

### Two-stage Ranking Strategy for Ad-hoc Video Search

Fangming Zhou<sup>1</sup>, Changqiao Wu<sup>2</sup>, Xiaofeng Guo<sup>2</sup>, Haofan Wang<sup>2</sup>, Jincan Deng<sup>2</sup>, Debing Zhang<sup>2</sup>

<sup>1</sup>Renmin University of China <sup>2</sup>MMU, Kuaishou Technology

TRECVID 2021 workshop

2021-12

### **Ad-hoc Video Search**

• ESSENCE: text-video matching in retrieval scenario (a hot topic)





Leaderboard of video retrieval on MSR-VTT-1kA [1]



[1] MSR-VTT-1kA Benchmark (Video Retrieval) | Papers With Code

### **Previous works**

- RUC\_AIMC3 at TRECVID 2020 <sup>[1]</sup>:
  - Two-branch model
  - Addition of irCSN feature
- RUCMM at TRECVID 2020 <sup>[2]</sup>:
  - Multi-space & multi-loss strategy
  - Addition of C3d feature



RUC\_AIMC3 at TRECVID 2020



RUCMM at TRECVID 2020

[1] Zhao et al., RUC\_AIM3 at TRECVID 2020: Ad-hoc video search & video to text description. TRECVID, 2020.
 [2] Li et al., Renmin University of China at TRECVID 2020: Sentence encoder assembly for ad-hoc video search. TRECVID, 2020.



#### • Stage - I : Keyword-based Rank

- Model architecture
- Training data
- Visual features and textual encoders

#### • Stage - II : Fine-grained Re-rank

- Frame-level matching
- Weighted sum of two stages as final similarity
- Reasonableness of re-ranking



- Stage I : Keyword-based Rank
  - Model architecture

SEA: sentence encoder assembly



- ✓ Multi-space architecture
- Learning k common space for k sentence encoders
- ✓ Combined loss:

$$loss = \sum_{i=1}^{k} loss_{i}(sentence, video)$$



[1] Li et al., SEA: Sentence encoder assembly for video retrieval by textual queries. TMM, 2021.

- Stage I : Keyword-based Rank
  - Model architecture



- Stage I : Keyword-based Rank
  - Model architecture SEA++ model
    - ✓ Individual common space for each combination.
    - ✓ First stage similarity:

 $S_{first} = \sum_{i=1}^{m} \sum_{j=1}^{n} S_{space,i,j}$ 

✓ Only the first **K** videos sorted according to  $S_{first}$  will be passed to the Stage – II.





- Stage I : Keyword-based Rank
  - Model architecture SEA++ model







- Stage I : Keyword-based Rank
  - Training data
    - ✓ Concepts in different datasets should be complementary.
    - ✓ Training data should be similar to V3C1.

Training data	Num of video/image	Num of sentence
MSR-VTT	10k	200k
TGIF	100k	124k
VATEX	32k	349k
MSCOCO	123k	616k



- Stage I : Keyword-based Rank
  - Visual features and textual encoders
    - ✓ Visual features: ResNeXt101, irCSN, CLIP, timesformer
    - ✓ Textual encoders: Bag-of-word, word2vec (keyword-based)

Model	Feature	TV19	TV20
	Resnext+irCSN	0.167	0.316
SEA(BoW, w2v)	Resnext+irCSN+CLIP	0.185	0.327
	Resnext+irCSN+CLIP+timesformer	0.191	0.332

Ablation experiment on visual feature



#### • Stage - II : Fine-grained Re-rank

- Frame-level matching
- Weighted sum of two stages as final similarity
- Reasonableness of re-ranking



- Stage II : Fine-grained Re-rank
  - Frame-level matching
    - ✓ We use out-of-box CLIP <sup>[1]</sup> as framelevel matching model.
    - ✓ Only Top-K videos sorted in Stage I are considered.
    - ✓ Second stage similarity:

 $S_{second} = \max(I_1 \cdot T, \cdots, I_n \cdot T)$ 

 $I_i = ImageEncoder(frame_i)$ 

T = TextEncoder(query)





[1] Radford et al,. Learning transferable visual models from natural language supervision, arxiv, 2021.

- Stage II : Fine-grained Re-rank
  - Weighted sum of two stages as final similarity

$$S(Q,V) = \begin{cases} 0, \\ w_1 \cdot S_{first} + w_2 \cdot S_{second}, \end{cases}$$

$$if S_{first} < S_{threshold}$$
$$if S_{first} \ge S_{threshold}^{*}$$

- Reasonableness of re-ranking
  - ✓ Keyframe is enough in most cases.
  - ✓ The sentence semantics is considered (versus previous keyword-based method).

 Models
 TV19
 TV20
 TV21

 baseline
 0.211
 0.362
 0.340

 Baseline + re-rank
 0.241
 0.360
 0.349



[\*] We set  $w_1=0.2 w_2=0.8$  in our experiments.

### **Overall Pipeline**





- Our final submitted runs as followed:
  - run 3: single SEA++ model

Submissions	TV19	TV20	TV21
Winner in 2019	0.163	-	-
Winner in 2020	-	0.359	_
run3	0.206	0.354	0.332



• Our final submitted runs as followed:

- run 3: single SEA++ model
- run 2: model ensemble<sup>1</sup>

Submissions	TV19	TV20	TV21
Winner in 2019	0.163	-	_
Winner in 2020	-	0.359	-
run3	0.206	0.354	0.332
run2	0.211	0.362	0.340



• Our final submitted runs as followed:

- run 3: single SEA++ model
- run 2: model ensemble<sup>1</sup>
- run 1: model ensemble<sup>1</sup> + re-rank

Submissions	TV19	TV20	TV21
Winner in 2019	0.163	-	_
Winner in 2020	-	0.359	-
run3	0.206	0.354	0.332
run2	0.211	0.362	0.340
run1(primary)	0.241	0.360	0.349



• Our final submitted runs as followed:

- run 3: single SEA++ model
- run 2: model ensemble<sup>1</sup>
- run 1: model ensemble<sup>1</sup> + re-rank
- run 4: model ensemble<sup>2</sup> + re-rank

Submissions	TV19	TV20	TV21
Winner in 2019	0.163	-	_
Winner in 2020	-	0.359	_
run3	0.206	0.354	0.332
run2	0.211	0.362	0.340
run1(primary)	0.241	0.360	0.349
run4	0.239	0.358	0.349



### **Take-home Message**

 1) We propose an improved video retrieval model, namely SEA++, which built a solid backbone for our best run.



SEA++ model



### **Take-home Message**

- 1) We propose an improved video retrieval model, namely SEA++, which built a solid backbone for our best run.
- 2) Re-ranking by CLIP is an effective method to gain higher performance.



# THANKS!

#### Contact with us:

fangming\_zhou@ruc.edu.cn, wuchangqiao@kuaishou.com, zhangdebing@kuaishou.com

