

# Harnessing Object Grounding for Time-Sensitive Video Understanding

Tz-Ying Wu

Sharath Nittur Sridhar

Subarna Tripathi

Intel



# Time-Sensitive Video Understanding (TSV)

- TSV tasks are challenging as they require **fine-grained** scene comprehension with temporal reasoning.

## Reasoning Temporal Localization (RTL)



**Q:** when does the woman's dance become the most **energetic** in the video?

**A:** Her dance is the most energetic between 7.37s and 12.81s when she does a **handspring**, which is more energetic than **standing** and **lying** on the floor.

[1] Huang et al., "LITA: Language Instructed Temporal-Localization Assistant". In ECCV 2024.

**Video:** [Timeline of video frames from 02:42 to 07:28]

**Timestamps:** 02:42, 03:17, 03:36, 03:56, 04:40, 05:05, 05:54, 06:14, 06:34, 06:53, 07:28

Task	User Query	TimeChat Response
Dense Video Captioning	Watch the video and <b>extract a maximum of 10 cooking steps</b> . For each step, determine the starting and ending times and provide a concise description.	<b>135 - 175 seconds</b> , trim fat from pork and cut it to slices. <b>180 - 209 seconds</b> , cut one green onion into 1 inch size. <b>214 - 225 seconds</b> , cut a red onion into strips. <b>226 - 242 seconds</b> , chop a chinese broccoli. <b>245 - 267 seconds</b> , cut the mushroom into thin slices. <b>299 - 323 seconds</b> , boil some udon noodles in water. <b>345 - 362 seconds</b> , boil the chinese broccoli and drain it. <b>412 - 419 seconds</b> , add the udon noodles and cook. <b>425 - 449 seconds</b> , add the premade sauce and stir.
Temporal Video Grounding	In which video clip does the <b>sauce mixing step</b> occur? Please provide the starting and ending times for that step.	The specific instruction for this step is to <b>combine soy sauce, water, and sugar in a bowl</b> from <b>274 - 290 seconds</b> .
Video Highlight Detection	Find the <b>highlight contents</b> in the video corresponding to " <b>frying the pork</b> ", determining the highlight timestamps and its saliency score on a scale from 1 to 5.	The highlight timestamps are in the <b>369, 371, 373, 375, 377, 379, 381, 383, 385, 387, 389 seconds</b> . Their saliency scores are <b>1.4, 2.8, 3.8, 4.0, 4.0, 4.0, 3.7, 4.0, 3.5, 4.0, 3.3</b> .

[2] Ren et al., "TimeChat: A Time-sensitive Multimodal Large Language Model for Long Video Understanding". In CVPR 2024.

# Motivation

- Video-LLMs compress frames into **highly condensed video tokens**.
- This works for global understanding but is **insufficient for TSV**.
- Objects provide key temporal cues yet are not modeled explicitly.
- We ask: *Can **grounded object (GO)** information improve Video-LLMs on TSV tasks?*

# Probing Video-LLMs with GO Information

*(i) Can providing GO information in Video-LLMs enhance TSV?*

*(ii) What levels of GO information is needed?*

As a preliminary experiment, we probe existing Video-LLMs with object cues by injecting the GO information into the text prompt:

- **Class:** “*<Obj> Objects in this video are: man, window, ...</Obj>*”.
- **Class+Time:** “*<Obj> Each object is provided with its timestamp and class label in the format of <time, class label>. Here are the objects: <91.2 second, man>, ...</Obj>*”.
- **Class+Time+Bbox:** “*<Obj> Each object bounding box is provided with its timestamp and class label in the format of <time, (x1,y1,x2,y2), class label>. Here are the objects: <91.2 second, (0.0001, 0.1715, 0.0806, 0.3784), man>, ...</Obj>*”.

# Probing Video-LLMs with GO Information

- **Object cues help**, and more detailed cues yield larger gains.
- However, this naïve approach:
  - Is **not robust** to object grounding errors.
  - Introduces a **large** number of **extra tokens** (~42 tokens per object).

GO at Inference	mIOU	P@0.5
n/a	23.04	18.78
Text (Class)	23.65	18.08
Text (Class+Time)	26.47	20.47
Text (Class+Time+Bbox)	27.63	23.71

**Table 1:** Probing LITA-13B on the RTL task with different levels of grounded object information.

Flip Ratio	mIOU	P@0.5
10%	27.07	23.25
20%	25.89	22.63
50%	24.63	20.13

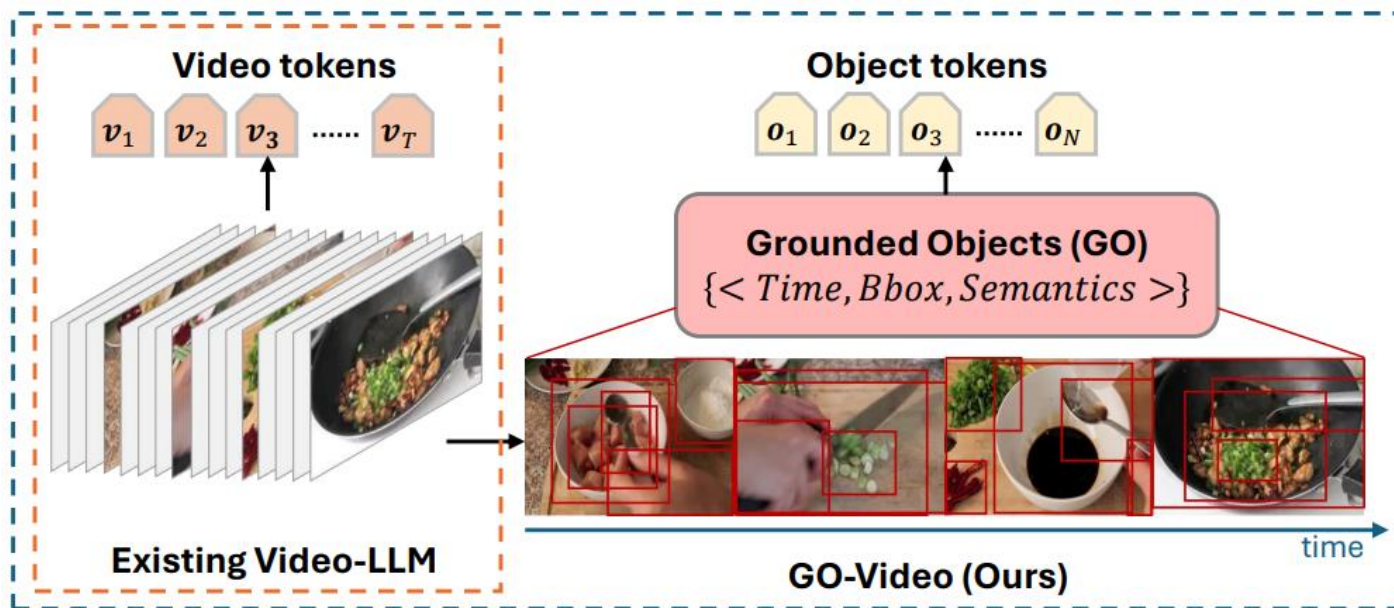
**Table 2:** Flipping a fraction of object class labels with random classes.

Shift	mIOU	P@0.5
$0.01 \times (H, W)$	27.37	21.13
$0.02 \times (H, W)$	26.42	19.01
$0.05 \times (H, W)$	24.71	17.14

**Table 3:** Jittering bounding box locations wrt. the image height ( $H$ ) and width ( $W$ ) in pixels.

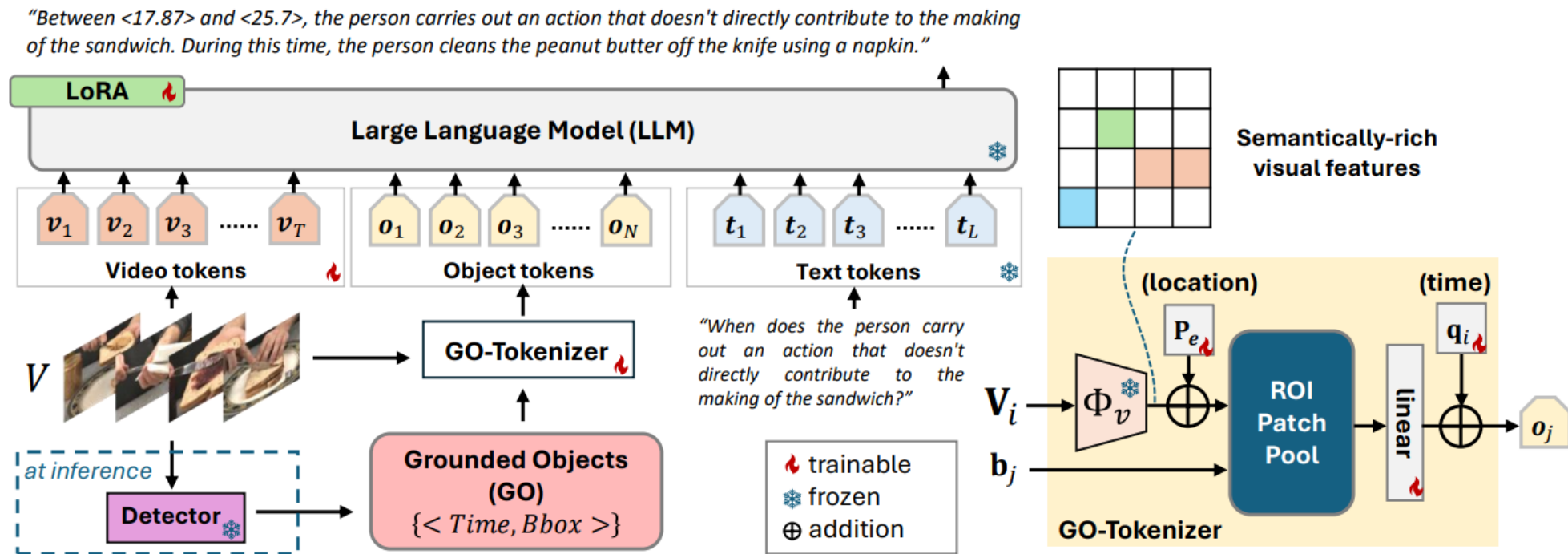
# Augmenting Video-LLMs with Object Tokens

- We introduce **object tokens** to complement **video tokens**.
- They are **compact** representations encapsulating each object's *spatial*, *semantics*, and *temporal* information into a **single** token.



# Grounded Object Tokenization

- To extract **object tokens**, we propose **Grounded Object Tokenizer (GO-Tokenizer)**, a lightweight add-on for existing Video-LLMs.



# Results

- **Text** does not yield consistent gain when introducing GO in prompt.
- In contrast, **GO-(LITA/TimeChat)** consistently outperforms.

Model	GO Training	mIOU	P@0.5	Rel. Score
LITA	✗	28.28	26.40	54.7
LITA + Text ( <i>Class+Time+Bbox</i> )	✗	28.07	25.83	54.3
<i>Finetuned on ActivityNet-Captions &amp; ActivityNet-RTL</i>				
LITA	✗	28.72	24.53	54.7
LITA + Text ( <i>Class+Time+Bbox</i> )	✗	28.82	25.10	55.6
GO-LITA	✓	<b>31.52</b>	<b>28.07</b>	<b>59.0</b>

**Table 4:** Evaluating LITA-13B on the ActivityNet-RTL dataset with/without GO information.

Model	GO Training	CIDEr	SODAc	F1
Valley [16]	✗	0.0	0.1	1.5
Video-LLaMA [27]	✗	0.0	0.0	0.1
TimeChat [22]	✗	3.4	1.2	12.6
<i>Finetuned on ActivityNet-Captions</i>				
TimeChat	✗	3.4	1.1	11.3
TimeChat + Text ( <i>Class+Time+Bbox</i> )	✗	2.4	0.8	10.6
GO-TimeChat	✓	<b>3.9</b>	<b>1.4</b>	<b>18.5</b>

**Table 5:** Zero-shot evaluation on Youcook2 dataset for dense video captioning.

# Results

- Text with **GO Training** still underperforms **GO-(LITA/TimeChat)**.

Model	Go Training	Tokens/GO	CIDEr	SODAc	F1
TimeChat	✗	0	3.4	1.1	11.3
TimeChat + Text ( <i>Class</i> )	✓	4	1.8	0.8	16.2
TimeChat + Text ( <i>Class+Time</i> )	✓	12	1.8	0.9	16.8
TimeChat + Text ( <i>Class+Time+Bbox</i> )	✓	42	1.6	0.8	13.0
GO-TimeChat	✓	1	<b>3.9</b>	<b>1.4</b>	<b>18.5</b>

**Table 6:** Comparison to Text with GO Training using TimeChat as the base model.

Model	Go Training	mIOU	P@0.5	Rel. Score
LITA	✗	28.72	24.53	54.7
LITA + Text ( <i>Class+Time+Bbox</i> )	✓	30.05	27.97	55.5
GO-LITA	✓	<b>31.52</b>	<b>28.07</b>	<b>59.0</b>

**Table 7:** Comparison to Text with GO Training using LITA as the base model.

# Visualization

- Video-LLMs with **GO-Tokenizer** generates captions with higher quality compared to **Text**.



***Thank you!***

**Harnessing Object Grounding for Time-Sensitive  
Video Understanding**

Tz-Ying Wu   Sharath Nittur Sridhar   Subarna Tripathi

Intel

