# Generating Question Relevant Captions to Aid Visual Question Answering

Presenter: Xiangjue Dong

# Visual Question Answering

- Top-down (Ren et al., 2015a; Fukui et al., 2016; Wu et al., 2016; Goyal et al., 2017; Li et al., 2018a)

- Bottom-up attentions (Anderson et al., 2018; Li et al., 2018b; Wu and Mooney, 2019)

- Enriching knowledge base (Li et al., 2018a; Narasimhan et al. 2018)

Jialin Wu, Zeyuan Hu, and Raymond J. Mooney. 2019.Generating Question Relevant Captions to Aid Vi-sual Question Answering.CoRR, abs/1906.00513.

# Image Captioning

- http://dbs.cloudcv.org/captioning
- Attention-based deep-learning models (Donahue et al., 2015; Karpathy and Fei-Fei, 2015; Vinyals et al., 2015; Luo et al., 2018; Liu et al., 2018)
- Encode the image using a CNN
- Build an attentional RNN, LSTM on top of the image features

Jialin Wu, Zeyuan Hu, and Raymond J. Mooney. 2019.Generating Question Relevant Captions to Aid Vi-sual Question Answering.CoRR, abs/1906.00513.

**Human Captions :**
1) A man on a blue surfboard on top of some rough water.
2) A young surfer in a wetsuit surfs a small wave.
3) A young man rides a surf board on a small wave while
a man swims in the background.
4) A young man is on his surf board with someone in the background.
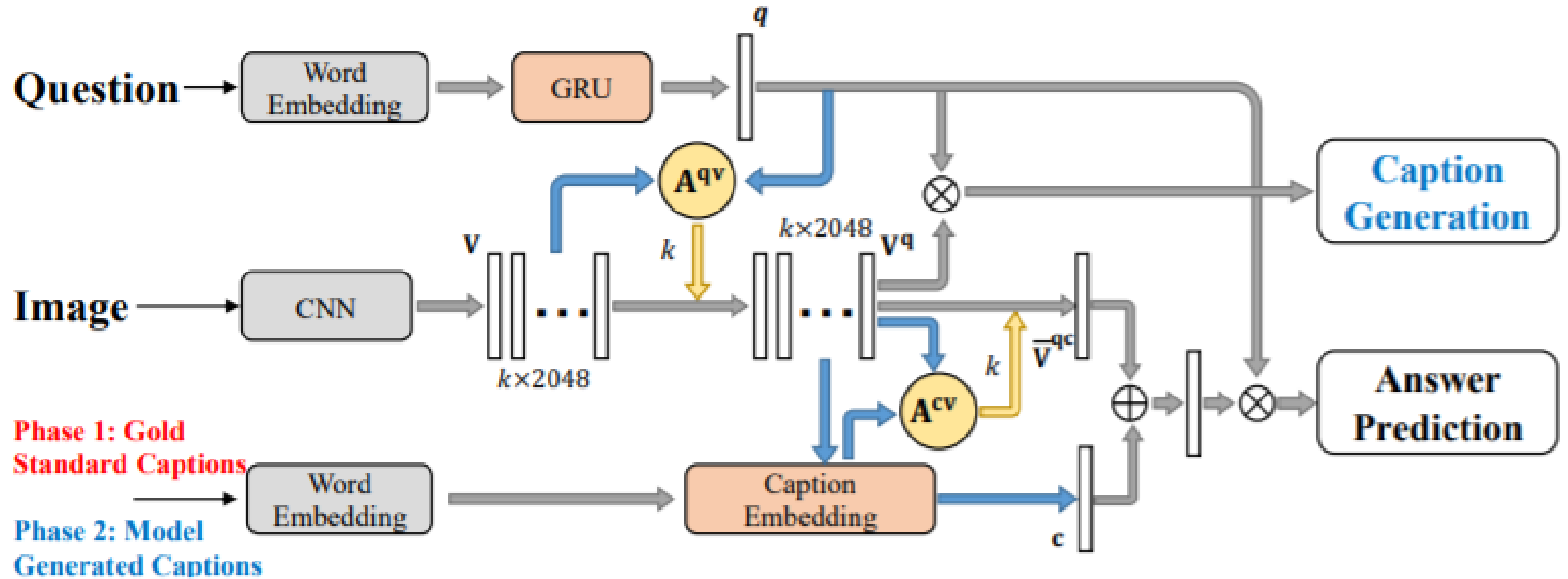5) A boy riding waves on his surf board in the ocean.

**Question 1: Does this boy have a full wetsuit on?**
Caption: A young man wearing **wetsuit** surfing on a wave.
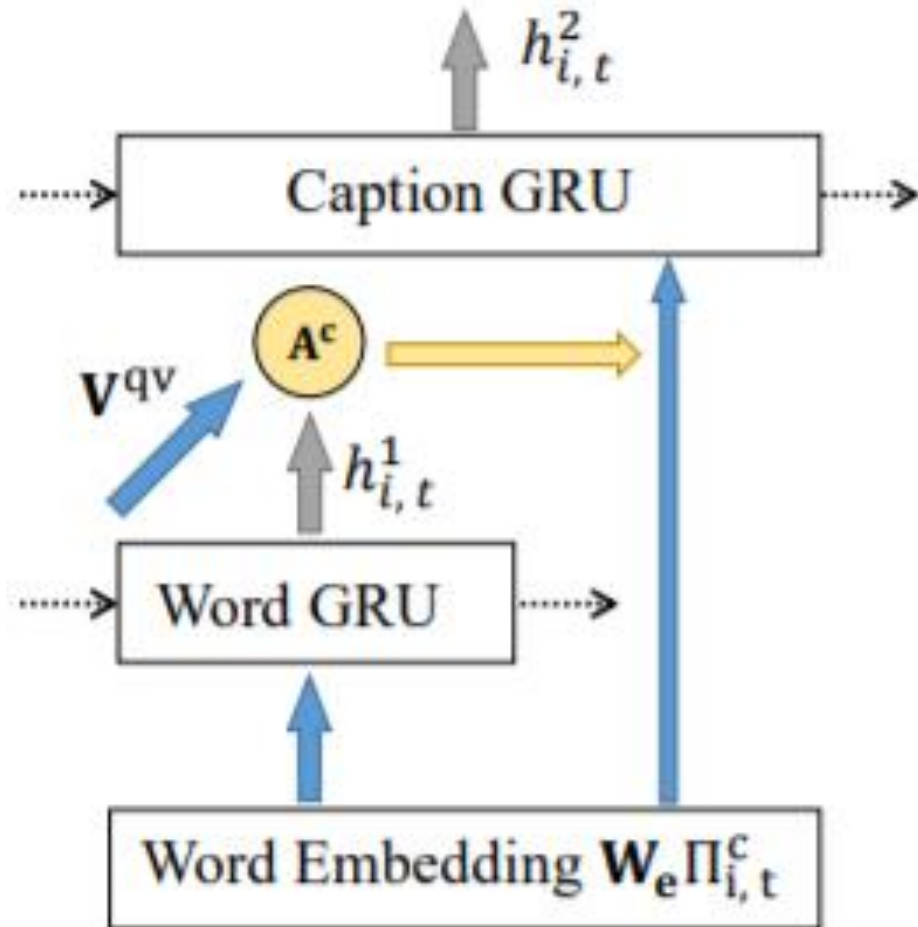**Question 2: What color is the board?**
Caption: A young man riding a wave on a **blue surfboard**.

Jialin Wu, Zeyuan Hu, and Raymond J. Mooney. 2019.Generating Question Relevant Captions to Aid Vi-sual Question Answering.CoRR, abs/1906.00513.

# Approach



Jialin Wu, Zeyuan Hu, and Raymond J. Mooney. 2019.Generating Question Relevant Captions to Aid Vi-sual Question Answering.CoRR, abs/1906.00513.

# Caption Embedding



$$h_{i,t}^1 = \text{GRU}(\mathbf{W}_e \Pi_{i,t}^c,\ h_{i,t-1}^1) \qquad (1)$$

$$h_{i,t}^2 = \text{GRU}(\alpha_{i,t}^c \mathbf{W}_e \Pi_{i,t}^c,\ h_{i,t-1}^2) \qquad (5)$$

Jialin Wu, Zeyuan Hu, and Raymond J. Mooney. 2019.Generating Question Relevant Captions to Aid Vi-sual Question Answering.CoRR, abs/1906.00513.

# Other Modules

- VQA Module

$$\mathbf{h} = \mathbf{q} \circ \left( f(\overline{\mathbf{v}}^{qc}) + f(\mathbf{c}) \right)$$

- Image Captioning Module

Jialin Wu, Zeyuan Hu, and Raymond J. Mooney. 2019.Generating Question Relevant Captions to Aid Vi-sual Question Answering.CoRR, abs/1906.00513.

# Datasets and Evaluation Metrics

- VQA Dataset:
  - VQA v2.0



What color are her eyes?
What is the mustache made of?

How many slices of pizza are there?
Is this a vegetarian pizza?

Is this person expecting company?
What is just under the tree?

Does it appear to be rainy?
Does this person have 20/20 vision?

Fig. 1: Examples of free-form, open-ended questions collected for images via Amazon Mechanical Turk. Note that commonsense knowledge is needed along with a visual understanding of the scene to answer many questions.

Jialin Wu, Zeyuan Hu, and Raymond J. Mooney. 2019.Generating Question Relevant Captions to Aid Vi-sual Question Answering.CoRR, abs/1906.00513.

# Datasets and Evaluation Metrics

- Image Captioning Dataset
  - MSCOCO 2014 dataset



The man at bat readies to swing at the pitch while the umpire looks on.

A large bus sitting next to a very tall building.

A horse carrying a large load of hay and two people sitting on it.

Bunk bed with a narrow shelf sitting underneath it.

Fig. 1: Example images and captions from the Microsoft COCO Caption dataset.

Jialin Wu, Zeyuan Hu, and Raymond J. Mooney. 2019.Generating Question Relevant Captions to Aid Vi-sual Question Answering.CoRR, abs/1906.00513.

# Results

| | Test-standard | | | |
|---|---|---|---|---|
| | Yes/No | Num | Other | All |
| Prior (Goyal et al., 2017) | 61.20 | 0.36 | 1.17 | 25.98 |
| Language-only (Goyal et al., 2017) | 67.01 | 31.55 | 27.37 | 44.26 |
| MCB (Fukui et al., 2016) | 78.82 | 38.28 | 53.36 | 62.27 |
| Up-Down (Anderson et al., 2018) | 82.20 | 43.90 | 56.26 | 65.32 |
| VQA-E (Li et al., 2018b) | 83.22 | 43.58 | 56.79 | 66.31 |
| **Ours**(single) | **84.69** | **46.75** | **59.30** | **68.37** |
| **Ours**(Ensemble-10) | **86.15** | **47.41** | **60.41** | **69.66** |

Jialin Wu, Zeyuan Hu, and Raymond J. Mooney. 2019.Generating Question Relevant Captions to Aid Vi-sual Question Answering.CoRR, abs/1906.00513.

# Results



Q: What is he doing?
Caption: A man is taking a picture of himself with a phone.
A: Taking picture.

Q: Is the cat watching TV?
Caption: A cat is watching a bird on the **screen.**
A: Yes.

Q: What colors are on the couch?
Caption: A living **room** with a blue and **white** bed.
A: Purple and white.

Q: What color is the vase?
Caption: A **white vase** filled with lots of **flowers.**
A: White.

Q: Is he wearing a hat?
Caption: A man with **glasses** and a **hat** on.
A: Yes.

Q: Is the tv on?
Caption: A **bird** flying on a large **television** screen.
A: Yes.

Q: Is there a picture on the wall?
Caption: A bedroom with **pictures** on the **wall.**
A: Yes.

Q: What color are the flowers?
Caption: A vase filled with lots of **red roses.**
A: Red.

Jialin Wu, Zeyuan Hu, and Raymond J. Mooney. 2019.Generating Question Relevant Captions to Aid Vi-sual Question Answering.CoRR, abs/1906.00513.

# Results



Question: What colors is the surfboard?        Answer: yellow and red

Caption: A group of people standing next to yellow board.

Visual attention                    Caption adjusted visual attention

Answer: Yellow and blue            Answer: Yellow and red

Jialin Wu, Zeyuan Hu, and Raymond J. Mooney. 2019.Generating Question Relevant Captions to Aid Vi-sual Question Answering.CoRR, abs/1906.00513.