



# Modeling Fine-Grained Hand-Object Dynamics for Egocentric Video Representation Learning

Baoqi Pei<sup>1,2</sup>, Yifei Huang<sup>3,2</sup>, Jilan Xu<sup>4,2</sup>, Guo Chen<sup>4,2</sup>, Yuping He<sup>5</sup>, Lijin Yang<sup>3</sup>,  
Yali Wang<sup>6,2</sup>, Weidi Xie<sup>7,2</sup>, Yu Qiao<sup>2</sup>, Fei Wu<sup>1</sup>, Limin Wang<sup>5,2</sup>

<sup>1</sup> Zhejiang University <sup>2</sup> Shanghai Artificial Intelligence Laboratory <sup>3</sup> The University of Tokyo  
<sup>4</sup> Fudan University <sup>5</sup> Nanjing University <sup>6</sup> SIAT <sup>7</sup> Shanghai Jiao Tong University



## I. Introduction

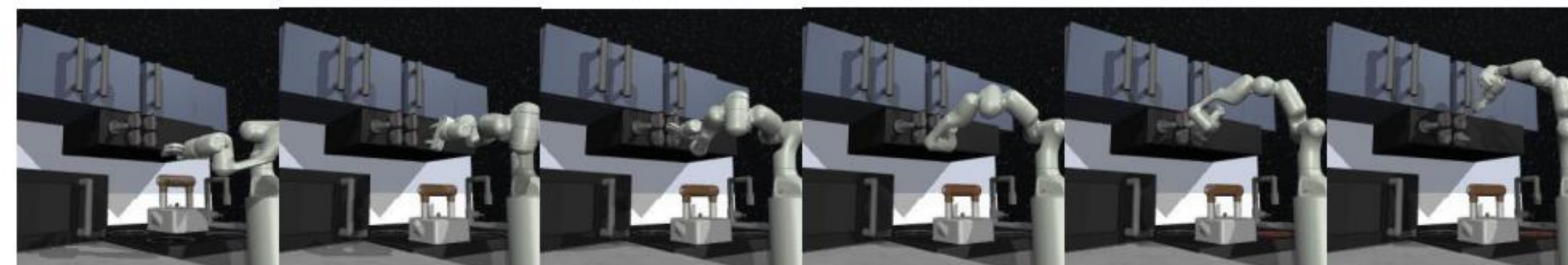
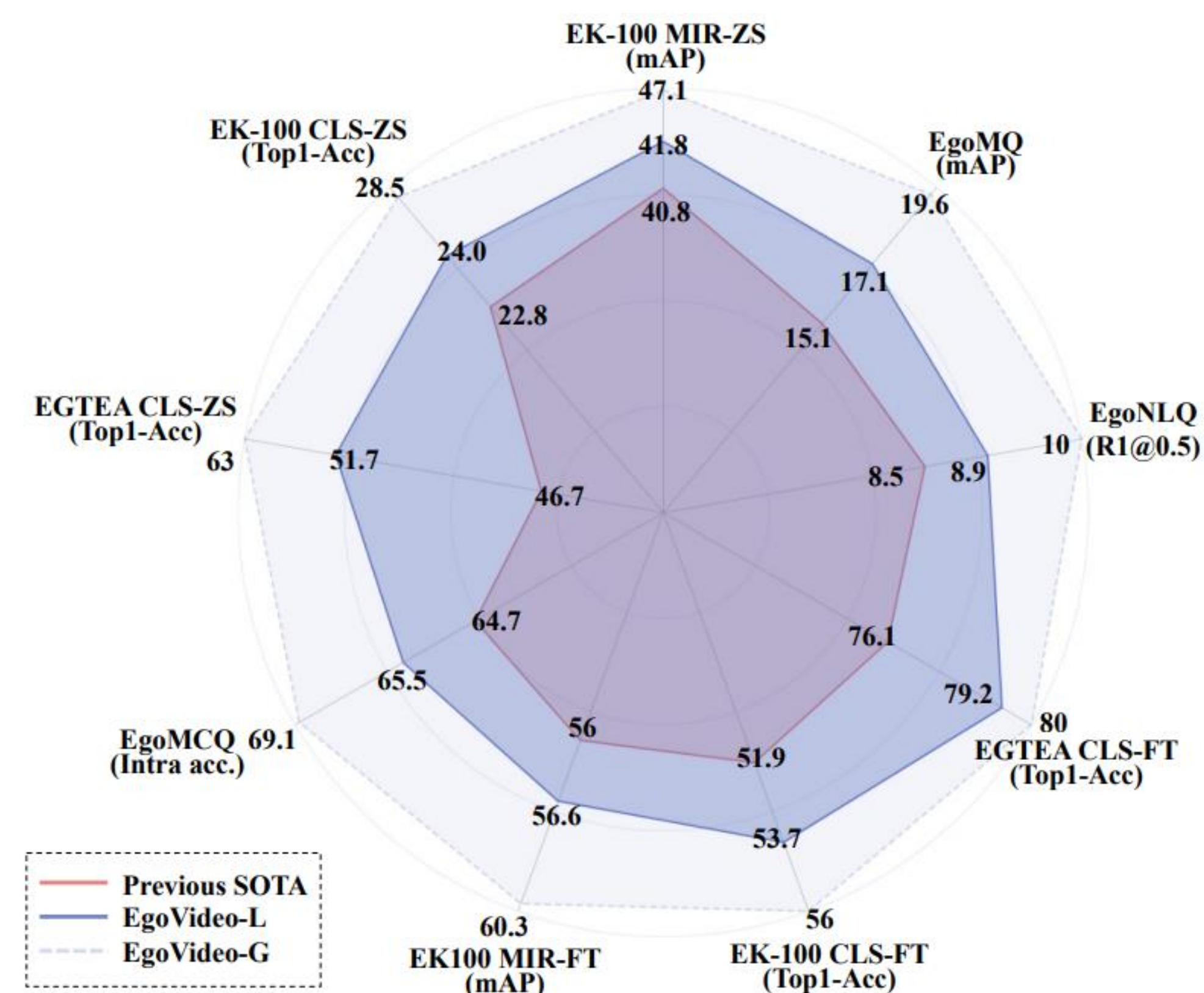
**Motivation:** Original annotations in egocentric video datasets are highly condensed, neglecting a crucial aspect — the fine-grained dynamics of hands and objects.

**Contribution:** We propose a new framework to integrate the modeling of fine-grained hand-object dynamics into the video representation learning process:

- **HOD data pipeline** to generate captions that describe fine-grained hand-object dynamics.
- **EgoVideo model** with a novel lightweight motion adapter and a co-training strategy.
- Get SOTA performance on **12** downstream tasks and robot manipulation tasks.

## III. Experiments

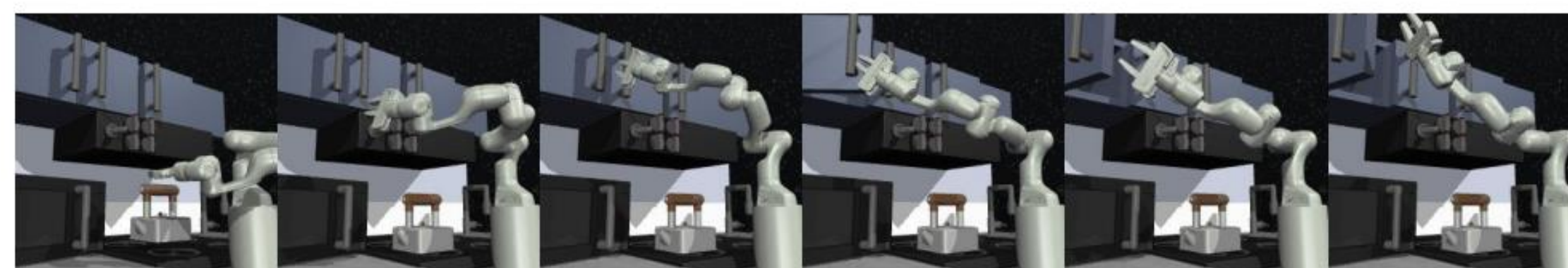
- Zero-shot and fine-tune results on multiple egocentric benchmarks.



Turn on the knob



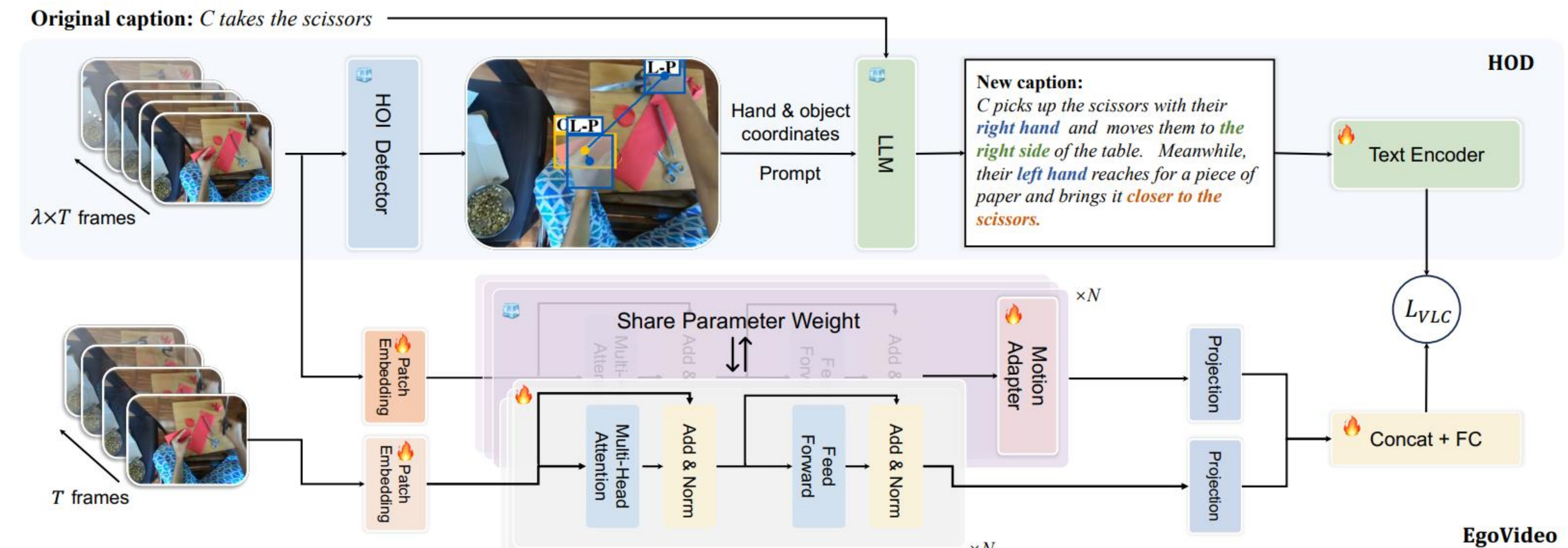
Open microwave



Open door

- Qualitative results on the Franka Kitchen dataset.

## II. Framework



### Data pipeline HOD:

- Use a hand object detector to generate bounding boxes for hands and objects.
- Modified the boxes to a format that LLM understands.
- Use a LLM as Rephraser to enrich the original video captions.

### Egocentric Representation Learning Model: Egovideo

- Comprising a backbone and a motion adapter.
- Adopt a co-training strategy to obtain richer video representations



We got 7 champions in the First Joint  
Egocentric Vision (EgoVis) Workshop!

Github Link:

Data & Training &  
Evaluation

