

1



2



3



4



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

Highlight | WED-PM-249 | Jun 21, 2023



Zhuowan Li¹



Xingrui Wang²



Elias Stengel-Eskin¹



Adam Kortylewski^{3,4}



Wufei Ma¹



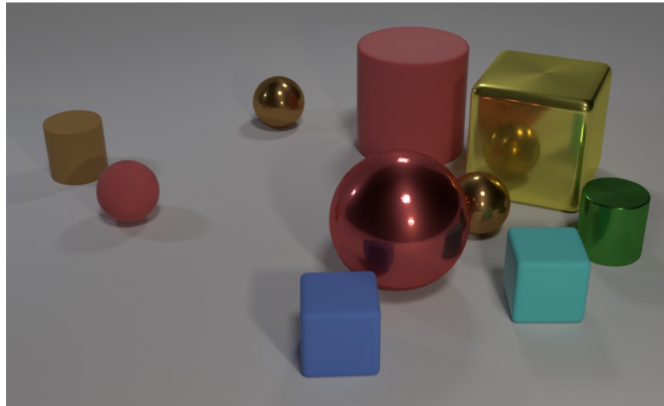
Benjamin Van Durme¹



Alan Yuille¹

Visual Reasoning

CLEVR



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

VQA2.0



What color are her eyes?

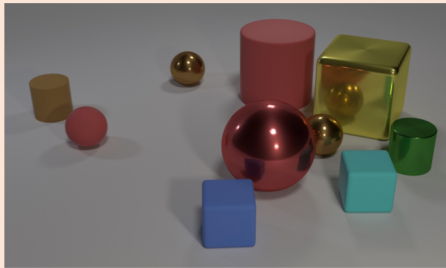
GQA



What color is the food on the red object left of the small girl that is holding a hamburger, yellow or brown?

Visual reasoning models suffer on out-of-domain testing

Domain-1



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

Training



poor performance

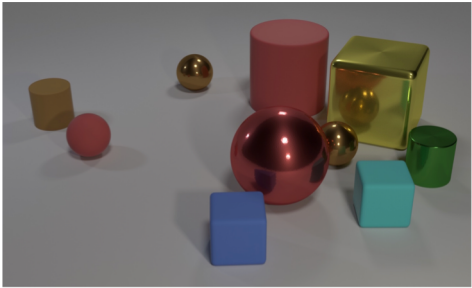
Domain-2



What color is the food on the red object left of the small girl that is holding a hamburger, yellow or brown?

Testing

Domain gaps contain multiple factors



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



What color are her eyes?



What color is the food on the red object left of the small girl that is holding a hamburger, yellow or brown?

How are they different from each other?

Multiple factors:

- Image styles?
- Question lengths?
- Concepts?
- ...

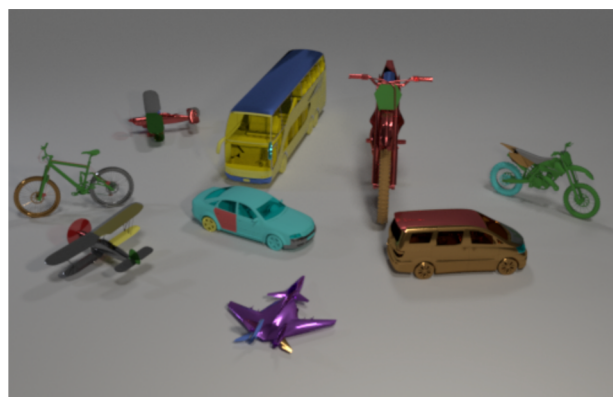


Super-CLEVR:

study each domain shift factor separately

Domain A

Super-CLEVR



“What color is the bus?”

Domain B

Visual Complexity



easy



middle



hard

Question Redundancy

- redundancy

standard

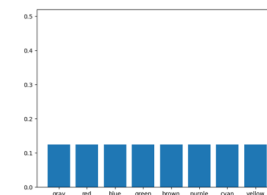
+ redundancy

“What color is the bus?”

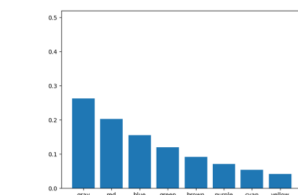
“What color is the large bus?”

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

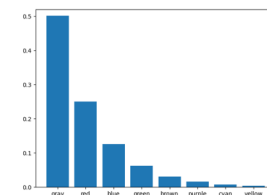
Concept Distribution



balanced

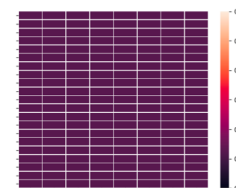


unbalanced

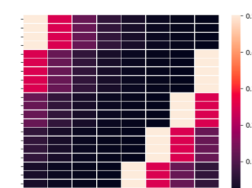


Long-tail

Concept Compositionality



well-composed

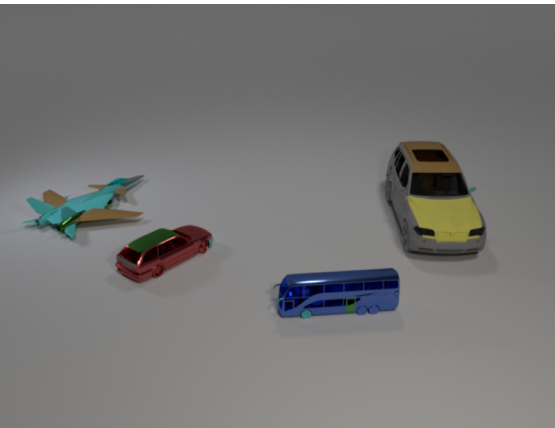


correlated

Super-CLEVR: a controllable dataset



The dirtbike that is the same size as the brown motorbike is what color?

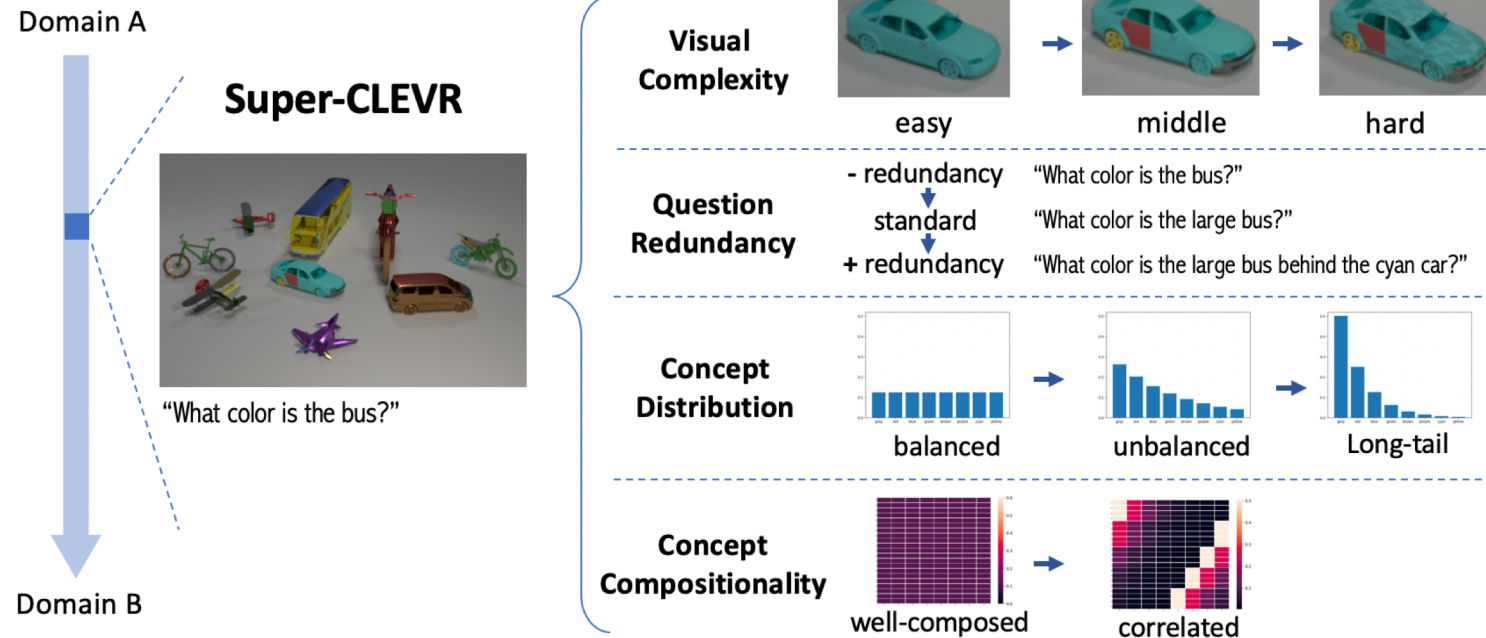


What is the color of the small car?

- Separately study the 4 domain shift factors
- 5 models are studied
- Analysis finding:
modular training + probabilistic execution → best model

In more detail...

- 4 domain shift factors
- 5 models we studied
- Analysis Results



Study domain shift by decomposition

Factors contributes to domain shifts in visual reasoning:

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

Decompose VQA domain shifts into 4 factors

- **visual complexity**
how hard is the image



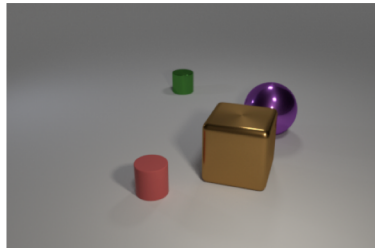
easy



middle



hard



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

- redundancy "What color is the bus?"
↓
standard "What color is the large bus?"
↓
+ redundancy "What color is the large bus behind the cyan car?"



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



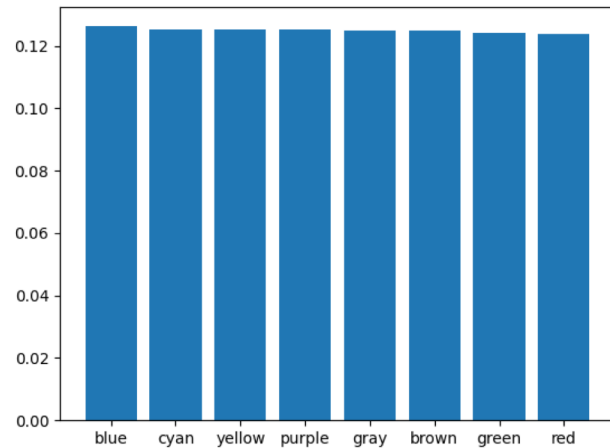
What is feeding the **large** animal ~~behind the fence~~?

- concept distribution
- concept compositionality

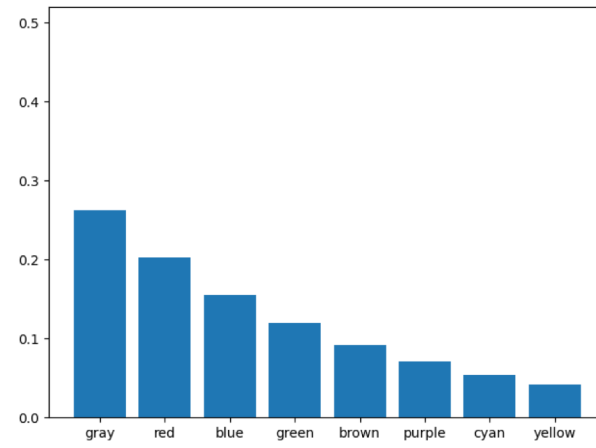
Decompose VQA domain shifts into 4 factors

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

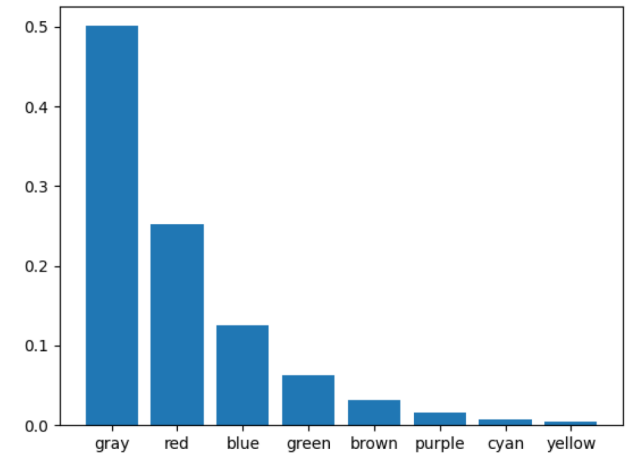
The distribution the concepts (objects names and attributes)



Well-balanced



Slightly unbalanced



Long-tail distributed

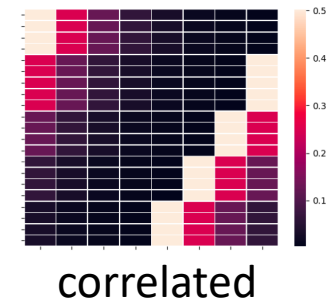
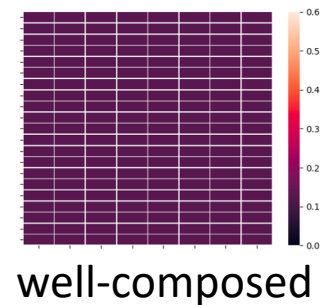
- concept compositionality

Decompose VQA domain shifts into 4 factors

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

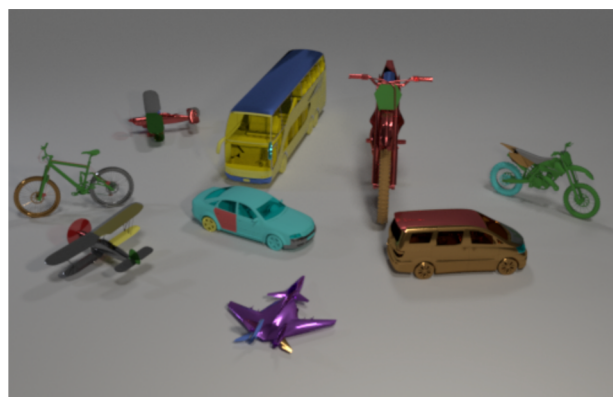
Some concepts always co-occur with another:

- *Bananas* are usually *yellow bananas*
- *Boys* are usually *little boys*
- *Skies* are usually *blue skies*
- ...



Domain A

Super-CLEVR



“What color is the bus?”

Domain B

Visual Complexity



easy



middle



hard

Question Redundancy

- redundancy

standard

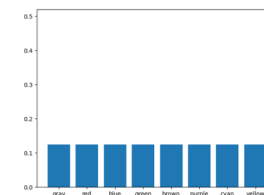
+ redundancy

“What color is the bus?”

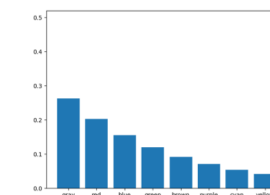
“What color is the large bus?”

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

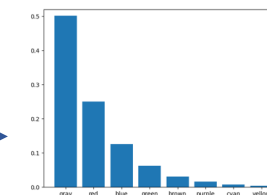
Concept Distribution



balanced

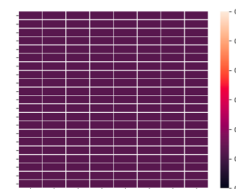


unbalanced

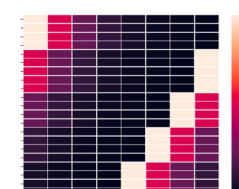


Long-tail

Concept Compositionality



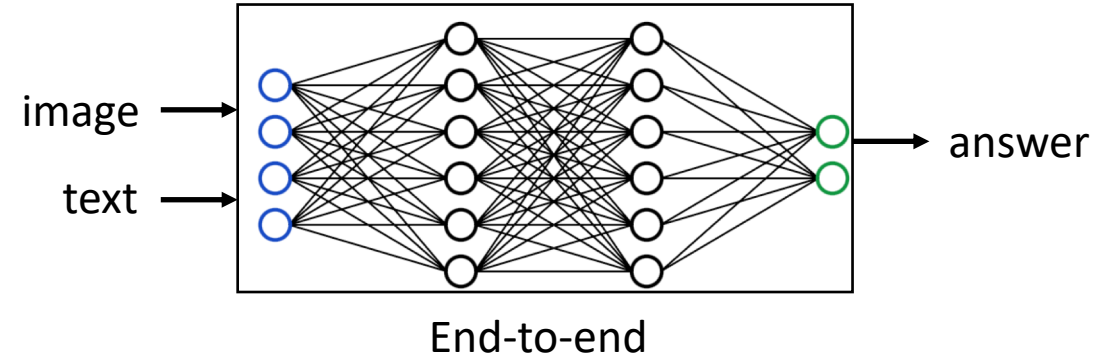
well-composed



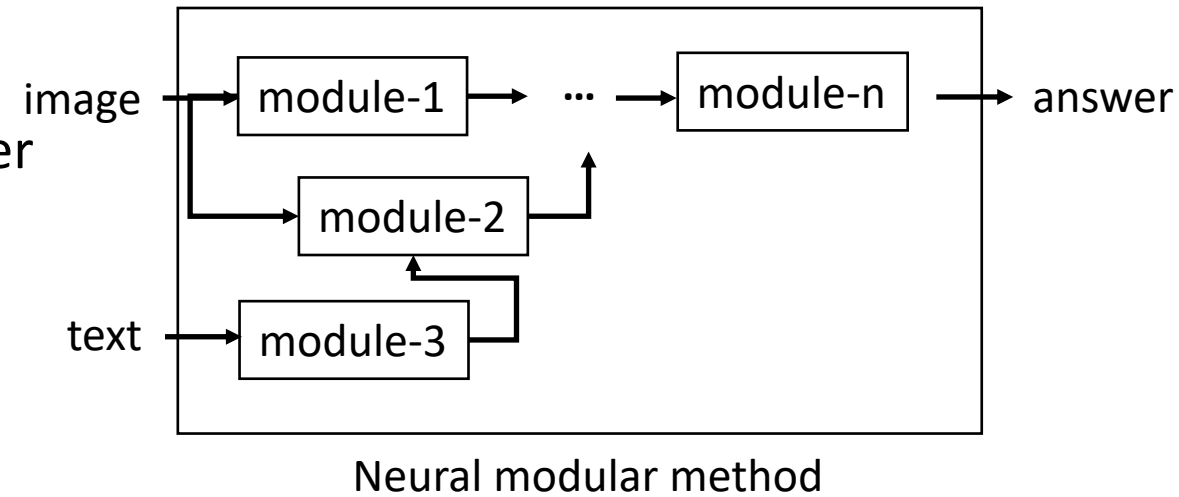
correlated

5 models are studied

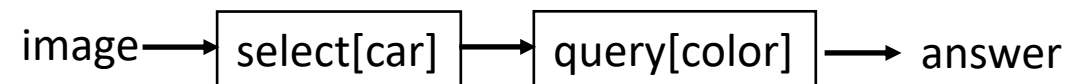
- non-modular {
- FiLM
two-stream feature merging
 - mDETR
pretrained transformer model



- modular {
- NSCL
neural symbolic concept learner
 - NSVQA
neural symbolic VQA
 - **Probabilistic NSVQA**
our method

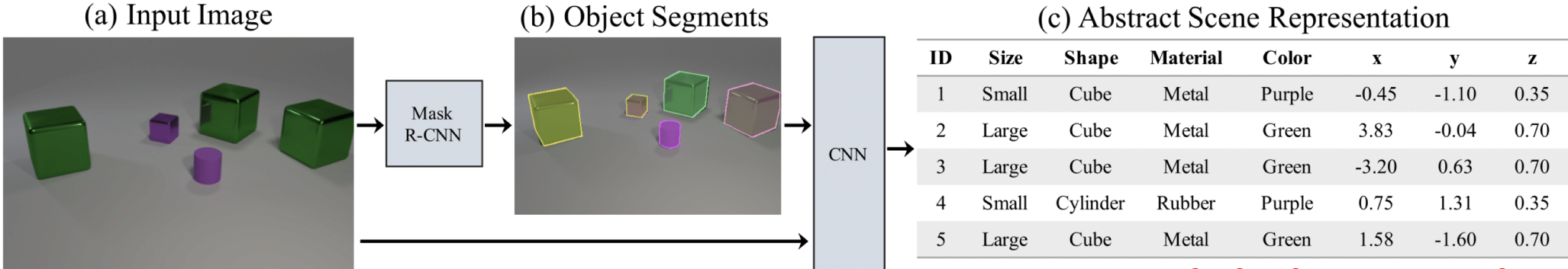


What color is the car?



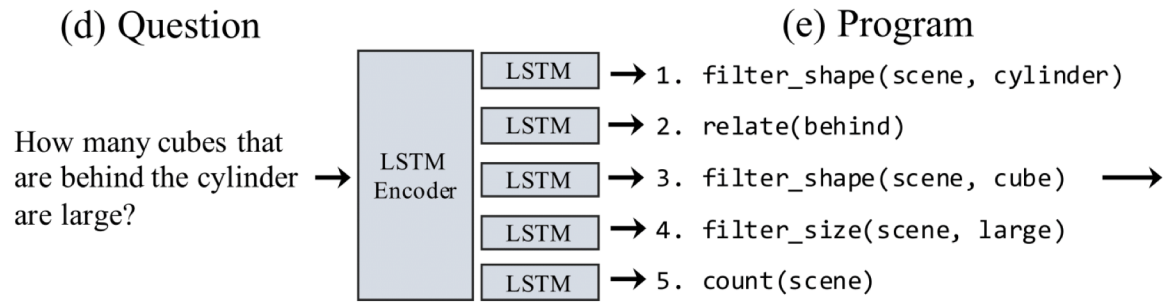
Probabilistic NSVQA (our model)

Recall NSVQA



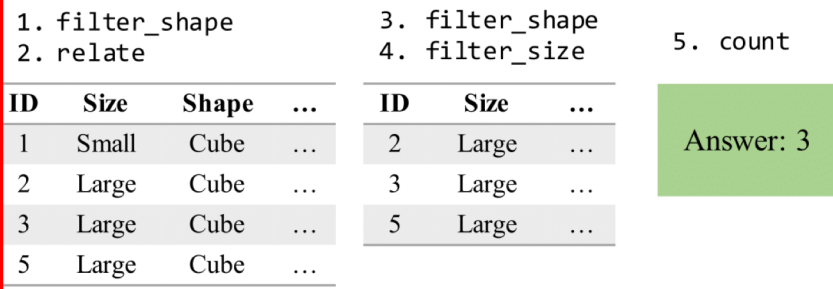
I. Neural Scene Parsing

II. Neural Question Parsing



Deterministic Execution

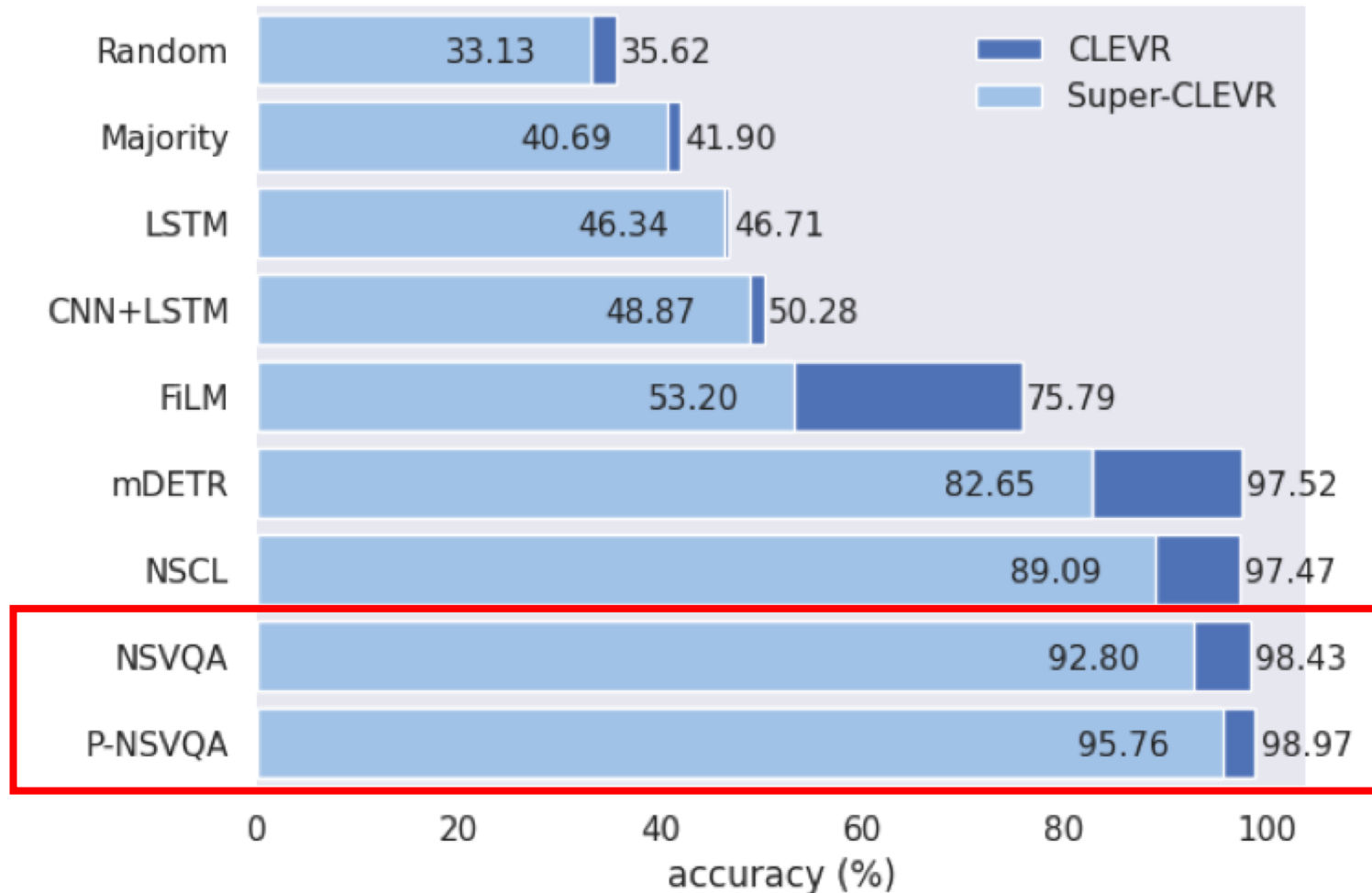
III. Symbolic Program Execution



Make it probabilistic!

Considering the uncertainty of scene parsing

In-domain Results



- Super-CLEVR is harder than CLEVR
- P-NSVQA is the best

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	92.95	<u>92.94</u>	<u>92.67</u>	<u>95.66</u>	95.72	<u>95.43</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>91.19</u>	92.78	<u>92.14</u>	<u>94.87</u>	95.72	<u>95.43</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>91.38</u>	<u>91.96</u>	92.80	<u>94.88</u>	<u>95.47</u>	95.72
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

Summary OOD testing results

Relative Degrade

- the percentage of accuracy decrease when the model is tested with domain that differs with training

$$RD = \frac{Acc_{iid} - Acc_{ood}}{Acc_{iid}}$$

		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	1.72	20.92	11.44
	Prob NSVQA	12.88	0.84	13.72	7.00

Table. *Relative Degrade* of models. Smaller is better.

Finding-1: modular models are very robust on redundancy

	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	1.72	20.92	11.44
Prob NSVQA	12.88	0.84	13.72	7.00

- Modular training is important
- Question component shouldn't compensate for visual understanding

Finding-2: P-NSVQA is the best on 3 out of 4 factors

	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	1.72	20.92	11.44
Prob NSVQA	12.88	0.84	13.72	7.00

- Modularity + Probabilistic execution -> best model

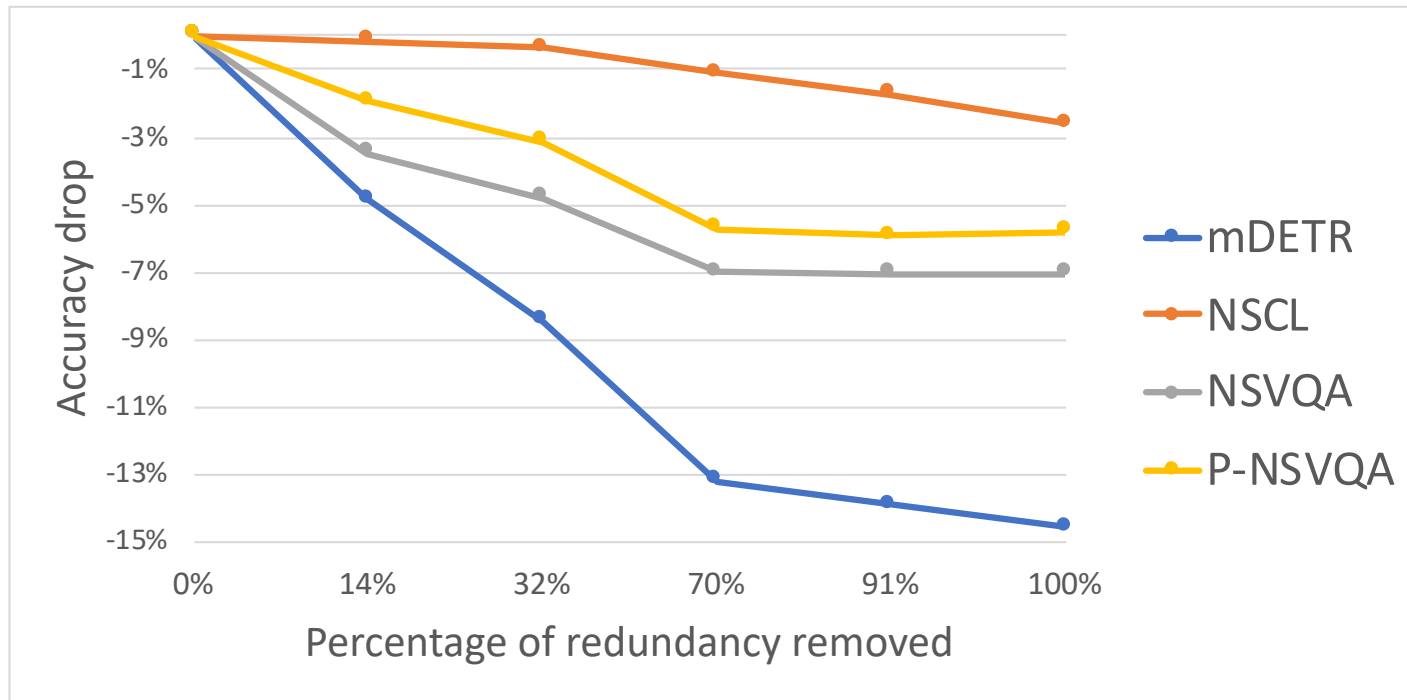
Finding-3: non-modular methods win on visual complexity

	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	1.72	20.92	11.44
Prob NSVQA	12.88	0.84	13.72	7.00

- mDETR has a more powerful pretrained visual component
- Visual scene parser (MaskRCNN) in modular methods can be improved.

Will the findings generalize to real data?

We verify findings on question redundancy on the real GQA dataset



When the redundant operations are progressively removed, the performance drops of modular methods are smaller than non-modular methods.

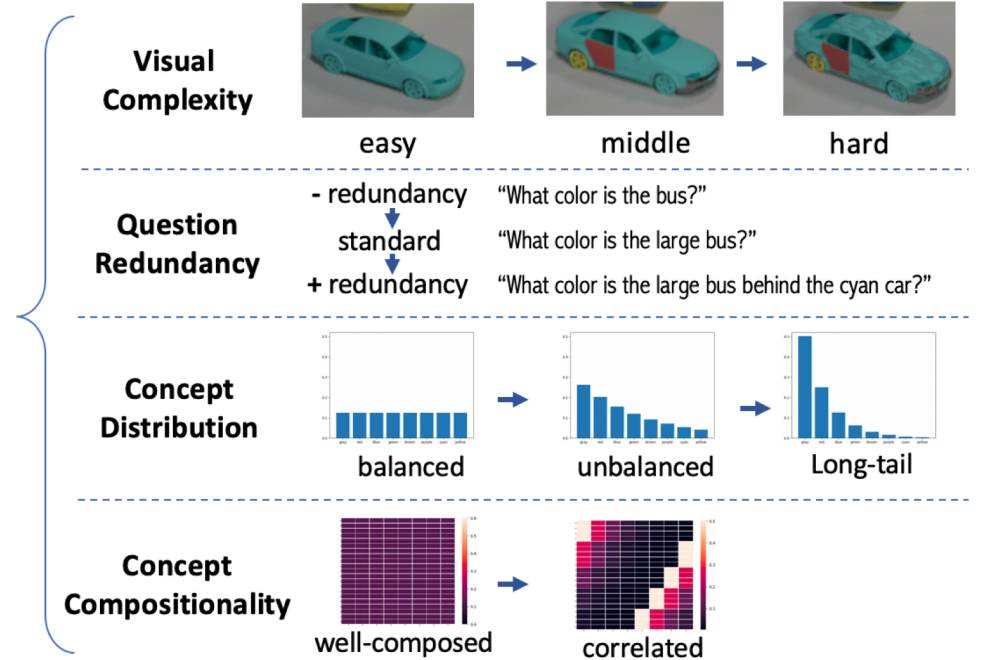
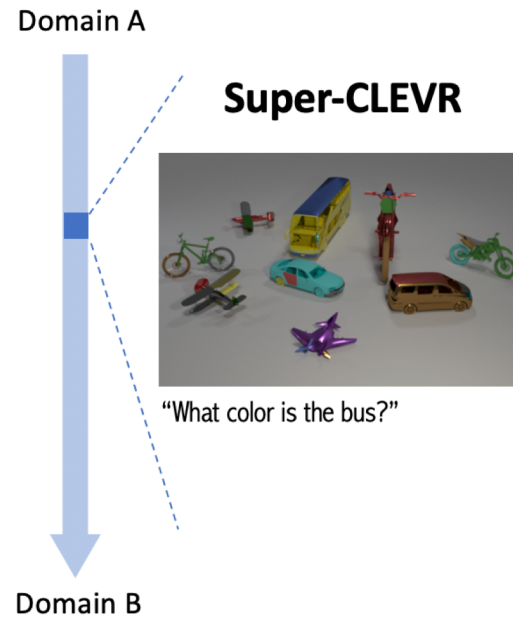
Take-home messages

- Super-CLEVR dataset

- Analysis findings:

1. Modular models are very robust on question redundancy.
2. P-NSVQA is the best on 3 out of 4 factors.
3. Non-modular methods win on visual complexity.

➔ Modularity and probabilistic execution are important; we need better visual modules.



Welcome to our session!

WED-PM-249 (Highlight)

Jun 21, 2023