

T2SGrid: Temporal-to-Spatial Gridification for Video Temporal Grounding

Anonymous CVPR submission

Abstract

001 Video Temporal Grounding (VTG) aims to localize the
002 video segment that corresponds to a natural language
003 query, which requires a comprehensive understanding of
004 complex temporal dynamics. Existing Vision-LMMs typi-
005 cally perceive temporal dynamics via positional encoding,
006 text-based timestamps, or visual frame numbering. How-
007 ever, these approaches exhibit notable limitations: assign-
008 ing each frame a text-based timestamp token introduces ad-
009 ditional computational overhead and leads to sparsity in vi-
010 sual attention, positional encoding struggles to capture ab-
011 solute temporal information, and visual frame numbering
012 often compromises spatial detail. To address these issues,
013 we propose Temporal to Spatial Gridification (T2SGrid), a
014 novel framework that reformulates video temporal under-
015 standing as a spatial understanding task. The core idea
016 of T2SGrid is to process video content in clips rather than
017 individual frames. we employ a overlapping sliding win-
018 dows mechanism to segment the video into temporal clips.
019 Within each window, frames are arranged chronologically
020 in a row-major order into a composite grid image, effec-
021 tively transforming temporal sequences into structured 2D
022 layouts. The gridification not only encodes temporal infor-
023 mation but also enhances local attention within each grid.
024 Furthermore, T2SGrid enables the use of composite text
025 timestamps to establish global temporal awareness. Ex-
026 periments on standard VTG benchmarks demonstrate that
027 T2SGrid achieves superior performance.

028 1. Introduction

029 Video content has become a ubiquitous medium for infor-
030 mation dissemination, yet efficiently localizing specific mo-
031 ments within vast, unstructured video streams remains a
032 fundamental challenge. Video Temporal Grounding (VTG)
033 [4, 14, 32, 42, 43], the task of identifying the precise video
034 segment that semantically corresponds to a natural language
035 query, serves as a key step toward bridging this gap. Success
036 in VTG hinges on a model’s ability to comprehend not only
037 static visual content but also intricate temporal dynamics,

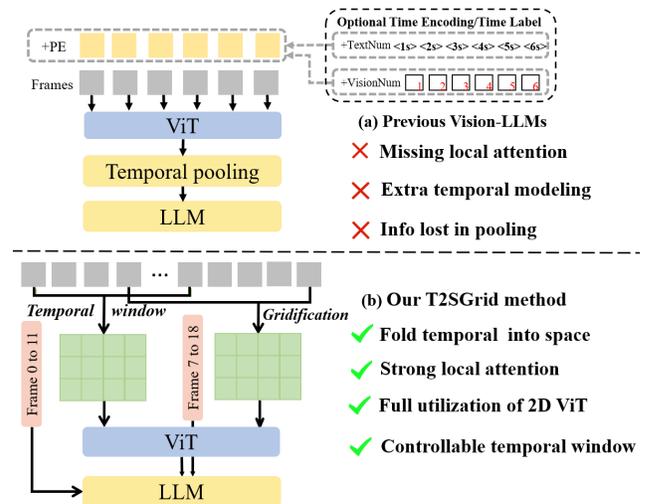


Figure 1. Comparison between T2SGrid and previous Vision-LLMs. (a) Traditional methods process frames sequentially or apply temporal pooling to capture information at multiple scales. However, sequential processing and pooling can cause information loss, obscure local temporal details within a window, and often require additional time encoding/label for temporal modeling. (b) Our T2SGrid method folds multiple frames within a temporal window into the spatial dimension via gridification, allowing direct processing a temporal window by the standard 2D ViT. This leverages the model’s strong spatial reasoning capability for temporal understanding.

including action sequences, event duration, and long-range dependencies.

The emergence of Vision-Large Language Models (Vision-LMMs) [2, 17, 22, 51, 53] has revolutionized visual understanding [50], exhibiting remarkable zero-shot reasoning [21, 25, 27, 33, 43, 56, 63, 69, 75, 79, 80] and multi-modal comprehension [13, 30, 38, 50, 76], primarily on static images. However, extending these spatially-oriented architectures to the temporal domain remains non-trivial. As shown in Figure 1 (a), Current approaches to incorporate temporal awareness [9, 18, 24, 29] include adding positional encoding (+PE) [15, 23, 35, 48, 57], as seen in Qwen2.5-VL [2] models, and using text-based timestamps

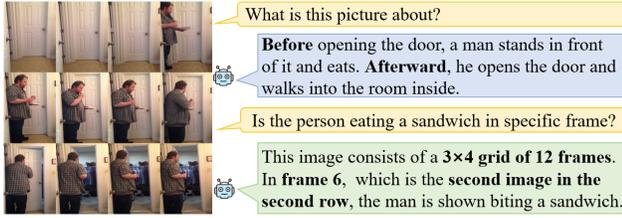


Figure 2. **Illustration of Qwen2-VL temporal reasoning on 2D grid layouts.** The model correctly infers the temporal order of events (before and after) and accurately identifies the biting action in frame 6.

(+TextNum) [31, 81, 82], as in Qwen3-VL [51] models, and visual numbering (+VisualNum) [5, 60, 61]. However, these approaches faces inherent limitations. Positional encoding, though effective for sequence modeling, fail to capture absolute temporal positions essential for grounding specific events and require additional encoding modules. Text-based timestamps (e.g., prompting with “Frame 1” or “1 seconds”) inevitably introduce a rapidly growing amount of textual tokens that leads to increasingly sparse visual attention as the video length increases. Visual frame numbering, which overlays timestamps directly on frames, degrades spatial detail and undermines the very visual features that Vision-LMMs rely on for semantic understanding.

To address these challenges, we propose Temporal to Spatial Gridification (T2SGrid), a novel framework that partially reformulates temporal reasoning as a spatial problem. Previously, Vision-LLMs treat a video as a linear sequence of frames. T2SGrid takes a different approach, which partitions the video into configurable temporal windows and arranges the frames within each window into a composite image with grid-structured layout (Figure 1 (b)), which we call “gridification” of the temporal window. Then, the Vision-LLM processes the grid images rather than the original frames.

The gridification process is central to our approach, establishing a novel paradigm for video input representation. T2SGrid offers several distinct advantages. First, it transforms temporal dynamics within a window into a spatial layout, leveraging spatial attention to enhance the understanding of local temporal dynamics. Second, the row-major frame arrangement in the grid image and the partial overlap between adjacent grid images implicitly function as a form of positional encoding. Crucially, the effectiveness of these advantages is fundamentally enabled by the demonstrated capability of modern Vision-LLMs to interpret grid-based imagery. For instance, as shown in Figure 2, these models can infer relative temporal relationships (e.g., “before” and “after”) by reading the spatial configuration of the grid from top-left to bottom-right, and can also identify the specific order of frames within the grid.

As shown in Figure 1 (b), beyond the input gridification, to enhance the model’s awareness of global time, we interleave the grid images with text-based timestamps. Unlike existing approaches that assign a timestamp to each individual frame, we associate a single composite text-based timestamp (e.g., “Frame 0 to 11”) with each grid image. By grouping multiple frames under one temporal descriptor, the model learns to associate a local window of visual content with a unified time interval, further improving temporal understanding.

In summary, our primary contributions are as follows. (1) We introduce T2SGrid, a novel paradigm that shifts video processing from individual frames to local temporal clips by transforming sequences of frames within a sliding window into a single, composite grid image. (2) Instead of assigning a timestamp to each frame, we use a single composite text timestamp for each grid image, enhancing global temporal awareness. (3) Extensive experiments on standard VTG and VQA benchmarks demonstrate that T2SGrid achieves superior performance.

2. Related Work

2.1. Video Temporal Grounding By Vision-LLMs

Video temporal grounding (VTG) [6, 8, 34, 36, 40, 66, 74, 77] aims to locate the precise time of specific actions or events within a video. For current Vision-LLMs [28, 38, 47, 51, 53], VTG is crucial for both temporal and spatial understanding. Many Vision-LLM-based approaches [12, 20, 32, 41, 43, 52, 62, 65, 80] have been proposed to tackle the VTG task. Some methods focus on fine-tuning existing models by providing textual timestamps [19, 36, 55] and constructing large temporally annotated datasets [45, 49]. Others design specialized temporal-aware modules using pooling operations [24, 38, 64, 72, 73] to enhance temporal reasoning. There are also approaches [59, 68, 70, 78, 83] that overlay frame indices [11, 37, 46, 58, 67] onto the original images to provide temporal information and feed the frames sequentially into the model. Although these methods achieve some success, they either require designing task-specific modules for temporal reasoning or compromise spatial information to preserve temporal cues. To address these limitations, we propose T2SGrid, which employs a gridification strategy along with a lightweight textual timestamps to provide global temporal context, eliminating the need for specialized module design or extensive dataset construction, and achieving superior performance.

2.2. Grid-Based Video Representation

Although several studies [3, 7, 10, 26, 71] have explored merging multiple video frames into a single composite image for video understanding, most of them primarily target

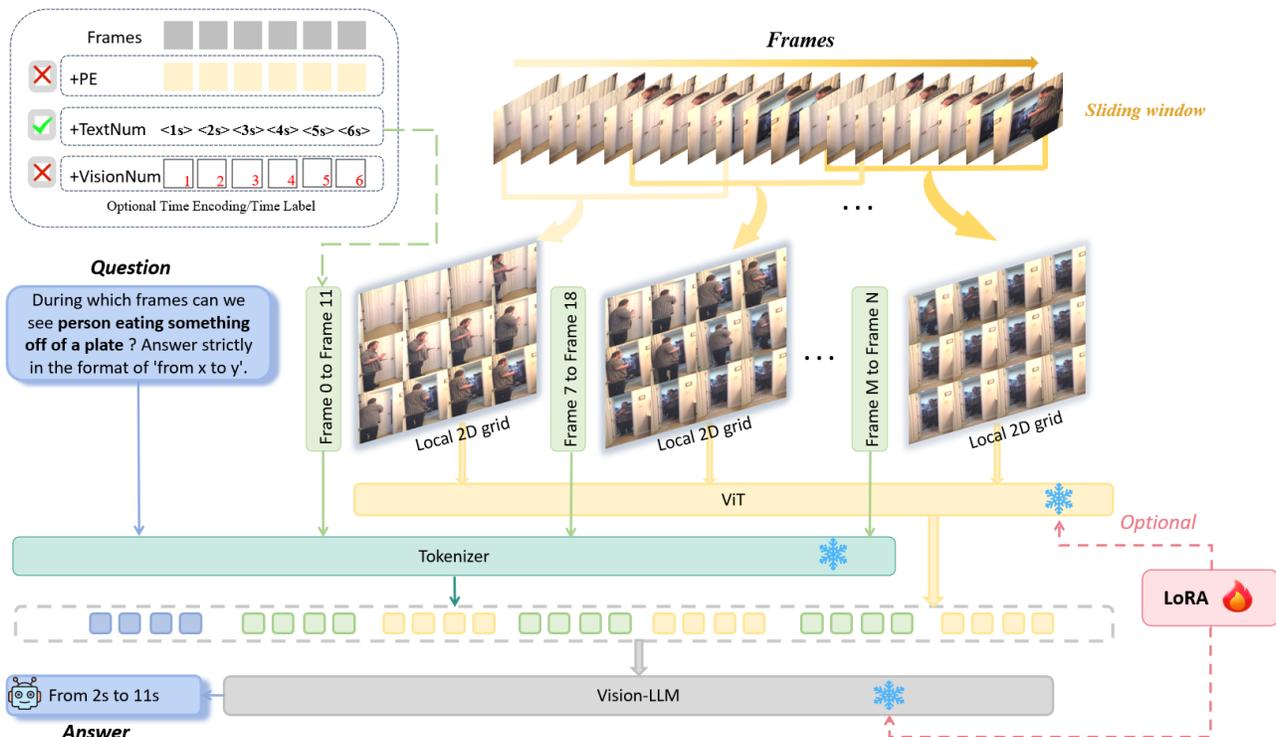


Figure 3. **Overview of our T2SGrid framework.** The original video frames are first arranged into a 2D grid in a row-major order (*gridification*) to enable spatialized temporal reasoning. A lightweight composite textual timestamp is incorporated to provide global temporal awareness. Our framework can operate in a training-free manner or be further enhanced via LoRA fine-tuning.

141 video question answering (VQA) [1, 39, 47]. For example,
 142 IG-VLM [26] selects a small set of keyframes and
 143 concatenates them into a single large image to facilitate VQA,
 144 yet such a keyframe grid can only model coarse-grained
 145 events. DynImg [3] enlarges a representative keyframe and
 146 appends several smaller frames below it as temporal cues
 147 to enhance video comprehension; however, such heuristic
 148 designs still fail to preserve fine-grained temporal order,
 149 thereby limiting the potential of grid-based representations
 150 for temporal reasoning. In contrast, our work is the first to
 151 reveal and leverage the intrinsic temporal information en-
 152 coded within the spatial grid itself, and to apply this prop-
 153 erty to video temporal grounding, thereby enhancing the
 154 generalization capability of gridification for temporal un-
 155 derstanding.

156 3. Method

157 We propose an innovative temporal-to-spatial gridification
 158 (T2SGrid) method. In Section 3.1, we first analyze how
 159 the gridification mechanism enhances local attention from a
 160 spatial perspective and improves temporal perception from
 161 the viewpoint of temporal attention. As illustrated in Fig-
 162 ure 3, our approach consists of two main stages: (1) Sliding-

Window Spatiotemporal Gridification, and (2) Temporal
 Modeling with T2SGrid, which are detailed in Section 3.2
 and Section 3.3, respectively.

3.1. Attention Analysis

163 Current Vision-LLMs process videos by encoding frames
 164 into a sequential feature representation. This sequence,
 165 concatenated over time, is then aligned with a language query
 166 to facilitate understanding. While this approach is effective
 167 for content recognition, it presents inherent challenges for
 168 tasks requiring fine-grained temporal reasoning. To investi-
 169 gate the underlying limitations of this sequential method-
 170 ology, we conducted an in-depth analysis of the model’s in-
 171 ternal attention mechanisms. Our central hypothesis is that
 172 the sequential processing of individual frame features biases
 173 Vision-LLMs towards recognizing a series of static spatial
 174 configurations rather than understanding the dynamic tem-
 175 poral evolution between them. To test this, we used Qwen2-
 176 VL-7B as case studies.

177 We visualize the cross-attention maps between the visual
 178 tokens and the query text token, and project them back to the
 179 original image or the grid representation. As shown in Fig-
 180 ure 4, when given the query “a person is putting a picture
 181 onto the wall”, the attention maps under sequential-frame
 182
 183
 184
 185

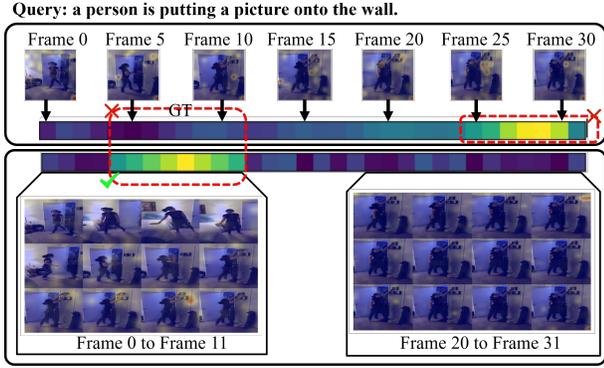


Figure 4. **Top:** spatial and temporal attention with sequential frame input. **Bottom:** spatial and temporal attention with grid-based input. The grid method captures dynamic actions and maintain focus on the correct temporal intervals.

input appear scattered or focus mainly on the person or the picture in a specific frame. This pattern indicates that sequential processing encourages the model to recognize what objects are present, but fails to capture how these objects change over time, resulting in weak sensitivity to motion-induced spatial variations. In contrast, our T2SGrid representation reorganizes the video into a structured 2D layout that places temporally adjacent patches into spatially coherent neighborhoods. This gridified structure strengthens local spatial attention and enables the model to capture subtle structural changes in the picture frame across neighboring cells. As shown in the Figure 4, the model now attends to how the person moves the picture onto the wall, demonstrating improved temporal understanding through spatial reasoning. We further average the attention over all visual tokens within each frame to obtain frame-level temporal attention. As shown in the figure, the temporal attention derived from sequential-frame input exhibits peaks that deviate substantially from the Ground Truth (GT), indicating that its temporal modeling is primarily driven by static object saliency rather than the actual progression of the action. In contrast, T2SGrid makes temporal cues more explicit: by embedding temporal adjacency into a 2D spatial layout and combining it with lightweight text-based temporal anchors, the model produces temporally consistent activation patterns that closely align with the GT. Further illustrative examples can be found in the Appendix.

In conclusion, our analyses reveal a fundamental limitation in sequential-frame processing: although the model is proficient at spatial recognition, this ability does not naturally extend to capturing fine-grained temporal progression. Sequential alignment encourages frame-by-frame object matching, causing the model to rely on static object saliency while overlooking subtle inter-frame motion cues. By contrast, the proposed gridification mechanism strength-

ens local spatial attention, enhances sensitivity to dynamic changes, and produces temporally consistent activation patterns that align closely with the GT. By reframing temporal reasoning into a structured 2D spatial layout, our approach enables the model to leverage its inherent spatial priors to capture genuine temporal cues and achieve more accurate action localization.

3.2. Sliding Window Spatiotemporal Gridification

We first employ a sliding window mechanism to obtain frame sequences, after which we apply spatiotemporal gridification to convert them into a coherent 2D spatial layout. Concretely, given a video with T frames, we define a temporal window size k and stride s . The i -th window W_i consists of:

$$W_i = \{f_{i \times s}, f_{i \times s + 1}, \dots, f_{i \times s + k - 1}\}. \quad (1)$$

For each W_i , we spatially rearrange its k **original resolution frames** into a single composite grid G_i . Notably, our method does not compromise spatial resolution; it only performs gridification of the original frames. The layout of this grid is flexible (e.g., $M \times N$) as long as $M \times N = k$. For instance, nine frames can be arranged in a 3×3 grid. The frames are consistently ordered in a row-major fashion (left-to-right, top-to-bottom).

Temporal localization methods processing input videos that uniformly sample frame sequences are fundamentally equivalent to sliding window processing with stride $s = 1$ and window size $k = 1$. While capturing per-frame static information, this approach poorly perceives dynamic changes between adjacent frames. Our method overcomes this limitation by setting $k > 1$ (e.g., 3×3 grid), where each window contains k consecutive frames forming a local temporal segment. When processing center frame f_t , the model simultaneously receives its neighborhood f_{t-4}, \dots, f_{t+4} , establishing complete dynamic context. To avoid splitting critical actions in short videos, we introduce overlap by setting the sliding-window stride $s < k$, preserving temporal continuity and allowing the model to capture both local dynamics and maintain global consistency. For longer videos, we set $s = k$ to avoid excessive overlap and computational redundancy. This design makes our approach suitable for long videos and robust to a wide range of frame rates (FPS), while maintaining the original video resolution throughout the spatiotemporal gridification process.

3.3. Temporal Modeling with T2SGrid

3.3.1. Implicit Temporal Encoding via Grid Layout

Although the T2SGrid is inherently spatial, arranging frames in a **row-major order** induces a deterministic mapping that implicitly encodes temporal information. Let N_c denote the number of frames per row. The temporal index t_f of a frame with row-column coordinates (r_f, c_f) is uniquely

271 defined as
 272
$$t_f = r_f \times N_c + c_f, \quad (2)$$

273 reflecting the linear progression along the row-major layout.
 274 The self-attention mechanism, however, operates on
 275 patch-level coordinates (r_p, c_p) and their corresponding 2D
 276 positional embeddings $E(r_p, c_p)$. A frame-level coordinate
 277 can be recovered from patch-level coordinates using

278
$$r_f = \lfloor r_p / h_{\text{patch}} \rfloor, \quad c_f = \lfloor c_p / w_{\text{patch}} \rfloor, \quad (3)$$

279 where h_{patch} and w_{patch} are the height and width of a frame
 280 in patches. Substituting Equation (3) into Equation (2) gives

281
$$t_f = \lfloor r_p / h_{\text{patch}} \rfloor \times N_c + \lfloor c_p / w_{\text{patch}} \rfloor, \quad (4)$$

282 showing that the temporal index is a well-defined function
 283 of patch coordinates.

284 Therefore, the 2D positional embeddings $E(r_p, c_p)$ con-
 285 tain sufficient information for the model to infer the **tempo-**
 286 **ral order** of frames, enabling implicit temporal reasoning
 287 without requiring explicit timestamps or frame identifiers.

288 Figure 2 shows concrete examples of how Qwen2-
 289 VL [53] interpret the grid structure, providing empirical ev-
 290 idence that the grid layout enables implicit temporal encod-
 291 ing via spatial understanding.

292 3.3.2. Absolute Global Temporal Awareness.

293 Although the sliding window strategy enhances local tem-
 294 poral reasoning through spatiotemporal grid representations
 295 of videos, it loses the absolute temporal position of each
 296 segment. This limitation arises from its capture of only rel-
 297 ative temporal relationships, such as action continuity be-
 298 tween adjacent frames. This limitation prevents the estab-
 299 lishment of absolute positioning along the global timeline,
 300 as the model loses perception of corresponding absolute
 301 time intervals in the original video (e.g., “From Frame 0 to
 302 Frame 8”) while comprehending dynamic evolution within
 303 grid G_i . Such deficiency impedes video temporal localiza-
 304 tion tasks demanding precise timestamp outputs (e.g., “Xs
 305 to Ys”) or global narrative comprehension (e.g., “What hap-
 306 pened at 7 seconds in the video?”), where absolute timing
 307 remains essential to satisfy rigid requirements for a tempo-
 308 ral reference system.

309 To preserve global time awareness, we introduce Absolu-
 310 te Time as solution. Before passing each grid image G_i
 311 to the LMM, we prepend a textual timestamp to the input
 312 prompt:

313
$$\text{Prompt}_i = [\text{“ from } T_{\text{start}} \text{ to } T_{\text{end}} \text{.”}]; [\text{Image: } G_i] \quad (5)$$

314 where T_{start} and T_{end} correspond to the start and end time of
 315 the i -th grid.

For holistic video understanding, multiple grids are or- 316
 ganized in an interleaved text-image sequence: 317

$$\begin{aligned} & [\text{Text: } T_{\text{start}}^{(1)} \text{ to } T_{\text{end}}^{(1)}] \rightarrow [\text{Image: } G_1]; \\ & [\text{Text: } T_{\text{start}}^{(2)} \text{ to } T_{\text{end}}^{(2)}] \rightarrow [\text{Image: } G_2]; \\ & \vdots \\ & [\text{Text: } T_{\text{start}}^{(n)} \text{ to } T_{\text{end}}^{(n)}] \rightarrow [\text{Image: } G_n]. \end{aligned} \quad (6) \quad 318$$

These absolute timestamps form a continuous temporal 319
 chain across the video, enabling the model to establish se- 320
 quential relationships while preserving temporal coherence. 321
 By interleaving grids with global textual time and leverag- 322
 ing cross-attention mechanisms, the model can reason not 323
 only about the dynamics within each grid but also about 324
 their placement along the full video timeline. 325

326 4. Experiment

We evaluate our model on standard Video Temporal 327
 Grounding benchmarks. Following prior work [19, 42, 44, 328
 54], we evaluate our method on the test sets of Charades- 329
 STA [14] and ActivityNet [4]. The evaluation metrics 330
 include mean Intersection over Union (mIoU) and Rec- 331
 all@1 at multiple IoU thresholds ($R@m$), where $m \in$ 332
 $\{0.3, 0.5, 0.7\}$, consistent with previous studies [19, 44]. 333

334 4.1. Experiment setup

Training dataset and benchmark We adopt a mixed 335
 training set, combining the training splits of ActivityNet- 336
 Captions [4] and Charades [14], containing 18K videos and 337
 50K question-answer pairs. Each video in our dataset is 338
 augmented and arranged using our T2SGrid method. The 339
 question-answer pairs follow a consistent template: ques- 340
 tions are formatted as “During which frames can we see 341
 query?”, and answers are formatted as “From x to y”, where 342
 x and y denote the start and end frame numbers of the 343
 queried event. 344

Training Detail We use Qwen2-VL-7B as the base 345
 model, which lacks temporal encoding and therefore en- 346
 ables a fair comparison of temporal strategies. The model 347
 is trained for 3 epochs with a batch size of 32 and a learn- 348
 ing rate of 2×10^{-5} . We adopt LoRA fine-tuning (rank 349
 64, $\alpha = 128$) on all linear layers. All experiments are 350
 conducted on $4 \times A100$ GPUs. Charades-STA is trained and 351
 evaluated under the overlap configuration (g43_s7), 352
 whereas ActivityNet-Captions uses the no-overlap config- 353
 uration (g43_s12). Additional details are provided in the 354
 Appendix. 355

356 4.2. Main result

357 4.2.1. Superior performance on VTG datasets

Table 1 reports the performance of our T2SGrid and 358
 T2SGrid-FT methods compared with previous approaches 359

Table 1. **Performance comparison on the video temporal grounding task with prior state-of-the-art methods.** *T2SGrid* denotes the use of spatiotemporal gridification with global time awareness, while *T2SGrid-FT* indicates fine-tuning with the instruction dataset augmented using our T2SGrid method.

Model	Charades-STA				ActivityNet			
	R@0.3	R@0.5	R@0.7	mIoU	R@0.3	R@0.5	R@0.7	mIoU
<i>VTG-Tuned Vision-LLMs</i>								
GroundingGPT [32]	-	29.6	11.9	-	-	-	-	-
LITA [20]	-	-	-	-	-	25.9	-	28.6
VTG-LLM [17]	52.0	33.8	15.7	-	-	-	-	-
TimeChat [44]	47.7	22.9	12.5	30.6	30.2	16.9	8.2	21.8
VTimeLLM [19]	51.0	27.5	11.4	31.2	44.0	27.8	14.3	30.4
Momentor [42]	42.9	23.0	12.4	29.3	42.6	26.6	11.6	28.5
HawkEye [54]	50.6	31.4	14.5	33.7	49.1	29.3	10.7	32.7
TimeSuite [72]	<u>69.9</u>	<u>48.7</u>	<u>24.0</u>	-	-	-	-	-
TRACE [16]	-	40.3	19.4	-	-	<u>37.7</u>	<u>24.0</u>	<u>39.0</u>
NumPro [61]	63.8	42.0	20.6	<u>41.4</u>	<u>55.6</u>	37.5	20.6	38.8
<i>General Vision-LLMs</i>								
GPT-4o [22]	55.0	32.0	11.5	35.4	33.3	21.2	10.4	23.7
+ <i>T2SGrid</i>	57.3 (+2.3)	36.7 (+4.7)	14.8 (+3.3)	36.9 (+1.5)	47.5 (+14.2)	32.1 (+10.9)	19.7 (+9.3)	35.6 (+11.9)
Qwen3-VL-8B [51]	69.3	43.4	17.5	43.1	39.0	26.2	15.8	29.3
+ <i>T2SGrid</i>	71.4 (+2.1)	47.0 (+3.6)	20.7 (+3.2)	44.9 (+1.8)	44.8 (+5.8)	26.8 (+0.6)	16.1 (+0.3)	32.5 (+3.2)
LLaVA-OneVision-1.5-8B [28]	19.8	6.7	2.3	14.5	9.2	4.9	2.1	7.5
+ <i>T2SGrid</i>	45.0 (+25.2)	26.3 (+19.6)	11.9 (+9.6)	28.8 (+14.3)	31.5 (+22.3)	17.4 (+12.5)	10.4 (+8.3)	23.6 (+16.1)
Qwen2-VL-7B [53]	8.7	5.4	2.4	7.9	17.0	9.4	3.9	12.5
+ <i>T2SGrid</i>	70.1 (+61.4)	46.7 (+41.3)	20.1 (+17.7)	44.3 (+36.4)	46.2 (+29.2)	27.2 (+17.8)	15.4 (+11.5)	33.3 (+20.8)
+ <i>T2SGrid-FT</i>	76.9 (+68.2)	60.6 (+55.2)	35.9 (+33.5)	53.2 (+45.3)	64.4 (+47.4)	48.4 (+39.0)	29.5 (+25.6)	46.7 (+34.2)

360 on video temporal grounding benchmarks. After integrating
 361 the proposed T2SGrid temporal encoding, both open-source
 362 and closed-source models exhibit consistent improvements.

363 Specifically, GPT-4o, which already possesses strong
 364 temporal reasoning capability, achieves further gains with
 365 T2SGrid, showing that our method can enhance well-
 366 trained multimodal LLMs. Qwen3-VL-8B, a state-of-
 367 the-art vision-LLM utilizing textual timestamps for tem-
 368 poral modeling, also benefits modestly; for instance, its
 369 Charades-STA R@0.3 increases from 69.3 to 71.4. The
 370 smaller gain mainly arises from a slight conflict between our
 371 local temporal encoding and its original timestamp-based
 372 scheme. In contrast, Qwen2-VL-7B, which lacks explicit
 373 temporal encoding, exhibits a significant performance boost
 374 when combined with T2SGrid, reaching 70.1 R@0.3 and
 375 44.3 mIoU on Charades-STA that surpassing several VTG-
 376 tuned video LLMs.

377 Finally, LLaVA-OneVision1.5-8B, a model trained
 378 solely on static images, achieves remarkable improvements
 379 with T2SGrid: absolute gains of 25.2, 19.6, 9.6, and 14.3 on
 380 R@0.3, R@0.5, R@0.7, and mIoU, respectively. On Activ-
 381 ityNet, the corresponding improvements are 22.3, 12.5, 8.3,
 382 and 16.1. These results confirm that T2SGrid effectively
 383 leverages the spatial reasoning capability of image-based
 384 models to enable temporal understanding.

385 To further validate the effectiveness of our approach, We

386 fine-tune Qwen2-VL-7B on a dataset augmented and ar-
 387 ranged using our T2SGrid framework, resulting in the best
 388 mIoU scores on Charades-STA and ActivityNet, reaching
 389 53.2 and 46.7, respectively.

4.2.2. Training Effectiveness 390

391 In Table 5, we observe consistent performance gains when
 392 incorporating T2SGrid during fine-tuning (+T2SGrid-FT).
 393 Both the standard fine-tuning (FT) and T2SGrid-FT exper-
 394 iments are conducted on the same training data, ensuring a
 395 fair comparison. For LLaVA-OneVision-1.5-8B, T2SGrid
 396 yields clear improvements of +4.8 R@0.5 and +2.4 mIoU
 397 on Charades-STA, and large boosts of +13.4 R@0.3 and
 398 +10.2 mIoU on ActivityNet. For Qwen2-VL-7B, it fur-
 399 ther improves results by +3.8 R@0.5 and +2.8 mIoU on
 400 Charades-STA, and by +9.8 R@0.5, +10.3 R@0.7, and +9.0
 401 mIoU on ActivityNet. Overall, models fine-tuned with our
 402 T2SGrid framework not only achieve consistently higher
 403 scores across all metrics and datasets but also demonstrate
 404 stronger temporal grounding capability, confirming the su-
 405 perior training effectiveness of our approach.

4.2.3. Qualitative Results 406

407 Figure 5 illustrates a visualization on the Charades dataset,
 408 comparing our method with previous approaches. Given the
 409 query “person throwing a blanket onto the vacuum”, other
 410 methods predict incorrect temporal intervals, whereas our
 411

Table 2. **Training effectiveness of our method.** T2SGrid-FT consistently outperforms standard fine-tuning (FT) when trained on the same data with identical training settings, demonstrating the benefits of our approach.

Model	Charades-STA				ActivityNet			
	R@0.3	R@0.5	R@0.7	mIoU	R@0.3	R@0.5	R@0.7	mIoU
LLaVA-OneVision-1.5-8B [28]	19.8	6.7	2.3	14.5	9.2	4.9	2.1	7.5
+FT	70.6	52.9	29.8	48.4	50.1	35.1	18.5	33.4
+T2SGrid-FT	73.9 (+3.3)	57.7 (+4.8)	32.7 (+2.9)	50.8 (+2.4)	63.5 (+13.4)	45.4 (+10.3)	26.2 (+7.7)	44.6 (+10.2)
Qwen2-VL-7B [53]	8.7	5.4	2.4	7.9	17.0	9.4	3.9	12.5
+FT	73.1	56.8	32.7	50.4	54.9	37.2	19.2	37.7
+T2SGrid-FT	76.8 (+3.7)	60.6 (+3.8)	35.9 (+3.2)	53.2 (+2.8)	64.4 (+9.5)	48.4 (+9.8)	29.5 (+10.3)	46.7 (+9.0)

Table 3. **Performance comparison on Video-MME, MVBench, and VideoInstruct benchmarks.** Our T2SGrid method consistently improves temporal perception and action understanding, demonstrating strong generalization QA tasks.

Model	Video MME			MVBench				VideoInstruct		
	Temporal Perception	Temporal Reasoning	All	Action Sequence	Fine-grained Action	State Change	Scene Transition	All	Temporal Understanding	Detail Orientation
Qwen2-VL-7B	60.0	41.7	63.3	60.5	46.5	44.0	79.5	51.7	2.47	2.57
+T2SGrid	74.5	50.2	64.1	71.5	49.5	51.5	89.0	58.3	2.52	2.69

411 method accurately captures the start and end times of the
 412 action, demonstrating the effectiveness of T2SGrid. Addi-
 413 tional examples are provided in the Appendix.

4.2.4. Effectiveness on Long-Video and VQA Tasks

415 To further evaluate the effectiveness of our T2SGrid method
 416 in enhancing temporal understanding and action reason-
 417 ing, we conduct experiments on three video understand-
 418 ing benchmarks: VideoMME [13], MVBench [30], and
 419 VideoInstruct [38]. Among them, VideoMME is a long-
 420 form video QA benchmark that further demonstrates our
 421 method’s ability to reason over long-duration videos.

422 As shown in Table 3, the T2SGrid-enhanced model con-
 423 sistentlly outperforms the base Qwen2-VL-7B model across
 424 all metrics. On VideoMME, it improves temporal per-
 425 ception and temporal reasoning by 14.5 and 8.5, respec-
 426 tively. Similarly, on MVBench, TGrid yields consistent
 427 gains across all subcategories including action sequence,
 428 fine-grained action, state change, and scene transition, rais-
 429 ing the overall score from 51.7 to 58.3. On the VideoInstruct
 430 dataset, our method also achieves higher scores in both tem-
 431 poral understanding and detail orientation, confirming that
 432 T2SGrid effectively enhances temporal reasoning and fine-
 433 grained perception across diverse video QA tasks.

434 These results demonstrate the generality and robustness
 435 of T2SGrid across both long and short video QA tasks,
 436 highlighting its strong ability to improve temporal percep-
 437 tion and question answering performance.

4.3. Ablation studies

439 To isolate the contributions of each component of our
 440 method, we conduct a series of ablation studies on the Cha-
 441 rades dataset using Qwen2-VL-7B.

Table 4. **Ablation study of key components:** composite tex-
 tual timestamps (ComTextNum), sliding window, and grid layout.
 Each component contributes to improved temporal modeling and
 overall grounding performance.

ComTextNum	Sliding Window	Grid	R1@0.3	R1@0.5	R1@0.7	mIoU
✗	✗	✗	8.7	5.4	2.4	7.9
✓	✗	✗	53.5	23.2	7.9	32.9
✓	✓	✗	<u>58.3</u>	<u>35.1</u>	<u>13.6</u>	<u>36.5</u>
✓	✓	✓	70.1	46.7	20.1	44.3

Component Ablation. We conduct a comprehensive ab-
 442 lation study to evaluate the contribution of each component
 443 in our temporal encoding framework. ComTextNum de-
 444 notes the use of a textual timestamp inserted before each
 445 grid element, providing an explicit temporal anchor for the
 446 model. Sliding Window refers to partitioning the sequen-
 447 tial frames into temporally localized windows, while Grid rep-
 448 resents the transformation of frames within each window
 449 into a 2D spatial grid.
 450

451 As shown in Table 4, introducing ComTextNum substan-
 452 tially improves temporal information, boosting mIoU from
 453 7.9 to 32.9. Incorporating the sliding window mechanism
 454 further enhances local spatiotemporal attention, raising
 455 mIoU to 36.5. When we additionally reformulate each tem-
 456 poral window into a 2D grid structure, the model achieves
 457 the largest performance gain, reaching 70.1 R1@0.3 and
 458 44.3 mIoU. These improvements demonstrate that local im-
 459 plicit temporal encoding within the grid provides advan-
 460 tages over TextNum, validating the effectiveness of our core
 461 grid-based approach.

Temporal Modeling Method Ablation. We conduct an
 462 ablation study to evaluate four temporal encoding strategies
 463 on the Qwen2-VL-7B baseline: PE, TextNum, VisualNum,
 464

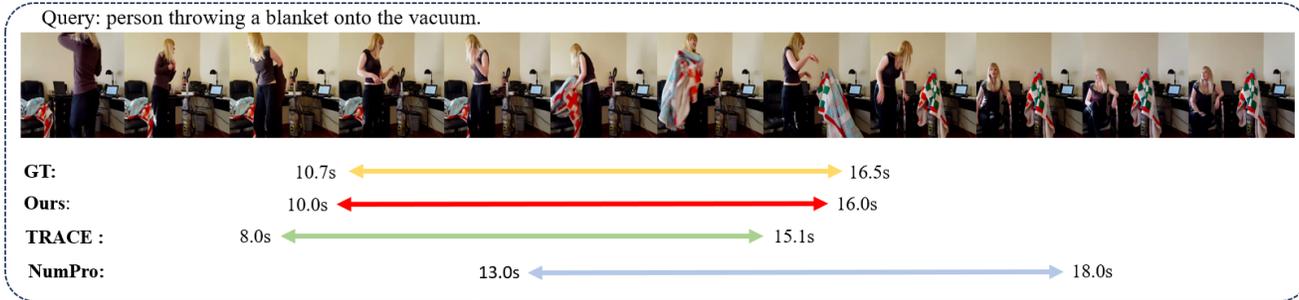


Figure 5. Qualitative Comparison with State-of-the-Art. Our method outperforms TRACE [16] and Numpro [61] on Charades by accurately identifying event boundaries in challenging scenes.

Table 5. Comparison of different temporal modeling strategies applied to the Qwen2-VL-7B baseline, including position encoding (PE), textual timestamping (TextNum), visual frame-level numeric labeling (VisualNum), and our proposed T2SGrid framework. mToken and mTime/s denote the average token count and inference time per sample, respectively.

Temporal Modeling	R1@0.3	R1@0.5	R1@0.7	mIoU	mToken	mTime/s
PE	53.1	28.4	10.6	33.8	5760.4	1.69
TextNum	53.5	23.2	7.9	32.5	5791.2	1.45
VisualNum	60.7	36.8	15.9	38.5	5760.4	2.17
Ours (wo/ overlap)	64.5	42.9	18.1	41.2	5766.1	1.43
Ours (w/ overlap)	70.2	46.7	20.1	44.3	7909.7	2.31

465 and Ours. PE encodes temporal order through positional
 466 encoding; TextNum adds textual timestamps before each
 467 frame or grid element; VisualNum draws numeric indices
 468 directly on frames. As shown in Table 5, Under the non-
 469 overlap setting, our method delivers the best overall perform-
 470 ance while reducing inference time by 34.1% over Visu-
 471 alNum. With overlap, it trades only 6% additional time for
 472 a 14% performance gain.

473 **Ablation Study on Grid Configuration** We conduct
 474 ablations over different grid configurations, denoted as
 475 $g\{col\}\{row\}_s\{stride\}$ where $k = col \times row$.

476 In Stage 1, we vary both grid size and stride (with $k = s$).
 477 As shown in Table 6, performance improves as the win-
 478 dow size increases: R1@0.3 rises from 53.5 (g11_s1) to
 479 63.7 (g33_s9), and mIoU increases from 32.9 to 41.2. The
 480 best configuration in this stage is g43_s12, reaching 64.5
 481 R1@0.3 and 41.2 mIoU. However, excessively large grids
 482 (e.g., g44_s16) lead to degraded performance, with mIoU
 483 dropping to 35.9. Similarly, g43 is also the optimal config-
 484 uration on ActivityNet.

485 Notably, under the no-overlap setting, our method does
 486 not reduce the resolution of the original frames, with the
 487 average token count (mToken) essentially the same as se-
 488 quential frame input (g11_s1), since we simply concatenate
 489 full resolution frames into a single grid. This design already
 490 yields strong gains.

491 In Stage 2, we fix the grid size to g43 and adjust the stride

Table 6. Ablation results of grid configurations on Qwen2-VL-7B. $g\{col\}\{row\}_s\{stride\}$ denoting grid size and stride. Stage 1 varies size and stride, and Stage 2 refines stride for g43.

Config	R1@0.3	R1@0.5	R1@0.7	mIoU	mToken	mTime/s
Stage 1: Exploring Grid size and stride						
g11_s1	53.5	23.2	7.9	32.9	5791.2	1.45
g22_s4	56.5	26.9	8.7	34.2	5789.0	1.43
g23_s6	61.6	38.7	16.2	38.8	5784.3	1.42
g32_s6	61.8	38.8	16.9	39.0	5784.3	1.43
g33_s9	63.7	42.8	18.4	41.1	5782.1	1.43
g34_s12	64.1	42.7	18.0	41.1	5781.4	1.42
g43_s12	64.5	42.9	18.1	41.2	5781.4	1.42
g44_s16	61.6	29.1	10.4	35.9	5779.7	1.41
Stage 2: Fine-tuning stride for optimal Grid size (g43)						
g43_s5	68.9	45.1	18.2	43.2	8971.4	2.39
g43_s6	69.5	46.4	19.6	43.7	8357.2	2.34
g43_s7	70.2	46.7	20.1	44.3	7909.6	2.31
g43_s12	64.5	42.9	18.1	41.2	5781.4	1.42

to introduce overlap. Overlapping windows yield consistent
 improvements: as the stride decreases from 12 to 7, R1@0.3
 increases from 64.5 to 70.2, and mIoU improves from 41.2
 to 44.3. The best configuration is g43_s7, which increases
 the number of tokens modestly but brings best performance.

5. Conclusion

In this work, we introduced T2SGrid, a novel frame-
 work designed to tackle the challenge of video temporal
 grounding (VTG) for Vision-LLMs. Our core idea is to
 reformulate temporal understanding as a spatial reasoning
 problem. By “gridifying” frames from a sliding temporal
 window into a single composite 2D image, T2SGrid directly
 leverages the powerful, pre-existing spatial attention mech-
 anisms of 2D Vision Transformers. This approach enhances
 the model’s ability to capture fine-grained local temporal
 dynamics without the need for task-specific temporal
 modules. We further combined this local gridification with
 a global absolute time modeling strategy, using composite
 text timestamps for each grid. Together, these components
 provide a stronger temporal encoding mechanism. Extensive
 experiments on standard VTG and VQA benchmarks demon-
 strate that T2SGrid achieves superior performance.

515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570

References

- [1] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015. 3
- [2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 1
- [3] Xiaoyi Bao, Chenwei Xie, Hao Tang, Tingyu Weng, Xiaofeng Wang, Yun Zheng, and Xingang Wang. Dyming: Key frames with visual prompts are good representation for multi-modal video understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23678–23688, 2025. 2, 3
- [4] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–970, 2015. 1, 5
- [5] Mu Cai, Haotian Liu, Siva Karthik Mustikovela, Gregory P Meyer, Yuning Chai, Dennis Park, and Yong Jae Lee. Vip-llava: Making large multimodal models understand arbitrary visual prompts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12914–12923, 2024. 2
- [6] Jingyuan Chen, Xinpeng Chen, Lin Ma, Zequn Jie, and Tat-Seng Chua. Temporally grounding natural sentence in video. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 162–171, Brussels, Belgium, 2018. Association for Computational Linguistics. 2
- [7] Lin Chen, Xilin Wei, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Zhenyu Tang, Li Yuan, et al. Sharegpt4video: Improving video understanding and generation with better captions. *Advances in Neural Information Processing Systems*, 37:19472–19495, 2024. 2
- [8] Yi-Wen Chen, Yi-Hsuan Tsai, and Ming-Hsuan Yang. End-to-end multi-modal video temporal grounding. *Advances in Neural Information Processing Systems*, 34:28442–28453, 2021. 2
- [9] Andong Deng, Zhongpai Gao, Anwesa Choudhuri, Benjamin Planche, Meng Zheng, Bin Wang, Terrence Chen, Chen Chen, and Ziyang Wu. Seq2time: Sequential knowledge transfer for video llm temporal grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13766–13775, 2025. 1
- [10] Lars Doorenbos, Federico Spurio, and Juergen Gall. Video panels for long video understanding, 2025. 2
- [11] Zhizhao Duan, Hao Cheng, Duo Xu, Xi Wu, Xiangxie Zhang, Xi Ye, and Zhen Xie. Cityllava: Efficient fine-tuning for vlms in city scenario, 2024. 2
- [12] Zheyu Fan, Jiateng Liu, Yuji Zhang, Zihan Wang, Yi R. Fung, Manling Li, and Heng Ji. Video-llms with temporal visual screening, 2025. 2
- [13] Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24108–24118, 2025. 1, 7
- [14] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. Tall: Temporal activity localization via language query. In *Proceedings of the IEEE international conference on computer vision*, pages 5267–5275, 2017. 1, 5
- [15] Junqi Ge, Ziyi Chen, Jintao Lin, Jinguo Zhu, Xihui Liu, Jifeng Dai, and Xizhou Zhu. V2pe: Improving multi-modal long-context capability of vision-language models with variable visual position encoding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 21070–21084, 2025. 1
- [16] Yongxin Guo, Jingyu Liu, Mingda Li, Qingbin Liu, Xi Chen, and Xiaoying Tang. Trace: Temporal grounding video llm via causal event modeling. *arXiv preprint arXiv:2410.05643*, 2024. 6, 8
- [17] Yongxin Guo, Jingyu Liu, Mingda Li, Dingxin Cheng, Xiaoying Tang, Dianbo Sui, Qingbin Liu, Xi Chen, and Kevin Zhao. Vtg-llm: Integrating timestamp knowledge into video llms for enhanced video temporal grounding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3302–3310, 2025. 1, 6
- [18] Miran Heo, Min-Hung Chen, De-An Huang, Sifei Liu, Subhashree Radhakrishnan, Seon Joo Kim, Yu-Chiang Frank Wang, and Ryo Hachiuma. Omni-rgpt: Unifying image and video region-level understanding via token marks. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 3919–3930, 2025. 1
- [19] Bin Huang, Xin Wang, Hong Chen, Zihan Song, and Wenwu Zhu. Vtimellm: Empower llm to grasp video moments, 2023. 2, 5, 6
- [20] De-An Huang, Shijia Liao, Subhashree Radhakrishnan, Hongxu Yin, Pavlo Molchanov, Zhiding Yu, and Jan Kautz. Lita: Language instructed temporal-localization assistant. In *European Conference on Computer Vision*, pages 202–218. Springer, 2024. 2, 6
- [21] Zeyi Huang, Yuyang Ji, Xiaofang Wang, Nikhil Mehta, Tong Xiao, Donghyun Lee, Sigmund Vanvalkenburgh, Shengxin Zha, Bolin Lai, Licheng Yu, Ning Zhang, Yong Jae Lee, and Miao Liu. Building a mind palace: Structuring environment-grounded semantic graphs for effective long video analysis with llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 24169–24179, 2025. 1
- [22] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 1, 6
- [23] Habib Irani and Vangelis Metsis. Positional encoding in transformer-based time series models: a survey. *arXiv preprint arXiv:2502.12370*, 2025. 1
- [24] Jinhyun Jang, Jungin Park, Jin Kim, Hyeongjun Kwon, and Kwanghoon Sohn. Knowing where to focus: Event-aware

744	[51] Qwen Team. Qwen3-vl: Sharper vision, deeper thought, broader action. <i>Qwen Blog</i> . Accessed, pages 10–04, 2025. 1, 2, 6	
745		
746		
747	[52] Haibo Wang, Zhiyang Xu, Yu Cheng, Shizhe Diao, Yufan Zhou, Yixin Cao, Qifan Wang, Weifeng Ge, and Lifu Huang. Grounded-videollm: Sharpening fine-grained temporal grounding in video large language models, 2025. 2	
748		
749		
750		
751	[53] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. <i>arXiv preprint arXiv:2409.12191</i> , 2024. 1, 2, 5, 6, 7	
752		
753		
754		
755		
756	[54] Yueqian Wang, Xiaojun Meng, Jianxin Liang, Yuxuan Wang, Qun Liu, and Dongyan Zhao. Hawkeye: Training video-text llms for grounding text in videos, 2024. 5, 6	
757		
758		
759	[55] Ye Wang, Ziheng Wang, Boshen Xu, Yang Du, Kejun Lin, Zihan Xiao, Zihao Yue, Jianzhong Ju, Liang Zhang, Dingyi Yang, et al. Time-r1: Post-training large vision language model for temporal video grounding. <i>arXiv preprint arXiv:2503.13377</i> , 2025. 2	
760		
761		
762		
763		
764	[56] Ziyang Wang, Shoubin Yu, Elias Stengel-Eskin, Jaehong Yoon, Feng Cheng, Gedas Bertasius, and Mohit Bansal. Videotree: Adaptive tree-based video representation for llm reasoning on long videos. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 3272–3283, 2025. 1	
765		
766		
767		
768		
769		
770	[57] Xilin Wei, Xiaoran Liu, Yuhang Zang, Xiaoyi Dong, Pan Zhang, Yuhang Cao, Jian Tong, Haodong Duan, Qipeng Guo, Jiaqi Wang, Xipeng Qiu, and Dahua Lin. Videorope: What makes for good video rotary position embedding?, 2025. 1	
771		
772		
773		
774		
775	[58] Junda Wu, Zhehao Zhang, Yu Xia, Xintong Li, Zhaoyang Xia, Aaron Chang, Tong Yu, Sungchul Kim, Ryan A Rossi, Ruiyi Zhang, et al. Visual prompting in multimodal large language models: A survey. <i>arXiv preprint arXiv:2409.15310</i> , 2024. 2	
776		
777		
778		
779		
780	[59] Mingrui Wu, Xinyue Cai, Jiayi Ji, Jiale Li, Oucheng Huang, Gen Luo, Hao Fei, Guannan Jiang, Xiaoshuai Sun, and Rongrong Ji. Controlmllm: Training-free visual prompt learning for multimodal large language models. In <i>Advances in Neural Information Processing Systems</i> , pages 45206–45234. Curran Associates, Inc., 2024. 2	
781		
782		
783		
784		
785		
786	[60] Tung-Yu Wu, Sheng-Yu Huang, and Yu-Chiang Frank Wang. Data-efficient 3d visual grounding via order-aware referring. In <i>2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)</i> , pages 3107–3117. IEEE, 2025. 2	
787		
788		
789		
790	[61] Yongliang Wu, Xinting Hu, Yuyang Sun, Yizhou Zhou, Wenbo Zhu, Fengyun Rao, Bernt Schiele, and Xu Yang. Number it: Temporal grounding videos like flipping manga. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 13754–13765, 2025. 2, 6, 8	
791		
792		
793		
794		
795	[62] Huijuan Xu, Kun He, Leonid Sigal, Stan Sclaroff, and Kate Saenko. Text-to-clip video retrieval with early fusion and recaptioning. <i>arXiv preprint arXiv:1804.05113</i> , 2(6):7, 2018. 2	
796		
797		
798		
799	[63] Mingze Xu, Mingfei Gao, Zhe Gan, Hong-You Chen, Zhengfeng Lai, Haiming Gang, Kai Kang, and Afshin De-	
800		
	ghan. Slowfast-llava: A strong training-free baseline for video large language models, 2024. 1	801
		802
[64]	Weijie Xu, Jingwei Tan, Shulin Wang, and Sheng Yang. Temporal relation-aware global attention network for temporal action detection. In <i>International Conference on Intelligent Computing</i> , pages 257–269. Springer, 2024. 2	803
		804
		805
		806
[65]	Yifang Xu, Yunzhuo Sun, Zien Xie, Benxiang Zhai, and Sidan Du. Vtg-gpt: Tuning-free zero-shot video temporal grounding with gpt. <i>Applied Sciences</i> , 14(5):1894, 2024. 2	807
		808
		809
[66]	Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. Tubedetr: Spatio-temporal video grounding with transformers. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 16442–16453, 2022. 2	810
		811
		812
		813
[67]	Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. <i>arXiv preprint arXiv:2310.11441</i> , 2023. 2	814
		815
		816
		817
		818
[68]	Lingfeng Yang, Yueze Wang, Xiang Li, Xinlong Wang, and Jian Yang. Fine-grained visual prompting, 2023. 2	819
		820
[69]	Zaiquan Yang, Yuhao Liu, Gerhard Hancke, and Rynson W. H. Lau. Unleashing the potential of multimodal llms for zero-shot spatio-temporal video grounding, 2025. 1	821
		822
		823
[70]	Yuan Yao, Ao Zhang, Zhengyan Zhang, Zhiyuan Liu, Tat-Seng Chua, and Maosong Sun. Cpt: Colorful prompt tuning for pre-trained vision-language models. <i>AI Open</i> , 5:30–38, 2024. 2	824
		825
		826
		827
[71]	Jinhui Ye, Zihan Wang, Haosen Sun, Keshigeyan Chandrasegaran, Zane Durante, Cristobal Eyzaguirre, Yonatan Bisk, Juan Carlos Niebles, Ehsan Adeli, Li Fei-Fei, Jiajun Wu, and Manling Li. Re-thinking temporal search for long-form video understanding. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 8579–8591, 2025. 2	828
		829
		830
		831
		832
		833
		834
[72]	Xiangyu Zeng, Kunchang Li, Chenting Wang, Xinhao Li, Tianxiang Jiang, Ziang Yan, Songze Li, Yansong Shi, Zhengrong Yue, Yi Wang, Yali Wang, Yu Qiao, and Limin Wang. Timesuite: Improving mllms for long video understanding via grounded tuning, 2025. 2, 6	835
		836
		837
		838
		839
[73]	Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding, 2023. 2	840
		841
		842
[74]	Hao Zhang, Aixin Sun, Wei Jing, and Joey Tianyi Zhou. Temporal sentence grounding in videos: A survey and future directions. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 45(8):10443–10465, 2023. 2	843
		844
		845
		846
[75]	Jinglei Zhang, Yuanfan Guo, Rolandos Alexandros Potamias, Jiankang Deng, Hang Xu, and Chao Ma. Vtimecot: Thinking by drawing for video temporal grounding and reasoning. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 24203–24213, 2025. 1	847
		848
		849
		850
		851
		852
[76]	Tao Zhang, Xiangtai Li, Hao Fei, Haobo Yuan, Shengqiong Wu, Shunping Ji, Chen Change Loy, and Shuicheng Yan. Omg-llava: Bridging image-level, object-level, pixel-level reasoning and understanding. In <i>Advances in Neural Information Processing Systems</i> , pages 71737–71767. Curran Associates, Inc., 2024. 1	853
		854
		855
		856
		857
		858

- 859 [77] Zhu Zhang, Zhou Zhao, Yang Zhao, Qi Wang, Huasheng
860 Liu, and Lianli Gao. Where does it exist: Spatio-temporal
861 video grounding for multi-form sentences. In *Proceedings*
862 *of the IEEE/CVF Conference on Computer Vision and Pat-*
863 *tern Recognition (CVPR)*, 2020. 2
- 864 [78] Heng Zhao, Zhao Yinjie, Bihan Wen, Yew-Soon Ong, and
865 Joey Tianyi Zhou. Video-text prompting for weakly super-
866 vised spatio-temporal video grounding. In *Proceedings of*
867 *the 2024 Conference on Empirical Methods in Natural Lan-*
868 *guage Processing*, pages 19494–19505, 2024. 2
- 869 [79] Minghang Zheng, Xinhao Cai, Qingchao Chen, Yuxin Peng,
870 and Yang Liu. Training-free video temporal grounding using
871 large-scale pre-trained models. In *European Conference on*
872 *Computer Vision*, pages 20–37. Springer, 2024. 1
- 873 [80] Minghang Zheng, Xinhao Cai, Qingchao Chen, Yuxin Peng,
874 and Yang Liu. Training-free video temporal grounding using
875 large-scale pre-trained models. In *Computer Vision – ECCV*
876 *2024*, pages 20–37, Cham, 2025. Springer Nature Switzer-
877 land. 1, 2
- 878 [81] Sheng Zhou, Junbin Xiao, Qingyun Li, Yicong Li, Xun
879 Yang, Dan Guo, Meng Wang, Tat-Seng Chua, and Angela
880 Yao. Egotextvqa: Towards egocentric scene-text aware video
881 question answering. In *Proceedings of the Computer Vi-*
882 *sion and Pattern Recognition Conference*, pages 3363–3373,
883 2025. 2
- 884 [82] Yunsong Zhou, Linyan Huang, Qingwen Bu, Jia Zeng,
885 Tianyu Li, Hang Qiu, Hongzi Zhu, Minyi Guo, Yu Qiao,
886 and Hongyang Li. Embodied understanding of driving sce-
887 narios. In *European Conference on Computer Vision*, pages
888 129–148. Springer, 2024. 2
- 889 [83] Kangyu Zhu, Ziyuan Qin, Huahui Yi, Zekun Jiang, Qicheng
890 Lao, Shaoting Zhang, and Kang Li. Guiding medical vision-
891 language models with explicit visual prompts: Framework
892 design and comprehensive exploration of prompt variations,
893 2025. 2