

Video-LLMs with Temporal Visual Screening

Anonymous ACL submission

Abstract

Current Video-LLMs lack human-like temporal screening capabilities and their behavioral manifestations (e.g. scrubbing the seek bar), struggling with fine-grained semantics due to a static, sparse frame sampling paradigm. To address this, we introduce **Temporal Visual Screening (TVS)**, a cognitively inspired novel task formalized as an endomorphic transformation, $\mathcal{F}_{TVS} : (V, Q) \rightarrow (v', q')$, which maps a task instance to a computationally tractable and semantically coherent one within the same bimodal space. This is achieved by excising irrelevant segments while synchronously rewriting the query to ensure answer invariance. Designed as a modular front-end adapter under VideoQA setting, TVS can be seamlessly integrated into **both** Video Instruction Tuning (training) and Video Question Answering (inference) pipelines. To facilitate this new task, we introduce the first dedicated benchmark and propose ReSimplifyIt, a baseline that surpasses prior methods from similar tasks by 0.47 F-1 score with competitive query rewriting performance. Integrating TVS yields significant relative gains of 7.33% in training and 33.7% in inference, validating the efficacy of our approach for video-language understanding.¹

1 Introduction

Human cognition optimizes information processing via a two-stage control mechanism: a fast, pre-attentive screening to prune redundancy, followed by focused reasoning on a compacted, goal-aligned representation (Treisman and Gelade, 1980; Wolfe, 1994; Zacks et al., 2007; Hochstein and Ahissar, 2002). Behaviorally, this manifests as scrubbing the seek bar to locate regions of interest before detailed viewing (Wu and Xie, 2024), minimizing Extraneous Load to purify Germane Load (Sweller, 1988; Mayer and Moreno, 2003; Paas et al., 2003),

enabling **goal-oriented** problem reconstruction before further reasoning.

In contrast, current Video-LLMs (Chen et al., 2023; Li et al., 2023a; Xu et al., 2024a; Maaz et al., 2024; Li et al., 2024c; Ma et al., 2024) ingest frames statically and equally, forcing them to confront the unmitigated reasoning load of the raw $(Video, Query)$ instance. This lack of reasoning-guided temporal skimming prevents the elimination of superfluous perceptual burdens, diluting semantic focus on critical segments and introducing noise. These inefficiencies hamper both training (Lin et al., 2024; Zhang et al., 2025)—by weakening supervision signals—and inference—by reducing task-oriented focus.

The traditional grounding or grounded QA task (Xiao et al., 2024; Chen et al., 2024a) would fail to address these, as grounding relies on **depictive-level** similarity and **perceptual-level** alignment to enhance *visual cues*, arising from a deeper misalignment between the problem representation and the goal state: across time during reasoning, the query often entangles multi-hop temporal, relational, and causal dependencies with sparse, weakly localized visual evidence.

This necessitates an interpretable, controllable, goal-oriented mechanism that optimizes reasoning structure and cognitive load while treating input modalities **as a whole**. We propose **Temporal Visual Screening (TVS)**, a cognitively inspired **endomorphic** transformation:

$$\mathcal{F}_{TVS} : (V, Q) \rightarrow (v, q)$$

that reconstructs a computationally tractable, semantically coherent instance within the same multimodal task space. TVS performs in a meticulously designed goal-invariant manner, which aligns with the fixity of the Goal State from the Problem Space Theory (Newell and Simon, 1972), enhances granularity of multimodal alignment on the cognitive

¹Our code is available [here](#).

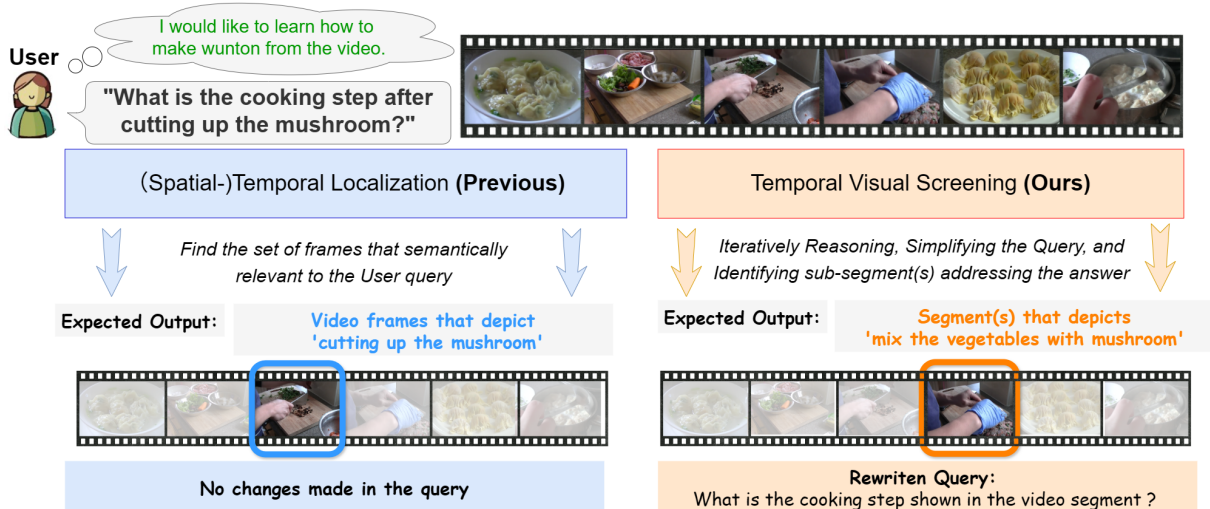


Figure 1: Input and output example of the TVS task as a side-by side comparison to a superficially similar task: temporal localization. TVS focuses on goal-driven cognitive alignment and performs reasoning-guided problem reconstruction through information screening, whereas grounding only emphasizes perceptual alignment and locates depictive level visual evidence by direct moment retrieval.

level, and facilitates pinpointed cognitive focus for detailed multimodal understanding.

The benefits of TVS land in four concrete areas: (i) Training alignment: Under fixed budgets, removing off-target supervision boosts gradient signal-to-noise ratios and reduces spurious correlations; (ii) Inference robustness: Offloading non-informative perceptual burdens shortens the implicit reasoning program, unleashing reasoning capabilities; (iii) Interpretability: The compact (v', q') serves as a controllable, faithful artifact exposing model bottlenecks; (iv) Generalizability: As a model-agnostic front-end, TVS benefits diverse VideoQA agents.

We present *ReSimplifyIt*, a plug-and-play multi-agent framework, with its key design drawing inspirations from cognitive science principles, being able to **ReS**iliently **im**provise and **qual**ify proposals **It**eratively. It chains iterative screening rounds where a language-only Launcher module hypothesizes trimming instruction and rewritten query while utilizing the intrinsic relation between the two input modalities as a prior, a Validator executes and critiques these plans through a Viewer module, and memory trackers to assist self-correction. This design operationalizes a competence-from-consequence (Brooks, 1991) principle where trial execution provides constructive feedback that tightens the coupling between hypothesized reasoning structure and available evidence.

We further construct YouCookII-TVS, the first benchmark dedicated to evaluating TVS as a joint transformation over video and query, enabling as-

essment of both screening quality and downstream VideoQA performance.

Through quantitatively benchmarking TVS and its impact on both downstream training and inference, we show that TVS is a non-trivial challenging task that remains far from being solved and reveals a key bottleneck for downstream VideoQA reasoning, and thus deserves to be formulated, benchmarked, and studied independently.

Our work is arguably the first to formally recognize that grounding and trimming a video creates a semantic mismatch with a complex, context-dependent query, to formalize the task upon it, and to extensively explore temporal filtering mechanism into the Video-LLM framework. Our results demonstrate that TVS not only improves inference-time performance but also facilitates more structured and effective training by finer-grained multimodal alignment, paving the way for more scalable, interpretable, and cognitively-inspired video-language understanding systems.

The main contributions of this work are summarized below:

- We propose a novel cognitively-inspired task, temporal visual screening (TVS), which plays a pivotal role in video understanding by guiding task-aware reasoning structure filtering, thereby addressing a key bottleneck in multi-modal alignment. In addition, we propose a benchmark YouCookII-TVS for the TVS task, which comprises 238.2 Hours of videos and 2754 questions.
- We introduce the ReSimplifyIt framework, the

first baseline for the TVS task, as a model-agnostic plug-in compatible with **any** existing videoQA and instruction tuning pipelines.

- We quantitatively evaluate the impact of TVS in VideoQA and Video Instruction Tuning. Our method brings considerable performance boost to multiple strong baselines over various VideoQA benchmarks (over 10% and 5% absolute gain on inference and training stage, respectively).

2 Formulation of Temporal Visual Screening (TVS)

2.1 Task Definition

TVS is elegantly designed to pinpoint critical video focus without altering data formats, ensuring full procedural compatibility with any existing VideoQA and Video Instruction Tuning pipelines as a front-end adapter. Refer to Figure 1 for illustration. Formally, given a full-length video $V = \{F_1, F_2, \dots, F_T\}$ consisting of T ordered frames and a natural language query Q , the **Temporal Visual Screening (TVS)** task aims to output:

- a reconstructed query q , derived from Q , and
- a new **continuous** video

$$v = \{F_{s_1}, F_{s_1+1}, \dots, F_{e_1}, F_{s_2}, F_{s_2+1}, \dots, F_{e_2}, \dots, F_{s_n}, F_{s_n+1}, \dots, F_{e_n}\} \subseteq V$$

being the concatenation of n non-overlapping segments of the original video, where $1 \leq s_1 \leq e_n \leq T$ and $n \geq 1$,

such that v contains the minimal but sufficient visual information to correctly answer q , and preserves *low-level temporal continuity* essential for motion, spatial, and object-level reasoning. Specifically, we define a function:

$$\mathcal{F}_{\text{TVS}} : (V, Q) \rightarrow (v, q)$$

subject to the efficacy conditions given next.

Efficacy conditions for TVS Maintaining both *uni-modal efficacy* and *cross-modal synergy*, TVS facilitates temporal visual focus by trimming out insignificant segments of the video, while *synchronously* reconstructing the text query to preserve original semantic alignment. For demonstration, refer to the Appendix. The transformations applied to q and v must be tightly coupled to preserve answer consistency despite input reduction, ensuring semantic alignment and outcome equivalence between the reduced input and the original task. Formally, It must meet the following:

1. **Continuity for Video Output:** v must be a temporally **contiguous** subsegment of V , or a concatenation of such segments, preserving low-level visual integrity necessary for tracking motion, appearance, and scene dynamics.
2. **Minimal Sufficiency for Answer:** A frame $F_t \in v$ is included *if and only if* removing it would alter the model output: $\exists \delta a$ such that $\mathcal{M}(q, v \setminus \{F_t\}) = a + \delta a$, $|\delta a| > \epsilon$ where ϵ is a task-specific confidence or semantic tolerance threshold. This guarantees that each included frame contributes non-trivially to the answer.
3. **Reasoning Minimality for Query Output:**

$$q^* = \underset{q}{\operatorname{argmin}} \operatorname{Infer}(q)$$

where $\operatorname{Infer}(q)$ denotes the number of reasoning steps needed to derive the answer from visual input. For example, reasoning on the phrase “within 5 seconds after event A” requires three reasoning steps on the temporal axis: localization on event A, relational reasoning on “after”, and temporal reasoning on “5 seconds”. This ensures that q focuses on a *minimal grounded fact* answerable within a concise visual span.

4. **Cross-modal Synergy:** For any reasonable VideoQA agent \mathcal{M} , the answer to q given v must match the answer given the full video:

$$\forall \mathcal{M}, \quad \mathcal{M}(q, v) = \mathcal{M}(Q, V).$$

This enforces that v retains all information necessary to support the same answer as the full input, regardless of model architecture.

This formulation enables precise, interpretable alignment between question, answer, and visual evidence, while ensuring temporal coherence and model-agnostic consistency. Refer to Figure 3 for more demonstrations of the TVS task.

2.2 Integrating TVS into Video-LLM Workflows

Temporal Visual Screening (TVS) can be effectively incorporated into both the training and inference workflows of video-language models, enabling cleaner supervision, stronger alignment, and more interpretable model behavior. Below, we describe two major modes of application.

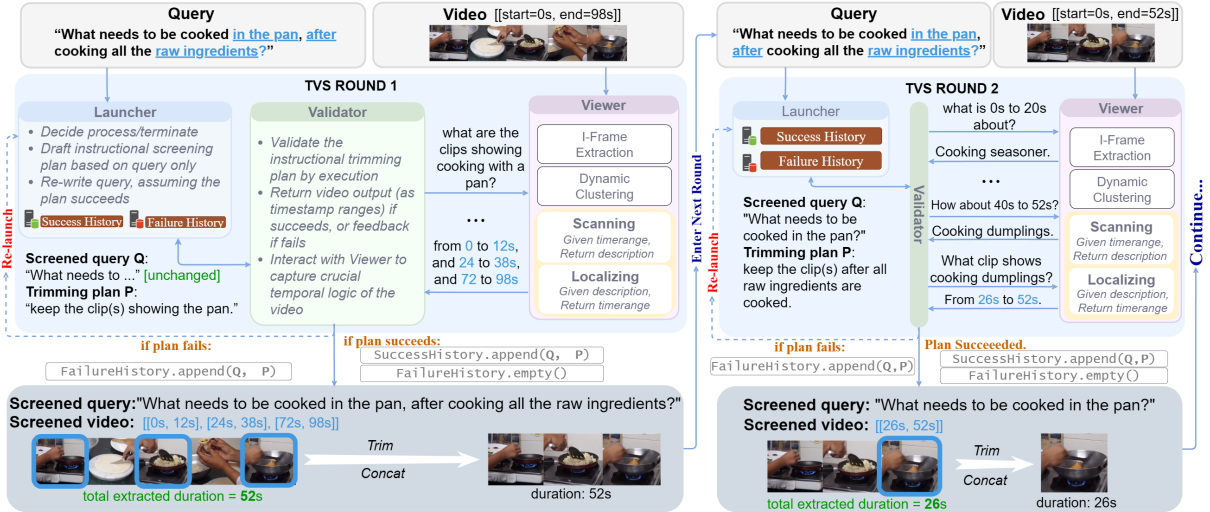


Figure 2: Snapshot examples of the workflow of our proposed ReSimplifyIt framework.

TVS for Inference-Time Simplification in Video-LLMs At inference, TVS acts as a lightweight pre-processing module. Given a video V and query Q , it identifies a salient segment v and reconstructs a simplified query q (see Section 3). This yields a new input pair (v', q') , reducing the temporal and semantic complexity the model must handle. Rather than processing the full video or ambiguous and overgeneralized instructions, the model concentrates on compact, relevant content, facilitating fine-grained temporal reasoning and direct visual grounding. This design aims to improve inference efficiency, robustness, and interpretability.

TVS for Improving Training-Time Visual-Language Alignment Training Video-LLMs often involves complex supervision tuples (V, Q, a) , where queries may be abstract or multi-hop. TVS mitigates this by transforming them into more compact triples (v, q, a) , where v is a high-utility segment and q is a grounded reformulation of the original query. This decomposition offloads off-target supervision to upstream modules, allowing the model to specialize in two core competencies: temporal understanding and multimodal alignment.

3 Method

3.1 ReSimplifyIt Framework

Drawing inspirations from *dynamic and interactive attention deployment* (Treisman and Gelade, 1980; Wolfe, 1994) and humans' information screening process by dragging the video's progress bar, we propose our multi-agent ReSimplifyIt framework. The framework is composed of three main

agentic components: the Launcher, the Validator, and the Viewer, accompanied by two extra helper modules as memory trackers: the Failure History and the Success History. For implementation details, please refer to Section 4. We denote $\{(v_r, q_r) | 1 \leq r \leq R\}$ as the whole TVS process, where v_r, q_r represents the output video clip and question of the r th round, respectively, and R sets the maximum number of rounds. Refer to Algorithm 2 for the complete algorithm of our TVS process performed by our ReSimplifyIt framework, and Figure 2 for a snapshot demonstration.

Initiating a TVS Round: Launcher drafting a video trimming instruction The TVS process unfolds over iterative rounds, akin to how humans drag the progress bar multiple times based solely on empirical surmise before actually viewing the video content. In each round r , a **Launcher** module generates a trimming instruction i_r and a revised query q_r as a trial, which trimming instruction may only contain *high-level* semantics yielding a declarative goal (e.g. "keep the clip after event X" instead of "clip([5s, 10s]"), based solely on the previous query q_{r-1} , without vision access. This trial-based methodology strikes an effective balance between performance and efficiency as the inherent semantic relevance between the query and video input can thereby serve as a prior to enabling plausible instruction generation without video access.

Akin to human landing on undesired segments after dragging the progress bar, we equip the Launcher module with self-correction based adaptive refinement mechanisms to tackle rare failure cases and ensure correctness. Specifically, we

maintain two memory tracker modules: **Success History (SH)** and **Failure History (FH)**. A successful trial is added to *SH* and terminates the round; otherwise, it is recorded in *FH* and re-triggers the Launcher with updated feedback.

We formalize the Launcher’s behavior as:

$$(i_r, q_r) = \text{Launcher}(q_{r-1}, FH, SH, t)$$

where t denotes the task prompt.

Our feedback-driven formulation achieves iterative self-correction, a key distinction from prior work (Wang et al., 2024b; Shang et al., 2024). It decomposes the task into progressive rounds to enable iterative self-correction and foster structured reasoning, which enhances multi-hop robustness while minimizing correction costs. The lightweight, video-free Launcher further boosts efficiency by performing abductive reasoning on the language query without sacrificing fidelity.

Validating a trial: Validator executing instruction The **Validator module** receives the high-level instruction i_r from the Launcher and determines its success through execution. The module performs two tightly coupled tasks: (i) *assessing the feasibility* of the instruction (i.e., whether it is able to succeed), and (ii) *executing* the instruction to obtain results or feedback. Unlike the Launcher, the Validator *has indirect access to visual semantics* by interacting with the Viewer module while not taking visual inputs or frame captions itself.

It returns a tuple (d_r, m_r) , where:

- d_r indicates whether the instruction is deemed succeeded (feasible) or failed (infeasible),
- m_r contains either the resulting trimmed video (in the form of timestamp ranges) if successfully executed the instruction, or an explanation if failed.

$$(d_r, m_r) = \text{Validator}(i_r, v_{r-1}, t)$$

Here, $\text{Validator}()$ denotes the Validator module, v_{r-1} is the previous video state, and t is the task prompt.

The Validator may also decide to call the Viewer module before returning to Launcher. In this scenario, d_r and m_r stand for this decision symbol and the message to the Viewer, respectively.

Our unified Validator assesses plan viability via direct execution rather than handcrafted rules, embodying a “competence from consequence” philosophy (Brooks, 1991). This integration simplifies

control flow and reduces inter-module communication. Crucially, failures are not terminal but provide constructive feedback for the Launcher to refine instructions, transforming execution into a mechanism for both validation and continual adaptation.

Scanning the video: Viewer scanning and localizing Inspired by the Bottom-Up (stimulus-driven) and Top-Down (goal-driven) Activation schema proposed by GS2.0 (Wolfe, 1994), which also guides humans’ scrubbing the seek bar when watching videos, we design the **Viewer** module to explicitly incorporate two complementary tasks:

(i) **Scanning**: summarize the content of a video snippet given a timestamp range;

(ii) **Localizing**: retrieve a timestamp range given a text summarization of a video snippet.

Following complementarily symmetric task structure, this design of the Viewer module achieves elegant and efficient bidirectional navigation by providing both top-down (content-driven) and bottom-up (time-driven) exploration, making the Viewer module highly flexible.

To support both tasks, we first perform a two-stage keyframe extraction and captioning process: (1) extract I-frames using MPEG-4 compression to capture key visual content, and (2) apply a dynamic frame clustering algorithm to select the final keyframes, which adaptively adjust the number of clusters without supervision. Compared to frame clustering techniques in previous work (Wang et al., 2024c), this approach demonstrates stronger generalizability and robustness. For the algorithm of the keyframe extraction process, please refer to algorithm 1. This pre-processing is denoted as:

$$C = \text{prep}(start, end).$$

Scanning executes by summarizing the snippet content via LLM reasoning over C , optionally querying additional frames, and is denoted as:

$$Cap = \text{Scanner}(\text{prep}(start, end), t).$$

Localizing follows a lightweight three-stage search: (1) locate top-k candidate timestamps, (2) select the best, and (3) expand it into a full range. The LLM may optionally query extra frames for confirmation, ensuring minimal frame access while maintaining accuracy. This process is denoted as:

$$(t_{start}, t_{end}) = \text{Localizer}(\text{prep}(0, d), q).$$

This lightweight yet effective three-stage design achieves fine-grained temporal grounding with minimal overhead, showcasing the Viewer’s adaptability and plug-and-play potential.

400 Algorithms of the keyframe extraction and local- 438
401 ization stages are provided in appendix A. 439

402 3.2 TVS-guided inference 440

403 Attributed to its endomorphic property, TVS can 441
404 serve as a plug-and-play adapter process for any 442
405 VideoQA pipeline, including both inference and 443
406 training stages. Given the input question q and 444
407 video V , TVS-guided inference is conducted as: 445

$$\begin{aligned} response &= VideoQA(V', q') \\ V', q' &= TVS(V, q) \end{aligned}$$

408 where $VideoQA$ can be any VideoQA pipeline, 450
409 including video-LLMs, LLM-assisted pipelines, or 451
410 any others alternatives, and $TVS()$ stands for the 452
411 Temporal Visual Screening process. 453

412 3.3 TVS-guided training 454

413 TVS can also be applied to training stage to purify 455
414 supervision signal. Given input question q , input 456
415 video V , and ground truth answer a , we compute 457
416 the likelihood during the TVS-guided training as: 458

$$p(\mathbf{X}_A | \mathbf{X}_V, \mathbf{X}_Q) = \prod_{i=1}^L p_{\theta}(\mathbf{X}_A^{[i]} | \mathbf{X}_V, \mathbf{X}_Q, \mathbf{X}_A^{[1:i-1]})$$

$$\mathbf{X}_A, \mathbf{X}_Q, \mathbf{X}_V = f_t(a), f_t(q'), f_v(v')$$

$$V', q' = TVS(V, q)$$

417 where f_t, f_v are the text and visual tokenizers. 461

418 4 Experiment Settings 462

419 4.1 YouCookII-TVS benchmark 463

420 As TVS is a new task, few benchmarks is capable 464
421 of performing the evaluation. While some studies 465
422 (Lei et al., 2018; Chen et al., 2024b; Lei et al., 466
423 2020) have explored grounded VideoQA which 467
424 may be mistakenly seen as equivalent, they mainly 468
425 focus on improving visual evidence. See Appendix 469
426 F for details. The synchronous update of video 470
427 and query inputs in TVS naturally resolves any 471
428 potential semantic mismatch, maximally enhancing 472
429 generalizability and adaptability of our work. 473

430 To bridge this gap, we introduce the YouCookII- 474
431 TVS benchmark, which provides both TVS and 475
432 standard VideoQA labels. Built upon the 476
433 YouCookII dataset (Zhou et al., 2018), it comprises 477
434 2,754 datapoints split into training, validation, and 478
435 test sets (roughly 7:2:1). This dual-task bench- 479
436 mark enables unified evaluation of both TVS and 480
437 VideoQA. Refer to Appendix B for more details. 481

4.2 Datasets and Benchmarks 438

439 We conduct the evaluation from three aspects. 440

441 Firstly, we evaluate the quality of the TVS pro- 442
443 cess on **YouCookII-TVS** (Section 4.1). Next, 444
445 we examine the impact of TVS on down- 446
447 stream VideoQA performance across the fol- 448
449 lowing benchmarks: **YouCookII-TVS**, which 450
451 supports evaluation of both the TVS process 452
453 and VideoQA performance; **ActivityNet-QA** (Yu 454
455 et al., 2019), **NExT-OE** (Xiao et al., 2021), 456
457 **Video-MME** (Fu et al., 2024), **MLVU** (Zhou 458
459 et al., 2024a), **LVBench**(Wang et al., 2024a), and 460
461 **EgoSchema** (Mangalam et al., 2023). Lastly, 462
463 we investigate the role of TVS in video instruc- 464
465 tion tuning by fine-tuning Video-LLM baselines 466
467 on **YouCookII-TVS**, and comparing their down- 468
469 stream performance across various benchmarks. 469
470 For **NExT-QA** (Xiao et al., 2021), we conduct 471
472 open-ended generation during inference and map 473
474 predictions to MCQ format using GPT-3.5, en- 474
475 suring consistent evaluation across baselines and 475
476 benchmarks (see the Appendix for details). 476

4.3 Implementation Details 460

461 We provide implementation details in Appendix D. 461

4.4 Baselines 462

463 **Temporal Visual Screening** Considering the 463
464 lack of baselines of the new TVS task, we adopt 464
465 two baselines from the training-free temporal local- 465
466 ization task which are considered to be robust and 466
467 generalizable to unseen datasets, to make a side- 467
468 by-side comparison of the video output alone of 468
469 TVS. Specifically, we adopt VTG-GPT (Xu et al., 469
470 2024b), a proposal-based method which made one 470
471 of the first attempts to training-free video temporal 471
472 grounding, and Zheng et al. (2024a)’s work which 472
473 comprehend candidate proposals based on both 473
474 static and dynamic matching scores. For the query 474
475 output of the TVS process, we report open-ended 475
476 evaluation result of our proposed framework. 476

477 **TVS-guided inference** Integrating TVS as a 477
478 front-end module into existing VideoQA agents 478
479 yields a novel VideoQA framework. We adopt sev- 479
480 eral Video-LLMs—Video-ChatGPT (Maaz et al., 480
481 2024), Video-LLaVA (Lin et al., 2023), ChatU- 481
482 niVi (Jin et al., 2023), LLaVA-NeXT (Liu 482
483 et al., 2024b), InternVL3.5 (Wang et al., 2025), 483
484 Qwen2.5VL (Bai et al., 2025b), Qwen3VL (Bai 484
485 et al., 2025a), LLaVA-OneVision (Li et al., 485
486 2024a)—as representative baselines. We also 486

Method	Temporal Relational		Timepoint Indexed		Multifaceted Integrative		Average	
	mIoU	F1	mIoU	F1	mIoU	F1	mIoU	F1
VTG-GPT (Xu et al., 2024b)	0.17	0.29	0.15	0.26	0.16	0.24	0.16	0.27
Zheng et al. (2024a)	0.11	0.19	0.04	0.08	0.06	0.10	0.07	0.12
ReSimplifyIt (Ours)	0.23	0.37	0.98	0.99	0.47	0.64	0.56	0.67

(a) Results on video output

Method	Temporal Relational	Timepoint Indexed	Multifaceted Integrative	Average
ReSimplifyIt (Ours)	66.8	78.5	72.8	72.7

(b) Results on query rewriting

Table 1: Stage-1 evaluation results on YouCookII-TVS dataset.

Method	Size	ActivityNetQA		YC2TVS		EgoSchema		LVBench		MLVU		Video-MME		NExT-OE	
		w/tvs	w/o	w/tvs	w/o	w/tvs	w/o	w/tvs	w/o	w/tvs	w/o	w/tvs	w/o	w/tvs	w/o
Video-ChatGPT (Maaz et al., 2024)	7B	58.5	50.5	38.9	28.91	29.2	23.0	23.5	22.9	22.5	18.6	31.1	29.5	49.8	41.5
Video-LLaVA (Lin et al., 2023)	7B	61.2	52.1	43.2	32.4	43.8	40.0	26.5	23.2	31.9	27.5	43.4	39.1	48.0	43.3
ChatUniVi (Jin et al., 2023)	7B	63.2	52.0	56.5	47.4	–	–	–	–	–	–	–	–	34.5	25.9
LLaVA-NExT (Liu et al., 2024b)	7B	67.8	62.5	58.2	48.6	38.6	38.2	31.4	23.8	34.2	28.3	43.8	40.3	52.3	43.7
InternVL3.5 (Wang et al., 2025)	8B	59.9	57.2	51.5	45.7	67.2	61.8	46.3	39.6	46.8	45.1	60.4	56.9	59.1	53.9
InternVL3.5 (Wang et al., 2025)	14B	60.1	58.9	52.1	48.8	70.6	67.6	47.8	42.9	46.5	45.2	62.2	61.1	58.2	55.6
Qwen2.5-VL (Bai et al., 2025b)	3B	57.6	55.8	47.9	40.3	61.8	57.2	44.2	36.3	45.3	41.3	59.0	57.7	54.8	51.6
Qwen2.5-VL (Bai et al., 2025b)	7B	58.0	55.2	48.0	41.1	68.2	66.4	43.0	34.8	43.1	37.6	60.0	58.5	55.2	53.4
Qwen3-VL (Bai et al., 2025a)	4B	59.9	57.7	48.0	45.6	71.2	70.2	43.2	37.5	49.0	42.8	60.5	60.6	61.1	58.2
Qwen3-VL (Bai et al., 2025a)	32B	60.9	60.5	47.8	45.2	72.8	72.6	–	–	44.9	38.7	64.0	61.2	61.4	58.9
LLaVA-OneVision (Li et al., 2024a)	4B	–	–	–	–	18.6	18.2	–	–	13.5	5.5	–	–	–	–
LLaVA-OneVision (Li et al., 2024a)	8B	32.4	27.5	–	–	38.4	32.5	–	–	14.1	12.2	51.1	51.3	–	–
VideoAgent (Wang et al., 2024b)	-	61.5	60.2	53.9	38.4	–	–	–	–	–	–	–	–	47.8	49.6
VideoTree (Wang et al., 2024c)	-	63.6	59.0	69.4	57.1	–	–	–	–	–	–	–	–	61.9	57.7
GPT-4o (OpenAI et al., 2024a)	-	75.65	72.2	73.84	63.59	72.38	71.17	46.22	35.04	48.92	44.73	66.91	61.63	62.25	54.9
GPT-4.1-mini (OpenAI et al., 2024b)	-	77.07	75.19	76.04	68.74	71.12	72.02	46.35	34.43	50.67	39.90	62.45	61.20	68.93	57.45
GPT-4-turbo (OpenAI et al., 2024b)	-	77.1	74.34	79.39	70.61	69.20	66.60	–	–	–	–	–	–	70.35	61.4

Table 2: Evaluation results of TVS-guided inference. **Bold** values indicate the better-performing result for each baseline, comparing the *with TVS* v.s. *without TVS* configurations on each benchmark.

examine LLM-assisted frameworks such as VideoTree (Wang et al., 2024c) and VideoAgent (Wang et al., 2024b), which, unlike the static encoding approaches of Video-LLMs, dynamically extract video frames based on the textual query.

TVS-guided training Following Section 3.3, we performed video instruction tuning on Video-ChatGPT (Maaz et al., 2024), Video-LLaVA (Lin et al., 2023), and ChatUniVi (Jin et al., 2023) on YouCookII-TVS. We kept the LLM backbone frozen, and only tuned their multimodal projectors.

5 Evaluation Results

We highlight the importance of the TVS task, as it addresses a core bottleneck in VideoQA by enabling models to screen for semantically aligned moments in long videos by abstracting task-level semantics and guiding modality-aware information filtering, which in turn enhances multi-modal alignment. In this section, we reveal its importance through: (1) using ground-truth TVS significantly

boosts performance (e.g., +7.3% relative gain in Table 3), showing its potential as a key supervision signal; (2) current methods fail to solve TVS effectively (Table 1), calling for dedicated solutions. These motivate our design of ReSimplifyIt as a plug-and-play approach to this essential task.

Temporal Visual Screening Our ReSimplifyIt framework outperform both baselines significantly, on every metric and subset of our YouCookII-TVS benchmark. On the average performance of the whole test set of YouCookII-TVS, ReSimplifyIt achieved an mIoU score higher than 300% of the baselines’ performance. For query reshaping, our proposed framework also achieved notably good performance. Refer to Table 1 and Table 4 for results. We provide the prompt in Appendix E.

Ablation Study of ReSimplifyIt Framework We conduct ablation study on ReSimplifyIt framework to evaluate the effectiveness its design. More details are provided in Appendix C.

Method		ActivityNetQA		YouCookII-TVS		NExT-QA		NExT-OE	
		w/tvs	w/o	w/tvs	w/o	w/tvs	w/o	w/tvs	w/o
Video-ChatGPT	w/o tvs	57.4	49.8	45.1	35.6	48.6	41.3	46.2	40.0
	w/ tvs (gt)	62.7	54.4	49.8	38.9	51.0	45.7	49.2	43.5
	w/ tvs (pred)	60.0	51.3	48.2	37.1	49.6	44.3	48.7	41.9
Video-LLaVA	w/o tvs	53.9	45.1	46.3	35.9	50.2	47.8	47.0	39.1
	w/ tvs (gt)	57.0	48.9	52.6	41.6	55.7	50.8	52.1	44.9
	w/ tvs (pred)	54.2	46.9	48.1	37.4	51.0	49.3	50.7	41.8
ChatUniVi	w/o tvs	63.2	50.0	62.4	57.4	–	–	29.0	22.5
	w/ tvs (gt)	65.6	54.4	67.8	63.8	–	–	34.3	26.0
	w/ tvs (pred)	64.9	52.6	62.8	61.6	–	–	32.1	24.0

Table 3: Evaluation results of TVS-guided training by downstream inference. On the rows, *w/o tvs*, *w/tvs (gt)*, and *w/tvs (pred)* indicates training with vanilla YouCookII-TVS dataset, training with ground truth TVS labels, and training with the predicted TVS results, respectively; on the columns, *w/o tvs* and *w/tvs* indicates the inference mode.

TVS guided inference As shown in Table 2 and Table 5, nearly all baselines - ranging from open source models of different sizes and architectures to proprietary models - saw solid absolute performance gain, suggesting that Video-LLMs largely suffer from superfluous cognition noise. The **consistent performance gains across all benchmarks** underscore the universality of this general reasoning load optimization problem across diverse tasks and scenarios. In contrast, both LLM-assisted reasoning frameworks saw smaller gain. We hypothesize that the design and implementation of these frameworks have intrinsically incorporated the TVS process, by leveraging strong reasoning ability of external LLMs or multi-turn interaction.

TVS guided training All checkpoints trained on ground truth TVS labels achieved stable performance gain over their counterparts trained on the vanilla YouCookII-TVS benchmark, which performance gain of the checkpoints trained on the predicted TVS output of ReSimplifyIt framework were weaker. We argue that TVS-guided training benefits Video-LLMs better, particularly when faced with untrimmed videos and complex text instructions requiring multi-hop reasoning. Refer to Table 3 for further details.

6 Related Work

We provide more details in appendix F.

Video-LLMs for VideoQA Video-LLMs have spurred a wave of models aimed at enhancing video understanding (Lin et al., 2023; Ma et al., 2024; Li et al., 2024b, 2023b; Liu et al., 2024b; Xu et al.,

2024a), while the sparsity and query-invariant nature of their encoding limits efficacy in capturing fine-grained spatial-temporal details.

LLM-assisted Agentic Reasoning for VideoQA Some other work proposing robust VideoQA baselines opt to explore pure-text LLM assisted frameworks or multi-agent systems for VideoQA (Wang et al., 2024c; Shang et al., 2024; Wang et al., 2024b). These methods adopt LLM-based methods to serve as a scheduler, which implicitly fulfills the TVS objective to a significant extent.

7 Conclusion

We introduce Temporal Visual Screening (TVS), a cognitively inspired task that reconstructs a VideoQA instance (V, Q) into an answer-preserving compact pair (v' , q') to purify cognition load. We realize TVS via ReSimplifyIt, a plug-and-play, model-agnostic, and multi-agent framework. To evaluate screening-centric reasoning, we synthesize YouCookII-TVS. Across models and datasets, TVS consistently strengthens training-time alignment and inference-time robustness, yielding notable accuracy gains and exposing bottlenecks through compact, auditable rationales. Our work highlights the importance of information screening and underscores the potential for significant advancements in video understanding and processing by incorporating dynamic information screening mechanisms into VideoQA agents.

Limitations

While applying the initiative of information screening brings considerable performance boost, there

are still limitations and room for improvements of our work. Firstly, as the name suggests, temporal visual screening is only applied on the temporal axis and is unable to filter redundant spatial visual information. We recognize spatial visual screening as the more difficult counterpart of our task. Secondly, although our proposed framework’s nature of being plug-and-play and adaptable to any VideoQA agent is considered as a strength, it’s end-to-end counterpart, which may mainly relate to Video-LLMs’ new visual encoder capable of inter-frame reasoning and query-adaptive frame sampling, is left unexplored. In fact, we believe that both extending plug-and-play external modules and improving end-to-end internal model structure can potentially pave the way to information screening.

References

Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhi-fang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, and 45 others. 2025a. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*.

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 others. 2025b. Qwen2.5-vl technical report. *Preprint*, arXiv:2502.13923.

Gedas Bertasius, Heng Wang, and Lorenzo Torresani. 2021. Is space-time attention all you need for video understanding? In *Proceedings of the International Conference on Machine Learning (ICML)*.

Rodney A Brooks. 1991. Intelligence without representation. *Artificial intelligence*, 47(1-3):139–159.

Guo Chen, Yicheng Liu, Yifei Huang, Yuping He, Baoqi Pei, Jilan Xu, Yali Wang, Tong Lu, and Limin Wang. 2024a. Cg-bench: Clue-grounded question answering benchmark for long video understanding. *Preprint*, arXiv:2412.12075.

Guo Chen, Yin-Dong Zheng, Jiahao Wang, Jilan Xu, Yifei Huang, Junting Pan, Yi Wang, Yali Wang, Yu Qiao, Tong Lu, and Limin Wang. 2023. Videollm: Modeling video sequence with large language models. *Preprint*, arXiv:2305.13292.

Qirui Chen, Shangzhe Di, and Weidi Xie. 2024b. Grounded multi-hop videoqa in long-form egocentric videos. *Preprint*, arXiv:2408.14469.

Shaoxiang Chen and Yu-Gang Jiang. 2019. Semantic proposal for activity localization in videos via sentence query. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence*, AAAI’19/IAAI’19/EAAI’19. AAAI Press.

Chaoyou Fu, Yuhan Dai, Yondong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, and 1 others. 2024. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*.

Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. 2017. Tall: Temporal activity localization via language query. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 5277–5285.

Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. 2017. Localizing moments in video with natural language. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 5804–5813.

Shaul Hochstein and Merav Ahissar. 2002. View from the top: Hierarchies and reverse hierarchies in the visual system. *Neuron*, 36(5):791–804.

Peng Jin, Ryuichi Takano, Caiwan Zhang, Xiaochun Cao, and Li Yuan. 2023. Chat-univi: Unified visual representation empowers large language models with image and video understanding. *arXiv preprint arXiv:2311.08046*.

Didier Le Gall. 1991. Mpeg: a video compression standard for multimedia applications. *Commun. ACM*, 34(4):46–58.

Jie Lei, Licheng Yu, Mohit Bansal, and Tamara Berg. 2018. TVQA: Localized, compositional video question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1369–1379, Brussels, Belgium. Association for Computational Linguistics.

Jie Lei, Licheng Yu, Tamara Berg, and Mohit Bansal. 2020. Tvqa+: Spatio-temporal grounding for video question answering. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8211–8225, Online. Association for Computational Linguistics.

Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024a. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.

Kunchang Li, Yanan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023a. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*.

697	Kunchang Li, Yali Wang, Yinan He, Yizhuo Li,	Richard E. Mayer and Roxana Moreno. 2003. Nine	754
698	Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo	ways to reduce cognitive load in multimedia learning.	755
699	Chen, Ping Lou, Limin Wang, and Yu Qiao. 2024b.	<i>Educational Psychologist</i> , 38(1):43–52.	756
700	Mvbench: A comprehensive multi-modal video un-		
701	derstanding benchmark . In <i>2024 IEEE/CVF Confer-</i>	Allen Newell and Herbert A. Simon. 1972. <i>Human</i>	757
702	ence on Computer Vision and Pattern Recognition	<i>Problem Solving</i> . Prentice-Hall, Englewood Cliffs,	758
703	(CVPR) , pages 22195–22206.	NJ.	759
704	Kunchang Li, Yali Wang, Yinan He, Yizhuo Li,	OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher,	760
705	Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo	Adam Perelman, Aditya Ramesh, Aidan Clark,	761
706	Chen, Ping Luo, Limin Wang, and Yu Qiao.	AJ Ostrow, Akila Welihinda, Alan Hayes, Alec	762
707	2024c. Mvbench: A comprehensive multi-	Radford, Aleksander Mađry, Alex Baker-Whitcomb,	763
708	modal video understanding benchmark . <i>Preprint</i> ,	Alex Beutel, Alex Borzunov, Alex Carney, Alex	764
709	arXiv:2311.17005.	Chow, Alex Kirillov, and 401 others. 2024a. Gpt-	765
710	Yanwei Li, Chengyao Wang, and Jiaya Jia. 2023b.	4o system card . <i>Preprint</i> , arXiv:2410.21276.	766
711	Llama-vid: An image is worth 2 tokens in large lan-		
712	guage models . In <i>European Conference on Computer</i>	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	767
713	Vision .	Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	768
714	Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and	man, Diogo Almeida, Janko Altmenschmidt, Sam Alt-	769
715	Li Yuan. 2023. Video-llava: Learning united visual	man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	770
716	representation by alignment before projection. <i>arXiv</i>	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming	771
717	<i>preprint arXiv:2311.10122</i> .	Bao, Mohammad Bavarian, Jeff Belgum, and	772
718	Yijie Lin, Jie Zhang, Zhenyu Huang, Jia Liu, Zujie Wen,	262 others. 2024b. Gpt-4 technical report . <i>Preprint</i> ,	773
719	and Xi Peng. 2024. Multi-granularity correspon-	arXiv:2303.08774.	774
720	dence learning from long-term noisy videos. In <i>Pro-</i>		
721	<i>ceedings of the International Conference on Learning</i>	Fred Paas, Alexander Renkl, and John Sweller. 2003.	775
722	<i>Representations</i> .	Cognitive load theory and instructional design:	776
723	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae	Recent developments . <i>Educational Psychologist</i> ,	777
724	Lee. 2024a. Improved baselines with visual instruc-	38(1):1–4.	778
725	tion tuning. In <i>Proceedings of the IEEE/CVF Con-</i>		
726	<i>ference on Computer Vision and Pattern Recognition</i>	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	779
727	(CVPR) , pages 26296–26306.	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	780
728	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan	try, Amanda Askell, Pamela Mishkin, Jack Clark,	781
729	Zhang, Sheng Shen, and Yong Jae Lee. 2024b. Llava-	Gretchen Krueger, and Ilya Sutskever. 2021. Learn-	782
730	next: Improved reasoning, ocr, and world knowledge .	ing transferable visual models from natural language	783
731	Ruyang Liu, Chen Li, Yixiao Ge, Thomas H. Li, Ying	supervision . <i>Preprint</i> , arXiv:2103.00020.	784
732	Shan, and Ge Li. 2024c. Bt-adapter: Video conversa-		
733	tion is feasible without video instruction tuning . In	Chuyi Shang, Amos You, Sanjay Subramanian, Trevor	785
734	<i>2024 IEEE/CVF Conference on Computer Vision and</i>	Darrell, and Roei Herzig. 2024. TravelER: A modu-	786
735	<i>Pattern Recognition (CVPR)</i> , pages 13658–13667.	lar multi-LMM agent framework for video question-	787
736	Fan Ma, Xiaojie Jin, Heng Wang, Yuchen Xian, Jiashi	answering . In <i>Proceedings of the 2024 Conference</i>	788
737	Feng, and Yi Yang. 2024. Vista-llama: Reducing	<i>on Empirical Methods in Natural Language Process-</i>	789
738	hallucination in video language models via equal	<i>ing</i> , pages 9740–9766, Miami, Florida, USA. Associ-	790
739	distance to visual tokens. In <i>Proceedings of the</i>	ation for Computational Linguistics.	791
740	<i>IEEE/CVF Conference on Computer Vision and Pat-</i>		
741	<i>tern Recognition (CVPR)</i> , pages 13151–13160.	John Sweller. 1988. Cognitive load during problem	792
742	Muhammad Maaz, Hanoona Rasheed, Salman Khan,	solving: Effects on learning. <i>Cognitive Science</i> ,	793
743	and Fahad Khan. 2024. Video-ChatGPT: Towards	12(2):257–285.	794
744	detailed video understanding via large vision and	Anne M. Treisman and Garry Gelade. 1980. A feature-	795
745	language models . In <i>Proceedings of the 62nd An-</i>	integration theory of attention . <i>Cognitive Psychology</i> ,	796
746	<i>nuual Meeting of the Association for Computational</i>	12(1):97–136.	797
747	<i>Linguistics (Volume 1: Long Papers)</i> , pages 12585–	Weihan Wang, Zehai He, Wenyi Hong, Yean Cheng,	798
748	12602, Bangkok, Thailand. Association for Compu-	Xiaohan Zhang, Ji Qi, Shiyu Huang, Bin Xu, Yuxiao	799
749	tational Linguistics.	Dong, Ming Ding, and Jie Tang. 2024a. Lvbench:	800
750	Karttikeya Mangalam, Raiymbek Akshulakov, and Ji-	An extreme long video understanding benchmark .	801
751	tendra Malik. 2023. Egoschema: A diagnostic bench-	<i>Preprint</i> , arXiv:2406.08035.	802
752	mark for very long-form video language understand-		
753	ing . <i>Preprint</i> , arXiv:2308.09126.	Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu,	803
		Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin	804
		Jing, Shenglong Ye, Jie Shao, Zhaokai Wang, Zhe	805
		Chen, Hongjie Zhang, Ganlin Yang, Haomin Wang,	806
		Qi Wei, Jinhui Yin, Wenhao Li, Erfei Cui, and 56	807
		others. 2025. InternV13.5: Advancing open-source	808
		multimodal models in versatility, reasoning, and effi-	809
		ciency . <i>Preprint</i> , arXiv:2508.18265.	810

811	Xiaohan Wang, Yuhui Zhang, Orr Zohar, and Serena	Jeffrey M. Zacks, Nicole K. Speer, Khená M. Swal-	864
812	Yeung-Levy. 2024b. Videoagent: Long-form video	low, Todd S. Braver, and Jeremy R. Reynolds. 2007.	865
813	understanding with large language model as agent.	Event perception: A mind/brain perspective . <i>Psycho-</i>	866
814	<i>European Conference on Computer Vision (ECCV)</i> .	<i>logical Bulletin</i> , 133(2):273–293.	867
815	Ziyang Wang, Shoubin Yu, Elias Stengel-Eskin, Jae-	Hao Zhang, Aixin Sun, Wei Jing, and Joey Tianyi Zhou.	868
816	hong Yoon, Feng Cheng, Gedas Bertasius, and Mo-	2020. Span-based localizing network for natural lan-	869
817	hit Bansal. 2024c. Videotree: Adaptive tree-based	guage video localization . In <i>Proceedings of the 58th</i>	870
818	video representation for llm reasoning on long videos .	<i>Annual Meeting of the Association for Computational</i>	871
819	<i>Preprint</i> , arXiv:2405.19209.	<i>Linguistics</i> , pages 6543–6554, Online. Association	872
820	Jeremy M. Wolfe. 1994. Guided search 2.0: A revised	for Computational Linguistics.	873
821	model of visual search. <i>Psychonomic Bulletin &</i>	Yuji Zhang, Sha Li, Cheng Qian, Jiateng Liu, Pengfei	874
822	<i>Review</i> , 1(2):202–238.	Yu, Chi Han, Yi R Fung, Kathleen McKeown,	875
823	Penghao Wu and Saining Xie. 2024. V?: Guided visual	Chengxiang Zhai, Manling Li, and 1 others. 2025.	876
824	search as a core mechanism in multimodal llms. In	The law of knowledge overshadowing: Towards un-	877
825	<i>Proceedings of the IEEE/CVF Conference on Com-</i>	derstanding, predicting, and preventing llm halluci-	878
826	<i>puter Vision and Pattern Recognition</i> , pages 13084–	nation. <i>arXiv preprint arXiv:2502.16143</i> .	879
827	13094.	Minghang Zheng, Xinhao Cai, Qingchao Chen, Yuxin	880
828	Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng	Peng, and Yang Liu. 2024a. Training-free video tem-	881
829	Chua. 2021. Next-qa: Next phase of question-	poral grounding using large-scale pre-trained models .	882
830	answering to explaining temporal actions. In <i>Pro-</i>	<i>Preprint</i> , arXiv:2408.16219.	883
831	<i>ceedings of the IEEE/CVF Conference on Computer</i>	Minghang Zheng, Xinhao Cai, Qingchao Chen, Yuxin	884
832	<i>Vision and Pattern Recognition (CVPR)</i> , pages 9777–	Peng, and Yang Liu. 2024b. Training-free video tem-	885
833	9786.	poral grounding using large-scale pre-trained models .	886
834	Junbin Xiao, Angela Yao, Yicong Li, and Tat-Seng	<i>Preprint</i> , arXiv:2408.16219.	887
835	Chua. 2024. Can i trust your answer? visually	Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao,	888
836	grounded video question answering. In <i>Proceedings</i>	Xi Yang, Yongping Xiong, Bo Zhang, Tiejun Huang,	889
837	<i>of the IEEE/CVF Conference on Computer Vision</i>	and Zheng Liu. 2024a. Mlvu: A comprehensive	890
838	<i>and Pattern Recognition</i> , pages 13204–13214.	benchmark for multi-task long video understanding.	891
839	Huijuan Xu, Kun He, Bryan A. Plummer, Leonid Sigal,	<i>arXiv preprint arXiv:2406.04264</i> .	892
840	Stan Sclaroff, and Kate Saenko. 2019. Multilevel lan-	Luowei Zhou, Chenliang Xu, and Jason J Corso. 2018.	893
841	guage and vision integration for text-to-clip retrieval.	Towards automatic learning of procedures from web	894
842	In <i>AAAI</i> .	instructional videos . In <i>AAAI Conference on Artifi-</i>	895
843	Lin Xu, Yilin Zhao, Daquan Zhou, Zhijie Lin,	<i>cial Intelligence</i> , pages 7590–7598.	896
844	See Kiong Ng, and Jiashi Feng. 2024a. Pllava :	Xingyi Zhou, Anurag Arnab, Shyamal Buch, Shen Yan,	897
845	Parameter-free llava extension from images to videos	Austin Myers, Xuehan Xiong, Arsha Nagrani, and	898
846	for video dense captioning . <i>ArXiv</i> , abs/2404.16994.	Cordelia Schmid. 2024b. Streaming dense video cap-	899
847	Yifang Xu, Yunzhuo Sun, Zien Xie, Benxiang Zhai,	tioning. In <i>Proceedings of the IEEE/CVF Conference</i>	900
848	and Sidan Du. 2024b. Vtg-gpt: Tuning-free zero-	<i>on Computer Vision and Pattern Recognition (CVPR)</i> ,	901
849	shot video temporal grounding with gpt. <i>Applied</i>	pages 18243–18252.	902
850	<i>Sciences</i> , 14(5):1894.	A Viewer Implementation Details	903
851	Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yuet-	Keyframe Extraction and Captioning. We uti-	904
852	ing Zhuang, and Dacheng Tao. 2019. Activitynet-qa:	lize a two-stage keyframe extraction process com-	905
853	A dataset for understanding complex web videos via	bined with frame captioning. In the first stage,	906
854	question answering. In <i>AAAI</i> , pages 9127–9134.	we implement the MPEG-4 compression technique	907
855	Yitian Yuan, Tao Mei, and Wenwu Zhu. 2019. To find	(Le Gall, 1991) to extract all I-frames as candidate	908
856	where you talk: temporal sentence localization in	keyframes. I-frames typically contain rich visual	909
857	video with attention based location regression . In	content and clarity, or represent scene transitions.	910
858	<i>Proceedings of the Thirty-Third AAAI Conference</i>	In the second stage, we apply a modified Isodata	911
859	<i>on Artificial Intelligence and Thirty-First Innovative</i>	clustering algorithm to the visual features of these I-	912
860	<i>Applications of Artificial Intelligence Conference and</i>	frames, selecting cluster centers as final keyframes.	913
861	<i>Ninth AAAI Symposium on Educational Advances in</i>	This algorithm adaptively determines the number	914
862	<i>Artificial Intelligence</i> , AAAI’19/IAAI’19/EAAI’19.	of clusters, unlike KNN-based clustering (Wang	915
863	AAAI Press.		

Method	Temporal Relational				Timepoint Indexed				Multifaceted Integrative				Average			
	mIoU	Pre.	Cov.	F1	mIoU	Pre.	Cov.	F1	mIoU	Pre.	Cov.	F1	mIoU	Pre.	Cov.	F1
VTG-GPT (Xu et al., 2024b)	0.17	0.26	0.32	0.29	0.15	0.22	0.30	0.26	0.16	0.23	0.31	0.24	0.16	0.27	0.31	0.27
Zheng et al. (2024a)	0.11	0.18	0.21	0.19	0.04	0.08	0.09	0.08	0.06	0.10	0.11	0.10	0.07	0.12	0.13	0.12
ReSimplifyIt (Ours)	0.23	0.36	0.39	0.37	0.98	1.0	0.98	0.99	0.47	0.60	0.69	0.64	0.56	0.65	0.69	0.67

(a) more results on stage 1 video output

Table 4: More results on Stage-1 evaluation on YouCookII-TVS dataset. 'Pre.' and 'Cov.' stands for 'precision' and 'coverage'.

et al., 2024c; Zhou et al., 2024b), which requires a pre-defined number of clusters or external supervision. Our approach ensures generalizability and robustness for diverse video types.

After clustering, we utilize an off-the-shelf vision-language model (VLM) to generate frame captions for the selected keyframes. The overall process is denoted by:

$$C = \text{prep}(start, end)$$

where C is the set of frame captions between the given timestamps.

Scanning. Given a timestamp range $(start, end)$, the Viewer calls an external LLM with C to produce an overview caption. It may query additional frames via a captioning tool to refine its summary, mimicking how users drag to specific timestamps for clarification. The process is:

$$Cap = \text{Scanner}(\text{prep}(start, end), t)$$

where t is the task prompt.

Localizing. To identify the timestamp range matching a textual description, we propose a lightweight three-stage search process:

Stage 1: We feed an external LLM the set of keyframe captions to let it acquire an overall capture of the video content. At the same time, we instruct the LLM to output five most possible timestamps that depicts the texture query. The LLM is able to call frame caption tool to acquire extra captions at arbitrary timestamps, before it's confident enough to output the answer. Formally:

$$P = \text{stage}_1(\text{prep}(0, d), e, t)$$

where e are extra captions and t is the prompt. P gives the resulting list of five timestamps which depicts the language query the best.

Stage 2: We initialize the conversation and feed it the frame captions at these five timestamps, while instructing it to pick the one timestamp that best depicts the language query, out of the five candidates. Also, the LLM is able to call frame caption tool

to acquire extra captions at arbitrary timestamps before it's confident enough to output the answer.

$$t_{best} = \text{stage}_2(P, e, t)$$

where P is the output of the previous step, and t is the task prompt.

Stage 3: Last, we initialize the conversation again, and instruct the external LLM to expand the single timestamp t_{best} from the last step to a timestamp range. The frame caption at t_{best} is provided, and the LLM is still able to acquire more captions by tool calling to confirm the boundary. Considering the difficulty in dealing with dynamic transitions of events, as explored by some previous work (Zheng et al., 2024b), we simply apply a hard value of 5 seconds on the output. Formally:

$$(t'_{start}, t'_{end}) = \text{stage}_3(t_{best}, e, t)$$

$$(t_{start}, t_{end}) = (t'_{start} - 5, t'_{end} + 5)$$

This three-stage process balances precision and efficiency by minimizing frame access while ensuring robust temporal grounding. For visual clarity of the Isodata clustering, refer to Algorithm 1 for details.

B Details on the YouCookII-TVS benchmark

We constructed a synthetic question-answering dataset named **YouCookII-TVS**, based on the YouCookII dataset, to support fine-grained temporal and semantic understanding of cooking videos.

B.1 Source Data Preparation

The original YouCookII dataset (Zhou et al., 2018) contains temporally annotated instructional videos. Each annotation includes a segment $[s_i, e_i]$ representing start and end times (in seconds), along with a natural language description of the cooking step.

To ensure video consistency and avoid duplications, we have verified that each video clip name is unique across the dataset source.

Method	ActivityNetQA	
	w/tvs	w/o
Video-ChatGPT (Maaz et al., 2024)	58.5	50.5
Video-LLaVA (Lin et al., 2023)	61.2	52.1
ChatUniVi (Jin et al., 2023)	63.2	52.0
LLaVA-NExT (Liu et al., 2024b)	67.8	62.5
VideoAgent (Wang et al., 2024b)	61.5	60.2
VideoTree (Wang et al., 2024c)	59.0	63.6

Method	YouCookII-TVS							
	TRR.		TIR.		MIR.		Avg.	
	w/tvs	w/o	w/tvs	w/o	w/tvs	w/o	w/tvs	w/o
Video-ChatGPT (Maaz et al., 2024)	31.7	30.5	45.2	26.5	39.8	29.73	38.9	28.91
Video-LLaVA (Lin et al., 2023)	39.3	37.4	46.8	28.2	43.5	31.6	43.2	32.4
ChatUniVi (Jin et al., 2023)	54.4	39.8	57.4	52.8	58.7	49.6	56.5	47.4
LLaVA-NExT (Liu et al., 2024b)	46.5	44.1	65.9	33.4	62.2	38.3	58.2	48.6
VideoAgent (Wang et al., 2024b)	30.2	48.7	46.1	55.2	38.9	57.8	38.4	53.9
VideoTree (Wang et al., 2024c)	59.3	64.0	59.2	72.2	52.8	72.0	57.1	69.4

Method	NExT-QA							
	Tem.		Cau.		Des.		Avg.	
	w/tvs	w/o	w/tvs	w/o	w/tvs	w/o	w/tvs	w/o
Video-ChatGPT (Maaz et al., 2024)	45.5	23.7	45.5	56.7	45.6	43.0	45.5	44.2
Video-LLaVA (Lin et al., 2023)	42.8	42.8	50.4	48.9	55.1	44.5	52.2	46.3
ChatUniVi (Jin et al., 2023) -	-	-	-	-	-	-	5	28
LLaVA-NExT (Liu et al., 2024b)	57.5	52.3	62.0	59.0	61.4	56.5	60.5	56.6
VideoAgent (Wang et al., 2024b)	49.2	47.3	43.6	41.9	51.7	51.1	47.0	45.0
VideoTree (Wang et al., 2024c)	62.4	59.9	57.7	66.1	66.8	72.3	60.0	65.2

Method	NExT-OE							
	Tem.		Des.		Cau.		Avg.	
	w/tvs	w/o	w/tvs	w/o	w/tvs	w/o	w/tvs	w/o
Video-ChatGPT (Maaz et al., 2024)	49.6	40.8	51.2	43.2	46.9	38.6	49.8	41.5
Video-LLaVA (Lin et al., 2023)	37.5	37.5	53.2	46.7	47.6	43.4	48.0	43.3
ChatUniVi (Jin et al., 2023)	30.8	30.8	36.9	23.1	34.0	26.4	34.5	25.9
LLaVA-NExT (Liu et al., 2024b)	41.4	39.11	59.7	46.0	50.1	44.6	52.3	43.7
VideoAgent (Wang et al., 2024b)	44.1	53.4	48.8	46.0	50.4	52.6	47.8	49.6
VideoTree (Wang et al., 2024c)	52.1	61.1	58.0	60.7	64.4	65.6	57.7	61.9

Table 5: Evaluation results of TVS-guided inference on four benchmark: ActivityNetQA, YouCookII-TVS, NExT-QA, and NExT-OE, with the sub-categories presented in each benchmark.

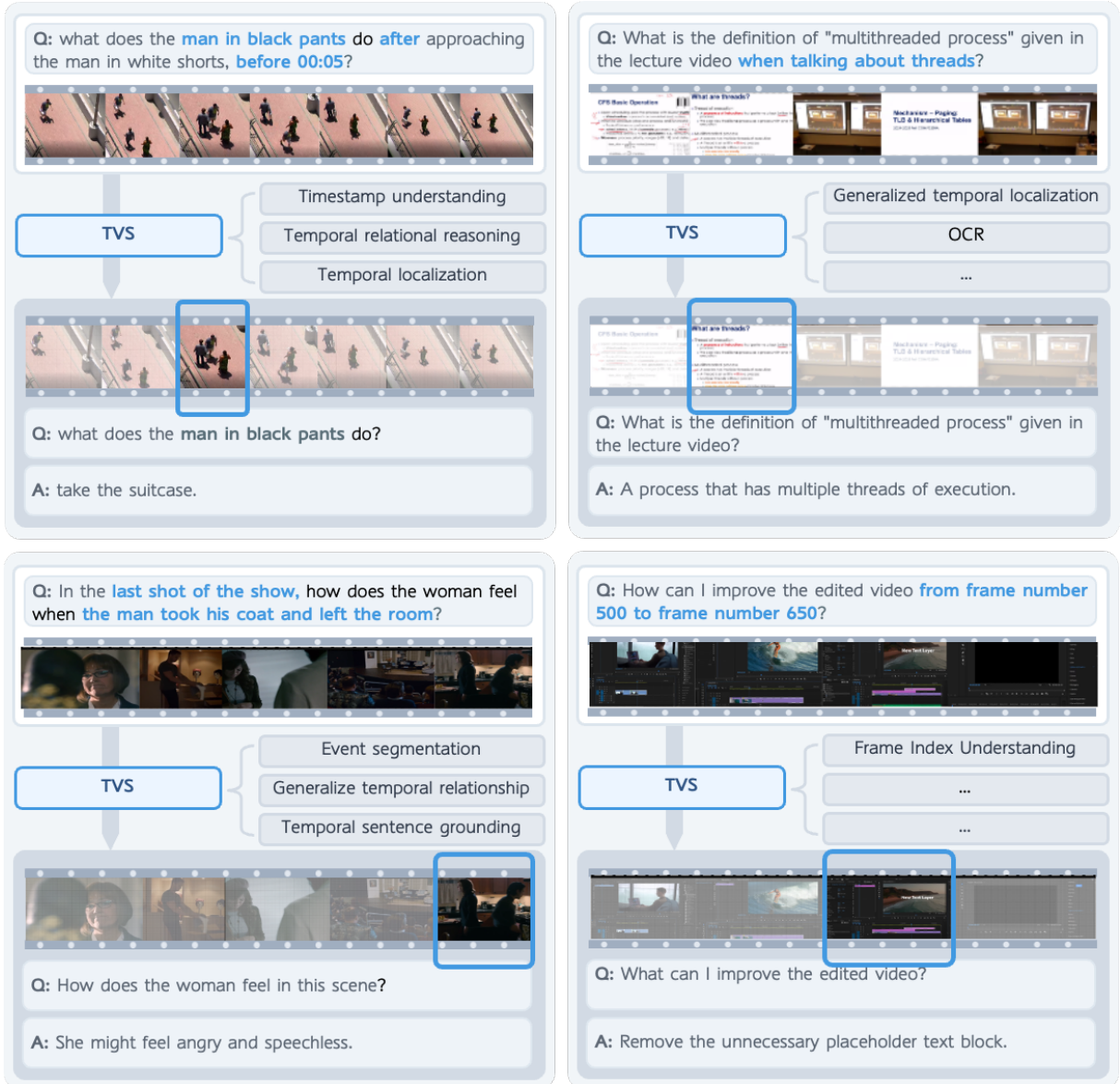


Figure 3: More examples of the TVS task.

B.2 Annotation Grouping via Temporal Connectivity

We define a temporal connectivity criterion to group sequential cooking steps into higher-level event triplets. Given two segments $[s_1, e_1]$ and $[s_2, e_2]$, we define their overlap ratio as:

$$\text{overlap_ratio} = \frac{|\min(e_1, e_2) - \max(s_1, s_2)|}{\max(e_1, e_2) - \min(s_1, s_2)}$$

Two segments are considered *connectable* if:

$$s_2 > s_1, \quad e_2 > e_1, \quad \text{and} \quad \text{overlap_ratio} \leq \theta$$

We set $\theta = 0.1$ in our experiments. Using this rule, we perform a greedy grouping of annotations into connected segments, and extract all valid length-3 subsequences (triplets) from each group.

B.3 Triplet-Based Question Generation

Each triplet $T = \{t_1, t_2, t_3\}$ consists of three temporally ordered steps. For each T , we generate nine different types of question-answer pairs by instantiating predefined templates. The question types are categorized into three groups:

Temporal Relational Reasoning (TRR)

- trr1: What is the cooking step after of $[description]$?
- trr2: What is the cooking step before $[description]$?
- trr3: What is the cooking step between $[description]$ and $[description]$?

Algorithm 1 ISODATA Clustering

Require: Data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, Candidate frames $\mathbf{F} \in \mathbb{R}^{n \times c \times w \times h}$, frame feature extractor $Enc()$, initial cluster count k , max iterations T , minimum intra-cluster similarity θ_{split} , maximum inter-cluster similarity θ_{merge} , max clusters k_{max} , min clusters k_{min} , max center shift δ_{max} , min elements per cluster n_{min}

Ensure: Cluster assignments \mathbf{C} , Final cluster centers \mathbf{M}

```
1:  $\mathbf{X} \leftarrow Enc(\mathbf{F})$ 
2:  $\mathbf{X}_{norm} \leftarrow \text{normalize}(\mathbf{X}, \ell_2)$ 
3:  $\mathbf{M} \leftarrow \text{random\_sample}(\mathbf{X}, k)$ 
4:  $t \leftarrow 0$ 
5: repeat
6:   Compute cosine similarity matrix  $\mathbf{S} = \mathbf{X}_{norm} \mathbf{M}_{norm}^T$ 
7:    $\mathbf{C} \leftarrow \arg \max(\mathbf{S}, \text{axis} = 1)$  ▷ Assign points to nearest clusters
8:   for each cluster  $i$  do
9:     Update center:  $\mathbf{m}_i \leftarrow \text{mean}(\mathbf{X}[\mathbf{C} == i])$ 
10:  end for
11:  if any cluster size  $< n_{min}$  then
12:    Merge smallest cluster with nearest neighbor ▷ Minimum elements enforcement
13:  end if
14:  for each cluster  $i$  do
15:    if intra-cluster similarity( $\mathbf{X}[\mathbf{C} == i], \mathbf{m}_i$ )  $< \theta_{split}$  and  $k < k_{max}$  then
16:      Split cluster  $i$  into two new clusters ▷ Splitting phase
17:       $k \leftarrow k + 1$ 
18:    end if
19:  end for
20:  for all cluster pairs  $(i, j)$  do
21:    if inter-cluster similarity( $\mathbf{m}_i, \mathbf{m}_j$ )  $> \theta_{merge}$  and  $k > k_{min}$  then
22:      Merge clusters  $i$  and  $j$  ▷ Merging phase
23:       $k \leftarrow k - 1$ 
24:    end if
25:  end for
26:  Compute center shifts  $\Delta \mathbf{M} = \|\mathbf{M}_{new} - \mathbf{M}_{old}\|$ 
27:   $t \leftarrow t + 1$ 
28: until  $t \geq T$  or  $\max(\Delta \mathbf{M}) < \delta_{max}$ 
29: return  $\mathbf{C}, \mathbf{M}$ 
```

Algorithm 2 *ReSimplifyIt*

Require: Video v , question q

Ensure: Screened video V' , screened question q'

```
1:  $V\_copy, q\_copy \leftarrow V, q$ 
2:  $success\_history, failure\_history \leftarrow SuccessHistory(), FailureHistory()$ 
3: while true do
4:    $launcher, validator, viewer \leftarrow Launcher(), Validator(), Viewer()$ 
5:    $decision, q', trimming\_instruction \leftarrow launcher(q\_copy, success\_history, failure\_history)$ 
6:   if  $decision == "proceed"$  then
7:      $judgement, request, result, reason \leftarrow validator(q\_copy, q', trimming\_instruction)$ 
8:     while  $judgement == "view"$  do
9:        $response \leftarrow viewer(V', request)$ 
10:       $validator.read\_response(response)$ 
11:    end while
12:    if  $judgement == "succeeded"$  then
13:       $V' \leftarrow result$ 
14:       $success\_history.append([q\_copy, q', trimming\_instruction])$ 
15:       $failure\_history.empty()$ 
16:       $V\_copy, q\_copy \leftarrow V', q'$ 
17:    else
18:       $failure\_history.append([q\_copy, q', trimming\_instruction])$ 
19:    end if
20:  else
21:    return  $V\_copy, q\_copy$ 
22:  end if
23: end while
```

1019	Timepoint Indexed Reasoning (TIR)	B.5 Dataset Splitting	1059
1020	• tir1: What is the step between timestamps	To support evaluation, we partition the dataset into	1060
1021	s_2 and e_2 ?	training, validation, and test splits. For each ques-	1061
1022	• tir2: What is the step between frame indices	tion type, we allocated:	1062
1023	$f_{s_2} = s_2 \cdot r$ and $f_{e_2} = e_2 \cdot r$?	train: 1926, val: 270, test: 558	1063
1024	• tir3: What step appears within $f_d = (e_2 -$	This roughly follows 7:1:2.	1064
1025	$s_2) \cdot r$ frames after s_2 seconds?	More details about the dataset statistics can be	1065
1026	Here, r denotes the video frame rate.	found in Figure 4.	1066
1027	Multifaceted Integrative Reasoning (MIR)	C Ablation Study	1067
1028	• mir1: What is the first step after timestamp	C.1 Experiment Settings	1068
1029	s_2 ?	We first ablate the modular design—responsible	1069
1030	• mir2: What is the last step before timestamp	for facilitating structured reasoning—by proposing	1070
1031	e_2 ?	the ReSimplifyIt-simple variant. In this setting,	1071
1032	• mir3: Within s_1 and e_3 , what is (are) the	all underlying interfaces for handling textual and	1072
1033	cooking step(s) apart from $[description]$ and	visual inputs, such as ISODATA-clustering (Algo-	1073
1034	$[description]$?	rithm 1) and frame-caption querying, are preserved.	1074
1035	Template instantiation is performed by replacing	However, the entire reasoning pipeline is collapsed	1075
1036	placeholders with actual sentences and timestamps	into a single, universal agent. Specifically, we ini-	1076
1037	(framestamps) from the triplet.	tialize an external agentic LLM with task descrip-	1077
1038	B.4 Data Structuring and Metadata	tions and operational instructions, and expose the	1078
1039	Each generated data point is stored with the follow-	mentioned interfaces either through conversa-	1079
1040	ing fields:	tional context or tool invocation. This unified	1080
1041	• vid_name, vid_fname: Video ID and file-	agent is then responsible for all reasoning pro-	1081
1042	name.	cedures—functionally covering the roles of the	1082
1043	• vid_duration, vid_frame_rate: Metadata	Launcher, Validator, and Viewer modules, as well	1083
1044	from video parsing.	as memory tracking—in the original ReSimplifyIt	1084
1045	• type: One of the nine QA types.	framework, and ultimately produces the final out-	1085
1046	• question: Instantiated natural language	put.	1086
1047	query.	Then, we further ablate the video access com-	1087
1048	• answer: Corresponding ground-truth step de-	pletely, i.e. No access to the video content is pro-	1088
1049	scription, serving as the ground-truth label for	vided throughout the <i>entire</i> reasoning process of the	1089
1050	the VideoQA task.	external agent. As the multi-turn interaction (feed-	1090
1051	• gt_timestamp: Temporal segment(s) serving	back) are all essentially from the reference to video	1091
1052	as the ground-truth label for our TVS task on	information to refine text input, disabling video	1092
1053	video trimming.	access also renders multi-turn interactions unneces-	1093
1054	• gt_rewritten_query: Natural language	sary. To reflect this, we introduce the ReSimplifyIt-	1094
1055	query serving as the ground-truth label for	blind variant, in which a tool-calling LLM gener-	1095
1056	our TVS task on query re-writing.	ates a rewritten query and a fixed sequence of tool	1096
1057	We generated a total of $N = 2754$ QA samples,	invocations within a single-turn conversation. This	1097
1058	covering all types evenly.	sequence is subsequently executed by a separate	1098
		executor module to produce the video output.	1099
		Refer to Figure 5 for an illustrations of these	1100
		frameworks.	1101
		C.2 Evaluation Results	1102
		Evaluation results are presented in Table 6. Al-	1103
		though the removal of modular design yields com-	1104
		parable video outputs, the modular architecture still	1105

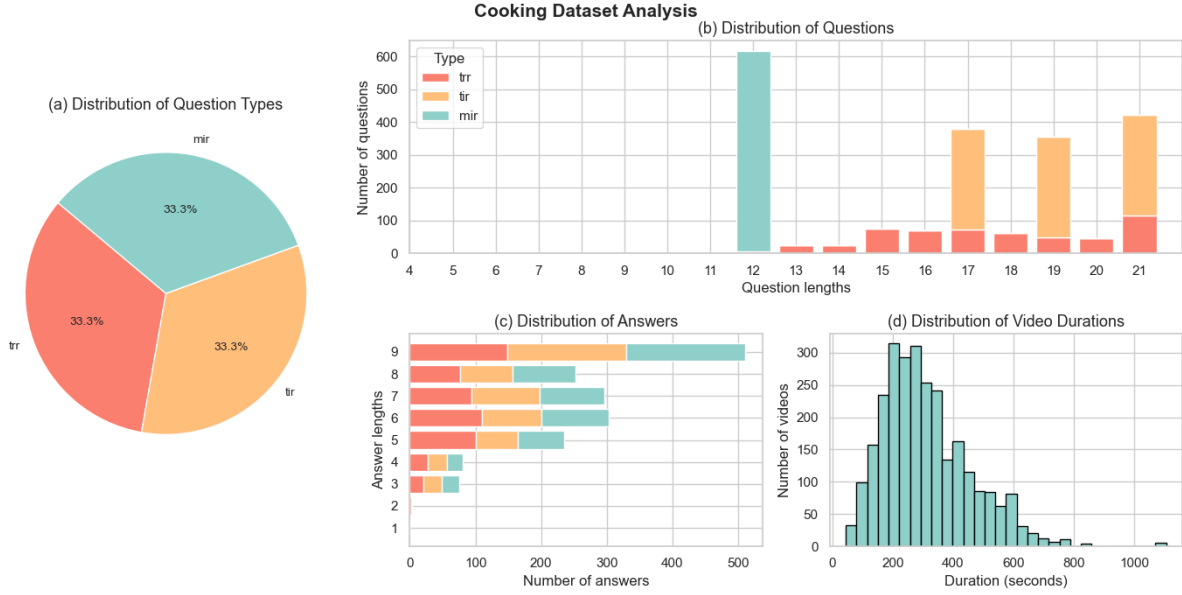


Figure 4: YouCookII-TVS statistics.

Method	Temporal Relational		Timepoint Indexed		Multifaceted Integrative		Average	
	mIoU	F1	mIoU	F1	mIoU	F1	mIoU	F1
ReSimplifyIt (Ours)	<u>0.23</u>	0.37	0.98	0.99	0.47	0.64	0.56	0.67
ReSimplifyIt-simple (Ours)	0.24	0.37	0.98	0.99	<u>0.42</u>	<u>0.57</u>	<u>0.55</u>	<u>0.64</u>
ReSimplifyIt-blind (Ours)	0.12	0.20	0.97	0.99	0.38	0.55	<u>0.49</u>	0.58

(a) Results on video output.

Method	Temporal Relational	Timepoint Indexed	Multifaceted Integrative	Average
ReSimplifyIt (Ours)	66.8	78.5	72.8	72.7
ReSimplifyIt-simple (Ours)	66.2	81.9	68.1	72.0
ReSimplifyIt-blind (Ours)	65.0	80.5	70.7	<u>72.1</u>

(b) Results on query rewriting.

Table 6: Ablation studies on our ReSimplifyIt framework. ReSimplifyIt-simple and ReSimplifyIt-blind represent ablations on modular design and feedback from video access, respectively.

demonstrates advantages—particularly on the *Multifaceted Integrative* subset, which involves more complex multi-hop reasoning, demonstrating the effectiveness of modular design and structured reasoning in more complex and intricate reasoning scenario. In comparison, the ablation of video access brings notable performance drop, underscoring the critical role of cooperative, feedback-driven multi-turn interaction. This aligns with the interdependency of the vision and text modality, which lies at the core of our TVS task formulation, reflecting the essence of the task initiative.

C.3 Full list of tools for ReSimplifyIt-blind framework

```
1. get_duration():
    Return the duration of the video as a
    floating point value.
```

```
2. get_resolution():
    Return the resolution of the video, as a
    tuple.

3. get_total_frame_num():
    Return total number of the frames of the
    video, as an integer.

4. grounding_select(obj_name,
    concerned_indices_input):
    Return, in the form of a list of integers,
    the indices of all frames
    containing the object given by obj_name,
    after taking the intersection
    of indices provided by
    concerned_indices_input. If None is passed,
    selects all frames.

5. indices_list_intersect(list1, list2):
    Return the intersection of two lists of
    indices.
```

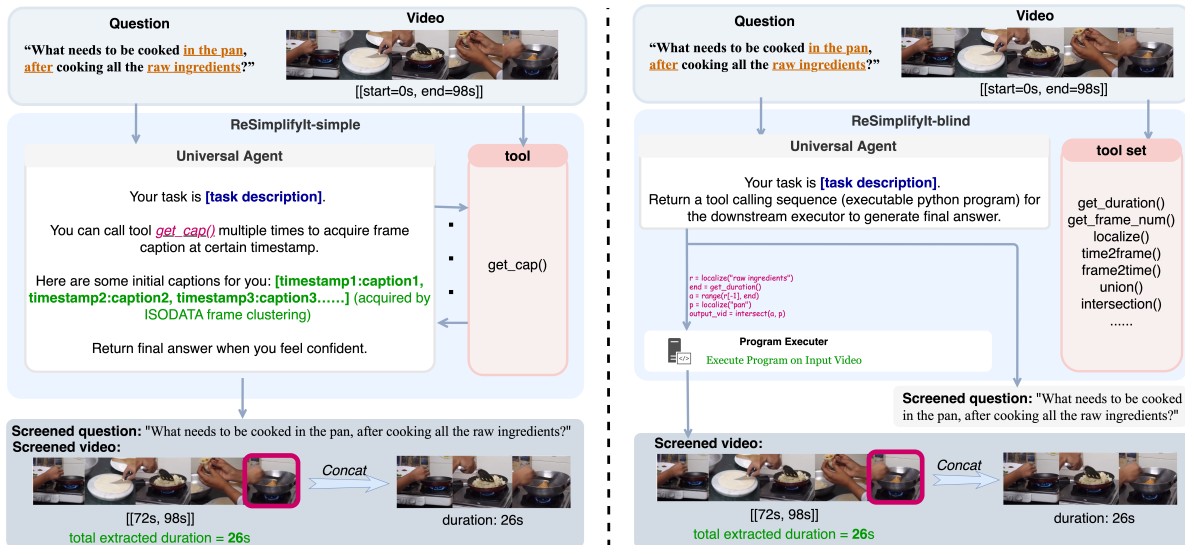


Figure 5: Workflow of *ReSimplifyIt-simple* (left) and *ReSimplifyIt-blind* (right).

- 1147 6. `indices_list_union(list1, list2)`:
1148 Return the union of two lists of indices.
- 1149 7. `indices_concat_and_fill(list1, list2)`:
1150 Return the sorted union of list1 and list2,
1151 then fill in missing
1152 values to make the sequence continuous.
- 1153 8. `indices_concat(list1, list2)`:
1154 Return the concatenation of the two lists.
- 1155 9. `timestamp_to_single_index(timestamp)`:
1156 Return the frame index corresponding to the
1157 given timestamp (in seconds).
- 1158 10. `single_timestamp_to_index_range(timestamp)`:
1159 Return indices of 60 consecutive frames
1160 centered at the timestamp.
- 1161 11. `range_timestamp_to_index_range(start, end)`:
1162 Return all frame indices between the start
1163 and end timestamps.

1170 D Implementation details of ReSimplifyIt 1171 framework

1172 The Launcher module is based on single-turn con-
1173 versation with GPT-4o. The Validator and Viewer
1174 module, including the scanner and localizer, are
1175 primarily based on multi-turn conversation with
1176 GPT-4o. For ISODATA frame clustering, we use
1177 CLIP (Radford et al., 2021) to obtain the visual
1178 features of the I-frames. We adopt the off-the-shelf
1179 LLaVA-1.5 (Liu et al., 2024a) for frame captioning.

1180 E Prompts

1181 E.1 Prompt used for evaluation of query 1182 rewriting of TVS

1183 You are a helpful assistant to evaluate the
1184 quality of an output of a special type of
1185 sentence compression, which sentence is in
1186 the form of a question.
1187
1188 You will be given an output of such compression
1189 process and the ground truth answer, and the
1190 original input question sentence before the
1191 compression. Please evaluate the output
1192 based on the following three criterias:
1193
1194 1. Relevance: to minimize unwanted
1195 information
1196 - in this criteria, a candidate output gets
1197 full mark if it doesn't contain any
1198 information (phrases, concepts, etc.) out of
1199 the scope of the original input question.
1200 - the more information it contains that is
1201 not mentioned in the original input question
1202 , the more marks are deducted.
1203
1204 2. Simplicity: to minimize tangential
1205 information
1206 - in this criteria, a candidate output gets
1207 full mark if doesn't contain any information
1208 (phrases, concepts, etc.) that is included
1209 in the original input question but not
1210 included in the ground truth compressed
1211 question. In other words, whether the
1212 question sentence is fully compressed.
1213 - the more information it contains that is
1214 included in the original input question but
1215 not included in the ground truth question,
1216 the more marks will be deducted.
1217
1218 3. Completeness: to minimize over-
1219 compression
1220 - in this criteria, a candidate output gets
1221 full mark if it contains all information
1222 included in the ground truth compressed
1223 question.
1224 - the more information contained in the
1225 ground truth output question is found
1226 missing in the output compressed question,
1227

1228 the more marks will be deducted.
 1229
 1230 Here is the original input question: [
 1231 original_question_flag], and
 1232 here is the ground truth compressed question:
 1233 [ground_truth_screened_question_flag], and
 1234 here is the output compressed question: [
 1235 output_screened_question_flag].
 1236
 1237 Please rate the output compressed question
 1238 on a scale from 0 to 100, with 0 being the
 1239 worst and 100 being the best (full mark).
 1240 Now, rate the quality of the output
 1241 compressed question based on all the
 1242 information above.
 1243 Return your answer in this json format: {"
 1244 score": [your score, from 0 to 100]}.

1246 **E.2 Prompts used for ReSimplifyIt**
 1247 **framework**

1248 Prompt for the Launcher module:

1249 Given a video question answering case, i.e. a
 1250 video, a question, and an answer, sometimes
 1251 it is possible to cut the video by only
 1252 keeping a sub-clip or several sub-clips (
 1253 concatenate if so) to be the video output
 1254 and simultaneously modify the text question
 1255 to be the paired text output while keeping
 1256 the answer consistent and unchanged, so that
 1257 the answer is still completely compatible
 1258 to the modified question and the obtained
 1259 video clip.
 1260 We define this action as "screening". By doing
 1261 screening, the video becomes shorter, and
 1262 therefore introduces lower cost to the
 1263 downstream video question answering model.
 1264 When modifying the question, note that the
 1265 downstream video question answering model
 1266 will only be seeing the sub-clip and think
 1267 that the shown sub-clip is the whole video,
 1268 so make sure the modified question is
 1269 perfectly paired and aligned with and
 1270 adapted to the modification.
 1271
 1272 [In context examples flag 1]
 1273
 1274 Sometimes, such screening can be repeatedly
 1275 applied sequentially by several rounds,
 1276 where the input in each round is the output
 1277 video clip and the output modified question
 1278 of its preceeding round. As mentioned, the
 1279 answer consistency should be always ensured
 1280 throughout the whole process, and is always
 1281 completely compatible to the resulting video
 1282 clip and modified question of every round.
 1283 If succeeded, each round would make the
 1284 video shorter and making the situation
 1285 closer to the optimal case.
 1286
 1287 [In context examples flag 2]
 1288
 1289 Your duty is to help complete the screening. Due
 1290 to some limit, you are only provided with
 1291 the text question but not the video input.
 1292 Your task is to initiate an immediate plan for
 1293 the screening operation, including a natural

1296 language instruction telling which video
 1297 clip(s) to keep (similar to examples above)
 1298 and the paired modified text question. If
 1299 you think the screening may compose more
 1300 than one round, you only need to perform one
 1301 round next.
 1302 As you cannot see the actual video, your plan
 1303 might fail, if the downstream video editing
 1304 agent taking and executing your instruction
 1305 on video clipping finds it infeasible.
 1306 Therefore, your plan is referred to as a '
 1307 trial'.
 1308 In your current situation, some rounds or trials
 1309 of the screening process might have already
 1310 been conducted, and provided as following
 1311 information:
 1312 1, the 'failure history': it is about the
 1313 history of failed previous trials of the
 1314 screening of your current round. Here is the
 1315 failure history for you (an empty indicates
 1316 that you are making the first trial of this
 1317 round): [failure_history_flag];
 1318 2, the 'success history': it is about the
 1319 history of the one success trial of all
 1320 previous round. Here is the failure history
 1321 for you (an empty indicates that you are at
 1322 the first round): [success_history_flag].
 1323
 1324 Now, here is the original question of this round
 1325 for video question answering: "[
 1326 question_flag]"
 1327 Again, please tell your plan which describes
 1328 what part of the video should be kept. Also,
 1329 give the modified question under the
 1330 assumption that the video is processed
 1331 smoothly according to your plan.
 1332
 1333 If you feel that there is no room to make such
 1334 screening (e.g. when the question is being
 1335 general like "what is this video about") so
 1336 you feel that you shouldn't make any plan,
 1337 you should decide to terminate the process.
 1338
 1339 Hints:
 1340 1. Before processing, remember to take a look in
 1341 the 'failure history' and 'success history'
 1342 information;
 1343 2. If you find that the success history contains
 1344 a case whose modified_question and the
 1345 description both significantly overlap with
 1346 the ones you are about to make, then you
 1347 should avoid making the same plan again. In
 1348 this case, you should switch to a clear
 1349 reasonable sensible alternative plan, or
 1350 make a decision to terminate the
 1351 modification process if you can't
 1352 confidently find one;
 1353 3. If you find that the failure history contains
 1354 a case whose modified_question and the
 1355 description both significantly overlap with
 1356 the ones you are about to make, or that any
 1357 of the "reason" in the failure history
 1358 records is going to make your attempt fail,
 1359 then you should avoid making the same plan
 1360 again. In this case, you should switch to a
 1361 clear reasonable sensible alternative plan,
 1362 or make a decision to terminate the
 1363 modification process if you can't
 1364 confidently find one.
 1365

1366 Return your plan in this json format (keep in
 1367 mind here, that your response should be in
 1368 json format):
 1369 {"decision": [your decision, either "process" or
 1370 "terminate"], "modified_question": [the
 1371 modified question, or "N/A" if your previous
 1372 decision is "terminate"], "description": [
 1373 Description of what part of the video should
 1374 be kept as wanted sub-clip. The description
 1375 will be passed to downstream processor to
 1376 validate. Return "N/A" if your decision is "
 1377 terminate".]}
 1378 """"
 1379
 1380

1382 **Prompt for the Validator module:**

1383 You will be given a natural language instruction
 1384 telling you to trim a video, which
 1385 instruction itself might be infeasible. The
 1386 reason it might be infeasible is because the
 1387 agent who gave the instruction had no
 1388 access to the actual video content, so it
 1389 might be infeasible if take the actual video
 1390 content into consideration.
 1391
 1392 Therefore, I need you to be a helpful assistant
 1393 to confirm if the trimming plan is feasible.
 1394 Specifically, your job is to act as a validator
 1395 to validate whether it is feasible (whether
 1396 the video content really supports the plan).
 1397 If it is feasible, you will need to implement
 1398 the plan and return the resulting sub-clip
 1399 of the plan in the form of a two-layer list.
 1400 [In context example flag 1]
 1401 If it is not feasible, you will need to tell the
 1402 reason.
 1403
 1404 Here are some examples and explanations for
 1405 infeasible trimming instructions:
 1406
 1407 [start of examples]
 1408
 1409 [In context examples flag 2]
 1410
 1411 [end of examples]
 1412
 1413 You can invoke a viewer multiple times to
 1414 acquire the video content (partial content
 1415 each time), which viewer is a downstream
 1416 module prepared to assist you.
 1417
 1418 Note that the viewer is capable to deal with two
 1419 type of questions:
 1420
 1421 1. snippet rough scanning, e.g. "what is the
 1422 video about from xxx second to xxx second?"
 1423 2. localizing, e.g. "which segment contains xxx
 1424 (event/object)?"
 1425
 1426 Therefore, if [your decision] is "view", [your
 1427 message] should follow one of the above two
 1428 example templates.
 1429
 1430 Here is the trimming instruction you need to
 1431 validate: [plan_flag].
 1432 The video length is [video_length_flag] seconds,
 1433 and its frame rate is [frame_rate_flag] fps
 1434
 1435

1436
 1437
 1438 Now, in each following turn of this conversation,
 1439 you need to give your response in this json
 1440 format: {"decision": [your decision], "
 1441 message": [your message]}. This works as
 1442 follows:
 1443 [your decision]: either "succeeded", "failed",
 1444 or "view". "succeeded" means that the plan
 1445 is successfully implemented, "failed" means
 1446 that it is not, and by "view" you invoke the
 1447 viewer to provide partial video content for
 1448 you. If you choose "view", the viewer will
 1449 take your message and return as you
 1450 requested, and the conversation will
 1451 continue. If you choose "succeeded" or "
 1452 failed", it will be your final decision and
 1453 the conversation will end.
 1454 [your message]: if you choose "succeeded" as
 1455 your decision, then it should be the two-
 1456 layer list as mentioned before, as the video
 1457 edit result of the plan. If you choose "
 1458 failed", this should be a brief reason on
 1459 the failure (e.g. requested timestamp
 1460 exceeds video length, video doesnt have the
 1461 object/event needed, etc.). If you choose "
 1462 view", this should be the question to ask
 1463 the viewer about the video content.
 1464
 1465 Hint: it is not always necessary to invoke the
 1466 viewer. [In context example flags 3]
 1467
 1468 For your ease of decision, here are some initial
 1469 frame captions and their timestamps for you
 1470 (in the form of key value pairs, where the
 1471 value is the frame caption and the key is
 1472 its corresponding timestamp, in the unit of
 1473 second): [initial_captions_flag].
 1474 [sys_usr_split]
 1475
 1476 Now let's start!
 1477

1478 **Prompt for Viewer module:**

1479 You are a helpful and smart assistant that can
 1480 respond to an upstream request a video by
 1481 invoking tools. The length of this video is
 1482 [video_length_flag].
 1483
 1484 Here is the content of the request: [
 1485 validator_request_flag].
 1486
 1487 Here are the tools you can access (you might
 1488 access them multiple times if you want):
 1489
 1490 1. scan(start, end): Return the overall caption
 1491 of the video snippet (clip) between start
 1492 and end timestamp, which are the parameters
 1493 with the unit of second.
 1494 2. localize(query): Return the video location,
 1495 in the form of a timestamp range given by
 1496 the start and end timestamp, which contains
 1497 the visual content of the query (which query
 1498 might be an object, event, etc.)
 1499 3. get_image_cap(timestamp): parameter "
 1500 timestamp" is an integer in the unit of
 1501 second. Return the caption of the video
 1502 frame at the given timestamp (regard the
 1503 frame as an image).
 1504
 1505

1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528

Now, in each turn of the following conversation, your response should be in the following json format: {"decision": [your decision], "message": [your message]}. This works as follows:
[your decision]: either "tool" (if you wish to call tool in this round) or "respond" (if you feel that you are able to respond to the upstream module's request by comprehending your current knowledge acquired about the video.).
[your message]: if your decision is "tool", then this should only contain tool calling following given format, e.g. 'get_image_cap(10)', 'scan(21, 35)', 'localize("kicking the ball")'. If your decision is "respond", then this should be your response to the upstream request.
[sys_usr_split]
Now let's start!

E.3 Prompt used for ReSimplifyIt-simple framework

1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574

You are a assistant for the video question answering process, in which a candidate is presented with a video and a question for them to answer.\
Your objective is to help the candidate so that they will be able to give the answer with watching the shortest possible sub-clip(s) of the video. \
Your task is to cut the video to acquire this sub-clip(s) and also to modify the question, so that the candidate directly answering your modified question with presented only this sub-clip(s) of the video would be equivalent to answering the original question with presented the original whole uncut video. \
[In context examples flag 1]
You will need to cut the video in the form of providing me the timestamps, which is a list of [start, end] unit clips in the unit of second. \
A tool (python function) will be helping you to get the frame caption of at a certain timestamp (in the unit of second). Whenever you need to call this tool, send a message in this json format: {"decision": "tool", "parameter": [timestamp you need]}. [In context example flag 2] \
Whenever you think you are confident enough to provide the timestamp, return {"decision": "end", "timestamps": [your result timestamps], "revised_question": [your revised question]}. [In context example flag 3].
Before we formally begin, here is a set of original captions with their timestamps provided for you to have an overall rough understanding of the video: [initial_captions_flag].

Also, the frame rate of this video is [frame_rate_flag] frames per second, and the total duration is [duration_flag].
[sys_usr_split]Now let's begin! and the original question is "[original_question_flag]".

E.4 Prompt used for ReSimplifyIt-blind framework

You are a assistant for the video question answering process, in which a candidate is presented with a video and a question for them to answer.\
Your objective is to help the candidate so that they will be able to give the answer with watching the shortest possible sub-clip(s) of the video. \
Your task is to cut the video to acquire this sub-clip(s) and also to modify the question, so that the candidate directly answering your modified question with presented only this sub-clip(s) of the video would be equivalent to answering the original question with presented the original whole uncut video. \
[In context examples flag 1]
You will be provided with a list of tools to process the video, and the original question to be answered by the candidate based on which to select the frames. Here is the list of tools you have access to, with the description (content in the brackets are the arguments needed):
[1]: get_duration(): return the duration of the video as a floating point value.
[2]: get_resolution(): return the resolution of the video, as a tuple.
[3]: get_total_frame_num(): return total number of the frames of the video, as an integer.
[4]: grounding_select(obj_name, concerned_indices_input): return, in the form of a list of integers, the indices of all frames containing the object given by obj_name, after taking the intersection of indices provided by the argument 'concerned_indices_input'. 'concerned_indices_input' is also a list of indices, and will be set to indices of all frames in the video if 'None' is passed.
[5]: indices_list_intersect(frame_indices_list_1, frame_indices_list_2): return, in the form of a list of integers, the intersection of the two arguments as list. Both argument 'frame_indices_list_1' and 'frame_indices_list_2' are a list of indices.
[6]: indices_list_union(frame_indices_list_1, frame_indices_list_2): return, in the form of a list of integers, the union of the two arguments as list. Both argument 'frame_indices_list_1' and 'frame_indices_list_2' are a list of indices.
[7]: indices_concat_and_fill(frame_indices_list_1, frame_indices_list_2): first take the sorted union of the two lists given by the arguments, and then fill in all the missing values so that every two

1575
1576
1577
1578
1579
1580
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643

adjacent element only differ by 1. Both argument 'frame_indices_list_1' and 'frame_indices_list_2' are a list of indices.

[8]: indices_concat(frame_indices_list_1, frame_indices_list_2): return, in the form of a list, the concatenation of the two lists provided by the arguments.

[9]: timestamp_to_single_index(timestamp): return a list with a single integer, which integer is the index of the frame at the given timestamp. The argument timestamp is a floating point value, whose unit is second.

[10]: single_timestamp_to_index_range(timestamp): return, in the form of a list, the indices of 60 consecutive frames, the midpoint of which is at the given timestamp. The argument timestamp is a floating point value, whose unit is second.

[11]: range_timestamp_to_index_range(start, end): return, in the form of a list, the indices of all frames which are between the two timestamps which are provided by the arguments. The argument start and end are both floating point values, whose unit are both second.

Above are all the tools to have access to. Please note that selecting frames out of all the frames of the original video is being cut and clipped, therefore you will also need to modify the aforementioned prompt, to make it align well with the reduced video frames.

[sys_usr_split]

Now the original question is: [question_flag]. Having access to the information of all the tools mentioned above, provide me the python code which could achieve the selection of frames. You may define variables to store intermediate result, and determine the value of some arguments when necessary, but you should not require the downstream task operator to replace any of your assumption on arguments, as no more information but the original video is provided to the downstream task. Please use the variable name 'final_frames' to store your final list of frame index. Please only provide the code and revised question in this format: {"Code:[your whole paragraph of code] Revised question:[your revised prompt]"}, where [your whole paragraph of code] should be an empty string if you think no tools need to be called and the whole original video should be passed to the downstream task.

1700 F More on Related Work

1701 **Video-LLMs for VideoQA** Video-LLMs have
 1702 spurred a wave of models aimed at enhancing video
 1703 understanding by leveraging the language capabilities
 1704 of large language models (LLMs) (Lin et al.,
 1705 2023; Ma et al., 2024; Li et al., 2024b, 2023b;
 1706 Liu et al., 2024b; Xu et al., 2024a; Wang et al.,
 1707 2025; Bai et al., 2025b,a; Li et al., 2024a). Some
 1708 approaches (Chen et al., 2023; Li et al., 2024c,

2023a) utilize dedicated video encoders such as
 video transformers (Bertasius et al., 2021) or con-
 volution networks. However, these designs often
 demand large-scale annotated video-text data and
 significant computational resources. To address
 this, alternative methods adapt pre-trained image-
 domain MLLMs to video inputs (Liu et al., 2024c;
 Maaz et al., 2024), offering improved practical-
 ity. Nonetheless, the sparsity and query-invariant
 nature of video encoding in existing models limits
 their ability to capture fine-grained spatial-temporal
 details effectively, especially under token budget
 constraints. This work addresses such inefficien-
 cies by introducing a query-adaptive processing
 mechanism inspired by the principle of informa-
 tion screening, aiming to reconcile the trade-off
 between token efficiency and representational fi-
 delity.

LLM-assisted Agentic Reasoning for VideoQA

Another line of research for video question an-
 swering lies in building pure-text LLM assisted
 frameworks or multi-agent systems for VideoQA
 (Wang et al., 2024c; Shang et al., 2024; Wang et al.,
 2024b). Compared to Video-LLMs, which rep-
 resents end-to-end single-pass approaches, these
 methods adopt traditional or LLM-based meth-
 ods to proactively sample relevant video frames.
 VideoAgent (Wang et al., 2024b) and TravLER
 (Shang et al., 2024) utilized LLM’s planning abil-
 ity to conduct iterative keyframe searching, while
 VideoTree (Wang et al., 2024c) presented query-
 adaptive hierarchical tree-based keyframe selec-
 tion.

Temporal Sentence Grounding for Videos

Temporal sentence grounding (TSG) aims to lo-
 calize the video segment best matching a lan-
 guage query. Early sliding-window methods (e.g.,
 TALL (Gao et al., 2017), MCN (Hendricks et al.,
 2017)) were costly, while later proposal-based (Xu
 et al., 2019; Chen and Jiang, 2019) and proposal-
 free methods (Yuan et al., 2019; Zhang et al., 2020)
 improved efficiency via query-guided proposals or
 direct boundary prediction. All these methods take
 a video and a query as input and predict a match-
 ing temporal span. In contrast, our proposed TVS
 generalizes beyond TSG by supporting inter-frame
 reasoning and joint adaptation over both video and
 query, rather than retrieving only frames directly
 mentioned in the query.

1758 **Grounded videoQA** Another seemingly-similar
1759 line of research lies in integrating grounding tech-
1760 niques into VideoQA pipelines (Xiao et al., 2024;
1761 Chen et al., 2024a). Grounded VideoQA seeks
1762 to enhance model faithfulness by explicitly link-
1763 ing a model’s textual output to “visual cues” or
1764 “evidence” within the video. This approach is ef-
1765 fective at mitigating hallucinations for descriptive
1766 queries by enforcing *perceptual* alignment. descrip-
1767 tive tasks, the approach suffers from a fundamental
1768 conceptual incongruity with the nature of complex
1769 video reasoning. The semantics of a query and
1770 its corresponding video segment are often holisti-
1771 cally entangled; for example, a high-level question
1772 about intent, causality, procedure, or even a simple
1773 temporal relational reasoning does not map to an
1774 atomic piece of visual evidence but is inferred from
1775 a continuous temporal context. The attempt to im-
1776 pose a discrete justification framework onto this
1777 intertwined semantic space leads to an inherently
1778 brittle and ill-defined notion of a “visual cue”.

1779 TVS naturally sidesteps this ambiguity by fun-
1780 damentally reframing the objective, turning from
1781 “*what to answer*” to “*what to ask*”. More fundamen-
1782 tally, TVS introduces a more principled paradigm
1783 that respects this intrinsic entanglement. Rather
1784 than seeking to atomize evidence for a given an-
1785 swer, TVS reformulates the task itself through a
1786 priori contextual simplification. It isolates the mini-
1787 mal, self-contained (*video, query*) sub-problem
1788 through an **endomorph**ic transformation that
1789 keeps with the original task space unchanged and
1790 preserves the necessary holistic context for rea-
1791 soning. This approach is not merely a circumven-
1792 tion of the definitional challenge; it establishes a
1793 more fundamentally sound and cognitively aligned
1794 task. By first reducing the problem space to a man-
1795 ageable and semantically coherent unit—mirroring
1796 the human strategy of simplifying a problem be-
1797 fore attempting to solve it—TVS offers a more
1798 robust foundation for complex video-language un-
1799 derstanding.