

# Reasoning with Heterogeneous Graph Alignment for Video Question Answering

Jiang, Pin, and Yahong Han (AAAI 2020)

Ahjeong Seo

Seoul National University

# Contents

1. 논문 소개
2. Reproduce 결과
3. 코드 실행 방법

# 1. 논문 소개

# 1. Introduction

## Recent efforts towards VideoQA (AS – IS)

### ■ Inter-modality correlations

- Uncover latent correlations between video content and words' semantics

- Beyond rnns: Positional self-attention with co-attention for video question answering (AAAI 2019, Li et al. )
- Progressive attention memory network for movie story question answering (CVPR 2019, Kim et al. )

### ■ Intra-modality correlations

- Incorporating correlations inside videos or dependencies

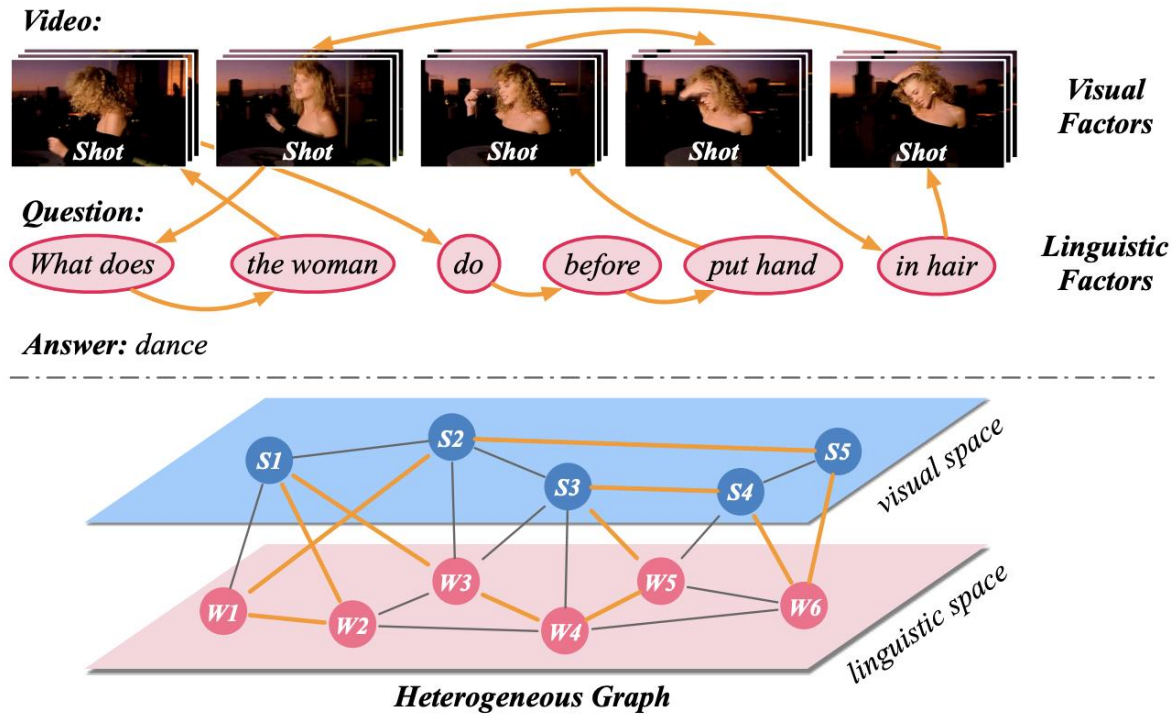
- Tgif-qa: Toward spatio-temporal reasoning in visual question answering (CVPR 2017, Jang et al. )
- Heterogeneous memory enhanced multimodal attention model for video question answering (CVPR 2019, Fan et al. )

### ■ Integrating correlations of both inter- and intra-modality

- Current methods lack of simultaneously reasoning inter- and intra-modality relations

# 1. Introduction

## Heterogeneous Graph Alignment (TO – BE)



## Contributions

- Novel heterogeneous graph over different modality factors
  - Reason inside one or inter-modality  
 $"s_1 \rightleftharpoons s_3"$ ,  $"s_1 \rightleftharpoons w_2"$
- Variety of modular co-attention embedding
  - Important role in cross-modal fusion and alignment before further aligned in GCN

## 2. Related Work

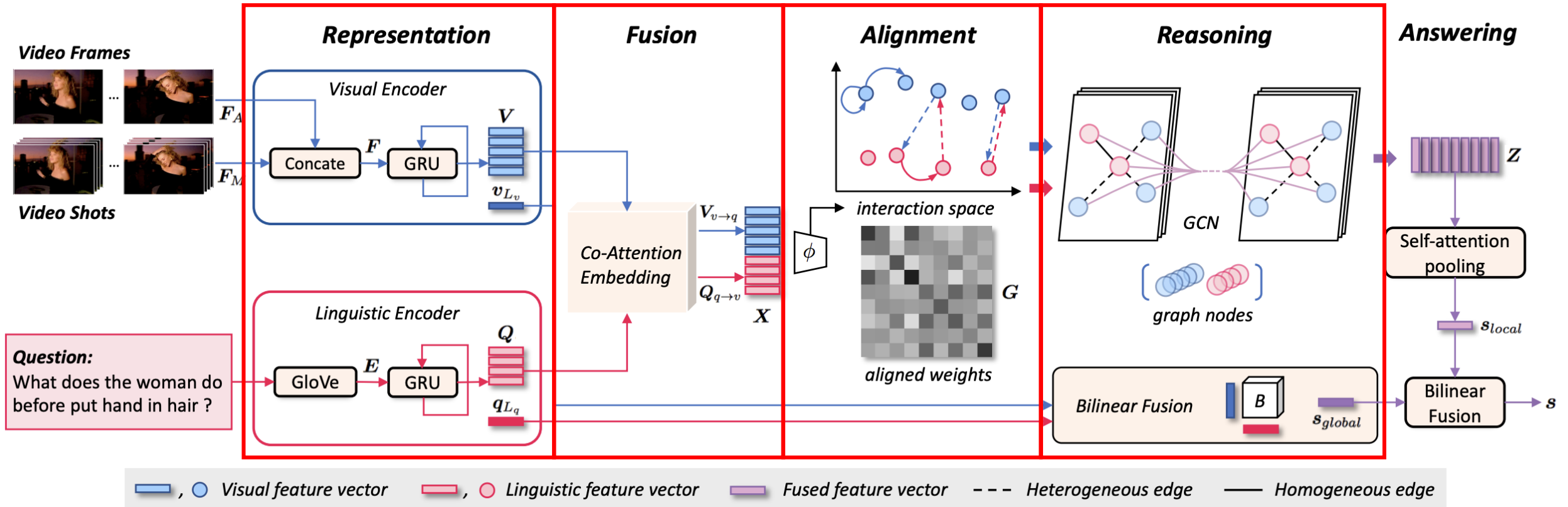
- Video Question Answering

- Extends VisualQA to video domain
- higher demands on **spatio-temporal** understanding and reasoning

- Multimodal Information Fusion

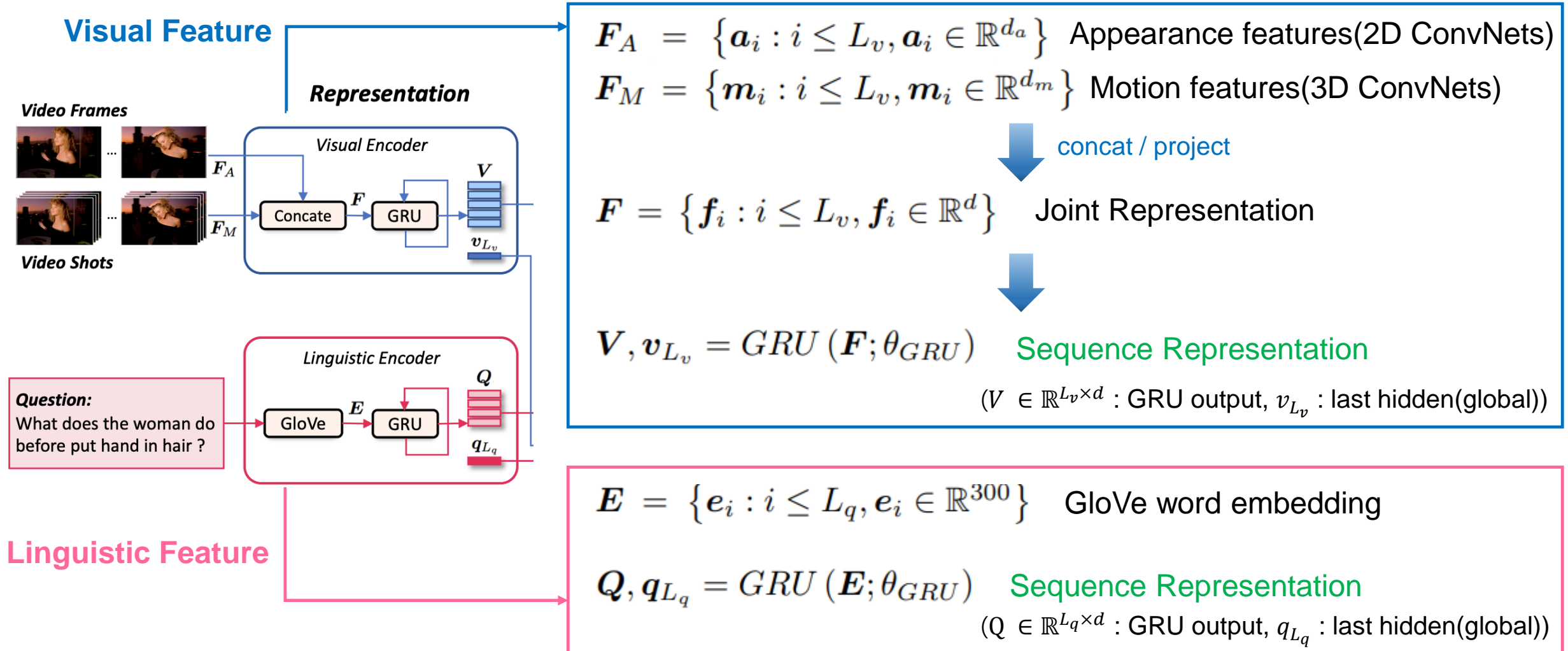
- Early fusion(vector operation: concatenation, element-wise addition and multiplication) and projected in joint space by neural network
- **Attention mechanism** enhance interaction between modalities
- **Co-attention**
- **Bilinear pooling** : fuse multimodal vectors by computing outer product

# 3. Method



# 3. Method

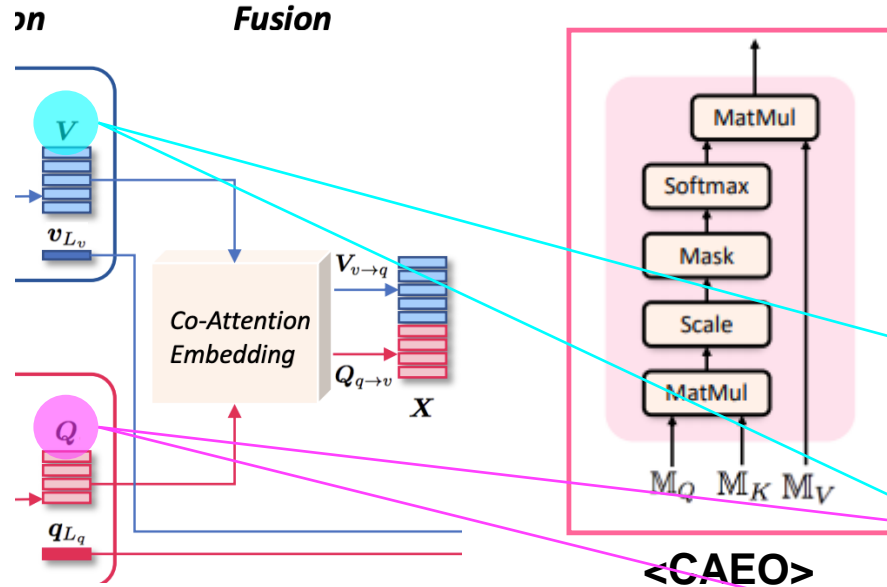
## 3.1 Visual and Linguistic Contextual Representation





# 3. Method

## 3.2 Cross-Modal Joint Fusion and Alignment

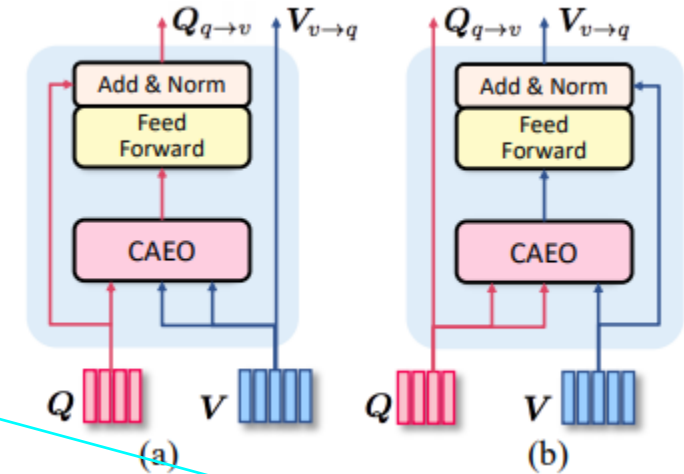


<CAEO>

$$CAEO(M_Q, M_K, M_V) = \text{softmax} \left( \frac{M_Q M_K^T}{\sqrt{d}} \right) M_V$$

- Query usually indicate different modality from Key, Value
- CAEO embeds the information of query into key's feature space

### 3.2.1 Attention-Based Structure



- From Linguistic to Visual space

$$Q_{CAEO} = CAEO(W_q^q Q, W_k^q V, W_v^q V),$$

$$Q_{q \to v} = \text{LayerNorm}(FF_q(Q_{CAEO}) + Q).$$

- From Visual to Linguistic space

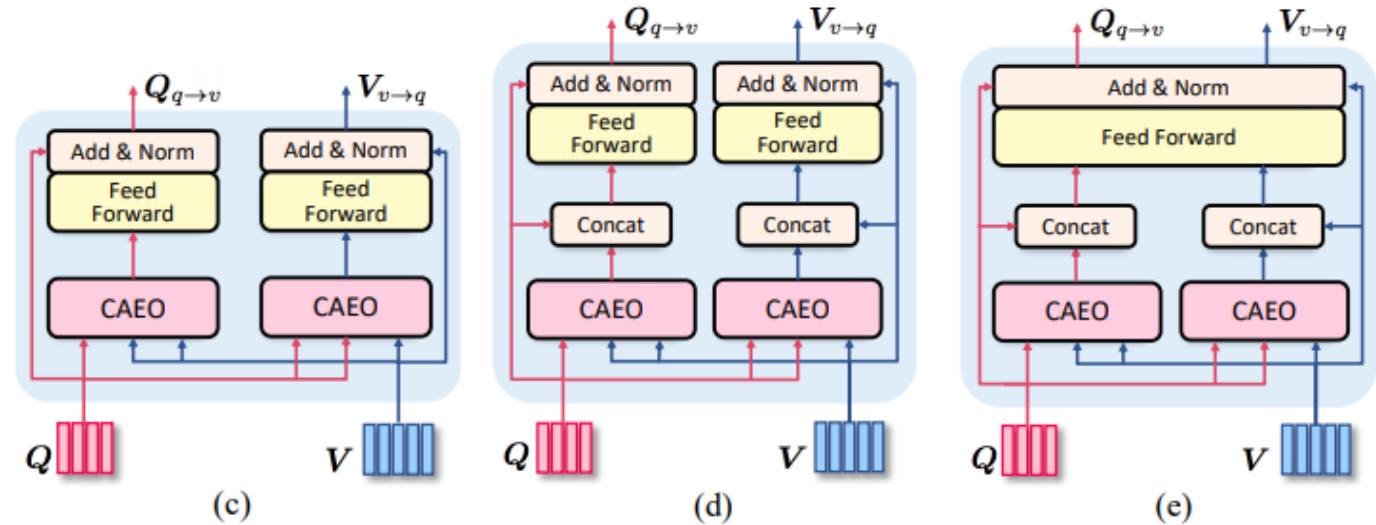
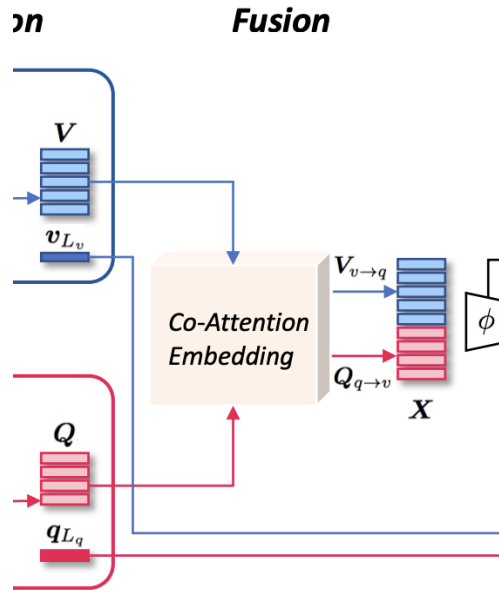
$$V_{CAEO} = CAEO(W_q^v V, W_k^v Q, W_v^v Q),$$

$$V_{v \to q} = \text{LayerNorm}(FF_v(V_{CAEO}) + V).$$

# 3. Method

## 3.2 Cross-Modal Joint Fusion and Alignment

### 3.2.2 Co-Attention-Based Structure



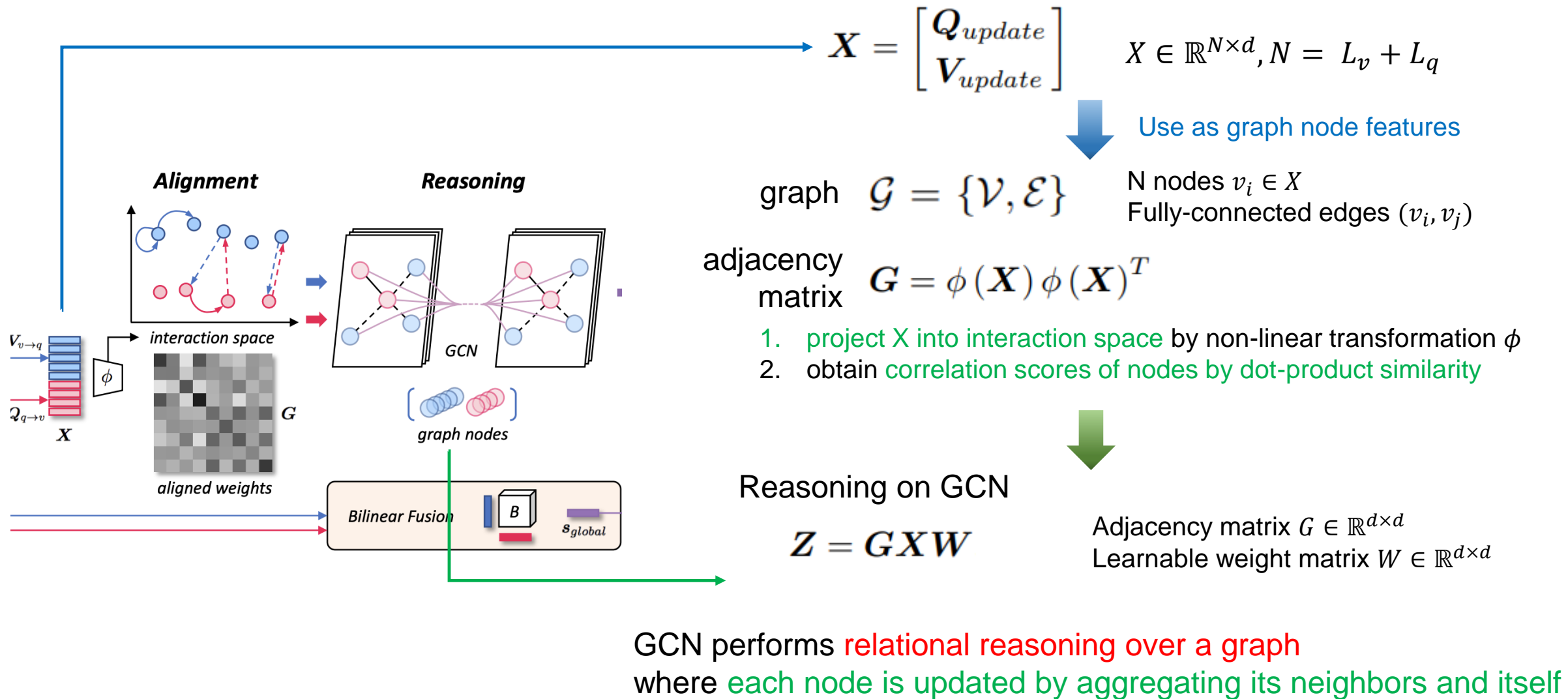
$$Q_{q \rightarrow v} = \text{LayerNorm} (FF_q ([Q_{CAEO}; Q]) + Q)$$

$$V_{v \rightarrow q} = \text{LayerNorm} (FF_v ([V_{CAEO}; V]) + V)$$

- Crossover transformation is crucial to fuse and align information of different modalities
- Combine two transformation and introduce a symmetrical co-attention operation

# 3. Method

## 3.3 Heterogeneous Graph Reasoning

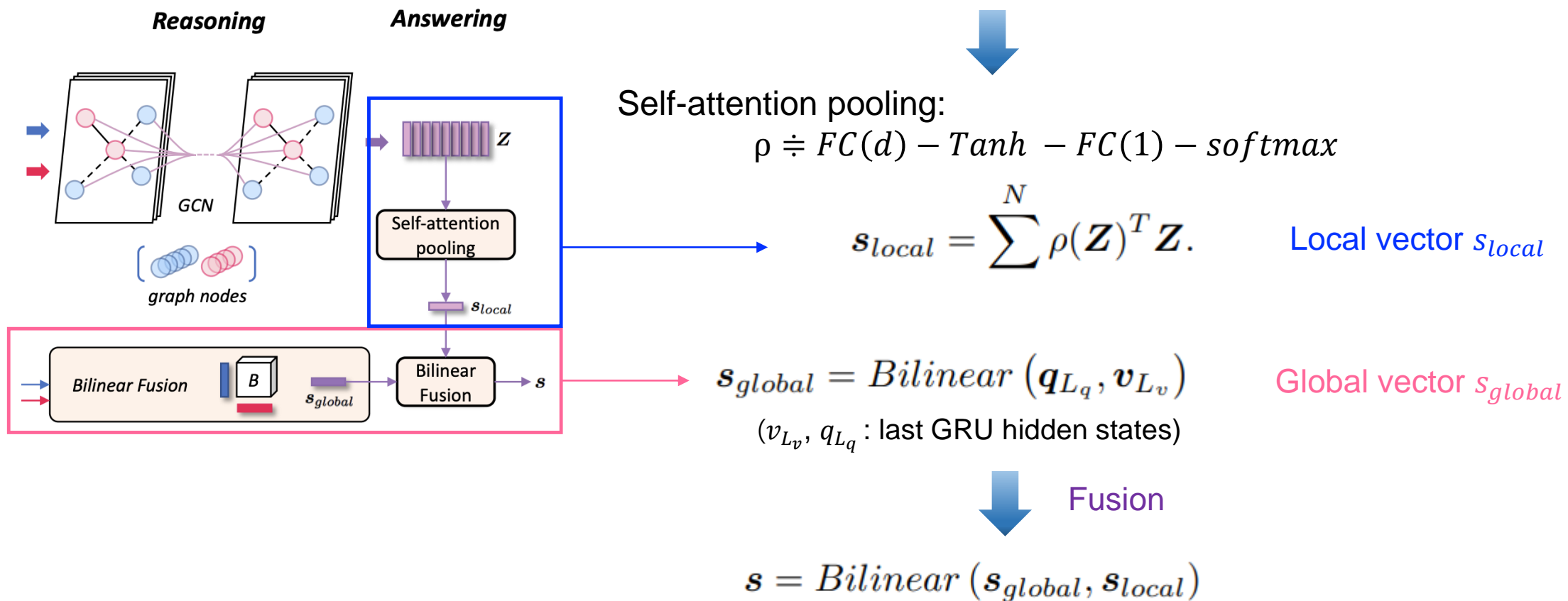


# 3. Method

## 3.3 Heterogeneous Graph Reasoning

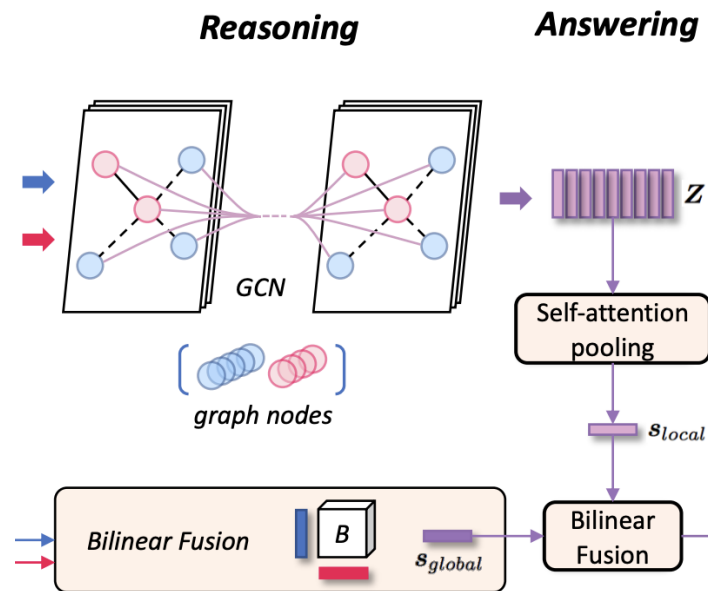
### ~ 3.4 Global and Local Information Fusion

GCN output  $Z = GXW$



## 3. Method

### 3.5 Answer Prediction and Evaluation



Input video clip  $v \in V$ , related question  $q \in Q$ , correct answer  $a^*$  and candidate set  $\{a_i\}_{i=1}^K$

1. Concentrate the question with each answer candidates



Get K candidate sequences

2. Linear regression function:

inputs the final output  $s$   
outputs K scores for all candidates

$$Loss = \sum_{i=1}^{K-1} \max(0, 1 + s_i^n - s^p)$$

incorrect answer scores  $s_i^n$   
correct answer score  $s^p$

## 4. Experiments

### 4.1 State of the Art Comparison

- Dataset : **TGIF-QA**

- 165K Q&A pairs from 72K animated GIFs
- Has four task types

- 1) Repetition Count(*Count*) :

open-ended numbers to count the number of repetition in action

- 2) Repeating Action(*Action*) :

multiple-choice task to identifying a repetitive action

- 3) State Transition(*Trans.*) :

5-options multiple-choice which identifies the transition of two states

- 4) FrameQA(*FrameQA*) :

open-ended task that can be answered from single frame

Table 1: State-of-the-art comparison on TGIF-QA dataset. Mean  $\ell_2$  loss for Count, and accuracy (%) for others.

Methods	Count	Action	Trans.	FrameQA
Random	19.62	20.00	20.00	0.06
ST-VQA-Sp.	4.28	57.3	63.7	45.5
ST-VQA-Tp.	4.40	60.8	67.1	49.3
ST-VQA-Sp.Tp.	4.56	57.0	59.6	47.8
CT-SAN	5.14	56.1	64.0	39.6
Co-Mem	<u>4.10</u>	68.2	74.3	51.5
PSAC	4.27	70.4	76.9	<b>55.7</b>
Fan et al.	<u>4.10</u> <sup>1</sup>	<u>73.9</u>	77.8	53.8
ST-VQA★	4.22	73.5	<u>79.7</u>	52.0
Ours HGA	<b>4.09</b>	<b>75.4</b>	<b>81.0</b>	<u>55.1</u>

## 4. Experiments

### 4.1 State of the Art Comparison

- Dataset : **MSVD-QA and MSRVTT-QA**
  - Two datasets generated from video descriptions
  - 50K and 243K Q&A pairs
  - Consists of different types of questions : *what, who, how, when, where*

Table 2: State-of-the-art comparison on MSVD-QA dataset. Mean  $\ell_2$  loss for Count, and accuracy (%) for others.

Methods	What (8,149)	Who (4,552)	How (370)	When (58)	Where (28)	All (13,157)
ST-VQA	18.1	50.0	<b>83.8</b>	72.4	28.6	31.3
Co-Mem	19.6	48.7	81.6	<b>74.1</b>	31.7	31.7
AMU	20.6	47.5	83.5	72.4	<b>53.6</b>	32.0
Fan et al.	22.4	50.1	73.0	70.7	42.9	33.7
Ours HGA	<b>23.5</b>	<b>50.4</b>	83.0	72.4	46.4	<b>34.7</b>

Table 3: State-of-the-art comparison on MSRVTT-QA dataset. Mean  $\ell_2$  loss for Count, and accuracy (%) for others.

Methods	What (49,869)	Who (20,385)	How (1,640)	When (677)	Where (250)	All (72,821)
ST-VQA	24.5	41.2	78.0	<b>76.5</b>	34.9	30.9
Co-Mem	23.9	42.5	74.1	69.0	<b>42.9</b>	32.0
AMU	26.2	43.0	80.2	72.5	30.0	32.5
Fan et al.	26.5	43.6	82.4	76.0	28.6	33.0
Ours HGA	<b>29.2</b>	<b>45.7</b>	<b>83.5</b>	75.2	34.0	<b>35.5</b>

## 4. Experiments

### 4.2 Ablation Study

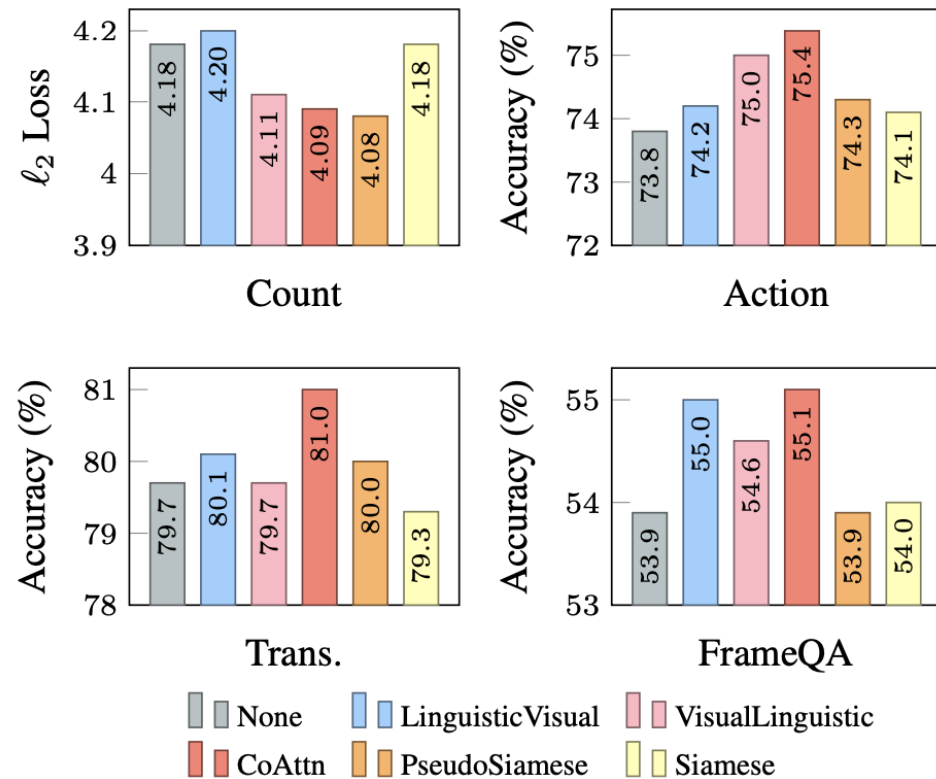


Figure 4: Experimental results of *CAEO* variants.

- On TGIF-QA dataset
- Appending CAEO improve the performance to some extent
- CoAttn achieves the best on *Action*, *Trans.*, and *FrameQA*
- But sub-optimal in *Count* task
- PseudoSiamese and Siamese hurt model quality



## 4. Experiments

### 4.2 Ablation Study

Table 4: Ablation study on TGIF-QA dataset. Mean  $\ell_2$  loss for *Count*, and accuracy (%) for others.

Methods	Count	Action	Trans.	FrameQA
GRU (w/ local fusion)	4.31	55.4	69.8	53.0
+ global fusion (baseline)	4.19	73.4	78.1	55.7
+ GCN §3.3	4.18	74.5	79.7	53.9
+ CoAttn §3.2	4.17	73.9	78.8	55.4
+ GCN §3.3	4.09	75.4	81.0	55.1
HGA w/o global fusion	4.25	73.6	79.6	53.6
HGA w/o local fusion	4.24	71.3	77.8	53.9

- On TGIF-QA dataset
- Methods
  - GRU : vanilla parallel architecture
  - + global fusion : in section 3.4 (add local feature with global feature)
  - + GCN : add GCN
  - + CoAttn : add CAEO module
  - + GCN : add GCN and CAEO
- Combining “CoAttn” and “GCN”
  - Two modules promote each other

## 4. Experiments

### 4.3 Examples of successful cases

*Q: What does the woman do before put hand in hair ?*

*A1: dance*

*A2: point gun*

*A3: close*

*A4: pick up head*

*A5: grab hair*



*Q: What does the man do before brush hair ?*

*A1: stare down*

*A2: hug a boy*

*A3: emerge from a robot*

*A4: contemplate*

*A5: fix jacket*



*Q: How many times does the man kick soccer ball ?*

*A: 6*



*Q: What jumps up at itself in the mirror ?*

*A: cat*



Figure 5: Typical examples of successful cases on TGIF-QA dataset.

## 2. Reproduce 결과

## ■ TGIF-QA Dataset에 대해 성능 비교 결과

### 1. 논문에 report된 성능

Methods	Count	Action	Trans.	FrameQA
HGA	4.09	75.4	81.0	55.1

### 2. Reproduce 성능

Methods	Count	Action	Trans.	FrameQA
Reproduced HGA	4.05	75.50	80.95	55.24

- Count와 Action, FrameQA의 경우 report된 성능보다 오히려 높은 accuracy를 얻었습니다.
- 주어지는 test dataset에 대해 실험하였기 때문에 실제 leaderboard상에서의 성능과 차이를 보이는 것으로 예상됩니다.
- Training log는 github code의 log 폴더에 저장되어 있습니다.

## ▪ Inference 결과 화면 캡처

### - Count

```
Test|Epoch: 0, Acc: 23.747890=844/3554, Test Loss: 3.965284  
Test|Count Real Loss: 4.052
```

### - Action

```
Test|Epoch: 0, Acc: 75.505717=1717/2274, Test Loss: 0.424502  
Save model at ./saved_models/MMModel/
```

### - Trans.

```
Test|Epoch: 0, Acc: 80.953145=5045/6232, Test Loss: 0.211618  
Save model at ./saved_models/MMModel/
```

### - FrameQA

```
Frame Frame  
Test|Epoch: 0, Acc: 55.240669=7563/13691, Test Loss: 2.224057  
Save model at ./saved_models/MMModel/
```

### 3. 코드 실행 방법

- Github Link: <https://github.com/ahjeongseo/HGA>

## 1. Requirements

- python 3.6
- pytorch 1.1
- colorlog
- bootstrap.pytorch
- block.bootstrap.pytorch
- numpy==1.18.5

(설명)

1. 'conda create -n final\_project python=3.6' 명령어로 python 3.6 가상환경을 만듭니다.
2. 'conda install pytorch==1.1.0 torchvision==0.3.0 cudatoolkit=9.0 -c pytorch' 로 pytorch 1.1을 설치합니다.  
저의 경우 cuda 9.0 version을 사용하였습니다. Cuda version에 따라 <https://pytorch.org/get-started/previous-versions/>를 참고할 수 있습니다.
3. 'pip install -r requirements.txt'를 통해 필요한 패키지를 다운 받습니다.

## 2. Data Download

- C3D feature, ImageNet feature, GloVe feature:

<https://github.com/fanchenyou/HME-VideoQA/tree/master/gif-qa>에서 다운로드

- TGIF QA question csv file:

<https://github.com/YunseokJANG/tgif-qa/tree/master/dataset>에서 다운로드

## 3. Data path 설정

- main.py에서 주석처리한 부분 path설정

- 1) feat\_dir : ImageNet 및 C3D feature를 저장한 path
- 2) vc\_dir : GloVe feature를 저장한 path, word embedding관련 파일이 dataloader를 실행하면 저장됨
- 3) df\_dir : TGIF question csv file을 저장한 path
- 4) checkpoint : model inference시 model file명
- 5) save\_path : training시 model 저장 path 및 inference시 모델 불러올 path



## 4. Training

**CUDA\_VISIBLE\_DEVICES=0 python main.py --test --task Count --num\_workers 2 --batch\_size 64**

- Task : Count, Action, Trans, FrameQA

## 5. Test (Inference)

**CUDA\_VISIBLE\_DEVICES=0 python main.py --test --task Count --num\_workers 2 --batch\_size 64**

- Task : Count, Action, Trans, FrameQA

- checkpoint option추가 가능 : inference 할 모델 파일 명 추가 (ex. --checkpoint Count\_4.092.pth)

## 6. 중요 코드 설명

- /data\_utils/dataset.py : TGIFQA dataloader를 구성하기 위한 class 파일
- /data\_utils/.. : 나머지 파일은 dataset.py를 구성하기 위한 submodule
- /models/lstm\_cross\_cycle\_gcn\_dropout.py : HGA 모델이 구현된 class 파일
- lstm\_cross\_cycle\_gcn\_dropout.py의 submodule들
  - /models/rnn\_encoder.py : video 및 text의 sequence representation을 위한 rnn 구현
  - /models/q\_v\_transformer.py : video 및 text의 attention, co-attention 구현
  - /models/gcn.py : co-attention된 video, text feature를 gcn으로 학습하기 위한 코드 구현
- main.py : 모델 학습을 위한 training 코드

*\* 모델 부분 코드 구현(models/...py)에 대한 자세한 설명은 주석으로 추가하였습니다.*