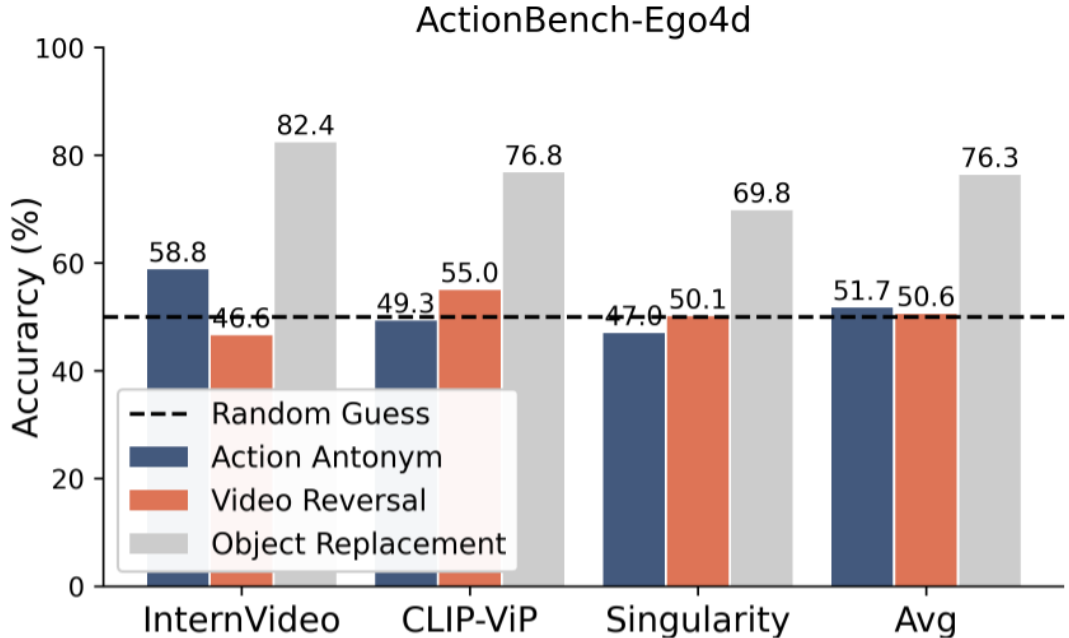


VidLM Result  
**77.9%**



**22.1%**

# Results



# Qwen Series



<https://qwen.ai/>

**Qwen3.5: Towards Native Multimodal Agents**

**Qwen3.5: Towards Native Multimodal Agents**

We are delighted to announce the official release of Qwen3.5, introducing the open-weight of the first model in the Qwen3.5 series, namely Qwen3.5-397B-A17B. As a native vision-language model,...

Open-Source | 2026/02/15

**Qwen-Image-2.0: Professional infographics, exquisite photorealism**

**Qwen-Image-2.0: Professional infographics, exquisite p...**

We are launching Qwen-Image-2.0, a next-generation foundational image generation model. The key highlights of Qwen-Image-2.0 include: Professional Typography Rendering: Supports 1k-token instructions for...

Open-Source | 2026/02/09

**Qwen3-Coder-Next: Pushing Small Hybrid Models on Agentic Coding**

**Qwen3-Coder-Next: Pushing Small Hybrid Models on A...**

--- We introduce Qwen3-Coder-Next, an open-weight language model designed specifically for coding agents and local development. Built on top of Qwen3-Next-80B-A3B-Base, which adopts a novel architecture...

Open-Source | 2026/02/02

**Qwen3-ASR & Qwen3-ForcedAligner is Now Open Sourced: Robust, Streaming and Multilingual!**

**Qwen3-ASR & Qwen3-ForcedAligner is Now Open Sour...**

Qwen3-ASR family includes two powerful all-in-one speech recognition models and a novel non-autoregressive speech forced alignment model. Qwen3-ASR-1.7B and Qwen3-ASR-0.6B are ASR models that support...

Open-Source | 2026/01/28

Widely used in academic research  
Multilingual and multimodal support

My personal take: Probably the strongest  
open-weight LLM right not