Ask Your Neurons: A Neural-based Approach to Answering Questions about Images

Author: Mateusz Malinowski, Marcus Rohrbach, Mario Fritz Presenter: Hooman Shariati, Wen Xiao

1.Introduction



Q: How many chairs are on the right side of the table in the image ?

A: 3

Q: What is in front of the door and on the right of the table in the image ?

A: chair

Q: What is in front of the white board or in front of the door in the image ?

A: table, chair

2

Visual Turing Test

- It is a system that generates a random sequence of binary questions specific to the test image, such that the answer to any question k is unpredictable given the true answers to the previous k-1 questions (also known as history of questions).
- Aim: evaluate the Image understanding of a computer system, and an important part of image understanding is the *story line of the image*.
- Give a new direction to the computer vision research which would lead to the introduction of systems that will be one step closer to *understanding images the way humans do*.

Wikipedia

Image credit to Gemana et al., PNAS 2015



1. Q: Is there a person in the blue region?	A: yes
2. Q: Is there a unique person in the blue region?	A: yes
(Label this person 1)	
3. Q: Is person 1 carrying something?	A: yes
4. Q: Is person 1 female?	A: yes
5. Q: Is person 1 walking on a sidewalk?	A: yes
6. Q: Is person 1 interacting with any other object?	A: no
1	
9. Q: Is there a unique vehicle in the yellow region?	A: yes
(Label this vehicle 1)	•
10. Q: Is vehicle 1 light-colored?	A: yes
11. Q: Is vehicle 1 moving?	A: no
12. Q: Is vehicle 1 parked and a car?	A: yes
14. O: Does vehicle 1 have exactly one visible tire?	A: no
15. O: Is vehicle 1 interacting with any other object?	A: no
17. Q: Is there a unique person in the red region?	A: no
18. Q: Is there a unique person that is female in the red region?	A: no
19. Q: Is there a person that is standing still in the red region?	A: yes
20. Q: Is there a unique person standing still in the red region?	A: yes
(Label this person 2)	
1	
23. Q: Is person 2 interacting with any other object?	A: yes
24. Q: Is person 1 taller than person 2?	A: amb
25. Q: Is person 1 closer (to the camera) than person 2?	A: no
26. Q: Is there a person in the red region?	A: yes
27. Q: Is there a unique person in the red region?	A: yes
(Label this person 3)	
1	
36. Q: Is there an interaction between person 2 and person 3?	A: yes
37. Q: Are person 2 and person 3 talking?	A: yes

3

Related Work (other tasks)

1. How can the computer 'see' the image?

CNN for Visual Recognition(as we did in A2)

2. How can the computer 'understand' the question and 'answer' the question?

RNN/LSTM for sequence modeling(as we did in A3)

3. Other task to combine the image and text?

Image caption(as we did in A4), Image grounding

Model: Ask Your Neurons(End-to-end)



Model: Ask Your Neurons



Predict the answer:

- $\hat{\boldsymbol{a}}_t = \operatorname*{arg\,max}_{\boldsymbol{a} \in \mathcal{V}} p(\boldsymbol{a} | \boldsymbol{x}, \boldsymbol{q}, \hat{A}_{t-1}; \boldsymbol{\theta})$
- $oldsymbol{x}$ image representation

$$oldsymbol{q} = egin{bmatrix} oldsymbol{q}_1, \dots, oldsymbol{q}_{n-1}, \llbracket ? \rrbracket \end{bmatrix}$$

- question word sequence

$$\hat{A}_{t-1} = \{ \hat{a}_1, \dots, \hat{a}_{t-1} \}$$

- the set of previous words

Notice: Loss only at answer words

LSTM



Non-linearity:

$$\sigma(v) = (1 + e^{-v})^{-1}$$

$$\phi(v) = \frac{e^v - e^{-v}}{e^v + e^{-v}} = 2\sigma(2v) - 1$$

LSTM equations:

$$i_t = \sigma(W_{vi}v_t + W_{hi}h_{t-1} + b_i)$$

$$f_t = \sigma(W_{vf}v_t + W_{hf}h_{t-1} + b_f)$$

$$o_t = \sigma(W_{vo}v_t + W_{ho}h_{t-1} + b_o)$$

$$g_t = \phi(W_{vg}v_t + W_{hg}h_{t-1} + b_g)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t$$

$$h_t = o_t \odot \phi(c_t)$$

Loss: Cross-entropy

CNN

- 1. CNN models are pre-trained on ImageNet dataset
- 2. GoogleNet consistently outperforms AlexNet





8

Neural model vs Symbolic model

Symbolic approach (NIPS'14)

- Explicit representation
- Independent components
- Detectors, Semantic Parser, Database
 - Segmented pictures
- Components trained separately
- Many 'hard' design decisions



Ask Your Neurons

- Implicit representation
- End-to-end formula
- From images and questions to answers
- Joint training
- Fewer design decisions



Neural Visual QA vs Neural Image Caption

Neural Image Description

- Conditions on an image
- Generates a description
- Sequence of words
- Loss at every step
- Hard to validate
 - Diversity of description



Ask Your Neurons (Our)

- Conditions on an image and a question
- Generates an answer
 - Sequence of answer words
- Loss only at answer words
- Easy to validate
 - Generally, questions has unique answers



10

Training:

GoogleNet pretrained on ImageNet.

Default hyper parameters for LSTM and CNN.

Randomly initialized the last FC layer of CNN, trained together with LSTM.

Train, validate and test on the same dataset DAQUAR as their previous work

No information on the training/validation/test set split of their data, or on the definition of their accuracy metrics.

We assumed it was the same as in their previous work.

Dataset

DAtaset for QUestion Answering on Real-world images (DAQUAR)

795 training images 6795 question-answer pairs

653 test images 5673 question-answer pairs

Asked 5 humans to provide questions and answers the only instructions were:

"Provide valid questions and answers related to basic colors, numbers, or types of both objects and sets of objects"

some biases showing humans tend to focus on a few prominent objects. For instance we have more than 400 occurrences of table and chair in the answers.



OA: (What is behind the table?, window) Spatial relation like 'behind' are dependent on the reference frame. Here the annotator uses observer-centric view.



OA: (what is beneath the candle holder, decorative plate) Some annotators use variations on spatial

cabinet)

interpretations.

relations that are similar, e.g. 'beneath' is closely related to 'below'. QA: (what is in front of the wall divider?.

Annotators use additional properties to clarify object references (i.e. wall divider). Moreover, the perspective plays an important role in these spatial relations

The annotators are using different names to call the same things. The names of the stand', 'stool', and 'cabinet'.

Some objects, like the table on the left of image, are severely occluded or truncated. Yet, the annotators refer to them in the questions.



OA: (What is in front of toilet?, door) Here the 'open door' to the restroom is not clearly visible, yet captured by the annotator.



OA: (what is behind the table?, sofa) Spatial relations exhibit different reference frames. Some annotations use observercentric, others object-centric view QA: (how many lights are on?, 6) Moreover, some questions require detection of states 'light on or off'



O: what is at the back side of the sofas? Annotators use wide range spatial relations, such as 'backside' which is object-centric.



QA1: (what is in front of the curtain behind the armchair?, guitar)

guitar)

Spatial relations matter more in complex environments where reference resolution becomes more relevant. In cluttered scenes, pragmatism starts playing a more important

OA2: (what is in front of the curtain?,

role

brown object near the bed include 'night



OA1: (How many doors are in the image?, 1 OA: (How many drawers are there?, 8) OA2:(How many doors are in the image?, 5 The annotators use their common-sense Different interpretation of 'door' results in different counts: 1 door at the end of the hall vs. 5 doors including lockers

the corner?, microwave)

frequently used by humans.

References like 'corner' are difficult to

resolve given current computer vision

models. Yet such scene features are

knowledge for amodal completion. Here the annotator infers the 8th drawer from the context

QA: (What is the shape of the green chair?, horse shaped) In this example, an annotator refers to a "horse shaped chair" which requires a quite abstract reasoning about the shapes.



OA: (How many doors are open?, 1) Notion of states of object (like open) is not well captured by current vision techniques. Annotators use such attributes frequently for disambiguation.



OA: (Where is oven?, on the right side of refrigerator) On some occasions, the annotators prefer to use more complex responses. With spatial relations, we can increase the answer's precision

Evaluation metrics:

- 1. Strict string matching: $\frac{1}{N} \sum_{i=1}^{N} \mathbf{1} \{A^i = T^i\} \cdot 100$
- 2. Semantic matching using WUP to account for word-level ambiguities:
 - a. I.e 'carton' and 'box' can be associated with similar concepts, so the model should not be strongly penalized for this type of mistakes.

Wu-Palmer (WUP) word similarity measure

- 1. Based on edge counting in a taxonomy like WorldNet or Ontology
- 2. WUP also weights the edges based on distance in the hierarchy.
 - a. Ex: Going from inanimate to animate is a larger distance than going from Felid to Canid.
- 3. WordNet:
 - Large lexical database of English words grouped into sets of cognitive synonyms (synsets),
 each expressing a distinct concept
 - b. Synsets are interlinked by means of conceptual-semantic and lexical relations



Figure credit to Malinowski et al., Multi Question ICCV,2014

WUPS (WUP Set)

$$\mathsf{WUPS}(A,T) = \frac{1}{N} \sum_{i=1}^{N} \min\{\prod_{a \in A^i} \max_{t \in T^i} \mathsf{WUP}(a,t), \prod_{t \in T^i} \max_{a \in A^i} \mathsf{WUP}(a,t)\} \cdot 100$$

Multiply WUP(a, b) with 0.1 whenever WUP(a, b) < t

$$WUPS(A,T) = \frac{1}{N} \sum_{i=1}^{N} \min\{\prod_{a \in A^i} \max_{t \in T^i} \mu(a,t), \prod_{t \in T^i} \max_{a \in A^i} \mu(a,t)\}$$

For precise answers, consider to words similar if WUP(a, b) > 0.9

t = 1, is same as string matching

Ground Truth	Predictions		
Armchair	Wardrobe	Chair	
Accuracy	0 =	= 0	
Wu-Palmer Similarity [1]	0.8 <	< 0.9	
WUPS @0.9 (NIPS'14)	≈0 <	< 0.9	

Evaluation

- Comparison with previous approach based on semantic parsing
- Comparison with how well questions can be answered without images
- Tried different subsets of the dataset and different accuracy metrics (to boost their score?)

Accu-	WUPS	WUPS
racy	@0.9	@0.0
7.86	11.86	38.79
17.49	23.28	57.76
19.43	25.28	62.00
50.20	50.82	67.27
17.06	22.30	56.53
17.15	22.80	58.42
7.34	13.17	35.56
	Accu- racy 7.86 17.49 19.43 50.20 17.06 17.15 7.34	Accu- racyWUPS @0.97.8611.8617.4923.2819.4325.2850.2050.8217.0622.3017.1522.807.3413.17

- Their performance drops dramatically with longer answers.
- They mention dataset bias:
 - 90% of the answers contain a single word



- 5 additional test answers for each image-question pair (by 5 additional people).
- Same directions as before.



 Their explanation as to why the benchmark performance of humans was 50%.

Two new scores to capture consensus

• Average Consensus Metric (ACM): Prefers mainstream answers.

$$\frac{1}{NK} \sum_{i=1}^{N} \sum_{k=1}^{K} \min\{\prod_{a \in A^{i}} \max_{t \in T_{k}^{i}} \mu(a, t), \prod_{t \in T_{k}^{i}} \max_{a \in A^{i}} \mu(a, t)\}$$

• Min Consensus Metric (MCM): Prefers closest matching answers.

$$\frac{1}{N}\sum_{i=1}^N \max_{k=1}^K \left(\min\{\prod_{a\in A^i} \max_{t\in T^i_k} \mu(a,t), \ \prod_{t\in T^i_k} \max_{a\in A^i} \mu(a,t)\} \right)$$

	Accu- racy	WUPS @0.9	WUPS @0.0			
Average Consensus Metric	33 .		1			
Language only (ours)						
- multiple words	11.60	18.24	52.68	17.06	22.30	56.53
- single word	11.57	18.97	54.39	17.15	22.80	58.42
Neural-Image-QA (ours)						
- multiple words	11.31	18.62	53.21	17.49	23.28	57.76
- single word	13.51	21.36	58.03	19.43	25.28	62.00
Min Consensus Metric						
Language only (ours)						
- multiple words	22.14	29.43	66.88			
- single word	22.56	30.93	69.82			
Neural-Image-QA (ours)						
- multiple words	22.74	30.54	68.17			
- single word	26.53	34.87	74.51			

	Accu-	Accu- WUPS	WUPS
	racy	@0.9	@0.0
Subset: No agreement	3 <u>5</u>		
Language only (ours)			
- multiple words	8.86	12.46	38.89
- single word	8.50	12.05	40.94
Neural-Image-QA (ours)			
- multiple words	10.31	13.39	40.05
- single word	9.13	13.06	43.48
Subset: $\geq 50\%$ agreement			
Language only (ours)			
- multiple words	21.17	27.43	66.68
- single word	20.73	27.38	67.69
Neural-Image-QA (ours)			
- multiple words	20.45	27.71	67.30
- single word	24.10	30.94	71.95
Subset: Full Agreement			
Language only (ours)			
- multiple words	27.86	35.26	78.83
- single word	25.26	32.89	79.08
Neural-Image-QA (ours)			
- multiple words	22.85	33.29	78.56
- single word	29.62	37.71	82.31

Some examples

Counting Questions:

What is on the right sid	e of the cabinet?	How many drawers are there?	What is the largest object?
Neural-Image-QA:	bed	3	bed
Language only:	bed	6	table

Color Questions

What is on the results	refrigerator?	What is the colour of the comforter?	What objects are found on the bed?
Neural-Image-QA:	magnet, paper	blue, white	bed sheets, pillow
Language only:	magnet, paper	blue, green, red, yellow	doll, pillow

Spatial Relationship Questions



Discussion

Strength and Weaknesses:

- Novel approach to an interesting problem.
 - But weak on evaluation. And low on implementation details
- Compared only to their own previous works (without even mentioning the architecture of their previous works or making a comparison).
- Unclear about training/test/validation splits and training parameters
- Modified both the dataset and evaluation metrics to boost their scores, to no avail.
 - Changed the number of answer words
 - Changed the number of provided ground truth answers for each question
 - Used several metrics

Difficulty with Spatial Relations

- Perform relatively well on "what color" and "how many" questions, but they have difficulty with questions like "what is to the left of the fridge".
 - Could be due to CNNs.
 - Providing more spatial information through an attention mechanism might help

- Also, difficulty with small objects, questions with negations, and shapes.
 - They attribute this to under-representation of these cases in training data

Doesn't Learn Enough From Images

- Humans answer 7.34% without images, and 50.20 % with images.
- Their system answers 17.06% without images, but only 17.49% with images.

• Our suggestions:

- Increase the learning capacity of the portion of the model that learns from images
- Encode the entire question first and pass it along with the image to a seperate LSTM
- Pre-train the LSTM on a different question/answers set to reduce dependence on the particular question/answers contained in the training set.

VQA: Visual Question Answering

www.visualqa.org

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

Abstract—We propose the task of *free-form* and *open-ended* Visual Question Answering (VQA). Given an image and a natural language question about the image, the task is to provide an accurate natural language answer. Mirroring real-world scenarios, such as helping the visually impaired, both the questions and answers are open-ended. Visual questions selectively target different areas of an image, including background details and underlying context. As a result, a system that succeeds at VQA typically needs a more detailed understanding of the image and complex reasoning than a system producing generic image captions. Moreover, VQA is amenable to automatic evaluation, since many open-ended answers contain only a few words or a closed set of answers that can be provided in a multiple-choice format. We provide a dataset containing ~0.25M images, ~0.76M questions, and ~10M answers (www.visualqa.org), and discuss the information it provides. Numerous baselines and methods for VQA are provided and compared with human performance.

1 INTRODUCTION

We are witnessing a renewed excitement in multi-discipline Artificial Intelligence (AI) research problems. In particular, research in image and video captioning that combines Computer Vision (CV), Natural Language Processing (NLP), and Knowledge Representation & Reasoning (KR) has dramatically increased in the past year [14], [7], [10], [36], [24], [22], [51]. Part of this excitement stems from a belief that multi-discipline tasks like image captioning are a step towards solving AI. However, the current state of the art demonstrates that a coarse scene-level understanding of an image paired with word n-gram statistics suffices to generate reasonable image captions, which suggests image captioning may not be as "AI-complete" as desired.

What makes for a compelling "AI-complete" task? We believe that in order to spawn the next generation of AI algorithms, an ideal task should (i) require *multi-modal knowledge* beyond a single sub-domain (such as CV) and (ii) have a well-defined *quantitative evaluation metric* to track progress. For some



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



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





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



Does it appear to be rainv

Does this person have 20/20 vision?

3

9



http://vqa.cloudcv.org/

Visual Question Answering(CVPR 2016)

- Visual Question Answering Dataset (VQA):
 - 250K images (COCO and abstract scenes)
 - 760K questions
 - 10M answers by multiple people
 - "yes/no", "number", and "object" answers; majority single word
 - Has confidence and Consensus measures (i.e. how many people agree on a given answer)
- Opens the way for automatic evaluation
 - many open-ended answers contain only a few words or a closed set of answers that can be provided in a multiple-choice format: http://visualqa.org/visualize/
- Adds human baseline performance and compares previous VQA methods



Is something under the sink broken?	yes yes	no
 What number do you see?	33 33 33	567



Does this man have children?	yes yes	ye ye
Is this man crying?	no	no
	no	ye
	no	ye



How many glasses are on the table?	3 3 3	2 2 6
What is the woman reaching for?	door handle glass wine	fruit glass remote



an you park ere?	ou park no no no	
What color is he hydrant?	white and orange white and orange white and orange	red red yello



as the pizza been	yes	yes
aked?	yes	yes
hat kind of cheese is opped on this pizza?	feta feta ricotta	mozzarella mozzarella mozzarella



boy on the ground	yes	no
has broken legs?	yes	yes
Why is the boy	his friend is hurt	ghost
on the right	other boy fell down	lightning
freaking out?	someone fell	sprayed by ho
	-	



What kind of store is this?	bakery bakery pastry	art supplies grocery grocery
Is the display case as full as it could be?	no	no
	no	yes
	10.00	11000



How many pickles are on the plate?	1 1	1
What is the shape of the plate?	circle round round	circle round round



Are the kids in the room the grandchildren of the adults?	yes yes	yes yes yes
What is on the bookshelf?	nothing nothing nothing	book book book



How many bikes are there?	2 2 2 2	3 4 12
What number is	48	4
the bus?	48	number 6



	What does the sign say?	stop	stop yield
	What shape is this sign?	octagon octagon	diamond octagon round



How many balls are there?	2222	1 2 3
What side of the	right	left
teeter totter is on the ground?	right right side	left right side

37 Credit to Agrawal et al., Visual Question answering, ICCV,2016

Yin and Yang: Balancing and Answering Binary Visual Questions (CVPR 2016)



complementary scenes

Tuple: <girl, walking, bike> Question: Is the girl walking the bike?

Credit to Zhanget al., Yin and Yang, ICCV,2016

Making the V in VQA Matter: (CVPR 2017)

- Answers why models ignore visual information
 - Inherent structure in our world and bias in language are easier signals to learn from

- Suggests a way to counter language priors
 - For each question, collect complementary images such that every question is associated with pair of similar images that result in two different answers to the question.

• Balanced VQA dataset

Where is the child sitting? fridge arms







Who is wearing glasses? man woman





Is the umbrella upside down? yes no





How many children are in the bed?



