

2A

Deep Compositional Question Answering With Neural Module Networks

GHAZAL SAHEBZAMANI, DELARAM BEHNAMI, PARDESS DANAEI, GOLARA JAVADI

19 MARCH 2019

Introduction

2

- ▶ Visual question answering based on neural module networks (NMNs)
- ▶ Understanding of both visual scenes and natural language

What color is the necktie?



Yellow

Acknowledgement

3

► PAPER 1

- Andreas J, Rohrbach M, Darrell T, Klein D. Neural module networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition 2016 (pp. 39-48).
- Andreas J, Rohrbach M, Darrell T, Klein D. Learning to compose neural networks for question answering. 2016. PowerPoint Presentation.

► PAPER 2

- Wu Q, Teney D, Wang P, Shen C, Dick A, van den Hengel A. Visual question answering: A survey of methods and datasets. Computer Vision and Image Understanding. 2017 Oct 1;163:21-40.

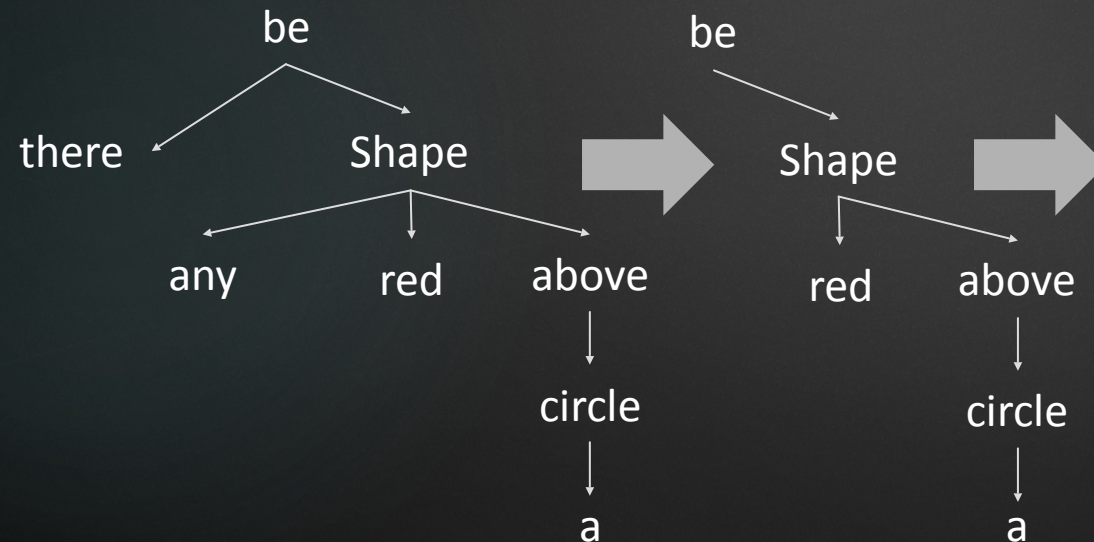
► PAPER 3

- Hu R, Andreas J, Darrell T, Saenko K. Explainable neural computation via stack neural module networks. In Proceedings of the European Conference on Computer Vision (ECCV) 2018 (pp. 53-69).

orks

NMN: Assemble a network on the fly from a collection of jointly-learned modules
 Analyze question with semantic parser
 Determine computational units based on the analysis
 Ground the

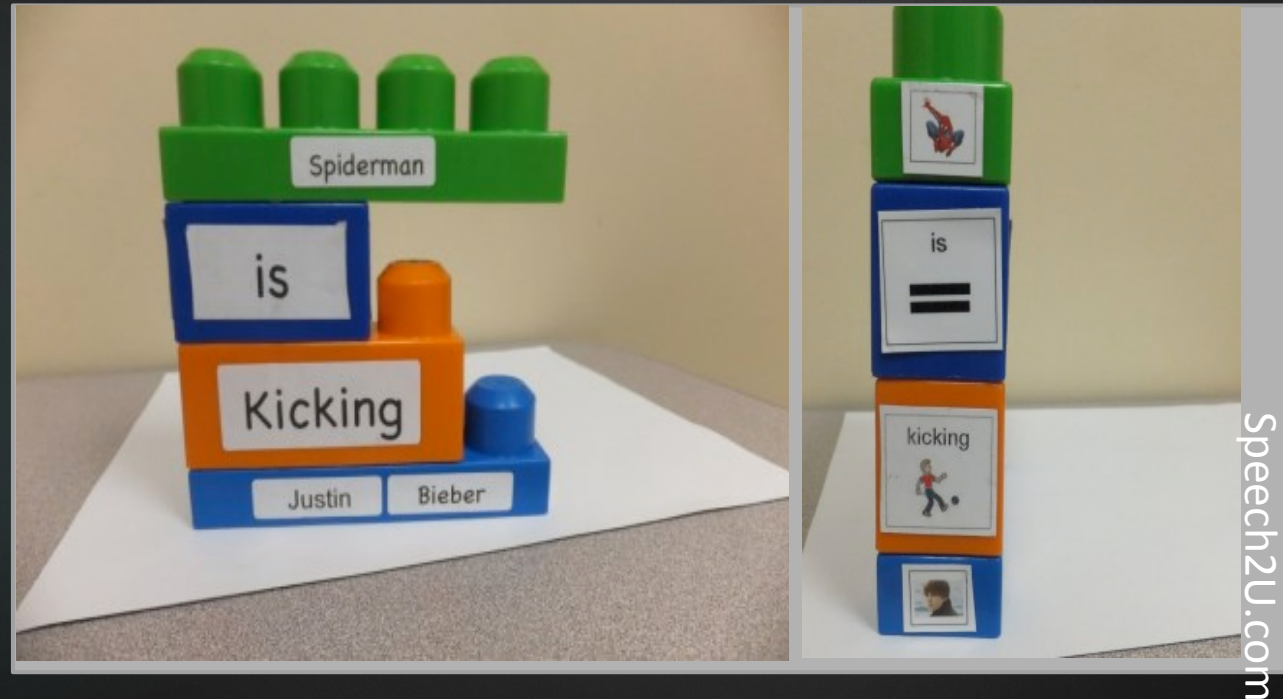
- ▶ Neural networks representation
 - ▶ Represent question as bags of words
 - ▶ Represent question using a recurrent neural network
 - ▶ Train a simple classifier on the encoded question and image
- ▶ Semantic Parsers



Approach

NMN: Assemble a network on the fly from a collection of jointly-learned modules
Common to initialize systems for vision tasks with a prefix of a network trained

- ▶ Assemble a network on the fly from a collection of jointly-learned modules
- ▶ No single “best network” for all tasks
- ▶ Question answering is a highly-multitask learning setting

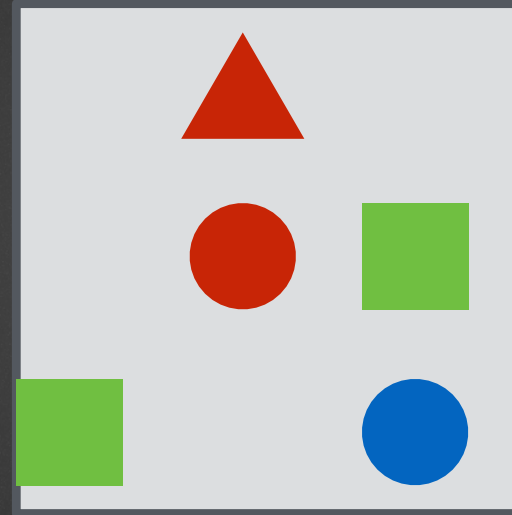


Approach

Ground the question on the image
Represent images with visual features
and attention
Messages passed between modules
may be raw image features, attentions,

- ▶ Ground the question on the image
 - ▶ Neural nets learn lexical groundings

Is there a red shape above a
circle?



Yes

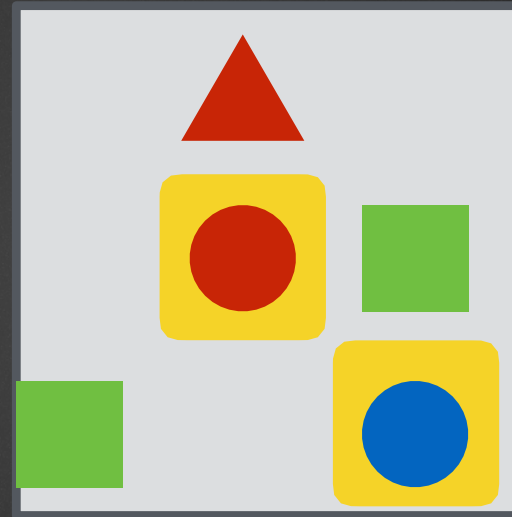
- ▶ Use a recurrent network (LSTM) to read the question for final answer

Approach

Ground the question on the image
Represent images with visual features
and attention
Messages passed between modules
may be raw image features, attentions,

- ▶ Ground the question on the image
 - ▶ Neural nets learn lexical groundings

Is there a red shape above a
circle?



Yes

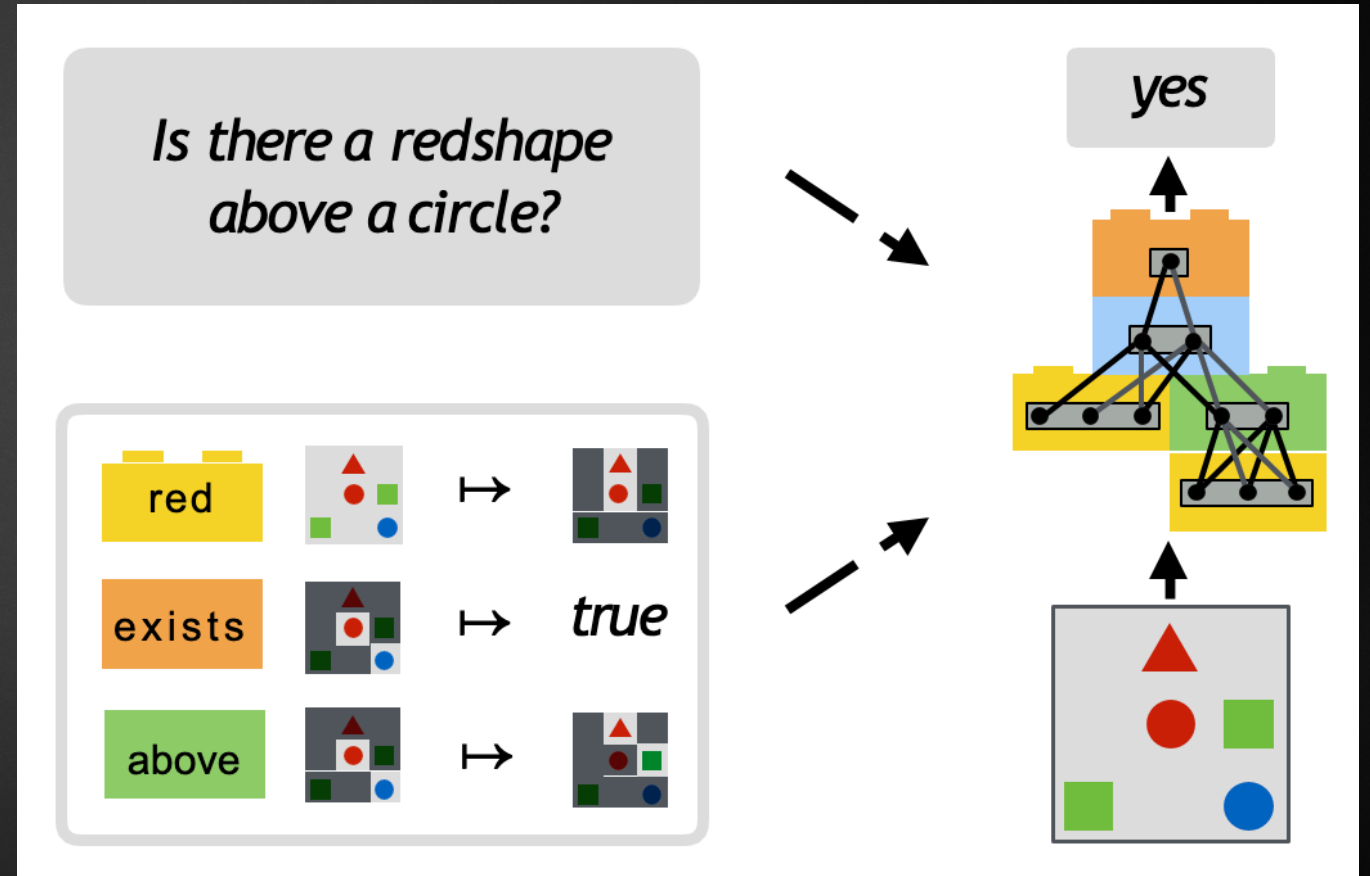
- ▶ Use a recurrent network (LSTM) to read the question for final answer

Neural Module Networks (NMNs)


each module has a set of parameters,
and a Network layout predictor

Instantiate network based on network
predictor

- ▶ Model is a collection of modules
- ▶ Training datum
 - ▶ Natural language question
 - ▶ Image
 - ▶ Answer
- ▶ Classifier

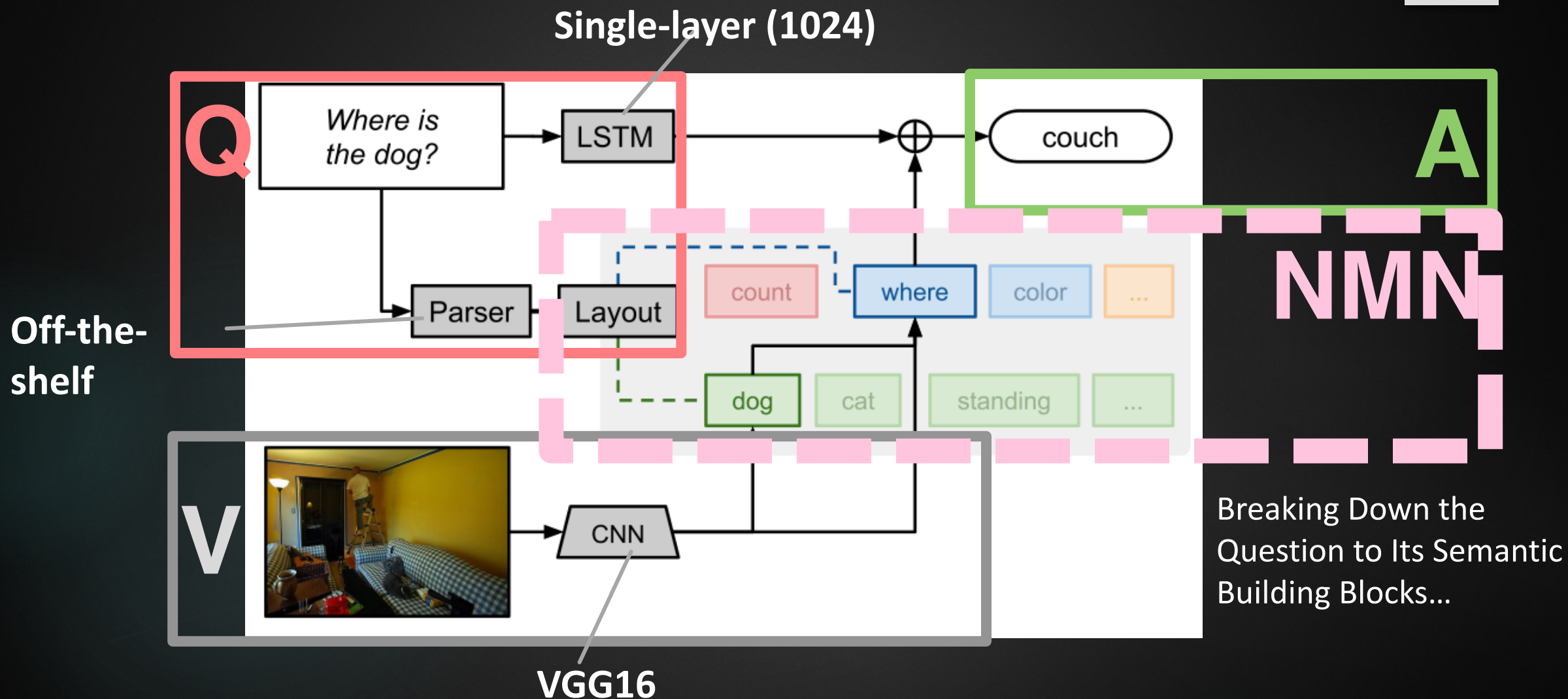


VQA System

V	
Q	<i>Is there a red shape above a circle?</i>
A	<i>Yes.</i>

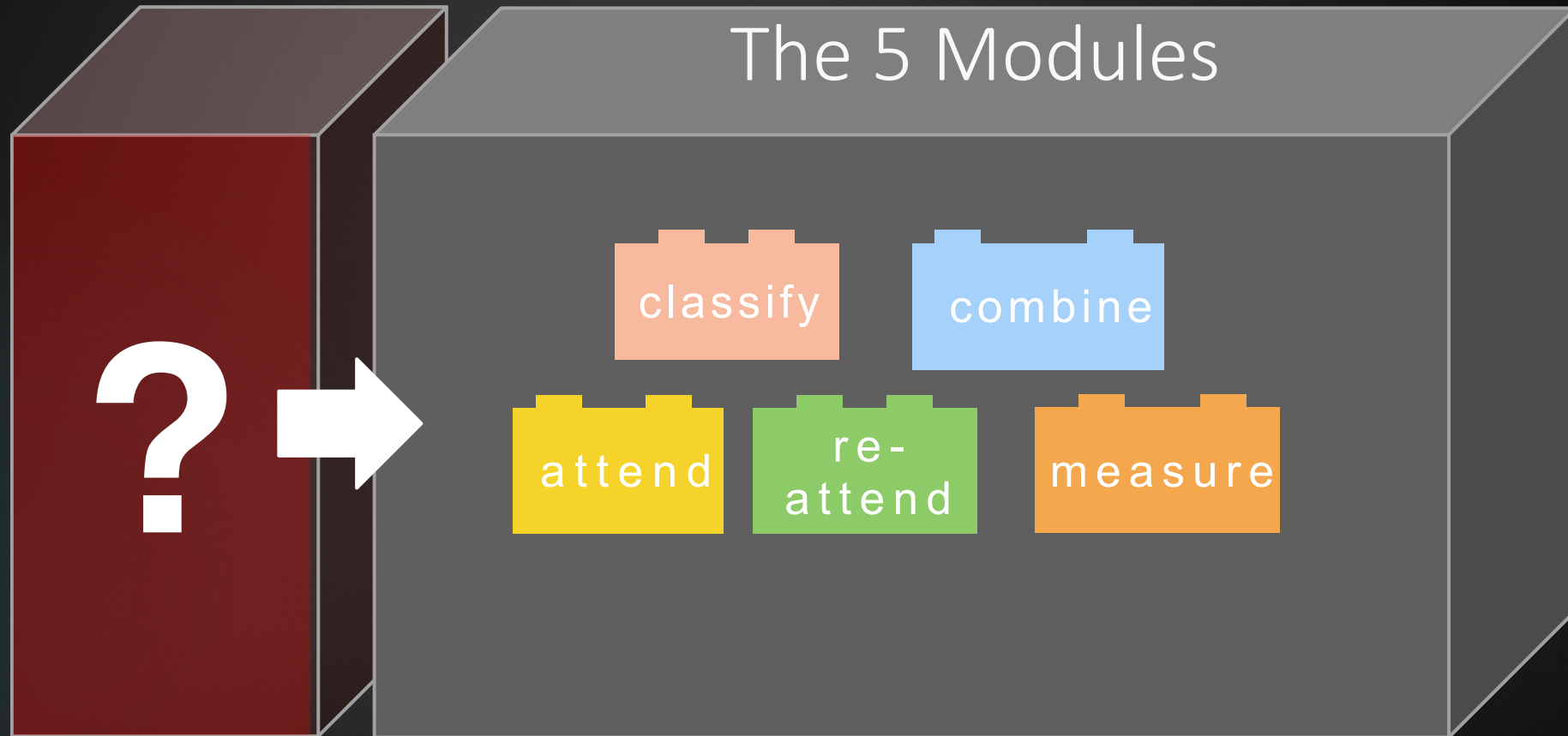
Compositional VQA with NMNs

10



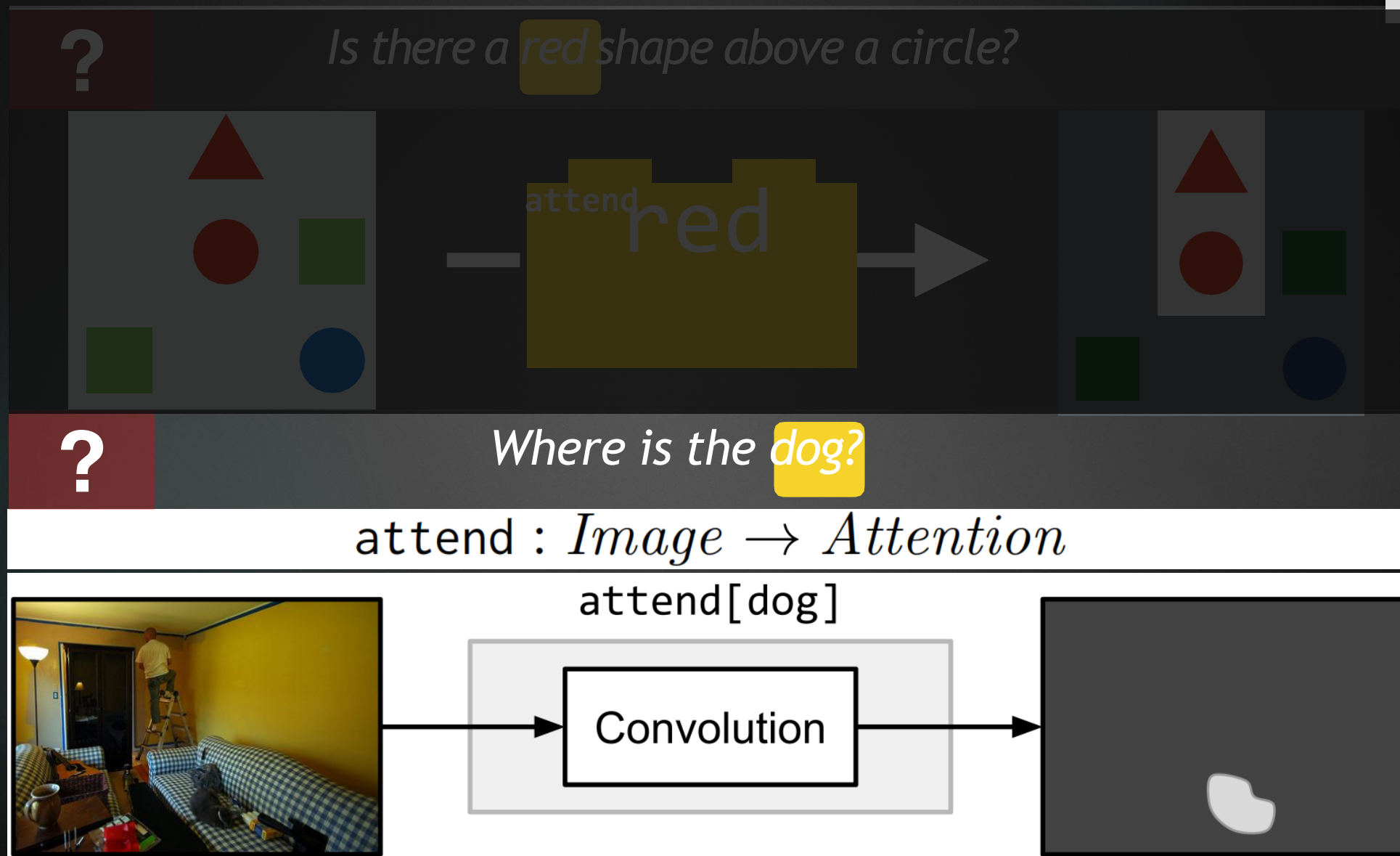
What Are the Building Blocks of Questions in VQA?

11



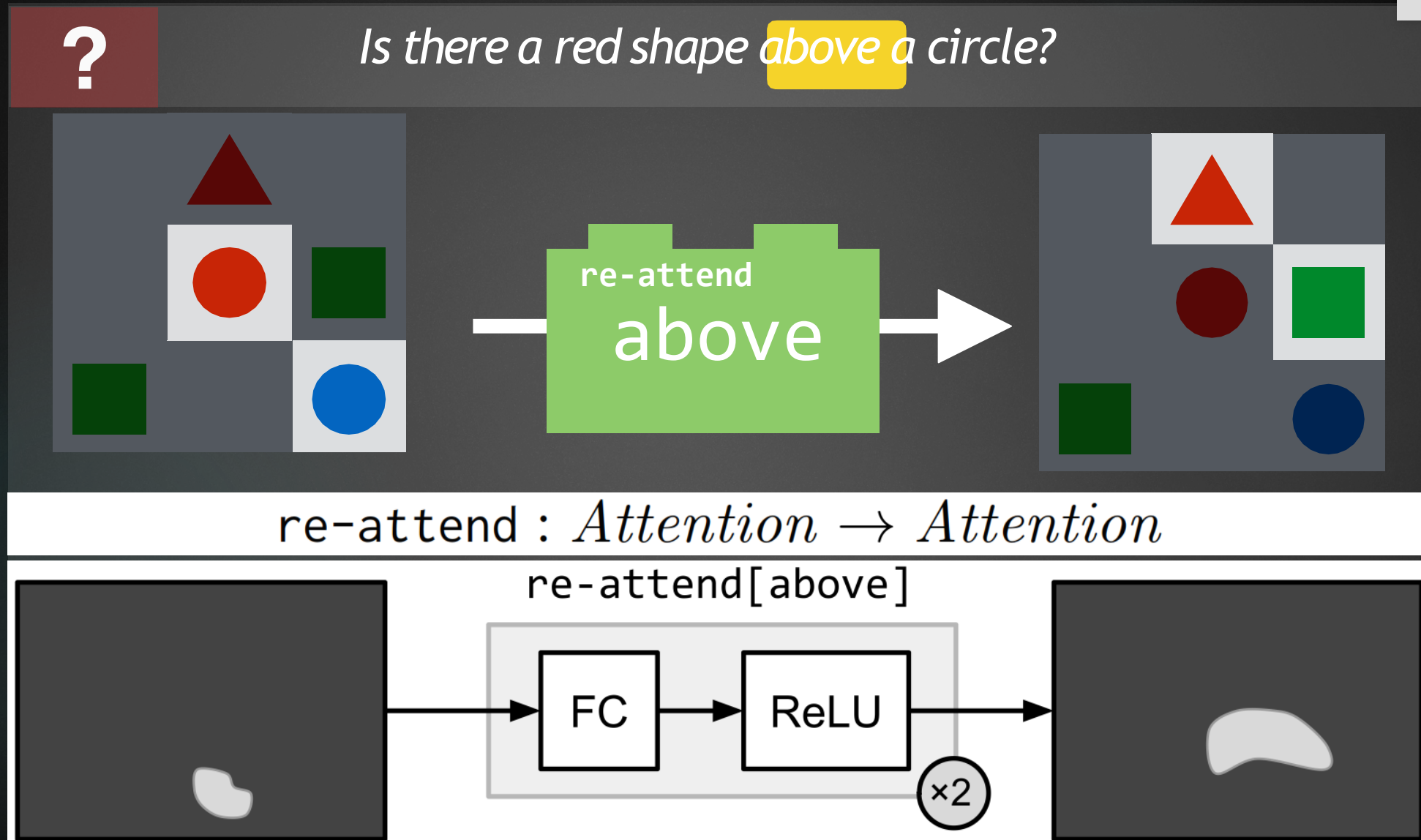
1. The **attend** Module for Attention to Predicates

12



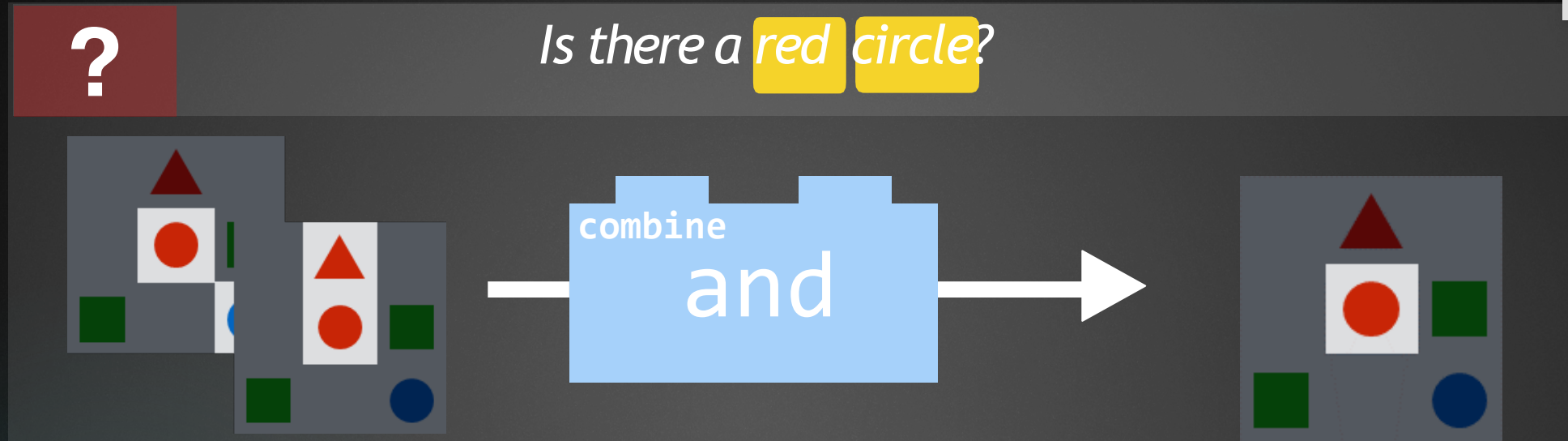
2. The **re-attend** Module for Relationships

13

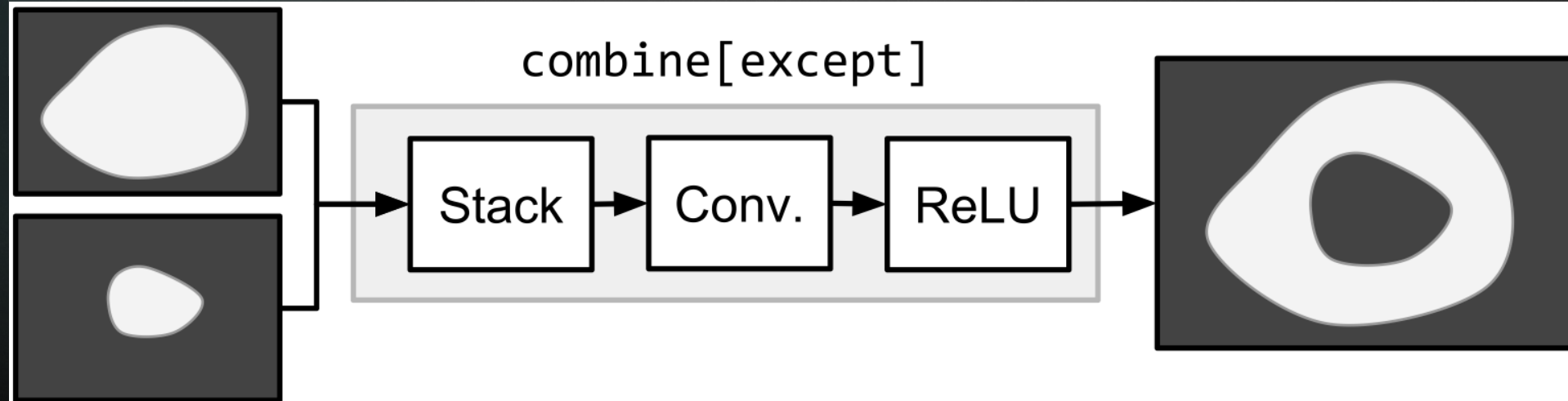


3. The **combine** Module for Logical Operations

14



$\text{combine} : \text{Attention} \times \text{Attention} \rightarrow \text{Attention}$

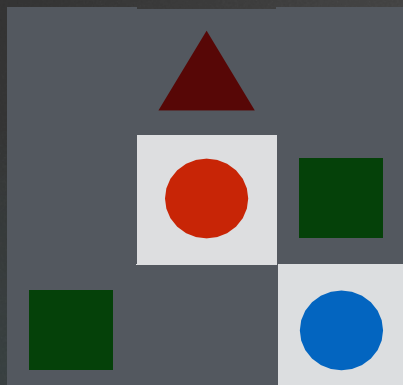


4. The **measure** Module

15

?

Is there a red shape above a circle?



measure
exists

True

$\text{measure} : \text{Attention} \rightarrow \text{Label}$

$\text{measure}[\text{exists}]$



FC

ReLU

FC

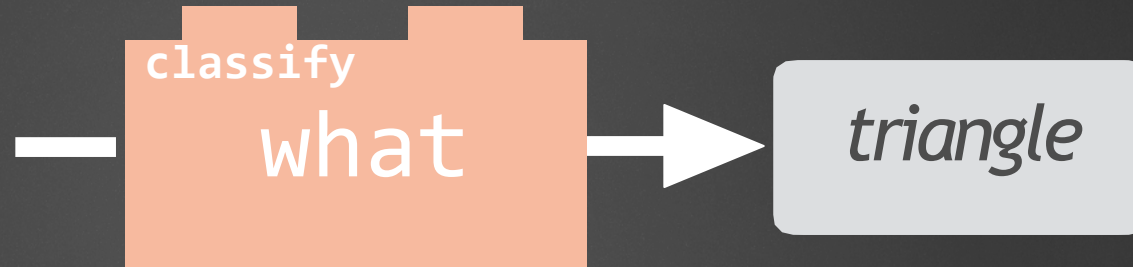
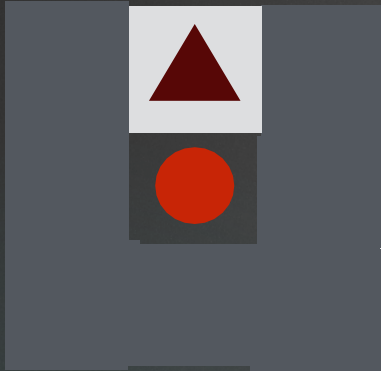
Softmax

yes

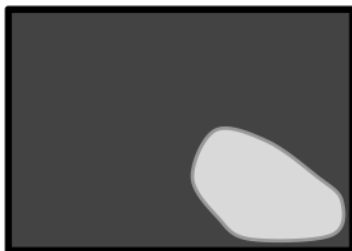
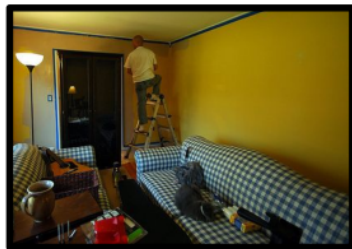
5. The **classify** Module

?

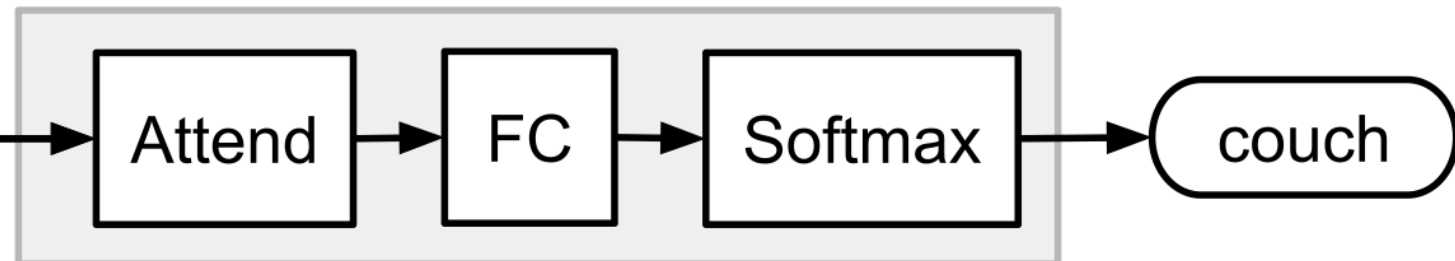
What is the a red shape above a circle?



$\text{classify} : \text{Image} \times \text{Attention} \rightarrow \text{Label}$



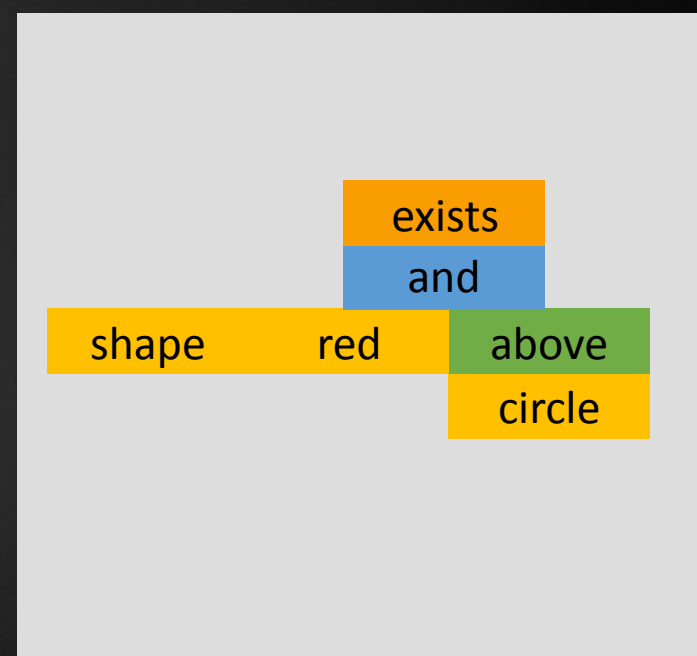
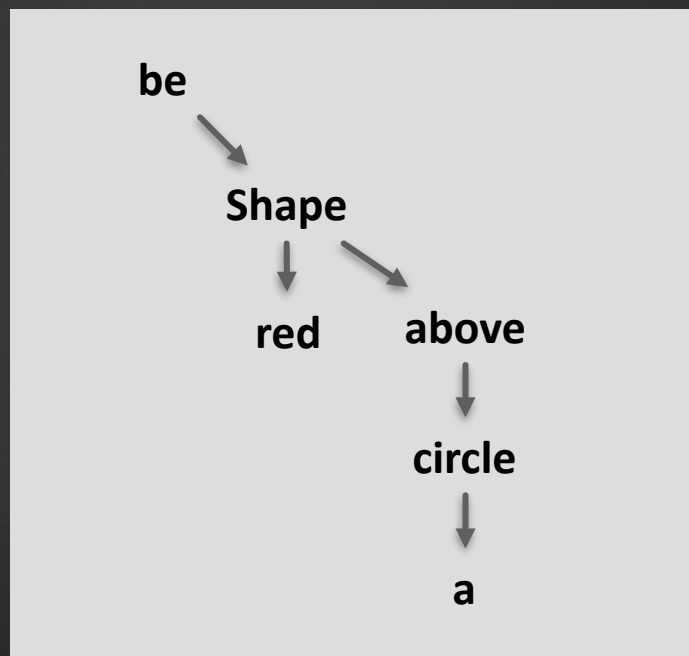
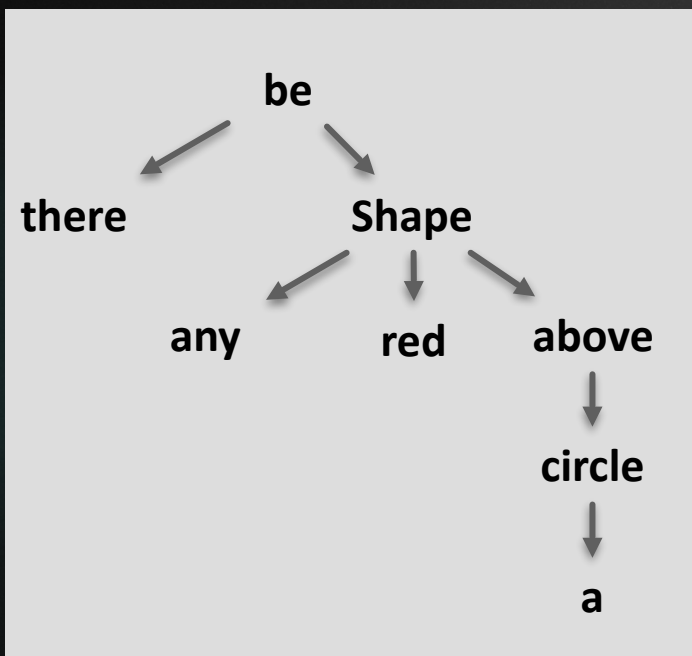
$\text{classify}[\text{where}]$



Training NMNs

17

Is there a red shape above a circle?



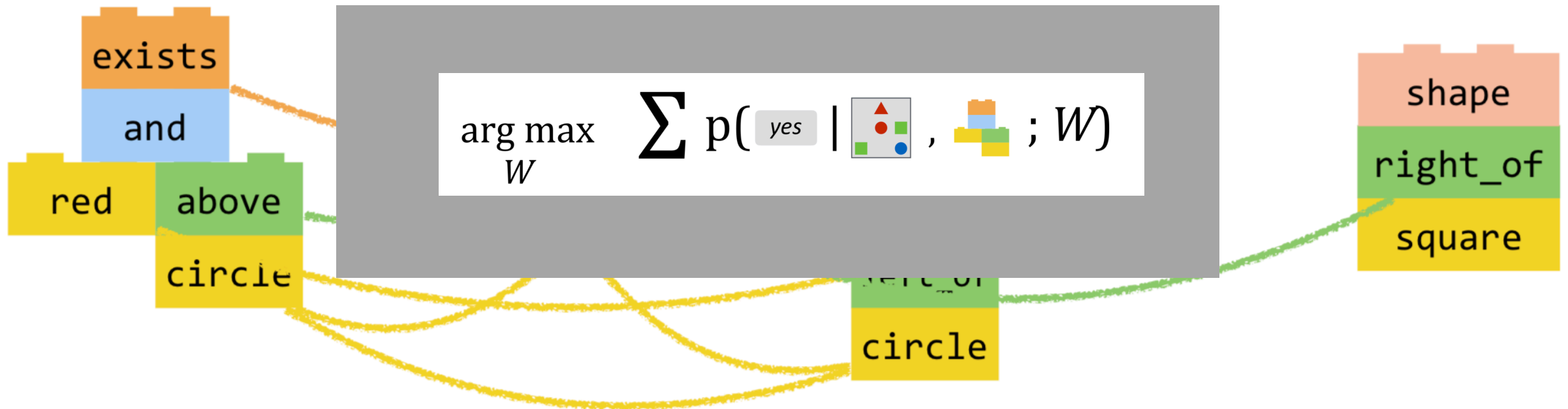
Training NMNs

18

- ▶ Learning parameters for individual module
- ▶ Supervised learning problem to maximize probability of output

Yes

Blue



VQA Dataset

19

*What color
is the necktie?*

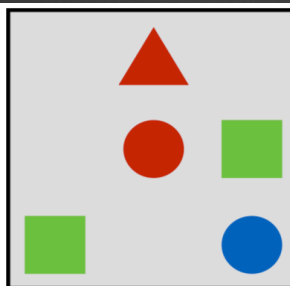


yellow

[Antol et al. 2015]

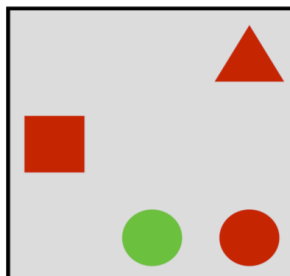
SHAPES Dataset - Novel

*Is there a red
shape above
a circle?*



yes

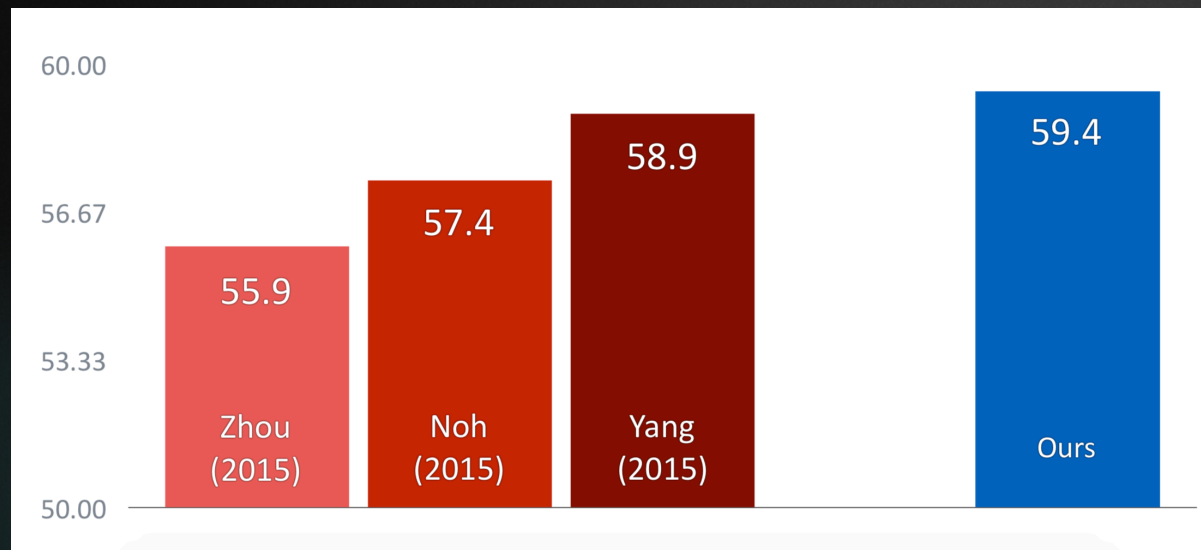
*Is a green shape
above left of a
red shape?*



no

Experiment Results - VQA

20



	test-dev				test
	Yes/No	Number	Other	All	All
LSTM [2]	78.20	35.7	26.6	48.8	–
VIS+LSTM [2]	78.9	35.2	36.4	53.7	54.1
NMN	69.38	30.7	22.7	42.7	–
NMN+LSTM	77.7	37.2	39.3	54.8	55.1

*What color is
she wearing?*



color

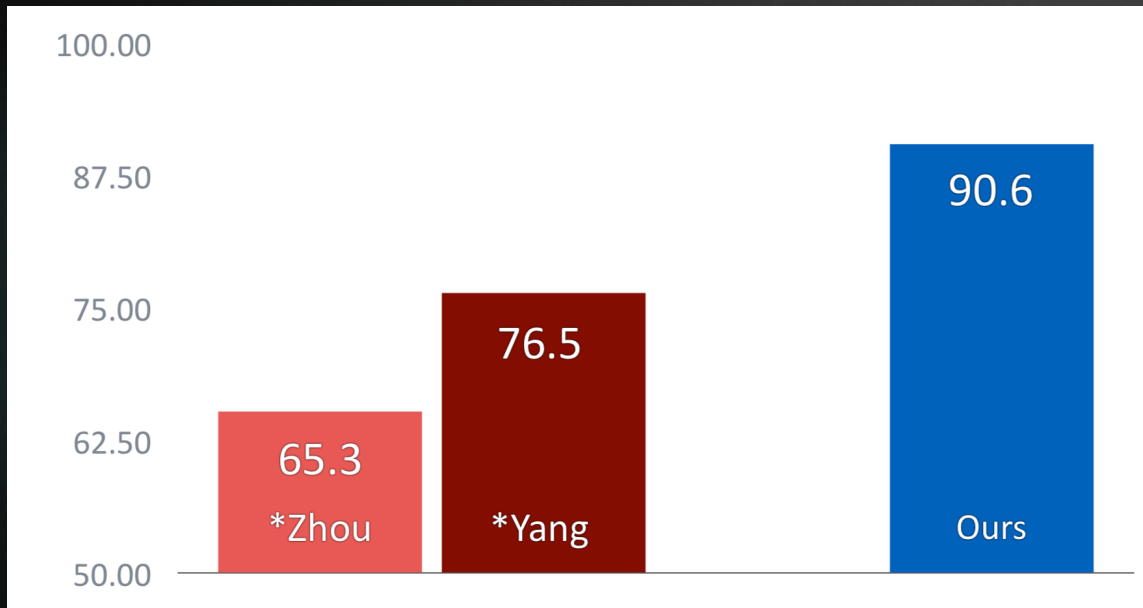
wear

white

Experiment Results - SHAPES

21



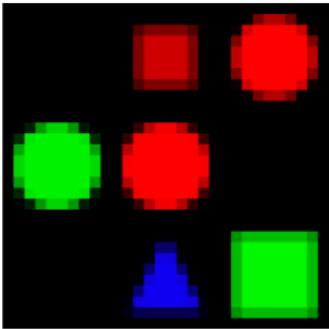
- ▶ Size: number of modules needed to instantiate NMN
- ▶ NMN (easy): modified training set with no size-6 questions



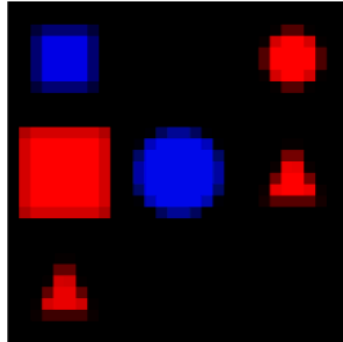


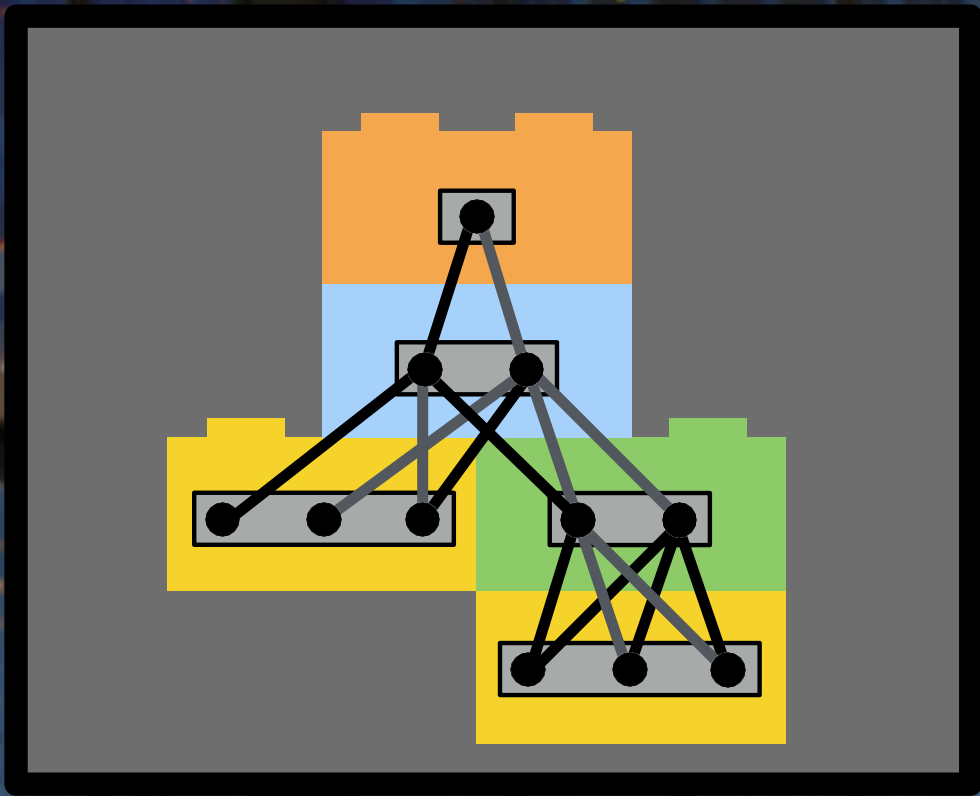
	size 4	size 5	size 6	All
Majority	64.4	62.5	61.7	63.0
VIS+LSTM	71.9	62.5	61.7	65.3
NMN	89.7	92.4	85.2	90.6
NMN (easy)	97.7	91.1	89.7	90.8

Discussion

22

		
<i>what color is the vase?</i>	<i>is the bus full of passengers?</i>	<i>is there a red shape above a circle?</i>
<code>classify[color](attend[vase])</code>	<code>measure[is](combine[and](attend[bus], attend[full])</code>	<code>measure[is](combine[and](attend[red], re-attend[above](attend[circle])))</code>
green (green)	yes (yes)	no (no)

		
<i>what material are the boxes made of?</i>	<i>is this a clock?</i>	<i>is a red shape blue?</i>
<code>classify[material](attend[box])</code>	<code>measure[is](attend[clock])</code>	<code>measure[is](combine[and](attend[red], attend[blue]))</code>
leather (cardboard)	yes (no)	yes (no)

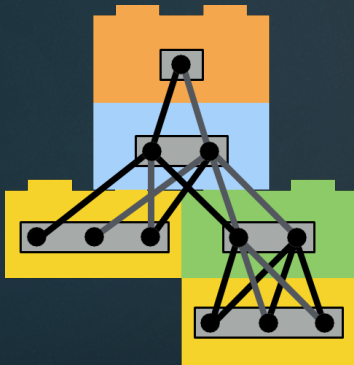


A Closer Look...

24

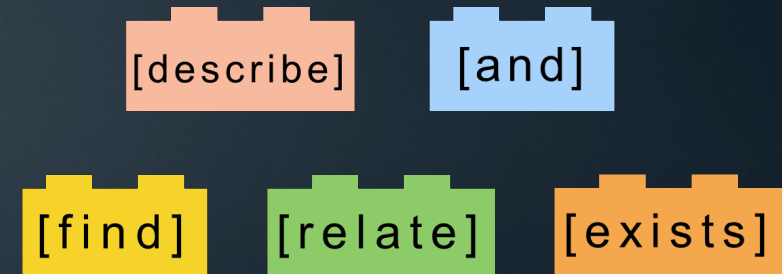
Pros:

- ▶ Combines advantages of:
 - ▶ Representation learning (like a neural net)
 - ▶ Compositionality (like a semantic parser)



Cons:

- ▶ Limited to produce arbitrary forms

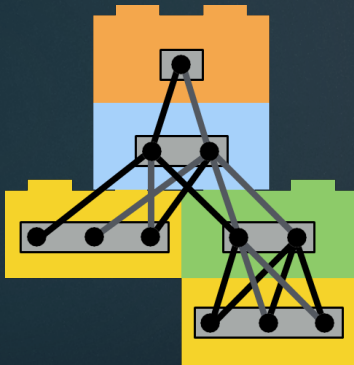


A Closer Look...

25

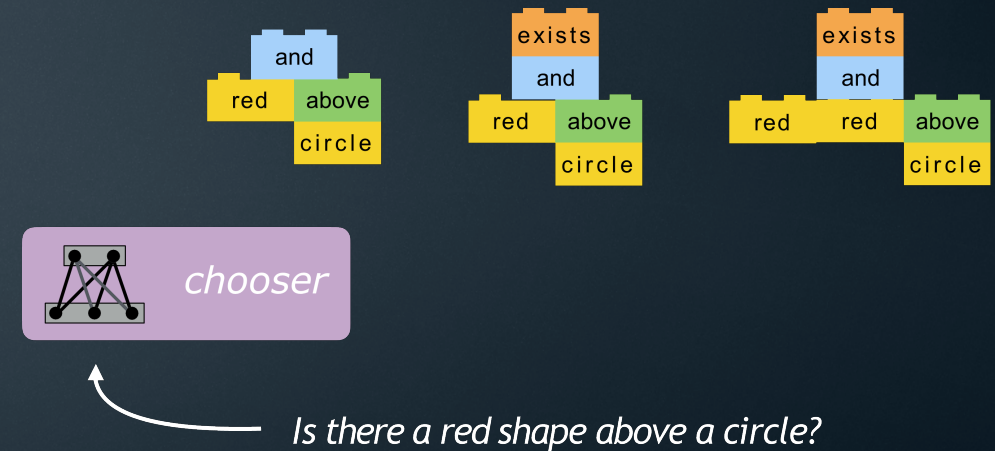
Pros:

- ▶ Combines advantages of:
 - ▶ Representation learning (like a neural net)
 - ▶ Compositionality (like a semantic parser)



Cons:

- ▶ Limited to produce arbitrary forms
- ▶ Structure selection

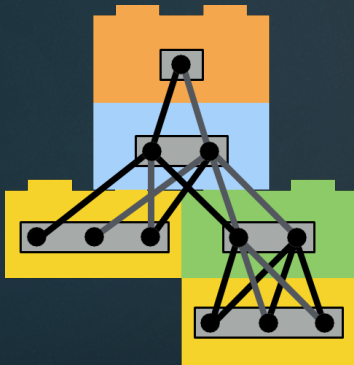


A Closer Look...

26

Pros:

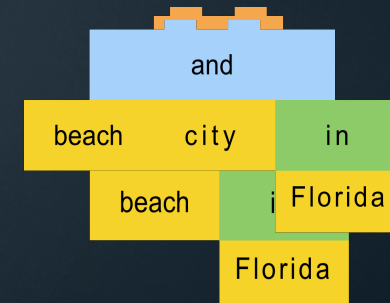
- ▶ Combines advantages of:
 - ▶ Representation learning (like a neural net)
 - ▶ Compositionality (like a semantic parser)



Cons:

- ▶ Limited to produce arbitrary forms
- ▶ Structure selection
- ▶ Potential error in generated model

What beach city is
there in Florida?



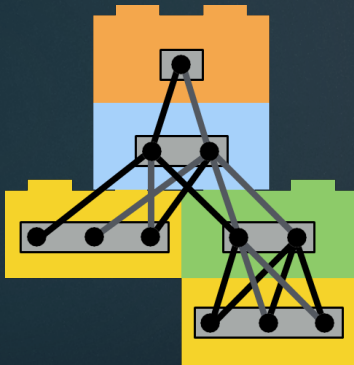
}
(wrong module
behavior)

A Closer Look...

27

Pros:

- ▶ Combines advantages of:
 - ▶ Representation learning (like a neural net)
 - ▶ Compositionality (like a semantic parser)



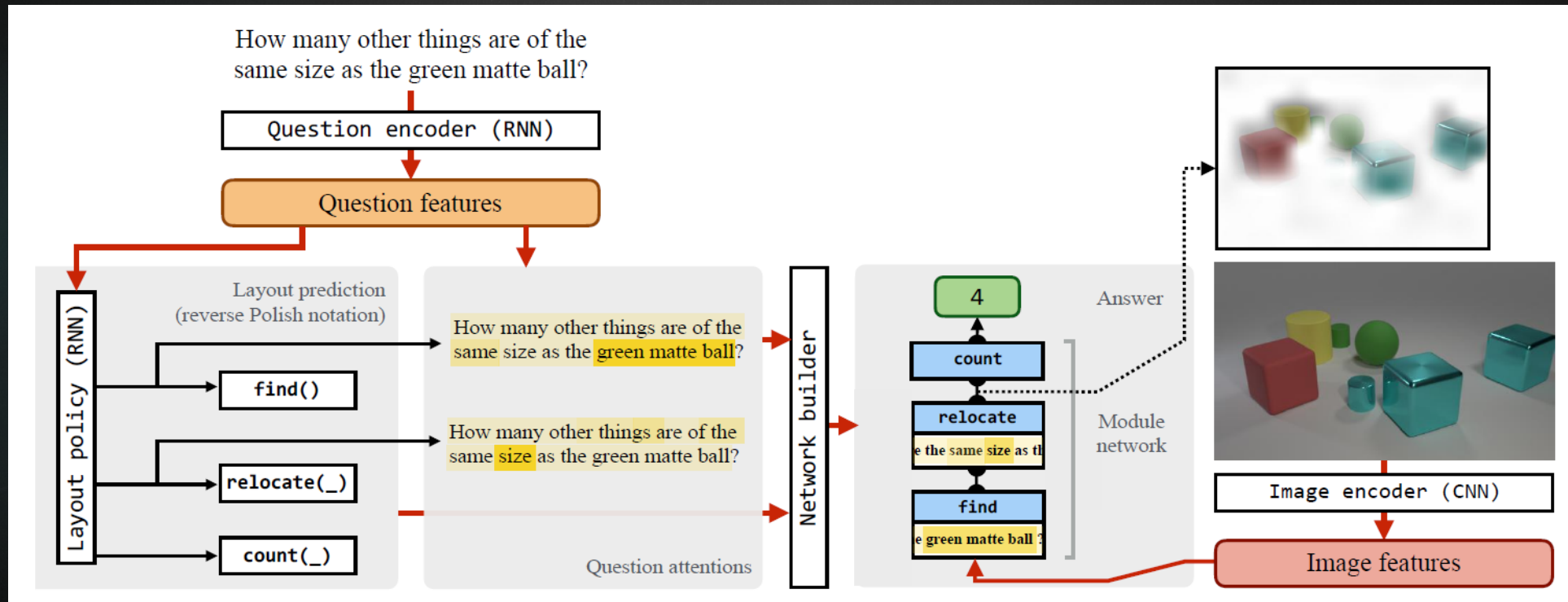
Cons:

- ▶ Limited to produce arbitrary forms
- ▶ Structure selection
- ▶ Potential error in generated model
- ▶ End-2-End structure learning

Extension of the Work

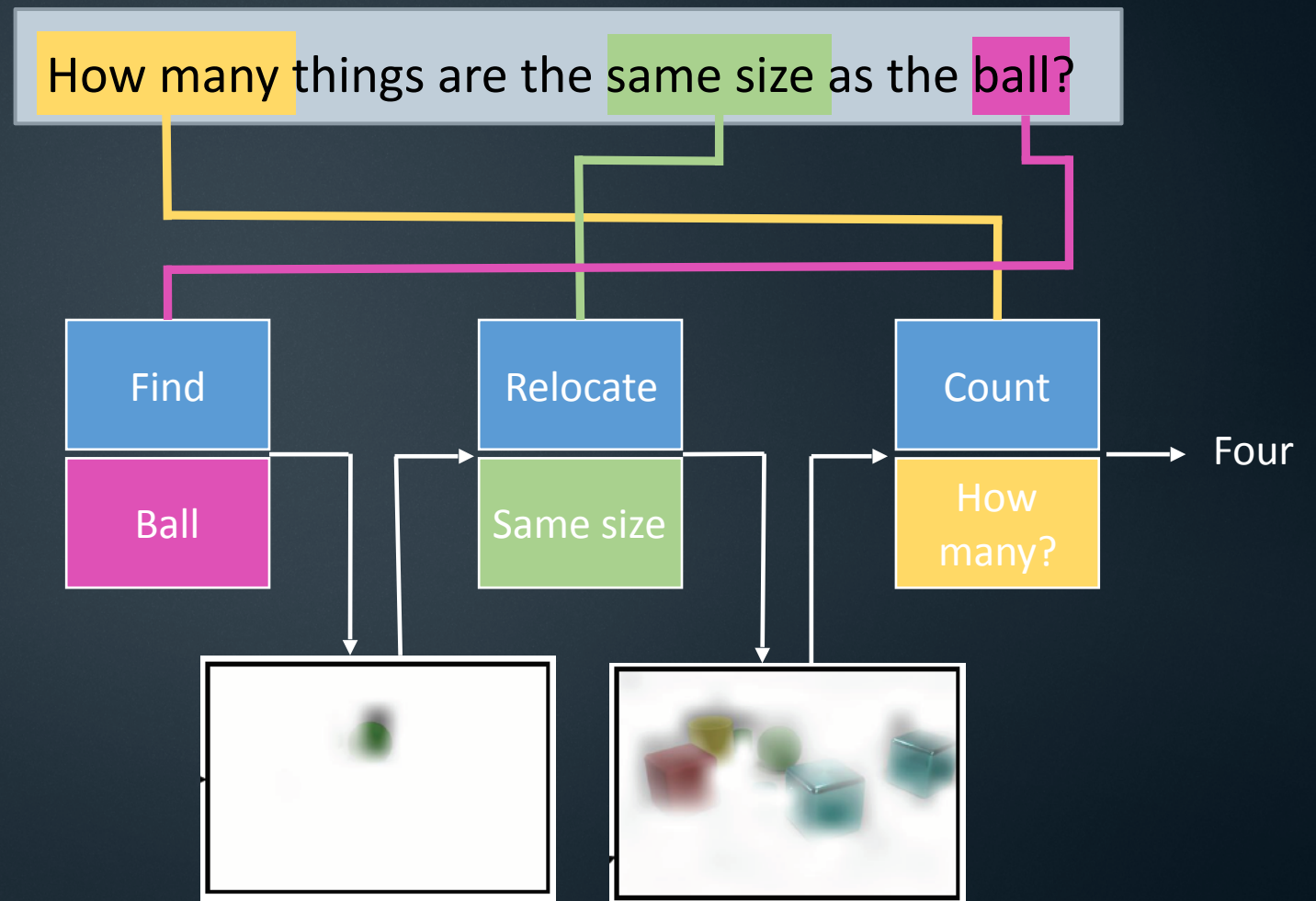
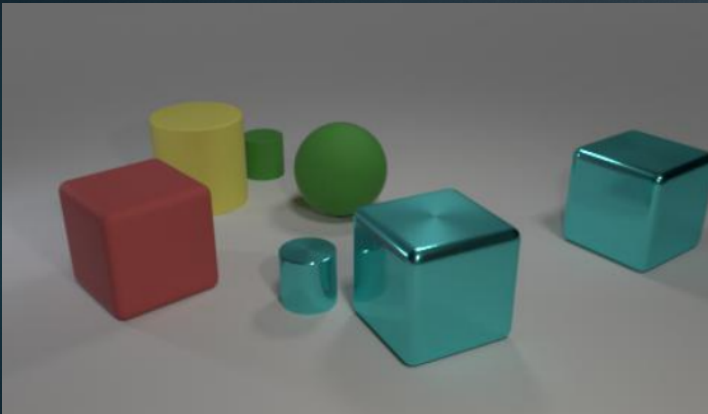
28

- Learning to Reason: End-to-End Module Networks for Visual Question Answering



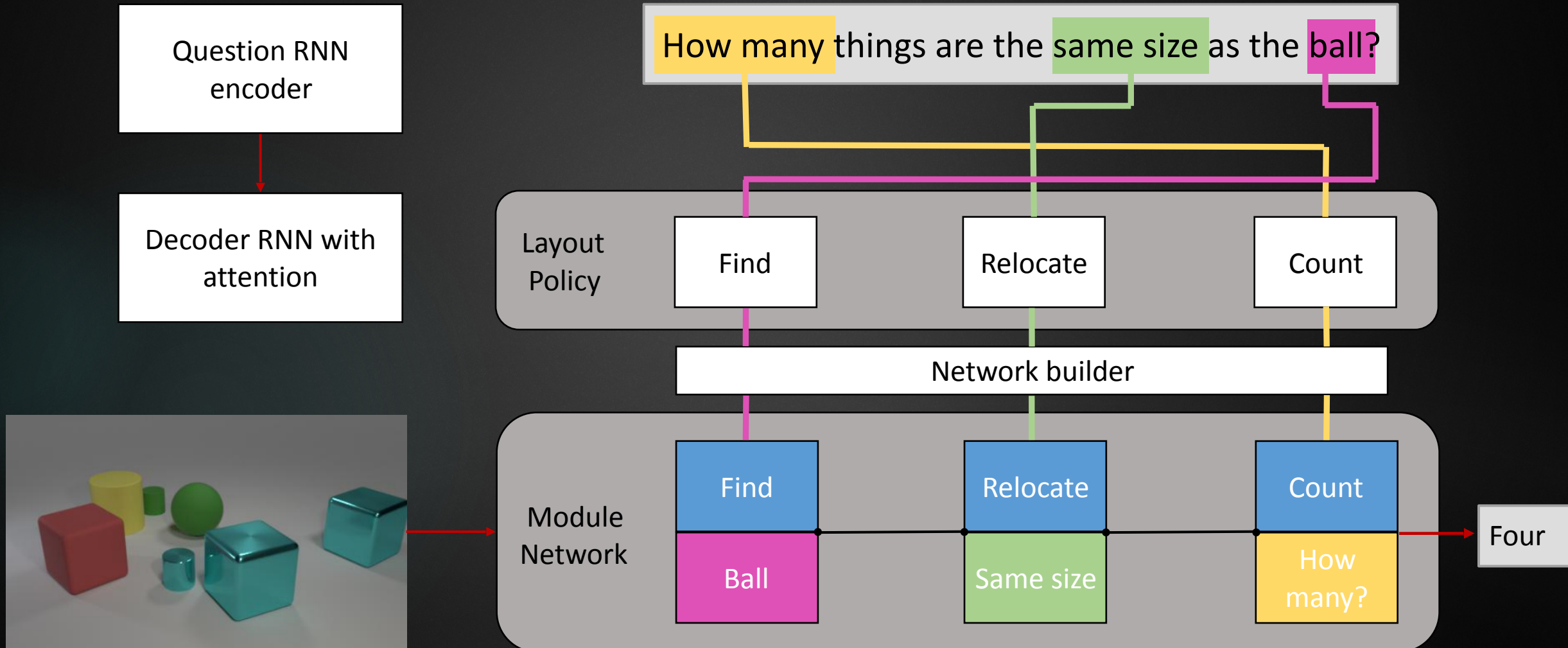
Extension of the Work

29



Extension of the Work

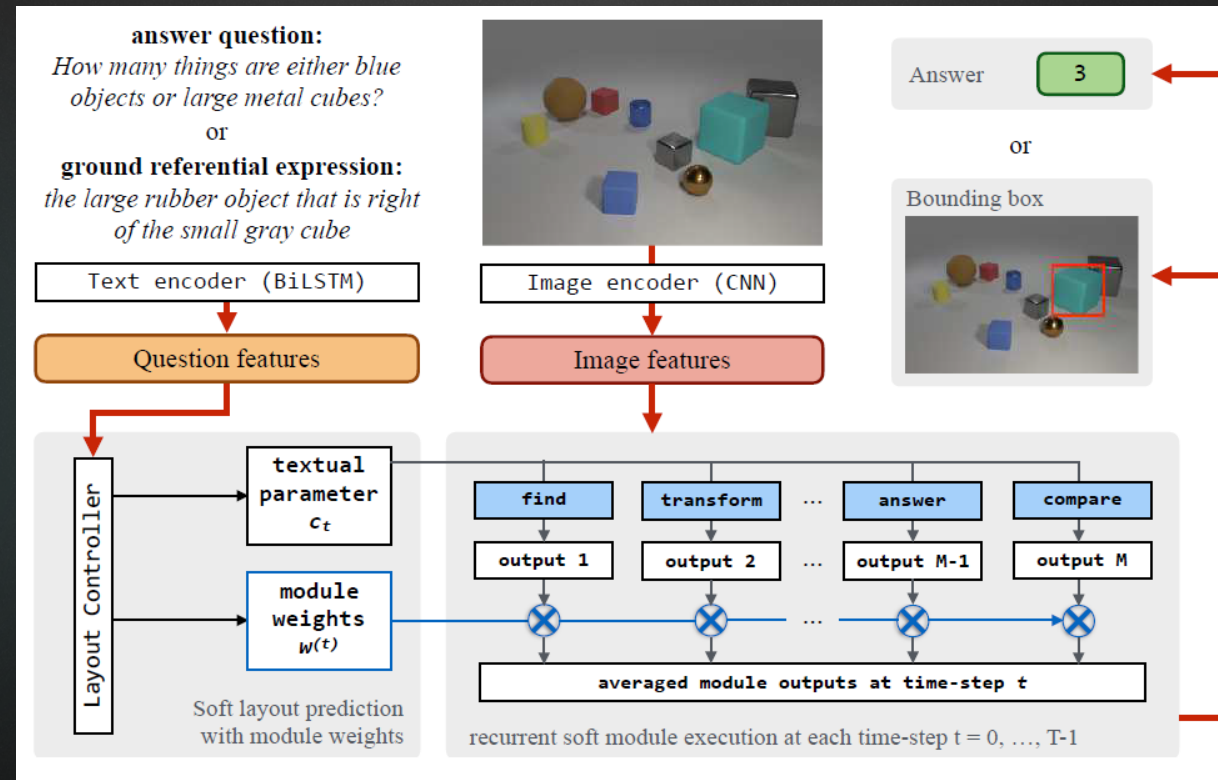
30



Extension of the Work

31

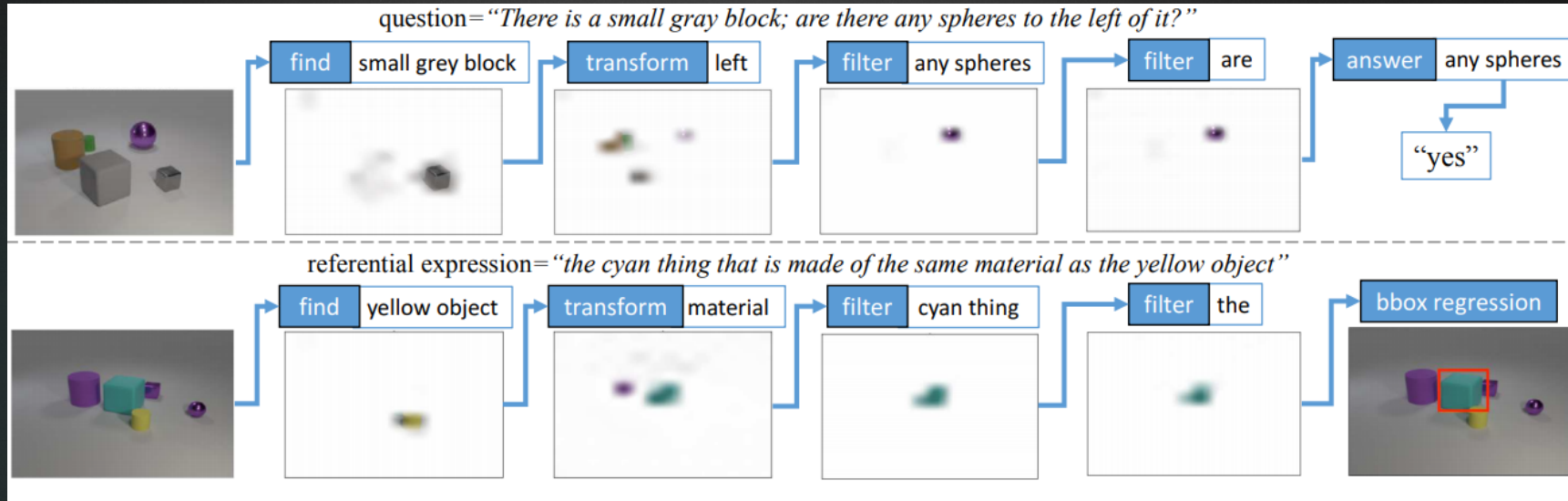
- Explainable Neural Computation via Stack Neural Module Networks



Extension of the Work

32

- Answer visual questions
- Ground referential expressions



- ▶ Andreas J, Rohrbach M, Darrell T, Klein D. Neural module networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition 2016 (pp. 39-48).
- ▶ Andreas J, Rohrbach M, Darrell T, Klein D. Learning to compose neural networks for question answering. 2016. PowerPoint Presentation.
- ▶ Chen L, Tang S. Visual Question Answering. 2016. PowerPoint Presentation.
- ▶ Wu Q, Teney D, Wang P, Shen C, Dick A, van den Hengel A. Visual question answering: A survey of methods and datasets. Computer Vision and Image Understanding. 2017 Oct 1;163:21-40.
- ▶ Hu R, Andreas J, Darrell T, Saenko K. Explainable neural computation via stack neural module networks. In Proceedings of the European Conference on Computer Vision (ECCV) 2018 (pp. 53-69).



Thank You
and Happy Spring