# DSTC7-AVSD: Scene-Aware Video Dialogue Systems with Dual Attention

Ramakanth Pasunuru

www.rama-kanth.com

Mohit Bansal
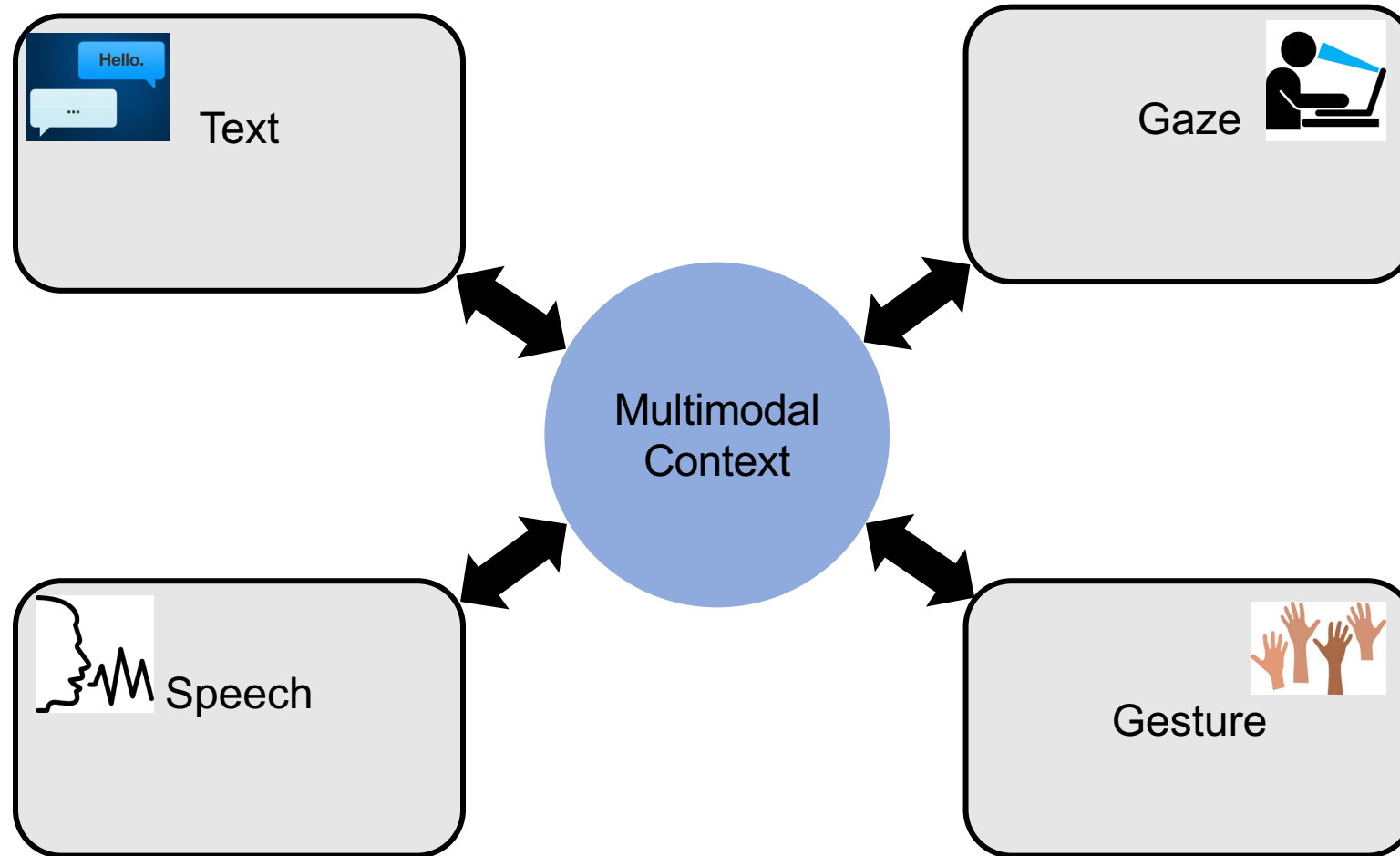
www.cs.unc.edu/~mbansal/
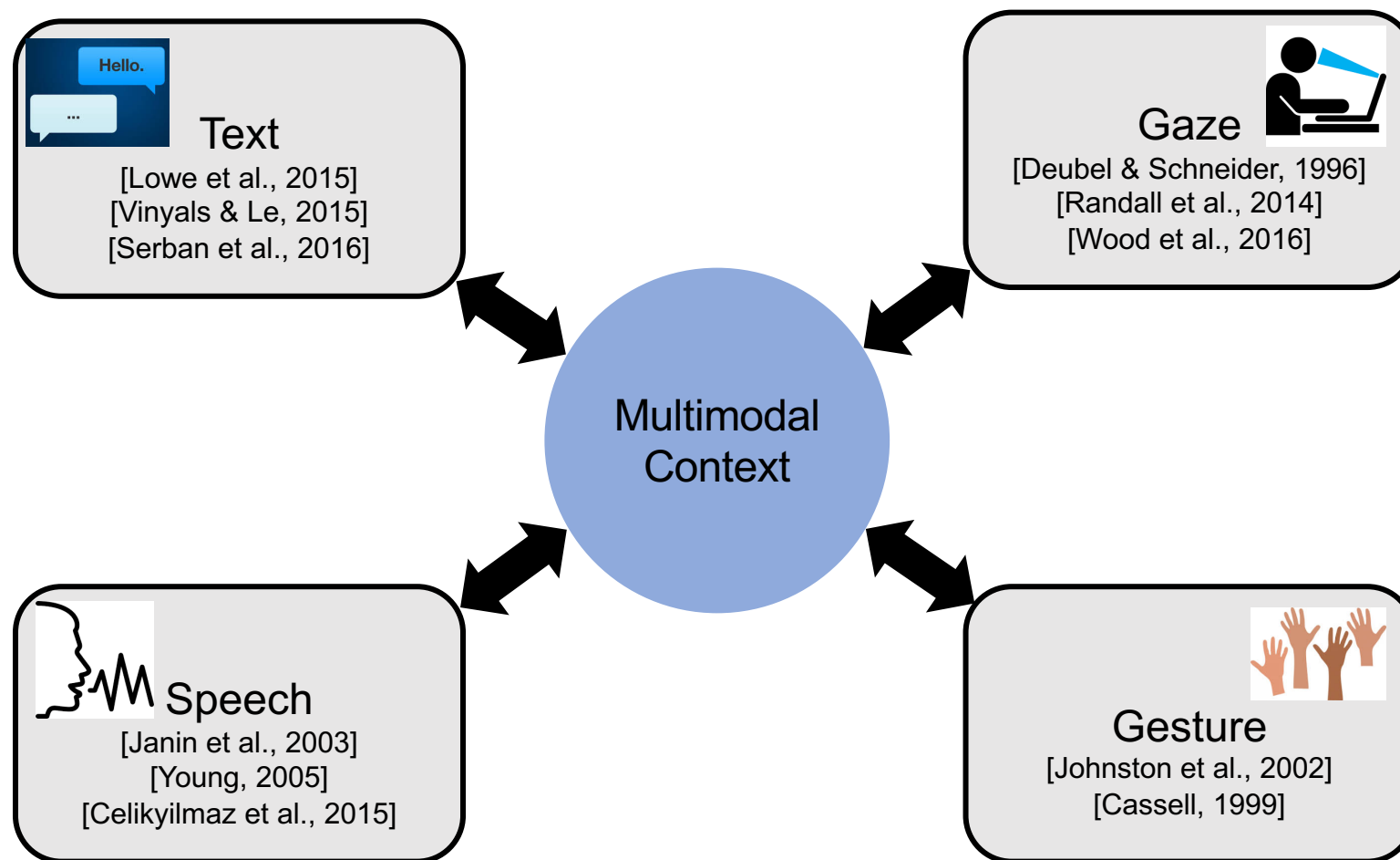
THE UNIVERSITY of NORTH CAROLINA at CHAPEL HILL

UNC NLP

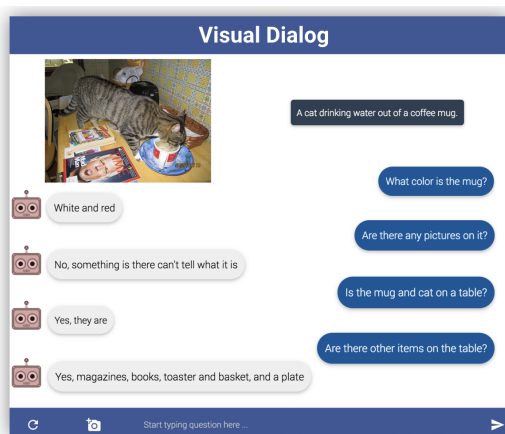# Dialogue Context

# Dialogue Context



Text
[Lowe et al., 2015]
[Vinyals & Le, 2015]
[Serban et al., 2016]

Gaze
[Deubel & Schneider, 1996]
[Randall et al., 2014]
[Wood et al., 2016]

Multimodal
Context

Speech
[Janin et al., 2003]
[Young, 2005]
[Celikyilmaz et al., 2015]

Gesture
[Johnston et al., 2002]
[Cassell, 1999]

# Visual Context



Image-based Context



[Das et al., 2017]



[De Vries et al., 2017]
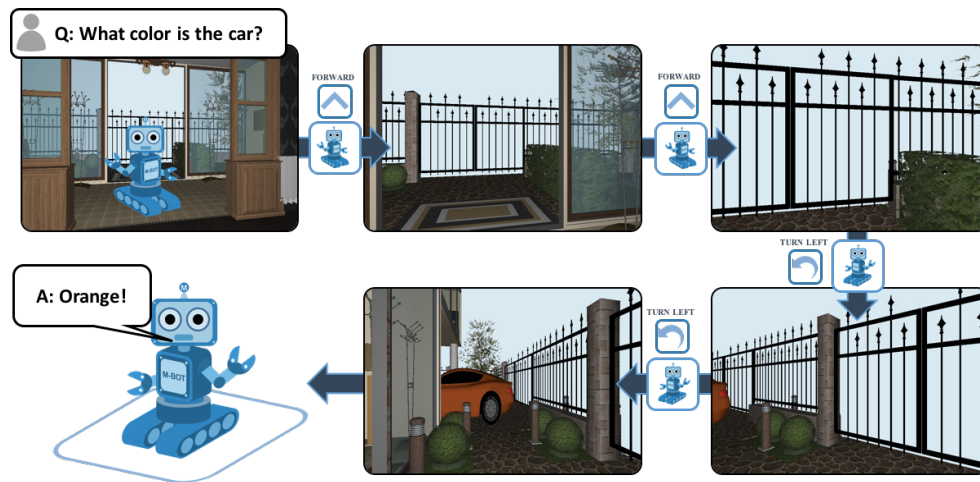


[Mostafazadeh et al., 2017]



[Celikyilmaz et al., 2014]

# Visual Context



Dynamic-Visual Context

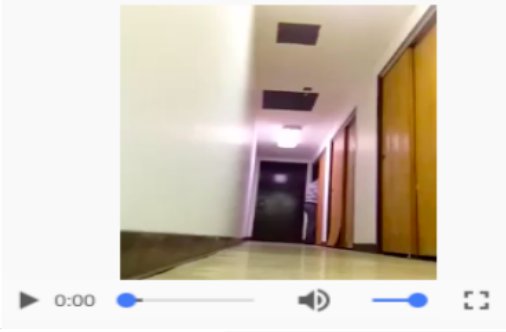Dynamic-Visual Context with Multi-Speaker

[Das et al., 2018]

[Pasunuru & Bansal, 2018]

# Visual+Audio Context



[Alamri et al., 2018]

# Visual+Audio Context

UNC NLP

**Task:**

**Input**

Question

Video
Chat History
Summary

**Output**

Answer



| | **Person A (Questioner)** |
|---|---|
| 1. | How many people are in the video? |
| 2. | Is he speaking with anyone? |
| 3. | What room is he in? |
| 4. | What is the man doing? |
| 5. | Does he start the video in the hallway? |
| 6. | Where does he put the tie and shirt? |
| 7. | Does he leave the hallway? |
| 8. | Does he open the closet door? |
| 9. | Can you tell what he grabs from the closet? |
| 10. | Is there anything else I should know? |

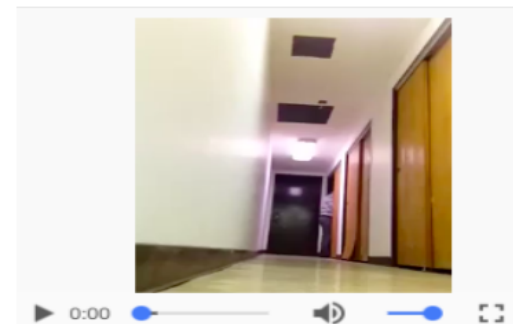| | **Persone B (Answerer)** |
|---|---|
| 1. | There is only one man in the video |
| 2. | No there is no sound |
| 3. | He is in a hallway |
| 4. | He is taking off his tie and shirt |
| 5. | Yes he does start in the hallway |
| 6. | He puts it in a closet |
| 7. | After he puts his stuff in the closet he grabs something out of the closet |
| 8. | No it is already open |
| 9. | He grabs a box and then starts walking toward the camera |
| 10. | No that is it from start to finish |

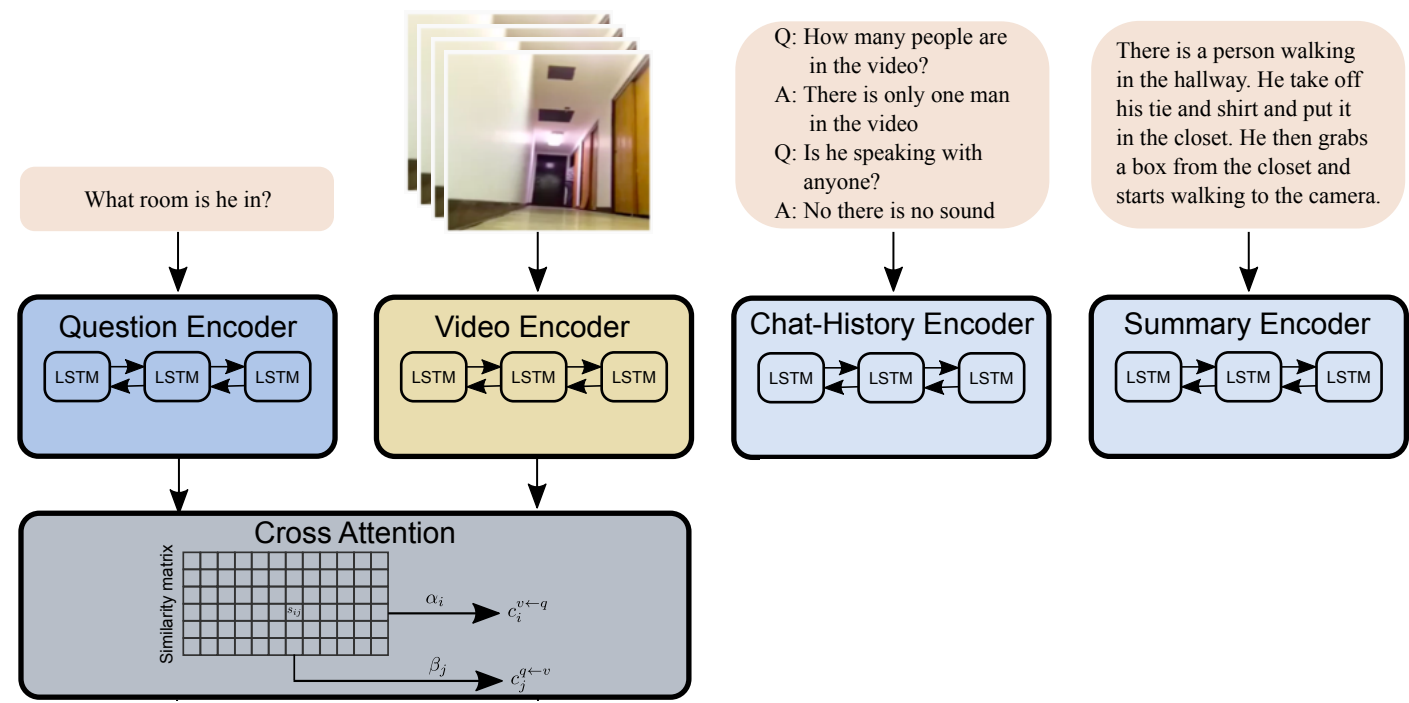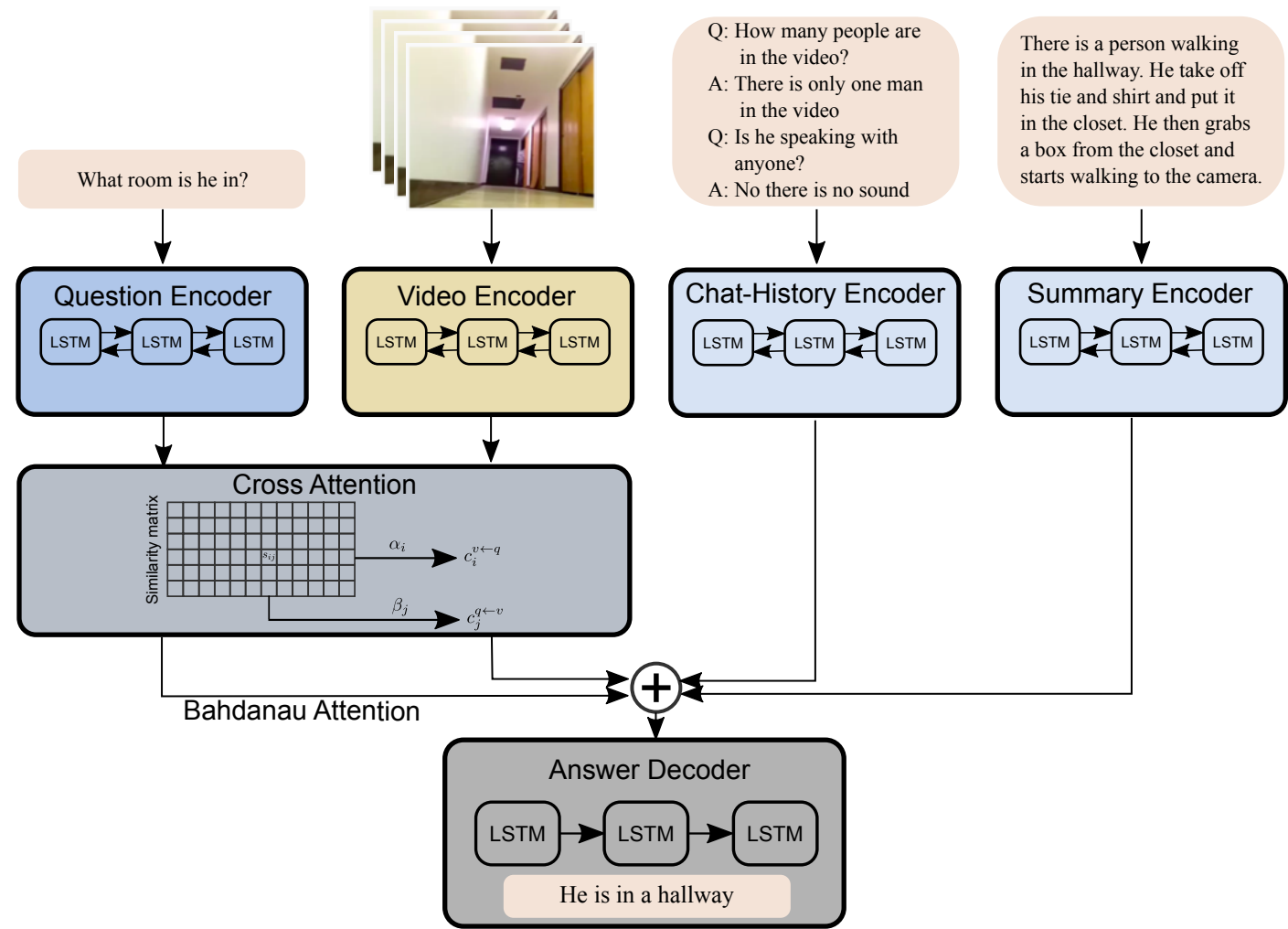[Alamri et al., 2018]

# Model

# Model

[Seo et al., 2017]

# Model



Note that we do not use audio features in our models

[Bahdanau et al., 2015; Seo et al., 2017]

# Results

| Model | METEOR | CIDEr | BLEU-4 | ROUGE-L |
|---|---|---|---|---|
| Video Only | 12.43 | 95.54 | 8.83 | 34.23 |
| Video + Chat History | 14.13 | 105.39 | 10.58 | 36.54 |
| Video + Chat History + Summary | 14.94 | 112.80 | 11.22 | 37.53 |
| Video + Chat History + Summary + Cross-attention | 14.95 | 115.82 | 11.38 | 37.87 |

Our models' performance on AVSD dataset's public test set. All of these models use the question information.

# Results

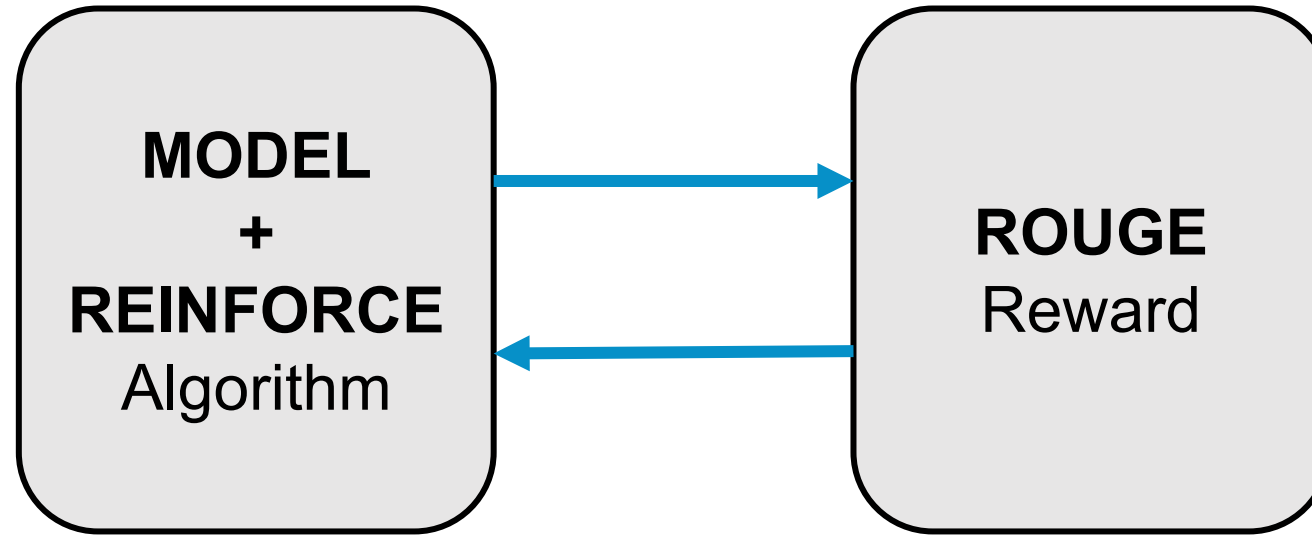| Model | METEOR | CIDEr | BLEU-4 | ROUGE-L |
|---|---|---|---|---|
| Video Only | 12.43 | 95.54 | 8.83 | 34.23 |
| Video + Chat History | 14.13 | 105.39 | 10.58 | 36.54 |
| Video + Chat History + Summary | 14.94 | 112.80 | 11.22 | 37.53 |
| Video + Chat History + Summary + Cross-attention | 14.95 | 115.82 | 11.38 | 37.87 |

Our models' performance on AVSD dataset's public test set. All of these models use the question information (no audio information).

# Other Methods

- Policy gradient based reinforcement learning

- Contextualized ELMo word embeddings

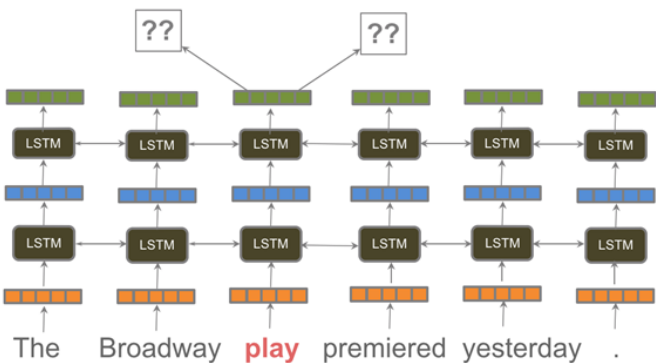- Using external data
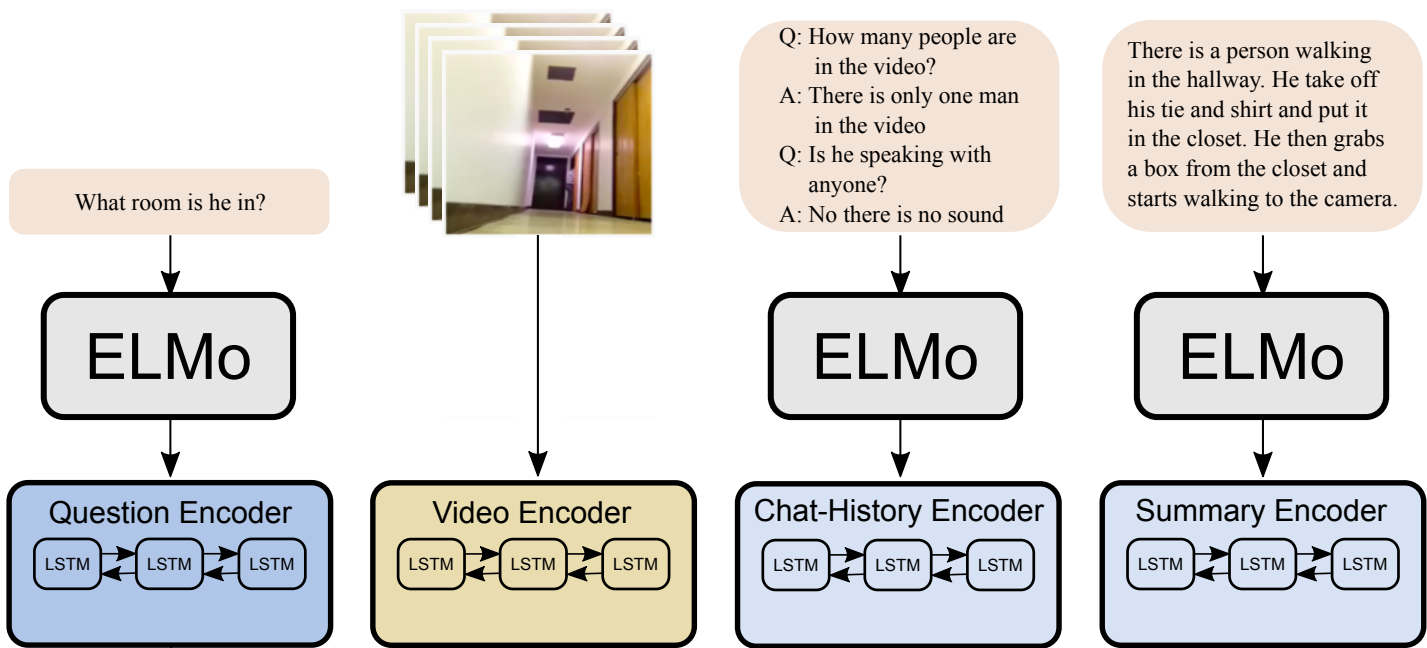
- Pointer-generator copy model

# Policy Gradients



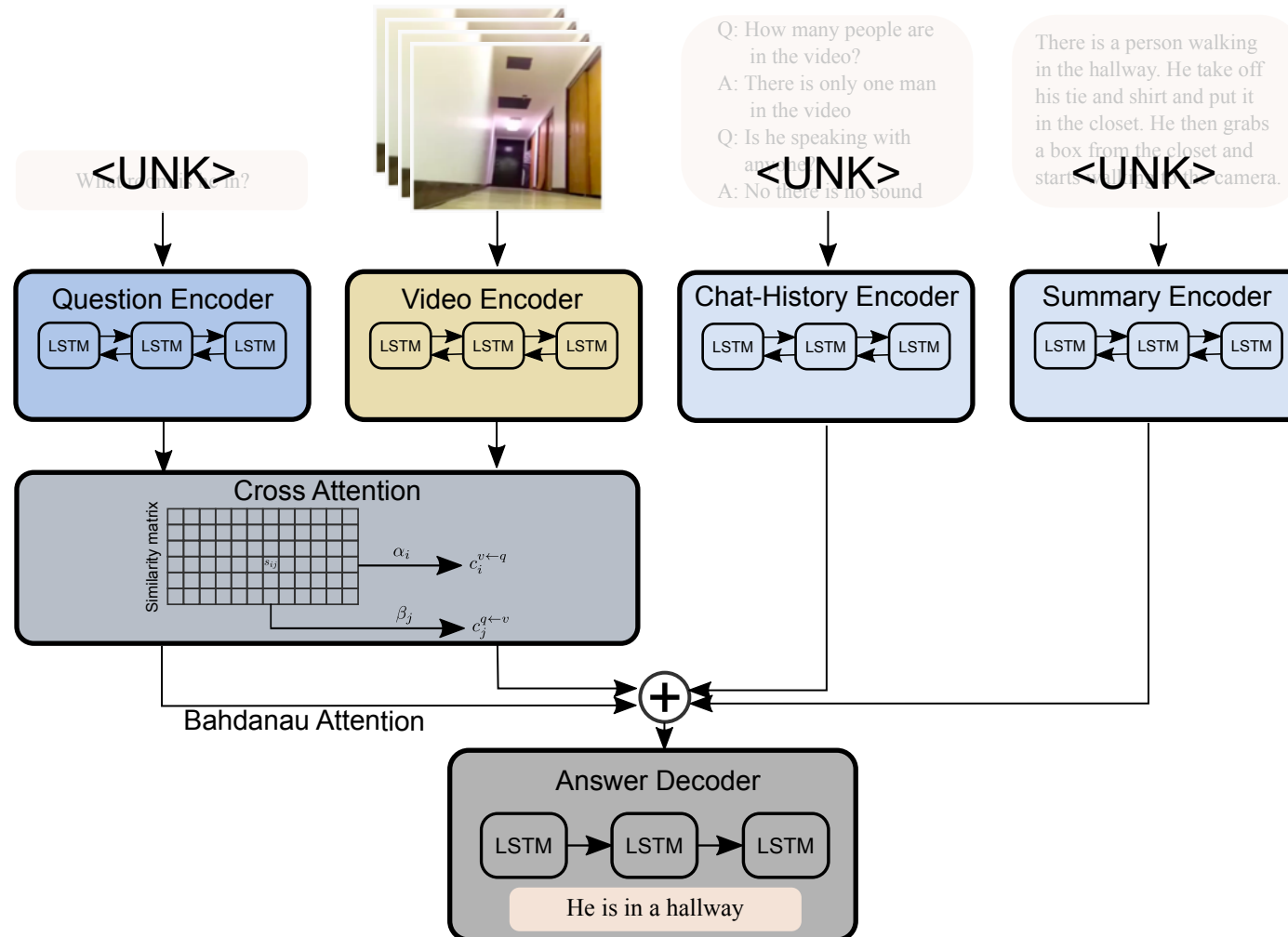$$\nabla_\theta L(\theta) = -\mathbb{E}_{w^s \sim p_\theta}\left[r(w^s) \cdot \nabla_\theta \log p_\theta(w^s)\right]$$

[Williams, 1992]

# Contextualized ELMo Word Embeddings



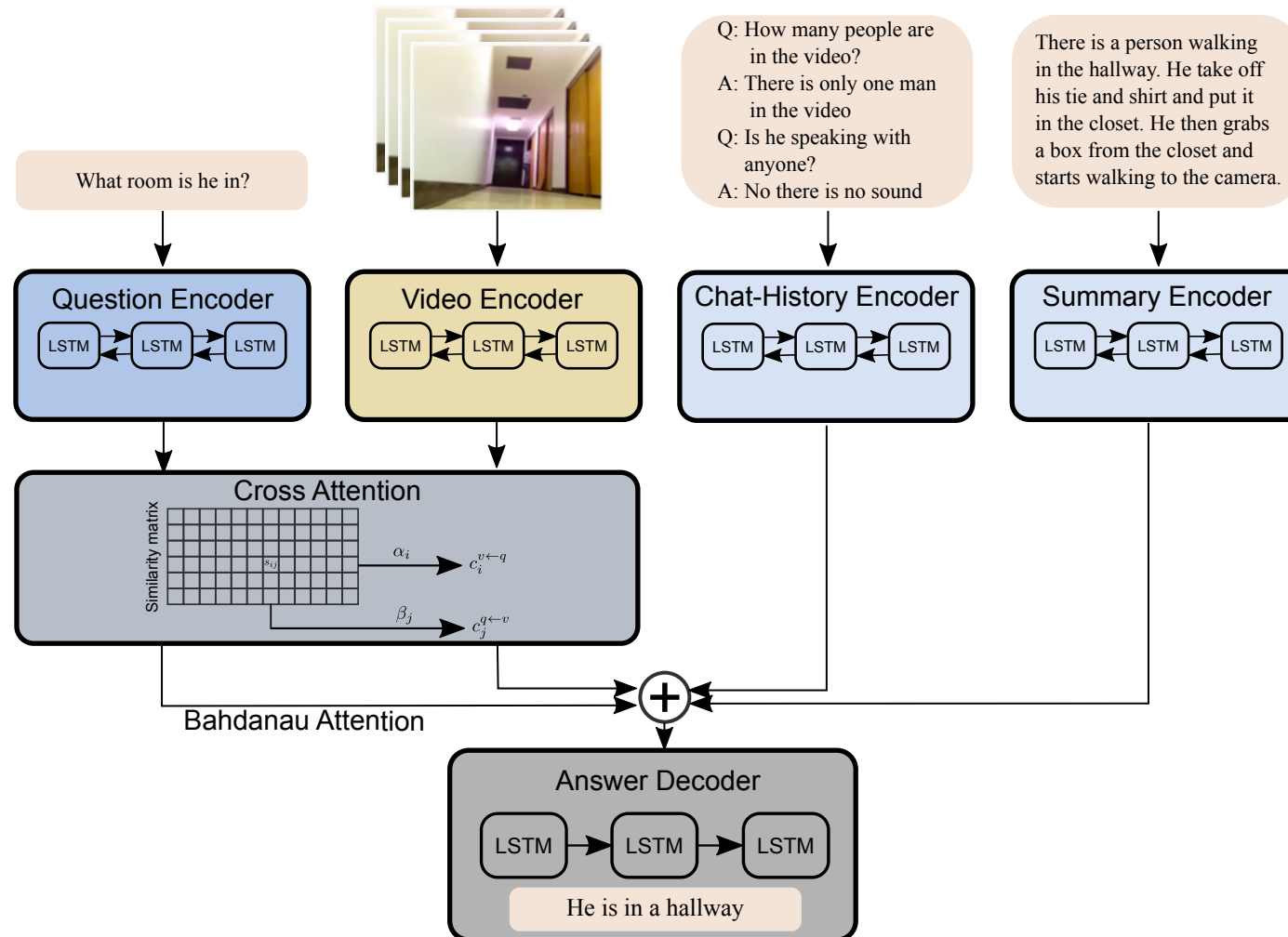[Peters et al., 2018]

# Using External Data (MSR-VTT)

[Xu et al., 2016]

# Using External Data (MSR-VTT)

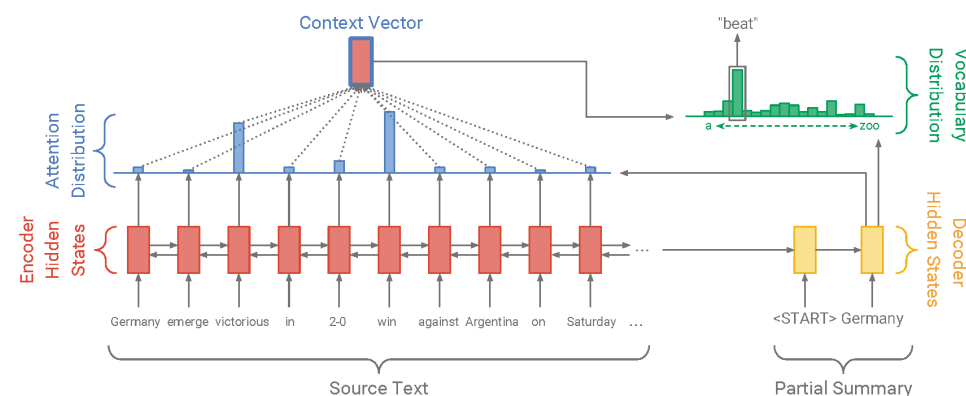[Xu et al., 2016]

# Pointer-generator Copy Model

- Lot of words in the question can also be present in the answer

- The final word distribution is a weighted combination of the vocab distribution and attention distribution

- Question-based pointer

- Joint question- and summary-based pointer



[See et al., 2017]

# Future Work

- Further analyze and improve these promising approaches with specific RL rewards, contextualized large language models, and joint copy models

- We will add Audio features to our final model

- Effective ways of extending cross-attention to multiple modalities (question+summary; question+chat-history)

# Thanks!