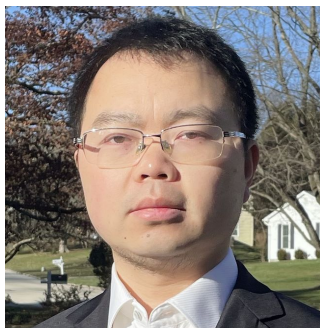
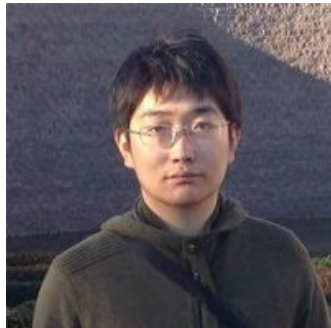


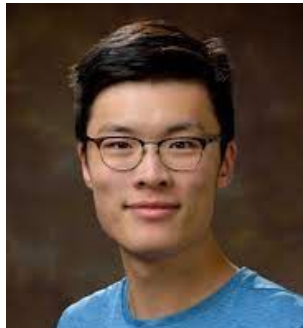
Indirectly Supervised Natural Language Processing



Wenpeng Yin



Muhao Chen



Ben Zhou



Qiang Ning



Kai-Wei Chang



Dan Roth

July 2023
ACL Tutorials
Indirectly Supervised Natural Language Processing



Indirect Supervision for Multi-modality Learning

Kai-Wei Chang

Department of Computer Science

University of California, Los Angeles

July 2023

ACL Tutorials

Indirectly Supervised Natural Language Processing

How to teach a model to locate “paramedics”



Visual Recognition with Language Descriptions



From ChatGPT:

Paramedic typically wears a uniform with **a patch or logo** identifying them as a member of an emergency medical services (EMS) team. The uniform may include a **shirt, pants, and jacket** with reflective strips for visibility. They may also wear protective gear such as **gloves, goggles, and a mask**. They often carry equipment such as a backpack with medical supplies, a radio, and a defibrillator. They may also wear a duty belt with a flashlight, scissors, and other tools. The appearance can vary depending on the agency.

Situated visual recognition

- ❖ How to teach a model to locate **the car in accident?**



Limitations of Supervised Data

- ❖ Detection data have **limited categories**
 - ❖ Objects365: 365 categories
 - ❖ LVIS: ~1,200 categories
 - ❖ Visual Genome: ~1,600 categories
- ❖ Detection data have **limited images**
 - ❖ Objects365 : < 1M
 - ❖ OpenImage : ~2M
- ❖ Hard to scale up because of annotation cost and long-tail distribution!

Learning Visual Concepts with Indirect Supervision

- Go beyond supervised object recognition
- Explicit labeling is expensive and incomplete
- How to leverage indirect supervision signals
 - Image/video caption
 - Rich descriptions in language
 - Unaligned text, images, and video

Outline

- Learning VL representation w/ indirect supervision
- Learning to recognize objects w/ image captions
- Learning to recognize objects w/ rich descriptions

Learning VL representation with indirect supervision

Motivating Example – Go Beyond Object Detection

Need to learn the association between **cake** with **birthday**



Q: What is this person doing?

A. He is celebrating birthday with his friends.

B. He is attending a formal banquet.

...



cake

birthday

Learning Associations from Image Captioning



People are making a **cake** for someone's **birthday** party.



cake

birthday



Q: What is this person doing?

A. He is celebrating birthday with his friends.

B. He is attending a formal banquet.

...

Learning Vision-and-Language Representation!



Harold (Liunian) Li



Mark Yatskar

Several people **walking** on a **sidewalk** in the **rain** with **umbrellas**.

Main training objective is to predict missing words.



VisualBERT

The model projects words and image regions into the same vector space and uses multiple Transformer layers to build joint representations.



Several people [MASK] on a [MASK] in the [MASK] with [MASK].



Input consists of an image and a caption with some masked words. Such data is easy to obtain from the internet.



Is it raining outside?

- a) Yes, it is snowing.
- b) Yes, [person8] and [person10] are outside.
- c) No, it looks to be fall.
- d) Yes, it is raining heavily.

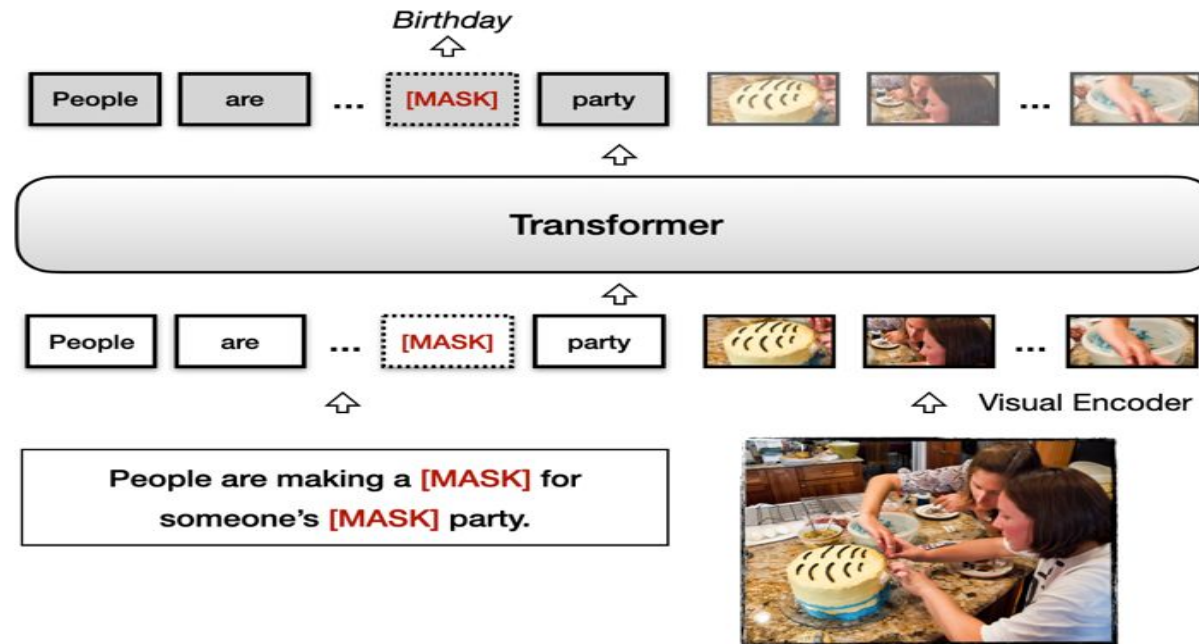
An example from the VCR dataset

Unsupervised pre-training on vision and language

Transfer to answering commonsense questions

Pre-training VisualBERT from Image Caption

Masked language modeling with the image



Sentence-Image Prediction



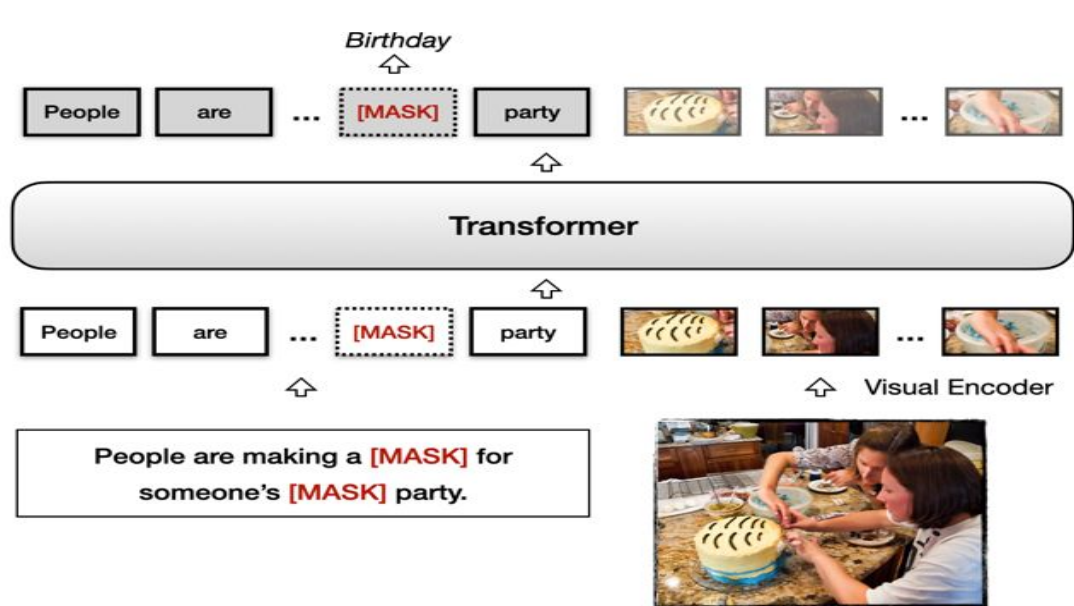
People are making a cake for someone's birthday party.

positive

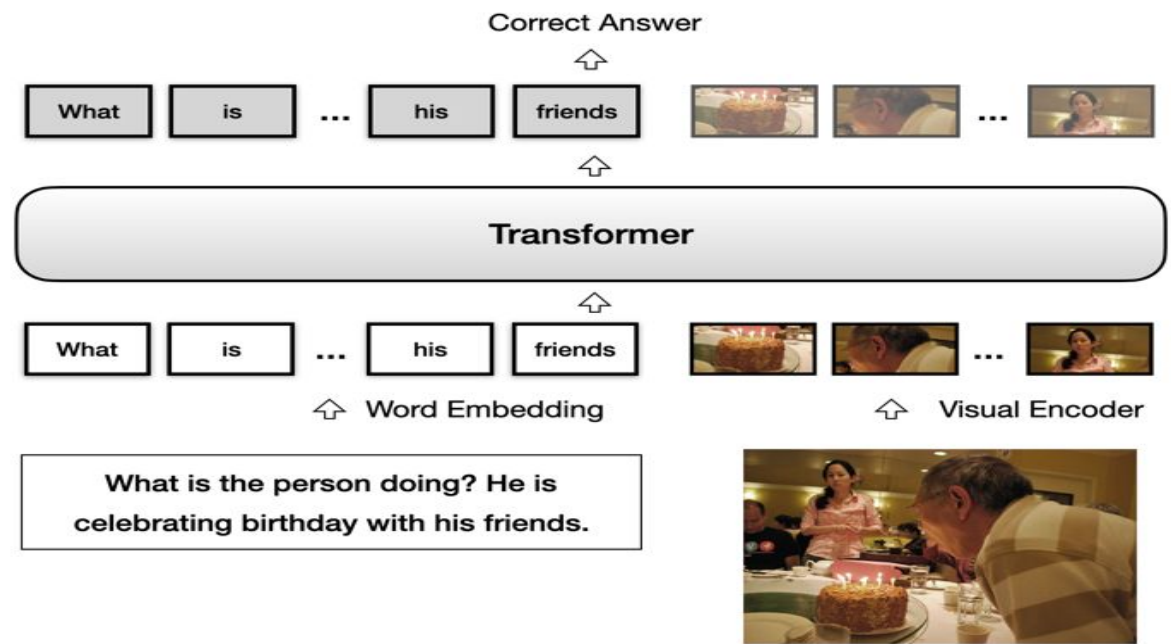
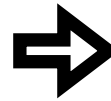
People are playing a ball in the park.

negative

Fine-Tuning VisualBERT

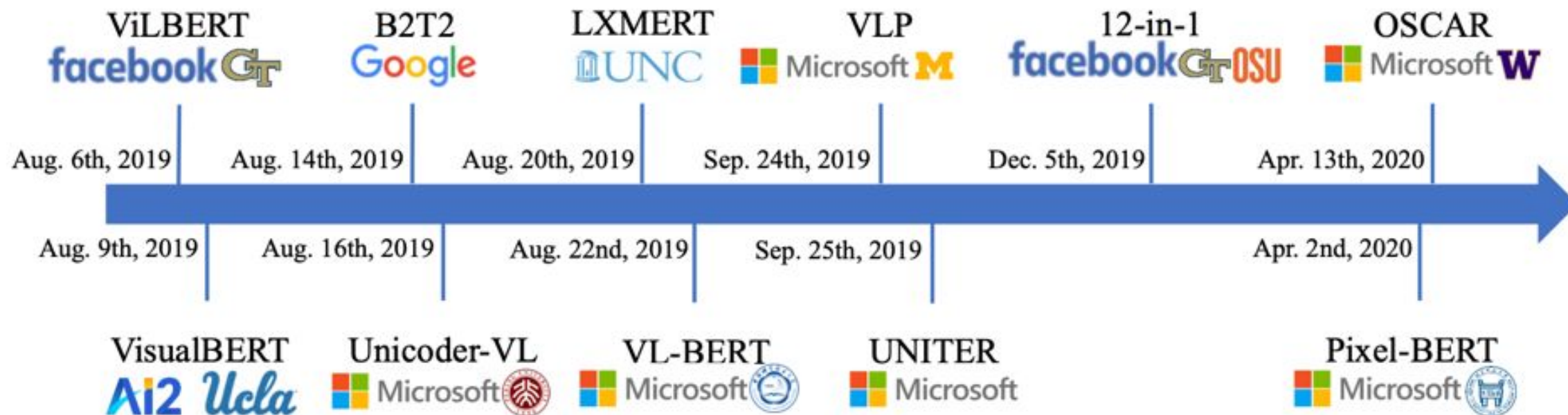


Pre-training (Representation Learning)



Fine-tuning (Solve a specific task)

Various Design Choices for VL Pre-Training



Learning VL representation with indirect supervision

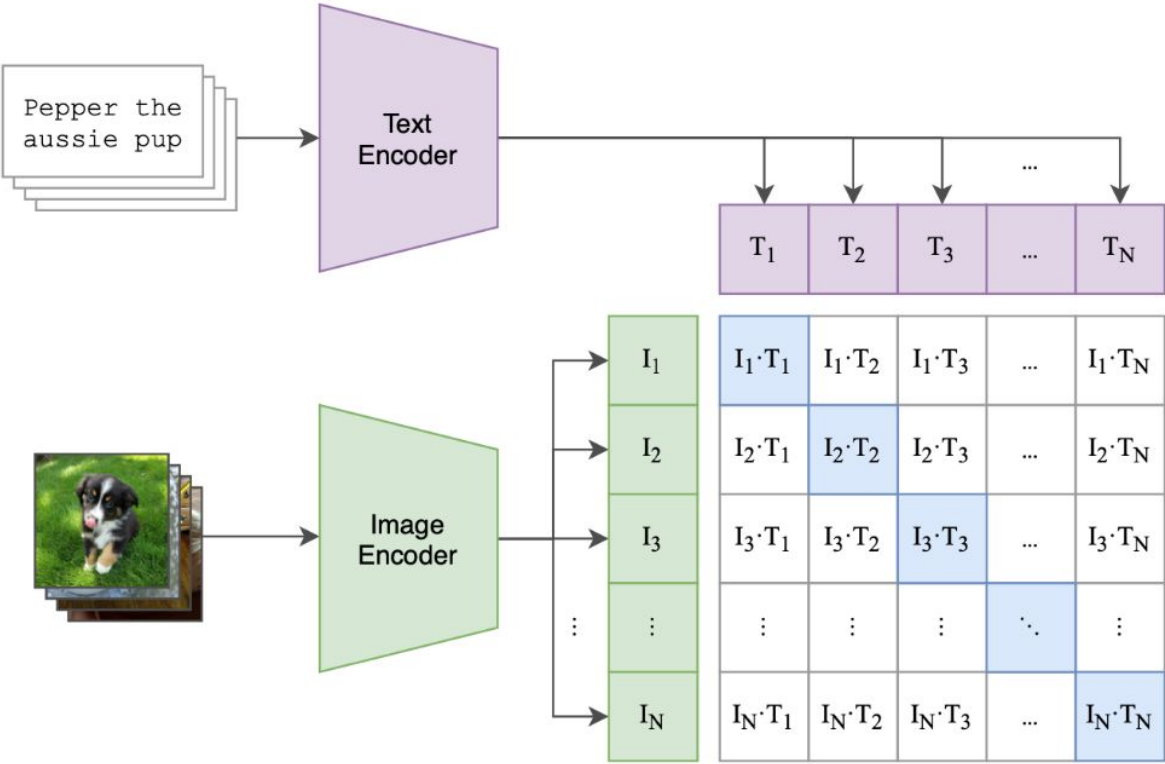
Image Classification w/ Vision&Langauge

- ❖ CLIP (Radford et al., 2021): image classification as image-text matching
- ❖ Leverage millions of semantic-rich image-caption data available on the web

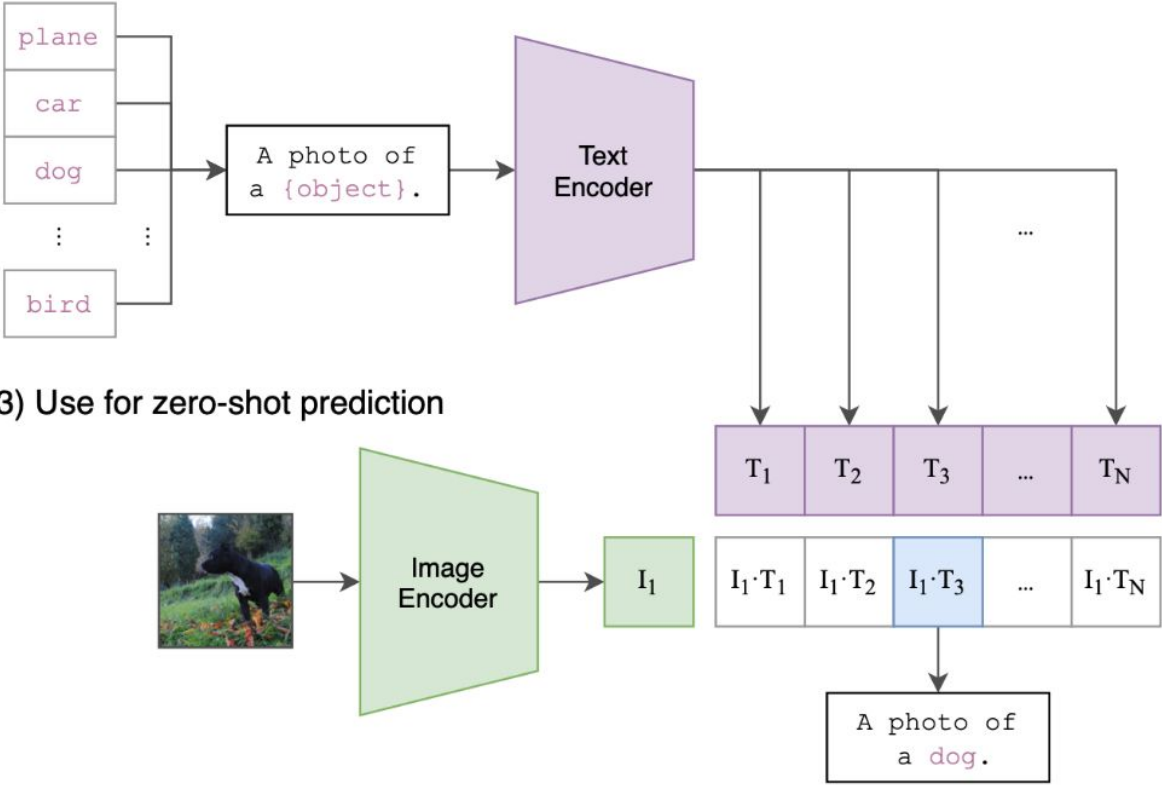


Constructive Language-Image Pre-training (CLIP) (Radford et al., 2021)

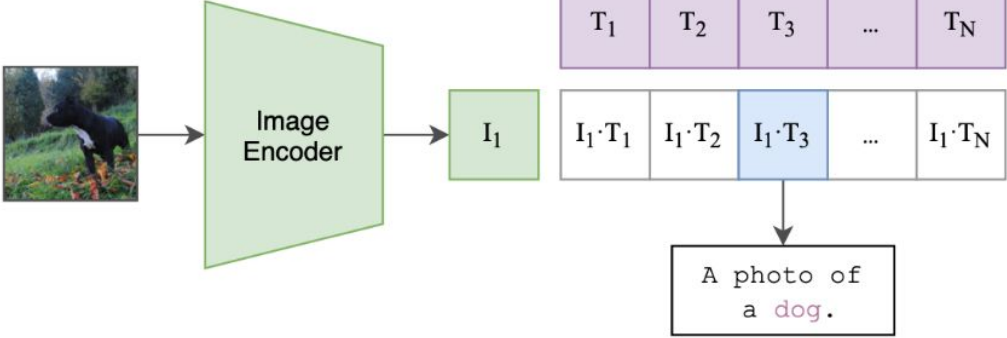
(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction



Language-Based Recognition Models: Image Classification

Food101

guacamole (90.1%) Ranked 1 out of 101 labels



✓ a photo of **guacamole**, a type of food.

✗ a photo of **ceviche**, a type of food.

✗ a photo of **edamame**, a type of food.

✗ a photo of **tuna tartare**, a type of food.

✗ a photo of **hummus**, a type of food.

FGVC Aircraft

Boeing 717 (45.6%) Ranked 2 out of 100 labels



✗ a photo of a **mcdonnell douglas md-90**, a type of aircraft.

✓ a photo of a **boeing 717**, a type of aircraft.

✗ a photo of a **fokker 100**, a type of aircraft.

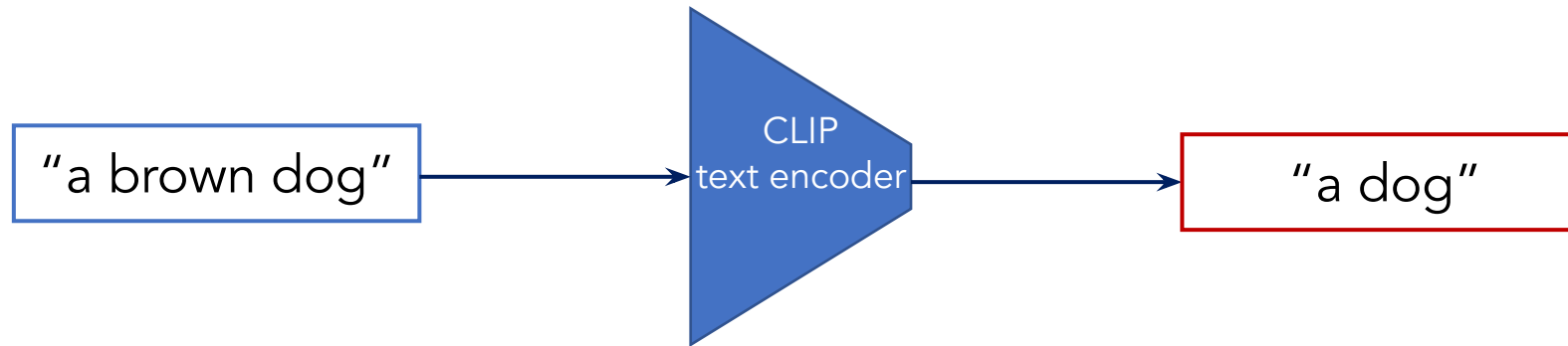
✗ a photo of a **mcdonnell douglas dc-9-30**, a type of aircraft.

✗ a photo of a **boeing 727-200**, a type of aircraft.

Use CLIP with a simple query: a photo of {class name}

Text encoders are performance bottlenecks

❖ What if the text encoder isn't perfect



Amita Kamath

You can't cram the meaning of a single \$&!#* sentence into a single \$!#&* vector!



Professor Raymond J. Mooney

Text Encoders are Performance Bottlenecks in Contrastive Vision-Language Models

Amita Kamath, Jack Hessel, and Kai-Wei Chang, in *Arxiv*, 2023.

Probing text encoder in CLIP

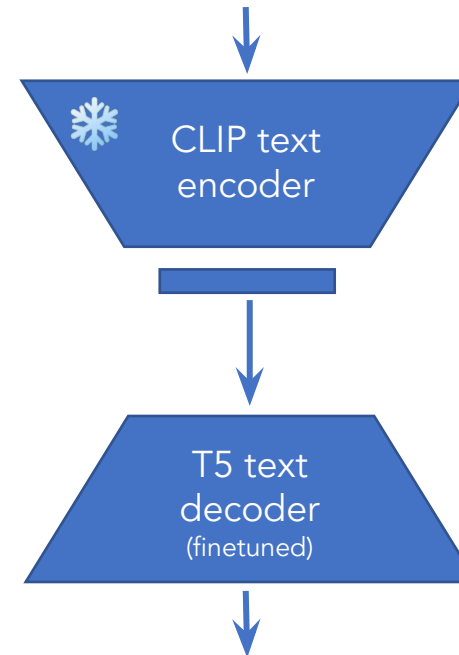
1. Create increasingly compositional text prompts

2. Feed them into CLIP's text encoder

3. Try to decode out the original prompt

Less than 30% captions
can be reconstructed

an iguana
a happy dinosaur
a surfer carrying a lifeguard
an orangutan eating and an officer flying



an iguana
an amusing dinosaur
a lifeguard carrying a surfer
an orangutan and an officer eating an orangutan

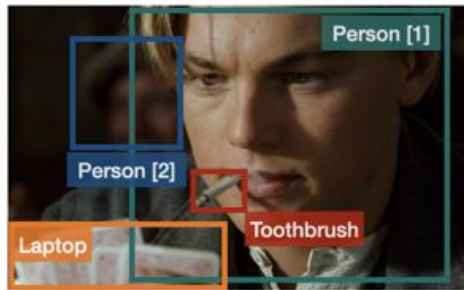
✓
✓
✗
✗21

GLIP: Object Detection as Phrase Grounding

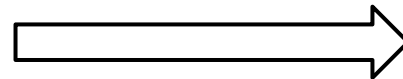
Phrase Grounding : Given a sentence and an image, locate the entities in the image



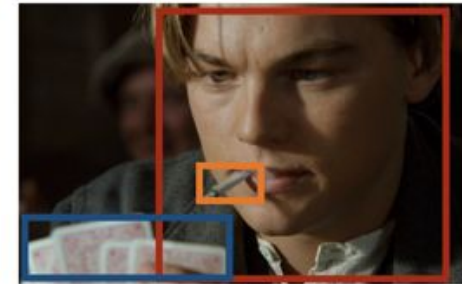
Harold (Liunian) Li



Detection Model
(Faster RCNN)



Object detection as grounding a language instruction



Unified Commonsense Model



Detect: Person, Cat, Dog, ...,
Poker, ..., Cigarette, ...

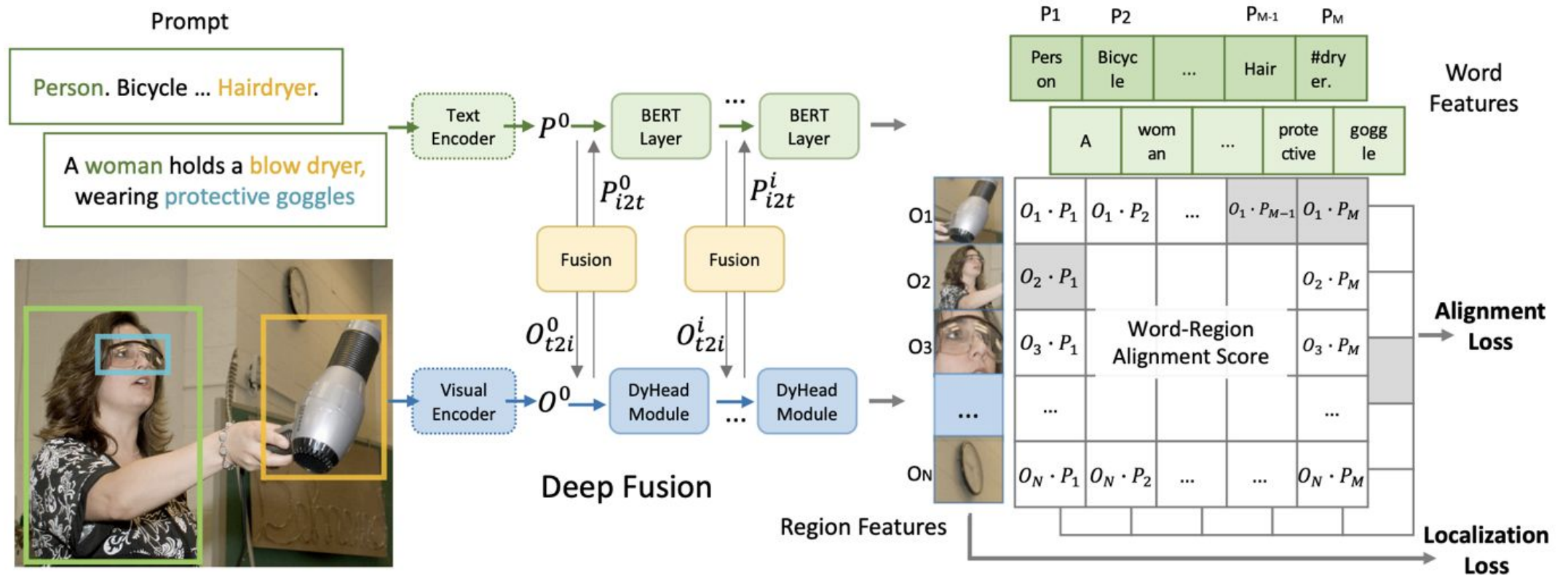


Grounded Language-Image Pre-training

Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, Kai-Wei Chang, and Jianfeng Gao, in CVPR, 2022.

GLIP: Overview

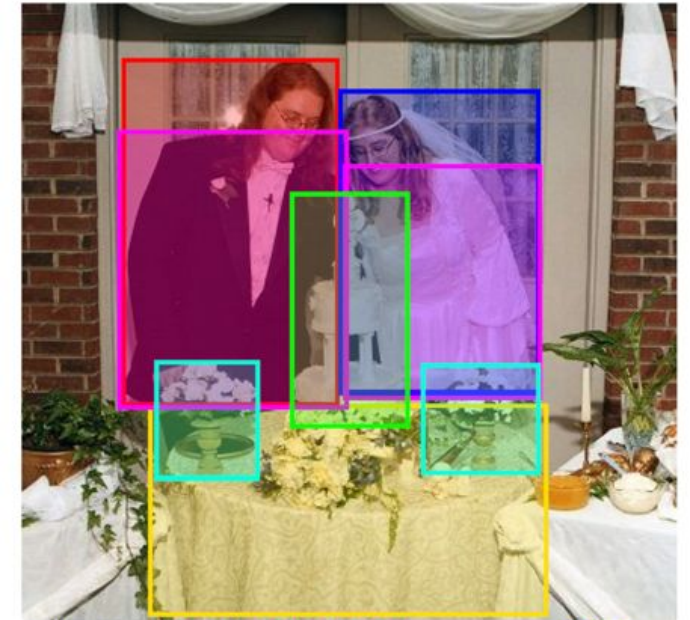
❖ Align objects to phrases in text



$$O = \text{Enc}_I(\text{Img}), P = \text{Enc}_L(\text{Prompt}), S_{\text{ground}} = OP^T,$$

Pre-training with Scalable Semantic-Rich Data

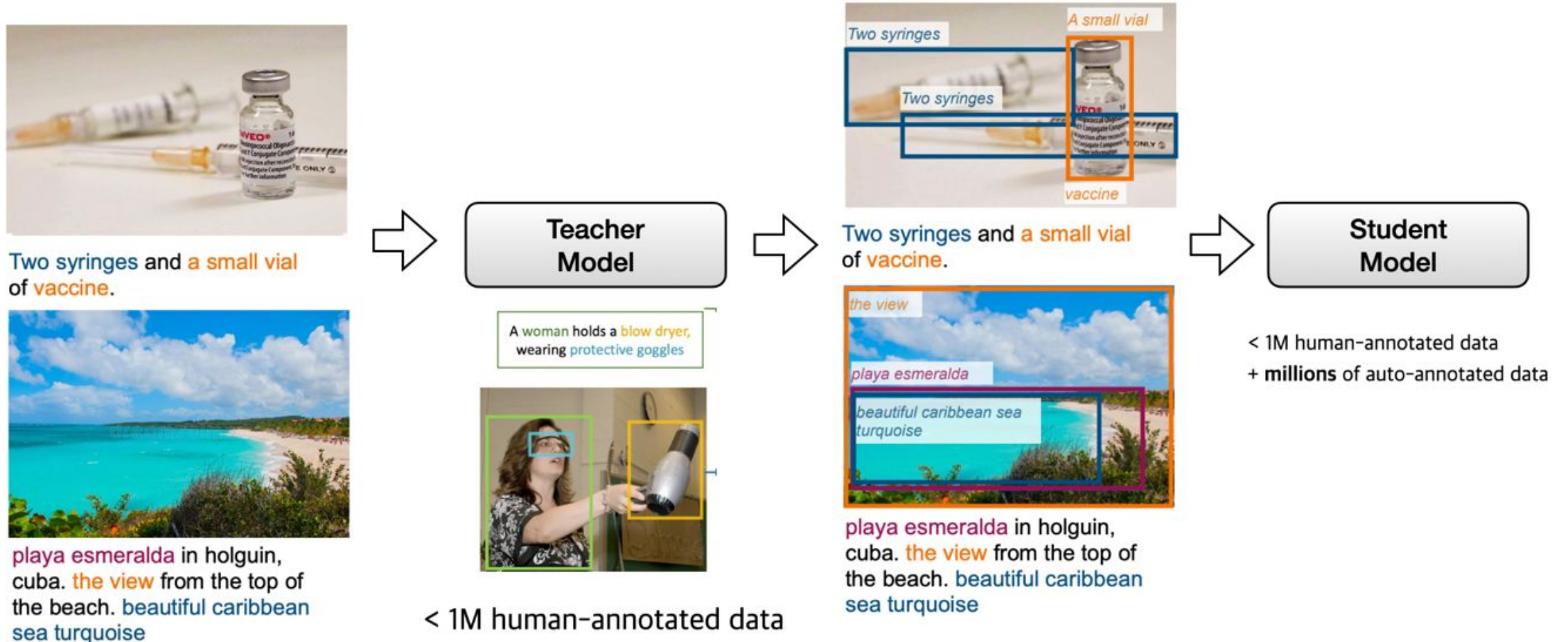
- ❖ Grounding data are **semantic-rich** and **scalable**
- ❖ **Gold grounding data:**
 - ❖ Flickr30K has **44,518** unique phrases
 - ❖ VG Caption has **110,689** unique phrases
 - ❖ 0.8M grounding data > 2M detection data



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**.

Scaling up with image-caption web data

Scalable training with 24M self-supervised web data

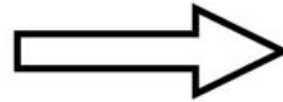


Pseudo grounding data (distant supervision)

- ❖ Train a teacher model with gold grounding data; produces boxes given image-caption data

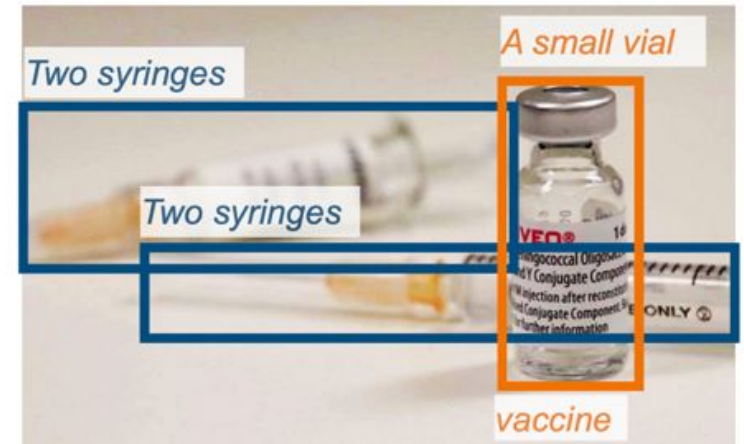


Two syringes and a small vial of vaccine.



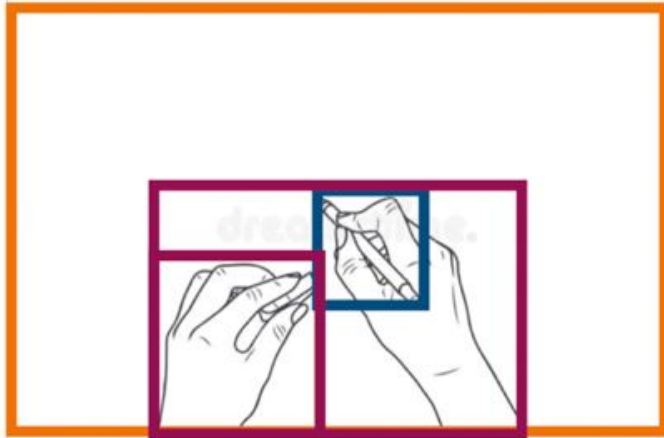
Teacher GLIP

Trained on gold detection & grounding data



Two syringes and a small vial of vaccine.

Pseudo grounding example



sketch illustration - female hands write with a pen. arm, art, background, black, care, concept, counting, design, drawing, finger, fingers, five, gesture royalty free illustration

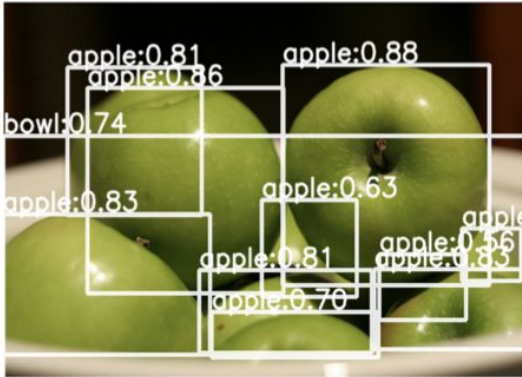


hard times teach us valuable lessons. handwriting on a napkin with a cup of coffee stock photos



save the straws classic t-shirt

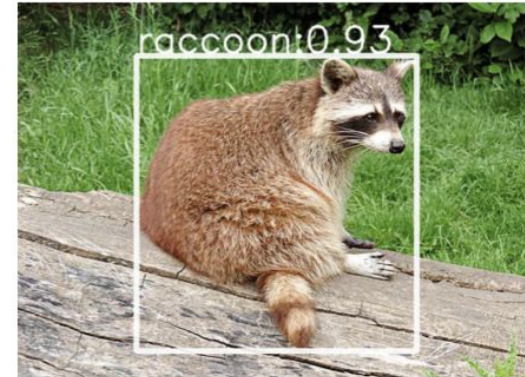
Language-Based Recognition Models: Object Detection



Prompt : person. bicycle.
car. motorcycle...



Prompt : aerosol can...
lollipop... pendulum...



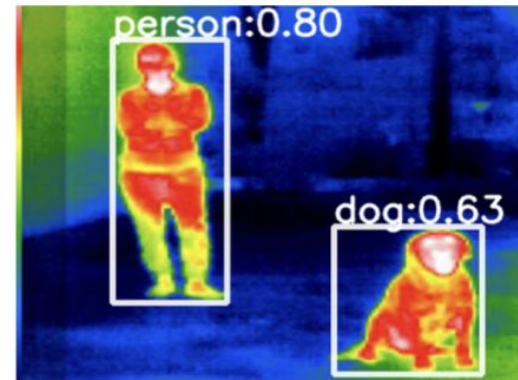
Prompt : raccoon



Prompt : pistol



Prompt : there are some
holes on the road

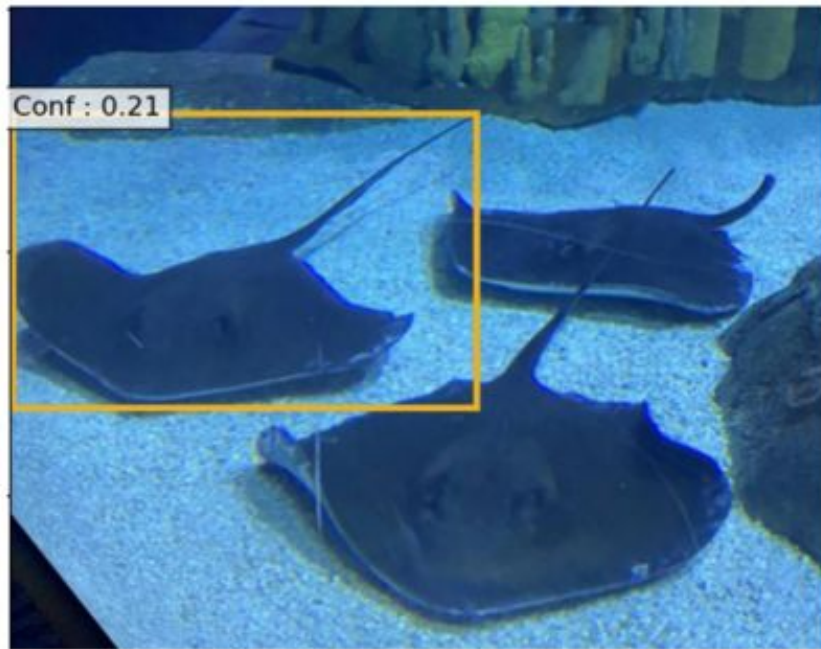


Prompt : person. dog.

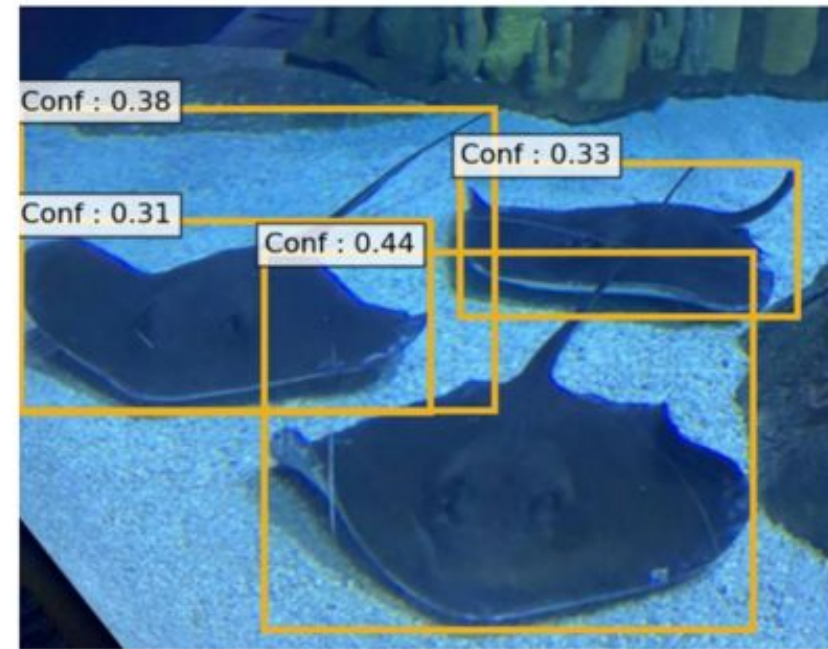
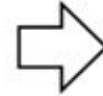
<https://huggingface.co/spaces/haotiz/glip-zeroshot-demo>

Object Detection with Instructions

- ❖ Learn from human instructions on the fly

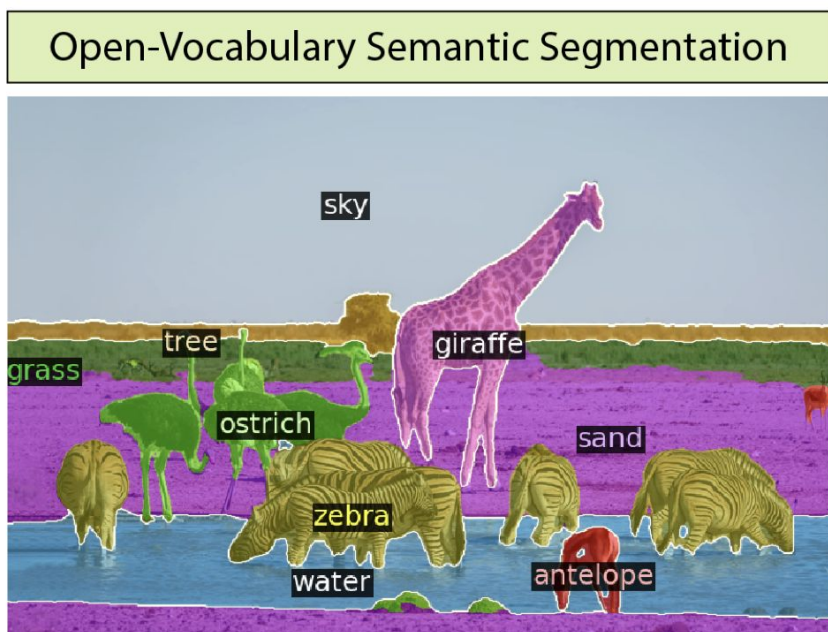


Prompt: ... **stingray** ...



Prompt: ... **stingray**, which is flat and round...

Language-Based Recognition Models: Segmentation

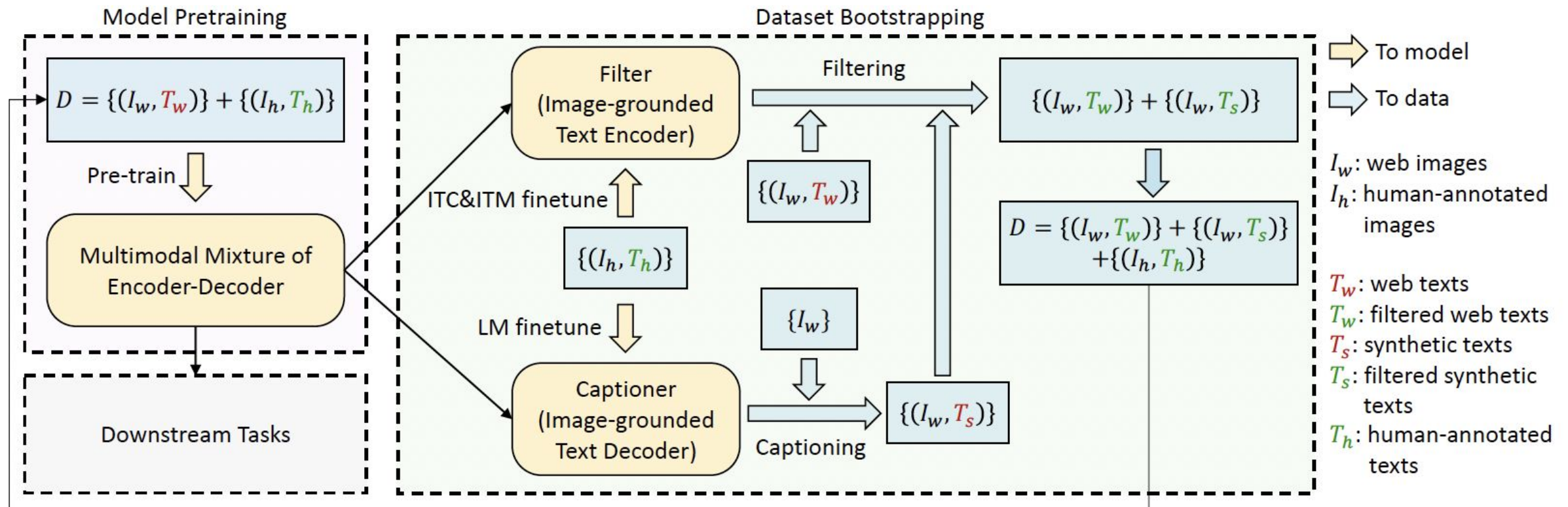


Segment: zebra. ostrich. water. ...

Generalized Decoding for Pixel, Image, and Language

Xueyan Zou, Zi-Yi Dou, Jianwei Yang, Zhe Gan, Linjie Li, Chunyuan Li, Xiyang Dai, Harkirat Behl, Jianfeng Wang, Lu Yuan, Nanyun Peng, Lijuan Wang, Yong Jae Lee, Jianfeng Gao

Bootstrapping Language-Image Pre-training (BLIP)



BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation

Junnan Li, Dongxu Li, Caiming Xiong, Steven Hoi

Learning to recognize objects w/ rich descriptions

How to Deal with Complex Language Queries?

1. Generalizing to novel categories
(object detection / segmentation in the wild)

Detect: **mallet**

v.s.

Detect: **mallet**, a kind of tool, wooden handle with a round head, used for pounding or hammering

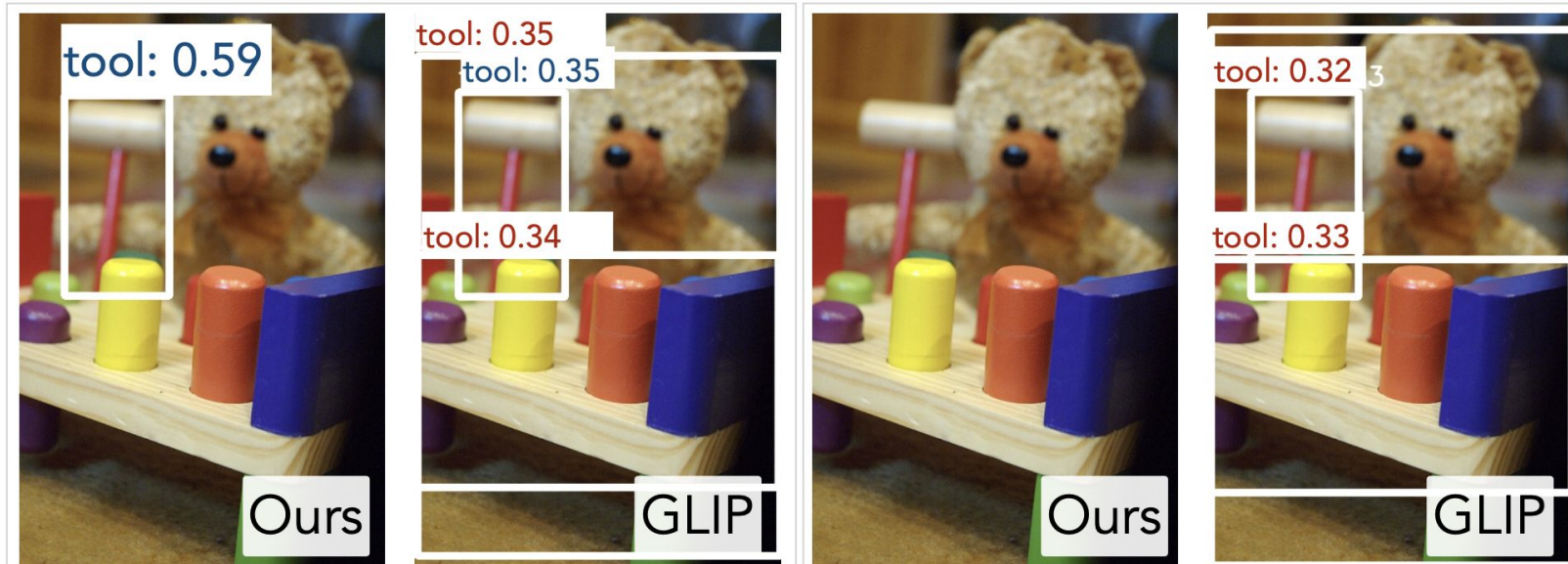


Limitation of VL Pre-Training Models

Detect with specifications for shape & subpart
(w/o object name)

Target Object

Confusable Object



A kind of tool, wooden handle with a round head, used for pounding or hammering





A kind of tool, long handle, sharp blade, could be used for chopping wood

Limitation of VL Pre-Training Models

Detect with specifications for relation

Target Object

Confusable Object

 <p>clown: 0.72</p> <p>Ours</p>	 <p>clown: 0.74</p> <p>GLIP</p>	 <p>clown: 0.71</p> <p>Ours</p>	 <p>clown: 0.33</p> <p>GLIP</p>
<p>A <u>clown</u> making a balloon animal for a pretty lady</p>		<p>A <u>clown</u> kicking a soccer ball for a pretty lady</p>	

Challenge #1: Fine-Grained Descriptions Rare in Pre-Training Data

A toy bear holding
a mallet.



Reporting bias: people do not write obvious things

When writing captions, we tend to directly use entity names rather than give descriptions for subparts, shapes, textures, etc.

The World of an Octopus: How Reporting Bias Influences a Language Model's Perception of Color

Cory Paik, Stéphane Aroca-Ouellette, Alessandro Roncone, Katharina Kann

Solution #1: Generating Descriptions from Large Language Models

User:

What are the useful features for identifying **mallet**?

A toy bear holding
a mallet.



GPT-3:

Mallet is a kind of tool, wooden handle, ...

Build a vocabulary of 10K noun phrases on Conceptual Captions and VG

Sample descriptions for each noun phrase (<1 day via API)

VISUAL CLASSIFICATION VIA DESCRIPTION FROM
LARGE LANGUAGE MODELS

Sachit Menon, Carl Vondrick

**ELEVATER: A Benchmark and Toolkit for Evaluating
Language-Augmented Visual Models**

Chunyuan Li^{1*}, Haotian Liu², Liunian Harold Li³, Pengchuan Zhang¹, Jyoti Aneja¹,
Jianwei Yang¹, Ping Jin¹, Houdong Hu¹, Zicheng Liu¹, Yong Jae Lee², Jianfeng Gao¹

Challenge #2: Model Might Ignore Description

The model is not incentivised to “read” the descriptions

It tends to learn using the following shortcuts:

- Entity shortcut

- Positive query shortcut

Entity Shortcut



Contrastive learning objective -> distinguishing two sub-queries:

Q1: **Mallet**, which has a wooden handle with a round head, used for pounding or hammering

Q2: **Ax**, which has a long handle and a sharp blade, could be used for chopping wood

Which query aligns to the image?

Entity Shortcut



Do you remember the descriptions?

Q1: **Mallet**, which has a wooden handle with a round head, used for pounding or hammering

Q2: **Ax**, which has a long handle and a sharp blade, could be used for chopping wood

Entity Shortcut: focusing on center entities is enough to distinguish Q+ and Q-

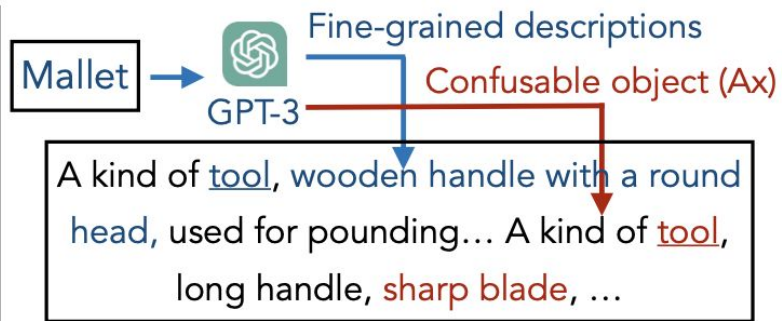
Solution for Entity Shortcut: Context-Sensitive Query





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

Description-rich and context-

Q1: a kind of **tool**, which has a wooden handle with a round head

Q2: a kind of **tool**, which has a long handle and a sharp blade

The labels of the word "tool" now depends on its language context

Positive Query Shortcut



Detect: Mallet.

Bear. Cat...

A toy bear holding
a mallet.

The imperfect unification of detection and grounding data

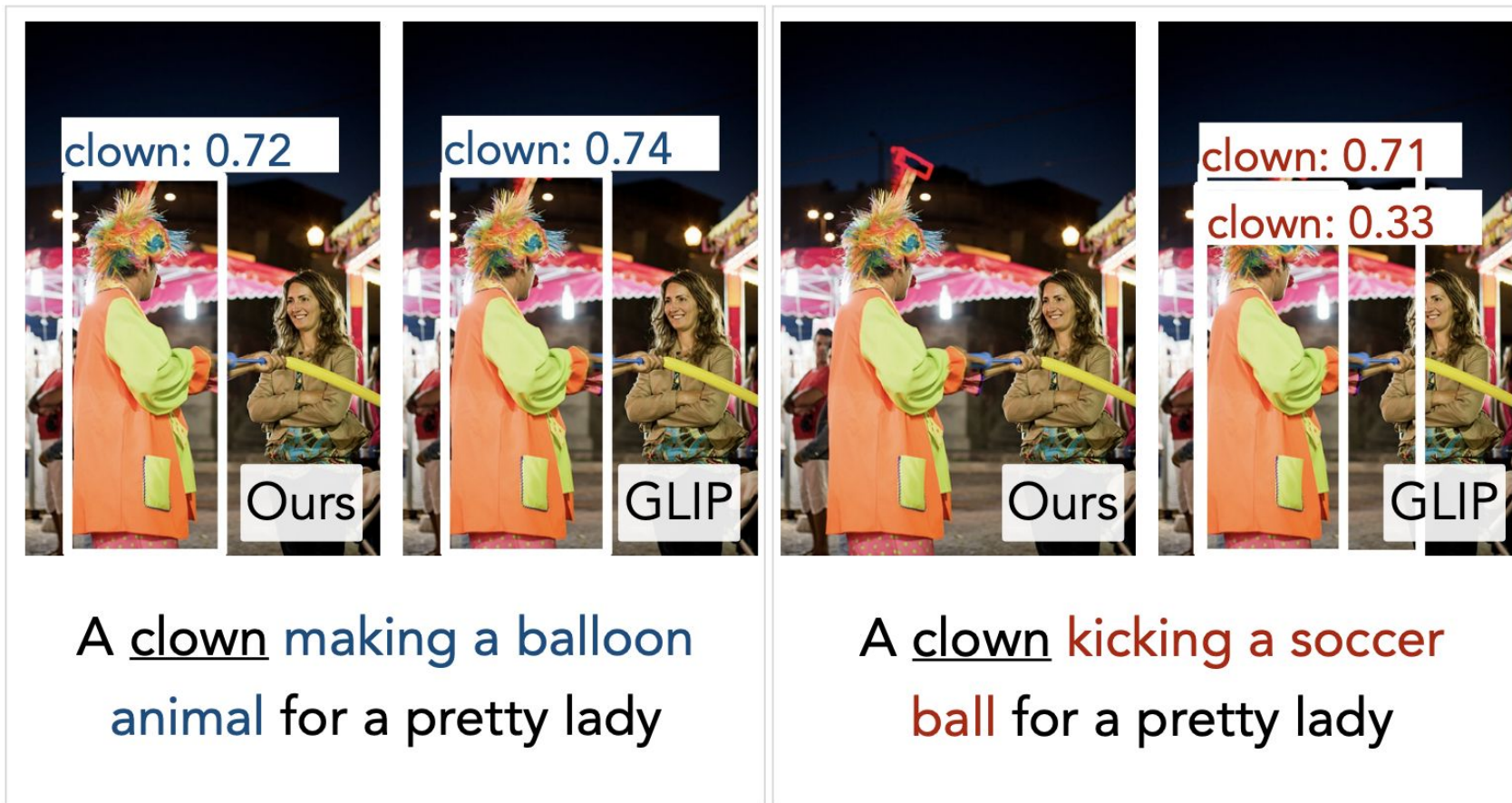
Detection task: target object may not in image

Grounding data: query is a caption that always corresponds to the image

Positive Query Shortcut: language-like query is always positive

Potential problem of many *deep-fused* VL models (MDETR, GLIP, FIBER, ...)

Positive Query Shortcut



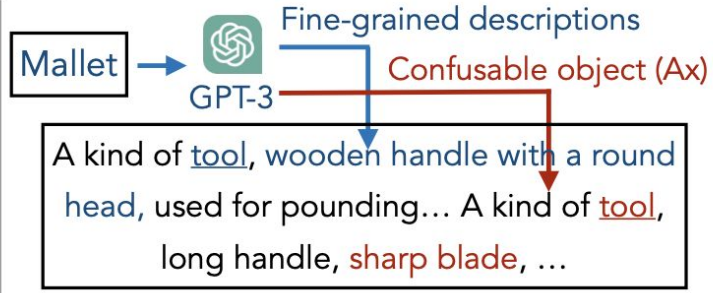
Solution #2: Context-Sensitive Query

Detect: Mallet.
Bear. Cat...

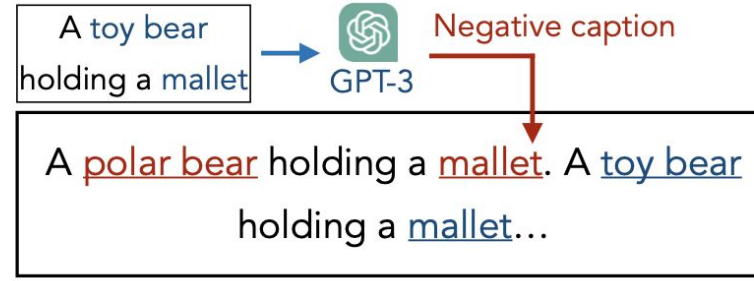
A toy bear holding a mallet.



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

Labels of "mallet" depends on context as well

In contrast to classical detection training (where we drop the examples with too few boxes), we need to have full-negative queries to teach the model to suppress false positives

Summary: Language Descriptions as a Supervision Signal

Label as supervision (ImageNet/COCO):

- Costly to annotate

Caption as supervision (CLIP/ALIGN/ViLD/GLIP):

- + Easy to scale
- Learns mostly object-entity alignment

Description as supervision:

- + Easy to scale
- + Decompose recognition of objects into grounding of attributes, parts, shapes, etc
- + More flexible language queries

otter

black: yes
white: no
brown: yes
stripes: no
water: yes
eats fish: yes



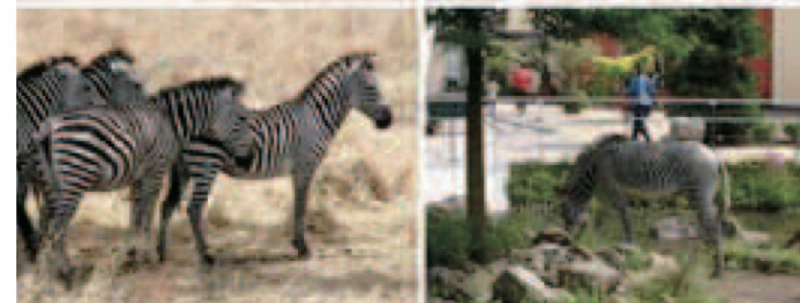
polar bear

black: no
white: yes
brown: no
stripes: no
water: yes
eats fish: yes



zebra

black: yes
white: yes
brown: no
stripes: yes
water: no
eats fish: no



Attribute-Based Classification for
Zero-Shot Visual Object Categorization