# Video Retrieval of Autonomous Driving Scenarios Based on Large Models

12310520 Rui Yuhan 12310437 Qiao Shihan

*Abstract*—Video retrieval in autonomous driving scenarios demands models to comprehend complex visual environments and align them with detailed textual descriptions. While multimodal large language models (MLLMs) have demonstrated potential in vision-language tasks, their application to domain-specific video retrieval remains underexplored. Furthermore, the scarcity of publicly available high-quality annotated video data in the autonomous driving domain makes it challenging for models to accurately learn video content semantics from limited data. Additionally, significant differences exist between autonomous driving videos and those in general video datasets, leading to suboptimal performance of generic large models on autonomous driving datasets.

To address these challenges, we evaluate the performance of two distinct pretrained MLLMs on a carefully curated autonomous driving dataset. This dataset features high-quality captions that have undergone multi-round optimization, combining human annotation expertise with advanced automated labeling techniques, while also serving as a benchmark platform for exploring more efficient annotation strategies. Experimental results demonstrate that our approach effectively retrieves relevant video segments based on textual queries, highlighting the potential of MLLMs to enhance retrieval accuracy and scalability.

*Index Terms*—Autonomous Driving Video Retrieval, Multimodal Large Language Models (MLLMs), Video Understanding, Automated Annotation Optimization

## I. INTRODUCTION

## II. TWO KNOWLEDGE DISTILLATION METHODS USED
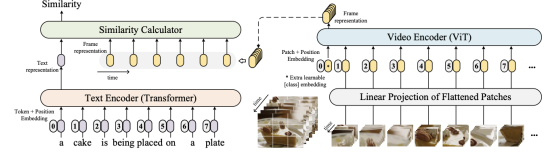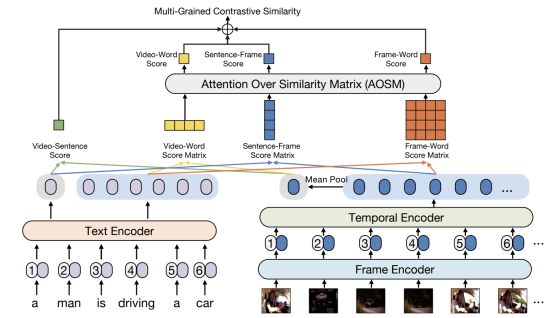
### A. DLDKD

### B. TeachCLIP

*1) Inspiration:* Existing VTR models generally fall into two categories: lightweight global feature models (e.g., CLIP4Clip) and heavy fine-grained models (e.g., X-CLIP). The former is computationally efficient but often lacks frame-level details, resulting in suboptimal retrieval accuracy. The latter leverages frame-level interactions to enhance performance but incurs high computational costs.

In terms of vision and text encoding, both CLIP4Clip and X-CLIP share nearly identical bottom-level encoding mechanisms (Frame → Patch → ViT, Word → Embedding → Transformer). The primary distinction lies in how they utilize these features in the upper layers: CLIP4Clip performs pooling to retain only global video vectors, whereas X-CLIP preserves fine-grained features and performs multi-granularity alignment.

TeachCLIP aims to bridge this gap by distilling the fine-grained alignment capability of heavy models into a lightweight student model.



(a) CLIP4Clip Framework



(b) X-CLIP Framework

Fig. 1. Comparison of CLIP4Clip and X-CLIP frameworks. Both use similar bottom-up encoding but differ in upper-layer feature usage.
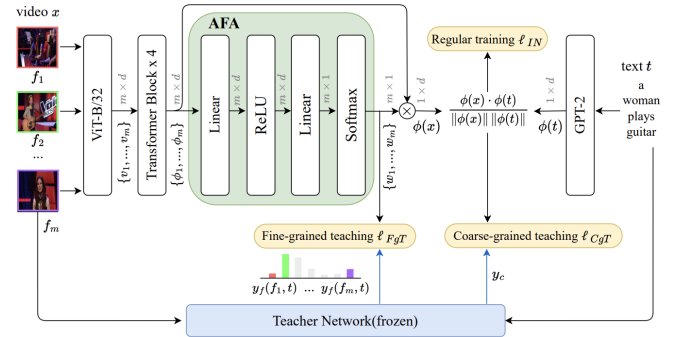


Fig. 2. The framework of TeachCLIP.

*2) Framework:* The TeachCLIP framework is designed to transfer the fine-grained knowledge from a teacher model to a student model while maintaining high inference efficiency.

**Student Model:** The student model is built upon CLIP4Clip. It samples the video into frames, processes each frame through the CLIP ViT to obtain frame features, and enhances them via a temporal Transformer. A key innovation is replacing the original mean pooling with Attentional Frame-Feature Aggregation (AFA). AFA generates frame weights $\{w_i\}$ to compute a weighted sum of frame features:

$$\phi(x) = \sum_{i=1}^{m} w_i \cdot \phi_i \qquad (1)$$

The AFA module consists of a lightweight structure (Linear → ReLU → Linear → Softmax) and introduces negligible parameters.

**Teacher Model:** The teacher model provides two types of supervision:

- **Video-level soft labels** ($y_c$)**:** The teacher provides the distribution of video-text correlations to guide the student's ranking.
- **Frame-level soft labels** ($y_f$)**:** The teacher provides the distribution of frame-text correlations to guide the student's frame weight allocation in AFA.

**Learning Objectives:** The training involves three losses:

1) *Frame-level Distillation* ($\ell_{FgT}$)*:* Optimizes the AFA weights to match the teacher's frame-text similarity distribution. The loss function is defined as:

$$\ell_{FgT} = -\frac{1}{b} \sum_{i=1}^{b} \sum_{k=1}^{m} y_f(f_{i,k}, t_i) \log w_{i,k} \qquad (2)$$

where $b$ is the batch size, $m$ is the number of frames, and $y_f$ represents the frame-text similarity from the teacher.

2) *Video-level Distillation* ($\ell_{CgT}$)*:* Aligns the student's video-text similarity matrix with the teacher's multi-grained similarity matrix using the Pearson distance $d_p$:

$$
\begin{aligned}
\ell_{CgT} = &\frac{1}{b} \sum_i d_p(\sigma(B_{i,\cdot}), \sigma(y_c(v_i, \cdot))) \\
&+ \frac{1}{b} \sum_j d_p(\sigma(B_{\cdot,j}), \sigma(y_c(\cdot, t_j)))
\end{aligned}
\qquad (3)
$$

where $\sigma$ denotes the softmax function and $B$ is the similarity matrix of the student.

3) *Contrastive Learning* ($\ell_{IN}$)*:* Standard maximization of similarity for positive pairs and minimization for negative pairs, using symmetric InfoNCE loss:

$$\ell_{IN} = \frac{1}{2} \left[ \ell_{NCE}^{v \to t} + \ell_{NCE}^{t \to v} \right] \qquad (4)$$

where

$$\ell_{NCE}^{v \to t} = -\frac{1}{b} \sum_{i=1}^{b} \log \frac{\exp(B_{ii}/\tau)}{\sum_{j=1}^{b} \exp(B_{ij}/\tau)} \qquad (5)$$

and $\ell_{NCE}^{t \to v}$ is defined symmetrically.

The total loss is $\ell = \ell_{CgT} + \ell_{FgT} + \ell_{IN}$. During inference, the teacher model is discarded, and only the lightweight student model with AFA is used.

*3) TeachCLIP vs X-CLIP Parameters and FLOPs Comparison:* Table I presents the comparison of parameters and FLOPs between TeachCLIP and the teacher model (X-CLIP).

## III. ALTERATIONS ON TEACHCLIP

### A. On student model

### B. On teacher model

## IV. CONCLUSION

Autonomous driving video retrieval presents unique challenges due to the complexity of real-world driving scenarios and the scarcity of high-quality annotated datasets. In this work, we explored the potential of multimodal large language models (MLLMs) to bridge the gap between visual perception and textual understanding in autonomous driving contexts. By leveraging a carefully curated dataset with optimized annotations, combining human expertise and automated refinement using Video-LLama2, we demonstrated that pretrained MLLMs can be effectively adapted for domain-specific video retrieval tasks.

Our experiments with two state-of-the-art MLLMs (Vast and Clip-Vip) revealed that fine-tuning on high-quality annotated data significantly improves retrieval accuracy, validating the importance of robust dataset construction. Additionally, our work establishes a benchmark for future research in autonomous driving video understanding, providing a foundation for developing more interpretable and scalable retrieval systems.

Moving forward, we anticipate that further advancements in MLLMs, combined with more sophisticated annotation pipelines, will enhance the generalization capabilities of autonomous driving systems. Future research could explore dynamic annotation strategies, real-time retrieval optimization, and the integration of multi-sensor data (e.g., LiDAR and radar) to further improve model performance in complex driving environments.

Ultimately, this study highlights the critical role of high-quality data and domain-specific adaptation in advancing autonomous driving technologies, paving the way for safer and more intelligent transportation systems.

TABLE I
COMPARISON OF PARAMETERS AND FLOPs

| Model | Parameters | FLOPs (Inference) | Description |
|---|---|---|---|
| **TeachCLIP** | $\approx 200$M | 53.65G (12 frames) | Student model, based on CLIP4Clip + AFA, maintains lightweight inference |
| **X-CLIP (Teacher)** | $\approx 220$M | 145G (8 frames) / 287G (16 frames) | Teacher model, introduces multi-grained contrastive similarity, higher inference cost |