



JOHNS HOPKINS
UNIVERSITY

Towards Generalizable Visual Reasoning

Zhuowan Li

The Department of Computer Science

Johns Hopkins University

Advisors: Alan Yuille, Benjamin Van Durme

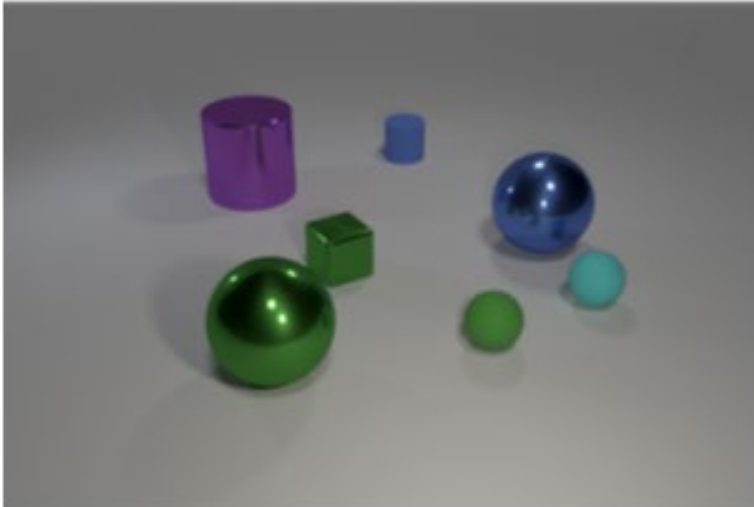
Humans reason with vision and language



Leaves turn
yellow in autumn.



Visual Question Answering is a challenging task



Q: The cylinder that is the same size as the metallic sphere is what color?

A: Purple

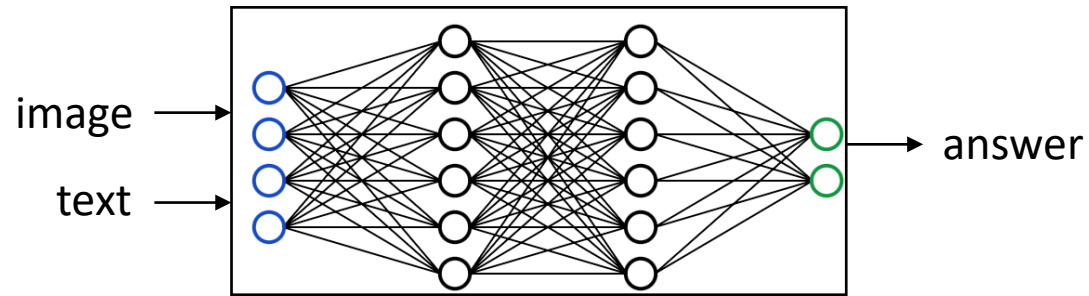


Q: What does the little boy in front of the table hold?

A: Toothbrush

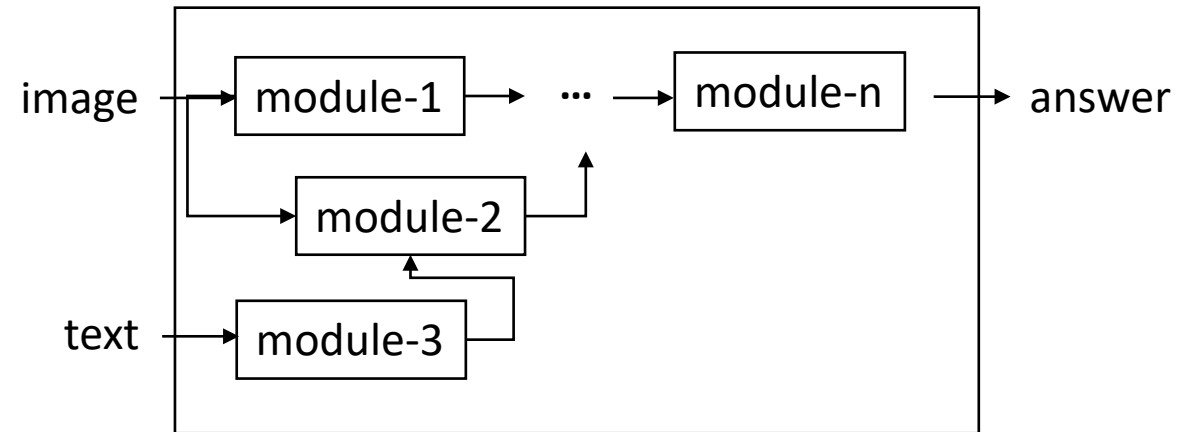
A paradox in VQA

Standard end-to-end models



- Perform well on IID setting
- Not robust to distribution shifts

Neural modular methods



- More robust
- Not perform well on standard real datasets

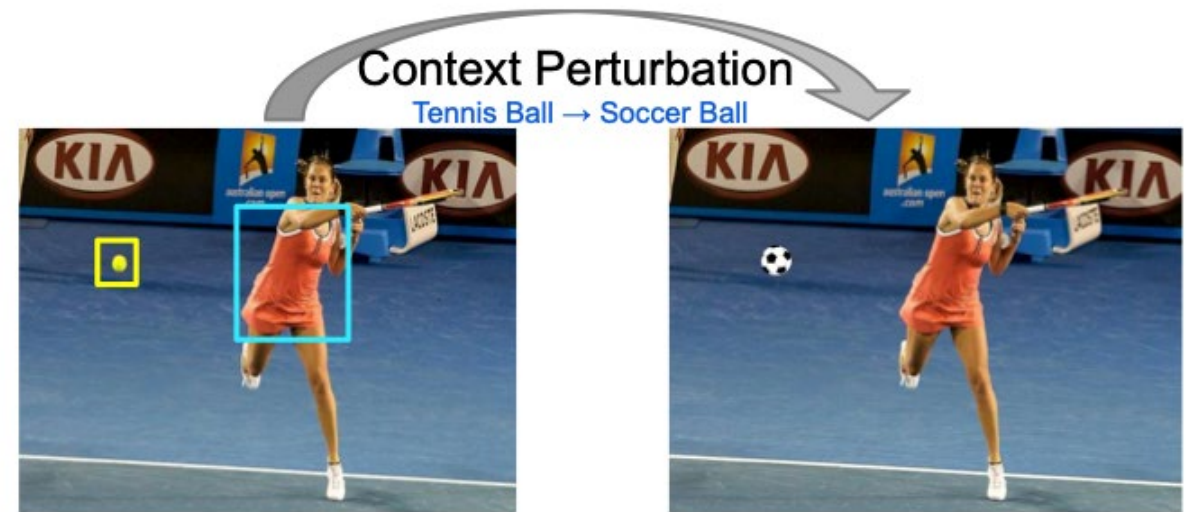
We want to understand this better

Standard VQA models are not robust

Text shortcuts

- “what color..” → white
- “is there..” → yes
- “how many..” → 2
- “What color is the banana” → yellow

Visual contexts



Question: What color is the woman’s dress?

Ground-truth answer: Orange.

Model prediction: Orange. ✓

Model prediction: White. ✗

Swapmix: Diagnosing and regularizing the over-reliance on visual context in visual question answering.

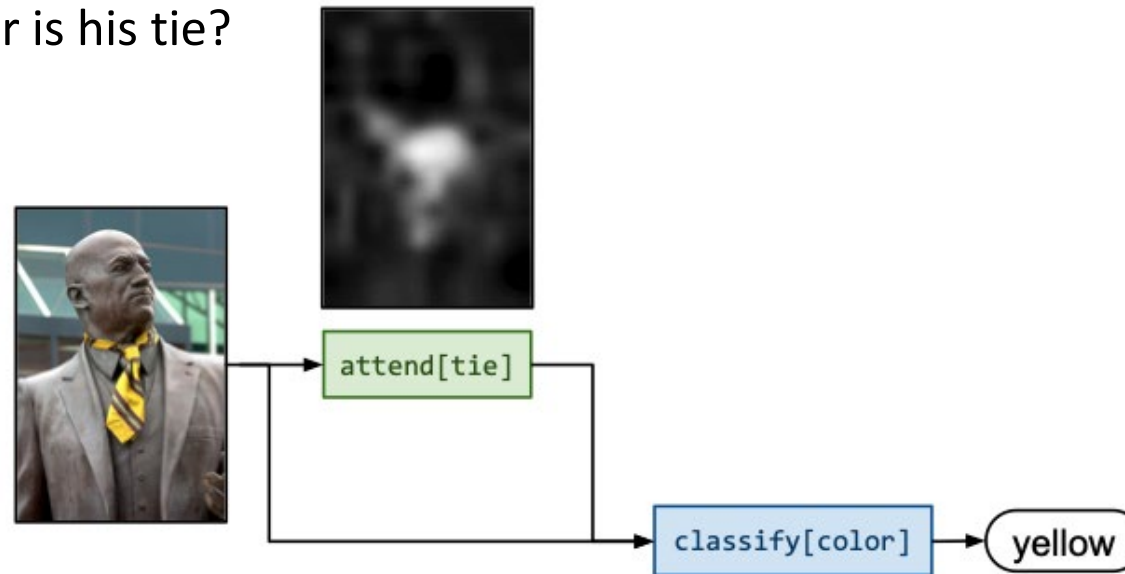
Vipul Gupta, **Zhuowan Li**, Adam Kortylewski, Chenyu Zhang, Yingwei Li, Alan Yuille.

In CVPR 2021.

An alternative: Neural Modular Methods

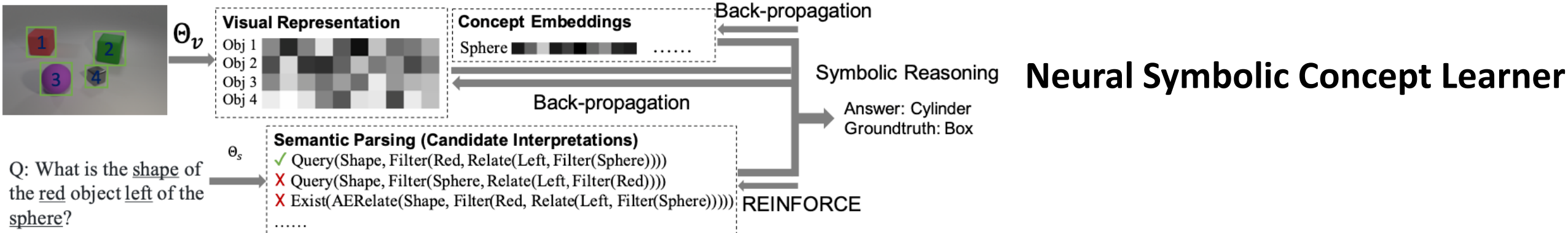
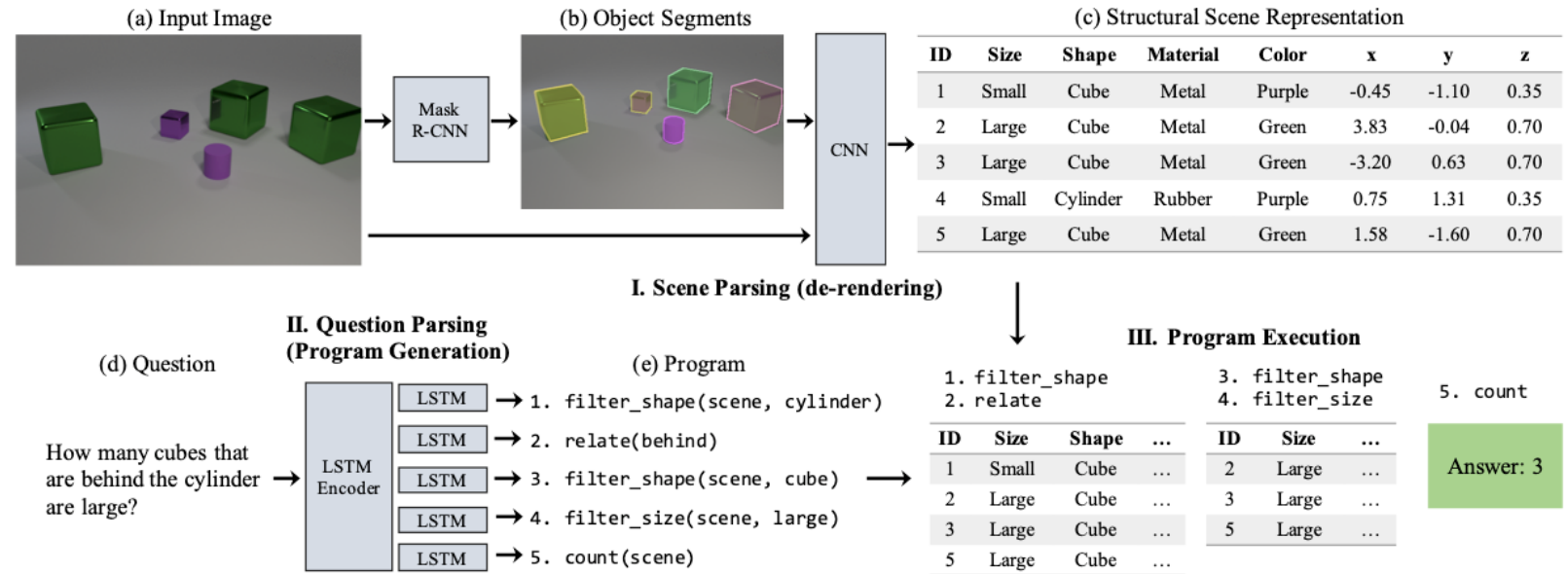
- Parse a question into a series of operations
- Each operation is implemented as a separate neural module

What color is his tie?



An alternative: Neural Modular Methods

Neural Symbolic VQA



Yi, Kexin, et al. "Neural-symbolic vqa: Disentangling reasoning from vision and language understanding." *NeurIPS* 2018.

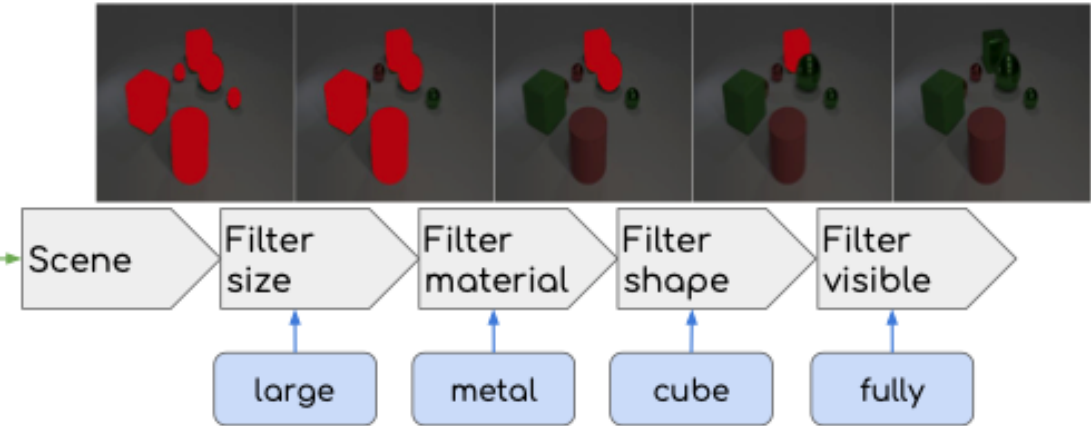
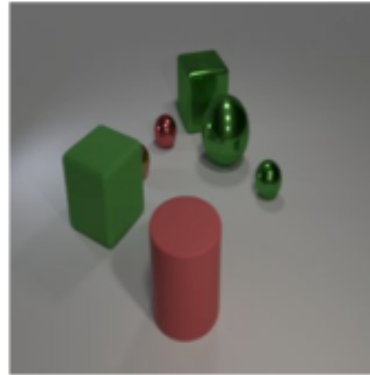
Mao, Jiayuan, et al. "The Neuro-Symbolic Concept Learner: Interpreting Scenes, Words, and Sentences From Natural Supervision." *ICLR* 2019.

Neural modular methods

Pros:

- Interpretable

CLEVR-Ref+ (by CCVL)
CVPR 2019



The fully visible big shiny block(s)

- Data-efficient
- Robust
- SoTA performance on CLEVR
- ...

Neural modular methods

Pros:

- Interpretable
- Data-efficient
- Robust (Need to be verified)
- SoTA performance on CLEVR
- ...

Cons:

- Need explicit reasoning programs (partly addressed by NSCL)
- Low performance on real images

Neural modular methods

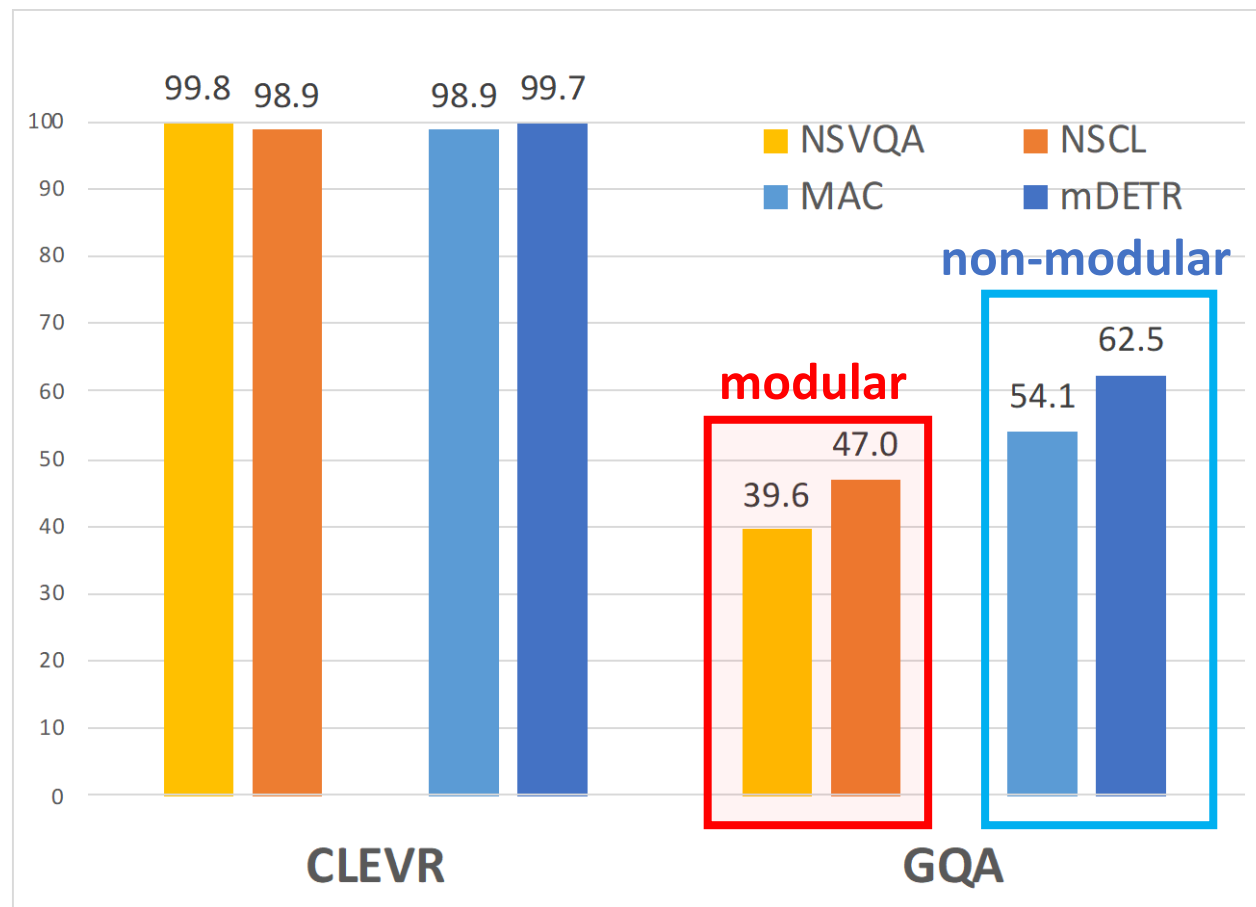
Pros:

- Interpretable
- Data-efficient
- Robust (Need to be verified)
- SoTA performance on CLEVR
- ...

Cons:

- Need explicit reasoning programs (partly addressed by NSCL)
- **Low performance on real images**

Modular methods suffer on real data



Outline

- Why do symbolic methods suffer on real images? How to improve them? [ICCV 2021]
- Super-CLEVR: How to study domain robustness more systematically? [CVPR 2023 Highlight]
- Extension of Super-CLEVR with Part, Pose, Occlusion [Ongoing]

Why do Symbolic Methods suffer on Real Images? How to improve them?

Calibrating Concepts and Operations: Towards Symbolic Reasoning on Real Images.

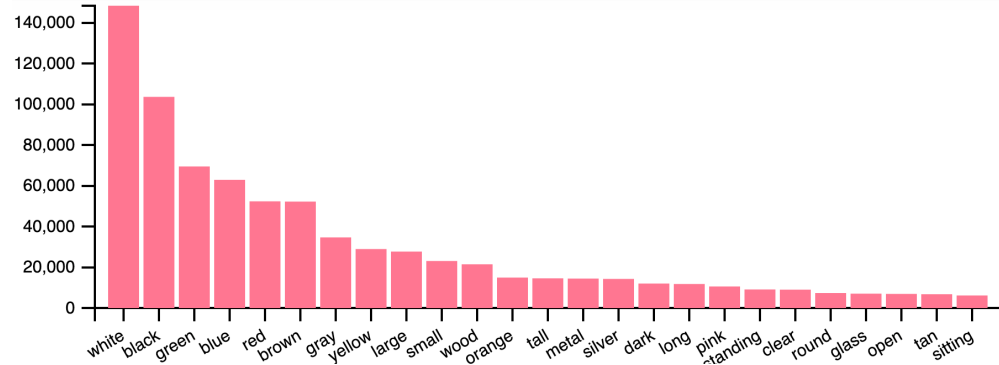
Zhuowan Li, Elias Stengel-Eskin, Yixiao Zhang, Cihang Xie, Quan Hung Tran, Benjamin Van Durme, Alan Yuille

In ICCV 2021.

Overview: two reasons, and solutions

Reason-1:

Long-tail distribution

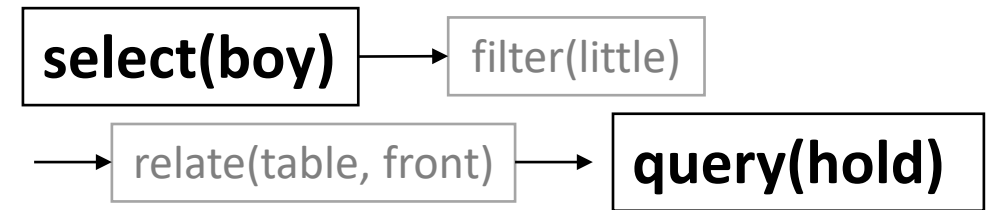


Solution-1:

Calibrating concepts

Reason-2:

Unequal importance of reasoning steps

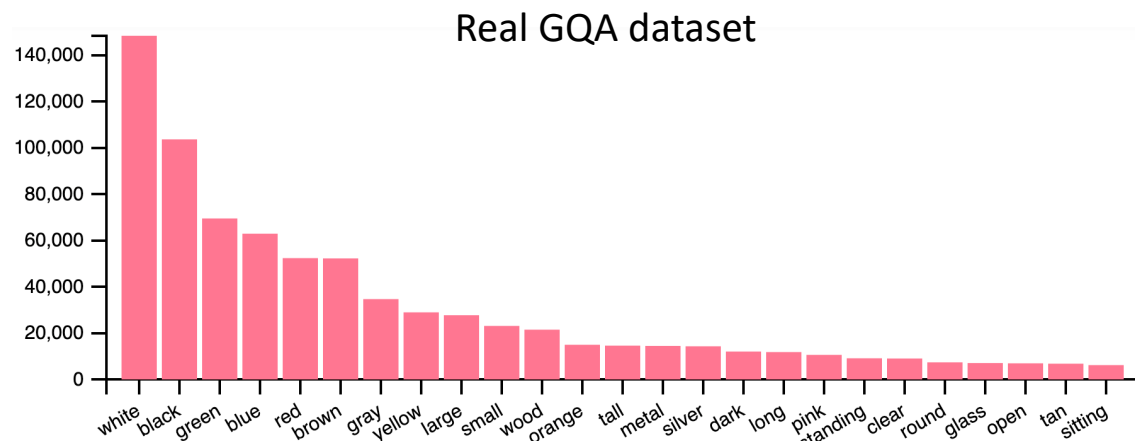
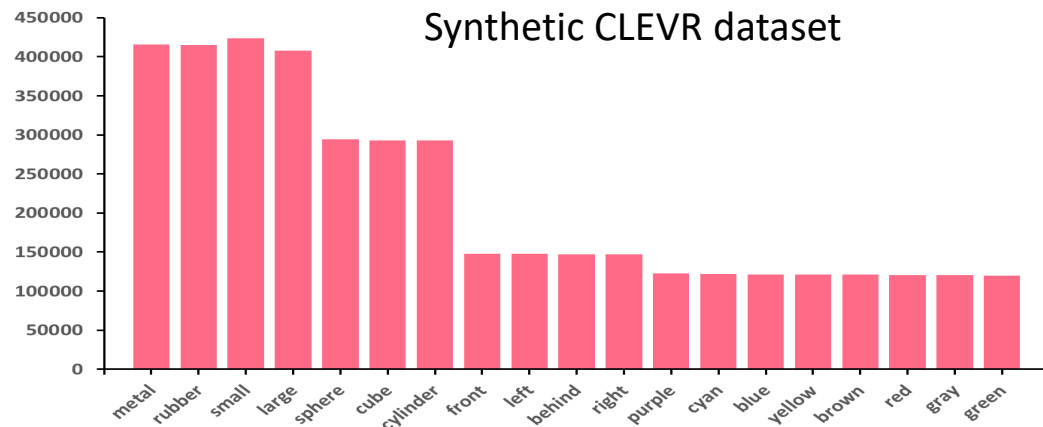


What does the little boy in front of the table hold?

Solution-2:

Calibrating operations

1. Real data suffers from long-tailed distribution

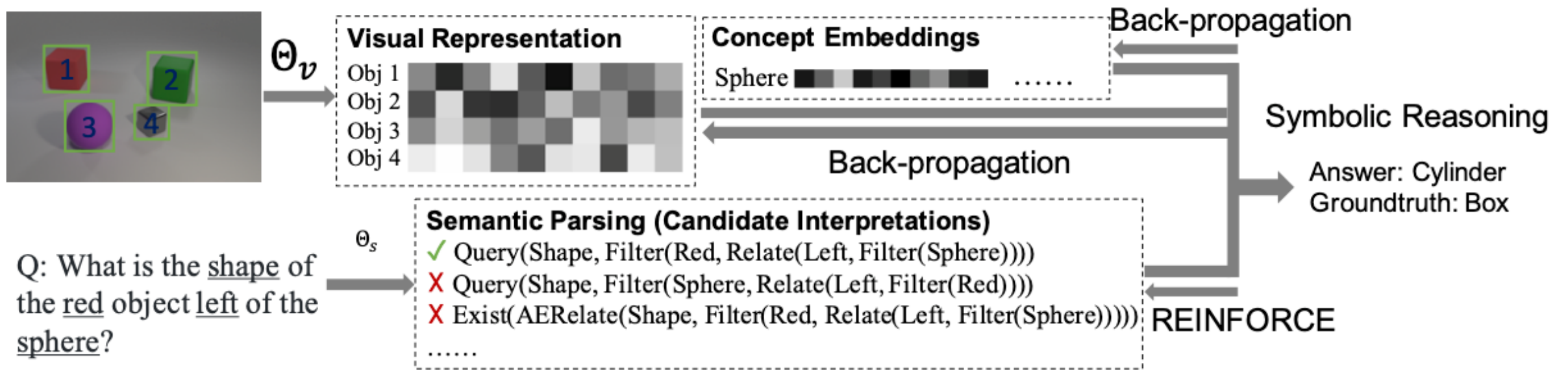


white, black, green, ...

tan, sitting, open, ...

How do long tails affect concept learning?

Recall NSCL:



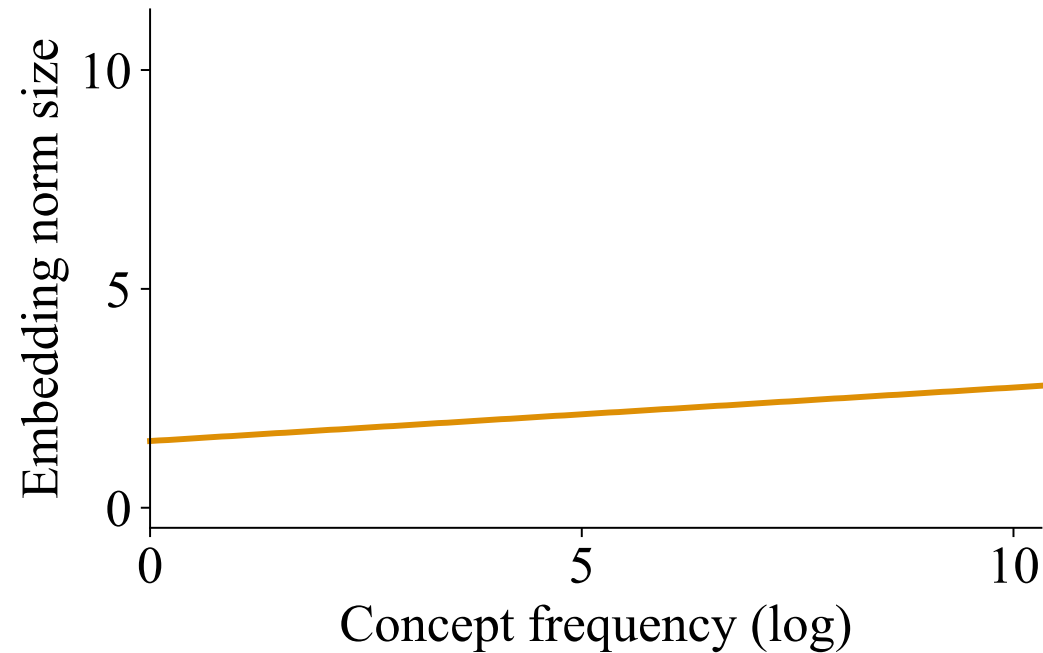
Normalized Concept embeddings hinder the learning the concept distributions

How does this distribution affect concept learning?

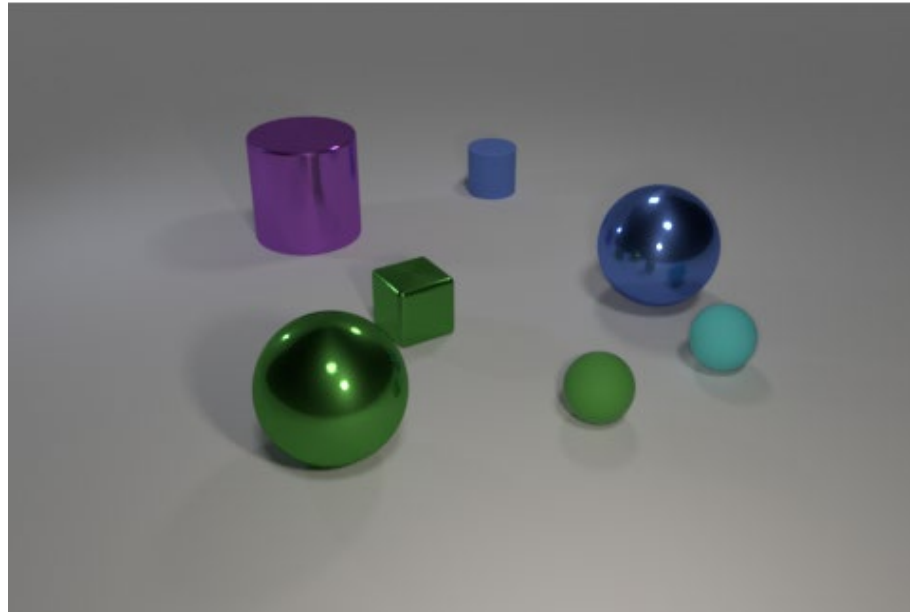
- By default, concept embeddings are normalized.
- Simply removing the normalization on concept embedding yields substantially better performance (+3.4%)
- **The magnitudes of the embeddings matter!**

Unnormalizing concept weights

Positive correlation between concept frequency and embedding norm size



2. Reasoning steps are of unequal importance



Q: The cylinder that is the same size as the metallic sphere is what color?

A: Purple

Prog: select(sphere) filter(→metallic) samesize() →
filter(cylinder) querycolor() →



Q: What does the little boy in front of the table hold?

A: Toothbrush

Prog: select(boy) filter(little) relate_s(table, front)
query_rel_o(hold)

→

→

→

Intuition: what if we adjusting operation weights?

Question: Is there a bag in this image that is not black?

Groundtruth: No



(1) Select(bag) scores:

$[-7.0, -6.0, 2.1, -9.9]$

(2) Filter(not black) scores:

$[0.8, -0.7, -1.7, 2.1]$

Merge: (1) + (2):

$[-6.2, -5.3, 0.4, -7.8]$

Exist?

\rightarrow ~~X~~ Answer: Yes

With weight: (1) + 2*(2)

$[-5.4, -4.6, -1.3, -5.7]$

Exist?

\rightarrow \checkmark Answer: No

Method: Calibrating Concepts and Operations

- Calibrating concepts
 - Explicitly learnable norm size for each concept:

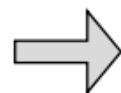
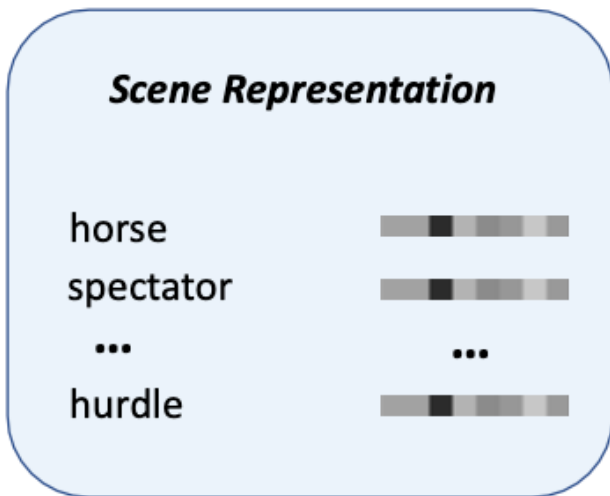
$$\mathbf{c}_{concept} = w_{concept}^{type} \mathbf{c}_{concept}$$

- Calibrating operations
 - Train a LSTM weight predictor to predict weights for each module using context
 - Merge operation results with learned weights:

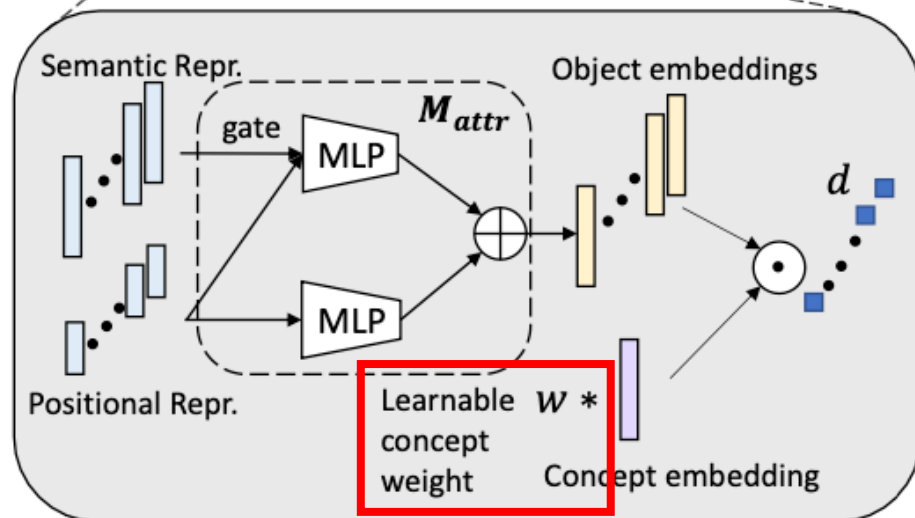
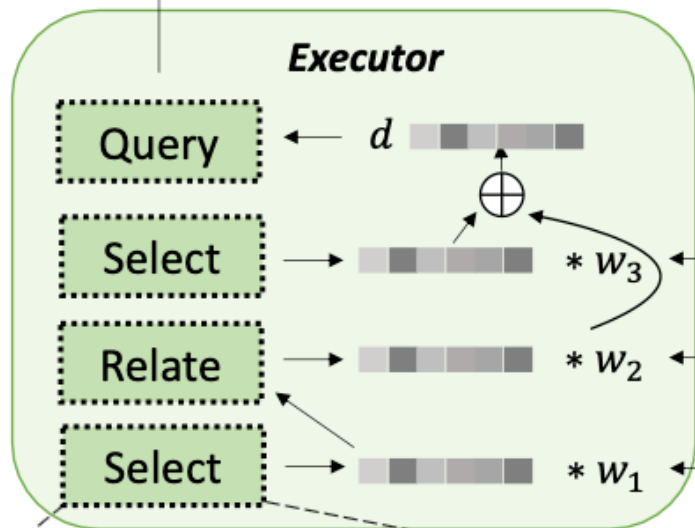
$$\mathbf{d}_i = \sum_{j \in \mathcal{D}(p_i)} w_j \mathbf{d}_j$$



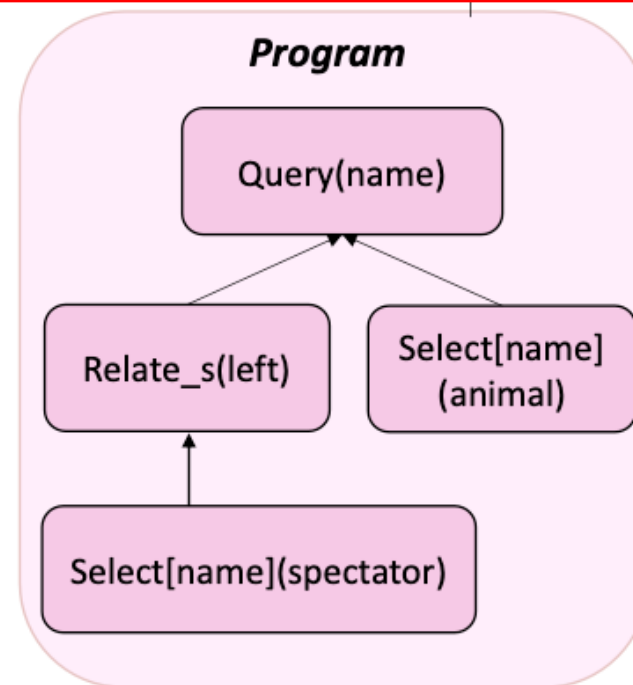
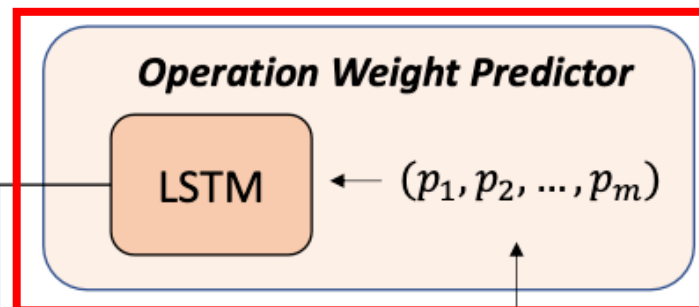
Faster-RCNN



Answer: horse



**Calibrating operations:
predict weights for each operation**



Question: How is the animal to the left of the spectators called?

**Calibrating Concepts:
Learn magnitudes for each concept in each module**

Both the concept and operation calibrations help!

Model Level	Concept	Operation	Accuracy
1 (Baseline)	Normalized	Average	47.01
2	Normalized	Calibrated	51.03
3	Unnormalized	Calibrated	54.65
4 (Ours)	Calibrated	Calibrated	56.13

Our method helps bridge the gap

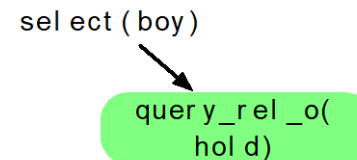
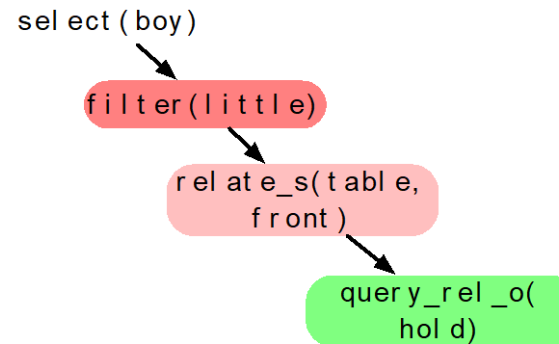
	Model	Accuracy
non-symbolic	LXMERT (Tan and Bansal, 2019)	60.33
	NSM (Hudson and Manning, 2019)	63.17
	MNM (Chen et al., 2021)	60.83
symbolic	∇ -FOL (Amizdeh et al., 2020)	54.76
	CCO (Ours, 2021)	56.38

Analysis of operation weights

- Prune low-weight modules progressively from the question
- The proposed perturb split can be used to analyze model behaviors

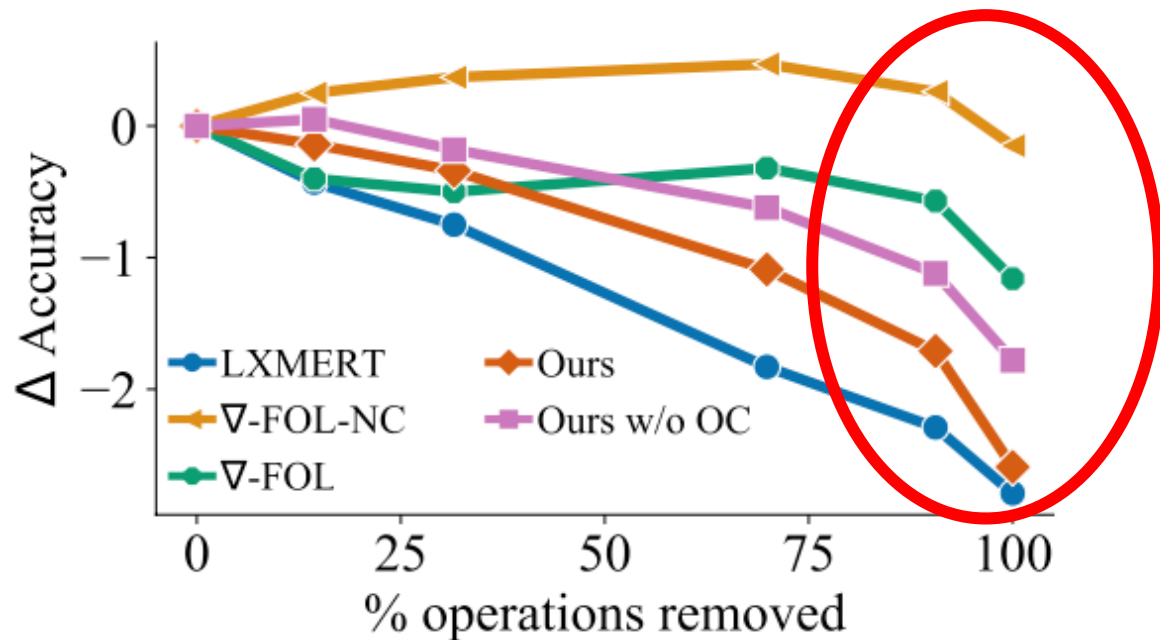


What does the little boy in front of the table hold?



What does the boy hold?

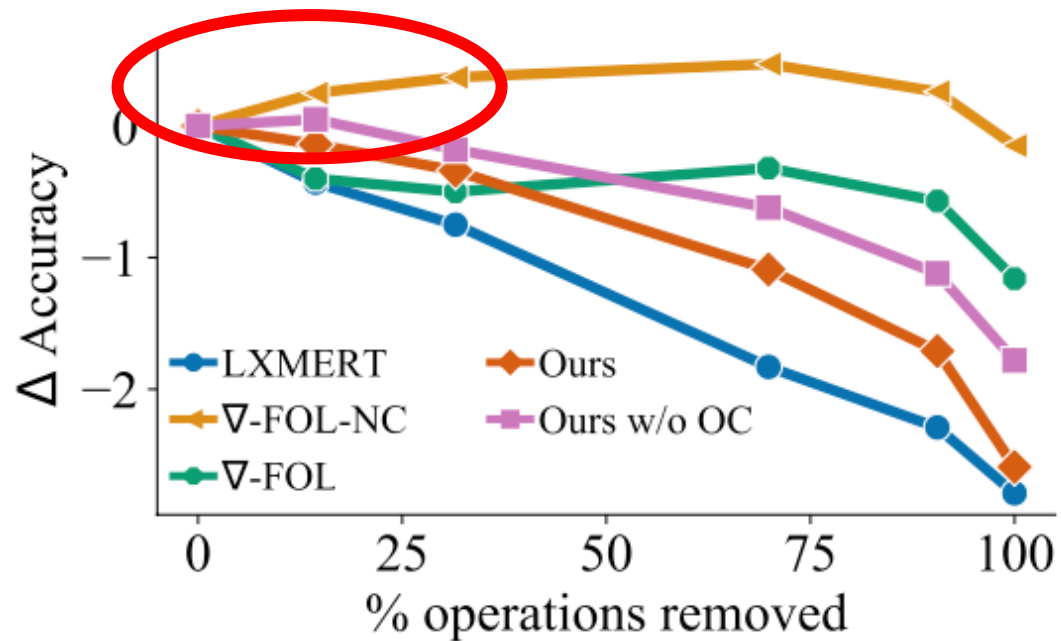
Analysis: Perturbed Test Split



The weight predictor assigns higher weights to more important operations.

Figure 5: Accuracy drop of different models when the testing questions are progressively perturbed by removing reasoning operations with low weights.

Analysis: Perturbed Test Split



Symbolic methods even increase performance when low-weighted operations are removed.

Figure 5: Accuracy drop of different models when the testing questions are progressively perturbed by removing reasoning operations with low weights.

Summary of Neuromodular on Real Data

- The performance of Neuromodular methods is improved by the two methods described above. But this does not solve the problem. Why can't Neuromodular or standard methods get 100% success?
- What else is going on?
- Conjectures:
 - (1) The standard end-to-trained methods can exploit the biases of the datasets, but the Neuromodular approaches are less effective at this.
 - (2) The Neuromodular methods use deep networks as their vision modules. They need better vision modules.
 - (3) The training confounds the vision and the language by training them together.
- How to study this? Better controlled datasets. Out-of-distribution testing.

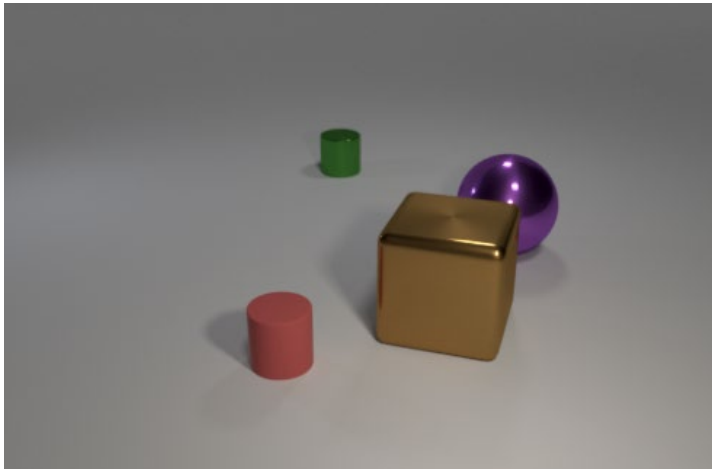
How to study the domain robustness more systematically?

Super-CLEVR: A Virtual Benchmark to Diagnose Domain Robustness in Visual Reasoning?

Zhuowan Li, Xingrui Wang, Elias Stengel-Eskin, Adam Kortylewski, Wufei Ma, Benjamin Van Durme, Alan Yuille.

CVPR 2023 Highlight.

Super-CLEVR dataset



CLEVR



Super-CLEVR

- Create the Super-CLRVR dataset using more realistic objects.
- This is more challenging
- It controls domain shifts to study robustness.

Super-CLEVR

- Super-CLEVR is more complex than CLEVR. It contains classes of objects – vehicles – from the ShapeNet repository. These are rendered to generate semi-realistic images.
- Super-CLEVR is controllable. We can systematically vary factors like the numbers of objects in the images, their poses in 3D, the occlusion, and so on.
- Super-CLEVR can be used to test VQA algorithms on data on which it has been trained. But it can also be used to test how VQA tests algorithms to generalize to out-of-distribution domains.
- Super-CLEVR can be extended to test VQA with adversarial examiners (but this has not been done).

Generalization: Decompose and analyze

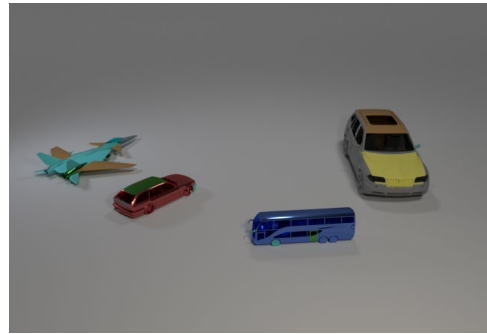
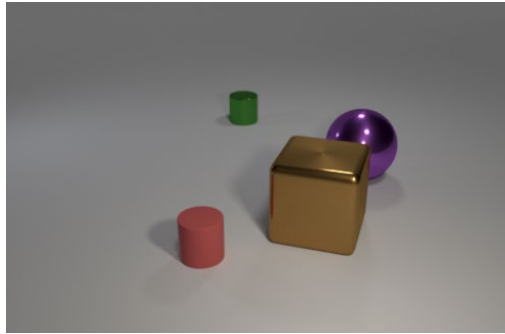
Four robustness factors in VQA domain shifts:

- visual complexity
- question redundancy
- concept distribution
- concept compositionality

Decompose VQA domain shifts into 4 factors

- **visual complexity**

how hard is the image



Easy



Hard

- question redundancy
- concept distribution
- concept compositionality

Decompose VQA domain shifts into 4 factors

- visual complexity
- **question redundancy**
the question may contain unnecessary information



What does the **little** boy ~~in front of the table~~ hold?



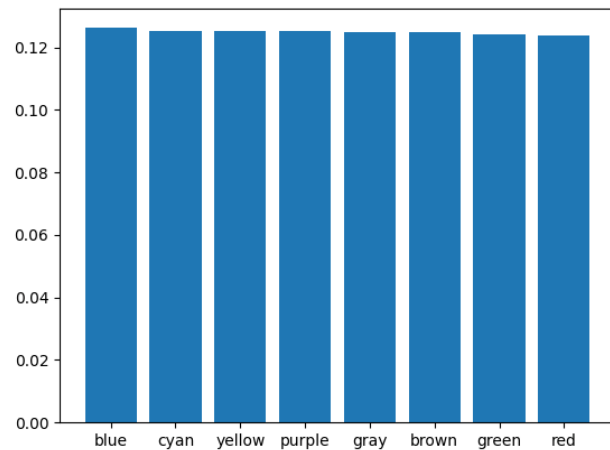
What is feeding the **large** animal ~~behind the fence~~?

- concept distribution
- concept compositionality

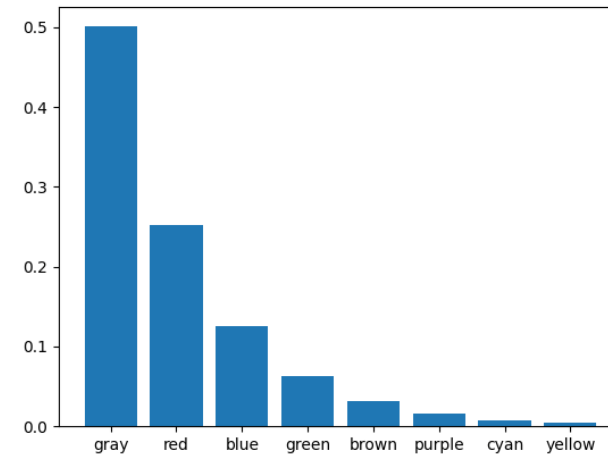
Decompose VQA domain shifts into 4 factors

- visual complexity
- question redundancy
- **concept distribution**

The distribution the concepts (objects names and attributes)



Well-balanced



Long-tail distributed

- concept compositionality

Domain A

Super-CLEVR



“What color is the bus?”

Domain B

Visual Complexity



easy



middle



hard

Question Redundancy

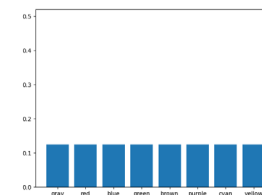
- redundancy
↓
standard
↓
+ redundancy

“What color is the bus?”

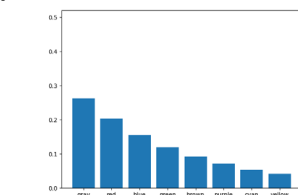
“What color is the large bus?”

“What color is the large bus behind the cyan car?”

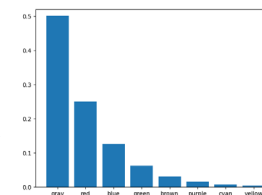
Concept Distribution



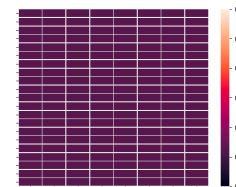
balanced



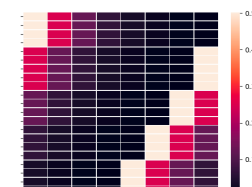
unbalanced



Concept Compositionality



well-composed



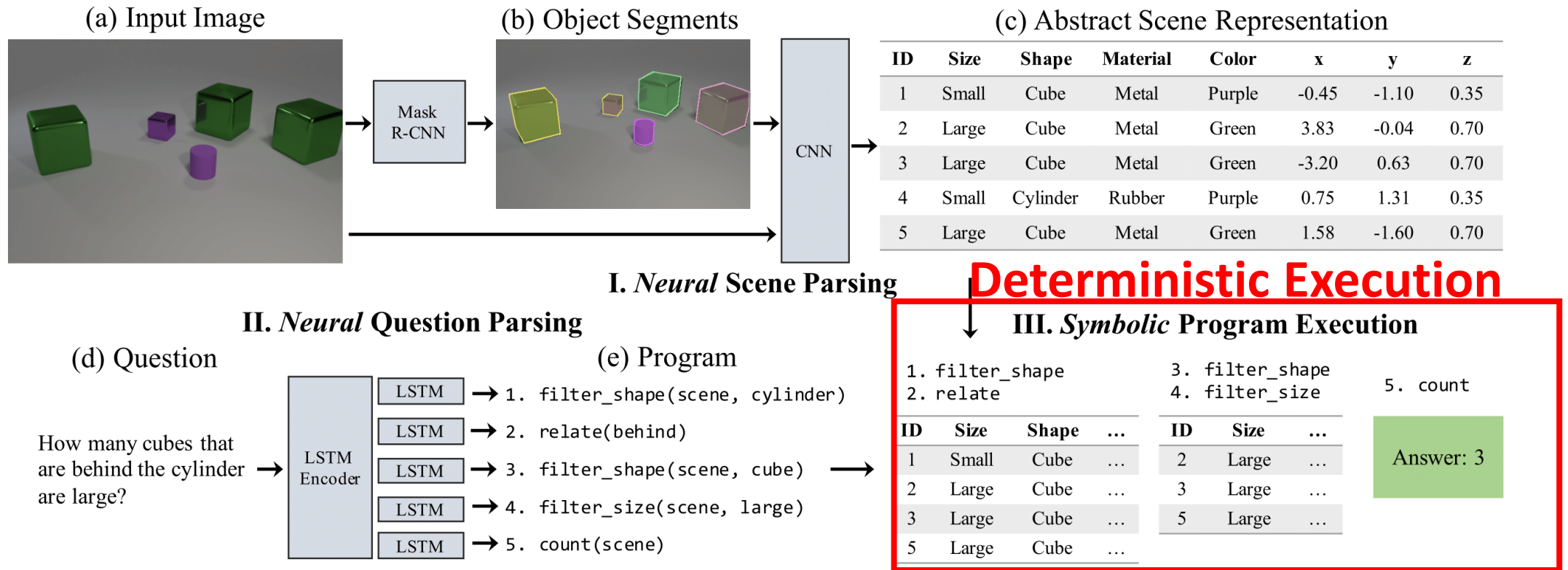
correlated

Five models are studied

- non-modular {
- FiLM
two-stream feature merging method
 - mDETR
pretrained transformer model

- modular {
- NSCL
neural modular method
 - NSVQA
neural modular method
 - **Probabilistic NSVQA**
our method

Recall Neural-Symbolic VQA (NSVQA)



Important – not emphasized in their paper – the training is modular. The language and the vision components are trained separately. No joint training.

Yi, Kexin, et al. "Neural-symbolic vqa: Disentangling reasoning from vision and language understanding." *NeurIPS* 2018.

Prob-NSVQA considers the confidence of the scene parser predictions.

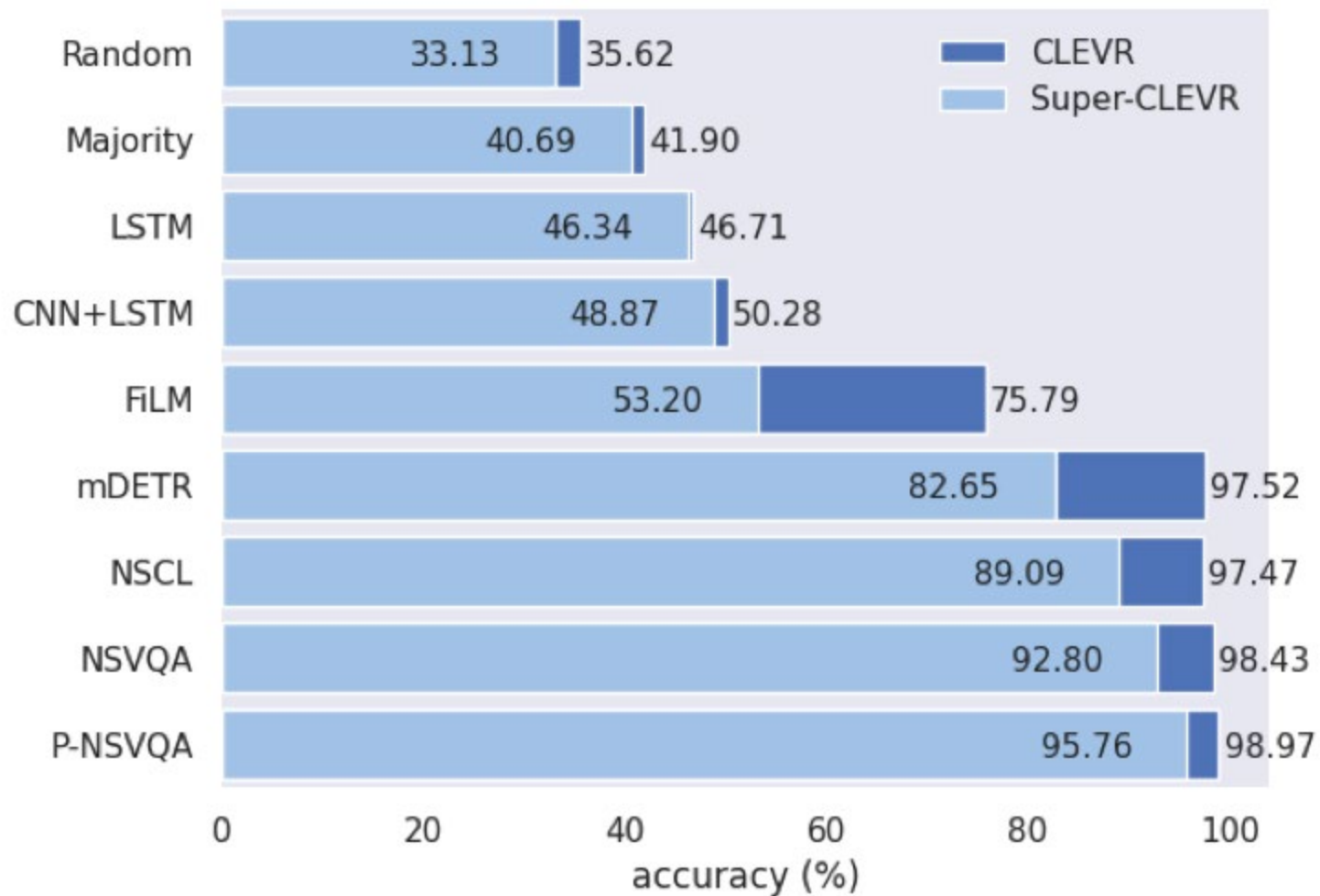
- We introduce probabilities to Prob-NSVQA. E.g., the vision modules output the probabilities that they have detected and classified an object.
- Given an image containing n objects, we maintain a vector of probs:

$$\mathbf{p} = [p^1, p^2, \dots, p^n]:$$

- *filter_{identifier}[attribute]* (e.g. *filter_{color}[red]*)

e.g. for the filter operation we set $p^k = p^k * P_{attribute}^k$

Super-CLEVR is harder than CLEVR (I.I.D. testing)



Out-of-domain testing: Complete Results

	FiLM			mDETR			NSCL			NSVQA			Prob NSVQA		
Visual Complexity															
	easy	mid	hard	easy	mid	hard	easy	mid	hard	easy	mid	hard	easy	mid	hard
easy	59.96	<u>53.95</u>	<u>50.66</u>	93.36	<u>84.30</u>	<u>82.97</u>	95.13	<u>92.31</u>	<u>90.81</u>	95.19	<u>94.19</u>	<u>94.09</u>	96.76	<u>95.98</u>	<u>96.37</u>
mid	57.41	<u>53.28</u>	<u>50.18</u>	83.34	<u>82.36</u>	<u>81.27</u>	<u>84.5</u>	89.10	<u>86.33</u>	<u>81.99</u>	<u>92.80</u>	93.78	<u>86.25</u>	95.76	<u>95.11</u>
hard	55.95	<u>53.11</u>	<u>50.47</u>	<u>79.71</u>	<u>79.94</u>	80.71	<u>76.85</u>	<u>78.66</u>	85.08	<u>73.11</u>	<u>79.71</u>	92.65	<u>79.81</u>	<u>86.47</u>	95.36
Question Redundancy															
	rd-	rd	rd+	rd-	rd	rd+	rd-	rd	rd+	rd-	rd	rd+	rd-	rd	rd+
rd-	<u>51.42</u>	<u>52.54</u>	53.51	83.94	<u>80.37</u>	<u>66.28</u>	<u>88.64</u>	<u>88.82</u>	90.33	-	<u>92.98</u>	-	-	<u>95.71</u>	-
rd	<u>50.39</u>	<u>53.28</u>	54.78	82.77	<u>82.36</u>	<u>70.36</u>	<u>88.45</u>	<u>89.10</u>	91.45	-	<u>92.80</u>	-	-	<u>95.76</u>	-
rd+	<u>46.14</u>	<u>52.30</u>	71.47	<u>78.48</u>	<u>84.05</u>	90.42	<u>87.94</u>	<u>88.34</u>	91.16	-	<u>92.69</u>	-	-	<u>95.73</u>	-
Concept Distribution															
	bal	slt	long	bal	slt	long	bal	slt	long	bal	slt	long	bal	slt	long
bal	<u>50.47</u>	<u>53.04</u>	54.35	80.71	<u>75.79</u>	<u>74.54</u>	85.08	<u>83.79</u>	<u>75.10</u>	92.65	<u>90.82</u>	<u>83.74</u>	95.36	<u>94.89</u>	<u>89.88</u>
long	<u>49.43</u>	<u>54.75</u>	62.96	<u>79.06</u>	<u>80.29</u>	90.66	<u>85.33</u>	<u>89.42</u>	91.10	<u>92.73</u>	93.38	<u>92.53</u>	<u>96.31</u>	96.32	<u>95.25</u>
head	<u>48.60</u>	<u>58.06</u>	61.60	<u>80.75</u>	<u>79.60</u>	87.46	<u>84.58</u>	<u>88.39</u>	90.19	<u>93.87</u>	94.82	<u>92.48</u>	<u>96.42</u>	96.80	<u>95.92</u>
tail	51.80	<u>48.70</u>	<u>50.08</u>	81.50	<u>70.88</u>	<u>60.94</u>	86.10	<u>80.27</u>	<u>60.55</u>	90.26	<u>89.20</u>	<u>75.32</u>	94.08	<u>93.20</u>	<u>82.68</u>
oppo	49.06	<u>48.93</u>	<u>46.68</u>	79.13	<u>68.37</u>	<u>56.98</u>	85.07	<u>77.86</u>	<u>55.14</u>	91.22	<u>88.65</u>	<u>71.32</u>	95.76	<u>94.09</u>	<u>79.74</u>
Concept Compositionality															
	co-0	co-1	co-2	co-0	co-1	co-2	co-0	co-1	co-2	co-0	co-1	co-2	co-0	co-1	co-2
co-0	<u>53.28</u>	57.00	<u>56.1</u>	83.36	<u>77.03</u>	<u>82.43</u>	89.1	<u>82.52</u>	<u>83.77</u>	92.80	<u>90.11</u>	<u>91.59</u>	95.76	<u>94.02</u>	<u>95.12</u>
co-1	<u>52.41</u>	60.57	<u>56.67</u>	<u>79.46</u>	<u>82.45</u>	83.93	<u>78.89</u>	87.18	<u>84.2</u>	<u>78.74</u>	<u>89.99</u>	90.67	<u>87.12</u>	<u>94.53</u>	94.78
co-2	<u>52.96</u>	<u>57.37</u>	60.53	<u>80.03</u>	<u>77.41</u>	87.24	<u>78.40</u>	<u>81.55</u>	88.84	<u>77.85</u>	<u>89.28</u>	92.23	<u>87.19</u>	<u>93.49</u>	95.61

Out-of-Domain Testing: Summary Results

		Visual	Redund.	Dist.	Comp.
non-modular	FiLM	4.03	21.33	28.46	9.04
	mDETR	9.81	19.05	36.34	9.45
modular	NSCL	15.57	0.92	37.44	15.40
	NSVQA	17.48	0.08	20.92	11.44
	NSVQA + Prob	12.88	0.08	13.72	7.00

Table 2. *Relative Degrade* under domain shifts, *i.e.* the percentage of accuracy decrease when the model is tested with domain that differs with training. Lower *RD* means better robustness.

Findings: comparison between models

	Visual	Redund.	Dist.	Comp.
FiLM	4.03	21.33	28.46	9.04
mDETR	9.81	19.05	36.34	9.45
NSCL	15.57	0.92	37.44	15.40
NSVQA	17.48	0.08	20.92	11.44
NSVQA + Prob	12.88	0.08	13.72	7.00

- **Neural symbolic methods are robust on question redundancy**
Question parsing is easy, and trained separately

Findings: comparison between models

	Visual	Redund.	Dist.	Comp.
FiLM	4.03	21.33	28.46	9.04
mDETR	9.81	19.05	36.34	9.45
NSCL	15.57	0.92	37.44	15.40
NSVQA	17.48	0.08	20.92	11.44
NSVQA + Prob	12.88	0.08	13.72	7.00

- Neural symbolic methods are robust on question redundancy
- Neural symbolic methods are **only** robust on question redundancy
Why?]
Maybe we need not only modular network, but also modular training

Findings: comparison between models

	Visual	Redund.	Dist.	Comp.
FiLM	4.03	21.33	28.46	9.04
mDETR	9.81	19.05	36.34	9.45
NSCL	15.57	0.92	37.44	15.40
NSVQA	17.48	0.08	20.92	11.44
NSVQA + Prob	12.88	0.08	13.72	7.00

- P-NSVQA is the most robust on 3 out of 4 factors
 - Probabilistic + symbolic -> best model

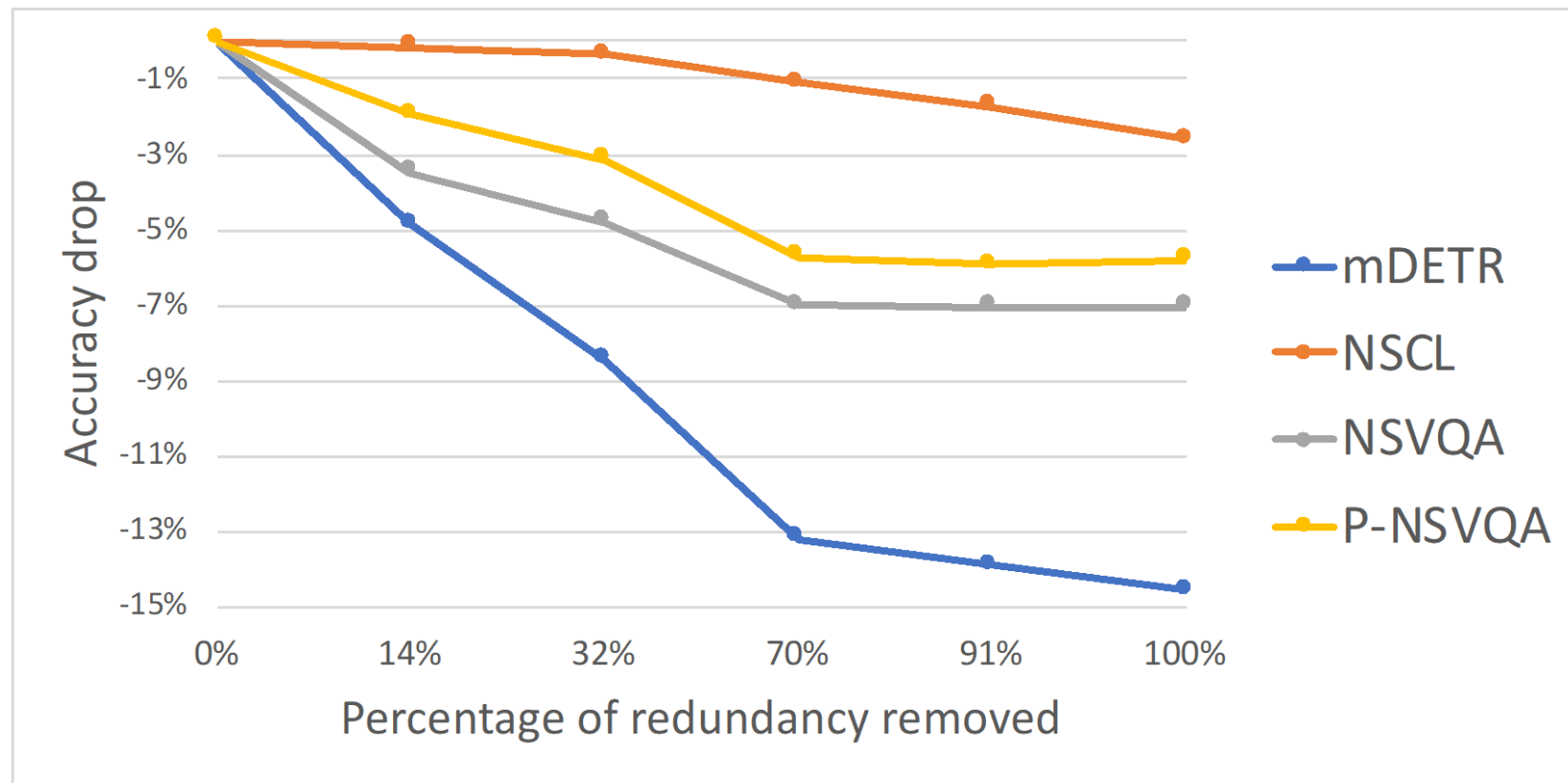
Findings: comparison between models

	Visual	Redund.	Dist.	Comp.
FiLM	4.03	21.33	28.46	9.04
mDETR	9.81	19.05	36.34	9.45
NSCL	15.57	0.92	37.44	15.40
NSVQA	17.48	0.08	20.92	11.44
NSVQA + Prob	12.88	0.08	13.72	7.00

- On visual complexity, end-to-end methods are more robust
mDETR has a more powerful visual component

Will the findings generalize to real data?

For question redundancy:



Progressively remove the redundant operations from questions in GQA dataset

Extension: Super-CLEVR with Parts, 3D Pose, Occlusion

Ongoing work

Part questions

Objects from UDA-Part dataset

Object with parts:



car



bus



airplane



bike



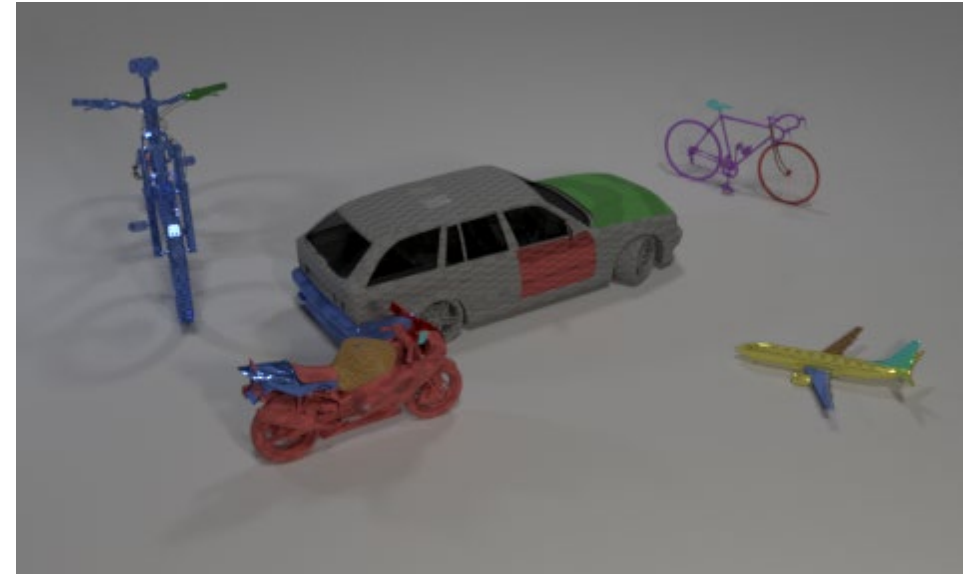
motorcycle

Texture: dotted, checkered, striped, none

Color: green, gray, brown, yellow, red, purple, cyan, blue

Size: large, small

Material: rubber, metal



Q: What is the color of the **front wheel** of the small purple bicycle?

Q: What is the material of the yellow **object that has** a blue part?

Q: What is the color of the front wheel that **belongs to the same object as** the cyan seat?

Figure 2. Super-CLEVR contains 21 vehicle models belonging to 5 categories, with controllable attributes.

Pose questions



Q: Which direction is the tiny blue school bus facing?

Q: What is the brown thing facing in the same direction as the tiny blue school bus?

Q: Is the plane and the tiny blue school bus facing in the opposite direction?

Occlusion questions

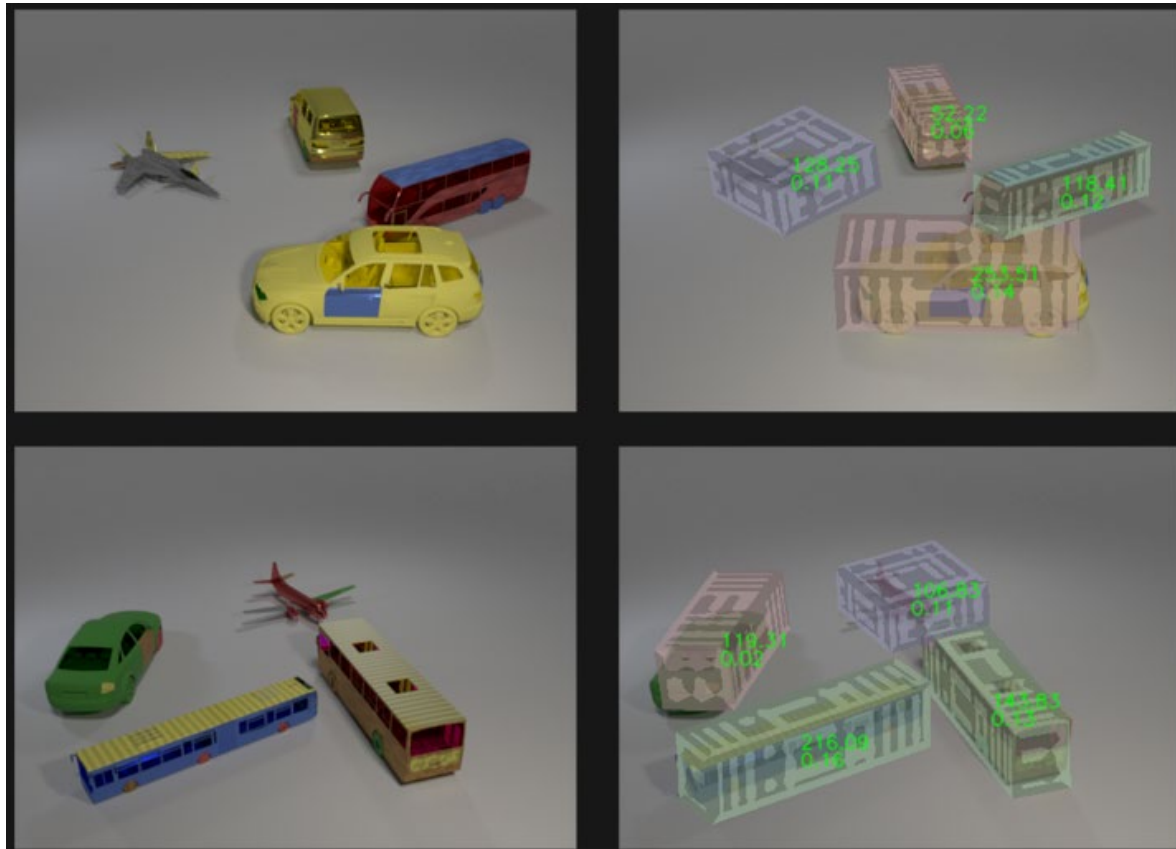


Q: What is the size of the purple object **that is occluded**?

Q: What part of the small rubber object is **occluded**?

Why part, pose, occlusion?

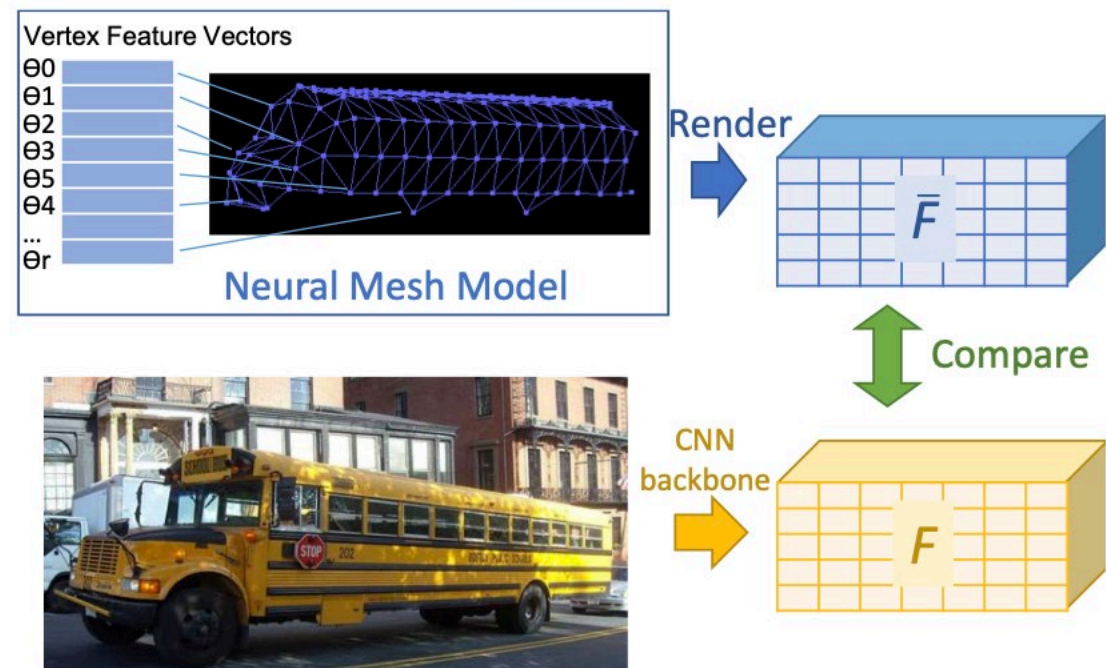
3D knowledge is necessary to answer those questions.



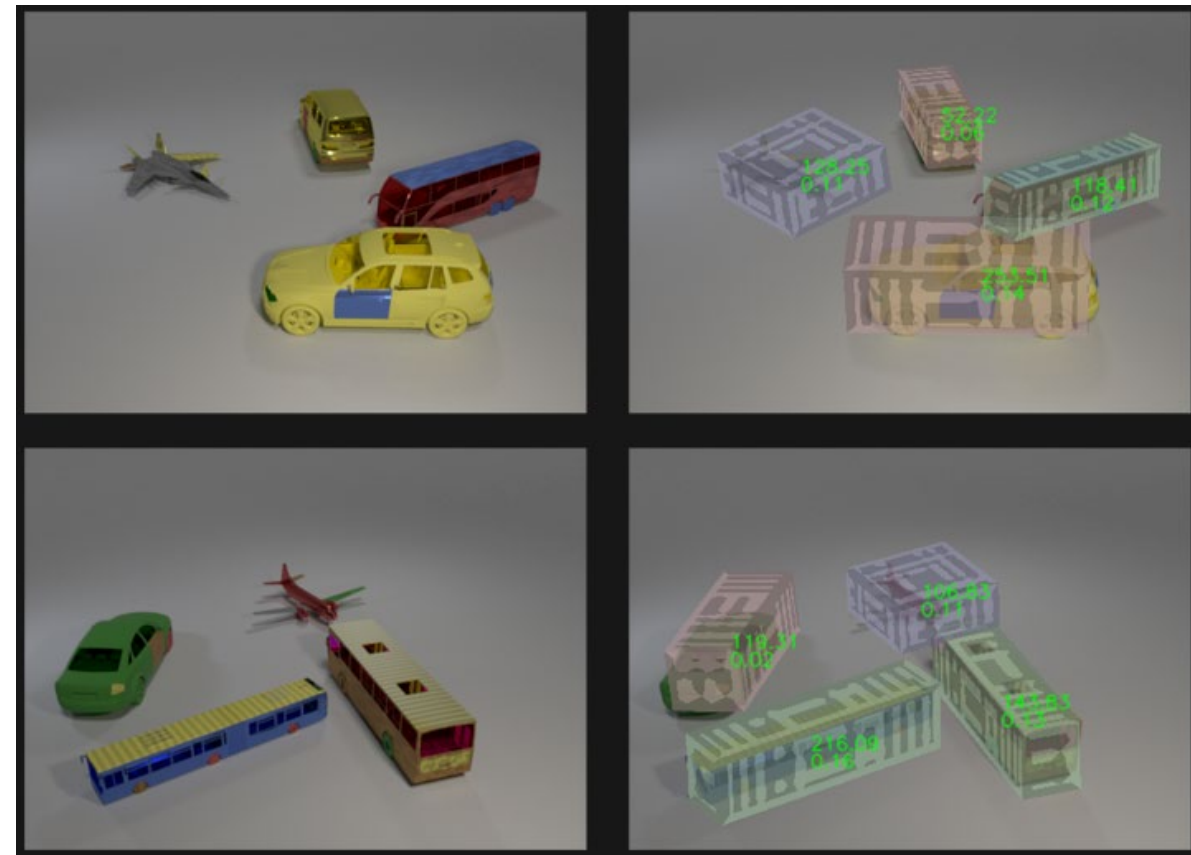
Introduce 3D Pose estimation in VQA

NeMo (Wang et al. ICLR 2021)

- A generative pose estimation model
- Robust to occlusion



NeMo Render-and-Compare: Feature Map



Ongoing results

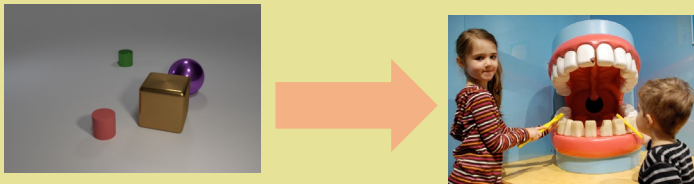
	Part	Pose	Occlusion	Part + Occlusion
mDETR	41.52	71.76	64.99	44.60
P-NSVQA with FRCNN	-	87.78	-	-
P-NSVQA with NeMO	-	86.40	-	-
P-NSVQA with NeMO & GT		91.34		

Table being completed.

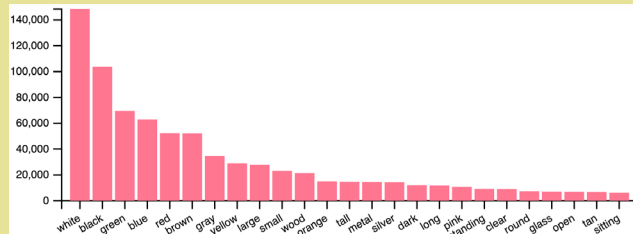
Towards Generalizable Visual Reasoning

Calibrating concepts and operations

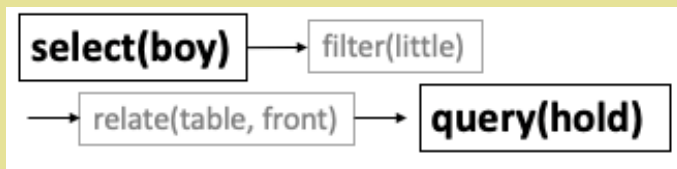
Synthetic-real gaps, and solutions



Long tail distribution



Unequal importance of operations



Super-CLEVR

systematically study domain robustness



Domain shift factors

- Visual complexity
- Question redundancy
- Concept distribution
- Concept compositionality
- ...

3D-aware VQA

- Object parts
- 3D poses
- Occlusions
- ...

Probabilistic + Symbolic → robust model

Future direction?

SS-CLEVR – more Realistic images.

Adversarial Testing

Stronger Vision Modules.

Thanks! Questions?

zli110@jhu.edu

