

# **Weak Multimodal Supervision for Robust Visual Recognition**

**Adriana Kovashka**

Associate Professor

Department of Computer Science  
University of Pittsburgh



# (Original) Research Motivation

- Media content affects public opinion and societal outcomes (e.g. via elections)
- We want to understand agenda/intent in the media
  - Images in political articles, visual advertisements
- Data is limited, supervision is noisy
  - Images with intent not as abundant as general images
  - Expensive to carefully annotate all knowledge required
  - Text in ads/articles helps understand images, but relation is indirect, abstract
- Can we learn useful models from noisy data?
- How to transfer knowledge from other domains?

# Outline

1. Visual reasoning for advertisements
  - Visual/language reasoning about ad strategies (ICCV 2021, ACL 2023)
2. Learning concepts from language supervision
  - Modeling complementary image-text relationships (BMVC 2018, ECCV 2020, TPAMI 2022)
  - Modeling variance and noise in object naming (EACL 2024, NAACL Findings 2024)
  - Investigating the role of language context (WACV 2024)
3. Recognizing objects with unexpected appearance
  - Across modalities (ICMR 2023)
  - Across geographies (CVPR 2024)

# Decoding image advertisements



## Recognized Concepts

people      commerce      street  
business      stock      city

## What Does It Mean?

Food at Burger King must taste really good since even competitor's employee secretly buys it.

## Image Captioning

A man standing in front of a display of food.

A man standing in front of a display of a store.

- State-of-the-art vision systems (were) inadequate to describe the messages hidden behind purposefully designed advertisements

# Challenges: Human associations (symbols)



# Challenges: Complex relation of image/text

- Is it an advertisement for ...



← Public Service Announcement  
Against Human Trafficking

“Winter Collection”

Commercial Advertisement →  
Premium Women Clothing



# Challenges: Atypical objects



# Challenges: Implied physical processes



# Advertisement dataset

We collected **64,832** advertisement images with:

- 204,340 **topic** annotations
- 102,340 **sentiment** annotations
- 202,090 **action/reason (what/why)** annotations
- 64,131 **symbolism** annotations
- 20,000 **strategy** annotations
- 11,130 **slogan** annotations

<http://cs.pitt.edu/~kovashka/ads>



# Decoding image advertisements

- What message does the ad convey (*action*), and what arguments does it provide for taking the suggested action (*reason*)?
- Multiple-choice task: Given  $k$  options for action-reason statements, pick one that matches the image



Correct action-reason statements:

- I should drink evian because it helps you recover
- I should buy Evian because it keeps us young
- I should drink Evian because it will keep me like a baby



# Shortcuts in Visual Commonsense Reasoning

**Shortcut effects: Example and definition.** Consider the following example from the VCR dataset (Zellers et al. 2019). In the figure, [person1] (male) is on the right and [person2] (female) is on the left.



Question: What does [person1] think of [person2]'s dress?

Correct answer: [person1] thinks [person2] looks stunning in her dress.

Incorrect #1: She does not approve.  
Incorrect #2: [person2] is a girl and girls like to wear makeup.

Incorrect #3: [person1] is confused and annoyed by [person2] following her in the store.

Table 5: Our method enables the most robust training. All results show Q→A except for the bottom two which show QA→R. The best method per group on Q→A is **bolded**, and the best method per task is underlined.

Method	STD VAL	RULE- SINGULAR	RULE- PLURAL	ADVTOP- 1
BASELINE (B2T2)	68.5	63.3	65.3	37.0
MASKING 0.05	<b>69.3</b>	<b>63.9</b>	<b>66.0</b>	48.8
MASKING 0.10	68.7	62.8	64.7	50.1
MASKING 0.15	68.2	62.0	63.3	<b>50.6</b>
MASKING 0.30	64.1	56.6	56.8	47.5
MASKING 0.05 + MLM	68.5	62.9	64.8	47.3
MASKING 0.10 + MLM	69.1	63.8	65.0	<b>50.6</b>
OURS-CL INIT0.30 DECAY1E-4	69.6	64.5	64.7	51.7
OURS-CL INIT0.30 DECAY5E-5	<b>69.9</b>	<b>65.9</b>	<b>66.8</b>	54.5
OURS-CL INIT0.50 DECAY1E-4	69.4	65.0	65.0	53.0
OURS-CL INIT0.50 DECAY5E-5	69.8	65.4	66.3	<b>54.9</b>
BASELINE (B2T2)	68.5	64.9	67.9	34.7
OURS-CL INIT0.30 DECAY5E-5	<b>70.6</b>	<b>66.6</b>	<b>70.4</b>	<b>47.9</b>

[val-54]

Original Val data

Q: Where is [2] going ?  
A0 [2] is going into the store .  
A1 [2] is getting into a carriage .

Modified by rule  
(A single person)

A0 He is going into the store .  
A1 [2] is getting into a carriage .

Modified by an  
adversarial model

A0 [MASK] is going into the store .  
A1 [2] is getting into a [MASK] .

A2 [1] is going to the bathroom .  
A3 [1] is going outside to play after the conversation with [2] is over .

X A2 [2] is going to the bathroom .  
A3 [1] is going outside to play after the conversation with [2] is over .

X A2 [MASK] is going to the bathroom  
A3 [1] is [MASK] outside to play after the conversation with [2] is over .

[val-270]

Original Val data

Q: What are [1, 2] feeling ?  
A0 [1, 2] do not like the restaurant .  
A1 They are apprehensive .

Modified by rule  
(A group of people)

A0 [1, 2] do not like the restaurant .  
A1 [1, 2] are apprehensive .

Modified by an  
adversarial model

X A0 [1, 2] do not like the [MASK] .  
A1 They are apprehensive [MASK]

A2 They are both feeling happy .  
A3 [1, 2] are feeling drunk .

A2 They are both feeling happy .  
A3 [1, 2] are feeling drunk .

A2 They are [MASK] feeling happy .  
A3 [1, 2] are feeling [MASK] .

# Domain-Robust Visual Question Answering

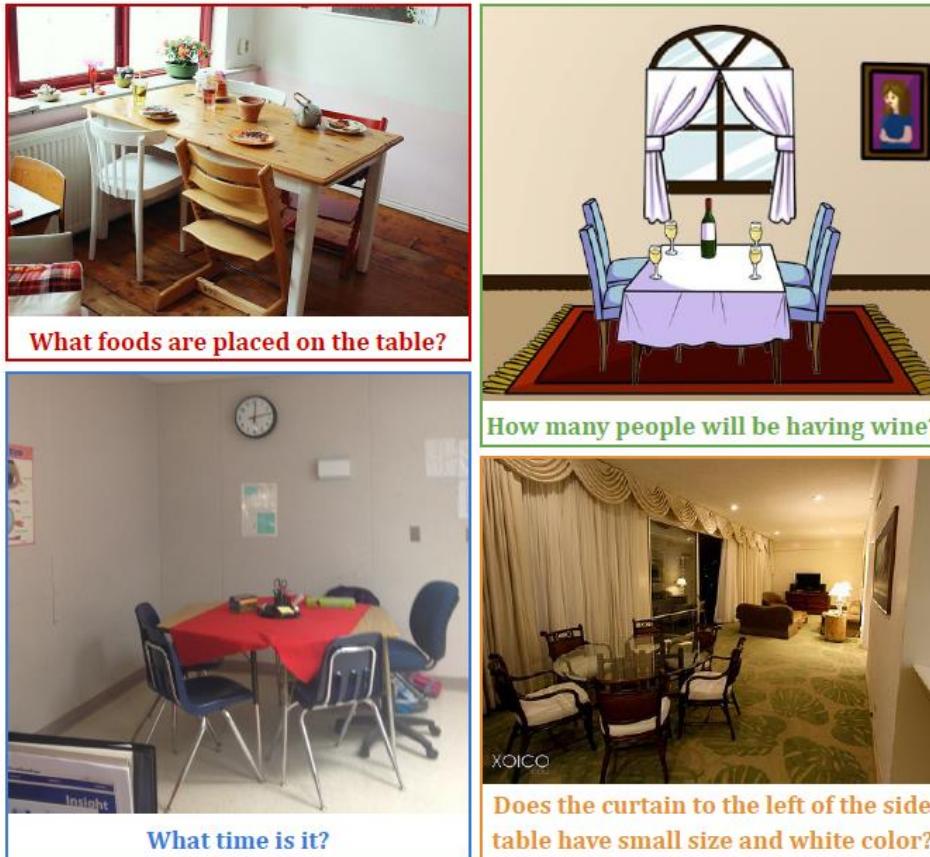


Figure 1. The same visual setting can be captured in different ways in VQA datasets, and paired with different information needs (questions). They may require deduction using visual contents, reading from a specific region of the image, or reasoning about complex spatial relationships. All examples are selected from real VQA datasets, *i.e.* **VQA v2**, **VQA Abstract**, **VizWiz** and **GQA**.

# Atypical images with persuasive intent

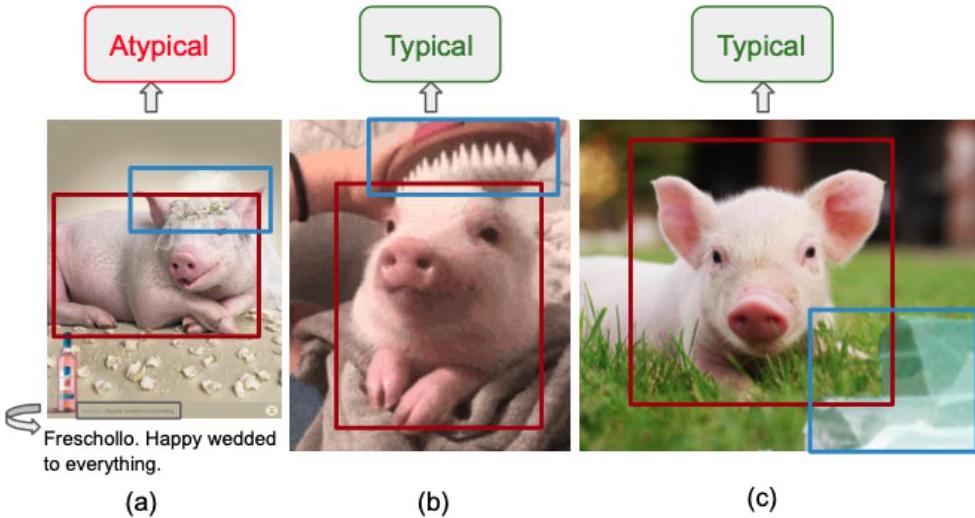


Figure 1. These images illustrate the importance of object interactions and their spatial relative position for atypicality detection. (a) Pig wearing a bridal veil is atypical; (b) If a handled brush instead of a veil is on top of the pig's head, then the image is typical; (c) If the veil's location is different, the image may also be typical.

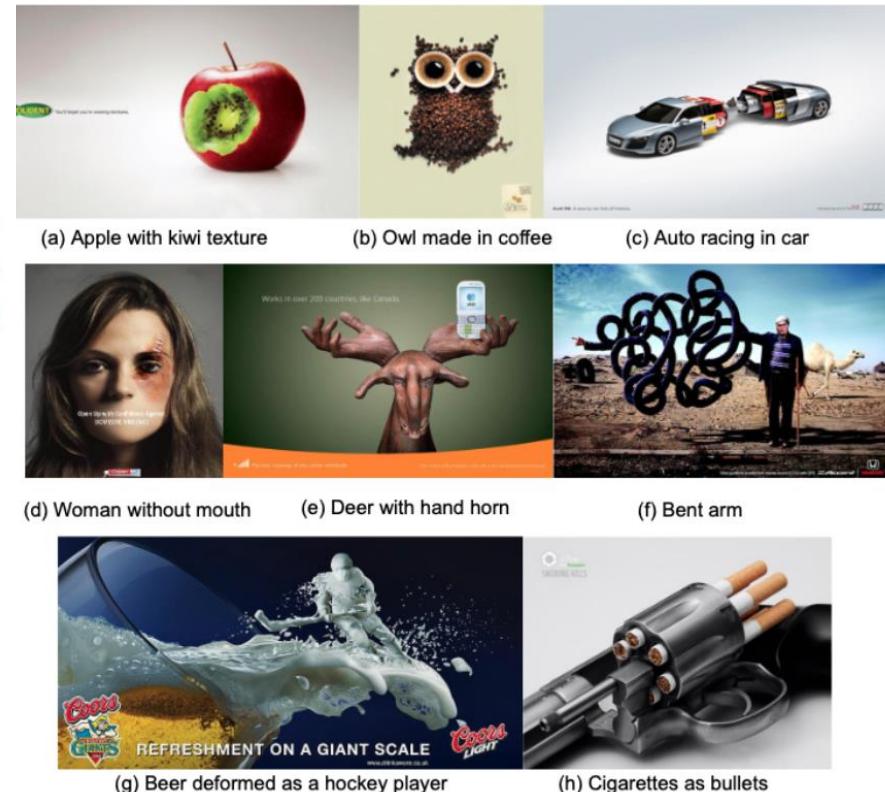


Figure 2. Atypical object transformations in ads; [48]'s dataset.

# Decoding Symbolism in Language Models

Relationship Type	Count	Example (signifier - signified)	Example (situated signifier - signified)
UsedFor	52	makeup - beauty	cartoon candy running on a treadmill - health
HasProperty	46	child - youth	workers sitting closely in a sofa - comfort
RelatedTo	47	mountain - adventure	cigarette smoke in the shape of mushroom cloud - danger
Others	94	chocolate - love	foot stepping on tombstone - death
Indirect	116	giraffe - love	shoes made out of red bull cans - strong

Table 2: Relationship types of *signifier-signified* in the set of advertising symbolism.

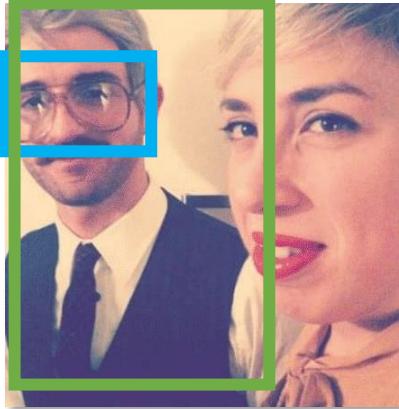
# Understanding the relation of image/text

Image Captioning



man in black shirt is  
playing guitar.

Visual Question Answering



Who is wearing glasses?  
man

Persuasive Ads using Visual Rhetoric

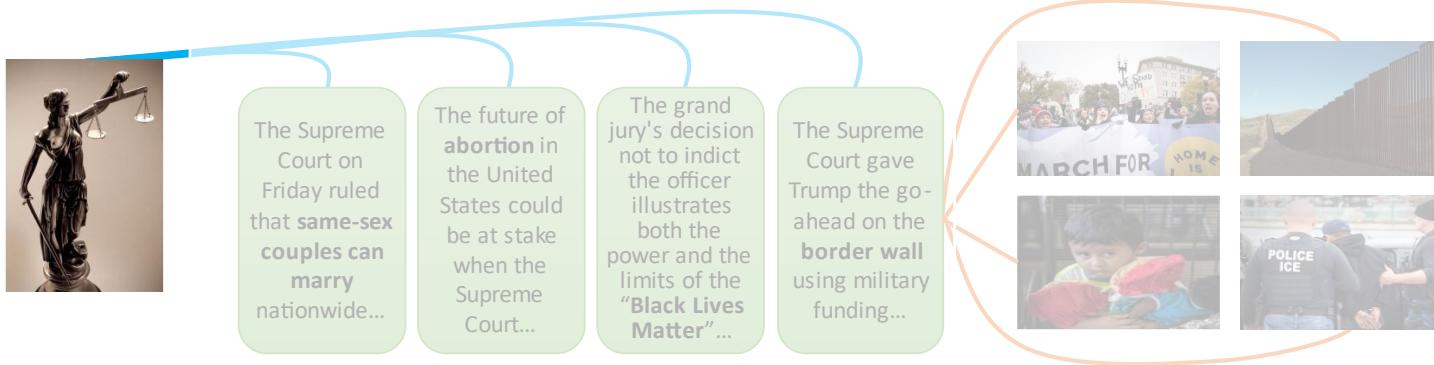


They can't afford to wait for evolution.  
Dolphin? Earphone?

Redundant messages are transmitted  
independently across channels.

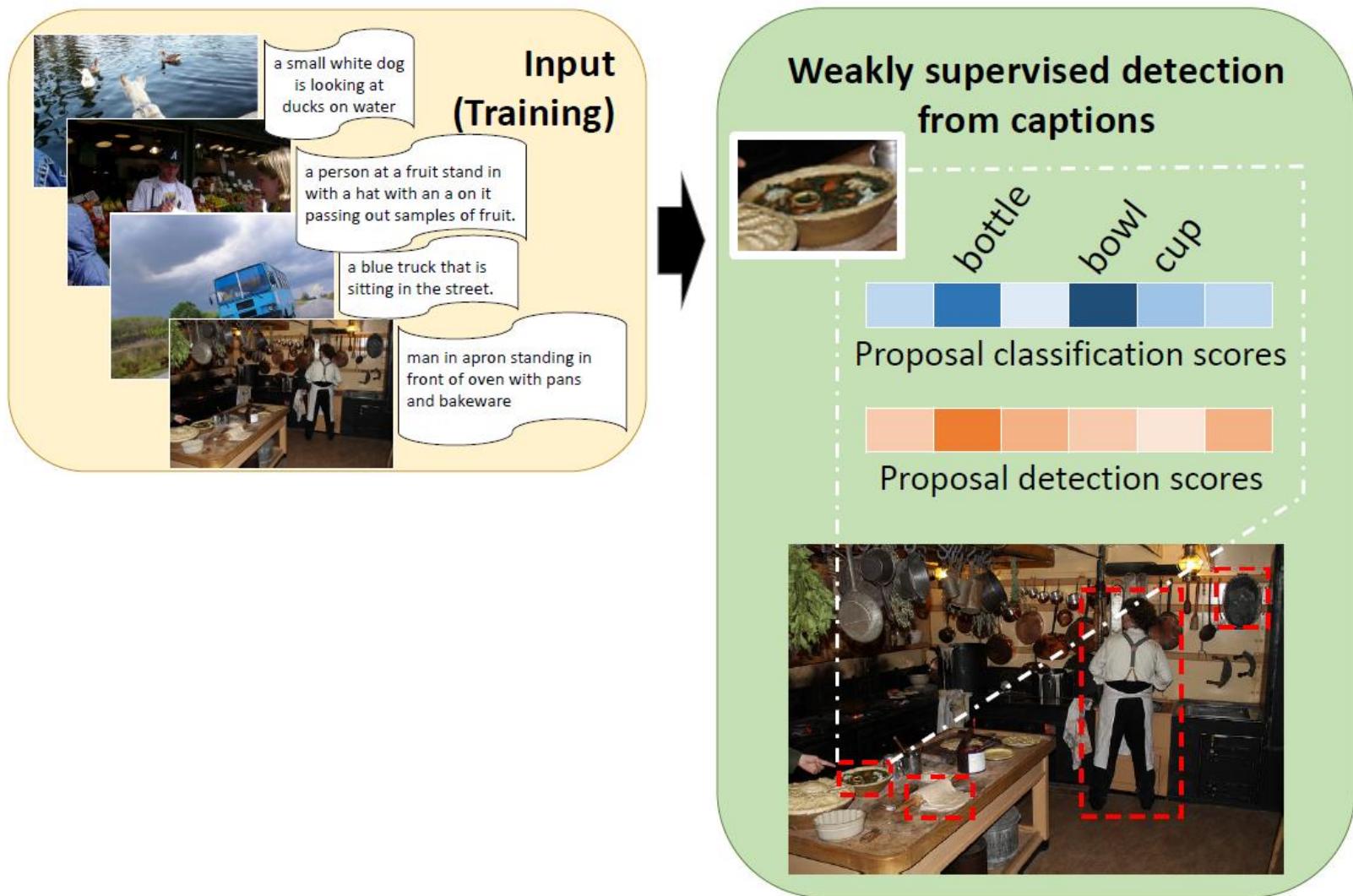
Elements across channels do not align  
or are connected in non-literal way.

# Multimodal complementarity



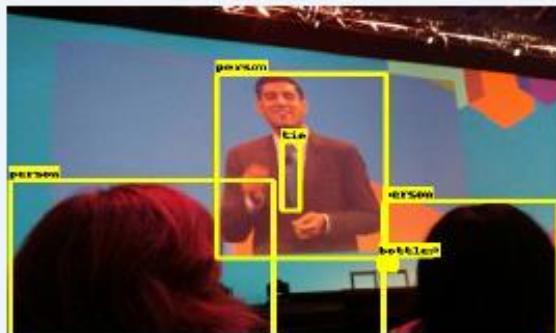
- Modalities complementary; each contributes to overall message the multimedia document seeks to convey
  - Figurative / symbolic usages (e.g. Themis, justice)
  - Abstract semantic concepts (e.g. immigration) can have diverse visual expressions

# Learning object detectors from captions



# Captions don't mention all objects present

- This causes a challenge for standard weakly-supervised detection methods

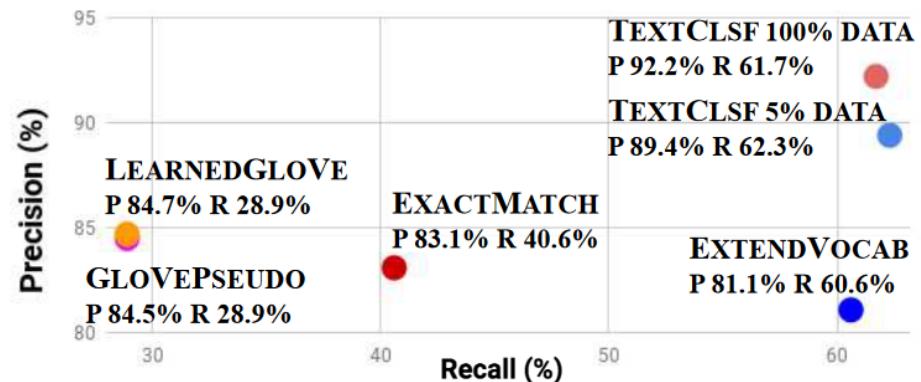


- People watch a man **delivering a lecture** on a screen
- A large screen showing a person **wearing a suit**
- An audience is looking at a film of a man talking that is projected onto a wall

**GROUNDTRUTH objects:** person, tie, bottle

In this example, the object “tie” cannot be extracted using a lexical matching method, but it can be inferred through reasoning (ties are worn at formal events)

- Can learn to map from captions to labels efficiently



# Dealing with noise in extracted labels

High



a person riding a wave on top of a surfboard



a pizza sitting on top of a white plate next to a bowl of fries



a model walks down the runway during the fall fashion show



athletic young woman training with green ball in the gym



biscuit cat kitten cute 2008



cat dog shelter death kill euthanize euthanasia

Low



a traffic light showing the symbols for stop and go



a display of fabrics in different colors and patterns



the man took to facebook to show what a bottle of water can do if left in a hot car



an extra double ensuite room for extra or visiting guests



akshay marc matthew joyson nisha srikanth jothi christmas 2007 mckees rocks img8687 diabetes diabetes365



karmann ghia volkswagen car automobile chrome vintage retro signage logo font typography red shiny

COCO

Conceptual Captions

MIRFlickr1M

Fig. 11: Image-caption pairs with high homogeneity scores on the top, and low scores on the bottom.

# VEIL: Vetting Extracted Image Labels from In-the-Wild Captions

	<p>Big <u>Bear</u> Lake in WV. near Bruceton Mills.</p> <p><b>Similar Context:</b> Co-Occuring Context <b>Visual Defects:</b> - <b>Linguistic Indicators:</b> Named Entity, Noun Modifer</p>		<p>A day on the <u>boat</u></p> <p><b>Similar Context:</b> Co-Occuring Context <b>Visual Defects:</b> Occlusion <b>Linguistic Indicators:</b> Prepositional Phrase</p>
	<p>This piece of a 7 up <u>bottle</u> from 1958 I found buried a foot deep in my back yard.</p> <p><b>Similar Context:</b> - <b>Visual Defects:</b> Key Parts Missing, Atypical <b>Linguistic Indicators:</b> -</p>		<p>Spinning in my desk <u>chair</u></p> <p><b>Similar Context:</b> Co-Occuring Context <b>Visual Defects:</b> - <b>Linguistic Indicators:</b> Beyond the Image, Prepositional Phrase</p>

Figure 1: Examples of noisy extracted labels (underlined) from our **Caption Label Noise** dataset. We categorize types of similar context present instead of the underlined object, as well as types of visual defects and linguistic indicators that are useful for detecting noise.

# VEIL: Vetting Extracted Image Labels from In-the-Wild Captions

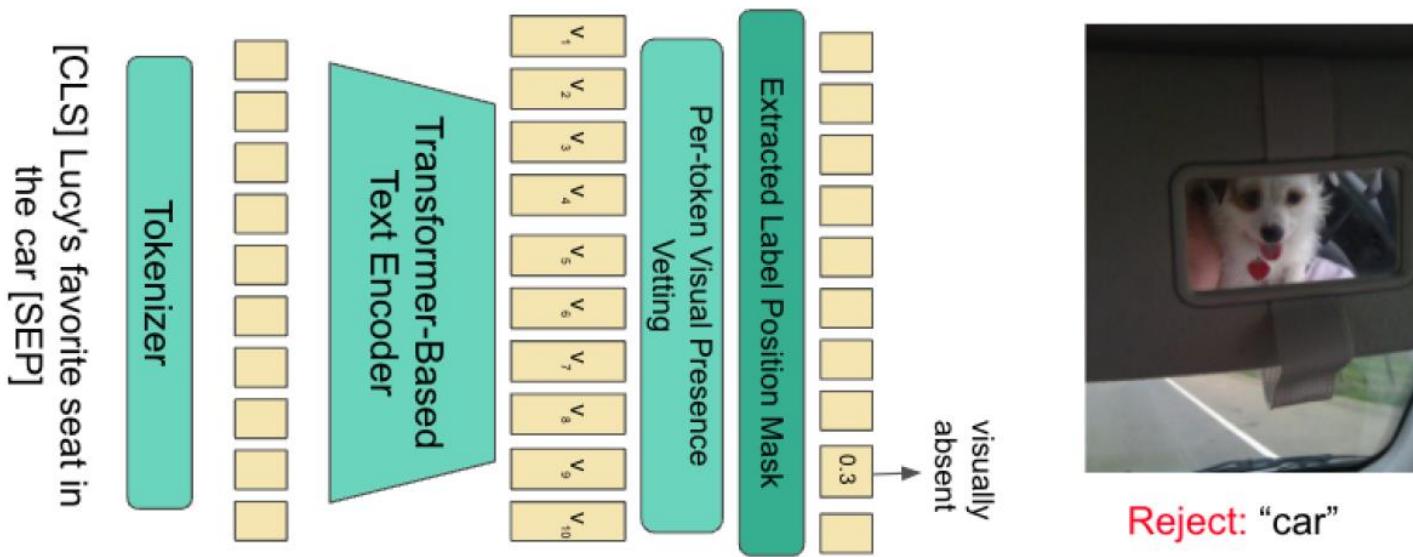
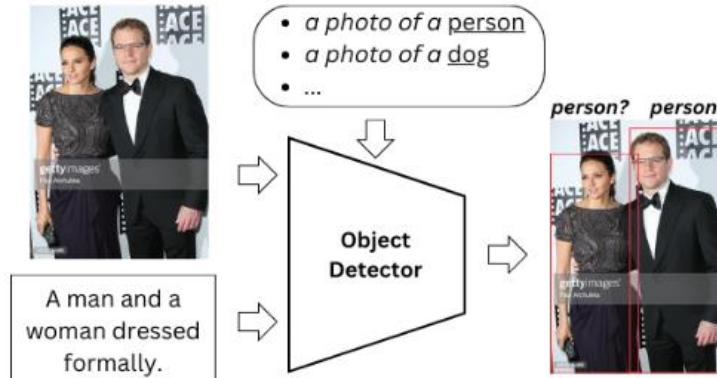


Figure 2: VEIL architecture. In this example, only “dog” is an extracted label and it fails the vetting process. The masking layer masks visual presence predictions for text tokens not corresponding to an extracted label.

# Synonym relations affect object detection learned on vision-language data



## 1) Data augmentation by back-translation

A man and a woman  
dressed formally. + En->De->En:  
A boy\_and a girl\_formally  
dressed.

## 2) Input class embeddings enrichment

• a photo of a person  
• a photo of a dog  
• ... { • a photo of a person  
• a photo of a man  
• a photo of a woman  
• ... }

Figure 1: Top: input to an open-vocabulary object detector: images, class embeddings, and captions; and its output: bounding boxes with associated labels. Bottom: our approaches. 1) Data augmentation by back-translation: add captions back-translated from a foreign language; 2) Class embeddings enrichment: consider synonyms when extracting class embeddings.

# Synonym relations affect object detection learned on vision-language data

Original	German	Russian
A <b>skate board rider</b> does a trick in front of a building.	A <b>skateboarder</b> does a trick in front of a building.	A <b>skater</b> does a trick in front of the building.
Three adults help a <b>youngster</b> follow a sheet of instructions.	Three adults help a <b>teenager</b> follow a sheet of instructions.	Three adults help the <b>teenager</b> follow instructions.

Table 1: Examples of (left) original COCO captions, (middle) captions back-translated from German, and (right) captions back-translated from Russian.

Captions	COCO names	Synonyms mean (std)	Avg.
<b>Class embeddings: COCO names</b>			
Original	<b>44.45</b>	33.87 (5.94)	35.63
BT: German	44.23	34.25 (5.32)	35.91
<b>Class embeddings: enriched</b>			
Original	43.58	37.25 (4.56)	38.31
BT: German	37.48	36.75 (4.56)	36.87
Curriculum	43.49	<b>37.93</b> (3.22)	<b>38.85</b>

Table 3: Class embedding enrichment: mAP@0.5 (as %) evaluated on COCO class embeddings (“COCO names”) and on synonyms embeddings (“Synonyms”).

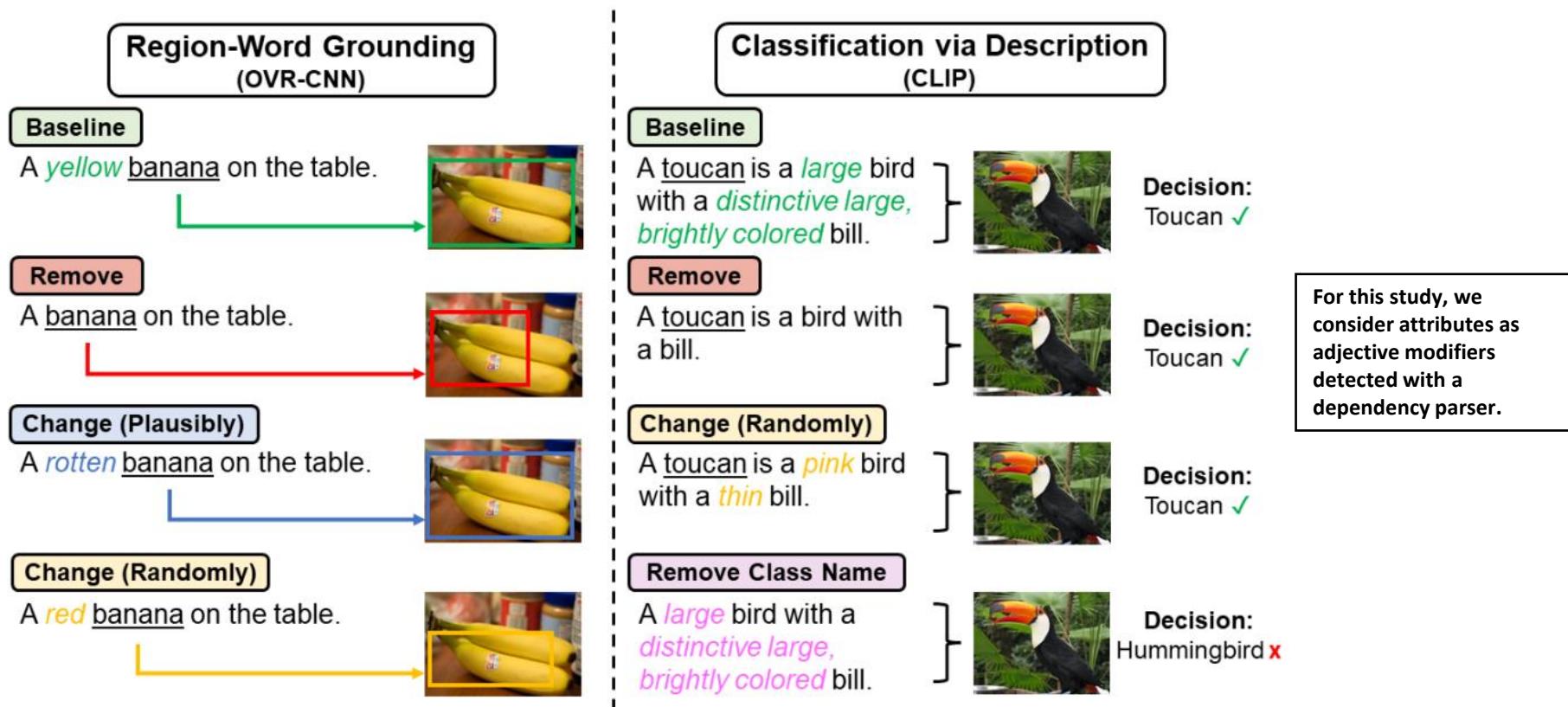
# Investigating the Role of Attribute Context in Vision-Language Models for Object Detection

- **RQ1:** Does attribute context play a role in VL pretraining for **object detection**?
  - Not significantly by default
- **RQ2:** Does learning to ground objects to **contextualized word embeddings** utilize attribute meaning?
  - Not significantly, unless adjective-perturbed caption negatives are used
- **RQ3:** Do VL models perform well at tasks where objects are described with/in terms of attributes?
  - No, there is room for improvement
  - Using class descriptions without class names and with only attributes is notably not effective
- **RQ4:** Can **contrastive negative sampling of captions** increase a model's ability to use attribute context?
  - Yes, in both region-word alignment pretraining and CLIP finetuning
- **RQ5:** What sampling mechanisms are most effective?
  - In OVR-CNN pretraining, both plausible and random sampling are effective
  - In CLIP finetuning, random sampling is most effective

# Measuring attribute sensitivity

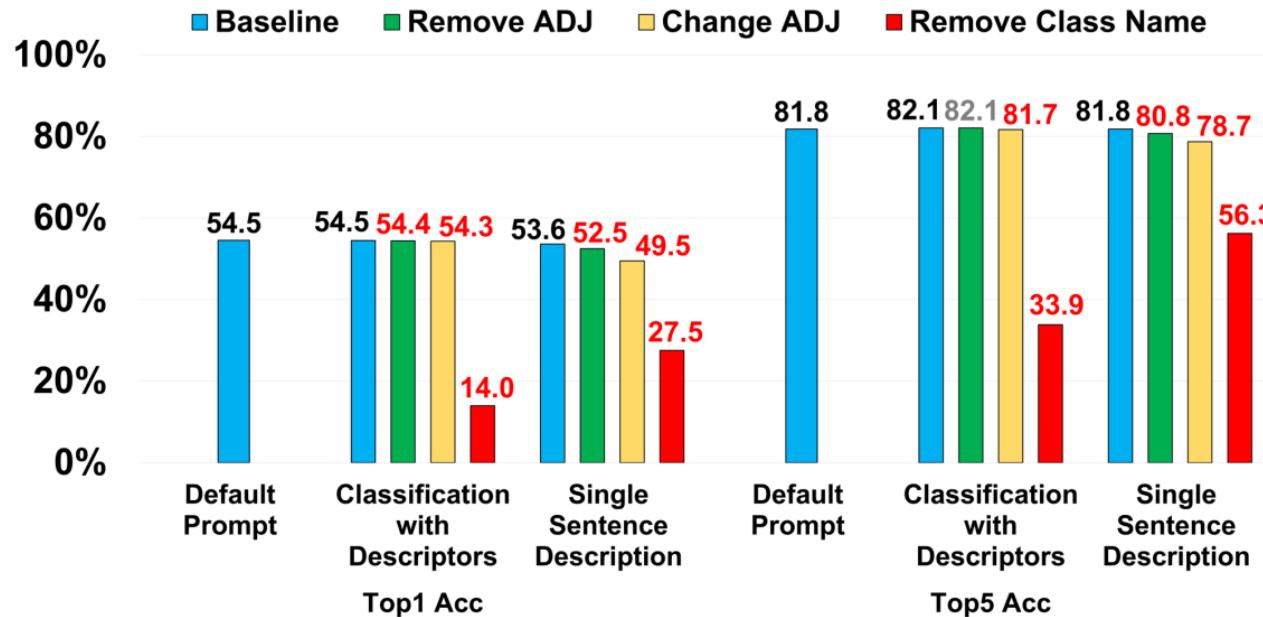
- **Our method:** Measure performance under various text changes

- Baseline – original caption/description
- Removing attributes – losing helpful information – *does performance drop?*
- Changing attributes (**plausibly/randomly**) – making text incorrect – *does performance drop?*
- Removing class names - Is describing an object by attributes alone effective?



# To what extent can CLIP effectively use attribute information?

- We test classification via description<sub>[1,2]</sub> with ImageNet-V2<sub>[3]</sub>
  - We find limited sensitivity in removing/changing attribute settings
    - More sensitivity observed in single-sentence case
  - We find **removing class names** to significantly reduce performance
    - Attribute-only descriptions are not sufficient for class differentiation



# Towards Shape-regularized Learning for Mitigating Texture Bias in CNNs

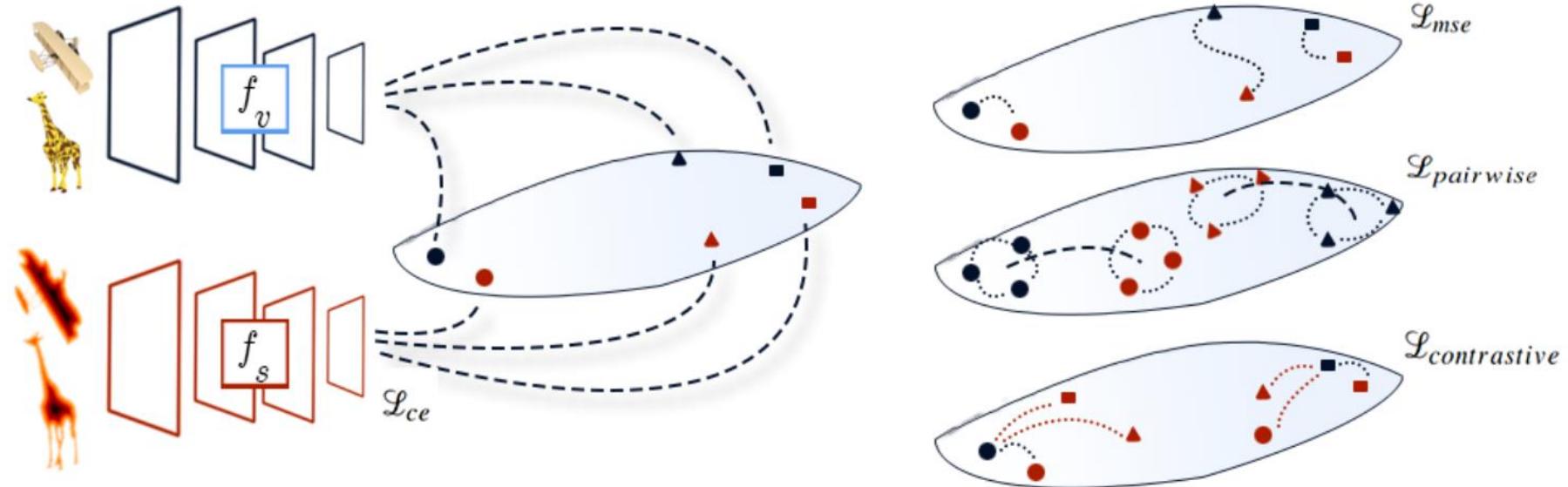


Figure 3: (Left) ResNet18 ( $f_v$ ) is trained to learn image features using the shape-bias loss. A pretrained ResNet18 ( $f_s$ ) is used as feature extractor for 2D distance transforms. Different categories are depicted by distinct shapes ( $\blacktriangle$ ,  $\bullet$ ,  $\blacksquare$ ) in the latent space. (Red) denotes shape features. Blue denotes image features. (Right) The shape-bias loss is composed of 3 losses, (a) mean squared error (b) pairwise loss and (c) contrastive loss. In addition, we use cross-entropy loss (for supervised learning). Blue dotted lines denote distance minimization. Red dotted lines denote distance maximization.

# Recognizing objects across geography

- Objects (e.g. cleaning equipment, plates, beds) look different in different countries; can knowledge from LLMs help us bridge the gaps in visual appearance?



**cleaning equipment:** ['bucket shape', 'handles for carrying', 'bright colors such as blue or green', 'mop or broom attachments', 'cleaning solution bottles', 'sponges or scrub brushes', 'dustpans or trash bags']



**bed:** ['rectangular shape', 'wooden frame', 'raised platform', 'woven mats or fabrics on top', 'often decorated with shells or feathers', 'may have mosquito nets or curtains']



**roof:** ['pointed or sloping shape', 'made of thatch or palm leaves', 'different colors and textures', 'often raised on stilts or posts', 'may have decorative elements such as carvings or patterns']



**plate:** ['round or oval shape', 'ceramic or clay material', 'intricate designs or patterns', 'earthy colors such as brown or beige', 'varying sizes and depths', 'may be used for serving food or as decorative objects']



**pet:** ['furry or feathered bodies', 'wagging tails or flapping wings', 'domesticated breeds such as cats or dogs', 'collars or leashes', 'obedience to human commands', 'friendly or playful behavior', 'common pets such as parrots or pigs.]



**salt:** ['white color', 'crystalline texture', 'irregular shapes', 'small granules', 'may appear in piles or mounds']

Figure 1. Examples of images from Papua New Guinea from DollarStreet [62] for which our proposed method allows us to confidently recognize the right category, while the baselines produce incorrect answers. We include information about the object's appearance in this country from a large language model, shown on the right of each image.

# Geo-Diverse Knowledge into Prompting for Increased Geographical Robustness in Object Recognition

## (1) Acquire Target Knowledge

Internal CLIP Knowledge  
+  
External LLM Knowledge

Target: Americas, Asia, Africa

A photo of a stove/hob in Mexico...  
A photo of a stove/hob in Vietnam...  
A photo of a stove/hob in Burkina Faso...

What are useful features for distinguishing a stove in a photo that I took in <country>?

LLM (davinci-003)

A stove/hob in Burkina Faso may be/have...

- Stone/mud material
- 3-4 burners
- Wide rectangular shape
- Metal/ceramic material
- Black/brown color
- Charcoal or wood fuel
- Metal grate on top

A stove/hob in Vietnam may be/have...

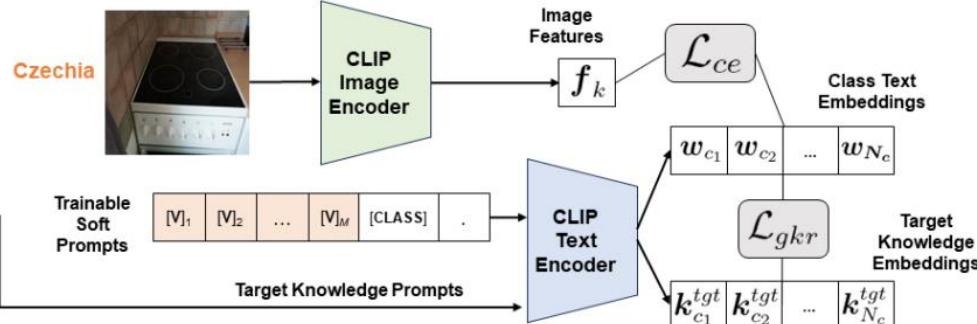
- Rectangular or circular shape
- Flat or slightly raised surface
- One, two, or four burners
- Black or white color
- Knobs, dials, and/or buttons
- Stainless steel or ceramic material

A stove/hob in Mexico may be/have...

- Four or more burners
- Metal or ceramic material
- White, black, or stainless-steel color
- Flat or slightly angled surface
- Knobs for setting temperature
- Over or broiler below burners
- Gas or electric powered
- Vent above hood

## (2) Optimize Soft Prompts on Source Data While Regularizing Towards Target Knowledge

Source: Europe



## (3) Recognize Objects in Target Countries



# Geo-Diverse Knowledge into Prompting for Increased Geographical Robustness in Object Recognition

Encoder	Prompting Method	Top-1 Accuracy										Top-3 Accuracy									
		Europe		Africa		Asia		Americas		Total		Europe		Africa		Asia		Americas		Total	
		Acc	Δ	Acc	Δ	Acc	Δ	Acc	Δ	Acc	Δ	Acc	Δ	Acc	Δ	Acc	Δ	Acc	Δ	Acc	Δ
ViT-B/32	Zero-Shot CLIP [36]	59.1	-	43.7	-	50.8	-	<b>55.3</b>	-	51.7	-	81.1	-	64.8	-	72.3	-	<b>77.4</b>	-	73.7	-
	GeneralLLM [30]	57.3	<b>-1.8</b>	44.3	<b>+0.6</b>	50.9	<b>+0.1</b>	54.6	<b>-0.7</b>	51.4	<b>-0.3</b>	78.8	<b>-2.3</b>	64.5	<b>-0.3</b>	72.1	<b>-0.2</b>	75.7	<b>-1.7</b>	73.0	<b>-0.7</b>
	CountryInPrompt	57.5	<b>-1.6</b>	45.2	<b>+1.5</b>	51.9	<b>+1.1</b>	55.0	<b>-0.3</b>	52.1	<b>+0.4</b>	80.2	<b>-0.9</b>	65.5	<b>+0.7</b>	73.3	<b>+1.0</b>	76.9	<b>-0.5</b>	73.9	<b>+0.2</b>
	CountryLLM	59.4	<b>+0.3</b>	45.2	<b>+1.5</b>	52.1	<b>+1.3</b>	<b>55.3</b>	0.0	52.6	<b>+0.9</b>	80.9	<b>-0.2</b>	66.4	<b>+1.6</b>	<b>73.6</b>	<b>+1.3</b>	<b>77.4</b>	0.0	74.6	<b>+0.9</b>
	CountryInPrompt+LLM	<b>60.8</b>	<b>+1.7</b>	<b>45.3</b>	<b>+1.6</b>	<b>52.2</b>	<b>+1.4</b>	55.0	<b>-0.3</b>	<b>52.8</b>	<b>+1.1</b>	<b>81.5</b>	<b>+0.4</b>	<b>67.4</b>	<b>+2.6</b>	<b>73.6</b>	<b>+1.3</b>	76.7	<b>-0.7</b>	<b>74.7</b>	<b>+1.0</b>
ViT-B/16	Zero-Shot CLIP [36]	64.3	-	46.9	-	53.9	-	<b>60.1</b>	-	55.5	-	84.3	-	69.3	-	75.9	-	81.1	-	77.2	-
	GeneralLLM [30]	64.2	<b>-0.1</b>	48.8	<b>+1.9</b>	<b>56.0</b>	<b>+2.1</b>	58.5	<b>-1.6</b>	56.8	<b>+1.3</b>	83.9	<b>-0.4</b>	71.1	<b>+1.8</b>	76.3	<b>+0.4</b>	80.4	<b>-0.7</b>	77.9	<b>+0.7</b>
	CountryInPrompt	63.9	<b>-0.4</b>	49.6	<b>+2.7</b>	55.7	<b>+1.8</b>	59.3	<b>-0.8</b>	56.6	<b>+1.1</b>	84.0	<b>-0.3</b>	71.3	<b>+2.0</b>	76.5	<b>+0.6</b>	80.0	<b>-1.1</b>	77.7	<b>+0.5</b>
	CountryLLM	65.2	<b>+0.9</b>	49.6	<b>+2.7</b>	55.6	<b>+1.7</b>	59.7	<b>-0.4</b>	57.0	<b>+1.5</b>	84.3	0.0	71.8	<b>+2.5</b>	<b>77.5</b>	<b>+1.6</b>	<b>81.5</b>	<b>+0.4</b>	<b>78.8</b>	<b>+1.6</b>
	CountryInPrompt+LLM	<b>65.5</b>	<b>+1.2</b>	<b>50.8</b>	<b>+3.9</b>	<b>56.0</b>	<b>+2.1</b>	59.7	<b>-0.4</b>	<b>57.4</b>	<b>+1.9</b>	<b>85.5</b>	<b>+1.2</b>	<b>72.5</b>	<b>+3.2</b>	77.0	<b>+1.1</b>	80.9	<b>-0.2</b>	78.7	<b>+1.5</b>
RN50	Zero-Shot CLIP [36]	53.0	-	38.0	-	44.4	-	49.8	-	45.7	-	76.5	-	60.2	-	66.4	-	<b>72.7</b>	-	68.1	-
	GeneralLLM [30]	55.5	<b>+2.5</b>	40.9	<b>+2.9</b>	46.9	<b>+2.5</b>	50.3	<b>+0.5</b>	47.9	<b>+2.2</b>	76.0	<b>-0.5</b>	61.2	<b>+1.0</b>	67.7	<b>+1.3</b>	71.1	<b>-1.6</b>	68.6	<b>+0.5</b>
	CountryInPrompt	54.5	<b>+1.5</b>	<b>43.4</b>	<b>+5.4</b>	47.0	<b>+2.6</b>	50.8	<b>+1.0</b>	48.4	<b>+2.7</b>	76.0	<b>-0.5</b>	<b>64.0</b>	<b>+3.8</b>	68.7	<b>+2.3</b>	<b>72.7</b>	0.0	<b>70.0</b>	<b>+1.9</b>
	CountryLLM	56.2	<b>+3.2</b>	41.1	<b>+3.1</b>	47.3	<b>+2.9</b>	50.4	<b>+0.6</b>	48.3	<b>+2.6</b>	<b>77.2</b>	<b>+0.7</b>	62.5	<b>+2.3</b>	<b>68.8</b>	<b>+2.4</b>	72.4	<b>-0.3</b>	<b>70.0</b>	<b>+1.9</b>
	CountryInPrompt+LLM	<b>56.4</b>	<b>+3.4</b>	43.0	<b>+5.0</b>	<b>48.0</b>	<b>+3.6</b>	<b>50.9</b>	<b>+1.1</b>	<b>49.1</b>	<b>+3.4</b>	76.7	<b>+0.2</b>	63.1	<b>+2.9</b>	68.3	<b>+1.9</b>	71.1	<b>-1.6</b>	69.4	<b>+1.3</b>

Table 1. **Zero-shot CLIP inference with descriptive knowledge prompts, top-1/3 balanced accuracy (Acc) on DollarStreet.** Strategies to capture CLIP’s internal country knowledge (CountryInPrompt), external LLM country knowledge (CountryLLM), and their combination (CountryInPrompt+LLM), often improve vs. the zero-shot CLIP baseline (prompt “a photo of a/an <object>”), especially on Africa and Asia; gains in green, drops in red. CountryLLM notably outperforms the GeneralLLM [30] baseline.

# Geo-Diverse Knowledge into Prompting for Increased Geographical Robustness in Object Recognition

Encoder	Prompting Method	Source Europe		Africa		Asia		Americas		Total	
		Acc	$\Delta$	Acc	$\Delta$	Acc	$\Delta$	Acc	$\Delta$	Acc	$\Delta$
ViT-B/16	CoOp [52]	72.2	-	53.9	-	61.5	-	68.6	-	61.7	-
	CoCoOp [51]	73.2	-	54.3	-	61.2	-	68.3	-	61.4	-
	KgCoOp [47]	73.1	-	54.4	-	62.6	-	68.7	-	62.4	-
	CountryInPrompt Reg	71.8	<i>-1.4</i>	56.8	<b>+2.4</b>	63.0	<b>+0.4</b>	69.8	<b>+1.1</b>	63.5	<b>+1.1</b>
	CountryLLM Reg	73.2	<i>0.0</i>	55.6	<b>+1.2</b>	63.0	<b>+0.4</b>	70.0	<b>+1.3</b>	63.2	<b>+0.8</b>
	CountryInPrompt+LLM Reg	<b>73.6</b>	<b>+0.4</b>	<b>57.2</b>	<b>+2.8</b>	<b>63.8</b>	<b>+1.2</b>	<b>70.3</b>	<b>+1.6</b>	<b>64.0</b>	<b>+1.6</b>
RN50	CoOp [52]	64.6	-	45.2	-	51.6	-	59.5	-	52.2	-
	CoCoOp [51]	62.9	-	44.5	-	51.0	-	58.3	-	51.4	-
	KgCoOp [47]	63.5	-	46.3	-	53.9	-	<b>60.5</b>	-	53.9	-
	CountryInPrompt Reg	63.5	<i>-1.1</i>	48.0	<b>+1.7</b>	53.9	0.0	60.3	<i>-0.2</i>	54.3	<b>+0.4</b>
	CountryLLM Reg	64.5	<i>-0.1</i>	47.4	<b>+1.1</b>	54.2	<b>+0.3</b>	59.9	<i>-0.6</i>	54.3	<b>+0.4</b>
	CountryInPrompt+LLM Reg	<b>65.5</b>	<b>+0.9</b>	<b>48.1</b>	<b>+1.8</b>	<b>54.5</b>	<b>+0.6</b>	60.4	<i>-0.1</i>	<b>54.8</b>	<b>+0.9</b>

Table 2. **Regularizing soft prompts with geographical knowledge, top-1 bal. acc. on DollarStreet.** We emphasize that our regularization aims to improve **target** performance, rather than source (gray, *italicized*). **Gains/drops** are shown vs. the *best* of soft prompt baselines (shaded). CountryInPrompt+LLM Reg achieves notable gains in target, especially on Africa. Methods use 16 shots per class.

# Recognizing objects across geography

- Language supervision has different structure across languages; not all languages mention all objects

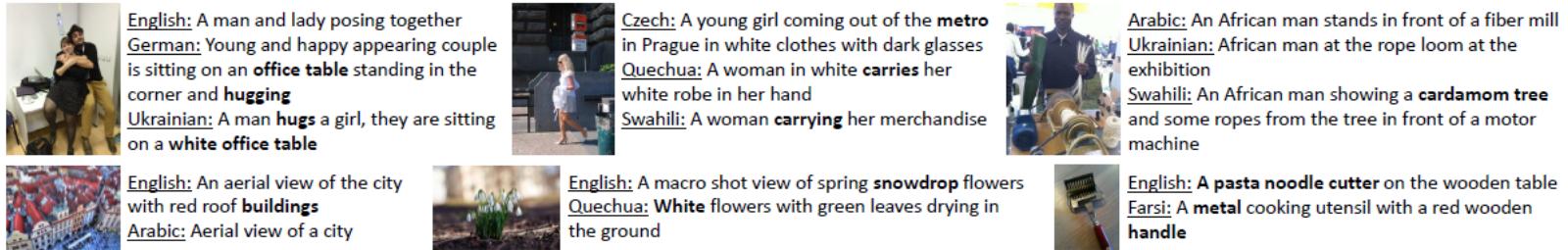


Figure 3. Representative examples of language differences in captions from the XM3600 dataset [68]. Noteworthy differences in mentioned objects and actions are highlighted.

	eu-mean	eu-stdev	ar-mean	ar-stdev	hi-mean	hi-stdev	id-mean	id-stdev	asia-mean	asia-stdev	sw
tree(s)	270.5---	92.9	349	19.799	<b>581.5+++</b>	214.3	286	49.5	274.7	63.52	<b>383</b>
mountain(s)	171.1---	47.78	183	24.0416	185.5	31.82	173	42.43	<b>218+++</b>	16.52	<b>208</b>
street	<b>100.9</b>	30.22	<b>124+++</b>	50.9117	61	7.071	38.5---	10.61	76.67	19.01	82
car(s)	207.3	19.96	235	24.0416	<b>239</b>	50.91	204	11.31	220	17.78	<b>270+++</b>
building(s)	244.8---	69.3	281.5	40.3051	329	108.9	<b>383.5</b>	84.15	253.3	49.9	<b>502+++</b>
restaurant	45.82	13.71	<b>54</b>	7.07107	19---	5.657	<b>50.5+++</b>	13.44	42.67	6.11	21
table	156.7	52.82	162.5	58.6899	<b>240+++</b>	12.73	<b>228</b>	93.34	185.3	43.66	121---
plate	112.5	25.93	90---	12.7279	105.5	10.61	109.5	33.23	<b>119.3+++</b>	5.132	<b>113</b>
box	18.14	4.454	15.5---	0.70711	15.5	2.121	<b>28+++</b>	4.243	<b>24</b>	2.646	18
bottle	10.23---	2.654	12	0	10.5	2.121	11	4.243	<b>14.67+++</b>	0.577	<b>18</b>
dog	26.23	5.108	28	1.41421	<b>31.5+++</b>	4.95	29.5	0.707	20.67---	5.508	<b>34</b>
woman	135.5	23.72	127	5.65685	114---	31.11	<b>164.5+++</b>	20.51	133.3	27.65	<b>160</b>

Table 1. Language shifts in terms of concept mentions in different languages. We grouped European languages (eu), Arabic and Farsi (ar), Hindi and Bengali (hi), Indonesian and Thai (id), Asian languages (asia), and report Swahili on its own (sw). The largest two numbers per row are bolded. Observe the differences between the language with highest (+++) and lowest (---) counts, which are significantly larger than the within-group standard deviations.

# Recognizing objects across geography

- Need to recognize rare, culture-specific objects (with limited data)



martenitsa: braided (doll)



pisanka: ornate (egg)



aguayo: vibrant (blanket)



kikombe: carved (cup)



sugar apple: bumpy (fruit)



ofuro: wooden (bath)

Figure 4. Rare objects unique to specific cultures, with their most distinctive visual attribute (and a related category).

# Acknowledgements

- Funding



- Students involved in highlighted projects



Keren Ye



Mingda  
Zhang



Chris  
Thomas



Erhan  
Unal



Kyle  
Buettner



Arushi Rai



Giacomo  
Nebbia