

A Survey on Video Temporal Grounding with Multimodal Large Language Model

Jianlong Wu, *Member, IEEE*, Wei Liu, Ye Liu, Meng Liu, *Member, IEEE*,
Liqiang Nie, *Senior Member, IEEE*, Zhouchen Lin, *Fellow, IEEE*, and Chang Wen Chen, *Fellow, IEEE*

Abstract—The recent advancement in video temporal grounding (VTG) has significantly enhanced fine-grained video understanding, primarily driven by multimodal large language models (MLLMs). With superior multimodal comprehension and reasoning abilities, VTG approaches based on MLLMs (VTG-MLLMs) are gradually surpassing traditional fine-tuned methods. They not only achieve competitive performance but also excel in generalization across zero-shot, multi-task, and multi-domain settings. Despite extensive surveys on general video-language understanding, comprehensive reviews specifically addressing VTG-MLLMs remain scarce. To fill this gap, this survey systematically examines current research on VTG-MLLMs through a three-dimensional taxonomy: 1) the functional roles of MLLMs, highlighting their architectural significance; 2) training paradigms, analyzing strategies for temporal reasoning and task adaptation; and 3) video feature processing techniques, which determine spatiotemporal representation effectiveness. We further discuss benchmark datasets, evaluation protocols, and summarize empirical findings. Finally, we identify existing limitations and propose promising research directions. For additional resources and details, readers are encouraged to visit our repository at <https://github.com/ki-lw/Awesome-MLLMs-for-Video-Temporal-Grounding>.

Index Terms—video-language understanding, video temporal grounding, fine-grained temporal understanding, vision-language model, large language model, multimodal learning.

1 INTRODUCTION

THE proliferation of untrimmed video content across domains such as surveillance, entertainment, and autonomous systems has created an urgent need for systems capable of precise temporal understanding. Real-world applications, including moment retrieval, scene editing, and temporal question answering, demand accurate identification of not only what events occur but precisely when they occur. Existing video-language models primarily focus on global or coarse-level video comprehension [1, 2, 3, 4], making them inadequate for tasks requiring fine-grained temporal grounding of events described by natural language. To address this capability gap, *Video Temporal Grounding* (VTG) has emerged as a pivotal research area. VTG involves localizing video segments that correspond specifically to given textual queries, enabling detailed interaction with video content. The core challenge of VTG lies in precisely aligning complex linguistic semantics with temporally distributed visual information, while simultaneously handling complex temporal relationships within the video.

As illustrated in Fig. 2, VTG encompasses several closely related but distinct tasks: (a) *Video Moment Retrieval* [5, 6], where the goal is to identify video segments matching natural language descriptions; (b) *Dense Video Captioning* [7, 8],

which requires generating temporally aligned captions for multiple events; (c) *Video Highlight Detection* [9, 10], aims at selecting segments most relevant to a given query; and (d) *Temporally Grounded Video Question Answering* [11, 12], which involves pinpointing the temporal evidence needed to accurately answer questions. Collectively, these tasks define the contemporary scope of VTG research and highlight the necessity of sophisticated temporal reasoning.

Although substantial progress has been achieved, early VTG methods based on traditional deep learning architectures [13, 14] continue to face significant limitations. These include challenges in bridging semantic gaps between visual and linguistic modalities, inadequate temporal context modeling, and limited generalization capabilities. Previous methods often relied on manually designed proposal-generation mechanisms [15, 16] or simple temporal boundary regression [17, 18], which lacked flexibility and interpretability. Recently, the advent of Large Language Models (LLMs) [19, 20, 21] and their multimodal variants, i.e., Multimodal Large Language Models (MLLMs) [22, 23, 24, 25], has dramatically reshaped the field of video-language understanding. These models provide powerful cross-modal reasoning, instruction-following capabilities, and robust zero-shot generalization, significantly enhancing the potential for effective VTG.

Motivated by these advancements, a rapidly growing research direction termed *VTG-MLLMs* [26, 27, 28] has emerged, leveraging MLLMs for temporal grounding tasks. The swift evolution of this subfield is visually charted in Fig. 1, which presents a chronological overview of representative VTG-MLLM approaches. Unlike traditional approaches relying solely on visual backbones with task-specific heads, VTG-MLLMs utilize general-purpose

• Jianlong Wu, Wei Liu, and Liqiang Nie are with the School of Computer Science and Technology, Harbin Institute of Technology, Shenzhen 518055, China (e-mail: wujianlong@hit.edu.cn; liuwei030224@gmail.com; nieliqiang@gmail.com).

• Meng Liu is with the School of Computer Science and Technology, Shandong Jianzhu University, Jinan 250101, China (e-mail: mengliu.sdu@gmail.com).

• Zhouchen Lin is with the School of Intelligence Science and Technology, Peking University, Beijing 100871, China (e-mail: zlin@pku.edu.cn).

• Ye Liu and Chang Wen Chen are with the Department of Computing, The Hong Kong Polytechnic University, Hong Kong, China (e-mail: coco.ye.liu@connect.polyu.hk; changwen.chen@polyu.edu.hk).

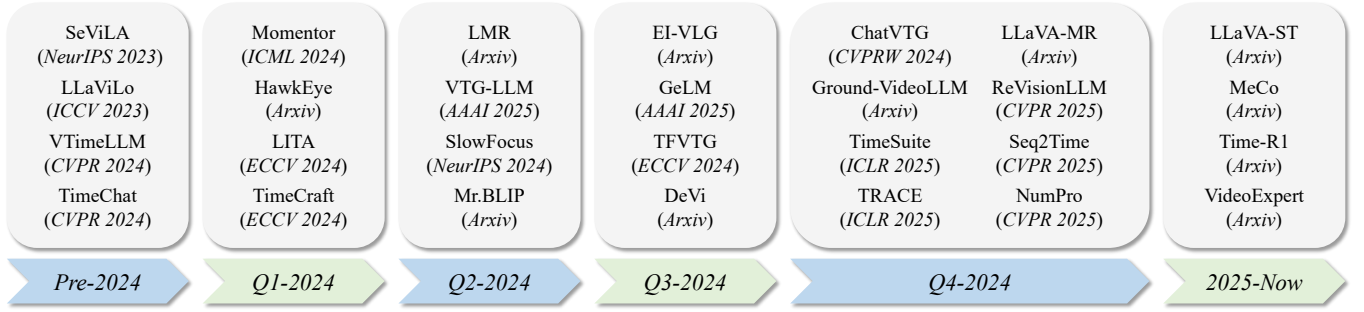


Fig. 1: An overview timeline of representative VTG-MLLMs. This timeline is organized according to the initial arXiv release dates of each work, with corresponding conference acceptance details included where applicable.

MLLMs to reason about temporal relationships, align semantics, and localize relevant video segments either directly or indirectly. VTG-MLLMs adopt diverse architectural strategies, with some methods employing MLLMs as high-level facilitators for semantic grounding [29, 30], and others using them explicitly for boundary prediction [31, 32]. Consequently, the VTG-MLLM field now encompasses varied architectural innovations, training paradigms, and representation techniques [33, 34, 35]. However, the rapid evolution and complexity of VTG-MLLM research present challenges for navigating the current literature. Existing surveys predominantly focus on general video-language modeling [36, 37, 38, 39, 40] or cover VTG from a pre-LLM perspective [41, 42, 43, 44], leaving a notable gap in systematic analyses of VTG in the LLM era.

To fill this gap, we present the first comprehensive survey exclusively dedicated to VTG-MLLMs, covering literature up to May 2025. This survey systematically organizes recent developments, identifies emerging technical trends, and outlines future research opportunities. Specifically, we introduce a structured three-dimensional taxonomy:

- **Functional Roles of MLLMs:** Classifying models based on whether MLLMs act as *Facilitators* assisting downstream grounding tasks, or as *Executors* directly predicting temporal boundaries.
- **Training Paradigms:** Distinguishing between *pretraining*, *fine-tuning*, and *training-free* approaches, each with unique trade-offs in generalization, task specialization, and supervision requirements.
- **Video Feature Processing Techniques:** Examining strategies for representing and integrating video inputs, including spatiotemporal tokenization and temporal modeling techniques.

Our taxonomy provides a progressive analytical framework, guiding readers from the high-level roles of MLLMs through training paradigms and down to video feature processing methods. By structuring our survey in this layered manner (see Fig. 3), we aim to offer clarity, advance comparative analyses, and identify underexplored avenues in VTG-MLLMs.

The remainder of this survey is organized as follows. Section 2 introduces the preliminaries of VTG-MLLMs, including an overview of VTG tasks and the foundational concepts behind MLLMs. Section 3 presents a detailed taxonomy of recent VTG-MLLM research, categorizing methods along three key dimensions: the *functional roles* of MLLMs,

training paradigms, and *video feature processing* strategies. These categories and their subtypes are illustrated in Fig. 3. Section 4 provides an overview of benchmark datasets, evaluation protocols, and comparative analysis of empirical results across existing VTG-MLLMs. Section 5 discusses open challenges and future research directions. Finally, Section 6 concludes the survey.

2 PRELIMINARIES

In this section, we provide an overview of core VTG tasks, as well as the foundational background on MLLMs.

2.1 Video Temporal Grounding

In this survey, we provide a comprehensive overview of four primary VTG tasks: video moment retrieval, dense video captioning, video highlight detection, and temporally grounded video question answering. In the following subsections, we briefly describe each of these tasks.

2.1.1 Video Moment Retrieval

Video moment retrieval (MR) [45, 46], also referred to as *temporal sentence grounding* [47, 48], *video moment localization* [49, 50] or *temporal video grounding* [27, 51], aims to identify and localize temporal segments within untrimmed videos based on natural language queries (see Fig. 2 (a)). This task represents the most direct and fundamental benchmark for evaluating the temporal grounding capabilities of VTG-MLLMs. It requires not only accurate alignment of textual descriptions with specific video segments, but also the ability to map video content to precise temporal boundaries, testing a model’s understanding of fine-grained temporal relationships and semantic coherence in untrimmed videos.

2.1.2 Dense Video Captioning

Dense video captioning (DC) [52, 53] involves generating detailed, temporally grounded descriptions for all significant events or actions occurring in an untrimmed video, along with their corresponding start and end timestamps (see Fig. 2 (b)). Unlike MR, which aims at localizing a single specific moment given a textual query, DC captures the complete narrative by identifying multiple events and their intricate temporal dependencies. This task assesses a model’s proficiency in comprehending extended temporal contexts and nuanced interactions within videos. Additionally, DC explicitly challenges models to manage overlapping events, a capability essential for achieving comprehensive fine-grained video understanding using large language models.

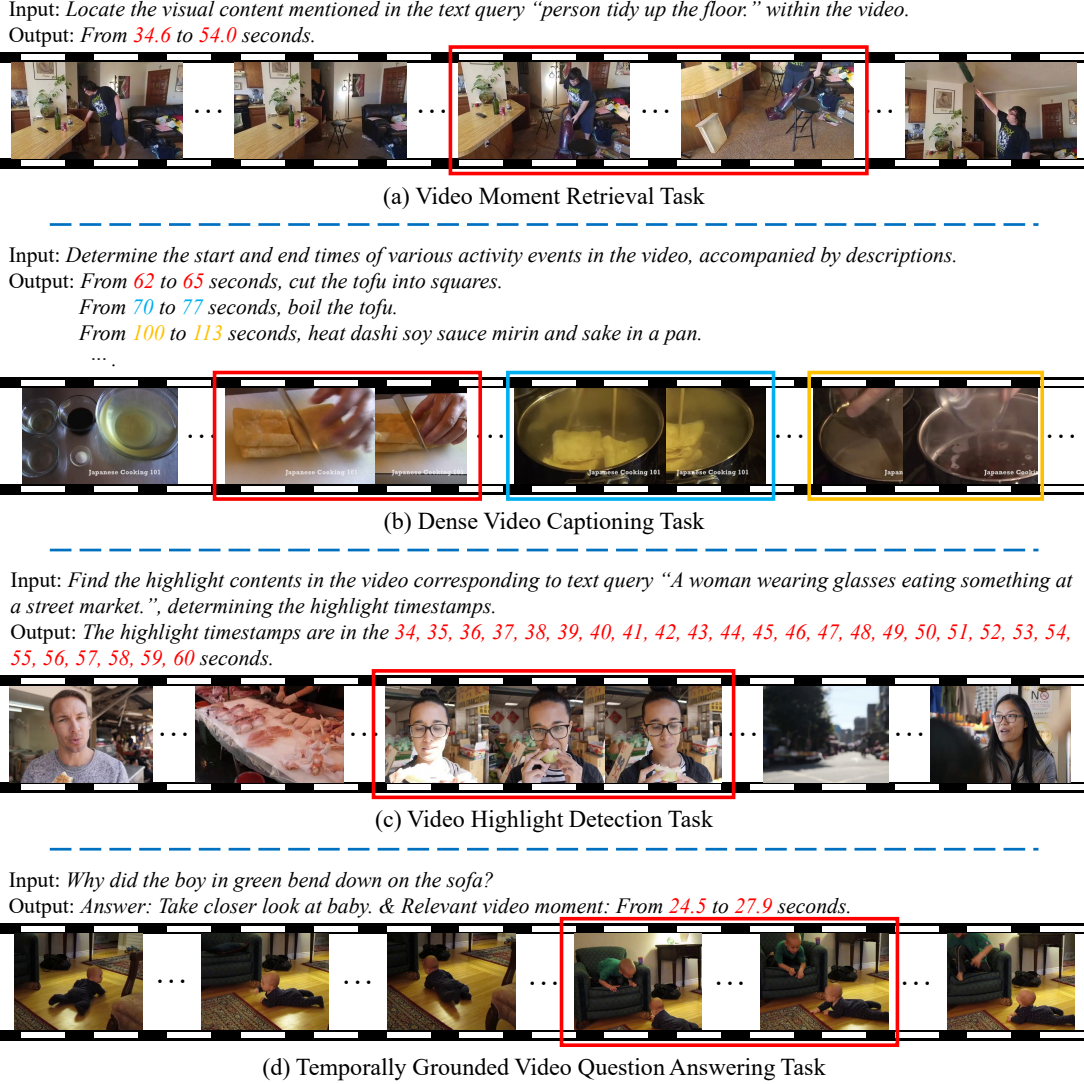


Fig. 2: Illustration of four core tasks in VTG: video moment retrieval, dense video captioning, video highlight detection, and temporally grounded video question answering. These tasks encompass a range of temporal reasoning challenges, including precise segment localization, multi-event description generation, highlight identification, and time-sensitive question answering, each requiring fine-grained temporal understanding.

2.1.3 Video Highlight Detection

Video highlight detection (HD) [9, 54] aims to identify keyframes or short segments within an untrimmed video that best match a given textual query, typically by assigning an importance or relevance score to these moments (see Fig. 2 (c)). Unlike MR and DC, which primarily operate at the event level, HD emphasizes frame-level precision. This task evaluates the ability of model to accurately pinpoint salient video clips that closely correspond to textual prompts and to assess their contextual significance. Such fine-grained alignment is essential for applications that require high temporal precision in identifying critical events.

2.1.4 Temporally Grounded Video Question Answering

Temporally grounded video question answering (GQA) [51, 55], also known as *grounded video QA* [11, 56] or *temporal video grounding of questions* [32], extends traditional video QA by requiring models to not only answer questions but also identify and localize the precise temporal intervals containing relevant visual evidence (see Fig. 2 (d)). Unlike

MR, GQA introduces the added complexity of integrating temporal localization with multimodal reasoning. This task is particularly critical for developing explainable video QA systems, as it demands explicit and interpretable connections between textual answers and visual evidence within the video content.

2.2 Multimodal Large Language Models

MLLMs [57, 58, 59] extend traditional LLMs by integrating multimodal encoders, such as image encoders [60, 61, 62, 63], video encoders [64, 65, 66, 67], and specialized cross-modal adapters [23, 68, 69]. Taking Video-LLM as an illustrative example, a video encoder processes a sequence of downsampled frames V , converting them into visual tokens $F^v = E_v(V)$. These visual tokens are then projected by a visual adapter to align with the embedding space of the language model, yielding aligned visual tokens $X^v = Q(F^v)$. Concurrently, an input textual query q , which may include instructions, prompts, or other textual elements, is encoded into linguistic tokens X^t via a textual encoder. The visual

and textual tokens are concatenated into a unified input sequence $[X^v, X^t]$, which is subsequently processed by the LLM to generate the appropriate inference.

Current research efforts on MLLMs emphasize maximizing the efficiency of leveraging LLMs’ advanced capabilities. Initial studies predominantly focus on designing cross-modality adapters aimed at mapping features from non-linguistic modalities into the semantic embedding space of language models. Flamingo [70], as a pioneering model, integrates visual and linguistic modalities through a gated cross-attention mechanism. Subsequently, various approaches, such as BLIP [71], mPLUG [72], and LanguageBind [73], adopt the Q-Former architecture to align visual representations, whereas models in the LLaVA series [22, 74, 75] introduce multilayer perceptrons (MLPs) as simpler yet effective connectors for modality integration. Additionally, more recent works have proposed lightweight and efficient alignment modules [76, 77, 78], continuing to enhance performance and model compactness.

Alongside architectural developments, training strategies constitute another critical research direction. Researchers have compiled large-scale multimodal pretraining datasets [79, 80, 81] to facilitate robust and diverse representation learning. Instruction-tuned datasets [82, 83] and specialized datasets tailored for chain-of-thought (CoT) reasoning [84, 85, 86] have been developed to improve task comprehension and generalization capabilities of MLLMs. Furthermore, parameter-efficient fine-tuning methods such as LoRA [87], LISA [88], and DoRA [89] have emerged, enabling efficient task-specific or domain-specific adaptation without extensive retraining.

3 A MULTI-FACETED TAXONOMY OF VTG-MLLMS

As established in the Introduction (Section 1), we employ a three-dimensional taxonomy to deconstruct the complexities of VTG-MLLMs. This section delves into the specifics of this classification. Our taxonomy (visualized in Fig. 3), which progresses from high-level architectural considerations to fine-grained processing techniques, will be explored in detail through the following dimensions:

- **The Functional Roles of MLLMs** (Section 3.1): We will analyze how the architectural positioning of MLLMs—whether as *Facilitators* aiding downstream tasks or as *Executors* directly undertaking temporal prediction—shapes their overall design and impact on temporal perception.
- **The Training Paradigms** (Section 3.2): This subsection will differentiate among the *pretraining*, *fine-tuning*, and *training-free* paradigms. The analysis will center on the inherent trade-offs each strategy presents in terms of generalization capability, task-specific adaptation, and overall resource demands.
- **The Video Feature Processing Techniques** (Section 3.3): Here, we will systematically examine the diverse methodologies for representing and integrating video inputs. This includes a closer look at spatiotemporal tokenization mechanisms within the token budget, and various temporal modeling approaches that enable models to effectively process and reason about dynamic visual content.

This structured examination will provide the foundation for our detailed review of specific methodologies and trends within the VTG-MLLM field in the sections that follow.

3.1 Functional Roles of MLLMs in VTG-MLLMs

The functional role of MLLMs characterizes their architectural integration within VTG pipelines, determining whether they function primarily as auxiliary modules facilitating cross-modal understanding or as central reasoning engines directly conducting temporal grounding. Accordingly, existing VTG-MLLMs can be categorized into two paradigms: 1) *Facilitators*, where MLLMs generate structured textual representations from video content to support downstream modules; and 2) *Executors*, where MLLMs directly perform temporal boundary prediction via integrated multimodal reasoning.

3.1.1 Facilitators

In the Facilitator role, MLLMs act as intermediaries by transforming complex video data into structured textual forms, as depicted in Fig. 4 (a). We formalize this process as a conditional generation problem:

$$T = \text{MLLM}_{\text{facilitator}}(V), \quad (1)$$

where V denotes the input video and T represents the generated textual descriptions. The generated textual outputs can either directly facilitate dataset construction or serve as semantic aids within dedicated downstream modules. Two primary application areas arise under this paradigm: *Dataset Construction* and *Expert Module Integration*.

Dataset Construction: MLLMs are extensively utilized to synthesize textual annotations, significantly enhancing the efficiency of dataset creation and expansion for model training and evaluation. For instance, Di and Xie [91] leverages Llama2 [128] to transform timestamped narrations from Ego4D [129] into temporally grounded QA pairs. Similarly, GPT-4o [130] was employed by Bao et al. [100] in VidMorp to automatically generate pseudo-labeled sentences aligned with video frames. Other works [95, 97] have similarly utilized advanced models like BLIP-2 [71], LLaVA [22], and Gemini-1.5 [131] to automate the annotation process and enrich datasets for VTG tasks.

Expert Module Integration: Beyond dataset generation, textual outputs from MLLMs can serve directly within VTG systems, either as semantic inputs in similarity-based grounding methods or as additional signals enhancing visual representations through cross-modal integration. For instance, Qu et al. [29] employs Video-ChatGPT [132] to generate multi-granularity clip captions, facilitating iterative query matching using Sentence-BERT [133]. Similarly, Xu et al. [94] utilizes MiniGPT-v2 [134] for caption generation and Baichuan2 [135] for query rewriting, reducing linguistic biases. Additionally, Cai et al. [90] leverages LLaVA-1.5 [74] to generate paragraph-level narrations, aligning them temporally with video features via cross-attention mechanisms [14, 17], thereby enhancing contextual understanding and robustness.

Summary: The Facilitator framework is advantageous due to its computational efficiency, ease of deployment, and inherent scalability, requiring minimal adaptation of pretrained MLLMs. However, reliance on static pretrained

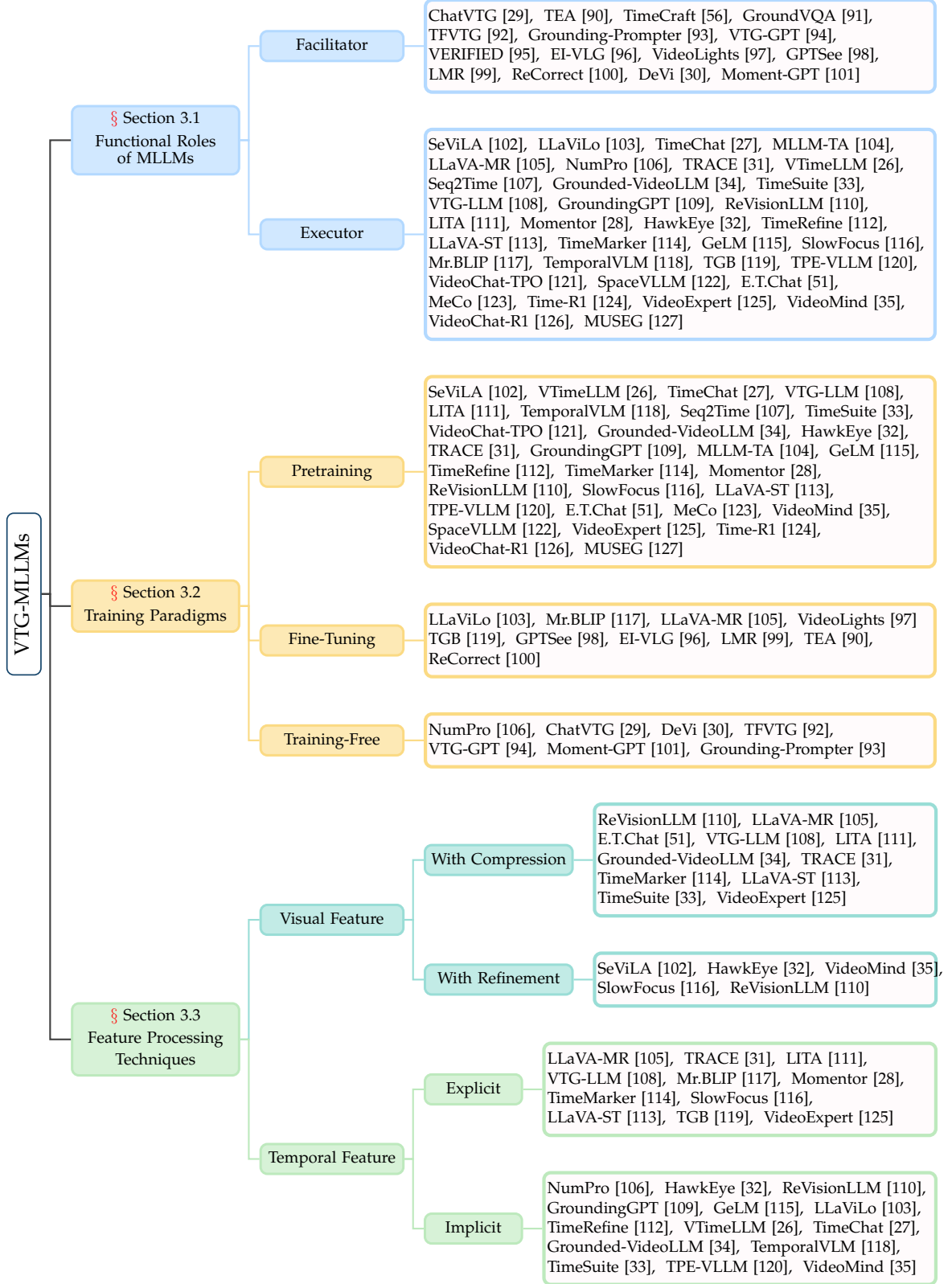


Fig. 3: Taxonomy of VTG-MLLMs, encompassing three primary dimensions: *Functional Roles of MLLMs* (Facilitators / Executors), *Training Paradigms* (Pretraining / Fine-Tuning / Training-Free), and *Feature Processing Techniques* (Visual Features / Temporal Features), each with distinct sub-categories reflecting the diverse strategies employed in this field.

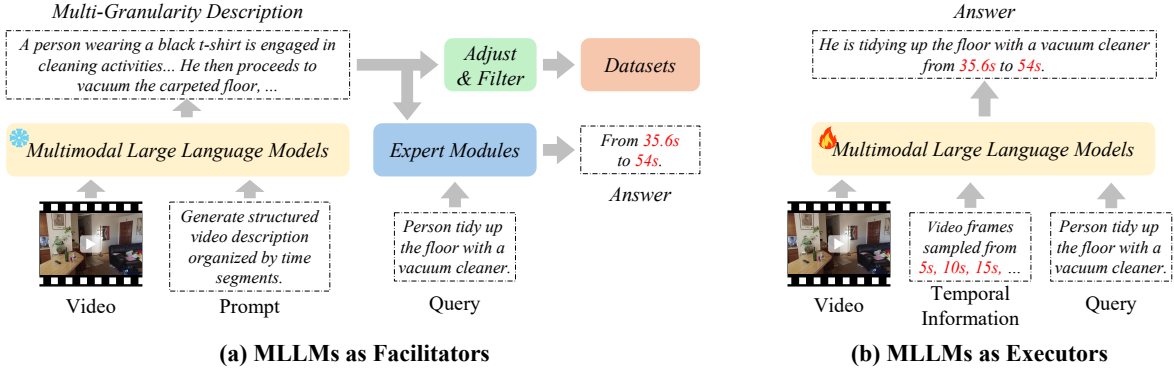


Fig. 4: Visualization of the distinct functional roles of MLLMs in VTG-MLLMs. The left panel depicts MLLMs as *Facilitators*, where they generate textual descriptions from video inputs, serving either as training data for downstream tasks or as auxiliary signals within expert modules. The right panel illustrates MLLMs as *Executors*, directly processing queries, video content, and temporal information (e.g., timestamps) to produce grounded outputs through a sequence-to-sequence prediction framework.

models carries limitations such as propagation of inherent biases [136] and constraints from original training data, potentially affecting the reliability of textual outputs and downstream performance. Furthermore, the fixed nature of off-the-shelf MLLMs inherently limits their capability for complex temporal reasoning, creating performance bottlenecks that sophisticated pipeline designs may not fully address. Nonetheless, as *Facilitators*, MLLMs remain valuable for efficient dataset curation, providing abundant task-specific data crucial for advancing VTG research.

3.1.2 Executors

When functioning as *Executors*, MLLMs directly perform the core tasks of VTG, formulating the problem as an end-to-end sequence-to-sequence (seq-to-seq) prediction challenge. In this setting, illustrated in Fig. 4 (b), the model jointly consumes raw video input and task-specific textual prompts to generate a temporally aligned output:

$$Y = \text{MLLM}_{\text{executor}}(V, Q, \tau), \quad (2)$$

where V denotes the input video stream, Q is the textual query, and τ optionally encapsulates temporal priors. The output Y represents predicted answers, timestamps, or task-specific tokens aligned with the video timeline.

This paradigm holds the potential to unify video understanding within a generative framework. However, it faces a significant hurdle: standard MLLMs [59, 132] often struggle to capture fine-grained temporal dependencies. This challenge largely stems from their vision encoders treating video as an unordered “bag-of-features,” which discards the crucial sequential information essential for precise event localization. To address these challenges, research has focused on two primary directions: *Architectural Enhancement* and *Training Optimization*.

Architectural Enhancement: Architectural innovations are designed to improve the temporal perception and reasoning capabilities of MLLMs. These enhancements typically involve either modifying the input feature processing pipeline or integrating temporal awareness directly into the LLM’s internal structure.

One line of research focuses on enhancing video feature representations before they are processed by the LLM. For instance, the method in Momentor [28] injects explicit

temporal position encodings into frame-level features to improve temporal localization. Grounded-VideoLLM [34] adopts a dual-stream architecture to separately capture spatial and temporal dynamics, while LLaVA-MR [105] introduces components for reducing redundancy and emphasizing critical dynamic moments. These strategies aim to provide the LLM with richer temporal context, laying a stronger foundation for subsequent reasoning.

A complementary approach modifies the internal architecture or output mechanisms of LLMs to better handle temporal cues. For example, GeLM [115] incorporates flexible grounding tokens for temporal evidence retrieval, and TRACE [31] adds task-specific decoding heads for structured temporal output. VideoExpert [125] integrates parallel reasoning and generation modules for specialized processing. Similarly, VideoMind [35] decomposes complex tasks into specialized roles, like a *Grounder* with a timestamp-decoder and a *Planner* to coordinate them, using a Chain-of-LoRA strategy for seamless collaboration. While these approaches can significantly enhance temporal understanding, they may introduce trade-offs, such as increased computational overhead or reduced general-purpose flexibility.

Training Optimization: Effective optimization strategies are critical for equipping MLLMs with robust temporal understanding, transforming them into capable *Executors*. These strategies typically form a holistic pipeline combining novel training curricula with temporal tasks and datasets.

A prevalent approach involves multi-stage training frameworks, as seen in VTimeLLM [26] and SlowFocus [116], which progressively refine the model’s temporal localization abilities. Another emerging direction is optimization via reinforcement learning (RL). Time-R1 [124] adapts RL-based strategies specifically for temporal reasoning, and VideoChat-R1 [126] explores the effectiveness of Group Relative Policy Optimization (GRPO) [137].

To further strengthen temporal perception, many methods incorporate explicit reasoning tasks. For example, TPE-VLLM [120] introduces novel training objectives targeting boundary detection and duration reasoning, improving its handling of complex temporal relationships.

Underpinning all these strategies is the reliance on high-quality, time-annotated datasets. For instance, the TimeIT dataset from TimeChat [27] provides rich timestamp anno-

tations essential for instruction tuning, while TimeSuite [33] offers a unified collection of diverse datasets to facilitate more comprehensive temporal learning.

Summary: The Executor paradigm represents a pivotal shift towards unified, end-to-end temporal grounding, allowing MLLMs to simultaneously process video content and textual queries within a tightly coupled seq-to-seq framework. This approach supports flexible input and output formats, capturing complex visual-textual correlations without relying on modular, cascaded architectures. However, this flexibility comes at a cost, often requiring extensive annotated datasets, significant computational resources, and complex training procedures. Despite these challenges, the Executor approach remains a promising direction for advancing fine-grained video understanding, with the potential to fundamentally reshape the field by integrating deeper temporal reasoning into multimodal models.

3.2 Training Paradigms of VTG-MLLMs

Building on the functional differentiation of MLLMs, this subsection examines the training paradigms used to adapt these models for effective video temporal grounding. The choice of training approach reflects not only the system’s design goals, i.e., whether to build a domain-generalist or task-specific model, but also the trade-offs in supervision, resource efficiency, and scalability. We categorize current VTG-MLLM approaches into three main paradigms: *Pre-training*, *Fine-Tuning*, and *Training-Free* pipelines.

3.2.1 Pretraining VTG-MLLMs

Pretraining in VTG-MLLMs aims to equip models with robust temporal reasoning capabilities via large-scale supervised learning. Fundamentally, like most generative multimodal approaches, the pretraining process involves training the model to generate a target output T , typically encompassing temporal annotations such as event order, timestamps, or durations, conditioned on an input video V and textual prompt P . Formally, this objective is represented as minimizing the pretraining loss over a dataset \mathcal{D} :

$$\mathcal{L}_{\text{pretrain}} = \sum_{(V,T) \in \mathcal{D}} \mathcal{L}_{\text{gen}}(T | V, P). \quad (3)$$

However, unlike general-purpose MLLMs, the primary innovation in pretraining-based VTG-MLLMs lies not in architectural design but in the creation of sophisticated training strategies and specialized pretraining datasets that tailor the optimization of Eqn. (3) for complex temporal understanding. We focus on two critical aspects of this approach: *Prevalent Pretraining Strategies* and *High-Quality Temporal Datasets*.

Prevalent Pretraining Strategies: The cornerstone of VTG pretraining is the multi-stage, progressive supervised learning pipeline. As introduced in the context of Executor models (Section 3.1.2), this strategy is founded on the principle of incremental learning, guiding the model from coarse-grained understanding to fine-grained localization. For instance, VTimeLLM [26] exemplifies this with its boundary-aware three-stage process, which sequentially tackles feature alignment, instruction tuning, and precise boundary optimization. Similarly, SlowFocus [116] integrates mixed-frequency sampling in its final training stages to enhance

temporal resolution. This multi-stage philosophy has become a de facto standard, with numerous other models like TimeMarker [114], GroundingGPT [109], and LLaVA-ST [113] adopting similar hierarchical frameworks to progressively refine temporal perception.

While multi-stage supervision remains the mainstream, a new wave of research is exploring innovative refinements, aiming for greater precision and efficiency. A particularly prominent direction is the application of RL. Inspired by the success of RL and techniques like GRPO [137] in complex reasoning domains such as code generation and mathematics, researchers have begun adapting these methods for VTG. RL allows for the direct optimization of task-specific metrics like IoU, achieved by designing a composite reward function that encourages both a structured reasoning process and high prediction accuracy. Pioneering this direction, Time-R1 [124] introduces a reasoning-guided framework with a novel reward mechanism; VideoChat-R1 [126] further offers a systematic exploration with GRPO, while MUSEG [127] addresses the single-segment limitation by enabling reasoning over multiple distributed events.

Beyond these, other novel strategies are also emerging. For example, Seq2Time [107] adopts a data-centric strategy, synthesizing sequential training data with self-generated temporal cues, while TimeRefine [112] reformulates temporal grounding as an iterative refinement task, allowing model to self-improve its localization accuracy.

High-Quality Temporal Datasets: High-quality, temporally annotated multimodal datasets are crucial for pretraining and instruction-tuning VTG models, providing the diverse contexts required for robust generalization. Building upon early efforts like TimeIT [27] and VTimeLLM [26], subsequent work has evolved along several key directions. One major line of work aims to enhance data scale and diversity for more effective pretraining. For instance, VTG-IT-120K [108] expands on TimeIT by incorporating annotations from YT-Temporal-180M [138], InternVid-G [32] enriches InternVid10M-FLT [139] with segment-level captions and hard negative samples for more precise grounding, and Vid-Morp [100] leverages pseudo-labeling on real-world videos to scale data creation. A parallel direction develops specialized instruction-tuning datasets to align models with complex temporal reasoning. Moment-10M [28], sampled from YT-Temporal-1B [140], was designed for this purpose, while E.T. Instruct 164K [51] complements this by providing a curated dataset across nine distinct tasks, specifically tailored for multi-event and time-sensitive scenarios. More recently, frontier datasets have begun to broaden the scope of VTG by integrating spatial dimensions [122] or advancing spatiotemporal understanding [113]. Collectively, these datasets, along with numerous other contributions [33, 34, 124], underscore a clear trajectory toward building more comprehensive and fine-grained data resources to tackle the full spectrum of temporal reasoning challenges.

Summary: The pretraining paradigm empowers MLLMs with robust temporal grounding capabilities, supporting generalization across diverse downstream VTG tasks. However, this approach also presents significant challenges, including the high computational cost of training and the substantial effort required to construct large, high-quality temporal datasets. Effective pretraining strategies must care-

fully balance task complexity and learning progression to maximize temporal understanding, making this a critical area for ongoing research.

3.2.2 Fine-Tuning VTG-MLLMs

In contrast to the resource-intensive pretraining paradigm (Section 3.2.1), fine-tuning VTG-MLLMs provides a more computationally efficient approach, requiring smaller, task-specific datasets. Research within this paradigm can be broadly divided into two main directions: *Direct Fine-Tuning of MLLMs* and *Offline Textualization with MLLMs*, closely aligning with the functional roles discussed in Section 3.1.

Direct Fine-Tuning of MLLMs: This approach directly fine-tunes general-purpose pre-trained MLLMs while maintaining their original architecture, adapting them for VTG tasks through task-specific training objectives. These methods typically reframe VTG as a seq-to-seq prediction problem, leveraging the contextual understanding capabilities already embedded in the models.

For instance, SeViLA [102] adapts BLIP-2 [71] into two interconnected components—a localizer and an answerer—strategically linking the outputs of the localizer to guide the answerer, thereby enhancing temporal precision. Similarly, LLaViLo [103] incorporates lightweight adapters to integrate video-text features, utilizing a multi-objective loss function for more refined temporal grounding. TGB [119] introduces a unique approach by leveraging CNN-extracted optical flow features as low-dimensional motion cues, enhancing temporal awareness without substantially increasing model complexity. Additionally, recent innovations have focused on optimizing BLIP-2 for fine-grained temporal understanding. Notably, Mr.BLIP [117] explores novel multimodal input sequences to improve temporal understanding of events, while LLaVA-MR [105] introduces a dynamic token compression strategy, reducing redundancy in spatiotemporal features and capturing more fine-grained event cues.

Despite their efficiency, these methods face the *catastrophic forgetting* challenge [141]. As the models adapt to fine-tuning data, they often lose the general-purpose capabilities acquired during pretraining, leading to performance degradation on broader video understanding tasks.

Offline Textualization with MLLMs: An alternative fine-tuning strategy employs MLLMs in a static capacity to convert raw video inputs into textual descriptions, which then guide downstream modules. This approach effectively bridges the gap between unstructured visual data and language-conditioned learning tasks, often integrating components from traditional VTG methods.

For example, GPTSee [98] generates detailed video descriptions, subsequently matched with textual queries to support moment localization. EI-VLG [96] incorporates these descriptions as environmental cues in a contrastive learning framework, refining the temporal precision of candidate segments. To address redundancy, LMR [99] uses cross-attention to highlight query-relevant segments, improving contextual alignment. Similarly, TEA [90] integrates these textual outputs with visual features to enhance semantic discriminability and temporal precision.

While this approach can significantly improve temporal grounding accuracy, it inherits certain limitations from tra-

ditional VTG methods. For instance, early two-stage matching approaches [142, 143] often rely on pre-defined segment boundaries, limiting global context modeling. Meanwhile, direct regression methods [144, 145] can misinterpret visually similar segments as semantically identical due to attention biases. In contrast, fine-tuned MLLMs provide richer, more context-aware embeddings, aligning visually similar events with distinct textual semantics and improving overall robustness.

Summary: Fine-tuning VTG-MLLMs offers a practical compromise between pretraining and fully training-free approaches, significantly reducing computational overhead while enhancing task-specific temporal alignment. However, these methods are inherently task-optimized, limiting their generalization across broader video understanding domains. As a result, fine-tuned models often excel in narrow, well-defined tasks but struggle with broader generalization—a critical challenge for future research.

3.2.3 Training-Free VTG-MLLMs

Training-free approaches represent a rapidly emerging paradigm in VTG-MLLMs, notable for their low computational overhead and zero-shot nature that eliminates the need for labeled supervision. These methods bypass the need for end-to-end training by leveraging pre-trained foundation models (e.g., MLLMs and LLMs) and specialized expert tools [71, 133, 146], enabling temporal grounding through purely inference-based pipelines. While they share architectural similarities with fine-tuning approaches (Section 3.2.2), training-free methods distinguish themselves by substituting trainable components with off-the-shelf models, significantly reducing the need for task-specific parameter updates. Current training-free VTG-MLLMs generally adopt one of two principal strategies: *Feature Similarity Matching* and *LLM-Driven Reasoning*, distinguished by their approach to leveraging MLLM-generated textualizations for temporal localization.

Feature Similarity Matching: This strategy relies on extracting high-dimensional semantic representations from both natural language queries and textualized video content using pre-trained encoders. Temporal grounding is then achieved by identifying the video span s^* that maximizes a similarity score with the query \mathbf{q} , formalized as:

$$s^* = \arg \max_{s_i \in \mathcal{S}} \text{sim}(E_Q(\mathbf{q}), E_V(\mathbf{v}_{s_i})), \quad (4)$$

where $\mathcal{S} = \{s_1, s_2, \dots, s_N\}$ denotes the set of candidate video spans, $E_Q(\cdot)$ and $E_V(\cdot)$ are frozen encoders for the query and video, respectively, and $\text{sim}(\cdot, \cdot)$ denotes a similarity function such as cosine similarity.

For instance, Moment-GPT [101] extends the VTG-GPT [94] framework by combining frame-level captions from MiniGPT-v2 [134] with segment-level captions from Video-ChatGPT [132], matching them to textual queries using a similarity-based retrieval approach. TFGVTG [92] further refines this paradigm by decomposing complex queries into sub-events using an LLM, followed by segment matching through the BLIP-2 Q-Former [71]. The final predictions are derived by integrating localized spans via temporally-aware filtering that accounts for the sub-event order and relations, improving the accuracy of temporal localization.

LLM-Driven Reasoning: This alternative strategy treats VTG as a high-level textual inference task, leveraging the reasoning capabilities of LLMs to comprehend and localize temporal segments based on enriched video descriptions.

For example, Grounding-prompter [93] reformulates VTG as a long-text comprehension task, aligning speech transcriptions and visual captions with timestamp annotations. It employs a four-step multiscale denoising chain-of-thought approach, progressively refining coarse temporal predictions through iterative prompts. In a similar vein, DeVl [30] uses Video-LLaVA [58] to perform hierarchical, multi-scale captioning, followed by query-driven refinement using GPT-4o [130]. This multi-stage reasoning process allows the model to better capture event dependencies and fine-grained temporal structures, leading to more accurate localization without additional training.

Beyond these dominant strategies, emerging methods are exploring novel approaches to training-free temporal grounding. For instance, NumPro [106] introduces a unique numbering scheme inspired by manga panel sequencing, inserting numerical identifiers into each video frame to enhance temporal traceability. This subtle form of visual embedding enables LLMs to track frame sequences more effectively without modifying the model architecture or requiring fine-tuning, preserving general video comprehension while improving temporal precision.

Summary: Training-free VTG-MLLMs provide a lightweight, modular alternative to conventional fine-tuning approaches, effectively decomposing VTG into manageable subtasks, i.e., captioning, matching, and reasoning, without the computational overhead of extensive training. By leveraging powerful pre-trained models and off-the-shelf components, these methods reduce the cost and complexity of domain-specific adaptation, making them a compelling choice for scenarios where data availability and computational resources are limited. However, their reliance on predefined embeddings and static representations can introduce challenges in capturing fine-grained temporal dependencies, presenting an ongoing area for innovation.

3.3 Video Feature Processing in VTG-MLLMs

At the most fine-grained level of our taxonomy, we examine the video feature extraction strategies that underpin VTG-MLLMs. As discussed in the taxonomy of functional roles, Facilitator-based methods often rely on pre-trained, frozen modules to provide high-level video embeddings. In contrast, Executor-oriented designs require more sophisticated mechanisms to handle raw video inputs, reflecting their more direct involvement in temporal reasoning and event localization. This subsection focuses on the critical strategies for extracting and processing visual and temporal features within Executor-based VTG-MLLMs.

3.3.1 Efficient Visual Feature Handling

Given the dense nature of frame-level information and the constrained input size of most LLMs, efficient visual feature handling is essential for capturing fine-grained temporal cues without overwhelming the model. These techniques can be broadly categorized into three main approaches, as illustrated in Fig. 5: *Learnable Token-Based Compression*, *Pooling-Based Compression*, and *Coarse-to-Fine Progressive Refinement*.

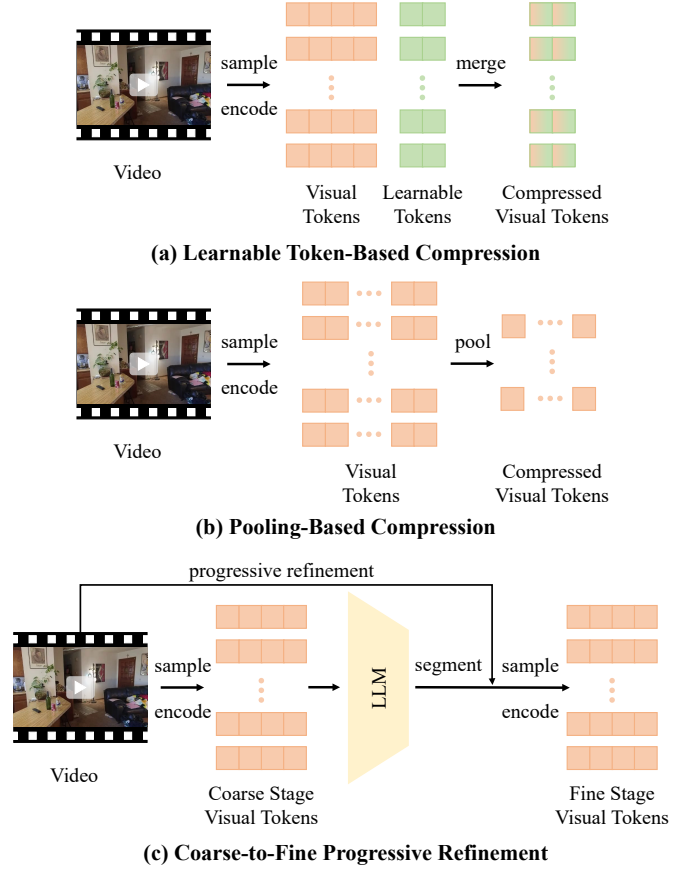


Fig. 5: Representative strategies for maximizing information utility within limited input token budgets: (a) Learnable Token-Based Compression, (b) Pooling-Based Compression, and (c) Coarse-to-Fine Progressive Refinement.

Learnable Token-Based Compression: This strategy employs learnable tokens to compress high-dimensional visual features into a concise and manageable representation, as depicted in Fig. 5 (a). Formally, given a set of trainable tokens \mathbf{Q} and raw visual features \mathbf{V} , the compressed representation \mathbf{C} is obtained by a parameterized function f_{compress} :

$$\mathbf{C} = f_{\text{compress}}(\mathbf{Q}, \mathbf{V}; \theta) \quad (5)$$

where θ represents learnable parameters. For example, VTG-LLM [108] introduces Slot-Based Token Compression, where a set of learnable slot embeddings aggregate information from raw visual tokens based on similarity. TRACE [31] adopts a similar strategy, compressing dense visual inputs into compact summaries using slot embeddings. ReVision-LLM [110] employs a [CLS]-like token [147] to aggregate segment features through self-attention, providing compact yet semantically rich representations, effectively acting as another instantiation of Eqn. (5).

Pooling-Based Compression: Pooling techniques aggregate local or global visual features to reduce dimensionality yet retain key semantic information, as shown in Fig. 5 (b). For instance, LITA [111] applies multi-granularity pooling across spatial and temporal dimensions, while Grounded-VideoLLM [34] and TimeMarker [114] use dynamically adjustable pooling kernels to capture hierarchical visual cues. LLaVA-MR [105] introduces token variance-based selection, prioritizing high-variance tokens to capture dynamic content more effectively. Alternatively, TimeSuite [33] reduces

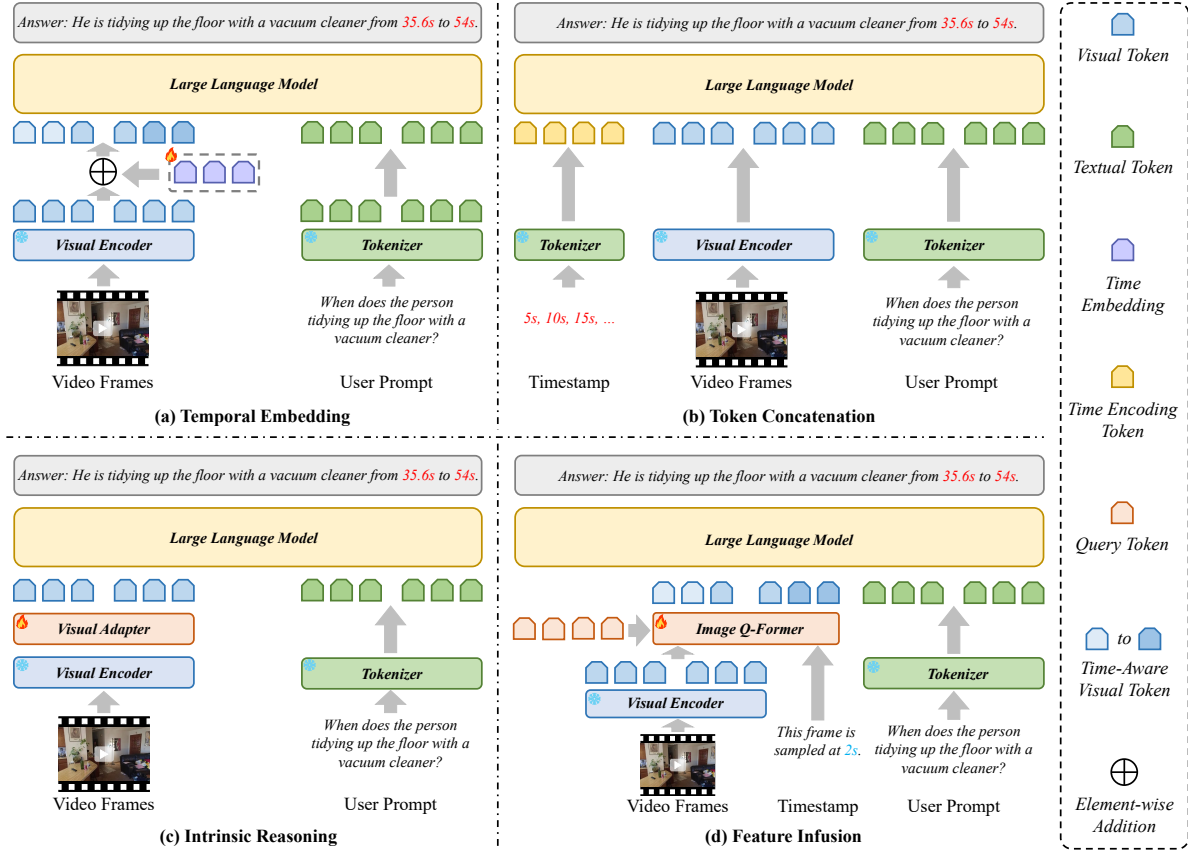


Fig. 6: Illustration of temporal feature processing strategies in VTG-MLLMs, categorized into (a)-(b) *Explicit Modeling*, which directly integrates timestamp information through methods like temporal embedding and token concatenation, and (c)-(d) *Implicit Modeling*, which relies on the reasoning capabilities of LLMs to infer temporal relationships through intrinsic reasoning and feature infusion.

token overhead through token shuffling and projection, achieving compression without extra parameters.

Coarse-to-Fine Progressive Refinement: Rather than compressing dense frame-level features upfront, these methods progressively refine temporal predictions to improve efficiency under strict token constraints, as outlined in Fig. 5 (c). SeViLA [102] is an early example that adopts a coarse-to-fine localization strategy by selecting language-aware keyframes before answer prediction. Similarly, HawkEye [32] adopts a recursive grounding approach, narrowing the temporal search space in iterative stages. Another prominent strategy involves a multi-stage refinement—first identifying a coarse temporal segment and then adjusting its boundaries—a method effectively employed by ReVisionLLM [110], SlowFocus [116], and VideoMind [35] to reduce token overhead while preserving accuracy.

Summary: Efficient visual feature handling is paramount in VTG-MLLMs to reconcile the richness of dense video data with the input token limitations. The above-discussed strategies represent distinct philosophies for information reduction. Token-based and pooling methods achieve upfront compression. In contrast, progressive refinement adopts an iterative approach, selectively focusing computational resources on temporally relevant segments. Collectively, these techniques are crucial for enabling MLLMs to process detailed video sequences, striking a vital balance between capturing fine-grained temporal nuances necessary for precise grounding and maintaining computational tractability.

3.3.2 Temporal Representation and Modeling

Unlike global video understanding tasks, fine-grained temporal grounding demands precise reasoning about temporal relationships to align video frames with timestamp intervals. This requirement is critical for VTG tasks such as moment retrieval and dense video captioning, which rely on accurate temporal boundary predictions. Robust timestamp representation mechanisms are essential for achieving this level of temporal precision. To address these challenges, temporal feature modeling in VTG-MLLMs can be broadly categorized into *Explicit* and *Implicit* modeling strategies, distinguished by whether temporal cues are directly injected into the model’s input stream or contextually assimilated through its architectural design and reasoning capabilities.

Explicit Modeling: Explicit modeling strategies directly embed temporal information into the input or feature representations of MLLMs, as illustrated in Fig. 6 (a)-(b). These approaches aim to provide precise temporal context by incorporating explicit time markers within the visual feature space, enhancing the model’s ability to align video frames with timestamps. Broadly, these methods can be categorized into two main approaches: *Temporal Embedding* and *Token Concatenation*, each with distinct mechanisms for integrating temporal cues.

Temporal Embedding. One common approach involves augmenting visual tokens with dedicated temporal embeddings, effectively integrating time information into the

sequence of input tokens. If \mathbf{v}_i represents the i -th visual token and \mathbf{e}_{t_i} is its corresponding temporal embedding, the augmented visual token \mathbf{v}_i' can be formed as:

$$\mathbf{v}_i' = \mathbf{v}_i + \mathbf{e}_{t_i}. \quad (6)$$

For instance, VTG-LLM [108] introduces learnable absolute time embeddings initialized to zero, preserving the original semantic integrity of visual tokens generated by pre-trained encoders. In contrast, LITA [111] adopts a relative time representation, segmenting videos into T equal-length chunks and assigning unique temporal tokens (e.g., $\langle 1 \rangle$ to $\langle T \rangle$) to each segment, providing a coarse but computationally efficient temporal structure. Momentor [28] takes a more granular approach by defining N learnable anchor points, each representing a specific temporal position within the video. These anchors define a continuous temporal feature space, allowing for more precise localization through interpolation. Other methods also leverage explicit time encodings. For example, TGB [119], SlowFocus [116], and LLaVA-ST [113] incorporate temporal position embeddings to enhance temporal awareness during fine-tuning. These embeddings provide clear temporal context, improving temporal alignment without notably altering the underlying architecture of MLLMs.

Token Concatenation. An alternative approach involves tokenizing timestamps directly from sampled frames, integrating these temporal markers with visual and textual tokens to form a unified input sequence. If \mathbf{S}_P , \mathbf{S}_V , and \mathbf{S}_T represent sequences of prompt tokens, visual tokens, and tokenized timestamps, respectively, the final input sequence \mathbf{S}_{input} fed to the MLLM can be a concatenation:

$$\mathbf{S}_{input} = \text{Concat}(\mathbf{S}_P, \mathbf{S}_V, \mathbf{S}_T). \quad (7)$$

The specific order and interleaving strategy can vary. For instance, LLaVA-MR [105] dynamically selects relative frame indices or absolute timestamps based on the frame sampling rate, interleaving these markers with special tokens like $\langle time_begin \rangle$ and $\langle time_end \rangle$ to denote temporal boundaries. TimeMarker [114] adopts a similar strategy, inserting explicit temporal separators (e.g., “second2.0”) into the input sequence to enhance temporal context. Additionally, Mr.BLIP [117] systematically explores various design choices for time representation, including relative versus absolute time, decimal versus integer formats, and different token ordering schemes, evaluating their impact on model performance. TRACE [31] further extends this approach by integrating temporal tokenization with visual feature embeddings, providing a tightly coupled representation of spatial and temporal information.

Implicit Modeling: Implicit modeling strategies aim to capture temporal relationships within video data through latent representations, leveraging the inherent reasoning and contextual understanding capabilities of large language models. Unlike explicit methods, which directly associate timestamps with visual inputs, implicit approaches integrate temporal cues more fluidly, embedding temporal knowledge without requiring explicit time markers. These strategies generally fall into two main categories: *Intrinsic Reasoning* and *Feature Infusion*, each employing distinct techniques to embed temporal context into visual representations, as illustrated in Fig. 6 (c)-(d).

Intrinsic Reasoning. This approach relies on the LLM’s

inherent ability to infer temporal relationships indirectly from the interplay between visual features and time-related language prompts. Rather than embedding explicit timestamps, these methods leverage numerical cues, iterative refinement, and boundary-aware reasoning to capture temporal dynamics. For instance, NumPro [106] introduces numerical indices directly into video frames, allowing the LLM to infer sequence order through positional awareness. Grounded-VideoLLM [34] adopts a similar strategy, introducing specialized temporal tokens into the LLM’s vocabulary, enabling unified modeling of time and semantics. TimeRefine [112] reframes temporal grounding as a progressive refinement task, where the model first predicts coarse intervals (e.g., “15.0s to 27.5s”) and subsequently refines these estimates by predicting offset adjustments (e.g., “+4.0s and -1.5s”), achieving fine-grained localization through iterative reasoning. Other models, such as VTimeLLM [26] and TPE-VLLM [120], incorporate boundary-aware tasks during pretraining, explicitly teaching the model to reason about event durations and transitions, thereby enhancing temporal precision without the need for explicit time tokens.

Feature Infusion. Feature infusion techniques integrate temporal context directly into visual feature representations by conditioning the feature extraction process. This is often achieved using architectures like Q-Formers, designed to jointly learn spatiotemporal embeddings. Formally, given raw visual features \mathbf{V}_{raw} and a temporal descriptor \mathbf{T}_{desc} (e.g., “This frame is sampled at 2s”), the infused features $\mathbf{F}_{infused}$ can be generated as:

$$\mathbf{F}_{infused} = \text{Extractor}(\mathbf{V}_{raw}, \mathbf{T}_{desc}; \theta), \quad (8)$$

where *Extractor* (e.g., a Q-Former) processes \mathbf{V}_{raw} conditioned on \mathbf{T}_{desc} , allowing the model to capture subtle, context-dependent temporal cues without explicit time tokenization. For example, TimeChat [27] and TemporalVLM [118] leverage this by providing such temporal descriptors as conditional inputs to their Q-Former, guiding the model to incorporate temporal context. Similarly, TimeSuite [33] generates segment-level features that capture temporal dynamics across longer intervals, enabling more comprehensive temporal reasoning.

Summary: Temporal representation and modeling are foundational for endowing VTG-MLLMs with the capacity for precise temporal localization. Explicit modeling strategies directly furnish MLLMs with unambiguous temporal information, offering direct control and interpretability over temporal cues. Implicit modeling, on the other hand, leverages the inherent sequential processing and reasoning capabilities of LLMs or integrates temporal context more subtly during feature extraction. These approaches reflect ongoing exploration into how best to integrate the continuous nature of time with the discrete, symbolic processing of LLMs, ultimately shaping the model’s ability to perform nuanced temporal reasoning and accurate boundary prediction.

4 EXPERIMENTAL EVALUATION

In this section, we provide a comprehensive analysis of the performance of various VTG-MLLMs across four core tasks. We begin by establishing the experimental settings, outlining the benchmark datasets and evaluation metrics that form the basis for fair comparison. Subsequently, we present

TABLE 1: Comparative statistics of major datasets commonly used in VTG. These datasets span a range of tasks, including MR, DC, HD, and GQA, as described in Section 2.1. Key dataset attributes include the number of videos, average video duration, average moment duration, total number of queries, and average query length, providing a comprehensive overview of their scale and complexity.

Dataset	Year	Task	#Videos	Duration / Video	Duration / Moment	#Queries	#Words / Query
Charades-STA [148]	2017	MR	6,672	0.5min	8.1s	16,128	7.2
ActivityNet-Captions [7]	2017	MR & DC	14,926	1.96min	36.2s	71,957	14.8
YouCook2 [149]	2017	DC	2,000	5.26min	-	15,433	8.8
QVHighlights [54]	2021	MR & HD	10,148	2.5min	24.6s	10,310	11.3
NExT-GQA [11]	2023	GQA	1,557	0.70min	6.92s	8,911	-

a detailed performance comparison, considering both zero-shot and fine-tuning scenarios to assess their generalization capabilities and task-specific effectiveness.

4.1 Experimental Settings

This subsection outlines the fundamental experimental settings employed in this survey. We first introduce the *Benchmark Datasets* used for evaluation, followed by a description of the *Evaluation Metrics* employed to assess model performance across different VTG tasks.

4.1.1 Benchmark Datasets

Benchmark datasets are foundational for the development and rigorous evaluation of VTG research, providing essential video content and corresponding natural language annotations necessary for model training and standardized comparisons. The current VTG dataset landscape can be broadly categorized into two groups: *Standard Benchmarks*, which are well-established and widely adopted by the research community, and *Emerging Benchmarks*, designed to address gaps in existing datasets by targeting complex scenarios such as long-term temporal reasoning and fine-grained event distinctions.

Given space limitations, we focus this subsection on key characteristics of select standard benchmarks that are utilized in our experimental evaluations (Sections 4.2 and 4.3). Comparative statistics for these datasets are summarized in Table 1. A comprehensive review, including additional standard and emerging benchmarks, is provided in the supplementary material.

Charades-STA [148], derived from the Charades dataset [150], is widely employed for MR. It includes approximately 6.7k videos depicting indoor daily activities, annotated with about 16.1k query-moment pairs. Videos average 30 seconds, moments average 8 seconds, and textual queries are concise (7.2 words on average), making it suitable for action-centric grounding tasks.

ActivityNet-Captions [7], an extension of the human activity dataset ActivityNet [151], is primarily utilized for DC, though also popular in MR tasks. It features around 15k videos annotated with multiple temporally localized events per video (avg. 3.65 events per video). With an average video duration of 1.96 minutes, average moment length of 36.2 seconds, and longer descriptive queries (14.8 words), making it suitable for evaluating models on complex temporal structures and narratives.

YouCook2 [149] offers a large-scale instructional video dataset tailored for procedural activity understanding, primarily applied in DC. It comprises 2,000 long, untrimmed

cooking videos, collectively spanning 176 hours and covering 89 distinct recipes, resulting in approximately 15.4k query-moment pairs. Each video is densely annotated with precise temporal boundaries for step-level imperative sentences (avg. 8.8 words).

QVHighlights [54] is specifically designed for query-based HD and also finds application in MR. The dataset includes approximately 10.1k YouTube videos, each cropped to a fixed duration of 150 seconds, and features around 10.3k query-moment pairs. A distinctive aspect of QVHighlights is its allowance for multiple disjoint moments per query, reflecting the non-linear nature of salient video highlights, which average about 24.6 seconds. Evaluation of its test set is managed via an official server¹.

NExT-GQA [11] extends the NExT-QA dataset [152], designed to support research on QA tasks that require temporal grounding as evidence. It adds precise temporal annotations for approximately 8.9k QA pairs across 1.6k videos, pinpointing critical segments (avg. 6.9 seconds) needed to answer questions. A notable feature is its weakly-supervised setup where only validation and test sets provide ground-truth temporal labels.

4.1.2 Evaluation Metrics

Effective evaluation of VTG-MLLMs relies on well-defined metrics that assess both the temporal localization accuracy and semantic alignment of predicted moments with ground-truth segments.

Video Moment Retrieval is quantified by two key metrics for localization accuracy: *Mean Intersection over Union (mIoU)* [148] quantifies the average overlap between predicted and ground-truth segments, indicating localization accuracy; and *Recall at Rank n with an IoU threshold of m ($R@n$ ($IoU=m$))* [153] quantifies the percentage of cases where the top- n retrieved segments have an IoU with the ground truth greater than or equal to the threshold m , under varying IoU thresholds (e.g., $R@1$ at $IoU=0.3/0.5/0.7$).

Dense Video Captioning hinges on both temporal precision and linguistic quality. Key metrics include: *SODA_c* [154] assesses structural alignment between captions and temporal segments, emphasizing event segmentation coherence; *METEOR* [155] and *CIDEr* [156] evaluate linguistic similarity between generated and reference captions, considering language precision and recall; and *F1 Score* measures the harmonic mean of precision and recall for instances within a video, assessing coverage and relevance.

Video Highlight Detection commonly relies on metrics: *Mean Average Precision (mAP)* measures the area under the precision-recall curve at various IoU thresholds, such as

1. Evaluation: <https://codalab.lisn.upsaclay.fr/competitions/6937>

TABLE 2: Zero-shot performance comparison on video moment retrieval benchmarks (Charades-STA [148] and ActivityNet-Captions [7]). Methods are grouped by training paradigm (PT: pretraining, TF: training-free). The best results are **boldfaced**, and the second-best results are underlined. This formatting convention is uniformly adopted across all subsequent tables unless explicitly noted otherwise.

Method	Paradigm	Charades-STA [148]				ActivityNet-Captions [7]			
		mIoU	R@0.3	R@0.5	R@0.7	mIoU	R@0.3	R@0.5	R@0.7
Time-R1 [124]	PT	-	78.1	60.8	35.3	-	-	-	-
VideoMind [35]	PT	50.2	<u>73.5</u>	<u>59.1</u>	<u>31.2</u>	33.3	48.4	30.3	15.7
TimeMarker [114]	PT	<u>48.4</u>	<u>73.5</u>	51.9	26.9	-	-	-	-
VideoChat-T [33]	PT	-	69.9	48.7	24.0	-	-	-	-
LLaVA-ST [113]	PT	42.4	63.1	44.8	23.4	-	-	-	-
MeCo [123]	PT	-	-	44.4	17.5	-	-	-	-
E.T.Chat [51]	PT	-	-	43.2	19.4	-	-	-	-
VideoExpert [125]	PT	41.1	61.5	40.3	20.9	-	-	-	-
TRACE [31]	PT	-	-	40.3	19.4	39.0	-	37.7	24.0
VideoChat-TPO [121]	PT	38.1	58.3	40.2	18.4	-	-	-	-
MLLM-TA [104]	PT	-	-	37.9	18.1	-	-	27.6	18.3
Ground-VideoLLM [34]	PT	36.8	54.2	36.4	19.7	<u>36.1</u>	46.2	30.3	19.0
VTG-LLM [108]	PT	-	-	33.8	15.7	-	-	-	-
TPE-VLLM [120]	PT	34.7	55.5	33.1	14.7	<u>36.1</u>	50.4	<u>35.4</u>	<u>19.2</u>
TimeChat [27]	PT	-	-	32.2	13.4	-	-	-	-
HawkEye [32]	PT	33.7	50.6	31.4	14.5	32.7	<u>49.1</u>	29.3	10.7
TemporalVLM [118]	PT	-	-	30.1	13.2	-	-	-	-
GroundingGPT [109]	PT	-	-	29.6	11.9	-	-	-	-
VTimeLLM [26]	PT	31.2	51.0	27.5	11.4	30.4	44.0	27.8	14.3
Momentor [28]	PT	28.5	42.6	26.6	11.6	29.3	42.9	23.0	12.4
TFVTG [92]	TF	44.5	67.0	50.0	<u>24.3</u>	34.1	<u>49.3</u>	27.0	<u>13.4</u>
VTG-GPT [94]	TF	<u>39.8</u>	<u>59.5</u>	<u>43.7</u>	25.9	30.5	47.1	<u>28.3</u>	12.8
Moment-GPT [101]	TF	36.5	58.2	38.4	21.6	<u>30.8</u>	48.1	31.1	14.9
ChatVTG [29]	TF	34.9	52.7	33.0	15.9	27.2	40.7	22.5	9.4

0.5 and 0.75, reflecting the model’s ability to accurately rank and localize key moments; and *HIT@1* [157] calculates the proportion of instances where the top-ranked highlight prediction matches a ground-truth highlight, providing a straightforward measure of first-choice accuracy.

Temporally Grounded Video Question Answering synthesizes metrics for both temporal localization and question-answering accuracy: *mIoU* follows the same setting as in MR; *Mean Intersection over Prediction (mIoP)* [11] evaluates the proportion of the predicted segment overlapping with the ground-truth moment, as an alternative alignment metric; *Acc@QA* [158] measures the percentage of correctly answered questions, irrespective of temporal grounding; and *Acc@GQA* [11] extends *Acc@QA* by requiring correct answers with accurate temporal grounding, where the predicted segment must achieve an IoP score of at least 0.5.

4.2 Zero-Shot Performance Comparison

Zero-shot evaluation measures the ability of VTG-MLLMs to generalize to new datasets without any dataset-specific fine-tuning. This evaluation is crucial for understanding the inherent adaptability and robustness of different models, reflecting their capacity to transfer knowledge across diverse video grounding tasks. For clarity, models are categorized based on their training paradigms: *Pretraining (PT)* approaches (Section 3.2.1), which rely on large-scale multi-modal datasets to build general-purpose grounding capabilities, and *Training-Free (TF)* approaches (Section 3.2.3), which leverage pre-trained foundation models without requiring additional task-specific training.

Video Moment Retrieval. Table 2 presents a comparison of zero-shot performance on the Charades-STA [148] and ActivityNet-Captions [7] benchmarks.

PT Approaches: Among pretraining-based methods, models like Time-R1 [124], VideoMind [35], and TimeMarker [114] demonstrate state-of-the-art zero-shot performance, particularly on the Charades-STA dataset. Time-R1, in particular, success in achieving the highest recall scores likely stems from its innovative use of reinforcement learning. Similarly, TRACE [31] and TPE-VLLM [120] achieve competitive results on ActivityNet-Captions, benefiting from their fine-grained temporal modeling and boundary-aware pertaining.

TF Approaches: In the training-free category, methods like TFVTG [92] and VTG-GPT [94] also exhibit competitive performance, particularly on Charades-STA. Notably, TFVTG outperforms several pretraining-based models in terms of mIoU and R@0.5. These models benefit from flexible, modular designs that leverage powerful pre-trained language and vision models without requiring task-specific fine-tuning.

Dense Video Captioning. Table 3 presents zero-shot performance comparisons on the ActivityNet-Captions [7] and YouCook2 [149] benchmarks.

PT Approaches: Within the pretraining category, methods like TRACE [31], Grounded-VideoLLM [34], and VTimeLLM [26] consistently achieve strong results across key metrics, including SODA_c, CIDEr, and METEOR. TRACE stands out as a top performer on both ActivityNet-Captions and YouCook2, reflecting its effective integration of fine-grained temporal encoding and sequence generation.

TF Approaches: Training-free methods are generally not included in this evaluation, as they often rely on offline

TABLE 3: Zero-Shot performances on dense video captioning benchmarks (ActivityNet-Captions [7] and YouCook2 [149]).

Method	Paradigm	ActivityNet-Captions [7]				YouCook2 [149]		
		SODA_c	CIDEr	METEOR	F1 Score	SODA_c	CIDEr	F1 Score
TRACE [31]	PT	6.0	<u>25.9</u>	<u>6.4</u>	39.3	2.2	8.1	22.4
Ground-VideoLLM [34]	PT	6.0	-	6.8	-	-	-	-
VTimeLLM [26]	PT	<u>5.8</u>	27.6	6.8	-	0.9	3.4	-
VideoExpert [125]	PT	-	-	-	-	<u>2.1</u>	<u>6.0</u>	-
VTG-LLM [108]	PT	5.1	20.7	5.9	34.8	1.5	5.0	<u>17.5</u>
TimeChat [27]	PT	4.7	19.0	5.7	<u>36.9</u>	1.2	3.4	12.6
Momentor [28]	PT	2.3	14.9	4.7	-	-	-	-
TemporalVLM [118]	PT	-	-	-	-	1.2	3.7	13.1

TABLE 4: Zero-Shot performance comparison on video highlight detection benchmark QVHighlights [54].

Method	Paradigm	mAP	HIT@1
VideoExpert [125]	PT	35.8	<u>52.7</u>
TRACE [31]	PT	<u>26.8</u>	42.7
VideoChat-T [33]	PT	26.5	54.1
MLLM-TA [104]	PT	19.2	31.8
VTG-LLM [108]	PT	16.5	33.5
TemporalVLM [118]	PT	16.4	31.3
TimeChat [27]	PT	14.5	23.9
Momentor [28]	PT	7.6	-

MLLMs for generating descriptions, which do not directly align with the standard online evaluation protocols typically used for dense captioning.

Video Highlight Detection. Table 4 presents zero-shot performance comparisons on the QVHighlights [54] benchmark, a challenging dataset specifically designed for localizing query-relevant highlight moments within videos.

PT Approaches: Among pretraining-based methods, VideoExpert [125] stands out with the highest mAP scores, reflecting its strong ability to capture fine-grained temporal relationships and accurately rank highlight moments. In contrast, videoChat-T [33] leads in Hit@1, demonstrating superior single-shot localization accuracy, which highlights its effectiveness in identifying salient temporal regions.

TF Approaches: Training-free methods are typically excluded from direct comparison in this context, as they often rely on frame-query similarity scoring, which inherently aligns well with highlight detection tasks. These approaches, while effective, lack the structured, end-to-end temporal grounding capabilities required for competitive performance in strictly defined zero-shot scenarios.

Temporally Grounded Video Question Answering. Table 5 presents zero-shot performance on the NExT-GQA [11] benchmark, a relatively recent dataset designed to evaluate temporally grounded question answering in videos.

PT Approaches: Among pretraining-based methods, VideoMind [35] showcases exceptional strength in temporal localization, achieving the highest scores across all IoU-related metrics. This proficiency is complemented by its outstanding Acc@GQA performance, largely attributed to its methodical and step-by-step approach to solving GQA problems. In parallel, VideoChat-TPO [121] demonstrates robust and competitive results, underscoring the efficacy of its unique task preference optimization strategy.

TF Approaches: In the training-free category, DeVi [30] stands out, achieving the highest scores in mIoU, Acc@GQA, and Acc@QA. This indicates a strong capacity for fine-grained video temporal understanding, likely due to the

carefully designed pipeline that effectively combines the perceptual capabilities of the pretrained MLLMs with the reasoning ability of LLMs.

Summary: Experimental results across the four tasks suggest that *PT Approaches* generally exhibit strong zero-shot performance due to their explicit temporal learning and large-scale data exposure, particularly excelling in tasks like moment retrieval. While *TF Approaches* are not directly comparable in dense captioning and highlight detection due to architectural differences, models like DeVi [30] show that with well-designed pipelines, *TF Approaches* can still achieve strong temporal reasoning in suitable tasks.

4.3 Fine-Tuning Performance Comparison

Fine-tuning evaluation assesses the performance of VTG-MLLMs after they are trained on the train-split of specific benchmark datasets. Unlike zero-shot evaluation, which tests the generalization ability of models without exposure to task-specific data, this setting allows models to adapt to domain-specific distributions and task requirements, often resulting in substantial performance gains. Notably, the NExT-GQA [11] benchmark is excluded from this evaluation, as its dataset definition does not provide a dedicated training split, precluding fine-tuning.

To better reflect differences in training paradigms, we categorize models into two types: *Pretraining (PT)* (Section 3.2.1) and *Fine-Tuning (FT)* (Section 3.2.2). The *PT* refers to models that undergo large-scale multimodal pretraining prior to task-specific fine-tuning on downstream training splits. In contrast, the *FT* is specifically designed to be trained directly on the training split of the downstream task.

Video Moment Retrieval. Table 6 presents the fine-tuning performance on Charades-STA [148].

PT Approaches: Pretraining-based models exhibit significant performance enhancements when fine-tuned on the benchmark. Notably, reinforcement learning-based methods, Time-R1 [124] and VideoChat-R1 [126], emerge as the top performers, underscoring the profound potential of RL in the MLLM fine-tuning process. Closely following these leaders, models like VideoChat-T [33] and VideoChat-TPO [121] also deliver highly competitive outcomes. These results underscore the benefit of large-scale pretraining as a foundation for task-specific adaptation, allowing models to further refine their temporal grounding capabilities.

FT Approaches: Models designed for direct fine-tuning also exhibit remarkable performance, particularly on Charades-STA. LLaVA-MR [105] leads with the highest mIoU, R@0.5, and R@0.7, outperforming all PT methods on these metrics. Mr.BLIP [117] is a close second with mIoU and

TABLE 5: Zero-Shot performance comparison on temporally grounded video QA benchmark NExT-GQA [11].

Method	Paradigm	mIoP	IoP@0.3	IoP@0.5	mIoU	IoU@0.3	IoU@0.5	Acc@GQA	Acc@QA
DeVi [30]	TF	39.3	-	37.9	22.3	-	17.4	<u>28.0</u>	71.6
VideoMind [35]	PT	<u>39.0</u>	56.0	<u>35.3</u>	31.4	50.2	25.8	28.2	-
VideoChat-TPO [121]	PT	35.6	<u>47.5</u>	32.8	27.7	<u>41.2</u>	<u>23.4</u>	25.5	-
VideoExpert [125]	PT	34.6	45.3	29.3	<u>27.9</u>	41.0	22.4	21.6	<u>71.1</u>
HawkEye [32]	PT	-	-	-	25.7	37.0	19.5	-	-
Ground-VideoLLM [34]	PT	34.5	42.6	34.4	21.1	30.2	18.0	26.7	-
SeViLA [102]	PT	29.5	34.7	22.9	21.7	29.2	13.8	16.6	68.1
VTIMELLM [26]	PT	27.9	30.3	23.8	18.3	27.7	14.1	12.7	-

TABLE 6: Fine-Tuning performance comparison on video moment retrieval benchmark Charades-STA [148].

Charades-STA [148]				
Method	Paradigm	mIoU	R@0.5	R@0.7
LLaVA-MR [105]	FT	59.8	70.7	49.6
Mr.BLIP [117]	FT	<u>58.6</u>	<u>69.3</u>	<u>49.3</u>
VideoLights [97]	FT	52.9	62.0	41.1
ReCorrect [100]	FT	48.7	54.4	31.1
LMR [99]	FT	-	55.9	35.2
LLaViLo [103]	FT	-	55.7	33.4
Time-R1 [124]	PT	-	72.2	<u>50.1</u>
VideoChat-R1 [126]	PT	60.8	<u>71.7</u>	50.2
VideoChat-T [33]	PT	-	67.1	43.0
VideoChat-TPO [121]	PT	<u>55.0</u>	65.0	40.7
SpaceVLLM [122]	PT	-	63.6	38.5
TRACE [31]	PT	-	61.7	41.4
VideoExpert [125]	PT	52.2	60.8	36.5
HawkEye [32]	PT	49.3	58.3	28.8
VTG-LLM [108]	PT	-	57.2	33.4
TemporalVLM [118]	PT	-	54.5	29.0
MLLM-TA [104]	PT	-	48.9	25.3
TimeChat [27]	PT	-	46.7	23.7

TABLE 7: Fine-Tuning performance comparison on dense video captioning benchmark YouCook2 [149].

Method	Paradigm	SODA_c	CIDEr	F1 Score
TRACE [31]	PT	6.7	35.5	31.8
VideoExpert [125]	PT	<u>4.2</u>	<u>18.7</u>	-
VTG-LLM [108]	PT	3.6	13.4	<u>20.6</u>
TimeChat [27]	PT	3.4	11.0	19.5
TemporalVLM [118]	PT	3.4	13.2	20.0

strong recall scores. On ActivityNet-Captions, Mr.BLIP [117] achieves the best R@0.5 and R@0.7 across all reported models. This demonstrates that well-optimized fine-tuning can effectively exploit the strong generalization capabilities of existing large vision-language models without requiring extensive pretraining related to temporal understanding.

Dense Video Captioning and Video Highlight Detection. The fine-tuning performance for DC on YouCook2 [149] and HD on QVHighlights [54] are presented in Table 7 and Table 8, respectively.

PT Approaches: For DC on YouCook2, TRACE [31] stands out significantly among pretraining-based methods, achieving the highest SODA_c, CIDEr, and F1 Score. VideoExpert [125] also shows competitive results, particularly in SODA_c and CIDEr. Turning to HD on QVHighlights, VideoChat-TPO [121] demonstrates exceptional performance among PT approaches, leading with the highest mAP and Hit@1. VideoExpert [125] also achieves a competitive mAP, while VideoChat-T [33] shows a strong Hit@1 score. For both DC and HD, these results show that PT models can be effectively adapted to the specific nuances

TABLE 8: Fine-Tuning performance comparison on video highlight detection benchmark QVHighlights [54].

Method	Paradigm	mAP	HIT@1
VideoChat-TPO [121]	PT	38.8	66.2
VideoExpert [125]	PT	<u>36.1</u>	<u>61.0</u>
TRACE [31]	PT	31.8	51.5
VideoChat-T [33]	PT	27.0	55.3
TemporalVLM [118]	PT	25.1	43.0
VTG-LLM [108]	PT	24.1	41.3
MLLM-TA [104]	PT	23.9	40.1
TimeChat [27]	PT	21.7	37.9

of these tasks through fine-tuning.

FT Approaches: For both DC and HD, results for *FT approaches* are absent in the presented tables under standard fine-tuning protocols. This is primarily due to two factors. Firstly, *FT approaches* of the “Offline Textualization with MLLMs” type (as described in Section 3.2.2) adopt architectures incompatible with the evaluation settings of DC and HD, which is analogous to the exclusion of *TF approaches* from zero-shot evaluations on these benchmarks. Secondly, “Direct Fine-Tuning of MLLMs” for these specific tasks, while theoretically possible, remains unreported in current public research efforts. The complexity of DC (generating multiple, temporally ordered descriptions) and the nuanced requirements of HD (identifying subtle, query-relevant highlights) may necessitate more specialized architectural adaptations or fine-tuning strategies for direct application of MLLMs, which are less explored in current literature compared to their application in MR.

Summary: Fine-tuning enables VTG-MLLMs to better align with task-specific objectives and dataset characteristics, leading to notable performance gains. *PT Approaches* like TRACE [31] and VideoExpert [125] benefit significantly from this process, consistently ranking among top performers across tasks. Meanwhile, *FT Approaches* such as LLaVA-MR [105] and Mr.BLIP [117] achieve competitive or even superior results in MR, highlighting the effectiveness of direct supervision. The absence of *FT Approaches* in tasks like DC and HD reflects current architectural limitations, marking this as an underexplored yet promising direction for future work.

5 LIMITATIONS AND FUTURE DIRECTIONS

The integration of LLMs into VTG has led to significant advancements, showing remarkable potential in understanding and localizing temporal events within videos. However, several critical limitations remain. This section outlines key challenges and proposes future research directions, focusing on 1) training paradigms, 2) feature representation, 3) temporal modeling, and 4) multimodal integration.

Training Paradigms: While fine-tuning (Section 3.2.2) and training-free (Section 3.2.3) approaches offer practical advantages, their performance is fundamentally capped by the capabilities inherited from their pretrained foundations [58, 132]. These strategies are inherently constrained in their ability to instill the fine-grained temporal nuances essential for complex VTG, particularly if such understanding is lacking in the base model.

To address this limitation, future research should advance along two complementary paths. The primary path involves enhancing the pretraining paradigm (Section 3.2.1) itself, by developing more targeted objectives and datasets that foster a deeper, native temporal perception in foundation models. A parallel, more agile path lies in pioneering advanced post-training techniques such as RL, offering a powerful mechanism to sharpen a model’s specialized skills beyond the reach of conventional Supervised Fine-Tuning (SFT). The synergy between these techniques will likely define the next generation of state-of-the-art VTG-MLLMs.

Feature Representation: Efficient and effective feature representation remains a significant challenge in VTG-MLLMs. Accurate temporal grounding, particularly in long videos or tasks requiring fine-grained event distinction, demands the processing of dense, high-resolution visual inputs. However, this directly conflicts with the strict input token limitations of current LLM architectures. Existing feature compression strategies (Section 3.3.1) help mitigate this issue but often risk discarding critical temporal cues, potentially compromising grounding accuracy.

To overcome this, future work could explore adaptive token selection strategies that dynamically prioritize the most relevant visual tokens while aggressively filtering out redundant information. Techniques such as temporal saliency detection, attention-based token pruning, and hierarchical feature aggregation could significantly reduce the computational burden without sacrificing precision. Moreover, developing more compact yet semantically rich video embeddings could further enhance scalability across a broad range of video understanding tasks [159, 160].

Temporal Modeling: Accurate temporal reasoning is central to VTG but remains an open challenge. Current approaches vary widely in their representation of temporal information, including explicit (e.g., absolute timestamps [114], frame indices [105]) and implicit (e.g., relative positions [111]) modeling strategies (Section 3.3.2). However, there is no clear consensus on the most effective approach, nor a comprehensive understanding of the trade-offs between these strategies.

Future research should focus on developing unified and expressive temporal encoding mechanisms that effectively capture both fine-grained event details and long-term dependencies. This entails exploring novel representations that embed temporal attributes such as duration, order, and causal relationships directly into the latent space of MLLMs. Furthermore, integrating temporal reasoning tasks into pretraining, such as sequence prediction, duration estimation, and multi-step causal reasoning, could substantially improve grounding precision and robustness.

Multimodal Integration: While most current VTG-MLLMs rely on visual and textual inputs, real-world videos are inherently multimodal, often containing rich audio

signals that provide complementary temporal cues. For example, sounds like speech, footsteps, and background noise can serve as precise temporal markers, enhancing grounding accuracy, particularly in visually ambiguous scenes [161], [162]. However, effectively integrating these additional modalities presents several challenges, including the need for precise temporal alignment, the scarcity of large-scale, multimodal datasets with accurate temporal annotations, and the increased computational complexity.

To address these issues, future work should focus on developing multimodal fusion architectures that can jointly model audio-visual signals with high temporal precision. Potential approaches include leveraging self-supervised learning to pretrain multimodal encoders, designing cross-modal attention mechanisms that dynamically weight different input streams, and creating synthetic multimodal training data to reduce the reliance on expensive manual annotation. Additionally, integrating pre-trained audio-language models, such as Whisper [163], with visual grounding pipelines could further enhance temporal reasoning across complex, multi-channel video content.

6 CONCLUSION

In this survey, we provide a comprehensive review of VTG-MLLMs, a critical area in fine-grained video understanding. MLLMs have significantly transformed VTG by introducing advanced reasoning and cross-modal alignment capabilities that surpass traditional task-specific methods. We systematically categorize and analyze current approaches, distinguishing between the functional roles of MLLMs as facilitators and executors. Additionally, we examine diverse training paradigms, including pretraining, fine-tuning, and training-free methods, alongside cutting-edge video feature processing techniques, such as visual feature compression and explicit versus implicit temporal modeling. To support ongoing research, we provide a structured overview of widely used datasets and benchmark comparisons, highlighting key performance trends across different VTG-MLLM architectures. Finally, we identify current challenges and propose future research directions to address limitations in training efficiency, temporal representation, and multimodal integration.

REFERENCES

- [1] Y. Wang, K. Li, Y. Li, Y. He, B. Huang, Z. Zhao, H. Zhang, J. Xu, Y. Liu, Z. Wang *et al.*, “Internvideo: General video foundation models via generative and discriminative learning,” *arXiv:2212.03191*, 2022.
- [2] Y. Wang, K. Li, X. Li, J. Yu, Y. He, G. Chen, B. Pei, R. Zheng, Z. Wang, Y. Shi *et al.*, “Internvideo2: Scaling foundation models for multimodal video understanding,” in *ECCV*, 2024, pp. 396–416.
- [3] M. Wang, J. Xing, J. Mei, Y. Liu, and Y. Jiang, “Actionclip: Adapting language-image pretrained models for video action recognition,” *TNNLS*, vol. 36, no. 1, pp. 625–637, 2023.
- [4] A. Yang, A. Nagrani, P. H. Seo, A. Miech, J. Pont-Tuset, I. Laptev, J. Sivic, and C. Schmid, “Vid2seq: Large-scale pretraining of a visual language model for dense video captioning,” in *CVPR*, 2023, pp. 10714–10726.
- [5] M. Regneri, M. Rohrbach, D. Wetzel, S. Thater, B. Schiele, and M. Pinkal, “Grounding action descriptions in videos,” *TACL*, vol. 1, pp. 25–36, 2013.
- [6] L. Anne Hendricks, O. Wang, E. Shechtman, J. Sivic, T. Darrell, and B. Russell, “Localizing moments in video with natural language,” in *ICCV*, 2017, pp. 5803–5812.

- [7] R. Krishna, K. Hata, F. Ren, L. Fei-Fei, and J. Carlos Niebles, "Dense-captioning events in videos," in *ICCV*, 2017, pp. 706–715.
- [8] T. Wang, R. Zhang, Z. Lu, F. Zheng, R. Cheng, and P. Luo, "End-to-end dense video captioning with parallel decoding," in *ICCV*, 2021, pp. 6847–6857.
- [9] Y. Song, J. Vallmitjana, A. Stent, and A. Jaimes, "Tvsum: Summarizing web videos using titles," in *CVPR*, 2015, pp. 5179–5187.
- [10] B. Xiong, Y. Kalantidis, D. Ghadiyaram, and K. Grauman, "Less is more: Learning highlight detection from video duration," in *CVPR*, 2019, pp. 1258–1267.
- [11] J. Xiao, A. Yao, Y. Li, and T.-S. Chua, "Can i trust your answer? visually grounded video question answering," in *CVPR*, 2024, pp. 13204–13214.
- [12] J.-J. Chen, Y.-C. Liao, H.-C. Lin, Y.-C. Yu, Y.-C. Chen, and F. Wang, "Rextime: A benchmark suite for reasoning-across-time in videos," in *NeurIPS*, 2024, pp. 28662–28673.
- [13] D. Zhang, X. Dai, X. Wang, Y.-F. Wang, and L. S. Davis, "Man: Moment alignment network for natural language moment retrieval via iterative graph adjustment," in *CVPR*, 2019, pp. 1247–1257.
- [14] W. Moon, S. Hyun, S. Park, D. Park, and J.-P. Heo, "Query-dependent video representation for moment retrieval and highlight detection," in *CVPR*, 2023, pp. 23023–23033.
- [15] M. Ma, S. Yoon, J. Kim, Y. Lee, S. Kang, and C. D. Yoo, "Vlanet: Video-language alignment network for weakly-supervised video moment retrieval," in *ECCV*, 2020, pp. 156–171.
- [16] M. Zhang, Y. Yang, X. Chen, Y. Ji, X. Xu, J. Li, and H. T. Shen, "Multi-stage aggregated transformer network for temporal language localization in videos," in *CVPR*, 2021, pp. 12669–12678.
- [17] Y. Liu, S. Li, Y. Wu, C.-W. Chen, Y. Shan, and X. Qie, "Umt: Unified multi-modal transformers for joint video moment retrieval and highlight detection," in *CVPR*, 2022, pp. 3042–3051.
- [18] Y. Yuan, T. Mei, and W. Zhu, "To find where you talk: Temporal sentence localization in video with attention based location regression," in *AAAI*, 2019, pp. 9159–9166.
- [19] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin *et al.*, "Opt: Open pre-trained transformer language models," *arXiv:2205.01068*, 2022.
- [20] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, Y. Li, X. Wang, M. Dehghani, S. Brahma *et al.*, "Scaling instruction-finetuned language models," *JMLR*, vol. 25, no. 70, pp. 1–53, 2024.
- [21] D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi *et al.*, "Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning," *arXiv:2501.12948*, 2025.
- [22] H. Liu, C. Li, Q. Wu, and Y. J. Lee, "Visual instruction tuning," in *NeurIPS*, 2023, pp. 34892–34916.
- [23] W. Dai, J. Li, D. Li, A. Tiong, J. Zhao, W. Wang, B. Li, P. N. Fung, and S. Hoi, "Instructblip: Towards general-purpose vision-language models with instruction tuning," in *NeurIPS*, 2023, pp. 49250–49267.
- [24] Z. Cheng, S. Leng, H. Zhang, Y. Xin, X. Li, G. Chen, Y. Zhu, W. Zhang, Z. Luo, D. Zhao, and L. Bing, "Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms," *arXiv:2406.07476*, 2024.
- [25] L. Xu, Y. Zhao, D. Zhou, Z. Lin, S. K. Ng, and J. Feng, "Pllava: Parameter-free llava extension from images to videos for video dense captioning," *arXiv:2404.16994*, 2024.
- [26] B. Huang, X. Wang, H. Chen, Z. Song, and W. Zhu, "Vtimellm: Empower llm to grasp video moments," in *CVPR*, 2024, pp. 14271–14280.
- [27] S. Ren, L. Yao, S. Li, X. Sun, and L. Hou, "Timechat: A time-sensitive multimodal large language model for long video understanding," in *CVPR*, 2024, pp. 14313–14323.
- [28] L. Qian, J. Li, Y. Wu, Y. Ye, H. Fei, T.-S. Chua, Y. Zhuang, and S. Tang, "Momentor: Advancing video large language model with fine-grained temporal reasoning," in *ICML*, 2024, pp. 41340–41356.
- [29] M. Qu, X. Chen, W. Liu, A. Li, and Y. Zhao, "Chatvtg: Video temporal grounding via chat with video dialogue large language models," in *CVPR*, 2024, pp. 1847–1856.
- [30] H. Qin, J. Xiao, and A. Yao, "Question-answering dense video events," *arXiv:2409.04388*, 2024.
- [31] Y. Guo, J. Liu, M. Li, X. Tang, Q. Liu, and X. Chen, "Trace: Temporal grounding video llm via causal event modeling," in *ICLR*, 2025.
- [32] Y. Wang, X. Meng, J. Liang, Y. Wang, Q. Liu, and D. Zhao, "Hawkeye: Training video-text llms for grounding text in videos," *arXiv:2403.10228*, 2024.
- [33] X. Zeng, K. Li, C. Wang, X. Li, T. Jiang, Z. Yan, S. Li, Y. Shi, Z. Yue, Y. Wang *et al.*, "Timesuite: Improving mllms for long video understanding via grounded tuning," in *ICLR*, 2025.
- [34] H. Wang, Z. Xu, Y. Cheng, S. Diao, Y. Zhou, Y. Cao, Q. Wang, W. Ge, and L. Huang, "Grounded-videollm: Sharpening fine-grained temporal grounding in video large language models," *arXiv:2410.03290*, 2024.
- [35] Y. Liu, K. Q. Lin, C. W. Chen, and M. Z. Shou, "Videomind: A chain-of-lora agent for long video reasoning," *arXiv:2503.13444*, 2025.
- [36] C. Li, Z. Gan, Z. Yang, J. Yang, L. Li, L. Wang, J. Gao *et al.*, "Multimodal foundation models: From specialists to general-purpose assistants," *arXiv:2309.10020*, 2023.
- [37] Y. Zhu, X. Li, C. Liu, M. Zolfaghari, Y. Xiong, C. Wu, Z. Zhang, J. Tighe, R. Manmatha, and M. Li, "A comprehensive study of deep video action recognition," *arXiv:2012.06567*, 2020.
- [38] M. Abdar, M. Kollati, S. Kuraparthi, F. Pourpanah, D. McDuff, M. Ghavamzadeh, S. Yan, A. Mohamed, A. Khosravi, E. Cambria *et al.*, "A review of deep learning for video captioning," *TPAMI*, pp. 1–20, 2024.
- [39] Z. Xing, Q. Feng, H. Chen, Q. Dai, H. Hu, H. Xu, Z. Wu, and Y.-G. Jiang, "A survey on video diffusion models," *ACM Computing Surveys*, vol. 57, no. 2, pp. 1–42, 2024.
- [40] J. Zhang, J. Huang, S. Jin, and S. Lu, "Vision-language models for vision tasks: A survey," *TPAMI*, vol. 46, no. 8, pp. 5625–5644, 2024.
- [41] X. Liu, X. Nie, Z. Tan, J. Guo, and Y. Yin, "A survey on natural language video localization," *arXiv:2104.00234*, 2021.
- [42] Y. Yang, Z. Li, and G. Zeng, "A survey of temporal activity localization via language in untrimmed videos," in *ICCV*, 2020, pp. 596–601.
- [43] X. Lan, Y. Yuan, X. Wang, Z. Wang, and W. Zhu, "A survey on temporal sentence grounding in videos," *TOMCCAP*, vol. 19, no. 2, pp. 1–33, 2023.
- [44] H. Zhang, A. Sun, W. Jing, and J. T. Zhou, "Temporal sentence grounding in videos: A survey and future directions," *TPAMI*, vol. 45, no. 8, pp. 10443–10465, 2023.
- [45] S. Zhang, H. Peng, J. Fu, and J. Luo, "Learning 2d temporal adjacent networks for moment localization with natural language," in *AAAI*, 2020, pp. 12870–12877.
- [46] K. Q. Lin, P. Zhang, J. Chen, S. Pramanick, D. Gao, A. J. Wang, R. Yan, and M. Z. Shou, "Univtg: Towards unified video-language temporal grounding," in *ICCV*, 2023, pp. 2794–2804.
- [47] W. Yang, T. Zhang, Y. Zhang, and F. Wu, "Local correspondence network for weakly supervised temporal sentence grounding," *TIP*, vol. 30, pp. 3252–3262, 2021.
- [48] D. Liu, X. Qu, J. Dong, P. Zhou, Y. Cheng, W. Wei, Z. Xu, and Y. Xie, "Context-aware biaffine localizing network for temporal sentence grounding," in *CVPR*, 2021, pp. 11235–11244.
- [49] S. Xiao, L. Chen, S. Zhang, W. Ji, J. Shao, L. Ye, and J. Xiao, "Boundary proposal network for two-stage natural language video localization," in *AAAI*, 2021, pp. 2986–2994.
- [50] H. Zhang, A. Sun, W. Jing, L. Zhen, J. T. Zhou, and R. S. M. Goh, "Natural language video localization: A revisit in span-based question answering framework," *TPAMI*, vol. 44, no. 8, pp. 4252–4266, 2021.
- [51] Y. Liu, Z. Ma, Z. Qi, Y. Wu, Y. Shan, and C. W. Chen, "Et bench: Towards open-ended event-level video-language understanding," in *NeurIPS*, 2024, pp. 32076–32110.
- [52] L. Zhou, Y. Zhou, J. J. Corso, R. Socher, and C. Xiong, "End-to-end dense video captioning with masked transformer," in *CVPR*, 2018, pp. 8739–8748.
- [53] V. Iashin and E. Rahtu, "Multi-modal dense video captioning," in *CVPRW*, 2020, pp. 958–959.
- [54] J. Lei, T. L. Berg, and M. Bansal, "Detecting moments and highlights in videos via natural language queries," in *NeurIPS*, 2021, pp. 11846–11858.
- [55] G. Chen, Y. Liu, Y. Huang, Y. He, B. Pei, J. Xu, Y. Wang, T. Lu, and L. Wang, "Cg-bench: Clue-grounded question answering benchmark for long video understanding," *arXiv:2412.12075*, 2024.
- [56] H. Liu, X. Ma, C. Zhong, Y. Zhang, and W. Lin, "Timecraft: Navigate weakly-supervised temporal grounded video question answering via bi-directional reasoning," in *ECCV*, 2025, pp. 92–107.
- [57] H. Zhang, X. Li, and L. Bing, "Video-llama: An instruction-tuned audio-visual language model for video understanding,"

- in *EMNLP*, 2023, pp. 543–553.
- [58] B. Lin, Y. Ye, B. Zhu, J. Cui, M. Ning, P. Jin, and L. Yuan, “Video-llava: Learning united visual representation by alignment before projection,” *arXiv:2311.10122*, 2023.
- [59] K. Li, Y. He, Y. Wang, Y. Li, W. Wang, P. Luo, Y. Wang, L. Wang, and Y. Qiao, “Videocat: Chat-centric video understanding,” *arXiv:2305.06355*, 2023.
- [60] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *ICLR*, 2020.
- [61] Q. Sun, Y. Fang, L. Wu, X. Wang, and Y. Cao, “Eva-clip: Improved training techniques for clip at scale,” *arXiv:2303.15389*, 2023.
- [62] Y. Fang, W. Wang, B. Xie, Q. Sun, L. Wu, X. Wang, T. Huang, X. Wang, and Y. Cao, “Eva: Exploring the limits of masked visual representation learning at scale,” in *CVPR*, 2023, pp. 19358–19369.
- [63] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *ICML*, 2021, pp. 8748–8763.
- [64] C. Feichtenhofer, Y. Li, K. He *et al.*, “Masked autoencoders as spatiotemporal learners,” in *NeurIPS*, 2022, pp. 35946–35958.
- [65] Z. Tong, Y. Song, J. Wang, and L. Wang, “Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training,” in *NeurIPS*, 2022, pp. 10078–10093.
- [66] L. Wang, B. Huang, Z. Zhao, Z. Tong, Y. He, Y. Wang, Y. Wang, and Y. Qiao, “Videomae v2: Scaling video masked autoencoders with dual masking,” in *CVPR*, 2023, pp. 14549–14560.
- [67] J. Wang, Y. Ge, R. Yan, Y. Ge, K. Q. Lin, S. Tsutsui, X. Lin, G. Cai, J. Wu, Y. Shan *et al.*, “All in one: Exploring unified video-language pre-training,” in *CVPR*, 2023, pp. 6598–6608.
- [68] G. Sun, W. Yu, C. Tang, X. Chen, T. Tan, W. Li, L. Lu, Z. Ma, Y. Wang, and C. Zhang, “video-salmonn: Speech-enhanced audio-visual large language models,” in *ICML*, 2024, pp. 47198–47217.
- [69] S. Azad, V. Vineet, and Y. S. Rawat, “Hierarq: Task-aware hierarchical q-former for enhanced video understanding,” *arXiv:2503.08585*, 2025.
- [70] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds *et al.*, “Flamingo: a visual language model for few-shot learning,” in *NeurIPS*, 2022, pp. 23716–23736.
- [71] J. Li, D. Li, S. Savarese, and S. Hoi, “Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models,” in *ICML*, 2023, pp. 19730–19742.
- [72] Q. Ye, H. Xu, G. Xu, J. Ye, M. Yan, Y. Zhou, J. Wang, A. Hu, P. Shi, Y. Shi *et al.*, “mplug-owl: Modularization empowers large language models with multimodality,” *arXiv:2304.14178*, 2023.
- [73] S. Huang, L. Dong, W. Wang, Y. Hao, S. Singhal, S. Ma, T. Lv, L. Cui, O. K. Mohammed, B. Patra *et al.*, “Language is not all you need: Aligning perception with language models,” in *NeurIPS*, 2023, pp. 72096–72109.
- [74] H. Liu, C. Li, Y. Li, and Y. J. Lee, “Improved baselines with visual instruction tuning,” in *CVPR*, 2024, pp. 26296–26306.
- [75] F. Li, R. Zhang, H. Zhang, Y. Zhang, B. Li, W. Li, Z. Ma, and C. Li, “Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models,” *arXiv:2407.07895*, 2024.
- [76] J. Ye, H. Xu, H. Liu, A. Hu, M. Yan, Q. Qian, J. Zhang, F. Huang, and J. Zhou, “mplug-owl3: Towards long image-sequence understanding in multi-modal large language models,” *arXiv:2408.04840*, 2024.
- [77] K. Chen, D. Shen, H. Zhong, H. Zhong, K. Xia, D. Xu, W. Yuan, Y. Hu, B. Wen, T. Zhang *et al.*, “Evlm: An efficient vision-language model for visual understanding,” *arXiv:2407.14177*, 2024.
- [78] M. Shi, S. Wang, C.-Y. Chen, J. Jain, K. Wang, J. Xiong, G. Liu, Z. Yu, and H. Shi, “Slow-fast architecture for video multi-modal large language models,” *arXiv:2504.01328*, 2025.
- [79] A. Miech, D. Zhukov, J.-B. Alayrac, M. Tapaswi, I. Laptev, and J. Sivic, “Howto100m: Learning a text-video embedding by watching hundred million narrated video clips,” in *ICCV*, 2019, pp. 2630–2640.
- [80] M. Bain, A. Nagrani, G. Varol, and A. Zisserman, “Frozen in time: A joint video and image encoder for end-to-end retrieval,” in *ICCV*, 2021, pp. 1728–1738.
- [81] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis, M. Wortsman *et al.*, “Laion-5b: An open large-scale dataset for training next generation image-text models,” in *NeurIPS*, 2022, pp. 25278–25294.
- [82] L. Li, Y. Yin, S. Li, L. Chen, P. Wang, S. Ren, M. Li, Y. Yang, J. Xu, X. Sun *et al.*, “M³it: A large-scale dataset towards multi-modal multilingual instruction tuning,” *arXiv:2306.04387*, 2023.
- [83] Z. Yin, J. Wang, J. Cao, Z. Shi, D. Liu, M. Li, X. Huang, Z. Wang, L. Sheng, L. Bai *et al.*, “Lamm: Language-assisted multi-modal instruction-tuning dataset, framework, and benchmark,” in *NeurIPS*, 2023, pp. 26650–26685.
- [84] H. Shao, S. Qian, H. Xiao, G. Song, Z. Zong, L. Wang, Y. Liu, and H. Li, “Visual cot: Advancing multi-modal language models with a comprehensive dataset and benchmark for chain-of-thought reasoning,” in *NeurIPS*, 2024, pp. 8612–8642.
- [85] T. Gao, P. Chen, M. Zhang, C. Fu, Y. Shen, Y. Zhang, S. Zhang, X. Zheng, X. Sun, L. Cao *et al.*, “Cantor: Inspiring multimodal chain-of-thought of mllm,” in *ACM MM*, 2024, pp. 9096–9105.
- [86] J. Wu, Z. Zhang, Y. Xia, X. Li, Z. Xia, A. Chang, T. Yu, S. Kim, R. A. Rossi, R. Zhang *et al.*, “Visual prompting in multimodal large language models: A survey,” *arXiv:2409.15310*, 2024.
- [87] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen *et al.*, “Lora: Low-rank adaptation of large language models,” in *ICLR*, 2022.
- [88] R. Pan, X. Liu, S. Diao, R. Pi, J. Zhang, C. Han, and T. Zhang, “Lisa: layerwise importance sampling for memory-efficient large language model fine-tuning,” in *NeurIPS*, 2024, pp. 57018–57049.
- [89] S.-Y. Liu, C.-Y. Wang, H. Yin, P. Molchanov, Y.-C. F. Wang, K.-T. Cheng, and M.-H. Chen, “Dora: Weight-decomposed low-rank adaptation,” in *ICML*, 2024, pp. 32100–32121.
- [90] W. Cai, J. Huang, S. Gong, H. Jin, and Y. Liu, “Mllm as video narrator: Mitigating modality imbalance in video moment retrieval,” *arXiv:2406.17880*, 2024.
- [91] S. Di and W. Xie, “Grounded question-answering in long egocentric videos,” in *CVPR*, 2024, pp. 12934–12943.
- [92] M. Zheng, X. Cai, Q. Chen, Y. Peng, and Y. Liu, “Training-free video temporal grounding using large-scale pre-trained models,” in *ECCV*, 2025, pp. 20–37.
- [93] H. Chen, X. Wang, H. Chen, Z. Song, J. Jia, and W. Zhu, “Grounding-prompter: Prompting llm with multimodal information for temporal sentence grounding in long videos,” *arXiv:2312.17117*, 2023.
- [94] Y. Xu, Y. Sun, Z. Xie, B. Zhai, and S. Du, “Vtg-gpt: Tuning-free zero-shot video temporal grounding with gpt,” *Applied Sciences*, vol. 14, no. 5, p. 1894, 2024.
- [95] H. Chen, X. Wang, H. Chen, Z. Zhang, W. Feng, B. Huang, J. Jia, and W. Zhu, “Verified: A video corpus moment retrieval benchmark for fine-grained video understanding,” *arXiv:2410.08593*, 2024.
- [96] H. Lee, S. Hong, M. Sung, and J. Choi, “Infusing environmental captions for long-form video language grounding,” *arXiv:2408.02336*, 2024.
- [97] D. Paul, M. R. Parvez, N. Mohammed, and S. Rahman, “Video-lights: Feature refinement and cross-task alignment transformer for joint video highlight detection and moment retrieval,” *arXiv:2412.01558*, 2024.
- [98] Y. Sun, Y. Xu, Z. Xie, Y. Shu, and S. Du, “Gptsee: Enhancing moment retrieval and highlight detection via description-based similarity features,” *SPL*, vol. 31, pp. 521–525, 2024.
- [99] W. Liu, B. Miao, J. Cao, X. Zhu, B. Liu, M. Nasim, and A. Mian, “Context-enhanced video moment retrieval with large language models,” *arXiv:2405.12540*, 2024.
- [100] P. Bao, C. Kong, Z. Shao, B. P. Ng, M. H. Er, and A. C. Kot, “Vid-morp: Video moment retrieval pretraining from unlabeled videos in the wild,” *arXiv:2412.00811*, 2024.
- [101] Y. Xu, Y. Sun, B. Zhai, M. Li, W. Liang, Y. Li, and S. Du, “Zero-shot video moment retrieval via off-the-shelf multimodal large language models,” in *AAAI*, 2025, pp. 8978–8986.
- [102] S. Yu, J. Cho, P. Yadav, and M. Bansal, “Self-chained image-language model for video localization and question answering,” in *NeurIPS*, 2023, pp. 76749–76771.
- [103] K. Ma, X. Zang, Z. Feng, H. Fang, C. Ban, Y. Wei, Z. He, Y. Li, and H. Sun, “Llavilo: Boosting video moment retrieval via adapter-based multimodal modeling,” in *ICCV*, 2023, pp. 2798–2803.
- [104] Y. Liu, H. Hou, F. Ma, S. Ni, and F. R. Yu, “Mllm-ta: Leveraging multimodal large language models for precise temporal video grounding,” *SPL*, vol. 32, pp. 281–285, 2025.
- [105] W. Lu, J. Li, A. Yu, M.-C. Chang, S. Ji, and M. Xia, “Llava-mr: Large language-and-vision assistant for video moment retrieval,”

- arXiv:2411.14505*, 2024.
- [106] Y. Wu, X. Hu, Y. Sun, Y. Zhou, W. Zhu, F. Rao, B. Schiele, and X. Yang, "Number it: Temporal grounding videos like flipping manga," *arXiv:2411.10332*, 2024.
- [107] A. Deng, Z. Gao, A. Choudhuri, B. Planche, M. Zheng, B. Wang, T. Chen, C. Chen, and Z. Wu, "Seq2time: Sequential knowledge transfer for video llm temporal grounding," *arXiv:2411.16932*, 2024.
- [108] Y. Guo, J. Liu, M. Li, D. Cheng, X. Tang, D. Sui, Q. Liu, X. Chen, and K. Zhao, "Vtg-llm: Integrating timestamp knowledge into video llms for enhanced video temporal grounding," in *AAAI*, 2025, pp. 3302–3310.
- [109] Z. Li, Q. Xu, D. Zhang, H. Song, Y. Cai, Q. Qi, R. Zhou, J. Pan, Z. Li, V. Tu *et al.*, "Groundinggpt: Language enhanced multimodal grounding model," in *ACL*, 2024, pp. 6657–6678.
- [110] T. Hannan, M. M. Islam, J. Gu, T. Seidl, and G. Bertasius, "Revisionllm: Recursive vision-language model for temporal grounding in hour-long videos," *arXiv:2411.14901*, 2024.
- [111] D.-A. Huang, S. Liao, S. Radhakrishnan, H. Yin, P. Molchanov, Z. Yu, and J. Kautz, "Lita: Language instructed temporal-localization assistant," in *ECCV*, 2025, pp. 202–218.
- [112] X. Wang, F. Cheng, S. Wang, H. Wang, M. M. Islam, L. Torresani, M. Bansal, G. Bertasius, and D. Crandall, "Timerefine: Temporal grounding with time refining video llm," *arXiv:2412.09601*, 2024.
- [113] H. Li, J. Chen, Z. Wei, S. Huang, T. Hui, J. Gao, X. Wei, and S. Liu, "Llava-st: A multimodal large language model for fine-grained spatial-temporal understanding," *arXiv:2501.08282*, 2025.
- [114] S. Chen, X. Lan, Y. Yuan, Z. Jie, and L. Ma, "Timemarker: A versatile video-llm for long and short video understanding with superior temporal localization ability," *arXiv:2411.18211*, 2024.
- [115] Q. Chen, S. Di, and W. Xie, "Grounded multi-hop videoqa in long-form egocentric videos," in *AAAI*, 2025, pp. 2159–2167.
- [116] M. Nie, D. Ding, C. Wang, Y. Guo, J. Han, H. Xu, and L. Zhang, "Slowfocus: Enhancing fine-grained temporal understanding in video llm," in *NeurIPS*, 2024.
- [117] B. Meinardus, A. Batra, A. Rohrbach, and M. Rohrbach, "The surprising effectiveness of multimodal large language models for video moment retrieval," *arXiv:2406.18113*, 2024.
- [118] F. J. Fateh, U. Ahmed, H. Khan, M. Z. Zia, and Q.-H. Tran, "Video llms for temporal reasoning in long videos," *arXiv:2412.02930*, 2024.
- [119] Y. Wang, Y. Wang, P. Wu, J. Liang, D. Zhao, Y. Liu, and Z. Zheng, "Efficient temporal extrapolation of multimodal large language models with temporal grounding bridge," in *EMNLP*, 2024, pp. 9972–9987.
- [120] X. Li, B. Wang, G. Shi, C. Feng, and J. Teng, "Mitigating the discrepancy between video and text temporal sequences: A time-perception enhanced video grounding method for llm," in *COLING*, 2025, pp. 9804–9813.
- [121] Z. Yan, Z. Li, Y. He, C. Wang, K. Li, X. Li, X. Zeng, Z. Wang, Y. Wang, Y. Qiao *et al.*, "Task preference optimization: Improving multimodal large language models with vision task alignment," *arXiv:2412.19326*, 2024.
- [122] J. Wang, Z. Liu, Y. Li, J. Ge, H. Xie, Y. Zhang *et al.*, "Spacevllm: Endowing multimodal large language model with spatio-temporal video grounding capability," *arXiv:2503.13983*, 2025.
- [123] Z. Pang, M. Otani, and Y. Nakashima, "Measure twice, cut once: Grasping video structures and event semantics with llms for video temporal localization," *arXiv:2503.09027*, 2025.
- [124] Y. Wang, Z. Wang, B. Xu, Y. Du, K. Lin, Z. Xiao, Z. Yue, J. Ju, L. Zhang, D. Yang, X. Fang, Z. He, Z. Luo, W. Wang, J. Lin, J. Luan, and Q. Jin, "Time-r1: Post-training large vision language model for temporal video grounding," *arXiv:2503.13377*, 2025.
- [125] H. Zhao, G.-P. Ji, R. Yan, H. Xiong, and Z. Li, "Videoexpert: Augmented llm for temporal-sensitive video understanding," *arXiv:2504.07519*, 2025.
- [126] X. Li, Z. Yan, D. Meng, L. Dong, X. Zeng, Y. He, Y. Wang, Y. Qiao, Y. Wang, and L. Wang, "Videochat-r1: Enhancing spatio-temporal perception via reinforcement fine-tuning," *arXiv:2504.06958*, 2025.
- [127] F. Luo, S. Lou, C. Chen, Z. Wang, C. Li, W. Shen, J. Guo, P. Li, M. Yan, J. Zhang *et al.*, "Museg: Reinforcing video temporal understanding via timestamp-aware multi-segment grounding," *arXiv:2505.20715*, 2025.
- [128] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, "Llama: Open and efficient foundation language models," *arXiv:2302.13971*, 2023.
- [129] K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger, H. Jiang, M. Liu, X. Liu *et al.*, "Ego4d: Around the world in 3,000 hours of egocentric video," in *CVPR*, 2022, pp. 18995–19012.
- [130] OpenAI, "Gpt-4o," 2024. [Online]. Available: <https://openai.com/index/hello-gpt-4o/>
- [131] G. Team, P. Georgiev, V. I. Lei, R. Burnell, L. Bai, A. Gulati, G. Tanzer, D. Vincent, Z. Pan, S. Wang *et al.*, "Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context," *arXiv:2403.05530*, 2024.
- [132] M. Maaz, H. Rasheed, S. Khan, and F. S. Khan, "Video-chatgpt: Towards detailed video understanding via large vision and language models," in *ACL*, 2024.
- [133] N. Reimers and I. Gurevych, "Sentence-bert: Sentence embeddings using siamese bert-networks," in *EMNLP*, 2019, pp. 3982–3992.
- [134] J. Chen, D. Zhu, X. Shen, X. Li, Z. Liu, P. Zhang, R. Krishnamoorthi, V. Chandra, Y. Xiong, and M. Elhoseiny, "Minigpt-v2: large language model as a unified interface for vision-language multi-task learning," *arXiv:2310.09478*, 2023.
- [135] A. Yang, B. Xiao, B. Wang, B. Zhang, C. Bian, C. Yin, C. Lv, D. Pan, D. Wang, D. Yan *et al.*, "Baichuan 2: Open large-scale language models," *arXiv:2309.10305*, 2023.
- [136] K. C. Fraser and S. Kiritchenko, "Examining gender and racial bias in large vision-language models using a novel dataset of parallel images," in *EACL*, 2024, pp. 690–713.
- [137] Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, X. Bi, H. Zhang, M. Zhang, Y. Li, Y. Wu *et al.*, "Deepseekmath: Pushing the limits of mathematical reasoning in open language models," *arXiv:2402.03300*, 2024.
- [138] R. Zellers, X. Lu, J. Hessel, Y. Yu, J. S. Park, J. Cao, A. Farhadi, and Y. Choi, "Merlot: Multimodal neural script knowledge models," in *NeurIPS*, 2021, pp. 23634–23651.
- [139] Y. Wang, Y. He, Y. Li, K. Li, J. Yu, X. Ma, X. Li, G. Chen, X. Chen, Y. Wang *et al.*, "Internvid: A large-scale video-text dataset for multimodal understanding and generation," *arXiv:2307.06942*, 2023.
- [140] R. Zellers, J. Lu, X. Lu, Y. Yu, Y. Zhao, M. Salehi, A. Kusupati, J. Hessel, A. Farhadi, and Y. Choi, "Merlot reserve: Neural script knowledge through vision and language and sound," in *CVPR*, 2022, pp. 16375–16387.
- [141] Y. Zhai, S. Tong, X. Li, M. Cai, Q. Qu, Y. J. Lee, and Y. Ma, "Investigating the catastrophic forgetting in multimodal large language models," *arXiv:2309.10313*, 2023.
- [142] S. Chen, W. Jiang, W. Liu, and Y.-G. Jiang, "Learning modality interaction for temporal sentence localization and event captioning in videos," in *ECCV*, 2020, pp. 333–351.
- [143] H. Wang, Z.-J. Zha, L. Li, D. Liu, and J. Luo, "Structured multi-level interaction network for video moment localization via language query," in *CVPR*, 2021, pp. 7026–7035.
- [144] J. Wang, L. Ma, and W. Jiang, "Temporally grounding language queries in videos by contextual boundary-aware prediction," in *AAAI*, 2020, pp. 12168–12175.
- [145] Y.-W. Chen, Y.-H. Tsai, and M.-H. Yang, "End-to-end multi-modal video temporal grounding," in *NeurIPS*, 2021, pp. 28442–28453.
- [146] S. Chen, Y. Wu, C. Wang, S. Liu, D. Tompkins, Z. Chen, and F. Wei, "Beats: Audio pre-training with acoustic tokenizers," *arXiv:2212.09058*, 2022.
- [147] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *NAACL*, 2019, pp. 4171–4186.
- [148] J. Gao, C. Sun, Z. Yang, and R. Nevatia, "Tall: Temporal activity localization via language query," in *ICCV*, 2017, pp. 5267–5275.
- [149] L. Zhou, C. Xu, and J. Corso, "Towards automatic learning of procedures from web instructional videos," in *AAAI*, 2018, pp. 7590–7598.
- [150] G. A. Sigurdsson, G. Varol, X. Wang, A. Farhadi, I. Laptev, and A. Gupta, "Hollywood in homes: Crowdsourcing data collection for activity understanding," in *ECCV*, 2016, pp. 510–526.
- [151] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Nibbles, "Activitynet: A large-scale video benchmark for human activity understanding," in *CVPR*, 2015, pp. 961–970.
- [152] J. Xiao, X. Shang, A. Yao, and T.-S. Chua, "Next-qa: Next phase of question-answering to explaining temporal actions," in *CVPR*, 2021, pp. 9777–9786.
- [153] R. Hu, H. Xu, M. Rohrbach, J. Feng, K. Saenko, and T. Darrell,

- “Natural language object retrieval,” in *CVPR*, 2016, pp. 4555–4564.
- [154] S. Fujita, T. Hirao, H. Kamigaito, M. Okumura, and M. Nagata, “Soda: Story oriented dense video captioning evaluation framework,” in *ECCV*, 2020, pp. 517–531.
 - [155] S. Banerjee and A. Lavie, “Meteor: An automatic metric for mt evaluation with improved correlation with human judgments,” in *ACL Workshop*, 2005, pp. 65–72.
 - [156] R. Vedantam, C. Lawrence Zitnick, and D. Parikh, “Cider: Consensus-based image description evaluation,” in *CVPR*, 2015, pp. 4566–4575.
 - [157] W. Liu, T. Mei, Y. Zhang, C. Che, and J. Luo, “Multi-task deep visual-semantic embedding for video thumbnail selection,” in *CVPR*, 2015, pp. 3707–3715.
 - [158] Y. Zhong, W. Ji, J. Xiao, Y. Li, W. Deng, and T.-S. Chua, “Video question answering: Datasets, algorithms and challenges,” in *EMNLP*, 2022, pp. 6439–6455.
 - [159] X. Wang, Q. Si, J. Wu, S. Zhu, L. Cao, and L. Nie, “Retake: Reducing temporal and knowledge redundancy for long video understanding,” *arXiv:2412.20504*, 2024.
 - [160] Y. Li, C. Wang, and J. Jia, “Llama-vid: An image is worth 2 tokens in large language models,” in *ECCV*, 2024, pp. 323–340.
 - [161] I. Viertola, V. Iashin, and E. Rahtu, “Temporally aligned audio for video with autoregression,” in *ICASSP*, 2025, pp. 1–5.
 - [162] G. Yariv, I. Gat, S. Benaim, L. Wolf, I. Schwartz, and Y. Adi, “Diverse and aligned audio-to-video generation via text-to-video model adaptation,” in *AAAI*, 2024, pp. 6639–6647.
 - [163] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, “Robust speech recognition via large-scale weak supervision,” in *ICML*, 2023, pp. 28 492–28 518.