

# Grounding Semantic Roles in Images

authors: Carina Silberer, Manfred Pinkal [EMNLP' 18]

Presented by: Boxin Du

University of Illinois at Urbana-Champaign



# Roadmap

- Motivation
- Problem Definition
- Proposed Method
- Evaluations
- Conclusion

# Motivation

- Scene interpretation
- Example:

image



*Well, the fridge  
broke, so I had to  
eat everything.*

text

TARGET IMAGE

CONTEXT

- Q: Why there is so much food on the table?
- The interpretation of a (visual) scene is related to the determination of its events, their participants and the roles they play therein (i.e., distill who did what to whom, where, why and how)

# Motivation (cont'd)

- Traditional Semantic Role Labeling (SRL):
  - Extract interpretation in the form of shallow semantic structures from natural language texts.
  - Applications: Information extraction, question answering, etc.
- Visual Semantic Role Labeling (vSRL):
  - Transfer the use of semantic roles to produce similar structured meaning descriptions for visual scenes.
  - Induce representations of texts and visual scenes by joint processing over multiple sources



# Roadmap

- Motivation
- Problem Definition
- Proposed Method
- Evaluations
- Conclusion

# Problem Definition

- Goal:
  - learn frame–semantic representations of images (vSRL)
  - Specifically, learn distributed situation representations (for images and frames), and participant representations (for image regions and roles)
- Two subtasks:
  - Role Prediction: predict the role of an image region (object) under certain frame
  - Role Grounding: realize (i.e. map) a given role to a specific region (object) in an image under certain frame

# Problem Definition (cont'd)

- Role Prediction:

- Given an image  $i$ , its region set  $R_i$ , map the regions  $r \in R_i$  to the predicted role  $e \in E$  and the frame  $f \in F$  it is associated with.

$$L : \{i\} \times R_i \rightarrow F \times E$$
$$L(i, r) = \arg \max_{(f, e), f \in F, e \in E_f} s(i, r, f, e)$$

$s()$  quantifies the visual-frame-semantic similarity between the region  $r$  and the role  $e$  of  $f$

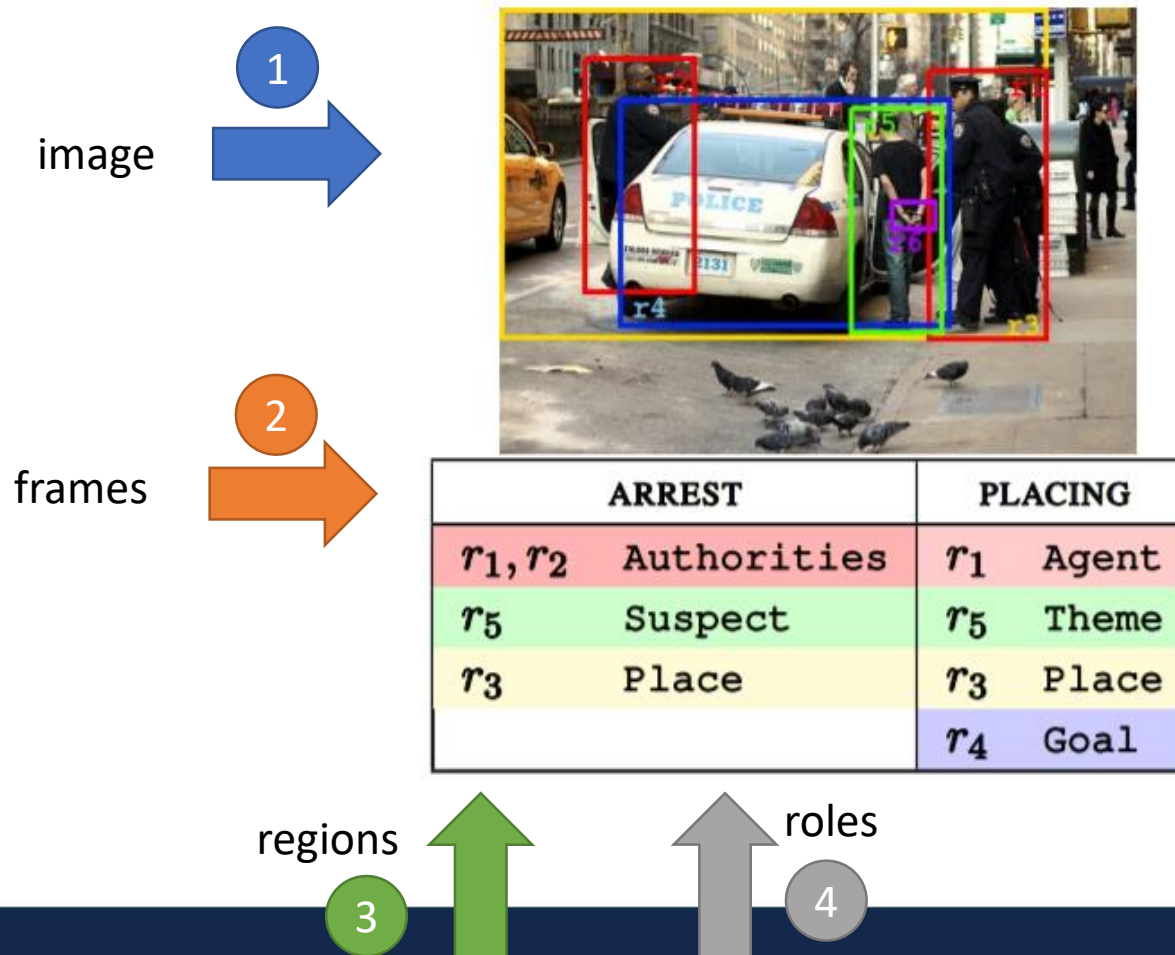
- Role Grounding:

- Given a frame  $f$  realized in  $i$ , ground each role  $e \in E_f$  in the region  $r \in R_i$  with the highest visual-frame semantic similarity to role  $e$ .

$$G : \{i\} \times \{f\} \times E_f \rightarrow R_i$$
$$G(i, f, e) = \arg \max_{r \in R_i} s(i, r, f, e)$$

# Problem Definition (cont'd)

- Example: given an image with annotations



- Role Prediction:

Given 1 3

Predict 2 4

- Role Grounding:

Given 1 2 4

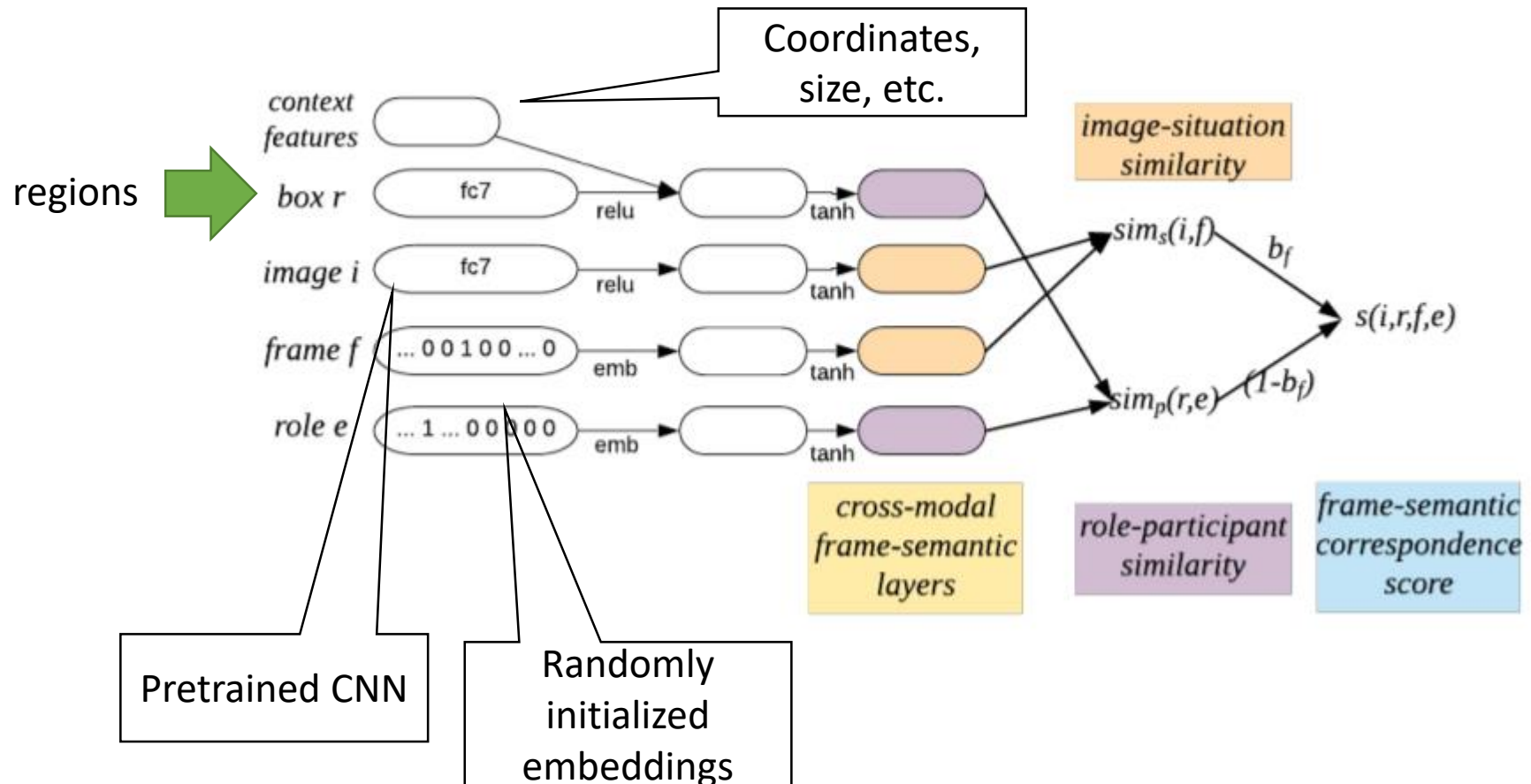
Predict 3

# Roadmap

- Motivation
- Problem Definition
- **Proposed Method**
- Evaluations
- Conclusion

# Proposed Method

- Overall architecture: Visual-Frame–Semantic Embedder



# Proposed Method

- Frame-semantic correspondence score:

$$s(q) = b_f \text{sim}_s(i, f) + (1 - b_f) \text{sim}_p(r, e)$$

- Training:

$$\theta = \arg \min_{\theta} \sum_{q \in Q} \frac{1}{K} \sum_{k=1}^K \max(0, M - s(q) + s(q'_k))$$

- Where the  $q = (i, r, f, e) \in Q$  and  $Q$  is the training set. For each positive example, the training stage samples  $K$  negative examples.

# Proposed Method

- Data:

- Apply PathLSTM [1] for extracting the grounded frame-semantic annotations

- E.g.



(img<sub>1</sub>, r<sub>5</sub>, PLACING, Theme)  
(img<sub>1</sub>, r<sub>1</sub>, PLACING, Agent)  
(img<sub>1</sub>, r<sub>4</sub>, PLACING, Goal)

(img<sub>1</sub>, r<sub>1</sub>-r<sub>2</sub>, ARREST, Authorities)  
(img<sub>1</sub>, r<sub>5</sub>, ARREST, Suspect)  
(img<sub>1</sub>, r<sub>3</sub>, ARREST, Place)

(1a) [r5 A man] is being placed in [r4 a police car] by [r1 a uniformed officer].

(1b) [r1,r2 The police] arresting [r5 someone] on [r3 a busy city street].

(1c) [r5 A young guy] is getting arrested.

(2a) PLACING (Theme:r5/A man,  
Goal:r4/a police car,  
Agent:r1/a uniformed officer )

(2b) ARREST (Authorities:r1,r2/The police,  
Suspect:r5/someone,  
Place:r3/on a busy city street )

(2c) ARREST ( Suspect:r5/A young guy )



# Roadmap

- Motivation
- Problem Definition
- Proposed Method
- **Evaluations**
- Conclusion

# Evaluations

- Role Prediction (dataset: Flickr30k):

		Fine-grained frame types							Coarse frame types			
		top-1-pred.			top-5 preds.			gt fr.	top-1-pred.		top-5 preds.	
		frame	fr.role	role	frame	fr.role	role	role	frame	fr.role	frame	fr.role
test set	Image-only	19.0	9.4	16.7	44.1	28.6	52.3	47.9	23.7	12.0	55.8	36.3
	ImgObject	18.7	12.8	24.1	44.9	33.8	61.2	64.3	22.6	15.5	55.5	41.4
	ImgObjLoc	18.6	13.5	25.9	46.8	35.7	62.2	65.7	23.0	16.7	56.5	43.2
reference	Image-only	<b>27.8</b>	13.2	17.2	55.2	39.3	57.3	50.2	<b>30.8</b>	14.6	67.8	46.6
	ImgObject	22.6	15.7	22.4	59.6	44.3	66.9	69.0	25.1	16.7	<b>68.8</b>	51.0
	ImgObjLoc	24.9	<b>17.4</b>	<b>23.6</b>	<b>60.2</b>	<b>47.3</b>	<b>68.6</b>	<b>70.3</b>	28.4	<b>19.7</b>	67.4	<b>53.3</b>

Image-only: a model that only uses the image as visual input

ImgObject: a model that does not use contextual box features

ImgObjLoc: the original model

- Obs.: horizontally the original model yields the overall best results; vertically the model is able to generalize over wrong role-filler pairs in the training data

# Evaluations

- Role Grounding (dataset: Flickr30k):

assigns each role  
randomly to a box in  
the image

			Fine-grained frame types						Fine-grained frame types						
			top-1 pred. filler			top-3 pred. fillers			top-1 pred. filler			top-3 pred. fillers			
			frame	fr.role	role	frame	fr.role	role	frame	fr.role	role	frame	fr.role	role	
test set	Random	$\tau_0$	37.7	23.6	25.3	70.8	56.5	59.4	props	5.5	3.7	4.1	15.7	10.6	11.6
	ImgObject		55.9	55.1	58.0	83.2	84.0	78.7		10.5	11.3	11.7	21.8	21.4	21.2
	ImgObjLoc		56.6	56.6	59.4	83.1	85.1	79.7		11.5	12.8	13.3	22.3	22.6	22.5
reference	Random	$\tau_0$	54.7	25.7	25.7	91.7	65.5	65.5	props	8.1	3.8	3.8	22.9	11.8	11.8
	ImgObject		78.9	62.1	62.1	95.8	88.2	83.6		13.7	12.8	12.8	39.6	30.9	28.2
	ImgObjLoc		<b>80.8</b>	<b>63.9</b>	<b>63.9</b>	<b>97.9</b>	<b>91.8</b>	<b>86.4</b>		<b>18.6</b>	<b>16.9</b>	<b>16.9</b>	<b>43.8</b>	<b>35.5</b>	<b>34.6</b>

Obs.: Horizontally ImgObjLoc is significantly more effective than ImgObject in all settings; vertically the models perform substantially better on the reference set than on the noisy test set (generalize over wrong role-filler pairs in the training data)

# Evaluations

- Visual Verb Sense Disambiguation (VerSe dataset):
  - The usefulness of the learned frame-semantic image representations on the task of visual verb disambiguation

Features	Motion	Non-motion
Random	$76.7 \pm 0.86$	$78.5 \pm 0.39$
MFS <sup>+</sup>	76.1	80.0
CNN <sup>+</sup>	82.3	80.0
Gella-CNN+O <sup>+</sup>	83.0	80.0
Gella-CNN+C <sup>+</sup>	82.3	80.3
CNN (reproduced)	83.1	$79.8 \pm 0.53$
<b>ImgObjLoc</b>	<b><math>84.8 \pm 0.69</math></b>	<b><math>80.4 \pm 0.57</math></b>

those which have at least 20 images and at least 2 senses

- Obs.: ImgObjLoc vectors outperform all comparison models on motion verbs; comparable with CNN on non-motion verbs.
- Reason: only frame-semantic embeddings are used?

# Roadmap

- Motivation
- Problem Definition
- Proposed Method
- Evaluations
- Conclusion

# Conclusion

- Goal:
  - grounding semantic roles of frames which an image evokes in the corresponding image regions of its fillers.
- Proposed method:
  - A model that learns distributed situation representations (for images and frames), and participant representations (for image regions and roles) which capture the visual–frame-semantic features of situations and participants, respectively.
- Results:
  - Promising results on role prediction, grounding (making correct predictions for erroneous data points)
  - It outperforms or is comparable to previous work on the supervised visual verb sense disambiguation task

# Thanks!



# VQA: Visual Question Answering

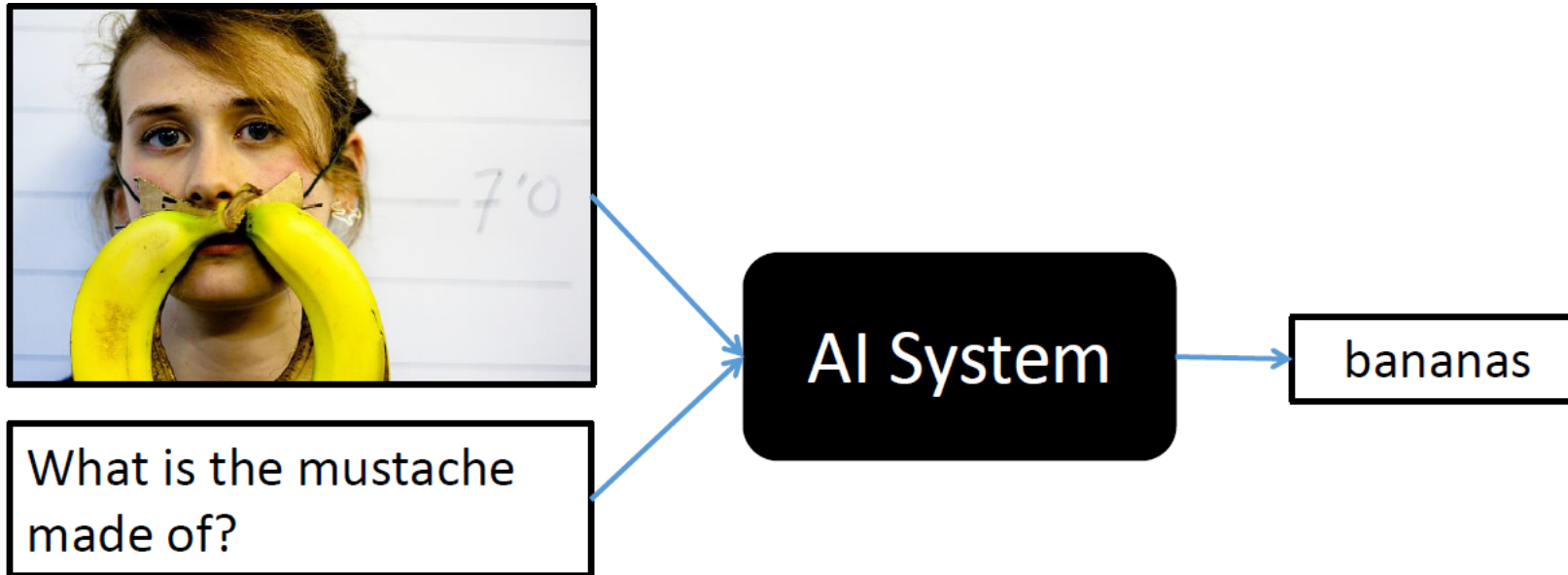
---

Aishwarya Agrawal, Jiasen Lu, Stanislaw Antol, Margaret Mitchell,  
C. Lawrence Zitnick, Dhruv Batra, Devi Parikh  
ICCV 2015

Presented by: Xinyang Zhang



# What is VQA?

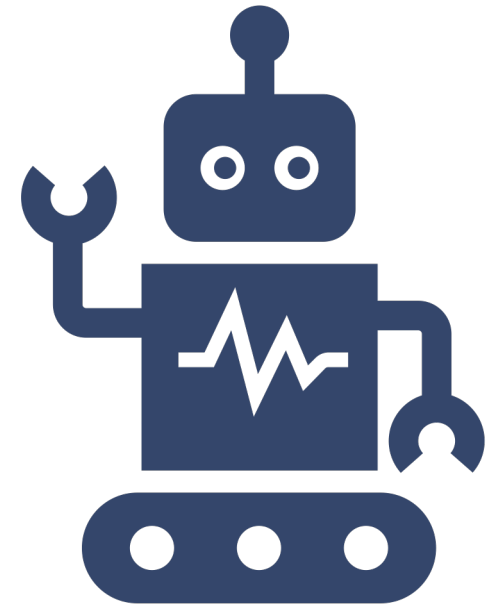


# Main contributions

- A new task
- A new dataset
- Baseline models

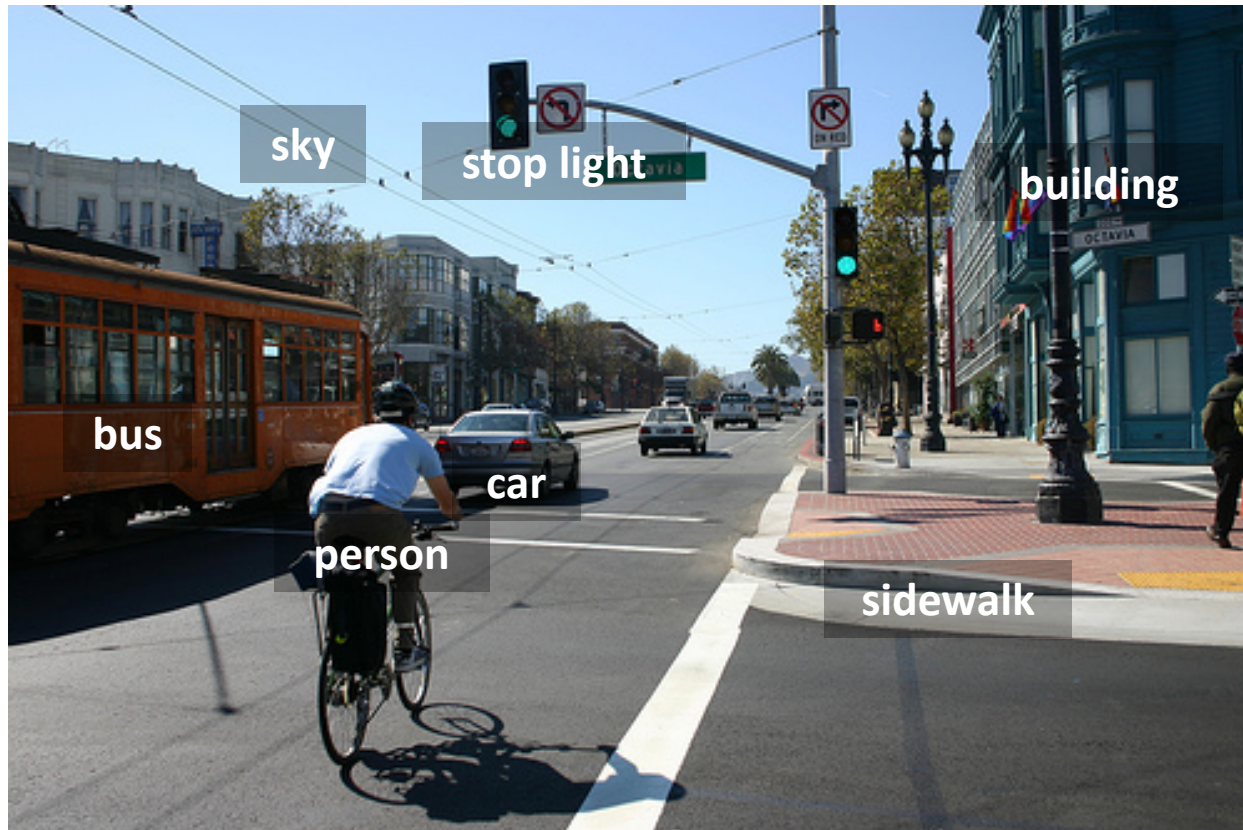
# Why VQA?

- Towards an “AI-complete” task

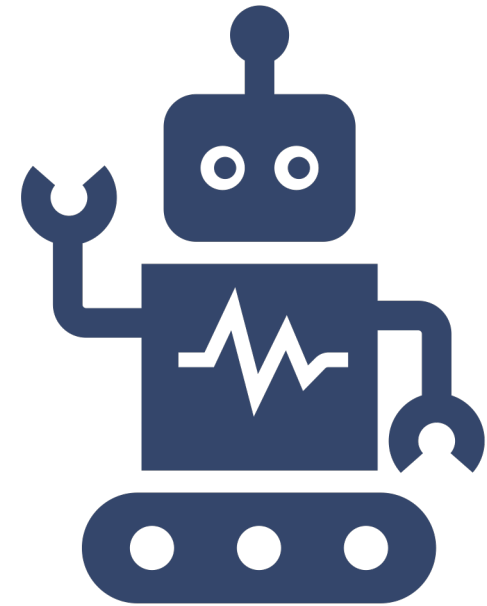


# Why VQA?

- Towards an “AI-complete” task



Object recognition?



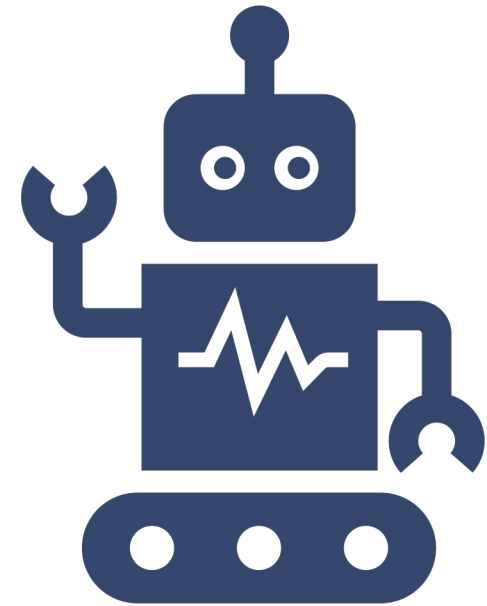


# Why VQA?

- Towards an “AI-complete” task



Scene recognition?

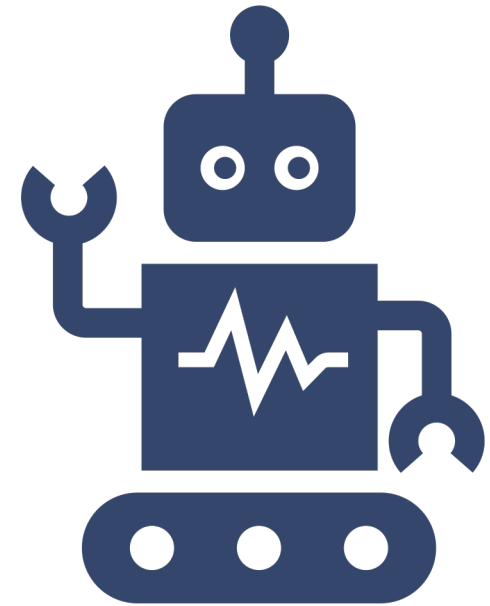


# Why VQA?

- Towards an “AI-complete” task



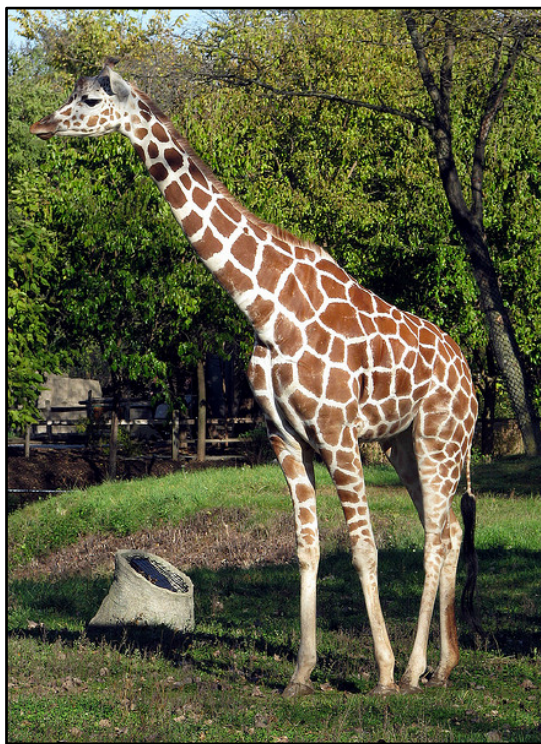
Image captioning?



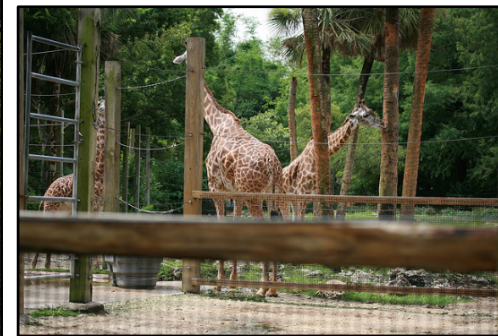
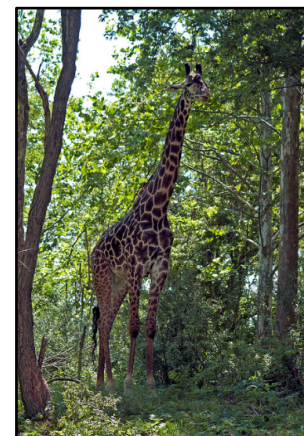
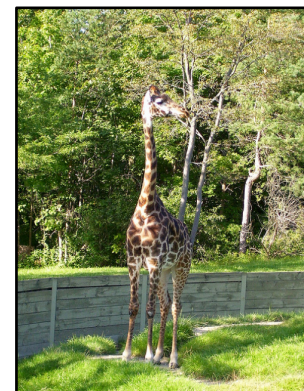
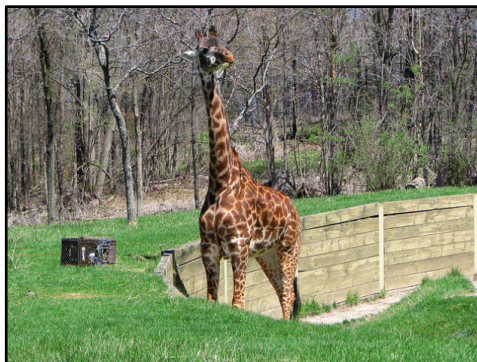
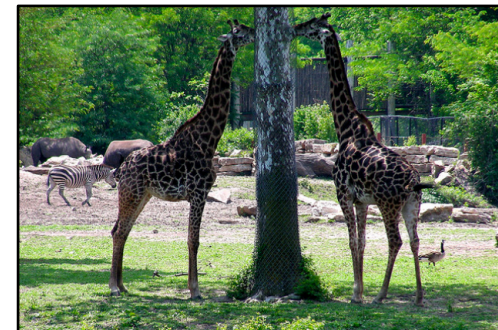
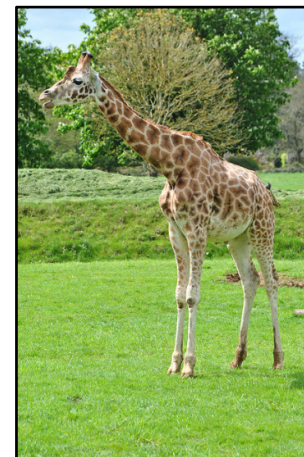


# Why VQA?

- Towards an “AI-complete” task



A giraffe standing in the grass next to a tree.





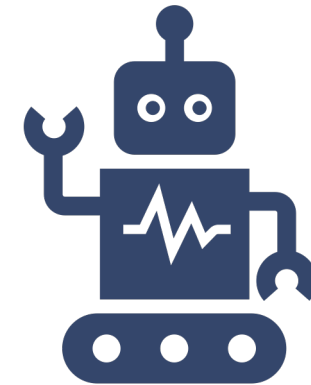
# Why VQA?

- Towards an “AI-complete” task



Answer questions about the scene

- Q: How many buses are there?
- Q: What is the name of the street?
- Q: Is the man on bicycle wearing a helmet?





# Why VQA?

- Towards an “AI-complete” task
  1. Multi-modal knowledge
  2. Quantitative evaluation

# Why VQA?

- Flexibility of VQA
  - Fine-grained recognition
    - “What kind of cheese is on the pizza?”
  - Object detection
    - “How many bikes are there?”
  - Knowledge base reasoning
    - “Is this a vegetarian pizza?”
  - Commonsense reasoning
    - “Does this person have 20/20 vision?”

# Why VQA?

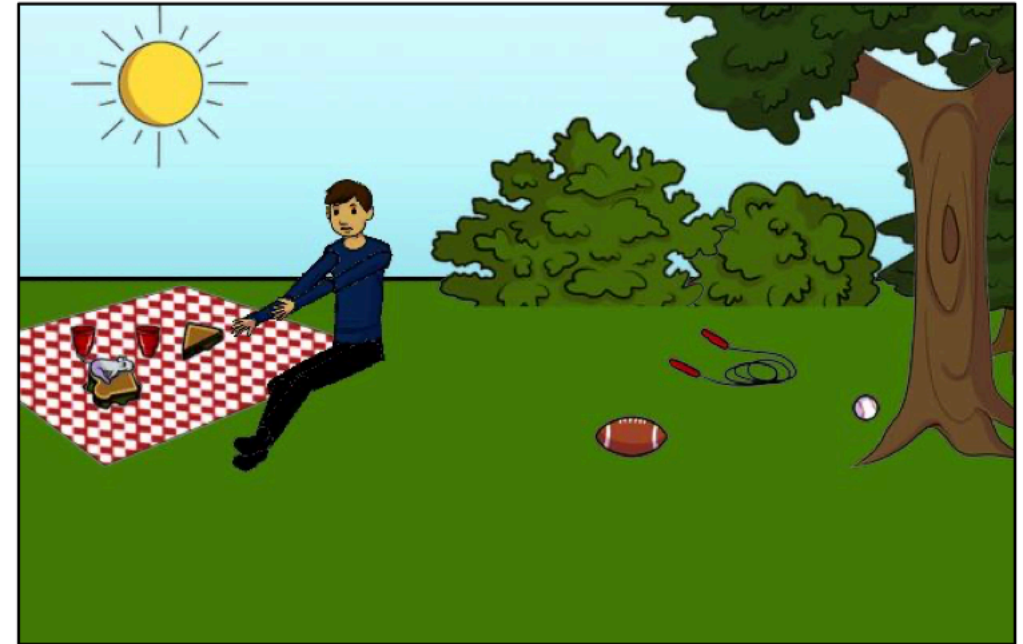
- Automatic quantitative evaluation possible
  - Multiple choice questions
  - “Yes” or “no” questions (~40%)
  - Numbers (~13%)
  - Short answers (one word 89.32%, two words 6.91%, three words 2.74%)

# How to collect a high-quality dataset?

- Images



Real Images  
(from MS COCO)



Abstract Scenes  
(curated)

# How to collect a high-quality dataset?

- Questions
  - Interesting and diverse
  - High-level image understanding
  - Require image to answer

*“We have built a **smart robot**. It understands a lot about images. **It can** recognize and name all the objects, it knows where the objects are, it can recognize the scene (e.g., kitchen, beach), people’s expressions and poses, and properties of objects (e.g., color of objects, their texture). Your task is to **stump this smart robot**!*

*Ask a question about this scene that this **smart robot probably can not answer**, but any **human can easily answer** while looking at the scene in the image.”*

“Smart robot” interface

# How to collect a high-quality dataset?

- Answers
  - 10 human answers
  - **Encourage short phrases** instead of long sentence
  - (1) Open-ended & (2) multiple-choice
- Evaluation
  - Exact match

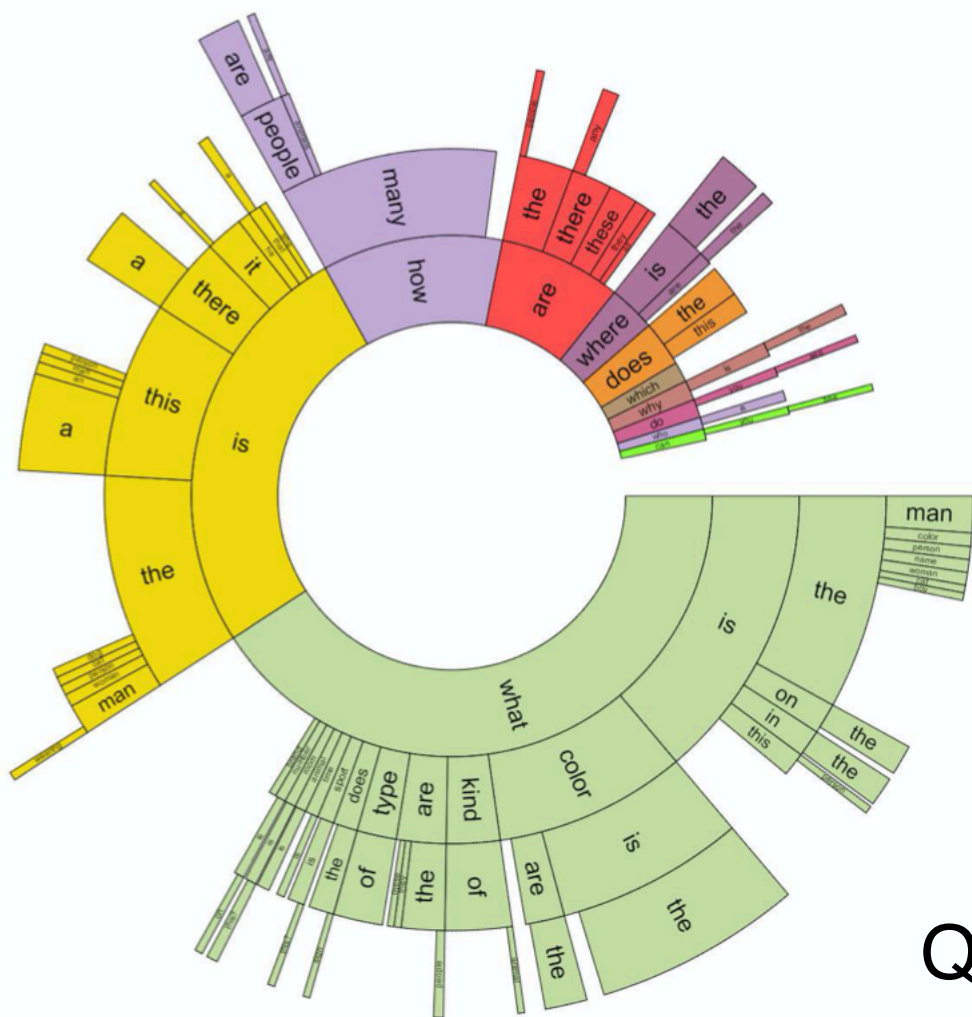
$$\text{accuracy} = \min\left(\frac{\# \text{ humans that provided that answer}}{3}, 1\right)$$

# Dataset Analysis

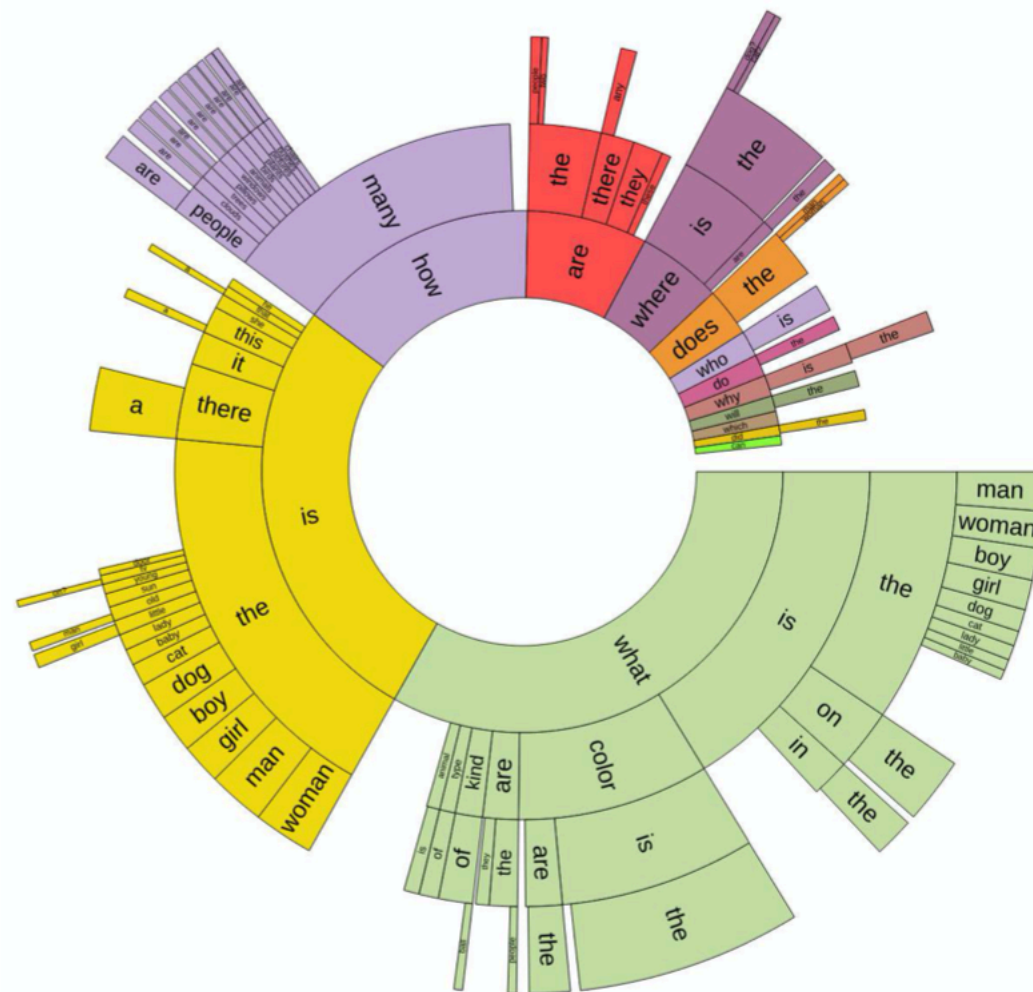
- ~0.25M images, ~0.76M questions, ~10M answers

# Dataset Analysis

## Real Images



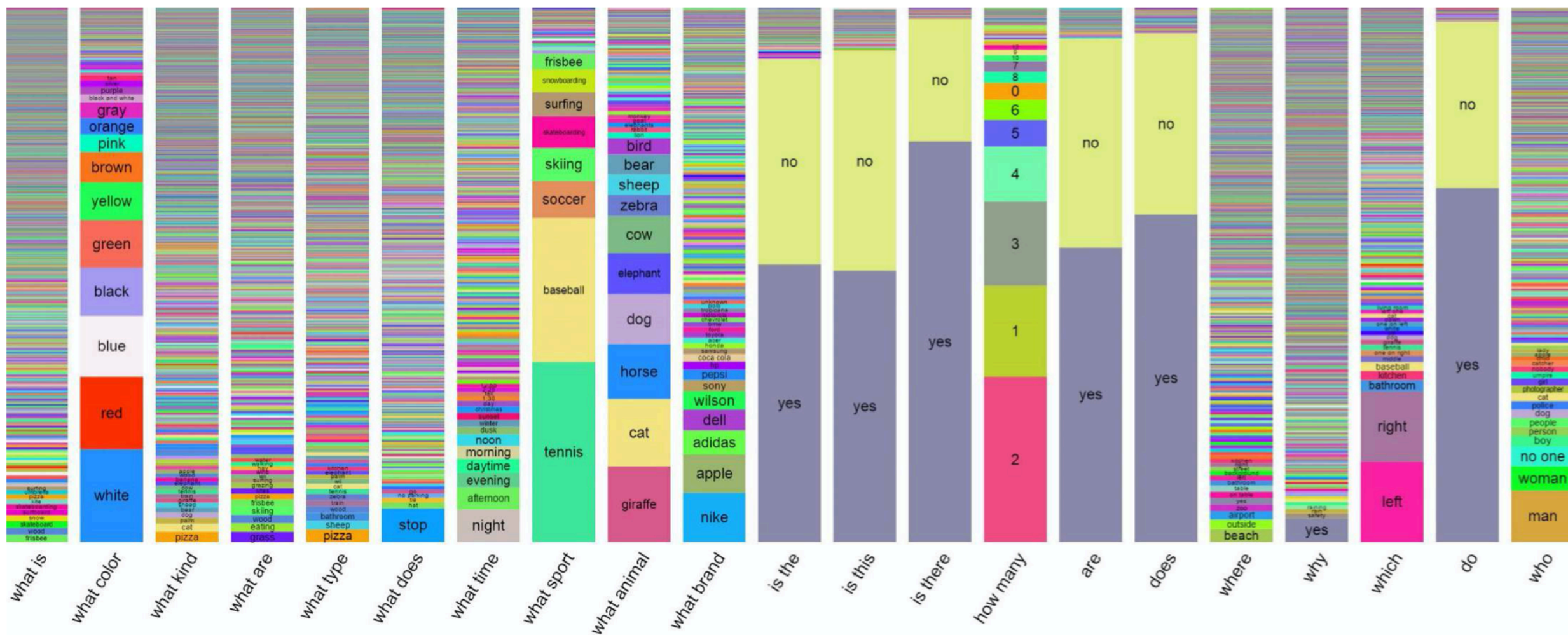
## Abstract Scenes



# Questions



# Dataset Analysis



# Answers

# Dataset Analysis

- Commonsense: Is image necessary?



Is something under	yes	no
the sink broken?	yes	no
	yes	no

What number do	33	5
you see?	33	6
	33	7

Dataset	Input	All	Yes/No	Number	Other
Real	Question	40.81	67.60	25.77	21.22
	Question + Caption*	57.47	78.97	39.68	44.41
	Question + Image	83.30	95.77	83.39	72.67
Abstract	Question	43.27	66.65	28.52	23.66
	Question + Caption*	54.34	74.70	41.19	40.18
	Question + Image	87.49	95.96	95.04	75.33

# Dataset Analysis

- Commonsense needed? Age group

## **3-4 (15.3%)**

Is that a bird in the sky?

What color is the shoe?

How many zebras are there?

Is there food on the table?

Is this man wearing shoes?

## **5-8 (39.7%)**

How many pizzas are shown?

What are the sheep eating?

What color is his hair?

What sport is being played?

Name one ingredient in the skillet.

## **9-12 (28.4%)**

Where was this picture taken?

What ceremony does the cake commemorate?

Are these boats too tall to fit under the bridge?

What is the name of the white shape under the batter?

Is this at the stadium?

## **13-17 (11.2%)**

Is he likely to get mugged if he walked down a dark alleyway like this?

Is this a vegetarian meal?

What type of beverage is in the glass?

Can you name the performer in the purple costume?

Besides these humans, what other animals eat here?

## **18+ (5.5%)**

What type of architecture is this?

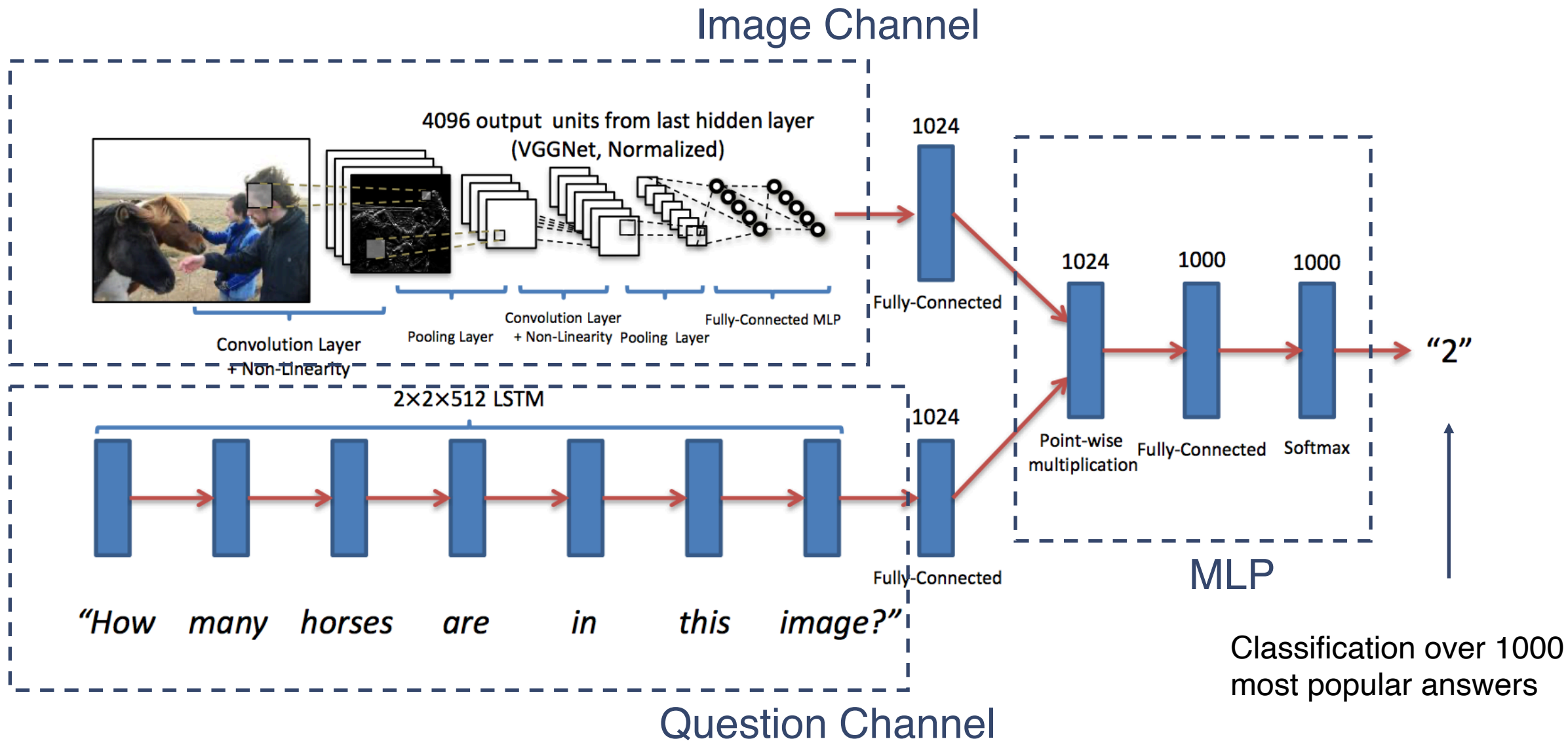
Is this a Flemish bricklaying pattern?

How many calories are in this pizza?

What government document is needed to partake in this activity?

What is the make and model of this vehicle?

# Model



# Results

	Open-Ended				Multiple-Choice			
	All	Yes/No	Number	Other	All	Yes/No	Number	Other
prior (“yes”)	29.66	70.81	00.39	01.15	29.66	70.81	00.39	01.15
per Q-type prior	37.54	71.03	35.77	09.38	39.45	71.02	35.86	13.34
nearest neighbor	42.70	71.89	24.36	21.94	48.49	71.94	26.00	33.56
BoW Q	48.09	75.66	36.70	27.14	53.68	75.71	37.05	38.64
I	28.13	64.01	00.42	03.77	30.53	69.87	00.45	03.76
BoW Q + I	52.64	75.55	33.67	37.37	58.97	75.59	34.35	50.33
LSTM Q	48.76	78.20	35.68	26.59	54.75	78.22	36.82	38.78
LSTM Q + I	53.74	78.94	35.24	36.42	57.17	78.95	35.80	43.41
deeper LSTM Q	50.39	78.41	34.68	30.03	55.88	78.45	35.91	41.13
deeper LSTM Q + norm I	<b>57.75</b>	<b>80.50</b>	<b>36.77</b>	<b>43.08</b>	<b>62.70</b>	<b>80.52</b>	<b>38.22</b>	<b>53.01</b>
Caption	26.70	65.50	02.03	03.86	28.29	69.79	02.06	03.82
BoW Q + C	54.70	75.82	40.12	42.56	59.85	75.89	41.16	52.53

Image alone  
performs poorly

# Results

	Open-Ended				Multiple-Choice			
	All	Yes/No	Number	Other	All	Yes/No	Number	Other
prior (“yes”)	29.66	70.81	00.39	01.15	29.66	70.81	00.39	01.15
per Q-type prior	37.54	71.03	35.77	09.38	39.45	71.02	35.86	13.34
nearest neighbor	42.70	71.89	24.36	21.94	48.49	71.94	26.00	33.56
BoW Q	48.09	75.66	36.70	27.14	53.68	75.71	37.05	38.64
I	28.13	64.01	00.42	03.77	30.53	69.87	00.45	03.76
BoW Q + I	52.64	75.55	33.67	37.37	58.97	75.59	34.35	50.33
LSTM Q	48.76	78.20	35.68	26.59	54.75	78.22	36.82	38.78
LSTM Q + I	53.74	78.94	35.24	36.42	57.17	78.95	35.80	43.41
deeper LSTM Q	50.39	78.41	34.68	30.03	55.88	78.45	35.91	41.13
deeper LSTM Q + norm I	<b>57.75</b>	<b>80.50</b>	<b>36.77</b>	<b>43.08</b>	<b>62.70</b>	<b>80.52</b>	<b>38.22</b>	<b>53.01</b>
Caption	26.70	65.50	02.03	03.86	28.29	69.79	02.06	03.82
BoW Q + C	54.70	75.82	40.12	42.56	59.85	75.89	41.16	52.53

Language-alone is surprisingly well

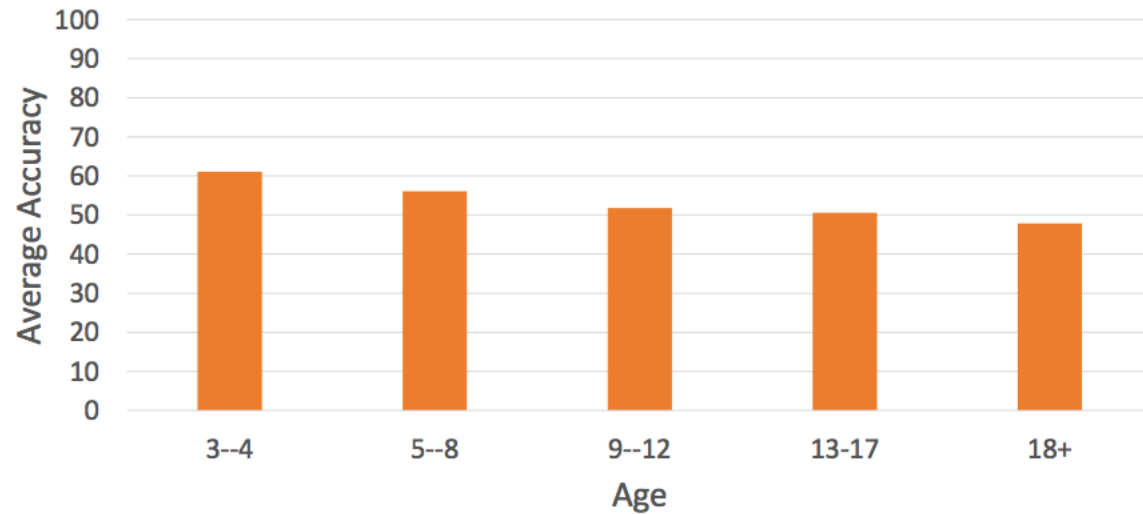


# Results

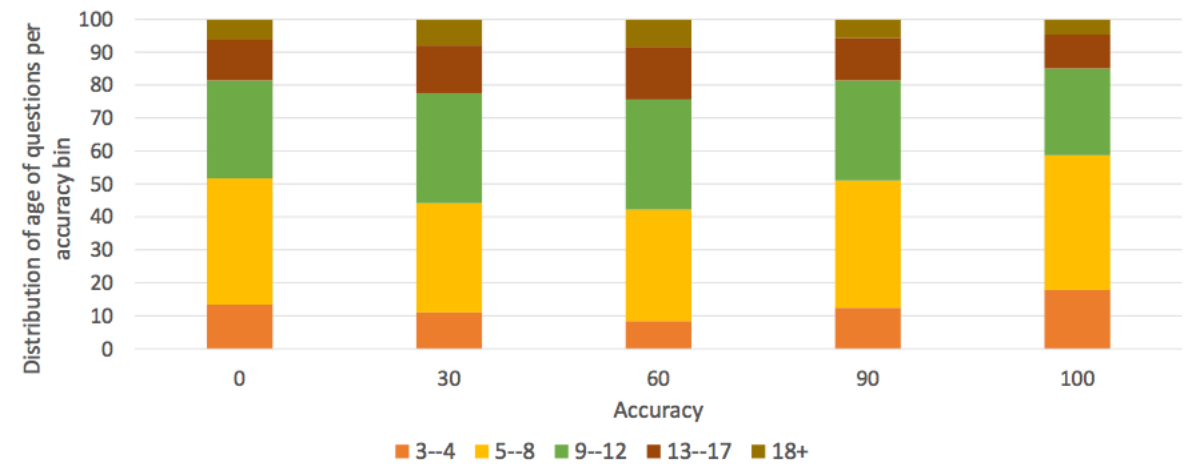
	Open-Ended				Multiple-Choice			
	All	Yes/No	Number	Other	All	Yes/No	Number	Other
prior (“yes”)	29.66	70.81	00.39	01.15	29.66	70.81	00.39	01.15
per Q-type prior	37.54	71.03	35.77	09.38	39.45	71.02	35.86	13.34
nearest neighbor	42.70	71.89	24.36	21.94	48.49	71.94	26.00	33.56
BoW Q	48.09	75.66	36.70	27.14	53.68	75.71	37.05	38.64
I	28.13	64.01	00.42	03.77	30.53	69.87	00.45	03.76
BoW Q + I	52.64	75.55	33.67	37.37	58.97	75.59	34.35	50.33
LSTM Q	48.76	78.20	35.68	26.59	54.75	78.22	36.82	38.78
LSTM Q + I	53.74	78.94	35.24	36.42	57.17	78.95	35.80	43.41
deeper LSTM Q	50.39	78.41	34.68	30.03	55.88	78.45	35.91	41.13
deeper LSTM Q + norm I	<b>57.75</b>	<b>80.50</b>	<b>36.77</b>	<b>43.08</b>	<b>62.70</b>	<b>80.52</b>	<b>38.22</b>	<b>53.01</b>
Caption	26.70	65.50	02.03	03.86	28.29	69.79	02.06	03.82
BoW Q + C	54.70	75.82	40.12	42.56	59.85	75.89	41.16	52.53

Combined sees  
significant gain

# Results



Accuracy by “age” of the question



“Age” of the question by accuracy

Model estimated to perform as well as a 4.74-year-old child



Thank you! Questions?

# The PhotoBook Dataset: Building Common Ground through Visually-Grounded Dialogue

Janosch Haber, Tim Baumgärtner, Ece Takmaz, Lieke Gelderloos, Elia Bruni, Raquel Fernández

<https://arxiv.org/pdf/1906.01530.pdf>

Presented By:  
Anant Dadu

# Contents

- Explanation of Visual Grounded Dialogue
- Shortcoming in Existing Works
- Task Setup
- Advantages
- Reference Chain
- Experiments
- Results

# Visual Grounded Dialogue

- The task of using natural language to communicate about visual input.
- The models developed for this task often focus on specific aspects such as image labelling, object reference, or question answering.

# Example

The little girl is standing with skis on her feet



## Human-Human Dialogue

what color are the skis ?  
Are there any other people?  
Is this outdoors?  
Do you see snow?  
Is it currently snowing?  
Is she on a slope or hill?  
Do you see trees?  
Do you see the sky?  
Is she wearing gloves?  
Is she wearing a hat?

A UNK color  
Not that i can see  
Yes  
Yes  
No I don't think so  
No i don't think so  
Yes  
No  
Yep  
yes

# Shortcoming in Existing Works

- Models fail to produce consistent outputs over a conversation.


**Reason:** It can be attributed to a **missing representation of the participant's shared common ground which** develops and extends during an interaction.

# Task Setup


- Two participants are paired for an online multi-round image identification game.
- Game Description:
  - Interface:**
    - page of a photo book (collection of 6 images)
    - some images are shown to both of them (common images) while other for each one of them are different
  - Task:**
    - mark these highlighted target images as either common or different by chatting with their partner.


# Screenshot of the Game Interface

Page 1 of 5





☒ Common ☐ Different






☐ Common ☐ Different







☐ Common ☐ Different

**YOU:** Do you have a man with two dogs on a bed?

**Robin:** With a purple wall in the background?

**YOU:** Yes

**Robin:** Then yes.

**Robin:** I have a little boy holding a phone to a teddy bear

**YOU:** I have that one as well

My next one is a boy sleeping with dolls

59 characters remaining.

Send

Figure 1: Screenshot of the Amazon Mechanical Turk user interface designed to collect the PhotoBook dataset.



# Advantages

- **Characteristic of dataset:** dialogues in the PhotoBook dataset contain multiple descriptions of each of the target images
- **Possible applications.:**
  - investigating participant cooperation
  - collaborative referring expression generation (single noun phrase for image)
  - description of image with respect to the conversation's common ground.

# Model

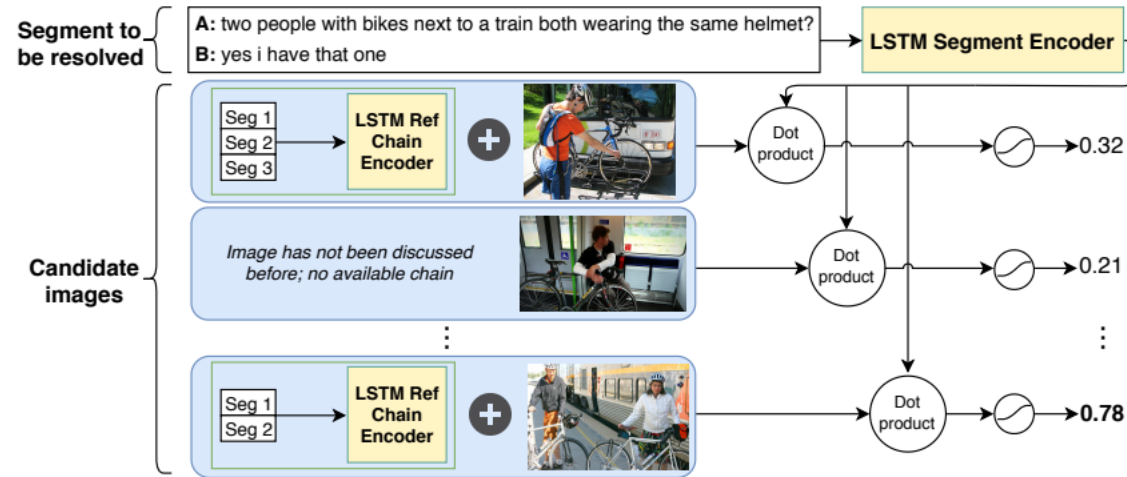


Figure 3: Diagram of the model in the HISTORY condition. For simplicity, we only show three candidate images. Some candidate images may not have a reference chain associated with them, while others may be linked to chains of different length, reflecting how many times an image has been referred to in the dialogue so far. In this example, the model predicts that the bottom candidate is the target referent of the segment to be resolved.

# Results

Model	Precision	Recall	F1
Random baseline	15.34	49.95	23.47
NO-HISTORY	56.65	75.86	64.86
HISTORY	56.66	77.41	65.43
HISTORY/No image	35.66	63.18	45.59

Table 3: Results for the target images in the test set.



**Reference chain with two segments:**  
 (1) A: a woman sitting in front of a monitor with a dog wallpaper while holding a plastic carrot  
 (2) B: carrot eating girl  
 A: no carrot eating girl on my end  
**Segment to be resolved:**  
 (4) B: I see the carrot lady again



**Reference chain with three segments:**  
 (1) A: I have a strange bike with two visible wheels in the back  
 (2) B: strange one  
 (3) A: strange bike again yes  
**Segment to be resolved:**  
 (4) B: strange

Figure 5: Reference chain for each of the two displayed images. The dialogue segments in the chains are slightly simplified for space reasons. **Left:** Both the HISTORY and the NO-HISTORY models succeed at identifying this image as the target of the segment to be resolved. **Right:** The NO-HISTORY model fails to recognise this image as the target of the segment to be resolved, while the HISTORY model succeeds. The distractor images for these two examples are available in Appendix E.

THANK YOU

# ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks



**Jiasen Lu<sup>1</sup>, Dhruv Batra<sup>1,2</sup>, Devi Parikh<sup>1,2</sup>, Stefan Lee<sup>1,3</sup>**

<sup>1</sup>Georgia Institute of Technology, <sup>2</sup>Facebook AI Research, <sup>3</sup>Oregon State University

# What is ViLBERT?

## Vision Language Tasks

Pretraining representation

ViLBERT

Finetuning

Pretrained on Conceptual caption dataset:  
(image, text) pairs



Is there something to cut the vegetables with?

VQA



Guy in yellow dribbling ball

Referring Expressions



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.

VCR Q→A

Rationale: a) is correct 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.

VCR QA→R

A large bus sitting next to a very tall building.

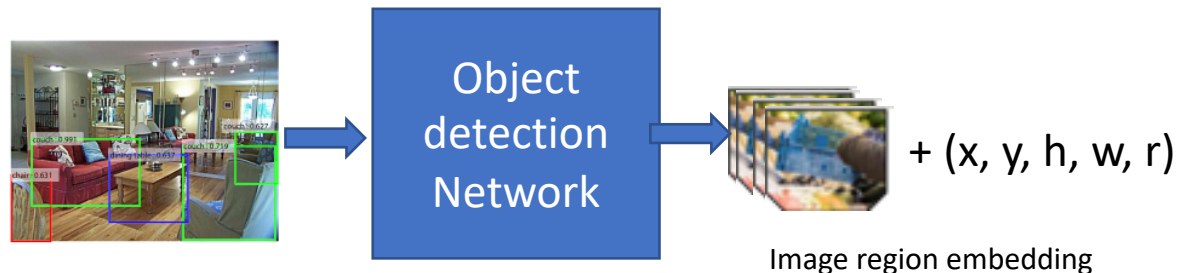
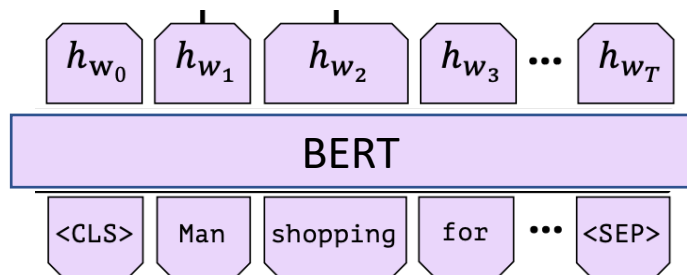


Caption-Based Image Retrieval

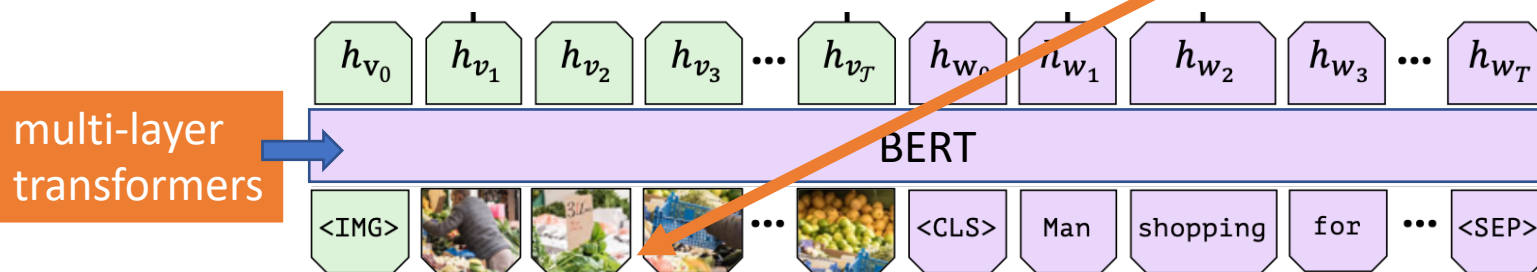


# From BERT to ViLBERT

- BERT



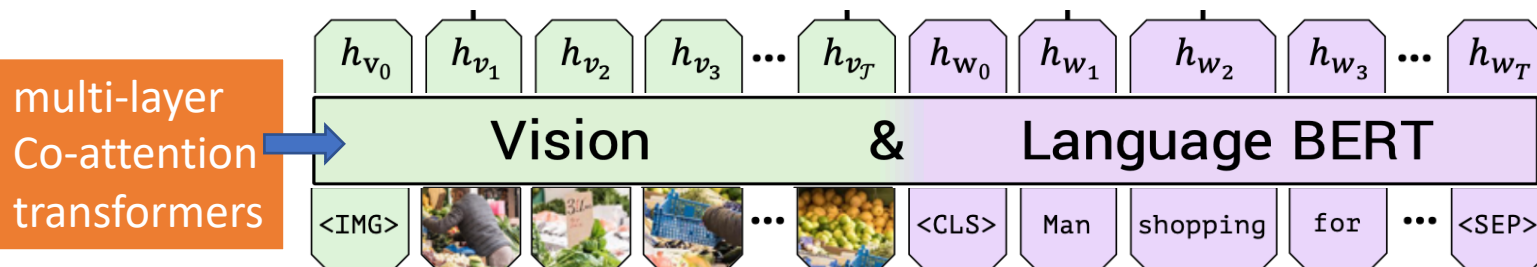
- Single Stream Vision Language BERT



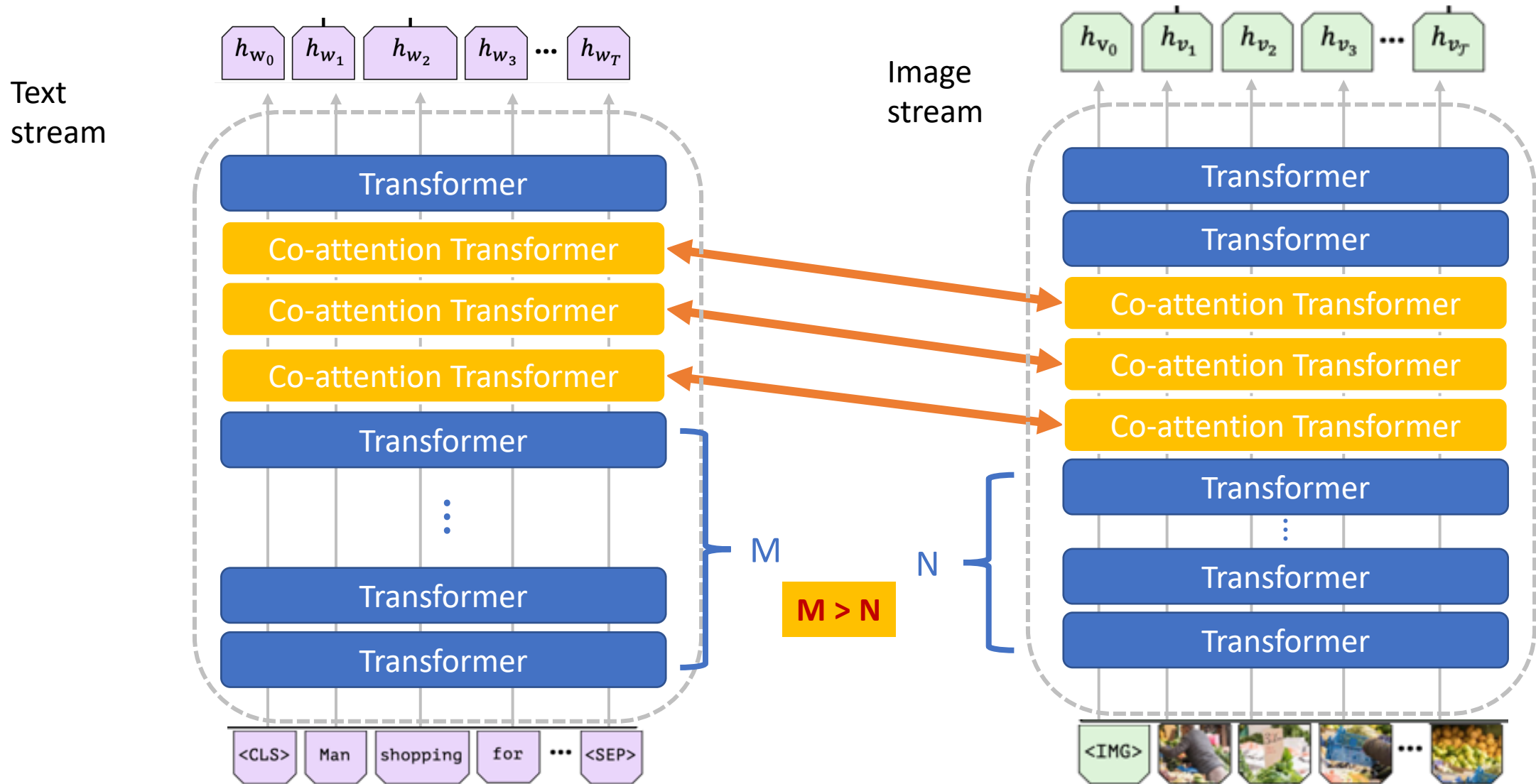
Problems:

Inputs from the two modalities are treated equally, but image region representation may be weaker as is already encoded by a deep network

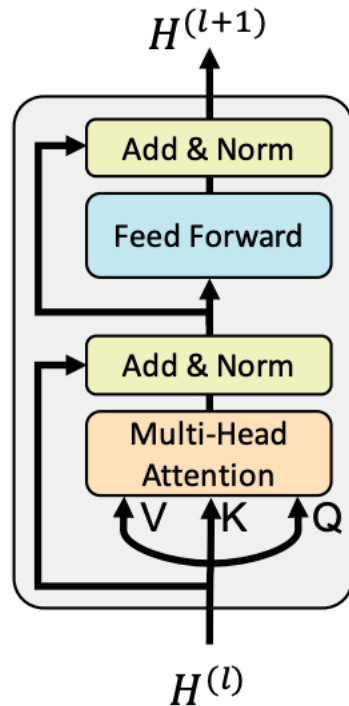
- ViLBERT (Co-Attention)



# The two streams model



# Transformer Layers



(a) Standard encoder transformer block

## Self-attention w/ queries, keys, values

Let's add learnable parameters ( $k \times k$  weight matrices), and turn each vector  $\mathbf{x}^{(i)}$  into three versions:

- **Query** vector  $\mathbf{q}^{(i)} = \mathbf{W}_q \mathbf{x}^{(i)}$
- **Key** vector:  $\mathbf{k}^{(i)} = \mathbf{W}_k \mathbf{x}^{(i)}$
- **Value** vector:  $\mathbf{v}^{(i)} = \mathbf{W}_v \mathbf{x}^{(i)}$

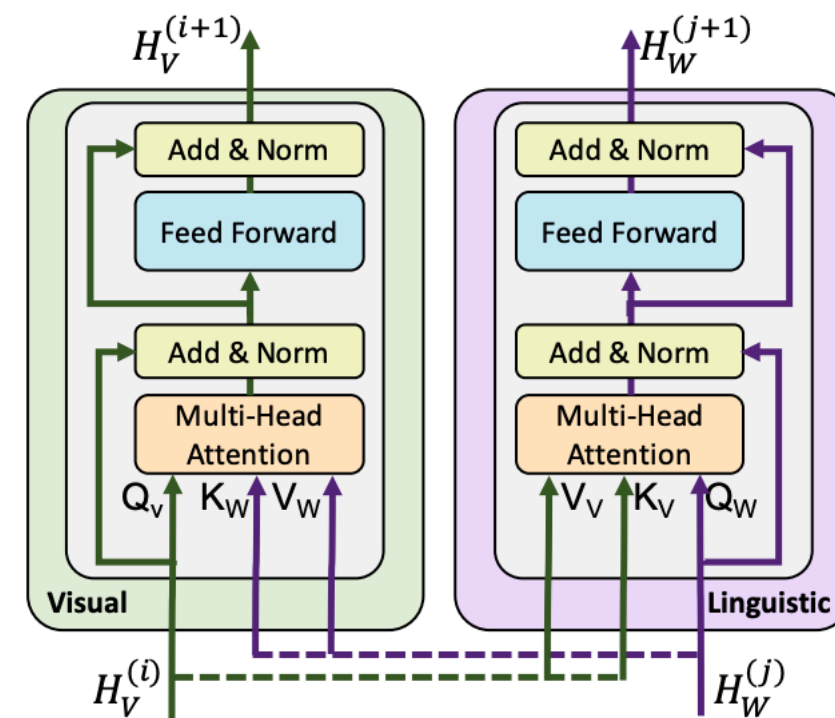
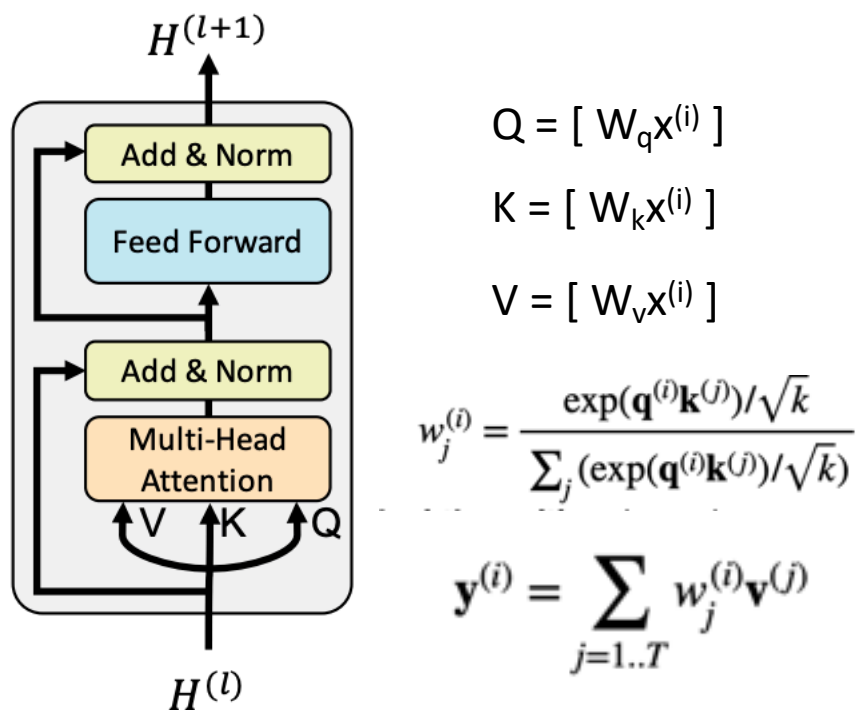
The **attention weight of the  $j$ -th position** to compute the **new output for the  $i$ -th position** depends on the **query of  $i$**  and the **key of  $j$  (scaled)**:

$$w_j^{(i)} = \frac{\exp(\mathbf{q}^{(i)} \mathbf{k}^{(j)} / \sqrt{k})}{\sum_j (\exp(\mathbf{q}^{(i)} \mathbf{k}^{(j)} / \sqrt{k}))}$$

The **new output vector for the  $i$ -th position** depends on the **attention weights** and **value** vectors of all **input positions  $j$** :

$$\mathbf{y}^{(i)} = \sum_{j=1..T} w_j^{(i)} \mathbf{v}^{(j)}$$

# Co-Attention Transformer Layers



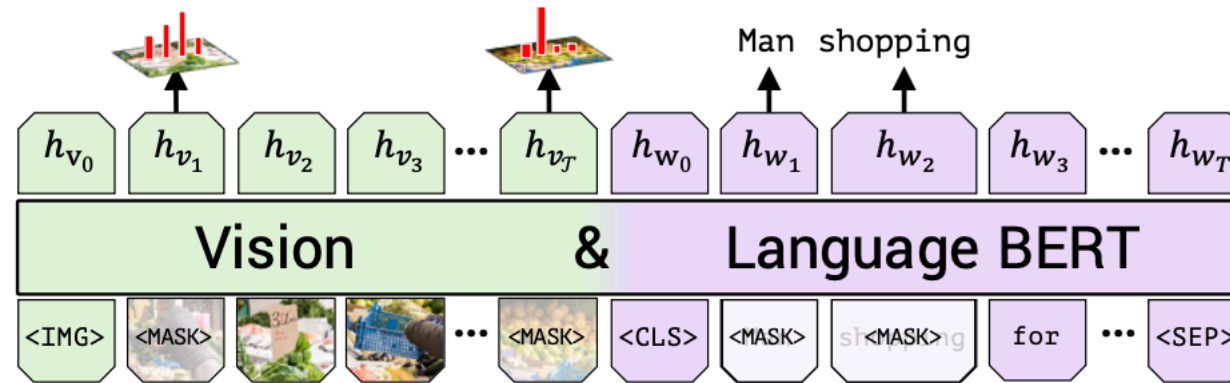
(a) Standard encoder transformer block

(b) Our co-attention transformer layer

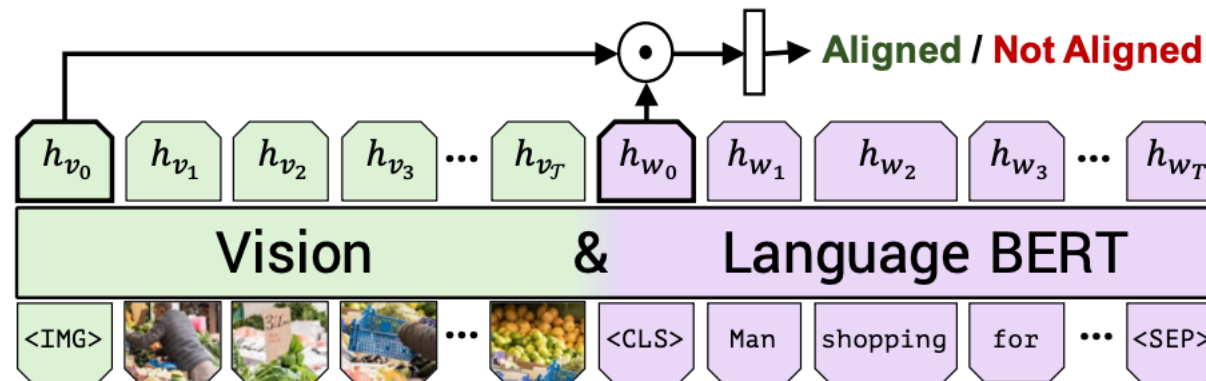
1. Two modalities have separate streams
2. Keys and values from each modality are passed as input to the other modality's multi-headed attention blocks.
3. The attention-pooled features for each modality conditioned on the other

# Training tasks (Objectives)

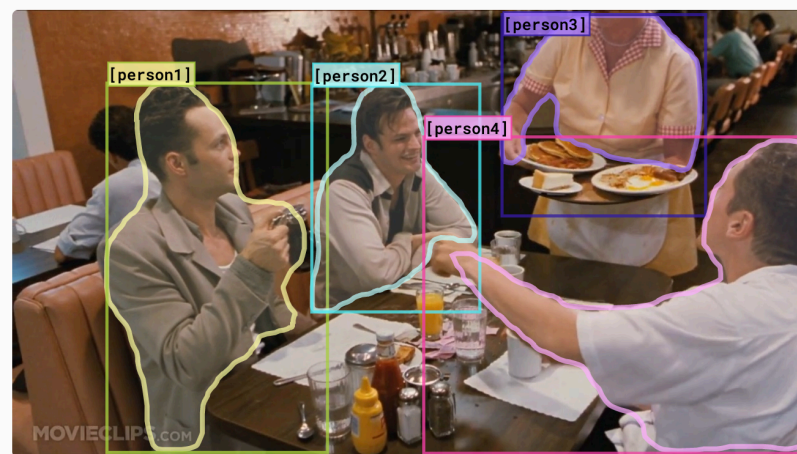
- Masked Multi-modal learning



- Multi-modal alignment prediction



# Finetuning – Visual Commonsense Reasoning



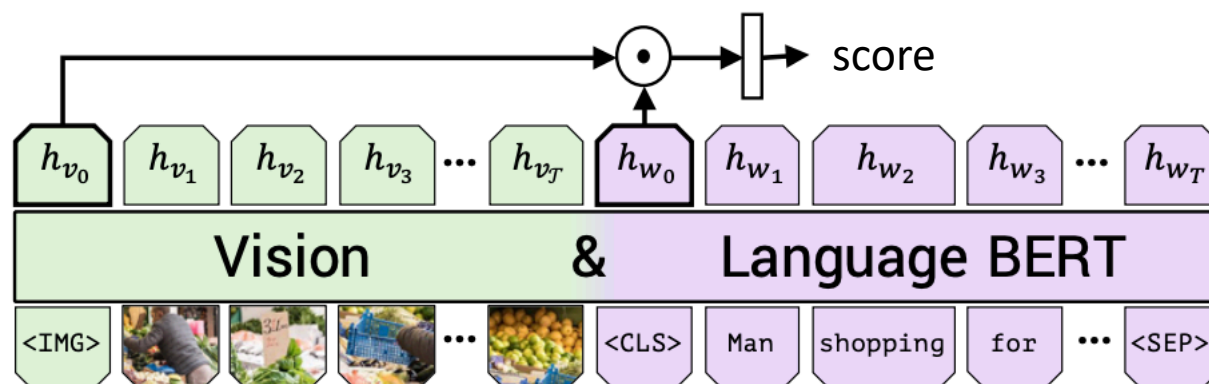
hide all show all [person1] [person2] [person3] [person4]  
more objects »

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.



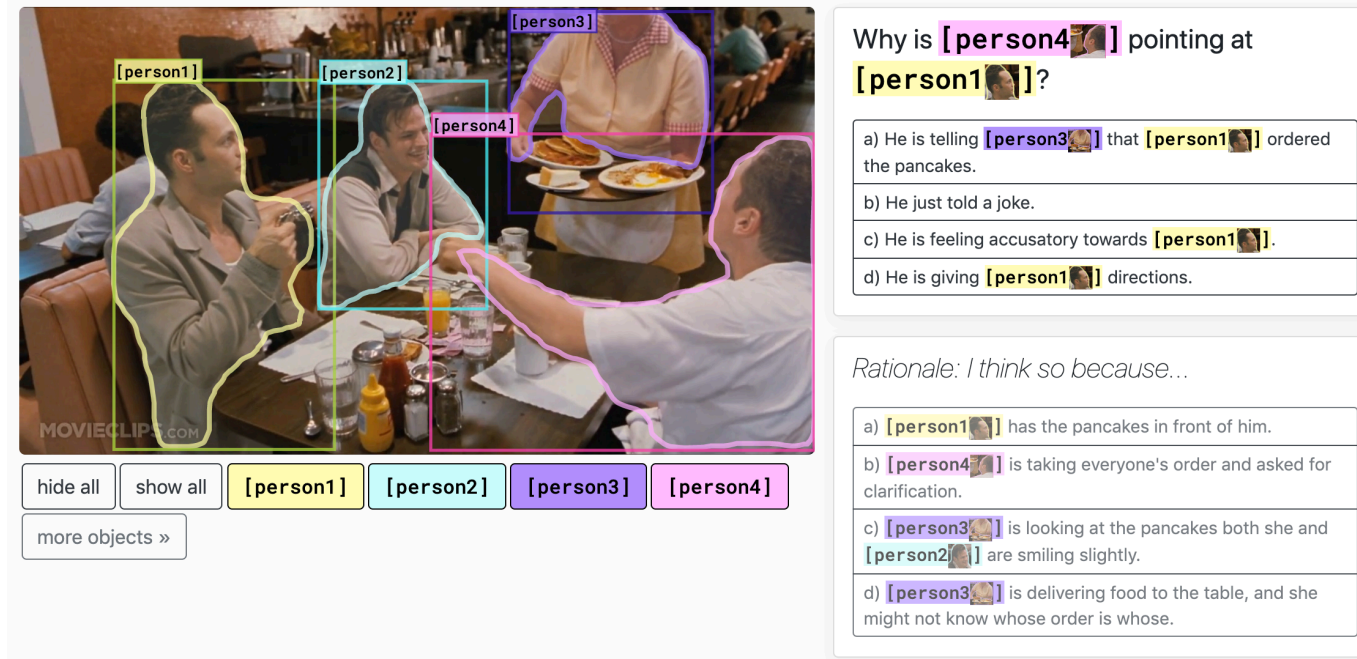
Question + One of the candidate answer

Train: **softmax** + **cross-entropy** (correct 1 /wrong 0)

Test: select the candidate answer with the **max predicted score**



# Finetuning – Visual Commonsense Reasoning



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.

	Q->A	QA->R	Q->AR
SOTA	63.8	67.2	43.1
<b>ViLBERT</b>	<b>72.42</b>	<b>74.47</b>	<b>54.04</b>

# Finetuning – Grounding Referring Expressions

RefCOCO+ testA

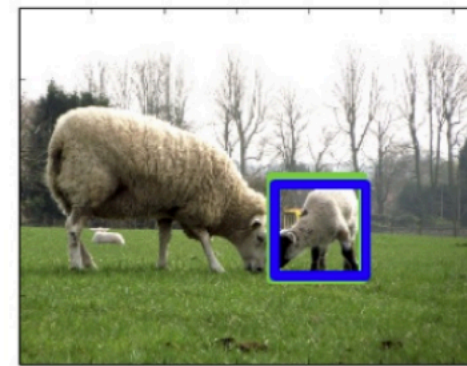


blurry person  
with sleeveless and sitting



man in full view in all black

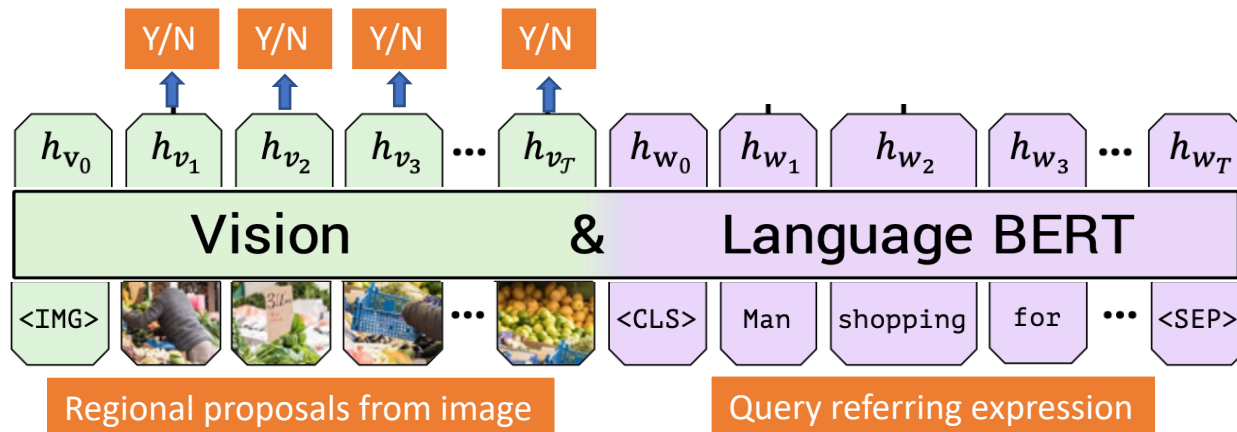
RefCOCO+ testB



small one grazing



books about bears



Train: **softmax** + **cross-entropy** (1 for correct; 0 for wrong)

Test : Select region with the **max predicted score**

# Finetuning – Grounding Referring Expressions

RefCOCO+ testA

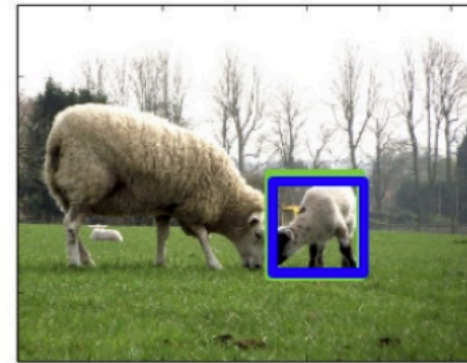


blurry person  
with sleeveless and sitting



man in full view in all black

RefCOCO+ testB



small one grazing

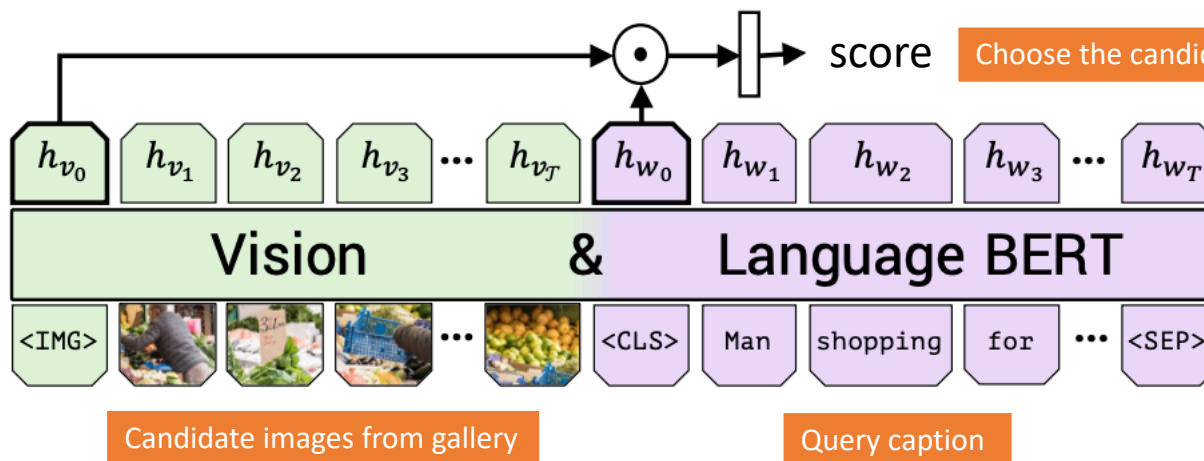


books about bears

	Val	testA	testB
SOTA	65.33	71.62	56.02
<b>ViLBERT</b>	<b>72.34</b>	<b>78.52</b>	<b>62.61</b>

# Finetuning – Caption-based Image Retrieval

- Query: A woman sings on stage as a man plays an instrument.
- Gallery:



Train: **softmax** + **cross-entropy** on each region embedding  
(1 for correct; 0 for wrong)

neg pairs: (rand img, cap) (img, rand cap) (hard img, cap)

Test : Select region with the **max predicted score**



# Finetuning – Caption-based Image Retrieval

- Query: A woman sings on stage as a man plays an instrument.
- Gallery:



	Q->A	QA->R	Q->AR
SOTA	48.60	77.70	85.20
<b>ViLBERT</b>	<b>58.20</b>	<b>84.90</b>	<b>91.52</b>

# References

- Lu, Jiasen, et al. "Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks." *Advances in Neural Information Processing Systems*. 2019.
- Devlin, Jacob, et al. "Bert: Pre-training of deep bidirectional transformers for language understanding." *arXiv preprint arXiv:1810.04805* (2018).
- Antol, Stanislaw, et al. "Vqa: Visual question answering." *Proceedings of the IEEE international conference on computer vision*. 2015.
- Zellers, Rowan, et al. "From recognition to cognition: Visual commonsense reasoning." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2019.
- Kazemzadeh, Sahar, et al. "Referitgame: Referring to objects in photographs of natural scenes." *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 2014.
- Young, Peter, et al. "From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions." *Transactions of the Association for Computational Linguistics* 2 (2014): 67-78.