

# Text-Visual Prompting for Efficient 2D Temporal Video Grounding

**Yimeng Zhang<sup>1,2</sup>, Xin Chen<sup>2</sup>, Jinghan Jia<sup>1</sup>, Sijia Liu<sup>1</sup>, Ke Ding<sup>2</sup>**

<sup>1</sup> OPTML Lab, Michigan State University

<sup>2</sup> Applied ML, Intel



# 1. What is Temporal Video Grounding (TVG)?

## 1. What is Temporal Video Grounding (TVG)?

Temporal Video Grounding is to **match a descriptive sentence** with one segment ( or moment ) in an untrimmed video that is of the same semantics.

## 1. What is Temporal Video Grounding (TVG)?

Temporal Video Grounding is to **match a descriptive sentence** with one segment ( or moment ) in an untrimmed video that is of the same semantics.

→ **Input:**    Descriptive Sentence    +    One Untrimmed Video

## 1. What is Temporal Video Grounding (TVG)?

Temporal Video Grounding is to **match a descriptive sentence** with one segment ( or moment ) in an untrimmed video that is of the same semantics.

- **Input:** Descriptive Sentence + One Untrimmed Video
- **Output:** Starting and ending points of the target segment

## 1. What is Temporal Video Grounding (TVG)?

Temporal Video Grounding is to **match a descriptive sentence** with one segment ( or moment ) in an untrimmed video that is of the same semantics.

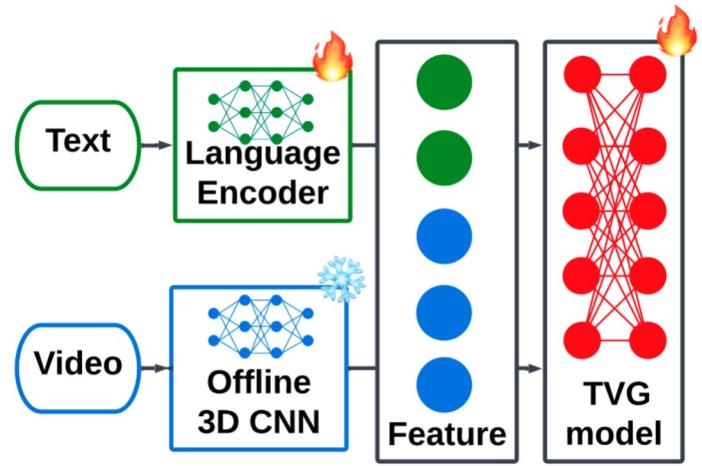
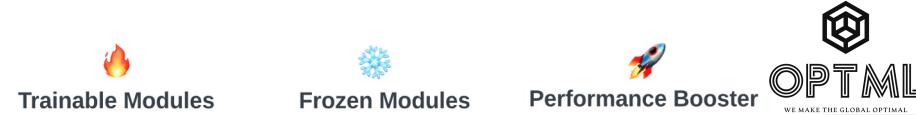
- **Input:** Descriptive Sentence + One Untrimmed Video
- **Output:** Starting and ending points of the target segment

### [ Existing Works ]

- Two-Stage “propose-and-rank” (Proposal-based)
- Regression-based (Proposal-free)
- Reinforcement learning-based

## 2. Temporal Video Grounding (TVG) Method Comparison

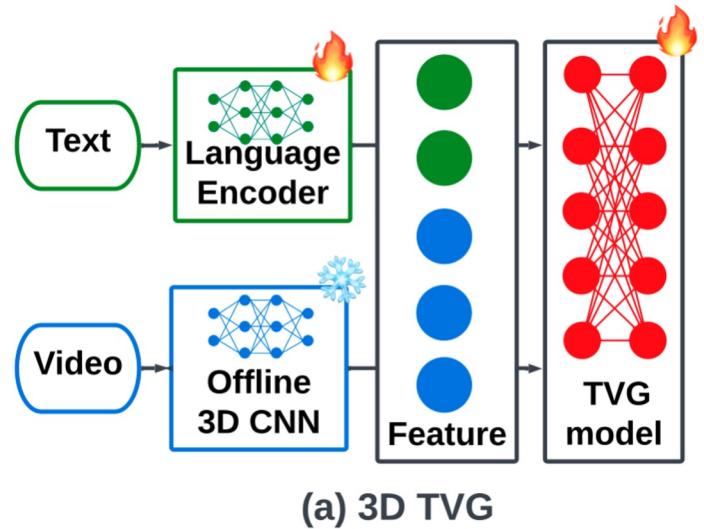
## 2. Temporal Video Grounding (TVG) Method Comparison



(a) 3D TVG

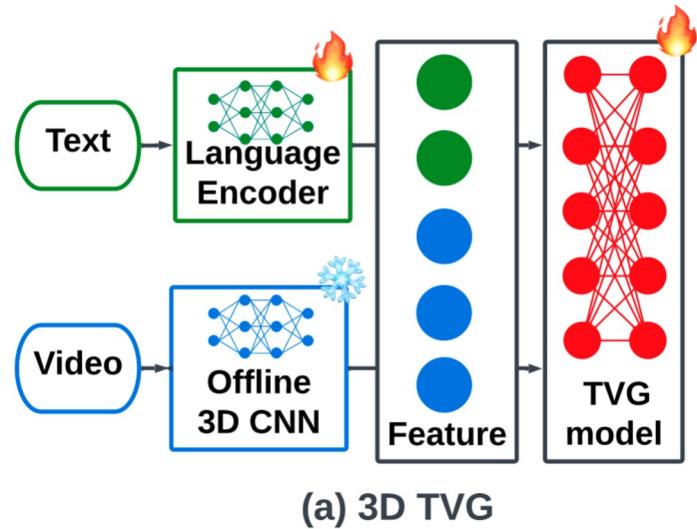
a) 3D TVG

## 2. Temporal Video Grounding (TVG) Method Comparison



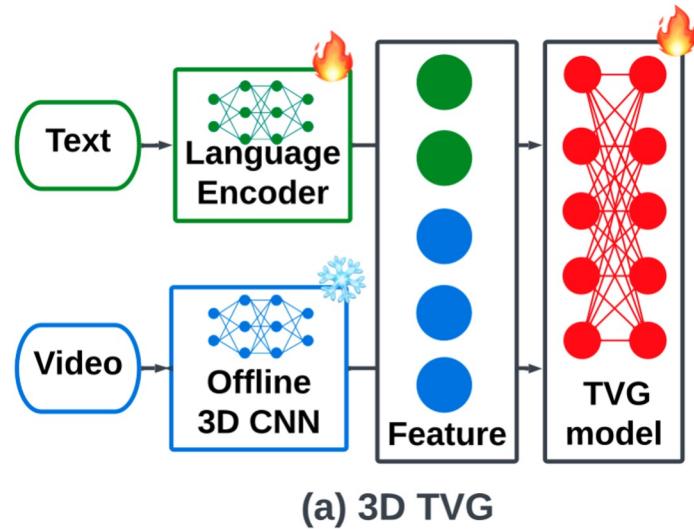
- a) 3D TVG
- Using offline 3D CNN as the video encoder.

## 2. Temporal Video Grounding (TVG) Method Comparison



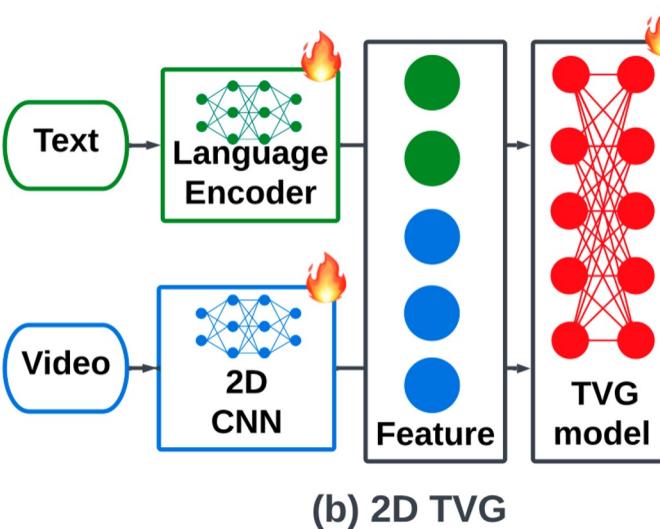
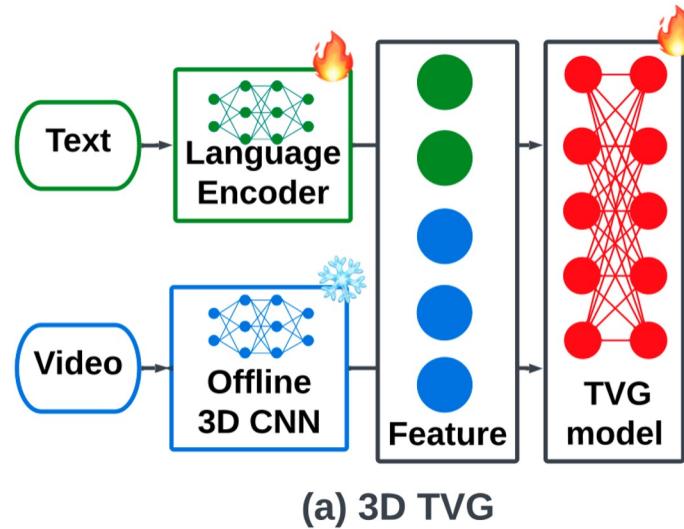
- a) 3D TVG
  - Using offline 3D CNN as the video encoder.
  - During training, **3D-CNN parameters are fixed**, which means modules for text and video processing **cannot be co-trained** for better multimodal feature fusion.

## 2. Temporal Video Grounding (TVG) Method Comparison



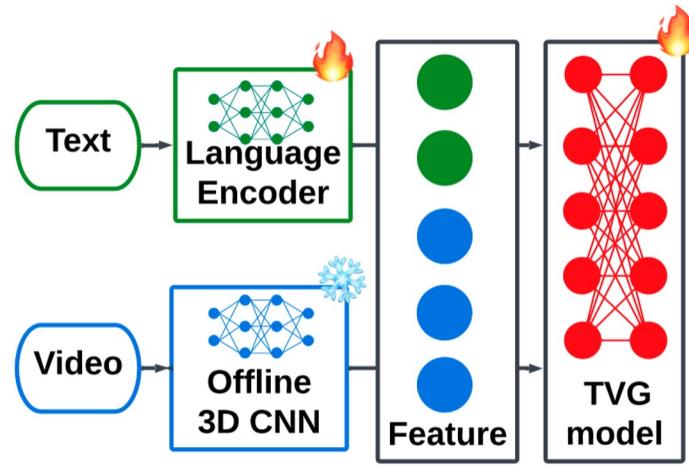
- a) 3D TVG
  - Using offline 3D CNN as the video encoder.
  - During training, **3D-CNN parameters are fixed**, which means modules for text and video processing **cannot be co-trained** for better multimodal feature fusion.
  - It is challenging to train 3D-CNNs, which is why most methods **do not involve 3D-CNNs during training and directly utilize the video features** extracted by offline 3D-CNNs as the video input.

## 2. Temporal Video Grounding (TVG) Method Comparison

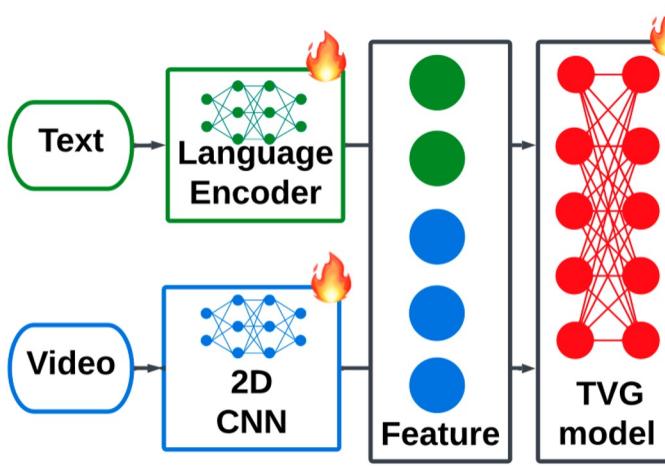


- a) 3D TVG
  - Using offline 3D CNN as the video encoder.
  - During training, **3D-CNN parameters are fixed**, which means modules for text and video processing **cannot be co-trained** for better multimodal feature fusion.
  - It is challenging to train 3D-CNNs, which is why most methods **do not involve 3D-CNNs during training and directly utilize the video features** extracted by offline 3D-CNNs as the video input.
  
- a) 2D TVG
  - Using 2D CNN as the video encoder.

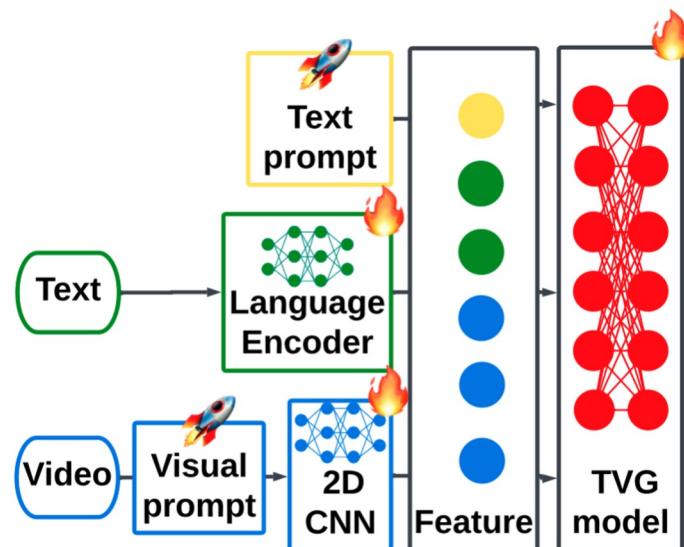
## 2. Temporal Video Grounding (TVG) Method Comparison



(a) 3D TVG



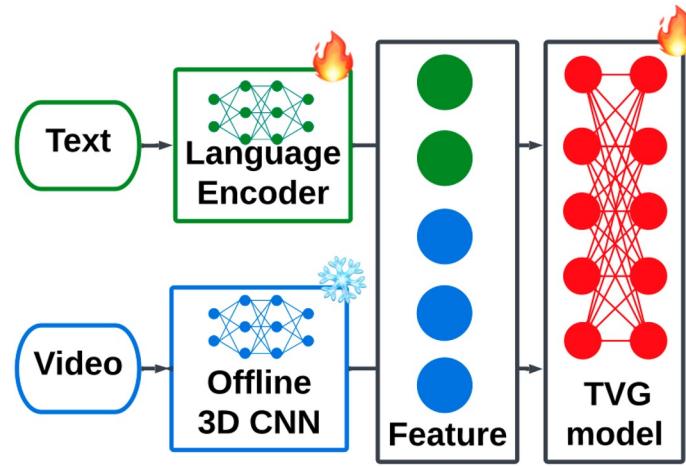
(b) 2D TVG



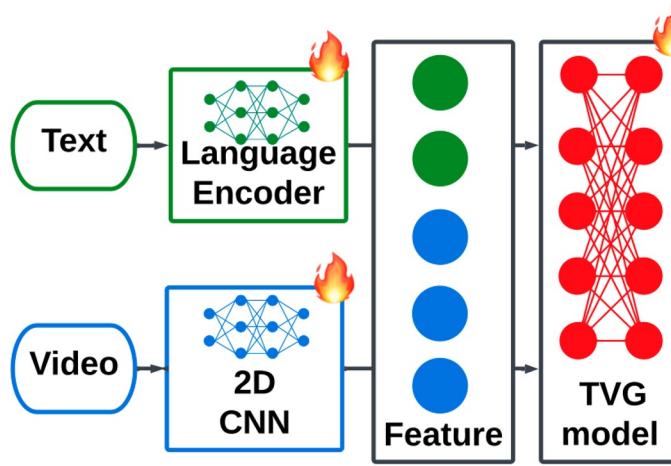
(c) TVP-based 2D TVG (Ours)

- a) **3D TVG**
  - Using offline 3D CNN as the video encoder.
  - During training, **3D-CNN parameters are fixed**, which means modules for text and video processing **cannot be co-trained** for better multimodal feature fusion.
  - It is challenging to train 3D-CNNs, which is why most methods **do not involve 3D-CNNs during training and directly utilize the video features** extracted by offline 3D-CNNs as the video input.
- a) **2D TVG**
  - Using 2D CNN as the video encoder.
- a) **TVP-Based 2D TVG**
  - The proposed **text-visual prompts (TVP)** compensate for the lack of spatiotemporal information in 2D CNNs for visual feature extraction.

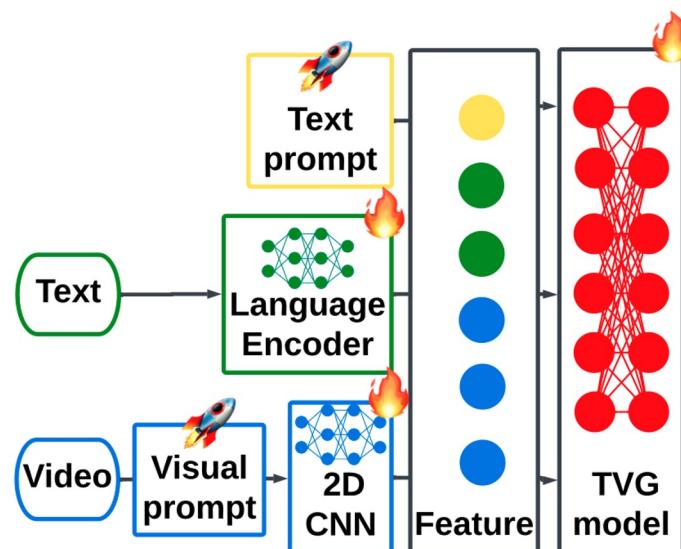
## 2. Temporal Video Grounding (TVG) Method Comparison



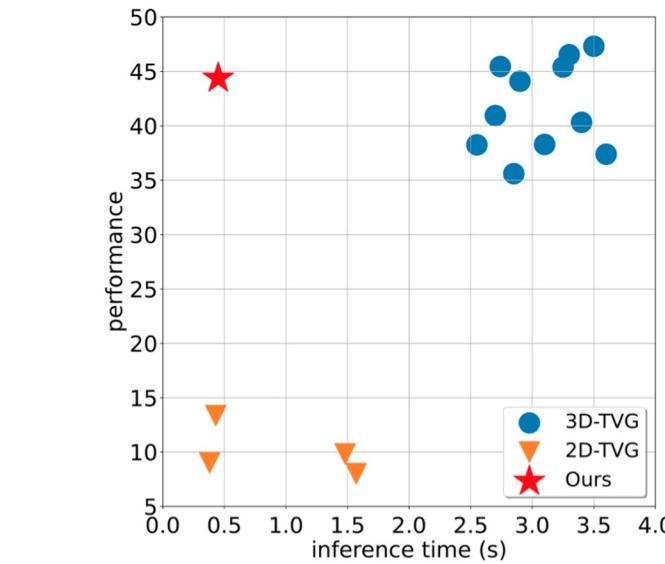
(a) 3D TVG



(b) 2D TVG



(c) TVP-based 2D TVG (Ours)



(d) Overall performance comparison

### a) 3D TVG

- Using offline 3D CNN as the video encoder
- During training, **3D-CNN parameters are fixed**, which means modules for text and video processing **cannot be co-trained** for better multimodal feature fusion.
- It is challenging to train 3D-CNNs, which is why most methods **do not involve 3D-CNNs during training and directly utilize the video features** extracted by offline 3D-CNNs as the video input.

### a) 2D TVG

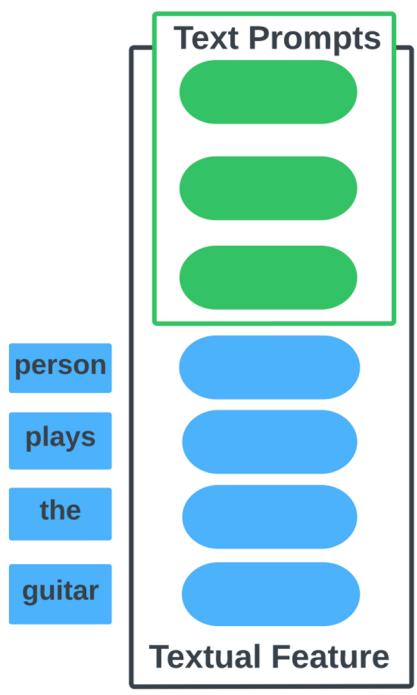
- Using 2D CNN as the video encoder.

### a) TVP-Based 2D TVG

- The proposed **text-visual prompts (TVP)** compensate for the lack of spatiotemporal information in 2D CNNs for visual feature extraction.

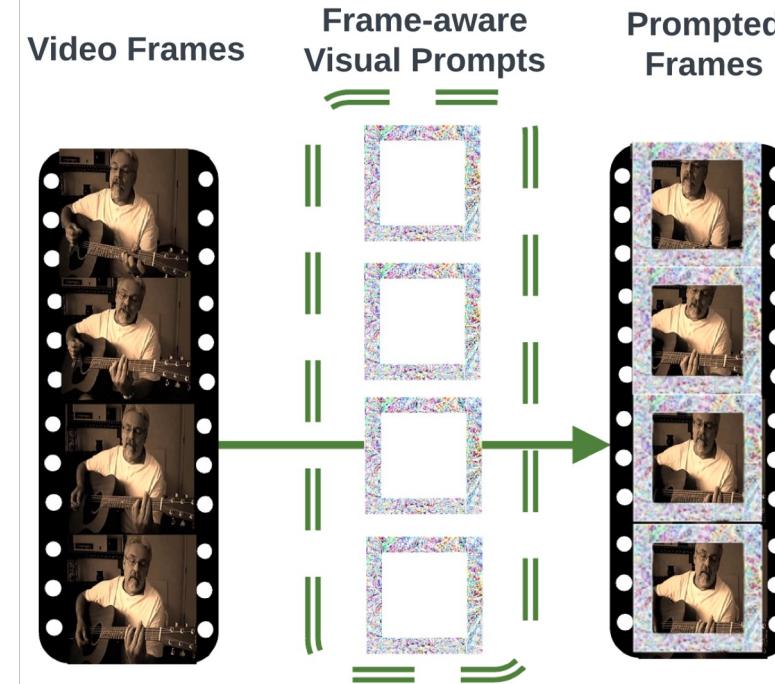
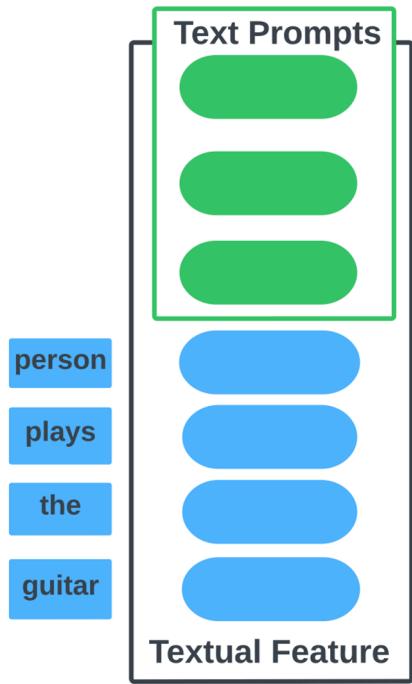
### 3. Text Prompts and Visual Prompts

### 3. Text Prompts and Visual Prompts



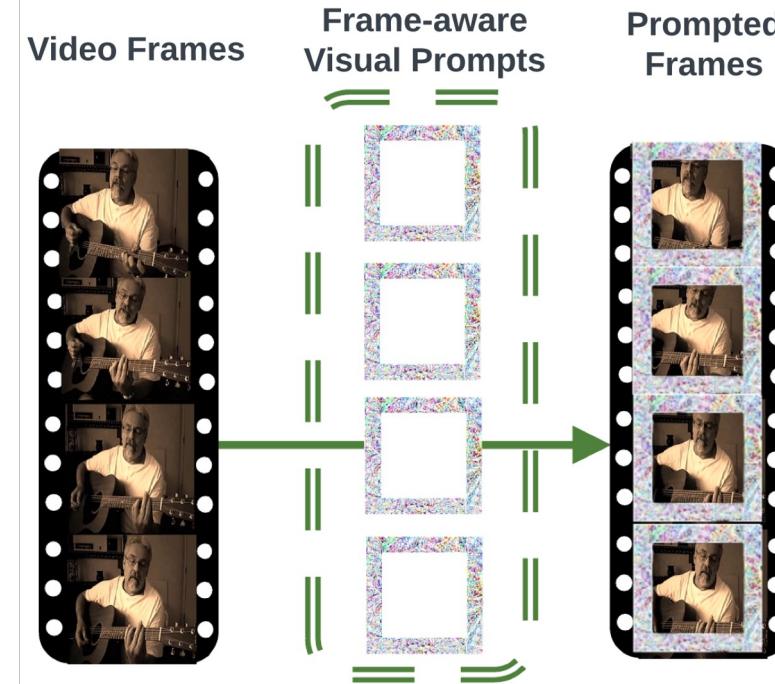
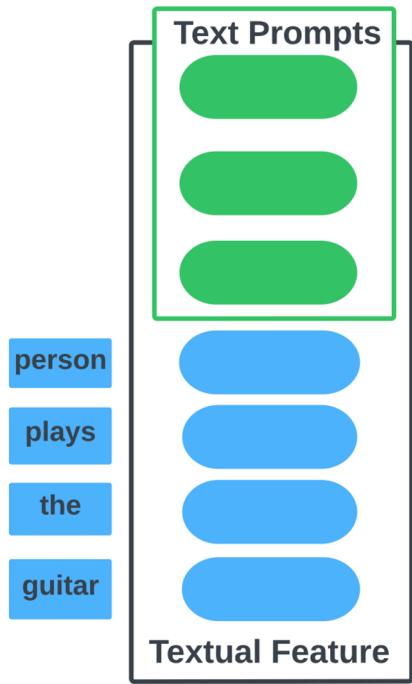
- Text prompts are directly applied in the feature space.

### 3. Text Prompts and Visual Prompts



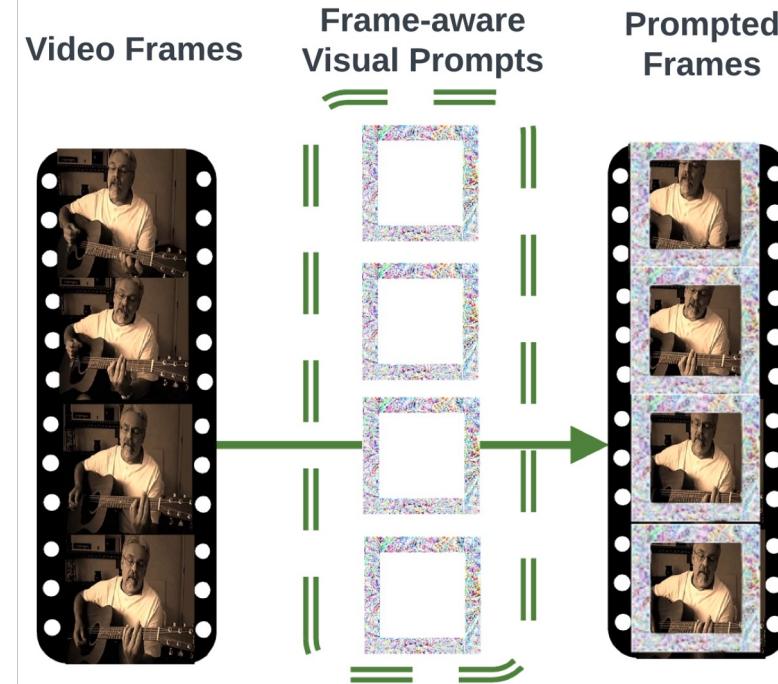
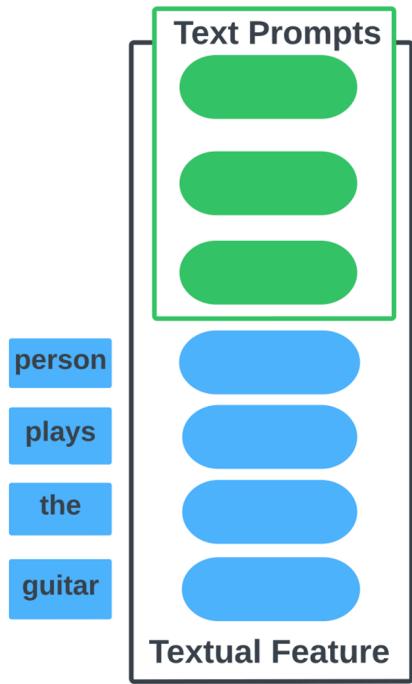
- Text prompts are directly applied in the feature space.
- A set of frame-aware visual prompts are applied to pixel space of video frames in order.

### 3. Text Prompts and Visual Prompts



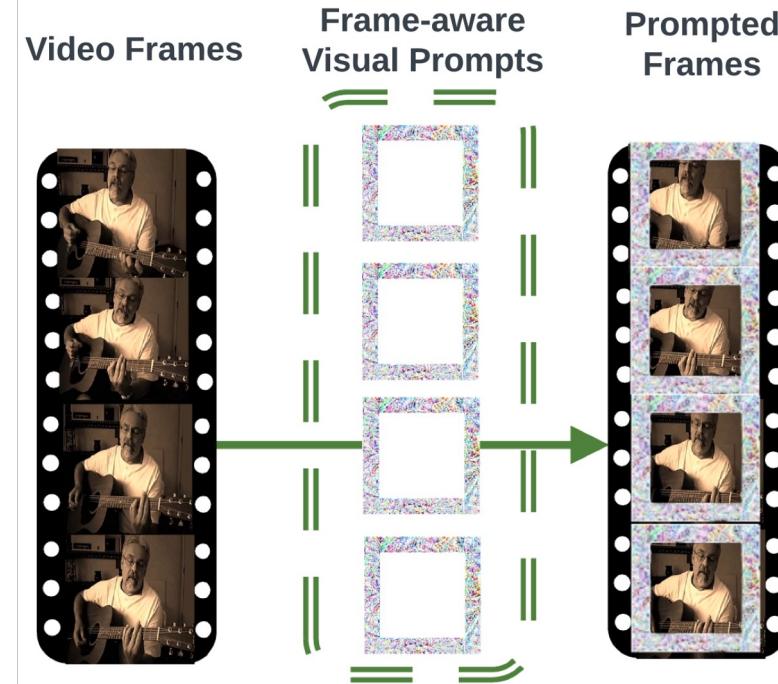
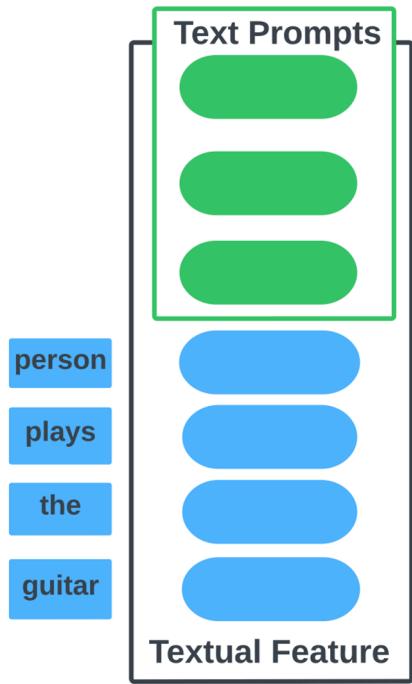
- Text prompts are directly applied in the feature space.
- A set of frame-aware visual prompts are applied to pixel space of video frames in order.
- During training, only the set of visual prompts and text prompts are updated through backpropagation.

### 3. Text Prompts and Visual Prompts



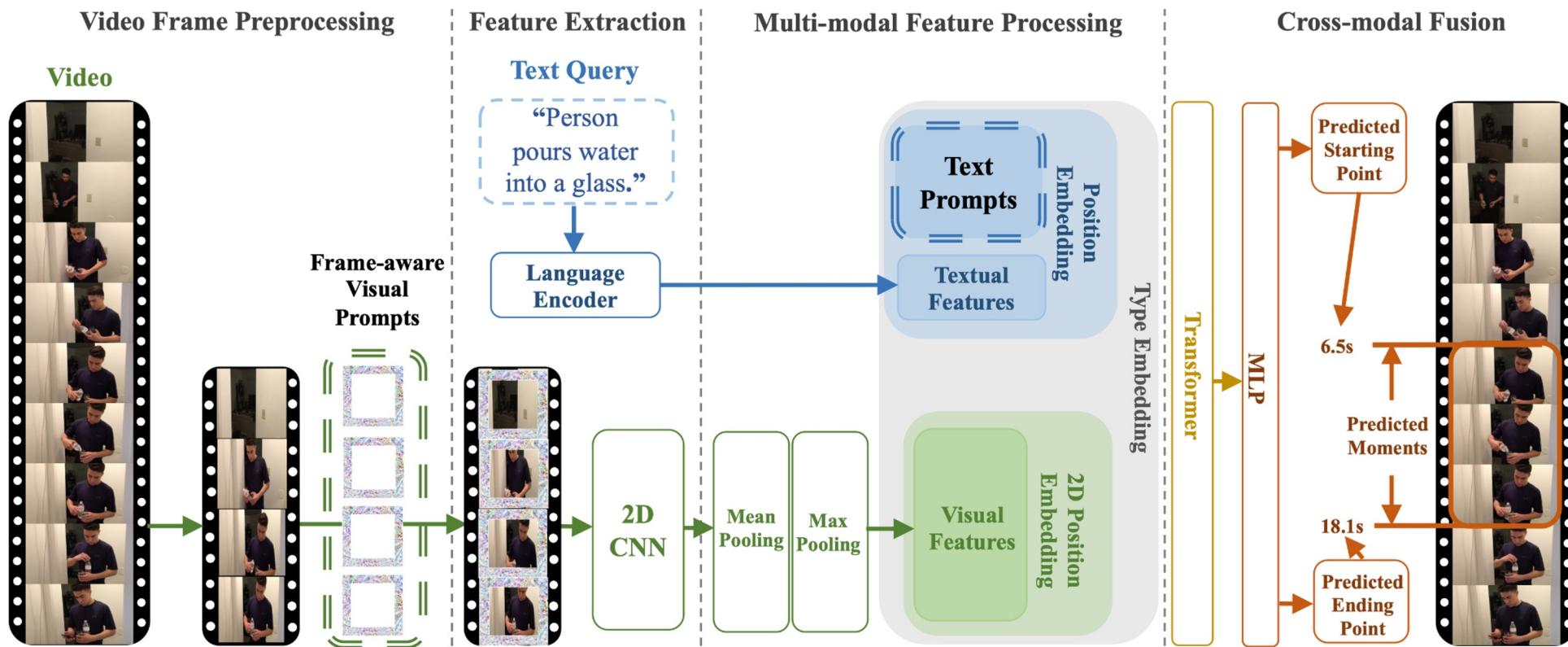
- Text prompts are directly applied in the feature space.
- During training, only the set of visual prompts and text prompts are updated through backpropagation.
- During finetuning, prompts are frozen, and the parameters of the TVG model and encoders are updated.

### 3. Text Prompts and Visual Prompts

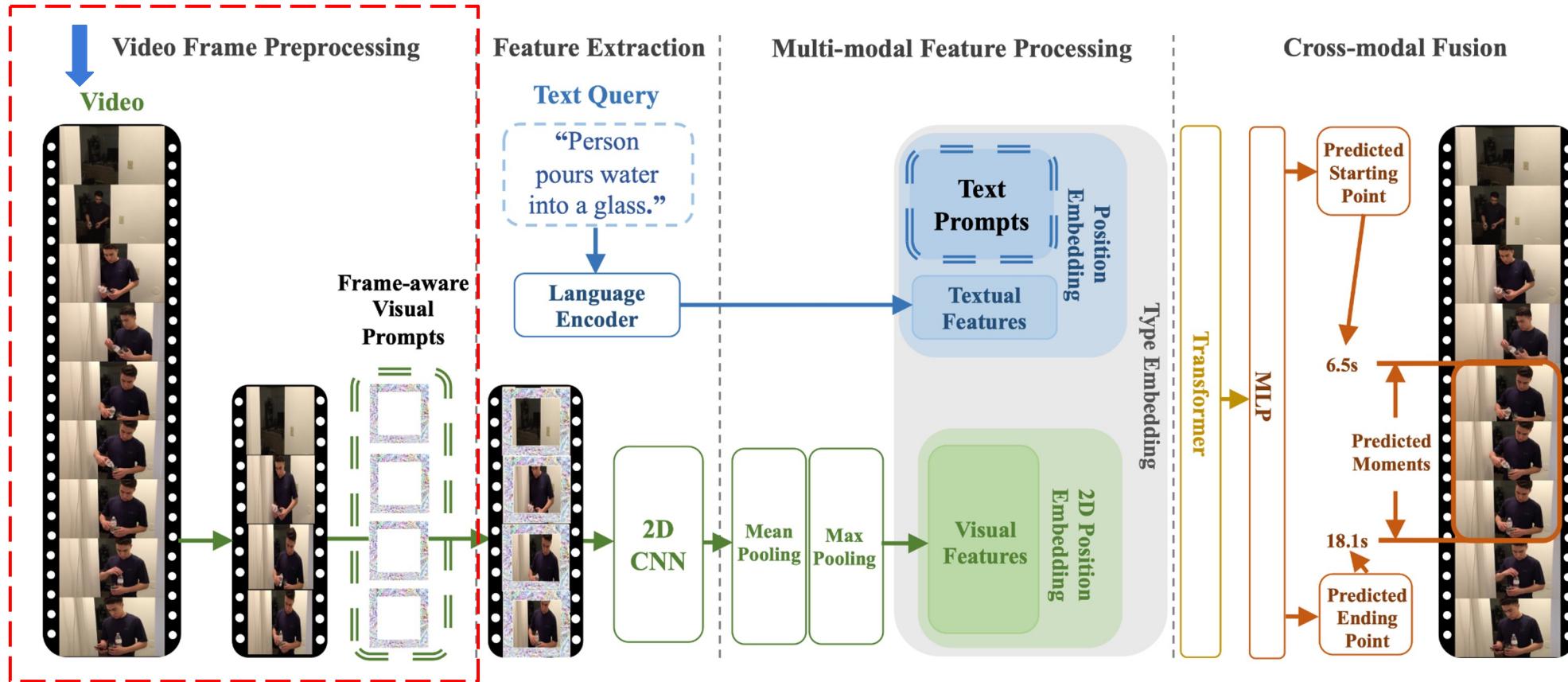


- Text prompts are directly applied in the feature space.
- a set of frame-aware visual prompts are applied to pixel space of video frames in order.
- During training, only the set of visual prompts and text prompts are updated through backpropagation.
- During finetuning, prompts are frozen, and the parameters of the TVG model and encoders are updated.
- During testing, the set of optimized visual prompts and the optimized text prompts are **applied to all test-time video-query pairs**.

## 4. Text-Visual Prompt for 2D TVG

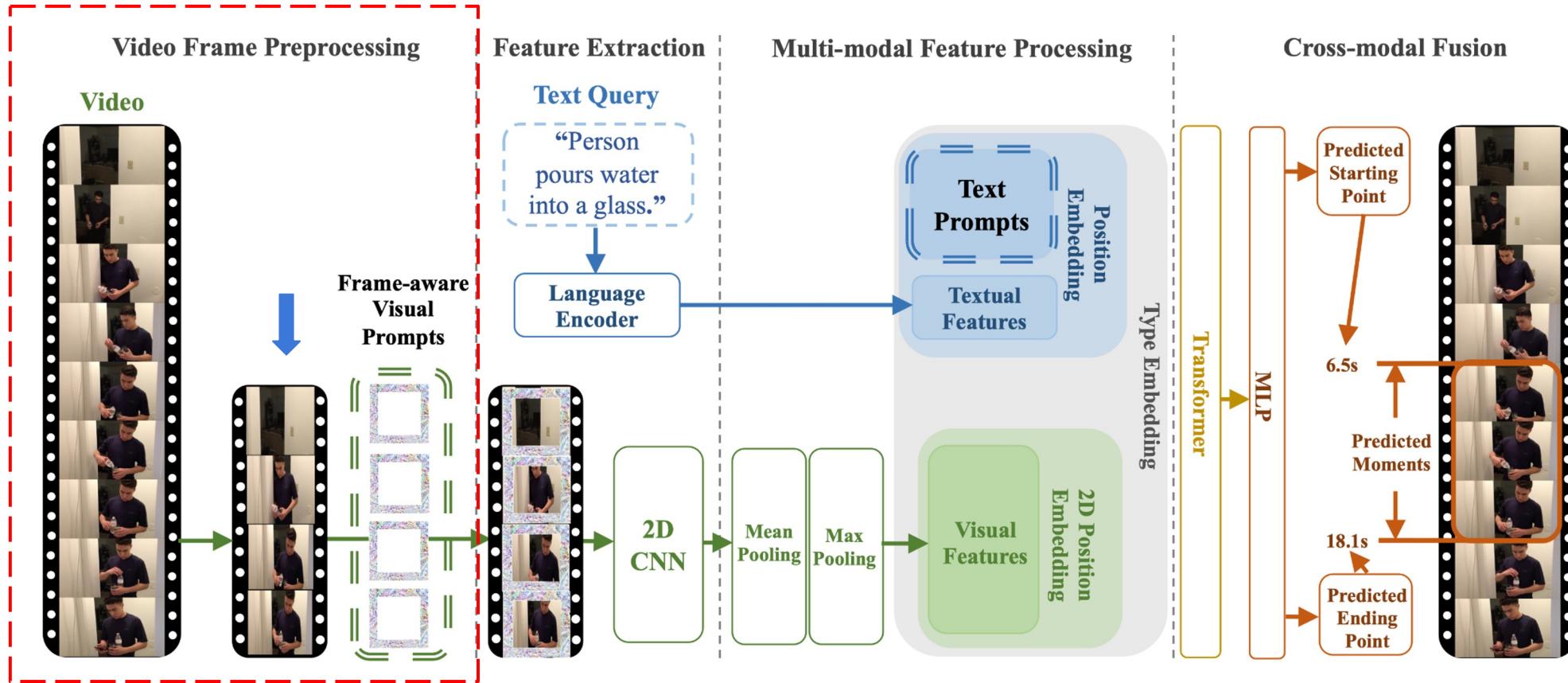


## 4. Text-Visual Prompt for 2D TVG



Video frame preprocessing

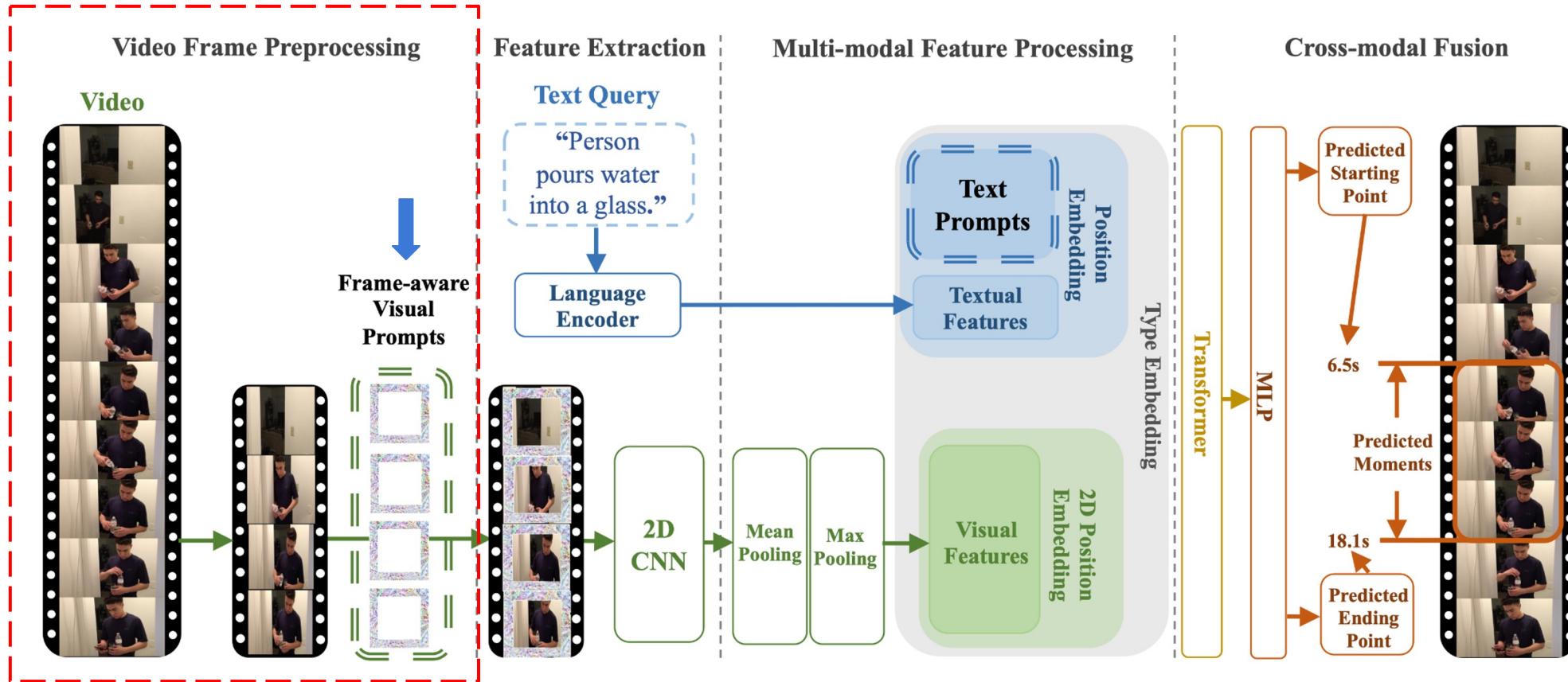
## 4. Text-Visual Prompt for 2D TVG



### Video frame preprocessing

- 1) Uniformly sample frames from input video.

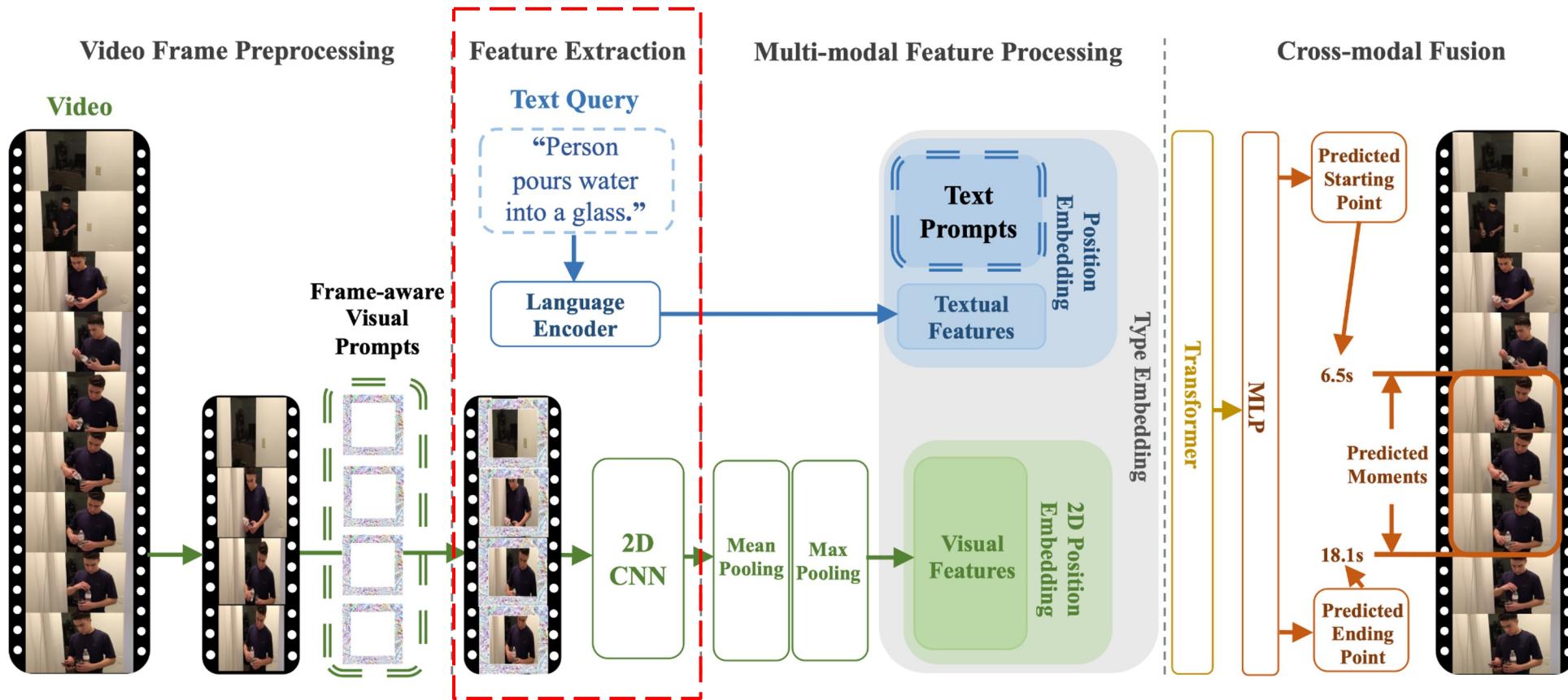
## 4. Text-Visual Prompt for 2D TVG



### Video frame preprocessing

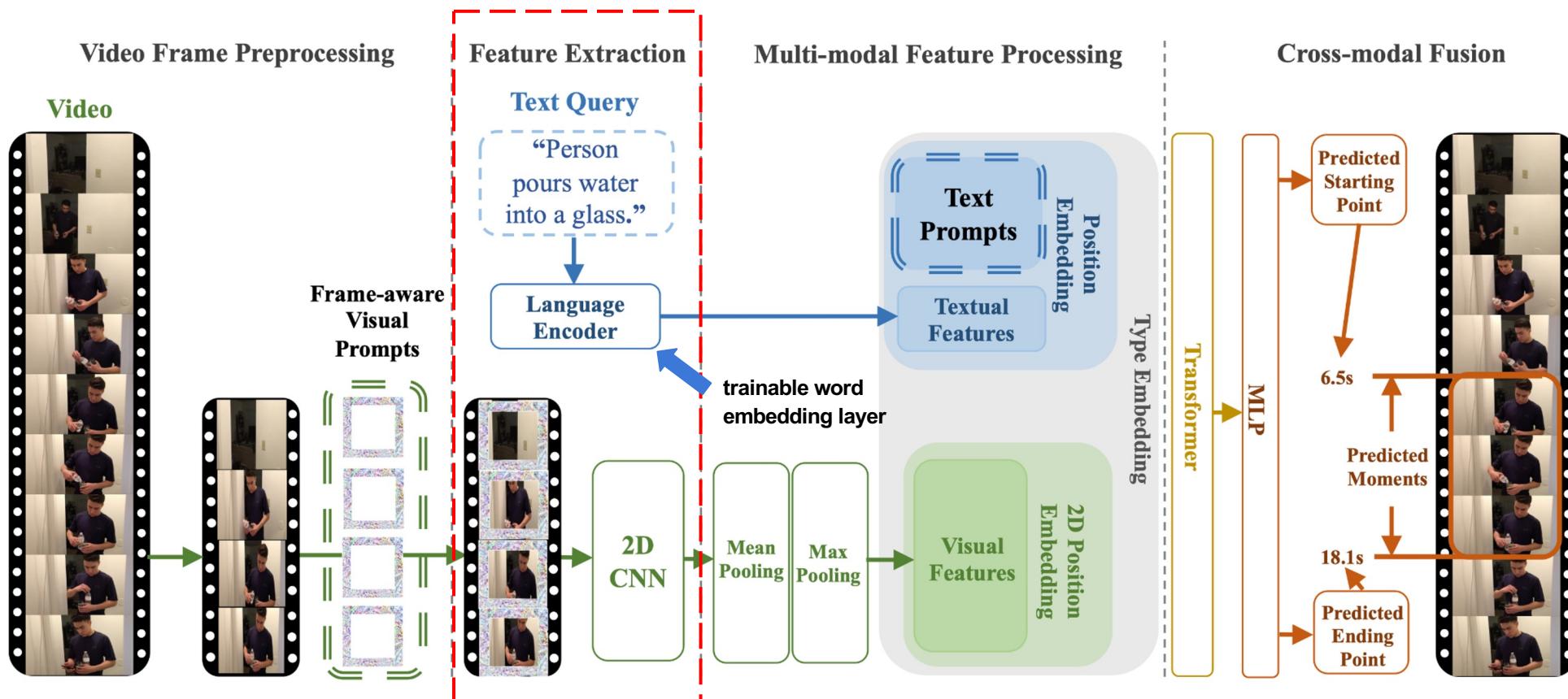
- 1) Uniformly sample frames from input video.
- 2) Apply an set of frame-aware visual prompts to the sampled frames in order.

## 4. Text-Visual Prompt for 2D TVG



**Feature extraction**

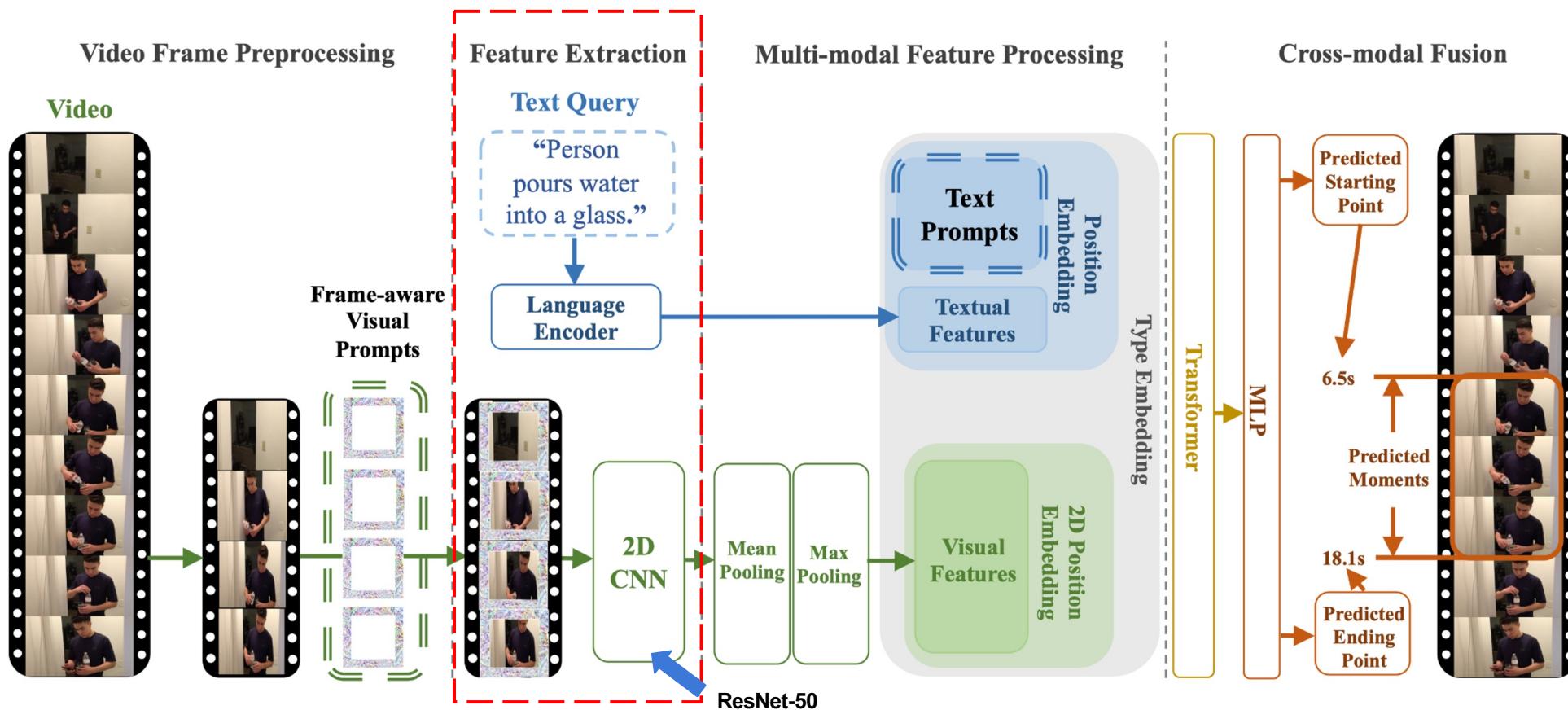
## 4. Text-Visual Prompt for 2D TVG



### Feature extraction

- 1) The language encoder extracts textual features.

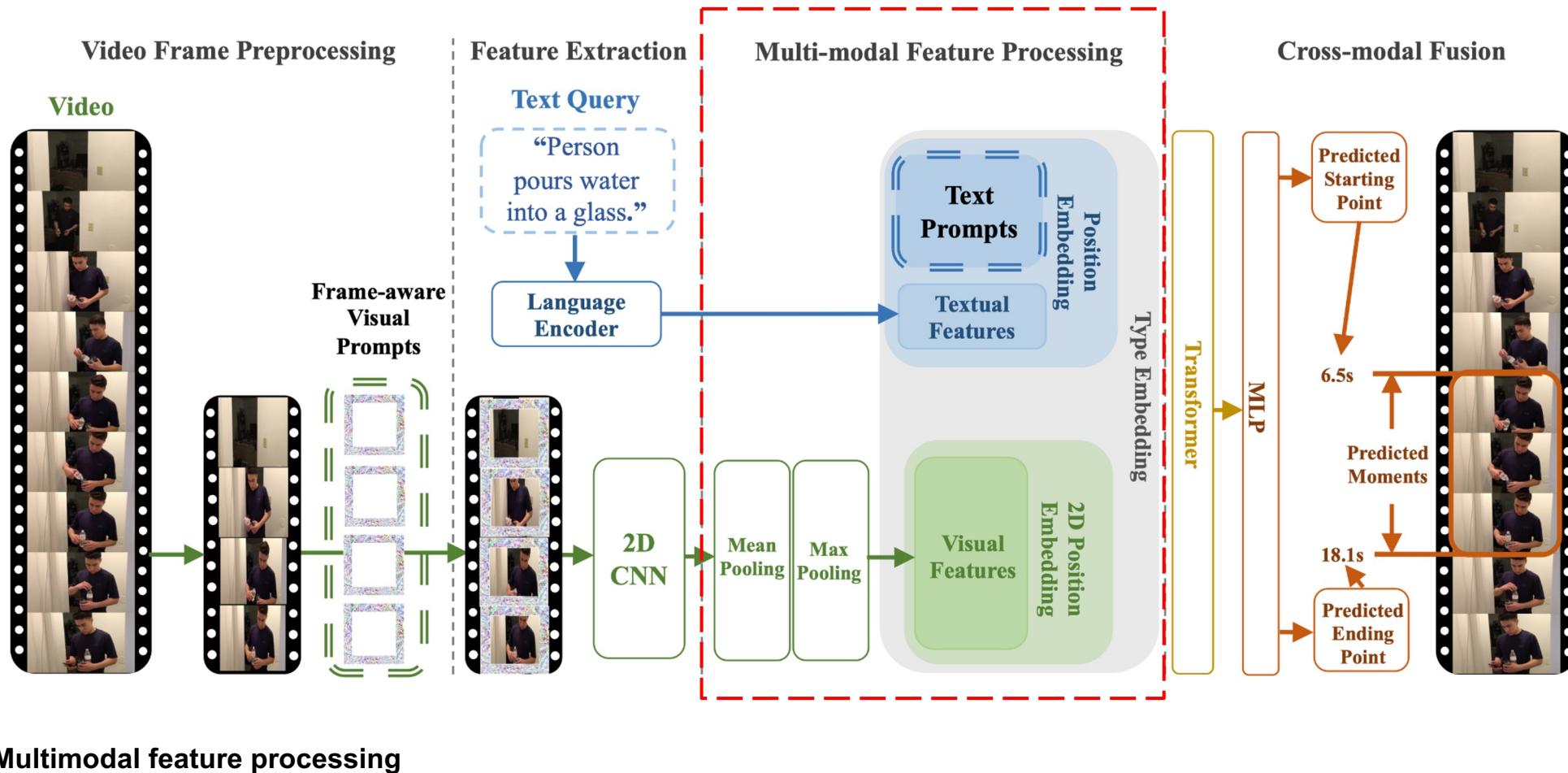
## 4. Text-Visual Prompt for 2D TVG



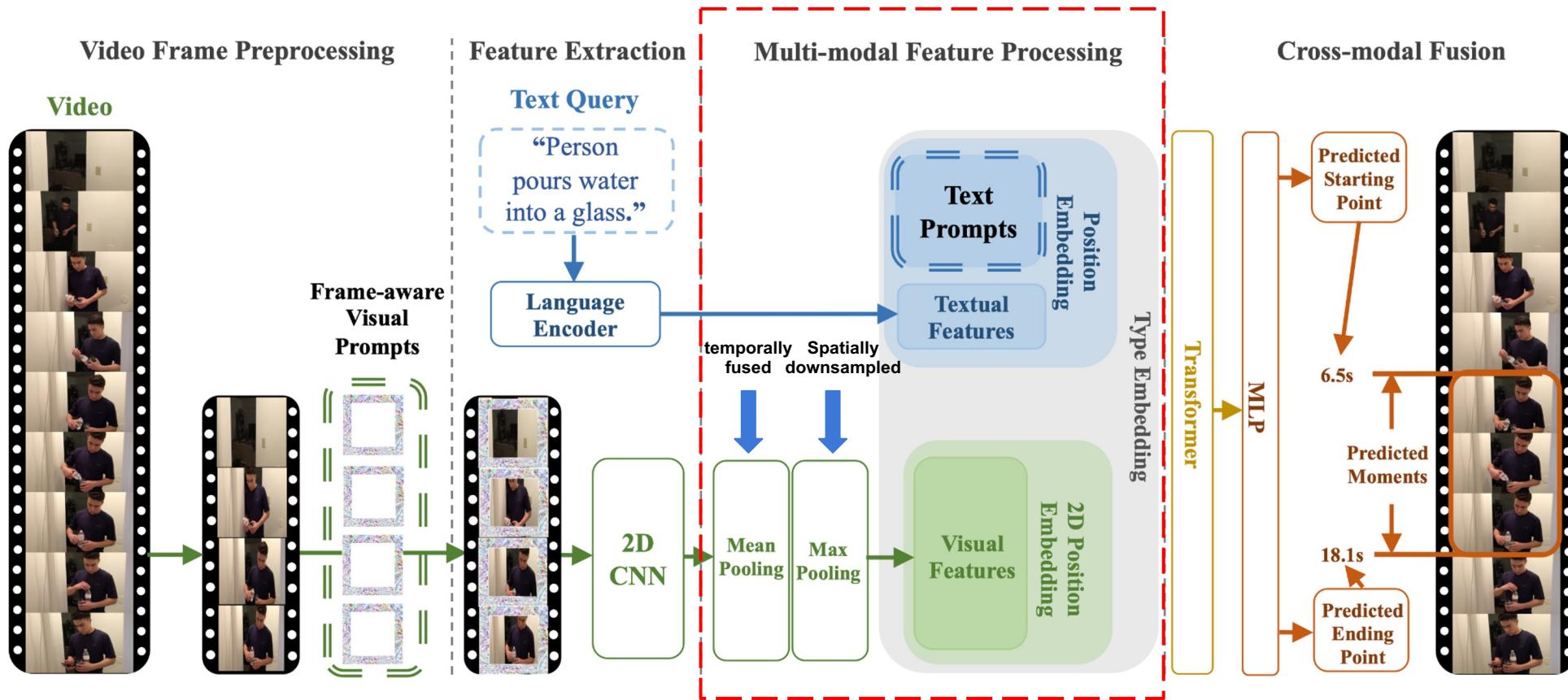
### Feature extraction

- 1) The language encoder extracts textual features.
- 2) 2D CNN extracts features from sampled video frames with visual prompts.

## 4. Text-Visual Prompt for 2D TVG



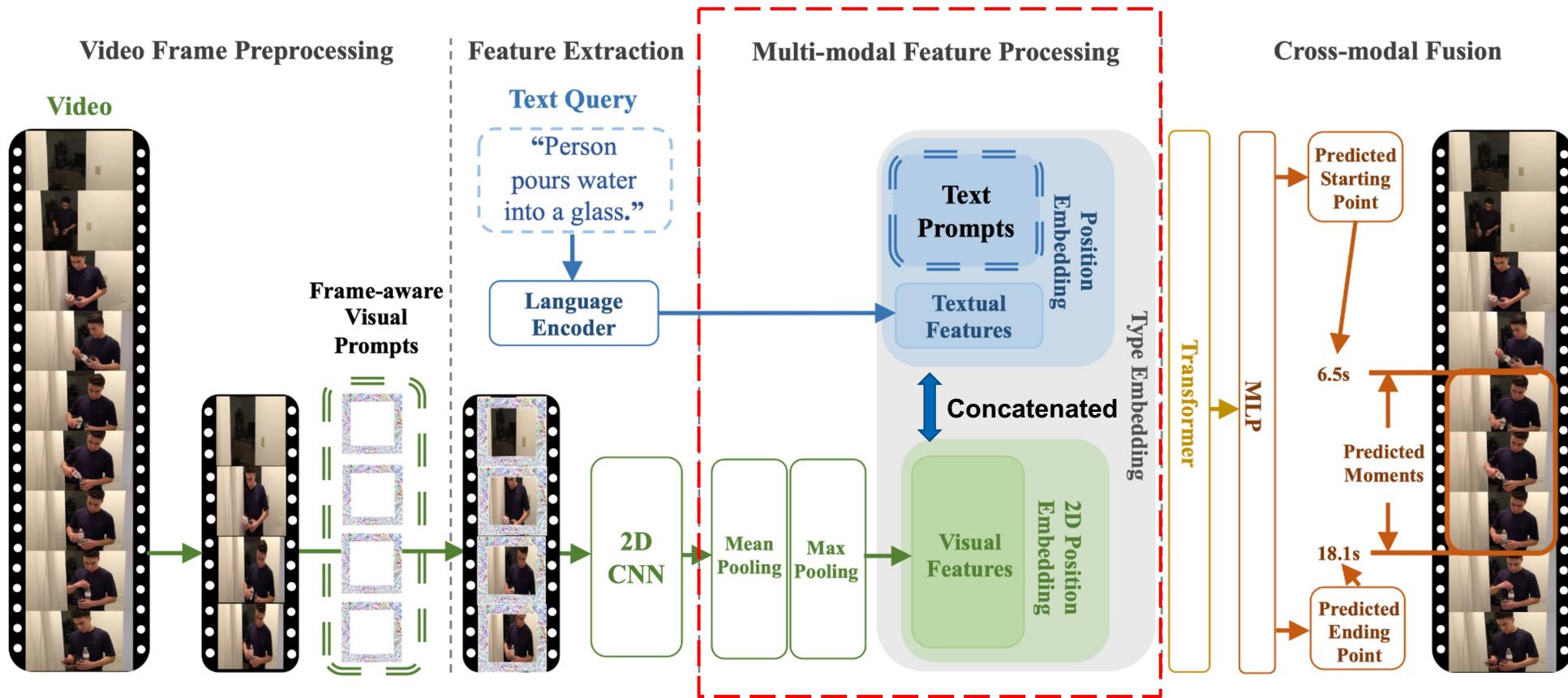
## 4. Text-Visual Prompt for 2D TVG



### Multimodal feature processing

- 1) Visual features would be temporally fused and spatially downsampled by mean pooling and max pooling, respectively.

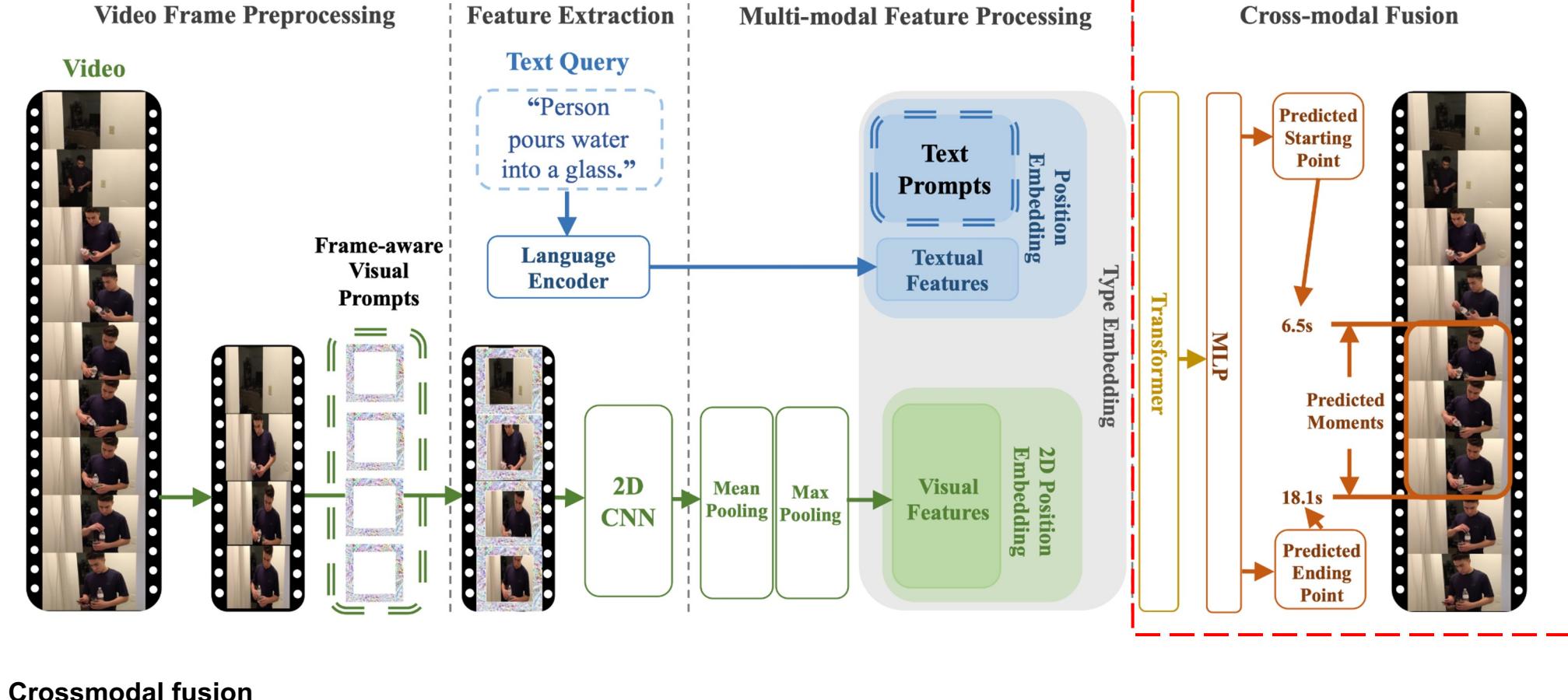
## 4. Text-Visual Prompt for 2D TVG



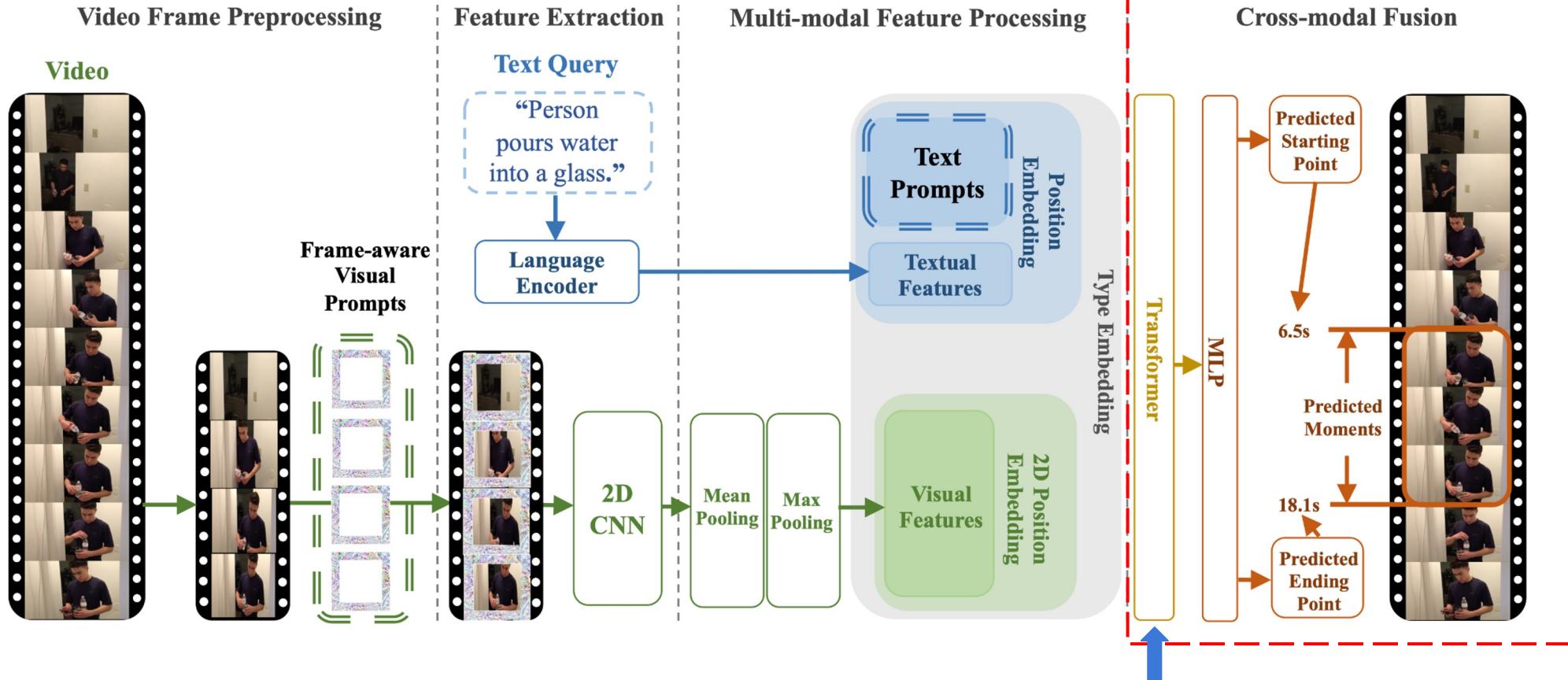
### Multimodal feature processing

- 1) Visual features would be temporally fused and spatially downsampled by mean pooling and max pooling, respectively.
- 2) The 2D visual features would be concatenated with textual features and text prompts.

## 4. Text-Visual Prompt for 2D TVG



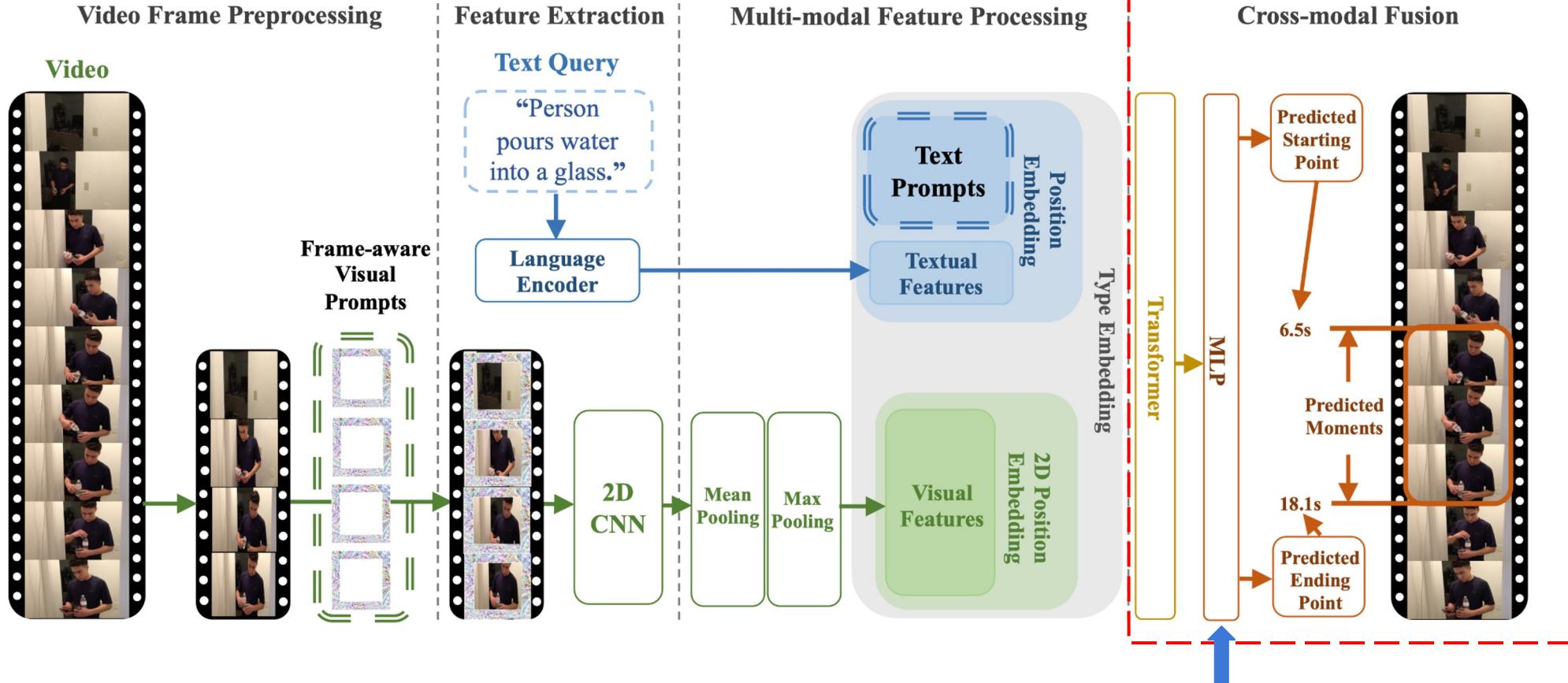
## 4. Text-Visual Prompt for 2D TVG



### Crossmodal fusion

- 1) The multimodal features would be processed by a 12-layer transformer encoder,

## 4. Text-Visual Prompt for 2D TVG



### Crossmodal fusion

- 1) The multimodal features would be processed by a 12-layer transformer encoder,
- 2) MLP would predict the starting/ending time points of the target moment.

## 6. Training Pipeline

## 6. Training Pipeline

- a) **Cross-modal pretraining** on large-scale image-text datasets.  
(COCO Captions and Visual Genome Captions)

## 6. Training Pipeline

- a) Cross-modal pretraining** on large-scale image-text datasets.  
(COCO Captions and Visual Genome Captions)
  
- b) Base model training** on the target dataset.

## 6. Training Pipeline

- a) **Cross-modal pretraining** on large-scale image-text datasets.  
(COCO Captions and Visual Genome Captions)
  
- b) **Base model training** on the target dataset.
  
- c) **Prompt training.**      ← Base model parameter are frozen !

## 6. Training Pipeline

- a) **Cross-modal pretraining** on large-scale image-text datasets.  
(COCO Captions and Visual Genome Captions)
  
- b) **Base model training** on the target dataset.
  
- c) **Prompt training.** ← Base model parameter are frozen !
  
- d) **Base model finetuning.** ← Text-Visual Prompts are frozen !

## 7. Dataset

Dataset	Charades-STA	ActivityNet Captions
Domain	Indoor Activity	Indoor/Outdoor Activity
# Videos	6,672	14,926
Avg. Video Length ( <i>second</i> )	30.6	117.6
# Moments	11,767	71,953
Avg. Moment Length ( <i>second</i> )	8.1	37.1
Vocabulary Size	1,303	15,505
# Queries	16,124	71,953
Avg. Query Length ( <i>word</i> )	7.2	14.4

Table 1. Statics of temporal video grounding benchmark datasets (Charades-STA and ActivityNet Captions datasets).

## 8. Evaluation Metric

**Acc(R@1, IoU=m)**

## 8. Evaluation Metric

$\text{Acc}(\text{R@1}, \text{IoU}=m)$



The percentage of predicted moments  
achieving IoU higher than  $m$  with the groundtruth moment.

## 8. Experimental Results

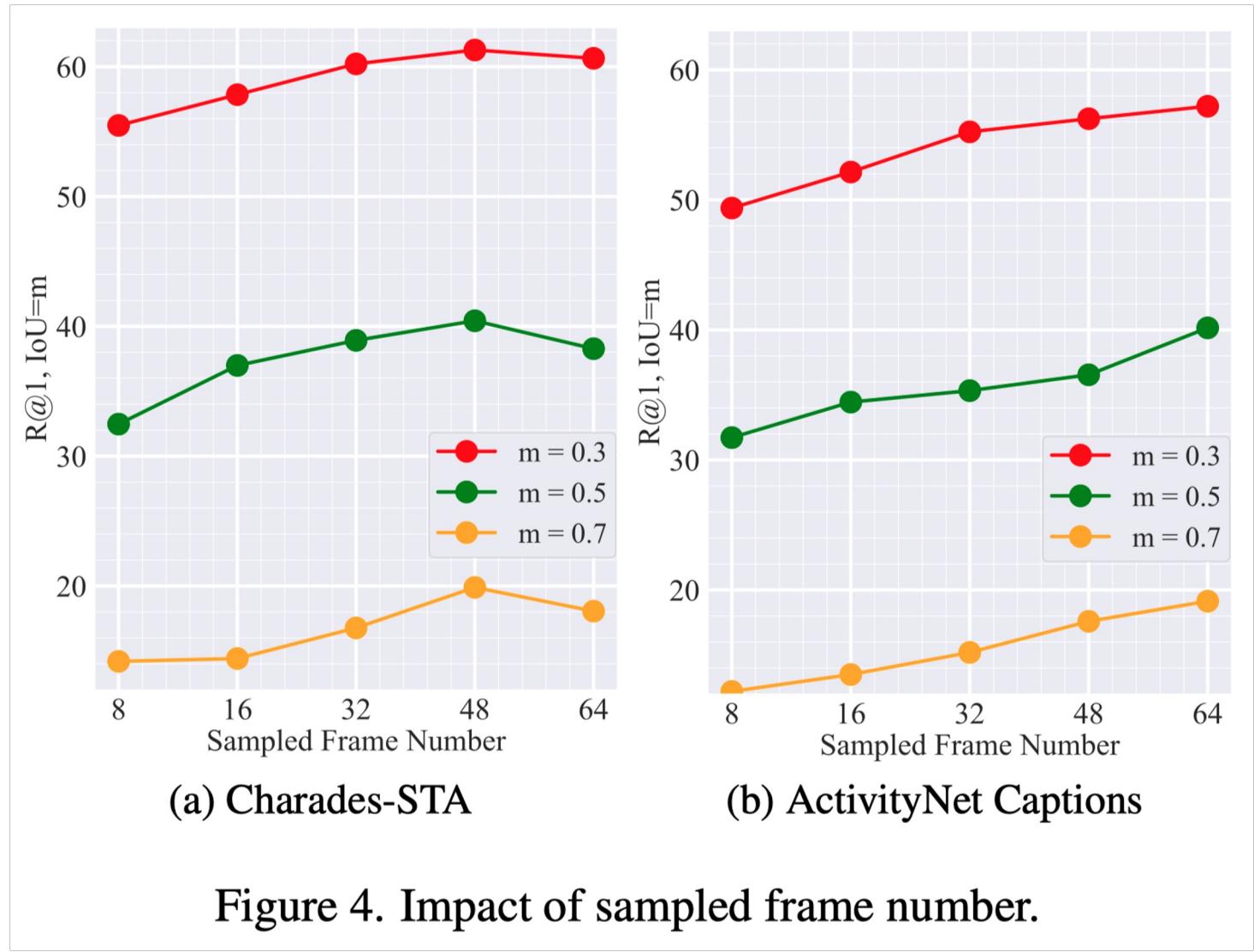
Table 2. Performance comparison of different thresholds  $m$  on the Charades-STA dataset.

Type	Method	Visual Feature	Acc(R@1, IoU= $m$ )		
			$m=0.3$	$m=0.5$	$m=0.7$
3D TVG	BPNet [53]	C3D	55.46	38.25	20.51
	LPNet [52]	C3D	59.14	40.94	21.13
	QSPN [55]	C3D	54.70	35.60	15.80
	TSP-PRL [51]	C3D	-	45.45	24.75
	TripNet [15]	C3D	54.64	38.29	16.07
	DRN [59]	C3D	-	45.40	26.40
	CPNet [28]	C3D	-	40.32	22.47
	DEBUG [34]	C3D	54.95	37.39	17.92
	ExCL [14]	I3D	61.50	44.1	22.40
	VSLNet [63]	I3D	64.30	<b>47.31</b>	<b>30.19</b>
2D TVG	MAN [61]	I3D	-	46.53	22.72
	CTRL [12]	VGG	13.5	9.82	-
	MCN [1]	VGG	17.46	8.01	-
	ABLR [58]	VGG	24.36	9.01	-
TVP-Based 2D TVG	SAP [5]	VGG	27.42	13.36	-
	<b>Ours</b>				
	Base w/o prompts	ResNet	61.29	40.43	19.89
	Base + Visual Prompts		65.38	44.31	20.22
	Base + Text Prompts		65.81	43.44	20.65
	Base + Both Prompts		<b>65.92</b>	44.39	21.51

Table 3. Performance comparison of different thresholds  $m$  on the ActivityNet Captions dataset.

Type	Method	Visual Feature	Acc(R@1, IoU= $m$ )		
			$m=0.3$	$m=0.5$	$m=0.7$
3D TVG	CTRL [12]	C3D	28.70	14.00	-
	BPNet [53]	C3D	59.98	42.07	24.69
	LPNet [52]	C3D	<b>64.29</b>	<b>45.92</b>	25.39
	QSPN [55]	C3D	45.30	27.70	13.60
	TSP-PRL [51]	C3D	56.02	38.83	-
	TripNet [15]	C3D	48.42	32.19	13.93
	DRN [59]	C3D	-	45.45	24.36
	CPNet [28]	C3D	-	40.56	21.63
	ABLR [58]	C3D	55.67	36.79	-
	DEBUG [34]	C3D	55.91	39.72	-
<b>Ours</b>	ExCL [14]	C3D	63.00	43.60	24.10
	VSLNet [63]	C3D	63.16	43.22	<b>26.16</b>
	Base w/o prompts	ResNet	57.20	40.16	19.14
	Base + Visual Prompts		60.12	43.39	23.71
TVP-Based 2D TVG	Base + Text Prompts		60.48	42.58	24.39
	Base + Both Prompts		60.71	43.44	25.03

## 8. Experimental Results



## 8. Experimental Results

Table 4. The performance comparison of different visual prompt sizes on Charades-STA dataset.

Visual Prompt Size	Acc(R@1, IoU= $m$ )			Prompt + Frame
	$m=0.3$	$m=0.5$	$m=0.7$	
0	61.29	40.43	19.89	
16	61.29	40.43	20.00	
32	61.94	39.78	19.35	
48	63.66	42.37	20.00	
72	63.87	43.66	19.78	
96	<b>65.38</b>	<b>44.31</b>	<b>20.22</b>	
128	64.73	43.66	19.78	

Table 5. The performance comparison of different text prompt sizes on Charades-STA dataset.

Text Prompt Size	Acc(R@1, IoU= $m$ )		
	$m=0.3$	$m=0.5$	$m=0.7$
0	57.20	40.16	19.14
5	65.38	41.94	20.43
10	<b>65.81</b>	43.44	20.65
15	65.59	43.23	21.29
20	64.95	<b>43.87</b>	<b>21.51</b>
25	63.66	42.80	20.65
30	64.46	42.63	20.51

## 8. Experimental Results

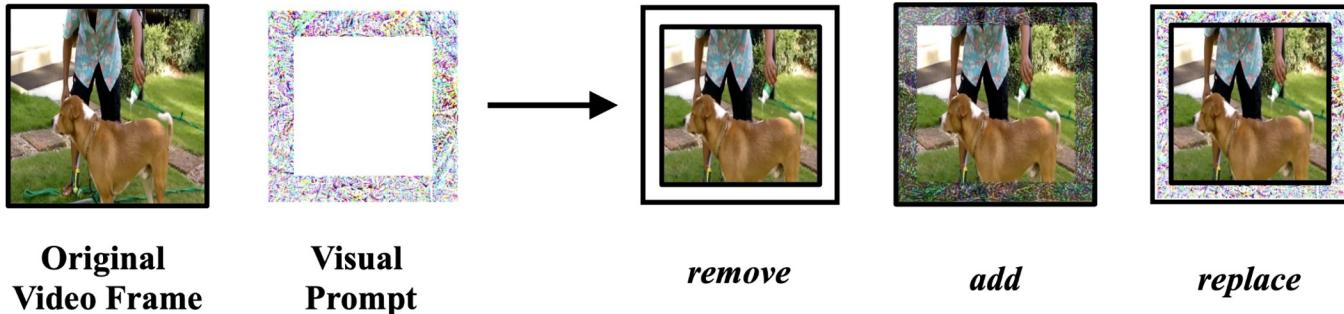


Table 6. The performance comparison of different visual prompt operations ('remove', 'add', 'replace') with fixed visual prompt size  $p = 96$  on Charades-STA and ActivityNet Captions datasets.

Operation	Charades-STA			ActivityNet Captions		
	R@1, IoU= $m$			R@1, IoU= $m$		
	$m=0.3$	$m=0.5$	$m=0.7$	$m=0.3$	$m=0.5$	$m=0.7$
Original	61.29	40.43	19.89	57.20	40.16	19.14
Remove	61.29	40.43	20.0	57.20	40.16	19.14
Add	61.08	39.57	20.22	57.15	40.16	19.27
Replace	<b>65.38</b>	<b>44.31</b>	<b>20.22</b>	<b>60.12</b>	<b>43.39</b>	<b>23.71</b>

## 8. Experimental Results

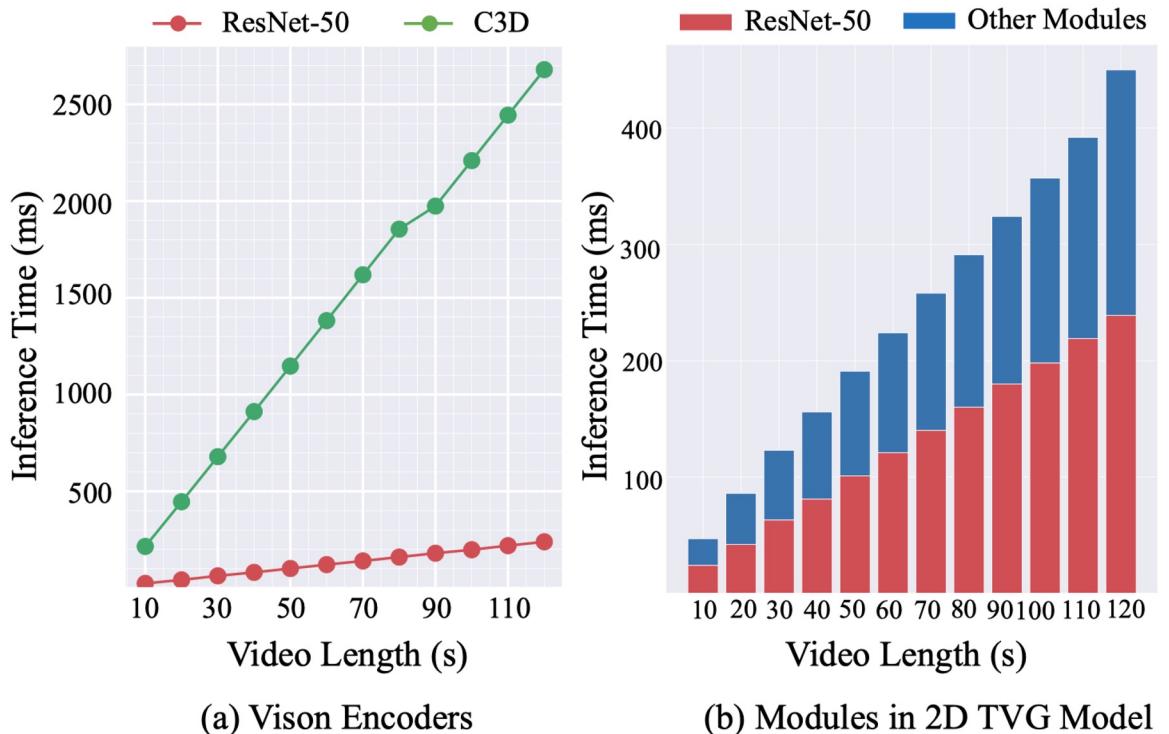


Figure 5. Inference time comparison. (a) inference time comparison between 2D vision encoder (ResNet-50) and 3D vision encoder (C3D). (b) inference time comparison between the vision encoder and the other modules of the 2D TVG model, where the sampled frame number for our TVP framework is  $1.2 \times$  the length of the video in seconds.

# THANKS for watching! !

## Any Questions?



**More Details on  
Project Website**