

# CRET: Cross-Modal Retrieval Transformer for Efficient Text-Video Retrieval

Authors: Kaixiang Ji, Jiajia Liu, Weixiang Kong, Liheng Zhong, Jian Wang, Jingdong Chen, Wei Chu

Siwen Tu and Shupef Li

Leiden Institute of Advanced Computer Science  
April 11, 2023



**Universiteit  
Leiden**  
The Netherlands

① Motivation

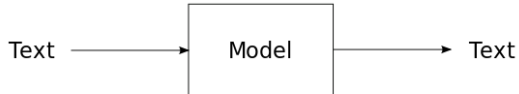
② Methodology

③ Experiments

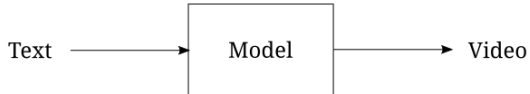
④ Critical Review

# Motivation

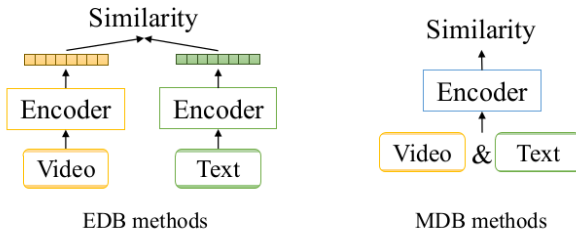
- Unimodal information retrieval task.



- Multimodal information retrieval task, e.g. text-to-video retrieval.



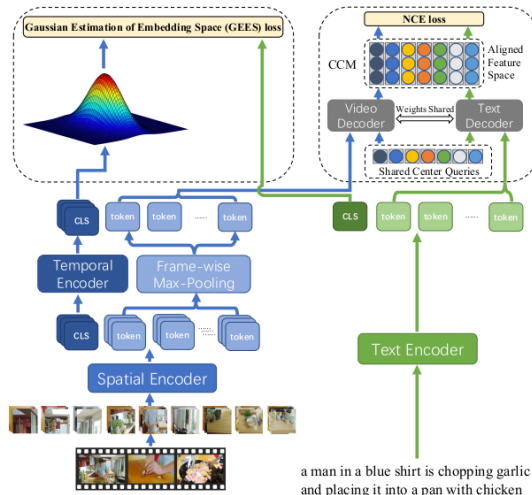
# EDB method versus MDB method



(a) EDB and MDB methods mainly differ at whether explicit embeddings of text/video are generated.

- The EDB method utilizes embeddings extracted from the text and videos to model the evaluation of distance.
- The MDB method iterates all text-video pairs to evaluate the distance without extracting embeddings.

# Overview of the CRET model



# Two main contributions: CCM & GEES

## Cross-modal correspondence modeling (CCM)

- Utilize transformer decoders to align the features from text and video modalities.
- Use queries as common centers of features from both modalities.
- Parameters are shared between decoders of two modalities.

## Gaussian estimation of embedding space (GEES)

$$\mathbf{Z}_{c,j} = \text{softmax} \left( \frac{(Q_c W_j^Q)(E W_j^K)^T}{\sqrt{d_k}} \right) (E W_j^V)$$

- ① Calculate the distance between token features and the query center. Regard the distance as the weight of corresponding features.
- ② Concatenate and project the aligned features.
- ③ Calculate the similarity score of features from text and video modalities.

# Experiments

- **Datasets:** MSRVT, LSMDC, MSVD, DiDeMo.
- **Metrics:** R@K, MdR.
- **Results:**
  - MSRVT

	Weight Initialization		E2E	Type	R@1↑	R@5↑	R@10↑	MdR↓
	Visual	Textual						
MIL-NCE [36]	K+H	G+H	✓	EDB	9.9	24.0	32.4	29.5
JSFusion [54]	I	N.A.	✓	MDB	10.2	31.2	43.2	13.0
HT [38]	I	G	✓	EDB	12.4	36.0	52.0	10.0
HT [38]	I+H	G+H	✓	EDB	14.9	40.2	52.8	9.0
ActBERT [59]	I+H	B+H		MDB	16.3	42.8	56.9	10.0
HiT(appearance-only) [30]	I+H	B+H	✓	EDB	18.2	41.9	55.5	5.0
TACo(R-152) [53]	I+H	B+H	✓	MDB	18.9	46.2	58.8	7.0
UniVL(FT-Joint) [33]	K+H	B+H		EDB	20.6	49.1	62.9	6.0
UniVL(FT-Align) [33]	K+H	B+H		MDB	21.2	49.6	63.1	6.0
ClipBERT [28]	I+C+V	C+V	✓	MDB	22.0	46.8	59.9	6.0
<b>Ours</b>	I	B	✓	EDB	<b>23.9</b>	<b>50.8</b>	<b>63.4</b>	<b>5.0</b>

# Experiments

- Results:
  - LSMDC

	E2E	Type	R@1↑	R@5↑	R@10↑	MdR↓
CT-SAN [55]	✓	MDB	5.1	16.3	25.2	46.0
HT [38]	✓	EDB	5.8	18.8	28.4	45.0
HT* [38]	✓	EDB	7.1	19.6	27.9	40.0
NoiseE* [1]		EDB	6.4	19.8	28.4	39.0
JSFusion [54]	✓	MDB	9.1	21.2	<b>34.1</b>	36.0
<b>Ours</b>	✓	EDB	<b>10.0</b>	<b>24.9</b>	33.4	<b>34.0</b>

- MSVD

	E2E	Type	R@1↑	R@5↑	R@10↑	MdR↓
HT [38]	✓	EDB	13.0	37.4	52.4	10.0
HT* [38]	✓	EDB	15.5	40.9	55.7	8.0
NoiseE* [1]		EDB	20.3	49.0	63.3	6.0
CLIP4Clip <sup>†</sup> [34]		EDB	46.2	76.1	84.6	2.0
<b>Ours</b>	✓	EDB	<b>49.0</b>	<b>87.0</b>	<b>95.0</b>	<b>2.0</b>



# Experiments

- **Results:**
  - DiDeMo

	E2E	Type	R@1↑	R@5↑	R@10↑	MdR↓
S2VT [49]	✓	EDB	11.9	33.6	-	13.0
FSE [57]	✓	EDB	13.9	36.0	-	11.0
ClipBERT <sup>‡</sup> [28]	✓	MDB	20.4	48.0	60.8	6.0
<b>Ours</b>	✓	EDB	<b>21.2</b>	<b>50.3</b>	<b>63.5</b>	<b>6.0</b>

- **Ablation studies.**
- **Validation of Gaussian assumption.**

# Critical review

- Readability and structure.
  - Illustrate CRET method clearly.
  - Satisfy requirements of the scientific paper.
- Reproducibility: Source code, the availability of data sets, experimental settings.
- Importance.
  - Theoretical contributions.
  - Practical applications.
- Summary of strong and weak points.