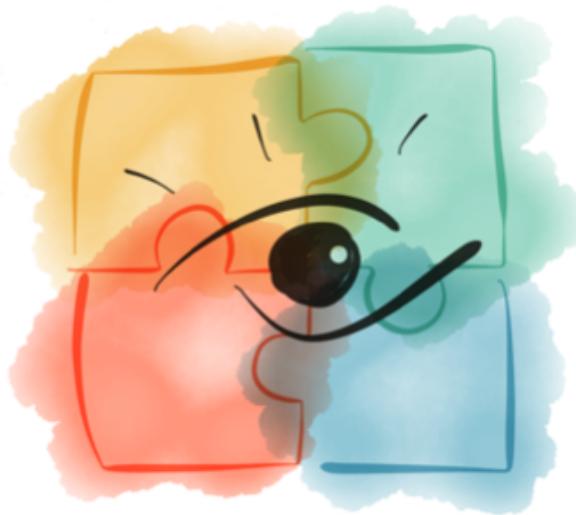


Forcing Vision and Language Models to Not Just Talk But Also Actually See



Devi Parikh

Pythia for Vizwiz



Vivek
Natarajan*



Tina
Jiang*



Meet
Shah



Xinlei
Chen



Dhruv
Batra



Devi
Parikh



Marcus
Rohrbach

* - indicates equal contribution

facebook Artificial Intelligence Research



facebook
Artificial Intelligence Research | A-STAR

Pythia

Pythia is a modular framework for Visual Question Answering research, which formed the basis for the winning entry to the VQA Challenge 2018 from Facebook AI Research (FAIR)'s A-STAR team. Please check our [paper](#) for more details.

(A-STAR: Agents that See, Talk, Act, and Reason.)



Q: What is the cat wearing?
A: Hat



Q: What is the weather like?
A: Rainy



Q: What surface is this?
A: Clay



Q: What is the weather like?
A: Sunny



Q: What color is the cat's eyes?
A: Green



Q: What toppings are on the pizza?
A: Mushrooms

VizWiz Challenge Entry: 10:50 am to 11:20 am

Pythia for Vizwiz



Vivek
Natarajan*



Tina
Jiang*



Meet
Shah



Xinlei
Chen



Dhruv
Batra



Devi
Parikh



Marcus
Rohrbach

* - indicates equal contribution

facebook Artificial Intelligence Research



facebook
Artificial Intelligence Research | A-STAR



People coloring a street in rural Virginia.



It was a great event! It brought families out, and the whole community together.



Q. What are they coloring the street with?

A. Chalk



An aerial photograph of a street festival. The street is covered with numerous colorful chalk drawings, including various shapes, letters, and figures. People are walking along the street, and there are some buildings and trees in the background.

AI: What a nice picture! What event was this?

User: *“Color College Avenue”. It was a lot of fun!*

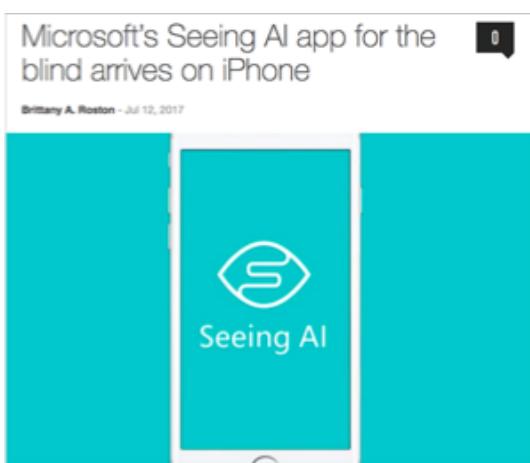
AI: I am sure it was! Do they do this every year?

User: *I wish they would. I don't think they've organized it again since 2012.*

...

Applications

FACEBOOK'S AI CAN CAPTION PHOTOS FOR THE BLIND ON ITS OWN



VQA Challenge @ CVPR16, 17, 18

Competition



VQA Real Image Challenge (Open-Ended)

Organized by vqateam - Current server time: March 22, 2016, 5 a.m. UTC

▶ Current

Next

Learn

Overview

Evaluation

Terms and Conditions

State-of-the-art: 54% → 72%

Visual Question Answering (VQA)



What color are her eyes?
What is the mustache made of?



How many slices of pizza are there?
Is this a vegetarian pizza?



Is this person expecting company?
What is just under the tree?



Does it appear to be rainy?
Does this person have 20/20 vision?

Recent progress in computer vision and natural language processing has demonstrated that lower-level tasks are much closer to being solved. We believe that the time is ripe to pursue

the challenging task of visual question answering.

Success Cases



Q: What is the woman holding?

GT A: laptop

Machine A: laptop



Q: What room is the cat located in?

GT A: kitchen

Machine A: kitchen



Q: Is it going to rain soon?

GT A: yes

Machine A: yes



Q: Is this a casino?

GT A: no

Machine A: no

Models affected by language priors

Test Sample



Q: What color
are the
safety cones?

Test Sample

Nearest Neighbor Training Samples



Q: What color
are the
safety cones?

GT Ans: green

Predicted Ans: orange

Test Sample



Q: What color
are the
safety cones?

GT Ans: green

Nearest Neighbor Training Samples



Q: What color
are the
cones?

GT Ans: orange

Predicted Ans: orange

Test Sample



Q: What color
are the
safety cones?

GT Ans: green

Nearest Neighbor Training Samples



Q: What color
are the
cones?

GT Ans: orange



Q: What color
is the
cone?

GT Ans: orange

Predicted Ans: orange

Test Sample



Q: What color
are the
safety cones?

GT Ans: green

Predicted Ans: orange

Nearest Neighbor Training Samples



Q: What color
are the
cones?

GT Ans: orange



Q: What color
is the
cone?

GT Ans: orange



Q: What color
are
the cones?

GT Ans: orange



Q: Are **A:** military

Q: Are they **A:** yes

Q: Are they playing **A:** yes

Q: Are they playing a **A:** yes

Q: Are they playing a game? **A:** yes

GT Ans: yes



Q: How **A:** no

Q: How many **A:** 2

Q: How many horses **A:** 2

Q: How many horses are **A:** 2

Q: How many horses are on **A:** 2

Q: How many horses are on the **A:** 2

Q: How many horses are on the beach? **A:** 2

GT Ans: 6



Q: Is **A:** kitchen

Q: Is the **A:** outside

Q: Is the bench **A:** no

Q: Is the bench made **A:** no

Q: Is the bench made of **A:** no

Q: Is the bench made of metal? **A:** no

GT Ans: yes

Q: What does the red sign say?

Predicted Ans: stop

Correct Response



Q: How many zebras?

Predicted Ans: 2

Correct Response



Q: What covers the ground?

Predicted Ans: snow

All Correct Responses



How do we force vision+language
models to also *look*?

Outline

- Datasets
 - E.g., VQA v2.0
- Evaluation metrics
 - E.g., # pairs of images correctly answered
- Novel problem spaces
 - E.g., changing priors (VQA-CP), novel object captioning, robust captioning, visual coreference resolution
- Inductive biases in models
 - E.g., G-VQA, Neural Baby Talk
- Human-in-the-loop evaluation
 - E.g., accuracy is not the final metric, performance of human-AI team is

Outline

- Datasets
 - E.g., VQA v2.0
- Evaluation metrics
 - E.g., # pairs of images correctly answered
- Novel problem spaces
 - E.g., changing priors (VQA-CP), novel object captioning, robust captioning, visual coreference resolution
- Inductive biases in models
 - E.g., G-VQA, Neural Baby Talk
- Human-in-the-loop evaluation
 - E.g., accuracy is not the final metric, performance of human-AI team is

Balancing the VQA dataset

What game is this?
Tennis



Balancing the VQA dataset



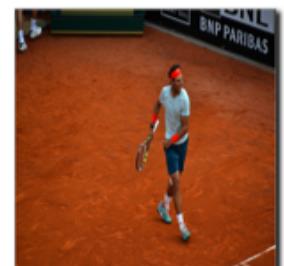
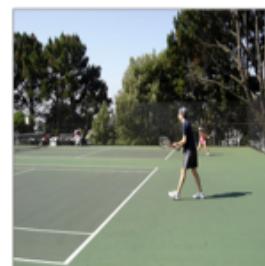
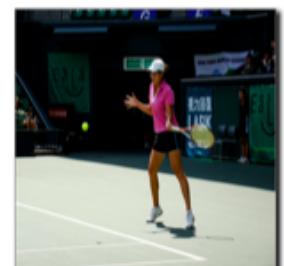
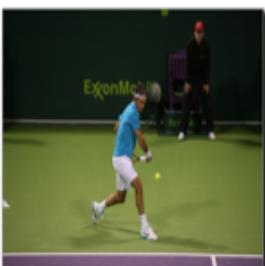
PAGE 1/5

SHOW INSTRUCTIONS

NOT POSSIBLE

PREVIOUS

NEXT



Balancing the VQA dataset

Is the TV on?

yes



no



Balancing the VQA dataset

How many pets are present?

2



1



Balancing the VQA dataset

What sign is this?

handicap



one way



Balancing the VQA dataset

Where is the child sitting?

fridge



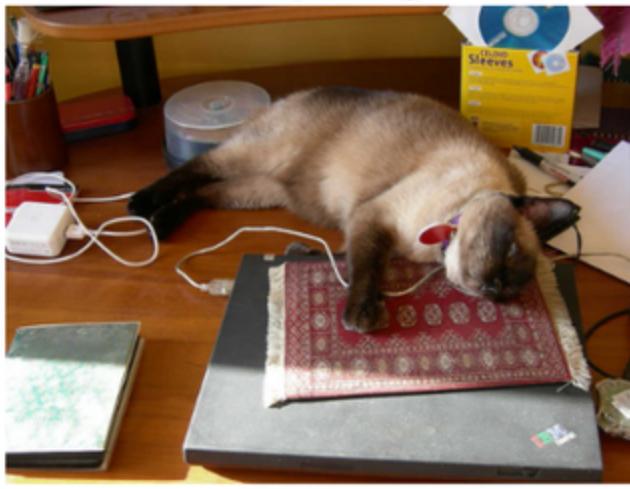
arms



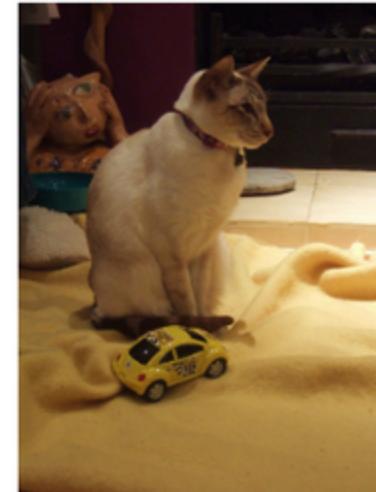
Balancing the VQA dataset

What is the cat doing on the rug?

sleeping



sitting



Balancing the VQA dataset

What color are the pants?

orange



brown



VQA v2.0

- More balanced than VQA v1.0
 - Entropy of answers increases by 56%
- Bigger than VQA v2.0
 - ~1.8 times image-question pairs

Benchmarking SOTA VQA models

- SOTA VQA models
 - Drop in performance by 7-8%
 - Gain 1-2% back when re-trained on balanced dataset
- By answer types
 - Biggest drop in performance in yes/no (10-12%)
 - Biggest improvement gained by re-training in yes/no (3-4%) and number (2-3%)

Trends

	By Answer Type			Overall
	Yes/No	Number	Other	
UC Berkeley & Sony ^[14]	83.79	38.9	58.64	66.9
Naver Labs ^[10]	83.78	37.67	54.74	64.89
DLAIT ^[5]	83.65	39.18	52.62	63.97
snubi-naverlabs ^[25]	83.64	38.43	51.61	63.4

Red arrows indicate changes from the previous row:

- UC Berkeley & Sony to Naver Labs: 0.15% increase in Yes/No
- Naver Labs to DLAIT: 1.51% increase in Number
- DLAIT to snubi-naverlabs: 7.03% decrease in Other
- UC Berkeley & Sony to snubi-naverlabs: 3.5% decrease in Overall

VQA v2.0

2nd and 3rd VQA Challenges @ CVPR17, 18.

Removing Language Priors

Scene 1/3 - Also need at least: 1 person

You must ACCEPT the HIT before you can start the real task.

Prev

Next

Question Is there a place to sit other than the floor?

Answer yes



Scene Depth



Flip



Type



Want to work on this HIT?

Accept HIT

Want to see other HITs?

Skip HIT

Removing Language Priors

Answer: No



Answer: Yes



complementary scenes

Question: Is the girl walking the bike?

Outline

- Datasets
 - E.g., VQA v2.0
- Evaluation metrics
 - E.g., # pairs of images correctly answered
- Novel problem spaces
 - E.g., changing priors (VQA-CP), novel object captioning, robust captioning, visual coreference resolution
- Inductive biases in models
 - E.g., G-VQA, Neural Baby Talk
- Human-in-the-loop evaluation
 - E.g., accuracy is not the final metric, performance of human-AI team is

Classifying a pair of complementary scenes

	Training set	
	Unbalanced	Balanced
Blind (no image features)		
Holistic image features		
Attention-based image features		

Outline

- Datasets
 - E.g., VQA v2.0
- Evaluation metrics
 - E.g., # pairs of images correctly answered
- Novel problem spaces
 - E.g., changing priors (VQA-CP), novel object captioning, robust captioning, visual coreference resolution
- Inductive biases in models
 - E.g., G-VQA, Neural Baby Talk
- Human-in-the-loop evaluation
 - E.g., accuracy is not the final metric, performance of human-AI team is

(A related) problem with existing setup

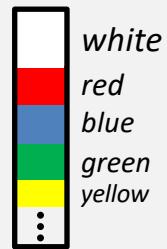
Train

Q: What color is the dog?

A: White



Training Prior



(A related) problem with existing setup

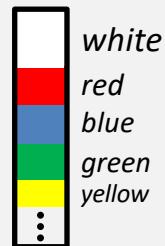
Train

Q: What color is the dog?

A: White



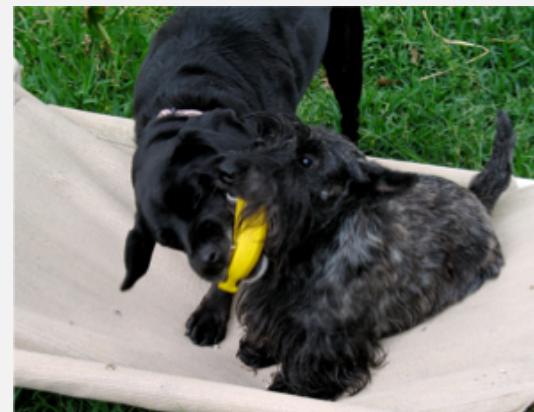
Training Prior



Test

Q: What color is the dog?

A: Black



(A related) problem with existing setup

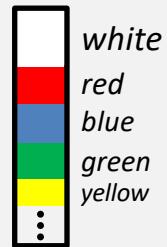
Train

Q: What color is the dog?

A: White



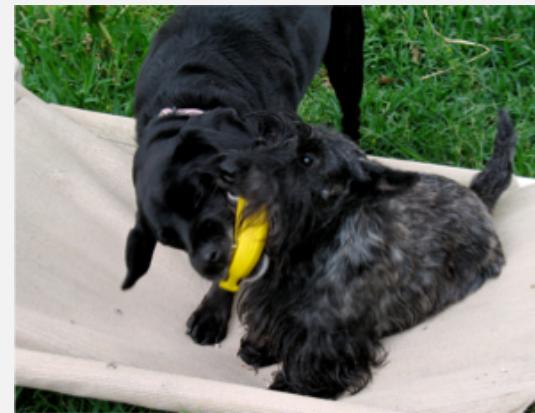
Training Prior



Test

Q: What color is the dog?

A: Black



Prediction:
White

(A related) problem with existing setup

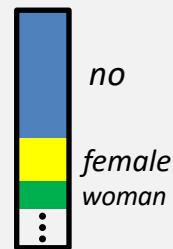
Train

Q: Is the person wearing shorts?

A: No



Training Prior



Test

Q: Is the person wearing shorts?

A: Yes



Prediction:

No

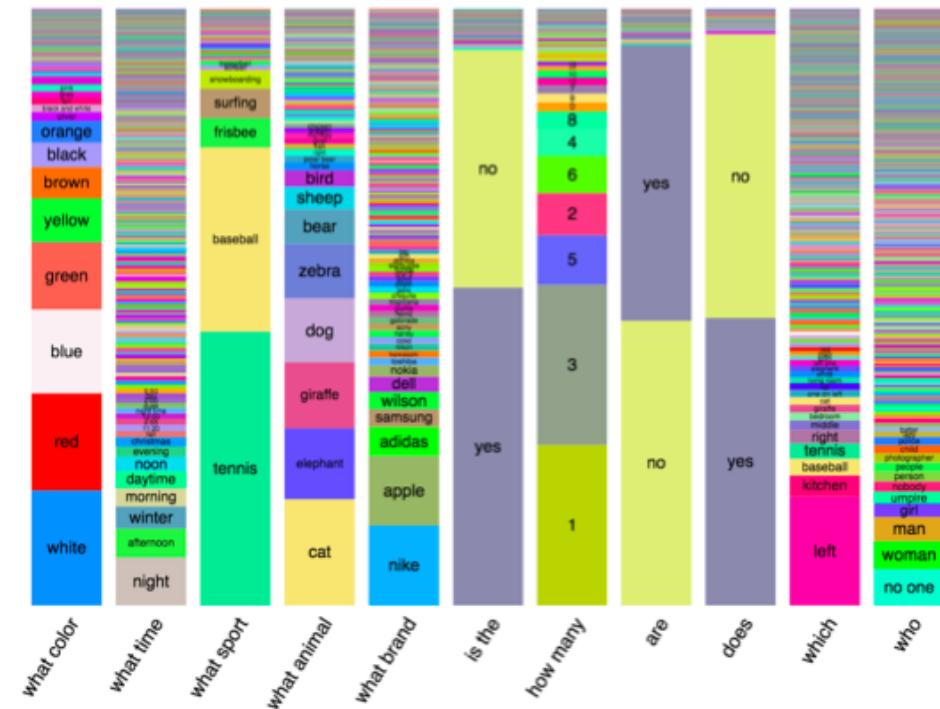
(A related) problem with existing setup

- Similar priors in train and test
- Memorization does not hurt as much
- Problematic for benchmarking progress

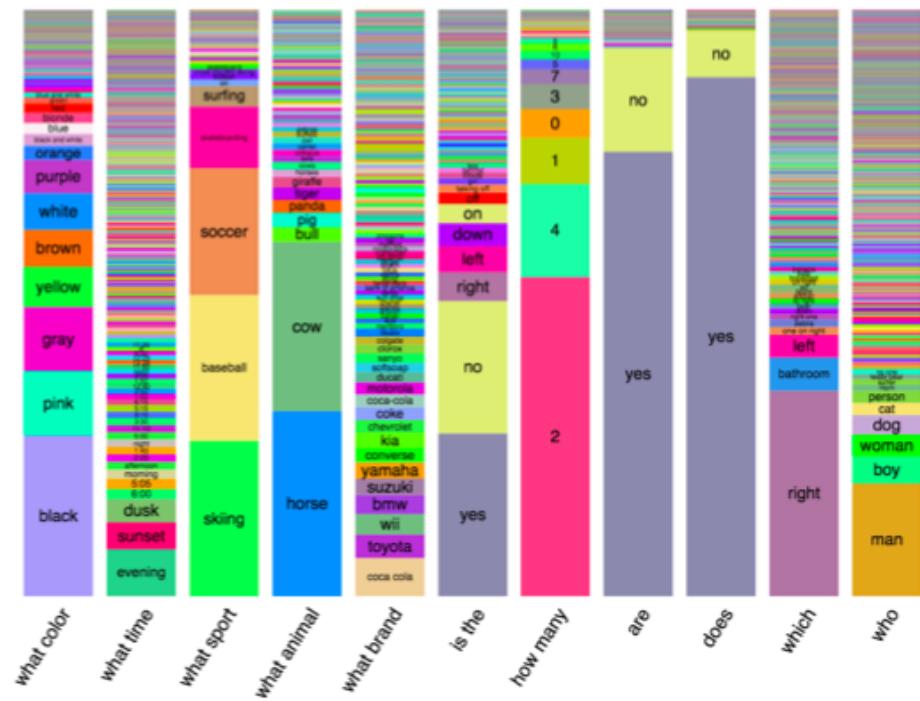
Meet VQA-CP!

- Visual Question Answering under Changing Priors
- A new split of the VQA v1.0 dataset (Antol et al., ICCV 2015)

VQA-CP Train Split



VQA-CP Test Split



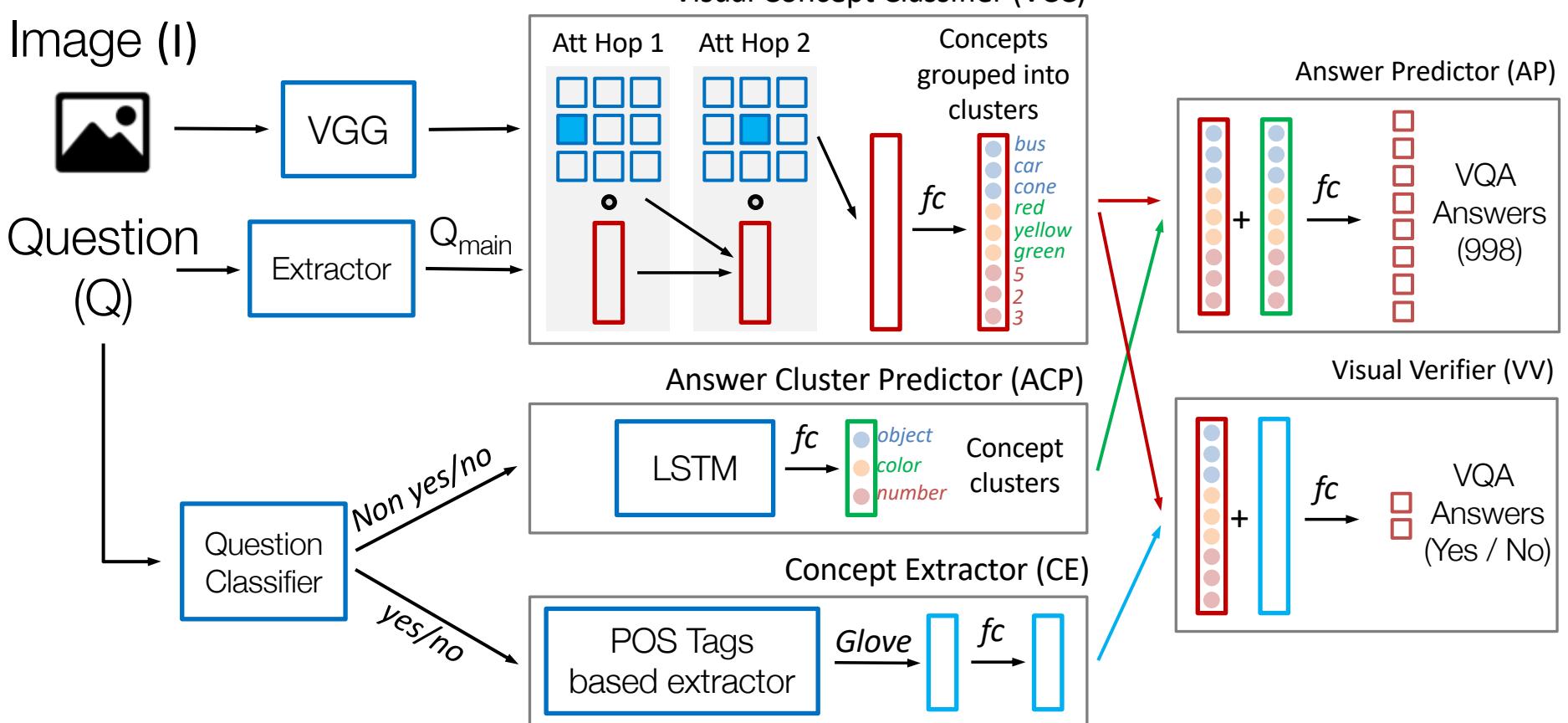
Performance of VQA models on VQA-CP

Model	Dataset	Overall	Yes/No	Number	Other	
d-LSTM Q + norm I (Antol et al. ICCV15)	VQA	54.23	79.81	33.26	40.35	31% drop
	VQA-CP	23.51	34.53	11.40	17.42	
NMN (Andreas et al. CVPR16)	VQA	54.83	80.39	33.45	41.07	25% drop
	VQA-CP	29.64	38.85	11.23	27.88	
SAN (Yang et al. CVPR16)	VQA	55.86	78.54	33.46	44.51	29% drop
	VQA-CP	26.88	35.34	11.34	24.70	
MCB (Fukui et al. EMNLP16)	VQA	60.97	81.62	34.56	52.16	27% drop
	VQA-CP	34.39	37.96	11.80	39.90	

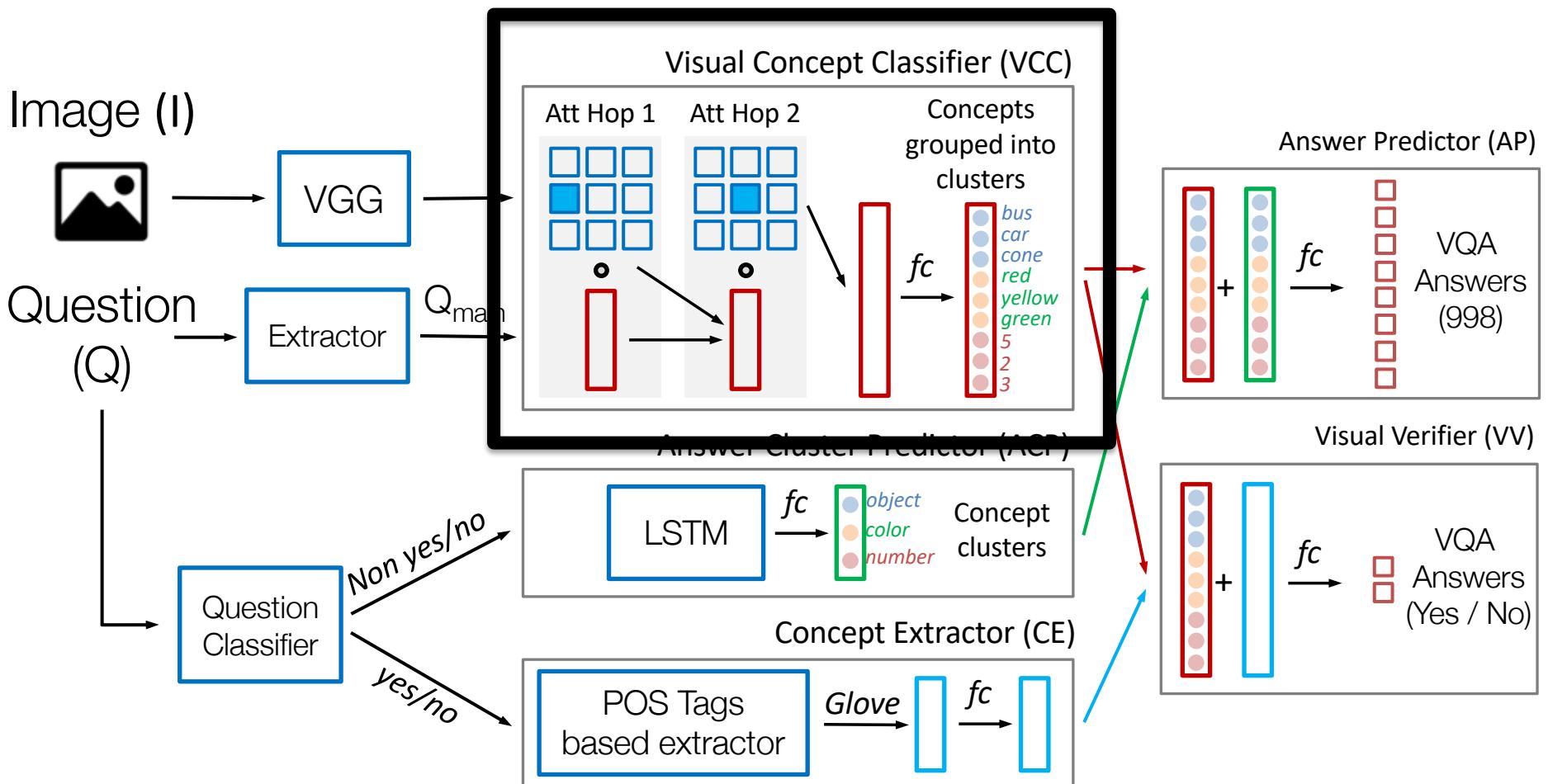
Outline

- Datasets
 - E.g., VQA v2.0
- Evaluation metrics
 - E.g., # pairs of images correctly answered
- Novel problem spaces
 - E.g., changing priors (VQA-CP), novel object captioning, robust captioning, visual coreference resolution
- Inductive biases in models
 - E.g., G-VQA, Neural Baby Talk
- Human-in-the-loop evaluation
 - E.g., accuracy is not the final metric, performance of human-AI team is

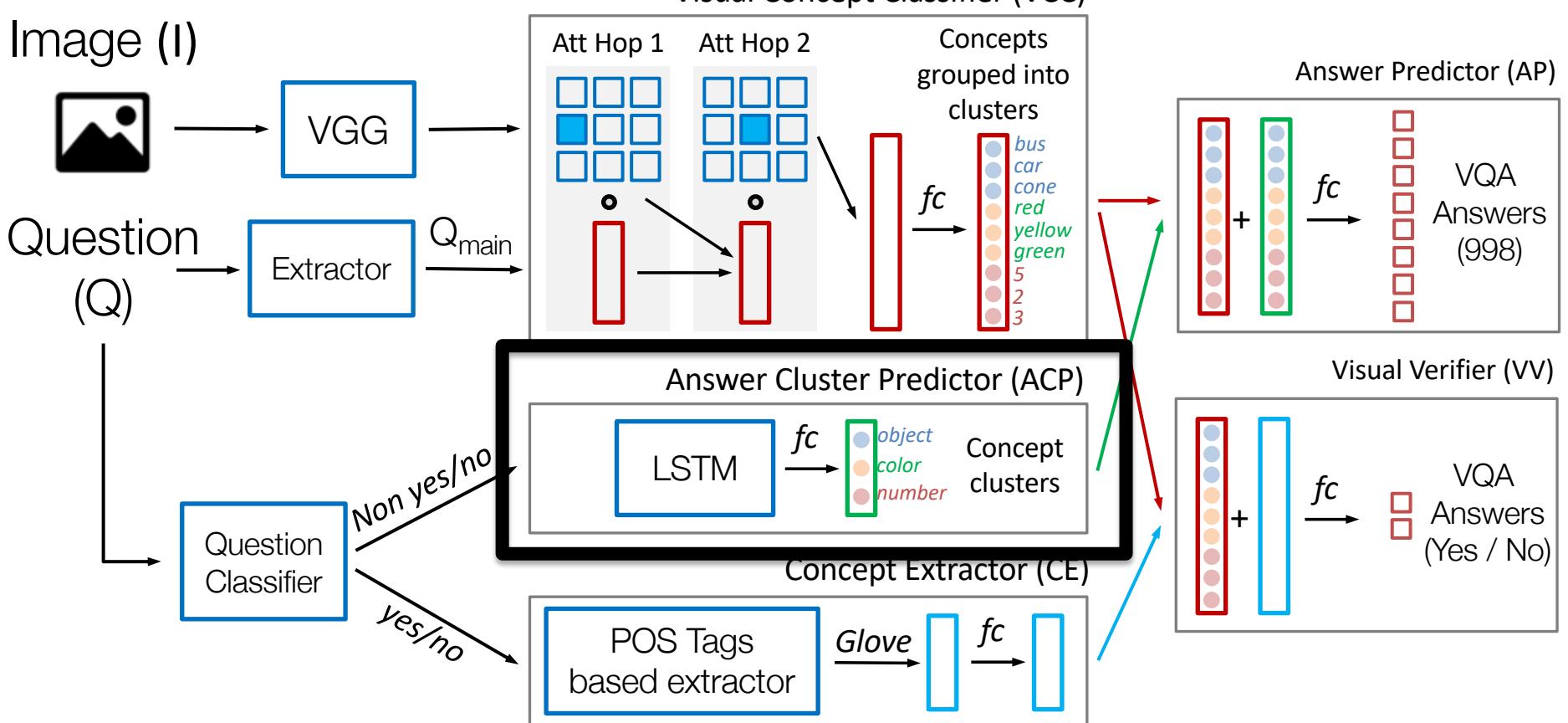
Grounded-VQA (GVQA)



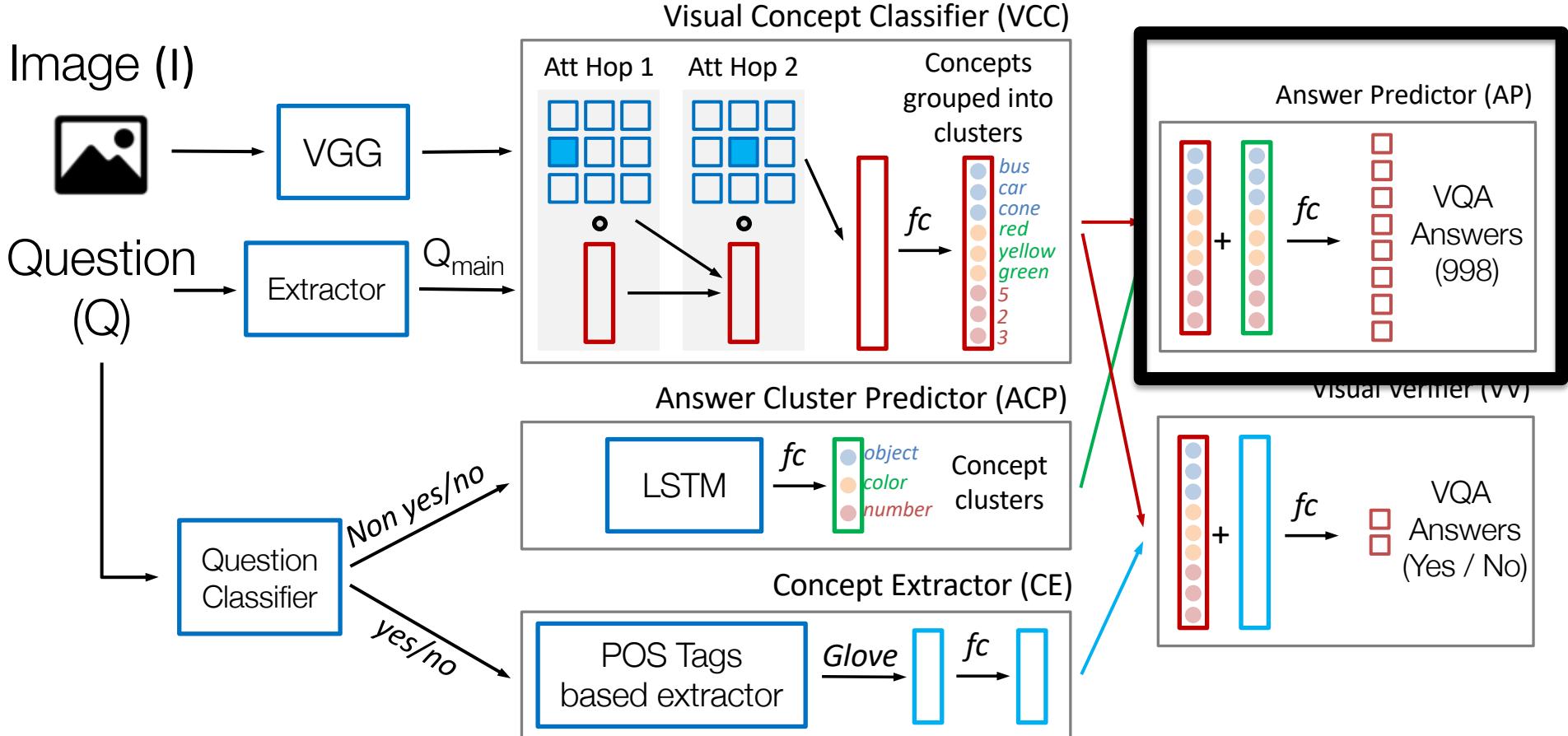
Grounded-VQA (GVQA)



Grounded-VQA (GVQA)



Grounded-VQA (GVQA)



Visual Coreference Resolution in Visual Dialog using Neural Module Networks



C : The **boat** has a **dragon head** on the

Q1 : Is the **boat** on **water** ?

A1 : Yes

Q2 : What color is **it** ?

A2 : White

Q3 : Does the **dragon** have a body ?

A3 : No, just the **head**

Q4 : What color is **it** ?

Known Entities
boat

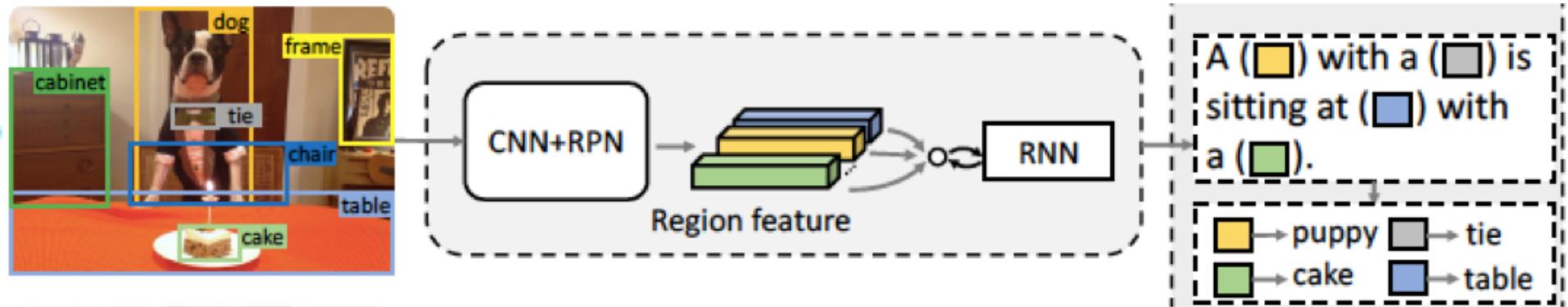


dragon head

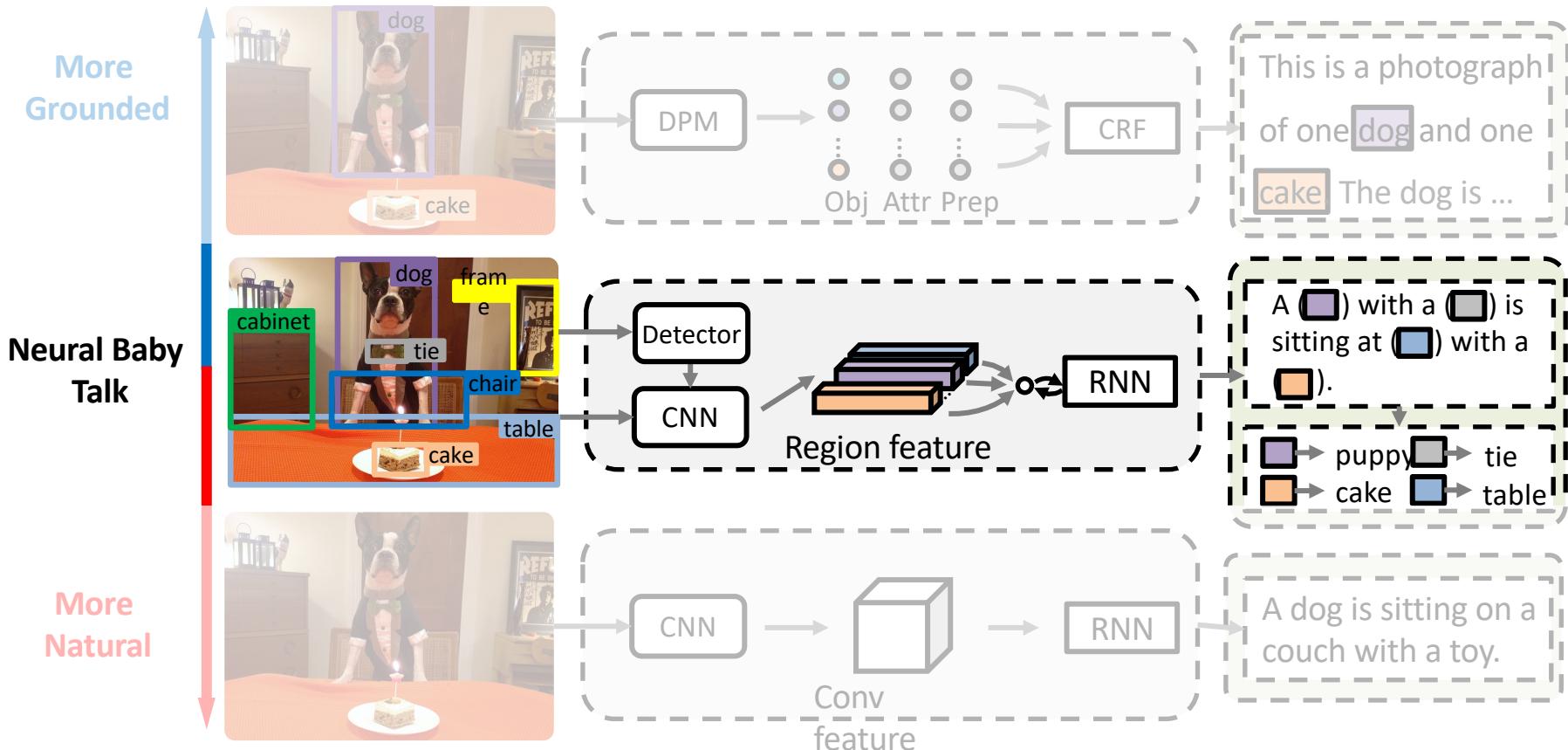


water

Neural Baby Talk



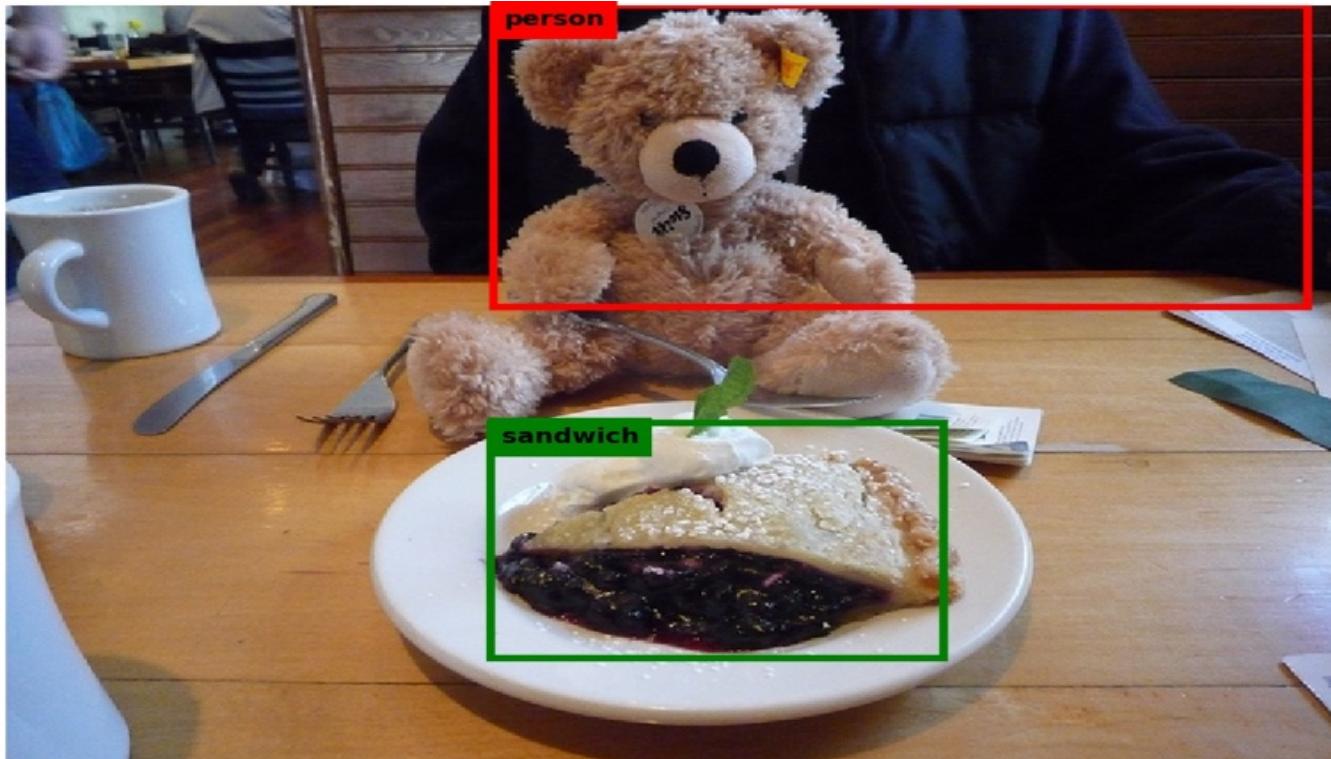
Framework



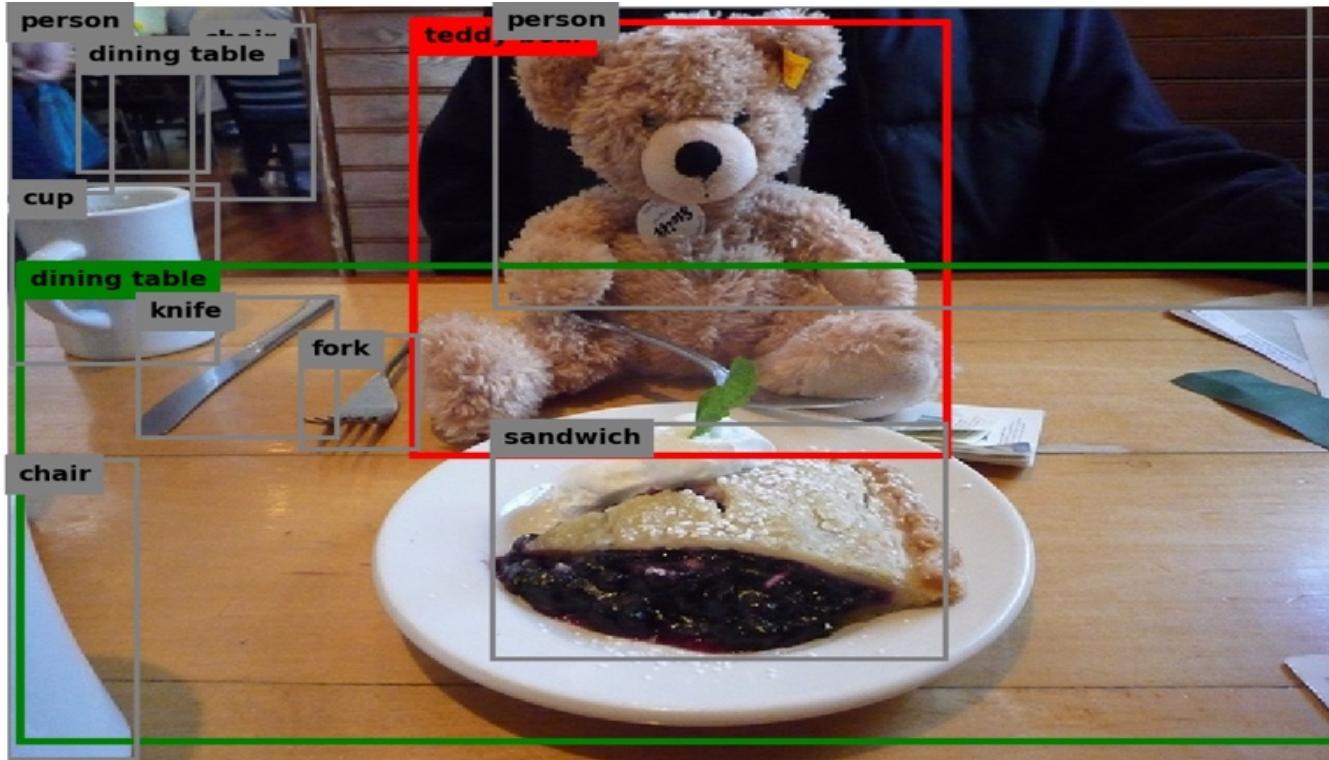
A close up of a stuffed animal on a plate.



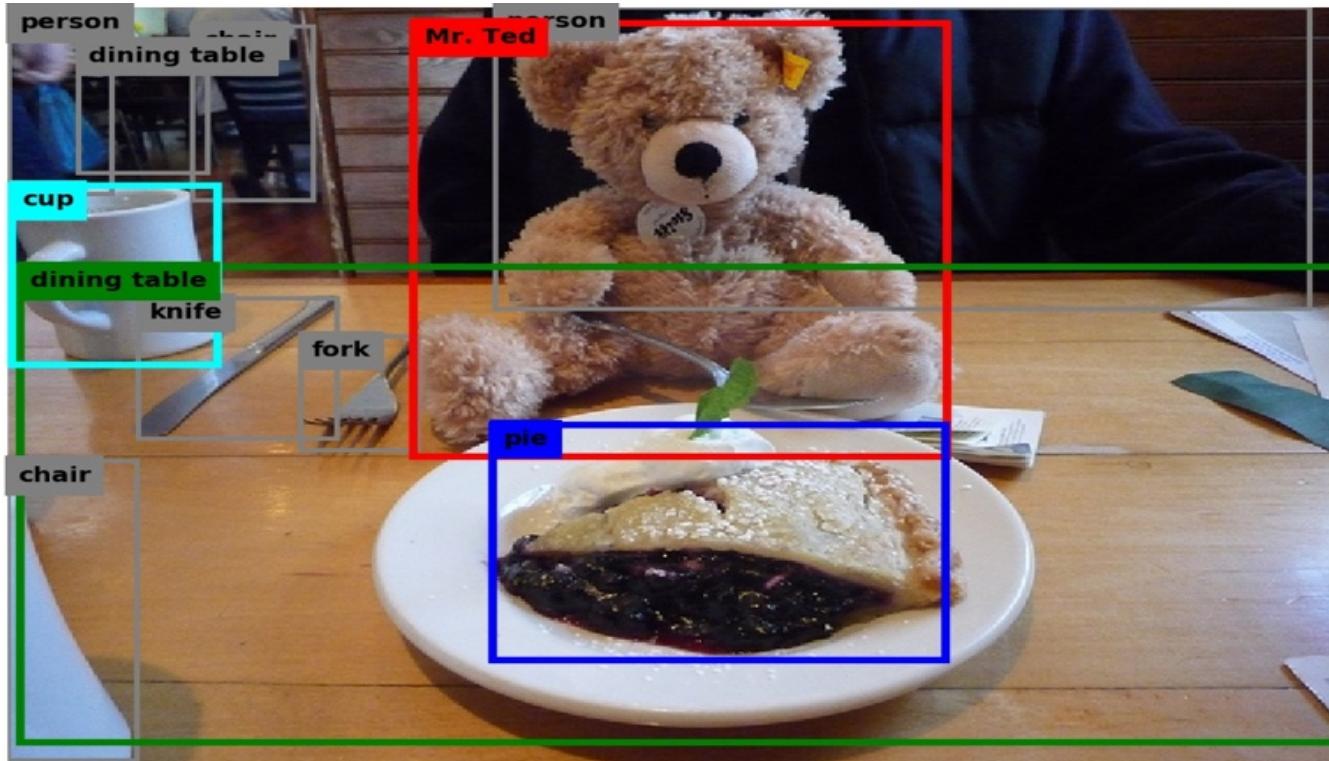
A **person** is sitting at a table with a **sandwich**.



A **teddy bear** sitting on a **table** with a **plate of food**.



A Mr. Ted sitting at a **table** with a **pie** and a **cup** of coffee.



Outline

- Datasets
 - E.g., VQA v2.0
- Evaluation metrics
 - E.g., # pairs of images correctly answered
- Novel problem spaces
 - E.g., changing priors (VQA-CP), novel object captioning, robust captioning, visual coreference resolution
- Inductive biases in models
 - E.g., G-VQA, Neural Baby Talk
- Human-in-the-loop evaluation
 - E.g., accuracy is not the final metric, performance of human-AI team is

GuessWhich Game

Total points earned: 0 Possible points for this game: 200

Hi, my name is Quincy. I am an Artificial Intelligence. You can see an image on your screen. I cannot see this image, but I can only see a one line description of this image. My task is to find this target image from a large pool of images. Please help me! I will ask you questions about the image, and your answers will help me find the right image. I will try to ask reasonable questions, but I am not perfect. I make quite a few mistakes. And my memory is a little weak. So please be patient with me, I may ask the same question more than once. I hope we can work together to find the image! Let's do this! Note: My knowledge of English is limited. Sometimes if I don't know the right word, I say UNK.



Target image description: Some boats are going around the river near the docks.

Q: 1/9 what color are boats ?

|Start typing answer here ...



GuessWhich Game

Total points earned: 0 Possible points for this game: 200



Hi, my name is Abot. I am an Artificial Intelligence. I have been assigned one of these images as the target image. I am not allowed to show you the image, but as a start, I will describe the image to you in a sentence. You can then ask me follow up questions about it. When ready, submit one of the images on the left as your best guess. I will try to describe the image and answer your questions, but I am not perfect. I make quite a few mistakes. I hope we can work together to find the image! Let's do this! Note: My knowledge of English is limited. Sometimes if I don't know the right word, I say UNK. You will win points based on how accurately you are able to guess.



Target image description: a man sitting on a couch with a laptop



Based on your understanding of the image description, pick the image that you think is the most relevant.

Click on the image that is your best guess

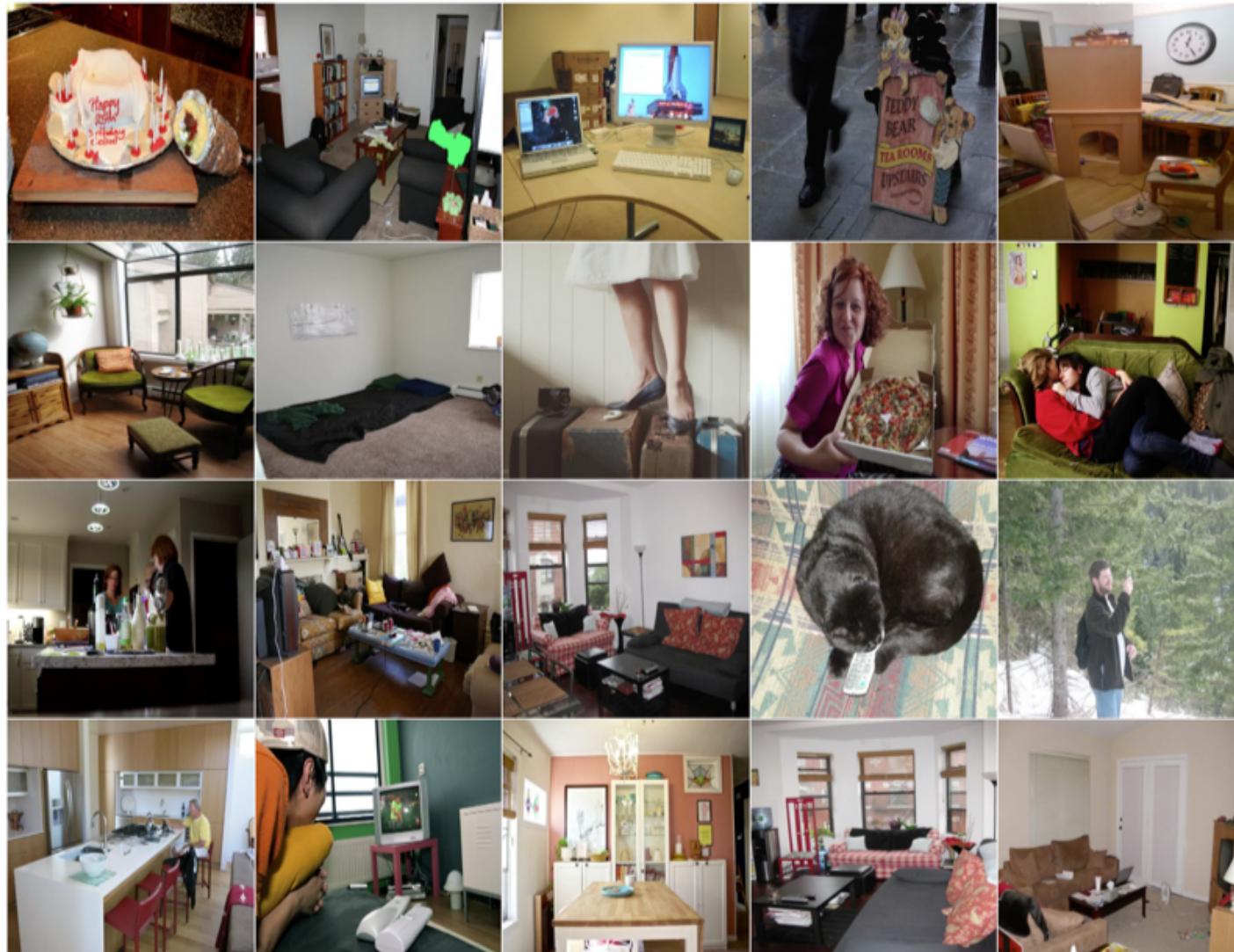


Image Guessing Game

- Compared two bots
- One did better at this game in bot-bot evaluation
- Trend did not generalize when evaluated with humans*

* caveats apply

Summary

- VQA is a rich problem space
 - Vision, language, attention, reasoning, external knowledge, HCI, ...
- Challenges:
 - Grounding and generalization
- Forcing models to look
 - Counter via datasets
 - E.g., VQA v2.0, VQA-CP
 - Counter via evaluation metrics
 - E.g., # pairs of images correctly answered
 - Counter via novel problem spaces
 - E.g., novel object captioning, robust captioning, visual coreference resolution
 - Counter via inductive biases in models
 - E.g., G-VQA, Neural Baby Talk
 - Counter via human-in-the-loop
 - E.g., accuracy is not the final metric, performance of human-AI team is

Thank you.