

The 17th Conference of the European Chapter  
of the Association for Computational Linguistics

## Towards a Unified Model for Generating Answers and Explanations in Visual Question Answering



Chenxi Whitehouse, Tillman Weyde and Pranava Madhyastha ([chenxi.whitehouse@city.ac.uk](mailto:chenxi.whitehouse@city.ac.uk))

City, University of London

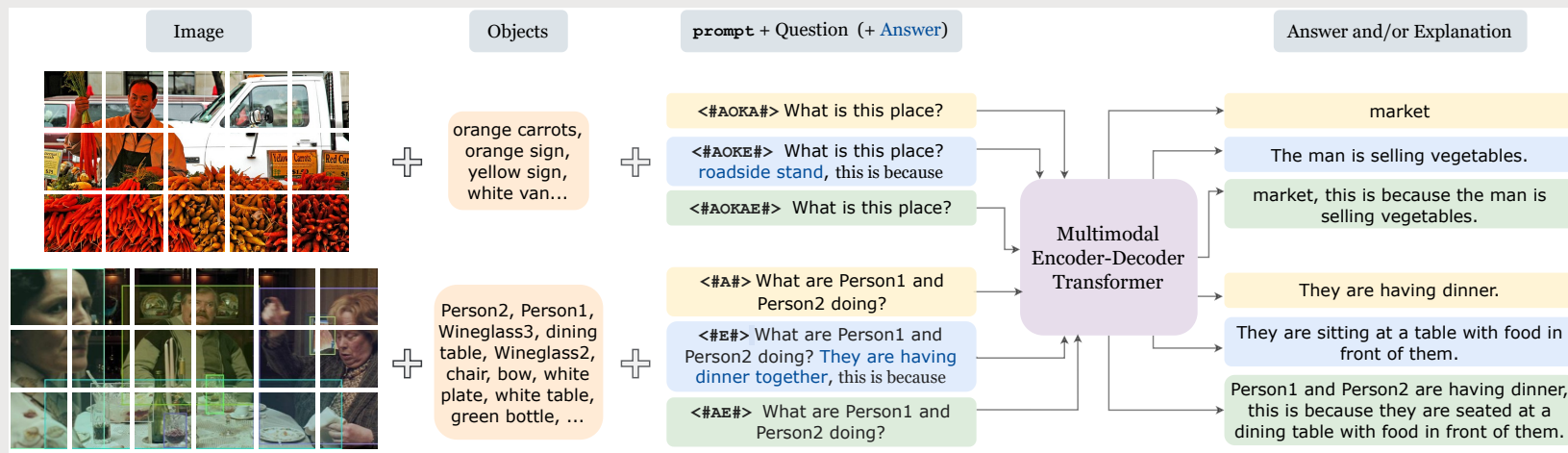
May 2023

# Motivation and Contribution

- Background:
  - Providing explanations for Visual Question Answering (VQA) tasks is desirable
  - Current explanation models for VQA generally trained separately from the QA model, resulting in less grounded answer and explanation
- Our proposal:
  - UMAE: A **unified model** for **answer** and **explanation** generation
  - Multitask learning with single artificial prompt tokens to distinguish tasks while joint training
  - Use perplexity as criteria to map open-ended generations to multiple-choice options
  - SOTA explanation generation scores and promising out-of-domain performance on VQA

# UMAE Illustration

- Train a multimodal encoder-decoder model on mix of VQA tasks for jointly optimizing answer & explanation
- Distinguish the training instances and target output with artificial prompt tokens (e. g.  $\langle \#A\# \rangle$ ,  $\langle \#E\# \rangle$ ).
- Top and bottom examples are from A-OKVQA ([Schwenk et al., 2022](#)) and VCR ([Zellers et al., 2019](#)), respectively



# Artificial Prompt Tokens

- Add a single artificial prompt token at the beginning of the textual input to
  - Distinguish different datasets and tasks
    - `<#A#>` for generating answer, `<#AE#>` for explanations, `<#AE#>` for both answer and explanation
  - Learn the shared semantics among different tasks
- These tokens are abstract, simple yet effective
- Different from natural language prompt commonly used in seq2seq models such as T5

## Artificial prompts

`<#A#>` What are Person1 and Person2 doing?

`<#E#>` What are Person1 and Person2 doing? **They are having dinner together**, this is because

`<#AE#>` What are Person1 and Person2 doing?

## Natural Language prompts

[Provide answer and explanation] What are Person1 and Person2 doing?

# Perplexity as Multiple-Choice Metric

- Map open-ended generated text to multiple-choice options
  - Limitation of existing methods using semantic embedding similarities such as Glove
  - We instead feed the same visual and textual input to the model and calculate the perplexity of each answer being generation
  - Choose the lowest perplexity option as the final answer
  - Results in better performance than mapping with generation metrics (BLEU, BERTScore)

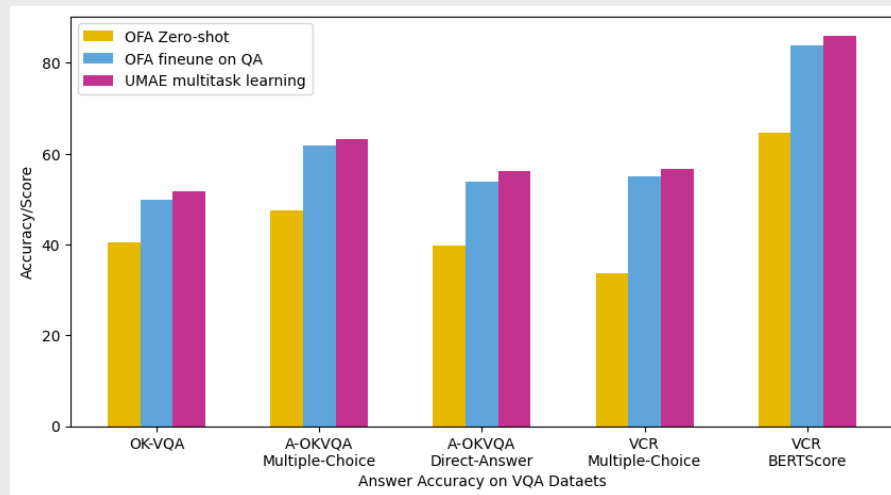
# Model and Datasets

- We built on OFA, a multimodal encoder-decoder model ([Wang et al., 2022](#))
  - Additionally extract bottom-up top-down features and attributes and feed to OFA
  - We do not use candidate answer set as OFA
- Datasets:
  - Train on three knowledge-intensive visual question answering tasks:
    - OKVQA ([Marino, et al., 2019](#), question and answer)
    - A-OKVQA (answer and explanation)
    - VCR (answer and explanation)
  - Out-of-domain evaluation on VQA-X ([Park, et al., 2018](#))

# Experimental Results

## ■ Answer Accuracy

- UMAE achieves better results than trained models separately
- Refer to the paper for more detailed scores



# Experimental Results

## ■ Explanation Performance

- OFA zero-shot not able to follow natural language instruction to generate explanations
- UMAE achieves SOTA explanation generation on A-OKVAQA, VCR and promising out-of-domain results on VQA-X

DATASET	MODEL	e-ViL SCORES			N-GRAM SCORES					LEARNT SCORE
		$S_O$	$S_T$	$S_E$	BLEU4	ROUGE-L	METEOR	CIDEr	SPICE	BERTSCORE
A-OKVQA	OFA*	4.44	56.19	7.90	0.30	4.45	3.26	4.82	4.62	68.64
	OFA <sub>Q→A</sub> +OFA <sub>QA→E</sub>	35.82	74.32	48.29	22.18	48.51	23.56	86.76	22.46	85.96
	UMAЕ <sub>A-OKVQA</sub>	37.10	73.97	50.15	<b>27.61</b>	52.23	24.06	<b>104.39</b>	22.88	87.86
	UMAЕ <sub>ALL</sub>	<b>37.91</b>	<b>74.59</b>	<b>50.82</b>	27.35	<b>52.56</b>	<b>24.83</b>	101.09	<b>23.33</b>	<b>88.21</b>
VCR	e-UG	19.30	<b>69.80</b>	27.60	4.30	22.50	11.80	32.70	12.60	79.00
	UMAЕ <sub>VCR</sub>	<b>22.57</b>	<b>56.68</b>	39.82	12.25	28.87	16.67	<b>48.14</b>	27.36	81.77
	UMAЕ <sub>ALL</sub>	<b>22.82</b>	56.66	<b>40.27</b>	<b>13.44</b>	<b>29.53</b>	<b>17.54</b>	47.33	<b>26.45</b>	<b>81.91</b>
VQA-X	e-UG	36.50	80.50	45.40	23.20	45.70	22.10	74.10	20.10	87.00
	UMAЕ <sub>ALL</sub>	31.58	77.65	40.67	14.63	35.12	20.29	50.35	19.13	85.40

Table 2: Explanation Scores. OFA\* is the pretrained OFA, showing the transferability of OFA for generating explanations with natural language instructions. Results with e-UG are from [Kayser et al. \(2021\)](#). We show the best results of A-OKVQA and VCR in bold. The last row in blue shade shows *out-of-domain* performance.



# Conclusion

- Jointly optimising answer and explanation improves quality in both in VQA
- Artificial prompt tokens is a simple and effect addition to the training data to boost multitask learning
- Perplexity as multiple-choice options metric outperform other metrics based on evaluating similarities
- We also discuss dataset quality limitation in the paper

An aerial photograph of London, England, with a semi-transparent red overlay. The city's dense urban landscape is visible, including the London Eye, the Shard, and various other skyscrapers and buildings. The text "Thank you!" is centered in white.

**Thank you!**