

# Strefer: Empowering Video LLMs with Space-Time Referring and Reasoning via Synthetic Instruction Data

Honglu Zhou, Xiangyu Peng, Shrikant Kendre, Michael S. Ryoo, Silvio Savarese,  
Caiming Xiong, Juan Carlos Nieves  
Salesforce AI Research

<https://strefer.github.io>

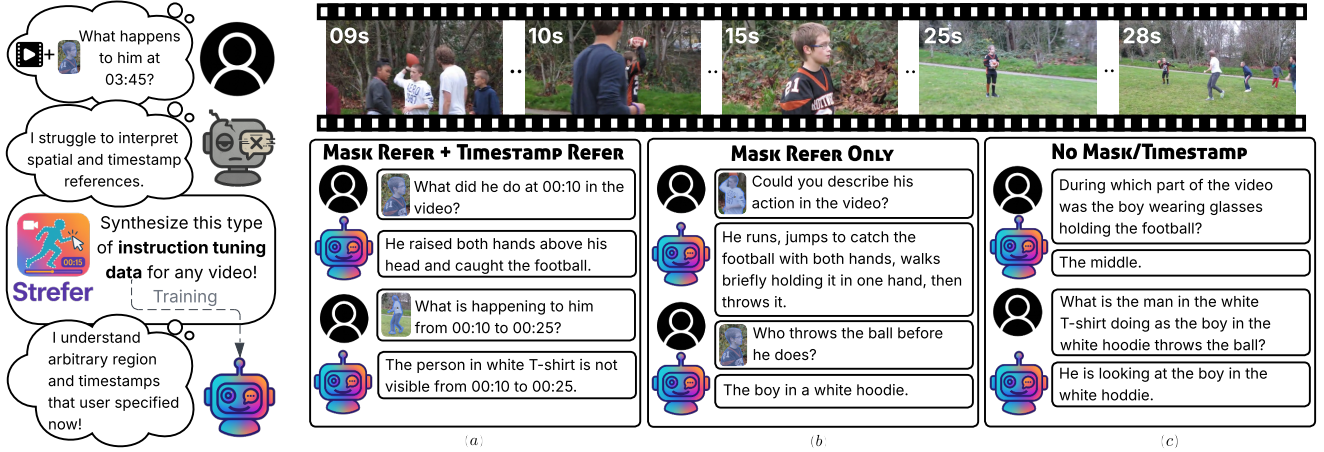


Figure 1. Our goal is to synthesize instruction-response pairs through a scalable, grounded approach that enhances fine-grained spatial and temporal perception and reasoning over videos for tuning Video LLMs. We hypothesize that the video space-time referring task, which requires understanding user-specified regions within a video at particular moments or over defined time intervals, can greatly enhance the effectiveness and versatility of video analysis with Video LLMs, facilitating advanced applications in areas including navigation, surveillance, and interactive robotics. To this end, we introduce **Strefer** (Figure 3), which synthesizes instruction data aligned with our goals. **Strefer** produces instruction-response pairs for: (a) queries with mask and timestamp references, (b) queries with mask (or timestamp) references only, and (c) queries without either. Though current implementation of **Strefer** does not any use proprietary models, without the need to annotate large volumes of new videos, instruction data from **Strefer** empowers models for space-time referring and spatiotemporal reasoning (ref. Table 2, 3, and 4). Using **Strefer**, we generated 947, 854 instruction data points from just 4, 253 NExT-QA [56] videos (test set excluded; up to 3 minutes long). While our final recipe, consisted exclusively of short videos a few minutes long and added only 545 more videos compared to the baseline, it led to performance gains across multiple benchmarks.

## Abstract

Next-generation AI companions must go beyond general video understanding to resolve spatial and temporal references in dynamic, real-world environments. Existing Video Large Language Models (Video LLMs), while capable of coarse-level comprehension, struggle with fine-grained, spatiotemporal reasoning, especially when user queries rely on time-based event references for temporal anchoring, or gestural cues for spatial anchoring to clarify object references and positions. To bridge this critical gap, we introduce **Strefer**, a synthetic instruction data generation framework designed to equip Video LLMs with spatiotemporal referring and reasoning capabilities. **Strefer** produces diverse instruction-tuning data using a data engine that pseudo-annotates temporally dense, fine-grained

video metadata, capturing rich spatial and temporal information in a structured manner, including subjects, objects, their locations as masklets, and their action descriptions and timelines. Our approach enhances the ability of Video LLMs to interpret spatial and temporal references, fostering more versatile, space-time-aware reasoning essential for real-world AI companions. Without using proprietary models, costly human annotation, or the need to annotate large volumes of new videos, experimental evaluations show that models trained with data produced by **Strefer** outperform baselines on tasks requiring spatial and temporal disambiguation. Additionally, these models exhibit enhanced space-time-aware reasoning, establishing a new foundation for perceptually grounded, instruction-tuned Video LLMs.

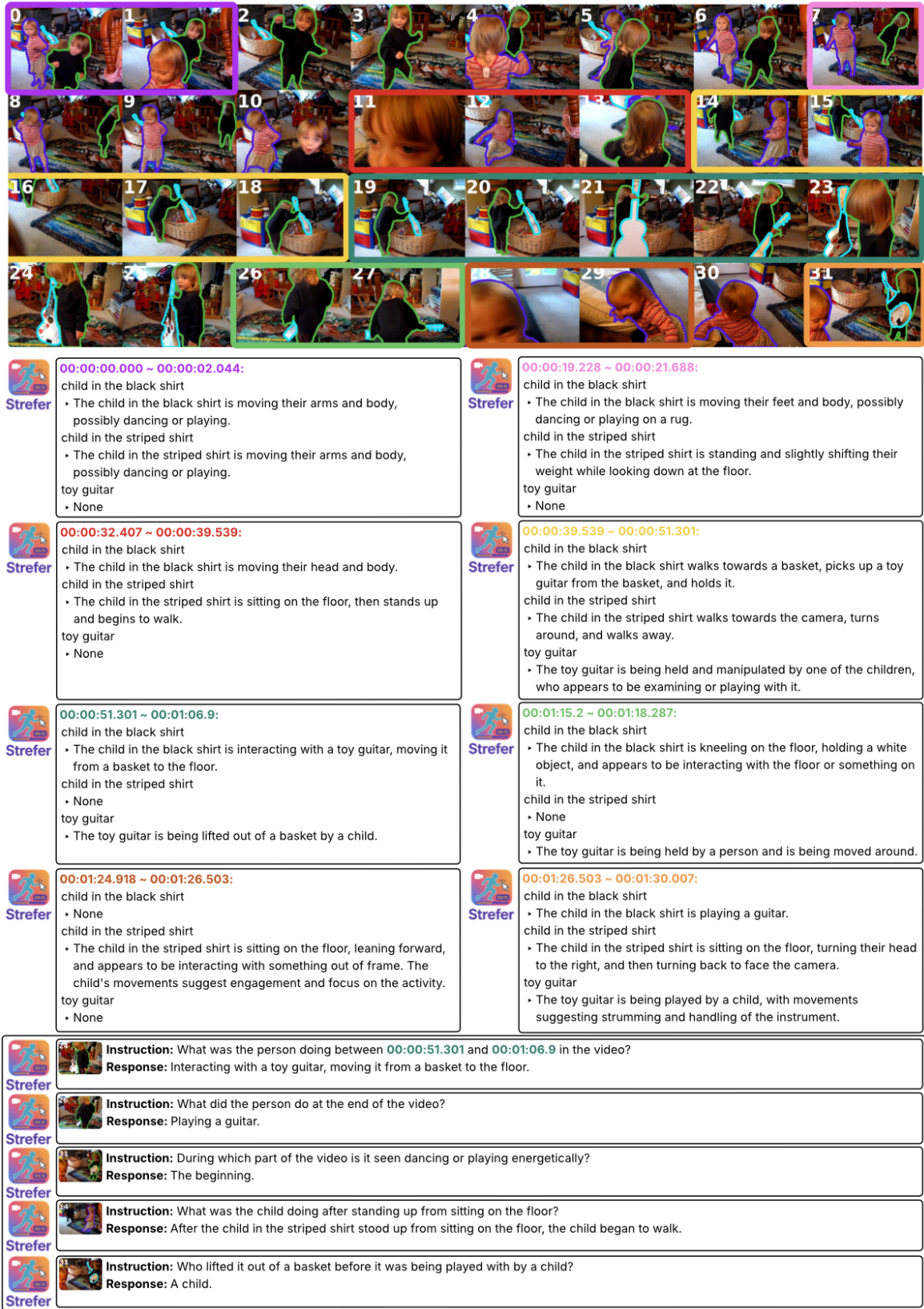


Figure 2. Example of **Strefer**-annotated instruction-response pairs (bottom) and video metadata (top). Each instruction begins with the prefix: *Please answer the following question about the <region>* (omitted in the figure). For each instruction-response pair, the boundary of object mask referred to by <region> is shown beside the pair. **Strefer** automatically clips the video into segments and pseudo-annotates the video metadata, including active entities, their locations (as masklets), and their action descriptions and timelines for complex video scenarios, such as scenes containing multiple entities of the same category, and cases where entities do not appear in the first frame, or temporarily exit and re-enter the frame. Based on the auto-generated video metadata, **Strefer** produces instruction-response pairs requiring no legacy annotations or manual efforts.



# Table of Contents

## Contents

<b>1. Introduction</b>	<b>4</b>
<b>2. Methodology</b>	<b>4</b>
2.1. <b>Strefer</b> : Automatic Data Engine	4
2.1.1 . Entity Recognizer	5
2.1.2 . Referring Parser	5
2.1.3 . Referring Masklet Generator	5
2.1.4 . Video Clipper	7
2.1.5 . Video Transcriber	7
2.1.6 . Video Instruction Data Generator	7
2.2. Synthesized Data	8
2.3. Modeling for Space-Time Referring	8
<b>3. Evaluation Details</b>	<b>8</b>
3.1. Training Details	8
3.2. Benchmark Details	9
<b>4. Experiments</b>	<b>9</b>
4.1. Quantitative Results	9
4.1.1 . Main Results and Findings	9
4.1.2 . Discussions and Insights	10
4.1.3 . Visual Prompting for Space-Time Referring	12
4.2. Qualitative Results	13
4.2.1 . <b>Strefer</b> -Synthesized Instruction Data	13
4.2.2 . Our Referring Masklet Generation	14
4.2.3 . <b>Strefer</b> -Trained Model	14
<b>5. Limitations and Future Directions</b>	<b>15</b>
<b>6. Related Work</b>	<b>17</b>
<b>7. Conclusion</b>	<b>18</b>
<b>A Appendix / Supplemental Material</b>	<b>22</b>
A.1. Additional Qualitative Results	22
A.1.1. Qualitative Results of <b>Strefer</b> -Synthesized Data	22
A.1.2. Qualitative Results of <b>Strefer</b> -Trained Model	22
A.2. Additional <b>Strefer</b> Details	22
A.3. Model Details	24
A.3.1. Architecture Overview	24
A.3.2. Video Token Representation	29
A.3.3. Masklet Reference Token Representation	30
A.3.4. Timestamp Reference Token Representation	31

## 1. Introduction

The vision of real-life AI assistants as everyday companions is rapidly becoming a reality. To function seamlessly in real-world environments, these agents must understand human queries that are grounded in both space and time. People often rely on non-verbal cues such as gestures that clarify which object is being referred to, or time-based references tied to specific moments in the past. For instance, a user might point at a cup and say, “Can you bring me that cup?”, a request that may be ambiguous without detailed verbal clarification if several similar looking cups are nearby. They may also ask questions using time-based references: “At 11 a.m. today, whom did I talk to at the market?” These scenarios highlight the critical need for space-time referring capabilities, where an AI agent must resolve queries with references that are not purely linguistic, but situated in dynamic spatio-temporal contexts.

Recent Video Large Language Models (Video LLMs) have shown promise in general video understanding [3, 7, 8, 10, 17, 24, 28, 45, 52, 59, 60, 70]. However, these models tend to operate at a coarse level, lacking the granularity required to track and reason about specific object states, movements, and temporal relations [9, 33]. Their limitations are especially evident in videos with complex spatial and temporal structures, where multiple entities interact and change over time. They are also incapable of handling user queries that are not purely verbal but contain specific space-time references (*ref.* Fig. 1). This shortcoming stems not only from architectural constraints, but more significantly, from the scarcity of fine-grained, object-level instruction tuning data focused on spatial and temporal understanding, referring and reasoning within complex videos.

To address the limitations of current video instruction datasets, we introduce **Strefer** (*ref.* Fig. 3), a novel data engine that systematically generates synthetic, fine-grained, spatiotemporally and semantically rich instruction data for training Video LLMs on space-time reference and reasoning tasks. Our approach begins with a modular framework that orchestrates multiple agents—including pre-trained Large Language Models (LLMs), Video LLMs, and Pixel-Level Multimodal Vision Foundation Models (e.g., RexSeek [20], GroundingDINO [32] and SAM2 [44])—to *pseudo-annotate* video metadata with temporally dense and object-centric space-time information. This metadata captures detailed spatial and temporal structures, such as subjects, objects, their locations as masklets (segmentation masks tracked over time), and action timelines. Building on this structured metadata, we leverage in-context learning and well-defined task schemas to guide LLMs in generating high-utility instruction data for tuning Video LLMs.

Unlike existing datasets and synthesis approaches, which often rely on legacy annotations [16, 25, 26, 35, 38, 47] or limited-scale but high-cost human labeling [2], **Strefer**

automatically produces instruction-response pairs grounded in object-centric, spatiotemporal video structures at scale. Furthermore, our framework supports the generation of multimodal user prompts that mimic realistic human-AI interactions in dynamic environments. These prompts encourage models to reason about spatial references<sup>1</sup> [36] (e.g., single frame or trajectories of masks<sup>2</sup>) and temporal dynamics (e.g., event sequences, temporal dependencies, and specific time-based moment anchoring). Crucially, our system is designed to handle complex scenarios that challenge existing data synthesis methods—such as scenes containing multiple entities of the same category and cases where entities do not appear in the first frame, or temporarily exit and re-enter the frame—and is capable of processing minutes-long videos rather than just short clips a few seconds long, thereby fostering more robust and generalizable spatiotemporal perception and reasoning in Video LLMs.

Through the development of **Strefer**, we make the following key contributions:

- We provide a scalable methodology for pseudo-labeling videos with temporally dense, object-centric, structured space-time metadata, as well as for generating synthetic instruction-response pairs that promote nuanced spatial and temporal perception and reasoning over videos.
- We operationalize data synthesis for an underexplored but essential capability of real-world AI agents—fine-grained, spatiotemporal, object-centric referencing in video-based queries—paving the way for more perceptually grounded and context-aware interactions.
- We demonstrate that Video LLMs trained on our data outperform baselines in tasks requiring spatial and temporal disambiguation. These models also show improved space-time-aware reasoning, marking a crucial step toward AI systems that can understand and act within the full space-time fabric of our visual world.

## 2. Methodology

In the following, Sec. 2.1 details the data synthesis process, Sec. 2.2 outlines the resulting dataset, and Sec. 2.3 presents modeling choices (w or w/o model architectural changes) for space-time referring in general-purpose Video LLMs.

### 2.1. Strefer: Automatic Data Engine

Figure 3 illustrates the **Strefer** framework for video pseudo-annotation and instruction data generation. All components in the **Strefer** framework are pre-trained, frozen, open-source models, and used off-the-shelf.

<sup>1</sup>Spatial visual prompts may be derived from user interactions (e.g., clicking) or from user gestures, including finger pointing and eye gaze.

<sup>2</sup>To support diverse, free-form spatial reference from users (e.g., points, scribbles, and boxes), we choose segmentation masks as the model input because masks contain richer information and many forms of spatial reference can be easily transformed into masks using off-the-shelf tools like SAM2 [43].

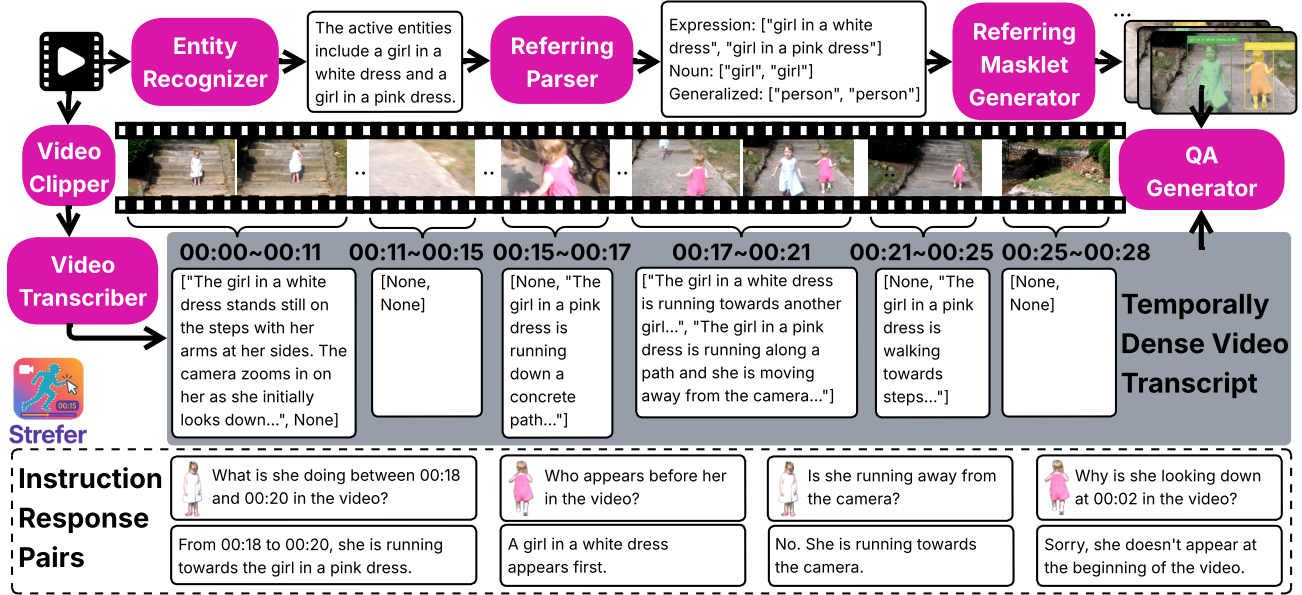


Figure 3. We introduce **Strefer**, a novel data engine that automatically generates synthetic instruction data—without manual effort or legacy annotation—featuring multimodal prompts grounded in complex spatiotemporal video structures, designed to train Video LLMs for space-time referring and reasoning tasks (Sec. 2.1). By design, **Strefer** handles challenging scenarios—such as scenes containing multiple entities of the same category and cases where entities do not appear in the first frame or temporarily exit and re-enter the frame—while scaling to minutes-long videos beyond the scope of existing auto-data engines in the video domain. The current implementation of **Strefer** does not use any proprietary models. The *Entity Recognizer* and *Video Transcriber* are Video LLMs (Tarsier-34b [52]), the *Referring Parser* is an LLM (Qwen2.5-32B-Instruct [48]), and the *QA Generator* uses either templates or an LLM (Qwen2.5-32B-Instruct [48]). The *Video Clipper* is based on PySceneDetect [4], SigLIP [69], and hierarchical clustering; and Fig. 4 illustrates our novel *Referring Masklet Generator*, leveraging GroundingDINO [32], SAM2 [43], and RexSeek [20]. Without the need to annotate large volumes of new videos, instruction data from **Strefer** empowers models for space-time referring and spatiotemporal reasoning (ref. Table 2, 3, and 4).

### 2.1.1. Entity Recognizer

Given an input video, a Video LLM is employed as an *Active Entity Recognizer*, tasked with identifying all active entities present throughout the video. Here, an entity refers to a person, an animal, or an object, while an active entity is defined as any entity that exhibits dynamic behavior—such as movement, interaction with other objects, or performing actions. To guide the model’s recognition process, we explicitly provide the definition of active entity in the prompt (see Appendix). We also instruct the model to use language that can clearly distinguish each entity from others.

### 2.1.2. Referring Parser

Since the Active Entity Recognizer outputs a descriptive paragraph, we use an LLM to extract three structured lists: (1) entity-referring expressions (e.g., “girl in a white dress”), (2) their noun categories (e.g., “girl”), and (3) generalized categories (e.g., “girl” → “person”, “parrot” → “bird”). Shortening referring expressions and extracting their broader concepts benefits the subsequent referring masklet generation module.

### 2.1.3. Referring Masklet Generator

The entity-referring expressions and their generalized categories, together with the video input, are passed to the

*Referring Masklet Generator* (Fig. 4) which generates masklets corresponding to each referring expression.

**Challenges:** Prior methods such as GroundedSAM2 [44] face several significant challenges when applied to the task of referring masklet generation in complex videos. While GroundedSAM2 supports both image and video input for referring segmentation and tracking, it struggles with multi-word referring expressions—especially in scenes with multiple entities sharing the same noun. This limitation stems from its use of spaCy [49] to extract a short noun phrase from complex expressions, which is then used as a prompt for GroundingDINO [32], an open-vocabulary object detector that expects simple object nouns. Once the bounding boxes of simple object nouns on the first frame are identified, GroundedSAM2 employs SAM2 [42] to propagate bounding boxes through subsequent frames of the video and generate segmentation masklets.

In contrast, RexSeek [20] is recently introduced which handles complex referring expressions and performs well in disambiguating entities with similar object nouns. However, it is limited to static images and cannot manage video-specific challenges such as motion blur, occlusion, or object re-identification across frames.



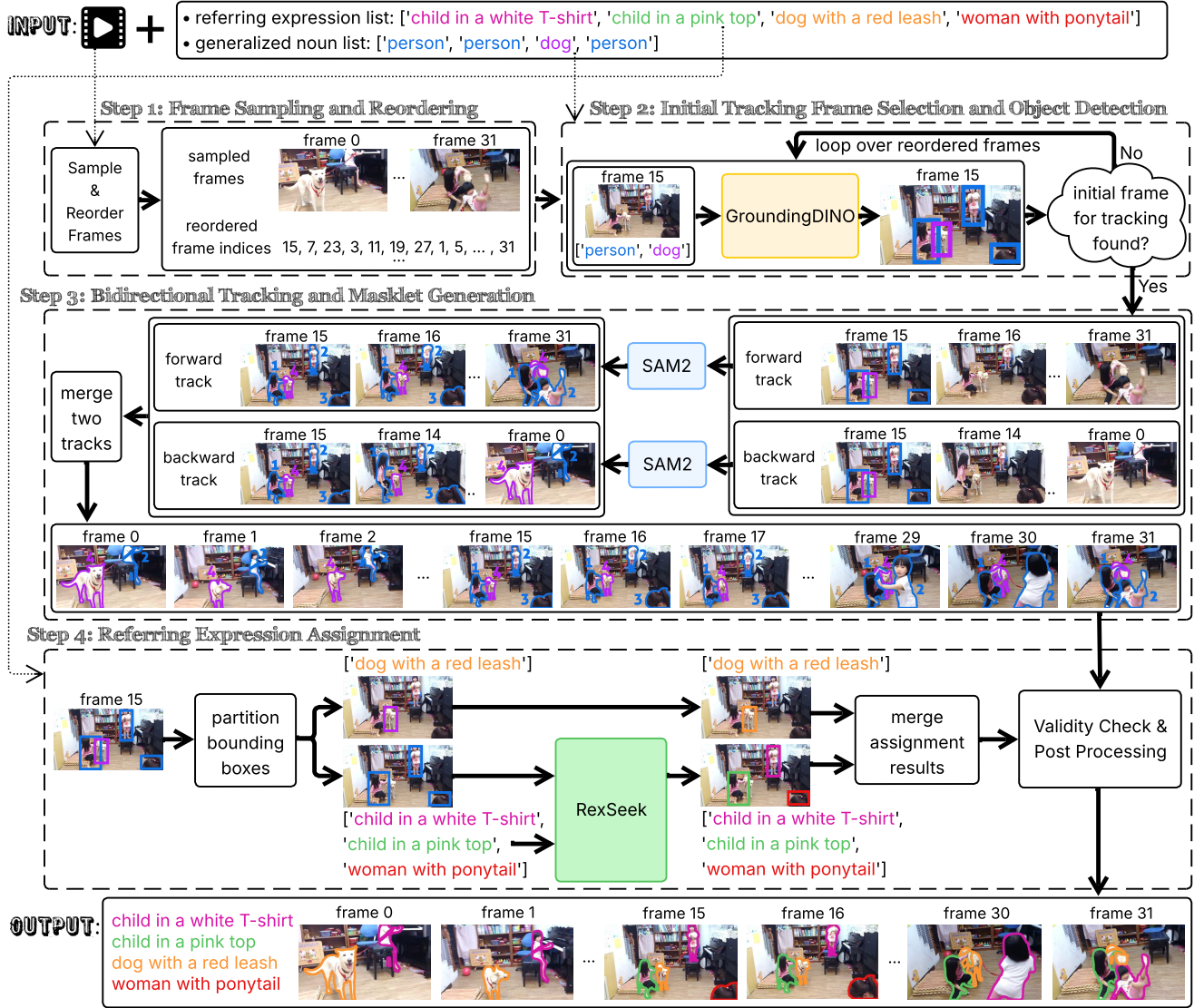


Figure 4. **Overview of the Referring Masklet Generation Pipeline within Strefer.** This pipeline produces tracked segmentation masks from videos with complex structures based on multi-word natural language referring expressions. Our masklet generator is carefully crafted to address key limitations overlooked by prior works [22, 37, 68] by orchestrating complementary strengths of the state-of-the-art pixel-level vision foundation models to achieve more effective results. It robustly handles challenging scenarios, including multiple same- or similar-category entities described differently, entities absent in the first frame, and entities that temporarily exit and re-enter the scene.

Moreover, GroundedSAM2’s performance degrades significantly when the target object is not visible in the first frame. Although SAM2 can handle scenarios where objects temporarily disappear and reappear, it assumes that tracking begins from an ideal, user-provided starting frame in which all target objects are visible. Consequently, selecting a suitable starting frame is critical for achieving effective tracking performance with SAM2. However, GroundedSAM2 automatically initializes SAM2’s tracking from the video’s first frame, without any mechanism to select a more appropriate frame or to handle complex video conditions.

**Our Referring Masklet Generator:** To address these challenges, we introduce a modular pipeline (Fig. 4).

**Step 1: Frame Sampling and Reordering.** The video is sampled into frames, which are then reordered using a heuristic that assumes important content typically occurs near the middle of the video. This reordering facilitates a more efficient search for a suitable tracking start frame.

**Step 2: Initial Tracking Frame Selection and Object Detection.** This step identifies an initial frame for tracking and simultaneously performs object detection using generalized entity nouns. Our key finding is that GroundingDINO performs more reliably when prompted with simple nouns (e.g., “person” instead of “bride”), and therefore, we use generalized entity nouns as prompts. The selected starting frame is the first in the reordered frame list where

the number of detected objects matches or exceeds the number of referring expressions. If no such frame exists, the frame with the highest number of detected objects is used.

**Step 3 : Bidirectional Tracking and Masklet Generation.** Using SAM2, tracking is performed both forward and backward from the selected initial frame, with object bounding boxes of generalized nouns detected by GroundingDINO on that frame as the input prompt. The resulting mask tracks from the two directions are then merged based on overlapping detections in the initial frame. Each generalized noun may have one or multiple masklets produced, depending on the specific video scenario.

**Step 4 : Referring Expression Assignment.** RexSeek is used to assign each referring expression to its corresponding masklet. For nouns with multiple candidate masklets (e.g., “person” referring to several individuals in the video), the full referring expressions are leveraged to disambiguate and establish the correct associations. This assignment is performed on the same frame used to initialize segmentation and tracking. Specifically, we first group masklets by their associated generalized nouns. Then, for each group requiring association resolution, given the frame and the bounding boxes of that noun in the frame, we prompt RexSeek with the prompt template: “Please detect {referring} in this image. Answer the question with object indexes.” The full referring expression replaces the word {referring} in the prompt template.

#### 2.1.4. Video Clipper

To capture temporally dense dynamics of each entity (e.g., movements, actions, and their coarse temporal boundaries), we segment the video into short clips using visual scene changes and semantic frame-level shifts. We first apply PySceneDetect’s ContentDetector [4] with a threshold of 20, empirically chosen based on qualitative assessment.

Many videos remained unclipped by PySceneDetect, despite clear action or event transitions, due to its reliance on HSV-based frame differences that miss semantic changes. To overcome this, we extract SigLIP [69] frame embeddings at 3 FPS and apply a custom clipping algorithm using Hierarchical Agglomerative Clustering to segment semantically distinct video segments.

#### 2.1.5. Video Transcriber

The *Video Transcriber* is a Video LLM that generates behavior-centric descriptions for entities, clip by clip. For each entity in the referring expression list, we iterate through the video clips and apply a two-step prompting process: (1) *Presence Check*: We first ask the model, “Is there a *entity*? Answer ‘Yes’ or ‘No’.” (2) *Behavior Description*: If the model responds with “Yes”, we follow up with: “What clearly happens to the *entity*? Describe only what is visibly happening to the *entity*, without inference or

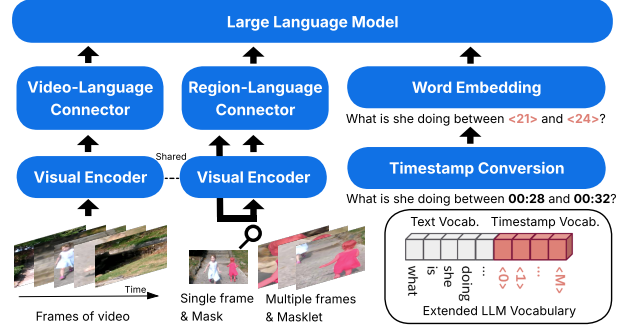


Figure 5. **Model Architecture:** Plug-and-play modules (Region-Language Connector, Timestamp Conversion) enhance general-purpose Video LLMs with space-time referring capabilities. It is worth noting that incorporating these modules for space-time referring is not strictly necessary (see Sec. 2.3).

assumptions.” The term “entity” is substituted with the actual referring expression (e.g., “girl in a white dress”).

This method yields a structured, temporally-dense video transcript focused on explicit, observable dynamics of individual entities.

#### 2.1.6. Video Instruction Data Generator

Finally, using the video metadata we have collected—such as structured video transcripts and language-described masklets—we generate instruction-style question-answer (QA) pairs either through templates or LLMs.

For example, given temporally-dense descriptions of entities across clips and knowledge of when each entity first appears, we can design templates to generate questions that test the understanding of temporal ordering—e.g., identifying the sequence in which entities appear. We can also ask about the actions of a specific entity within a given time interval, even if the entity does not appear during that interval but is present in other segments of the video.

Template-based question-answer generation typically involves designing a separate template for each task type (e.g., entity behavior captioning, temporal ordering). Crafting templates that effectively capture long-range dependencies—and implementing code to automatically generate QA pairs that conform to these templates—is difficult to scale. Alternatively, we can use the video transcript along with a set of in-context examples, and leverage an LLM to generate QA pairs that follow our predefined task definitions.

Moreover, to encourage reliance on masklets for resolving references, an LLM uses the entity referring expression list (which includes entities grounded with masklets) to replace specific language-based entity references with pronouns or generic terms (e.g., “she”, “he”, “the person”) in the generated questions.

We present the question task types and their definitions, used in both the template-based and LLM-based ap-

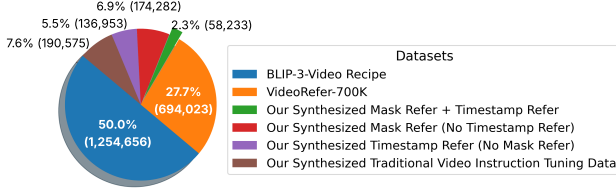


Figure 6. Data composition of our **final recipe** used in our experiments in Sec. 4.

proaches, in Table 7 of the Appendix. Each task is carefully designed to focus on space-time perception and reasoning.

## 2.2. Synthesized Data

We implement **Strefer** leveraging Tarsier-34b [52] and Qwen2.5-32B-Instruct [48] as the Video LLM and LLM resp. For the Masklet Generator, we use RexSeek-3B [20], SAM2 (sam2.1\_hiera\_large) [42] and GroundingDINO (grounding-dino-tiny) [32]. By the time of writing, we have generated video metadata and 947,854 instruction-response pairs using 4,253 NExT-QA [56] videos (excluding test set; avg. 40 seconds long and up to 3 minutes long). We illustrate the data composition of our final recipe in Fig. 6. Our synthesized instruction data is organized into 8 distinct groups based on whether masks or timestamps appear in the question and the QA-generation source (template or LLM), covering 11 different question task types. For details, please refer to Table 1 and Table 7 in the Appendix.

## 2.3. Modeling for Space-Time Referring

**Modeling with Architectural Enhancements.** We add plug-in-and-play modules to established general-purpose Video LLMs to unleash their fine-grained mask-level comprehension at any specific regions and any timestamps for a given video. To support detailed region-level understanding, we incorporate the spatiotemporal object encoder design from VideoRefer [68], which enables the model to understand fine-grained mask and masklet. For precise timestamp-level comprehension, we introduce learning special temporal tokens inspired by GroundedLLM [50], allowing the model to interpret specific moments in time properly. For details, please refer to Appendix. The resulting model (Fig. 5) is a next-token-prediction Video LLM with fine-grained, mask-level comprehension across arbitrary spatial regions and temporal segments of a video.

**Modeling without Architectural Changes.** It is worth noting that incorporating these modules for space-time referring is not strictly necessary. We also explore visual prompting approaches—SoM [63] for masklet comprehension and NumberIt [55] for timestamp understanding.

## 3. Evaluation Details

### 3.1. Training Details

To evaluate the quality of our synthesized instruction data, we integrate it into a base video instruction tuning recipe, which combines the video instruction-tuning data used by BLIP-3-Video [45] with VideoRefer-700K [68]. Our baseline is the model tuned on this base recipe. The video instruction-tuning data used by BLIP-3-Video comprises data from multiple sources, including Mira [21], VideoInstruct-100K [35], MSVD-QA [57], MSRVT-QA [57], ActivityNet-QA [66], TGIF-QA [19], and NExT-QA [56]. VideoRefer-700K is a recently released instruction-tuning dataset for video mask and masklet referring tasks, but it lacks timestamp-referring instructions.

The full model is tuned except the visual encoder. The visual encoder is not tuned due to insufficient data, which prevents effective tuning. Since the encoder is designed to extract complex visual patterns and features from raw RGB signals, it requires a large amount of data to generalize well. To be specific, we start from an image-comprehension vision LLM, the pre-trained BLIP-3 [61] model, with additional untrained architectural enhancements from BLIP-3-Video [45], as well as those described in Sec. 2.3. We adapt BLIP-3 for video and masklet comprehension by fine-tuning the full model illustrated in Fig. 5, except for the visual encoder, using 32 frames per video and 32 temporal tokens. Other hyperparameters, such as learning rate and batch size, were selected based on downstream evaluation results for the ‘Baseline’ model presented in the result tables. However, when tuning the model using recipes integrated with our data, we did not change any hyperparameters from those used in the ‘Baseline’ model. The training takes roughly 1 day and requires  $3 \times 8$  H200 GPUs. The resulting model has 4B parameters.

In the baseline and ablation models, if its training data lacks mask-referring instructions, the corresponding modules are excluded; likewise, timestamp-related modules are omitted if timestamp-referring instructions are not present in training. Therefore, the ‘Baseline Ablation’ model presented in the result tables shares the same architecture as BLIP-3-Video [45]; the model does not include the plug-and-play modules described in Sec. 2.3, as its training data lacks instructions that refer to masks or timestamps. The ‘Baseline’ model does not have ‘Timestamp Conversion’ or an extended LLM vocabulary for learning special temporal tokens (see Fig. 5) due to the lack of instruction data involving specific timestamps.

Our synthetic data includes full-length masklets per referring entity, but for efficient training, we sample a single mask on a random frame per instruction-response pair. At evaluation, we use the full masklet. Training with full masklets is expected to further improve performance.



### 3.2. Benchmark Details

We describe the evaluation benchmarks for Mask-Referred Regional Description, Mask-Referred Regional QA, and Timestamp-Referred Video QA below, as these represent less common evaluation settings for Video LLMs.

**VideoRefer-Bench<sup>D</sup>** [68] assesses the model’s ability to describe an entity across a video, given a mask or masklet of that entity. The benchmark comprises 400 videos from the test set of Panda-70M [6].

To evaluate performance on this benchmark, we use the following instruction template: “Please give a detailed description of the highlighted object `<region>` in the video.” The word `<region>` is substituted with model-extracted regional tokens if the model has built-in mechanisms to extract regional features.

The model evaluation is performed by GPT-4o by assigning scores to the generated predictions on a scale range from 0 to 5 across the following four dimensions [68]:

- **Subject Correspondence:** This dimension evaluates whether the subject of the generated description accurately corresponds to that specified in the ground truth.
- **Temporal Description:** This aspect analyzes whether the representation of the object’s motion is consistent with the actual movements.
- **Appearance Description:** This criterion assesses the accuracy of appearance-related details, including color, shape, texture, and other relevant visual attributes.
- **Hallucination Detection:** This facet identifies discrepancies by determining if the generated description includes any facts, actions, or elements absent from reality, like imaginative interpretations or incorrect inferences.

**VideoRefer-Bench<sup>Q</sup>** [68] evaluates a model’s ability to answer video entity-related questions, given one or more entities’ masks or masklets within a video. The benchmark includes 1,000 multiple-choice questions spanning 198 videos sourced from various datasets, including the test set of MeViS [12], A2D-Sentences [14], and Refer-YouTube-VOS [53]. Questions are crafted to assess different dimensions of understanding, including Basic Questions, Sequential Questions, Relationship Questions, Future Predictions and Complex/Reasoning Questions.

Sequential Questions typically ask about entity action and ordering; Basic Questions typically concern attributes like object color. Relationship Questions involve more than one object regions in the question. Future Predictions involve weakly grounded reasoning about forthcoming events. Notably, models generally perform best on Complex/Reasoning Questions, making this category the easiest despite its name.

We use the following instruction template: “Please answer the following question about the `<region>`. {question}”.

**Timestamp-based Yes/No QA on QVHighlights** is a task that repurposes existing annotations from the video highlight detection dataset, QVHighlights [23]. Specifically, for each annotated segment—defined by a start and end timestamp and an associated language description—we construct a question prompt in the following form:

‘Does the following description accurately reflect what happens in the video between `<start_time>` and `<end_time>`? Description: {description}. Respond with ‘Yes’ or ‘No’ only.’ Each of these prompts is assigned the ground truth answer “Yes”.

To generate negative (i.e., “No”) samples, we randomly select segments from the same video that do not overlap with any annotated intervals. To ensure this, we first expand each annotated timestamp by a buffer of 5 seconds on both sides, then merge overlapping intervals to form a set of excluded ranges. We then identify all remaining gaps in the video timeline that lie outside these excluded regions. From these valid gaps, we randomly select a new segment that satisfies a minimum duration of 10 seconds. A description from an annotated segment is then paired with this unrelated time window to form a mismatched QA example with the correct answer “No”.

We ensure a balanced answer distribution, with almost 50% of the samples labeled as “Yes” and 50% as “No”.

For each question, we substituted `<start_time>` and `<end_time>` with their corresponding timestamps. For models that do not learn temporal tokens, timestamps are represented by default in the `HH:MM:SS.xxx` format. For models that do learn temporal tokens, we use temporal tokens to substitute `<start_time>` and `<end_time>`. For example, if a model learns 32 temporal tokens and the video’s duration is 90 seconds, a timestamp like `00:00:19.228` is converted to `<7>` ( $\frac{19.228}{90} \times 32 \approx 7$ ).

## 4. Experiments

### 4.1. Quantitative Results

#### 4.1.1. Main Results and Findings

We explore the final data recipe and ablate groups of our synthesized instruction data, and results are listed in Tables 2, 3, and 4. We created model variants by incrementally adding our data groups to the base recipe. These groups were categorized based on whether the instructions involved mask- or timestamp-referring cues, and whether the QA pairs were generated by an LLM or derived from templates. We summarize our findings as follows:

**Finding 1:** The video space-time referring task, which requires understanding both the full video and the user-specified regions at specific times, enhances spatiotemporal understanding in Video LLMs.

Integrating the base recipe with our LLM-synthesized  $\mathcal{G}7$  data—featuring queries with mask and timestamp references—boosts performance across several tasks: mask-referred video regional description improves from 3.28 to 3.34, mask-referred video QA from 0.665 to 0.672, timestamp-based QA from 0.5288 to 0.5390, and TempCompass [33] from 60.100 to 60.650. Impressively,  $\mathcal{G}7$  contains just 27K samples (1.39% of the base recipe).

Incorporating additional data containing queries with only mask references ( $\mathcal{G}6$ ) leads to consistent performance gains across all tasks—except for QVHighlights [23]. On this benchmark, using both  $\mathcal{G}6$  and  $\mathcal{G}7$  results in a slightly lower score (0.5337) compared to using only  $\mathcal{G}7$  (0.5390). We hypothesize that this drop is due to the relatively smaller proportion of timestamp-referred queries, as the inclusion of 174K mask-referred instances may dilute the supervision of timestamp-based video understanding. Despite this, the combined data still outperforms the baseline, which lacked any timestamp-referred training examples. Notably, further augmenting training with  $\mathcal{G}6$  yields the highest score (4.4525) on the Subject Correspondence task, which evaluates the model’s ability to align the described subject with that specified in the ground truth, highlighting the benefits of our mask-referring data in enhancing model understanding of specific regions within the video content.

To further enhance the ability of model to interpret timestamp-based video content, we introduce  $\mathcal{G}8$ , which yields a QVHighlights score of 0.5672.  $\mathcal{G}8$  combines timestamp and mask references and is derived from a subset of  $\mathcal{G}5$ , a template-generated timestamp referring data. We sampled remaining  $\mathcal{G}5$  that were not converted into mask referring data and incorporate it into the training recipe, performance on QVHighlights improved into 0.5900.

**Finding 2: Traditional video instruction data—when designed with a focus on dynamics and enriched with both positive and negative questions—can further improve performance, even on space-time referring tasks.**

Our final recipe incorporates an additional component,  $\mathcal{G}1$ , which is generated from templates and specifically designed to focus on entity behavior. It includes negative questions—such as asking the model to describe the behavior of entities that are absent from a given video segment—to help the model avoid relying on shortcuts or making incorrect assumptions. Surprisingly, this addition not only enhances performance on video temporal perception and reasoning tasks (TempCompass and VideoMME [13]), but also improves results on space-time referring understanding tasks: Compared to the baseline, the average score on mask-referred video regional description increases to 3.39; mask-referred video QA accuracy improves to 0.688; timestamp-based QA rises to 0.6031; performance on TempCompass reaches 61.675; and performance on VideoMME improves to 37.70. The final recipe consistently and signifi-

cantly improves performance across benchmarks.

**Finding 3: Template-generated event sequencing data may enhance higher-level, long-term temporal reasoning skills, but this can come at the cost of precise, fine-grained spatial-temporal understanding. Therefore, a more balanced data mixture may be necessary to support both broad temporal abstractions and localized, detailed spatial-temporal comprehension.**

$\mathcal{G}2$  comprises multiple-choice traditional video instruction data that tasks the model with identifying the correct temporal order from the wrong ones in which entities first appear in the video. We hypothesize that this task format fosters long-range temporal reasoning abilities. Empirically, *before incorporating  $\mathcal{G}1$* , we find that incorporating  $\mathcal{G}2$  improves performance on temporally-focused benchmarks such as TempCompass with short videos (from 60.750 to 60.925), as well as VideoMME—a long-video understanding benchmark containing hour-long videos (from 35.90 to 41.65), supporting this hypothesis.

However, we observe a performance drop when  $\mathcal{G}2$  is added in tasks like mask-referred captioning, mask-referred QA, and timestamp-referred QA. We attribute this decline to a trade-off in the type of reasoning the model develops: while  $\mathcal{G}2$  promotes higher-level long-term temporal reasoning skills, the mask or timestamp-referred understanding tasks demand fine-grained, localized understanding of video content. Thus, the broader temporal abstractions encouraged by  $\mathcal{G}2$  may come at the expense of precision in spatial-temporal fine-grained, localized understanding.

After incorporating  $\mathcal{G}1$ , further adding  $\mathcal{G}2$  leads to performance drop across all benchmarks, even though  $\mathcal{G}2$  contains only 1K samples. This suggests that  $\mathcal{G}1$ , which incorporates negative questions and dynamic understanding, may already cultivate a certain degree of temporal reasoning. Consequently, the benefits of  $\mathcal{G}2$  may be diminished or overshadowed in the presence of  $\mathcal{G}1$ .

#### 4.1.2. Discussions and Insights

Using **Strefer**, we generated 947,854 instruction data points from just 4,253 NExT-QA videos. While our **final recipe** added only 545 more videos compared to the baseline, it led to performance gains across multiple benchmarks with both short and long videos.

This highlights a key insight: the value of video instruction-tuning data is not merely in its quantity, but in its quality and specificity. **Rather than indiscriminately scaling the dataset with more videos and generic questions, we found that carefully crafting meaningful, well-grounded questions and answers leads to significantly better training outcomes for video-language models.**

We identify the following three critical principles for curating effective video instruction data:

- **Video-grounded instruction-response pairs:** Both

Group	Description	Visual Input	Task Types (ref. Tab. 7)	# Samples
$\mathcal{G}_1$	Traditional data, template generated	Clip	1,2,3,4	190,575
$\mathcal{G}_2$	Traditional data, template generated	Video	5	1,316
$\mathcal{G}_3$	Traditional data, LLM generated	Video	7	288,630
$\mathcal{G}_4$	<b>Timestamp referring</b> data, LLM generated	Video	8,9,10,11	44,243
$\mathcal{G}_5$	<b>Timestamp referring</b> data, template generated	Video	6	190,575
$\mathcal{G}_6$	<b>Mask referring</b> instead of language referring, derived from $\mathcal{G}_3$	Video	7	174,282
$\mathcal{G}_7$	<b>Mask referring</b> instead of language referring, derived from $\mathcal{G}_4$	Video	8,9,10,11	27,092
$\mathcal{G}_8$	<b>Mask referring</b> instead of language referring, derived $\mathcal{G}_5$	Video	6	31,141

Table 1. **Overview of our synthesized instruction-tuning data.** Using **Strefer**, we generated 947,854 instruction data from only 4,253 NExT-QA videos. Our synthesized data includes a diverse range of temporally focused questions (i.e., tasks)—such as temporal relations, event sequencing, time-referencing, and coarse time localization. It incorporates both positive and negative questions to help the model learn to avoid relying on shortcuts or making incorrect assumptions. Although our **final recipe** introduced only 545 additional videos beyond the baseline, it yielded noticeably improved performance across multiple benchmarks (cf. Tab. 2, Tab. 3, and Tab. 4).

Mask-Referred Regional Description (VideoRefer-Bench <sup>D</sup> [68])	Samples Added (%)	Avg.	Subject Correspondence	Temporal Description	Appearance Description	Hallucination Detection
GPT-4o	N/A	3.25	4.15	3.11	3.31	2.43
GPT-4o-mini	N/A	3.05	3.89	2.62	3.18	2.50
Baseline Ablation: Video Instruction-Tuning Data [45]	N/A	2.7308	3.5200	2.4235	2.5639	2.4160
Baseline: Base Recipe (1,948,679 samples)	N/A	3.2837	4.3775	2.9523	3.1075	2.6975
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8 + \text{Remaining } \mathcal{G}_5 \text{ (Sampled)} + \mathcal{G}_1$	28.73%	<b>3.3947</b>	<u>4.4400</u>	3.0575	<b>3.2763</b>	2.8050
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8 + \text{Remaining } \mathcal{G}_5 \text{ (Sampled)} + \mathcal{G}_2$	19.02%	3.3537	4.3650	3.0150	<u>3.2675</u>	2.7675
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8 + \text{Remaining } \mathcal{G}_5 \text{ (Sampled)} + \mathcal{G}_1 + \mathcal{G}_2$	28.80%	3.3742	4.3975	<u>3.0676</u>	3.2493	2.7825
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8 + \text{Remaining } \mathcal{G}_5 \text{ (Sampled)}$	18.95%	3.3740	4.4025	3.0300	3.2512	<u>2.8125</u>
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8$	11.93%	<u>3.3821</u>	4.3625	<b>3.1150</b>	3.2317	<b>2.8195</b>
+ $\mathcal{G}_6 + \mathcal{G}_7$	10.33%	3.3710	<b>4.4525</b>	3.0579	3.2412	2.7325
+ $\mathcal{G}_7$	1.39%	3.3421	4.3825	3.0150	3.2311	2.7400

Table 2. **Regional Description** results on VideoRefer-Bench<sup>D</sup> (**Best/Second**). This evaluation assesses the model’s ability to describe an entity across a video, given a mask or masklet of that entity. If a video contains multiple entities, the model must leverage the region input to distinguish and describe the correct one. Across result tables: The **blue-highlighted** row represents our final recipe that performs well across all benchmarks. GPT-4o and GPT-4o-mini results are listed but excluded when highlighting the top two performers.

Mask-Referred Regional QA (VideoRefer-Bench <sup>Q</sup> [68])	Samples Added (%)	Avg.	Sequential Questions	Relationship Questions	Basic Questions	Future Predictions	Reasoning Questions
GPT-4o	N/A	0.713	0.745	0.660	0.623	0.737	0.880
GPT-4o-mini	N/A	0.658	0.671	0.565	0.576	0.754	0.859
Baseline Ablation: Video Instruction-Tuning Data [45]	N/A	0.621	0.6015	0.4920	0.6042	0.7192	0.8321
Baseline: Base Recipe (1,948,679 samples)	N/A	0.665	0.6367	0.5476	0.6382	0.7807	0.8741
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8 + \text{Remaining } \mathcal{G}_5 \text{ (Sampled)} + \mathcal{G}_1$	28.73%	<b>0.688</b>	<b>0.6640</b>	<u>0.5753</u>	0.6680	<b>0.8070</b>	0.8671
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8 + \text{Remaining } \mathcal{G}_5 \text{ (Sampled)} + \mathcal{G}_2$	19.02%	0.677	0.6250	0.5634	0.6808	0.7894	<u>0.8741</u>
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8 + \text{Remaining } \mathcal{G}_5 \text{ (Sampled)} + \mathcal{G}_1 + \mathcal{G}_2$	28.80%	0.683	0.6523	0.5555	<b>0.6936</b>	<u>0.7982</u>	0.8531
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8 + \text{Remaining } \mathcal{G}_5 \text{ (Sampled)}$	18.95%	0.681	<b>0.6640</b>	0.5674	0.6510	0.7894	<u>0.8741</u>
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8$	11.93%	0.678	0.6289	0.5714	<u>0.6851</u>	0.7807	0.8601
+ $\mathcal{G}_6 + \mathcal{G}_7$	10.33%	<u>0.685</u>	<u>0.6601</u>	<b>0.5952</b>	0.6553	0.7543	<b>0.8811</b>
+ $\mathcal{G}_7$	1.39%	0.672	0.6523	0.5674	0.6425	0.7631	0.8671

Table 3. **Regional QA** results on VideoRefer-Bench<sup>Q</sup> (**Best/Second**). This evaluation assesses the model’s ability to answer questions about one or more entities in a video, given a mask or masklet of the entity.

Timestamp-Referred QA (Yes/No) & Traditional QA	Samples Added (%)	QVHighlights [23]	TempCompass [33]					VideoMME [13]
			Avg.	Yes	No	MCQ	Caption Matching	
Baseline Ablation: Video Instruction-Tuning Data [45]	N/A	0.5297	58.525	57.3	59.3	67.5	50.0	35.85
Baseline: Base Recipe (1,948,679 samples)	N/A	0.5288	60.100	61.9	58.4	70.5	49.6	37.45
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8 + \text{Remaining } \mathcal{G}_5 \text{ (Sampled)} + \mathcal{G}_1$	28.73%	<b>0.6031</b>	<b>61.675</b>	<u>65.1</u>	59.1	70.9	<b>51.6</b>	37.70
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8 + \text{Remaining } \mathcal{G}_5 \text{ (Sampled)} + \mathcal{G}_2$	19.02%	0.5774	60.925	63.8	58.9	<u>71.1</u>	49.9	<b>41.65</b>
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8 + \text{Remaining } \mathcal{G}_5 \text{ (Sampled)} + \mathcal{G}_1 + \mathcal{G}_2$	28.80%	<u>0.5951</u>	61.250	64.5	<b>59.8</b>	70.5	50.2	34.70
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8 + \text{Remaining } \mathcal{G}_5 \text{ (Sampled)}$	18.95%	0.5900	60.750	63.4	58.8	70.3	50.5	35.90
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8$	11.93%	0.5672	<b>61.675</b>	<b>66.0</b>	59.3	<b>71.4</b>	50.0	38.15
+ $\mathcal{G}_6 + \mathcal{G}_7$	10.33%	0.5337	<u>61.550</u>	<u>65.1</u>	<u>59.4</u>	70.7	<u>51.0</u>	<u>39.15</u>
+ $\mathcal{G}_7$	1.39%	0.5390	60.650	63.0	58.7	70.1	50.8	33.90

Table 4. **Results of timestamp-based QA on QVHighlights and traditional QA focusing on temporal cues using TempCompass and VideoMME (Best/Second).** We transform annotations from QVHighlights into questions tied to specific timestamps, each expecting a ‘Yes’ or ‘No’ answer. For VideoMME, we report results on the Temporal Perception and Temporal Reasoning subsets (no subtitles). While VideoMME benchmarks long video understanding (up to an hour), our data of entirely short videos can help improve model performance.



<b>Mask-Referred Regional Description</b> (VideoRefer-Bench <sup>D</sup> [68])	Samples Added (%)	Avg.	Subject Correspondence	Temporal Description	Appearance Description	Hallucination Detection
Baseline Ablation Model	N/A	2.7308	3.5200	2.4235	2.5639	2.4160
Baseline Ablation Model + SoM [63]	N/A	2.6593	3.5600	2.1834	2.3576	2.5363

Table 5. **Regional Description** results on VideoRefer-Bench<sup>D</sup> before and after applying the visual prompting method, SoM [63].

Timestamp-Referred QA (Yes/No)	Samples Added (%)	QVHighlights [23]
Baseline Ablation: Video Instruction-Tuning Data [45]	N/A	0.5297
Baseline Ablation: Video Instruction-Tuning Data [45] + NumberIt [55]	N/A	0.5301
Baseline: Base Recipe (1,948,679 samples)	N/A	0.5288
Baseline: Base Recipe (1,948,679 samples) + NumberIt [55]	N/A	0.5311
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8$ + Remaining $\mathcal{G}_5$ (Sampled) + $\mathcal{G}_1$	28.73%	<b>0.6031</b>
+ $\mathcal{G}_6 + \mathcal{G}_7 + \mathcal{G}_8$ + Remaining $\mathcal{G}_5$ (Sampled) + $\mathcal{G}_1$ + NumberIt [55]	28.73%	<b>0.6041</b>

Table 6. **Timestamp-based Yes/No QA** results on QVHighlights before and after applying the visual prompting method, NumberIt [55].

questions and responses must be tightly linked to the video content. **Strefer** achieves this through spatiotemporally grounded metadata generated using multiple pixel-level vision foundation models.

- **Temporal and fine-grained reasoning emphasis:** Effective data should challenge models to focus on dynamics and reason about time, especially those involving fine-grained details and long-range dependencies in both space and time.
- **Diverse tasks and formats:** A mix of task types and instruction data formats ensures broader coverage and more robust video understanding and reasoning capabilities.

#### 4.1.3. Visual Prompting for Space-Time Referring

It is worth noting that incorporating the plug-and-play modules and modify the architecture of the pre-trained general-purpose Video LLM for space-time referring is not strictly necessary. We explore visual prompting approaches: SoM [63] for masklet comprehension and NumberIt [55] for timestamp understanding.

**SoM: Mask-Overlay-Frame Prompting.** We follow the implementation of the Set-of-Mark (SoM) method from VideoRefer [11] to apply masks to video frames, as originally proposed by [63]. We also changed the question prompt into: “I have outlined an object with a red contour in the video. Please describe the object in detail.” The results are presented in Table 5. After applying SoM, the average performance on the Mask-Referred Video Regional Description task decreases, but performance increases on certain metrics, e.g., Subject Correspondence.

Our analysis reveals that the effectiveness of SoM is highly sensitive to the way masks are rendered on the video frames. In our initial implementation, we used thicker mask boundaries and semi-transparent red fill color. This approach led to severe hallucinations by the model, which often misinterpreted masked regions as merely red-colored

objects. In contrast, the SoM implementation from VideoRefer [11] uses thinner boundaries and fully transparent fills, resulting in significantly improved performance over our version. Nevertheless, the performance remains lower than the baseline without any SoM prompting.

#### **NumberIt: FrameID-Overlay-Frame Prompting.**

Similar to SoM, NumberIt [55] overlays the frame ID at a specific location on each frame. We overlaid the frame ID in red, following the authors’ suggestion, and placed each ID in the top-left corner of the corresponding frame (the resulting rendering effect is similar to Fig. 20). We also modified each question as follows: The red numbers on each frame represent the frame number. Does the following description accurately reflect what happens in the video between <frame.start> and <frame.end>? Description: {description}. Respond with ‘Yes’ or ‘No’ only. For each question, we substituted <frame.start> and <frame.end> with their corresponding frame IDs. The results are listed in Table 6. The performance on timestamp-based Yes/No Video QA in QVHighlights shows a slight improvement after applying the visual prompting method, NumberIt.

**Summary:** Using a pretrained general-purpose Video LLM for space-time referring tasks does not necessarily require altering the model architecture. We found that visual prompting approaches, such as SoM and NumberIt, can help the model perform mask-referring or timestamp-referring tasks; however, their effectiveness appears limited in the absence of model tuning. We hypothesize that incorporating these techniques during model fine-tuning—while preserving the original architecture—may lead to more performance gains [55], which we leave for future work.

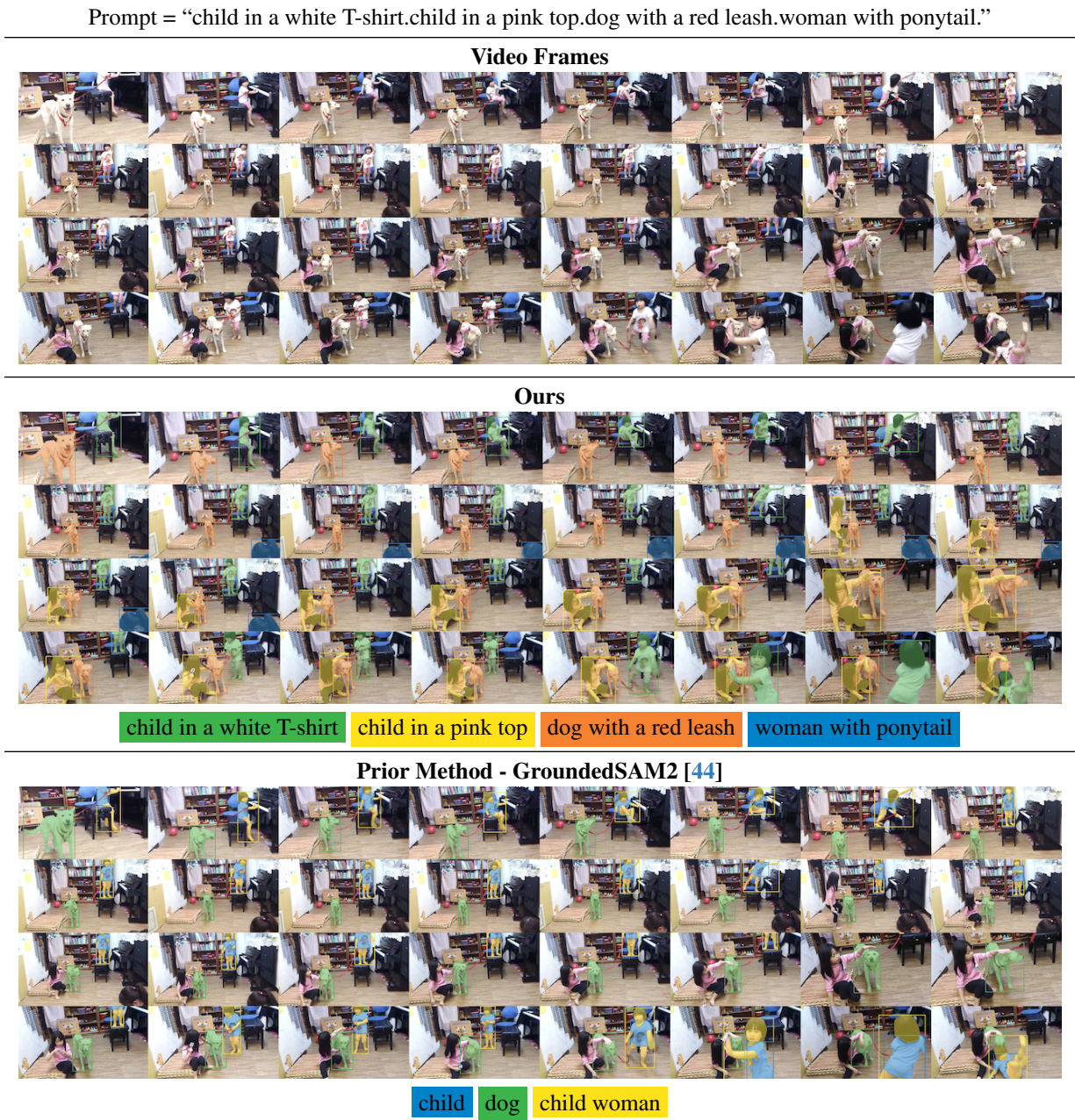


Figure 7. **Qualitative Results of Referring Masklet Generation.** In this video, our method accurately generates masklets corresponding to the input referring expressions. In contrast, GroundedSAM2 [44] fails to assign expressions to masklets and does not detect the woman and the child in a pink top, who appear midway through the video.

## 4.2. Qualitative Results

In this section, we present qualitative results that best illustrate the strengths and limitations of our method and model. We discuss limitations in detail in the next section.

### 4.2.1. Strefer-Synthesized Instruction Data

We present qualitative results of our **Strefer**-synthesized data in Fig. 2 and Fig. 11 in the Appendix. Despite some noise in the pseudo-annotated video metadata, the synthesized instruction-response pairs capture essential spatio-temporal dynamics from the videos and successfully enhance the model’s performance.



Prompt = “woman on a bicycle.man in a blue shirt.man in a white shirt.man in a black shirt.”

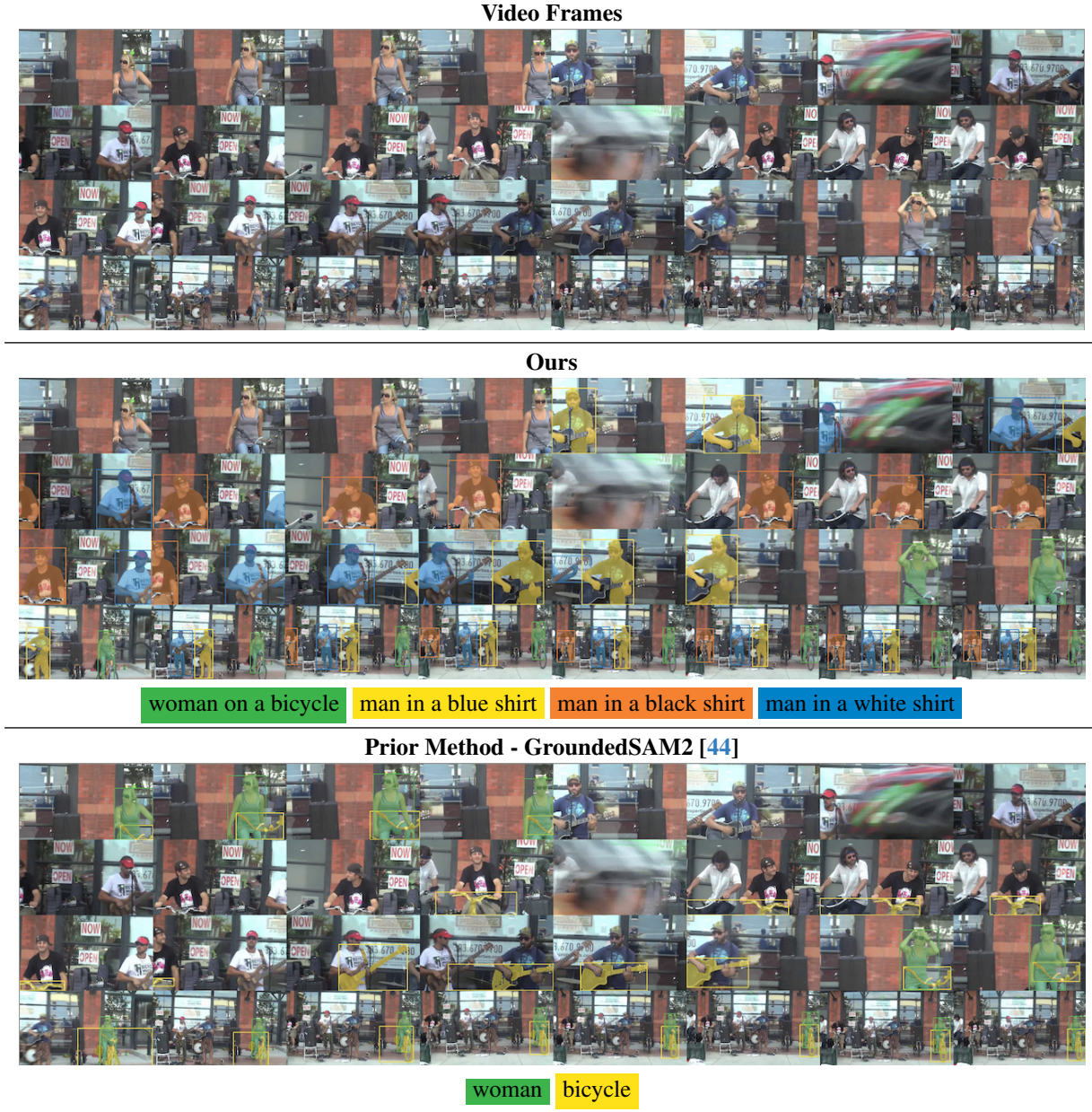


Figure 8. **Failure Results of Referring Masklet Generation.** Our method fails to consistently track the woman on a bicycle throughout the video, while GroundedSAM2 [44] fails to detect, track, and differentiate the individuals referenced in the input text prompt. Videos with heavy motion blur and long-range dependencies remain challenging to handle.

#### 4.2.2. Our Referring Masklet Generation

We present qualitative results of our novel referring masklet generation pipeline in Fig. 7, Fig. 8, Fig. 12 and Fig. 13 (some figures are in the Appendix). Our referring masklet generation pipeline demonstrates superior performance compared to the widely adopted prior method, GroundedSAM2 [44].

#### 4.2.3. Strefer-Trained Model

We also present comprehensive qualitative results of our **Strefer**-trained model in comparison to the ‘Baseline’ model. We observe that the models tend to exhibit a foreground bias and misinterprets the masklet. The model trained on our data mitigates this limitation (see Fig. 9, Fig. 14). In cases requiring entity disambiguation (e.g.,



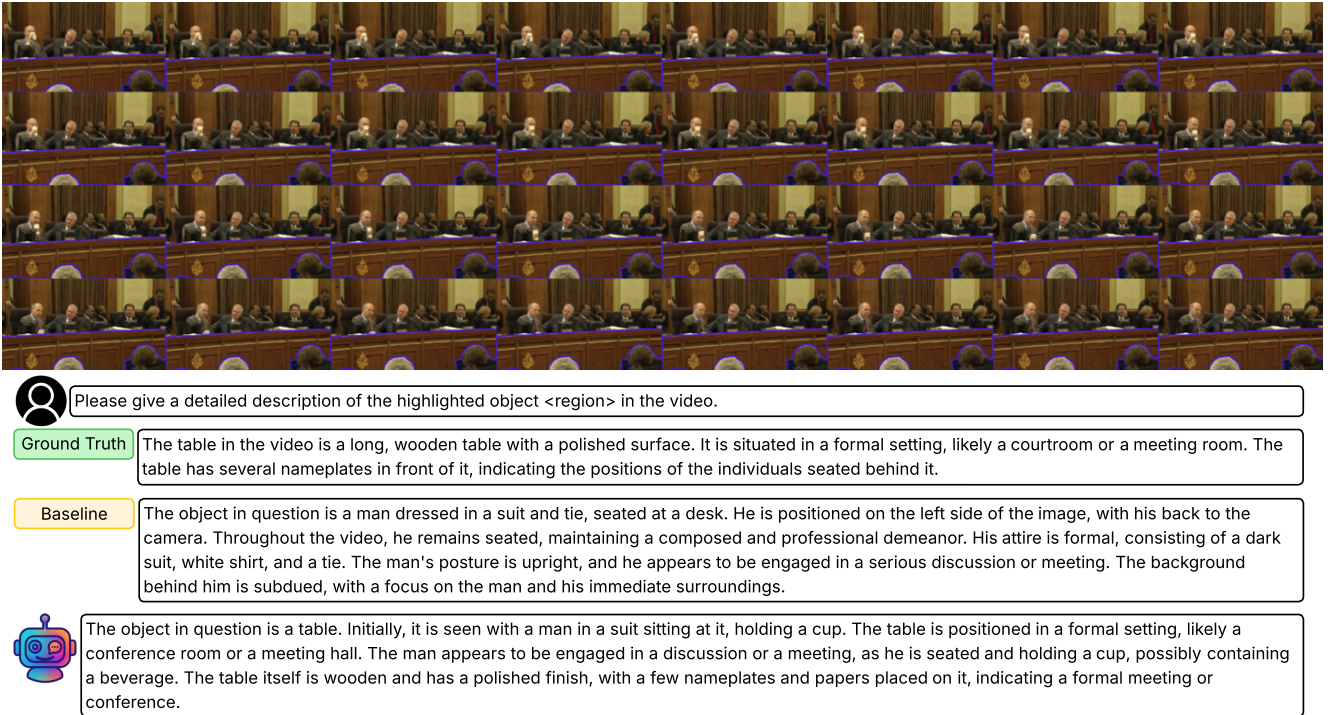


Figure 9. **Qualitative Results of Video LLMs Trained w/ and w/o Strefer-Synthesized Data.** This sample is drawn from VideoRefer-Bench<sup>D</sup>, designed to assess a model’s performance on the task of **Mask-Referred Regional Description**. The boundary of the region referred to by the mask in this sample is highlighted in purple. While the video includes several individuals as prominent foreground elements, the masklet specifically refers to the table, not the people. The baseline model, however, fails to interpret the mask correctly and mistakenly answers that the referred object is a man. In contrast, the model trained on **Strefer**-generated data accurately identifies the masklet-referred region as a table.

Fig. 15, Fig. 17, Fig. 18, and Fig. 19), our model produces more accurate results. These gains stem from two design choices: (1) our LLM-based question generation strategy ensures that questions explicitly reference masklet information, and (2) our carefully designed masklet generator addresses the limitations of GroundedSAM2 in handling such scenarios (in contrast, GroundedSAM2 is used by VideoRefer-700K [68] in its data generation process). Additionally, our model demonstrates superior fine-grained spatiotemporal action understanding (see Fig. 16) and exhibits superior performance in precise timestamp-based video comprehension (i.e., Fig. 10, Fig. 20, and Fig. 21).

## 5. Limitations and Future Directions

**Strefer** synthesized data is not error-free. For example, in the last blue-highlighted video segment shown in Fig. 11 of the Appendix, **Strefer** identifies the woman as not present. However, a human viewer would easily identify the woman in that segment while watching the video, despite the fact that she is largely occluded and the frames are mostly occupied by the child. This error also occurs because we did not employ a more complex, dense video captioning framework that leverages *inter-segment* information, such as a hierarchical [74] or differential video cap-

tioning [5] method. We actually tried these approaches, but they did not yield better results than clip-by-clip captioning using current open-source models. We also experimented with several alternative open-source mid-scale Video LLMs, including Qwen2.5-VL-7B-Instruct, LLaVA-NeXT-Video-34B, LLaVA-OneVision-7B, and Tarsier-7B. Ultimately, we selected Tarsier-34B, as it appeared to provide more accurate, action-centric descriptions.

Scenes characterized by high visual clutter and significant dynamic variations continue to pose substantial challenges. Fig. 8 illustrates a failure case of our referring masklet generation—the woman is not tracked or segmented in the first four frames. This highlights that videos with heavy motion blur and long-range dependencies remain challenging to handle, even for our method. The issue stems from the tracking limitations of SAM2, the tracking and segmentation model we employ which is sensitive to the selection of the start tracking frame. Fig. 2 presents another example of tracking and segmentation failure—the child in the black shirt lacks associated masks in frame 10 and frame 11, despite being clearly visible.

**Strefer** may inherit limitations from the underlying models used in its modular components, such as pixel-level foundation models, LLMs, and Video LLMs. For

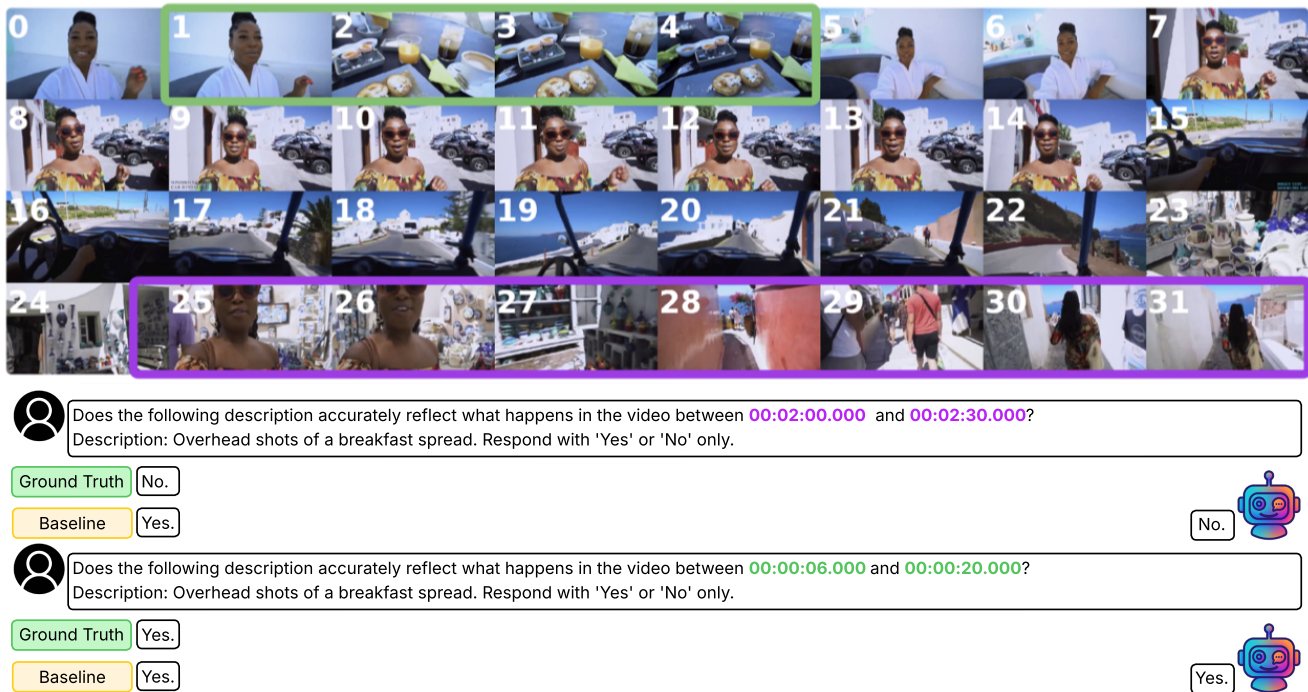


Figure 10. **Qualitative Results of Video LLMs Trained w/ and w/o Strefer-Synthethized Data.** This sample is drawn from QVHighlights, using our repurposed task designed to assess a model’s performance on **Timestamp-Referred Video QA**. The segment boundaries corresponding to the timestamps in the first and second questions are highlighted in purple and green, respectively. The model trained on our **Strefer**-generated data correctly answers both questions, demonstrating superior understanding of precise moments and segments in videos compared to the baseline.

instance, LLMs and Video LLMs are known to hallucinate, potentially introducing misleading information into annotated video metadata and synthesized question–answer pairs. Similarly, the extraction of pixel-level information may be less reliable in videos with highly similar individuals or densely populated scenes (e.g., crowded urban environments), which the pixel-level foundation models we used could not reliably handle. Despite these challenges, the modular nature of **Strefer** positions it well to benefit from future improvements in its underlying models—including LLMs, Video LLMs, and grounding vision foundation models—as they continue to evolve.

**Strefer** involves multiple models, which may hinder exact reproduction of its pseudo-annotation, data synthesis as well as our model training procedures for research groups with limited resources. However, **Strefer** is a multi-stage modular framework, wherein individual model components can be substituted with more computationally efficient alternatives, enabling flexible adaptation under varying resource constraints.

In terms of the limitations inherent to models trained on **Strefer**-synthesized data, Fig. 22 shows a typical failure case where the region indicated by the mask is neither the primary foreground object nor centrally positioned in the

frame. This reflects a common issue across all models, including baselines, which tend to exhibit both a foreground main object bias and a center bias. Furthermore, Fig. 23 illustrates another failure case, underscoring the persistent difficulty in capturing fine-grained action semantics.

Future work may further improve **Strefer** by refining individual modules—for example, improving the video clipping pipeline to produce entity-centric segments, and incorporating feedback-verification mechanisms to minimize hallucinated content in video metadata and instruction-following pairs. Additionally, given the potential for error propagation in our modular framework, as well as the nature of synthetic data, which may not fully match real human question distributions, future research is encouraged to develop effective filtering strategies for the synthetic instruction-tuning data as a final quality assurance or adjustment step in the data engine, thereby enabling more efficient, reliable, and robust model training.

For the development of the referring and reasoning video model, our current trained models are limited to mask-based spatial referring and is not trained for other types of spatial references such as points, boxes, and scribbles. However, since these other forms of spatial references are inherently sparser than masks and can be easily derived from

object segmentation masks, an avenue for future research is to explore transforming the current mask-level instruction-tuning data produced by **Strefer** into alternative data formats, and to train models that can comprehend more diverse forms of user spatial references.

Moreover, the LLM backbone of models we trained on our data is based on the pretrained microsoft/Phi-3-mini-4k-instruct [1]. As with many other Video LLMs, the performance of our model is heavily influenced by the capabilities of the underlying LLM. We encourage further research into training Video LLMs on **Strefer**-synthesized data using larger and more powerful models. However, this typically demands significantly more tuning data and greater computational resources.

In our experiments, we conducted rather limited exploration of the optimal training data mixture. As a result, the current composition may not represent the most effective setup for fostering broad, balanced and transferable skills. Future work could focus on systematically optimizing the data composition, which is likely to result in more substantial and consistent performance gains across diverse benchmarks and metrics.

Finally, our model is grounded at the perception level rather than at the output response generation level. While grounding at the output level offers a more direct path to interpretable video-language reasoning, it requires training data with high-fidelity spatiotemporal annotations. At present, the boundary of the fine-grained space-time information generated by **Strefer** may lack the precision required to reliably supervise such models.

Our work centers on *referring* understanding—where the model leverages fine-grained spatiotemporal cues as conditional input, and our models are trained using synthesized instruction-tuning data, rather than being directly supervised by pseudo-annotated dense video metadata. This setup is inherently more robust to moderate imperfections—such as missing entities or imprecise temporal boundaries. For example, as shown in Fig. 2, even when the temporal span of the masklet for the child in the black shirt is incomplete, the associated instruction-response pair remains accurate and meaningful. This resilience arises because referring understanding does not require an exhaustive coverage of the language-described pixel-level space-time information, making it more adaptable under the current data limitations.

Looking ahead, advancing output-level spatiotemporal grounding in Video LLMs holds significant promise for improving their generalization, reliability and fine-grained spatiotemporal reasoning skills. We encourage future work to pursue this direction by leveraging more accurate spatial-temporal annotations aligned with language, ideally enabled by an enhanced, scalable, and automated data generation pipeline.

## 6. Related Work

**Video Spatial Referring.** A growing body of research has enabled Multimodal Large Language Models (MLLMs) to perform regional understanding in static images [34, 41, 58, 62, 64, 67, 71–73], but spatial referring understanding in videos remains underexplored. Earlier efforts include Artemis [40] for single-object referencing using box-level representations, as well as Elysium [51] and Merlin [65] that transform box coordinates into textual prompts to help the LLM identify the referred objects. Nevertheless, these methods are plagued by imprecise regional understanding.

Recently, a new wave of concurrent work has shifted towards mask-based region-level understanding in videos. DAM [29] and PAM [31], while demonstrating extension to videos, primarily focus on regional captioning rather than general instruction-following. Both Omni-RGPT [16] and SAMA [47] introduced automatic annotation methods using GPT-4o and Gemini 1.5 Pro, resp., to transform existing video datasets with annotated regions into conversational data. VideoRefer [68] introduces a multi-agent annotation pipeline that operates without the need for pre-existing annotations. However, its design struggles with complex videos such as scenes involving multiple entities of the same category or cases where entities are temporarily absent. Additionally, it cannot produce temporal annotations, limiting its ability to support time-sensitive tasks.

**Video Timestamp Referring.** We posit that enabling models to comprehend timestamp-specific video segments may enhance their video temporal reasoning skills. Therefore, we propose **Strefer**, a framework that generates timestamp-based temporal annotations and instructions for videos to empower model to understand video timestamp-based content referencing. To the best of our knowledge, there is no existing work that focuses solely on this task. However, there has been a line of research on temporally grounded Video LLMs [15, 18, 39, 50, 55, 75], which are capable of both understanding and localizing temporal boundaries in response to user queries. These models require training on boundary-accurate temporal localization data, and we consider temporal or spatiotemporal localization a future work.

**Video Instruction-Following Data.** Early efforts such as VideoChat [25], Video-ChatGPT [35] and MVBench [26] convert existing video annotations into a conversational format. Later on, more studies employed proprietary models to produce training data such as ShareGPT4Video [5] and MiraData [21] for the video captioning task, but obtaining versatile and effective video instruction-following data has remained challenging, with earlier datasets criticized for their limited utility [74]. LLaVA-Video-178K [74] then takes a step forward by sourcing existing video datasets (with videos up to 3 minutes) and enriching them with open-ended and multiple-choice questions across diverse tasks,



developed through a combination of GPT-4o and human efforts. LongViTU [54] and VideoMarathon [30] were also introduced for hour-scale long video training.

Inspired by these efforts—as well as the TempCompass benchmark [33], which emphasizes the diversity of temporal aspects and task formats, we design instruction-following tasks and formats that capture a wide range of temporal variations using pseudo-annotated video metadata. We are also influenced by T3 [27], which shows that with proper temporally focused task design, textual temporal instruction data can enhance temporal reasoning without video-based training. However, we argue that true spatiotemporal reasoning of an AI agent requires the video input. Our synthetic instruction data target space-time referring understanding that encourages perceptually-grounded, instruction-tuned Video LLMs.

**Positioning of Our Work.** None of the above general video instruction data is designed to empower Video LLMs for the fine-grained, space-time referring tasks. Currently, the only available instruction tuning data for masklet referring comprehension in videos is VideoRefer-700K [68], which we added into our base training recipe but it lacks temporally focused tasks and annotations. Unlike prior work, **Strefer** pseudo-annotates spatiotemporally dense metadata—subjects, objects, their locations as masklets, as well as their action descriptions and timelines—for complex video scenarios. It generates spatially and temporally fine-grained, grounded instruction-response pairs, requiring no legacy annotations. Without proprietary models or the need to annotate large volumes of new videos, instruction data from **Strefer** empowers models for space-time referring and spatiotemporal reasoning.

## 7. Conclusion

In this work, we introduced **Strefer**, a novel synthetic instruction data generation framework designed to equip Video LLMs with fine-grained, space-time referring and reasoning capabilities. By leveraging temporally dense, pseudo-annotated video metadata, **Strefer** systematically produces instruction-response pairs that are richly grounded in space and time. This enables models to better resolve complex user queries involving object-centric events, temporal references, and gestural or spatial cues—challenges that existing Video LLMs struggle to address.

Our experiments demonstrate that models trained with **Strefer** data outperform baselines on tasks requiring spatiotemporal disambiguation, confirming the efficacy of our approach in fostering more perceptually grounded and context-aware video understanding. **Strefer** thus represents a foundational step toward building next-generation AI companions capable of realistic human-AI interactions and sophisticated spatiotemporal reasoning. Future work

may explore effective methods for filtering synthetic data and advancing the output-level spatiotemporal grounding capabilities of Video LLMs, further enhancing their generalization and reasoning performance.

## References

- [1] Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024. 17
- [2] Ali Athar, Xueqing Deng, and Liang-Chieh Chen. Vicas: A dataset for combining holistic and pixel-level video understanding using captions with grounded segmentation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 19023–19035, 2025. 4
- [3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 4
- [4] Brandon Castellano and contributors. PySceneDetect: Video Scene Cut Detection. <https://www.scenedetect.com/>, 2025. Version 0.6.6 (released March 9, 2025). 5, 7
- [5] Lin Chen, Xilin Wei, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Zhenyu Tang, Li Yuan, et al. Sharegpt4video: Improving video understanding and generation with better captions. *Advances in Neural Information Processing Systems*, 37:19472–19495, 2024. 15, 17
- [6] Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiang-wei Chao, Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, 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. 9
- [7] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024. 4
- [8] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, et al. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms. *arXiv preprint arXiv:2406.07476*, 2024. 4
- [9] Jang Hyun Cho, Andrea Madotto, Effrosyni Mavroudi, Triantafyllos Afouras, Tushar Nagarajan, Muhammad Maaz, Yale Song, Tengyu Ma, Shuming Hu, Suyog Jain, et al. Perceptionlm: Open-access data and models for detailed visual understanding. *arXiv preprint arXiv:2504.13180*, 2025. 4
- [10] Jihoon Chung, Tyler Zhu, Max Gonzalez Saez-Diez, Juan Carlos Niebles, Honglu Zhou, and Olga Russakovsky. Unifying specialized visual encoders for video language models. *arXiv preprint arXiv:2501.01426*, 2025. 4
- [11] DAMO-NLP-SG. Videorefer benchmark evaluation for general mllms. <https://github.com/DAMO->



NLP-SG/VideoRefer/blob/main/benchmark/evaluation\_general\_mllms.md, 2024. Accessed: 2025-06-29. 12

- [12] Henghui Ding, Chang Liu, Shuting He, Xudong Jiang, and Chen Change Loy. Mevis: A large-scale benchmark for video segmentation with motion expressions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 9
- [13] Chaoyou Fu, Yuhang Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24108–24118, 2025. 10, 11
- [14] Kirill Gavriluk, Amir Ghodrati, Zhenyang Li, and Cees G. M. Snoek. Actor and action video segmentation from a sentence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 9
- [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 *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3302–3310, 2025. 17
- [16] Miran Heo, Min-Hung Chen, De-An Huang, Sifei Liu, Subhashree Radhakrishnan, Seon Joo Kim, Yu-Chiang Frank Wang, and Ryo Hachiuma. Omni-rgpt: Unifying image and video region-level understanding via token marks. *arXiv preprint arXiv:2501.08326*, 2025. 4, 17
- [17] Wenyi Hong, Weihang Wang, Ming Ding, Wenmeng Yu, Qingsong Lv, Yan Wang, Yean Cheng, Shiyu Huang, Junhui Ji, Zhao Xue, et al. Cogvlm2: Visual language models for image and video understanding. *arXiv preprint arXiv:2408.16500*, 2024. 4
- [18] Bin Huang, Xin Wang, Hong Chen, Zihan Song, and Wenwu Zhu. Vtimellm: Empower llm to grasp video moments. In *CVPR*, 2024. 17
- [19] Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and Gunhee Kim. Tgif-qa: Toward spatio-temporal reasoning in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2758–2766, 2017. 8
- [20] Qing Jiang, Lin Wu, Zhaoyang Zeng, Tianhe Ren, Yuda Xiong, Yihao Chen, Qin Liu, and Lei Zhang. Referring to any person, 2025. 4, 5, 8
- [21] Xuan Ju, Yiming Gao, Zhaoyang Zhang, Ziyang Yuan, Xintao Wang, Ailing Zeng, Yu Xiong, Qiang Xu, and Ying Shan. Miradata: A large-scale video dataset with long durations and structured captions. *Advances in Neural Information Processing Systems*, 37:48955–48970, 2024. 8, 17
- [22] Evangelos Kazakos, Cordelia Schmid, and Josef Sivic. Large-scale pre-training for grounded video caption generation. *arXiv preprint arXiv:2503.10781*, 2025. 6
- [23] Jie Lei, Tamara L Berg, and Mohit Bansal. Detecting moments and highlights in videos via natural language queries. *Advances in Neural Information Processing Systems*, 34: 11846–11858, 2021. 9, 10, 11, 12
- [24] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 4
- [25] KunChang Li, Yanan He, Yi Wang, Yizhuo Li, Wenhui Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023. 4, 17
- [26] Kunchang Li, Yali Wang, Yanan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22195–22206, 2024. 4, 17
- [27] Lei Li, Yuanxin Liu, Linli Yao, Peiyuan Zhang, Chenxin An, Lean Wang, Xu Sun, Lingpeng Kong, and Qi Liu. Temporal reasoning transfer from text to video. In *ICLR 2025*. OpenReview.net, 2025. 18
- [28] Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language models. 2024. 4
- [29] Long Lian, Yifan Ding, Yunhao Ge, Sifei Liu, Hanzi Mao, Boyi Li, Marco Pavone, Ming-Yu Liu, Trevor Darrell, Adam Yala, et al. Describe anything: Detailed localized image and video captioning. *arXiv preprint arXiv:2504.16072*, 2025. 17
- [30] Jingyang Lin, Jialian Wu, Ximeng Sun, Ze Wang, Jiang Liu, Yusheng Su, Xiaodong Yu, Hao Chen, Jiebo Luo, Zicheng Liu, et al. Unleashing hour-scale video training for long video-language understanding. *arXiv preprint arXiv:2506.05332*, 2025. 18
- [31] Weifeng Lin, Xinyu Wei, Ruichuan An, Tianhe Ren, Tingwei Chen, Renrui Zhang, Ziyu Guo, Wentao Zhang, Lei Zhang, and Hongsheng Li. Perceive anything: Recognize, explain, caption, and segment anything in images and videos. *arXiv preprint arXiv:2506.05302*, 2025. 17
- [32] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 4, 5, 8
- [33] Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and Lu Hou. Tempcompass: Do video llms really understand videos? *arXiv preprint arXiv:2403.00476*, 2024. 4, 10, 11, 18
- [34] Chuofan Ma, Yi Jiang, Jiannan Wu, Zehuan Yuan, and Xiaojuan Qi. Groma: Localized visual tokenization for grounding multimodal large language models. In *European Conference on Computer Vision*, pages 417–435. Springer, 2024. 17
- [35] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*, 2023. 4, 8, 17
- [36] Arjun Mani, Nobline Yoo, Will Hinthorn, and Olga Russakovsky. Point and ask: Incorporating pointing into visual question answering. *arXiv preprint arXiv:2011.13681*, 2020. 4
- [37] Shehan Munasinghe, Rusiru Thushara, Muhammad Maaz, Hanoona Abdul Rasheed, Salman Khan, Mubarak Shah, and

- Fahad Khan. Pg-video-llava: Pixel grounding large video-language models. *arXiv preprint arXiv:2311.13435*, 2023. 6
- [38] Shehan Munasinghe, Hanan Gani, Wenqi Zhu, Jiale Cao, Eric Xing, Fahad Shahbaz Khan, and Salman Khan. Videoglm: A large multimodal model for pixel-level visual grounding in videos. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 19036–19046, 2025. 4
- [39] Long Qian, Juncheng Li, Yu Wu, Yaobo Ye, Hao Fei, Tat-Seng Chua, Yueting Zhuang, and Siliang Tang. Momen-tor: Advancing video large language model with fine-grained temporal reasoning. *arXiv preprint arXiv:2402.11435*, 2024. 17
- [40] Jihao Qiu, Yuan Zhang, Xi Tang, Lingxi Xie, Tianren Ma, Pengyu Yan, David Doermann, Qixiang Ye, and Yunjie Tian. Artemis: Towards referential understanding in complex videos. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 17
- [41] Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdel-rahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Eric Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel grounding large multimodal model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13009–13018, 2024. 17
- [42] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 5, 8
- [43] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 4, 5, 22, 30
- [44] Tianhe Ren, Shuo Shen, et al. Grounded sam 2: Ground and track anything in videos with grounding dino, florence-2 and sam 2. GitHub repository, 2025. <https://github.com/IDEA-Research/Grounded-SAM-2>. 4, 5, 13, 14, 24, 25
- [45] Michael S Ryoo, Honglu Zhou, Shrikant Kendre, Can Qin, Le Xue, Manli Shu, Jongwoo Park, Kanchana Ranasinghe, Silvio Savarese, Ran Xu, et al. xgen-mm-vid (blip-3-video): You only need 32 tokens to represent a video even in vlms. *arXiv preprint arXiv:2410.16267*, 2024. 4, 8, 11, 12, 29
- [46] Eli Schwartz, Leshem Choshen, Joseph Shtok, Sivan Doveh, Leonid Karlinsky, and Assaf Arbelle. Numerologic: Number encoding for enhanced llms’ numerical reasoning. *arXiv preprint arXiv:2404.00459*, 2024. 31
- [47] Ye Sun, Hao Zhang, Henghui Ding, Tiegua Zhang, Xingjun Ma, and Yu-Gang Jiang. Sama: Towards multi-turn referential grounded video chat with large language models. *arXiv preprint arXiv:2505.18812*, 2025. 4, 17
- [48] Qwen Team. Qwen2.5: A party of foundation models, 2024. 5, 8
- [49] Yuli Vasiliev. *Natural language processing with Python and spaCy: A practical introduction*. No Starch Press, 2020. 5
- [50] Haibo Wang, Zhiyang Xu, Yu Cheng, Shizhe Diao, Yufan Zhou, Yixin Cao, Qifan Wang, Weifeng Ge, and Lifu Huang. Grounded-videollm: Sharpening fine-grained temporal grounding in video large language models. *arXiv preprint arXiv:2410.03290*, 2024. 8, 17, 31
- [51] Han Wang, Yongjie Ye, Yanjie Wang, Yuxiang Nie, and Can Huang. Elysium: Exploring object-level perception in videos via mllm. In *European Conference on Computer Vision*, pages 166–185. Springer, 2024. 17
- [52] Jiawei Wang, Liping Yuan, Yuchen Zhang, and Haomiao Sun. Tarsier: Recipes for training and evaluating large video description models, 2024. 4, 5, 8
- [53] Jiannan Wu, Yi Jiang, Peize Sun, Zehuan Yuan, and Ping Luo. Language as queries for referring video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. Also explores Ref-Youtube-VOS, Ref-DAVIS17, A2D-Sentences, JHMD-B-Sentences. 9
- [54] Rujie Wu, Xiaojian Ma, Hai Ci, Yue Fan, Yuxuan Wang, Haozhe Zhao, Qing Li, and Yizhou Wang. Longvit: Instruction tuning for long-form video understanding. *arXiv preprint arXiv:2501.05037*, 2025. 18
- [55] Yongliang Wu, Xinting Hu, Yuyang Sun, Yizhou Zhou, Wenbo Zhu, Fengyun Rao, Bernt Schiele, and Xu Yang. Number it: Temporal grounding videos like flipping manga. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13754–13765, 2025. 8, 12, 17
- [56] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9777–9786, 2021. 1, 8
- [57] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 25th ACM international conference on Multimedia*, pages 1645–1653, 2017. 8
- [58] Jiarui Xu, Xingyi Zhou, Shen Yan, Xiuye Gu, Anurag Arnab, Chen Sun, Xiaolong Wang, and Cordelia Schmid. Pixel-aligned language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13030–13039, 2024. 17
- [59] Lin Xu, Yilin Zhao, Daquan Zhou, Zhijie Lin, See Kiong Ng, and Jiashi Feng. Pllava: Parameter-free llava extension from images to videos for video dense captioning. *arXiv preprint arXiv:2404.16994*, 2024. 4
- [60] Mingze Xu, Mingfei Gao, Zhe Gan, Hong-You Chen, Zhengfeng Lai, Haiming Gang, Kai Kang, and Afshin Dehghan. Slowfast-llava: A strong training-free baseline for video large language models. *arXiv preprint arXiv:2407.15841*, 2024. 4
- [61] Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Viraj Prabhu, Yutong Dai, Michael S Ryoo, et al. xgen-mm (blip-3): A family of open large multimodal models. *arXiv preprint arXiv:2408.08872*, 2024. 8

- [62] An Yan, Zhengyuan Yang, Junda Wu, Wanrong Zhu, Jianwei Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Julian McAuley, Jianfeng Gao, et al. List items one by one: A new data source and learning paradigm for multimodal llms. *arXiv preprint arXiv:2404.16375*, 2024. [17](#)
- [63] Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. *arXiv preprint arXiv:2310.11441*, 2023. [8](#), [12](#)
- [64] Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen Zhang, Zirui Wang, Liangliang Cao, Shih-Fu Chang, and Yinfei Yang. Ferret: Refer and ground anything anywhere at any granularity. *arXiv preprint arXiv:2310.07704*, 2023. [17](#)
- [65] En Yu, Liang Zhao, Yana Wei, Jinrong Yang, Dongming Wu, Lingyu Kong, Haoran Wei, Tiancai Wang, Zheng Ge, Xiangyu Zhang, et al. Merlin: Empowering multimodal llms with foresight minds. In *European Conference on Computer Vision*, pages 425–443. Springer, 2024. [17](#)
- [66] Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yuet-ing Zhuang, and Dacheng Tao. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 9127–9134, 2019. [8](#)
- [67] Yuqian Yuan, Wentong Li, Jian Liu, Dongqi Tang, Xinjie Luo, Chi Qin, Lei Zhang, and Jianke Zhu. Osprey: Pixel understanding with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 28202–28211, 2024. [17](#)
- [68] Yuqian Yuan, Hang Zhang, Wentong Li, Zesen Cheng, Boqiang Zhang, Long Li, Xin Li, Deli Zhao, Wenqiao Zhang, Yueting Zhuang, et al. Videorefer suite: Advancing spatial-temporal object understanding with video llm. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 18970–18980, 2025. [6](#), [8](#), [9](#), [11](#), [12](#), [15](#), [17](#), [18](#)
- [69] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986, 2023. [5](#), [7](#)
- [70] 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. [4](#)
- [71] Hao Zhang, Hongyang Li, Feng Li, Tianhe Ren, Xueyan Zou, Shilong Liu, Shijia Huang, Jianfeng Gao, Leizhang, Chunyuan Li, et al. Llava-grounding: Grounded visual chat with large multimodal models. In *European Conference on Computer Vision*, pages 19–35. Springer, 2024. [17](#)
- [72] Haotian Zhang, Haoxuan You, Philipp Dufter, Bowen Zhang, Chen Chen, Hong-You Chen, Tsu-Jui Fu, William Yang Wang, Shih-Fu Chang, Zhe Gan, et al. Ferret-v2: An improved baseline for referring and grounding with large language models. *arXiv preprint arXiv:2404.07973*, 2024.
- [73] Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Yu Liu, Kai Chen, and Ping Luo. Gpt4roi: Instruction tuning large language model on region-of-interest. In *European Conference on Computer Vision*, pages 52–70. Springer, 2025. [17](#)
- [74] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data, 2024. [15](#), [17](#)
- [75] Henghao Zhao, Ge-Peng Ji, Rui Yan, Huan Xiong, and Zechao Li. Videoexpert: Augmented llm for temporal-sensitive video understanding. *arXiv preprint arXiv:2504.07519*, 2025. [17](#)

## A. Appendix / Supplemental Material

### A.1. More Qualitative Results

#### A.1.2 Strefer-Synthesized Data

#### A.1.2 Strefer-Trained Model

### A.2. More Strefer Details

#### A.3. Model Details

##### A.3.1 Architecture Overview

##### A.3.2 Video Token Representation

##### A.3.3 Masklet Reference Token Representation

##### A.3.4 Timestamp Reference Token Representation

### A.1. Additional Qualitative Results

#### A.1.1. Qualitative Results of Strefer-Synthesized Data

We present qualitative results of our **Strefer**-synthesized data in Fig. 2 and Fig. 11, as well as qualitative results of our novel referring masklet generation pipeline within **Strefer** in Fig. 7, Fig. 8, Fig. 12 and Fig. 13. Detailed discussions of the qualitative results are presented in Sec. 4.2 and Sec. 5 of the main paper.

#### A.1.2. Qualitative Results of Strefer-Trained Model

We present qualitative results of our final **Strefer**-trained model in comparison to the ‘Baseline’ model (from our quantitative result tables). Specifically, results are shown in Fig. 9 and Fig. 14 for the task of **Video Regional Captioning/Description**; Fig. 15, Fig. 16, Fig. 17, Fig. 18 and Fig. 19 for the task of **Video Regional QA**; and Fig. 10, Fig. 20 and Fig. 21 for the task of **Timestamp-Referred Video QA**. Three failure cases are also shown in Fig. 22, Fig. 23 and Fig. 24. Observations are in the figure captions and in Sec. 4.2 and Sec. 5 of the main paper.

### A.2. Additional Strefer Details

Notably, in the design of **Strefer**, we choose masks to accommodate diverse, free-form spatial references from users (e.g., points, scribbles, etc.), which can be readily converted into masks using off-the-shelf tools like SAM2 [43].

We present our designed Referring Masklet Generation Pipeline in Algorithm 1 and the Video Clipper in Algorithm 2. On the right, we list the prompt used for the Video LLM-based Active Entity Recognizer.

### Algorithm 1 Referring Masklet Generation Pipeline

---

**Require:** Video  $\mathcal{V}$ , Referring Expressions  $\mathcal{R} = [r_1, \dots, r_n]$ , Generalized Nouns  $\mathcal{G} = [g_1, \dots, g_n]$

**Ensure:** Masklets aligned to each referring expression  $r_i \in \mathcal{R}$

```

1: procedure GENERATEMASKLETS( $\mathcal{V}, \mathcal{R}, \mathcal{G}$ )
2:    $\mathcal{F}, S \leftarrow \text{SAMPLEANDREORDERFRAMES}(\mathcal{V})$ 
3:    $f^*, D_{f^*} \leftarrow \text{SELECTINITIALFRAME}(S, \mathcal{G}, \mathcal{R})$ 
4:    $\mathcal{M} \leftarrow \text{BIDIRECTIONALSEGMENTATIONTRACKING}(\mathcal{F}, f^*, D_{f^*})$ 
5:    $\mathcal{M} \leftarrow \text{ASSIGNEXPRESSIONSTOMASKLETS}(\mathcal{M}, f^*, \mathcal{R}, \mathcal{G})$ 
6:   return  $\mathcal{M}$ 
7: end procedure

8: procedure SAMPLEANDREORDERFRAMES( $\mathcal{V}$ )
9:    $\mathcal{F} \leftarrow \text{sample frames from } \mathcal{V}$ 
10:   $S \leftarrow \text{reorder } \mathcal{F} \text{ using a middle-first recursive strategy}$ 
11:  return  $(\mathcal{F}, S)$ 
12: end procedure

13: procedure SELECTINITIALFRAME( $S, \mathcal{G}, \mathcal{R}$ )
14:   $\text{max\_count} \leftarrow -1, f^* \leftarrow S[-1], \text{best\_frame} \leftarrow S[0], \text{best\_detections} \leftarrow \emptyset$ 
15:  for each frame  $f \in S$  do
16:     $D_f \leftarrow \text{GROUNDINGDINO}(f, \text{set}(\mathcal{G}))$ 
17:    if  $|D_f| > \text{max\_count}$  then
18:       $\text{max\_count} \leftarrow |D_f|, \text{best\_frame} \leftarrow f, \text{best\_detections} \leftarrow D_f$ 
19:    end if
20:    if  $|D_f| \geq |\mathcal{R}|$  then
21:      return  $(f, D_f)$ 
22:    end if
23:  end for
24:  if  $f^* == S[-1]$  and  $\text{best\_frame} \neq f^*$  then
25:     $f^* \leftarrow \text{best\_frame}, D_{f^*} \leftarrow \text{best\_detections}$ 
26:  end if
27:  return  $(f^*, D_{f^*})$ 
28: end procedure

29: procedure BIDIRECTIONALSEGMENTATIONTRACKING( $\mathcal{F}, f^*, D_{f^*}$ )
30:  Define forward sequence:  $\mathcal{F}^{\rightarrow} = [f^*, f^*+1, \dots]$ 
31:  Define backward sequence:  $\mathcal{F}^{\leftarrow} = [f^*, f^*-1, \dots]$ 
32:  Initialize  $\mathcal{T} \leftarrow []$ 
33:  for each clip  $\mathcal{C} \in \{\mathcal{F}^{\rightarrow}, \mathcal{F}^{\leftarrow}\}$  do
34:     $\mathcal{M}_{\mathcal{C}} \leftarrow \text{SAM2}(\mathcal{C}, D_{f^*}, \text{video})$ 
35:    Append  $\mathcal{M}_{\mathcal{C}}$  to  $\mathcal{T}$ 
36:  end for
37:  return  $\text{MERGE\_TRACKING\_RESULTS}(\mathcal{T})$ 
38: end procedure

39: procedure ASSIGNEXPRESSIONSTOMASKLETS( $\mathcal{M}, f^*, \mathcal{R}, \mathcal{G}$ )
40:  Partition  $\mathcal{M}$  into groups based on  $\mathcal{G}$ 
41:  for each group  $G$  in partitioned  $\mathcal{M}$  do
42:     $\mathcal{B}_G \leftarrow \text{available bounding boxes on frame } f^* \text{ in group } G$ 
43:    for each referring expression  $r_i \in \mathcal{R}$  do
44:       $b_i \leftarrow \text{REXSEEK}(f^*, \mathcal{B}_G, r_i)$ 
45:      Assign  $r_i$  to mask in  $G$  corresponding to box  $b_i$ 
46:    end for
47:  Post-process assignments to ensure valid mapping
48:  end for
49:  return  $\mathcal{M}$ 
50: end procedure

```

---

#### Prompt P.1: Entity Recognizer

Prompt: What “active” scene entities can you identify from the video? An entity refers to an object, and “active” scene entities are scene objects that have any dynamic behaviors, such as actions, interactions with others, or movements. Please compile a list of clearly visible “active” scene entities from the video. Use entity appearance in concise description to distinguish one “active” scene entity from another if possible.





Figure 11. **Example of Strefer-Synthesized Instruction-Response Pairs (left) and Pseudo-Annotated Video Metadata (right).** Each instruction begins with the prefix: “Please answer the following question about the <region>” (and the prefix is omitted in the figure). For each instruction-response pair, the boundary of the object mask referred to by <region> is shown next to the pair and highlighted in color. **Strefer** automatically clips the video into segments and pseudo-annotates the video metadata—including active entities, their locations (as masklets), and their action descriptions and timelines—for complex video scenarios, such as scenes containing multiple entities of the same category, and cases where entities do not appear in the first frame, or temporarily exit and re-enter the frame; based on the auto-generated video metadata, it produces instruction-response pairs, requiring no legacy annotations or manual efforts. Though current implementation of **Strefer** does not any use proprietary models, without the need to annotate large volumes of new videos, instruction data from **Strefer** empowers models for space-time referring and spatiotemporal reasoning (ref. Table 2, 3, and 4).

---

Prompt = “man in the grey shirt.man in the black jacket.”

---

Video Frames



Ours



man in the grey shirt

man in the black jacket

Prior Method - GroundedSAM2 [44]



man

man

Figure 12. **Qualitative Results of Referring Masklet Generation.** In this video, our method accurately generates masklets corresponding to the input referring expressions. In contrast, GroundedSAM2 [44] fails to differentiate between the man in the grey shirt and the man in the black jacket.

### A.3. Model Details

#### A.3.1. Architecture Overview

The Video LLM processes a video and a user’s multimodal query to generate a textual response. A multimodal query consists of the textual component of the question, a masklet

along with its associated frames referring to a specific region within the video, and optionally, one or more specific timestamps within the video.

The architecture of the Video LLM is illustrated in Fig. 5. At a high level, the LLM processes four types of input tokens: (i) *visual tokens*, which encode the global context of

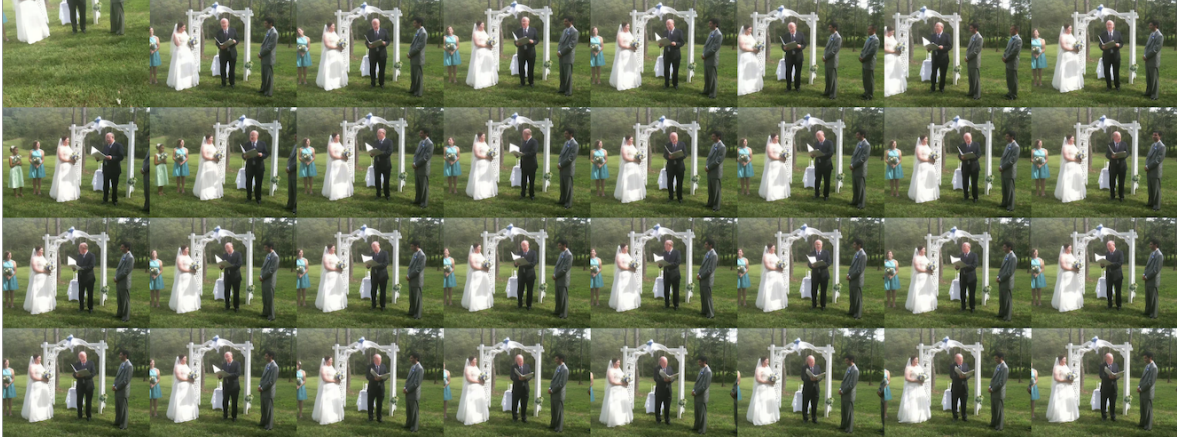


---

Prompt = “bride.groom.bridesmaid.officiant.”

---

Video Frames



Ours



bride groom officiant bridesmaid

Prior Method - GroundedSAM2 [44]



groom officiant groom bride bridesmaid

Figure 13. **Qualitative Results of Referring Masklet Generation.** Our method accurately generates masklets corresponding to the input referring expressions. In contrast, GroundedSAM2 [44] converts the input text prompt into the class names *groom officiant*, *groom*, *bride*, and *bridesmaid*. GroundedSAM2 then fails to detect the bridesmaid in the video, while incorrectly assigning the class name *bridesmaid* to the actual bride, *groom* to the officiant, and *groom officiant* to the groom.





Please give a detailed description of the highlighted object <region> in the video.

Ground Truth

The parking lot is a large, open area with multiple cars parked in designated spaces. The surface is smooth and appears to be made of asphalt. The parking lot is bordered by orange and white barriers on one side, and there are several orange cones placed throughout the area.

Baseline

A car is parked on the left side of the road. It is a white vehicle with a blue stripe running along its side. The car is stationary and does not appear to be in motion.



The object in question is a parking lot. Initially, it is depicted as a vast, open space with a smooth, dark asphalt surface. The parking lot is bordered by a series of orange and white barriers, which are arranged in a neat, parallel fashion, creating a clear demarcation for the area. The barriers are evenly spaced, ensuring a uniform appearance.

Figure 14. **Qualitative Results of Video LLMs Trained w/ and w/o Strefer-Synthethized Data.** This sample is drawn from VideoRefer-Bench<sup>D</sup>, designed to assess a model’s performance on the task of **Mask-Referred Regional Description**. The boundary of the region referred to by the mask in this sample is highlighted in purple. The baseline model tends to have a foreground bias, fails to interpret the masklet correctly, and mistakenly responds that the referred region is a car. In contrast, the model trained on **Strefer**-generated data accurately identifies the masklet-referred region as a parking lot.



What action is <object0><region> performing? Pick the best option:

- (A) Sitting and eating
- (B) Turning around and moving to the right
- (C) Lying down to eat
- (D) Walking to the leftmost part

Ground Truth

(B) Turning around and moving to the right

Baseline

(A) Sitting and eating

(B) Turning around and moving to the right



Figure 15. **Qualitative Results of Video LLMs Trained w/ and w/o Strefer-Synthethized Data.** This sample is drawn from VideoRefer-Bench<sup>Q</sup>, designed to assess a model’s performance on the task of **Mask-Referred Regional QA**. The boundary of the region referred to by the mask in this sample is highlighted in purple. The model trained on **Strefer**-generated data correctly identifies the masklet-referred region and action.

the video; (ii) *region tokens*, which represent specific visual regions referenced in the user query (e.g., a mask or masklet); (iii) *timestamp/temporal tokens*, which indicate

particular temporal locations within the video; and (iv) *text tokens*, which represent the textual content of the query itself. These tokens are jointly fed into the LLM, which then

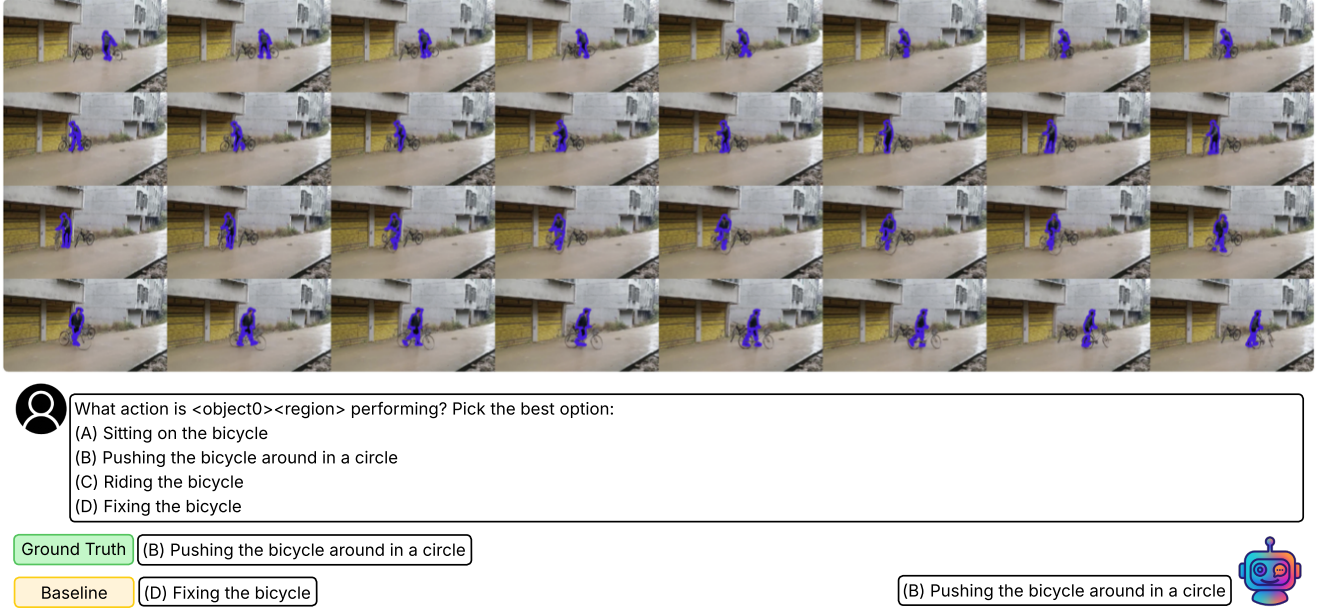


Figure 16. **Qualitative Results of Video LLMs Trained w/ and w/o Strefer-Synthethized Data.** This sample is drawn from VideoRefer-Bench<sup>Q</sup>, designed to assess a model’s performance on the task of **Mask-Referred Regional QA**. The boundary of the region referred to by the mask in this sample is highlighted in purple. In this sample, the model must demonstrate fine-grained spatiotemporal action understanding due to the small size of the mask and the subtle motion differences between the correct and negative options. The model trained on **Strefer**-generated data successfully identifies both the region referred to by the masklet and the corresponding action.

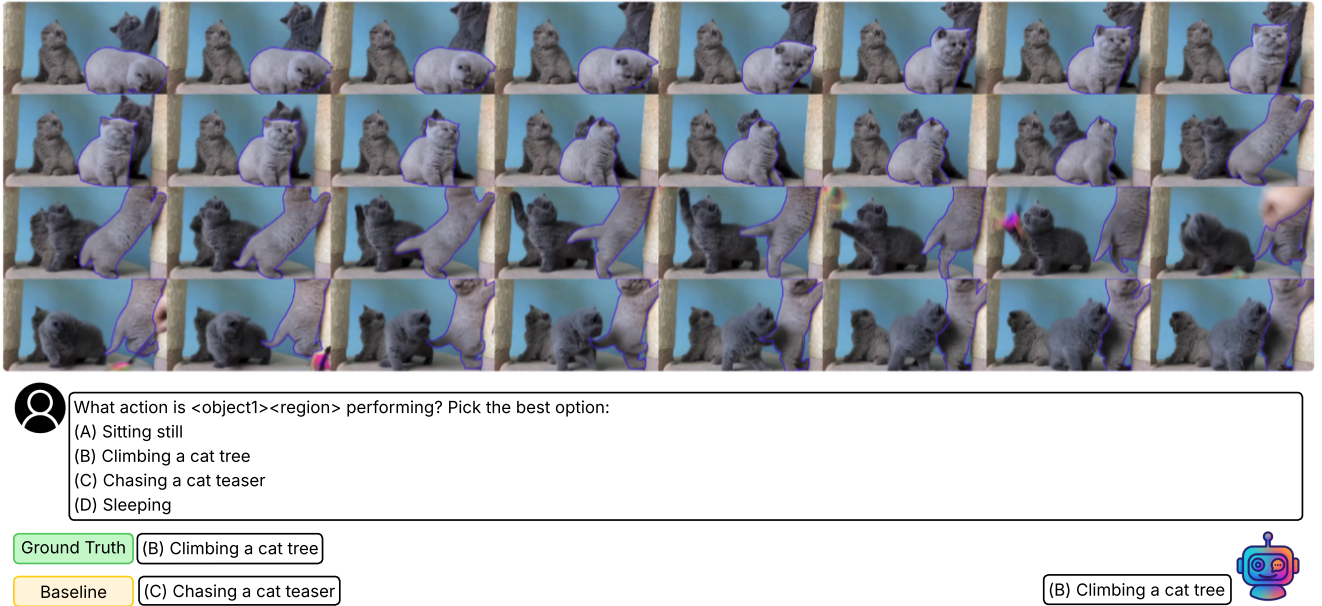


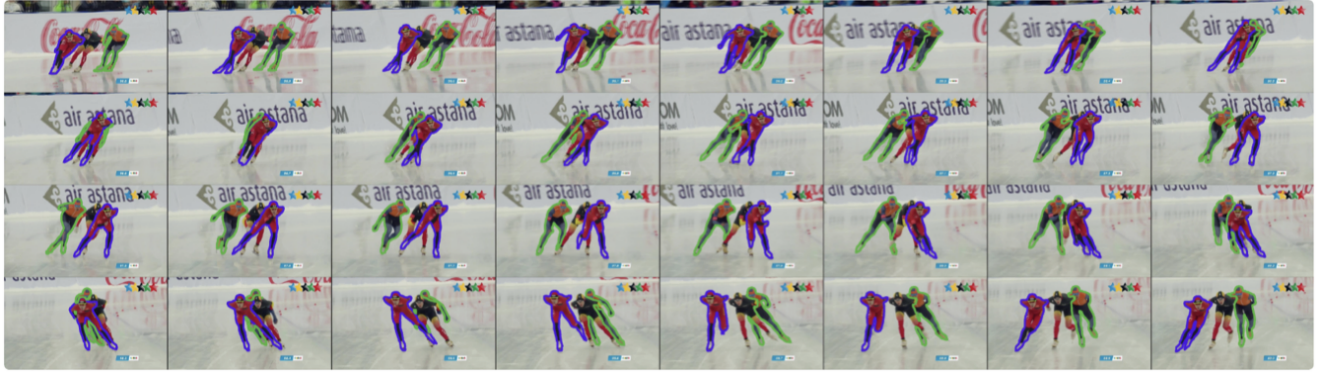
Figure 17. **Qualitative Results of Video LLMs Trained w/ and w/o Strefer-Synthethized Data.** This sample is drawn from VideoRefer-Bench<sup>Q</sup>, designed to assess a model’s performance on the task of **Mask-Referred Regional QA**. The boundary of the region referred to by the mask in this sample is highlighted in purple. The model trained on **Strefer**-generated data correctly identifies the masklet-referred region and action.

auto-regressively generates a textual response.

The construction of visual, region, and timestamp tokens from raw inputs—namely, the video and multimodal user

query—is detailed in Sec. A.3.2, Sec. A.3.3, and Sec. A.3.4, respectively.





How does <object1><region>'s position relate to <object3><region>'s position? Pick the best option:

- (A) <object1> is in front of <object3>
- (B) <object1> is behind <object3>
- (C) <object1> is beside <object3>
- (D) <object1> is not visible in relation to <object3>

Ground Truth

(A) <object1> is in front of <object3>

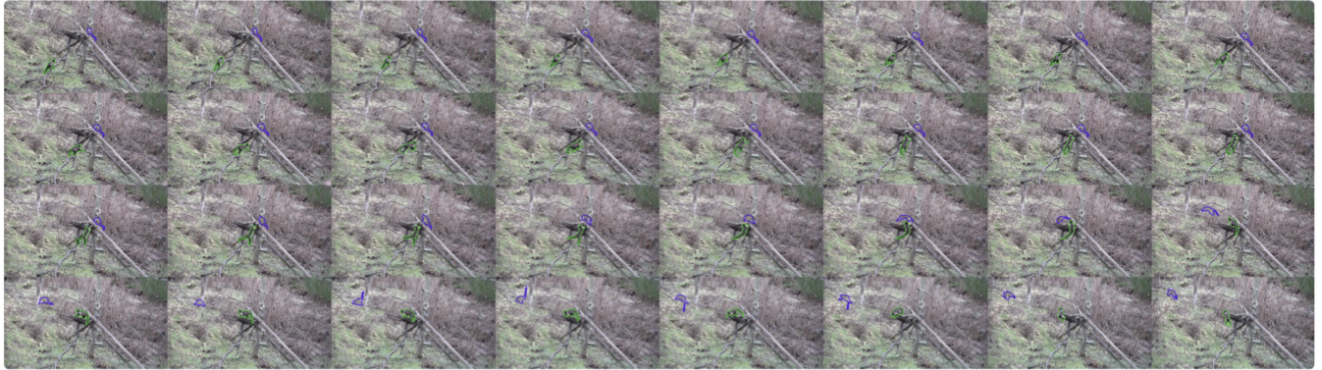
Baseline

(B) <object1> is behind <object3>

(A) <object1> is in front of <object3>



Figure 18. **Qualitative Results of Video LLMs Trained w/ and w/o Strefer-Synthethized Data.** This sample is drawn from VideoRefer-Bench<sup>Q</sup>, designed to assess a model's performance on the task of **Mask-Referred Regional QA**. This sample presents a multi-masklet scenario, with two masklets referring to two different individuals. The boundary of the <object1> region is highlighted in purple, and <object2> is highlighted in green. The model trained on **Strefer**-generated data correctly answers this multi-masklet reference question by effectively analyzing the relationship between the two masklets within the video context.



How does the position of <object1><region> relate to <object2><region>? Pick the best option:

- (A) <object1> is above <object2> on the tree
- (B) <object1> is below <object2> on the tree
- (C) <object1> and <object2> are on the same branch
- (D) <object1> is on a different tree than <object2>

Ground Truth

(A) <object1> is above <object2> on the tree

Baseline

(B) <object1> is below <object2> on the tree

(A) <object1> is above <object2> on the tree



Figure 19. **Qualitative Results of Video LLMs Trained w/ and w/o Strefer-Synthethized Data.** This sample is drawn from VideoRefer-Bench<sup>Q</sup>, designed to assess a model's performance on the task of **Mask-Referred Regional QA**. This sample presents a multi-masklet scenario, with two masklets referring to two different individuals. The boundary of the <object1> region is highlighted in purple, and <object2> is highlighted in green. Kindly zoom in, as the regions are relatively small and may be difficult to discern. The model trained on **Strefer**-generated data correctly answers this multi-masklet reference question by effectively analyzing the relationship between the two masklets within the video context.



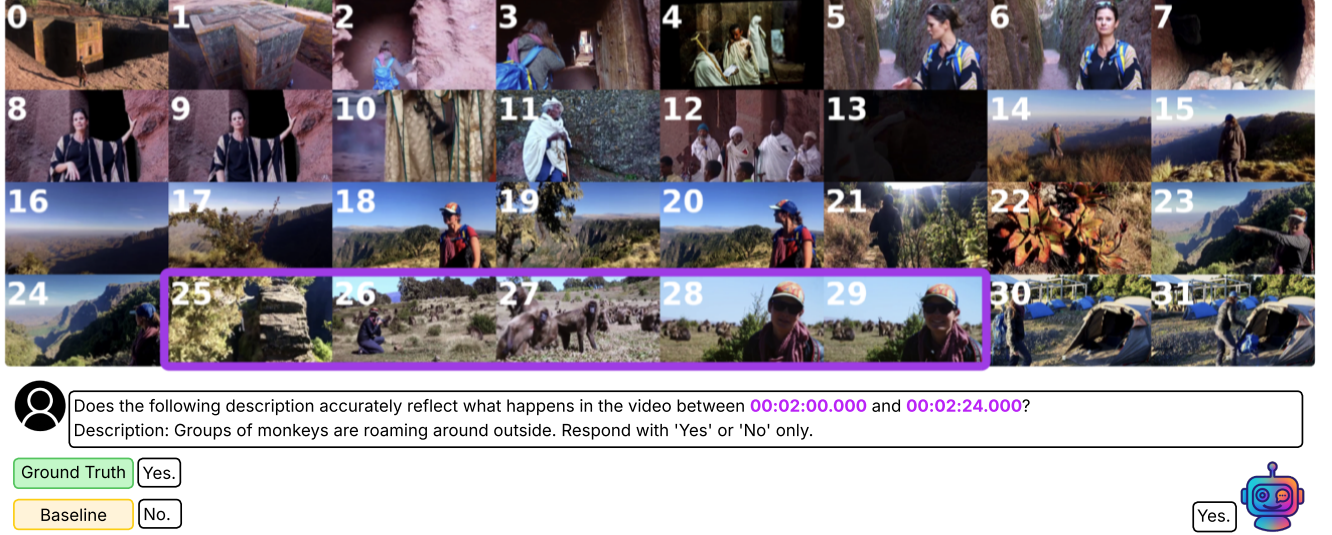


Figure 20. **Qualitative Results of Video LLMs Trained w/ and w/o Strefer-Synthethized Data.** This sample is drawn from QVHighlights, using our repurposed task designed to assess a model’s performance on **Timestamp-Referred Video QA**. The boundary of segment corresponding to the timestamps in the question is highlighted in purple. The model trained on our **Strefer**-generated data correctly answers the question, demonstrating superior understanding of precise moments and segments in videos compared to the baseline.

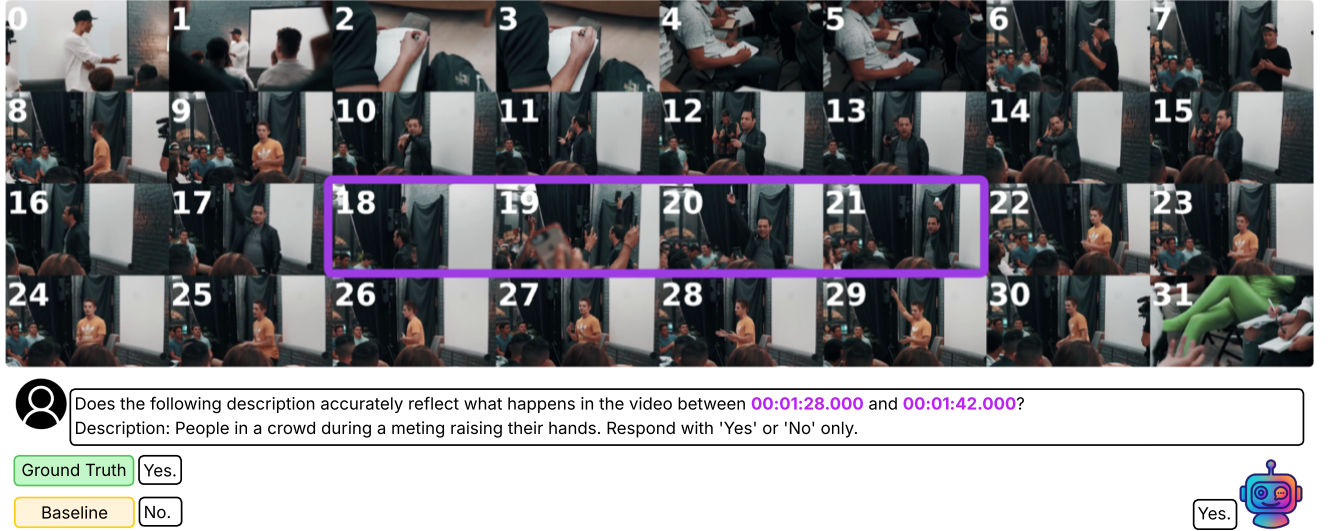


Figure 21. **Qualitative Results of Video LLMs Trained w/ and w/o Strefer-Synthethized Data.** This sample is drawn from QVHighlights, using our repurposed task designed to assess a model’s performance on **Timestamp-Referred Video QA**. The boundary of segment corresponding to the timestamps in the question is highlighted in purple. The model trained on our **Strefer**-generated data correctly answers the question, demonstrating superior understanding of precise moments and segments in videos compared to the baseline.

### A.3.2. Video Token Representation

Given an input video  $x_v \in \mathbb{R}^{T_v \times 3 \times H_v \times W_v}$ , where  $T_v$  is the number of frames and  $H_v, W_v$  are the height and width of the frames, a visual encoder extracts the video’s global visual features  $f_v \in \mathbb{R}^{t_v \times d_v \times h_v \times w_v}$ .

A Video-Language Connector is then applied on top of the visual encoder to project the global visual features into

a sequence of visual tokens  $e_v \in \mathbb{R}^{L_v \times d}$ , where  $d$  represents the dimensionality of the language model’s input token space, and  $L_v$  is the number of visual tokens of a video. This connector aligns the visual features to the input space of a language model while preserving semantics relevant for multimodal understanding. In some designs (e.g., BLIP-3-Video [45] that our model architecture is based on), the connector also incorporates a token compression module to

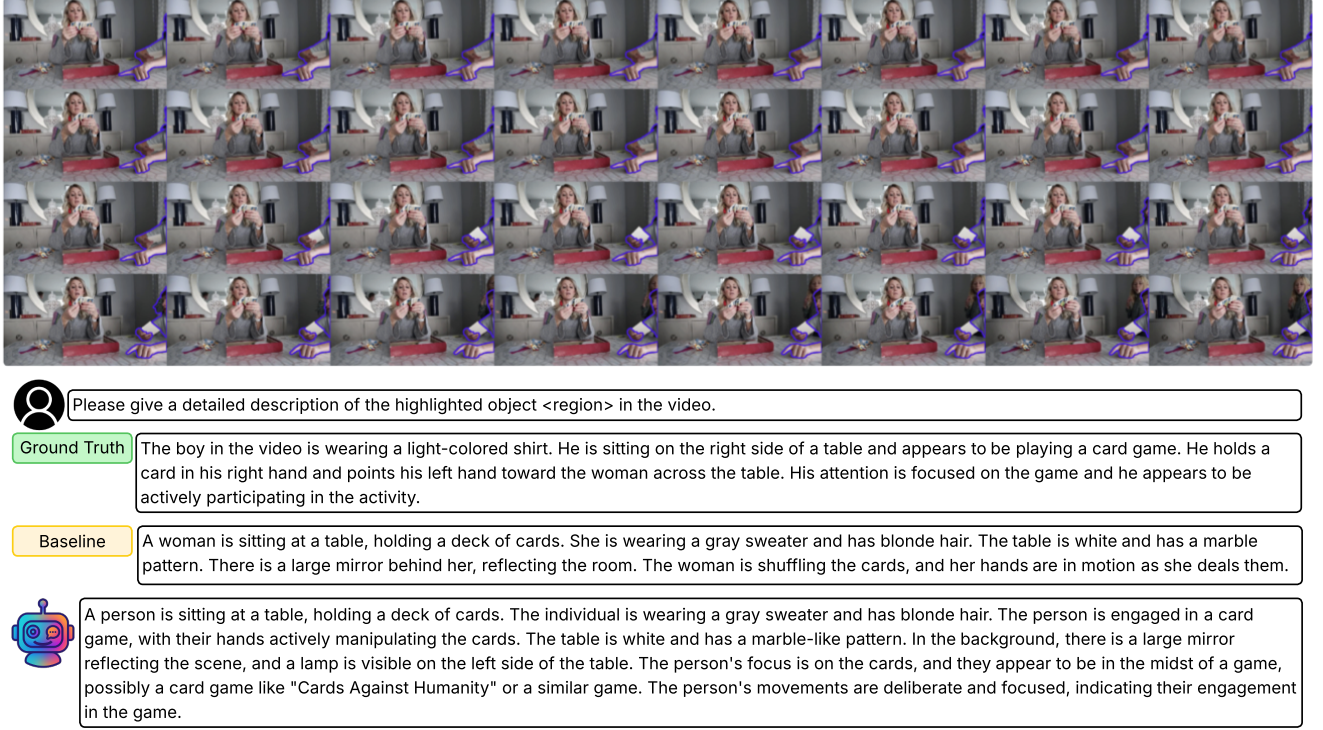


Figure 22. **Failure Results of Video LLMs Trained w/ and w/o Strefer-Synthethized Data.** This sample is drawn from VideoRefer-Bench<sup>D</sup>, designed to assess a model’s performance on the task of **Mask-Referred Regional Description**. The boundary of the region referred to by the mask in this sample is highlighted in purple. The masklet is intended to refer to the boy on the right, but he is mostly out of view, while a woman appears prominently in the center of the video. Both the baseline model and the model trained on **Strefer**-generated data fail to correctly interpret the masklet.

#### Algorithm 2 Video Clipping Pipeline

```

1: procedure CLIPVIDEO(video)
2:    $B \leftarrow \text{PYSCENEDETECT}(\text{video}, \text{threshold}=20)$ 
3:   if  $B = \emptyset$  and  $\text{DURATION}(\text{video}) \geq 3$  sec then
4:      $E \leftarrow \text{GETEMBEDDINGS}(\text{video})$ ;
5:      $D \leftarrow \text{PAIRWISEDISTANCES}(E)$ 
6:      $T \leftarrow \text{CLUSTERINGAUTOTHRESHOLD}(D, 1.7)$ ;
7:      $C \leftarrow \text{HIERARCHICALAGGLOMERATIVECLUSTERING}(E, T)$ 
8:      $B \leftarrow \text{EXTRACTCLIPTIMESTAMPBOUNDARIES}(C)$ 
9:   end if
10:  return  $B$ 
11: end procedure

12: procedure CLUSTERINGAUTOTHRESHOLD( $D, f$ )
13:   $m \leftarrow \text{mean}(D)$ ;  $s \leftarrow \text{std}(D)$ ;  $M \leftarrow \text{max}(D)$ 
14:  return  $\min(m + f \cdot s, M)$ 
15: end procedure

```

reduce the number of tokens, improving efficiency without sacrificing critical information.

#### A.3.3. Masklet Reference Token Representation

Our modified Video LLM is designed to understand user queries about videos that involve spatial or spatiotemporal, local regional references. To support diverse, free-form spatial reference from users (e.g., points, scribbles, etc.), we standardize them by converting these free-form spatial references into regional masks before processed by the model.

This approach is effective because many forms of spatial reference can be easily transformed into masks using off-the-shelf tools like SAM2 [43].

**Mask and Masklet.** A regional mask is represented as a 2D binary matrix  $\mathbb{R}^{H_m \times W_m}$ , where  $H_m$  and  $W_m$  are the height and width of the image containing the region of interest, with a value of 1 inside the region and 0 outside. When extended over time, a temporal sequence of such regional masks  $x_r \in \mathbb{R}^{T_m \times H_m \times W_m}$  is referred to as a masklet. Since a mask is special case of masklet with only one frame, we describe the masklet feature extraction process below.

**Masklet Token Representation.** Leveraging the same visual encoder, our model extracts image feature maps  $f_m \in \mathbb{R}^{t_m \times d_v \times h_m \times w_m}$  for the frames that contain the masklet  $x_r$ . The masklet  $x_r$  and its corresponding frames’ feature maps  $f_m$  are then processed by a Region-Language Connector, which outputs region tokens  $e_r \in \mathbb{R}^{L_r \times d}$  that are aligned to the language space, where  $L_r$  is a predefined number of region tokens.

The Region-Language Connector begins by resizing the binary masklet  $x_r$  via bilinear interpolation to match the spatial (and temporal if the visual encoder condenses the



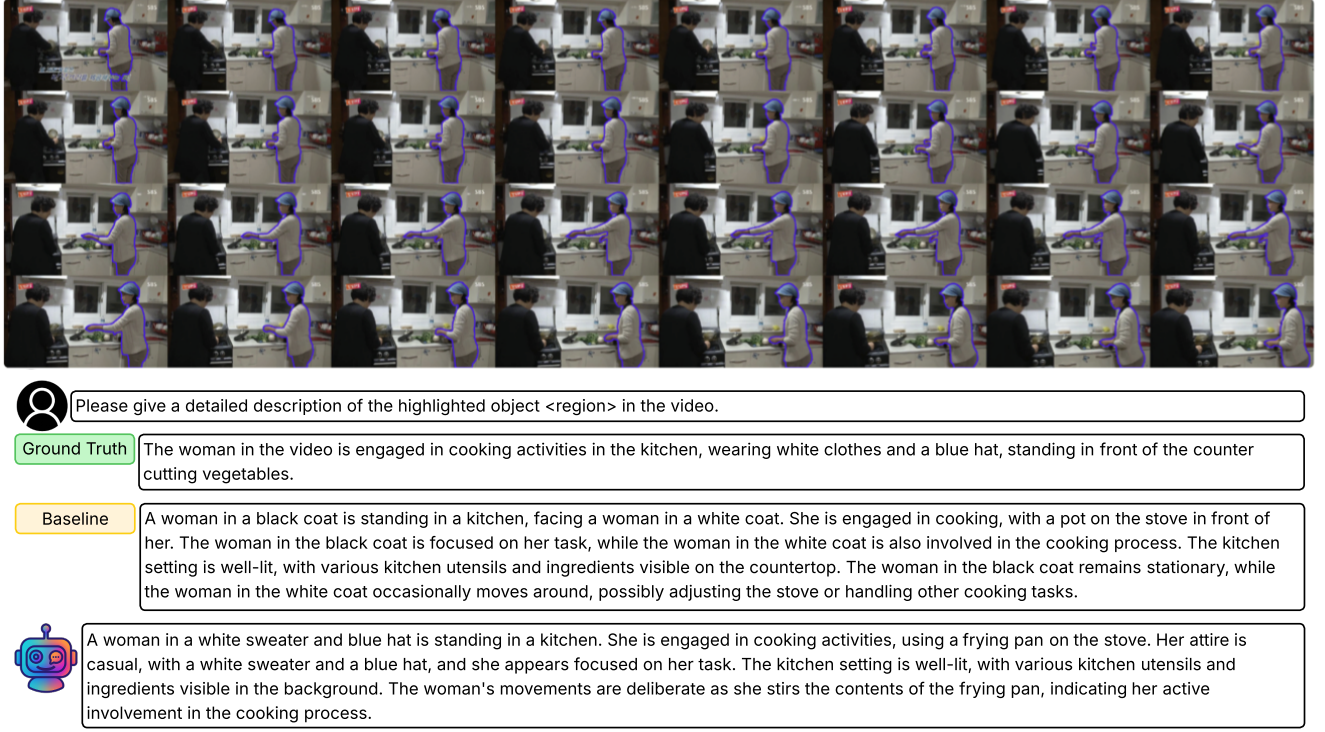


Figure 23. **Failure Results of Video LLMs Trained w/ and w/o Strefer-Synthethized Data.** This sample is drawn from VideoRefer-Bench<sup>D</sup>, designed to assess a model’s performance on the task of **Mask-Referred Regional Description**. The boundary of the region referred to by the mask in this sample is highlighted in purple. While the model trained on **Strefer**-generated data correctly identifies that the masklet refers specifically to the woman in the white sweater, it incorrectly responds that her action is “frying pan”.

time axis) dimensions of  $f_m$ , yielding a resized masklet of shape  $\mathbb{R}^{t_m \times h_m \times w_m}$ . A Mask Pooling operation is then applied: average pooling is performed over the spatial locations within the mask region for each frame, producing a pooled feature representation  $p \in \mathbb{R}^{t_m \times d_v}$ . This representation can be interpreted as a sequence of  $t_m$  region tokens, each of dimensionality  $d_v$ .

To reduce the temporal redundancy, a Temporal Token Merge module condenses the  $t_m$  tokens into  $L_r$  representative ones ( $L_r < t_m$ ). Specifically, for  $p \in \mathbb{R}^{t_m \times d_v}$ , cosine similarities are computed between each pair of temporally adjacent tokens:

$$s_{i,i+1} = \frac{p^i \cdot p^{i+1}}{\|p^i\| \cdot \|p^{i+1}\|}, \quad 0 \leq i < t_m - 1 \quad (1)$$

This yields a similarity vector  $s \in \mathbb{R}^{t_m-1}$ . A similarity threshold  $\theta$  is then selected as the  $L_r$ -th largest value in  $s$ . Next, the sequence  $p$  is processed sequentially from the beginning to the end to form token groups. An initially empty group is created and the first token in  $p$  is added to it. For each index  $i$  from 0 to  $t_m - 2$ , if  $s_{i,i+1} \geq \theta$ , then  $p^{i+1}$  is added to the current group. Otherwise, the current group is finalized, and a new group is initiated with  $p^{i+1}$ .

This process produces exactly  $L_r$  token groups. Each

group is finally merged into a single representative token by averaging the embeddings of all tokens within the group.

Finally, the resulting  $L_r$  tokens, each in  $\mathbb{R}^{d_v}$ , are projected into the language embedding space via an MLP, producing the final region tokens  $e_r \in \mathbb{R}^{L_r \times d}$ .

Note that the Temporal Token Merge module is bypassed when the user query involves only a single frame mask (as opposed to a masklet).

#### A.3.4. Timestamp Reference Token Representation

By design, our model is effective in scenarios where users may refer to specific times within videos in their queries. However, LLMs often struggle with interpreting numerical values [46]. To address this challenge, we adopt the Temporal Token Representation method introduced in Grounded-VideoLLM [50], which discretizes continuous time into a sequence of temporal tokens, making time-related reasoning more manageable for LLMs.

Suppose the video has a duration of  $L$  seconds. We divide it into  $M$  equal-length, non-overlapping, and non-spacing segments, resulting in  $M+1$  anchor points that span from the start to the end of the video. These anchor points, labeled from  $\langle 0 \rangle$  to  $\langle M \rangle$ , represent evenly spaced temporal positions throughout the video. Each specific times-



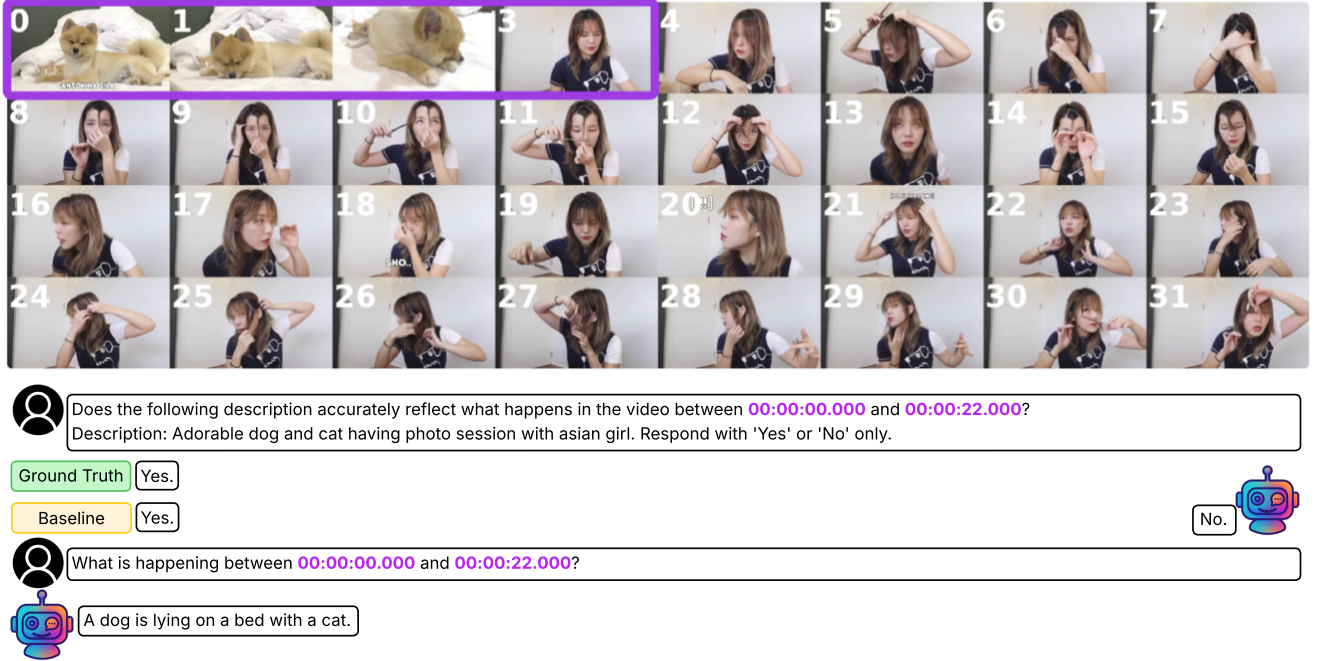


Figure 24. **Failure Results of Video LLMs Trained w/ and w/o Strefer-Synthethized Data.** This sample is drawn from QVHighlights, using our repurposed task designed to assess a model’s performance on **Timestamp-Referred Video QA**. The boundary of segment corresponding to the timestamps in the question is highlighted in purple. Although the model trained on our **Strefer**-generated data fails to answer the question correctly, it does accurately recognize that the segment shows a dog lying on a bed with a cat. We suspect the model’s failure stems from its disagreement with the description, which is not fully grounded in the visual content—for example, the video segment does not clearly depict a photo session involving the dog, the cat, and the girl.

tamp within the video is mapped to an anchor point and then encoded as a temporal token. For example,  $\langle 0 \rangle$  marks the beginning of the video, while  $\langle M \rangle$  represents the end. These  $M + 1$  anchor points are added to the LLM’s vocabulary (by expanding the LLM’s vocabulary), enabling unified modeling of time alongside text. Mathematically:

A specific continuous timestamp  $\tau$  can be easily converted to a temporal token  $\langle t \rangle$  and vice versa:

$$t = \text{Round} \left( M \cdot \frac{\tau}{L} \right), \quad \tau = L \cdot \frac{t}{M} \quad (2)$$

In this way, specific timestamps in the user query are converted into timestamp anchor points. Both text and timestamp anchor points are then mapped to embeddings through the extended word embedding layer of the LLM, forming interleaved text tokens and temporal tokens.

In our model tuning and evaluation experiments, since our Video LLM processes 32 input frames, we set  $M = 31$  to learn 32 temporal tokens.

Type ID & Task	Frames	Source	Format	Example Question-Answer Pair	Mask-Refer Version
1. Ask the model to describe the behavior of entities that are present in a segment of the video.	Frames are extracted only from the segment of the video.	Template	OE	Question: <video>What is happening to the woman? Answer: The woman is engaged in a dance with the man, involving spins and turns. She is lifted off the ground by the man during the dance.	Question:<video>Please answer the following question about the <region>. What is happening to her? Answer: The woman is engaged in a dance with the man, involving spins and turns. She is lifted off the ground by the man during the dance.
2. Ask the model to describe the behavior of entities that are not present in a segment of the video; the model should respond with uncertainty (e.g., "Sorry, I'm not sure").	Frames are extracted only from the segment of the video.	Template	OE	Question: <video>What is currently happening to the person in a green hoodie? Answer: The person in a green hoodie seems to be not clearly visible.	N/A
3. Ask a yes/no question about the presence of an entity in a segment of the video; if present, the model should describe its behavior; if absent, the model should respond with uncertainty.	Frames are extracted only from the segment of the video.	Template	OE	Question: <video>Were you able to see a woman in a black jacket? Answer: Yes. The woman walks towards the child seated on the sofa.	N/A
4. Ask a yes/no question about the presence of an entity in a segment of the video; the model should respond with a concise "Yes" or "No" only.	Frames are extracted only from the segment of the video.	Template	OE	Question:<video> Is there a woman in a black jacket? Answer only "Yes" or "No". Answer: Yes.	N/A
5. Ask the model to identify the correct temporal order in which entities first appear in the video from multiple choices.	Frames are extracted from the full video.	Template	MCQ	Question:<video> Which order shows their first appearance in the video? (A) child interacting with the plant bed, child holding a bag and a toy, child walking across the lawn (B) child holding a bag and a toy, child approaching a plant bed, child interacting with the plant bed (C) child approaching a plant bed, child holding a bag and a toy, child interacting with the plant bed (D) child walking across the lawn, child holding a bag and a toy, child interacting with the plant bed Answer: (B)	N/A
6. Ask the model to describe the behavior of entities that may or may not be present in a specific time range of the video; the question refers to a time range.	Frames are extracted from the full video.	Template	OE	Question:<video> Could you explain what the girl in the yellow coat is doing between 00:00:05 and 00:00:12.210? Answer: The girl in the yellow coat is carefully watering plants in a garden.	Question:<video>Please answer the following question about the <region>. Could you explain what she is doing between 00:00:05 and 00:00:12.210? Answer: The girl in the yellow coat is carefully watering plants in a garden.
7. Ask the model to describe what happened generally or to a specific entity during a specific time range in the video; the question refers to a time range.	Frames are extracted from the full video.	LLM	OE & MCQ	Question:<video> What else did the woman interviewing the man do between 00:00:00 and 00:00:07.007? Answer: The woman interviewing the man is talking as well.	Question:<video>Please answer the following question about the <region>. What else did she do between 00:00:00 and 00:00:07.007? Answer: The woman interviewing the man is talking as well.
8. Ask the model to describe what happened generally or to a specific entity during a coarse time range in the video (e.g., throughout the video, beginning, middle, or end).	Frames are extracted from the full video.	LLM	OE & MCQ	Question:<video> What else did the woman interviewing the man do in the beginning of the video? Answer: The woman interviewing the man is talking as well.	Question:<video>Please answer the following question about the <region>. What else did she do in the beginning of the video? Answer: The woman interviewing the man is talking as well.
9. Ask the model to identify when a specific behavior or event occurs within the video; expect the model to answer with a coarse time range in the video (e.g., throughout the video, beginning, middle, or end).	Frames are extracted from the full video.	LLM	OE & MCQ	Question:<video> During which part of the video was the child in pink dress riding a tricycle? Answer: The beginning.	Question:<video>Please answer the following question about the <region>. During which part of the video was this person riding a tricycle? Answer: The beginning.
10. Ask the model to describe the behavior of an entity before/during/after something else occurs.	Frames are extracted from the full video.	LLM	OE & MCQ	Question:<video> What is the adult doing while the child is riding a tricycle? Answer: The adult is watching and walking behind the child.	Question:<video>Please answer the following question about the <region>. What is he doing while the child is riding a tricycle? Answer: The adult is watching and walking behind the child. -
11. Ask the model to identify the entity involved before/during/after something else occurs.	Frames are extracted from the full video.	LLM	OE & MCQ	Question: Who is walking behind the child in blue while the child is riding a tricycle? Answer: An adult wearing a black shirt.	Question:<video>Please answer the following question about the <region>. Who is walking behind this child while the child is riding a tricycle? Answer: An adult wearing a black shirt.

Table 7. **Details of Strefer-synthesized video instruction data.** The table details the question task types, their visual inputs, QA generation sources, formats, examples, and the mask-referring versions of the QAs. 'OE' denotes open-ended QA, and 'MCQ' indicates multiple-choice QA.