# LRTA: A Transparent Neural-Symbolic Reasoning Framework with Modular Supervision for Visual Question Answering

Weixin Liang[1], Feiyang Niu[2], Aishwarya Reganti[2], Govind Thattai[2], Gokhan Tur[2]

[1]Department of Electrical Engineering, Stanford University, US
[2]Alexa AI, Amazon, US

## 1. Multimodal Question Answering

- **Problem**: answer free-form questions by reasoning about presented images
- **Dataset**: GQA
  - 113,018 images & 1.5M questions
  - 1702 object classes
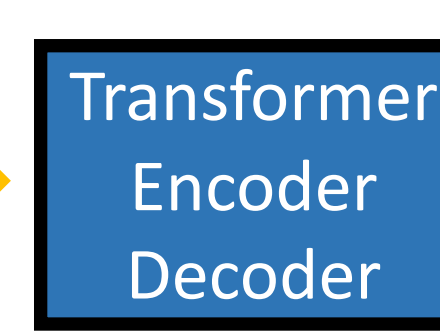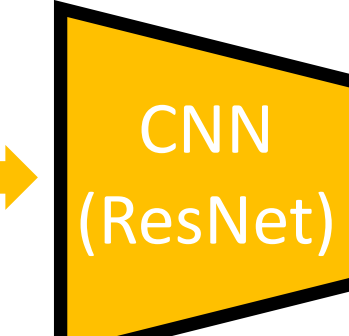
Scene graph / Question semantics

**Question:**
Is there any red object to the left of the small girl who is holding a hamburger?

**Short Answer**: Yes
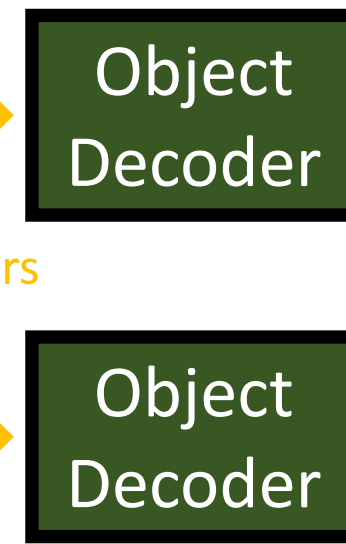**Long Answer**: Yes, there is a red tray

## 2. Motivation

**Question   Image**

Transformer Multimodal Encoder

Classifier

**Answer**

LXMERT

**Our Main Baseline in this talk**

**Limitations**
- LXMERT not really understands the question
  - Mask out relationship words (e.g. "to the left of") in the questions: 59.7% → 55.5% (less than 5% drop)[1]
- A "black-box" neural encoder without human readable justification

[1] Sanjay Subramanian, Sameer Singh, and Matt Gardner. 2019. Analyzing compositionality of visual question answering. Visually Grounded Interaction and Language Workshop.

## 3. LRTA: A More Explainable Approach

- **Question:** Is there any red object to the left of the small girl who is holding a hamburger?
- **LRTA:** Solving the problem step-by-step like humans
  - (1) **Look** at the image
  - (2) **Read** the question
  - (3) **Think** (Multi-Step)
    - (3.1) hamburger
    - (3.2) small girl
    - (3.3) tray
  - (4) **Answer**

## 4. Look: Scene Graph Generation

- **Extended Facebook DETR: Object + Bounding Box + Attributes + Relationships**

CNN (ResNet) → Transformer Encoder Decoder

**Decode Object Vectors**

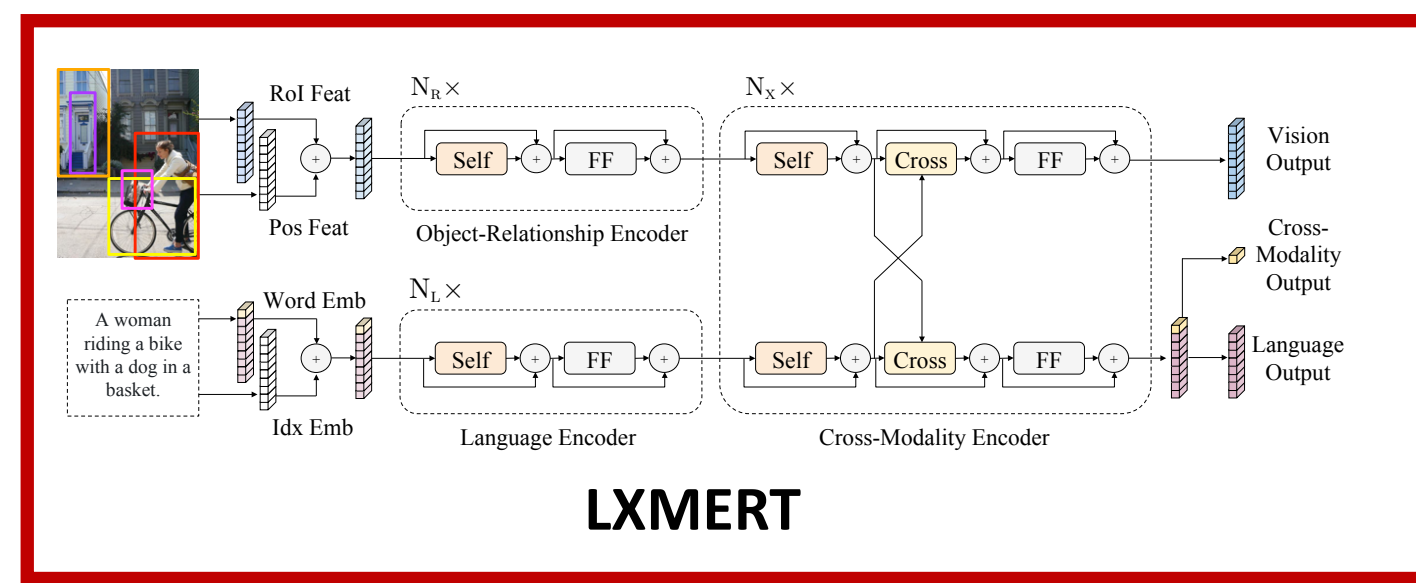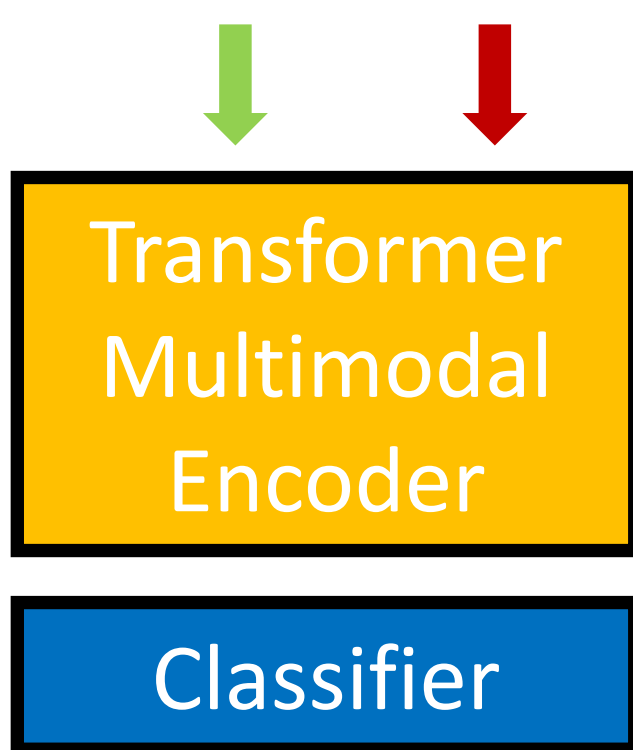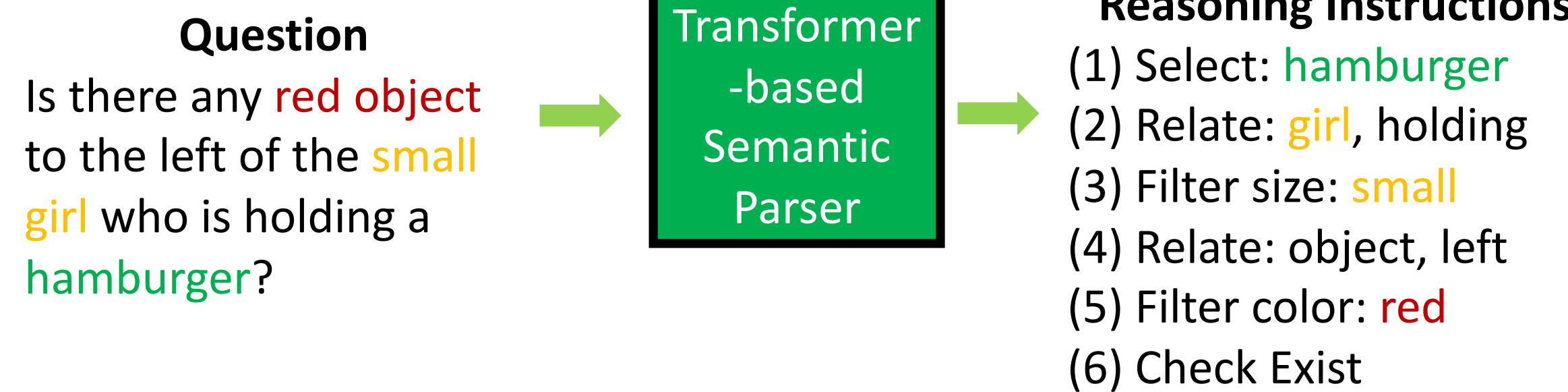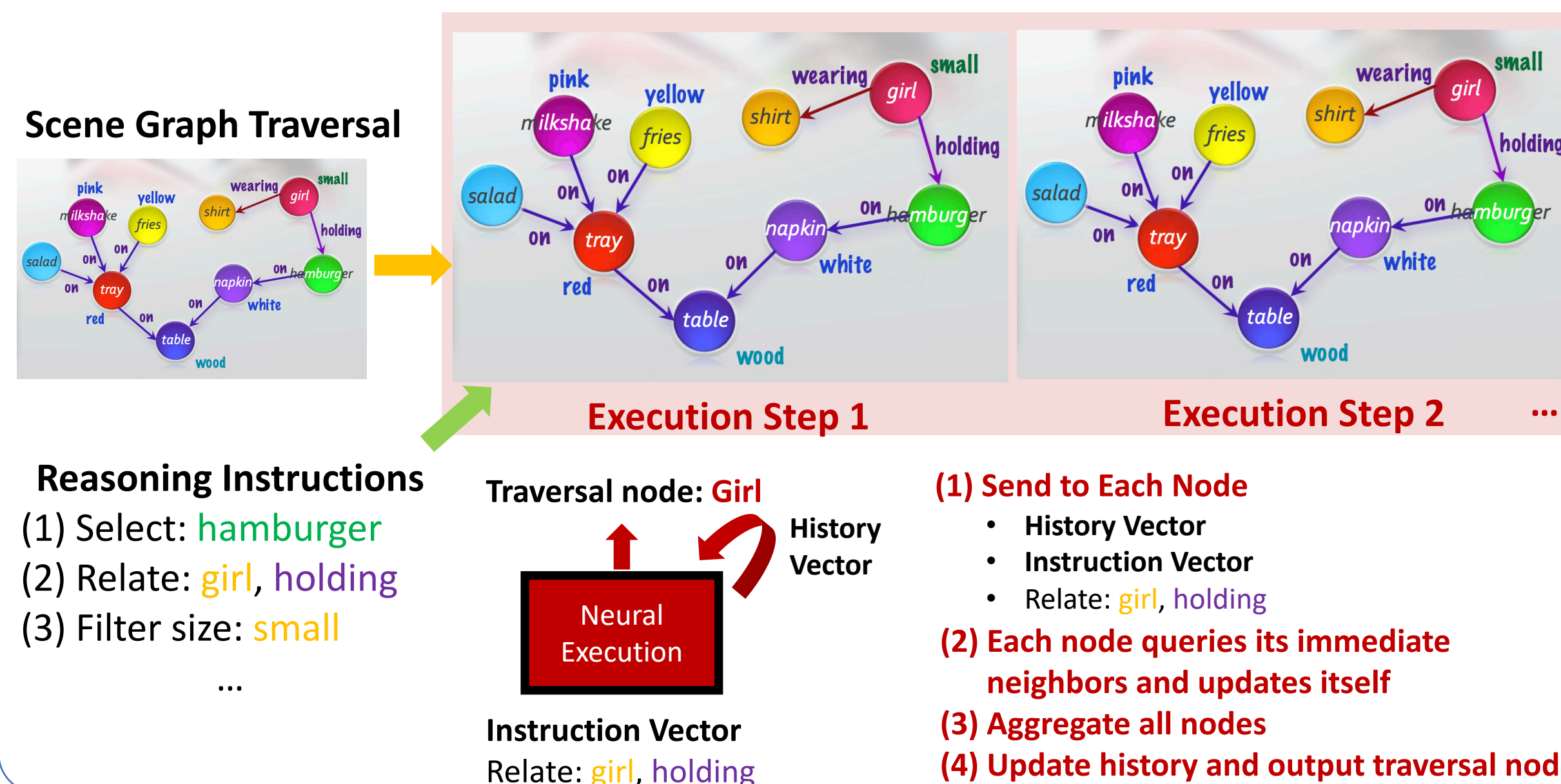Object Decoder → Object Class Bounding Box **+ Attributes (small, …) + Relationships (holding, …)**

Object Decoder

Small / Girl / Holding / Hamburger

## 5. Read: Semantic Parser Module

**Question**
Is there any red object to the left of the small girl who is holding a hamburger?

→ Transformer-based Semantic Parser →

**Reasoning Instructions**
(1) Select: hamburger
(2) Relate: girl, holding
(3) Filter size: small
(4) Relate: object, left
(5) Filter color: red
(6) Check Exist

## 6. Reason: Neural Execution Module

**Scene Graph Traversal**

**Reasoning Instructions**
(1) Select: hamburger
(2) Relate: girl, holding
(3) Filter size: small
...

Execution Step 1    Execution Step 2 …

**Traversal node: Girl**
History Vector
Neural Execution
Instruction Vector
**Instruction Vector**
Relate: girl, holding

**(1) Send to Each Node**
- History Vector
- Instruction Vector
- Relate: girl, holding

**(2) Each node queries its immediate neighbors and updates itself**
**(3) Aggregate all nodes**
**(4) Update history and output traversal node**

## 7. Answer: Natural Language Generation

Question → **Read** → Semantic Parsing → **Reasoning Instructions** → (Recurrent) **Reason** Neural Execution → **Answer** Answer Generation → **Full Answer**

Image → **Look** Scene Graph Generation → Scene Graph

**Transformer Decoder Architecture**

**Input: Neural Execution Module's History Vectors**

## 8. LRTA Overview

**Step 1: Look - Scene Graph Generation**
CNN → Transformer Encoder Decoder → Scene Graph Generation
N Object Vectors

**Step 2: Read - Semantic Parsing**
Question: Is there any red object left of the girl that is holding a hamburger?
Transformer Encoder Decoder → Instruction Vector Decoder
M Instruction Vectors

Reasoning Instructions
(1) Select: hamburger
(2) Relate: girl, holding
(3) Relate: object, left
(4) Filter color: red

**Step 3: Think - Visual Reasoning**
Execution Step 1 — Neural Execution Engine Execution Step 2 (Recurrent)
Instruction Vector #1 - Select: hamburger
Instruction Vector #2 - Relate: girl, holding

**Step 4: Answer - Full Answer Generation**
Recurrent History of Neural Execution Engine → Natural Language Generation → Yes, there is a red tray.
Natural Language Justification

## 9. Experiments

| Model | Long Acc | Short Acc |
|---|---|---|
| Human | - | 89.30% |
| Bottom-up | - | 49.74% |
| MAC | - | 54.06% |
| LXMERT | 28.00% | 56.20% |
| LRTA | 43.10% | 54.48% |

**Table 1: End-to-end training experiment on testdev set**

| Model | Long Acc | Short Acc |
|---|---|---|
| LRTA trained w/ visual oracle | | |
| Evaluated w/o attributes | 67.79% | 78.21% |
| Evaluated w/o relations | 67.95% | 75.47% |
| Evaluated w/o attributes & relations | 50.15% | 61.15% |
| Evaluated w/ visual oracle | 85.99% | 93.10% |
| LRTA trained w/ reading oracle | | |
| Evaluated w/ reading oracle | 55.45% | 64.36% |

**Table 2: Validation study on valid set**

| Model | Short Acc Drop (from → to) |
|---|---|
| VB & PRPN masked | |
| LXMERT | 19.43% (56.20% → 36.77%) |
| LRTA | 26.20% (54.48% → 28.28%) |
| Attributes masked | |
| LXMERT | 9.41% (56.20% → 46.79) |
| LRTA | 21.03% (54.48% → 33.45) |

**Table 3: Perturbation analysis on testdev set. The larger drop the better**

## 10. Conclusions

- **Contributions**
  - We propose LRTA, an end-to-end trainable, modular VQA framework facilitating explain-ability and enhanced error analysis as compared to contemporary black-box approaches.
  - We formulate VQA as a full answer generation problem to improve explainability and discourage superficial guess for answering the questions.
- **Future works**
  - Visual understanding poses as a bottleneck from our validation study and more model architectures should be explored and compared.
  - Scene graph data exhibit heavy long-tailed bias and an unbiased scene graph prediction needs to be explored, e.g. Tang et al 2020[1]

[1] K. Tang, Y. Niu, J. Huang, J. Shi, and H. Zhang. Unbiased scene graph generation from biased training. In 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR2020, Seattle, WA, USA, June 13-19, 2020, pages 3713–3722. IEEE, 2020.