# ChronusOmni: Improving Time Awareness of Omni Large Language Models

Yijing Chen[1†], Yihan Wu[1†], Kaisi Guan[1], Yuchen Ren[1], Yuyue Wang[1], Ruihua Song[1*], Liyun Ru[2*]

[1]Gaoling School of Artificial Intelligence, Renmin University of China  [2]Baichuan Inc.

## Abstract

*Time awareness is a fundamental ability of omni large language models, especially for understanding long videos and answering complex questions. Previous approaches mainly target vision-language scenarios and focus on the explicit temporal grounding questions, such as identifying when a visual event occurs or determining what event happens at a specific time. However, they often make insufficient use of the audio modality, and overlook implicit temporal grounding across modalities—for example, identifying what is visually present when a character speaks, or determining what is said when a visual event occurs—despite such cross-modal temporal relations being prevalent in real-world scenarios. In this paper, we propose **ChronusOmni**, an omni large language model designed to enhance temporal awareness for both explicit and implicit audiovisual temporal grounding. First, we interleave text-based timestamp tokens with visual and audio representations at each time unit, enabling unified temporal modeling across modalities. Second, to enforce correct temporal ordering and strengthen fine-grained temporal reasoning, we incorporate reinforcement learning with specially designed reward functions. Moreover, we construct **ChronusAV**, a temporally-accurate, modality-complete, and cross-modal-aligned dataset to support the training and evaluation on audiovisual temporal grounding task. Experimental results demonstrate that ChronusOmni achieves state-of-the-art performance on ChronusAV with more than 30% improvement and top results on most metrics upon other temporal grounding benchmarks. This highlights the strong temporal awareness of our model across modalities, while preserving general video and audio understanding capabilities. Code and dataset are available at* https://github.com/YJCX330/Chronus/.

## 1. Introduction

Video temporal grounding [18, 19, 55] aims to align video semantics with temporal segments, serving as a critical

---

†Equal contribution. Contact chenyijing@ruc.edu.cn
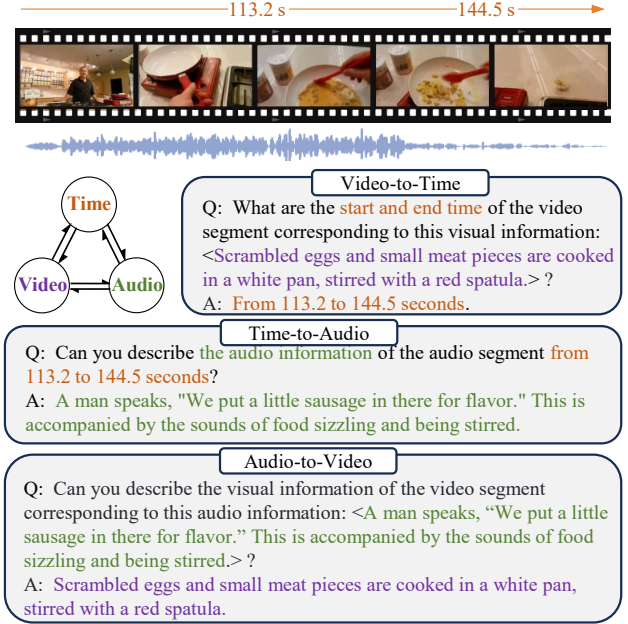*Corresponding authors.

Figure 1. Illustration of the audiovisual temporal grounding task. Three primary elements Video, Audio and Time are connected through six directional basic temporal grounding subtasks: Video-to-Time, Time-to-Video, Audio-to-Time, Time-to-Audio, Video-to-Audio, and Audio-to-Video.

capability for video understanding. Early studies in this field [7, 18, 19, 23, 24, 43, 53, 54] primarily rely solely on visual information to align video content with textual descriptions. However, as in-the-wild videos increasingly contain rich audio cues (dialogues, environmental sounds, and background music), visual information alone is insufficient for precise grounding. This motivates a shift toward audiovisual temporal grounding, which incorporates both visual and auditory modalities to capture complementary cues and improve time awareness in complex real-world scenarios.

Temporal grounding serves as a prerequisite for sophisticated multimodal reasoning. Consider a typical audiovisual scenario (e.g., in a movie dialogue scene, the camera first stays on Character A's expression; off-screen, Character B delivers the line "Let's go now." Then the shot cuts to B, and B's lip movements align with the dialogue.), where the semantic meaning of an event can only be accurately

understood when both visual and auditory cues are jointly analyzed and temporally aligned. In such cases, it requires model to understand not only the **localization information in video streams** and the **localization information in audio streams** (explicit temporal grounding), but also **the synchronization of audio-visual interactions** (implicit temporal grounding). However, most existing works focus solely on explicit temporal grounding [18, 23, 24, 49, 63], relying on either visual frames or audio features without capturing temporal dependencies cross audio and video. As a result, they struggle in complex real-world scenarios where auditory and visual cues complement each other. This limitation reveals a fundamental challenge for multimodal temporal grounding: models must be capable of jointly representing, aligning, and reasoning over temporal dynamics across both auditory and visual streams. To address the above challenge, in this work, we improve the audiovisual temporal grounding task in a systematic way. First, we formally define the audiovisual temporal grounding task and systematically analyze its core challenges. As illustrated in Figure 1, a model with audiovisual temporal grounding capability should support six directions of prediction across time, video, and audio: **video-to-time**, **time-to-video**, **audio-to-time**, **time-to-audio**, **audio-to-video**, and **video-to-audio**. Among them, the first four directions are explicit temporal grounding, while the last two directions constitute implicit temporal grounding. This requires three essential abilities: (1) video content localization, (2) audio content localization, and (3) synchronization modeling of audio-visual interactions. Second, we design audiovisual temporal-synchronized representation, improving multimodal joint temporal reasoning ability with temporal-aware training strategy. Furthermore, we construct a comprehensive audiovisual dateset to support the training and evaluation for audiovisual temporal understanding.

The main contributions can be summarized as follows.

- We propose **ChronusOmni**, an omni large language model designed for the audiovisual temporal grounding task. More specifically, we design an audiovisual temporal interleaved representation that aligns visual, audio, and temporal information effectively. We further design the temporal-aware supervised finetuning and reinforcement learning strategy upon the representation, showing remarkable performance on both audiovisual joint temporal understanding and unimodel understanding.
- We introduce **ChronusAV**, a temporally-accurate, modality-complete, and cross-modal-aligned dataset to support the training and evaluation on the audiovisual temporal grounding task, covering six cross-modal prediction directions among time, video, and audio.
- Experimental results show that ChronusOmni achieves the best performance in all metrics over ChronusAV (with more than 30% improvement) and most metrics over

other temporal grounding benchmarks, while preserving its general video and audio understanding capabilities. These results demonstrate the effectiveness of our proposed temporal modeling method and training strategy.

## 2. Related Work

**Temporal Grounding MLLMs.** Traditional video temporal grounding [55] aims to precisely align video semantics with corresponding time, including various tasks such as moment retrieval [4, 13, 22, 62], dense video caption [25, 60], and video highlight detection [57]. With the development of video large language models (LLMs), some recent works leverage the reasoning capabilities of LLMs to enhance temporal grounding. Some models directly infer frame indices relying on their ability to reason over frame order [23, 31], while others introduce absolute timestamp encoding to strengthen temporal localization [7, 18, 19, 41]. Beyond the video modality, several audiovisual LLMs [14, 17, 33, 48, 58, 59, 64] extend temporal perception capabilities to audiovisual scenarios. Many use learnable temporal positional embeddings [58, 59, 64] or specially designed time encoders [33], which require substantial amounts of temporal training data to develop fine-grained temporal sensitivity from scratch. ARC-Hunyuan-Video [14] instead embeds a watermark representing absolute time on each frame. It introduces additional complexity because the model must perform OCR-like operations to recover temporal information. To address these limitations, we propose an audiovisual temporal-synchronized representation. We encode time directly as a text modality and interleave audio, visual, and time tokens along the timeline, enabling explicit, fine-grained alignment across modalities. With our temporal-aware training strategy, the model acquires strong audiovisual temporal grounding ability using only a relatively small amount of training data.

**Temporal Grounding Datasets.** Most existing temporal grounding benchmarks [21, 26, 27, 46, 46, 65] rely on visual information solely, neglecting the audio modality entirely. AVEL [51] and UnAV-100 [16] incorporate both audio and video streams with coarse event-level labels without detailed captions, limiting the model's ability to understand fine-grained audiovisual temporal details. LongVALE [17] and TriSense-2M [33] offer more detailed temporal event descriptions that combine information from both modalities without separate audio and video captions. This entanglement makes it difficult to analyze or assess a model's temporal understanding capability within and across each modality independently. VUE-TR [49] provides separate temporal annotations for video and audio streams, but its audio annotations focus only on human speech, restricting its applicability to more diverse audio-visual scenarios. Furthermore, none of the existing temporal grounding benchmarks evaluate implicit temporal perception across audio
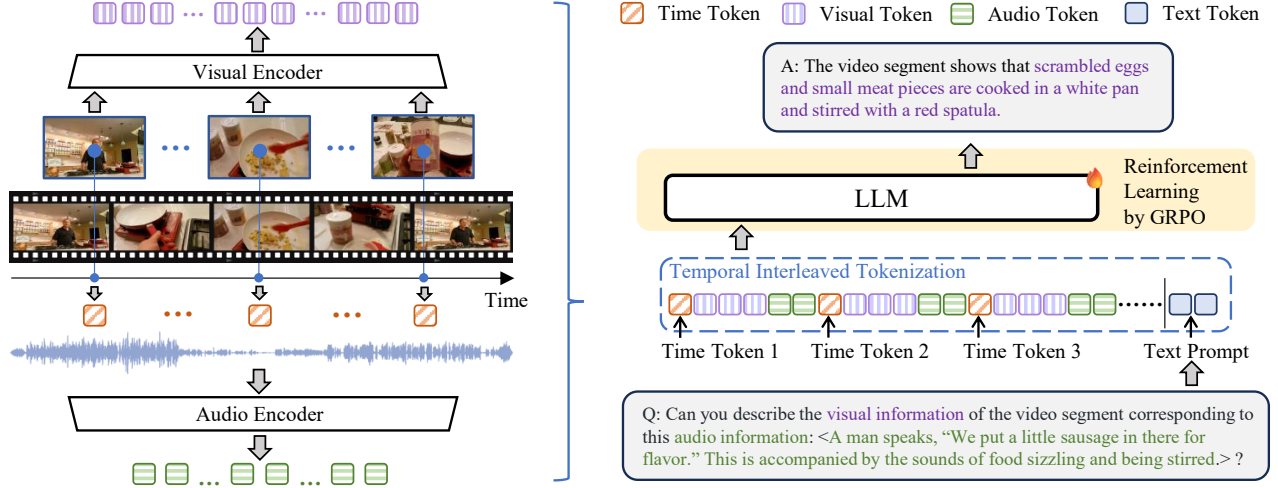
Figure 2. The architecture of ChronusOmni. Time, video, audio are tokenized and interleaved at each time step. The token sequence is along with text prompt is input into an LLM, which is supervised finetuned and further enhanced by reinforcement learning.

and video modalities, such as grounding what is heard when something is seen, or what is seen when something is heard. To address these limitations, we introduce ChronusAV, a temporally-accurate, modality-complete, and cross-modal-aligned dataset that provides time aligned triplets of (timestamp, video segment caption, audio segment caption) for each video. This design enables flexible construction of all six multimodal temporal grounding subtasks (as shown in Figure 1) and establishes a comprehensive foundation for advancing multimodal temporal grounding research.

## 3. Task Definition: Audiovisual Temporal Grounding

In this section, we formally define the audiovisual temporal grounding task and its six subtasks accordingly. Suppose we have a video $v$ and its corresponding audio $a$, both with a duration $t$. We denote a seqence of absolute timestamps as $t_0, t_1, t_2, ..., t_n$, which divide the video into $n$ temporal segments: $(t_0, t_1), (t_1, t_2), ..., (t_{n-1}, t_n)$. For each interval $(t_i, t_{i+1})$, we extract a video clip $v_i$ and an audio clip $a_i$. The clip $v_i$ (and similarly $a_i$) starts at $t_i$ and ends at $t_{i+1}$. Thus, we obtain $n$ aligned tuples:

$$D_i = \{v_i, a_i, (t_i, t_{i+1})\}, \tag{1}$$

where each tuple contains a video segment, an audio segment, and its corresponding timestamp.

Given a tuple $D_i$, the audiovisual temporal grounding task consists of querying one element to infer another. This give rise to six subtasks, as illustrated in Figure 1.

- **Video-to-Time** (V2T). Given the information of the video segment $v_i$, the goal is to ground its corresponding absolute time segment $(t_i, t_{i+1})$.
- **Time-to-Video** (T2V). Given the absolute time segment $(t_i, t_{i+1})$, the objective is to obtain the information of the associated video segment $v_i$.

- **Audio-to-Time** (A2T). Given the information of the audio segment $a_i$, the goal is to ground its corresponding absolute time segment $(t_i, t_{i+1})$.
- **Time-to-Audio** (T2A). Given the absolute time segment $(t_i, t_{i+1})$, the objective is to obtain the information of the associated audio segment $a_i$.
- **Video-to-Audio** (V2A). Given the visual content of $v_i$, V2A aims to identify the content of the time-synchronized audio segment $a_i$.
- **Audio-to-Video** (A2V). Given the audio information of $a_i$, the goal is to identify the time-synchronized video segment $v_i$.

## 4. Proposed Method: ChronusOmni

To equip the model with fine-grained audiovisual time awareness, we propose ChronusOmni, a unified framework that achieves precise alignment across video, audio, and time. Unlike prior MLLMs that process audio and video largely in isolation, ChronusOmni organize tokens in an interleaved manner to explicitly model cross-modal temporal dependencies. In addition, we design a two-stage, coarse-to-fine temporal optimization strategy to strengthen the model's temporal reasoning ability, resulting in robust performance on multimodal temporal grounding tasks.

### 4.1. Audiovisual Temporal Representation

To improve cross-modal synchronization and achieve fine-grained alignment of audiovisual information and absolute time, we propose a temporal interleaved tokenization method, explicitly aligning the three elements–video, audio, and time along the timeline, as depicted in Figure 2.

For visual input, we uniformly sample a fixed number of video frames. A visual encoder and audio encoders are used to transform the video frames and audio

into modality-specific tokens. Because the time intervals between frames are not fixed, explicit timestamps are essential for conveying absolute temporal information to the model. For each sampled frame, we extract its corresponding time point and convert it into text following the fixed format second{t} in [7]. For example, for a 126-second video sampled at 2-second intervals, we collect 64 timestamps: "second{0.0}", "second{2.0}", "second{4.0}", ..., "second{126.0}". These timestamps are encoded directly into text tokens, making them fully interpretable by the LLM. In contrast to methods that rely on learnable temporal embeddings, our explicit textual timestamps do not require aligning a learned embedding space with the language space, resulting in better adaptability, easy extensibility, and more stable temporal grounding performance.

Before feeding the tokens into the LLM, we interleave the token sequences from these three elements to achieve fine-grained temporal alignment. The interleaved tokens can be represented as:

$$I = [T_1, V_1, A_1, T_2, V_2, A_2, ..., T_i, V_i, A_i, ...], \quad (2)$$

where $T_i$ is the absolute timestamp tokens corresponding to the specific time point $t_i$, $V_i$ is the tokens encoded from the frame at $t_i$, and $A_i$ is the tokens encoded from the audio between $t_i$ and $t_{i+1}$. By using interleaved tokens, the model can establish a fine-grained connection among absolute time, visual information, and audio information.

### 4.2. Audiovisual Temporal Optimization

To enable the model to gradually acquire and robustly exploit temporal information, we further design a coarse-to-fine training strategy, including temporal-aware finetuning and temporal-aware reinforcement learning.

**Stage 1: Temporal-aware Supervised Finetuning.** In the first stage, we optimize the ChronusOmni via supervised fine-tuning (SFT). The model is trained on the video and audio dense captioning task, where for each video-audio pair input, the model learns to localize events and output timestamps, visual caption and audio caption for each event.

This training objective, together with our temporal-synchronized representation introduced in Section 4.1, guides the model to understand the alignment relationships among video, audio, and time.

However, SFT alone presents inherent limitations for temporal grounding [54]. First, although time is a continuous variable, SFT treats it as a categorical label, ignoring the metric structure—i.e., the magnitude and distance between time points. Second, as SFT exposes the model to gold prefixes throughout training, the model tends to rely on pattern memorization rather than learning to precisely localize boundaries. Consequently, SFT struggles to establish a robust mapping between audiovisual features and accurate temporal intervals.

**Stage 2: Reinforcement Learning with GRPO.** To overcome the limitations of SFT, we employ reinforcement learning in the second stage. Specifically, we adopt Group Relative Policy Optimization (GRPO) [10, 45] to train the model on the audiovisual temporal grounding task directly. This paradigm replaces the restrictive maximum likelihood objective with task-specific, outcome-driven reward functions, enabling the model to explore better strategies for aligning and integrating detailed audiovisual information with temporal boundaries.

For moment retrieval subtasks (V2T, A2T), the model outputs a temporal interval. We adopt the Intersection over Union (IoU) between the predicted interval and the ground-truth interval as the reward function:

$$R_{\text{IoU}} = \frac{|I_{\text{pred}} \cap I_{\text{gt}}|}{|I_{\text{pred}} \cup I_{\text{gt}}|}, \quad (3)$$

where $I_{\text{pred}}$ is the predicted time interval, and $I_{\text{gt}}$ is the ground-truth interval. To encourage consistent formatting that aligns with our timestamp design (Section 4.1), we also introduce a format reward enforcing outputs in the form "second{start}-second{end}":

$$R_{\text{format}} = \begin{cases} 0, & \text{if output format is wrong,} \\ 1, & \text{if output format is correct.} \end{cases} \quad (4)$$

For subtasks requiring the model to output a video or audio caption (T2A, V2A, T2V, A2V), we use Meteor [3] as the reward, which is a commonly used metric for evaluating video caption quality [1, 66]:

$$R_{\text{Meteor}} = \text{Meteor}(C_{\text{pred}}, C_{\text{gt}}), \quad (5)$$

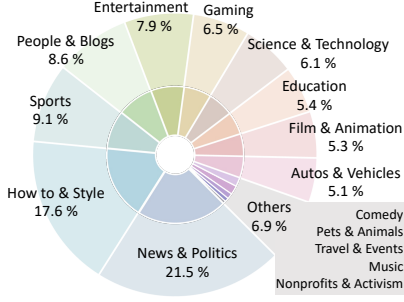where $C_{\text{pred}}$ and $C_{\text{gt}}$ is the predicted and ground-truth caption respectively.

This coarse-to-fine strategy effectively mitigates the discretization and exposure-bias issues of SFT by incorporating reward-based temporal supervision. As a result, the model achieves more accurate boundary localization and stronger audiovisual alignment, yielding substantially improved multimodal temporal grounding performance.

## 5. ChronusAV: Multimodal Dataset for Audiovisual Temporal Grounding

As discussed in Section 2, to bridge the gap of missing audiovisual temporal grounding datasets, we construct ChronusAV. ChronusAV is a large-scale dataset tailored for the audiovisual temporal grounding task, containing separate, detailed captions for audio (including sound, speech, and music) and visual information with precise temporal boundaries.

### 5.1. Construction of ChronusAV

To construct a temporally-accurate, modality-complete, and cross-modal-aligned dataset, we design a systematic dataset construction pipeline:

(a) Videos category distribution      (b) Videos duration distribution (seconds)

(c) The duration distribution of the video or audio segments corresponding to the QA in the test set (seconds)

Figure 3. Statistics of ChronusAV.

Table 1. Comparison of ChronusAV to existing related datasets. ChronusAV covers large-scale open-domain long videos and, among the listed benchmarks, is the only dataset that simultaneously provides precise timestamps, multimodal annotations across vision, sound/music, and speech, as well as separate captions for audio and video. A&V: audio and video.
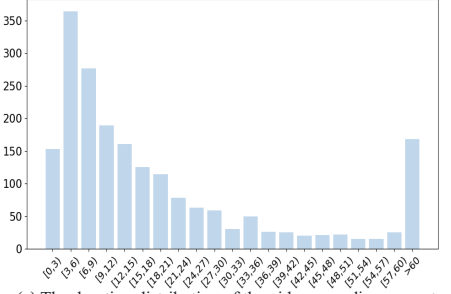
| Benchmark | Videos | Avg. video len | Annotations | Domain | Timestamps | Vision | Sound&Music | Speech | A&V seperate captions |
|---|---|---|---|---|---|---|---|---|---|
| Charades-STA[46] | 10K | 30s | 16K | daily activity | ✓ | ✓ | ✗ | ✗ | - |
| ActivityNet Caps[26] | 20K | 180s | 72K | daily activity | ✓ | ✓ | ✗ | ✗ | - |
| VALOR [35] | 1.18M | 10s | 1.18M | open | ✗ | ✓ | ✓ | ✗ | ✗ |
| VAST [6] | 27M | 30s | 27M | open | ✗ | ✓ | ✓ | ✓ | ✗ |
| AVEL [50] | 4K | 10s | 4K | open | ✓ | ✓ | ✓ | ✗ | ✗ |
| UnAV-100 [15] | 30K | 42.1s | 84K | open | ✓ | ✓ | ✓ | ✗ | ✗ |
| Shot2Story [20] | 43K | 17.1s | 181K | open | ✓ | ✓ | ✗ | ✓ | ✓ |
| VUE-TR [49] | 428 | 907s | 1598 | open | ✓ | ✓ | ✗ | ✓ | ✓ |
| LongVALE [17] | 8.4K | 235s | 108K | open | ✓ | ✓ | ✓ | ✓ | ✗ |
| ChronusAV | 47K | 226s | 677K | open | ✓ | ✓ | ✓ | ✓ | ✓ |

- **Video collection.** To construct a real-world, open-domain dataset, we select English videos from the large-scale, high-resolution Panda-70M dataset [8]. To cover videos of varying durations and reflect real-world temporal complexity, we select untrimmed long videos with audio tracks and restrict durations to 60–600 seconds, ensuring sufficient context for multimodal analysis.

- **Video segmentation.** To obtain precise event timestamps, we segment long, untrimmed videos into meaningful and coherent event units. Following Panda-70M, we first split basic visual scenes and then merge semantically similar ones, ensuring that each final segment represents a complete and coherent event. The corresponding audio tracks are segmented using the same timestamps to maintain audio–visual alignment. Finally, we retain only those audio–video pairs with 5 to 30 segments, ensuring that each sample contains rich temporal dynamics without being excessively fragmented.

- **Modality-specific annotation.** For each segment, we produce fine-grained, modality-aware captions by separately describing the visual and auditory streams. We use Gemini-2.5-Flash and Gemini-2.5-Pro [9] to annotate video and audio, respectively.

- **Human verification.** To ensure quality and modality separation of LLM-generated captions, we conduct a human study on 1,000 randomly sampled segments. Annotators rate accuracy and cross-modal leakage while viewing the corresponding video and audio. Results show that video captions are mostly accurate in 96.1% of cases; audio captions in 93.5%. And modality disentanglement is strong: 99.3% of video captions and 97.5% of audio captions show no or only minor cross-modal leakage. This high level of quality assurance validates the reliability of our automatically generated annotations for large-scale audiovisual temporal grounding training and evaluation.

- **Training set and benchmark construction.** Ultimately, we annotate 677K segments from 47K untrimmed videos. Each annotation is a (timestamp, video caption, audio caption) triplet, supporting all subtasks of audiovisual temporal grounding. We hold out 2,000 videos with corresponding audio to construct the test set and use the remaining 45K for training. For the test set, we randomly select one segment per video and generate six QA pairs covering the six subtasks, as shown in Figure 1, yielding 12K QA pairs. All test QAs are human-verified and refined to ensure accurate, unique answers.

As a result, we get the ChronusAV, a temporally-accurate, modality-complete, and cross-modal-aligned dataset for audiovisual temporal grounding.

## 5.2. Analysis of ChronusAV

As described in Section 5.1, we construct the ChronusAV to enable the model training and a comprehensive evaluation of audiovisual temporal grounding. In this section, we provide an overview of the dataset and highlight the key characteristics that distinguish ChronusAV from existing temporal

grounding datasets, which is also showed in Table 1.

- **Large-scale and long-duration videos.** ChronusAV contains 47K multi-shot videos totaling 2,922 hours, with an average duration of 226 seconds. As shown in Figure 3 (b), the videos range from 60 to 600 seconds, substantially longer than most prior benchmarks (e.g., ActivityNet Caps 180s, UnAV-100 42.1s). This long-form design provides rich temporal context and enables robust modeling of long-term event dependencies.
- **Diverse and open-domain coverage.** As shown in Figure 3 (a), ChronusAV spans 15 real-world domains, covering a wide spectrum of scenarios. This domain diversity ensures that ChronusAV aligns with the distributions of natural audiovisual content and supports strong generalization across contexts.
- **Fine-grained temporal annotations.** ChronusAV provides precise timestamps, which are essential for temporal grounding tasks. In contrast, large-scale datasets like VALOR and VAST do not provide temporal boundaries, limiting their suitability for fine-grained temporal alignment or localization. In the test set, QA-related segments peak at 3–6 seconds, as shown in Figure 3 (c), reflecting the benchmark's higher demands on a model's ability to capture subtle temporal cues across modalities.
- **Modality completeness.** ChronusAV is one of the few datasets that simultaneously provide visual content and comprehensive audio cues (speech, music, and sound). Many existing datasets are audio-incomplete—e.g., Charades-STA and ActivityNet Caps provide no audio annotation; VALOR and AVEL lack speech; Shot2Story and VUE-TR lack music and sound.
- **Separate audio and video captions.** ChronusAV includes modality-separated audio and video captions, enabling detailed and disentangled multimodal understanding. Most existing datasets (including LongVALE) do not offer this feature. Although Shot2Story and VUE-TR provide separate captions, their annotations lack sound-event descriptions, limiting their ability to capture real-world auditory complexity.

These features makes ChronusAV can serve as a solid foundation for training and assessing our model. More statistics and cases can be found in Appendix.

## 6. Experimental Settings

### 6.1. Implementation Details

We adopt multimodal LLM Ola [37] as our backbone, which is built upon Qwen-2.5-7B [61]. For visual processing, we use OryxViT [36] as the vision encoder and uniformly sample 64 frames from each input video. On the audio side, we employ Whisper-Large-V3 [42] as the speech encoder and BEATs [5] as the sound and music encoder. To achieve a comprehensive representation of the audio con-

tent, the embedding features from Whisper-Large-V3 and BEATs are concatenated across the channel dimension. The resulting audio features are then downsampled by a factor of 10 to reduce token length while preserving semantic information. Finally, both visual and audio features are projected through two-layer multi-layer perceptron adapters, which transform them into unified tokens for the LLM decoder.

### 6.2. Training Settings

We initialize ChronusOmni with the pretrained Ola model and only finetune the LLM parameters, keeping all encoders and adapters frozen throughout training.

- **SFT stage.** In this stage, we train on a mixture of How2 [44], AVSD [2], and ChronusAV datasets. We randomly sample 10K clips of How2-300h and train on the AVSR task. For AVSD, we split each dialogue into turn-level QA in every dialogue, randomly sample 30K QA pairs, and train on the AVQA task. For ChronusAV, we randomly sample 30K videos from the training set and train on the dense video and audio caption task. We train on all 70k samples for 1 epoch.
- **GRPO stage.** We construct the GRPO data using the ChronusAV training set. For each untrimmed video, we randomly select one segment and form a QA pair corresponding to one of the six subtasks using the segment's timestamp and captions. In this stage, we ultimately use 4000 pairs of QA data to train for 1000 steps.

### 6.3. Evaluation Settings

We evaluate ChronusOmni across a broad set of tasks and datasets to assess both its fine-grained temporal grounding ability and general audiovisual understanding.

- **ChronusAV.** We comprehensively assess ChronusOmni's temporal grounding performance using ChronusAV benchmark, and compare it with previous audiovisual LLMs. For V2T and A2T subtasks, we follow the evaluation metrics for the moment retrieval task used in previous works [18, 43], reporting Recall@1 at IoU thresholds of 0.5, 0.7. For another four subtasks, we use BLEU-4 [40], ROGUE-L [34], METEOR [3], and CIDEr [52] for standard caption quality evaluation.
- **LongVALE.** LongVALE [17] comprises three omni temporal tasks: Omni-TVG (predict the time segment corresponding to an omni event caption), Omni-DVC (predict timestamps and captions of all omni events), and Omni-SC (generate the omni caption for a specified time segment). We evaluate ChronusOmni on the LongVALE test set in a zero-shot setting.
- **Charades-STA and ActivityNet.** Charades-STA [46] and ActivityNet [26] are widely used video temporal grounding datasets. Because some of our training data (e.g., AVSD) contains Charades-STA videos, a fair zero-shot evaluation on Charades-STA is infeasible; thus, we

Table 2. Comparison with existing video-audio LLMs for multimodal temporal grounding task on ChronusAV benchmark. B: BLEU-4. R: ROUGE-L. M: METEOR. C: CIDEr. For Qwen3-Omni we use Qwen3-Omni-30B-A3B, while other models are all of size 7B.

| Model | V2T | | T2V | | | | A2T | | T2A | | | | V2A | | | | A2V | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | R@0.5 | R@0.7 | B | R | M | C | R@0.5 | R@0.7 | B | R | M | C | B | R | M | C | B | R | M | C |
| VideoLLaMA [64] | 2.00 | 0.80 | 0.10 | 1.50 | 0.93 | 0.03 | 1.70 | 0.80 | 0.05 | 1.13 | 0.63 | 0.06 | 0.08 | 1.11 | 0.81 | 0.06 | 0.07 | 1.22 | 0.93 | 0.00 |
| Ola [37] | 5.80 | 2.65 | 0.26 | 1.66 | 0.99 | 0.47 | 5.50 | 2.55 | 0.18 | 0.78 | 0.31 | 0.26 | 0.15 | 1.15 | 0.53 | 0.45 | 0.18 | 1.60 | 1.08 | 0.40 |
| AVicuna [48] | 10.75 | 4.70 | 0.04 | 1.03 | 0.32 | 0.12 | 8.20 | 3.45 | 0.01 | 0.51 | 0.13 | 0.02 | 0.02 | 0.54 | 0.23 | 0.10 | 0.11 | 1.45 | 0.65 | 0.38 |
| LongVALE-LLM [17] | 9.50 | 3.65 | 0.35 | 2.03 | 0.99 | 1.38 | 4.25 | 1.25 | 0.15 | 1.31 | 0.50 | 0.19 | 0.10 | 1.26 | 0.64 | 0.24 | 0.21 | 1.88 | 1.02 | 0.87 |
| Qwen2.5-Omni [58] | 7.05 | 3.00 | 0.35 | 1.89 | 1.02 | 0.79 | 10.10 | 3.70 | 0.67 | 1.01 | 0.64 | 0.26 | 0.54 | 0.91 | 0.61 | 0.30 | 0.19 | 1.61 | 0.97 | 0.58 |
| ARC-Hunyuan-Video [14] | 36.10 | 23.15 | 0.24 | 1.54 | 1.11 | 0.46 | 36.85 | 24.25 | 0.18 | 1.17 | 0.65 | 0.37 | 0.12 | 0.94 | 0.70 | 0.10 | 0.11 | 1.12 | 0.91 | 0.11 |
| Qwen3-Omni [59] | 37.85 | 21.80 | 0.37 | 2.13 | 1.62 | 0.92 | 46.70 | 33.10 | 0.94 | 2.26 | 1.20 | 2.18 | 0.35 | 1.58 | 1.25 | 0.74 | 0.22 | 1.72 | 1.49 | 0.39 |
| ChronusOmni | 63.15 | 45.95 | 1.16 | 3.37 | 2.12 | 5.07 | 90.50 | 79.85 | 6.78 | 6.86 | 4.50 | 34.30 | 3.61 | 4.90 | 3.27 | 13.60 | 1.01 | 3.17 | 2.12 | 3.03 |

Table 3. Comparison on LongVALE benchmark. "*" indicates that the model has been trained on the LongVALE training set.

| Model | Omni-TVG | | | | Omni-DVC | | | Omni-SC | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | R@0.3 | R@0.5 | R@0.7 | mIoU | SODA_c | CIDEr | METEOR | BLUE-4 | ROUGE-L | CIDEr | METEOR |
| VideoChat [30] | 2.2 | 0.9 | 0.4 | 3.0 | 0.7 | 0.2 | 0.9 | 0.5 | 9.6 | 0.0 | 8.2 |
| VideoChatGPT [38] | 4.9 | 2.0 | 0.9 | 5.0 | 0.7 | 0.1 | 0.9 | 0.4 | 14.0 | 0.9 | 5.9 |
| VideoLLaMA [64] | 2.5 | 1.1 | 0.3 | 1.9 | 0.6 | 0.6 | 0.9 | 0.9 | 11.5 | 0.1 | 8.9 |
| PandaGPT [47] | 2.5 | 1.0 | 0.3 | 2.2 | 0.5 | 0.0 | 0.6 | 0.6 | 14.9 | 0.3 | 8.9 |
| NExT-GPT [56] | 4.3 | 1.9 | 0.7 | 4.0 | 0.2 | 0.1 | 0.3 | 0.4 | 10.2 | 0.0 | 8.1 |
| TimeChat [43] | 5.8 | 2.6 | 1.1 | 5.2 | 1.6 | 0.1 | 1.4 | 1.2 | 16.1 | 1.6 | 10.0 |
| VTimeLLM [23] | 7.5 | 3.4 | 1.3 | 6.4 | 2.4 | 0.2 | 2.0 | 1.0 | 14.5 | 1.6 | 5.5 |
| TriSense [33] | 14.8 | 9.3 | 4.7 | 11.2 | - | - | - | 4.8 | 21.9 | 18.8 | 10.4 |
| LongVALE-LLM* [17] | 15.7 | 8.6 | 3.9 | 11.0 | 2.8 | 7.9 | 4.7 | 5.6 | 22.4 | 20.3 | 10.9 |
| ChronusOmni | 49.7 | 32.5 | 17.6 | 34.5 | 3.7 | 5.6 | 5.2 | 5.5 | 22.0 | 20.3 | 11.7 |

fine-tune for 1 epoch on Charades-STA using GRPO and report results under fine-tuning. For ActivityNet, we evaluate zero-shot on the test set. To enable fair comparison with prior VLMs, we use video-only inputs (no audio) for both datasets during training and testing.

- **General video and audio understanding benchmarks.** To evaluate whether our temporal modeling affects general audiovisual reasoning, we conduct zero-shot evaluations on four widely used benchmarks: Video-MME [11], Librispeech [39], VisSpeech [12], and MUSIC-AVQA [28]. Video-MME is vision-only, Librispeech is audio-only, while VisSpeech and MUSIC-AVQA utilize both visual and audio modalities.

## 7. Experimental Results

### 7.1. Results on Audiovisual Temporal Grounding

**Performance on ChronusAV.** We compare ChronusOmni with other omni LLMs on ChronusAV benchmark. As shown in Table 2, ChronusOmni consistently outperforms existing audiovisual LLMs across all subtasks by more than 30%. In V2T subtask, ChronusOmni reaches R@0.5/0.7 of 63.15/45.95, yielding 67%/98% gains over the next best model. In A2T subtask, the R@0.5/R@0.7 achieves 90.50/79.85, yielding 94%/142% improvements. For the remaining four subtasks, ChronusOmni delivers multiplicative gains across most captioning metrics, with particularly large magins on CIDEr, demonstrating stronger temporal audiovisual alignment. ChronusOmni also significantly outperforms recent omni-modal models such as ARC-Hunyuan-Video [14] and Qwen3-Omni [59], espe-

Table 4. Comparison with VLMs capable of temporal grounding on Charades-STA and ActivityNet. FT: fine-tuned setting; ZS: zero-shot setting. "*" indicates that this model is smaller in size than ChronusOmni, while the sizes of the other models are all 7B.

| Model | Charades-STA (FT) | | | ActivityNet (ZS) | | |
|---|---|---|---|---|---|---|
| | R@0.3 | R@0.5 | R@0.7 | R@0.3 | R@0.5 | R@0.7 |
| TimeChat [43] | - | 46.7 | 23.7 | 36.2 | 20.2 | 9.5 |
| Hawkeye [53] | 72.5 | 58.3 | 28.8 | 49.1 | 29.3 | 10.7 |
| VTimeLLM [23] | - | - | - | 44.0 | 27.8 | 14.3 |
| VTG-LLM [19] | - | 57.8 | 33.9 | - | - | - |
| TRACE [18] | - | 61.7 | 41.4 | - | - | - |
| TimeSuite [63] | 79.4 | 67.1 | 43.0 | - | - | - |
| iMOVE* [29] | 79.8 | 68.5 | 45.3 | 42.4 | 23.1 | 12.1 |
| VideoChat-R1 [32] | - | 71.7 | 50.2 | - | 33.4 | 17.7 |
| Time-R1[54] | 82.9 | 72.2 | 50.1 | 58.6 | 39.0 | 21.4 |
| ChronusOmni | 85.1 | 75.0 | 54.2 | 56.9 | 38.2 | 22.1 |

cially on audio-related subtasks, highlighting the limitations of prior models in handling complex auditory temporal cues. Notably, many baselines use more densely sampled video frames (e.g., 100 in LongVALE-LLM, 150 in ARC-Hunyuan-Video, 2 fps in Qwen3-Omni), yet ChronusOmni—using 64 frames—achieves superior temporal understanding, indicating both higher accuracy and greater efficiency.

**Performance on LongVALE.** Table 3 shows ChronusOmni's performance on LongVALE benchmark. As shown in the table, ChronusOmni delivers state-of-the-art performance on most evaluation metrics, while achieving the second-best results on the rest. In the Omni-TVG task, it achieves an mIoU of 34.5, nearly three times higher than the next best models—TriSense (11.2) and LongVALE-LLM (11.0)—demonstrating markedly stronger temporal local-

Table 5. Ablation study on our temporal interleaved tokenization method and training strategy. TIT: temporal interleaved tokenization.

| Model | V2T | | T2V | | | | A2T | | T2A | | | | V2A | | | | A2V | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | R@0.5 | R@0.7 | B | R | M | C | R@0.5 | R@0.7 | B | R | M | C | B | R | M | C | B | R | M | C |
| w/o TIT | 16.05 | 6.80 | 0.64 | 2.51 | 1.67 | 1.76 | 26.05 | 13.00 | 2.43 | 3.29 | 2.99 | 2.94 | 2.56 | 3.74 | 2.75 | 4.17 | 0.67 | 2.58 | 1.79 | 1.51 |
| w/o SFT | 59.40 | 37.20 | 0.53 | 2.39 | 1.81 | 0.69 | 84.45 | 69.25 | 2.45 | 3.29 | 1.89 | 6.91 | 2.12 | 3.15 | 1.96 | 5.32 | 0.49 | 2.34 | 1.84 | 0.43 |
| w/o GRPO | 30.00 | 15.45 | 0.46 | 2.53 | 1.07 | 2.37 | 34.55 | 19.00 | 0.04 | 1.48 | 0.35 | 1.85 | 0.54 | 1.40 | 0.73 | 1.24 | 0.12 | 1.81 | 0.66 | 0.76 |
| ChronusOmni | **63.15** | **45.95** | **1.16** | **3.37** | **2.12** | **5.07** | **90.50** | **79.85** | **6.78** | **6.86** | **4.50** | **34.30** | **3.61** | **4.90** | **3.27** | **13.60** | **1.01** | **3.17** | **2.12** | **3.03** |

ization ability. For Omni-DVC, ChronusOmni obtains the highest SODA_c (3.7) and METEOR (5.2), and ranks second in CIDEr (5.6), behind only LongVALE-LLM, which is specifically trained on the LongVALE training set. In Omni-SC, ChronusOmni achieves the best METEOR (11.7) and ties for the best CIDEr (20.3), while its BLEU-4 (5.5) and ROUGE-L (22.0) remain within negligible margins (0.1 and 0.4) of the leading scores.

## 7.2. Results on Visual-only Temporal Grounding

Besides audiovisual temporal grounding, we also evaluate ChronusOmni on visual-only temporal grounding. As shown in Table 4, ChronusOmni achieves state-of-the-art performance on Charades-STA (fine-tuned), obtaining 85.1/75.0/54.2 R@0.3/0.5/0.7. These results outperform the next best models by +2.2 (vs. Time-R1 at R@0.3), +2.8 (vs. Time-R1 at R@0.5), and +4.0 (vs. VideoChat-R1 at R@0.7). On ActivityNet (zero-shot), ChronusOmni ranks second at R@0.3 (56.9) and R@0.5 (38.2), only slightly behind Time-R1 by 1.7 and 0.8, respectively. It achieves the best R@0.7 (22.1), surpassing Time-R1 by +0.7, indicating stronger high-precision temporal localization under zero-shot transfer. Overall, ChronusOmni consistently outperforms or matches strong baselines on visual-only temporal grounding benchmarks, demonstrating robust and generalizable temporal understanding capabilities.

## 7.3. Results on General Video and Audio Understanding

We further evaluate whether audiovisual temporal-aware design affects the model's general multimodal reasoning ability. Figure 4 compares ChronusOmni with the base model across four standard video and audio understanding benchmarks. ChronusOmni matches the base model on general video QA (Video-MME) and shows only a minor degradation on pure speech recognition (LibriSpeech). In contrast, it yields substantial improvements on visual speech recognition (VisSpeech: 9.1 vs. 12.3 WER, a 26% reduction) and audio-visual reasoning (MUSIC-AVQA: +3.8 over the base model). It indicates that our temporal interleaved tokenization and temporal optimization do not impair general modality-specific abilities. Instead, they enhance the model's ability to perform joint audio–visual reasoning, demonstrating stronger multimodal integration without sacrificing overall understanding quality.
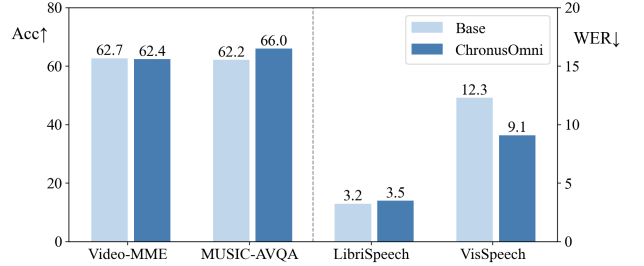


Figure 4. Evaluation on general video and audio understanding benchmarks. The evaluation metric for Video-MME and MUSIC-AVQA is Accuracy, with higher values being better. The evaluation metric for Librispeech and Visspeech is Word Error Rate (WER), with lower values being better. The "Base" is Ola.

## 7.4. Ablation Study

Table 5 shows the ablation results of our temporal interleaved tokenization method and training strategy. Removing temporal interleaved tokenization leads to severe degradation in moment retrieval: V2T R@0.7 drops from 45.95 to 6.80 and A2T R@0.7 from 79.85 to 13.00, highlighting its essential role in absolute time perception. SFT is crucial for caption quality: T2V CIDEr rises from 0.69 (w/o SFT) to 5.07, and T2A from 6.91 to 34.30. GRPO provides the largest single gains in video-audio alignment: removing it reduces all V2A and A2V metrics by over 3 times. With all components enabled, ChronusOmni achieves the best results across all six directions and all metrics, confirming that all three components are complementary and jointly necessary for strong audiovisual multimodal temporal grounding.

## 8. Conclusion

In this paper, we formally define audiovisual temporal grounding task and present ChronusOmni, a multimodal LLM tailored for this task. Our approach builds a temporally synchronized audio–visual representation and trains the model with temporal-aware supervised training and reinforcement learning to strengthen fine-grained temporal understanding. We also release ChronusAV, a comprehensive and standardized dataset for training and evaluation of audiovisual temporal grounding. Across ChronusAV and widely used public datasets, ChronusOmni sets a new state of the art, while maintaining robust video and audio understanding capabilities. These results demonstrate the effectiveness of our explicit cross-modal temporal alignment method and training strategy. Future work will focus on deploying our model in real-world interactive scenarios and scaling it to hour-long videos.

# References

[1] Moloud Abdar, Meenakshi Kollati, Swaraja Kuraparthi, et al. A review of deep learning for video captioning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 4

[2] Huda AlAmri, Vincent Cartillier, Abhishek Das, et al. Audio visual scene-aware dialog. In *CVPR*, pages 7558–7567, 2019. 6

[3] Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72, 2005. 4, 6

[4] Meinardus Boris, Batra Anil, Rohrbach Anna, and Marcus Rohrbach. The surprising effectiveness of multimodal large language models for video moment retrieval. *CoRR*, abs/2406.18113, 2024. 2

[5] Sanyuan Chen, Yu Wu, Chengyi Wang, et al. Beats: Audio pre-training with acoustic tokenizers. *arXiv preprint arXiv:2212.09058*, 2022. 6

[6] Sihan Chen, Handong Li, Qunbo Wang, Zijia Zhao, Mingzhen Sun, Xinxin Zhu, and Jing Liu. VAST: A vision-audio-subtitle-text omni-modality foundation model and dataset. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. 5

[7] Shimin Chen, Xiaohan Lan, Yitian Yuan, Zequn Jie, and Lin Ma. Timemarker: A versatile video-llm for long and short video understanding with superior temporal localization ability. *arXiv preprint arXiv:2411.18211*, 2024. 1, 2, 4

[8] Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, et al. Panda-70m: Captioning 70m videos with multiple cross-modality teachers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13320–13331, 2024. 5

[9] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, 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. 5

[10] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *CoRR*, abs/2501.12948, 2025. 4

[11] Chaoyou Fu, Yuhan Dai, Yondong Luo, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multimodal llms in video analysis. *CoRR*, abs/2405.21075, 2024. 7

[12] Valentin Gabeur, Paul Hongsuck Seo, Arsha Nagrani, Chen Sun, Karteek Alahari, and Cordelia Schmid. AVATAR: unconstrained audiovisual speech recognition. In *23rd Annual Conference of the International Speech Communication Association, Interspeech 2022, Incheon, Korea, September 18-22, 2022*, pages 2818–2822, 2022. 7

[13] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. TALL: temporal activity localization via language query. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 5277–5285, 2017. 2

[14] Yuying Ge, Yixiao Ge, Chen Li, et al. Arc-hunyuan-video-7b: Structured video comprehension of real-world shorts. *CoRR*, abs/2507.20939, 2025. 2, 7

[15] Tiantian Geng, Teng Wang, Jinming Duan, Runmin Cong, and Feng Zheng. Dense-localizing audio-visual events in untrimmed videos: A large-scale benchmark and baseline. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 22942–22951, 2023. 5

[16] Tiantian Geng, Teng Wang, Jinming Duan, Runmin Cong, and Feng Zheng. Dense-localizing audio-visual events in untrimmed videos: A large-scale benchmark and baseline. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 22942–22951. IEEE, 2023. 2

[17] Tiantian Geng, Jinrui Zhang, Qingni Wang, Teng Wang, Jinming Duan, and Feng Zheng. Longvale: Vision-audio-language-event benchmark towards time-aware omni-modal perception of long videos. In *CVPR*, pages 18959–18969, 2025. 2, 5, 6, 7

[18] Yongxin Guo, Jingyu Liu, Mingda Li, Qingbin Liu, Xi Chen, and Xiaoying Tang. TRACE: temporal grounding video LLM via causal event modeling. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025. 1, 2, 6, 7

[19] Yongxin Guo, Jingyu Liu, Mingda Li, et al. VTG-LLM: integrating timestamp knowledge into video llms for enhanced video temporal grounding. In *AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025, Philadelphia, PA, USA*, pages 3302–3310, 2025. 1, 2, 7

[20] Mingfei Han, Linjie Yang, Xiaojun Chang, Lina Yao, and Heng Wang. Shot2story: A new benchmark for comprehensive understanding of multi-shot videos. In *ICLR*, 2025. 5

[21] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, et al. Localizing moments in video with natural language. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 5804–5813. IEEE Computer Society, 2017. 2

[22] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, et al. Localizing moments in video with temporal language. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 1380–1390, 2018. 2

[23] Bin Huang, Xin Wang, Hong Chen, Zihan Song, and Wenwu Zhu. Vtimellm: Empower llm to grasp video moments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14271–14280, 2024. 1, 2, 7

[24] De-An Huang, Shijia Liao, Subhashree Radhakrishnan, et al. LITA: language instructed temporal-localization assistant. In

*Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LXIV*, pages 202–218, 2024. 1, 2

[25] Minkuk Kim, Hyeon Bae Kim, Jinyoung Moon, Jinwoo Choi, and Seong Tae Kim. Do you remember? dense video captioning with cross-modal memory retrieval. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 13894–13904, 2024. 2

[26] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *Proceedings of the IEEE international conference on computer vision*, pages 706–715, 2017. 2, 5, 6

[27] Jie Lei, Licheng Yu, Tamara L Berg, and Mohit Bansal. Tvr: A large-scale dataset for video-subtitle moment retrieval. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXI 16*, pages 447–463. Springer, 2020. 2

[28] Guangyao Li, Yake Wei, Yapeng Tian, Chenliang Xu, Ji-Rong Wen, and Di Hu. Learning to answer questions in dynamic audio-visual scenarios. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 19086–19096, 2022. 7

[29] Jiaze Li, Yaya Shi, Zongyang Ma, et al. imove: Instance-motion-aware video understanding. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 23959–23975, 2025. 7

[30] Kunchang Li, Yinan He, Yi Wang, et al. Videochat: Chat-centric video understanding. *CoRR*, abs/2305.06355, 2023. 7

[31] Xinhao Li, Yi Wang, Jiashuo Yu, et al. Videochat-flash: Hierarchical compression for long-context video modeling. *CoRR*, abs/2501.00574, 2025. 2

[32] Xinhao Li, Ziang Yan, Desen Meng, et al. Videochat-r1: Enhancing spatio-temporal perception via reinforcement fine-tuning. *arXiv preprint arXiv:2504.06958*, 2025. 7

[33] Zinuo Li, Xian Zhang, Yongxin Guo, et al. Watch and listen: Understanding audio-visual-speech moments with multimodal LLM. *CoRR*, abs/2505.18110, 2025. 2, 7

[34] Chin-Yew Lin and Franz Josef Och. Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics, 21-26 July, 2004, Barcelona, Spain*, pages 605–612, 2004. 6

[35] Jing Liu, Sihan Chen, Xingjian He, Longteng Guo, Xinxin Zhu, Weining Wang, and Jinhui Tang. VALOR: vision-audio-language omni-perception pretraining model and dataset. *IEEE Trans. Pattern Anal. Mach. Intell.*, 47(2): 708–724, 2025. 5

[36] Zuyan Liu, Yuhao Dong, Ziwei Liu, et al. Oryx mllm: On-demand spatial-temporal understanding at arbitrary resolution. In *The Thirteenth International Conference on Learning Representations*, 2025. 6

[37] Zuyan Liu, Yuhao Dong, Jiahui Wang, et al. Ola: Pushing the frontiers of omni-modal language model with progressive modality alignment. *CoRR*, 2025. 6, 7

[38] Muhammad Maaz, Hanoona Abdul Rasheed, Salman Khan, and Fahad Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 12585–12602, 2024. 7

[39] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: An ASR corpus based on public domain audio books. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2015, South Brisbane, Queensland, Australia, April 19-24, 2015*, pages 5206–5210, 2015. 7

[40] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA*, pages 311–318, 2002. 6

[41] Long Qian, Juncheng Li, Yu Wu, et al. Momentor: Advancing video large language model with fine-grained temporal reasoning. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*, 2024. 2

[42] Alec Radford, Jong Wook Kim, Tao Xu, et al. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518, 2023. 6

[43] Shuhuai Ren, Linli Yao, Shicheng Li, et al. Timechat: A time-sensitive multimodal large language model for long video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14313–14323, 2024. 1, 6, 7

[44] Ramon Sanabria, Ozan Caglayan, Shruti Palaskar, et al. How2: A large-scale dataset for multimodal language understanding. *CoRR*, abs/1811.00347, 2018. 6

[45] Zhihong Shao, Peiyi Wang, Qihao Zhu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *CoRR*, abs/2402.03300, 2024. 4

[46] Gunnar A. Sigurdsson, Gül Varol, Xiaolong Wang, et al. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part I*, pages 510–526, 2016. 2, 5, 6

[47] Yixuan Su, Tian Lan, Huayang Li, et al. Pandagpt: One model to instruction-follow them all. In *Proceedings of the 1st Workshop on Taming Large Language Models: Controllability in the era of Interactive Assistants!*, pages 11–23, 2023. 7

[48] Yunlong Tang, Daiki Shimada, Jing Bi, et al. Empowering llms with pseudo-untrimmed videos for audio-visual temporal understanding. In *AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025, Philadelphia, PA, USA*, pages 7293–7301, 2025. 2, 7

[49] Vidi Team, Celong Liu, Chia-Wen Kuo, et al. Vidi: Large multimodal models for video understanding and editing. *CoRR*, abs/2504.15681, 2025. 2, 5

[50] Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chenliang Xu. Audio-visual event localization in unconstrained videos. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part II*, pages 252–268. Springer, 2018. 5

[51] Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chenliang Xu. Audio-visual event localization in unconstrained videos. In *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part II*, pages 252–268. Springer, 2018. 2

[52] Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 4566–4575, 2015. 6

[53] Yueqian Wang, Xiaojun Meng, Jianxin Liang, Yuxuan Wang, Qun Liu, and Dongyan Zhao. Hawkeye: Training video-text llms for grounding text in videos. *CoRR*, abs/2403.10228, 2024. 1, 7

[54] Ye Wang, Ziheng Wang, Boshen Xu, et al. Time-r1: Post-training large vision language model for temporal video grounding. *arXiv preprint arXiv:2503.13377*, 2025. 1, 4, 7

[55] Jianlong Wu, Wei Liu, Ye Liu, et al. A survey on video temporal grounding with multimodal large language model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025. 1, 2

[56] Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal llm. In *Forty-first International Conference on Machine Learning*, 2024. 7

[57] Bo Xiong, Yannis Kalantidis, Deepti Ghadiyaram, and Kristen Grauman. Less is more: Learning highlight detection from video duration. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 1258–1267. Computer Vision Foundation / IEEE, 2019. 2

[58] Jin Xu, Zhifang Guo, Jinzheng He, et al. Qwen2.5-omni technical report. *CoRR*, abs/2503.20215, 2025. 2, 7

[59] Jin Xu, Zhifang Guo, Hangrui Hu, et al. Qwen3-omni technical report. *CoRR*, abs/2509.17765, 2025. 2, 7

[60] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10714–10726, 2023. 2

[61] An Yang, Baosong Yang, Beichen Zhang, et al. Qwen2.5 technical report. *CoRR*, abs/2412.15115, 2024. 6

[62] Abhay Zala, Jaemin Cho, Satwik Kottur, et al. Hierarchical video-moment retrieval and step-captioning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 23056–23065, 2023. 2

[63] Xiangyu Zeng, Kunchang Li, Chenting Wang, et al. Timesuite: Improving mllms for long video understanding via grounded tuning. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025. 2, 7

[64] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 543–553, 2023. 2, 7

[65] Luowei Zhou, Chenliang Xu, and Jason J. Corso. Towards automatic learning of procedures from web instructional videos. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 7590–7598, 2018. 2

[66] Xingyi Zhou, Anurag Arnab, Shyamal Buch, et al. Streaming dense video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18243–18252, 2024. 4

# ChronusOmni: Improving Time Awareness of Omni Large Language Models

## Supplementary Material

## 9. Additional Details of ChronusAV Dataset

### 9.1. Additional Details of ChronusAV Construction

Our dataset construction pipeline includes video collection, video segmentation, modality-specific annotation, human verification, and training set and benchmark construction, as introduced in Section 5.1. In this section, we provide further implementation details specifically focusing on the modality-specific annotation details and the rigorous human verification process employed to ensure data quality.

- **Details of modality-specific annotation.** As introduced in Section 5.1, we use Gemini-2.5-Flash and Gemini-2.5-Pro to annotate video and audio segment caption, respectively.

  In video segment caption annotation, our prompt is: "The user will input a video. Please provide a brief description of the main visual content of the video. Avoid specific timestamps and keep the content concise. Avoid using indicative phrases like 'in the video' and directly output the visual information." For video segments that are 16 seconds or shorter, we sample frames at 1 fps; for video segments longer than 16 seconds, we uniformly sample 16 frames. We then input the sampled video frames along with the prompt into Gemini-2.5-Flash to obtain the output caption for each video segment.

  In audio segment caption annotation, our prompt is: "The user will input a audio. Briefly describe the audio information, including the original text of the speech, audio events, etc. Avoid using indicative phrases like 'in the audio' and directly output the audio information. Do not interpret the meaning of the speech expressed in the audio; just record what you hear concisely. For audio events, only output those that you are very certain about, and disregard any uncertain sounds. Record speech and audio events in order, but avoid specific timestamps. ##For example: A train whistle blows followed by a character speaking in a crisp voice, 'Hello, my name is John. I would like to help you learn the numbers.' This is all accompanied by the sound of a train moving on its tracks." We input the audio segment with this audio caption prompt into Gemini-2.5-Pro to obtain the output caption for each audio segment.

- **Details of human verification.** To ensure quality of LLM-generated captions, we conduct a human study, as introduced in Section 5.1. Three professional annotators are recruited to evaluate the annotations. They are presented with the sampled 1000 video and audio segments and asked to assess the generated captions based on two distinct criteria: *Semantic Accuracy* and *Modality Independence*.



Figure 5. Distribution of visual segment caption length (words).



Figure 6. Distribution of audio segment caption length (words).

For *Semantic Accuracy*, we utilize a 3-point Likert scale: (1) *Accurate*: The caption perfectly describes the visual/audio content; (2) *Acceptable*: The caption is generally correct but misses minor details or contains slight hallucinations; (3) *Inaccurate*: The caption is missing important details or contains significant hallucinations.

For *Modality Independence* (Cross-modal Leakage), annotators checked for information leakage between modalities (e.g., visual captions describing auditory events like "loud explosion"). This was rated as: (1) *No Leakage*, (2) *Minor Leakage*, or (3) *Severe Leakage*.

To ensure consistency, we calculate the Fleiss' Kappa score among annotators, yielding $\kappa = 0.82$, indicating high inter-annotator agreement. Final labels are determined by majority voting.

Results show that video captions are mostly accurate (rated as Accurate or Acceptable) in 96.1% of cases, and audio captions are mostly accurate in 93.5%. Furthermore, 99.3% of video captions and 97.5% of audio captions show no or only minor cross-modal leakage. This high level of quality assurance validates the reliability of our automatically generated annotations for large-scale audiovisual temporal grounding training and evaluation.

```
# Segment 1
    Timestamp: ["0:00:00.0", "0:00:02.9"]
    Visual segment caption: The Universal Pictures logo appears with a globe and starfield background, followed by the MRC logo in
    blue text within a glowing rectangular frame against a black screen. The screen then turns completely black.
    Audio segment caption:  A low, resonant boom is followed by a gentle piano melody.

# Segment 2
    Timestamp: ["0:00:02.9", "0:00:06.0"]
    Visual segment caption: A view across dark blue ocean water towards a distant rugged coastline featuring cliffs, green hills, and a
    sandy beach.
    Audio segment caption: A character speaks in a soft, slightly trembling voice, "I was eight years old when my mother died,"
    accompanied by soft, somber music.

    ● ● ● ● ● ●

# Segment 11
    Timestamp: ["0:01:34.6", "0:01:45.4"]
    Visual segment caption: A large, shadowy figure with glowing green eyes and hands appears, creating a menacing presence. Several
    individuals, including men and women, react with expressions of fear and surprise to the imposing figure.
    Audio segment caption: Intense, dramatic music with metallic clanging sounds is followed by a monstrous roar. A female voice
    whispers, "He's here," as the dramatic music swells.

# Segment 12
    Timestamp: ["0:01:45.4", "0:02:17.2"]
    Visual segment caption: A sprawling, multi-tiered city atop a mountain, grand futuristic architecture, and large-scale mechanical
    structures. Scenes of intense action, explosions, and destruction are interspersed with close-ups of characters displaying fear,
    determination, and concern. Landscapes transition from cloudy skies to fiery, volcanic terrains and barren expanses.
    Audio segment caption: Dramatic, orchestral music plays. A male voice says, "In the great game of survival, this is checkmate."
    Another male voice says, "I knew you wouldn't leave me." A female voice responds, "Shut up and run." A male voice asks, "You sure
    you want to do this?" A female voice replies, "I have to. For my mother." A different female voice says, "You look at her and all you
    see are the jagged edges. But she is something quite different." A male voice adds, "She is beautiful and strange." The second female
    voice concludes, "And very, very rare."

# Segment 13
    Timestamp: ["0:02:17.3", "0:02:36.1"]
    Visual segment caption: Movie title card followed by film credits, concluding with a split black and white screen with text
    encouraging subscriptions and suggesting another video.
    Audio segment caption: Dramatic and suspenseful orchestral music.
```

Figure 7. Timestamps, visual captions and audio captions of a 157-second video with corresponding audio in ChronusAV.

## 9.2. Additional Analyses of ChronusAV Dataset

Figures 5 and Figures 6 illustrate the distribution of caption lengths (in words) for visual and audio segments, respectively. As shown in Figure 5, the visual segment captions exhibit a unimodal, quasi-normal distribution. The vast majority of visual captions fall within the range of 15 to 40 words, with a distinct peak in the [20, 25) interval. The distribution tails off smoothly, with negligible instances of captions exceeding 80 words, suggesting a high degree of consistency in the length of visual descriptions. Figure 6 show that the audio segment caption length peaks in the [10, 20) interval, and exhibits a clear long-tail distribution. The long-tail distribution is due to the fact that some audio segments contains rich speech information and the detailed speech transcriptions significantly increase the length of the audio captions. The average length of visual segment captions is 29.6 words, while the average length of audio segment captions is 50.6 words, indicating that our captions are quite detailed.

A case showing the annotations of a 157-second video in ChronusAV, is presented in Figures 7. For each untrimmed video with corresponding audio, we annotate the timestamps, visual segment captions and audio segment captions. Each (timestamp, visual segment caption, audio segment caption) tuple can be used to construct 6 types of QA in the audiovisual temporal grounding task. For example, we can use the timestamp, visual segment caption, audio segment caption of segment 11 in Figures 7 to construct 6 types of question-answer pairs, as shown in Figures 8.

Figure 8. A case of question and answers for 6 subtasks in the audiovisual temporal grounding task.

# 10. Additional Details and Results of Experiments

## 10.1. Training Hyperparameters

In supervised fine-tuning stage, we train on 70K data samples for 1 epoch, as introduces in Section 6.2. We employ a batch size of 16. Optimization is guided by a cosine learning rate scheduler with a peak learning rate of $5 \times 10^{-6}$, a warmup ratio of 0.05, a minimum learning rate ratio of 0.01, and no weight decay. The maximum context length is configured to 32K tokens. We trained for about 17 hours using $8\times$NVIDIA A800-80GB GPUs in this stage.

In reinforcement learning stage, we train on 4K data samples for 1 epoch, using a batch size of 4 for 1000 steps. During the rollout process within GRPO, the generation temperature is set to 1.0, and the maximum number of generated tokens is 1024. we sample 4 generations per prompt. The KL divergence penalty coefficient $\beta$ is set to 0.04 to ensure controlled deviation from the initial policy. Learning rate is set to $1 \times 10^{-6}$. We trained for about 27 hours using $4\times$NVIDIA A800-80GB GPUs in this stage.

## 10.2. Additional Details of Evaluation Settings

For all the models we evaluate, we use the officially recommended number of frames, for example, 8 frames for VideoLLaMA, 64 frames for Ola, 100 frames for Avicuna and LongVALE-LLM, 150 frames for ARC-Hunyuan-Video, 2 fps (for videos shorter than 384 seconds) or 768 frames (for videos longer than 384 seconds) for Qwen2.5-Omni and Qwen3-Omni. To ensure fairness and reproducibility, we use greedy sampling for all models on evaluation. Due to different models using various time formats during training (e.g. frame index for AVicuna and LongVALE-LLM, absolute time format "HH:MM:SS" for ARC-Hunyuan-Video), we convert the time format in the ChronusAV benchmark to match the format used by the models during training in order to ensure the best inference results.

## 10.3. Frame Sampling Ablations

We perform an ablation study on the number of sampled frames during the inference stage, with the model trained on 64 frames. As shown in Table 6, the model achieves op-

Table 6. Ablation study on sampled frame number of ChronusOmni. B: BLEU-4. R: ROUGE-L. M: METEOR. C: CIDEr.

| Model | Frame Number | V2T | | T2V | | | | A2T | | T2A | | | | V2A | | | | A2V | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | R@0.5 | R@0.7 | B | R | M | C | R@0.5 | R@0.7 | B | R | M | C | B | R | M | C | B | R | M | C |
| ChronusOmni | 32 | 51.40 | 32.95 | 1.02 | 3.12 | 1.93 | 4.44 | 83.85 | 66.90 | 6.32 | 6.36 | 4.22 | 27.02 | 3.15 | 4.46 | 3.14 | 10.73 | 0.49 | 2.34 | 1.84 | 0.43 |
| ChronusOmni | 64 | **63.15** | **45.95** | 1.16 | 3.37 | 2.12 | 5.07 | **90.50** | 79.85 | **6.78** | **6.86** | **4.50** | 34.30 | 3.61 | **4.90** | **3.27** | **13.60** | **1.01** | 3.17 | **2.12** | 3.03 |
| ChronusOmni | 128 | 60.20 | 43.85 | **1.20** | **3.42** | **2.14** | **5.24** | 90.30 | **82.90** | 6.38 | 6.76 | 4.23 | **34.62** | **3.62** | 4.76 | 3.00 | **13.60** | 1.00 | **3.19** | 2.11 | **3.43** |

Table 7. Comparison with base model on inference efficiency. ChronusOmni incurs only a slight increase in inference latency compared to the base model, striking a favorable balance between performance and efficiency.

| Model | Avg. CPU Time | Avg. GPU Time | Avg. Infer Time |
|---|---|---|---|
| Base model (Ola) | 1.51 s | 2.01 s | 3.52 s |
| ChronusOmni | 1.55 s | 2.18 s | 3.73 s |

timal performance across most subtasks when the inference frame number matches the training setting (64 frames). Reducing the frames to 32 results in a noticeable performance drop due to information loss. While increasing the sampling rate to 128 frames slightly improves some metrics by providing denser temporal guidance, it also degrades performance in some metrics. This decline is likely attributed to the gap between training and inference. Consequently, we maintain 64 frames as the default inference setting to ensure consistency and maximize performance.

## 10.4. Inference Efficiency Evaluation

Due to the introduction of additional time tokens and the token interleave process compared to the base model, We clearly compare the inference efficiency in Table 7 to analyze the computational burden brought by our approach. We test the average inference CPU and GPU time of the base model and ChronusOmni over 2000 samples of the A2T subtask in ChronusAV benchmark on a single A800 GPU. Compared with the base model, ChronusOmni introduces negligible computational overhead. Specifically, the average total inference time increases by only 0.21 seconds (from 3.52 s to 3.73 s), demonstrating that our method maintains the high efficiency of the backbone while incorporating additional capabilities.

## 10.5. Qualitative Results

We provide qualitative results of ChronusOmni and other two audiovisual LLMs in various subtasks of ChronusAV benchmark in Figures 9-14. As shown in Figures 9 and Figures 11, ChronusOmni can locate more precise time boundaries. Figures 10, 12, 13 and 14 shows ARC-Hunyuan-Video and Qwen3-Omni often describe information outside of the time periods in the time-related caption subtasks (T2V, T2A, V2A, A2V), or they may omit information within the time periods. In comparison, ChronusOmni's output demonstrates stronger temporal awareness.
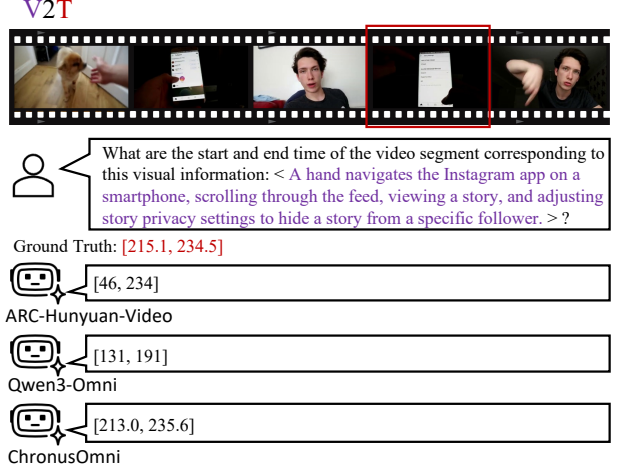


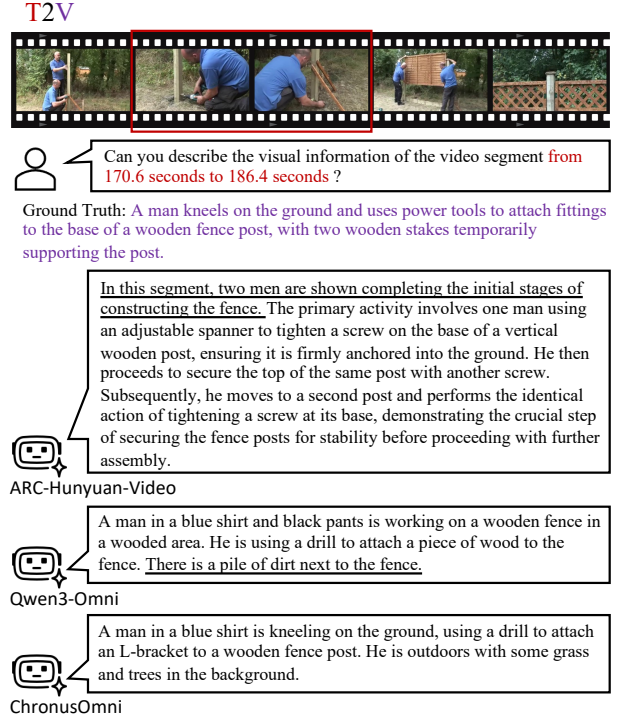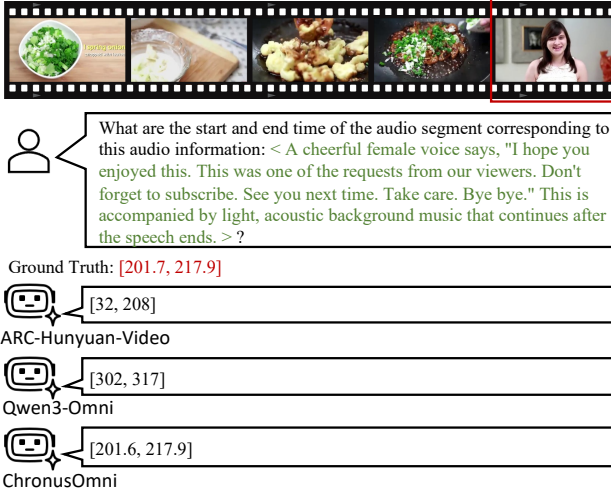Figure 9. Qualitative results on V2T subtask. The sample is from ChronusAV benchmark.



Figure 10. Qualitative results on T2V subtask. The sample is from ChronusAV benchmark. The underlined sentence indicates the incorrect part.

## A2T



**User:** What are the start and end time of the audio segment corresponding to this audio information: < A cheerful female voice says, "I hope you enjoyed this. This was one of the requests from our viewers. Don't forget to subscribe. See you next time. Take care. Bye bye." This is accompanied by light, acoustic background music that continues after the speech ends. > ?

Ground Truth: [201.7, 217.9]

**ARC-Hunyuan-Video:** [32, 208]

**Qwen3-Omni:** [302, 317]

**ChronusOmni:** [201.6, 217.9]

Figure 11. Qualitative results on A2T subtask. The sample is from ChronusAV benchmark.

## T2A



**User:** Can you describe the audio information of the audio segment from 42.7 seconds to 60.2 seconds ?

Ground Truth: A man says, "you could find out where a good food is." A soft click is heard, followed by another man asking, "Where can I get a falafel?" A beep sounds, and then a woman's voice states, "Looking for restaurants whose reviews mention falafel. I found quite a few restaurants a little ways from you." The first man continues, "and it'll instantly program it into your nav system and the car, it'll take you right there."

**ARC-Hunyuan-Video:** In this segment, a man asks the car's voice assistant where to find a restaurant with a review for falafel. The car's automated voice immediately responds, stating it has located several nearby restaurants that fit the criteria. The man concludes by mentioning that the system will automatically program the selected address into his navigation system.

**Qwen3-Omni:** The audio features a male speaker discussing the functionality of a car's navigation system. He explains how the system can find addresses, locate places to eat, and program destinations into the navigation system. The speaker also mentions that the system can handle text messages and music playback.

**ChronusOmni:** A man speaks, "you could find where you want to go. You could find out where good food is." A beep sounds, followed by a female voice speaking, "Where can I get a falafel?" Another beep sounds, and the female voice continues, "Looking for restaurants whose reviews mention falafel. I found quite a few restaurants a little ways from you." The man then says, "And it'll instantly program it into your nav system and the car, it'll take you right there."

Figure 12. Qualitative results on T2A subtask. The sample is from ChronusAV benchmark. The underlined sentence indicates the incorrect part.

## V2A



**User:** Can you describe the audio information of the audio segment corresponding to this visual information: < Person's hands attempting to open a plastic bag containing a brown substance, with a pot nearby. The word "STRUGGLES" is displayed on screen. > ?

Ground Truth: A person says, "It's a lot of oil," accompanied by the sound of a plastic wrapper crinkling. A bell rings, followed by a light laugh.

**ARC-Hunyuan-Video:** The audio segment corresponding to the visual of the person struggling to open a plastic bag is a loud, hissing or crinkling sound.

**Qwen3-Omni:** The audio features a person struggling to open a plastic bag, accompanied by the sound of crinkling plastic. The person says, "It's a lot of oil," and then, "Alright, now pour that in." The audio also includes the sound of a whisk being used to mix ingredients in a bowl.

**ChronusOmni:** A person speaks, "It's a lot of oil." This is followed by the sound of plastic crinkling and a loud bang. The person then laughs.

Figure 13. Qualitative results on V2A subtask. The sample is from ChronusAV benchmark. The underlined sentence indicates the incorrect part.

## A2V



**User:** Can you describe the visual information of the video segment corresponding to this audio information: < A clinking sound is heard, followed by a person speaking, "Next, add a half a cup of water." This is accompanied by the sound of liquid being poured and a thud. The person continues, "Four teaspoons of pure vanilla extract," as more liquid is poured, followed by another clinking sound. The person then says, "Four tablespoons of dark cocoa powder," which is accompanied by tapping sounds. > ?

Ground Truth: A woman in a kitchen adds ingredients to a blender on a countertop. She has short, curly dark hair and wears a grey top with a light-colored apron. The kitchen features wooden cabinets and a white stove.

**ARC-Hunyuan-Video:** A woman introduces her recipe for a healthy chocolate pudding and adds two peeled avocados into a blender. She adds maple syrup, water, vanilla extract, cocoa powder, and sea salt to the blender. The woman puts the lid on the blender, inserts the tamper, and blends all the ingredients together. She pours the finished chocolate pudding into a glass, tastes it, and shows the final product along with the recipe.

**Qwen3-Omni:** A woman in a kitchen prepares a chocolate pudding recipe using a blender. She adds ingredients such as avocado, maple syrup, water, vanilla extract, cocoa powder, and sea salt into the blender. After blending the mixture, she pours it into a glass bowl and garnishes it with a raspberry and mint leaf before taking a bite.

**ChronusOmni:** The video shows a woman in a kitchen, wearing an apron and a grey top, standing next to a blender. She is adding ingredients to the blender, which already contains some green ingredients. She pours water from a glass into the blender, followed by liquid from another bowl, then adds vanilla extract and cocoa powder.
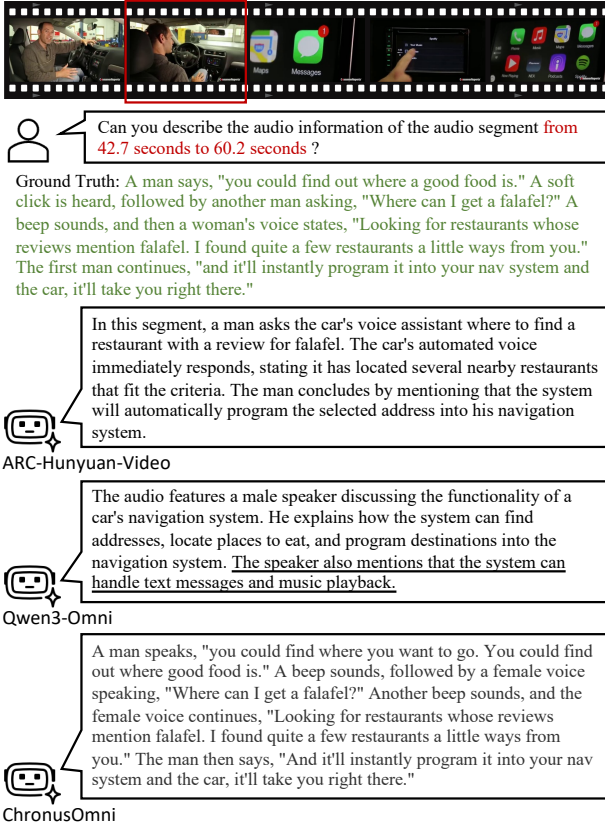
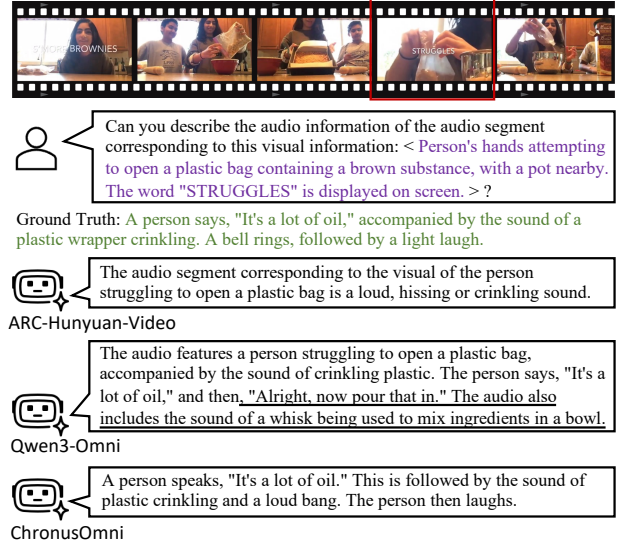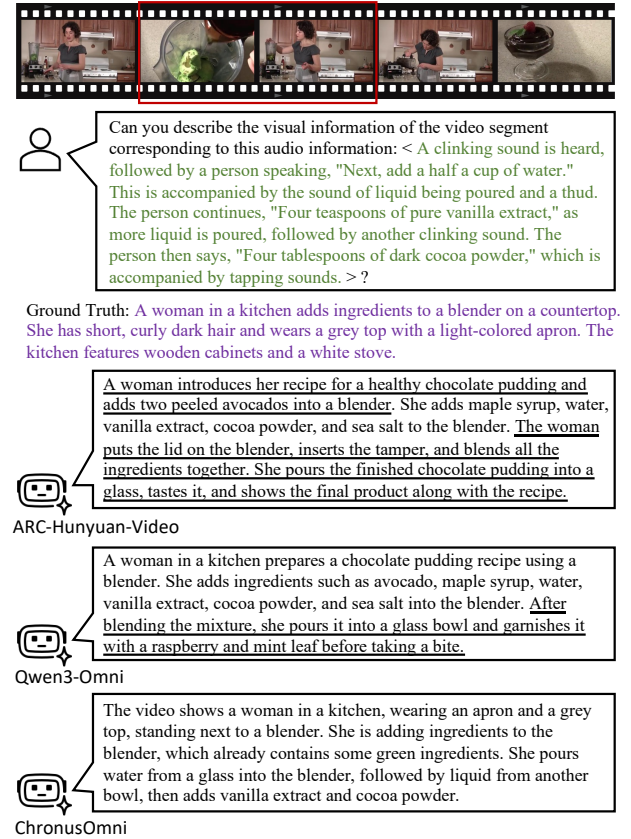Figure 14. Qualitative results on A2V subtask. The sample is from ChronusAV benchmark. The underlined sentence indicates the incorrect part.