# Multi-scale 2D Representation Learning for weakly-supervised moment retrieval

Ding Li \*, Rui Wu \*\*, Yongqiang Tang \*, Zhizhong Zhang \*, Wensheng Zhang \*

\* Institute of Automation, Chinese Academy of Sciences
\*\* Horizon Robotics, Beijing, China

- Fully-supervised moment retrieval
  - Input:
    - Untrimmed video & Text query
    - Temporal boundary annotations: [start, end]
  - Output: Corresponding segment: [start, end]

- Weakly-supervised moment retrieval
  - Input: Untrimmed video & Text query
  - Output: Corresponding segment: [start, end]

#### Motivation

- Moment candidate relations
- Variance of temporal scale of video moments



- Related work
  - 2D TAN
  - 2D temporal map
  - Temporal Adjacent Net
- Pros
  - Candidates Relations
- Cons
- ns



- Variance of temporal scale
- No supervision signals when start&end labels are missing

### Method

- Basic Video and Text Representation
  - 3D CNN spatio-temporal feature
  - Word-embedding & LSTM
- Multi-scale 2D Temporal Network
  - Multi-scale 2D temporal map of moment candidates
- Moments Evaluation Module
  - Moments Caption & Evaluation
- Loss functions
  - Reconstruction loss
  - Reconstruction-guided binary cross-entropy loss

# Framework of our method



- Multi-scale 2D temporal map
  - Basic video segment feature  $f^S \in R^{d^v}$
  - Single-scale 2D temporal map

$$F_{x,y}^{j} = \begin{cases} \max \operatorname{pool}\left(f_{x}^{S}, f_{x+1}^{S}, \cdots, f_{y}^{S}\right), & if 0 < x \le y < N_{j} \\ 0^{S}, & else \end{cases} \qquad F^{j} \in R^{N_{j} \times N_{j} \times d^{v}}$$

– Multi-scale temporal sampling on untrimmed video  $F^{M} = \left\{F^{j}\right\}_{i=1}^{N_{s}}$ 

 $N_j$  is the number of sampled segments in the j-th temporal scale,  $N_s$  is the number of scales

• Multi-scale 2D temporal network

- Alignment score  $P = \text{sigmoid} (W^F \cdot F_{cro})$ 

#### Cross-modal fusion

 Cross-modal feature map are fused by the Hadamard product and I2 normalization

$$F_{cro} = \left\| \begin{pmatrix} W^T \cdot F^T \cdot 1^T \end{pmatrix} \odot \begin{pmatrix} W^M \cdot F^M \end{pmatrix} \right\|_F$$
  
Multi-scale Text feature



- Moment caption
  - No temporal boundary annotations, we design the caption module and make model trainable.



Caption Module for text query reconstruction

#### Loss function

Reconstruction loss

$$L_{rec} = \frac{-1}{N_s K L} \sum_{k=1}^{N_s K} \sum_{l=1}^{L} \log P\left(w_l^* \left| F_{cro}^k, h_{l-1}^2, w_1, \cdots, w_{l-1} \right. \right)$$

Reconstruction-guided binary cross-entropy loss

Pseudo label

$$l_{k}^{j} = \frac{\sum_{l=1}^{L} \log \left( w_{l}^{*} \left| F_{cro}^{k}, h_{l-1}^{2}, w_{1}, w_{2}, \cdots, w_{l-1} \right) \right.}{\sum_{k=1}^{K} \sum_{l=1}^{L} \log \left( w_{l}^{*} \left| F_{cro}^{k}, h_{l-1}^{2}, w_{1}, w_{2}, \cdots, w_{l-1} \right) \right.}$$

$$y_k^j = \begin{cases} 0, & l_k^j \ge l_{\max} \\ 1 - l_k^j, & l_{\min} \le l_k^j < l_{\max} \\ 1, & l_k^j < l_{\min} \end{cases}$$

• RG-BCE Loss

$$L_{rg-bce} = \frac{1}{N_s K} \sum_{j=1}^{N_s} \sum_{k=1}^{K} y_k^j \log p_k^j + \left(1 - y_k^j\right) \log\left(1 - p_k^j\right)$$

#### • Experiments

- Dataset
  - Charades-STA & ActivityNet Captions
- Metric: R@K (Recall at K)
- Ablation

T-Scale	Multi-scale	IoU0.3		IoU0.5	
		R@1	R@5	R@1	R@5
64	×	44.25	63.66	27.07	51.79
64-24-8	$\checkmark$	44.52	63.70	25.00	52.05
64-24-6	$\checkmark$	47.99	66.41	21.09	44.30
64-24-4	$\checkmark$	49.79	72.57	29.68	58.66

TABLE III EXPERIMENT RESULTS WITH MULTIPLE TEMPORAL SCALES (T-SCALE).

#### • Experiments

#### Comparasion with SOTA

TERFORMANCE COMPARISON RESULTS ON CHARADES-STA DATASET.					
Method	Training	IoU0.5		IoU0.7	
		<b>R@</b> 1	R@5	R@1	R@5
Random	-	8.61	37.57	3.39	14.98
VSA-RNN	Full	10.50	48.43	4.32	20.21
VSA-STV	Full	16.91	53.89	5.81	23.58
CTRL [4]	Full	23.63	58.92	8.89	29.52
2D-TAN [27]	Full	39.70	80.32	23.31	51.26
TGA [14]	Weak	19.94	65.52	8.84	33.51
LoGAN [21]	Weak	34.68	74.30	14.54	39.11
SCN [13]	Weak	23.58	71.80	9.97	38.87
Ours	Weak	30.38	69.60	17.31	34.92

TABLE I

PERFORMANCE COMPARISON RESULTS ON CHARADES-STA DATASET

TABLE II Performance comparison results on ActivityNet Captions Dataset.

Method	Training	IoU0.3		IoU0.5	
		R@1	R@5	R@1	R@5
Random	-	18.64	52.78	7.63	29.49
VSA-RNN	Full	39.28	70.84	23.43	55.52
VSA-STV	Full	41.71	71.05	24.01	56.62
CTRL [4]	Full	47.43	75.32	29.01	59.17
2D-TAN [27]	Full	59.45	85.53	44.51	77.13
WS-DEC [3]	Weak	41.98	-	23.34	-
WSLLN [5]	Weak	42.80	-	22.70	-
SCN [13]	Weak	47.23	71.45	29.22	55.69
Ours	Weak	49.79	72.57	29.68	58.66

• Disscusion

TABLE IV EXPERIMENT RESULTS WITH DIFFERENT LOSS WEIGHT.

Loss weight	IoU0.3		IoU0.5	
Loss weight	R@1	R@5	R@1	R@5
$\lambda = 0.5$	47.09	74.62	21.63	54.01
$\lambda = 1.0$	49.79	72.57	29.68	58.66
$\lambda = 2.0$	43.85	<b>79.98</b>	24.68	59.67

• Visualization



### Conclusion

- Considering the various temporal scale of moment candidates as well as the temporal relations between them in weakly-supervised setting
- Generate more precise moment candidates with various temporal scales, facilitate weakly-supervised moments evaluation
- Future work
  - Action focus (person-centric)
  - Action re-identification



Feel free to contact us. Wenshengia.hotmail.com liding2019@ia.ac.cn