



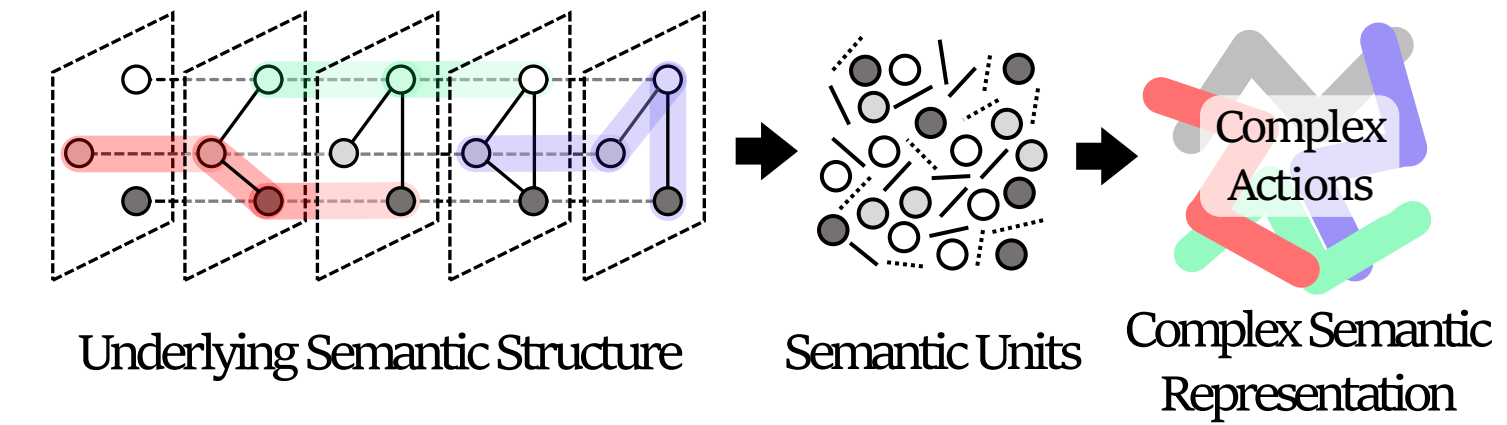
Introduction

Research Problem

- ▶ How to understand **multi-granular semantic structures** encompassing objects, scenes and video-wide contexts?
- ▶ How to understand video composed of **complex semantic structures**?

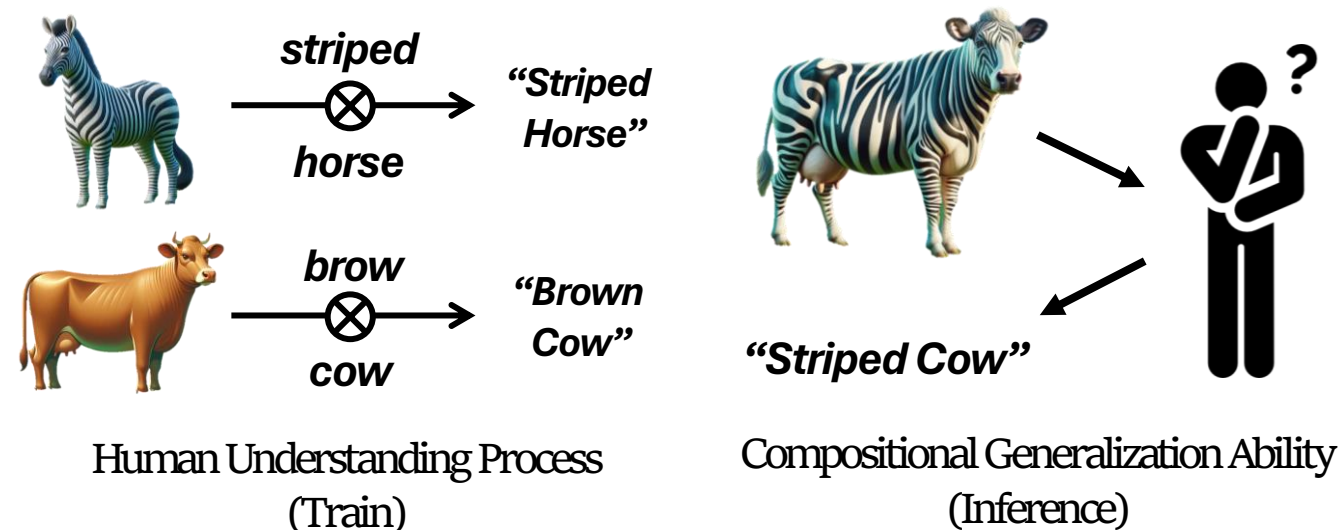
Key Insights

- ▶ **Spatiotemporal symbolic graph** as an input representation of a given video
- ▶ **Tokenizing** spatiotemporal graph and getting **semantic unit representation**
- ▶ Semantic units are then **combined to represent more complex semantics**



Compositional Generalization Ability

- ▶ Ability to understand conceptual combinations unseen in the training process.



Contributions

- ▶ A novel **object-oriented video understanding method** to learn the **multi-granular semantic structure of long videos** is suggested
- ▶ A novel **data split for compositional generalization test** of video understanding algorithms is proposed
- ▶ In experiments with both synthetic and real-world videos, we achieve new state-of-the-art performances

Proposed Method

Compositional Learning Framework

Spatiotemporal Graph Construction → Spatiotemporal Graph Transformer → Object-oriented Video Encoder → Embedding Disentangling Module

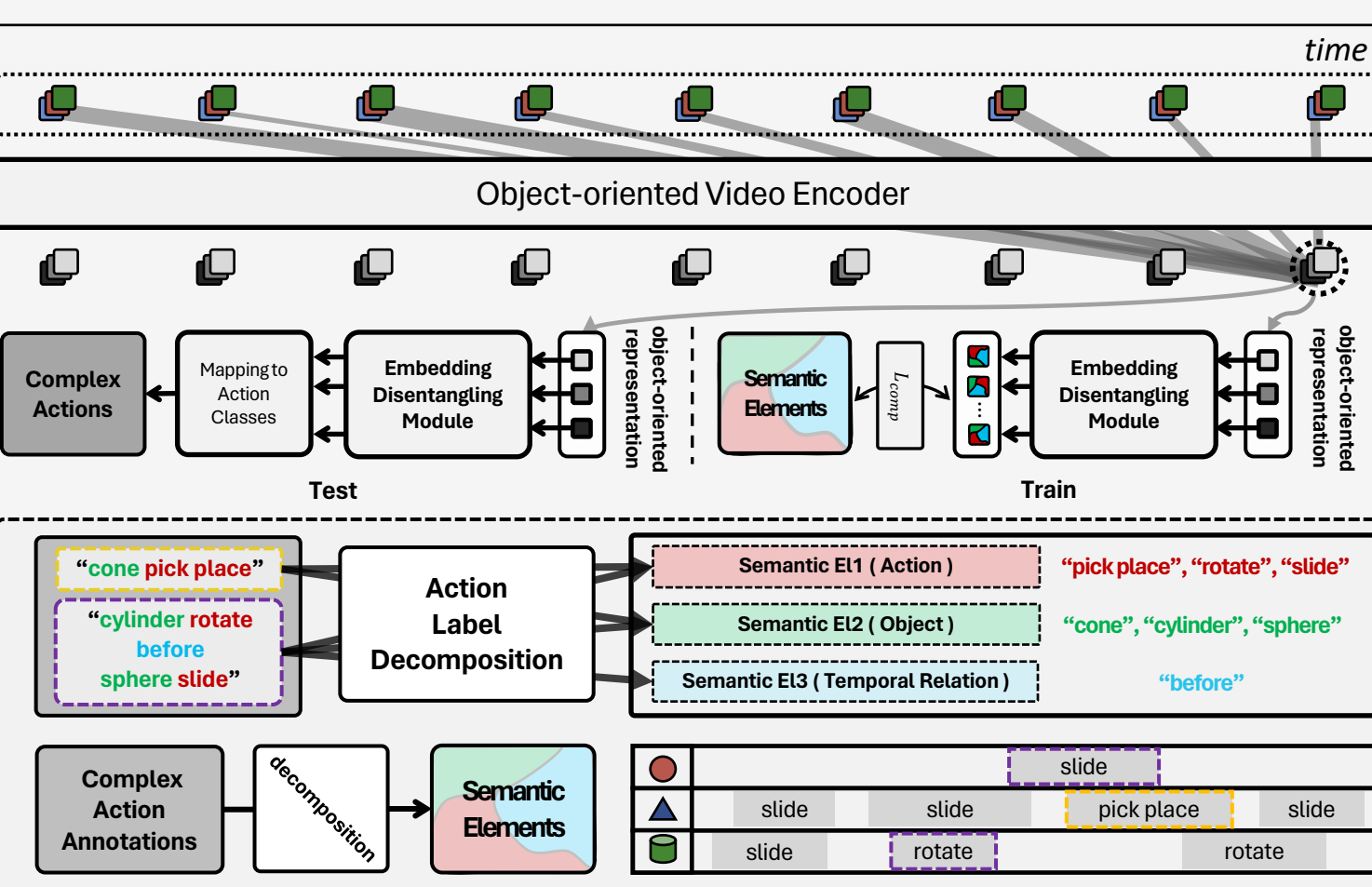
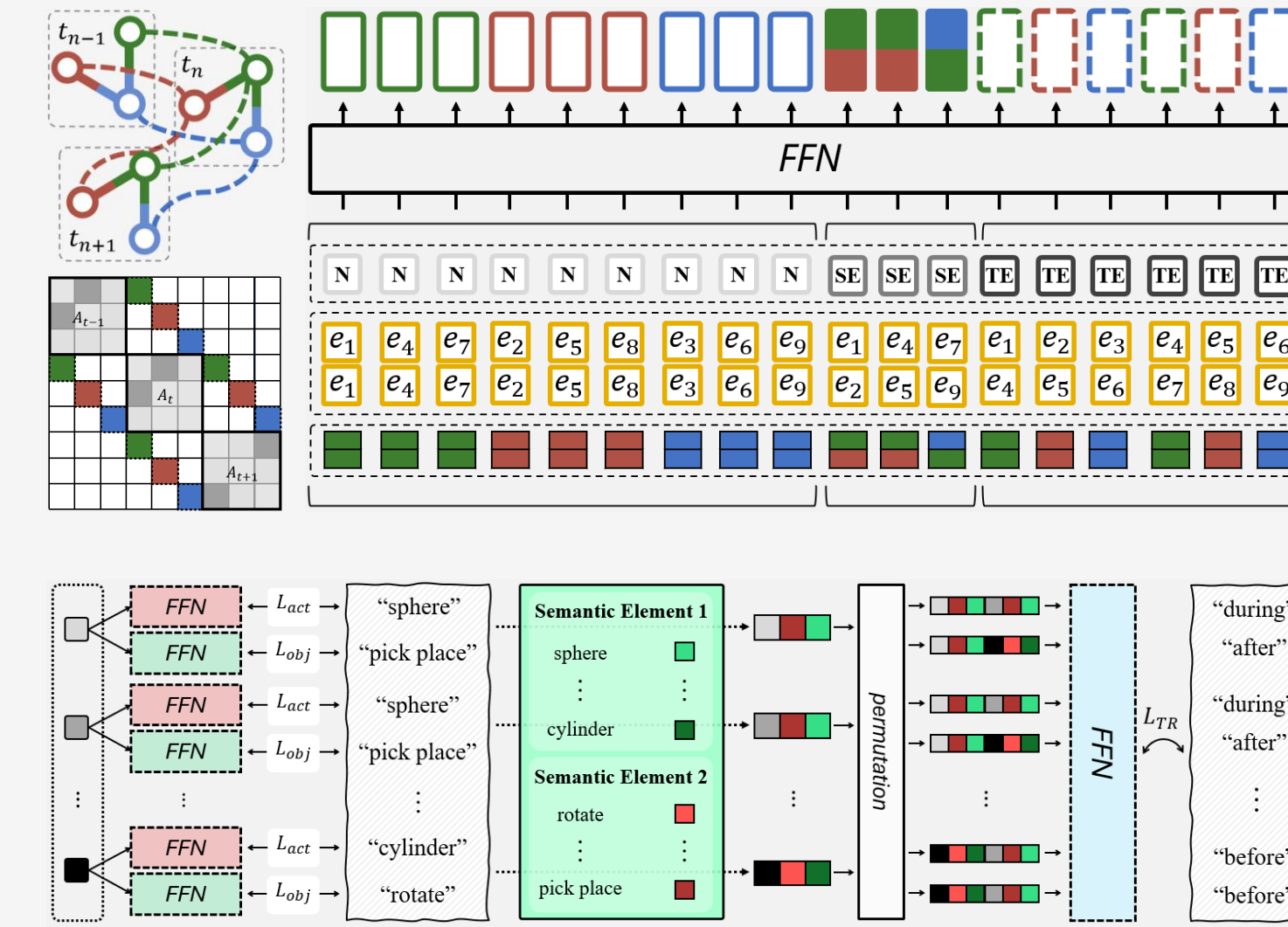
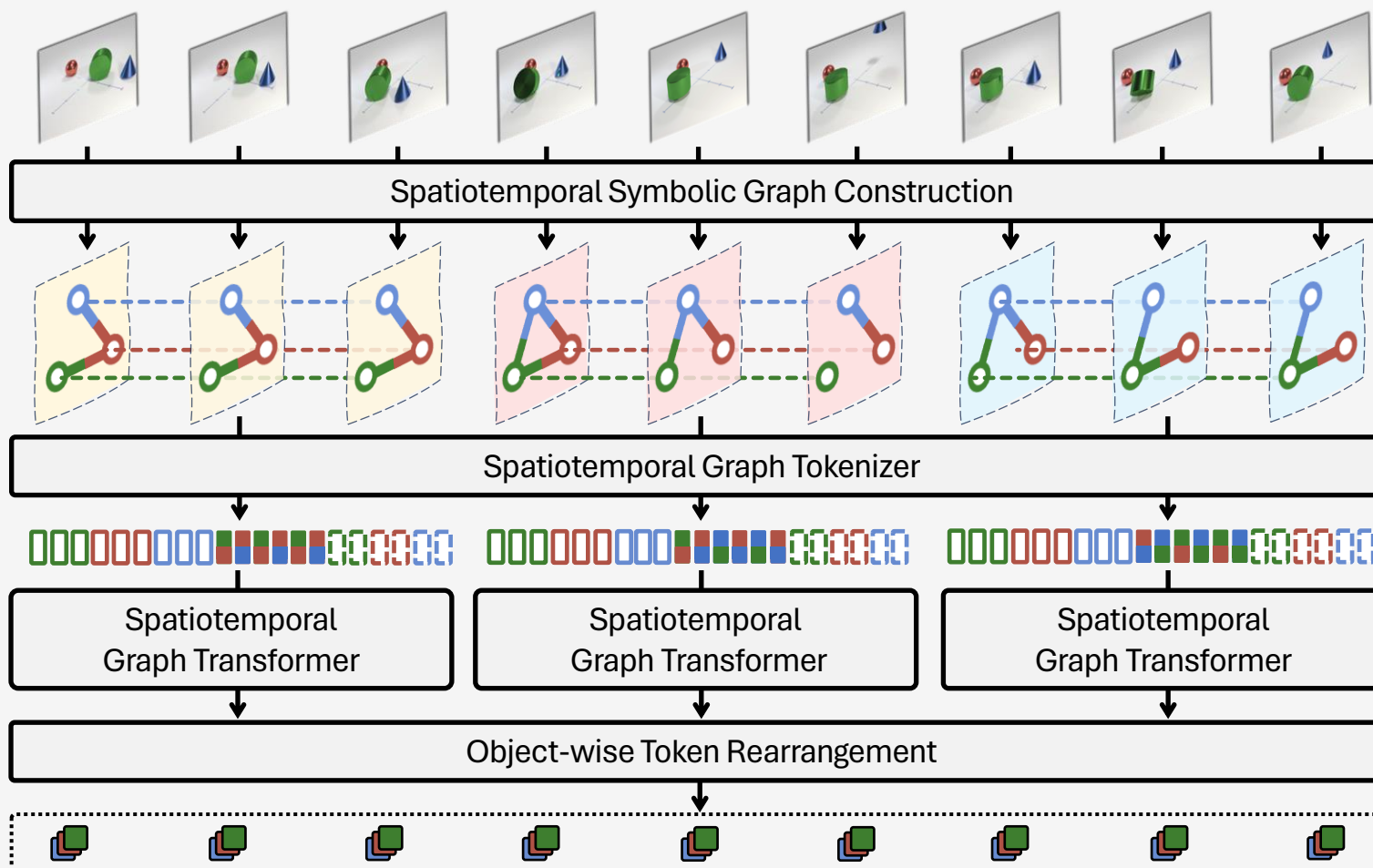
Spatiotemporal Graph Tokenizer & Graph Transformer

- ✓ All elements (**node**, **spatial edge**, **temporal edge**) are **tokenized** to be fed into the transformer as inputs
- ✓ We introduce **total adjacency matrix** which involves **temporal auxiliary edges** describing the temporal connection between two adjacent time
- ✓ As a positional embedding vector, **graph Laplacian of total adjacency matrix** is adapted

Embedding Disentangling Module

- ✓ To achieve compositional generalization ability, we introduce a method where **predictions are made for each semantic elements**. (e.g. object, action, temporal relation)
- ✓ To resolve this, **entangled feature embeddings are projected onto each independent subspace**. (implemented by learnable embedding)
- ✓ By concatenating each embedding, classification head is applied to **predict semantic labels that was previously decomposed**.

Overall Architecture



Experiments

Evaluation Methods

- ▶ Action Recognition with Synthetic Video Dataset (by CATER)
- ▶ Action Recognition with Real-world Video Dataset (by MOMA-LRG)
- ▶ **Compositional Generalization Test** (by newly suggested data splits)

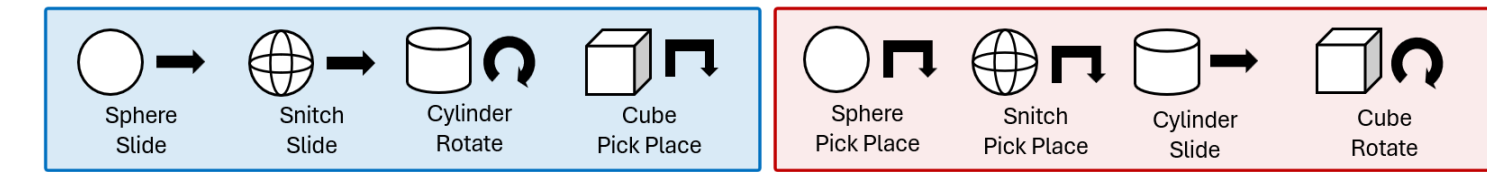
Compositional Generalization Test & Suggested Data Split

Train Phase

- Each of individual semantic elements is seen. (ex. “sphere”, “pick place” is seen)
- But specific combinations of semantic elements is unseen. (ex. “sphere pick place” is unseen)

Test Phase

- Input video containing “sphere pick place” is given



Trainset & Testset Example of Object + Action Combination

Quantitative Results

Method	mAP	Method	mAP	# params	Inference time (ms)
VIVIT [1]	25.45	VIVIT [1]	65.09	99M	23
R3D [4, 12]	44.2	FROZEN [2]	69.36	115M	40
R3D + NL [37]	45.9	Ours	72.83	7.3M	20
R3D + LSTM [10]	53.4				
R3D + NL + LSTM [10]	53.1				
SC3D + LSTM [29]	66.71				
Single stream SC3D + LSTM [29]	69.76				
VIVIT [1]	66.18				
FROZEN [2]	66.64				
Ours	75.40				

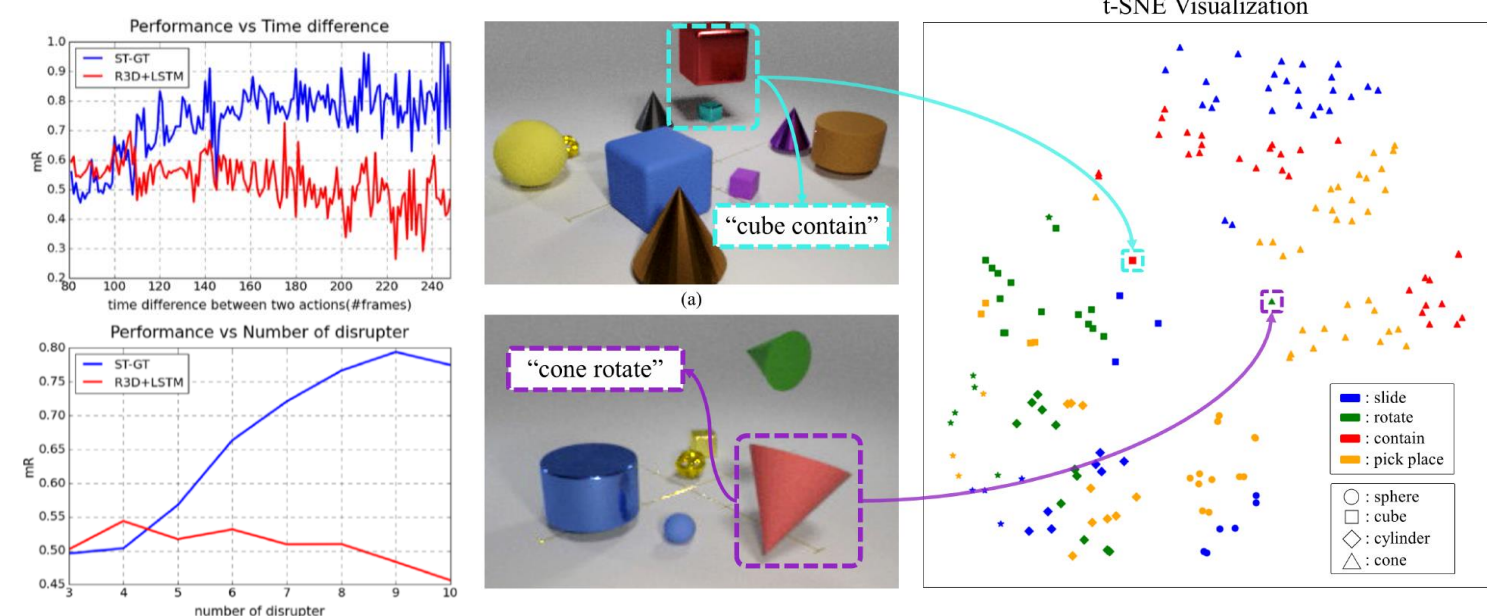
Method	mAP	Task 1	Task 2
R3D [4, 12]	98.8	val	test
R3D + NL [37]	98.9	90.42	90.38
Ours	99.88	99.89	99.88

Method	mAP	# params	Inference time (ms)
VIVIT [1]	65.09	99M	23
FROZEN [2]	69.36	115M	40
Ours	72.83	7.3M	20

Method	mAP	Task 1	Task 2
R3D [4, 12]	98.8	val	test
R3D + NL [37]	98.9	90.42	90.38
Ours	99.88	99.89	99.88

Action Recognition Performance & Ablation Study

Qualitative Results



Acknowledgement

This work is supported by IITP grant funded by MSIT (Grant No. 2022-0-00264/40%, 2022-0-00612/20%, 2022-0- 00951/20%), and IITP Artificial Intelligence Graduate School Program for Hanyang University funded by MSIT (Grant No. RS-2020-0-0120137/20%).