

Cross-Modal Interaction Networks for Query-Based Moment Retrieval in Videos

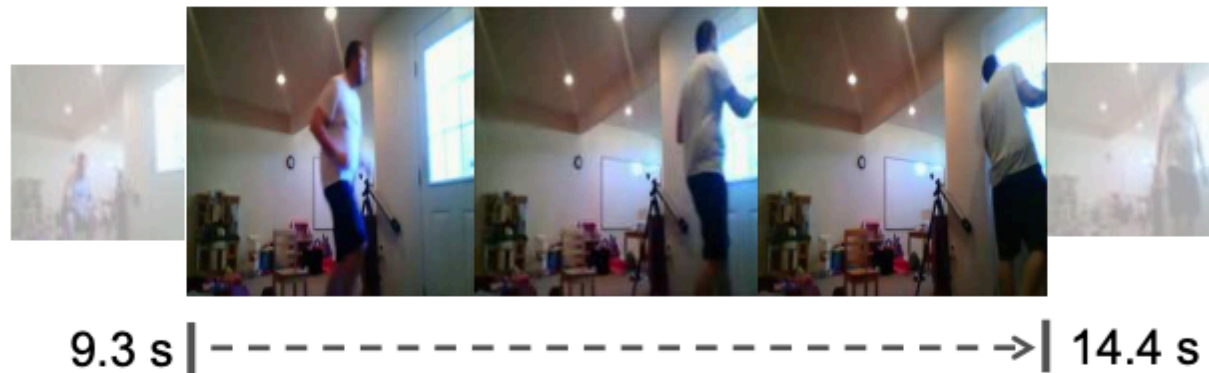
Zhu Zhang, Zhijie Lin, Zhou Zhao, Zhenxin Xiao
SIGIR 2019

Task

Query-Based Moment Retrieval

- 動画の一部を説明した文章を受けて、それに対応する場면을動画内からローカライズするタスク

Language Query:
A person runs to the window and then look out



[TALL](#)

Related works

- [MCN\(Moment Context networks\)](#) (Hendricks, ICCV2017)
動画特徴量(動画全体のglobal featureとタイムスタンプ内のlocal featureの合成)と言語特徴量の分布が最短距離となるタイムスタンプを推定
- [CTRL\(Cross-modal Temporal Regression Localizer\)](#) (Gao, ICCV2017)
クリップ動画(複数の固定サイズ)と言語の特徴量をcross-modalに合成
アライメントスコアを推定し, 最大スコアのモーメントの回帰を行って
ローカライズ
- [ACRN\(Attentive Cross-modal Retrieval Network\)](#) (Liu, SIGIR2018)
attentionを組み込んで認識により有効な特徴量を強調
各モーダル特徴量と融合特徴量を合成
回帰モデルからアライメントスコアとオフセットを算出し, ローカライズ
- [QSPN\(Query-Guided Proposal Network\)](#) (Xu, AAAI2019)
動画特徴量と言語特徴量をattentionを組み込んでearly fusionするマルチモデルからclassificationとregressionで学習
クエリの再生成を補助タスク

Problem

query representation → RNN

- query文の**文法的な文構造**を捉え切れていない

video representation → CNN+RNN

- 動画のsemanticな関係性の情報を**長期的**に捉えることができない

Cross-modal Interaction → attention (fusion)

- 一般的なattentionが**一層だけ**設けられているのみで十分な相互関係を見ることができない

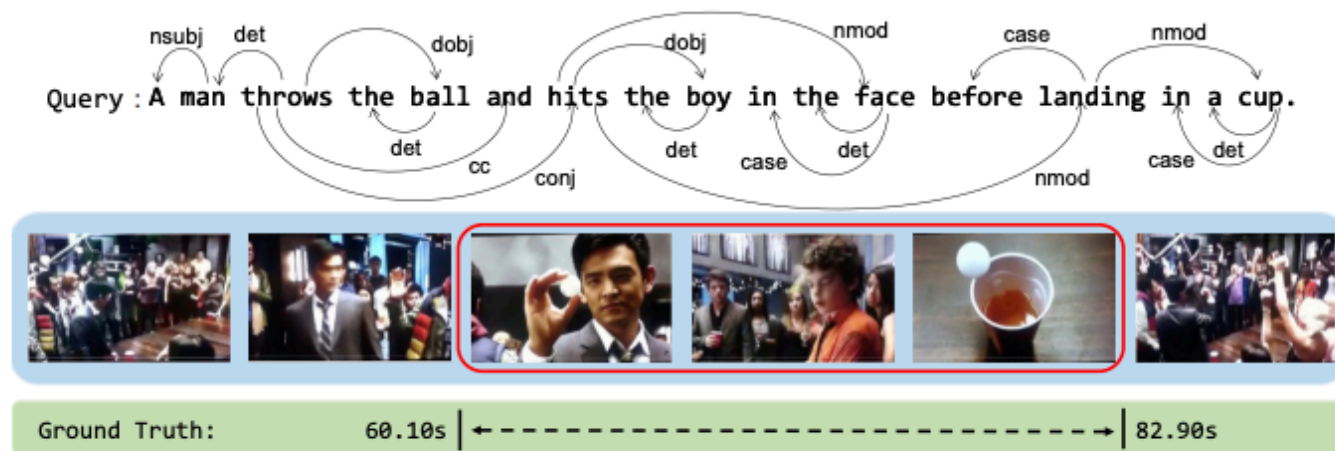


Figure 1: Query-Based Moment Retrieval in Video

Approach

query representation

→ GCN層によってquery文の文法的な構造を捉える

video representation

→ multi-head self-attention層によって動画の長期的でsemanticな関係性を捉える

Cross-modal Interaction → attention (fusion)

→ cross-modal interaction層の多層化によって効率的にモーダル間の統合的な情報を捉える

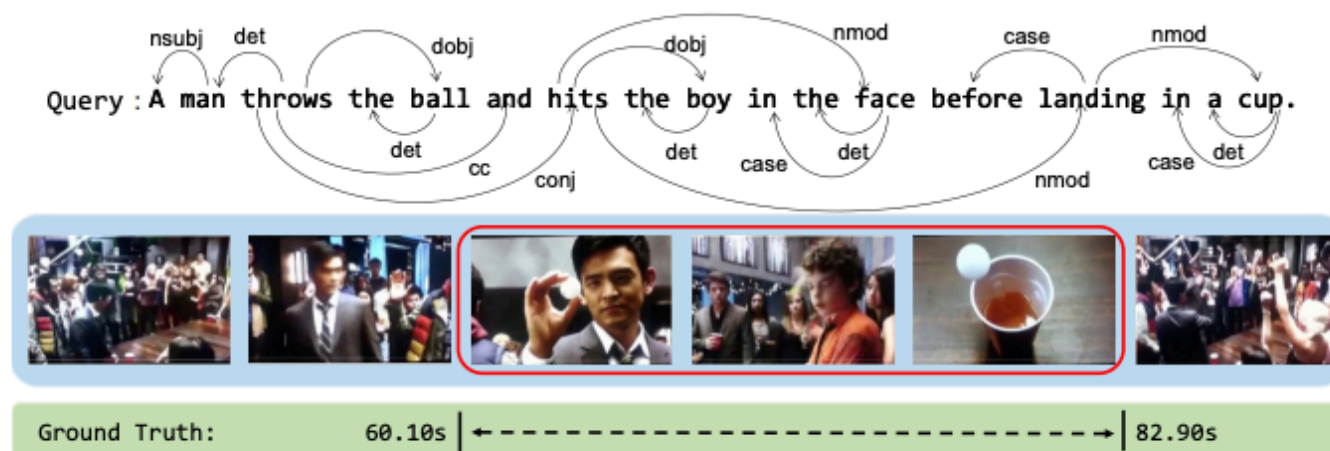
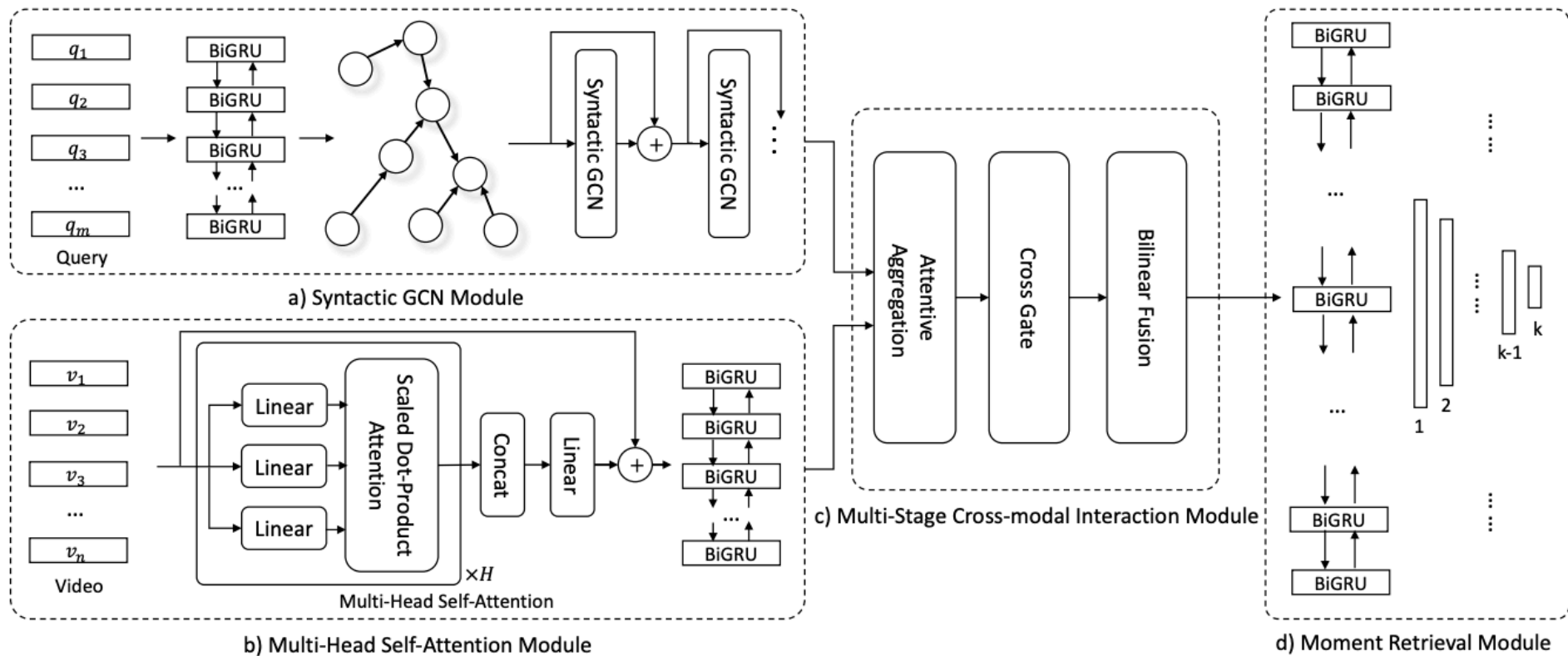


Figure 1: Query-Based Moment Retrieval in Video

Proposed Method



4モジュール構成

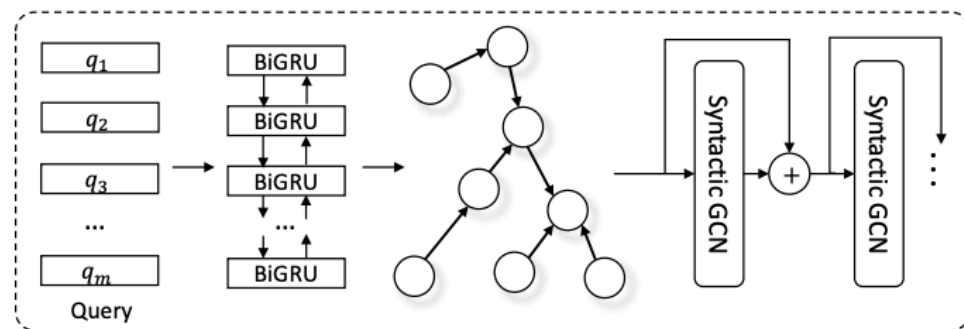
- a. syntactic GCNモジュール
- b. Multi-Head Self-Attentionモジュール
- c. Multi-Stage Cross-Modal Interactionモジュール
- d. Moment Retrievalモジュール

Proposed Method (Query)

文法的な文構造

入力: query

出力: syntactic-aware query representation



a) Syntactic GCN Module

- 文法的な構造関係の有向グラフを作成



$$\mathbf{g}_i^1 = \text{ReLU} \left(\sum_{j \in \mathcal{N}(i)} \mathbf{w}_{dir(i,j)}^g \mathbf{h}_j^q + \mathbf{b}_{lab(i,j)}^g \right)$$

- syntactic GCNを介して文法的な文構造を考慮した特徴量

$$\mathbf{o}_i^1 = \mathbf{g}_i^1 + \mathbf{h}_i^q$$

$$\begin{cases} \mathbf{g}^1 = \text{synGCN}(\mathbf{h}^q), \mathbf{o}^1 = \mathbf{g}^1 + \mathbf{h}^q \\ \mathbf{g}^l = \text{synGCN}(\mathbf{o}^{l-1}), \mathbf{o}^l = \mathbf{g}^l + \mathbf{o}^{l-1} \end{cases}$$

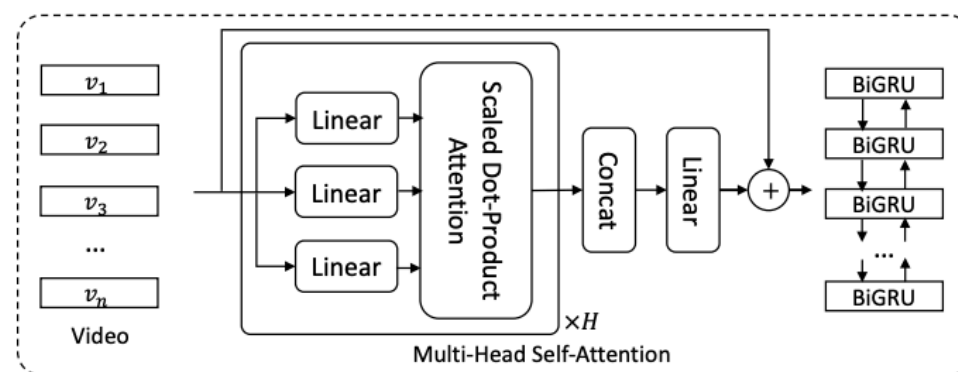
$$\mathbf{o}^l = (\mathbf{o}_1^l, \mathbf{o}_2^l, \dots, \mathbf{o}_m^l)$$

Proposed Method (Video)

動画の長期的でsemanticな関係性

入力: video

出力: video semantic representation



b) Multi-Head Self-Attention Module

- Multi-Head Self-Attentionによって近接フレーム間だけでなく、離れたフレーム間の関係性も考慮

$$\text{Attention}(\bar{Q}, \bar{K}, \bar{V}) = \text{Softmax}\left(\frac{\bar{Q}^T \bar{K}}{\sqrt{d_k}}\right) \bar{V}$$

$$\text{MultiHead}(\bar{Q}, \bar{K}, \bar{V}) = \mathbf{W}^O \text{Concat}(\text{head}_1, \dots, \text{head}_H)$$

$$\text{where head}_i = \text{Attention}(\mathbf{W}_i^Q \bar{Q}, \mathbf{W}_i^K \bar{K}, \mathbf{W}_i^V \bar{V})$$

$$\mathbf{V}^S = \text{MultiHead}(\mathbf{V}, \mathbf{V}, \mathbf{V}) + \mathbf{V}$$

Proposed Method (Cross-Modal)

入力: query and video semantic representation

出力: cross-modal representation

- Attentive Aggregation

attentionベースで各単語 j と各フレーム i との関係性 M_{ij}^{row}
フレーム毎でのqueryの集約表現 \mathbf{h}_i^s

$$\mathbf{h}_i^s = \sum_{j=1}^m M_{ij}^{row} \mathbf{o}_j^l$$

- Cross Gate

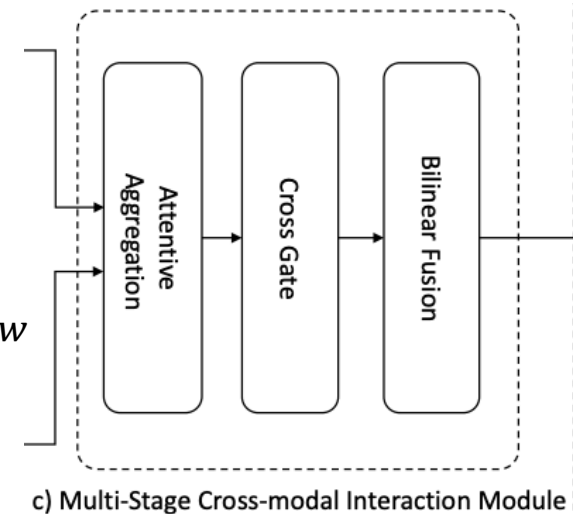
queryの集約表現 \mathbf{h}_i^s とframe semantic表現 \mathbf{h}_i^v との関係性から重要度が重み付け

$$\tilde{\mathbf{h}}_i^s = \mathbf{h}_i^s \odot \sigma(\mathbf{W}^v \mathbf{h}_i^v + \mathbf{b}^v)$$

$$\tilde{\mathbf{h}}_i^v = \mathbf{h}_i^v \odot \sigma(\mathbf{W}^s \mathbf{h}_i^s + \mathbf{b}^s)$$

- Bilinear Fusion

cross-modal表現 $\mathbf{f} = (\mathbf{f}_1, \dots, \mathbf{f}_n)$ を獲得



Proposed Method (Localization)

入力: cross-modal representation

出力: moment

- 複数サイズ k の固定幅ウィンドウを候補として用意
- 各タイムステップ i においてウィンドウをセット
- 同時に confidence score cs_i を算出
- これらのモーメントのオフセット(差分) $\hat{\delta}_s, \hat{\delta}_e$ を推定

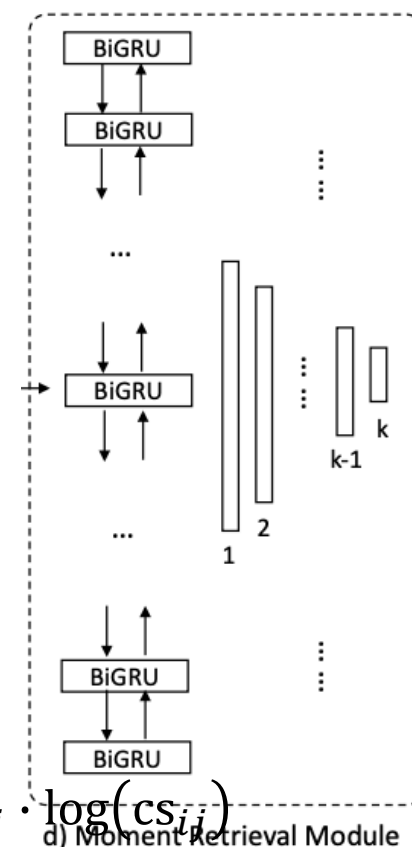
Loss

- **Alignment loss**
算出されるアライメントスコア confidence score に対して

$$L_{align} = -\frac{1}{nk} \sum_{i=1}^n \sum_{j=1}^k (1 - \text{IoU}_{ij}) \cdot \log(1 - cs_{ij}) + \text{IoU}_{ij} \cdot \log(cs_{ij})$$

- **Regression loss**
オフセットに対して (δ_s, δ_e) :GT, $(\hat{\delta}_s, \hat{\delta}_e)$:predicted

$$L_{reg} = -\frac{1}{N} \sum \left(\text{L1}(\delta_s - \hat{\delta}_s) + \text{L1}(\delta_e - \hat{\delta}_e) \right)$$



Results

Table 2: Performance Evaluation Results on the ActivityCaption Dataset ($n \in \{1, 5\}$ and $m \in \{0.3, 0.5, 0.7\}$).

Method	R@1	R@1	R@1	R@5	R@5	R@5
	IoU=0.3	IoU=0.5	IoU=0.7	IoU=0.3	IoU=0.5	IoU=0.7
MCN	39.35	21.36	6.43	68.12	53.23	29.70
VSA-RNN	39.28	23.43	9.01	70.84	55.52	32.12
VSA-STV	41.71	24.01	8.92	71.05	56.62	34.52
CTRL	47.43	29.01	10.34	75.32	59.17	37.54
ACRN	49.70	31.67	11.25	76.50	60.34	38.57
QSPN	52.13	33.26	13.43	77.72	62.39	40.78
CMIN	63.61	43.40	23.88	80.54	67.95	50.73

Table 3: Performance Evaluation Results on the TACoS Dataset ($n \in \{1, 5\}$ and $m \in \{0.1, 0.3, 0.5\}$).

Method	R@1	R@1	R@1	R@5	R@5	R@5
	IoU=0.1	IoU=0.3	IoU=0.5	IoU=0.1	IoU=0.3	IoU=0.5
MCN	3.11	1.64	1.25	3.11	2.03	1.25
VSA-RNN	8.84	10.77	4.78	19.05	13.90	9.10
VSA-STV	15.01	10.77	7.56	32.82	23.92	15.50
CTRL	24.32	18.32	13.30	48.73	36.69	25.42
ACRN	24.22	19.52	14.62	47.42	34.97	24.88
QSPN	25.31	20.15	15.23	53.21	36.72	25.30
CMIN	32.48	24.64	18.05	62.13	38.46	27.02

Results

Table 4: Performance Evaluation Results of Ablation Model on the ActivityCaption dataset.

Method	R@1	R@1	R@1	R@5	R@5	R@5
	IoU=0.3	IoU=0.5	IoU=0.7	IoU=0.3	IoU=0.5	IoU=0.7
w/o. GCN	60.12	40.84	21.79	78.23	65.67	45.43
w/o. SA	61.22	41.56	22.36	79.43	66.91	48.12
w/o. CG	60.57	41.21	22.01	78.62	65.99	46.89
w/o. BF	61.32	41.89	22.12	79.27	66.21	47.92
full	63.61	43.40	23.88	80.54	67.95	50.73

Table 5: Performance Evaluation Results of Ablation Model on the TACoS dataset.

Method	R@1	R@1	R@1	R@5	R@5	R@5
	IoU=0.1	IoU=0.3	IoU=0.5	IoU=0.1	IoU=0.3	IoU=0.5
w/o. GCN	30.54	23.22	10.03	57.69	37.12	26.16
w/o. SA	30.21	23.02	16.87	55.54	36.6	25.37
w/o. CG	31.96	23.59	17.47	61.87	38.11	26.79
w/o. BF	32.01	24.79	17.61	61.59	38.23	26.75
full	32.48	24.64	18.05	62.13	38.46	27.02

Results

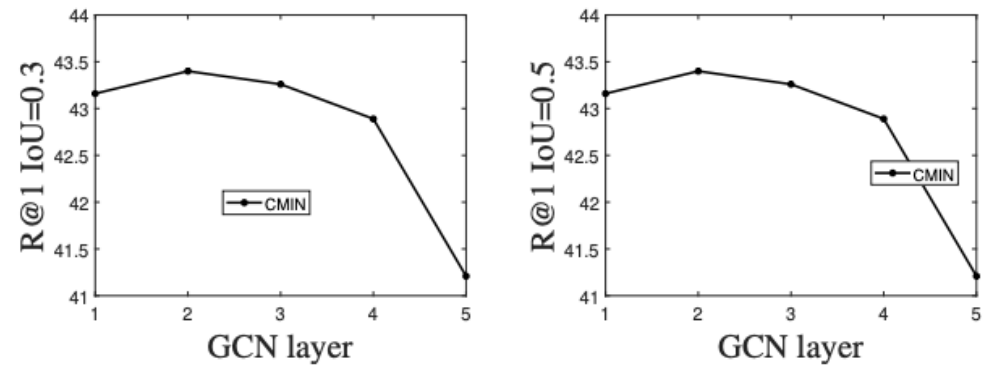


Figure 4: Effect of the Number of Stacked Syntactic GCN layers on the ActivityCaption Dataset.

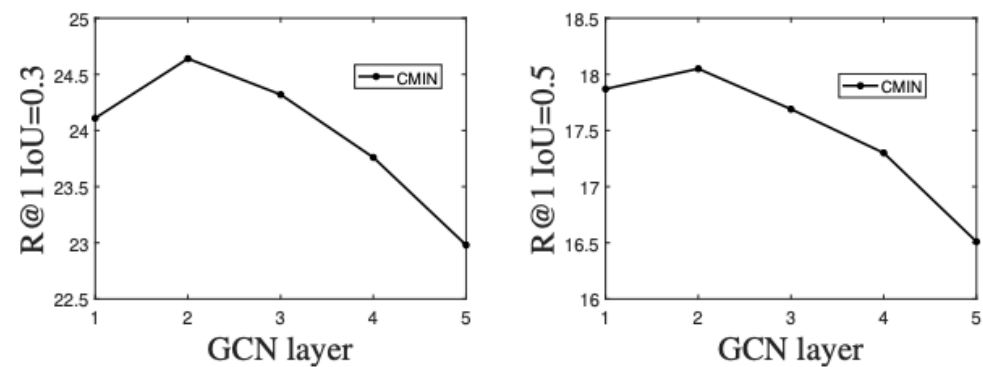
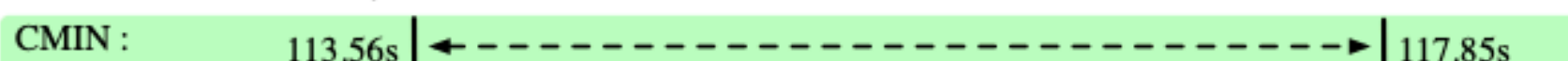
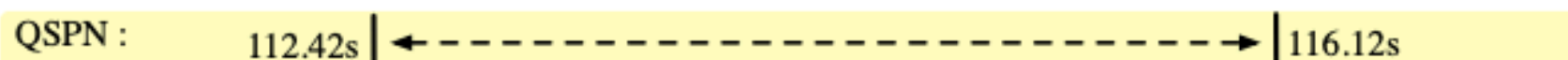
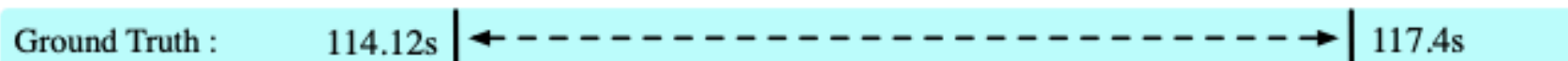
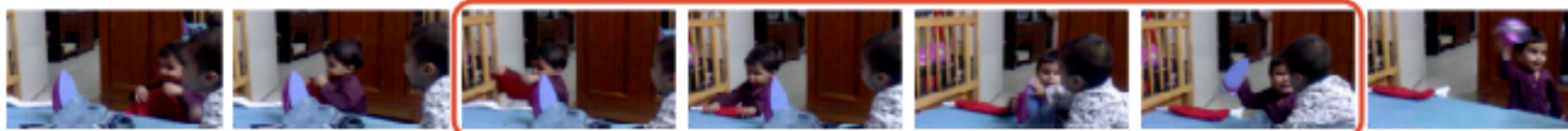


Figure 5: Effect of the Number of Stacked Syntactic GCN layers on the TACoS Dataset.

Syntactic GCN layerの層数による影響

Results

Query : The boy drops the cloths and takes the iron away **before** the baby can pick it up.



Query : The female athlete jumped over the pole **and** wave at everyone.

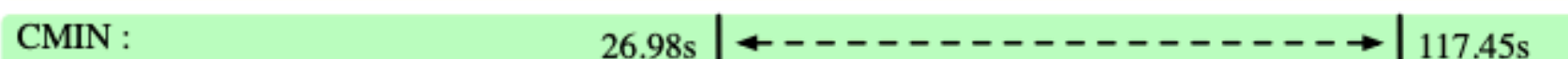
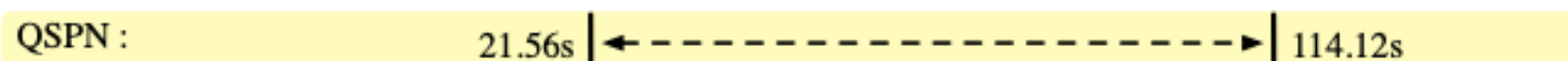
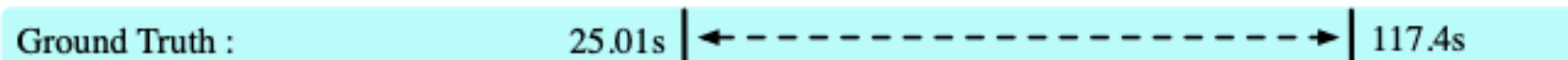


Figure 6: Examples on the ActivityCaption dataset.

Results

Query : **After** getting out the juicer, he juices the **first** orange half.



Ground Truth : 74.56s | ← ----- → | 120.34s

QSPN : 79.1s | ← ----- → | 129.31s

CMIN : 72.32s | ← ----- → | 117.78s

Query : She washes herb stems in the sink **before** placing them on the cuttingboard.



Ground Truth : 126.29s | ← ----- → | 127.62s

QSPN : 120.1s | ← ----- → | 133.34s

CMIN : 123.04s | ← ----- → | 130.72s

Figure 7: Examples on the TACoS dataset.

Results

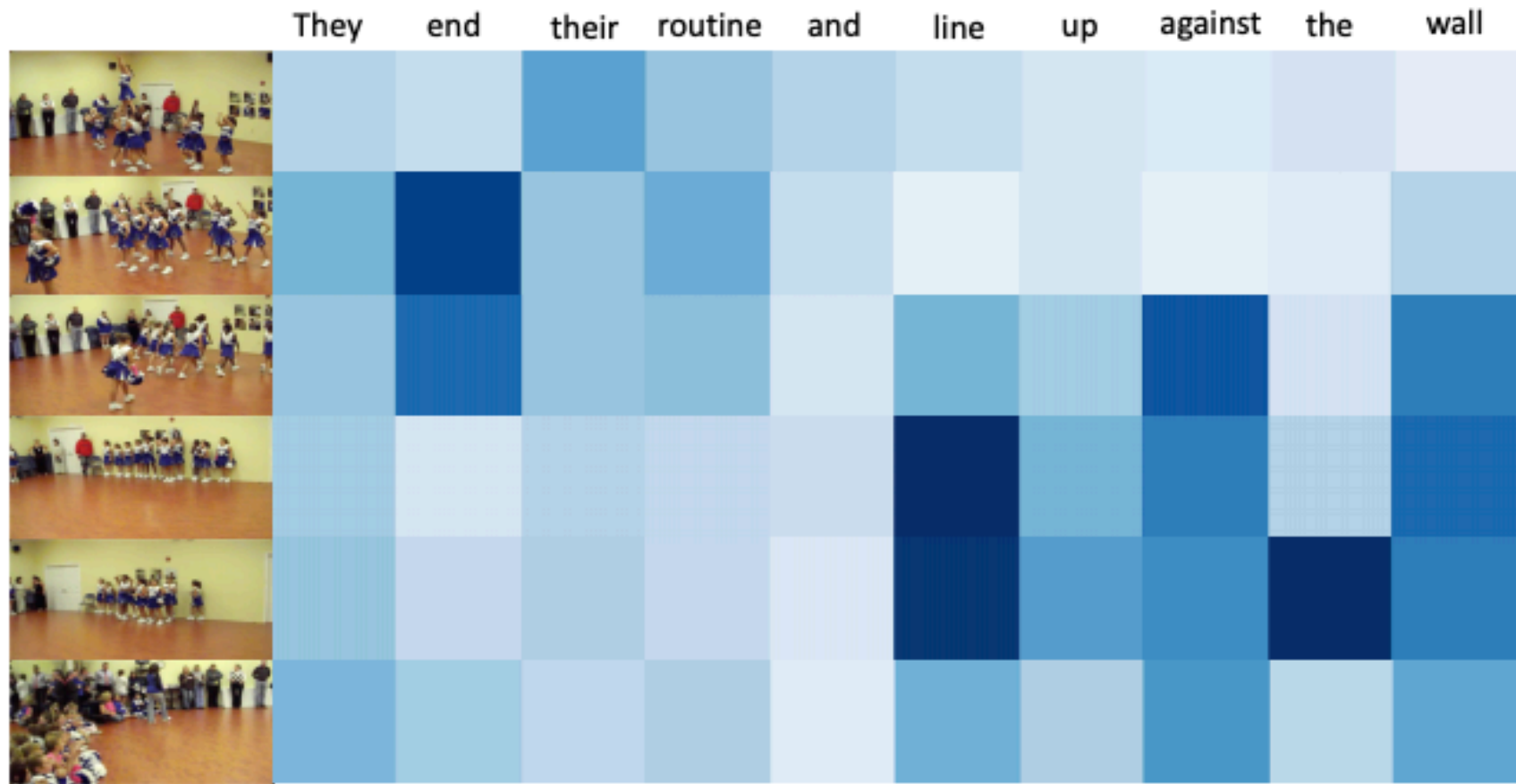


Figure 8: The Video-to-Query Attention Results in the Multi-Stage Cross-Modal Interaction Module