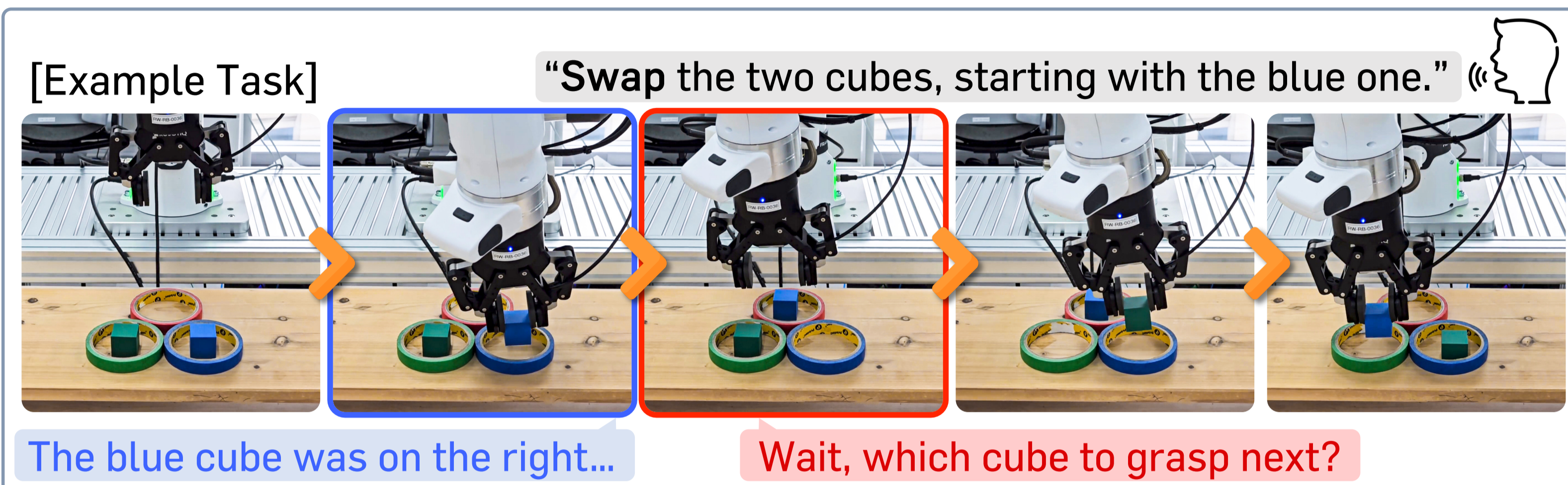


Myungkyu Koo^{*A}, Daewon Choi^{*A}, Taeyoung Kim^A,
Kyungmin Lee^A, Changyeon Kim^A, Younggyo Seo^{†B}, Jinwoo Shin^{†AC}
^AKAIST, ^BUC Berkeley, ^CRLWORLD, ^{*}Equal Contribution, [†]Equal Advising

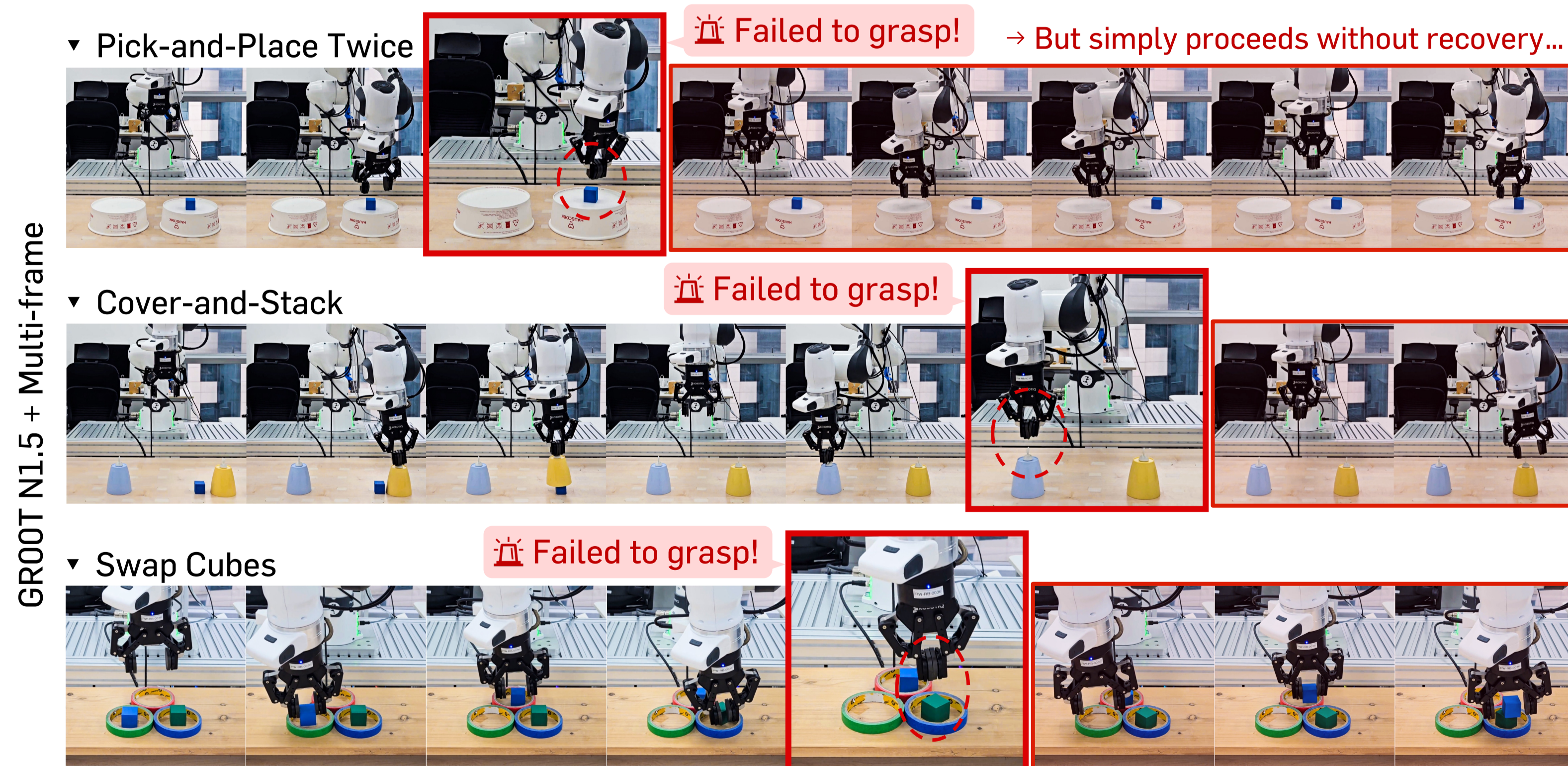
Motivation



- ✓ Most recent VLAs rely on a **single-frame assumption**, solely depending on current observations.
- ✓ Yet robotic manipulation is inherently **history-dependent!** (e.g., deciding which cube to grasp next; see example above)

Observation & Challenges

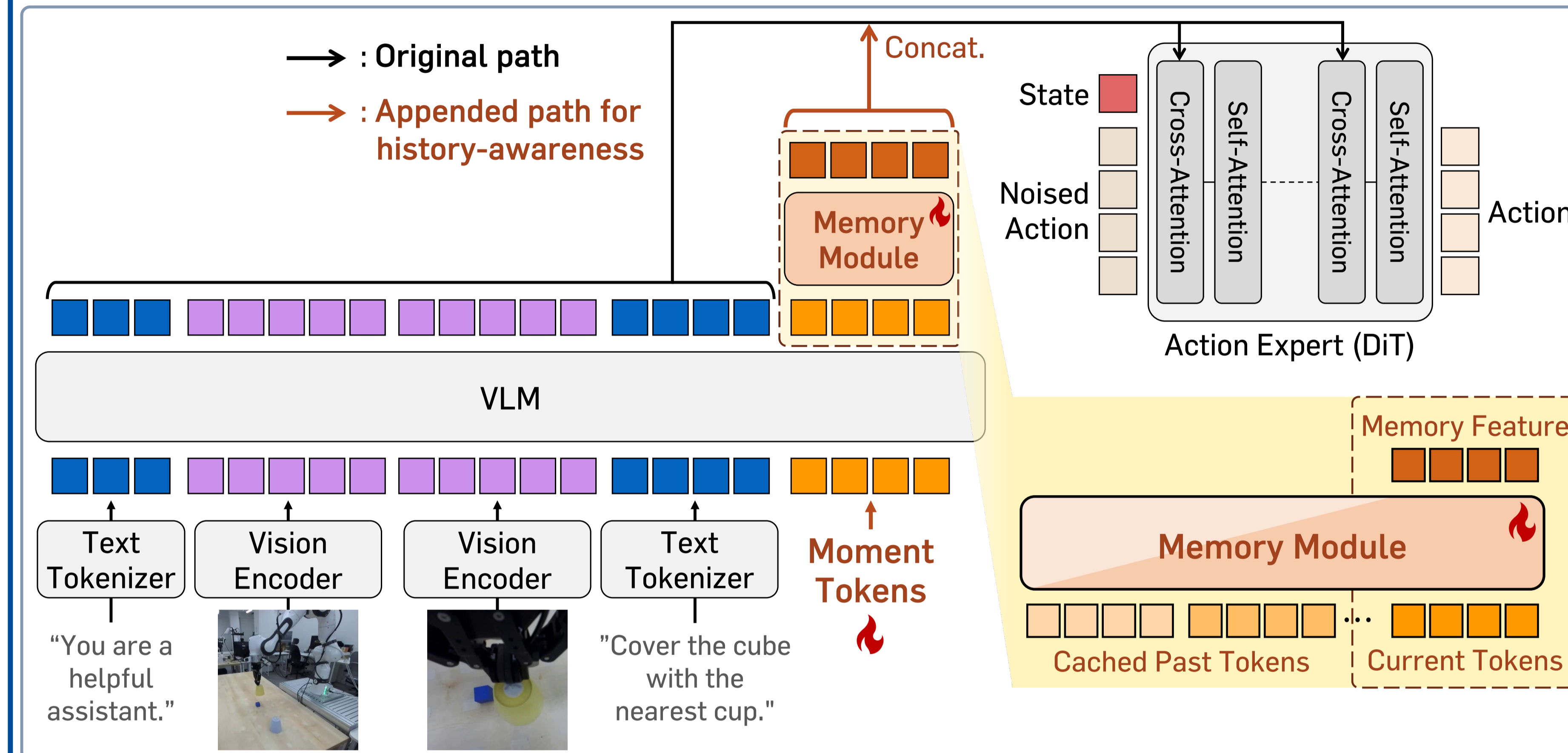
- ✓ When naively using multi-frame inputs for longer-context,
 - ① **Inefficient scaling:** Appending 4 past frames to VLA input results in **1.35x** latency, **3.64x** peak memory, lacking practicality.
 - ② **Poor generalization:** Easily suffers from causal confusion (e.g., overfits to past proprioceptive states, not current context) and fails to recover from intermediate failures (see below).



Our Answer: HAMLET

- ① **+47.2%** on real-world history-dependent tasks
- ② Consistent gains on public benchmarks (LIBERO, RoboCasa, Simpler)
- ③ Only **1.02x** latency, **1.96x** peak memory

Overview of HAMLET



- ✓ TL;DR: HAMLET is a plug-in framework to incorporate history-awareness into VLAs via i) **moment tokens** and ii) a **memory module**.

Method

1. Context compression via Moment Tokens

- **Function:** Appends **learnable tokens** to VLM input to compress per-timestep context into compact representations.

$$[\mathbf{h}_t; \mathbf{m}'_t] = \mathcal{F}_\theta([\mathbf{o}_t, \mathbf{c}; \mathbf{m}_t])$$

- **Initialization:** Utilizes **Time-Contrastive Learning** to capture task-relevant cues while suppressing static backgrounds.

$$\mathcal{L}_{\text{TCL}}(\mathbf{z}_t, \mathbf{z}_t^+) = - \sum_{t=1}^B \log \frac{\exp(\text{sim}(\mathbf{z}_t, \mathbf{z}_t^+)/\tau)}{\exp(\text{sim}(\mathbf{z}_t, \mathbf{z}_t^+)/\tau) + \exp(\text{sim}(\mathbf{z}_t, \mathbf{z}_t^-)/\tau)}, \quad \mathbf{z}_t = g(\mathbf{m}'_t)$$

\mathbf{z}_t^+ : Same timestep w/ image augmentation, \mathbf{z}_t^- : Different timestep

2. Memory consolidation via Memory Module

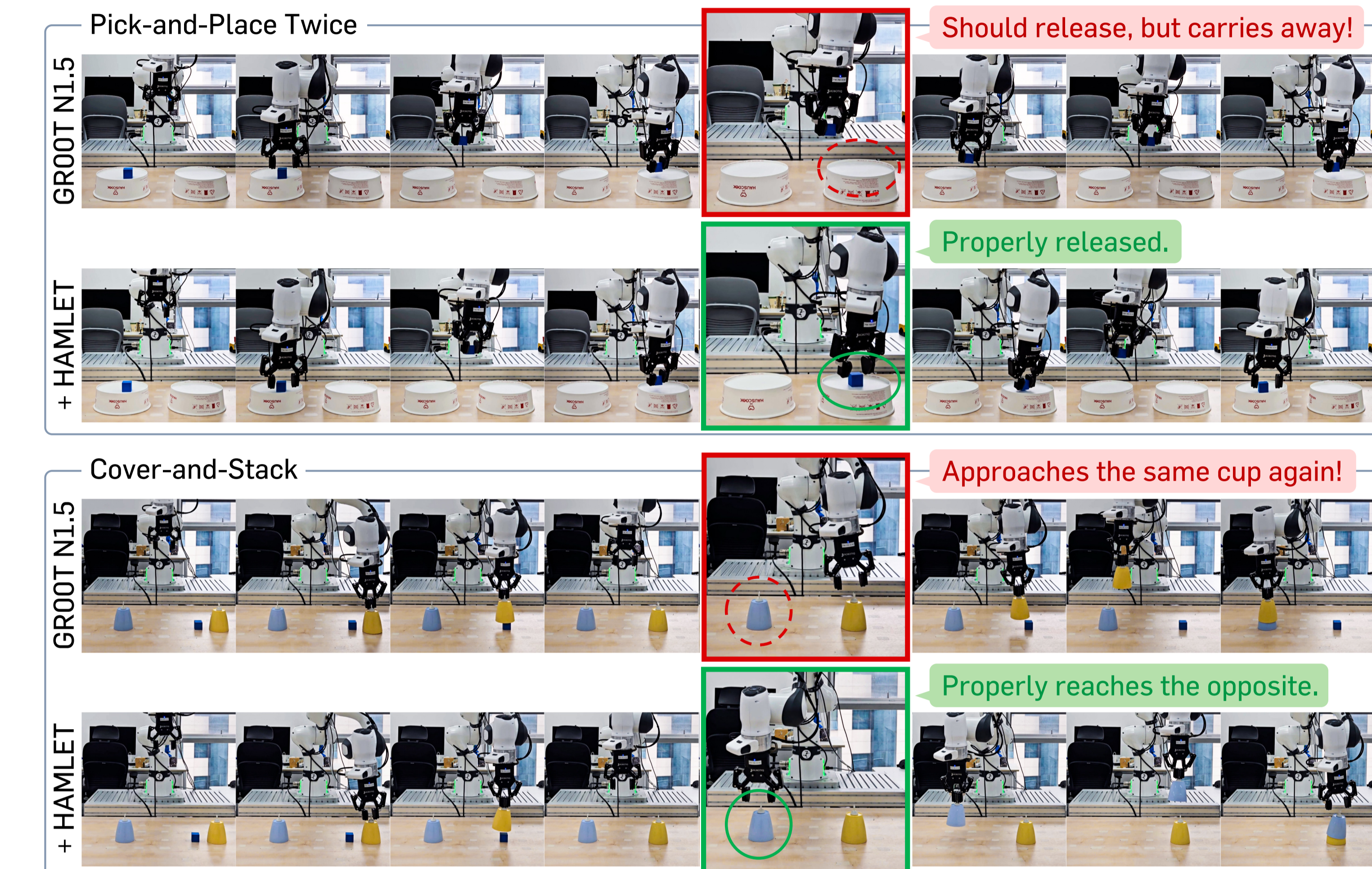
- **Architecture:** A lightweight Transformer that aggregates the cached moment tokens across the history horizon.
- **Mechanism:** Selectively attends to critical past moments via causal self-attention, producing a history-augmented feature.

$$\mathbf{M}' = [\mathbf{m}'_{t-k(T-1)}; \dots; \mathbf{m}'_{t-k}; \mathbf{m}'_t] \in \mathbb{R}^{L \times d}$$

$$[\mathbf{a}_t, \mathbf{a}_{t+1}, \dots, \mathbf{a}_{t+k-1}] = \mathcal{A}_\psi([\mathbf{h}_t; \tilde{\mathbf{m}}', \mathbf{s}_t]), \text{ where } \tilde{\mathbf{m}}' \text{ is the history-augmented feature}$$

Experiments

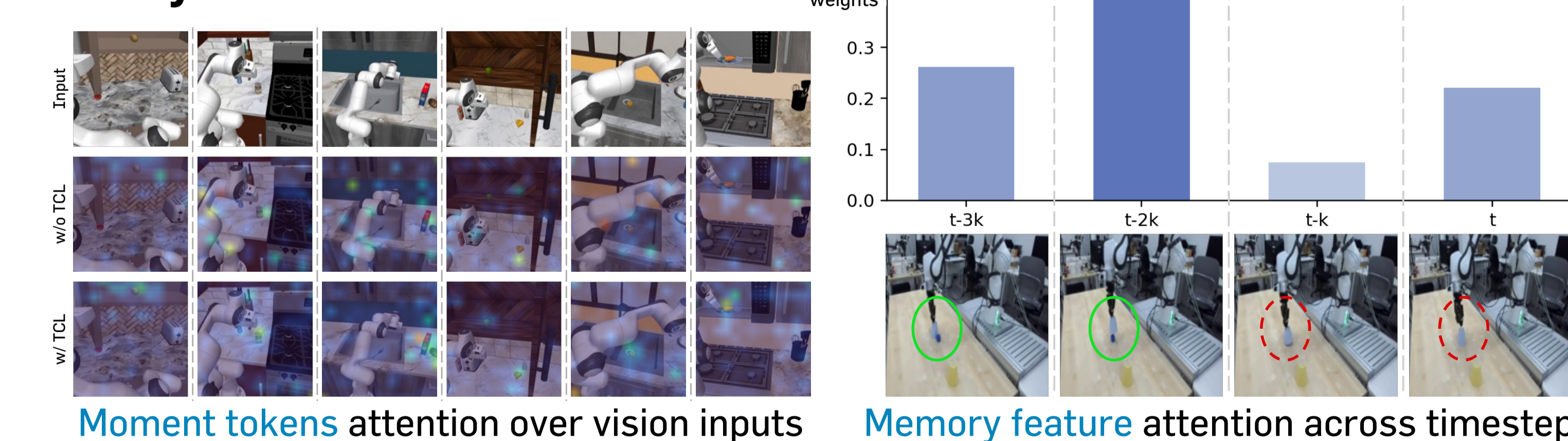
Qualitative Results



Quantitative Results

Method	History?	Pick-and-Place Twice		Cover-and-Stack		Swap Cubes		Avg.
		PnP Once	Success	Cover Cube	Success	Stage Cube	Success	
π_0	✗	54.2	25.0	87.5	58.3	83.3	12.5	31.9
π_0 -FAST	✗	37.5	20.8	54.2	12.5	66.7	4.2	12.5
GROOT N1	✗	54.2	25.0	79.2	33.3	75.0	33.3	30.6
GROOT N1.5	✗	54.2	12.5	62.5	37.5	87.5	37.5	29.2
+ Multi-frame	✓	79.2	45.8	70.8	33.3	91.7	58.3	45.8
+ HAMLET (Ours)	✓	91.7	66.7	95.8	79.2	95.8	83.3	76.4

Analysis



Moment Token	TCL	Memory Module	Avg.	RoboCasa Kitchen		LIBERO				
				100 demo	Spatial	Object	Goal	Long	Avg.	
✗	✗	✗	62.6	64.8	99.4	99.8	98.0	87.8	96.2	
✗	✗	✗	63.1	63.8	99.4	98.6	97.0	87.2	95.5	
✓	✓	✗	63.4	64.8	98.6	99.4	98.2	91.6	96.9	
✓	✓	✗	64.8	64.9	99.0	99.0	97.8	90.2	96.5	
✓	✓	✓	65.4	TCL (Ours)	65.4	99.0	100.0	99.2	92.2	97.7

Token Length	Avg.	Method	Avg.	Method	History Length	Latency (ms, ↓)	Peak memory (MB, ↓)
1	64.3	No Memory	62.6	GROOT N1.5	1	80.5 (1.00x)	289 (1.00x)
4	65.4	Moment Concat.	62.7	+ Multi-frame	4	108.5 (1.35x)	1051 (3.64x)
8	66.4	RNN	64.5	+ HAMLET (Ours)	4	82.4 (1.02x)	566 (1.96x)
16	65.9	LSTM	65.0	+ Multi-frame	8	193.0 (2.40x)	2023 (7.00x)
32	62.7	GRU	64.3	+ HAMLET (Ours)	8	85.8 (1.07x)	578 (2.00x)
64	62.5	Transformer	65.4				