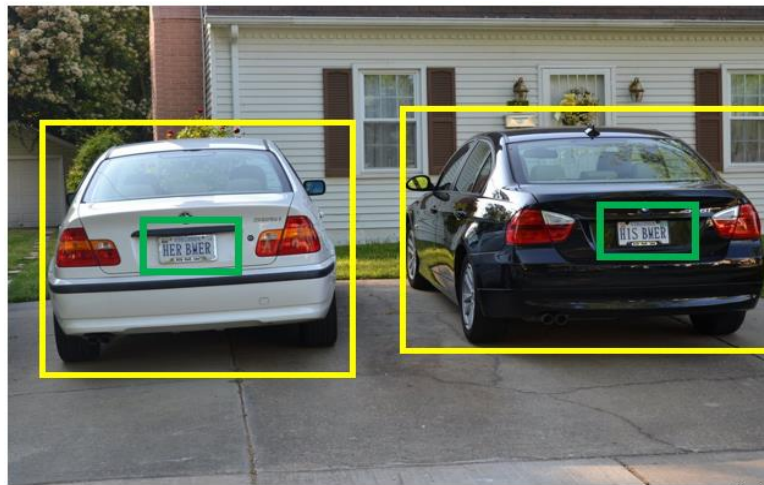# Multi-modal Contextual Graph Neural Network for TextVQA

Yaoyuan Liang

lyy8ztc@outlook.com

# Illustration of TextVQA and MCG Model



(a)

(b)

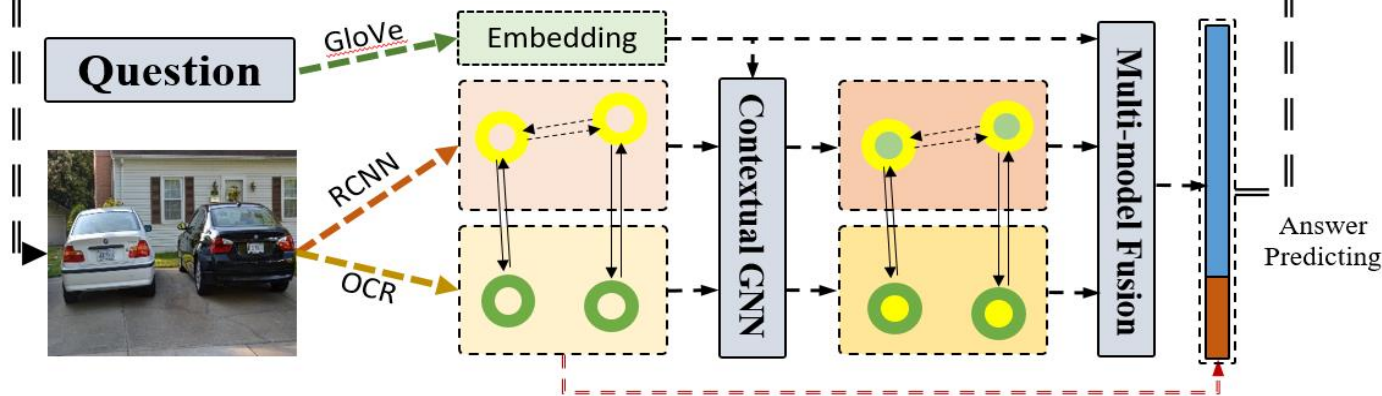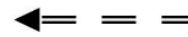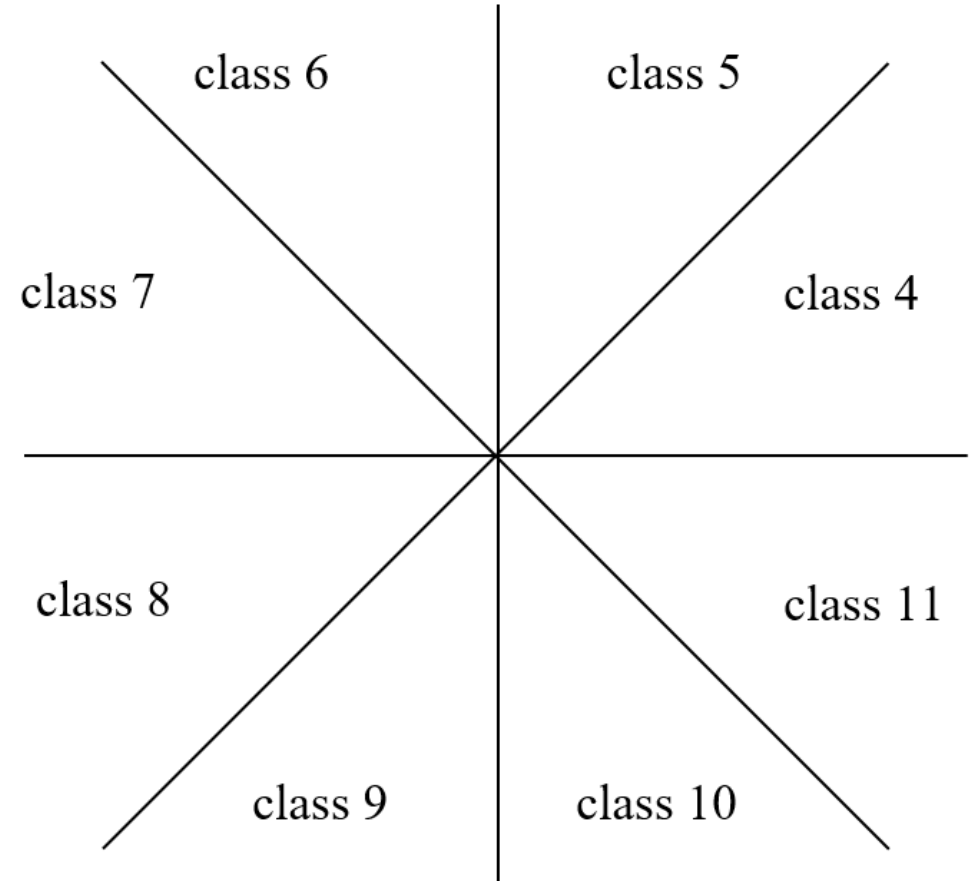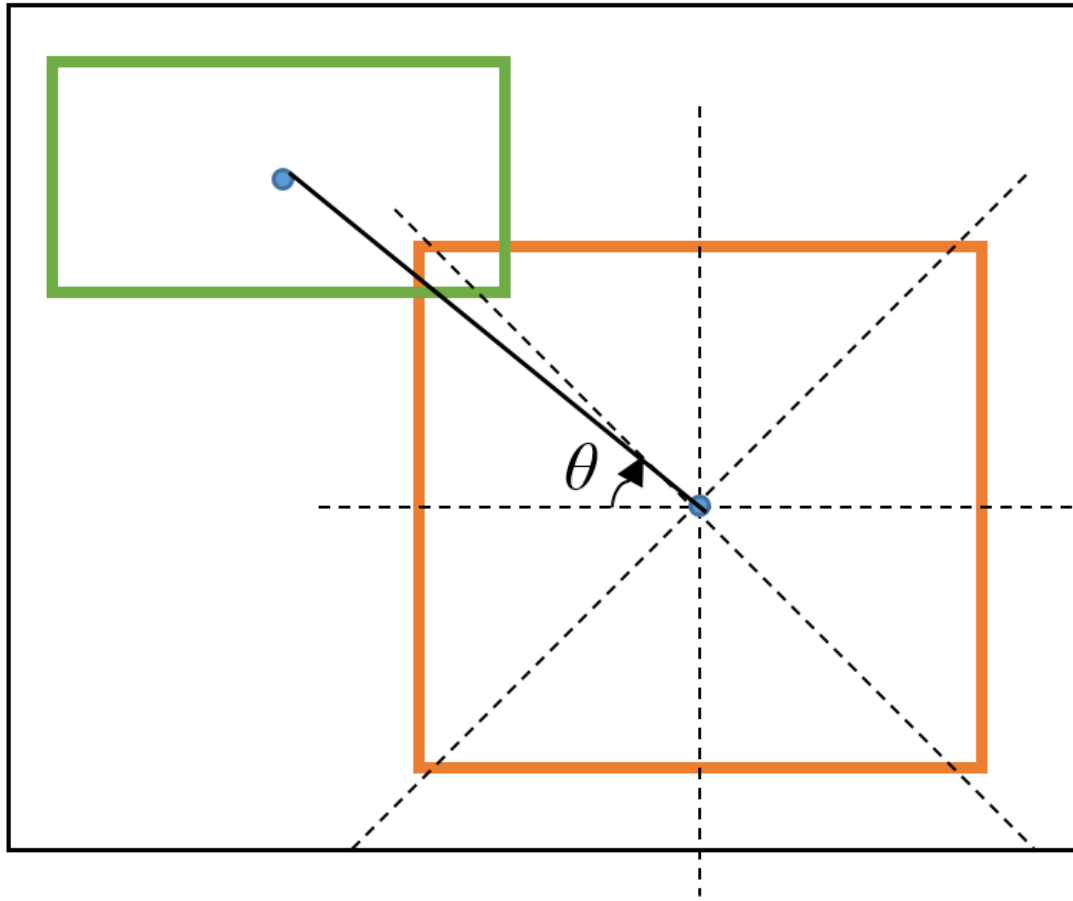Brief description of TextVQA problem, and an illustration of our MCG model structure, which contains a GNN-based contextual information propagation mechanism.

# Encoding Component

- Non-textual object features are extracted with a pre-trained Faster-RCNN model.

- For the scene texts in the image, we apply scene text detector Rosetta to identify tokens in the image. We get tokens, visual bounding box, and visual feature of scene texts. The visual feature is extracted trough feeding the bounding box into the Faster-RCNN model.

- For the question, we follow the common practice as in other VQA works.

# Relation Modeling Component

- Spatial Relationship modeling:
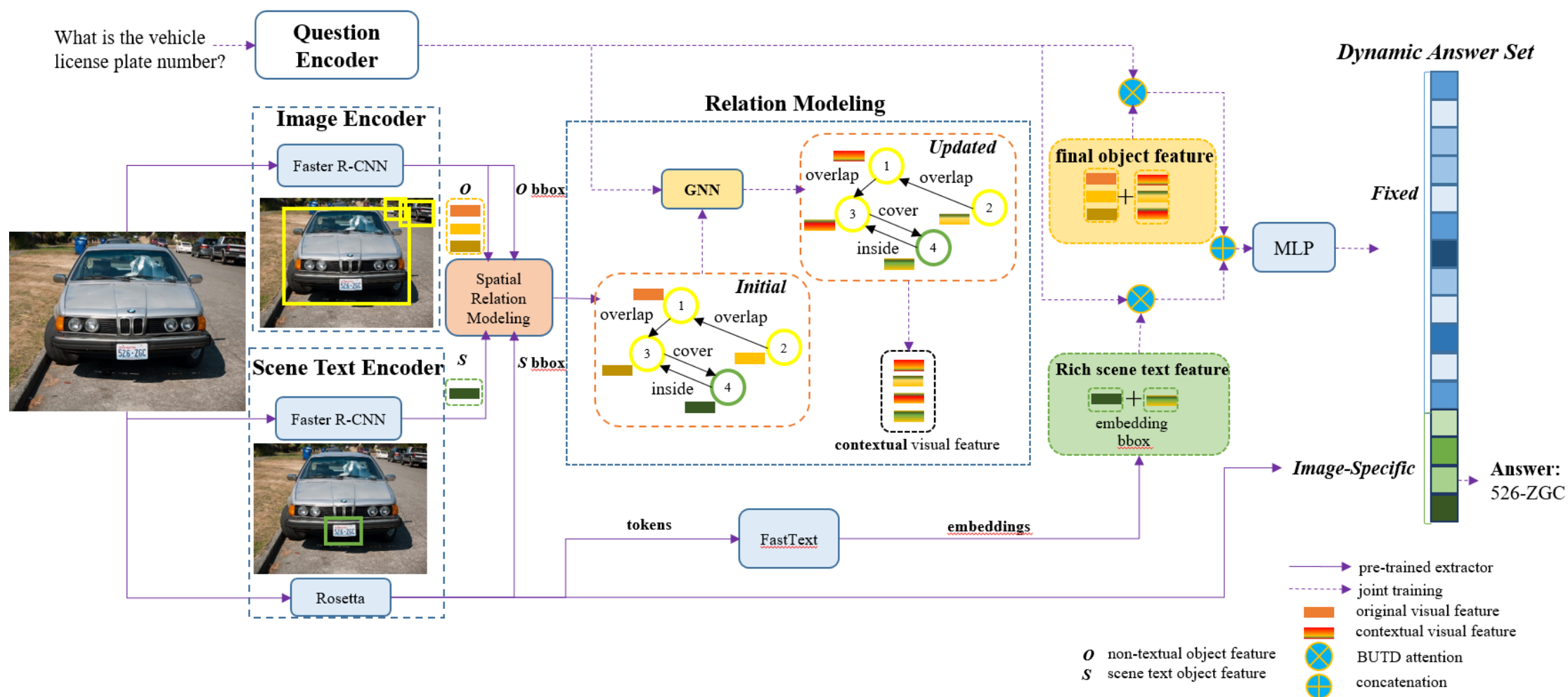
# Contextual GNN Propagation Mechansim

$$v_i^{(q)} = \sigma\left(q \cdot W_q v_i\right), \quad i = 1, 2, \cdots, K + M,$$

$$v_i^{h+1} = \sigma\left(\sum_{j \in \mathcal{N}_i} \alpha_{ij} \cdot W_h v_j^h\right),$$

$$\alpha_{ij}^l = \frac{\exp\left(\left(U^l v_i^h\right)^\top \cdot V^l v_j^h\right)}{\sum_{j \in \mathcal{N}_i} \exp\left(\left(U^l v_i^h\right)^\top \cdot V^l v_j^h\right)}, \qquad v_i^{h+1} = \Big\|_{l=1}^{L} \sigma\left(\sum_{j \in \mathcal{N}_i} \alpha_{ij}^l \cdot W_h^{\ l} v_j^h\right)$$

# Model Architecture

# Results

| Model | Object Combine | OCR Combine | No.of GNN Layer | Rich OCR Feature | Acc. on Val | Acc. on Test |
|---|---|---|---|---|---|---|
| LoRRA [29] | – | – | – | – | 26.56% | 27.63% |
| MCG(max-pooling) | – | – | 1 | yes | 17.85% | 17.34% |
| MCG | residual | residual | 1 | yes | 29.29% | 29.29% |
| MCG | 2 att. | concat. | 1 | yes | 27.68% | 27.91% |
| MCG | 2 att. | residual | 1 | no | 27.81% | 27.98% |
| MCG | 2 att. | residual | 2 | yes | 28.71% | 29.06% |
| MCG | 2 att. | residual | 1 | yes | **29.40%** | **29.61%** |

# Results-Qualitive



What is the **name** of the **hotspot**?
**LoRRA:** gates
**MCG:** vodafone

What **company** is on the **advert**?
**LoRRA:** zemel
**MCG:** nationwide

What kind of **gps logger** is it?
**LoRRA:** peceoi
**MCG:** wireless

What **brand** is the **yellow box**?
**LoRRA:** eauking
**MCG:** triscuit

How much **time** is left on the **washing machine**?
**LoRRA:** 0
**MCG:** 120

What **city** is **named**?
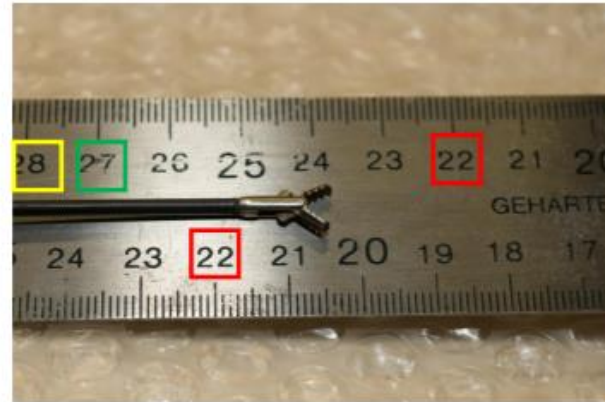**LoRRA:** new york
**MCG:** martinborough

# Results-Faulty



How many **way stop** is this **sign** for?

**LoRRA:** 3
**MCG:** all
**Human:** 4

What is the **largest number** on the **top row** of this **ruler**?

**LoRRA:** 22
**MCG:** 27
**Human:** 28

What does it say in **blue**?

**LoRRA:** kullik
**MCG:** ilihakvik
**Human:** kullik ilihakvik