# Multimodal Multihop Source Retrieval for WebQA

Nikhil Yadala[1], Paritosh Mittal[2], Saloni Mittal[1], Shubham Gupta[2]

[1]Language Technologies Institute, [2]Robotics Institute

**Group 10**

## Introduction

**Q:** At which festival can you see a castle in the background: Oktoberfest in Domplatz Austria or Tanabata festival in Hiratsuka, Japan?



- Information can be scattered across multiple sources; the proposed system should be capable of identifying and collating information critical to answering a question
- We aim to develop a system capable of selecting 'relevant' multimodal sources that can be combined to generate natural language answers to questions

## Motivation and Challenges

**Motivation**

- Information is rarely localized within individual sources
- Information can come from any combination of modalities
- Modality agnostic approach to generalize and scale with web data

**Challenges**

- Significant data imbalance between positive and negative sources
- Need for collective reasoning and 'smart' information aggregation

## Baselines

- **Lexical Overlap:** A trivial baseline that outputs the top 2 sources with the highest lexical overlap between question and caption
- **VLP:** A transformer-based model trained on MLM and VQA is used for source retrieval
  - ▶ Processes each source independently and hence poor in capturing multihop aspect of selection
  - ▶ Resource intensive and difficult to train

## Approach A: Dense Super-Node Graph
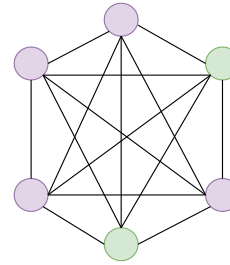
**Intuition**

- Unlike VLP, graphs can perform multihop reasoning on multiple sources
- It can learn meaningful connections between sources



**Message Passing**

- Super node contains all information about source and question
- All nodes pertaining to a question are connected together (dense)
- Source selection is reduced to node classification (+/-)
- Message passing mathematical formulation

$$x'_i = W_1 x_i + W_2 \cdot mean_{j \in \mathcal{N}(i)} x_j$$
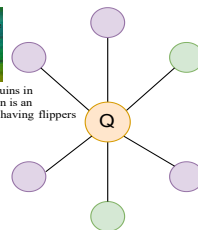
## Approach B: Star Graph

**Intuition**

- Dense graph has a large number of uninformative connections (90% negative sources)
- Dropping irrelevant connections can improve learning



- All sources for a question are connected to a central question node
- We use multiple layers of the GNNs to enable message passing through the question node
- Sparse graph leads to faster training and convergence

**Primitive Representations**

Sentence embeddings from BERT to represent textual modality and ResNet-152 features to represent image modality while SOTA uses VinVL, X101fpn and VLP based feature representations
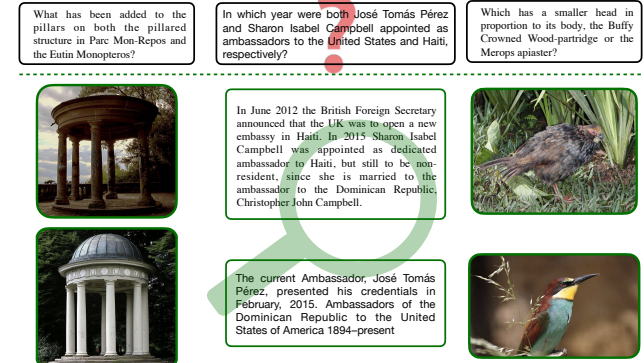
## Results

**Qualitative Results**



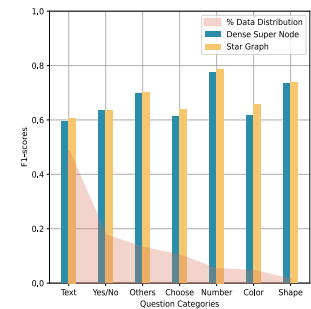Figure 1:Queries along with Retrieved Sources from Star Graph

**Quantitative Results**

| Modality | Lexical Overlap | VLP-VinVL | Super Node [†] | Star Graph [†] |
|----------|----------------|-----------|-----------|------------|
| Image | 44.83 | 68.13 | 65.59 | 66.58 |
| Text | 33.78 | 69.48 | 59.39 | 60.74 |

Table 1:F1-score comparison of baselines with our methods. [†] Ours

**Insights**

- Even with 'primitive' representations, graph based approach has comparable performance to SOTA due to inherent 'multihop' reasoning ability
- Intuition-based sparse connections are faster and improve the performance



## Ongoing Work

- Edge classification using graph attention networks
- Experiment with gated graphs for better information flow
- Using richer VinVL/CLIP features as node inputs