

“Unit-testing” deep learning with synthetic data for more informative evaluation

Alexander Kuhnle

Department of Computer Science and Technology
University of Cambridge (UK)

MILA – 10th October 2018

Overview

- ▶ Visual question answering
- ▶ Problems with the VQA Dataset
- ▶ Evaluation methodology
- ▶ ShapeWorld generation framework
- ▶ Evaluation of FiLM on ShapeWorld

Visual question answering

Examples



Where is this cat laying?
Is the cat awake?
What color is the cat?



Is the cat facing the computer?
Is the cat typing?
Is the cat playing with the mouse?



What object is shining on the animal?
What objects is the cat sitting behind?
How many cats?

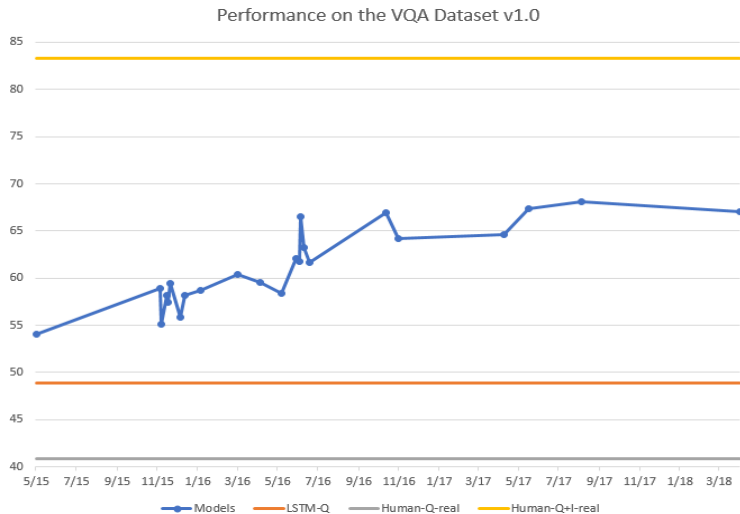


How many items are on the bookcase?
Are these two children related?
Is the dog begging for food?

⇒ **Visual Turing test?**

Visual question answering

Performance over time



Based on (incomplete) list of VQA papers with arXiv publication dates

Problems with the VQA Dataset

Question-answer biases

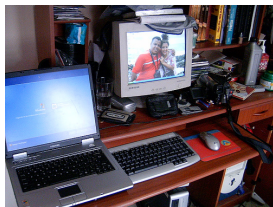


- ▶ What sport is...? \Rightarrow tennis (41%)

- ▶ How many...? \Rightarrow two (39%)



- ▶ Do you see a...? \Rightarrow yes (87%)



Problems with the VQA Dataset

Complete question/image understanding



- ▶ What...? ⇒ umbrella
- ▶ What season...? ⇒ summer
- ▶ What season of...? ⇒ summer
- ▶ ...
- ▶ What season of year was this photo taken in?
⇒ summer



- ▶ What does the red sign say? ⇒ stop



Examples from Agrawal et al. (<https://arxiv.org/abs/1606.07356>) and Devi Parikh's slides (<https://newgeneralization.github.io/>)

Problems with the VQA Dataset

Sensitivity to question words

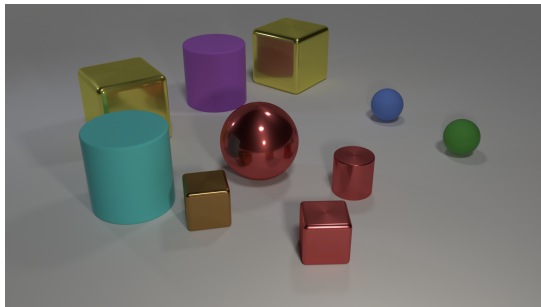


- ▶ How symmetrical are the white bricks on either side of the building? ⇒ very
- ▶ How **spherical** are the white bricks on either side of the building? ⇒ very
- ▶ How **soon** are the bricks **fading** on either side of the building? ⇒ very
- ▶ How **fast** are the bricks **speaking** on either side of the building? ⇒ very

Example from Mudrakarta et al. (<https://arxiv.org/abs/1805.05492>).

Problems with the VQA Dataset

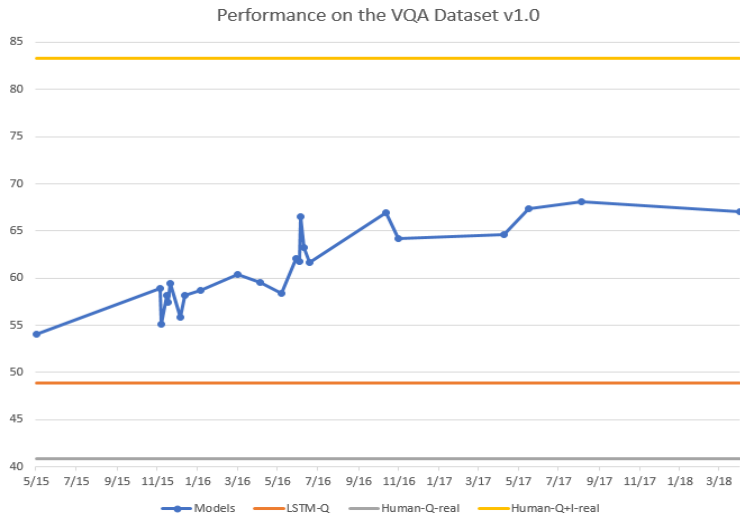
Low performance on CLEVR



- ▶ How many small spheres are there? $\Rightarrow 2$
- ▶ What number of cubes are small things or red metal objects? $\Rightarrow 2$
- ▶ Does the metal sphere have the same color as the metal cylinder? \Rightarrow Yes
- ▶ Are there more small cylinders than metal things? \Rightarrow No

Evaluation methodology

Meaningful progress?



Based on (incomplete) list of VQA papers with arXiv publication dates

Evaluation methodology

Pros and cons of crowd-sourced real-world datasets

Solve the problem/dataset?



Deep learning will find a way to make effective use of the data.

Evaluate model capabilities?



Are these datasets appropriate to investigate this question?

- ▶ Natural?
- ▶ Difficult?
- ▶ Specific?

⇒ **Synthetic data!**

Evaluation methodology

Other popular datasets with similar issues

SNLI – Stanford Natural Language Inference Corpus

C: A soccer game with multiple males playing.

H: Some men are playing a sport.

→ **entailment**

C: A smiling costumed woman is holding an umbrella.

H: A happy woman in a fairy costume holds an umbrella.

→ **neutral**

C: A man inspects the uniform of a figure in some East Asian country.

H: The man is sleeping

→ **contradiction**

SQuAD – Stanford Question Answering Dataset

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under **gravity**. The main forms of precipitation include drizzle, rain, sleet, snow, **graupel** and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals **within a cloud**. Short, intense periods of rain in scattered locations are called “showers”.

(1) What causes precipitation to fall?

⇒ **gravity**

(2) What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?

⇒ **graupel**

(3) Where do water droplets collide with ice crystals to form precipitation? ⇒ **within a cloud**

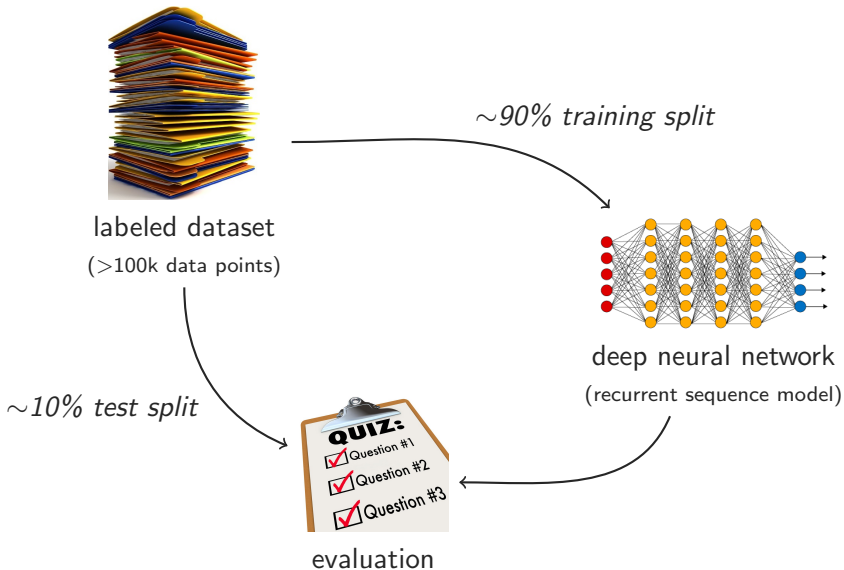
Evaluation methodology

“Growing pains” for deep learning evaluation

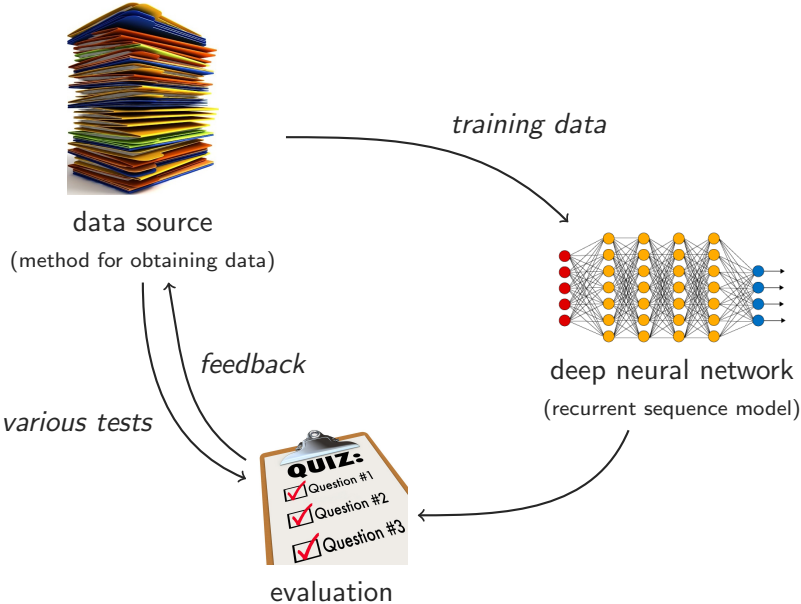
- ▶ Dataset bias and “cheating” models
 - ▶ Unexpectedly simple data and strong baselines
 - ▶ Adversarial examples with unintuitive model behavior
 - ▶ Replication and task/dataset transfer failure
- ⇒ **Symptoms of insufficient/inappropriate evaluation**

Evaluation methodology

Current approach

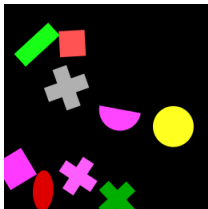


Evaluation methodology

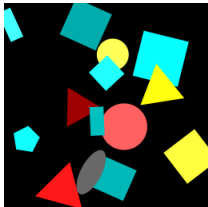


ShapeWorld generation framework

Examples: relations and quantifiers



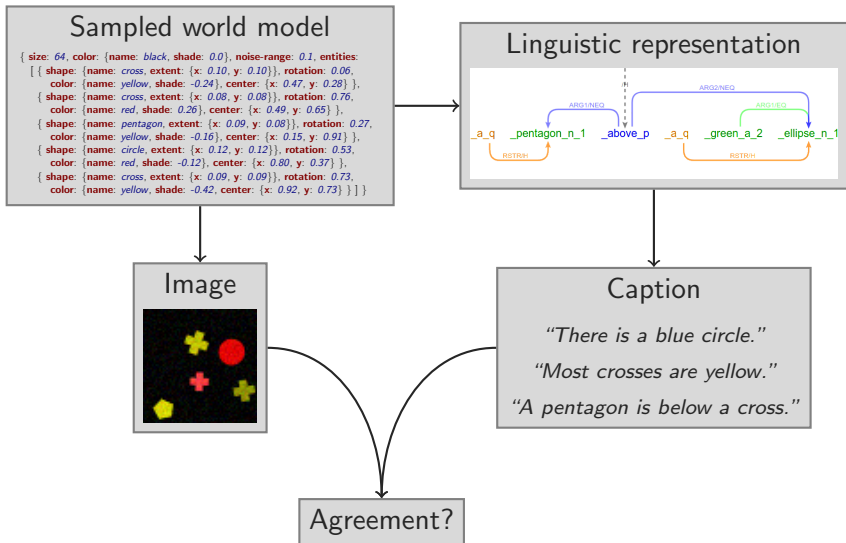
- ▶ A magenta square is to the right of a green shape.
- ▶ A yellow shape is not in front of a square.
- ▶ A circle is farther from an ellipse than a gray cross.
- ▶ A cross is not the same color as a green rectangle.
- ▶ The lowermost green shape is a cross.
- ▶ A red shape is the same shape as a green shape.



- ▶ Less than one triangle is cyan.
- ▶ At least half the triangles are red.
- ▶ More than a third of the shapes are cyan squares.
- ▶ Exactly all the five squares are red.
- ▶ More than one of the seven cyan shapes is a square.
- ▶ Twice as many red shapes as yellow shapes are circles.

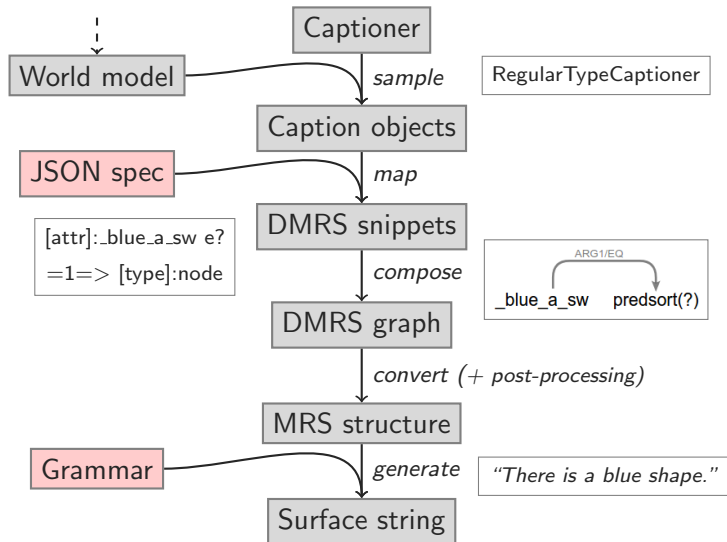
ShapeWorld generation framework

System overview



ShapeWorld generation framework

Language generation

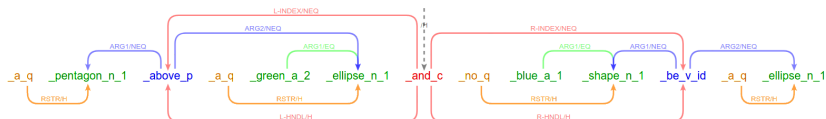


ShapeWorld generation framework

Compositionality

“A pentagon is above a green ellipse, and no blue shape is an ellipse.”

↑ ERG + ACE realization ↑



↑ Internal DMRS mapping ↑

$\exists a$	$a.shape=pg$	$a.y>b.y$	$\exists b$	$b.color=gr$	$b.shape=el$	\wedge	$\neg\exists c$	$c.color=bl$	true	$c=d$	$\exists d$	$d.shape=el$
$\exists a$	$a.shape=pg$	$a.y>b.y$	$\exists b$	$b.color=gr \wedge b.shape=el$		\wedge	$\neg\exists c$	$c.color=bl$		$c=d$	$\exists d$	$d.shape=el$
	$\exists a : a.shape=pg \wedge [\exists b : b.color=gr \wedge b.shape=el \wedge a.y>b.y]$					\wedge	$\neg\exists c : c.color=bl \wedge [\exists d : d.shape=el \wedge c=d]$					
$(\exists a : a.shape=pg \wedge [\exists b : b.color=gr \wedge b.shape=el \wedge a.y>b.y]) \wedge (\neg\exists c : c.color=bl \wedge [\exists d : d.shape=el \wedge c=d])$												

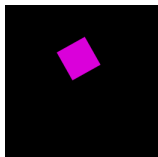
ShapeWorld generation framework

Design choices

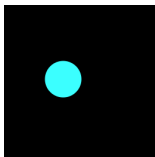
- ▶ Caption is extracted from image, i.e. world model
- ▶ Incorrect caption via minimal modification of correct one
- ▶ Three agreement values to avoid ambiguous cases
- ▶ Initialize generator/captioner values before sampling
- ▶ Various tautology/contradiction checks
- ▶ Modular and configurable

ShapeWorld generation framework

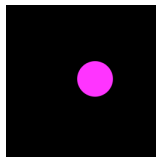
What type of generalization do we expect/desire?



magenta square



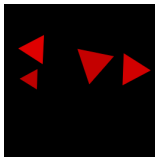
cyan circle



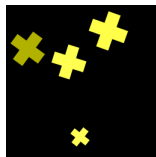
magenta circle



three crosses



four triangles

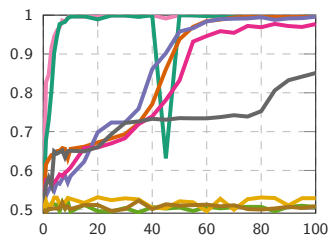


four crosses

Evaluation of FiLM on ShapeWorld

Results per instance type

Dataset	CNN-LSTM		CNN-LSTM-SA		FiLM	
(single-shape)	—	—	—	—	100.0	87.2
existential	100.0	81.1	100.0	99.7	100.0	99.9
logical	79.7	62.2	76.5	58.4	99.9	98.9
numbers	75.0	66.4	99.1	98.2	99.6	99.3
quantifiers	72.1	69.1	84.8	80.8	97.7	97.0
(simple-spatial)	81.4	64.8	81.9	57.7	85.1	61.3
relational	—	—	—	—	50.6	51.0
implicit-rel	—	—	—	—	52.9	53.2
superlatives	—	—	—	—	50.8	50.2

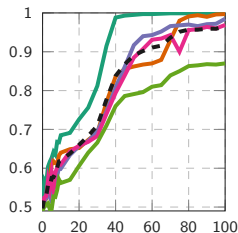


- ▶ Can relational-like instances implicitly be learned when training on a broader set of instances?
- ▶ Can relational-like instances be learned when (pre)training on simpler pedagogical instances?

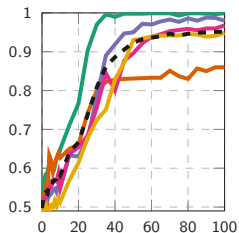
Evaluation of FiLM on ShapeWorld

Learning from a broader set of instances

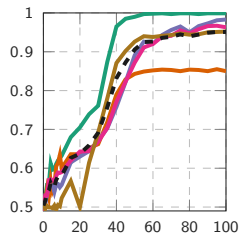
relational



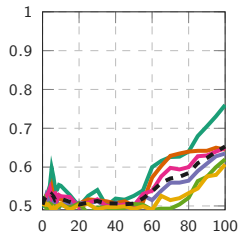
implicit-relational



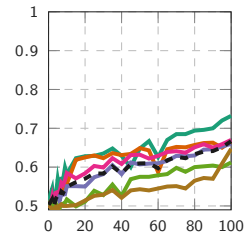
superlatives



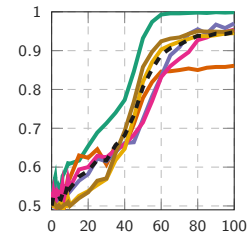
relation + implicit-rel



relational + superlat



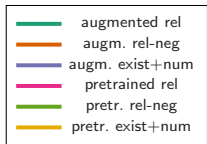
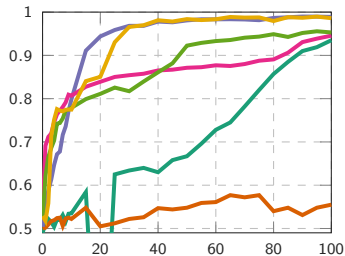
implicit-rel + superlat



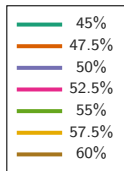
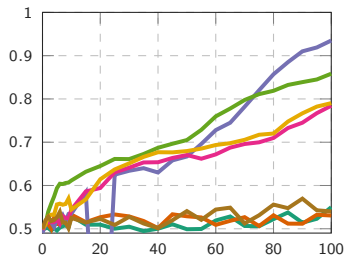
Evaluation of FiLM on ShapeWorld

Learning bootstrapped by simpler instances

augmentation vs pretraining



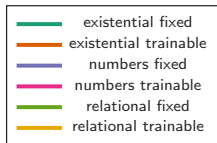
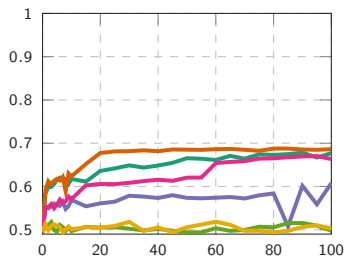
augmentation distributions



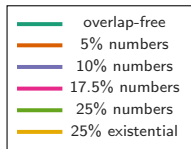
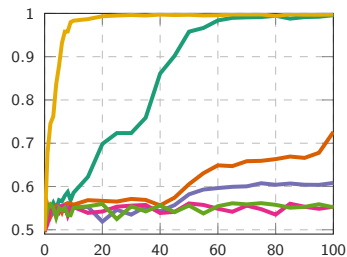
Evaluation of FiLM on ShapeWorld

Additional findings

pretrained ResNet doesn't work



overlapping objects impede learning



Conclusion

real-world data vs **synthetic data**

limited and expensive \longleftrightarrow unlimited amount

uncontrolled content \longleftrightarrow clean content

sparse instance coverage \longleftrightarrow targeted instance coverage

monolithic benchmark \longleftrightarrow tailored unit tests

test interpolation ability \longleftrightarrow test extrapolation ability

\Rightarrow **Complementary evaluation paradigms**



Thank you for your attention!

Questions?