



The **GUESSWHAT?!** and FiLM Stories

Presented by L. Celotti, for IGLU
Master Class HLU, Paris 11 Apr. 2018



Interactive Grounded Language Understanding

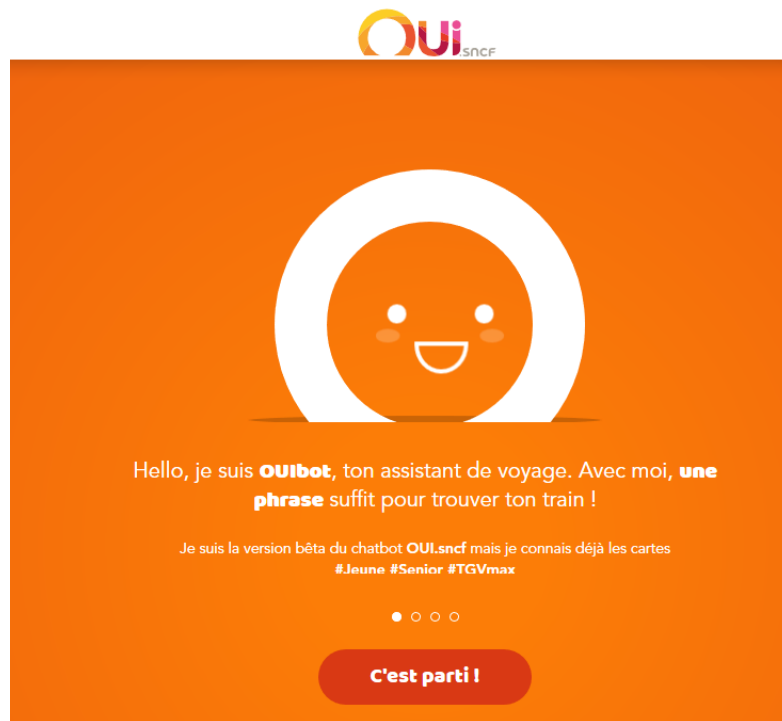


Motivation


Dialogue System



Dialogue System



OUI_{sncf}

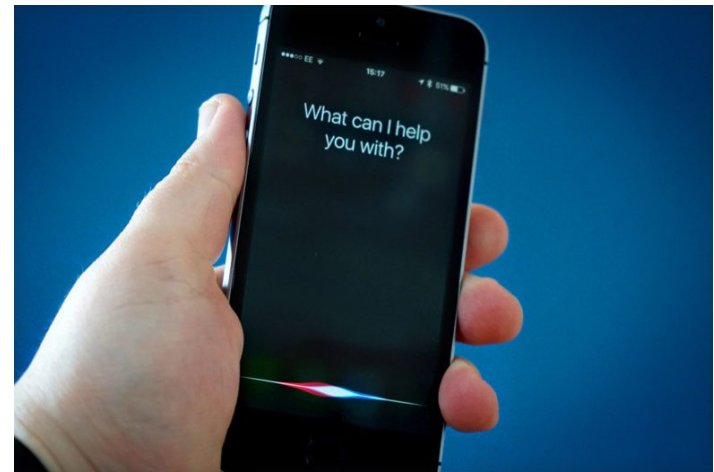


Hello, je suis **OUIbot**, ton assistant de voyage. Avec moi, **une phrase** suffit pour trouver ton train !

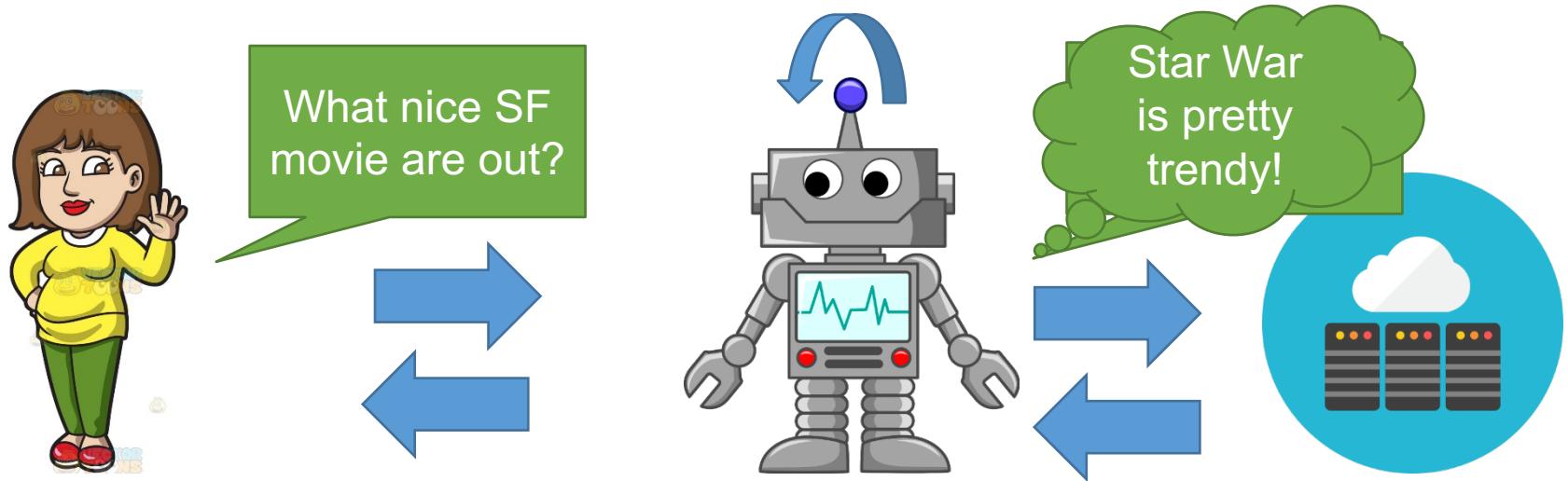
Je suis la version bêta du chatbot OUI.sncf mais je connais déjà les cartes
#Jeune #Senior #TGVmax

● ○ ○ ○ ○

C'est parti !



Dialogue System



Step 1 : Speech processing

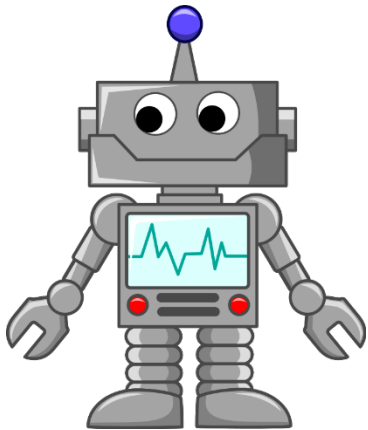
Step 2 : Language understanding

Step 3 : Query information

Step 4 : Language generation

Step 5 : Speech generation

Dialogue System



Expert system
(rule based)

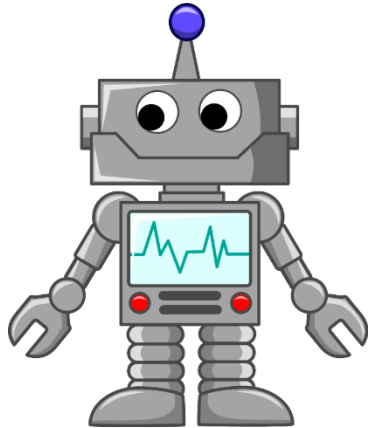
If `user.say() == "Hello"`:
-> "Say Hello"

if `user.say() == "What ${type} movie are out?"`:
-> Browse: `$type` in IMD
-> Say: "Do you know `${res}`?"

Generative model
(machine learning)

- Project question into high-dimensional space
- Generate word by word answer from this high-dimensional space

Dialogue System



Chit-chat dialogue

B: How are you today?

H: Fine! Thank you!

B: Did you hear the news? A new baby panda is born in the city zoo!

H: ohhhh, so lovely!

Goal-Oriented dialogue

H: I want to book a plane to London.

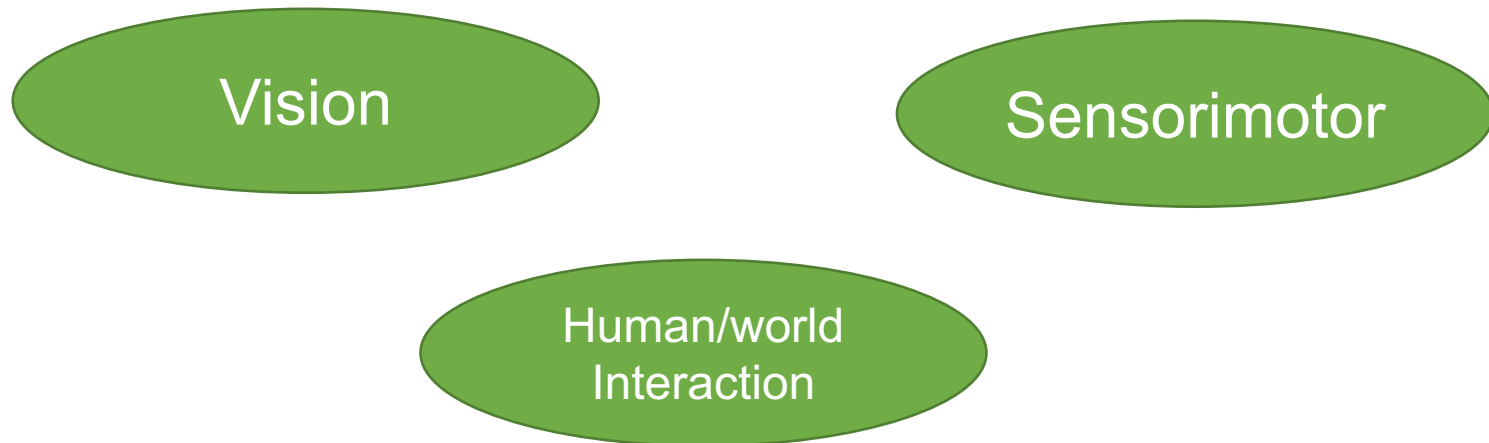
B: Sure, when do you want to leave?

H: Tomorrow, in the morning

B: There is plane at 9am etc...

Grounding Language

The grounding problem is related to the problem of how words (symbols) get their meanings. [1]



GuessWhat?!

The Game

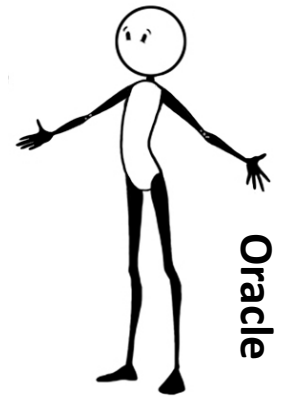
GuessWhat?! Game



Questionner



Oracle



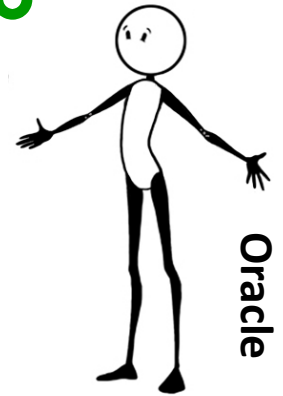
GuessWhat?! Game



Questionner



Oracle



GuessWhat?! Game

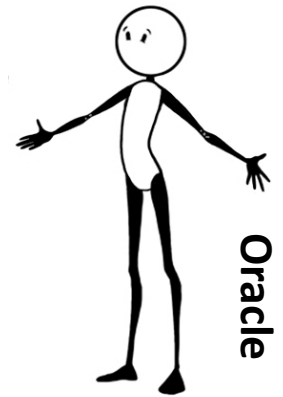


Questionner



Is it a vase ?

Oracle



GuessWhat?! Game

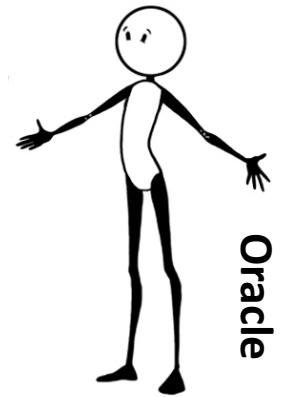


Questioner



Is it a vase ?

Yes



Oracle

GuessWhat?! Game

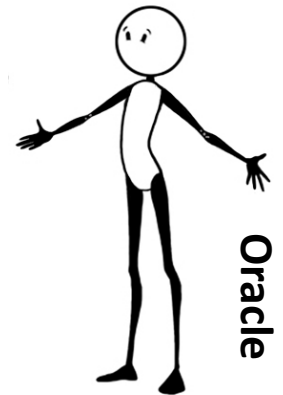


Questioner



Is it a vase ?
Is it in the front row?

Yes



Oracle

GuessWhat?! Game



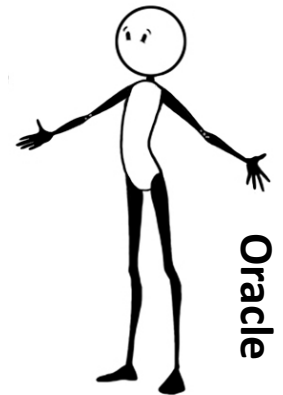
Questioner



Is it a vase ?
Is it in the front row?

Yes
Yes

Oracle



GuessWhat?! Game



Questioner



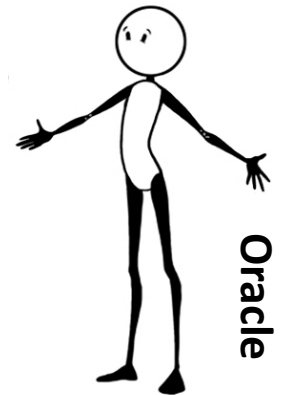
Is it a vase ?

Is it in the front row?

Does it have some red on it?

Yes

Yes



Oracle

GuessWhat?! Game



Questioner



Is it a vase ?

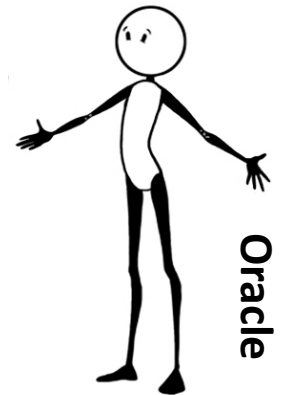
Is it in the front row?

Does it have some red on it?

Yes

Yes

No



Oracle

GuessWhat?! Game



Questioner



Is it a vase ?

Is it in the front row?

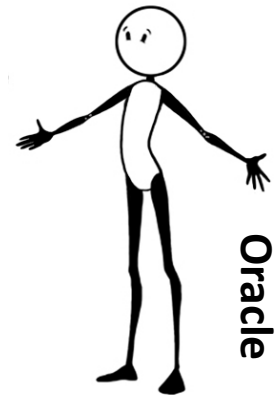
Does it have some red on it?

Is it the second vase from the right?

Yes

Yes

No



Oracle

GuessWhat?! Game



Questioner



Is it a vase ?

Is it in the front row?

Does it have some red on it?

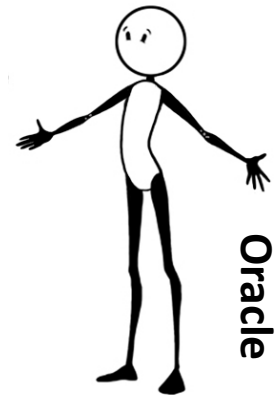
Is it the second vase from the right?

Yes

Yes

No

Yes



Oracle

GuessWhat?! Game



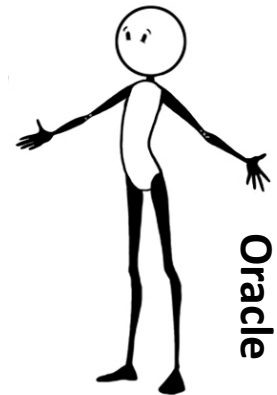
I found it!

Questioner



- Is it a vase ?
- Is it in the front row?
- Does it have some red on it?
- Is it the second vase from the right?

- Yes
- Yes
- No
- Yes



Oracle

GuessWhat?! Game



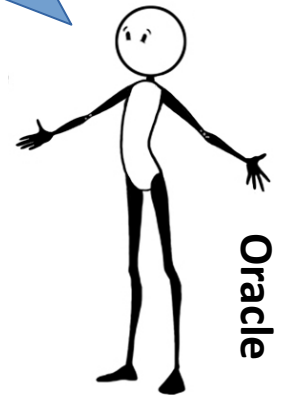
Correct!

Questioner



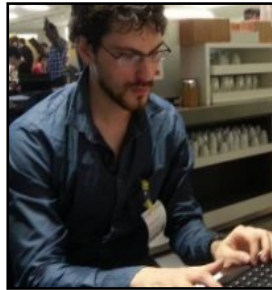
- Is it a vase ?
- Is it in the front row?
- Does it have some red on it?
- Is it the second vase from the right?

- Yes
- Yes
- No
- Yes



Oracle

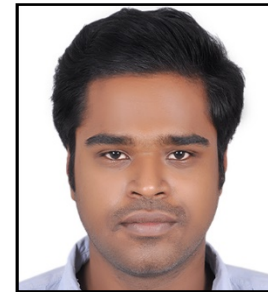
The Team



Florian Strub
Univ. of Lille, Inria



Harm de Vries
Univ. of Montreal



Sarath Chandar
Univ. of Montreal



Bilal Piot
Univ. of Lille, Inria



Hugo Larochelle
Google Brain



Jérémie Mary
Criteo



Olivier Pietquin
Univ. of Lille, Inria



Aaron Courville
Univ. of Montreal

GuessWhat?! Game

Wait! Is it really a difficult task?

- No rule-based
 - High dimensional input/output spaces
- Contextual
 - Each game has a new image
- End-2-end:
 - No image processing: input = raw pixel
 - No language processing: input/output = words
- Multi-objective
 - Asking meaningful questions
 - Asking coherent sequence of questions



GuessWhat?! Dataset

[Home](#)[People](#)[Download](#)[Explore](#)[Play the game](#)

Game instructions

This is a two-player game, yourself and a partner. You will be (randomly) assigned to play one of the two roles:



Questioner

Find the object

- You will be shown an image of a scene with multiple objects.
- One of the objects will be assigned as the target (but not visible to you).
- Your job is to locate that object by asking yes or no questions.
- You click on the GuessWhat! button once you are certain which object it is.
- All object segmentations are then shown in the image, and you click on the correct object.



Oracle

Answer the questions

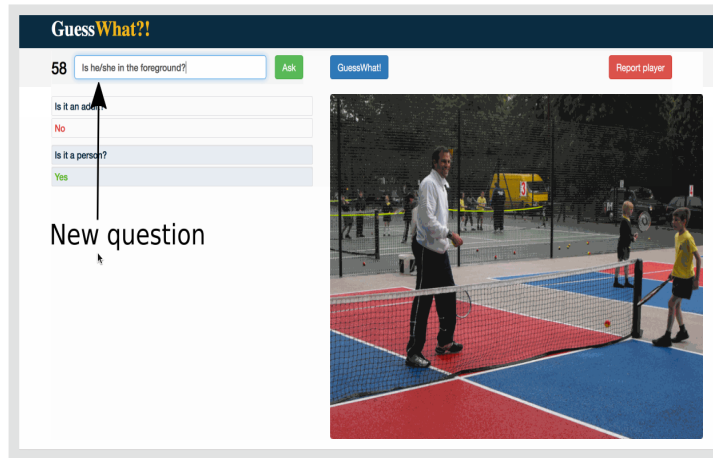
- You will be shown an image of a scene with multiple objects.
- One of the objects will be assigned as the target.
- Your partner will ask yes/no questions to locate this object.
- Your job is to answer their questions correctly.



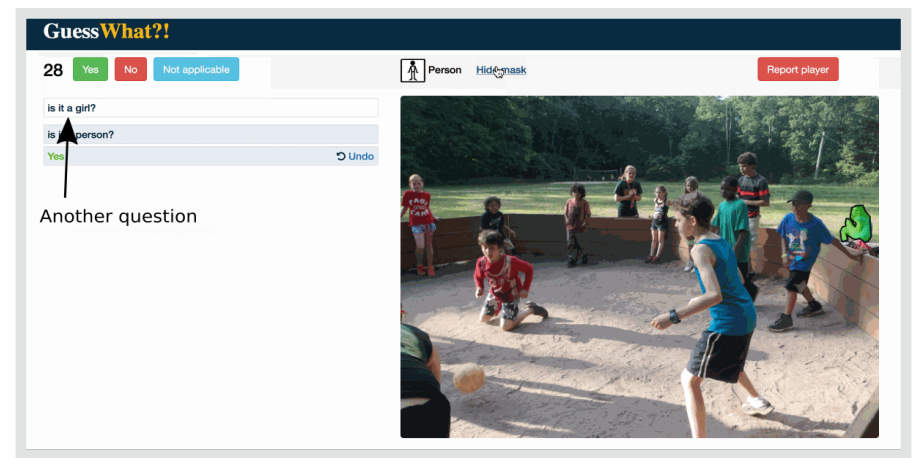
Ready to play?

[Start a game »](#)[Play with AI »](#)

GuessWhat?! Dataset



Questioner

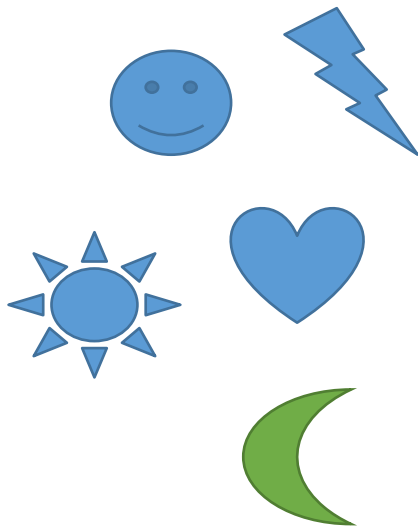


Oracle

- More than 10k players
- Top 10 players completed more than 2500 games
- Highest game time was 70+ hour

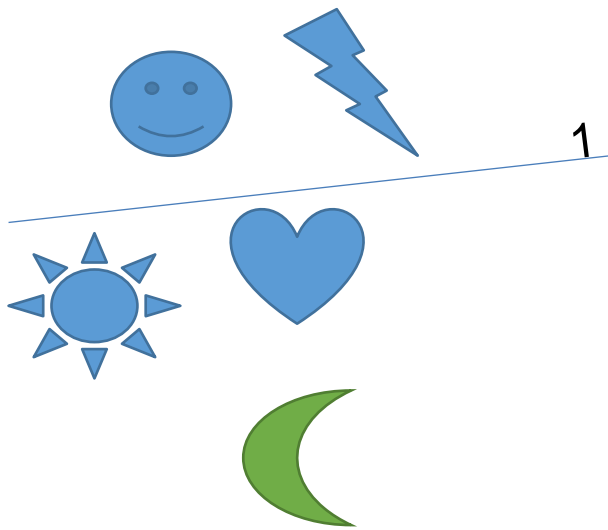


GuessWhat?! Dataset



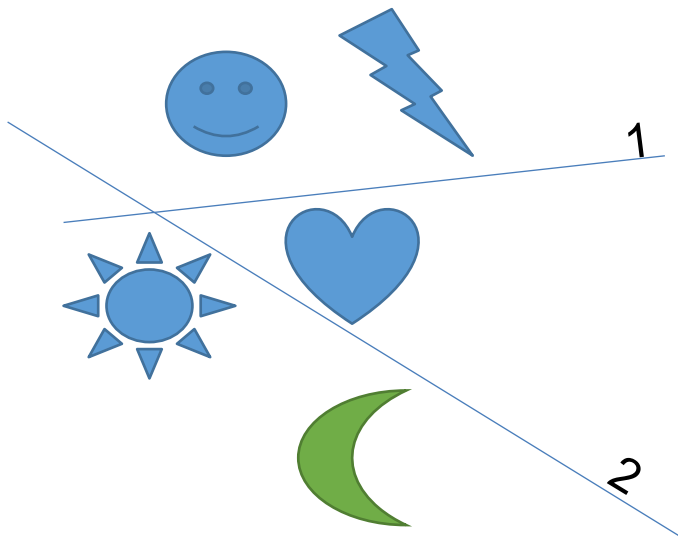
Optimal search in the object space is a binary search i.e. removing half of the objects at each step ($\log(n)$ complexity).

GuessWhat?! Dataset



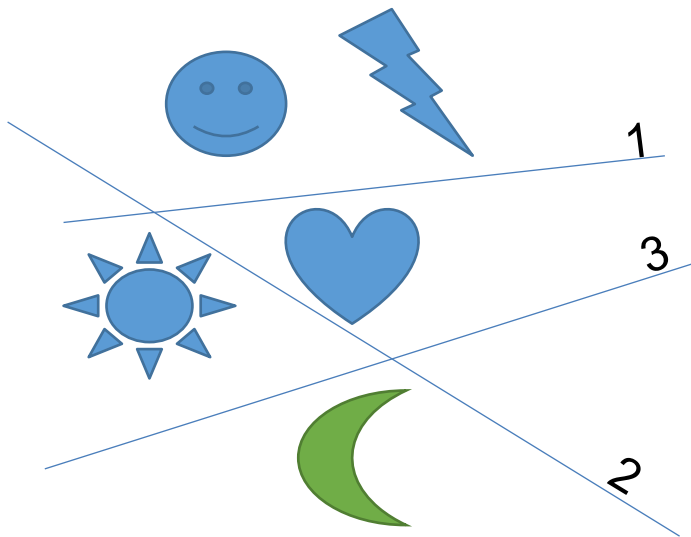
Optimal search in the object space is a binary search i.e. removing half of the objects at each step ($\log(n)$ complexity).

GuessWhat?! Dataset



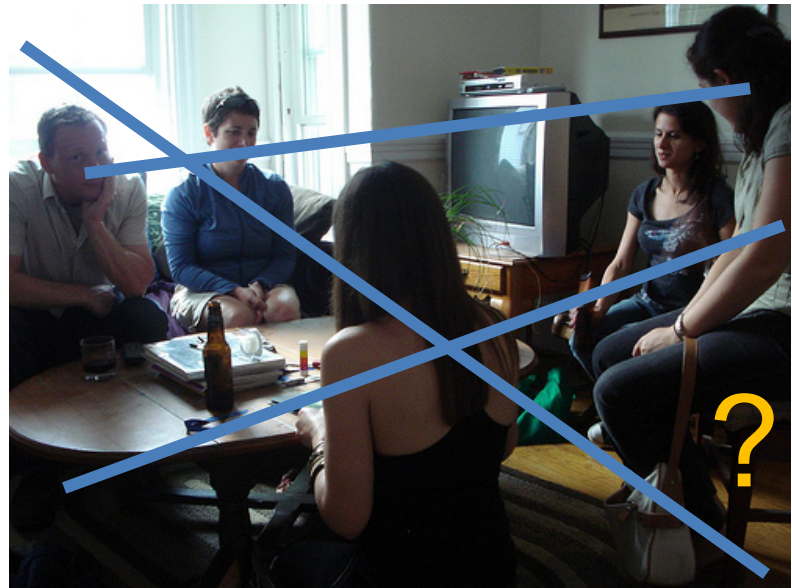
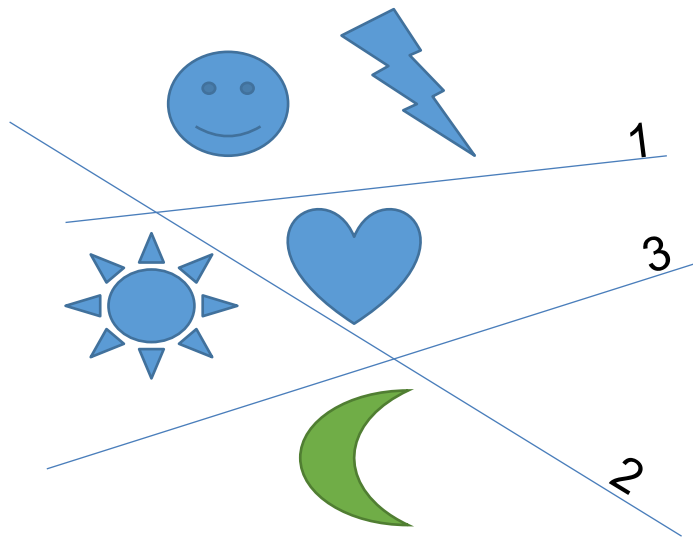
Optimal search in the object space is a binary search i.e. removing half of the objects at each step ($\log(n)$ complexity).

GuessWhat?! Dataset



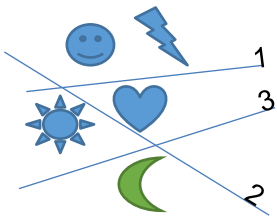
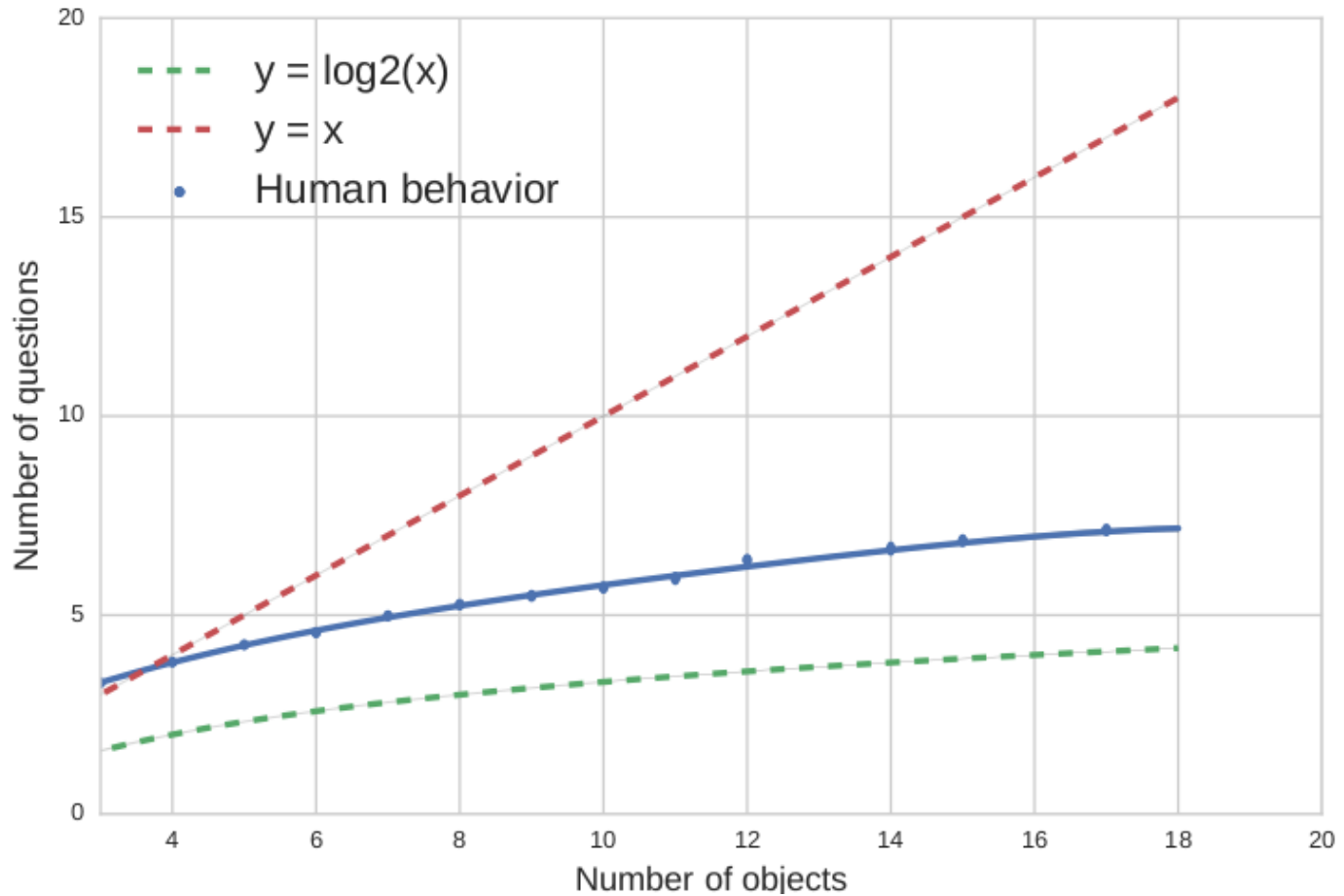
Optimal search in the object space is a binary search i.e. removing half of the objects at each step ($\log(n)$ complexity).

GuessWhat?! Dataset



Optimal search in the object space is a binary search i.e. removing half of the objects at each step ($\log(n)$ complexity).

GuessWhat?! Dataset



The Models

Models

Repeat until <stop_dialogue>



question



yes/no answer



Questioner

Oracle



Models

Repeat until <stop_dialogue>



question



yes/no answer



Questioner

Oracle

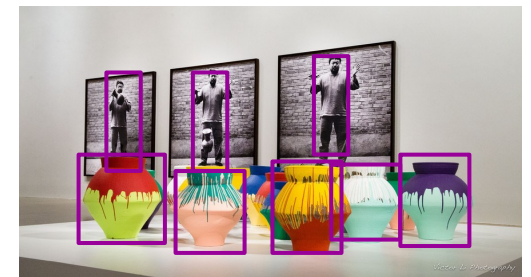
dialogue



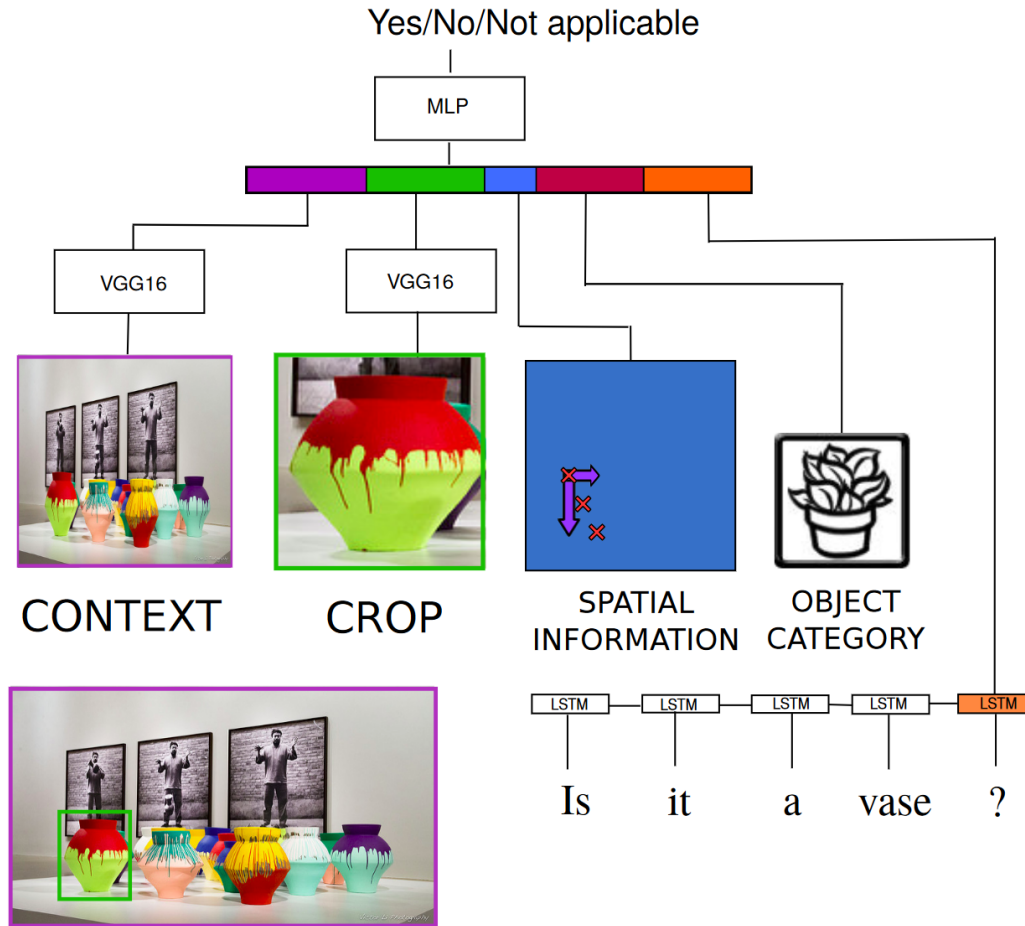
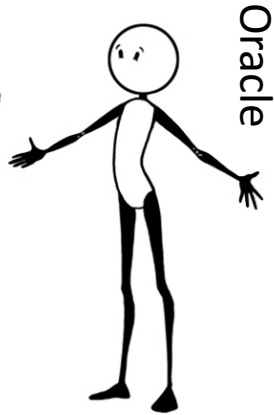
Find object?



Guesser

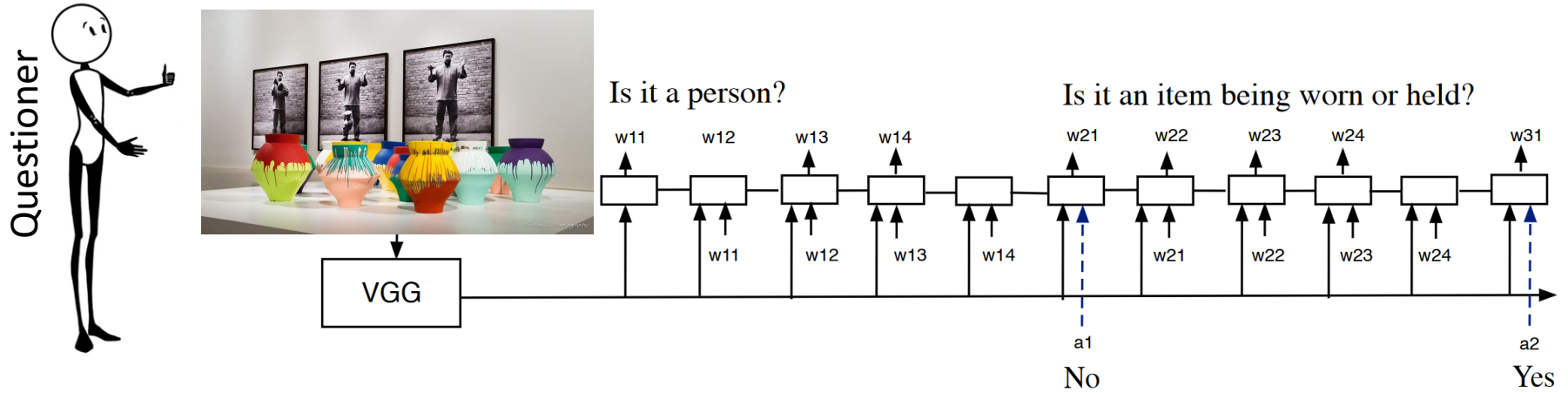


Models




79.5%
accuracy

Models



47.1%
accuracy

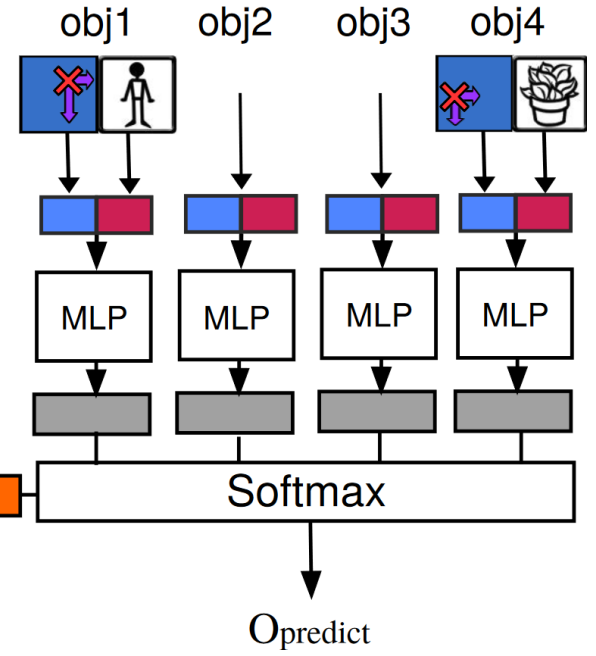
Models

Guesser 

Is it a vase? Yes
Is it partially visible? No
Is it in the left corner? No
Is it the turquoise and purple one? Yes





LSTM / HRED encoder



63.8%
accuracy

Where SL fails and RL not

Image	Beam Search	REINFORCE
	<p>Is it a person ? no Is it a ball ? no Is it a ball ? no Is it a ball ? no Is it a ball ? no</p> <p>Failure (blue bat)</p>	<p>Is it a person ? no Is a glove ? no Is an umbrella ? no Is in the middle ? no On a person? no is it on on far right? yes</p> <p>Success (red chair)</p>
	<p>Is it a person ? yes Is it the one in front ? yes Is it the one on the left ? no Is it the one in the middle with the red umbrella ? yes Is it the one to the right of the girl in ? no</p> <p>Failure (umbrella)</p>	<p>Is it a person ? yes Is it girl in white ? yes</p> <p>Success (girl)</p>

Where RL fails too

- Tend to lose grammatical structure: Is it a cat? → cat ?
- Environment use limited modalities → limits the number of strategy
- Stopping far from optimal !



Is it a suitcase? Yes
Is it in the left side? Yes



is it a person ? yes
is it in left ? no
is it in right ? no
is it in middle ? yes
is it in front ? yes
is it in middle ? yes

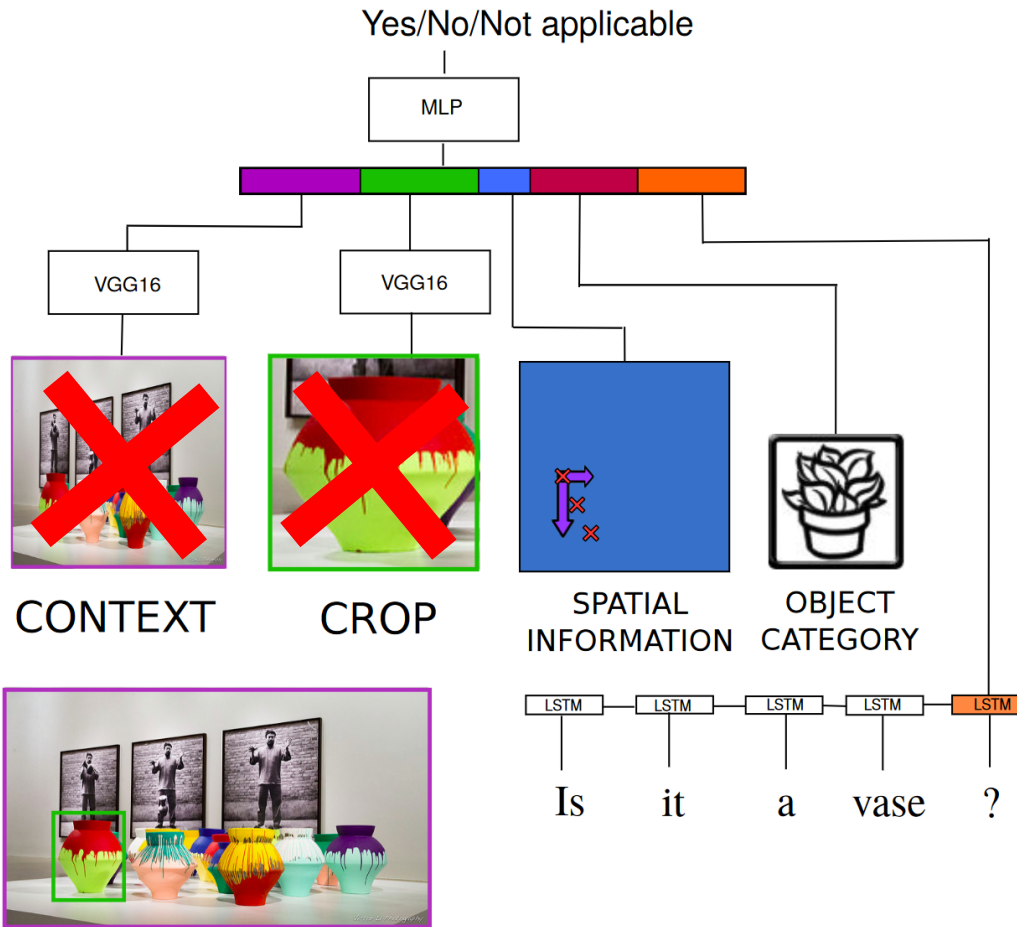
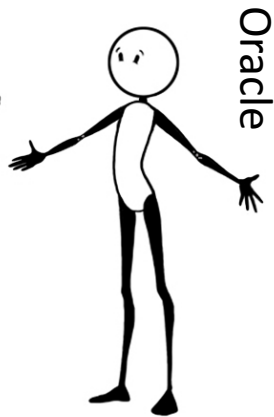
Language Bias

What is Language Bias?

- How many zebra are there in the image? 2 – 3
- Is the man wearing glasses? Yes
- What color is the banana? Yellow

Most of the SOTA models learn (with more or less success) language bias!!

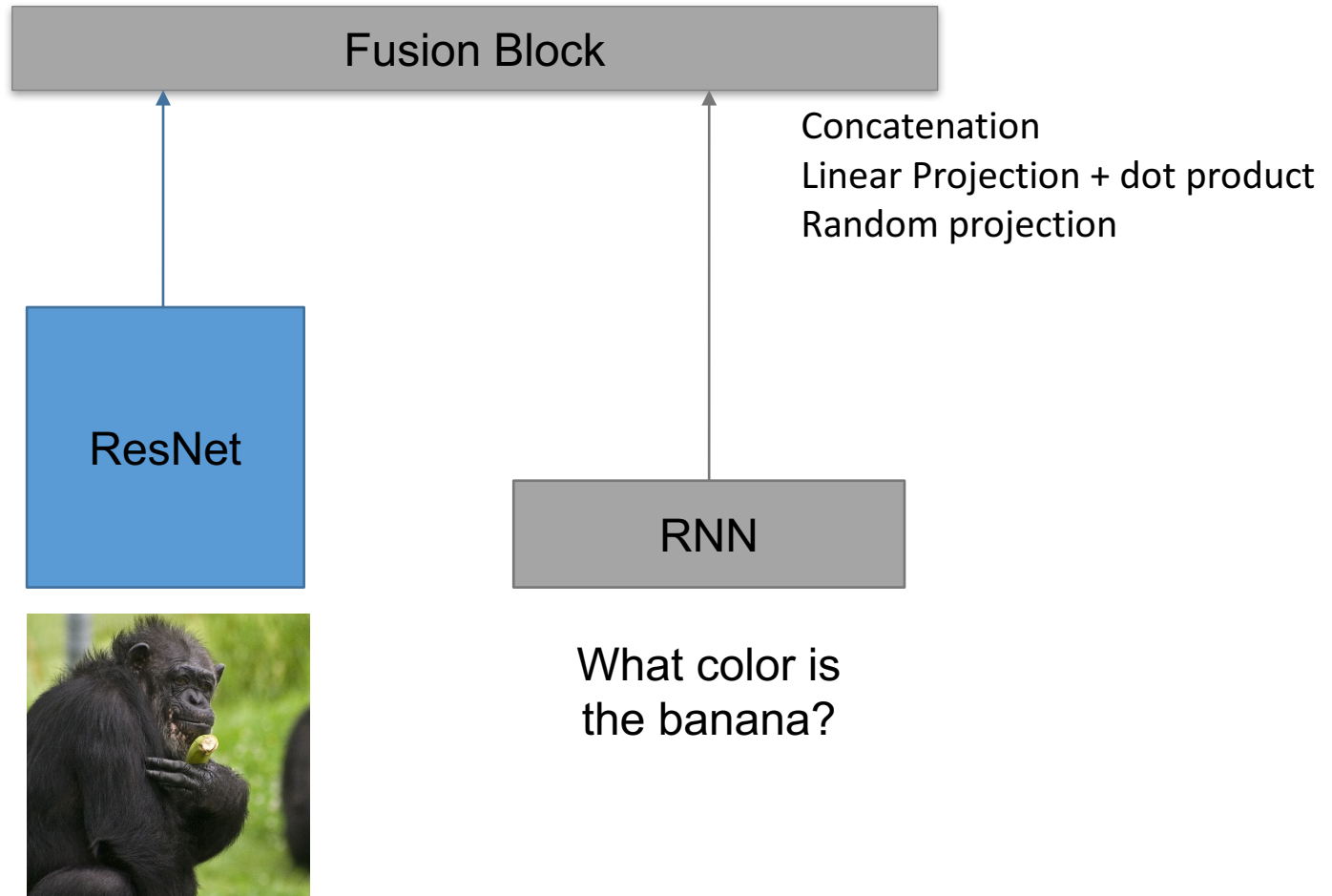
Visual Input worsened results



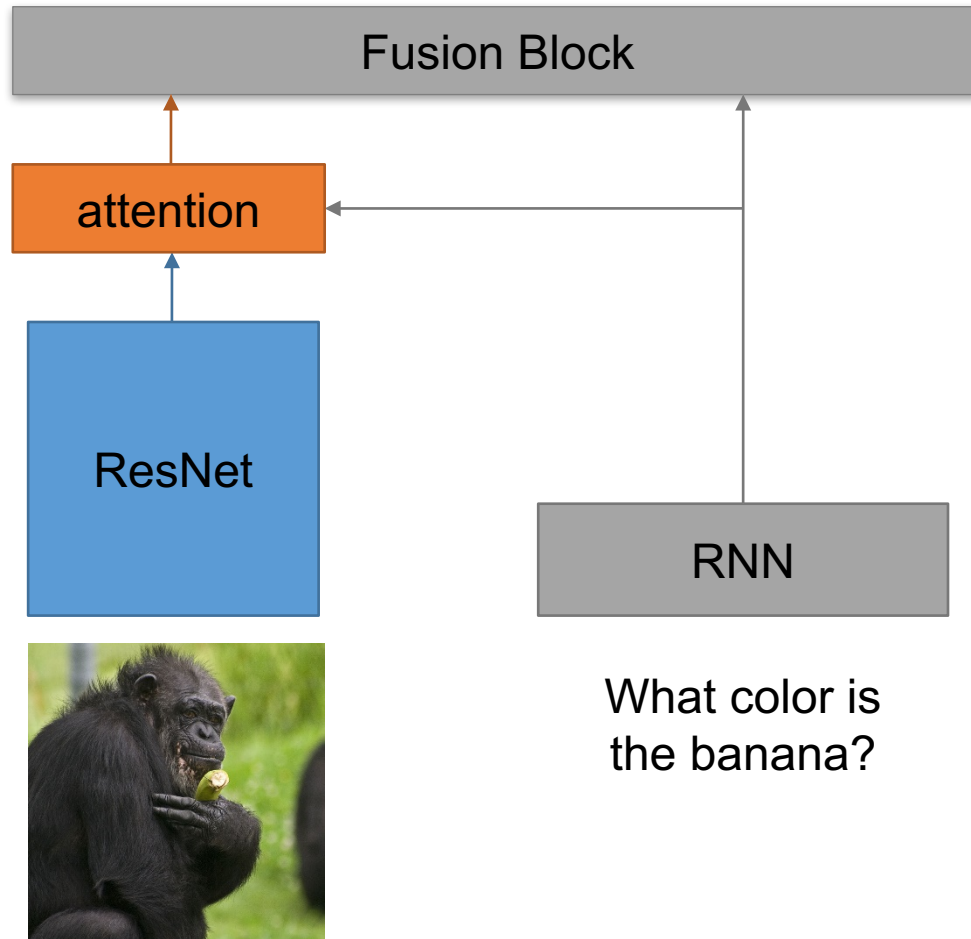
78.9%
accuracy

79.5%
accuracy

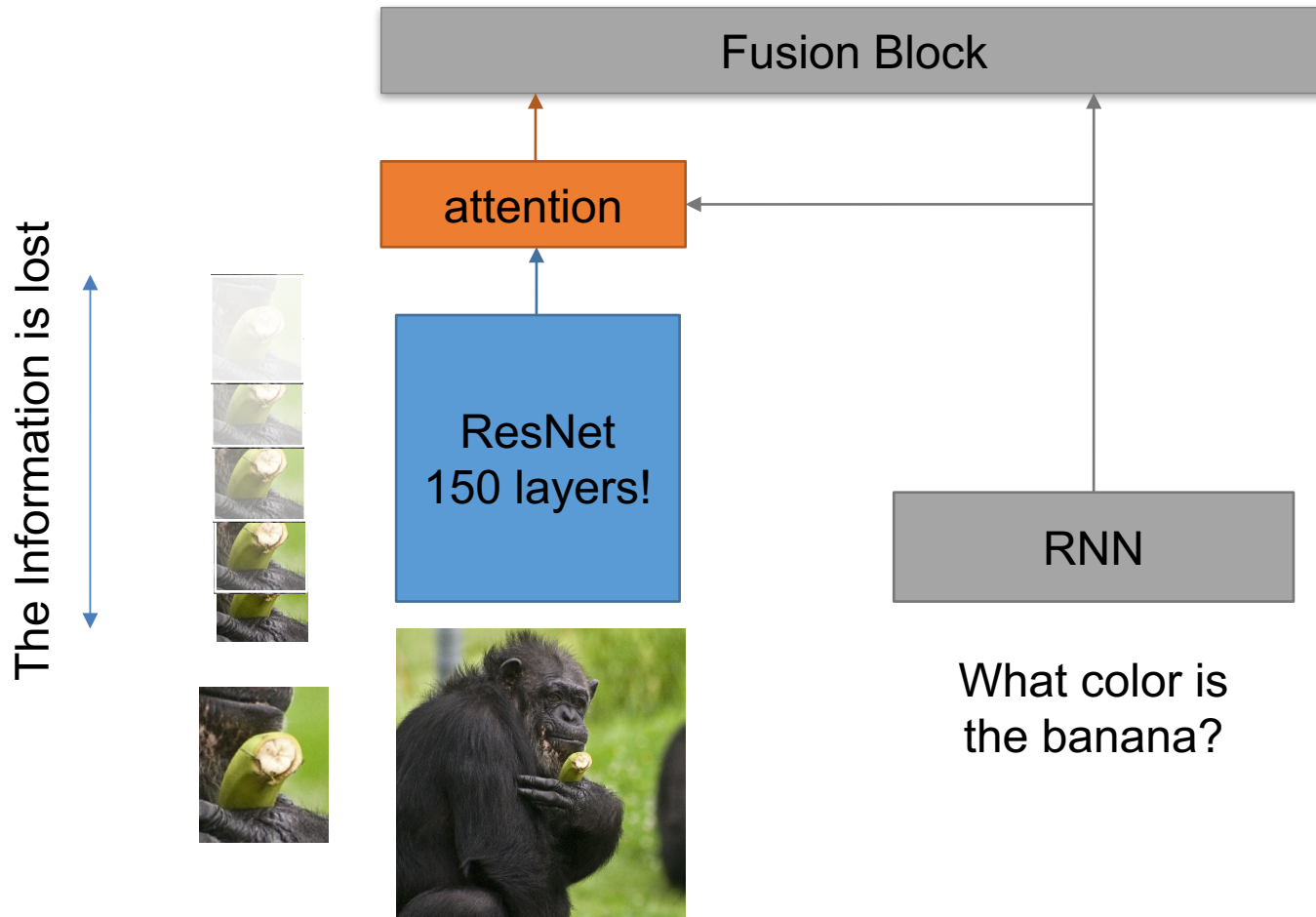
Concatenation wasn't enough



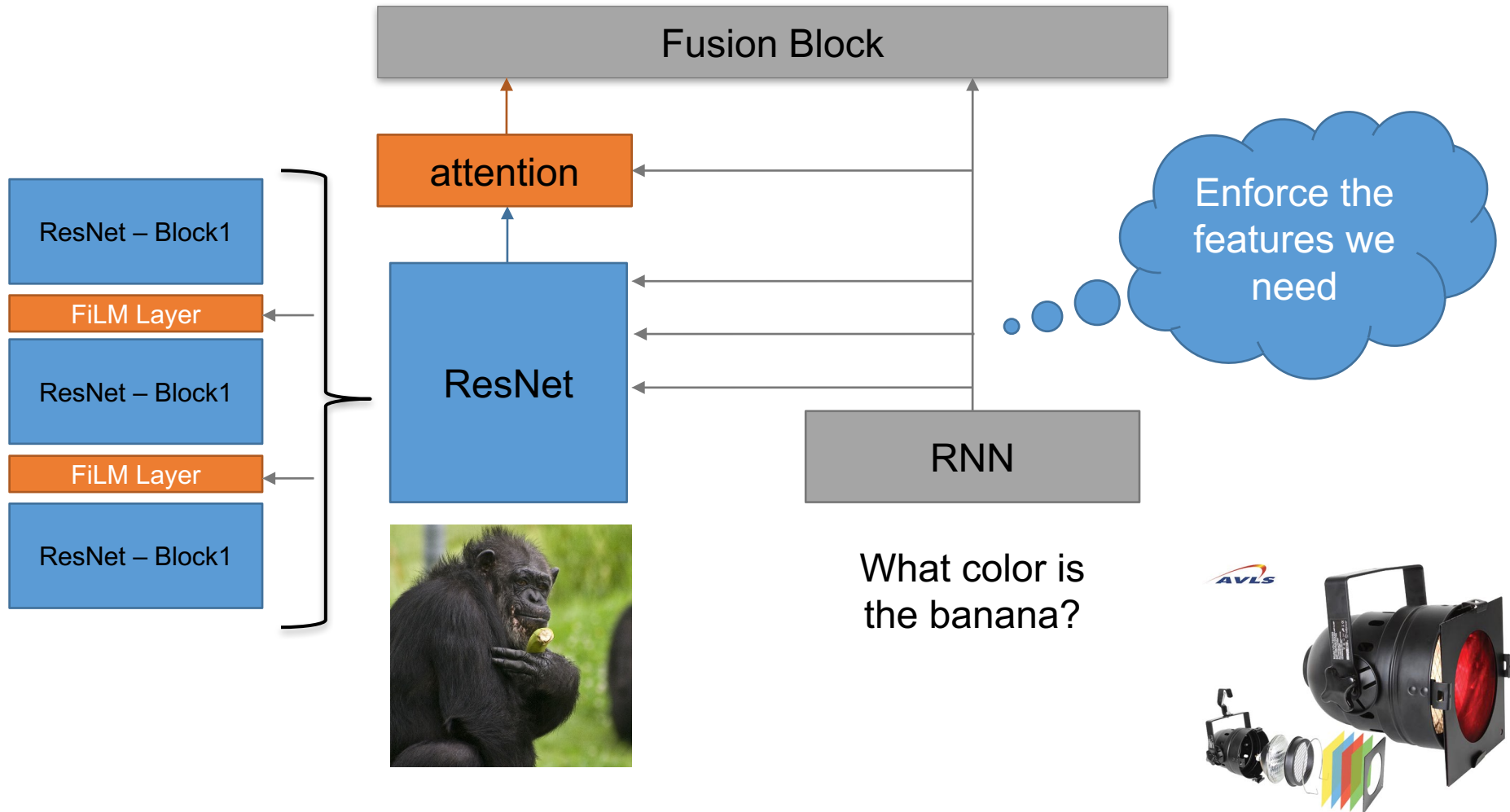
Attention isn't enough



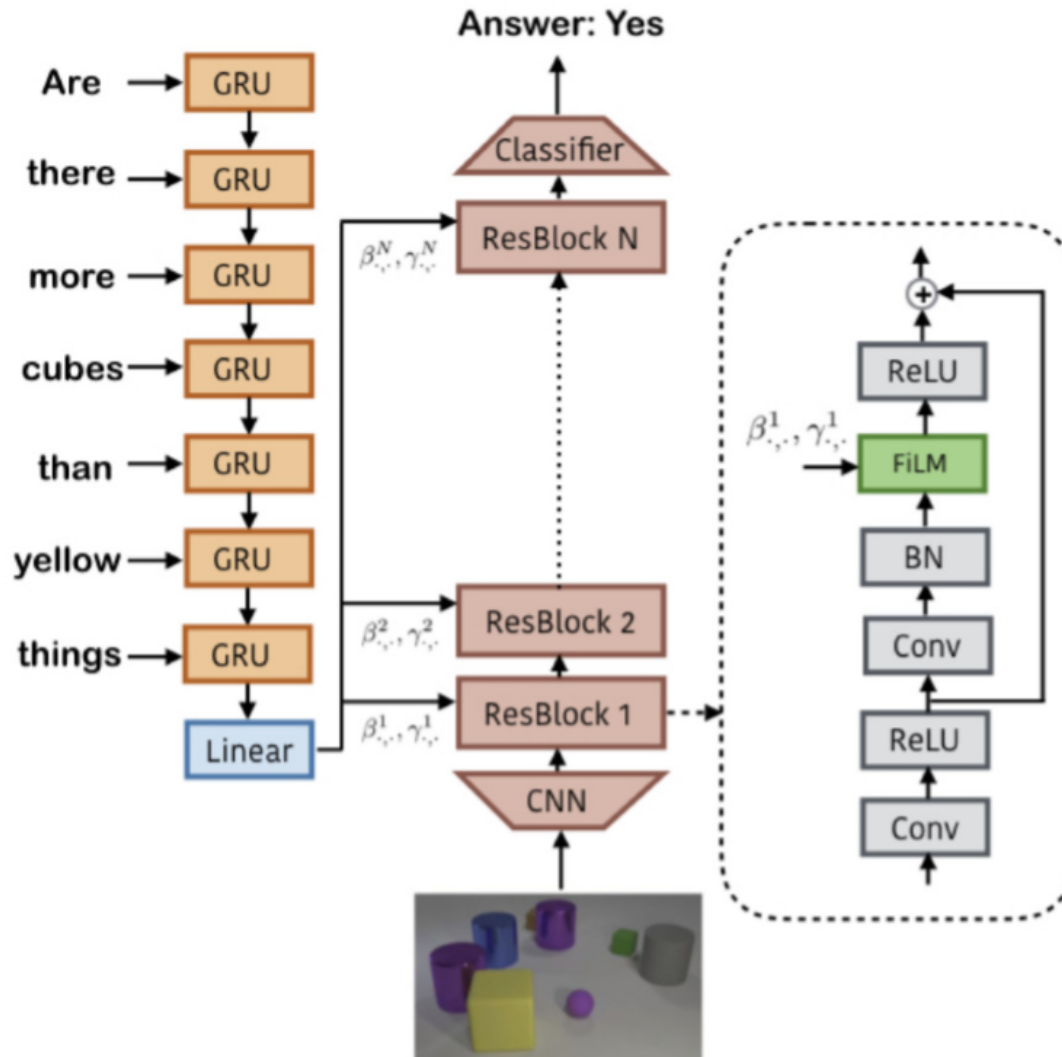
Information Lost



Distributed Attention and FiLM

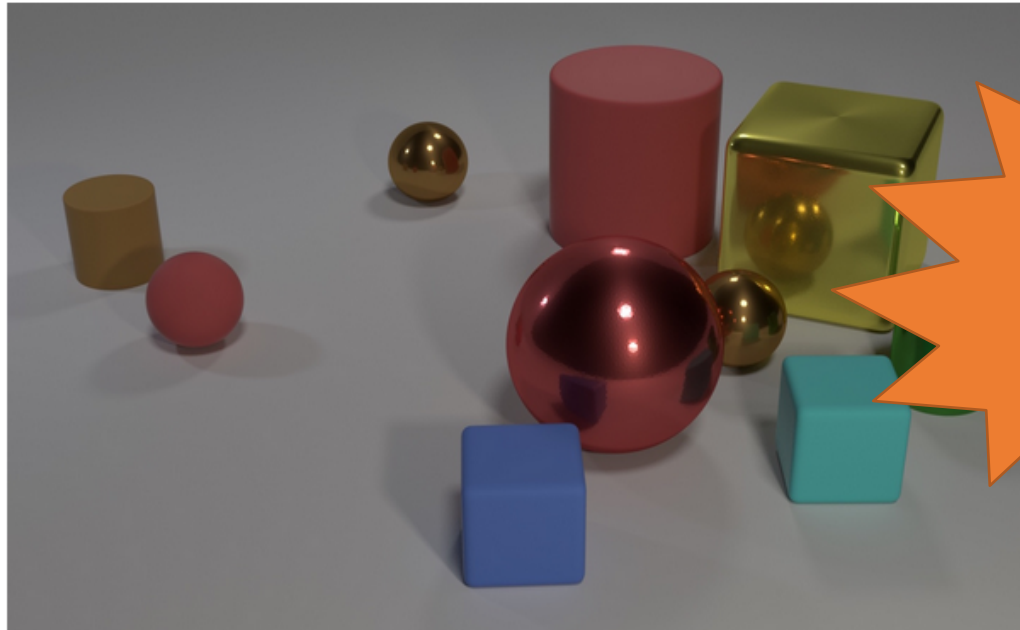


Distributed Attention and FiLM



Unbiased data: CLEVR

Questions in CLEVR test various aspects of visual reasoning including **attribute identification**, **counting**, **comparison**, **spatial relationships**, and **logical operations**.



Guaranteed
without
language
bias

Q: Are there an **equal number** of **large things** and **metal spheres**?

Q: **What size** is the **cylinder that is left of** the **brown metal** thing **that is left of** the **big sphere**?

Q: There is a **sphere** with the **same size as** the **metal cube**; is it **made of the same material as** the **small red sphere**?

Q: **How many** objects are **either small cylinders** or **red** things?

FiLM on CLEVR

Model	Overall
Human (Johnson et al.)	92.6
CNN+LSTM (Johnson et al.)	52.3
Stacked Attention (Santoro et al.)	76.6
End-to-End Module Networks* (Hu et al.)	83.7
Prog. Generator+Execution Engine* (Johnson et al.)	96.9
Relation Networks (Santoro et al.)	95.5
FiLM (Ours)	97.7

* = Uses additional program-level supervision.

Conclusions

Conclusions

We explained how

- **Guesswhat?!** proposes an interesting framework to study purposeful language interaction constrained by a visual scene
- To stay alert against **language bias**
- Benefit from placing **attention mechanisms all along** the convolutional pipeline
 - It will enable us to modulate from low level to high level features
- **FiLM** as a new performant attention mechanism

Thank you for Listening



guesswhat.ai



Supplementary Material

Language Grounding

How to train a computer to acquire *natural* language?



Language Grounding

How to train a computer to acquire *natural* language?



Language Grounding

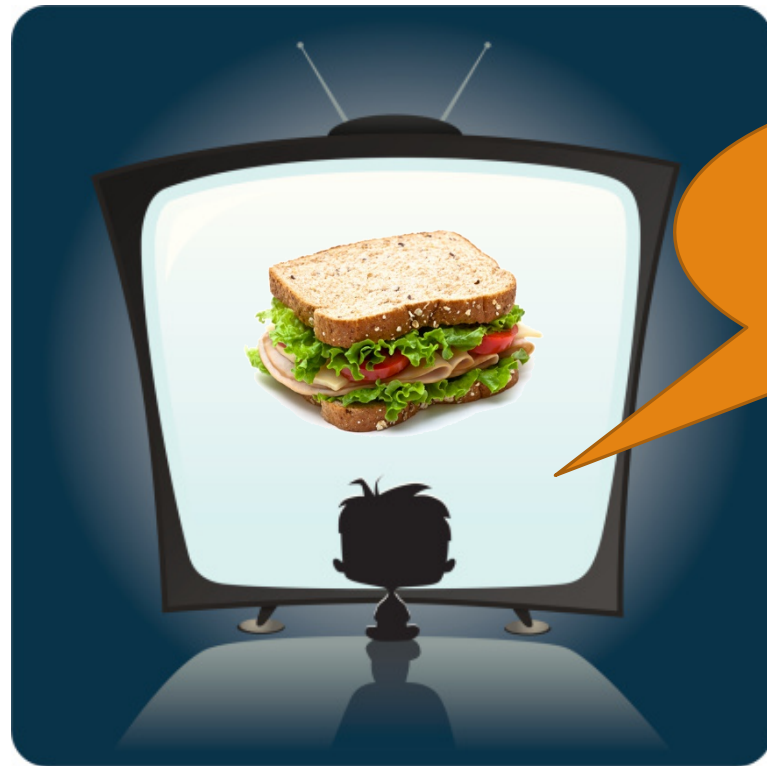
How to train a computer to acquire *natural* language?



Tramezino
サンドイッチ

Language Grounding

How to train a computer to acquire *natural* language?



Tramezino
サンドイッチ

Language Grounding

How to train a computer to acquire *natural* language?



Language Grounding

How to train a computer to acquire *natural* language?



Starting point

Issue : How to verify this hypothesis!

Game features:

- Dialogue (for interaction)
- Visually grounded
- Collaborative
- Goal-oriented with a clear reward

End-to-end Optimization of Goal-driven and Visually Grounded Dialogue Systems

GuessWhat?! Dataset



#64374

is it an animal? **Yes**

one of the two in the bottom right corner? **Yes**

the one most to the right? **No**

the one to the left of it? **Yes**

Success



Dataset
It's rich

- **155,280** played games
- **821,889** questions+answers
- **66,537** images
- **134,073** objects

[Download](#) the dataset.



#113037

is it a person? **Yes**

are they sitting in the front row? **No**

are they in the next row? **No**

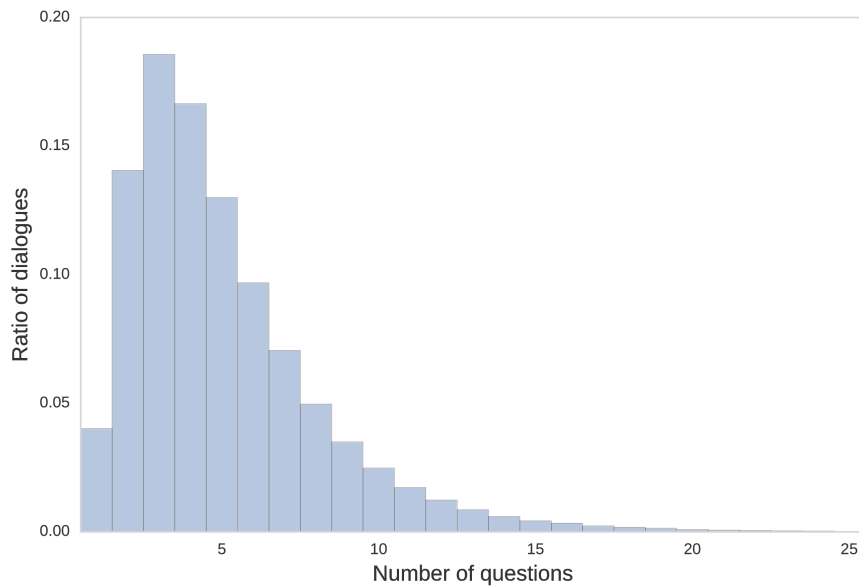
are they in the back row? **Yes**

are they on the left? **Yes**

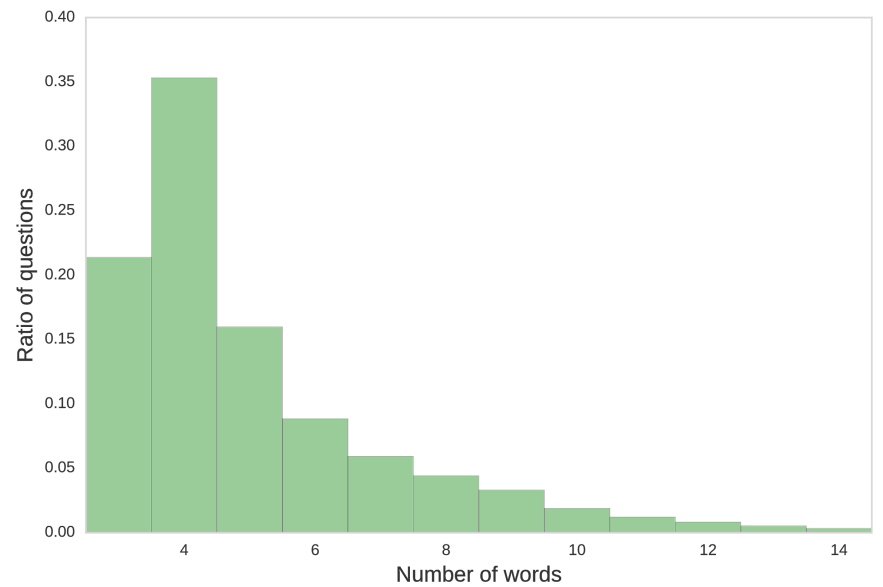
is it the guy with the pink shirt? **Yes**

Success

GuessWhat?! Dataset

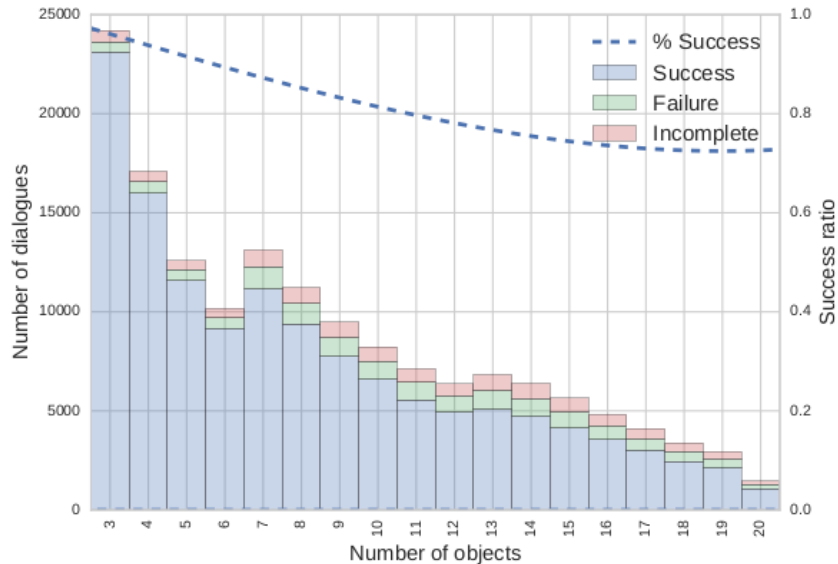


Average: 5+ questions

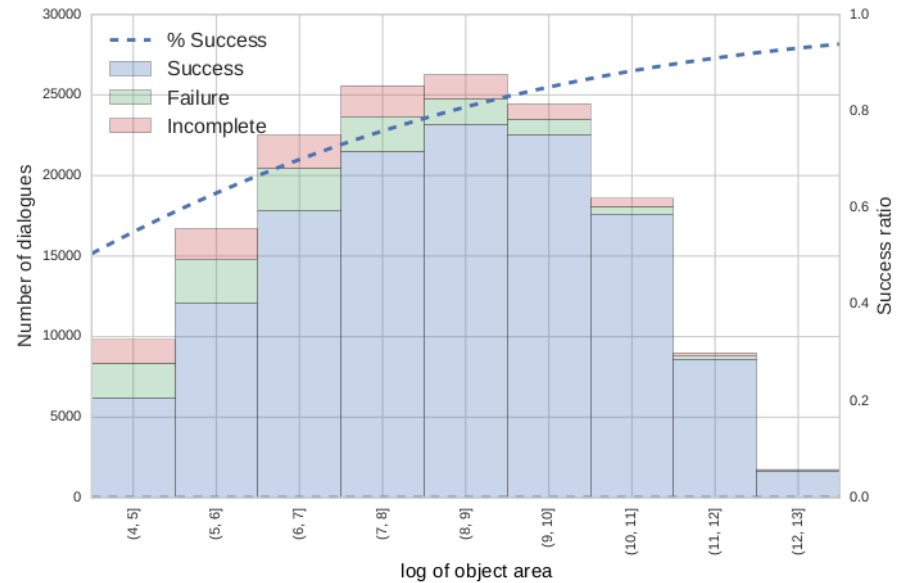


Average: ~5 words

GuessWhat?! Dataset



The more object there are, the lower is the success ratio



The bigger the object is, the higher is the success ratio

Models

Grounding the language requires to interact with the environment

Limitation of supervised learning:

- Action/state space of dialogues is large: Hard to generalize
- Imitation miss the planning aspect of dialogue
- Does supervised learning really manage to ground vision/language?



#113037

is it a person? Yes

are they sitting in the front row? No

are they in the next row? No

are they in the back row? Yes

are they on the left? Yes

is it the guy with the pink shirt? Yes

Success

What if the answer would have been no?

Is it the best question?

Can we phrase it differently?

Models

Repeat until <stop_dialogue>



question



yes/no answer



Questioner

Oracle

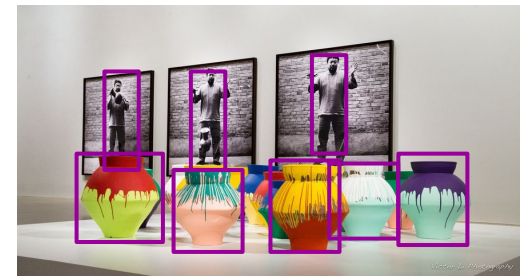
dialogue



Find object?



Guesseur



Models

Repeat until <stop_dialogue>



question



yes/no answer



Questioner

Oracle

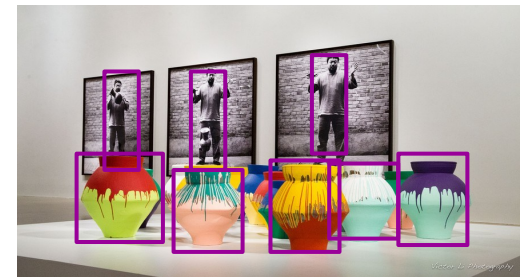
dialogue



Reward



Guesseur



Models

Repeat until <stop_dialogue>



question



yes/no answer



Questioner

Oracle

dialogue

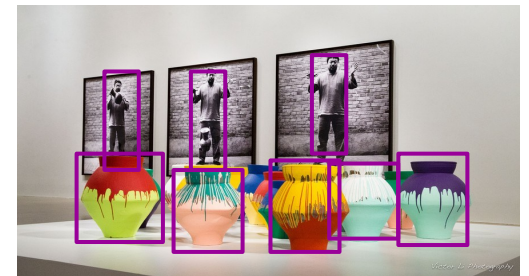


$\nabla J(\theta_h)$

Reward

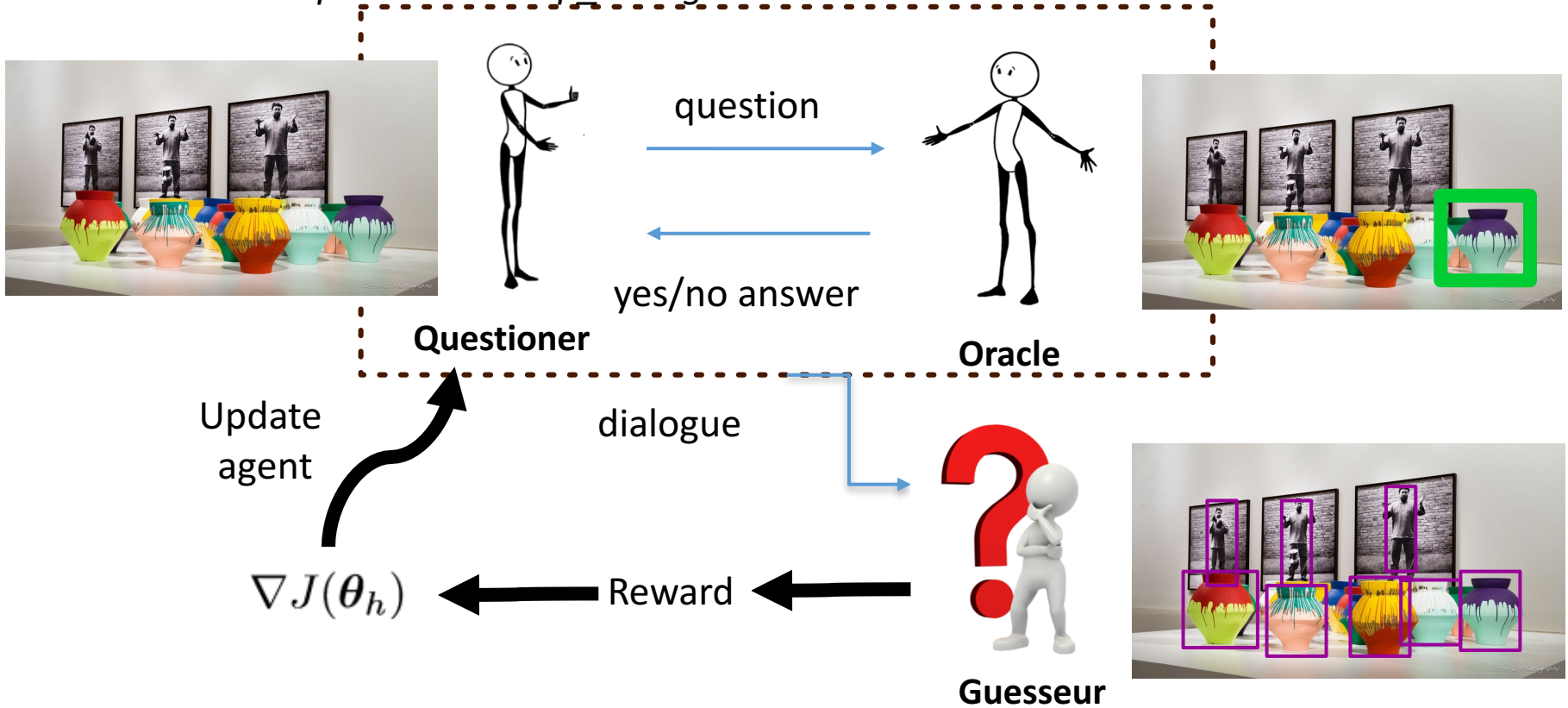


Guesseur



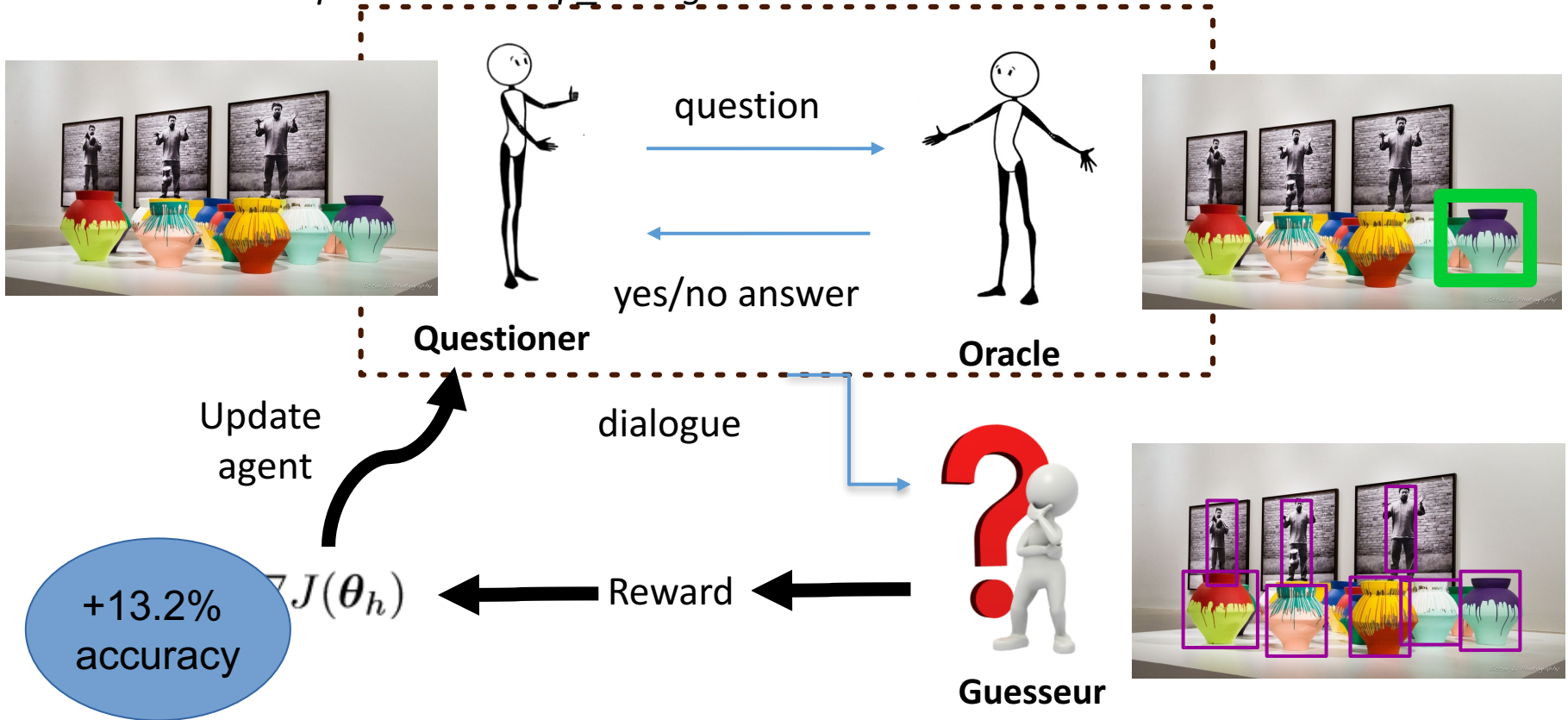
Models

Repeat until <stop_dialogue>

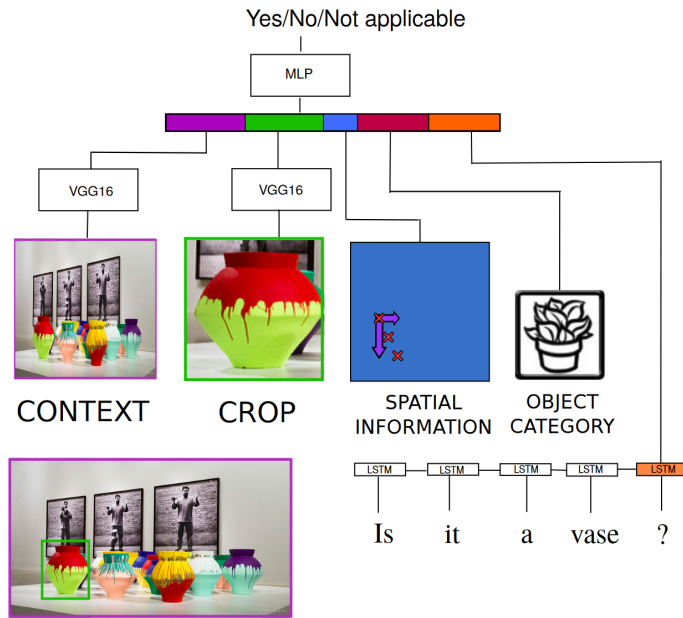


Models

Repeat until <stop_dialogue>



Multimodal Learning



Model	Test err
Dominant class (no)	50.9%
Question	41.2%
Image	46.7%
Crop	43.0%
Question + Crop	29.2%
Question + Image	39.8%
Question + Category	25.7%
Question + Spatial	31.3%
Question + Category + Spatial	21.5%
Question + Category + Crop	24.7%
Question + Spatial + Crop	26.2%
Question + Category + Spatial + Crop	22.1%
Question + Category + Spatial + Image	23.5%

- Image features are not helpful !!!

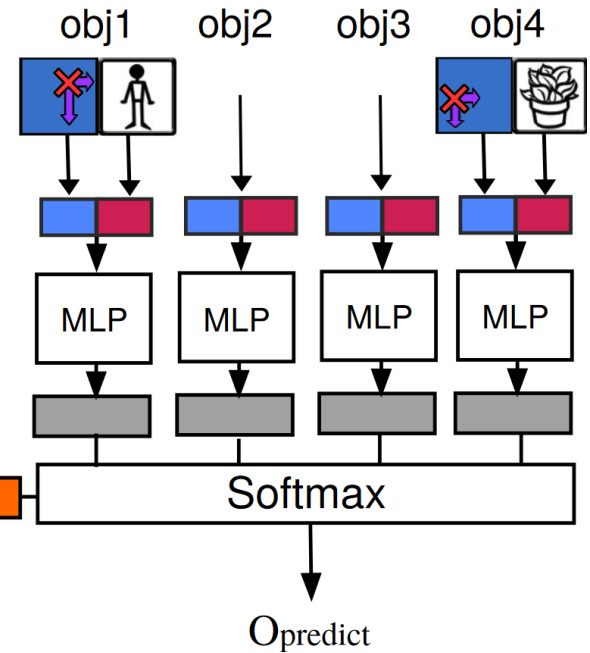
Multimodal Learning

Guesser ?

Is it a vase? Yes
Is it partially visible? No
Is it in the left corner? No
Is it the turquoise and purple one? Yes



LSTM / HRED encoder



63.8%
accuracy

Multimodal Learning

The RL quickly find the oracle/guesser limitations and focus and the question they can use



is it a person ? yes
is it in left ? no
is it in right ? no
is it in middle ? yes
is it in front ? yes
is it in middle ? yes

We need to improve the Oracle/Guesser to improve the language!

Multimodal Learning

Our models fail to fuse two modalities:

- First modality : Vision
- Second Modality: Language



Multimodal Learning

MS COCO [1]



The man at bat readies to swing at the pitch while the umpire looks on.



A large bus sitting next to a very tall building.

VQA [2]



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



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

RefeRit [3]

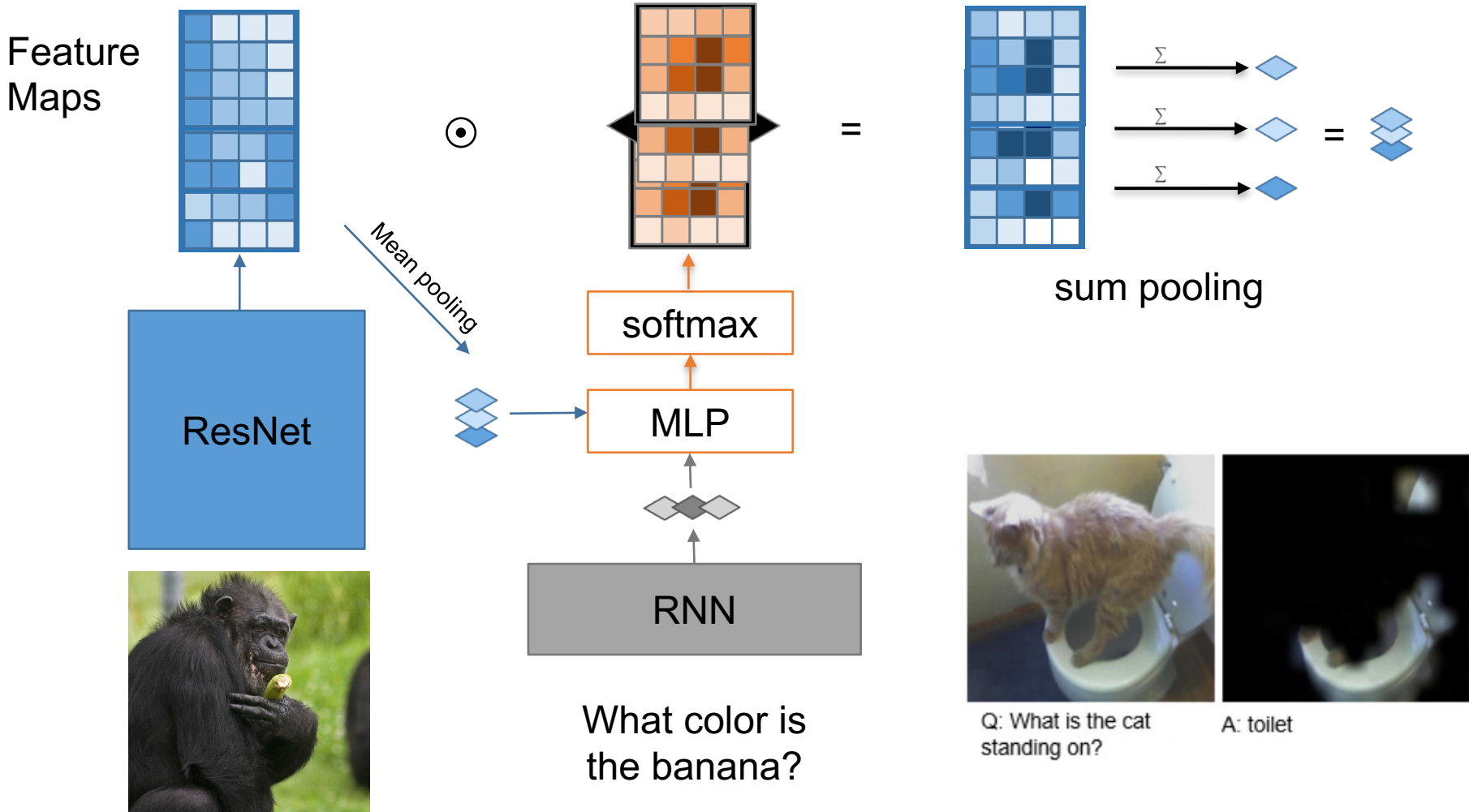


- The blue truck in the bottom right corner
- The light blue truck
- The blue truck on the right

Sample referring expressions for an object in a natural scene.

[1] T. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollar, and L. Zitnick. Microsoft coco: Common objects in context. In Proc of ECCV, 2014.
[2] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, Z. Lawrence, and D. Parikh. Vqa: Visual question answering. In Proc. of ICCV, 2015
[3] S. Kazemzadeh, V. Ordonez, M. Matten, and T. Berg. ReferItGame: Referring to Objects in Photographs of Natural Scenes. In Proc. of EMNLP, 2014.

Multimodal Learning



Multimodal Learning

Attention Mechanism



$$\xi_{w,h} = MLP([\mathbf{F}_{i,\cdot,w,h}; \mathbf{e}_q])$$

$$\alpha_{w,h} = \frac{\exp(\xi_{w,h})}{\sum_{w,h} \exp(\xi_{w,h})}$$

$$\mathbf{e}_v = \sum_{w,h} \alpha_{w,h} \mathbf{F}_{i,\cdot,w,h}$$

Concatenation

Dot product

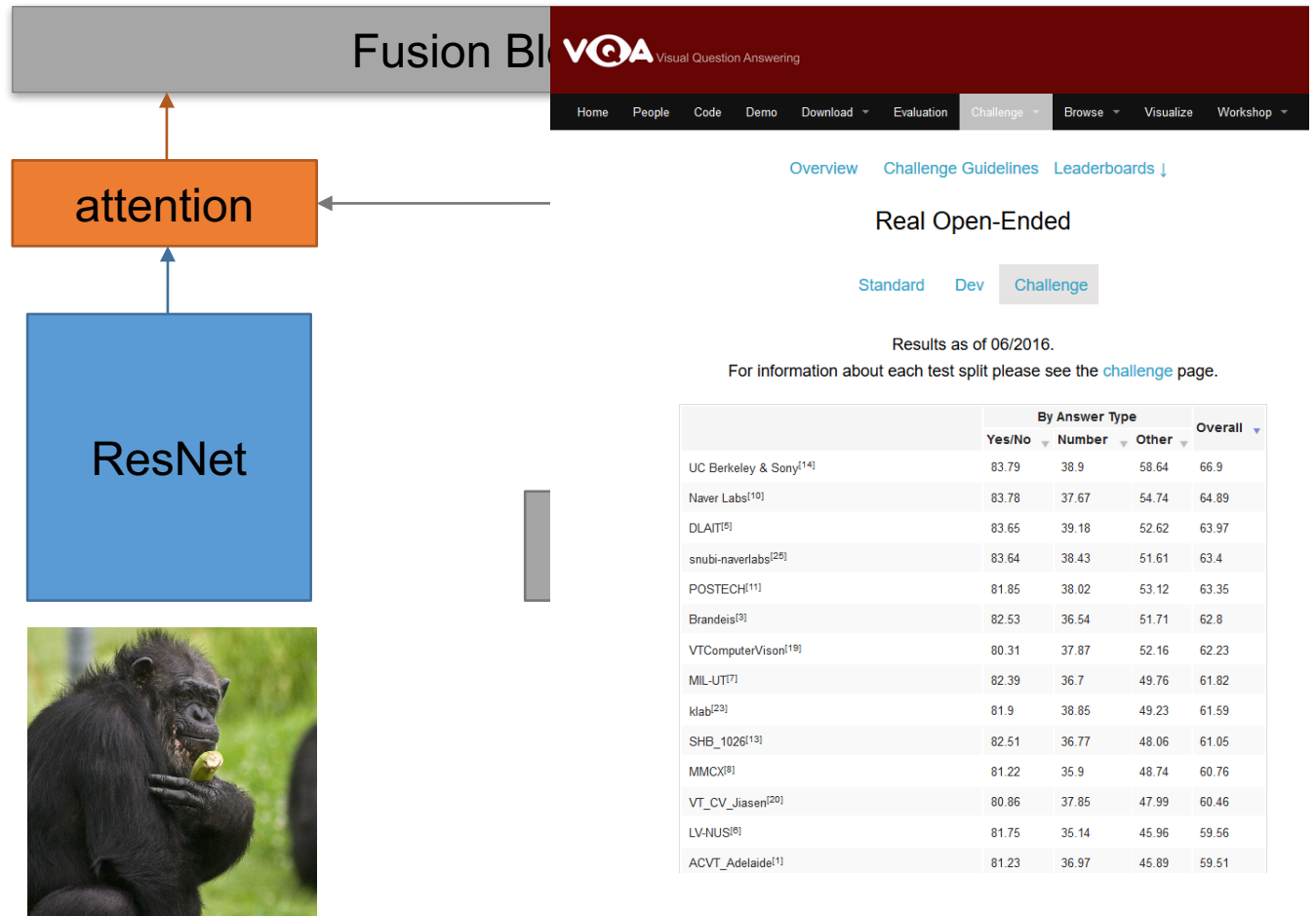
Random projection (MCB)

Linear projection + Dot Product (MLB)

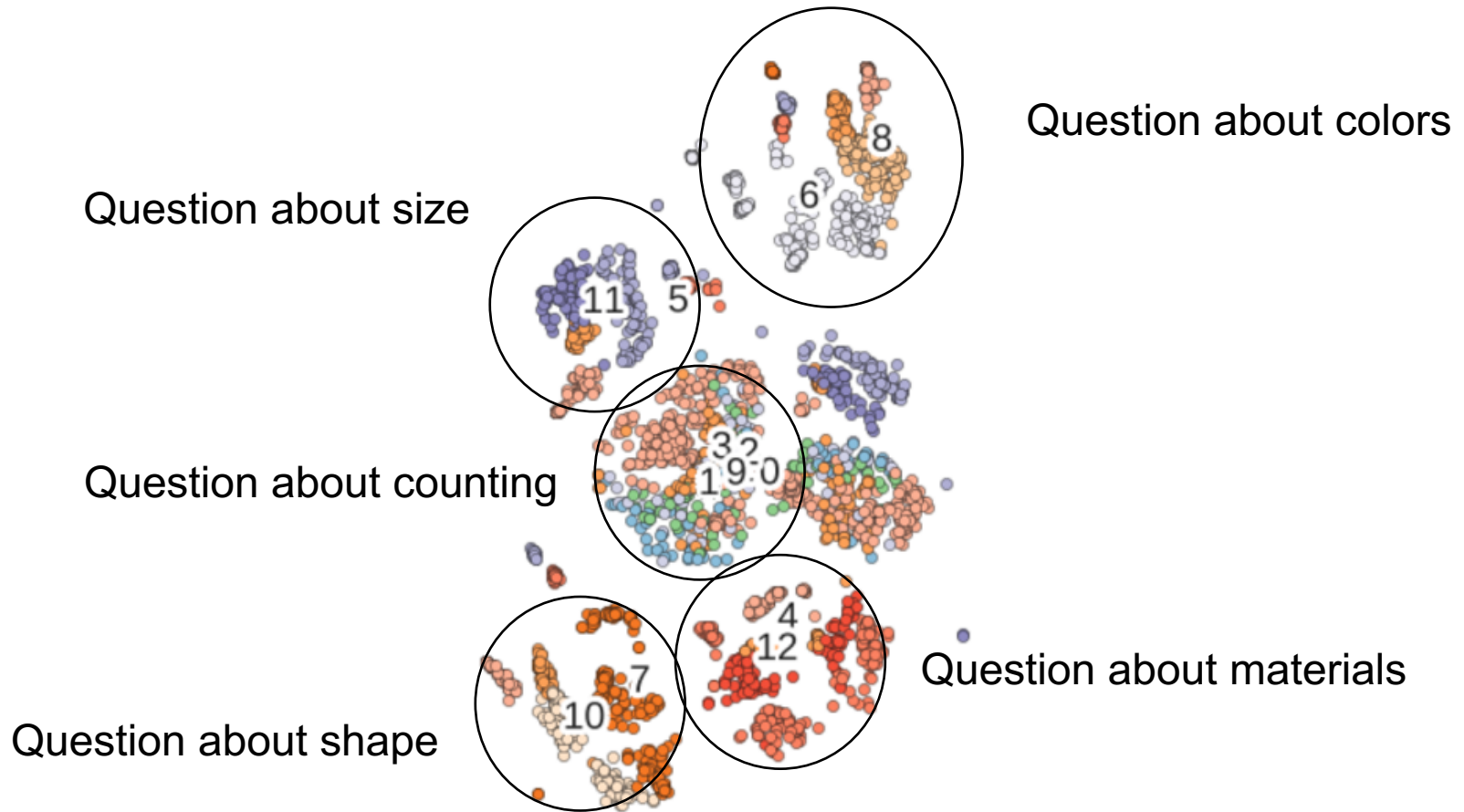
Tucker Decomposition

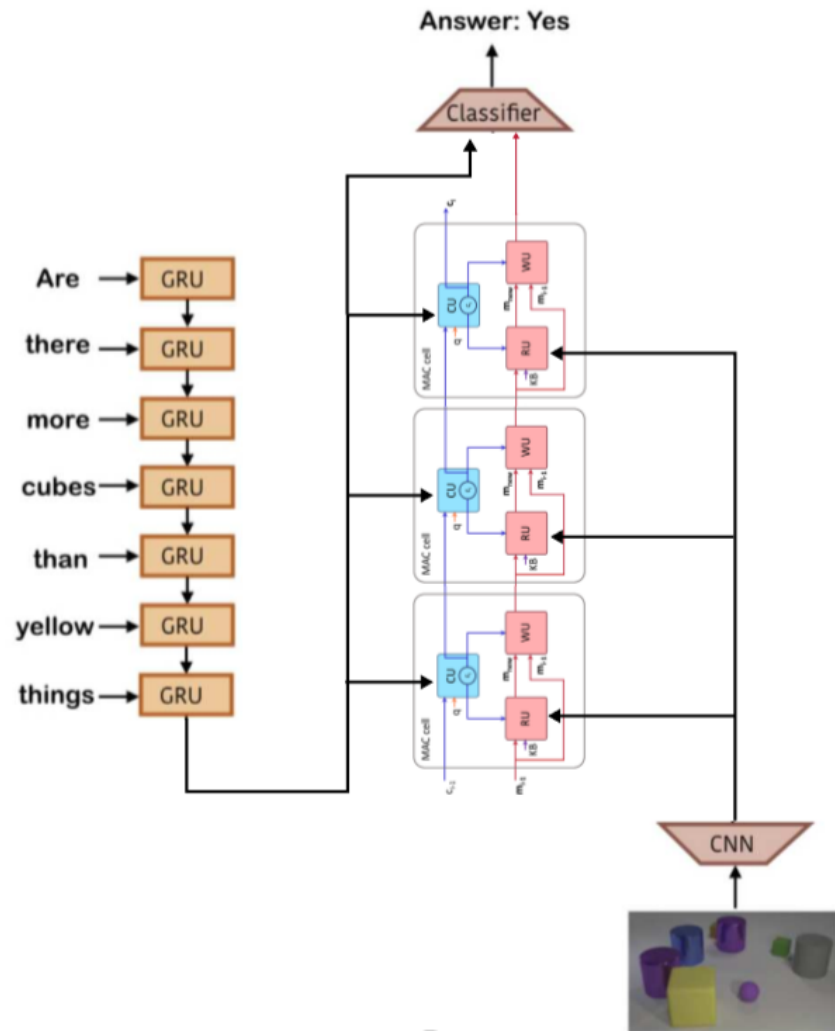
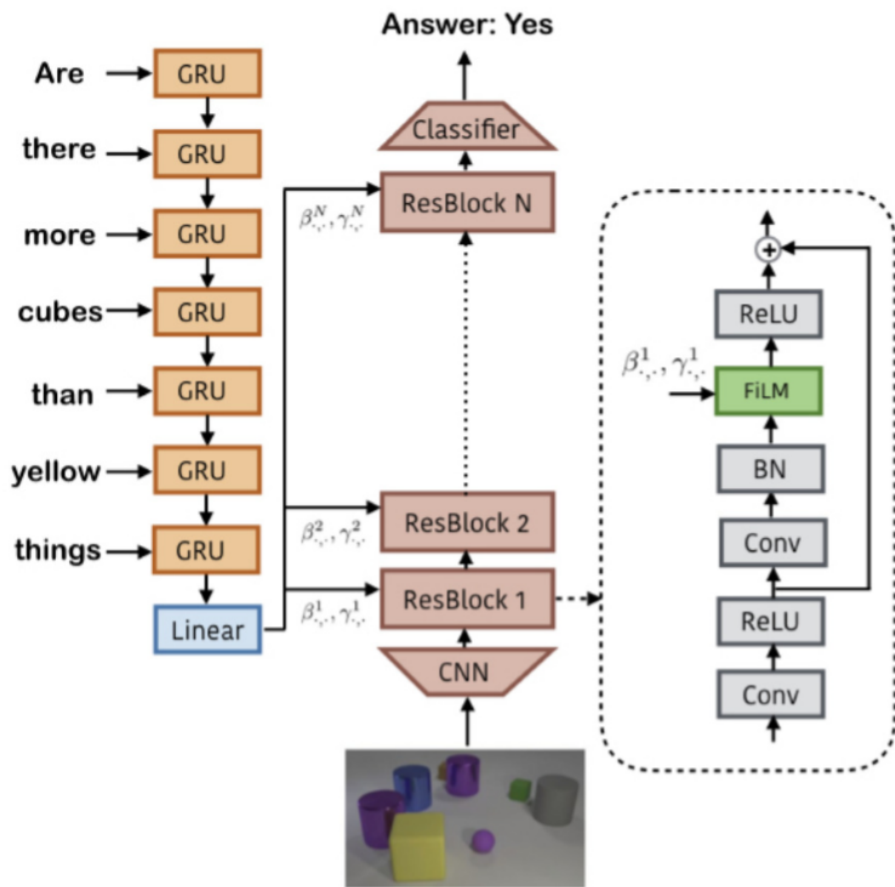
...

Multimodal Learning

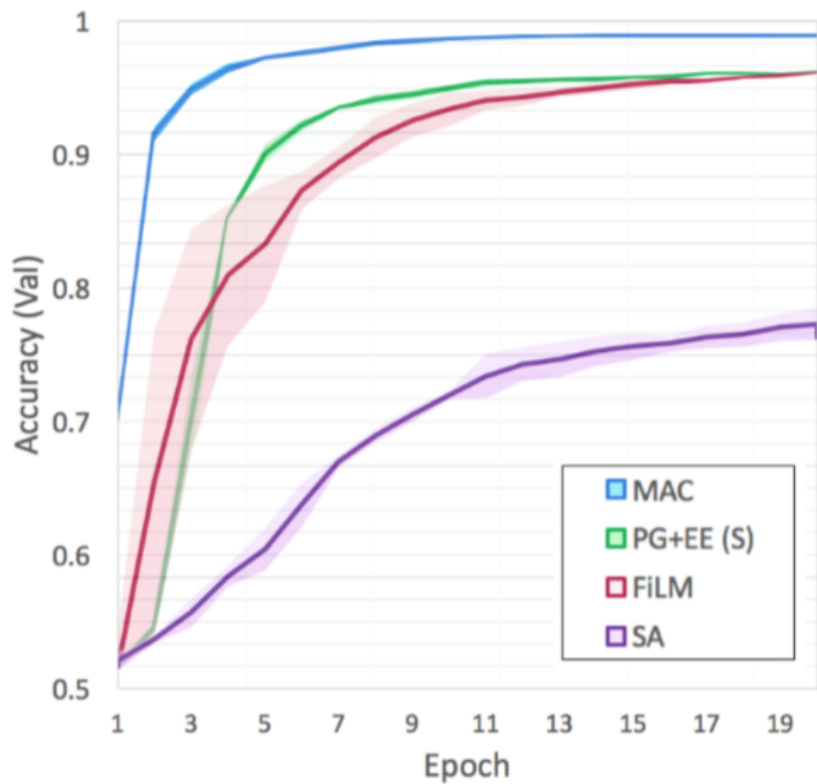


Multimodal Learning





Training Curve



Accuracy / Dataset Size (out of 700k)

