

ARC-Chapter: Structuring Hour-Long Videos into Navigable Chapters and Hierarchical Summaries

Junfu Pu*, Teng Wang*, Yixiao Ge†, Yuying Ge, Chen Li, Ying Shan

ARC Lab, Tencent PCG

*Core contributors, †Project lead

The proliferation of hour-long videos (e.g., lectures, podcasts, documentaries) has intensified demand for efficient content structuring. However, existing approaches are constrained by small-scale training with annotations that are typically short and coarse, restricting generalization to nuanced transitions in long videos. We introduce ARC-Chapter, the first large-scale video chaptering model trained on over million-level long video chapters, featuring bilingual, temporally grounded, and hierarchical chapter annotations. To achieve this goal, we curated a bilingual English-Chinese chapter dataset via a structured pipeline that unifies ASR transcripts, scene texts, visual captions into multi-level annotations, from short title to long summaries. We demonstrate clear performance improvements with data scaling, both in data volume and label intensity. Moreover, we design a new evaluation metric termed GRACE, which incorporates many-to-one segment overlaps and semantic similarity, better reflecting real-world chaptering flexibility. Extensive experiments demonstrate that ARC-Chapter establishes a new state-of-the-art by a significant margin, outperforming the previous best by 14.0% in F1 score and 11.3% in SODA score. Moreover, ARC-Chapter shows excellent transferability, improving the state-of-the-art on downstream tasks like dense video captioning on YouCook2.

Date: November 18, 2025

Github: <https://github.com/TencentARC/ARC-Chapter>

1 Introduction

The exponential proliferation of long-form video content, including educational lectures, vlogs, live streams, and meeting recordings—poses significant challenges for automatic content understanding. Video chaptering [35; 44] has emerged as a promising solution, segmenting videos into navigable and semantically coherent chapters. This enables efficient content retrieval, summarization, and enhanced user interaction, which are critical for managing and consuming large-scale video data.

Despite notable advances in segmenting short videos (usually within five minutes) for tasks such as action segmentation [8; 22; 27; 32; 39], temporal event localization [16; 54], and dense video captioning [19; 38; 46], the structuring of hour-long videos remains a formidable challenge. First, modeling sophisticated semantics across multimodal inputs, including visual and audio streams—over extended temporal horizons requires robust and scalable architectures. Second, the scarcity of large-scale datasets with fine-grained annotations hinders the development and evaluation of effective chaptering models. Third, existing evaluation metrics [10; 19] often fail to capture the semantic granularity of chapter boundaries, leading to suboptimal matching and similarity scoring between predicted and ground-truth segments [10].

In this technical report, we introduce ARC-Chapter, a comprehensive framework designed to address the unique challenges of long-form video structuring. As illustrated in Fig. 1, ARC-Chapter enables the segmentation of lengthy videos into navigable chapters and generates hierarchical summaries that capture both coarse and fine-grained content structure. Our work makes three primary contributions. First, we advance the scalability of video chaptering by developing the first large-scale model trained on one million long videos, totaling 400,000 hours of content. This dataset is fifty times larger than those used in previous studies [35], allowing our model to generalize across diverse video domains and formats. Second, we propose a semi-automatic annotation pipeline for hierarchical summaries, which leverages easily accessible human-annotated coarse labels. This pipeline integrates automatic speech recognition (ASR) derived transcripts with timestamped visual

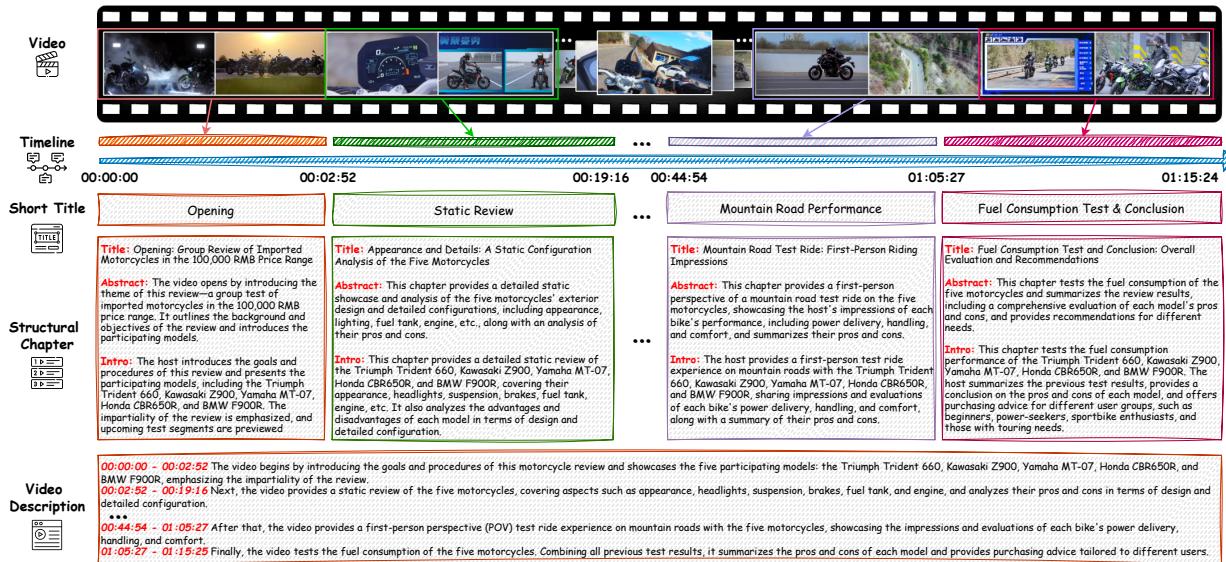


Figure 1 An illustration of the capabilities of our video chaptering model. Given a video, our model is able to generate timestamped chapters with three-level structured output: 1) **Short Title** - a concise label summarizing each chapter; 2) **Structural Chapter** - a detailed, structured annotation for each chapter, including a rewritten comprehensive Title, an abstract summarizing the core content, and an introduction describing key details and highlights; and 3) **Timestamp-Aligned Video Description** - fine-grained descriptions aligned with precise temporal boundaries. This hierarchical structure facilitates an efficient and precise understanding of video content.

elements, enabling a holistic and multimodal understanding of video content. Third, we introduce GRACE, a novel granularity-robust evaluation metric designed to address the semantic misalignment issues prevalent in existing chaptering benchmarks. GRACE provides a more accurate assessment of chapter boundary quality by accounting for varying levels of semantic granularity.

Our extensive experiments demonstrate the effectiveness of ARC-Chapter, which establishes a new state-of-the-art on both Chinese and English long-form video chaptering benchmarks. Specifically, ARC-Chapter substantially outperforms previous methods on the VidChapters-7M test sets (e.g., CIDEr: $100.9 \rightarrow 186.6$; F1: $45.3 \rightarrow 59.3$; SODA: $19.3 \rightarrow 30.6$). We validate the importance of multimodality, showing that our full model surpasses video-only and audio-only variants by 7.7 and 5.3 points on SODA, respectively. Furthermore, pretraining on our large-scale dataset significantly enhances transferability, evidenced by notable performance gains on downstream tasks like YouCook2 and ActivityNet Captions. Crucially, our work is the first to identify a clear scaling law in video chaptering: model performance consistently improves with increased training data and label density. This finding refutes previous observations that performance saturates on smaller datasets ($\sim 20k$ samples) [35] and suggests a promising direction for future research.

The remainder of this report is structured as follows: Section 2 reviews related works; Section 3 describes the dataset and annotation pipeline; Section 4 details our methodology and model architecture; Section 5 presents experimental results and analysis; Section 6 concludes.

2 Related Works

Global Video Understanding. Early video understanding [1; 7; 13; 23; 26; 33; 37; 41; 42; 49; 52; 53; 57] research primarily targeted global comprehension tasks, such as video question answering, video captioning, and video classification. These methods treat entire videos as holistic units, extracting global representations to predict semantic labels or generate summaries. While effective for short videos, they often fail to capture complex temporal dynamics and hierarchical structures of long-form content [24; 30].

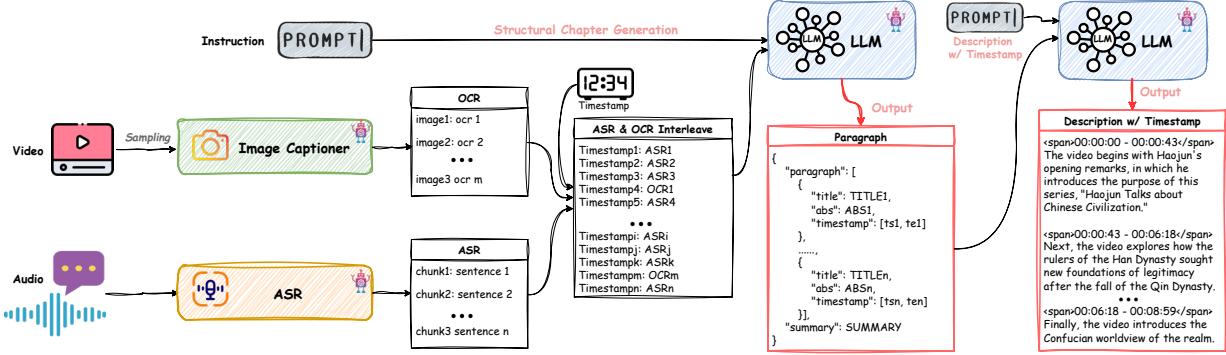


Figure 2 Overview of our automatic video annotation pipeline for hierarchical chaptering and summarization. We extract visual captions (OCR included) from sampled video frames and ASR transcripts from audio. These outputs are temporally aligned and interleaved into a unified multimodal transcript. This transcript, together with original chapter markers, is processed by an LLM to produce structured chapters and timestamp-aligned video descriptions.

Temporal Segmentation for Short Videos. To address the limitations of global approaches, recent works [14; 15; 17; 28; 30; 40; 47; 50; 56] have shifted towards modeling the temporal structure of videos. Datasets like ActivityNet Captions [19], Charades-STA [11], YouCook2 [55] and Breakfast [21] provide timestamped event annotations, enabling tasks such as temporal event localization, action segmentation, and dense video captioning. These approaches move beyond global representations to identify and describe fine-grained events and local temporal dependencies. However, most temporally-structured datasets [25; 48] are limited to short clips, typically under several minutes, and thus do not capture the challenges of ultra-long videos found in lectures, podcasts, or livestreams. The lack of large-scale, long-duration datasets with fine-grained temporal annotations remains a major bottleneck.

Long-Form Video Structuring. A few efforts [35; 45] have explored the structuring of hour-long videos. The VidChapters-7M dataset [45] provides a large-scale benchmark for video chaptering, with millions of videos and annotated chapter boundaries, better reflecting real-world scenarios such as vlogs, podcasts, and meetings where long-term temporal reasoning is essential.

Despite these advances, significant challenges remain. Existing chaptering models often rely on limited modalities, such as automatic speech recognition, are trained on small-scale datasets, and produce coarse, uninformative descriptions, which limits their scalability across diverse video domains. To address these issues, we propose a scalable, multimodal framework for long-form video chaptering, supported by a large-scale dataset with detailed chapter descriptions.

3 Data Collection and Annotation

A significant challenge in developing strong video chaptering models is the scarcity of publicly available datasets with detailed, multi-level annotations. Existing datasets typically provide only sparse labels, such as video-level categories for video classification or coarse temporal segments with brief titles such as VidChapters-7M. To address this limitation and to facilitate research on hierarchical video chaptering and summarization, we introduce a new, richly annotated video chaptering dataset. This section details our data curation and annotation pipeline.

3.1 Data Curation

One of the key contributions of our work is the introduction of a new large-scale dataset, named VidAtlas, which is designed for the task of hierarchical video chaptering and summarization. Our primary goal is to construct a dataset that not only provides accurate chapter boundaries but also offers dense, multi-granularity textual descriptions for both individual chapters and the entire video.

Data Sourcing. We begin by sourcing videos from the video platform. The primary selection criterion is the presence of author-provided chapter markers. These markers, which include the start/end timestamps and a short title for each chapter, are manually defined by the video uploader. This approach provides us with a highly accurate human-verified ground truth for the temporal segmentation of videos, which is a significant foundation for our subsequent annotation efforts. The collected videos, which are long, well-structured, and information-dense, are ideal candidates for video chaptering.

Filtering and Refinement. Starting with this initial collection, we apply several filtering criteria to guarantee the quality and diversity of our dataset for video understanding and chaptering. First, we retain videos whose durations lie between 2 minutes and 3 hours. This range excludes trivial short clips, which are unnecessary for chaptering, as well as overly long videos, which are often unstructured (e.g., live streams) and difficult to process due to the context-length limitations of our model. Second, we curate videos across a wide range of domains, including educational lectures, DIY tutorials, reviews & unboxings, interviews & podcasts, webinars & presentations, gaming & music albums, fitness & cooking and documentaries. This wide distribution of domains ensures that the dataset is not biased towards any specific genre and supports the development of more generalizable models.

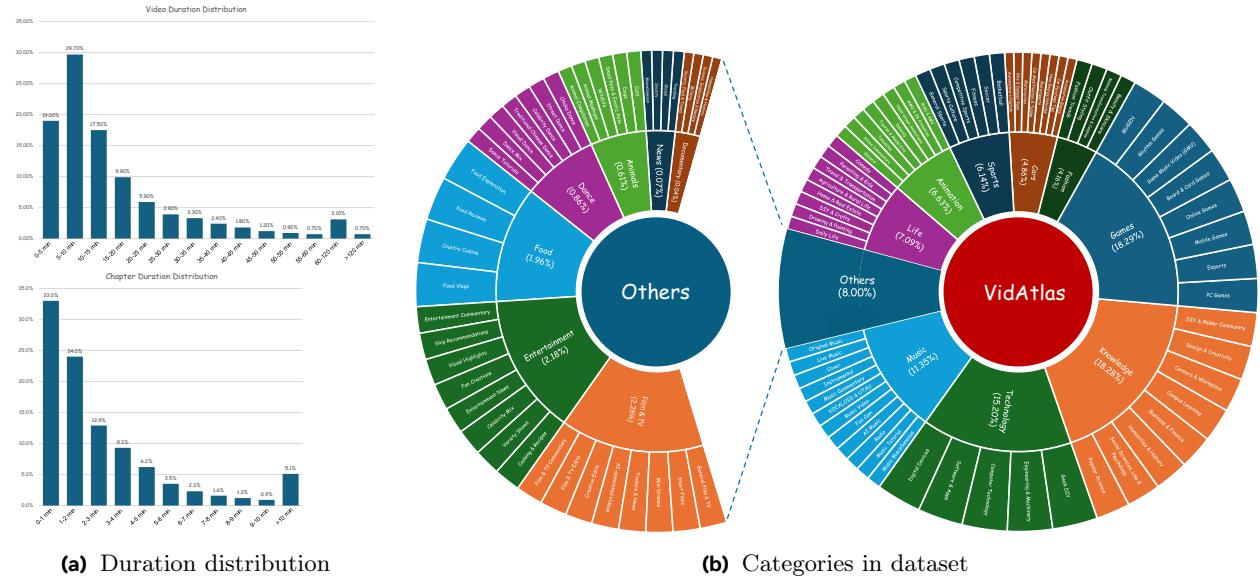


Figure 3 Dataset statistics: (a) Distribution of video durations (top) and chapter durations (bottom) in the VidAtlas dataset. (b) Distribution of video topics in VidAtlas.

3.2 Hierarchical Annotation

To generate high-quality video chaptering annotations, we design an automated annotation pipeline that leverages both multimodal content extraction and large language model (LLM)-based reasoning based on the videos with user-provided chapter markers, *i.e.* timestamps and brief title of each chapter. The illustration of our annotation pipeline is shown in Fig. 2.

Multimodal Information Extraction. Considering efficiency and cost, we avoid directly using multimodal large language models (MLLMs) for video annotation. Instead, we first extract multimodal information from video frames and audio, integrate this content, and then feed the result into text-only LLM for reasoning and annotation. Specifically, we use Whisper-v3 [29] to transcribe speech into text, segmented into sentences with the corresponding timestamps. In parallel, we uniformly sample video frames with a fixed sampling frame rate and employ Qwen2.5-VL-7B [4] to extract visual captions and on-screen text (OCR) for better understanding of the video content. Subsequently, the visual captions and ASR transcripts are temporally aligned based on their respective timestamps. This process allows us to interleave the textual content from both modalities into a unified chronologically ordered sequence. This multimodal transcript, together with

the original user-provided chapter timestamps and short titles, is fed into LLM for reasoning and structural segmentation.

LLM Reasoning and Chaptering. The LLM is prompted to analyze the transcript and reorganize the content into a structured set of chapters, each containing a comprehensive title, an abstract, an introduction, and precise temporal boundaries. Following this, we perform a verification step on the LLM’s output to ensure that the generated chapter boundaries strictly adhere to the original timestamps. Building upon the verified structured chapter information, we further prompt the LLM to produce a comprehensive, timestamped narrative description for the entire video. Through this annotation pipeline, we can efficiently obtain accurate, multi-level video chapter segmentation and descriptive annotations. The resulting annotations form a dense, hierarchically organized representation of long-form videos, supporting a wide range of research tasks in video understanding, temporal reasoning, chaptering, and summarization.

3.3 Dataset Statistics

We summarize the key statistics of our VidAtlas dataset and highlight the properties that make it suited for research on video chaptering and summarization. The dataset comprises 410k+ videos with an average duration of 16.8 minutes, amounting to more than 115k hours of diverse content. On average, each video is segmented into 5.5 chapters, with an average chapter duration of 182 seconds (approximately 3 minutes). Fig. 3a provides a detailed statistic of the duration distributions for both videos and chapters. Our dataset contains a wide spectrum of video and chapter lengths to ensure models are trained on a diverse temporal structures. This comprehensive video/chapter length distribution makes the models exposed to a variety of content length, from concise segments to hour-long narratives, forcing models to resolve both rapid topic shifts and sustained thematic segments. To mitigate genre bias, VidAtlas covers a wide array of topics, including 16 primary categories with over 100 subcategories, as shown in Fig. 3b. The categories of VidAtlas include Games, Knowledge, Technology, Music, Life, Animation, and Sports, together with other variety that captures long-tail topics. Videos in these categories are typically well-structured and information-dense, making them ideal for chaptering.

4 ARC-Chapter

4.1 Overall Framework

We leverage Qwen2.5-VL-7B [5] as our base model, enhancing its capabilities to process and structure video content into chapters. The architecture of our model is illustrated in Fig. 4. The model unifies three inputs: 1) an instruction prompt that specifies the task of input modalities and output schema. 2) a sequence of sampled video frames that provide appearance, layout and on-screen text (including subtitles which often align with the ASR transcript), and 3) a timestamp-aligned ASR transcript from audio. While both the video and ASR transcript inputs are optional, the model requires at least one modality to be provided. Frames are embedded with Qwen2.5-VL vision encoder and translated into visual tokens, while ASR transcript is tokenized as plain text with explicit timestamps. The vision encoder is kept frozen and the language model is instruction tuned on VidAtlas to specialize in video chaptering.

Prompt Design. The model’s behavior is guided by carefully designed prompts that specify the desired task and output format. To handle the diverse requirements of different inputs and outputs of the model, we design a set of 18 distinct prompt templates. These prompts are constructed based on three axes: language in source video, input modality, and desired output format.

- **Language:** We support English and Chinese to match the language of the source video.
- **Input Modality:** The prompt specifies whether the model should rely on ASR-only, video-only, or both video and ASR inputs. This allows for ablation studies and adaptation to scenarios where one modality may be absent or noisy.

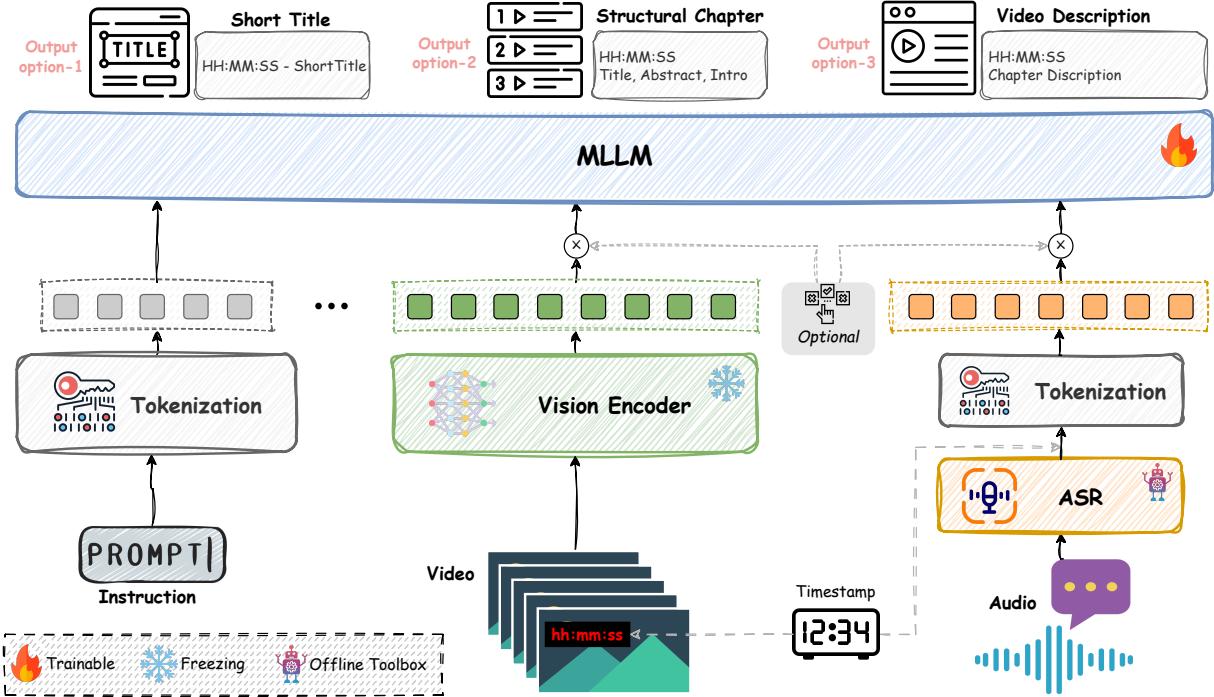


Figure 4 Overview of the model architecture for video chaptering. The model inputs include a task-specific prompt, sampled video frames, and timestamped ASR transcripts. Video frames are processed with a frozen vision encoder. The resulting visual features, along with the tokenized prompt and ASR text, are fed into a trainable multimodal large language model (MLLM). Based on the inputs, the model is able to generate chapters in various formats, including timestamped concise title, detailed structural chapters, or comprehensive video description with timestamps.

- **Output Format:** We define three distinct output structures: (a) *Short Titles* for concise chapter markers, (b) *Structured Chapters* that include a title, abstract, and introduction for each chapter, and (c) *Video Descriptions* that provide a dense, timestamp-aligned summary of the entire video.

Video Input. To balance temporal coverage and context budget, we follow the setup of Qwen2.5-VL and cap the visual stream at 768 frames sampled at up to 1 fps. That is to say, videos shorter than 12.8 minutes are sampled with 1 fps, while longer videos are uniformly down-sampled to 768 frames with a lower fps. The sampling strategy retains coarse global coverage for hour-long content, ensuring sufficient representation to capture the high-level semantic shifts necessary for the chaptering task. Since the model context length is shared across modalities, we dynamically adjust the per-frame token allowance according to the input of ASR transcript. For video-only inputs we use a higher frame resolution (higher token budget per frame) so that small text (OCR and subtitles) and fine-grained visual cues are preserved. When ASR is provided alongside video, we reduce frame resolution (thus reducing the number of visual tokens) so that the combined input of visual tokens and ASR text fits the maximum context length of MLLM. This dynamic allocation is implemented by adjusting image scaling and patch-tokenization parameters at preprocessing time. Moreover, to enhance temporal awareness, we randomly overlay timestamps onto the video frames, making the model more sensitive to the video timeline.

ASR Input. Although integrating raw audio features or learned audio embeddings from pretrained ASR models (e.g. Whisper [29]) is attractive, it presents severe scalability challenges for long-form video. For example, while Whisper-style audio encoder produces 50 audio tokens per second, a 60-minute audio therefore produces 180k tokens, far exceeding feasible LLM context budgets without aggressive compression or specialized audio-to-token aggregation. Furthermore, synchronizing fixed-rate audio features with dynamically sampled video frames poses an additional alignment problem. To address these practical constraints, we opt to use ASR transcripts as a highly effective proxy for the audio modality. Text is significantly more information-dense. Therefore,

the ASR transcript of a long audio segment occupies far fewer tokens than its raw feature representation. This makes processing hour-long videos computationally feasible for both training and inference. Although such a paradigm introduces an extra step for offline ASR transcription, we believe that trading a modest amount of offline processing time for the ability to handle long-form audio under strict context-length budgets is worthwhile. In our implementation, we use Whisper-large-v3 [29] to generate timestamped ASR transcripts. The model provides sentence-level segments with corresponding start timestamps. We formulate the ASR text and timestamp of each segment as *start time (hh:mm:ss)*: $\langle \text{ASR text} \rangle$. The normalized ASR transcript is then passed to the model either alone (ASR-only) or together with visual tokens (ASR+Video), providing dense semantic information that is particularly useful for temporal boundary detection and chaptering.

4.2 Training Strategy

Training Objective. We perform supervised instruction tuning on VidAtlas and VidChapter-7M using all prompt templates. The training objective is the standard autoregressive next-token prediction loss over the target sequence. Given a multimodal input sequence consisting of a prompt X_{prompt} , video frames X_{video} , and an ASR transcript X_{asr} (video stream X_{video} and ASR streams X_{asr} are optional), the model is trained to maximize the log-likelihood of the target output sequence $Y = (y_1, y_2, \dots, y_n)$ (e.g., a list of chapter titles, a structured chapter object, or a timestamped description):

$$\mathcal{L} = - \sum_{i=1}^n \log P(y_i | y_{<i}, X_{\text{prompt}}, X_{\text{video}}, X_{\text{asr}}),$$

where $y_{<i}$ represents the preceding ground-truth tokens. During training, the vision encoder is frozen to enable a larger context length, while all parameters of the large language model are optimized with the training objective.

Adaptive Modality Dropping. To enable a single model to perform well under various deployment conditions, we adopt an adaptive modality dropping strategy during training. For each training sample, we randomly configure the input with a certain probability to be one of three types: 1) **Video + ASR**: Both modalities are provided to the model. 2) **Video-only**: The ASR transcript is omitted, forcing the model to rely solely on visual information. and 3) **ASR-only**: The video frames are omitted, requiring the model to understand the content based on the transcript alone. This strategy prevents the model from becoming overly reliant on a single modality and ensures it develops a comprehensive understanding from all available input modalities. Consequently, a single trained model can be deployed to handle videos under various conditions during inference (whether only a video is available, only transcript is provided, or both are present), without requiring specialized models for each scenario.

4.3 Evaluation Metrics

Evaluation metrics can be divided into two aspects: (1) the accuracy of segmentation (e.g., Precision, Recall, and tIOU [20]), and (2) joint metrics that assess both segmentation and chapter captioning (e.g., CIDEr [20], SODA [10]). However, we observe that the primary metrics such as SODA, originally developed for dense video captioning, are not well-suited for the video chaptering task. While SODA enforces a one-to-one matching between predicted and ground-truth events to suppress redundancy in overlapping event detection, video chaptering requires segmenting videos into sequential, non-overlapping chapters. Furthermore, chaptering annotations often exhibit granularity ambiguity: different annotators may segment the same video at varying levels of detail—some may annotate coarse-grained chapters (e.g., by day in a travel vlog), while others may provide fine-grained chapters (e.g., by each visited site within a day). This results in multiple valid annotation granularities for the same content.

To address these challenges, we propose GRACE, a metric tailored for video chaptering. It introduces a many-to-one (set-to-one) matching paradigm, allowing each ground-truth (predicted) chapter to be matched with a set of predicted (ground-truth) chapters. As illustrated in Fig. 5, for each ground-truth chapter, GRACE evaluates the temporal overlap and semantic similarity between the chapter and its matched prediction set, using established language similarity metrics (e.g., BERTscore [51]) for textual comparison. Specifically, we aim to find a best many-to-one mapping M which splits both ground-truth set G and prediction set P into several pairs of groups $\{(P_i, G_i)\}_{i=1}^K$, followed by group-based similarity calculation:

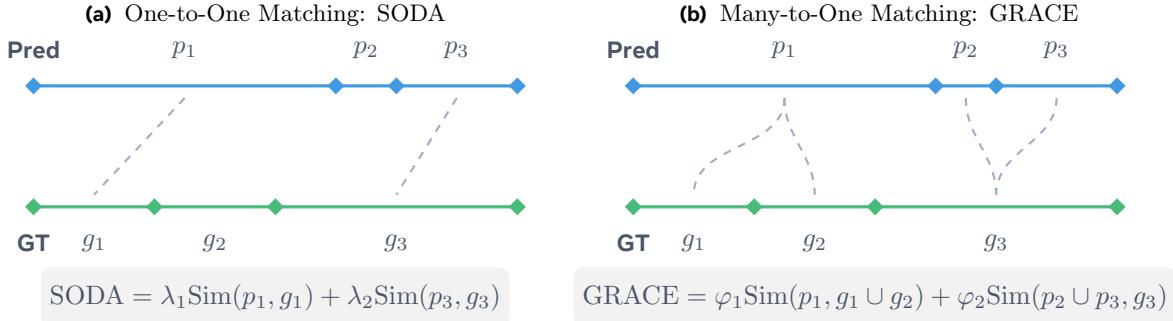


Figure 5 Comparison of one-to-one (SODA) and many-to-one (GRACE) matching strategies. The one-to-one matching can fail to account for important events like p_2 and g_2 , whereas the many-to-one strategy considers all predicted and ground-truth events for a more robust, overall assessment.

$$\text{GRACE} = \sum_{(P_i, G_i) \in M(P, G)} \varphi(P_i, G_i) \cdot \text{BERTscore}(P_i, G_i) \quad (1)$$

$$\varphi(P_i, G_i) = \frac{1}{|P_i||G_i|} \sum_{p \in P_i, g \in G_i} \text{IOU}(p, g) \quad (2)$$

$$\text{s.t. } P_i \cap P_j = \emptyset, \cup(P_i) = P, G_i \cap G_j = \emptyset, \cup(G_i) = G, \min(|P_i|, |G_i|) = 1 \quad (3)$$

where P_i and G_i represent groups of chapters. When calculating the BERTScore between two groups, we first concatenate all captions within each group into a single sentence, then compute the BERTScore between the two merged sentences. We adopt the dynamic time warping algorithm (DTW) [6; 31] to achieve the optimal matching $M(P, G)$, with IOU between two chapters being used as the matching criteria.

GRACE provides a more accurate and human-aligned assessment of chaptering models. This design confers several advantages: (1) robustness to annotation granularity, enabling fair evaluation across diverse annotation styles; (2) improved semantic fidelity, rewarding models that capture the full scope of ground-truth chapters; and (3) closer alignment with human judgment of chapter boundaries and content.

4.4 Reinforcement Learning with GRPO

While supervised fine-tuning (SFT) achieves strong performance, the standard cross-entropy loss does not directly optimize for the primary objective of video chaptering: temporal accuracy. To further enhance the model’s temporal localization capabilities, we introduce a subsequent reinforcement learning phase using the GRPO algorithm [12].

The core of this phase is a reward function designed to directly incentivize precise chapter boundary prediction. We leverage our proposed GRACE metric, which holistically evaluates both temporal alignment and semantic content. However, to specifically sharpen the model’s ability to predict accurate timestamps of segmented chapters, we formulate a simplified, temporal-only reward by omitting the semantic BERTscore component from Equation (1). For a given ground-truth chapter set G and a model-generated set P , the reward R is calculated by summing the temporal alignment scores φ over the optimal matching $M(P, G)$ found via DTW:

$$R = \sum_{(P_i, G_i) \in M(P, G)} \varphi(P_i, G_i). \quad (4)$$

This reward directly reflects the quality of the temporal segmentation, providing a clear and targeted optimization objective.

Due to the significant context length required for multimodal inputs, and to specifically bolster the model’s ability to reason from visual cues, we conduct this RL training phase using only the video modality. We select a diverse subset of 90k videos from both Chinese and English SFT data, ensuring that training samples cover all three output formats: short titles, structural chapters, and timestamped video description. We initialize the model with the weights from our best-performing SFT model and further optimize it using GRPO. The KL divergence coefficient is set to 0.01 to ensure that the policy does not stray far from the robust language generation capabilities learned during SFT, thereby balancing temporal refinement with descriptive quality.

Table 1 Comparison to the state of the art on VidChapters7M-test set: The results of compared methods are evaluated in the ASR-only setting from Chapter-Llama [35]. We evaluate ARC-Chapter with different input modalities: **-vid** for video, **-asr** for ASR, and **-vidasr** for both. “**Ft.**” indicates whether the model is finetuned for chaptering task. †denotes LLM-API results reported from Chapter-Llama. Our model, ARC-Cchapter, achieves the best performance across all metrics and video durations.

Backbone	Ft.	Short				Medium				Long				All			
		F1	tIoU	S	C												
GPT-4o-mini [18]†	X	32.1	64.5	7.2	42.4	30.5	62.3	6.1	30.6	28.0	61.0	6.0	27.3	31.2	63.6	6.8	37.8
GPT-4o [18]†	X	37.7	68.0	8.4	53.8	38.1	68.8	8.1	51.4	36.5	66.2	6.6	34.8	37.6	68.0	8.1	51.0
Gemini-2.0-Flash [34]†	X	39.9	69.2	12.0	72.8	43.8	71.4	11.2	70.3	34.9	66.2	9.0	51.6	40.2	69.3	11.4	69.7
Gemini-1.5-Pro [34]†	X	41.7	70.6	11.7	65.3	43.8	71.8	11.2	61.4	41.3	70.6	10.1	55.3	42.2	70.9	11.4	63.2
Vid2Seq [45; 46]	X	2.5	28.6	0.3	0.3	3.2	29.7	0.3	0.4	4.6	32.0	0.3	0.5	3.0	29.3	0.3	0.4
Llama 3.1-8B [9]	X	29.9	63.4	7.1	34.5	30.6	62.7	5.4	28.1	26.6	59.3	3.6	18.9	29.5	62.5	6.2	30.7
Vid2Seq [45; 46]	✓	33.4	63.7	15.2	74.9	19.0	53.3	7.5	31.9	16.7	50.8	5.9	28.4	26.7	58.6	11.6	55.8
Chapter-Llama [35]	✓	45.5	72.2	20.2	103.5	46.7	72.3	18.8	98.7	41.3	69.2	15.8	91.2	45.3	71.8	19.3	100.9
ARCChapter-asr¹	✓	54.5	76.7	26.3	144.1	55.9	77.5	25.1	143.0	55.1	77.0	24.8	158.0	54.5	76.7	25.3	144.0
ARCChapter-vid	✓	52.6	75.8	26.0	156.8	51.4	75.3	20.6	124.0	47.3	72.3	19.2	119.8	50.2	74.3	22.9	138.3
ARCChapter-vidasr	✓	60.0	80.1	32.5	195.7	59.2	79.4	29.6	177.3	60.2	79.9	29.2	190.3	59.3	79.6	30.6	186.6

Table 2 Comparison to the state of the art on VidChapter7M-sml300 with different input modalities. Our method, ARC-Chapter, demonstrates superior performance on VidChapter-sml300 by effectively integrating both speech and video information. The modalities of “Embed” and “Caption” in LLaMA and Chapter-LLaMA models play the same role as “Video” in ARC-Chapter model.

Method	Ft.	Modalities			Segmentation		Titles	
		Speech	Embed.	Caption	F1	tIoU	S	C
LLaMA 3.1-8B	X	X	X	✓	12.6	48.6	1.9	6.4
	X	✓	X	X	22.7	57.3	4.4	19.7
	X	✓	X	✓	29.9	63.0	6.9	33.7
Chapter-LLaMA	✓	✓	X	X	38.5	68.1	13.9	67.3
	✓	X	✓	X	38.4	66.5	3.4	7.3
	✓	X	X	✓	39.1	67.7	5.9	20.2
	✓	✓	✓	X	40.4	68.2	15.3	74.9
	✓	✓	X	✓	42.6	70.6	16.4	82.4
	✓	✓	✓	✓	44.4	71.5	16.3	84.2
		Speech	Video		F1	tIoU	S	C
ARCChapter	✓	✓	X		56.5	78.1	25.9	148.5
	✓	X	✓		50.0	74.3	21.6	130.8
	✓	✓	✓		62.4	81.6	30.1	190.7

5 Experiments

In this section, we conduct a series of experiments to thoroughly evaluate our video chaptering model. We first introduce the evaluation benchmarks, then present the main results and detailed ablation studies.

5.1 Evaluation Benchmark

To comprehensively assess our model’s capabilities in video chaptering, we evaluate it on three distinct benchmarks covering different languages, scales, and data modalities. The evaluation targets two key criteria: the precision of temporal boundary localization and semantic relevance of the generated chapter titles/descriptions. VidChapters7M is a large-scale English chaptering dataset. We use two of its standard splits for evaluation, i.e., VidChapters7M-test and VidChapters7M-sml300val. VidChapters7M-test is a large-scale test set comprising 8.2k samples. For this split, the compared methods are only based on ASR

¹For convenience, “ARC-Chapter” in the main text is abbreviated as “ARCChapter” in all experimental result tables.

Table 3 Comparison to the state of the art on VidAtlas-test set: “Ft.” indicates whether the model is finetuned for chaptering task. Modality‡ specifies which inputs are provided: A for ASR and V for video. † denotes LLM-API results. For API-base models, the video is converted into a textual description, which is then provided as input for LLM.

Backbone	Ft.	Modality‡		Short				Medium				Long				All				
		A	V	F1	tIoU	S	C	G												
Claude-Sonnet [3]†	X	✓	X	39.2	69.8	7.6	38.8	34.7	66.3	6.5	33.8	36.6	66.9	5.8	33.5	37.8	68.6	7.1	36.9	11.1
Doubaο-1.5-Pro [13]†	X	✓	X	38.8	70.4	7.4	40.6	35.8	68.4	6.9	38.3	36.1	67.1	3.2	17.4	37.7	69.5	6.7	36.4	9.8
DeepSeek-R1 [12]†	X	✓	X	40.0	71.1	11.0	48.8	37.9	69.5	9.6	45.2	35.7	66.8	6.3	28.3	38.9	70.1	10.0	44.8	13.4
Gemini-2.5-Pro [7]†	X	✓	X	39.6	68.3	8.1	44.6	30.6	60.1	6.3	37.4	34.0	60.2	9.9	54.0	45.2	73.2	9.7	53.5	14.9
GPT-4.1 [2]†	X	✓	X	36.5	68.6	6.6	34.6	33.0	66.1	5.8	32.4	36.0	66.3	5.9	33.0	35.7	67.7	6.3	33.9	-
Qwen-3-235B [43]†	X	✓	X	36.7	67.7	7.7	36.9	33.5	65.6	6.7	33.9	26.6	61.0	3.8	18.7	34.4	66.2	6.9	33.4	10.2
Claude-Sonnet [3]†	X	✓	✓	36.8	68.2	7.9	42.4	32.0	65.2	8.0	45.0	40.8	68.2	16.8	110.4	36.4	67.5	9.3	53.6	13.2
Doubaο-1.5-Pro [13]†	X	✓	✓	39.5	70.0	7.7	43.3	35.5	67.6	7.6	45.2	44.4	69.8	14.9	109.0	39.5	69.4	8.8	54.1	12.6
DeepSeek-R1 [12]†	X	✓	✓	39.4	69.9	10.5	50.0	38.0	68.7	10.8	54.9	62.2	80.3	48.2	264.4	41.1	70.5	13.9	69.7	17.1
Gemini-2.5-Pro [7]†	X	✓	✓	48.3	73.1	9.8	54.9	45.4	70.1	11.8	66.1	54.8	75.3	30.6	172.5	48.7	72.8	13.5	75.8	19.8
GPT-4.1 [2]†	X	✓	✓	35.3	67.2	6.3	34.2	30.8	64.2	6.2	34.8	43.9	69.2	19.1	120.2	35.8	66.9	8.3	47.9	11.7
Qwen-3-235B [43]†	X	✓	✓	24.8	59.2	6.5	31.9	19.5	52.9	5.6	28.2	27.5	57.8	16.0	92.9	24.1	57.7	7.8	40.7	9.6
ARCChapter-asr	✓	✓	X	57.3	79.3	24.1	103.3	60.1	80.8	24.5	113.5	63.2	79.5	28.1	140.6	58.8	79.7	24.8	111.3	28.0
ARCChapter-vid	✓	X	✓	57.1	79.1	21.2	91.5	55.9	78.2	18.4	88.2	62.0	79.4	27.9	137.8	57.6	78.9	21.6	98.1	25.0
ARCChapter-vidasr	✓	✓	✓	65.5	83.8	28.5	129.2	65.7	84.2	29.0	140.0	69.6	84.2	38.5	192.3	66.2	84.0	30.2	141.5	34.1

transcripts, while ARC-Chapter is evaluated with different input modalities. VidChapters7M-sml300val is a smaller validation set of 300 samples, which includes both the original videos and their corresponding ASR transcripts. This subset is ideal for fast evaluation and conducting modality ablation studies. To assess generalization beyond English, we additionally report experimental results on VidAtlas-test, a Chinese test set with more than 1.5k videos together with ASR transcripts and original videos.

5.2 Comparison with the State of the Art

Performance on VidChapters7M. As shown in Tab. 1, our ARC-Chapter significantly outperforms all existing methods on VidChapters7M-test benchmark. Our model achieves a new state-of-the-art result in the ASR-only regime, with an overall F1 score of 54.5, tIoU of 76.7, SODA of 23.5, and a CIDEr of 144.0. This represents a substantial improvement over the previous SOTA model, Chapter-Llama, with absolute gains of +9.2 in F1, +4.9 in tIoU, and +6.0 in the SODA score. Notably, the performance gain enlarges as video duration increases. For long videos (30-60 min), the evaluation metrics of SODA and CIDEr for ARC-Chapter are remarkably higher than which in Chapter-LLama, demonstrating the superior capability of our model in processing long videos. Even when compared against powerful general models like GPT-4o and Gemini-1.5-Pro, which are not finetuned on this task, ARC-Chapter perform much better. The experiments conducted on VidChapter7M-sml300 show more comparisons for different input modalities, shown in Tab. 2.

Performance on VidAtlas. As detailed in Tab. 3, we evaluate our model on the VidAtlas benchmark under three settings: ASR-only, video-only, and ASR+video. ARC-Chapter consistently establish a new state-of-the-art across all settings. Our full multimodal model, ARCChapter-vidasr, which leverages both ASR and video inputs, achieves an overall F1 score of 66.2, tIoU of 84.0, SODA of 30.2, CIDEr of 141.5, and GRACE of 34.1. This marks a significant leap over the strongest LLM, Gemini-2.5-Pro, with an absolute improvement of +17.5 in F1 score and more than doubling the SODA score (+16.7). Furthermore, our single-modality versions also demonstrate superior performance. The ASR-only model, ARCChapter-asr, achieves an F1 of 58.8, and the video-only model, ARCChapter-vid, scores an F1 of 57.6. From shot-to-long videos, our model consistently outperforms other models, demonstrating its robustness in handling extended content.

5.3 Transferability

To evaluate transferability, we pre-trained ARC-Chapter on our dataset before fine-tuning and testing it on the dense video captioning benchmarks, i.e., Youcook2 and ActivityNet Captions. As shown in Table 4, our model establishes a new state-of-the-art, significantly outperforming all prior MLLM-based methods.

Notably, for event segmentation ability, ARC-Chapter achieves an F1/SODA Score of 37.9/12.5 on YouCook2, a substantial improvement over the previous best of 33.5/7.9. This demonstrates that the knowledge acquired during pre-training effectively transfers and enhances performance on downstream tasks.

Table 4 Transferability Performance on YouCook2 and ActivityNet Captions [20] for Dense Video Captioning. All methods use visual modality as inputs without ASR. The **Rank(\downarrow)** column represents the overall performance, calculated as the arithmetic mean of a method’s rank across all reported metrics (M, S, C, and F1) for that dataset. Some results for ActivityNet Captions are sourced from [14] and [46]. * indicates zero-shot evaluation. The best results on each dataset are in **bold** and the second-best are underlined.

Method	YouCook2					ActivityNet Captions				
	M	S	C	F1	Rank \downarrow	M	S	C	F1	Rank \downarrow
GIT [36]	3.4	3.1	12.1	17.7	7.5	7.8	5.7	29.8	50.6	4.3
ECHR [46]	3.8	-	-	-	4.0	7.2	3.2	14.7	-	8.6
PDVC [46]	4.7	4.4	22.7	-	5.0	8.0	5.4	29.0	56.7	3.8
Vid2Seq [46]	<u>9.3</u>	<u>7.9</u>	<u>47.1</u>	27.3	2.8	<u>8.5</u>	5.8	<u>30.1</u>	52.4	<u>2.6</u>
CM ²	-	5.3	31.7	28.4	4.7	-	-	-	-	-
TimeChat [30]	-	3.4	11.0	19.5	8.0	5.7	4.7	19.0	36.9	8.8
VTimeLLM [17]	-	-	-	-	-	6.8	5.8	27.6	-	5.8
Momentor* [28]	-	-	-	-	-	4.7	2.3	14.9	-	10.7
TRACE [14]	-	6.7	35.5	31.8	3.7	6.4	<u>6.0</u>	25.9	39.3	5.8
VTG-LLM [15]	-	3.6	13.4	20.6	6.7	5.9	5.1	20.7	34.8	8.3
TimeExpert [47]	-	7.2	39.0	33.5	<u>2.7</u>	7.0	6.5	28.4	40.5	4.3
ARC-Chapter	9.6	12.5	69.4	37.9	1.0	<u>8.1</u>	5.9	35.4	55.9	2.0

5.4 Ablation Studies

5.4.1 Scaling Property

We analyze how ARC-Chapter scales with the amount of training data. Concretely, we subsample the training set at 20%, 40%, 60%, 80%, and 100% and keep the model architecture and prompt templates fixed. We evaluate three inference modalities, i.e. ASR-only, Video-only, and ASR+Video, on two benchmarks: VidChapters-7M (sml300val) and a sampled subset of the VidAtlas-testset for efficiency. As illustrated in Fig. 6, the performance across all metrics (F1, tIOU, SODA, and CIDEr) and input modalities (ASR-only, Video-only, Video+ASR) demonstrates a clear positive correlation with the amount of training data. Specifically, the full multimodal model (Video+ASR) consistently achieves the best performance. ARC-Chapter is highly data-efficient, achieving strong performance with as little as 20% of the training data. Furthermore, it is data-scalable, continuing to benefit from larger corpora for even better results.

5.4.2 Hierarchical Annotations

A core contribution of our work is the VidAtlas dataset, which features rich, hierarchical annotations. To validate the effectiveness of this data structure, we evaluate our model’s capability to generate outputs of varying complexity, from simple *Short Title* to detailed *Structural Info* which comprising a title, abstract and introduction for each chapter. The results are presented in Table 5. From the experimental results, our model successfully learns to generate these complex, structured outputs, achieving strong performance across all generated components (title, abstract, introduction) on both VidChapter-sml300 and VidAtlas-testset benchmarks, particularly when using both video and ASR inputs. This demonstrates a high degree of semantic understanding.

More importantly, the capability for detailed generation does not come at the cost of performance on the fundamental chaptering task. When comparing the segmentation metrics (temporal evaluation score F1 and tIoU) for the *Short Title* task versus the more demanding *Structural Info* task, we observe only a negligible difference. For example, on VidChapter-sml300, the multimodal model achieved an F1 score of 62.4 and a tIoU of 81.6 for *Short Title* generation, compared to slightly lower scores of 61.4 and 80.6 for *Structural Info* generation. Notably, this small margin represents the largest performance gap observed across all modality inputs on both benchmarks, indicating that the model can perform complex, multi-part generation in a single forward pass without compromising its core ability to accurately segment the video. This result strongly validates our hierarchical annotation strategy, demonstrating that training on such rich data endows the model with advanced structural reasoning capabilities.

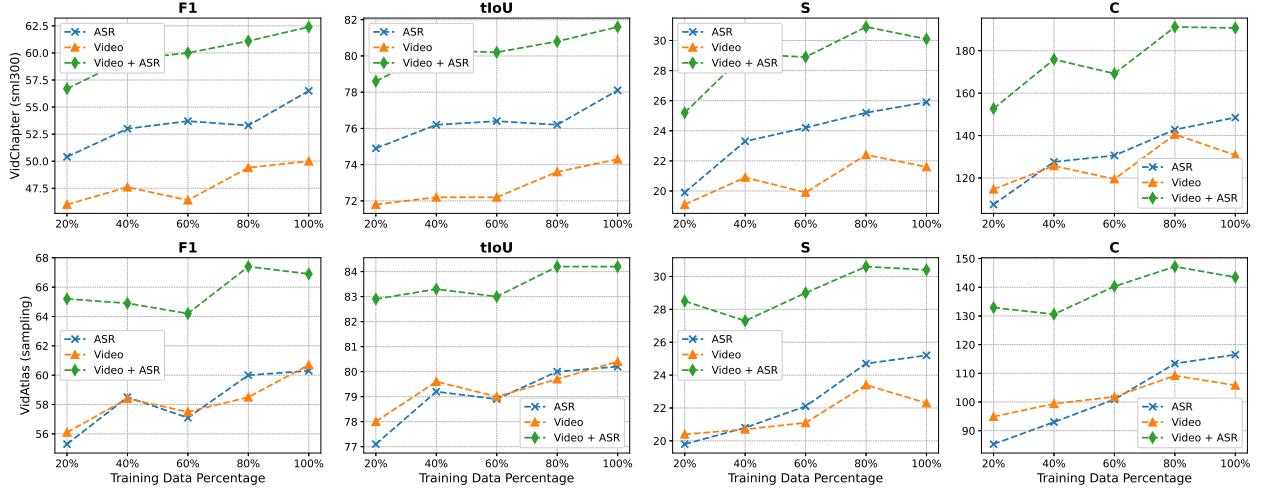


Figure 6 Data Scaling property of ARC-Chapter. We report the performance on VidChapter (a sampled subset) and VidAtlas test set with respect to different percentage of training samples.

Table 5 Ablation study on the model’s capability to generate hierarchical annotations. We compare models trained with *Short Title* and *Structural Info* (structured chapters with short title, title, abstract, and introduction) across different input modalities (A for ASR, and V for Video) on both English (VidChapter-sml300) and Chinese (VidAtlas-testset) benchmarks. Metrics include F1 and tIoU for boundary quality evaluation, and SODA(S), CIDEr(C), as well as our proposed GRACE(G) for semantic quality evaluation.

Dataset	Modality	Short Title						Structural Info														
		Segmentation		Short Title				Segmentation		Short Title				Title			Abstract			Intro		
		A	V	F1	tIoU	S	C	G	F1	tIoU	S	C	G	S	C	G	S	C	G	S	C	G
VidChapter-sml300 (English)	✓ X	56.5	78.1	25.9	148.5	33.0		54.8	77.1	25.5	147.9	32.5	12.8	91.2	25.6	12.3	14.5	25.1	11.8	11.9	24.6	
	X ✓	50.0	74.3	21.6	130.8	27.9		50.4	74.4	22.3	136.4	28.7	8.6	57.7	19.8	8.5	6.4	19.7	8.2	5.2	19.4	
	✓ ✓	62.4	81.6	30.1	190.7	38.4		61.4	80.6	30.8	194.5	38.4	14.6	107.2	28.6	13.4	14.5	27.4	13.0	10.2	27.0	
VidAtlas-testset (Chinese)	✓ X	58.8	79.7	24.8	111.3	28.0		59.1	79.8	25.5	112.8	28.6	16.2	101.7	27.0	17.5	57.8	31.8	16.4	36.0	29.6	
	X ✓	57.6	78.9	21.6	98.1	25.0		56.8	78.7	22.0	97.8	25.1	12.7	67.4	21.7	14.5	37.5	27.3	13.8	22.2	25.2	
	✓ ✓	66.2	84.0	30.2	141.5	34.1		65.9	83.8	30.8	143.5	34.6	18.5	119.8	30.7	19.1	66.3	35.3	18.2	39.8	33.0	

5.4.3 Performance with GRPO

To validate the effectiveness of our GRPO-based reinforcement learning stage, we compare the performance of our models before (SFT-base) and after (+RL) this optimization. The results, detailed in Table 6, confirm that GRPO serves as a powerful fine-tuning method for enhancing temporal precision in video chaptering. From the experimental results, we draw three key conclusions.

First, GRPO directly and consistently improves metrics correlated with temporal segmentation accuracy. As hypothesized, by optimizing with a reward focused on temporal alignment, we observe a clear performance boost in F1 and tIoU scores across all configurations. For instance, on the VidAtlas-test set, the GRPO model with video input achieves a notable gain of +0.8 in F1 and +0.7 in tIoU over its SFT baseline. This empirically validates that GRPO effectively sharpens the model’s ability to predict precise chapter boundaries.

Second, we observe a significant degree of cross-modal transferability from the RL training. Notably, despite the GRPO training being conducted exclusively on the video modality, the temporal localization performance of the ASR and Video+ASR inputs also improves. The GRPO model with Video+ASR input, for example, achieves a +1.5 F1 and +1.1 tIoU gain on VidChapter7M-test. This suggests that the optimization is not merely learning a superficial visual-to-temporal mapping but is refining a more abstract, modality-agnostic representation of temporal structure within the language model’s parameters.

Finally, these enhancements in temporal precision are achieved without sacrificing semantic quality. Crucially, although our reward function is agnostic to content, semantic metric such as CIDEr remain highly comparable to the SFT baseline, and in some cases even improve (e.g., +1.1 CIDEr for video input on VidChapters7M-

Table 6 Effectiveness of Reinforcement Learning with GRPO. We compare the performance of our models before (SFT) and after applying reinforcement learning (+RL) with GRPO. The evaluation is conducted on two benchmarks across different input modalities (A: ASR, V: Video). The results show that GRPO consistently improves temporal segmentation metrics (F1, tIoU) while maintaining or slightly improving semantic quality metrics (S: SODA, C: CIDEr). Bold numbers indicate the best performance between the base model and GRPO-enhanced model for each metric.

Method	Stage	Modality		VidChapters7M-test					VidAtlas-test				
		A	V	F1	tIoU	S	C	G	F1	tIoU	S	C	G
Base-asr GRPO-asr	sft	✓	X	54.5	76.7	26.3	144.0	28.9	58.8	79.7	24.8	111.3	28.0
	+rl	✓	X	54.8(+0.3↑)	77.2(+0.5↑)	25.3(-1.0↓)	143.7(-0.3↓)	28.8 (-0.1↓)	59.6(+0.8↑)	80.2(+0.5↑)	24.7(-0.1↓)	109.9(-1.4↓)	28.0(↑↓)
Base-vid GRPO-vid	sft	X	✓	50.2	74.3	22.9	138.3	25.4	57.6	78.9	21.6	98.1	25.0
	+rl	X	✓	50.6(+0.4↑)	74.8(+0.5↑)	22.9(↑↓)	139.4(+1.1↑)	25.4(↑↓)	58.4(+0.8↑)	79.6(+0.7↑)	21.9(+0.3↑)	98.2(+0.1↑)	25.0(↑↓)
Base-vidasr GRPO-vidasr	sft	✓	✓	59.3	79.6	30.6	186.6	34.3	66.2	84.0	30.2	141.5	34.1
	+rl	✓	✓	60.8(+1.5↑)	80.7(+1.1↑)	31.0(+0.4↑)	190.7(+4.1↑)	34.6(+0.3↑)	66.8(+0.6↑)	84.3(+0.3↑)	30.4(+0.2↑)	141.7(+0.2↑)	34.4(+0.3↑)

test.). Composite metrics like SODA and GRACE, which balance segmentation and description, also maintain their performance or exhibit slight gains. This indicates that the KL-regularized optimization successfully avoids policy degradation, suggesting a positive effect where more accurate segmentation enables the model to generate more focused and relevant content. In summary, GRPO acts as a critical fine-tuning step, effectively sharpening the model’s temporal acuity while preserving its descriptive capabilities.

5.5 Qualitative Visualization

To provide a more intuitive understanding of our model’s capabilities beyond quantitative metrics, we present qualitative examples on both English and Chinese videos. These visualizations showcase ARC-Chapter’s ability to generate accurate, coherent, and hierarchically structured outputs in multiple formats and languages.

Fig. 7 illustrates the model’s performance on a challenging English video discussing US debt and the role of stablecoins. The topic is dense with financial terminology and complex arguments. Our model successfully navigates this complexity across all output formats. The *Short Title* accurately segments the video into logical thematic units, such as "Intro", "Stablecoin Regulation". The *Video Description with Timestamp* summarizes the video content for each chapter. More impressively, the *Structural Chapters* demonstrates the model’s advanced capability for hierarchical chaptering. The generated title, abstract, and introduction for each chapter are distinct yet complementary, providing a rich, layered understanding of the content that mirrors human-authored summaries.

To showcase the multilingual performance of our model, Fig. 8 presents the results for a Chinese video on a similar topic. The model exhibits a comparable level of understanding and generation quality in Chinese. The generated *Short Titles* are precise. The detailed *Description* and *Structural Chapters* are fluent and contextually appropriate. This strong cross-lingual performance underscores the model’s ability to generalize the learned chaptering and summarization skills, rather than merely memorizing patterns in a single language.

Together, these qualitative examples confirm that ARC-Chapter is not only a powerful chaptering tool but also a versatile video understanding model capable of producing rich, structured, and multilingual summaries that are both accurate and useful for end-users.

6 Conclusion

In this report, we introduced ARC-Chapter, a scalable and robust framework for structuring long-form videos into semantically coherent chapters and hierarchical summaries. ARC-Chapter leverages a large-scale dataset of millions of long video chapters and employs a semi-automatic annotation pipeline. These innovations advance the state of the art in video chaptering and summary generation. We also proposed the GRACE metric, which addresses the limitations of existing evaluation methods by providing a granularity-robust assessment of chapter boundaries. Experimental results show that ARC-Chapter achieves superior performance across multiple benchmarks, video durations, and languages. These findings demonstrate the framework’s effectiveness and generalizability. ARC-Chapter has strong potential to facilitate efficient content navigation, retrieval, and understanding as long-form video content continues to grow rapidly.



Video Description with Timestamp

(00:00:00 - 00:00:48) The video begins by introducing the pressing issue of the US national debt, highlighting the potential consequences of a default and the Trump administration's interest in utilizing stablecoins as a potential solution. This sets the stage for a deeper dive into the mechanics of the debt ceiling and the role stablecoins might play. (00:00:48 - 00:05:35) Next, the video explains the debt ceiling and its implications, including the possibility of market stimulation due to reduced government debt issuance. The discussion covers the Federal Reserve's quantitative tightening (QT) program and its impact on commercial banks' cash reserves, ultimately leading to the exploration of how stablecoins could contribute to refilling the Treasury General Account (TGA). (00:05:35 - 00:09:23) Moving on, the video defines stablecoins and their connection to US dollar-backed assets. It examines the role of stablecoins like USDT and USDC in various market segments, such as crypto trading and DeFi, and discusses the potential for significant growth driven by increased demand from retail investors and PayPal's expanding use of stablecoins. (00:09:23 - 00:14:24) The video then delves into the potential of stablecoins to address the TGA funding gap, drawing parallels to the 2023 refill process. It analyzes the potential growth of USDT, USDC, and PayPal's PYUSD, considering the influence of retail investors, DeFi adoption, and PayPal's merchant network. The analysis suggests that stablecoins could play a crucial role in meeting the \$800 billion funding requirement. (00:14:24 - 00:17:08) Shifting focus to regulatory developments, the video discusses the two proposed stablecoin bills in Congress: the Senate's Genius Act and the House's Stable Act. It compares the two bills, highlighting their differing approaches to regulation, particularly concerning decentralized stablecoins, and emphasizes the importance of swift passage of these regulations. (00:17:08 - 00:20:57) Finally, the video explores the potential benefits of stablecoins for various cryptocurrencies, focusing on the blockchains where major stablecoins are active. It examines the growth of USDT, USDC, and PYUSD on platforms like Ethereum, Solana, and others, and discusses the potential impact of XRP's EVM-compatible sidechain on the XRP ecosystem and the broader crypto market.

Structural Chapter

Video Chapters

► Intro [00:00:00]

Title: US Debt Bubble & Stablecoins: An Overview

Intro: The video begins by discussing the US national debt exceeding \$36 trillion and the potential for a bubble burst. It introduces the idea that stablecoins could be a solution, noting that stablecoin issuers have already purchased over \$60 billion in US debt. The presenter, Nick, sets the stage for explaining the debt problem, the role of stablecoins, and which cryptocurrencies might benefit.

► US Debt Problem [00:00:48]

Title: The US Debt Ceiling Crisis and Market Impact

Intro: This section details the US debt ceiling situation, explaining that the US government hit its debt ceiling in January and cannot issue more debt. The presenter explains that this can paradoxically stimulate the market by reducing debt issuance, leading to more money flowing into other assets. The chapter highlights the need for the US government to refill the Treasury General Account (TGA) with \$800 billion in bonds, which presents a challenge for the market.

► Stablecoins & US Debt [00:05:35]

Title: Stablecoins as Potential Solutions for US Debt

Intro: This segment focuses on how stablecoins could be a solution to the US debt problem. It explains that stablecoins are crypto tokens pegged to fiat currencies, often the US dollar, and are backed by US bonds. The presenter notes that stablecoin issuers have purchased over \$60 billion in US bonds. The chapter highlights the potential of stablecoins to help refill the TGA, with the presenter emphasizing the importance of stablecoins in the market.

► How Much Stablecoin Growth Needed [00:09:22]

Title: Estimating Stablecoin Growth to Refill the TGA

Intro: This section examines how much stablecoin growth is needed to refill the TGA. The presenter uses the 2023 refill as a reference point, noting that the RRP facility was used to refill the TGA. The presenter then discusses the impact of Quantitative Tightening (QT) on the TGA refill. The presenter estimates that stablecoins could be the primary buyers of the \$800 billion needed to refill the TGA.

► Stablecoin Regulations [00:14:24]

Title: Stablecoin Regulation and Congressional Action

Intro: This segment covers the stablecoin bills being considered in Congress. The presenter mentions the Genius Act in the Senate and the Stable Act in the House. The presenter notes that the Genius Act is more favorable for the TGA refill because it has fewer restrictions. The presenter also discusses the potential for these bills to be incompatible and the possibility of a pro-crypto supermajority in Congress. The presenter notes that stablecoin regulations are expected to be passed by August.

► Which Cryptos Will Benefit [00:17:08]

Title: Cryptocurrencies Set to Benefit from Stablecoin Growth

Intro: This section identifies which cryptocurrencies will benefit from the stablecoin frenzy. The presenter suggests looking at the blockchains where the biggest stablecoins are active and growing. The presenter notes that USDT is growing the fastest on Tron, APTOS, and Ethereum's layer twos. The presenter also notes that USDC is growing the most on Hyperliquid, Solana, SUI, APTOS, and Sonic. The presenter also mentions that PayPal's PYUSD is expanding to Solana and Ethereum, and Ripple's XRP is launching an EVM-compatible sidechain.

Video Summary

This video analyzes the potential of stablecoins to address the US national debt bubble. It begins by outlining the US debt ceiling crisis and its market implications, then explores how stablecoins, particularly those backed by US bonds, could be a solution. The video estimates the necessary stablecoin growth to refill the Treasury General Account (TGA), discusses the regulatory landscape in Congress, and identifies cryptocurrencies poised to benefit from the anticipated stablecoin growth, focusing on the blockchains where major stablecoins are active. The video concludes by highlighting the potential for a significant injection of crypto-native liquidity and the impact on the broader crypto market, particularly altcoins on layer-one blockchains.

Figure 7 Qualitative results on an English video about finance and cryptocurrency.



Figure 8 Qualitative results on a Chinese video discussing stablecoins.

References

- [1] Abdelrahman Abouelenin, Atabak Ashfaq, Adam Atkinson, Hany Awadalla, Nguyen Bach, Jianmin Bao, Alon Benhaim, Martin Cai, Vishrav Chaudhary, Congcong Chen, et al. Phi-4-mini technical report: Compact yet powerful multimodal language models via mixture-of-loras. [arXiv preprint arXiv:2503.01743](#), 2025.
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. [arXiv preprint arXiv:2303.08774](#), 2023.
- [3] Anthropic. The claudie 3 model family: Opus, sonnet, haiku. 2024.
- [4] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-VL: A frontier large vision-language model with versatile abilities. [arXiv preprint arXiv:2308.12966](#), 2023.
- [5] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. [arXiv preprint arXiv:2502.13923](#), 2025.
- [6] Richard Bellman and Robert Kalaba. On adaptive control processes. *IRE Transactions on Automatic Control*, 4(2):1–9, 2003.
- [7] Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. [arXiv preprint arXiv:2507.06261](#), 2025.
- [8] Guodong Ding, Fadime Sener, and Angela Yao. Temporal action segmentation: An analysis of modern techniques. *TPAMI*, 46(2):1011–1030, 2023.
- [9] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. [arXiv preprint arXiv:2407.21783](#), 2024.
- [10] Soichiro Fujita, Tsutomu Hirao, Hidetaka Kamigaito, Manabu Okumura, and Masaaki Nagata. Soda: Story oriented dense video captioning evaluation framework. In *ECCV*, pages 517–531, 2020.
- [11] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. Tall: Temporal activity localization via language query. In *ICCV*, 2017.
- [12] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. [arXiv preprint arXiv:2501.12948](#), 2025.
- [13] Dong Guo, Faming Wu, Feida Zhu, Fuxing Leng, Guang Shi, Haobin Chen, Haoqi Fan, Jian Wang, Jianyu Jiang, Jiawei Wang, et al. Seed1. 5-vl technical report. [arXiv preprint arXiv:2505.07062](#), 2025.
- [14] Yongxin Guo, Jingyu Liu, Mingda Li, Qingbin Liu, Xi Chen, and Xiaoying Tang. Trace: Temporal grounding video llm via causal event modeling. [arXiv preprint arXiv:2410.05643](#), 2024.
- [15] Yongxin Guo, Jingyu Liu, Mingda Li, Dingxin Cheng, Xiaoying Tang, Dianbo Sui, Qingbin Liu, Xi Chen, and Kevin Zhao. Vtg-llm: Integrating timestamp knowledge into video llms for enhanced video temporal grounding. In *AAAI*, pages 3302–3310, 2025.
- [16] Fabian Caba Heilbron, Juan Carlos Niebles, and Bernard Ghanem. Fast temporal activity proposals for efficient detection of human actions in untrimmed videos. In *CVPR*, pages 1914–1923, 2016.
- [17] Bin Huang, Xin Wang, Hong Chen, Zihan Song, and Wenwu Zhu. Vtimellm: Empower llm to grasp video moments. In *CVPR*, pages 14271–14280, 2024.
- [18] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. [arXiv preprint arXiv:2410.21276](#), 2024.
- [19] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *ICCV*, pages 706–715, 2017.
- [20] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *ICCV*, pages 706–715, 2017.

- [21] Hilde Kuehne, Ali Arslan, and Thomas Serre. The language of actions: Recovering the syntax and semantics of goal-directed human activities. In CVPR, 2014.
- [22] Colin Lea, Michael D Flynn, Rene Vidal, Austin Reiter, and Gregory D Hager. Temporal convolutional networks for action segmentation and detection. In CVPR, pages 156–165, 2017.
- [23] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. arXiv preprint arXiv:2305.03726, 2023.
- [24] Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language models. In ECCV, pages 323–340, 2024.
- [25] Ye Liu, Kevin Qinghong Lin, Chang Wen Chen, and Mike Zheng Shou. Videomind: A chain-of-lora agent for long video reasoning. arXiv preprint arXiv:2503.13444, 2025.
- [26] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, et al. DeepSeek-VL: Towards real-world vision-language understanding. arXiv preprint arXiv:2403.05525, 2024.
- [27] Zijia Lu and Ehsan Elhamifar. Fact: Frame-action cross-attention temporal modeling for efficient action segmentation. In CVPR, pages 18175–18185, 2024.
- [28] Long Qian, Juncheng Li, Yu Wu, Yaobo Ye, Hao Fei, Tat-Seng Chua, Yueting Zhuang, and Siliang Tang. Momentor: Advancing video large language model with fine-grained temporal reasoning. In ICML, 2024.
- [29] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In ICML, pages 28492–28518, 2023.
- [30] Shuhuai Ren, Linli Yao, Shicheng Li, Xu Sun, and Lu Hou. Timechat: A time-sensitive multimodal large language model for long video understanding. In CVPR, pages 14313–14323, 2024.
- [31] Hiroaki Sakoe and Seibi Chiba. Dynamic programming algorithm optimization for spoken word recognition. IEEE Transactions on Acoustics, Speech, and Signal Processing, 26(1):43–49, 2003.
- [32] Yuhan Shen and Ehsan Elhamifar. Progress-aware online action segmentation for egocentric procedural task videos. In CVPR, pages 18186–18197, 2024.
- [33] Fangxun Shu, Lei Zhang, Hao Jiang, and Cihang Xie. Audio-visual llm for video understanding. arXiv preprint arXiv:2312.06720, 2023.
- [34] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805, 2023.
- [35] Lucas Ventura, Antoine Yang, Cordelia Schmid, and Güл Varol. Chapter-llama: Efficient chaptering in hour-long videos with llms. In CVPR, pages 18947–18958, 2025.
- [36] Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. Git: A generative image-to-text transformer for vision and language. arXiv preprint arXiv:2205.14100, 2022.
- [37] Peiyu Wang, Yichen Wei, Yi Peng, Xiaokun Wang, Weijie Qiu, Wei Shen, Tianyidan Xie, Jiangbo Pei, Jianhao Zhang, Yunzhuo Hao, et al. Skywork r1v2: Multimodal hybrid reinforcement learning for reasoning. arXiv preprint arXiv:2504.16656, 2025.
- [38] Teng Wang, Ruimao Zhang, Zhichao Lu, Feng Zheng, Ran Cheng, and Ping Luo. End-to-end dense video captioning with parallel decoding. In ICCV, pages 6847–6857, 2021.
- [39] Zhenzhi Wang, Ziteng Gao, Limin Wang, Zhifeng Li, and Gangshan Wu. Boundary-aware cascade networks for temporal action segmentation. In ECCV, pages 34–51, 2020.
- [40] Hao Wu, Huabin Liu, Yu Qiao, and Xiao Sun. Dibs: Enhancing dense video captioning with unlabeled videos via pseudo boundary enrichment and online refinement. In CVPR, pages 18699–18708, 2024.
- [41] Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. NExT-GPT: Any-to-any multimodal llm. arXiv preprint arXiv:2309.05519, 2023.
- [42] Cilin Yan, Haochen Wang, Shilin Yan, Xiaolong Jiang, Yao Hu, Guoliang Kang, Weidi Xie, and Efstratios Gavves. Visa: Reasoning video object segmentation via large language models. In ECCV, pages 98–115, 2024.

- [43] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. [arXiv preprint arXiv:2505.09388](#), 2025.
- [44] Antoine Yang, Arsha Nagrani, Ivan Laptev, Josef Sivic, and Cordelia Schmid. Vidchapters-7m: Video chapters at scale. [NeurIPS](#), 36:49428–49444, 2023.
- [45] Antoine Yang, Arsha Nagrani, Ivan Laptev, Josef Sivic, and Cordelia Schmid. Vidchapters-7m: Video chapters at scale. In [NeurIPS](#), 2023.
- [46] Antoine Yang, Arsha Nagrani, Paul Hongsuck Seo, Antoine Miech, Jordi Pont-Tuset, Ivan Laptev, Josef Sivic, and Cordelia Schmid. Vid2seq: Large-scale pretraining of a visual language model for dense video captioning. In [CVPR](#), 2023.
- [47] Zuhao Yang, Yingchen Yu, Yunqing Zhao, Shijian Lu, and Song Bai. Timeexpert: An expert-guided video llm for video temporal grounding. [arXiv preprint arXiv:2508.01699](#), 2025.
- [48] Xiangyu Zeng, Kunchang Li, Chenting Wang, Xinhao Li, Tianxiang Jiang, Ziang Yan, Songze Li, Yansong Shi, Zhengrong Yue, Yi Wang, et al. Timesuite: Improving mllms for long video understanding via grounded tuning. [arXiv preprint arXiv:2410.19702](#), 2024.
- [49] Boqiang Zhang, Kehan Li, Zesen Cheng, Zhiqiang Hu, Yuqian Yuan, Guanzheng Chen, Sicong Leng, Yuming Jiang, Hang Zhang, Xin Li, et al. Videollama 3: Frontier multimodal foundation models for image and video understanding. [arXiv preprint arXiv:2501.13106](#), 2025.
- [50] Haoji Zhang, Xin Gu, Jiawen Li, Chixiang Ma, Sule Bai, Chubin Zhang, Bowen Zhang, Zhichao Zhou, Dongliang He, and Yansong Tang. Thinking with videos: Multimodal tool-augmented reinforcement learning for long video reasoning. [arXiv preprint arXiv:2508.04416](#), 2025.
- [51] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. [arXiv preprint arXiv:1904.09675](#), 2019.
- [52] Xiaoying Zhang, Da Peng, Yipeng Zhang, Zonghao Guo, Chengyue Wu, Chi Chen, Wei Ke, Helen Meng, and Maosong Sun. Towards self-improving systematic cognition for next-generation foundation mllms. [arXiv preprint arXiv:2503.12303](#), 2025.
- [53] Yipeng Zhang, Yifan Liu, Zonghao Guo, Yidan Zhang, Xuesong Yang, Chi Chen, Jun Song, Bo Zheng, Yuan Yao, Zhiyuan Liu, et al. Llava-uhd v2: an mllm integrating high-resolution feature pyramid via hierarchical window transformer. [arXiv preprint arXiv:2412.13871](#), 2024.
- [54] Yue Zhao, Yuanjun Xiong, Limin Wang, Zhirong Wu, Xiaou Tang, and Dahua Lin. Temporal action detection with structured segment networks. In [ICCV](#), pages 2914–2923, 2017.
- [55] Luwei Zhou, Chenliang Xu, and Jason Corso. Towards automatic learning of procedures from web instructional videos. In [AAAI](#), 2018.
- [56] Xingyi Zhou, Anurag Arnab, Shyamal Buch, Shen Yan, Austin Myers, Xuehan Xiong, Arsha Nagrani, and Cordelia Schmid. Streaming dense video captioning. In [CVPR](#), pages 18243–18252, 2024.
- [57] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. [arXiv preprint arXiv:2504.10479](#), 2025.