## **Paper ID:1025**

# Multi-modal Contextual Graph Neural Network for Text Visual Question Answering

Yaoyuan Liang, Xin Wang, Xuguang Duan\*, Wenwu Zhu\*



## **Illustration of TextVQA and MCG MADA**

Embedding



Question

**Question:** What does the license plate say on the white car?

Scene Texts: Left: HER BWER Right: HIS BWER

Answer: HER BWER



## **Visual and Textual Feature Extractor**

Textual Object: Rosetta and Faster-RCNN ROI pooling. Non-textual Object: Faster-RCNN



# **Spatial Relationship Modelling**

#### class 6

class 5



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



**Overall Model Architecture and GNN Propagation Mechanism** 





OVERALL MODEL PERFORMANCE COMPARISION. THE VALIDATION SET ACCURACY (VAL)IS COMPUTED LOCALLY, WHILE THE TEST SET ACCURACY (TEST) IS OBTAINED THROUGH THE ONLINE JUDGING SYSTEM.

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%

Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vga models that can read. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 8317–8326, 2019.

Reference

Linjie Li, Zhe Gan, Yu Cheng, and Jingjing Liu. Relation-aware graph attention network for visual question answering. In Proceedings of



LoRRA: gates

MCG: vodafone

What kind of gps logger is it?

LoRRA: peceoi

MCG: wireless

MARTINBOROUGH WIN

What company is on the advert? What is the name of the hotspot? LoRRA: zemel MCG: nationwide



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



Qualitative examples from our MCG model on TextVQA test set

How many way stop is this sign for? LoRRA: 3 MCG: all Human: 4



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

LoRRA: kullik MCG: ilihakvik Human: kullik ilihakvik

Faulty examples from our MCG model on TextVQA test set. We can infer that previous work LoRRA and our MCG model are weak in 2 aspects:

- 1) do not have the ability in predicting answers that require more than 1 token.
- 2) unable to split extracted OCR tokens according to the semantic clue given in question.

the IEEE International Conference on Computer Vision, pages 10313–10322, 2019.



![](_page_0_Picture_43.jpeg)

lyy8ztc@outlook.com