# Efficient Visual Reasoning with Language Grounding

April 11th ,2022
Presenter: Sihui Wang

# Problem Statement

- Motivation

  - What are the trends and challenges for VQA?

  - Make a retrospective comparison for current methods

  - Find ideas for the Next Generation VQA models

# Problem Statement

- Long Term Goals
  - Find ideas for Next Generation VQA models
  - Reimplementation/Modification of Current Methods
  - In-depth Experiments and Analysis; Experiments on GQA Dataset
- Short Term Goals
  - Make a Survey of the Current Methods; Trend Analysis
  - Test and Comparison of Current Methods on CLEVR Dataset

# Problem Statement

- Input
  - CLEVR Dataset
  - Codes for PG + EE, NS-VQA
  - Codes for Baseline Methods (LSTM + CNN + Attention + MLP)
- Output
  - Trend Analysis
  - Comparison of Performance of PG + EE & NS-VQA against Baselines

# Related Work

- 3 Current Approaches

  - (PG + EE) Inferring and Executing Programs for Visual Reasoning

  - FiLM: Visual Reasoning with a General Conditioning Layer

  - (NS-VQA) Neural-Symbolic VQA: Disentangling Reasoning from Vision and Language Understanding

- Baseline Methods

  - (LSTM-Q) Visual Question Answering

  - (LSTM + CNN + Self Attention) Stacked Attention Networks for Image Question Answering

# Related Work

- Comparison:
  - Baseline Models: Combinations of LSTM, CNN, Attention...
  - PG + EE: Explicit Modelling of Reasoning Process (to avoid learning the biases)
  - FiLM: Refined Design of Feature Interaction Layer
  - NS-VQA: Question Parser Guided Program Generator

# Related Work

| | PG + EE | FiLM | NS-VQA |
|---|---|---|---|
| Features | CNN + LSTM | CNN + LSTM | CNN +LSTM |
| Reasoning | Dynamic Combination of Module Networks (RL) | Feature-wise Linear Modulation | Dynamic Combination of Module Networks (Symbolic Parsing) |
| High Level Ideas | Explicit Modelling of Reasoning | Efficient and Expressive Feature Interactions | Symbolic Structure as Prior Knowledge |
| Interpretability | High | Moderate (Visualization of Attention) | High |

# Related Work

- Trend:
    - From Basic Building Blocks to High-Level Ideas
    - Interpretability
    - Data Efficiency
    - Capacity
- Why This Project:
    - Retrospective Perspective and Evaluation

# Approach

- Literature Review and Trend Analysis

- Compare Current Methods' Performance on CLEVR Dataset

- Models to Use:

  - PG + EE

  - NS – VQA

  - Baseline Models: LSTM, LSTM + CNN, LSTM + CNN + Self Attention, LSTM + CNN + Self Attention + MLP
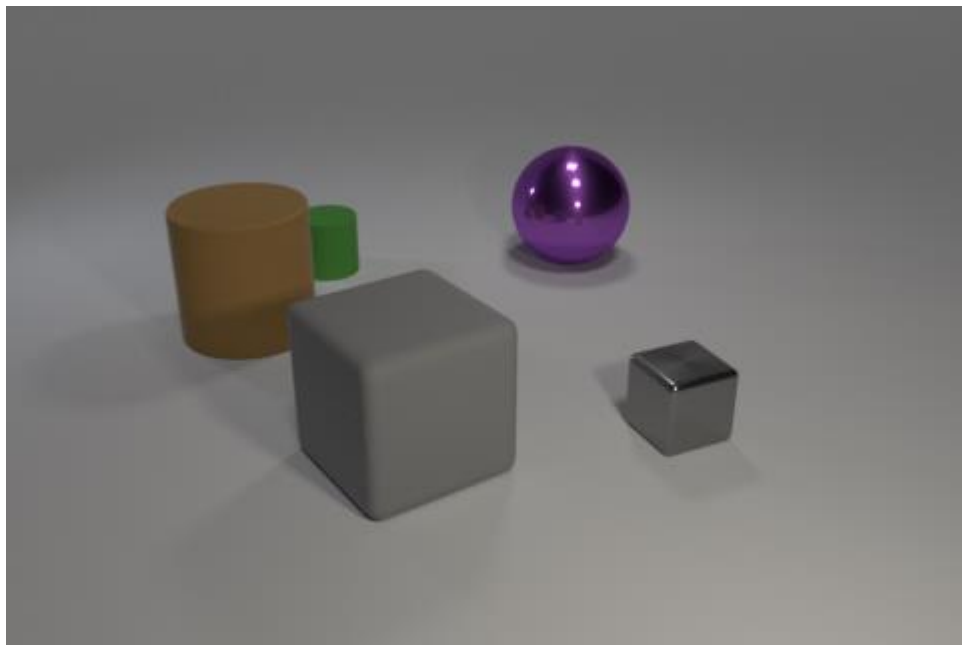
# Results and Analysis: Quantitative

- Accuracy on CLEVR Validation Dataset

|  | Accuracy | Accuracy (Previously Reported) |
|---|---|---|
| LSTM | 47.05 % | 46.8 - 47.0 % |
| LSTM + CNN | 54.19 % | 52.3 - 54.3 % |
| LSTM + CNN + SA | 69.48 % | 68.5 – 76.6 % |
| LSTM + CNN + SA + MLP | 73.15 % | 73.2 % |
| PG + EE (18K Prog) | 95.19 % | 95.4 % |
| NS - VQA | 99.8 % | 99.8 % |

# Results and Analysis: Qualitative

- Case Study: Why does PG + EE fail?



Typical Scenes:
Complex Questions: Logical Composition + Multiple Attributes (Color, Location, Texture,...) + Ask Number
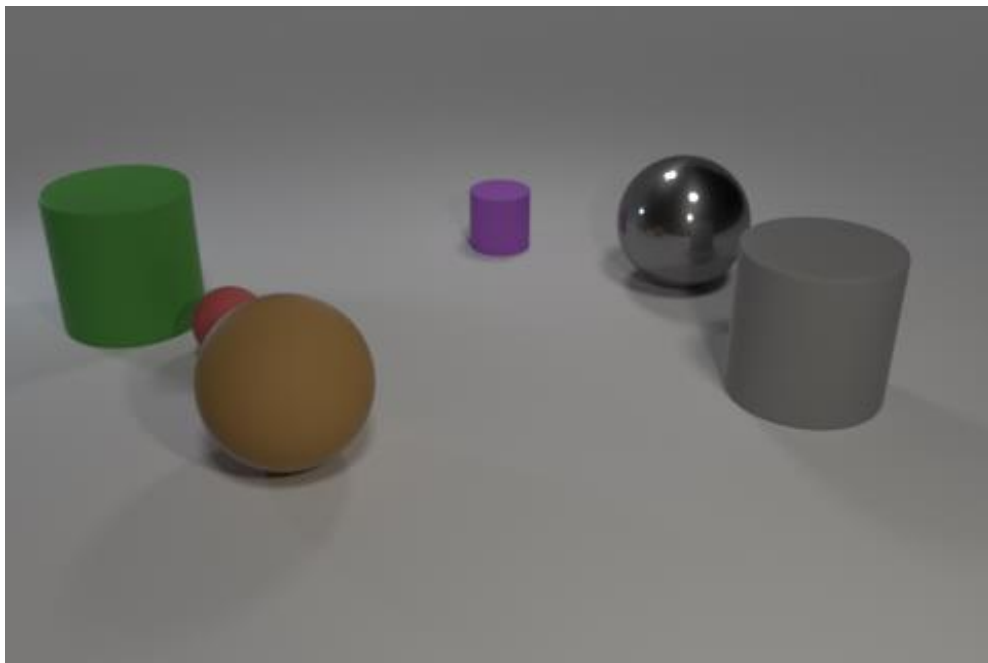
How many objects are either metal things behind the small green rubber cylinder or small green rubber objects?

Correct Answer: 2
Predicted Answer: 1

# Results and Analysis: Qualitative

- Case Study: Why does PG + EE fail?



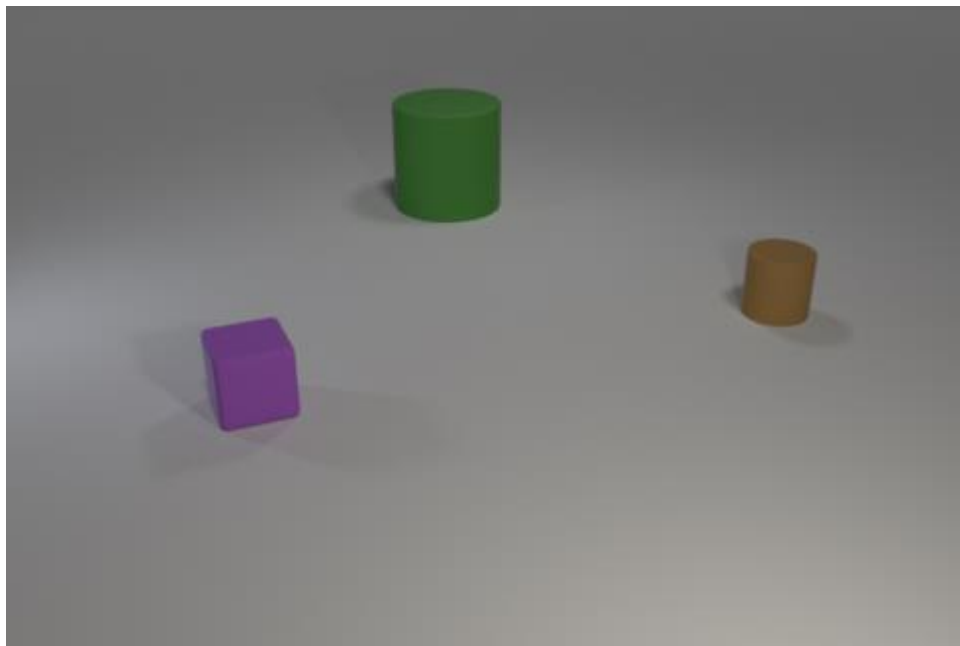What color is the small ball?
Correct Answer: red
Predicted Answer: cyan

Where does cyan come from?
Biases?

# Results and Analysis: Qualitative

- Case Study: Why does PG + EE fail?



What shape is the brown thing?
Correct Answer: cylinder
Predicted Answer: no

Totally irrelevant

# Conclusion

- All models' performance on CLEVR validation dataset match the claims in the papers

- Top-Peformance VQA models are guided by high-level insights and refined architecture design

- Even though PG + EE achieves good performance on CLEVR dataset, qualitative results show that it still doen't fully understand the questions.

# Limitations

- CLEVR's synthetic nature provides us the opportunity to perform breakdown analysis why certain method fails. However, in-depth comparison is not presented in this work.

- Comparisons are made only on synthetic datasets. Real world dataset should also be included in the study.

- Top performance on CLEVR doesn't mean current methods have learned reasoning abilities.

# Future works

- Include more breakdown analysis: which building block contributes to performance? What are the reasons for the failures?

- Make comparisons on real-world datasets, such as GQA

- Current Methods still lacks the real reasoning abilities. Future work should focus on improvements/mitigations.

# Takeaway

Current VQA models still don't fully understand the questions and images. Substantial work is required towards learning real reasoning abilities.

# Q & A