

---

# A Study on Multimodal and Interactive Explanations for Visual Question Answering

Kamran Alipour,<sup>1</sup> Jurgen P. Schulze,<sup>1</sup> Yi Yao,<sup>2</sup> Avi Ziskind,<sup>2</sup> and Giedrius Burachas<sup>2</sup>

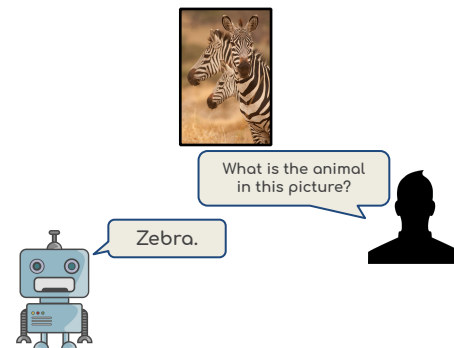
<sup>1</sup>UC San Diego

<sup>2</sup>SRI International

# Introduction

## Visual Question Answering

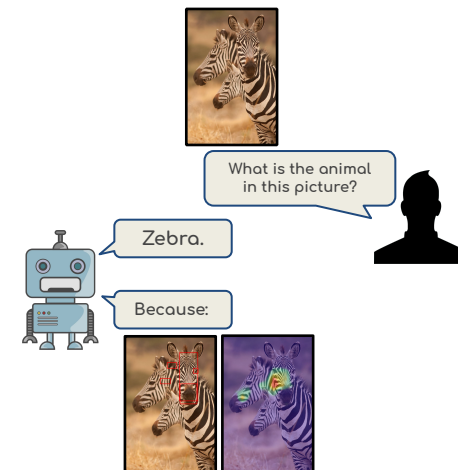
- Answering a question in natural language regarding a given image
- Attention-based models on image features influenced by the question to produce answers



## Introduction

### eXplainable AI (XAI)

- Used for years in medicine, robotics, ...
- Mainly attention-based
- Various formats: visual, textual, ...

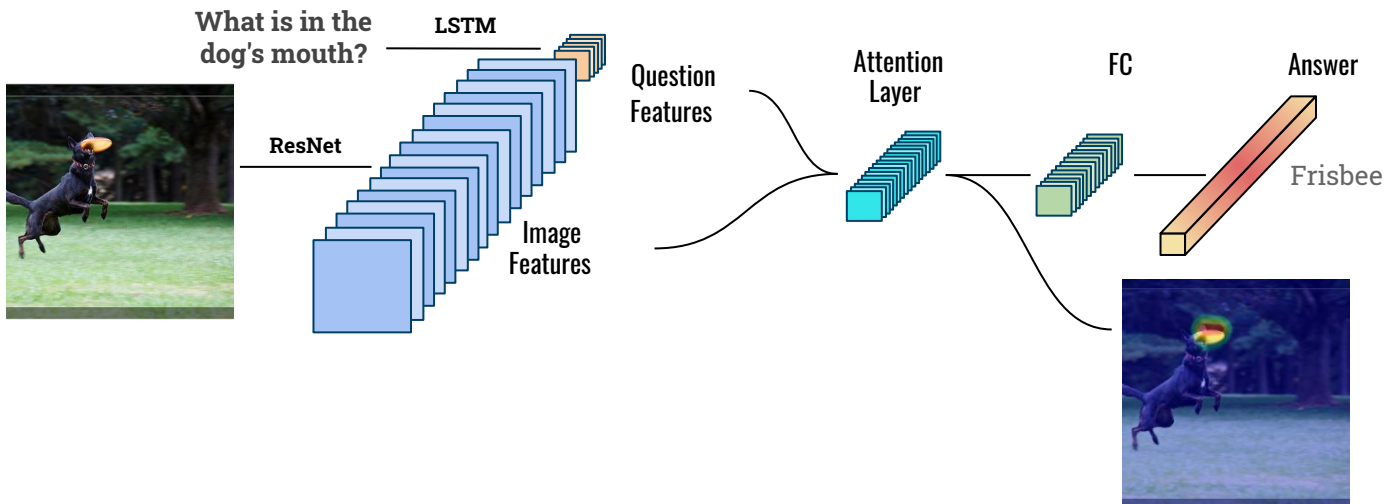


A decorative horizontal bar consisting of a teal segment on the left and an orange segment on the right.

## Contributions

- Generated multiple explanation modes from attention features and annotations
- Introduced an interactive explanation mode
- Prediction task: A novel assessment of the efficacy of explanations
- A user study to evaluate our XAI system

# eXplainable Visual Question Answering (XVQA)



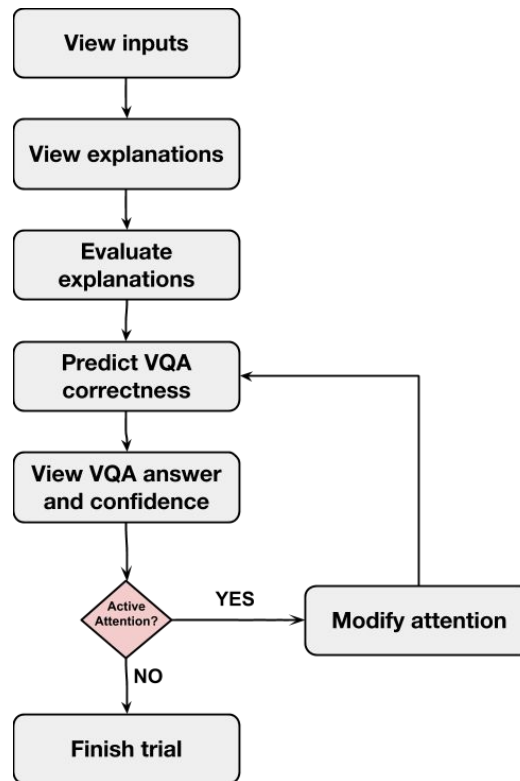
# User Study

## Input

- Images and Questions
- Explanations
- VQA correctness and confidence

## TASK:

- Rank explanations
- Predict if VQA answers correctly based on explanations.





## User Study

<b>Subjects:</b>	90
<b>Trials:</b>	10513

<b>Group</b>	
<b>NE</b>	Control group
<b>SP</b>	Spatial attention
<b>SA</b>	Active (steerable) attention
<b>SE</b>	Semantic
<b>OA</b>	Object attention
<b>AL</b>	All explanations

---

# Explanation modes

## Spatial attention

Shows the parts of the image the model is focused on while preparing the answer.

**QUESTION:** Where is this place?

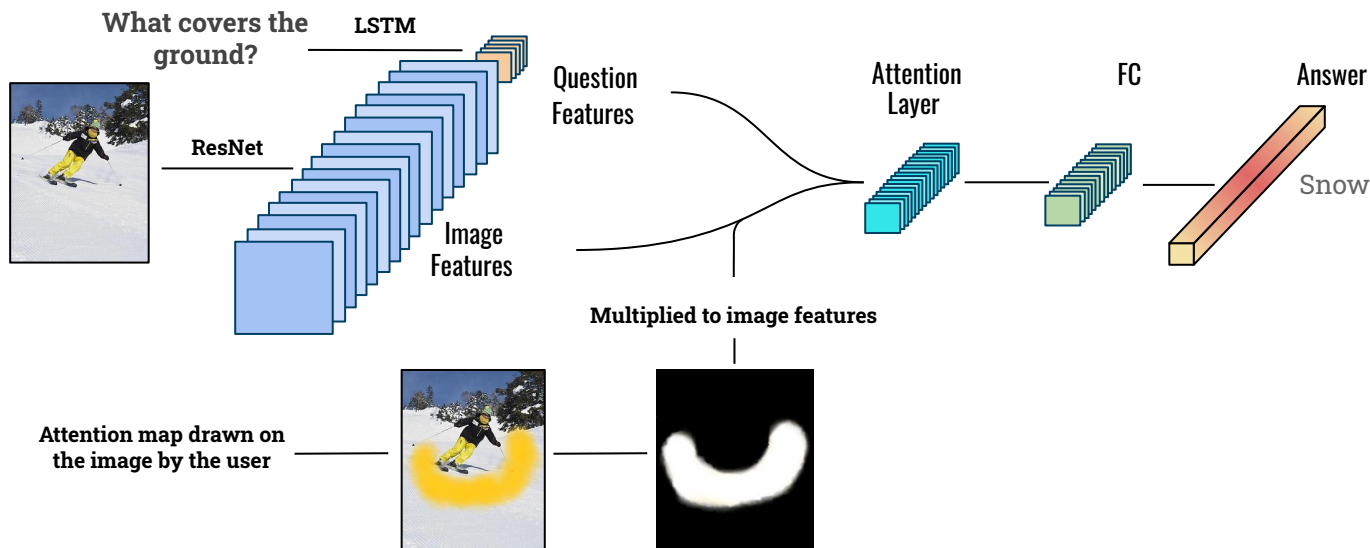
**ANSWER:** Airport.





# Explanation modes

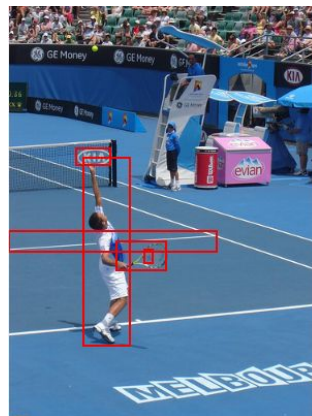
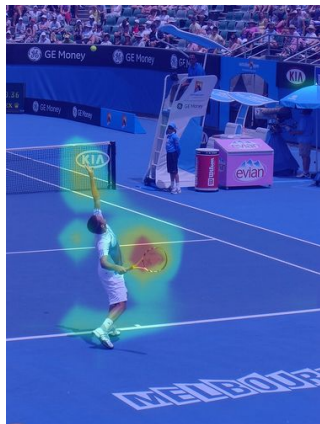
## Active (steerable) attention



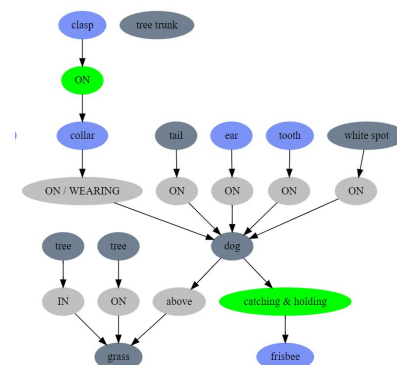
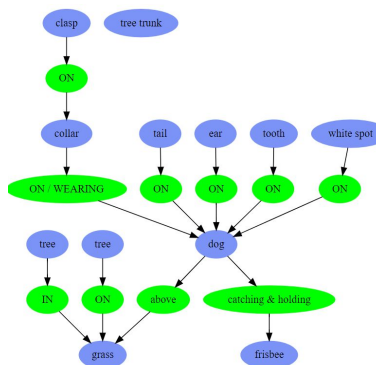
# Explanation modes

Semantic:

Bounding box



Scene graph



# Explanation modes

## Semantic:

**Textual Explanation** [Ghosh et al. 2019]

**QUESTION:** What is this sport?

**ANSWER:** Soccer.

Explanations:

Because the image contains: **white line** on **soccer field**

Because the image contains: **man** on **soccer field**

Because the image contains: **soccer shorts** with **numbers**

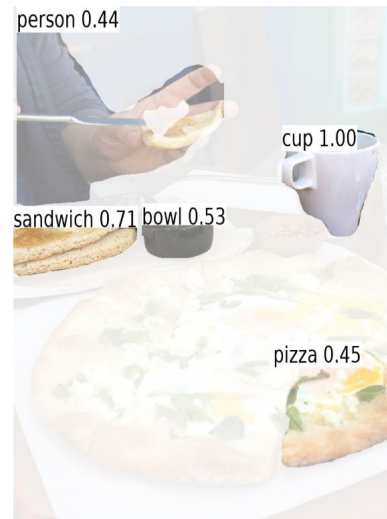


# Explanation modes

**Object attention** [Ray et al. 2019]

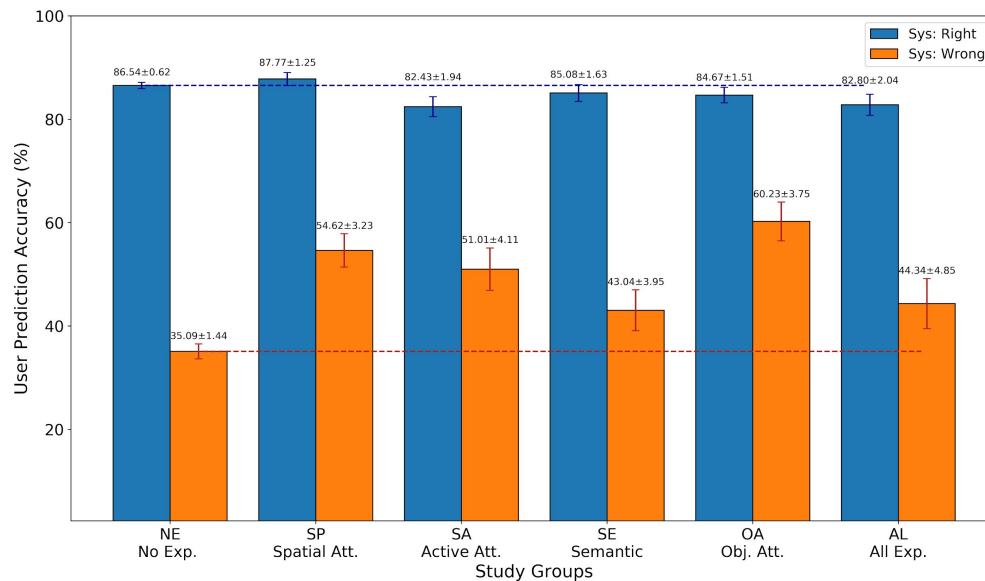
**QUESTION:** What food is he eating?

**ANSWER:** Sandwich.





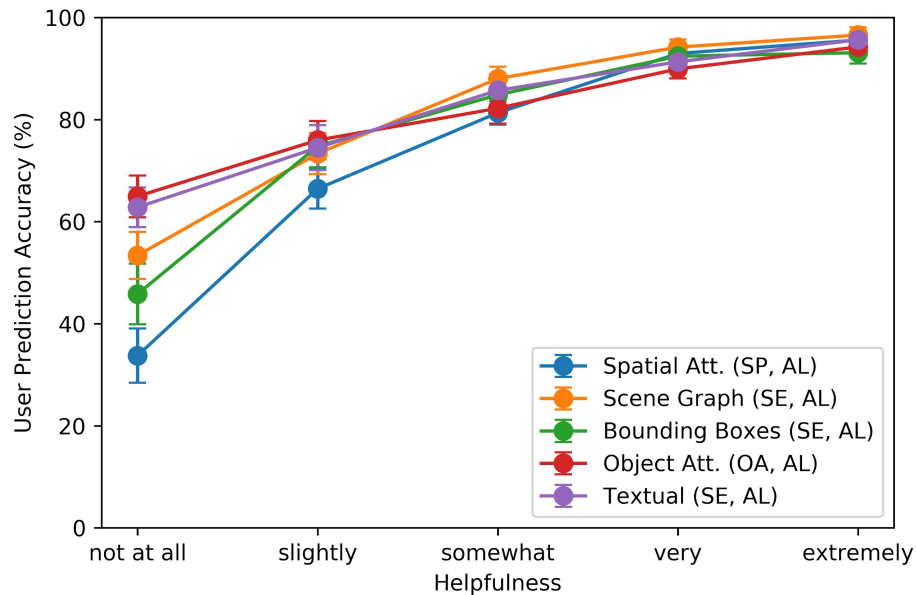
## Results



Prediction accuracies in different study groups



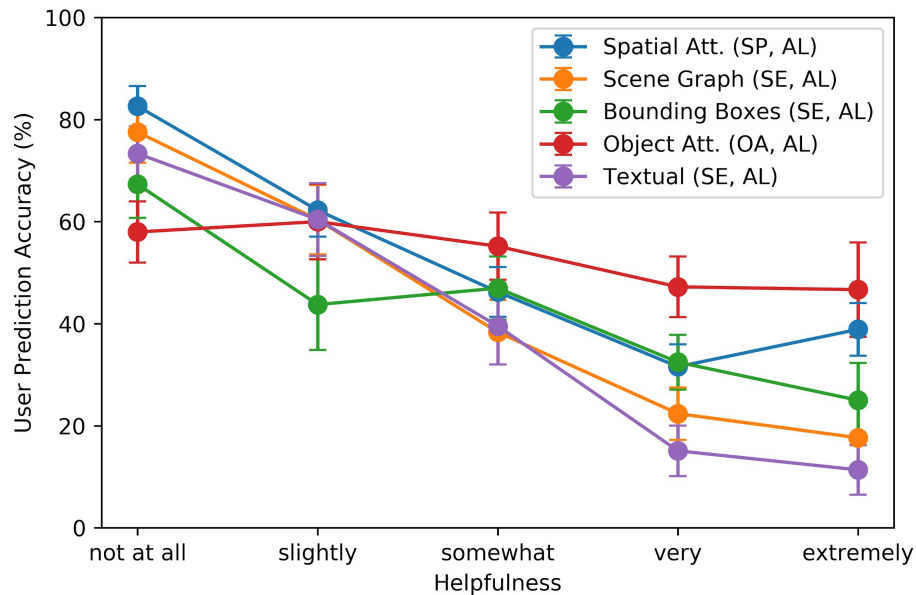
## Results



Prediction accuracy vs. explanation ratings when system **accurate**



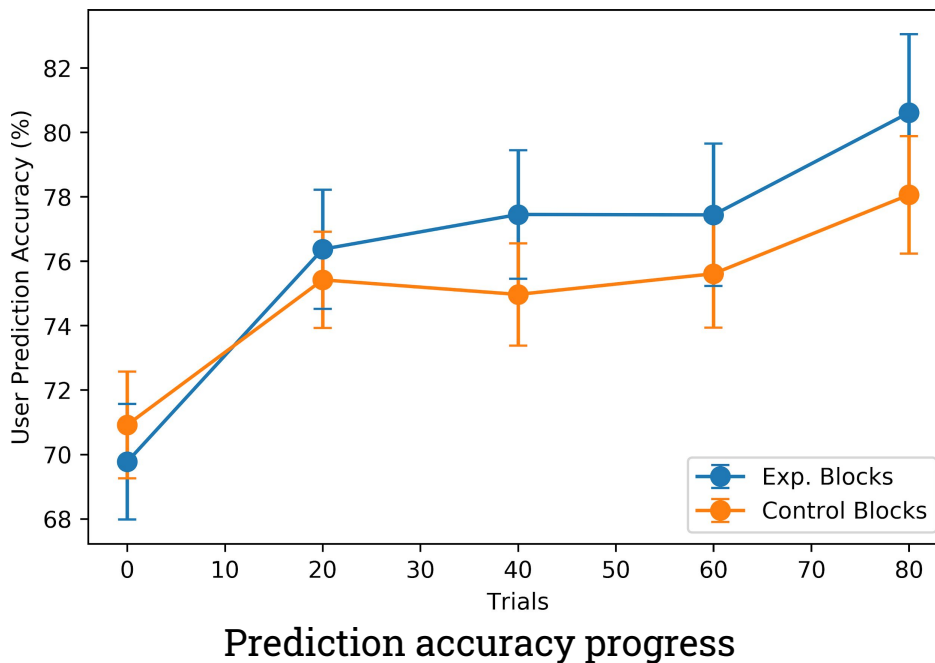
## Results



Prediction accuracy vs. explanation ratings when system **inaccurate**



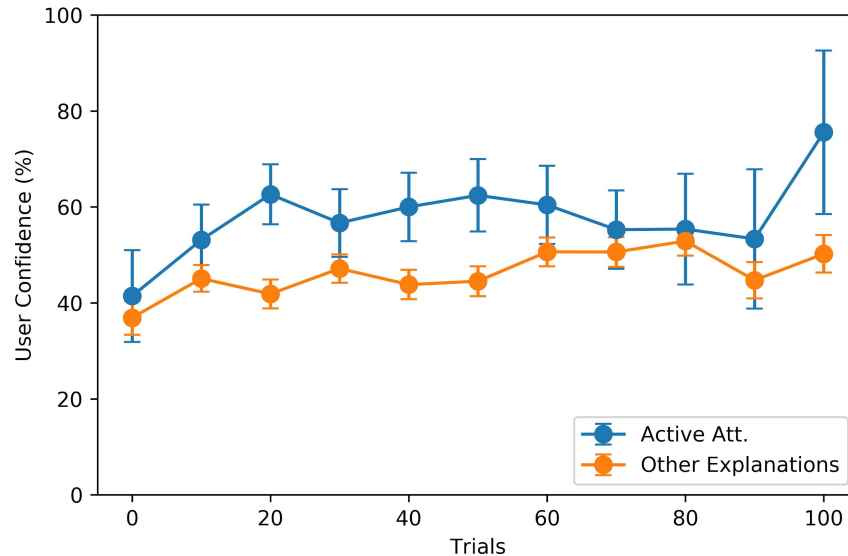
# Results







## Results



User confidence progress: Active attention vs. other explanations

A decorative horizontal bar consisting of a teal segment on the left and an orange segment on the right.

## Discussion

- Interactive explanations improved users confidence compared to other explanations
- When AI is inaccurate, explanations are more helpful
- Getting exposed to multiple explanation modes can be conflicting/overwhelming and reduce prediction accuracy

A decorative horizontal bar consisting of a teal segment on the left and an orange segment on the right.

## Conclusion and Future Work

- Interactive experiment to probe explanation effectiveness
- Explanations help accuracy prediction
- Users confidence improved when exposed to explanations



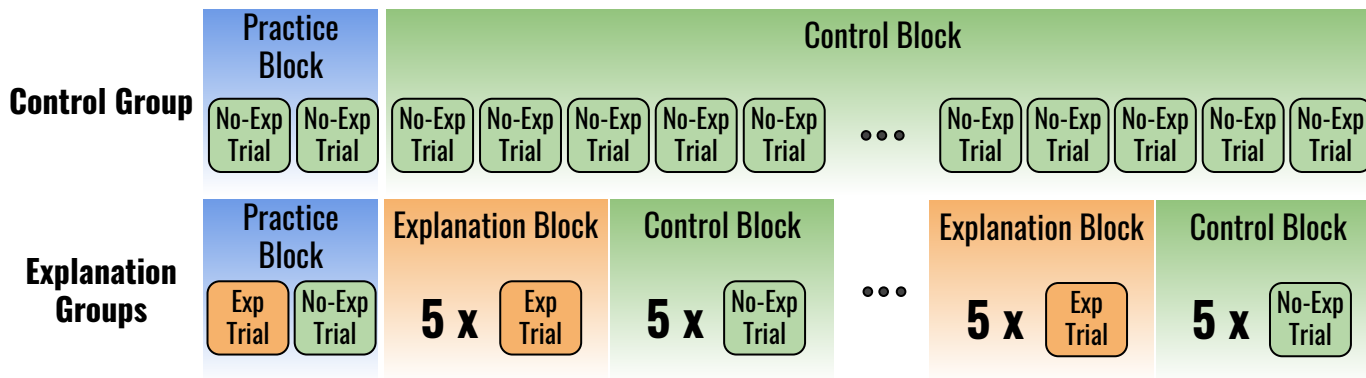
**Thank You!**

A decorative horizontal bar consisting of a teal segment on the left and an orange segment on the right.

## User Study Statistics

	<b>Group</b>	<b>Subjects</b>	<b>Trials</b>
<b>NE</b>	Control group	15	4124
<b>SP</b>	Spatial attention	15	1826
<b>SA</b>	Active (steerable) attention	15	1021
<b>SE</b>	Semantic	15	1261
<b>OA</b>	Object attention	15	1435
<b>AL</b>	All explanations	15	846
	<b>Total :</b>	90	10513

# User Study Structure



# User Study Structure

## Input

- Images and Questions
- Explanations
- VQA correctness and confidence

## TASK:

- Rank explanations
- Predict if VQA answers correctly based on explanations.

