

# RTQ: Rethinking Video-language Understanding Based on Image-text Model

Xiao Wang  
scz.wangxiao@gmail.com  
Harbin Institute of Technology,  
Shenzhen & JD.com Inc.

Yaoyu Li  
liyaoyu2014@gmail.com  
JD.com Inc.

Tian Gan\*  
gantian@sdu.edu.cn  
Shandong University

Zheng Zhang  
zhangzheng11@jd.com  
JD.com Inc.

Jingjing Lv  
lvjinghit@163.com  
JD.com Inc.

Liqiang Nie  
nieliqiang@gmail.com  
Harbin Institute of Technology,  
Shenzhen

## ABSTRACT

Recent advancements in video-language understanding have been established on the foundation of image-text models, resulting in promising outcomes due to the shared knowledge between images and videos. However, video-language understanding presents unique challenges due to the inclusion of highly complex semantic details, which result in *information redundancy*, *temporal dependency*, and *scene complexity*. Current techniques have only partially tackled these issues, and our quantitative analysis indicates that some of these methods are complementary. In light of this, we propose a novel framework called RTQ (Refine, Temporal model, and Query), which addresses these challenges simultaneously. The approach involves refining redundant information within frames, modeling temporal relations among frames, and querying task-specific information from the videos. Remarkably, our model demonstrates outstanding performance even in the absence of video-language pre-training, and the results are comparable with or superior to those achieved by state-of-the-art pre-training methods.

## CCS CONCEPTS

• Information systems → Multimedia and multimodal retrieval; • Computing methodologies → Computer vision.

## KEYWORDS

Video Retrieval; Video Caption; Video Question Answering

### ACM Reference Format:

Xiao Wang, Yaoyu Li, Tian Gan, Zheng Zhang, Jingjing Lv, and Liqiang Nie. 2023. RTQ: Rethinking Video-language Understanding Based on Image-text Model. In *Proceedings of the 30th ACM International Conference on Multimedia (MM '23)*, October 29–November 3, 2023, Ottawa, Canada. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/nmmnnnnn.nmmnnnnn>

\*Corresponding author: Tian Gan

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](mailto:permissions@acm.org).

MM'23, October 29–November 3, 2023, Ottawa, Canada

© 2023 Association for Computing Machinery.  
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00  
<https://doi.org/10.1145/nmmnnnnn.nmmnnnnn>

(a) Challenge 1: Information Redundancy



Caption: A person is pulling on a parachute.

(b) Challenge 2: Temporal Dependency



Question: What did the dog do after the girl puts her palm out?  
Candidates: A1: place paw on her palm, A2: kiss her.

(c) Challenge 3: Scene Complexity



Text query: The woman first spreads butter on one side of the bread slices, then places the bread butter onto a hot pan.

**Figure 1: Video language understanding requires dynamic perception and interpretation of complex semantics, which can be further decomposed into three challenges.**

## 1 INTRODUCTION

Video-language understanding ability can reflect the proficiency of intelligent agents in perceiving and interpreting visual and textual cues in the real world. This ability is evaluated through a range of tasks, including text-to-video retrieval, video captioning, and video question answering, etc. Recent approaches in this area generally modify pre-trained image-text models for video-language understanding, owing to the transferable knowledge acquired by these models [12, 17, 46]. Since image-text models have learned a lot of transferable vision-language knowledge, these approaches generally perform better and are becoming the *de facto* paradigm. However, such an approach has limitations in handling situations beyond the shared knowledge between images and videos.

Video-language understanding involves dynamic perception and interpretation of complex semantics. This can be further decomposed into three challenges as depicted in Fig. 1. The first challenge is **Information redundancy** which arises due to the presence of information that is duplicated or lacks semantic meaning. It hinders the model's ability to accurately recognize essential cues. For example, frames in Fig. 1(a) are quite similar, some of which can be removed without affecting the interpretation. The second

challenge is **Temporal dependency** which requires the model to identify and understand the relationships between video frames. Consider video in Fig. 1(b), multiple actions occur between the girl and the dog. To correctly answer the given question, the model has to recognize the order of these actions. The third challenge is **Scene complexity**, where a video depicts multiple concepts, but only some of them are task-relevant. In such cases, the model needs to prioritize task-relevant information to achieve better performance. If a model fails to consider this, it may become overwhelmed by the complexity of the scene, leading to poor performance. For instance, in video of Fig. 1(c), numerous cooking actions and ingredients are presented. To achieve better performance, the model must concentrate on crucial pieces of information such as objects and actions mentioned in text queries, captions, or questions.

Current approaches have mainly focused on addressing one or two of the aforementioned challenges. However, it is essential to consider these challenges jointly since they address different facets that could complement each other. For example, eliminating redundant information could benefit the model when solving the scene complexity problem [8]. Nonetheless, it is crucial to carry out quantitative analysis to ascertain the extent of the complementarity between these challenges. Moreover, joint modeling is not a straightforward task due to the challenge of cooperative design [11, 35]. Existing methods [25, 52] have addressed the information redundancy challenge by selecting meaningful tokens or frames. However, this selection process breaks the spatial consistency between frames, making it intractable for temporal modeling techniques. In this regard, TS2-Net [25] set the selection module as the final layer to enable temporal modeling. Nevertheless, this approach is sub-optimal since a lot of redundant information could be reduced in the shallower layers. On the other hand, some methods choose to ignore the information redundancy challenge, and only address the temporal dependency and scene complexity problems [34, 44].

This paper begins with a quantitative analysis of existing approaches, wherein we cluster them based on their predictions. Our analysis reveals that the models within each cluster address the same challenges and confirms their complementarity, as previously discussed. Based on this observation, we propose the **Refine**, **Temporal** model, and **Query (RTQ)** framework to jointly tackle the aforementioned challenges. The RTQ framework consists of three key components, each of which is designed to address a specific challenge. The first refinement component eliminates redundant patches across adjacent video frames using clustering, followed by the selection of representative patches. The second temporal modeling component uses an image backbone augmented with a temporal module. The module is designed to perceive and interpret temporal patterns in the video, without requiring spatial consistency between patches across frames. In the third query component, we adopt language query (text query, question, or generated caption) to accumulate task-relevant information gradually. The aforementioned three components can be realized through any appropriate methods. In this study, we select the simplest modules available. Specifically, we employ the non-parametric *k-medoids++* method for clustering, the message token mechanism for the temporal module, and cross-attention for the query. Despite their simplicity and the absence of pretraining, our approach achieves superior (or comparable)

performance to the pre-training based methods in text-to-video retrieval, video captioning, and video question answering.

In summary, our contributions are threefold:

- Our systemic analysis reveals that current methods focus only on restricted aspects of video-language understanding, and they are complementary.
- We propose the RTQ framework to jointly model information redundancy, temporal dependency, and scene complexity in video-language understanding.
- We demonstrate that, even without pre-training on video-language data, our method can achieve superior (or comparable) performance with state-of-the-art pre-training methods. We will make our code publicly available for further research<sup>1</sup>.

## 2 RELATED WORKS

### 2.1 Video-language Pre-training

The dominant pre-training methods for video-language understanding fall into two categories. The first category focus on data curation and refinement. For example, Zellers *et al.* [51] collect a diverse corpus of frames/ASR, named YT-Temporal-180M, from videos covering authentic situations, which improves downstream performance compared to curated instructional video corpora. However, video-language pre-training also suffers from incoherence and misalignment between ASR/subtitle and video [12]. To overcome this issue, Bain *et al.* [3] curate a video dataset, WebVid, with well-aligned textual description annotations. CLIP-ViP [46] uses an image captioning model to generate captions for the middle frame of videos in HD-VILA-100M to obtain more aligned video-text annotations. The texts are not so well aligned compared with videos in WebVid, but the scale of data is much bigger.

The second category focus on improving pre-training strategy. To bridge the modality gap between video and text, BridgeFormer [13] proposed a novel multiple-choice questions task to achieve fine-grained video-text interactions. OmniVL [34] and mPLUG-2 [44] explore a universal paradigm that benefit from joint modality learning. OmniVL adopts a unified transformer-based visual encoder for both image and video inputs, facilitating joint image-language and video-language pre-training. On the other hand, mPLUG-2 introduces a multi-module composition network that shares common universal modules for modality collaboration and disentangles different modality modules to handle modality entanglement.

Despite their contributions, all of the above methods mainly focus on pre-training data or strategy, neglecting architecture design for modeling information redundancy, temporal dependency, and scene complexity in video-language understanding, which prevents the full realization of the model's potential.

### 2.2 Video-language Model Architecture

Classical architectures utilize separately pre-trained vision and language backbones, which remain static during training [37]. However, recent studies have identified limitations in these approaches related to modality and domain gaps [17], whereas newer architectures based on image-text pre-training models show more promising results due to their ability to bridge the modality gap and

<sup>1</sup>See our GitHub repository <https://github.com/SCZwangxiao/RTQ-MM2023>.

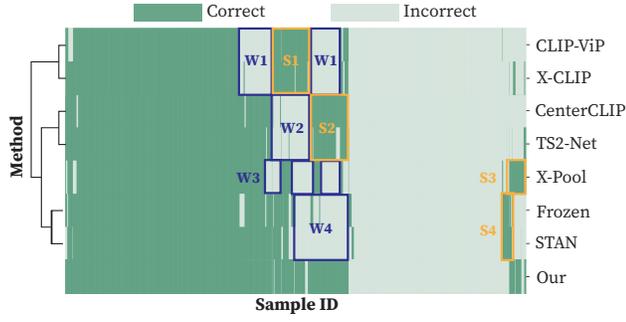


Figure 2: Clustering results of existing methods in the NeXt-QA [40] dataset.

enhance transferability [31]. As videos are essentially composed of image sequences, the insights from image-text models can also be applied to video-language understanding.

Recent architectures can be categorized into three groups: refinement, temporal modeling, and information query. **Refinement methods** aim to identify and eliminate irrelevant image patches or frames. In the case of TS2-Net [25], a scoring network is initially employed to assess the significance of each patch in the frames, followed by the application of a differentiable TopK function for patch selection. Zhao *et al.* [52] proposed to first cluster all patches in adjacent frames and then retain only the patches that are closest to each cluster centroid. **Temporal modeling methods** focus on modeling the temporal dependencies between video frames. Frozen [3] and STAN [24] insert a temporal attention layer in their image model. Their difference is that STAN [24] inserts the layer in parallel instead of sequential with the original model, which empirically performs better. X-CLIP [28] and CLIP-ViP [46] adopt the message tokens mechanism, where each frame has message token(s) to communicate with other frames. They differ in the construction approach and perception range of message tokens. CLIP4clip [27] employs a transformer to model temporal dependencies, while TS2-Net shifts spatial token features across adjacent frames to capture local movement. **Information query methods** are capable of mining task-relevant information from the whole video. Gorti *et al.* [14] proposed to use text as the query to guide the aggregation of useful information among the entire video. VideoCoCa [47] leverages a generative pooling mechanism that gradually accumulates relevant information for caption generation.

However, these methods only address one or two aspects of video understanding and fail to consider the complementarity of different architectures.

### 3 PRELIMINARY ANALYSIS

In Section 2.2, we highlighted that existing video-language models do not simultaneously address all three challenges and that we believe a joint approach is beneficial as the challenges may be complementary. To gain more objective insights, we conducted a quantitative analysis by clustering the models based on their successful and failed cases. Our rationale was to examine their complementarity by comparing representative cases within each cluster. In the following sections, we explain our clustering methodology and present our findings.

### 3.1 Clustering Method

To conduct a clustering analysis, it is necessary to establish a vector representation of a model to serve as the basis for measuring the distance between models [53]. Subsequently, the appropriate clustering methods are selected for analysis. In this study, we define a method’s representation as  $\mathbf{m} \in \mathbb{R}^N$ , where  $N$  is the number of samples in the validation set, and  $m_i \in \{0, 1\}$  indicates whether the  $i$ -th sample is correctly predicted. Hamming distance  $d(\mathbf{m}, \mathbf{n})$  is used to assess the similarity of two methods  $\mathbf{m}$  and  $\mathbf{n}$ , given by:

$$d(\mathbf{m}, \mathbf{n}) = \sum_{i=1}^N \mathbb{I}(m_i = n_i), \quad (1)$$

where  $\mathbb{I}(m_i = n_i) = 1$  only if  $m_i = n_i$ . Finally, we employ hierarchical clustering to explore the relationships among each method.

Our analysis was conducted on the validation split of the NeXt-QA dataset [40] to perform our analysis because it contains a considerable number of descriptive, temporal, and causal questions that thoroughly evaluate a model’s video understanding capabilities [4].

### 3.2 Result Analysis

The clustering results are depicted in Fig. 2, with each row representing a distinct method and each column a unique sample ID (approximately 5k IDs). The results revealed four clusters of methods, namely: (1) temporal modeling methods that leverage message tokens mechanism, such as CLIP-ViP [46], X-CLIP [28]; (2) refinement methods, including CenterClip [52], TS2-Net [25]; (3) information query approach XPool [14]; and (4) temporal modeling methods based on temporal attention, such as Frozen [3] and STAN [24]. The four clusters broadly align with the three categories of methods discussed in Section 2, namely refinement, temporal modeling, and information query.

Upon closer examination of the results, we find that each cluster of methods has distinct advantages in handling certain types of samples that other methods may struggle with (as indicated by the Area S1-S4 in Fig. 2). Meanwhile, these methods also exhibit weaknesses in handling particular samples, but which can be overcome effectively by other methods (as indicated by the Area W1-W4). Our findings indicate that current approaches tend to focus on specific aspects of video-language understanding, but exhibit some degree of complementarity. Therefore, it is possible and beneficial to jointly consider all three challenges to harness the strengths of these methods while mitigating their respective weaknesses. Based on these insights, we propose the RTQ framework, which we elaborate in the next section.

## 4 METHOD

### 4.1 Overview

Based on observations in Section 3, we propose addressing the three challenges of video-language understanding collaboratively in the RTQ framework, outlined in Fig. 3 (a). The framework leverages a video encoder with refinement and temporal modules, and a query component composed of a Mixture of Encoder-Decoder (MoED).

Specifically, the video encoder first encodes individual video frames utilizing  $K$  Vision Transformer (ViT) [9] layers, thereby generating image patch embeddings with semantic meanings. Then,

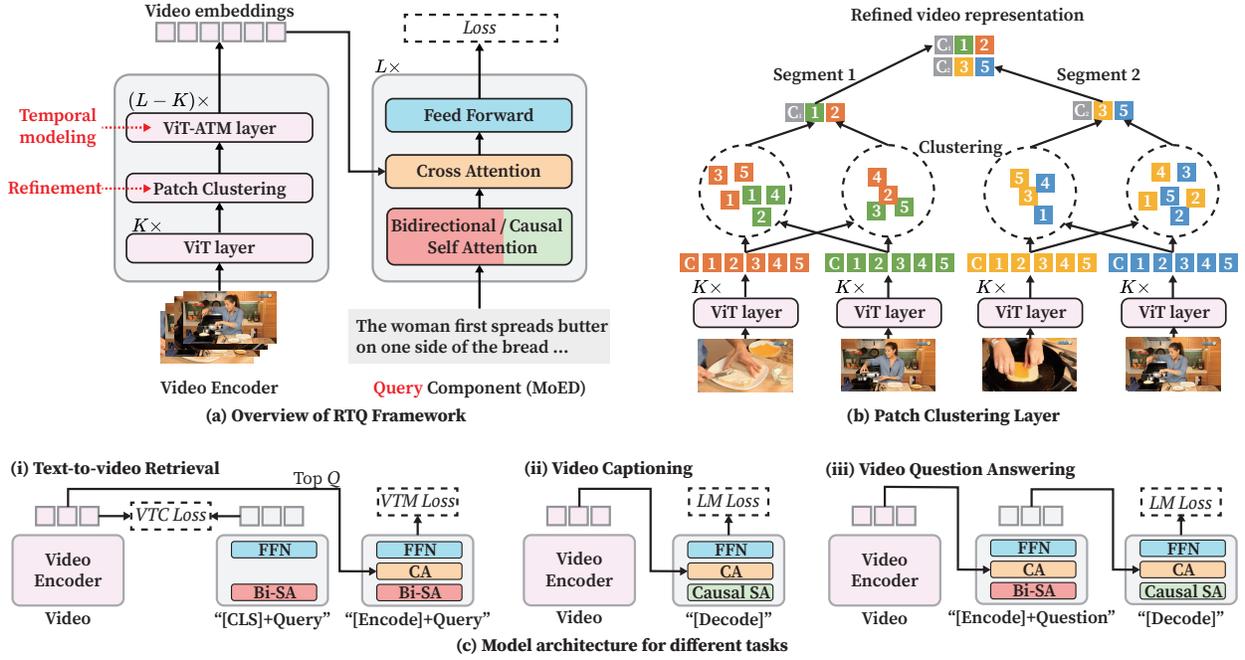


Figure 3: Method overview.

a patch clustering layer is employed to eliminate redundant patches by retaining only representative patches from the clustering results. Subsequently, the remaining patches are fed into  $(L - K)$  ViT Augmented with Temporal Module (ViT-ATM) layers, where  $L$  is the total number of video encoder layers. These layers enable the model to capture temporal dependencies between frames and produce video embeddings as the output. Finally, in the query component of the framework, task-specific text queries in each layer are employed to progressively collect task-relevant information from the video embeddings. This process varies depending on tasks, which we elaborate in Section 4.4.

## 4.2 Refinement Component

As illustrated in Fig. 3 (b), patch clustering layer is employed to refine the video embeddings. This process entails grouping patches in adjacent frames and selecting the representative ones to generate refined video embeddings.

The output frame embeddings of the ViT layers are denoted as  $\hat{\mathbf{V}}^K \in \mathbb{R}^{F \times (1+P_{in}) \times d}$ , where  $F$  is the number of input frames, “1” represents the [CLS] token,  $P_{in}$  is the number of patches per image divided by ViT, and  $d$  is the hidden dimension. Initially, frames are grouped into  $S$  segments, with each containing  $F/S$  frames. Next, the  $k$ -medoids++ method [52] is applied to cluster all  $F/S(1+P_{in})$  patches within each segment, resulting in  $(1+P_{out})$  clusters. Finally,  $(1+P_{out})$  patches that are closest to each cluster centroid are selected to form the refined video embeddings  $\mathbf{V}^K \in \mathbb{R}^{S \times (1+P_{out}) \times d}$ .

It is worth noting that the clustering process is not limited to the  $k$ -medoids++ method, and may be implemented with other clustering methods. The straightforward and non-parametric  $k$ -medoids++ method is employed in this study to demonstrate the effectiveness of our proposed framework.

## 4.3 Temporal Modeling Component

The temporal modeling is accomplished by incorporating temporal modules into the ViT layers. The temporal module generally consists of a modified or inserted layer in the original ViT layer. Temporal modules come in three varieties: temporal attention [3], message tokens [46], and temporal shifting [25]. We adopt the message tokens mechanism, the details and design rationales are explained below.

Both temporal attention and temporal shifting mechanisms necessitate spatial consistency among patch positions of input frames. This is because they perform temporal reasoning on patches that are located at the same positions across different frames. However, clustering disrupts spatial consistency in our framework, as illustrated in Fig. 3 (b). Segment 1 has patches from location ID “1, 2”, while segment 2 has patches from location ID “3, 5”. Furthermore, temporal reasoning is built upon the underlying assumption that content within each patch changes gradually, enabling the identification of temporal semantics by observing patches across different frames. This assumption is not applicable in the video-language understanding task because it predominantly involves processing lengthy, untrimmed videos. Our ablation studies in Section 5.4 have confirmed this view.

Different from these two mechanisms, the message token mechanism generally inserts several learnable embeddings along with patch embeddings, and utilize these learnable embeddings as an intermediary for temporal reasoning between frames. In our framework, we utilize a straightforward message token mechanism similar to that in X-CLIP [28]. Specifically, in the  $k$ -th layer, we first select and concatenate the [CLS] token in the refined video embedding to produce message token embeddings  $\hat{\mathbf{m}}^k \in \mathbb{R}^{S \times d}$ . We then perform Self-Attention (SA) along these tokens to learn the

temporal dependencies between frames:

$$\mathbf{m}^k = \hat{\mathbf{m}}^k + \text{SA}(\hat{\mathbf{m}}^k). \quad (2)$$

Finally,  $\mathbf{m}^k$  and other patch tokens are fed into the original ViT layer. After  $(L - K)$  ViT-ATM layer, our method produces the video embeddings video embeddings  $\mathbf{V}^L \in \mathbb{R}^{S \times (1+P_{\text{out}}) \times d}$ .

#### 4.4 Query Component

After the previous refinement and temporal modeling components, videos have been encoded into a temporal-aware representation with high information density. Nevertheless, a significant amount of task-irrelevant information remains due to the challenge of scene complexity. For instance, in retrieval and question answering tasks, the model should focus on information about the query or question. Similarly, in captioning, the model should emphasize previously undescribed objects and events in the current generated caption.

To address this issue, we introduce a query component that regards video embeddings as memories and use task-specific queries to gradually gather relevant details for the final results. To reuse common structures among various tasks, we integrate an  $L$  layer MoED [18] as the fundamental module of our query component. As highlighted with different colors in Fig. 3 (a), MoED consists of four modules, including Bi-directional Self Attention (BiSA), Causal Self Attention (Causal SA), Cross Attention (CA), and Feed Forward Network (FFN). They can be assembled into three variants for various tasks:

- **Text encoder** is the same as BERT [7], which encodes text by using a BiSA and an FFN in each layer. A [CLS] token is appended to the beginning of the text input to summarize it.
- **Video-grounded text encoder** gathers task-relevant visual information by incorporating one additional cross attention between the BiSA and the FFN in each layer of the text encoder. Thereinto, the text input (query or question) serves as the query and the flattened video embeddings serve as key and value. To accommodate the task requirements, a task-specific [Encode] token is appended to the text, and the resulting embedding of [Encode] contains the multimodal representation of the video-text pair.
- **Video-grounded text decoder** is responsible for collecting task-specific visual information to generate the desired text output. It replaces the BiSA layers of the video-grounded text encoder with Causal SA layers. The start of a sequence is identified using a [Decode] token, whereas the conclusion of the sequence is identified using an end-of-sequence token.

The subsequent discourse elaborates on how these three variations exhibit distinct approaches towards video-language understanding.

**Text-to-video retrieval.** As summarized in Fig. 3 (c)(i), we perform text-to-video retrieval in two stages: recall and re-rank. They are completed by the text encoder and video-grounded text encoder, respectively. During inference, we first recall Top  $Q$  videos by calculating the cosine similarity of the [CLS] token of video and text embeddings. Then, we re-rank the recalled videos by feeding the video-text query pair into the video-grounded text encoder, and use the output embedding of [Encode] with a fully-connected layer and sigmoid function to get the matching scores. The matching score and cosine similarity are added to get the final scores. The

self-attention and feedforward share parameters between the text encoder and video-grounded text encoder.

**Video captioning.** As summarized in Fig. 3 (c)(ii), our methods apply the video-grounded text decoder to generate captions based on video embeddings.

**Video question answering.** As depicted in Fig. 3 (c)(iii), for *open-ended QA*, our method first encodes video and question text into multimodal embeddings using the video-grounded text encoder. Then, we feed the multimodal embeddings into the video-grounded text decoder to generate answers. The encoder and decoder share parameters. For *multiple-choice QA*, we formulate it as a classification problem. Specifically, we concatenate the question and answer into a whole sentence, then we apply the video-grounded text encoder to encode the video and question-answer pair into multimodal embedding. We finally apply Softmax after a linear layer to get the score of the best answers.

#### 4.5 Training Objectives

We detail the loss functions and training strategies.

**Text-to-video retrieval.** We jointly train the text encoder and video-grounded text encoder. For the text encoder, we apply Video-Text Contrastive (VTC) loss, which aligns the video and text feature space by encouraging the [CLS] tokens of matched video-text pairs to have similar representations against the unmatched pairs. Formally, for the  $i$ -th video-text pair, given their [CLS] embeddings, we follow CLIP [31] to apply a linear projection and  $L_2$  normalization layer on them to obtain the hidden video vector  $\mathbf{v}_i \in \mathbb{R}^d$  and text vector  $\mathbf{t}_i \in \mathbb{R}^d$ . First of all, to maximize the benefits of contrastive learning with a larger batch size (which makes it more accurate theoretically) [38], we maintain three memory banks to store most recent  $M$  video vectors  $\{\mathbf{v}_m\}_{m=1}^M$  and text vectors  $\{\mathbf{t}_m\}_{m=1}^M$  from the momentum encoders, and the corresponding video/clip IDs  $\{y_m\}_{m=1}^M$ . Then we calculate the text-to-video contrastive loss  $\mathcal{L}_{t2v}$  and video-to-text loss  $\mathcal{L}_{v2t}$  as:

$$\begin{cases} \mathcal{L}_{t2v}(\mathbf{t}_i) = - \sum_{k \in \mathcal{P}(i)} \log \frac{\exp(\mathbf{t}_i^T \mathbf{v}_k / \tau)}{\sum_{m=1}^M \exp(\mathbf{t}_i^T \mathbf{v}_m / \tau)}, & (3) \\ \mathcal{L}_{v2t}(\mathbf{v}_i) = - \sum_{k \in \mathcal{P}(i)} \log \frac{\exp(\mathbf{v}_i^T \mathbf{t}_k / \tau)}{\sum_{m=1}^M \exp(\mathbf{v}_i^T \mathbf{t}_m / \tau)}, & (4) \end{cases}$$

where  $\mathcal{P}(i) = \{k | k \in M, y_k = y_i\}$  is the positive sample set, and  $\tau$  is the learnable temperature parameter. Finally we combine the above two losses for the VTC loss  $\mathcal{L}_{VTC}$ :

$$\mathcal{L}_{VTC} = \frac{1}{2} (\mathcal{L}_{t2v} + \mathcal{L}_{v2t}). \quad (5)$$

Note that to compensate for potential false negatives in the momentum encoder, we apply the momentum distillation strategy in ALBEF [19] to generate soft labels.

For the video-grounded text encoder, we apply Video-Text Matching (VTM) loss, which aims to learn a video-text multimodal representation that captures the fine-grained alignment between video and language. VTM corresponds to a binary classification task, where the model uses a VTM head (a linear layer) to predict whether a video-text pair is positive (matched) or negative (unmatched) given its multimodal feature of [Encode] token. Formally, for the

$i$ -th video-text pair, we first calculate their positive matching score  $p_i^+ \in \mathbb{R}$ . Then we randomly sample a video/text to replace the matched video/text to get the negative matching score  $p_i^{-v}/p_i^{-t} \in \mathbb{R}$ . Finally, we calculate the video-text matching loss  $\mathcal{L}_{\text{VTM}}$ :

$$\mathcal{L}_{\text{VTM}} = - \left[ \log(p_i^+) + \log(1 - p_i^{-v}) + \log(1 - p_i^{-t}) \right]. \quad (6)$$

Note that, to make the VTM loss more informative, we sample the negative samples using the hard negative mining strategy [19]. Specifically, we use the contrastive similarity distribution from Eqn. (3) and Eqn. (4) to sample hard negatives, where similar samples have a higher chance to be sampled. We sum  $\mathcal{L}_{\text{VTM}}$  and  $\mathcal{L}_{\text{VTC}}$  to get the final loss.

**Video Captioning.** During training, we apply Language Modeling (LM) Loss for the decoder, which optimizes a cross-entropy loss that trains the model to maximize the likelihood of the text in an auto-regressive manner. Formally, for each video-text pair  $(v, t)$ :

$$\mathcal{L}_{\text{LM}} = - \sum_{l=1}^L \log \left( P(t^l | t^{<l}, v) \right), \quad (7)$$

where  $L$  is the total length of the sentence. We apply a label smoothing of 0.1 when computing the loss. Compared to the masked language modeling loss that has been widely used for video-language pretraining, LM enables the model with the generalization capability to convert visual information into coherent captions.

**Video Question Answering.** *Open-ended QA* adopts LM loss, while *multiple-choice QA* employs VTM loss. Different from that in text-to-video retrieval, the negative samples come from false question-answer pairs instead of sampling.

## 5 EXPERIMENTS

### 5.1 Datasets

To assess the efficacy of video-language models, we conducted experiments on three distinct tasks: text-to-video retrieval, video caption, and video question answering. Each task was evaluated on its corresponding dataset.

**Text-to-Video Retrieval.** (i) **MSRVTT** [45] contains 10K YouTube videos with 200K descriptions, which is split into 9K videos for training and 1K videos for test. (ii) **DiDemo** [2] contains 10K Flickr videos with 40K sentences, where the test set contains 1,000 videos. We follow the standard setting to concatenate all descriptions in a video as a single query, and further evaluate paragraph-to-video retrieval. (iii) **ActivityNet Captions** [16] contains 20K YouTube videos annotated with 100K sentences. We follow the paragraph-to-video retrieval setting to train models on 10K videos and report results on the val1 set with 4.9K videos.

**Video Caption.** (i) **MSRVTT** [45] consists of 10K open-domain video clips, and each clip has 20 ground-truth captions. We use the standard captioning split, which has 6.5K training, 500 validation, and 2.9K testing videos. (ii) **MSVD** [5] contains 1,970 videos from YouTube with 80K descriptions, which is split into 1200, 100 and 670 videos for training, validation and testing, respectively.

**Video Question Answering.** (i) **NeXt-QA** [40] contains about 47.7K manually annotated questions for multi-choice QA collected from 5.4K videos. There are three types of questions: descriptive, temporal, and causal. (ii) **MSRVTT QA** [43] contain 50K QA pairs that focus on the description of video elements. For both two

datasets, we use the val split for model selection, and report the final results on the test split.

### 5.2 Experimental Settings

**5.2.1 Experimental Details.** We present the key implementation details shared across all tasks. Our models were trained on 8 GPUs. On the model architecture, we used ViT-B/16 [9] as the video backbone, and BERT-base [7] as the backbone of the query component. The number of layers  $L$  was 12. Their parameters were initialized from BLIP-base [18] without extra image-text datasets. On optimization techniques, we used an AdamW optimizer with a weight decay of 0.04. The learning rate was scheduled with a linear warm-up with 1,000 iterations, and cosine annealing starting at 10% of training following [36]. For each frame, we took random crops of resolution  $224 \times 224$  as inputs (thus  $P_{\text{in}}$  is 196) and applied RandAugment. We set  $P_{\text{out}}$  to be also 196 for all models. For task-specific hyperparameters, we set the momentum distillation [19] weight in VTC to 0.4. During inference, we used beam search with a beam size of 3. The number of frames  $F$ , segments  $S$ , learning rate, and batch size differ among datasets (see Appendix A).

**5.2.2 Evaluation Metrics.** For text-to-video retrieval, the metrics [10, 30] recall at rank  $K$  (**R@K**, higher is better) calculate the percentage of test samples with the correct result found in the Top- $K$  retrieved points to the query sample. We report the  $K = 1, 5, 10$ . For video captioning, we adopted BLEU-4 (**B-4**), ROUGE-L (**R-L**), and CIDEr (**C**) metrics [29] (higher is better for all). For video question answering, we reported accuracy (**Acc**). Additionally, we reported the accuracy of each subset in NeXt-QA [40] dataset.

### 5.3 Performance Comparison

In this section, we compared our method with recent state-of-the-art methods. Given that some baselines have several variants (differences in model capacity, training data, and post-processing). We ensured a fair comparison by selecting baselines with similar model capacity (ViT/B-16) if multiple variants were present. We grouped all baselines into with/without video-language pretraining methods (denoted as PT in all tables). For text-to-video retrieval, we did NOT use the post-processing trick Dual Softmax Loss (DSL) [6], and selected the baseline variants in the same manner. This is because DSL requires one-one mapping of text-video pairs, which is not a feasible approach in real-world scenarios. Furthermore, for the HiTeA [50] model in the NeXt-QA dataset, we executed the model once again to report its results in the test split rather than the original validation split.

**5.3.1 Text-to-Video Retrieval.** Results on text-to-video retrieval are presented in Table 1. Upon careful examination of the results, three key observations can be made.

- Our model consistently outperforms all methods without video-language pre-training. Furthermore, our method exhibits comparable performance to pre-training based models on MSR-VTT datasets and even surpasses their performance on DiDemo and ActivityNet-Captions datasets. These findings provide compelling evidence for the superiority of our RTQ framework, particularly considering that pre-training based approaches are trained with significantly larger amounts of data. For example, the CLIP-ViP

**Table 1: Comparison of text-to-video retrieval. "PT" stands for video-language pretraining.**

Model	PT	MSR-VTT				DiDemo				ActivityNet-Captions			
		R@1	R@5	R@10	MdR	R@1	R@5	R@10	MdR	R@1	R@5	R@10	MdR
BridgeFormer [13]	✓	37.6	64.8	75.1	3.0	37.0	62.2	73.9	3.0	-	-	-	-
OmniVL [34]		47.8	74.2	83.8	-	52.4	79.5	85.4	-	-	-	-	-
HiTeA [50]		46.8	71.2	81.9	-	56.5	81.7	89.7	-	49.7	77.1	86.7	-
mPLUG2-B [44]		48.3	75.0	83.2	-	52.3	80.8	87.5	-	-	-	-	-
STOA-VLP [54]		50.1	75.5	83.8	-	51.1	76.4	84.0	-	-	-	-	-
CLIP-ViP [46]		<b>54.2</b>	<b>77.2</b>	<b>84.8</b>	1.0	50.5	78.4	87.1	1.0	<u>53.4</u>	<u>81.4</u>	<u>90.0</u>	1.0
ClipBERT [17]	✗	22.0	46.8	59.9	6.0	20.4	48.0	60.8	6.0	21.3	49.0	63.5	6.0
X-Pool [14]		46.9	72.8	82.2	2.0	-	-	-	-	-	-	-	-
CenterCLIP [52]		48.4	73.8	82.0	1.0	-	-	-	-	46.2	77.0	87.6	2.0
STAN [24]		50.0	75.2	84.1	1.5	49.4	74.9	84.5	1.0	-	-	-	-
Cap4Video [39]		51.4	75.7	83.9	1.0	52.0	79.4	87.5	1.0	-	-	-	-
Ours		<u>53.4</u>	<u>76.1</u>	<u>84.4</u>	1.0	<b>57.6</b>	<b>84.1</b>	<b>89.8</b>	1.0	<b>53.5</b>	<b>81.4</b>	<b>91.9</b>	1.0

**Table 2: Comparison of video captioning.**

Model	PT	MSR-VTT			MSVD		
		B-4	R-L	CIDEr	B-4	R-L	CIDEr
UniVL [26]	✓	41.8	60.8	50.0	-	-	-
MV-GPT [32]		48.9	64.0	60.0	-	-	-
STOA-VLP [54]		45.8	<b>68.4</b>	60.2	<u>64.4</u>	<b>83.9</b>	<b>131.8</b>
Clip4Cap [33]	✗	46.1	63.7	57.7	-	-	-
HMN [49]		43.5	62.7	51.5	59.2	75.1	104.0
SwinBERT [23]		45.4	64.1	55.9	58.2	77.5	120.6
CMVC [48]		48.2	64.8	58.7	-	-	-
TextKG [15]		46.6	64.8	<u>60.8</u>	60.8	75.1	105.2
Ours		<b>49.6</b>	<u>66.1</u>	<b>69.3</b>	<b>66.9</b>	<u>82.8</u>	<u>123.4</u>

**Table 3: Comparison of video question answering.**

Model	PT	MSR-VTT	NExT-QA			
		Acc	Acc	Acc-D	Acc-T	Acc-C
MV-GPT [32]	✓	41.7	-	-	-	-
Flamingo [1]		<b>47.4</b>	-	-	-	-
HiTeA [50]		45.4	<u>62.4</u>	<u>75.5</u>	<u>58.7</u>	<u>60.6</u>
mPLUG2-B [44]		<u>46.3</u>	-	-	-	-
CoVGT (PT) [42]		40.0	59.7	68.4	58.0	58.0
ATP [4]		-	54.3	66.8	50.2	53.1
IGV [21]	✗	38.3	51.3	59.6	51.7	48.6
HQGA [41]		38.6	51.8	59.4	52.3	49.0
CoVGT [42]		38.3	59.4	66.8	57.0	58.5
Ours		42.1	<b>63.2</b>	<b>75.6</b>	<b>59.6</b>	<b>61.4</b>

model [46] is trained on a dataset comprising 100 million video-text pairs, and initialized from CLIP [31] model.

- Our method surpasses OmniVL [34] and mPLUG [44], despite their inclusion of temporal modeling and query structures. This result emphasizes the criticality of our refinement module.
- Our method achieves the best performance in DiDemo and ActivityNet Captions datasets, and comparable performance in the MSR-VTT dataset. This is because videos in the former two datasets are all untrimmed videos with longer average duration and richer content, making our method more effective as it suffers less from information redundancy and scene complexity.

**5.3.2 Video Captioning.** Outcomes on captioning are illustrated in Table 2. These results have led to three notable observations.

- Our model is consistently superior to all non-pretraining methods, and comparable with pretraining methods, demonstrating the superiority of our RTQ framework in generation tasks.

- In contrast to pretraining methods, our model exhibits superior performance in the Bleu metric but inferior results in the Rouge metric. Given that Bleu is typically indicative of precision and Rouge measures recall, a plausible explanation is that our RTQ framework possesses superior learning capabilities, resulting in lower training error on the dataset and hence, a higher degree of precision. Conversely, pretraining methods incorporate more external knowledge and thus exhibit better recall.
- Compared to pretraining methods on CIDEr metric, our model exhibits superiority in the MSR-VTT dataset, albeit inferiority in the MSVD dataset. As CIDEr metric considers both semantic similarity and diversity, it provides a more comprehensive evaluation for video captioning. Given the fact that MSVD has 5 times fewer videos than MSR-VTT, it is likely that additional knowledge acquired through pre-training could improve performance in the smaller MSVD dataset.

**5.3.3 Video Question Answering.** Results of video question answering are illustrated in Table 3. There are two observations.

- Our model is consistently superior to all non-pretraining methods, which demonstrates the superiority of our RTQ framework in complex video language tasks.
- Our model is superior to pretraining methods in multiple-choice QA (NExT-QA), while inferior in open-ended QA (MSR-VTT). We attribute the strong open-ended ability of pretraining models to the modal capacity and large pretraining corpus. Specifically, Flamingo has 80B parameters in total, while ours has only around 400M. Besides, Flamingo uses 2 billion image-text pairs and 27 million video-text pairs, while mPLUG-B uses 2 million video-text pairs and a large natural language corpus (WikiCorpus and common crawl) as the pre-training data.

It is noteworthy that HiTeA [50] formulates MSR-VTT QA as a multiple-choice problem, where the best matching word is selected as the answer. Since all answers in MSR-VTT consist of only a single word in a limited vocabulary (1.5K), such a setting leads to an overestimation of its open-ended ability.

## 5.4 Ablation Studies

This section presents a series of experiments aimed at analyzing the efficacy of our model, which are reported in Table 4.

**Table 4: Ablation studies of our RTQ framework. We report accuracies on the NeXt-QA [40] testset with descriptive (D), temporal (T), and causal (C) splits.**

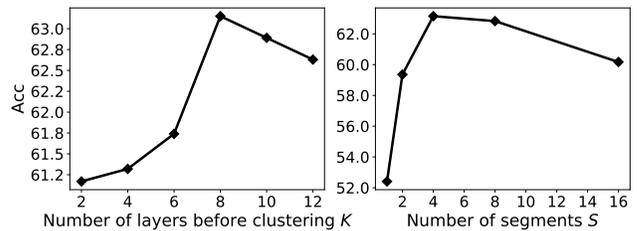
Model	Acc-C	Acc-D	Acc-T	Acc
Ours w/o R	58.01	71.95	57.55	60.17
Ours w/o T	58.92	75.27	57.10	61.04
Ours w/o Q	59.78	61.42	57.92	59.43
Ours w/ divST [3]	58.64	74.80	55.55	60.34
Ours w/ STAN [24]	58.49	73.95	55.56	60.12
Ours w/ Uniformer [20]	55.23	67.88	52.56	56.47
Ours w/ TSM [22]	57.60	74.45	55.55	59.73
Ours w/ token select [25]	55.92	70.08	56.11	58.32
<b>Ours</b>	<b>61.39</b>	<b>75.58</b>	<b>59.57</b>	<b>63.15</b>

**5.4.1 Ablation of the RTQ modules.** We first examined the distinct impacts of the R, T, and Q components on our model. We conducted three variations, namely: 1) **w/o R**, which removes the clustering layer; 2) **w/o T**, which removes the temporal modeling layer; and 3) **w/o Q**, which deletes the query layer. The experiments are conducted in the NeXt-QA [40] dataset to illustrate the differences in functionality among components. For the **w/o Q** implementation, we concatenated the question and answer candidate to form a sentence, and then formulated a matching problem between videos and question-answer pairs by computing the cosine similarity of their [CLS] tokens. Our observations are presented below.

- The refinement module has demonstrated consistent contributions across all question types, particularly for casual ones. This outcome can be attributed to the elimination of redundant information, which enhances the model’s comprehension of crucial objects, actions, and events depicted in the videos. Such improvement is beneficial for answering all categories of questions.
- The temporal modeling module exhibits a significant impact on temporal questions while exhibiting minimal influence on descriptive ones. This phenomenon is expected as the temporal module is intended to capture temporal cues.
- The query module has substantially contributed to descriptive questions, whilst demonstrating relatively minor effects on other question types. This outcome can be explained by the fact that descriptive questions demand a highly detailed understanding of the video content, which is facilitated by the query module. Conversely, temporal and causal questions require a more holistic understanding of the overall video content.

**5.4.2 Ablation of other designing choices.** We then replace modules in our framework for other design choices. Both **w/ divST** and **w/ STAN** replace the message token mechanism with temporal attention [3, 24]; **w/ Uniformer** and **w/ TSM** replaces the message token mechanism with unified transformer (temporal convolution and attention) [20] and temporal shift module [22], respectively. **w/ token select** replace the clustering module with token selection module [25]. We gain two observations.

- Methods of temporal attention and temporal shifting mechanism generally harm the performance. As analyzed in Section 4.3, these methods necessitate the spatial consistency in the patches of input frames, which is not guaranteed in our framework.



**Figure 4: Sensitivity analysis.**

- Token selection method also harms the performance. Considering that token selection method is built without any prior information. It may over-select patches, resulting fewer useful information for the query module.

## 5.5 Sensitivity Analysis

Our study aimed to investigate the impact of hyper-parameters on the performance of our model, as illustrated in Fig. 4. Specifically, we focused on two primary hyper-parameters: the number of layers preceding the cluster layer  $K$  and the number of segments of clustering  $S$ . We conducted experiments on the NeXt-QA dataset using 16 input frames ( $F$ ). Our investigation leads to two key observations.

- The clustering layer should be positioned within the deep layers of ViT. Specifically, the optimal number of layers before clustering is around 8. This is consistent with the design rationale outlined in Section 4. Specifically, the model should prioritize the semantic meaning of patches over their appearance in order to eliminate truly redundant patches. As deep layers contain more semantic information, they prove more effective in this regard.
- Secondly, we found that the optimal number of segments is approximately half of the input frames. This is a reasonable number, as reducing the number of segments increases patch purity but reduces patch integrity. Therefore, a trade-off must be struck between them.

## 6 CONCLUSION

Our systemic analysis reveals that current video-language understanding methods focus on limited aspects of the task, and methods targeting different challenges can complement each other. In light of this, we propose a framework integrating the refinement, temporal modeling, and query components to jointly tackle information redundancy, temporal dependency, and scene complexity, respectively. Remarkably, our method achieves superior (or comparable) performance to state-of-the-art pretraining methods without requiring video-language pretraining.

To further enhance the performance of our model, there are several potential directions. Firstly, video-language pretraining could be used to acquire more world knowledge, which is especially helpful for open-ended QA tasks. Secondly, developing more effective refinement, temporal modeling, and query modules could elevate the overall performance of our approach.

## 7 ACKNOWLEDGEMENTS

This work is supported by the National Natural Science Foundation of China, No.: 62176137, and No.: 62006140; the Shandong Provincial Natural Science and Foundation, No.: ZR2020QF106.

## REFERENCES

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L. Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. 2022. Flamingo: a Visual Language Model for Few-Shot Learning. In *Advances in Neural Information Processing Systems*. Curran Associates, 23716–23736.
- [2] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. 2017. Localizing moments in video with natural language. In *International Conference on Computer Vision*. IEEE, 5803–5812.
- [3] 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 *Conference on Computer Vision and Pattern Recognition*. IEEE, 1708–1718.
- [4] Shyamal Buch, Cristóbal Eyzaguirre, Adrien Gaidon, Jiajun Wu, Li Fei-Fei, and Juan Carlos Niebles. 2022. Revisiting the "Video" in Video-Language Understanding. In *Conference on Computer Vision and Pattern Recognition*. IEEE, 2907–2917.
- [5] David Chen and William B Dolan. 2011. Collecting highly parallel data for paraphrase evaluation. In *Annual Meeting of the Association for Computational Linguistics*. The Association for Computational Linguistics, 190–200.
- [6] Xing Cheng, Hezheng Lin, Xiangyu Wu, Fan Yang, and Dong Shen. 2021. Improving Video-Text Retrieval by Multi-Stream Corpus Alignment and Dual Softmax Loss. arXiv:2109.04290
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Conference of the North American Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics, 4171–4186.
- [8] Xingning Dong, Tian Gan, Xueming Song, Jianlong Wu, Yuan Cheng, and Liqiang Nie. 2022. Stacked Hybrid-Attention and Group Collaborative Learning for Unbiased Scene Graph Generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, 19405–19414.
- [9] 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. In *International Conference on Learning Representations*. OpenReview.net.
- [10] Yali Du, Yinwei Wei, Wei Ji, Fan Liu, Xin Luo, and Liqiang Nie. 2023. Multi-queue Momentum Contrast for Microvideo-Product Retrieval. In *Proceedings of International Conference on Web Search and Data Mining*. ACM, 1003–1011.
- [11] Tian Gan, Junnan Li, Yongkang Wong, and Mohan S. Kankanhalli. 2019. A Multi-sensor Framework for Personal Presentation Analytics. *ACM Transactions on Multimedia Computing, Communications, and Applications* 15, 2 (2019), 30:1–30:21.
- [12] Tian Gan, Qing Wang, Xingning Dong, Xiangyuan Ren, Liqiang Nie, and Qingpei Guo. 2023. CNVid-3.5M: Build, Filter, and Pre-Train the Large-Scale Public Chinese Video-Text Dataset. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, 14815–14824.
- [13] Yuying Ge, Yixiao Ge, Xihui Liu, Dian Li, Ying Shan, Xiaohu Qie, and Ping Luo. 2022. Bridging Video-Text Retrieval With Multiple Choice Questions. In *Conference on Computer Vision and Pattern Recognition*. IEEE, 16167–16176.
- [14] Satya Krishna Gorti, Noël Vouitsis, Junwei Ma, Keyvan Golestan, Maksims Volkovs, Animesh Garg, and Guangwei Yu. 2022. X-Pool: Cross-Modal Language-Video Attention for Text-Video Retrieval. In *Conference on Computer Vision and Pattern Recognition*. IEEE, 4996–5005.
- [15] Xin Gu, Guang Chen, Yufei Wang, Libo Zhang, Tiejian Luo, and Longyin Wen. 2023. Text with Knowledge Graph Augmented Transformer for Video Captioning. In *Conference on Computer Vision and Pattern Recognition*. IEEE, 1175–1175.
- [16] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. 2017. Dense-captioning events in videos. In *International Conference on Computer Vision*. IEEE, 706–715.
- [17] Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L. Berg, Mohit Bansal, and Jingjing Liu. 2021. Less Is More: ClipBERT for Video-and-Language Learning via Sparse Sampling. In *Conference on Computer Vision and Pattern Recognition*. IEEE, 7331–7341.
- [18] Junnan Li, Dongxu Li, Caimeing Xiong, and Steven C. H. Hoi. 2022. BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation. In *International Conference on Machine Learning*. PMLR, 12888–12900.
- [19] Junnan Li, Ramprasaath R. Selvaraju, Akhilesh Gotmare, Shafiq R. Joty, Caimeing Xiong, and Steven Chu-Hong Hoi. 2021. Align before Fuse: Vision and Language Representation Learning with Momentum Distillation. In *Advances in Neural Information Processing Systems*. Curran Associates, 9694–9705.
- [20] Kunlun Li, Yali Wang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Qiao. 2022. UniFormer: Unified Transformer for Efficient Spatial-Temporal Representation Learning. In *International Conference on Learning Representations*. OpenReview.net.
- [21] Yicong Li, Xiang Wang, Junbin Xiao, Wei Ji, and Tat-Seng Chua. 2022. Invariant Grounding for Video Question Answering. In *Conference on Computer Vision and Pattern Recognition*. IEEE, 2918–2927.
- [22] Ji Lin, Chuang Gan, and Song Han. 2019. TSM: Temporal Shift Module for Efficient Video Understanding. In *International Conference on Computer Vision*. IEEE, 7082–7092.
- [23] Kevin Lin, Linjie Li, Chung-Ching Lin, Faisal Ahmed, Zhe Gan, Zicheng Liu, Yumao Lu, and Lijuan Wang. 2022. SwinBERT: End-to-End Transformers with Sparse Attention for Video Captioning. In *Conference on Computer Vision and Pattern Recognition*. IEEE, 17928–17937.
- [24] Ruyang Liu, Jingjia Huang, Ge Li, Jiashi Feng, Xinglong Wu, and Thomas H. Li. 2023. Revisiting Temporal Modeling for CLIP-based Image-to-Video Knowledge Transferring. In *Conference on Computer Vision and Pattern Recognition*. IEEE, 6422–6431.
- [25] Yuqi Liu, Pengfei Xiong, Luhui Xu, Shengming Cao, and Qin Jin. 2022. TS2-Net: Token Shift and Selection Transformer for Text-Video Retrieval. In *European Conference on Computer Vision*. Springer, 319–335.
- [26] Huaishao Luo, Lei Ji, Botian Shi, Haoyang Huang, Nan Duan, Tianrui Li, Jason Li, Taroon Bharti, and Ming Zhou. 2020. UniVL: A Unified Video and Language Pre-Training Model for Multimodal Understanding and Generation. arXiv:2002.06353
- [27] Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. 2022. CLIP4Clip: An empirical study of CLIP for end to end video clip retrieval and captioning. *Neurocomputing* 508 (2022), 293–304.
- [28] Bolin Li, Houwen Peng, Minghao Chen, Songyang Zhang, Gaofeng Meng, Jianlong Fu, Shiming Xiang, and Haibin Ling. 2022. Expanding Language-Image Pretrained Models for General Video Recognition. In *European Conference on Computer Vision*, Vol. 13664. Springer, 1–18.
- [29] Liqiang Nie, Leigang Qu, Dai Meng, Min Zhang, Qi Tian, and Alberto Del Bimbo. 2022. Search-oriented Micro-video Captioning. In *International Conference on Multimedia*. ACM, 3234–3243.
- [30] Leigang Qu, Meng Liu, Jianlong Wu, Zan Gao, and Liqiang Nie. 2021. Dynamic Modality Interaction Modeling for Image-Text Retrieval. In *SIGIR Conference on Research and Development in Information Retrieval*. ACM, 1104–1113.
- [31] 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 *International Conference on Machine Learning*. PMLR, 8748–8763.
- [32] Paul Hongsuck Seo, Arsha Nagrani, Anurag Arnab, and Cordelia Schmid. 2022. End-to-end Generative Pretraining for Multimodal Video Captioning. In *Conference on Computer Vision and Pattern Recognition*. IEEE, 17938–17947.
- [33] Mingkang Tang, Zhanyu Wang, Zhenhua Liu, Fengyun Rao, Dian Li, and Xiu Li. 2021. CLIP4Caption: CLIP for Video Caption. In *International Conference on Multimedia*. ACM, 4858–4862.
- [34] Junke Wang, Dongdong Chen, Zuxuan Wu, Chong Luo, Luowei Zhou, Yucheng Zhao, Yujia Xie, Ce Liu, Yu-Gang Jiang, and Lu Yuan. 2022. OmniVL: One Foundation Model for Image-Language and Video-Language Tasks. In *Advances in Neural Information Processing Systems*. Curran Associates.
- [35] Shaokun Wang, Tian Gan, Yuan Liu, Li Zhang, Jianlong Wu, and Liqiang Nie. 2022. Discover Micro-Influencers for Brands via Better Understanding. *IEEE Transactions on Multimedia* 24 (2022), 2595–2605.
- [36] Xiao Wang, Tian Gan, Yinwei Wei, Jianlong Wu, Dai Meng, and Liqiang Nie. 2022. Micro-video Tagging via Jointly Modeling Social Influence and Tag Relation. In *MM '22: The 30th ACM International Conference on Multimedia, Lisboa, Portugal, October 10 - 14, 2022*. ACM, 4478–4486.
- [37] Yinwei Wei, Xiang Wang, Weili Guan, Liqiang Nie, Zhouchen Lin, and Baoquan Chen. 2019. Neural multimodal cooperative learning toward micro-video understanding. *Transactions on Image Processing* 29 (2019), 1–14.
- [38] Jianlong Wu, Wei Sun, Tian Gan, Ning Ding, Feijun Jiang, Jialie Shen, and Liqiang Nie. 2023. Neighbor-Guided Consistent and Contrastive Learning for Semi-Supervised Action Recognition. *IEEE Transactions on Image Processing* 32 (2023), 2215–2227.
- [39] Wenhao Wu, Haipeng Luo, Bo Fang, Jingdong Wang, and Wanli Ouyang. 2023. Cap4Video: What Can Auxiliary Captions Do for Text-Video Retrieval?. In *Conference on Computer Vision and Pattern Recognition*. IEEE, 1618–1628.
- [40] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. 2021. NEX-T-QA: Next Phase of Question-Answering to Explaining Temporal Actions. In *Conference on Computer Vision and Pattern Recognition*. IEEE, 9777–9786.
- [41] Junbin Xiao, Angela Yao, Zhiyuan Liu, Yicong Li, Wei Ji, and Tat-Seng Chua. 2022. Video as Conditional Graph Hierarchy for Multi-Granular Question Answering. In *AAAI Conference on Artificial Intelligence*. AAAI Press, 2804–2812.
- [42] Junbin Xiao, Pan Zhou, Angela Yao, Yicong Li, Richang Hong, Shuicheng Yan, and Tat-Seng Chua. 2023. Contrastive Video Question Answering via Video Graph Transformer. arXiv:2302.13668
- [43] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. 2017. Video question answering via gradually refined attention over appearance and motion. In *International Conference on Multimedia*. ACM, 1645–1653.

- [44] Haiyang Xu, Qinghao Ye, Ming Yan, Yaya Shi, Jiabo Ye, Yuanhong Xu, Chenliang Li, Bin Bi, Qi Qian, Wei Wang, Guohai Xu, Ji Zhang, Songfang Huang, Fei Huang, and Jingren Zhou. 2023. mPLUG-2: A Modularized Multi-modal Foundation Model Across Text, Image and Video. arXiv:2302.00402
- [45] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. Msr-vtt: A large video description dataset for bridging video and language. In *Conference on Computer Vision and Pattern Recognition*. IEEE, 5288–5296.
- [46] Hongwei Xue, Yuchong Sun, Bei Liu, Jianlong Fu, Ruihua Song, Houqiang Li, and Jiebo Luo. 2023. CLIP-ViP: Adapting Pre-trained Image-Text Model to Video-Language Representation Alignment. In *International Conference on Learning Representations*. OpenReview.net.
- [47] Shen Yan, Tao Zhu, Zirui Wang, Yuan Cao, Mi Zhang, Soham Ghosh, Yonghui Wu, and Jiahui Yu. 2023. VideoCoCa: Video-Text Modeling with Zero-Shot Transfer from Contrastive Captioners. arXiv:2212.04979
- [48] Bang Yang, Tong Zhang, and Yuexian Zou. 2022. CLIP Meets Video Captioning: Concept-Aware Representation Learning Does Matter. In *Chinese Conference of Pattern Recognition and Computer Vision*. Springer, 368–381.
- [49] Hanhua Ye, Guorong Li, Yuankai Qi, Shuhui Wang, Qingming Huang, and Ming-Hsuan Yang. 2022. Hierarchical Modular Network for Video Captioning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18–24, 2022*. IEEE, 17918–17927.
- [50] Qinghao Ye, Guohai Xu, Ming Yan, Haiyang Xu, Qi Qian, Ji Zhang, and Fei Huang. 2022. HiTeA: Hierarchical Temporal-Aware Video-Language Pre-training. arXiv:2212.14546
- [51] Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, and Yejin Choi. 2021. MERLOT: Multimodal Neural Script Knowledge Models. In *Advances in Neural Information Processing Systems*. Curran Associates, 23634–23651.
- [52] Shuai Zhao, Linchao Zhu, Xiaohan Wang, and Yi Yang. 2022. CenterCLIP: Token Clustering for Efficient Text-Video Retrieval. In *International SIGIR Conference on Research and Development in Information Retrieval*. ACM, 970–981.
- [53] Huasong Zhong, Jianlong Wu, Chong Chen, Jianqiang Huang, Minghua Deng, Liqiang Nie, Zhouchen Lin, and Xian-Sheng Hua. 2021. Graph Contrastive Clustering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. IEEE, 9204–9213.
- [54] Weihong Zhong, Mao Zheng, Duyu Tang, Xuan Luo, Heng Gong, Xiaocheng Feng, and Bing Qin. 2023. STOA-VLP: Spatial-Temporal Modeling of Object and Action for Video-Language Pre-training. arXiv:2302.09736

## Appendix A EXPERIMENTAL DETAILS

The number of frames  $F$ , segments  $S$ , top  $Q$  in the recall stage, learning rate, and batch size differ among datasets. We will detail them in Table A.1 to Table A.3.

**Table A.1: Experimental Details of text-to-video retrieval.**

Dataset	$F$	$S$	lr	bs	epoch	$Q$
MSR-VTT	12	6	2e-6	64	6	128
DiDemo	12	6	2e-5	128	10	128
ActivityNet-Captions	32	16	2e-5	64	10	512

**Table A.2: Experimental Details of video captioning.**

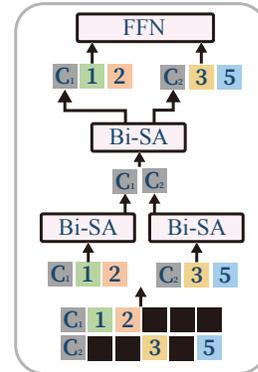
Dataset	$F$	$S$	lr	bs	epoch
MSR-VTT	12	6	5e-6	128	10
MSVD	12	6	2e-6	64	10

**Table A.3: Experimental Details of video question answering.**

Dataset	$F$	$S$	lr	bs	epoch
MSR-VTT	12	6	2e-5	128	6
NeXt-QA	16	4	1e-5	64	6

## Appendix B ADDITIONAL ILLUSTRATION

We illustrate the temporal modeling module of our model in Fig. B.1.



**Figure B.1: Illustration of the Temporal Modeling component.**

## Appendix C ADDITIONAL EXPLANATION ON CLUSTERING

On the logical relationship between the clustering results and three challenges: The logical relationship stems from the fact that each method within the clustering graph addresses a particular challenge discussed in the introduction (analyzed in Section 2.2). Consequently, by clustering their prediction results, we can identify shared advantages and disadvantages among these methods.