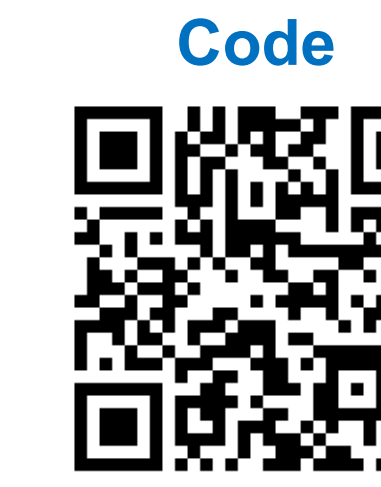


# Leveraging Temporal Contextualization for Video Action Recognition

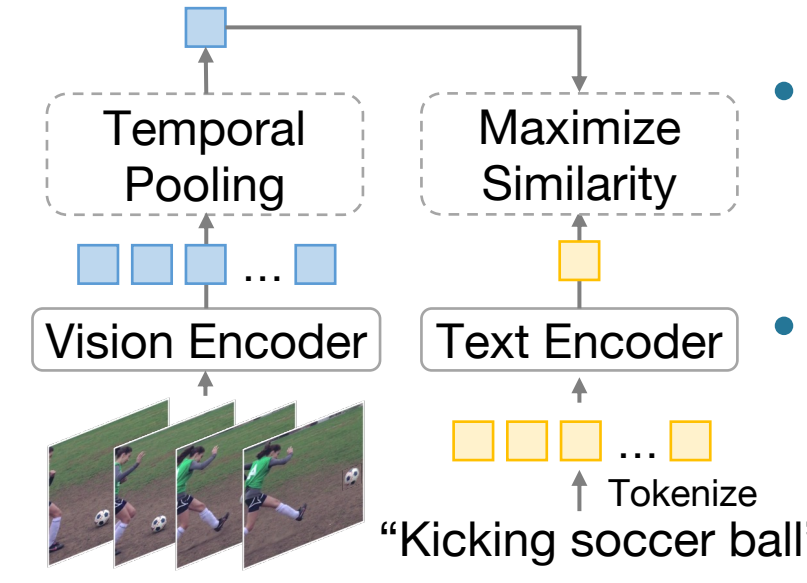
Minji Kim<sup>1†</sup> Dongyoon Han<sup>2</sup> Taekyung Kim<sup>2★</sup> Bohyung Han<sup>1★</sup>

<sup>†</sup> Work done during an internship at NAVER AI Lab    <sup>★</sup> Corresponding authors



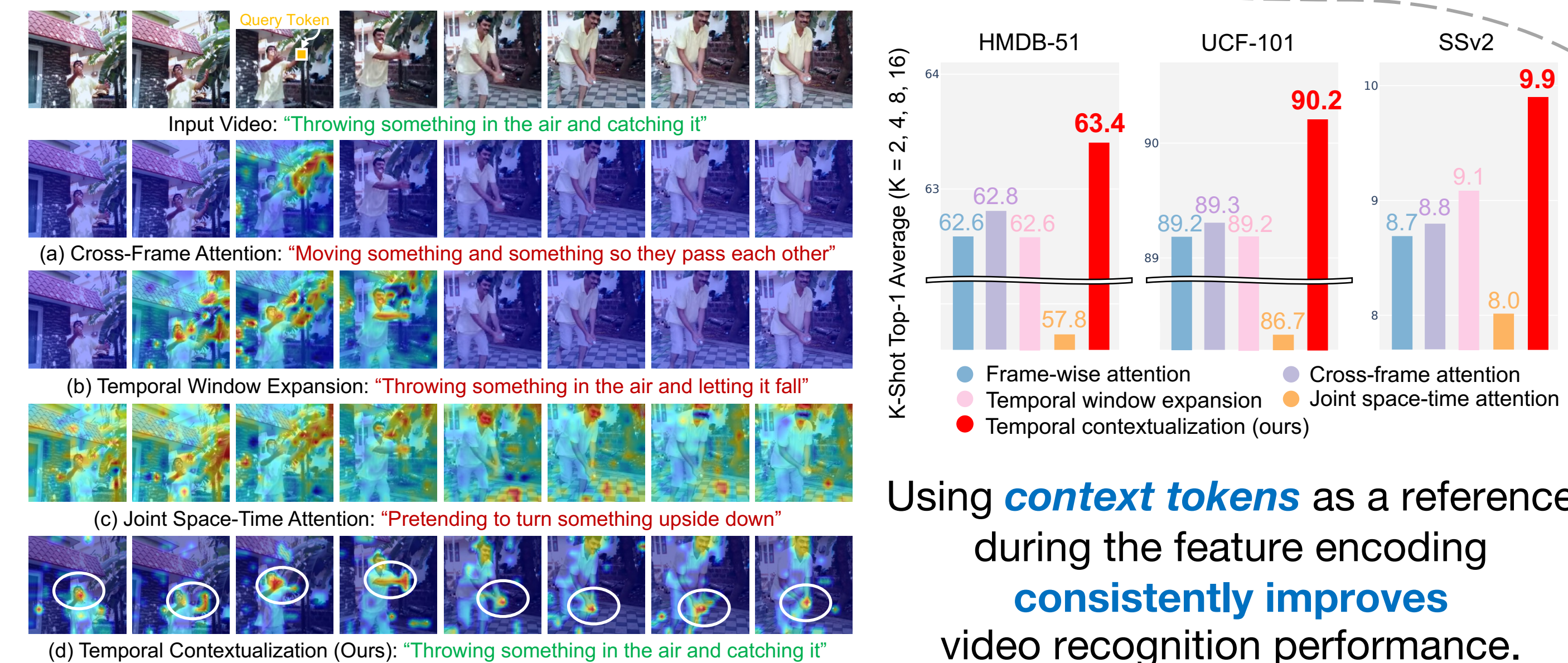
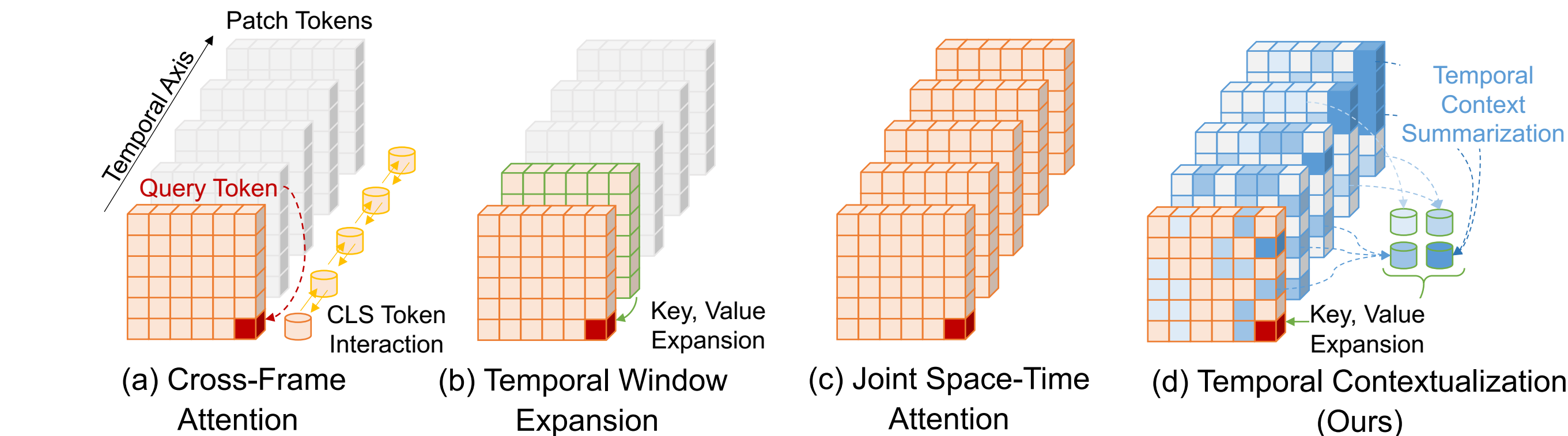
( See our paper for more results! )

## Background



- Tuning **CLIP for video recognition** enables open-vocab generalization without expensive video-text pretraining.
- A naïve baseline: **frame-wise attention**  
→ **Limitation: no token interactions in the temporal axis**
- To consider temporal cues, prior works additionally incorporate **reference tokens**:  $z_t^l = f_{\theta_v}^l(z_t^{l-1}, s^{l-1})$

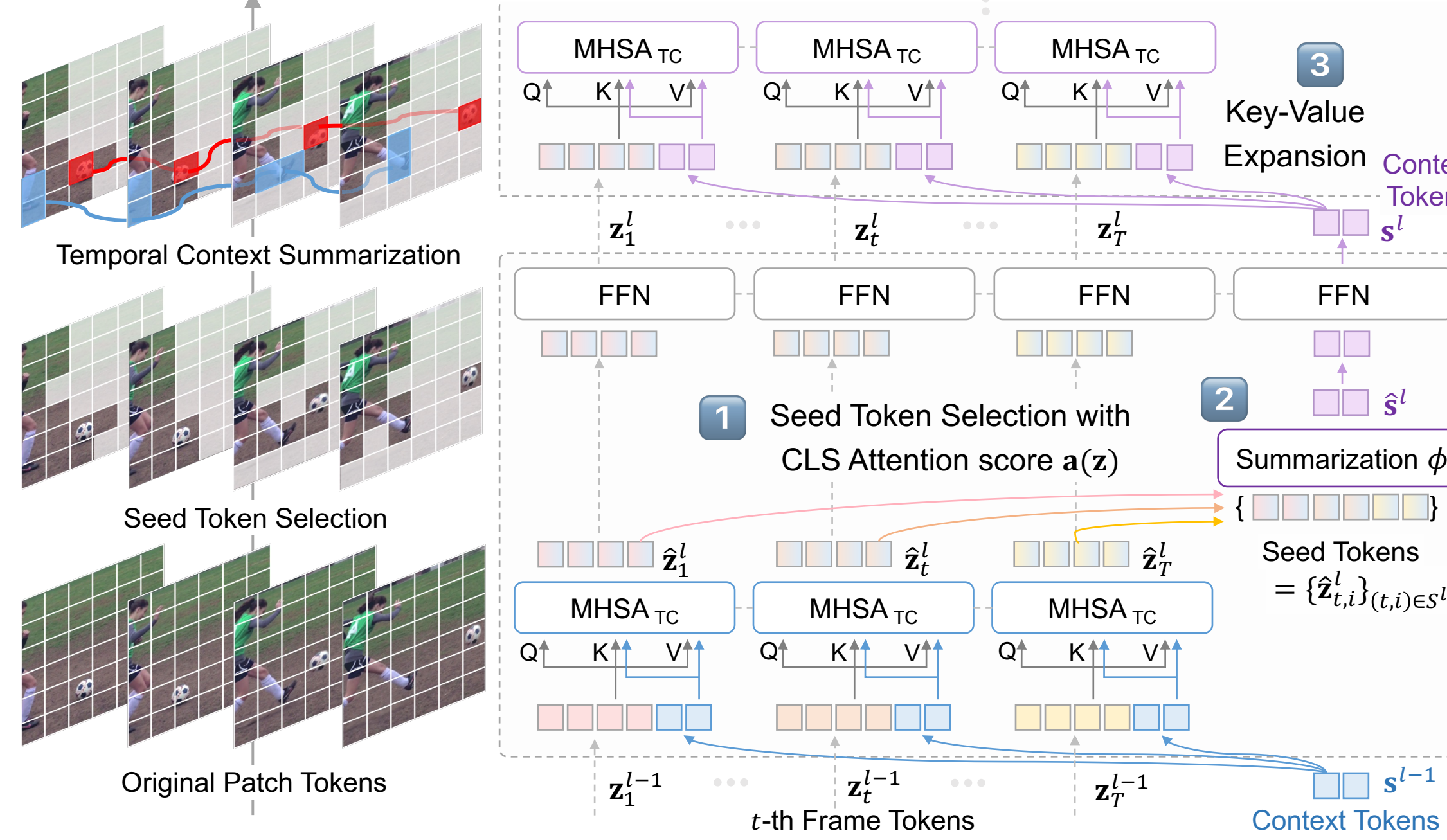
## Problem: Insufficient Token Interactions in Temporal Modeling



Using **context tokens** as a reference during the feature encoding **consistently improves** video recognition performance.

## Temporal Contextualization (TC)

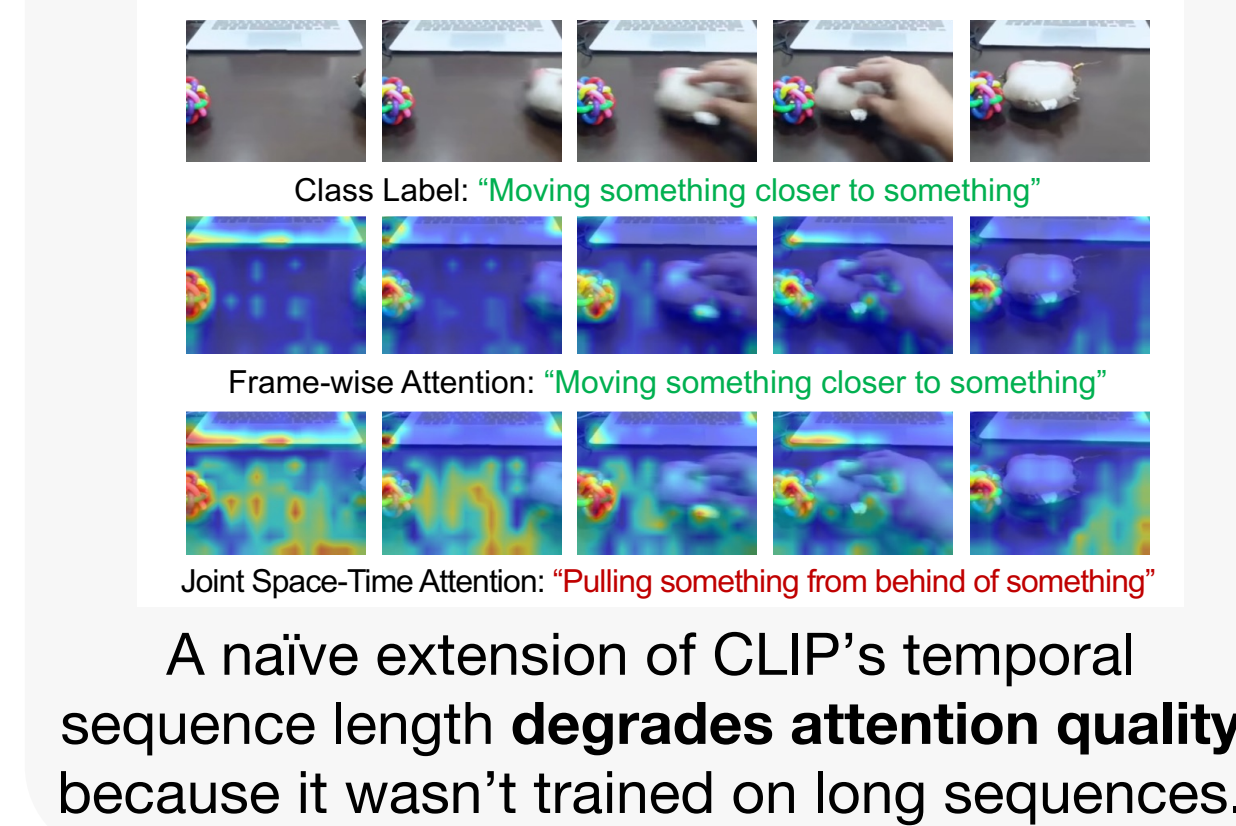
- Key Idea: Summarize informative tokens from the entire video** into a small set of tokens, called **context tokens**, and **reference** them during feature encoding.



## Video-conditional Prompting (VP)

- VP generates **instance-level prompts** to compensate for the **lack of textual semantics**.
- Video information from context tokens is injected to text prompt vectors** using a cross-attention mechanism.

### ✖ Pitfall of Joint Space-Time Attention

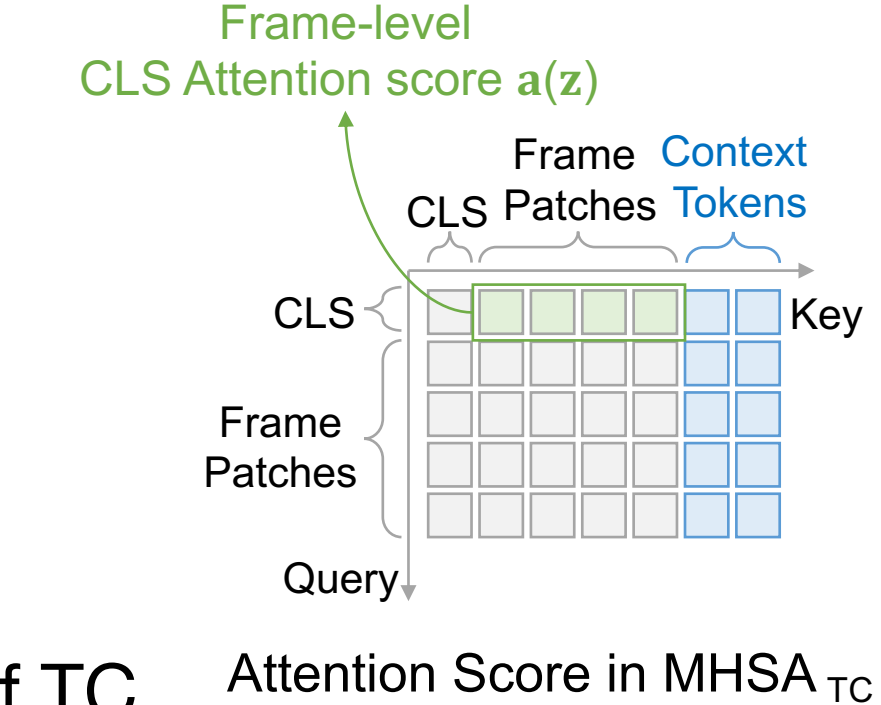
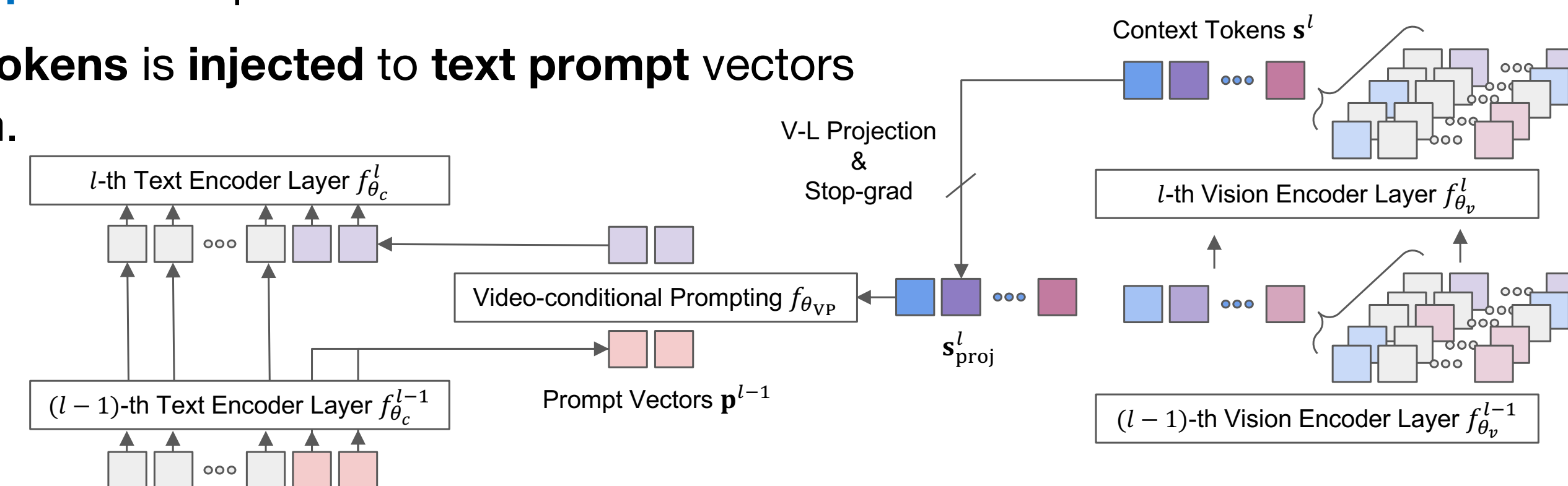


- Three steps of TC**

  - Seed token selection** in each frame
  - Spatio-temporal context summarization**
  - Temporal context infusion** to all tokens by **expanding key-value pairs**:

$$\text{Attention}_{TC}(z_t, s) = \text{Softmax} \left( \frac{Q_{z_t} [K_{z_t} | K_s]^T}{\sqrt{d}} + B \right) [V_{z_t} | V_s]$$

Learnable bias  $B_{ij} = \begin{cases} b_{\text{local}} & \text{if } j \leq N + 1 \\ b_{\text{global}} & \text{otherwise,} \end{cases}$



## Experiments

**SOTA performance in zero/few-shot, base2novel, fully-supervised video recognition**

- Results on **zero-shot action recognition**

Method	WE	HMDB-51	UCF-101	K600 (Top-1)	K600 (Top-5)	All (Top-1)
Vanilla CLIP [32]		40.8 ± 0.3	63.2 ± 0.2	59.8 ± 0.3	83.5 ± 0.2	54.6
ActionCLIP [39] <sup>†</sup>		49.1 ± 0.4	68.0 ± 0.9	56.1 ± 0.9	83.2 ± 0.2	57.7
A5 [14]		44.3 ± 2.2	69.3 ± 4.2	55.8 ± 0.7	81.4 ± 0.3	56.5
X-CLIP [29]		44.6 ± 5.2	72.0 ± 2.3	65.2 ± 0.4	86.1 ± 0.8	60.6
Vita-CLIP [41]		48.6 ± 0.6	75.0 ± 0.6	67.4 ± 0.5	-	63.7
ViFi-CLIP [34] <sup>†</sup>		52.3 ± 0.2	78.9 ± 1.1	70.7 ± 0.8	92.1 ± 0.3	67.3
<b>TC-CLIP (Ours)</b>		<b>53.7 ± 0.7</b>	<b>80.4 ± 0.9</b>	<b>72.7 ± 0.5</b>	<b>93.2 ± 0.2</b>	<b>68.9</b>
ActionCLIP [39] <sup>†</sup>	✓	51.9 ± 0.5	74.2 ± 1.0	67.5 ± 1.2	90.7 ± 0.1	64.5
ViFi-CLIP [34] <sup>†</sup>	✓	52.2 ± 0.7	81.0 ± 0.9	73.9 ± 0.5	93.3 ± 0.3	69.0
Open-VCLIP [42]	✓	53.9 ± 1.2	<b>83.4 ± 1.2</b>	73.0 ± 0.8	93.2 ± 0.1	<b>70.1</b>
<b>TC-CLIP (Ours)</b>	✓	<b>54.2 ± 0.7</b>	<b>82.9 ± 0.6</b>	<b>75.8 ± 0.5</b>	<b>94.4 ± 0.2</b>	<b>71.0</b>

Using LLM-based text augmentation

Case	WE	HMDB-51	UCF-101	K-600	All (Δ)	
MAXI [24]	✓	52.3 ± 0.7	78.2 ± 0.8	71.5 ± 0.8	92.5 ± 0.4	67.3
OST [4]	✓	55.9 ± 1.2	79.7 ± 1.1	75.1 ± 0.6	94.6 ± 0.2	70.2
FROSTER [10]	✓	54.8 ± 1.3	84.8 ± 1.1	74.8 ± 0.9	-	71.5
<b>TC-CLIP (Ours)</b>	✓	<b>56.0 ± 0.3</b>	<b>85.4 ± 0.8</b>	<b>78.1 ± 1.0</b>	<b>95.7 ± 0.3</b>	<b>73.2</b>

- Component-wise ablation: TC and VP are **both effective**.

Case	Without weight-space ensembling				With weight-space ensembling			
	HMDB-51	UCF-101	K-600	All (Δ)	HMDB-51	UCF-101	K-600	All (Δ)
Baseline	52.3 ± 0.2	78.9 ± 1.1	70.7 ± 0.8	67.3	52.2 ± 0.7	81.0 ± 0.9	73.9 ± 0.5	69.0
(a) +TC	53.6 ± 0.2	78.6 ± 1.0	71.8 ± 0.7	68.0 (+0.7)	54.3 ± 0.6	81.9 ± 1.0	75.5 ± 1.0	70.6 (+1.6)
(b) +VP	53.2 ± 0.8	80.5 ± 0.7	71.6 ± 0.9	68.4 (+1.1)	53.4 ± 0.8	82.0 ± 0.9	74.7 ± 0.7	70.0 (+1.0)
<b>(c) +TC+VP</b>	<b>53.7 ± 0.7</b>	<b>80.4 ± 0.9</b>	<b>72.7 ± 0.5</b>	<b>68.9 (+1.6)</b>	<b>54.2 ± 1.1</b>	<b>82.9 ± 0.9</b>	<b>75.8 ± 0.4</b>	<b>71.0 (+2.0)</b>

- TC is robust across diverse token aggregation strategies.**

Case	(a) Seed token selection strategy.				(b) Context token summarization strategy.				
	HMDB	UCF	SSv2	All (Δ)	HMDB	UCF	SSv2	All (Δ)	
Baseline	62.6	89.2	8.7	53.5	62.6	89.2	8.7	53.5	
No selection	62.8	89.8	9.7	54.1 (+0.6)	No merge	57.2	85.6	7.7	50.2 (-3.3)
Head-wise key norm	62.3	89.8	9.8	54.0 (+0.5)	Random merge	58.8	87.1	7.5	51.2 (-2.3)
Averaged key norm	62.5	89.4	9.3	53.7 (+0.2)	K-means [25]	62.1	89.7	9.0	53.6 (+0.1)
Head-wise CLS attn.	63.4	89.9	9.7	54.3 (+0.8)	DPC-KNN [13]	63.3	<b>90.2</b>	9.8	54.4 (+0.9)
Averaged CLS attn.	63.4	<b>90.2</b>	<b>9.9</b>	<b>54.5 (+1.0)</b>	Bipartite soft matching [1,15]	<b>63.4</b>	<b>90.2</b>	<b>9.9</b>	<b>54.5 (+1.0)</b>
Patch saliency [5]	62.9	<b>90.3</b>	9.6	54.2 (+0.7)	Bipartite w/ attention weights	62.9	89.8	<b>9.9</b>	54.2 (+0.7)
ATS [8]	<b>63.5</b>	<b>90.3</b>	9.8	<b>54.5 (+1.0)</b>	Bipartite w/ saliency weights [5]	62.4	89.9	9.6	54.0 (+0.5)

- Context-token-conditional text prompting is effective.**

Case	Use context tokens?	HMDB-51	UCF-101	K-600	All (Δ)
Baseline		52.3 ± 0.2	78.9 ± 1.1	70.7 ± 0.8	67.3
(a) Learnable prompt vectors		52.4 ± 0.4	78.4 ± 1.3	70.6 ± 0.7	67.1 (-0.2)
(b) Video-conditional prompting		53.2 ± 0.8	80.4 ± 0.7	71.6 ± 0.9	68.4 (+1.1)
<b>(c) Video-conditional prompting</b>	✓	<b>53.7 ± 0.7</b>	<b>80.4 ± 0.9</b>	<b>72.7 ± 0.5</b>	<b>68.9 (+1.6)</b>
(d) Vision-text late-fusion	✓	53.7 ± 0.7	79.0 ± 0.7	70.9 ± 0.6	67.9 (+0.6)

**Temporal Contextualization enhances CLIP's video understanding capability by infusing global information within the encoding process.**