

Boosting Video-Text Retrieval with Explicit High-Level Semantics

Haoran Wang*[†]

wanghaoran09@baidu.com
Department of Computer Vision
Technology (VIS), Baidu Inc.
Beijing, China

Fu Li

lifubaidu.com
Department of Computer Vision
Technology (VIS), Baidu Inc.
Beijing, China

Di Xu*[‡]

xudi20s@ict.ac.cn
Institute of Computing Technology,
Chinese Academy of Sciences
Beijing, China

Zhong Ji

jizhong@tju.edu.cn
School of Electrical and Information
Engineering, Tianjin University
Tianjin, China

Errui Ding

dingerrui@baidu.com
Department of Computer Vision
Technology (VIS), Baidu Inc.
Beijing, China

Dongliang He

hedongliang01@baidu.com
Department of Computer Vision
Technology (VIS), Baidu Inc.
Beijing, China

Jungong Han

jungonghan77@gmail.com
Computer Science Department,
Aberystwyth University
SY23 3FL, UK

ABSTRACT

Video-text retrieval (VTR) is an attractive yet challenging task for multi-modal understanding, which aims to search for relevant video (text) given a query (video). Existing methods typically employ completely heterogeneous visual-textual information to align video and text, whilst lacking the awareness of homogeneous high-level semantic information residing in both modalities. To fill this gap, in this work, we propose a novel visual-linguistic aligning model named HiSE for VTR, which improves the cross-modal representation by incorporating explicit *high-level semantics*. First, we explore the hierarchical property of explicit high-level semantics, and further decompose it into two levels, *i.e.* discrete semantics and holistic semantics. Specifically, for visual branch, we exploit an off-the-shelf semantic entity predictor to generate discrete high-level semantics. In parallel, a trained video captioning model is employed to output holistic high-level semantics. As for the textual modality, we parse the text into three parts including occurrence, action and entity. In particular, the occurrence corresponds to the holistic high-level semantics, meanwhile both action and entity represent the discrete ones. Then, different graph reasoning techniques

are utilized to promote the interaction between holistic and discrete high-level semantics. Extensive experiments demonstrate that, with the aid of explicit high-level semantics, our method achieves the superior performance over state-of-the-art methods on three benchmark datasets, including MSR-VTT, MSVD and DiDeMo.

CCS CONCEPTS

• **Information systems** → **Information retrieval; Retrieval models and ranking.**

KEYWORDS

Video-Text Retrieval, High-level Semantics, Vision-language Understanding

ACM Reference Format:

Haoran Wang, Di Xu, Dongliang He, Fu Li, Zhong Ji, Jungong Han, and Errui Ding. 2022. Boosting Video-Text Retrieval with Explicit High-Level Semantics. In *Proceedings of the 30th ACM International Conference on Multimedia (MM '22), October 10–14, 2022, Lisboa, Portugal*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3503161.3548010>

1 INTRODUCTION

With the exponentially increasing of on-line videos, video-text retrieval (VTR) is becoming an emerging requirement for people to perceive the world. It refers to searching for a video (text) when given a query text (video), which plays a fundamental role in vision and language understanding [3, 4, 18, 20, 24, 24, 29, 34, 40, 41, 43, 48]. Although remarkable progresses in this area have been made, it is still challenging to precisely match video and text since the raw input multi-modal data exist in heterogeneous spaces.

The main challenge for video-text alignment is how to narrow the gap between the heterogeneous representations from both modalities. To tackle this problem, previous solutions can be divided

*indicates equal contribution.

[†]indicates corresponding author.

[‡]Work done while Di Xu was a Research Intern at VIS, Baidu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '22, October 10–14, 2022, Lisboa, Portugal

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9203-7/22/10...\$15.00

<https://doi.org/10.1145/3503161.3548010>

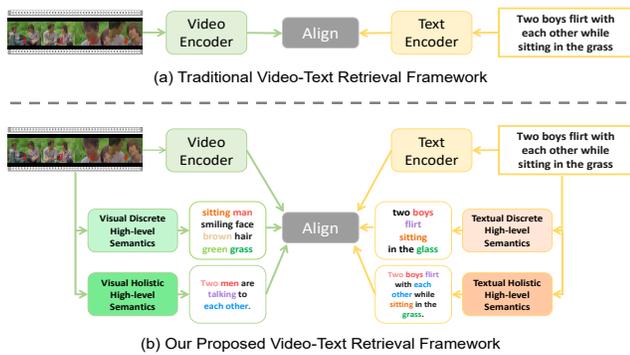


Figure 1: The conceptual comparison between traditional cross-modal aligning paradigm for video-text retrieval (VTR) (a) and our proposed VTR framework (b).

into two streams of works. One line of studies is built based on *single branch* representation architecture. These methods [31, 33, 50] typically employ an independent modality to represent video, *i.e.* appearance and utilize single encoder to represent it, whereas ignoring the that video is actually constituted by multiple constituent modalities. In contrast to these approaches, some multi-expert based studies [11, 14, 27] have been made to enhance video representation by introducing more complementary clues, such as *motion*, *audio* and *speech*, and achieve steady improvements.

While encouraging achievements have been made, most existing methods are restricted by relying on heterogeneous information to align video and text. But how to leverage homogeneous clues, *i.e.* **explicit high-level semantics (EHS)**, to perform cross-modal alignment is still unexplored. By comparison, in cognitive science [9], studies explain that humans can recognize worlds by extracting high-level descriptions, which have been shown effective in several vision-language understanding tasks [12, 17, 46, 49]. Therefore, it is worthy of exploring how to leverage EHS to improve VTR.

To fill this gap, we take a step towards exploiting homogeneous knowledge to improve VTR task, which is achieved by injecting EHS information into cross-modal representation learning. Firstly, we explore the hierarchical property of EHS, and further categorize it into two levels: *discrete semantics* and *holistic semantics*. Concretely, the former focuses on mining local and detailed semantics, which is represented by separate entities or phrases. Meanwhile, the latter is denoted by an entire describing sentence, which aims to extract more global and consecutive information. As illustrated in Figure 6, both types of EHS extracted from video are properly consistent with those from text, which are beneficial for bridging the modality gap.

In this work, to bridge the semantic discrepancy between video and text, we propose a high-level semantics aided (**HiSE**) visual-linguistic embedding model for VTR application, which injects the EHS into cross-modal representation learning. To realize it, we first elaborately design different architectures to collect hierarchical EHS information for video and text, respectively. As for the visual modality, we employ two separate components to generate the discrete and holistic EHS. Specifically, to extract discrete EHS from video, we leverage an off-the-shelf semantic entity predictor [1] to detect semantic entities for all sampled frames. Then, the obtained

entities are filtered according to their confidence scores and the retained part of them is taken as discrete EHS. Meanwhile, we employ a video captioning model to generate describing sentence for the given video. The output is taken as visual holistic EHS that contains more coherent information.

On the other side, to extract the hierarchical EHS from text, we utilize a semantic role parser to decompose it into a role graph, which includes three kinds of nodes: *occurrence*, *action* and *entity*. The occurrence refers to the whole sentence, which consists of more context information and is taken as textual holistic EHS. Besides, action is denoted by verbs and entity is represented by noun & noun phrases. Considering action and entity only describe local textual information, we name them as textual discrete EHS.

In addition to EHS acquisition, how to organize and integrate them also plays important role in cross-modal representation learning. For the video part, given the video discrete EHS, we use its textual embeddings in conjunction of a graph convolution operation to obtain the compact vectors, termed video discrete high-level semantic (VDS) representation. As for video holistic EHS representing, we employ another sentence-level text encoder to process it, which outputs the video holistic high-level semantic (VHS) representation. By contrast, for the textual branch, we uniformly employ the corresponding text embeddings to denote both discrete EHS and holistic EHS by compact vectors. Then, based on the role graph structure, we apply multi-graph reasoning to promote the interaction between the output embeddings from both levels. Lastly, we fuse the generated EHS representations with the raw outputs of modality-specific encoders by using their convex combination. Through integrating the EHS knowledge into the cross-modal representation framework, the bidirectional VTR performance can be remarkably boosted. On the whole, our main contributions lie in three-fold.

- We present a **High-level Semantic (HiSE)** aided visual-linguistic embedding framework for VTR. To the best of our knowledge, this is the first work integrating explicit high-level semantics into video-text retrieval, which leverages the unified abstract information to strength the semantic relationship between two modalities.
- We decompose high-level semantics into two levels: *i.e.* discrete semantics and holistic semantics, which are responsible to capture local and global information respectively. Moreover, we design elaborated architectures to represent and tie up them to improve the cross-modal representation.
- The extensive experiments on benchmark datasets not only demonstrate the superiority of our approach by outperforming state-of-the-art methods for VTR, but also exhibits the interpretability.

2 RELATED WORK

2.1 Video-Text Retrieval

Accompanied by the renaissance of deep learning, there have been growing interest in research for VTR [5, 14, 22, 27, 31, 34, 53]. A majority of early works tackle this task from perspective of representation architecture. On one side, some works [11, 14, 27, 31] introduce multiple experts to enhance video representation. For

instance, MoEE [31] jointly combined three sources of expert components for video encoding, including videos, motion and audio feature. On the other side, several studies [21, 22] focused on text encoding designing. Li *et al.* [22] introduced multiple sentence encoders and combining similarities from all text encoder-specific joint space. Based on the multi-expert representations, more works [8, 45] devoted to performing hierarchical alignments for VTR, and achieved continuous performance improvements.

Recently, since the prevalence of large-scale pre-training technique, a flurry of works [6, 13, 28, 35] leveraging this technique to promote video-text representation have emerged. As a representative study, Luo *et al.* [28] proposed a CLIP4CLIP method, which transfers the knowledge of the CLIP model to VTR application via an end-to-end fine-tuning. Then, Fang *et al.* [13] presented to plug a temporal information capturing module in CLIP4CLIP for video representation enhancing. Cao *et al.* [6] proposed a framework modelling visual consensus, which exploits commonsense information in both vision and language domain to improve VTR. To sum up, all above approaches perform cross-modal alignment based on completely heterogeneous video and text representations. Distinct from them, our HiSE additionally introduces homogeneous high-level semantic clues to narrow the modality gap.

2.2 Multi-modal Understanding Using High-level Semantics

High-level semantics plays critical role in multi-modal data understanding [12, 17, 17, 42, 46, 49, 51]. For instance, Wu *et al.* [46] proposed a method that incorporates semantic concepts into the CNN-RNN architecture to improve image captioning. Similarly, Yao *et al.* [49] integrated attributes into the CNN-RNN image captioning framework, followed by training in an end-to-end manner. For VQA, Hudson *et al.* [17] presented to exploit high-level scene-graph knowledge to transform both image and text into semantic concept-based representations. In contrast to previous works, to the best of our knowledge, we make first attempt to simultaneously integrate explicit high-level semantics into video and text representations for VTR. The most related work to ours is SCO [16], which leveraged concepts to enhance image representation for image-text retrieval. On the contrast, we further extend high-level semantics of video to discrete semantics and holistic one, and validate the effectiveness of the hierarchical information on VTR task.

3 METHODOLOGY

In this section, we detailedly introduce our proposed high-level semantic (HiSE) aided visual-linguistic embedding model for video-text retrieval (see Figure 2). Firstly, we illustrate the modality-specific encoders to represent video and text along with their corresponding memory banks. Then, we illustrate the concrete architecture of two explicit high-level semantics representation modules, *i.e.* VSE module and TSE module, respectively. Afterwards, the cross-modal representation aggregating manner, inference method and alignment objectives are introduced sequentially.

3.1 Video and Text Encoder

3.1.1 Video Encoder. To obtain the video representation, we first extract the frames from the video clip, and then employ a video

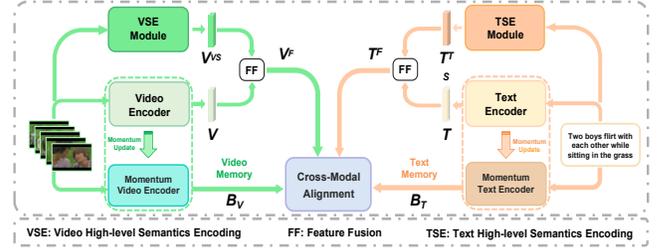


Figure 2: The overall architecture of proposed HiSE model for video-text retrieval. Taking the input video and text, on one hand, it simultaneously employ modality-specific encoders and dynamic memory banks to generate cross-modal representations. On the other hand, it additionally introduces two explicit high-level semantics encoding modules, *i.e.* VSE module and TSE module, to improve cross-modal alignment.

encoder to turn them into a sequence of features, followed by fusing them by a feature aggregator. To realize it, we adopt pre-trained CLIP [36] (ViT-B/32) as video encoder, which utilizes the ViT [10] backbone pre-trained based on 400 million image-text pairs. Specifically, we split an image into non-overlapping patches, then use a linear projection to project them into 1-D tokens. The output vectors are taken as input of ViT, which leverages the transformer architecture to model the interaction between image patches. Following CLIP, we adopt the output from the [class] token as the image representation. Consequently, given the input sequence of video frames $O = \{o_1, \dots, o_N\}$, the generated video features are denoted as $\bar{V} = \{v_1, \dots, v_N\}$. Afterwards, to further capture the temporal information between frames, we use another Transformer encoder [39] with position embedding to combine frames into one global video representation. Formally, the global video representation is calculated by $V = f_{video}(O)$, in which $f_{video}(\cdot)$ represents the video encoder.

3.1.2 Text Encoder. For caption encoding, we directly employ the text encoder from the CLIP to extract the textual representation. In particular, it refers to a Transformer with architecture modifications according to [36]. Following CLIP [36], we use the activations of the [EOS] token from the most top layer of the transformer as the global representation of the caption. Given the caption S , the global textual representation is computed according to $T = f_{text}(S)$, where $f_{text}(\cdot)$ denotes the text encoder.

3.1.3 Modality-specific Memory Banks Building. To further enhance the negative interactions for both modalities in contrastive learning, we propose to leverage two dynamic modality-specific memory banks B_V and B_T to store additional video and text representations. Particularly, we follow MoCo [15] to obtain momentum video encoder and text encoder by momentum updating their weights according to the corresponding modality-specific encoders, whose architectures are totally same as $f_{video}(\cdot)$ and $f_{text}(\cdot)$, respectively. As such, video or text samples from the latest training iterations are fed to the momentum encoders, which output video

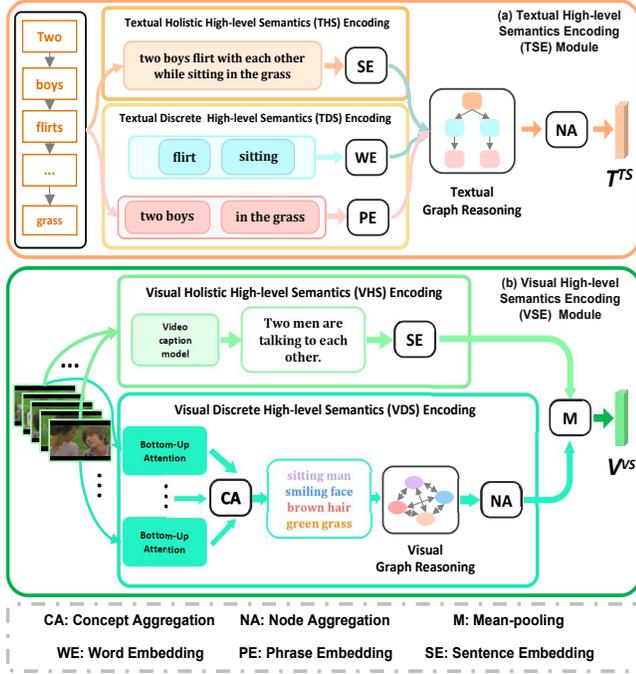


Figure 3: The concrete architecture of VSE and TSE modules in our HiSE model. The VSE module aims at encoding high-level semantic information from video, and the TSE module is designed to represent high-level semantics of text.

and text embeddings and restored them in coupled modality-specific memory banks. This process is implemented via deploying queues.

3.2 Textual High-level Semantics Encoding

In this part, we describe the detailed structure of Textual High-level Semantics Encoding (TSE) module (see Figure. 3(a)). We know that, video captioning sentence is naturally constructed based on hierarchical semantic structure, which can be expressed by the relationships between the *actions* and *entities* to describe *occurrence* (see Figure 3(a)). Such hierarchical structure is beneficial to realize comprehensive understanding for video captions. Considering these relationships can be conveyed by graph structures, we exploit graph convolution to encode the explicit high-level semantics for text. The processing pipeline includes three steps: 1) Initializing semantic role graph node; 2) Graph building; 3) Graph Reasoning.

Initializing semantic role graph node. First, the graph nodes are divided into three classes: *occurrence* nodes, *action* nodes, and *entity* nodes. Specifically, the occurrence node representations are initialized by employing the pre-trained CLIP text encoder, which shares the same architecture with the encoder in 3.1.2 (not sharing weights). To initialize action and entity nodes, we use tokenizer of CLIP text encoder to split and map them into token embeddings, then aggregate the token embedding vectors by mean-polling operation to denote phrase-level contents (entity). Finally, the initialized occurrence/action/entity nodes are denoted by $E_o/E_a/E_e$, respectively.

Role graph building. Given a video description S , we take it as the occurrence node in semantic role graph. Then, an off-the-shelf semantic role parser [38] is leveraged to extract noun phrases and verbs from S along with their semantic relationships. In our role graph, the noun phrases and verbs act as action nodes and entity nodes, respectively. The entity nodes are connected with different action nodes, where the edge type between them is defined by the semantic role of the entity in reference to the action. The action nodes all connect to the occurrence node with direct edges, in which the temporal clues between different actions can be captured.

Graph Reasoning. To effectively exploit the multiple roles of relational information contained in our graph, we adopt the relational graph convolutional network (R-GCN) [37] to model interactions in graph between nodes. Specifically, given the initialized nodes $E_i = \{E_o, E_a, E_e\}$ and role graphs $G = \{G^r\}$, through one-layer graph convolution with residual connection, the node embedding feature is computed as:

$$E_i^t = \rho \left(\sum_{r \in R} G^r E_i^{t,0} W^r + E_i^0 \right) \quad (1)$$

where $G^r \in G$ ($r \in R$) denotes the semantic role matrix under relation of r and R is the number of relation roles. W^r represents the learnable weight matrix under relation r . ρ is a ReLU function. To distinguish them, we term E_o^t and the sum vector of E_a^t , E_e^t as Textual Holistic High-level Semantic (THS) representation and Discrete High-level Semantic (TDS) representation. Finally, THS and TDS representations are fused by mean-polling to generate vector T^{TS} , namely Textual High-level Semantic (TS) representation, which simultaneously encodes holistic and discrete semantics from three types of nodes.

3.3 Visual High-level Semantics Encoding

In this section, we elaborate on the details of visual high-level semantics encoding (VSE) module (see Figure. 3(b)), which consists of two parallel sub-modules to encode discrete and holistic high-level semantics into video representation, respectively.

3.3.1 Visual Discrete Semantics Encoding. Video is made up sequence of consecutive frames, in which each frame composes of multiple subjects and background objects. These entities can express most informative semantics from perspective of local information acquiring. Inspired by these observations, we take these entities as discrete semantics of video and propose to incorporate it into the video representation learning.

Given sampled video frames $O = \{o_1, \dots, o_N\}$, for each sampled frame o_i , we adopt the bottom-up attention toolkit [1] to extract visual entities $E_i^v = \{E_{i,1}^v, E_{i,2}^v, \dots, E_{i,M}^v\}$ ($i \in [1, N]$). Specifically, the visual entity E^v includes three properties, including concept representation C^v , appearance representation A^v and position representation P^v . Firstly, the concept property is made up of an object name C_o^v and its decorated attribute C_a^v , which is completely homogeneous to the textual high-level semantics. Accordingly, we generate concept representation C^v as follows:

$$C^v = MLP(S_o + S_a) \quad (2)$$

where $MLP(\cdot)$ denotes a fully-connected layer.

Then, we introduce two additional information to be complementary for concept representation, *i.e.* appearance representation and position representation. As for appearance representation, we employ a fully-connected layer to project ROI (Region of Interest) feature into the joint space, obtaining the appearance representation \mathbf{A}^v . Similarly, another fully-connected layer is used to produce position representation \mathbf{P}^v based on the spatial coordinate of the entity \mathbf{E}^v . Afterwards, given three types of embeddings $\{\mathbf{C}^v, \mathbf{A}^v, \mathbf{P}^v\}$, similar to TSE module, we also utilize graph convolution operation to encode high-level semantics for video, including three steps: 1) Initializing semantic relation graph node; 2) Graph building; 3) Graph Reasoning.

Initializing semantic relation graph node. Different with textual semantic role graph node, visual relation graph only contain one type of nodes, *i.e.* entity nodes. Given the entity set $\mathbf{E}_{i,j}^v, i \in [1, N], j \in [1, M]$ extracted from video frame sequence, we only select top- K entities $\{\mathbf{E}_1^v, \mathbf{E}_2^v, \dots, \mathbf{E}_k^v\}$ according to their appearing frequency in video. Then, the entity representation can be defined by aggregating its three types of embeddings as follows:

$$\mathbf{E}_e = MLP([\mathbf{C}^v, (\mathbf{A}^v + \mathbf{P}^v)]) \quad (3)$$

where $[\cdot]$ represents the concatenate operation. The output \mathbf{E}_e serves as the initial node representation in graph.

Relation graph building. To further capture the associations between entities, we first build up an entity affinity graph. Concretely, given a set of visual entities $\mathbf{E}_i^v = \{\mathbf{E}_{i,1}^v, \mathbf{E}_{i,2}^v, \dots, \mathbf{E}_{i,M}^v\}, i \in [1, N]$, we measure the affinity between pairwise entities as follows:

$$\mathbf{h}(E_i, E_j) = \frac{\varphi(E_i) \phi(E_j)^T}{\sqrt{D}} \quad (4)$$

where $\varphi(E_i) = W_\varphi E_i$ and $\phi(E_j) = W_\phi E_j$ are two node embeddings, and \sqrt{D} is the dimension of graph nodes. W_φ and W_ϕ both represent embedding matrix. Then, we obtain a fully connected graph $\mathbf{H} = (\mathbf{V}, \mathbf{E})$, where \mathbf{V} is graph nodes initialized by visual entities and \mathbf{E} denotes the affinity matrix \mathbf{H} calculated from Eq. 4. $\mathbf{h}(\cdot, \cdot)$ indicates the strength of the affinity between two nodes in graph \mathbf{H} .

Graph Reasoning. Given the relation graph \mathbf{H} , we adopt one-layer Graph Convolutional Network(GCN) [19] with residue connection to promote the node embedding learning, which is defined as follows:

$$\mathbf{E}_e^v = \rho \left(\mathbf{H}^0 \mathbf{E}_e^{v,0} \mathbf{W}^v + \mathbf{E}_e^0 \right) \quad (5)$$

where \mathbf{W}^v is the weight matrix of the GCN layer and ρ is a ReLU function. Consequently, we can obtain the output node embedding \mathbf{E}_e^v , dubbed Visual Discrete High-level Semantic (VDS) representation.

3.3.2 Visual Holistic Semantics Encoding. Distinct from VDS encoding, extracting visual holistic semantics requires more comprehensive and generalized video understanding. Accordingly, the research target of video captioning [40, 55] is consistent with this goal. It refers to automatically generating natural language descriptions of videos, which is tightly associated with VTR task. Inspired by this observation, we attempt to transfer the prior knowledge restored in video captioning model to encode visual holistic semantics. In particular, we directly adopt the off-the-shelf video captioning

model [23] to produce the natural language description for video. Given an output captioning sentence \mathbf{S}^H , we also use the CLIP text encoder described in section 3.1.2 to encode it. The generated feature vector \mathbf{V}^{VHS} is termed as Visual Holistic Semantic (VHS) representation. Lastly, we aggregate VDS representation and VHS representation by mean-pooling, and thus obtain the final visual high-level semantic (VS) representation \mathbf{V}^{VS} .

3.4 Cross-Modal Representations Fusion

Given the original video-text representation $\mathbf{V}(\mathbf{T})$ and high-level semantic representation $\mathbf{V}^{VS}(\mathbf{T}^{VS})$, we use a simple convex combination operation to aggregate them as the final cross-modal representations, which can be defined as:

$$\begin{aligned} \mathbf{V}^F &= \alpha \mathbf{V} + (1 - \alpha) \mathbf{V}^{VS}, \\ \mathbf{T}^F &= \alpha \mathbf{T} + (1 - \alpha) \mathbf{T}^{VS}, \end{aligned} \quad (6)$$

where α is a tuning parameter balancing two types of representations. And \mathbf{v}^F and \mathbf{t}^F respectively denote the fused video and text representations.

3.5 Training and Inference

As for training objective, we employ the hubness-aware contrastive loss (HAL) [25] for aligning video and text : Given video representation $\mathbf{V} = \{\mathbf{V}_1, \dots, \mathbf{V}_Q\}$ and text representation $\mathbf{T} = \{\mathbf{T}_1, \dots, \mathbf{T}_R\}$, loss function $L_{HAL}(\mathbf{V}, \mathbf{T})$ can be formulated as:

$$\begin{aligned} L_{HAL} &= \frac{\mu}{Q} \sum_{q=1}^Q [\log(\sum_{r \neq q} \exp(\frac{(S_{qr} - \gamma)}{\mu}) + 1) - \log(S_{qq} + 1)] + \\ &\quad \frac{\mu}{R} \sum_{r=1}^R [\log(\sum_{q \neq r} \exp(\frac{(S_{rq} - \gamma)}{\mu}) + 1) - \log(S_{rr} + 1)]; \end{aligned} \quad (7)$$

where γ is a margin parameter; μ is a temperature parameter; N denotes the number of samples within the mini-batch; $S_{qr} = \cos(\mathbf{V}_q, \mathbf{T}_r)$, $S_{rq} = \cos(\mathbf{T}_r, \mathbf{V}_q)$, $S_{qq} = \cos(\mathbf{V}_q, \mathbf{T}_q)$ and $S_{rr} = \cos(\mathbf{T}_r, \mathbf{V}_r)$, where $\cos(\cdot, \cdot)$ represents similarity function calculating cosine distance.

To enhance cross-modal learning, two types of HAL loss are utilized. First, it is imposed on mini-batch data. Secondly, it is imposed on anchor sample in mini-batch and items from modality-specific memory banks. Formally, the final objective of our HiSE model is defined as:

$$L = \lambda_1 L_{HAL}(\mathbf{V}^F, \mathbf{T}^F) + \lambda_2 L_{HAL}(\mathbf{V}, \mathbf{B}_T) + \lambda_2 L_{HAL}(\mathbf{T}, \mathbf{B}_V). \quad (8)$$

where balancing parameters are set to $\lambda_1 = 10$ and $\lambda_2 = 0.1$.

For inference, we use the cosine distance between fused representations \mathbf{V}^F and (\mathbf{T}^F) to measure the cross-modal relevance.

4 EXPERIMENTS

4.1 Dataset and Settings

4.1.1 Datasets. We conduct experiments on three benchmark datasets for bidirectional video-text retrieval, including MSR-VTT [47], MSVD [7] and DiDeMo [2].

- **MSR-VTT** dataset contains 10K videos from YouTube website, with each video annotated with five 20 captions. We follow the 1k-A protocol in [50] and report the experimental

Table 1: Comparisons of Experimental Results on MSR-VTT 1k-A Testing Set. ★: The results of [54] are reported by the model with video encoder of ViT-B/32 for fair comparison.

Approach	Text-to-Video Retrieval				Video-to-Text Retrieval				R@Sum
	R@1	R@5	R@10	MdR	R@1	R@5	R@10	MdR	
Non-CLIP Based									
JSFusion[50]	10.2	31.2	43.2	13.0	-	-	-	-	-
CE[27]	20.9	48.8	62.4	6.0	20.6	50.3	64.0	5.3	267.0
MMT[14]	24.6	54.0	67.1	4.0	24.4	56.0	67.8	4.0	293.9
Support-Set (pretrained)[34]	30.1	58.5	69.3	3.0	28.5	58.6	71.6	3.0	316.6
HIT (pretrained)[26]	30.7	60.9	73.2	2.6	32.1	62.7	74.1	3.0	333.7
FROZEN[5]	31.0	59.5	70.5	3.0	-	-	-	-	-
CLIP Based									
CLIP[35]	31.2	53.7	64.2	4.0	27.2	51.7	62.6	5.0	290.6
CLIP4Clip-meanP[28]	43.1	70.4	80.8	2.0	43.1	70.5	81.2	2.0	389.1
CLIP4Clip-seqTransf[28]	44.5	71.4	81.6	2.0	42.7	70.9	80.6	2.0	391.7
VCM[6]	43.8	71.0	80.9	2.0	45.1	72.3	82.3	2.0	395.4
CenterCLIP [54] ★	44.2	71.6	82.1	2.0	42.8	71.7	82.2	2.0	394.6
HiSE	45.0	72.7	81.3	2.0	46.6	73.3	82.3	2.0	401.2

Table 2: Comparisons of Experimental Results on MSVD Testing Set.

Approach	Text-to-Video Retrieval				Video-to-Text Retrieval				R@Sum
	R@1	R@5	R@10	MdR	R@1	R@5	R@10	MdR	
Non-CLIP Based									
VSE[33]	12.3	30.1	42.3	14	34.7	59.9	70.0	33	249.3
CE[27]	19.8	49.0	63.8	-	-	-	-	-	-
MoEE[31]	21.1	52.0	66.7	5.0	27.3	55.1	65.0	4.3	287.2
TT-CE+[9]	25.4	56.9	71.3	4.0	27.1	55.3	67.1	4.0	303.1
Support-Set (pretrained)[34]	28.4	60.0	72.9	4.0	34.7	59.9	70.0	3.0	325.9
FROZEN[5]	31.0	59.5	70.5	3.0	-	-	-	-	-
CLIP Based									
CLIP[35]	37.0	64.1	73.8	3.0	59.9	85.2	90.7	1.0	410.7
CLIP4Clip-meanP[28]	46.2	76.1	84.6	2.0	56.6	79.7	84.3	1.0	427.5
CLIP4Clip-seqTransf[28]	45.2	75.5	84.3	2.0	62.0	87.3	92.6	1.0	446.9
HiSE	45.9	76.2	84.6	2.0	66.3	90.7	95.2	1.0	458.9

results. Specifically, 1k-A protocol adopts 9,000 videos for training and utilizes rest 1,000 video-text pairs for testing.

- **MSVD** dataset includes 1,970 videos with approximately 80,000 captions. It is split into 1,200 training, 100 validation, and 1000 testing videos. We report the performance of video-text retrieval on testing set with multiple captions per video.
- **DiDeMo** contains 10,000 videos, with each video annotated by 40 sentences. We follow [27] to conduct video-paragraph retrieval, where all the descriptions of one video are combined into a single text.

4.1.2 Evaluation Metrics. For evaluation, we employ three standard retrieval criteria: recall at rank K ($R@K$, higher is better), median rank (MdR, lower is better) and sum of recall ($R@sum$, higher is better). In particular, $R@K$ measures the the fraction of queries for which the matched item is found among the top K retrieved results. The MdR denotes the median rank of correct items in the retrieved

ranking list. Besides, $R@sum$ criterion is calculated by summing all metrics of $R@K$, which can better reflect the overall performance.

4.2 Implementation Details

For text encoding, the basic text transformer and occurrence node encoder in THS module are both initialized by CLIP text transformer. For video encoding, the spatial transformer (ViT) is initialized with CLIP (ViT/B-32). The dimension of the joint embedding space is set to 512. The caption token length is 32 and frame length is 12. Note that on DiDeMo dataset, the captioning sentences are contacted as one paragraph for video-paragraph retrieval, where the caption token length is set to 64. The size of memory banks is set to 4096 and the momentum coefficient is equal to 0.995. The fusion ratio parameter α in Eq.6 is empirically set to 0.9. In alignment objective, we follow [25] to set $\gamma = 0.3$ and $\mu = 0.1$ in Eq.7. Our model is fine-tuned by Adam optimizer with mini-batch size of 256. As for

Table 3: Performance of our HiSE method with different representation components on MSR-VTT 1k-A test set. Visual Discrete and Holistic High-level Semantics are abbreviated as “VDS” and “VHS”, respectively. Textual Discrete and Holistic Semantics are abbreviated as “VHS” and “THS”, respectively.

High-level Semantics Encoding Module				Text-to-Video Retrieval			Video-to-Text Retrieval		
VDS	VHS	TDS	THS	R@1	R@5	R@10	R@1	R@5	R@10
				43.6	72.1	80.9	44.5	72.2	82.0
			✓	43.6	72.4	81.2	44.7	72.3	81.3
		✓		43.6	72.5	81.2	44.6	72.2	81.3
		✓	✓	43.8	72.4	81.3	44.9	72.3	81.4
	✓	✓	✓	44.1	72.4	81.2	45.4	72.5	81.4
✓		✓	✓	44.7	72.6	81.3	46.0	72.8	81.9
✓	✓	✓	✓	44.6	72.8	81.2	46.4	73.1	82.2
✓	✓	✓	✓	45.0	72.7	81.3	46.6	73.3	82.3

the learning rate, we follow the CLIP [36] to decay it with a cosine schedule. The initial learning rate is $1e-7$ for basic text encoder and video encoder and $1e-4$ for other modules. All our experiments are implemented on 8 NVIDIA Tesla P40 GPUs.

4.3 Comparison to State-of-the-art Methods

The experimental results on MSR-VTT and MSVD datasets are listed in Table 4 and Table 5, respectively. Note that the results on DiDeMo dataset is placed in the supplementary materials due to limited space.

On MSR-VTT dataset, as shown in Table 4, we can see our HiSE outperforms the competitors in most evaluation metrics. For text retrieval, compared with the second best method, we achieve absolute boost (1.5%, 1.0%) on (R@1, R@5). For video retrieval, although R@1 of our method is slightly lower than that of [13], the summation of all our three criteria still outperforms it by 0.8%. Moreover, as the most comprehensive criteria, the R@sum of our model obviously surpasses other algorithms, which achieves 5.8% improvement in comparison to the best competitor.

The results on the MSVD dataset are presented in Table 5. It can be seen that our HiSE arrives at 458.9% on the criteria of “R@Sum”, which outperforms the second best method by 15.5%. Especially for text retrieval, the HiSE model surpasses the previous best method by (4.3%, 3.4%, 2.6%) on (R@1, R@5, R@10), respectively. These results substantially demonstrate the advance of our method. Moreover, it can be observed that the text retrieval performance is better than video retrieval, we conjecture the possible reason is that the high-level semantics can provide scarce and complementary information for visual modality. By contrast, it is naturally homogeneous to the textual modality, thus the improvement is not that striking.

4.4 Ablation Studies

In this section, we conduct a series of ablation experiments to explore the impacts of different components in our HiSE model. Note that all results are reported on MSR-VTT 1k-A test set.

4.4.1 Impact of Discrete High-level Semantics. To begin with, we investigate in how the incorporation of discrete high-level semantics affects the performance of HiSE model. From Table 3, we can

see that when TDS representation is employed, our HiSE can obtain 0.2% and 0.2% performance gain on R@1 for video retrieval and text retrieval, respectively. Besides, comparing line #3 with line #6, adopting VDS representation will further bring in 0.9% improvement for video retrieval and 1.3% improvement for text retrieval. These results verify the effectiveness of introducing discrete high-level semantics into video-text representation learning.

4.4.2 Impact of Holistic High-level Semantics. Here, we explore the impact of holistic high-level semantics extraction. From Table 3, comparing line #3 with line #4, we can see the VHS module collaborates well with THS module, which can bring about 0.2% performance boost on R@1 for video retrieval and 0.3% performance boost on R@1 for text retrieval, respectively. Moreover, comparing line #6 with line #8, it can be observed that the combination of all four types of high-level semantic representations will lead to the best performance. These results can verify two main points: 1) The holistic high-level semantics really contributes to improving video-text alignment. 2) The complementarity existing between discrete and holistic high-level semantics also matters for acquiring performance improvements.

4.5 Qualitative Analysis

4.5.1 High-level Semantics Generation Results Visualization. In this part, we display some visualization results of the detected discrete and holistic high-level semantic information. In Figure 4, the second column corresponds to the extracted video discrete semantics with their predicting confidences; the video holistic semantics. *i.e.* video captioning results, are listed in the the third column. For reference, several ground truth descriptions are randomly selected and presented in the last column. For example, as shown in Figure 4, the ground truth (GT) caption of video in first line is “a man walking down street in front of an official building”. As for discrete semantics, the predicting phrases, such as “walking man” and “white building”, can be properly aligned with the GT sentence. Moreover, the generated holistic semantics “a man walking on the street” is also very consistent with the GT one. These results can reflect additional interpretability conveyed by our model.

4.5.2 Bidirectional Cross-Modal Retrieval Results. We compare the bidirectional VTR results obtained by different models. From Figure

	Discrete high-level semantics: small window: 0.91 white building: 0.86 parked motorcycle: 0.86 blue sign: 0.76 shirt: 0.75 walking man: 0.66 white house: 0.58 standing person: 0.52	Holistic high-level semantics: a man is walking on the street	Ground Truth Captions: (1) a man walking down street in front of an office building (2) a large building with people walking around (3) people walk across a white buiding with a blue sign
	Discrete high-level semantics: brown hair: 0.93 brown shirt: 0.93 green grass: 0.89 white shirt: 0.84 sitting man: 0.74 smiling face: 0.66 open mouth: 0.62 young boy: 0.55	Holistic high-level semantics: two men are talking to each other	Ground Truth Captions: (1) 2 young men are laying in the grass and talking (2) guys are talking to each other (3) two boys flirt with each other while sitting in the grass
	Discrete high-level semantics: yellow food: 0.72 silver sink: 0.62 silver bowl: 0.86 brown food: 0.74 black pan: 0.59 frying food: 0.58 yellow oil: 0.55 metal pot: 0.50	Holistic high-level semantics: a person is cooking a dish in a pan	Ground Truth Captions: (1) there is a pan and a dish is frying on it (2) someone is adding oil and frying something in the big pan (3) a person places food in a hot pan and cooks it
	Discrete high-level semantics: brown horse: 0.96 green tree: 0.80 green shirt: 0.76 long tail: 0.74 brown sand: 0.61 running horse: 0.61 dirt road: 0.60 riding person: 0.53 runway: 0.49	Holistic high-level semantics: a woman is riding a horse	Ground Truth Captions: (1) a horse is racing (2) a few horses are riding down a track (3) a reporter discussing a horse race
	Discrete high-level semantics: blue shirt: 0.91 blue man: 0.86 sitting man: 0.85 smiling man: 0.83 wooden chair: 0.78 brown guitar: 0.67 metal basket: 0.55 wooden cabinet: 0.42	Holistic high-level semantics: two men are playing guitar	Ground Truth Captions: (1) men are playing the guitar together (2) two old men play a song on the guitar (3) two men are playing the guitar

Figure 4: Visualization results of video discrete and holistic high-level semantics extracted on MSR-VTT dataset. The listed ground truth sentences are randomly selected from MSR-VTT captions.

Query	HiSE	HiSE (w/o HS)
	<ol style="list-style-type: none"> 1. a computer generated cartoon figure operates a control panel while another character sleeps in the background 2. a squid is talking 3. it is the animation cartoon 	<ol style="list-style-type: none"> 1. a squid is talking 2. a computer generated cartoon figure operates a control panel while another character sleeps in the background 3. it is the animation cartoon
A computer generated cartoon figure operates a control panel while another character sleeps in the background		
	<ol style="list-style-type: none"> 1. band performing a hard rock song about diamonds in the sky 2. this is a rock band music video 3. a video of a rock group performing one of their songs 	<ol style="list-style-type: none"> 1. this is a rock band music video 2. band performing a hard rock song about diamonds in the sky 3. a video of a rock group performing one of their songs
Band performing a hard rock song about diamonds in the sky		
	<ol style="list-style-type: none"> 1. a lady named lizzy is speaking about movies she is wearing a very nice outfit 2. a girl is talking about a celebrity 3. interview with artist shanal twain 	<ol style="list-style-type: none"> 1. a girl is talking about a celebrity 2. a woman introducing someone 3. a girl talks about photos and her life
A lady named lizzy is speaking about movies she is wearing a very nice outfit		
	<ol style="list-style-type: none"> 1. a man is talking to an athlete 2. two people are preparing for sports 3. a football player with a football 	<ol style="list-style-type: none"> 1. two people are preparing for sports 2. a man is talking to an athlete 3. a football player with a football
A man is talking to an athlete		

Figure 5: Quantitative results of V2T and T2V retrieval on MSR-VTT dataset obtained by our model. The HiSE (w/o HS) model indicates the HiSE model without employing High-level Semantics (HS) for VTR. Each video is denoted by its single represented frame. For V2T direction, the ground-truth text are marked as red, while others text are in black. For T2V direction, the ground-truth frames are outlines in red rectangles.

5, it can be seen that the retrieval results listed in the left column is superior than those in the right column. These results further validate the adoption of explicit high-level semantics really contributes to improving the visual-linguistic embeddings, returning more reasonable retrieval results.

5 CONCLUSIONS

The ambiguous understanding of videos and texts impedes the ability of machine to build accurate cross-modal association. In this work, we proposed a explicit high-level semantics (HiSE) aided

visual-linguistic embedding model for improving video-text retrieval. Particularly, we study how to mine explicit high-level semantics from both texts and videos, and incorporate them into cross-modal representations learning. Doing so allows us to disentangle explainable information from raw data that supplies complementary knowledge for the traditional visual-linguistic aligning framework. The experiments conducted on three benchmark datasets validate our method achieves superior performance over the state-of-the-art solutions.

REFERENCES

- [1] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang. 2018. Bottom-up and top-down attention for image captioning and VQA. *CVPR*.
- [2] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. 2017. Localizing moments in video with natural language. In *Proceedings of the IEEE international conference on computer vision*. 5803–5812.
- [3] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, Z. C. Lawrence, and D. Parikh. 2015. VQA: Visual question answering. *ICCV*.
- [4] Yalong Bai, Jianlong Fu, Tiejun Zhao, and Tao Mei. 2018. Deep Attention Neural Tensor Network for Visual Question Answering. In *ECCV*.
- [5] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. 2021. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1728–1738.
- [6] Shuqiang Cao, Bairui Wang, Wei Zhang, and Lin Ma. 2022. Visual Consensus Modeling for Video-Text Retrieval. *AAAI*.
- [7] David Chen and William B Dolan. 2011. Collecting highly parallel data for paraphrase evaluation. In *Proceedings of the 49th annual meeting of the association for computational linguistics: human language technologies*. 190–200.
- [8] Shizhe Chen, Yida Zhao, Qin Jin, and Qi Wu. 2020. Fine-grained video-text retrieval with hierarchical graph reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10638–10647.
- [9] Ioana Croitoru, Simion-Vlad Bogolin, Marius Leordeanu, Hailin Jin, Andrew Zisserman, Samuel Albanie, and Yang Liu. 2021. Teactext: Crossmodal generalized distillation for text-video retrieval. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 11583–11593.
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xi-aohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *ArXiv abs/2010.11929* (2021).
- [11] Maksim Dzabraev, Maksim Kalashnikov, Stepan Komkov, and Aleksandr Petiushko. 2021. Mdmmt: Multidomain multimodal transformer for video retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3354–3363.
- [12] Hao Fang, Saurabh Gupta, Forrest Iandola, Rupesh K Srivastava, Li Deng, Piotr Dollár, Jianfeng Gao, Xiaodong He, Margaret Mitchell, John C Platt, et al. 2015. From captions to visual concepts and back. *CVPR*.
- [13] Han Fang, Pengfei Xiong, Luhui Xu, and Yu Chen. 2021. Clip2video: Mastering video-text retrieval via image clip. *arXiv preprint arXiv:2106.11097* (2021).
- [14] Valentin Gabeur, Chen Sun, Karteek Alahari, and Cordelia Schmid. 2020. Multimodal transformer for video retrieval. In *European Conference on Computer Vision*. Springer, 214–229.
- [15] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. , 9729–9738 pages.
- [16] Y. Huang, Q. Wu, C. Song, and L. Wang. 2018. Learning semantic concepts and order for image and sentence matching. *CVPR*.
- [17] Drew Hudson and Christopher D Manning. 2019. Learning by abstraction: The neural state machine. *Advances in Neural Information Processing Systems* 32 (2019).
- [18] Zhong Ji, Kexin Chen, and Haoran Wang. 2021. Step-Wise Hierarchical Alignment Network for Image-Text Matching. *IJCAI*.
- [19] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *ICLR*.
- [20] K.-H. Lee, X. Chen, G. Hua, H. Hu, and X. He. 2018. Stacked Cross Attention for Image-Text Matching. *ECCV*.
- [21] Xirong Li, Chaoxi Xu, Gang Yang, Zhineng Chen, and Jianfeng Dong. 2019. W2vv++ fully deep learning for ad-hoc search. In *Proceedings of the 27th ACM International Conference on Multimedia*. 1786–1794.
- [22] Xirong Li, Fangming Zhou, Chaoxi Xu, Jiaqi Ji, and Gang Yang. 2020. Sea: Sentence encoder assembly for video retrieval by textual queries. *IEEE Transactions on Multimedia* 23 (2020), 4351–4362.
- [23] Yehao Li, Yingwei Pan, Jingwen Chen, Ting Yao, and Tao Mei. 2021. X-modaler: A Versatile and High-performance Codebase for Cross-modal Analytics.
- [24] Yuanzhi Liang, Yalong Bai, Wei Zhang, Xueming Qian, Li Zhu, and Tao Mei. 2019. VrR-VG: Refocusing Visually-Relevant Relationships. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)* (2019), 10402–10411.
- [25] Fangyu Liu, Rongtian Ye, Xun Wang, and Shuaipeng Li. 2020. Hal: Improved text-image matching by mitigating visual semantic hubs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 11563–11571.
- [26] Song Liu, Haoqi Fan, Shengsheng Qian, Yiru Chen, Wenkui Ding, and Zhongyuan Wang. 2021. Hit: Hierarchical transformer with momentum contrast for video-text retrieval. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 11915–11925.
- [27] Yang Liu, Samuel Albanie, Arsha Nagrani, and Andrew Zisserman. 2019. Use what you have: Video retrieval using representations from collaborative experts. In *arXiv preprint arXiv:1907.13487*.
- [28] Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. 2021. Clip4clip: An empirical study of clip for end to end video clip retrieval. *arXiv preprint arXiv:2104.08860* (2021).
- [29] Lin Ma, Zhengdong Lu, and Hang Li. 2016. Learning to Answer Questions From Image Using Convolutional Neural Network. *AAAI*.
- [30] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9 (2008), 2579–2605.
- [31] Antoine Miech, Ivan Laptev, and Josef Sivic. 2018. Learning a text-video embedding from incomplete and heterogeneous data. *arXiv preprint arXiv:1804.02516* (2018).
- [32] Antoine Miech, Ivan Laptev, and Josef Sivic. 2018. Learning a text-video embedding from incomplete and heterogeneous data. *arXiv preprint arXiv:1804.02516* (2018).
- [33] Niluthpol Chowdhury Mithun, Juncheng Li, Florian Metzke, and Amit K Roy-Chowdhury. 2018. Learning joint embedding with multimodal cues for cross-modal video-text retrieval. In *Proceedings of the 2018 ACM on International Conference on Multimedia Retrieval*. 19–27.
- [34] Mandela Patrick, Po-Yao Huang, Yuki Asano, Florian Metzke, Alexander Hauptmann, Joao Henriques, and Andrea Vedaldi. 2020. Support-set bottlenecks for video-text representation learning. *arXiv preprint arXiv:2010.02824* (2020).
- [35] Jesús Andrés Portillo-Quintero, José Carlos Ortiz-Bayliss, and Hugo Terashima-Marin. 2021. A straightforward framework for video retrieval using clip. In *Mexican Conference on Pattern Recognition*. Springer, 3–12.
- [36] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *ICML*.
- [37] Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. 2018. Modeling relational data with graph convolutional networks. *ESWC*.
- [38] Peng Shi and Jimmy J. Lin. 2019. Simple BERT Models for Relation Extraction and Semantic Role Labeling. *ArXiv abs/1904.05255* (2019).
- [39] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, K. Lukasz, and I. Polosukhin. 2017. Attention is all you need. *NIPS*.
- [40] Bairui Wang, Lin Ma, Wei Zhang, and Wei Liu. 2018. Reconstruction Network for Video Captioning. *CVPR*.
- [41] Bokun Wang, Yang Yang, Xing Xu, Alan Hanjalic, and Heng Tao Shen. 2017. Adversarial Cross-Modal Retrieval. *ACM international conference on Multimedia*.
- [42] Haoran Wang, Ying Zhang, Zhong Ji, Yanwei Pang, and Lin Ma. 2020. Consensus-Aware Visual-Semantic Embedding for Image-Text Matching. *ECCV*.
- [43] Jingwen Wang, Wenhao Jiang, Lin Ma, Wei Liu, and Yong Xu. 2018. Bidirectional Attentive Fusion with Context Gating for Dense Video Captioning. *CVPR*.
- [44] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. 2019. Vatec: A large-scale, high-quality multilingual dataset for video-and-language research. In *ICCV*. 4581–4591.
- [45] Peng Wu, Xiangteng He, Mingqian Tang, Yiliang Lv, and Jing Liu. 2021. HANet: Hierarchical Alignment Networks for Video-Text Retrieval. In *Proceedings of the 29th ACM International Conference on Multimedia*. 3518–3527.
- [46] Qi Wu, Chunhua Shen, Lingqiao Liu, Anthony Dick, and Anton Van Den Hengel. 2016. What value do explicit high level concepts have in vision to language problems?. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 203–212.
- [47] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5288–5296.
- [48] Xing Xu, Fumin Shen, Yang Yang, D. Zhang, Heng Tao Shen, and Jingkuan Song. 2017. Matrix Tri-Factorization with Manifold Regularizations for Zero-Shot Learning. , 2007–2016 pages. *Computer Vision and Pattern Recognition (CVPR)*.
- [49] Ting Yao, Yingwei Pan, Yehao Li, Zhaofan Qiu, and Tao Mei. 2017. Boosting image captioning with attributes. In *Proceedings of the IEEE international conference on computer vision*. 4894–4902.
- [50] Youngjae Yu, Jongseok Kim, and Gunhee Kim. 2018. A joint sequence fusion model for video question answering and retrieval. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 471–487.

- [51] Alireza Zareian, S. Karaman, and Shih-Fu Chang. 2020. Bridging Knowledge Graphs to Generate Scene Graphs. *ECCV*.
- [52] Michael Zhang, James Lucas, Jimmy Ba, and Geoffrey E Hinton. 2019. Lookahead optimizer: k steps forward, 1 step back. *Neurips* 32.
- [53] Rui Zhao, Kecheng Zheng, Zheng-Jun Zha, Hongtao Xie, and Jiebo Luo. 2021. Memory Enhanced Embedding Learning for Cross-Modal Video-Text Retrieval. *arXiv preprint arXiv:2103.15686* (2021).
- [54] Shuai Zhao, Linchao Zhu, Xiaohan Wang, and Yi Yang. 2022. CenterCLIP: Token Clustering for Efficient Text-Video Retrieval. *SIGIR*.
- [55] Luowei Zhou, Yingbo Zhou, Jason J. Corso, Richard Socher, and Caiming Xiong. 2018. End-to-End Dense Video Captioning with Masked Transformer. , 8739-8748 pages.

A OVERVIEW OF APPENDIX

In this appendix, we give some details which were omitted in the main body of our manuscript owing to the limited space. Concerning the experiments, we report the experimental results on DiDeMo dataset, impact of representation fusion hyper-parameter α , impact of different alignment objectives, and data distribution visualization result in joint embedding space.

B EXPERIMENTS

B.1 Experimental Results on DiDeMo dataset

The experimental results on DiDeMo dataset are presented in Table 4. It can be seen that, in comparison to the method with the same raw video encoder as ours, *i.e.* CLIP4Clip-seqTransf, the overall retrieval quality reflected by the R@Sum metric is increased by a large margin (+ 8.5 %). Moreover, compared to the best competitor CLIP4Clip-meanP, our HiSE obtains 0.7 and 1.3 % performance gains on R@1 metric for video retrieval and text retrieval, respectively. We believe the major improvement derives from the exploitation of the introducing of explicit high-level semantics, which supplies more complementary information to narrow the modality gap.

B.2 Ablation Studies

Consistent with the main body of our paper, all our ablation experiments are conducted on MSR-VTT 1k-A test set.

B.2.1 Impact of Graph Reasoning in High-level Semantics Representation. In this part, we explore the affect of graph reasoning modules in high-level semantics representation components. As shown in Table 5, for visual semantics representation, when we replace mean-pooling with visual graph reasoning module, our model can obtain 0.1% and 0.4% performance gain on R@1 for video retrieval and text retrieval, respectively. As for textual branch, compared to adopting mean-pooling, employing textual graph reasoning to aggregate hierarchical information can result in 0.3% boost for video retrieval and 0.1% boost for text retrieval. Furthermore, the deployment of both visual and textual graph reasoning can lead to the best performance. These results is consistent with our designing aim of these modules, which leverage the graph reasoning technique to promote the interaction between hierarchies of high-level semantics.

B.2.2 Impact of Representation Fusion Parameter. To explore the impact of parameter α of fusing raw video-text representation $V(T)$ and high-level semantic representation $V^{VS}(V^{TS})$. In Figure 6, we can see that the bidirectional retrieval performances both decrease when α varies from 0.9 to 1. Considering $\alpha=1$ indicates the high-level semantics is removed from the video-text representation, these results validate the effectiveness of introducing high-level semantics to improve cross-modal discrimination.

B.2.3 Impact of Different Alignment Objectives. In this part, we analyze the impact of different alignment objectives in our HiSE method. Specifically, in Table 6, we present the retrieval result of HiSE replacing the M-HAL loss with the prevailing Bi-directional InfoNCE (B-InfoNCE) loss [28, 36] in recent comparison studies. From Table 6, it can be seen that although using the same objective, our model still outperforms the second best competitor by 0.6% improvement for video retrieval and 1.4% improvement for text

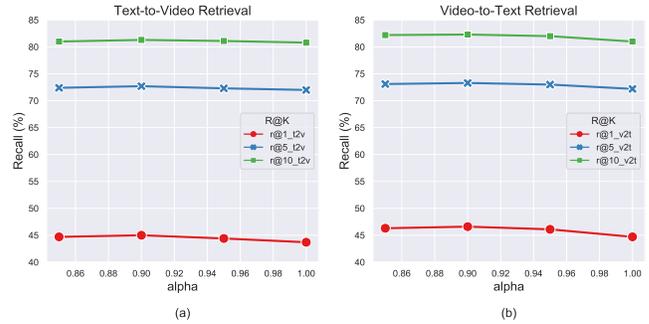


Figure 6: Impact of varied controlling parameters α on MSR-VTT 1k-A test set. Sub-figure (a) depicts how parameter α affects video retrieval performance, and Sub-figure (b) illustrates the corresponding text retrieval performance.

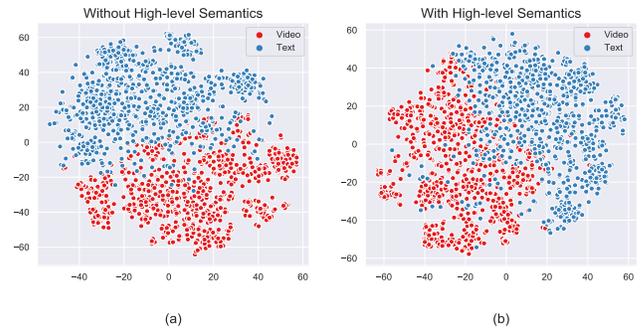


Figure 7: T-SNE visualization of the video-text representations generated by (a) baseline model with L_{HAL_L} loss and (b) full CODER model on MSR-VTT 1k-A test set (1000 videos and 1000 texts).

retrieval, respectively. It further validates the advance of our proposed explicit high-level semantics encoding modules. Besides, the two items in M-HAL loss can both bring about respective performance gain. These results confirm the effectiveness and rationality of our employed M-HAL loss for video-text retrieval.

B.3 Qualitative Analysis

B.3.1 T-SNE Visualization of Video-Text Representation. To further investigate how the explicit high-level semantics affects the learned joint embedding space, we adopt t-SNE [30] to visualize the learned cross-modal representations from MSR-VTT 1k-A test set, containing 1000 images and 1000 texts. In particular, the data distribution of the baseline model without utilizing explicit semantics and our full HiSE model are depicted in Figure 7(a) and Figure 7(b). In comparison to the former, it can be observed the distributions of videos and texts are further mixed by our proposed HiSE method. These results indicate that the proposed explicit semantics incorporating method contributes to reducing the distribution difference between two modalities.

Table 4: Comparisons of Experimental Results on DiDeMo Testing Set.

Approach	Text-to-Video Retrieval				Video-to-Text Retrieval				R@Sum
	R@1	R@5	R@10	MdR	R@1	R@5	R@10	MdR	
Non-CLIP Based									
S2VT[44]	11.9	33.6	-	13.0	13.2	33.6	-	15.0	-
FSE[52]	13.9	36.0	-	11.0	13.1	33.9	-	12.0	-
CE[27]	16.1	41.1	-	8.3	15.6	40.9	-	8.2	-
TT-CE+[9]	21.6	48.6	62.9	6.0	21.1	47.3	61.1	6.3	262.6
MoEE[32]	16.1	41.2	55.2	8.3	16.0	41.7	54.6	8.7	224.8
Frozen[5]	34.6	65.0	74.7	3.0	-	-	-	-	-
CLIP Based									
CLIP4Clip-seqTransf[28]	42.8	68.5	79.2	2.0	41.4	68.2	79.1	2.0	379.2
CLIP4Clip-meanP[28]	43.4	70.2	80.6	2.0	42.5	70.6	80.2	2.0	387.5
HiSE	44.1	69.9	80.3	2.0	43.8	70.4	79.2	2.0	387.7

Table 5: Performance of our HiSE method with different high-level semantics aggregating modules on MSR-VTT 1k-A test set. Visual Graph Reasoning and Textual Graph Reasoning are abbreviated as “VGR” and “TGR”, respectively. Mean-pooling operation for visual and textual modality are abbreviated as “VMP” and “TMP”, respectively.

High-level Semantics Aggregating Module				Text-to-Video Retrieval			Video-to-Text Retrieval		
VGR	TGR	VMP	TMP	R@1	R@5	R@10	R@1	R@5	R@10
-	-	✓	✓	44.5	72.2	81.0	46.0	73.0	82.1
✓	-	-	✓	44.6	72.3	81.1	46.4	73.2	82.2
-	✓	✓	-	44.8	72.5	81.3	46.1	73.1	82.1
✓	✓	-	-	45.0	72.7	81.3	46.6	73.3	82.3

Table 6: Performance of our HiSE with different alignment objectives on MSR-VTT 1k-A test set. Coupled Memory Banks are abbreviated as “CMB”.

Approach	Alignment Objective	Components		Text-to-Video Retrieval			Video-to-Text Retrieval		
		CMB	HAL	R@1	R@5	R@10	R@1	R@5	R@10
CLIP[35]	B-InfoNCE	-	-	31.2	53.7	64.2	27.2	51.7	62.6
CLIP4Clip-meanP[28]	B-InfoNCE	-	-	43.1	70.4	80.8	43.1	70.5	81.2
CLIP4Clip-seqTransf[28]	B-InfoNCE	-	-	44.5	71.4	81.6	42.7	70.9	80.6
HiSE	B-InfoNCE	-	-	45.1	72.1	80.9	44.1	72.8	81.5
HiSE	M-HAL	-	✓	44.6	72.5	81.1	46.3	73.0	82.2
HiSE	M-HAL	✓	-	45.0	72.7	81.3	46.6	73.3	82.3