

TwinVLA

Data-Efficient Bimanual Policy with Twin Single-Arm Vision-Language-Action Models

Hokyun Im^{1,2}, Euijin Jeong¹, Andrey Kolobov², Jianlong Fu², Youngwoon Lee¹

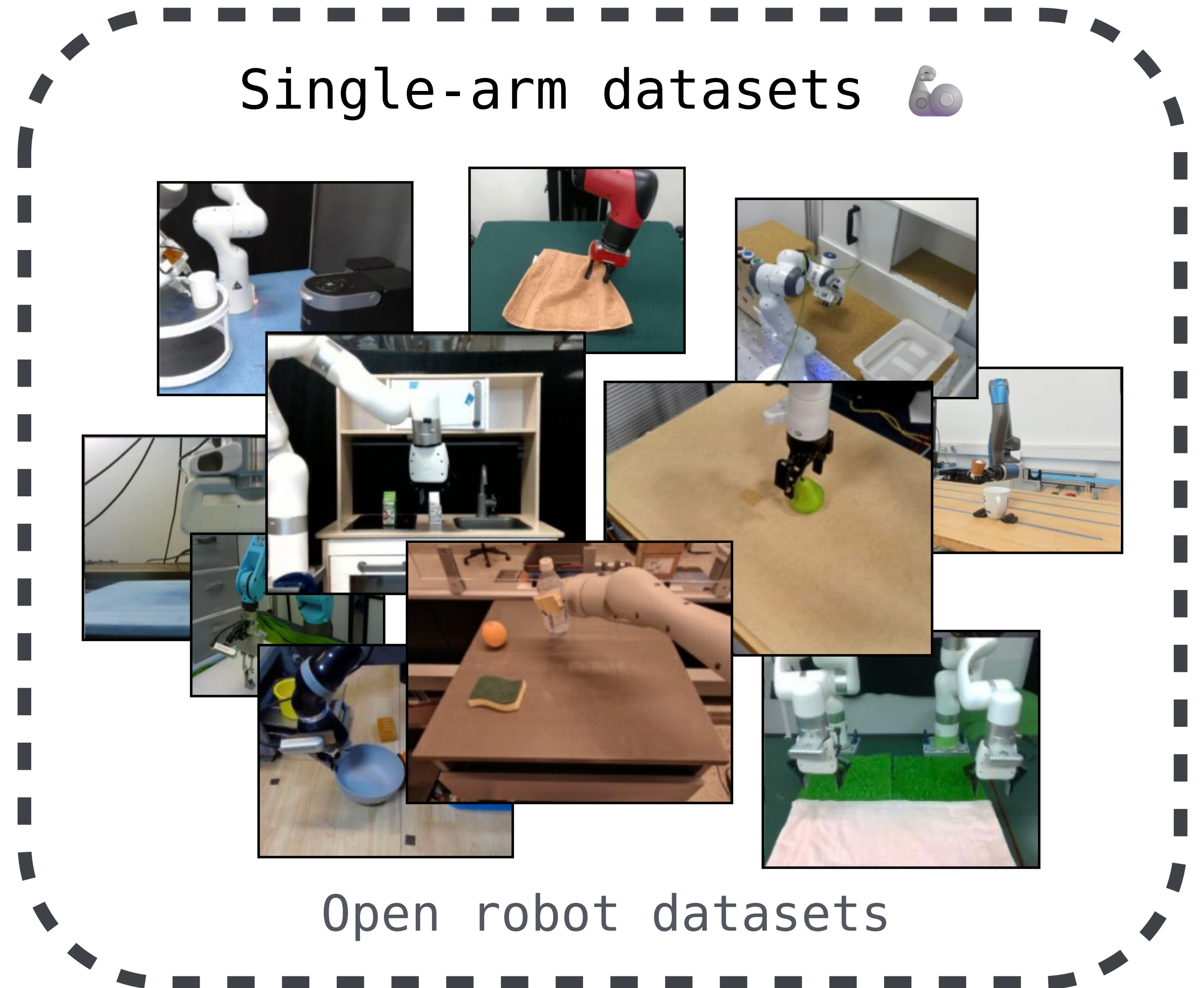
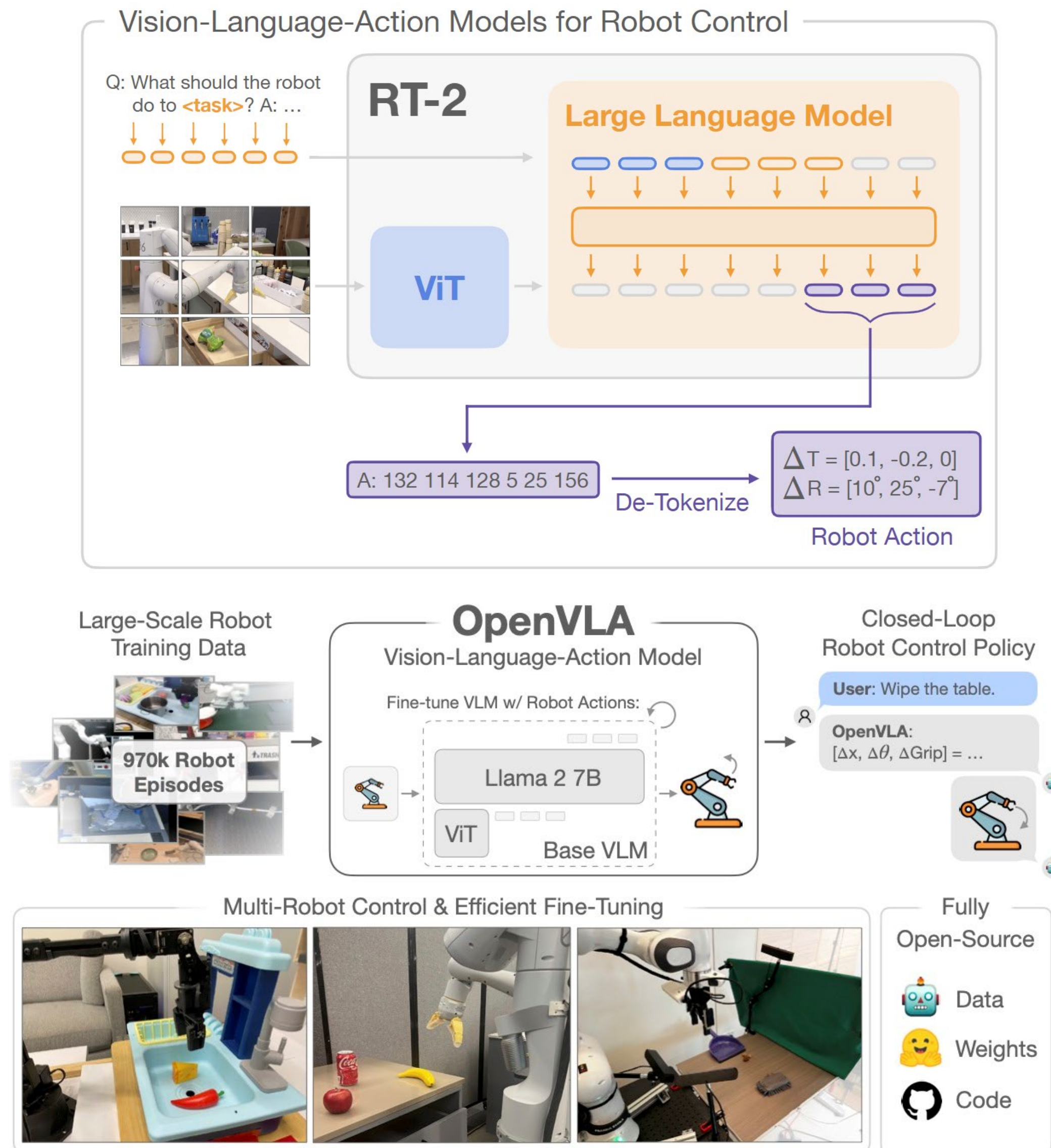


¹Yonsei University,



²Microsoft Research

Large, Diverse Robot Datasets Enable VLAs



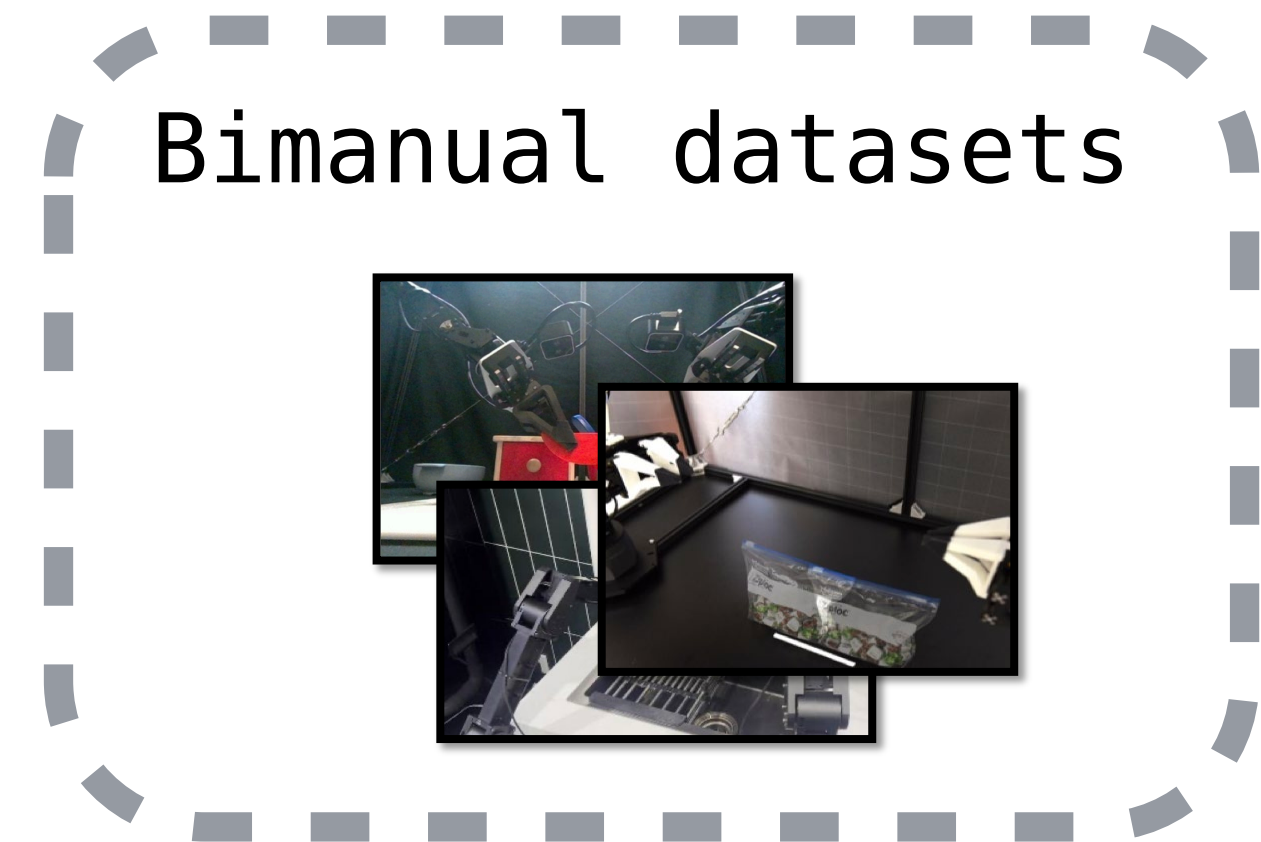
[1] Brohan et. al., "RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control", arXiv 2023

[2] Kim et. al., "OpenVLA: An Open - Source Vision - Language - Action Model", CoRL 2024

[3] OXE collaboration, "Open X - Embodiment: Robotic Learning Datasets and RT - X Models", ICRA 2024

Training Bimanual VLAs Remains Challenging

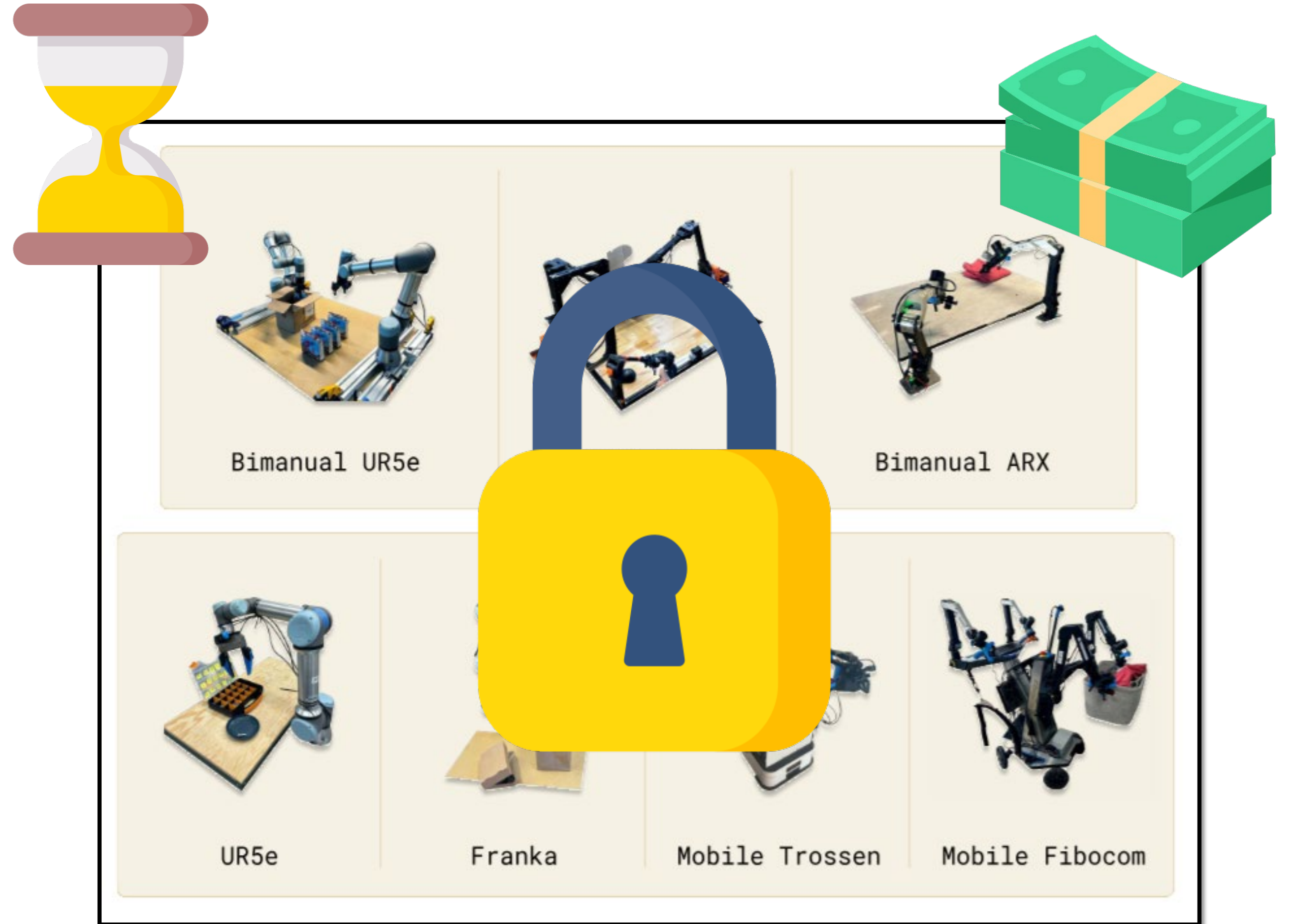
Bimanual VLA



Scarce!

Training Bimanual VLAs Remains Challenging

Bimanual VLA



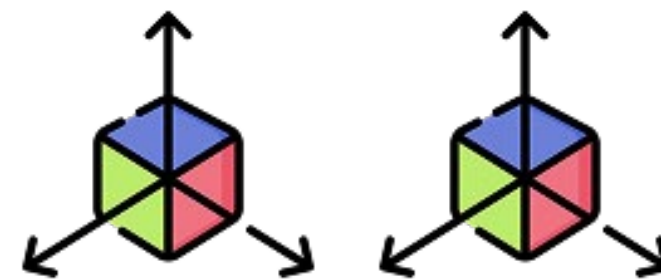
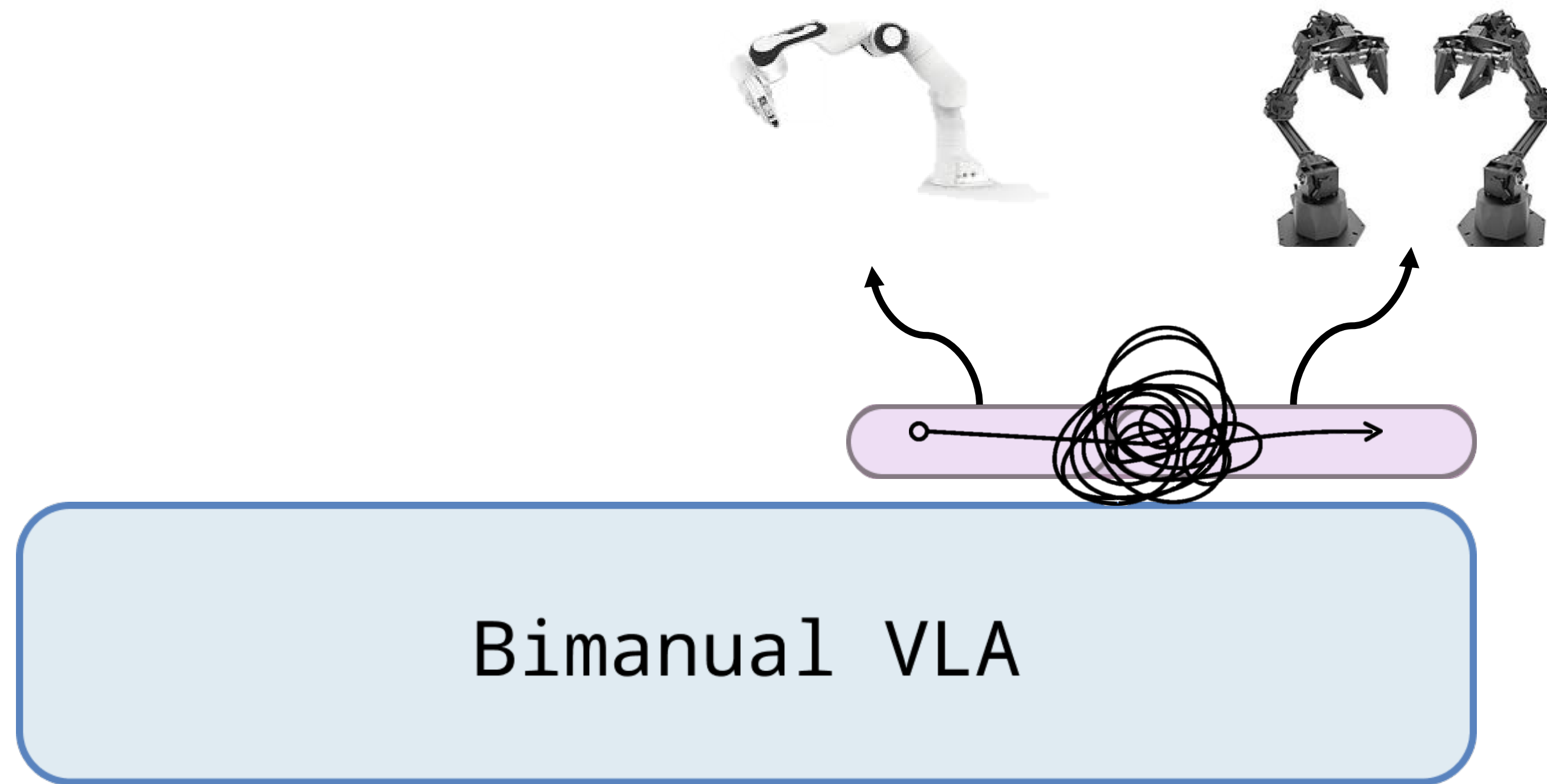
Proprietary datasets

Training Bimanual VLAs Remains Challenging

Bimanual VLA

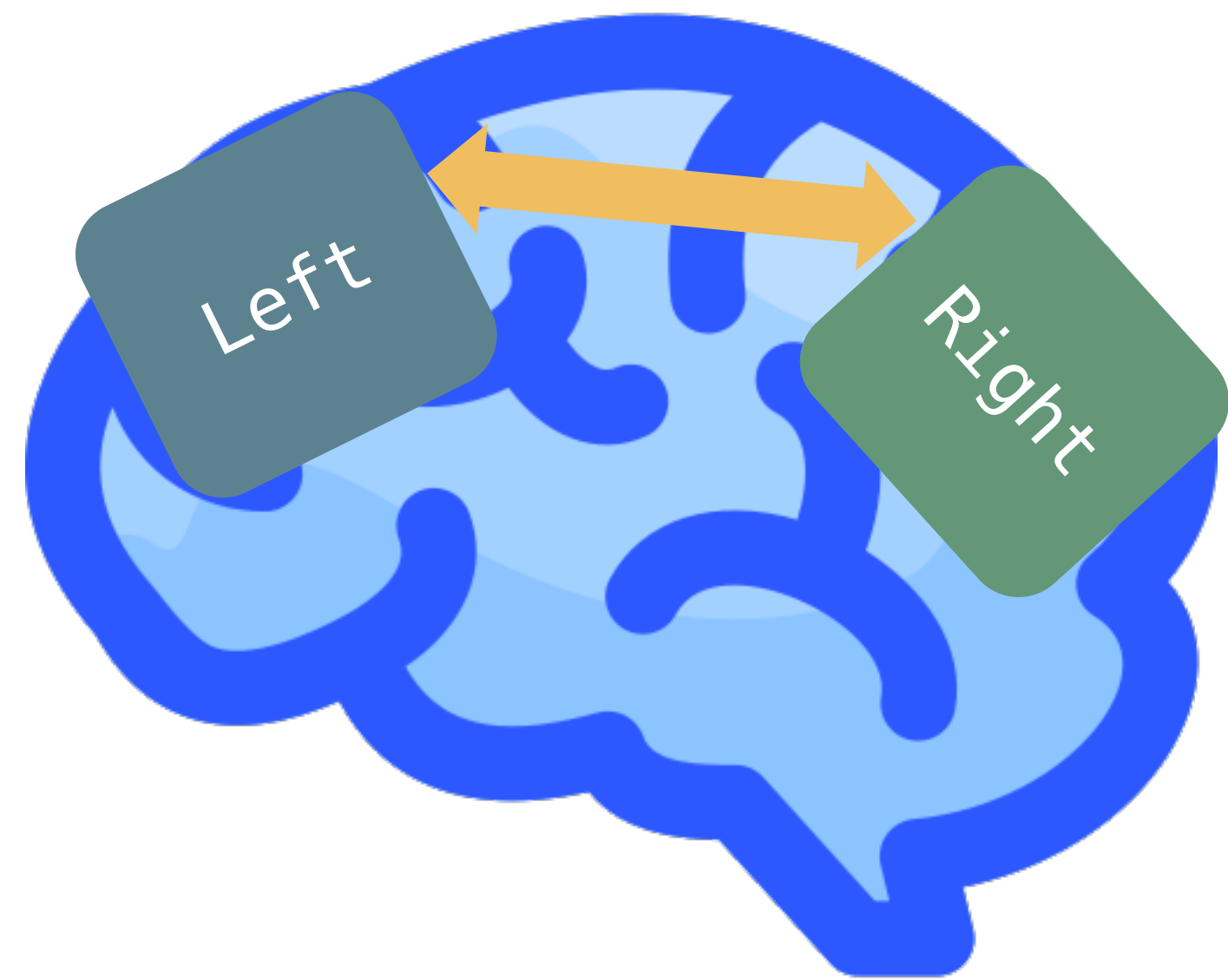
Not reproducible ✖

Training Bimanual VLAs Remains Challenging



Monolithic models

Our Intuition :
Bimanual Manipulation is a Coordination of Single-Arm Skills



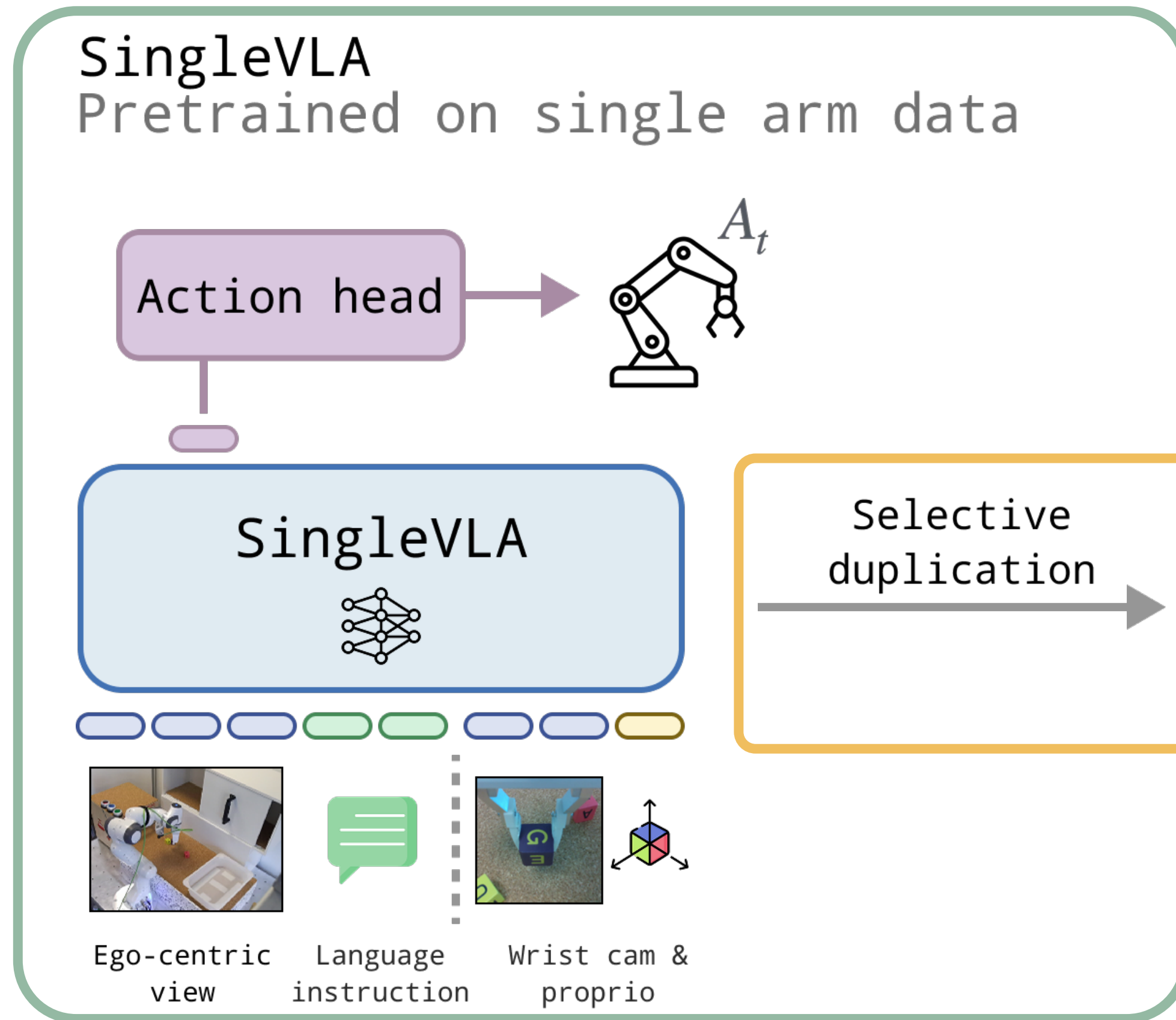
Human Brain



Bimanual Manipulation

Modular perspective on bimanual manipulation

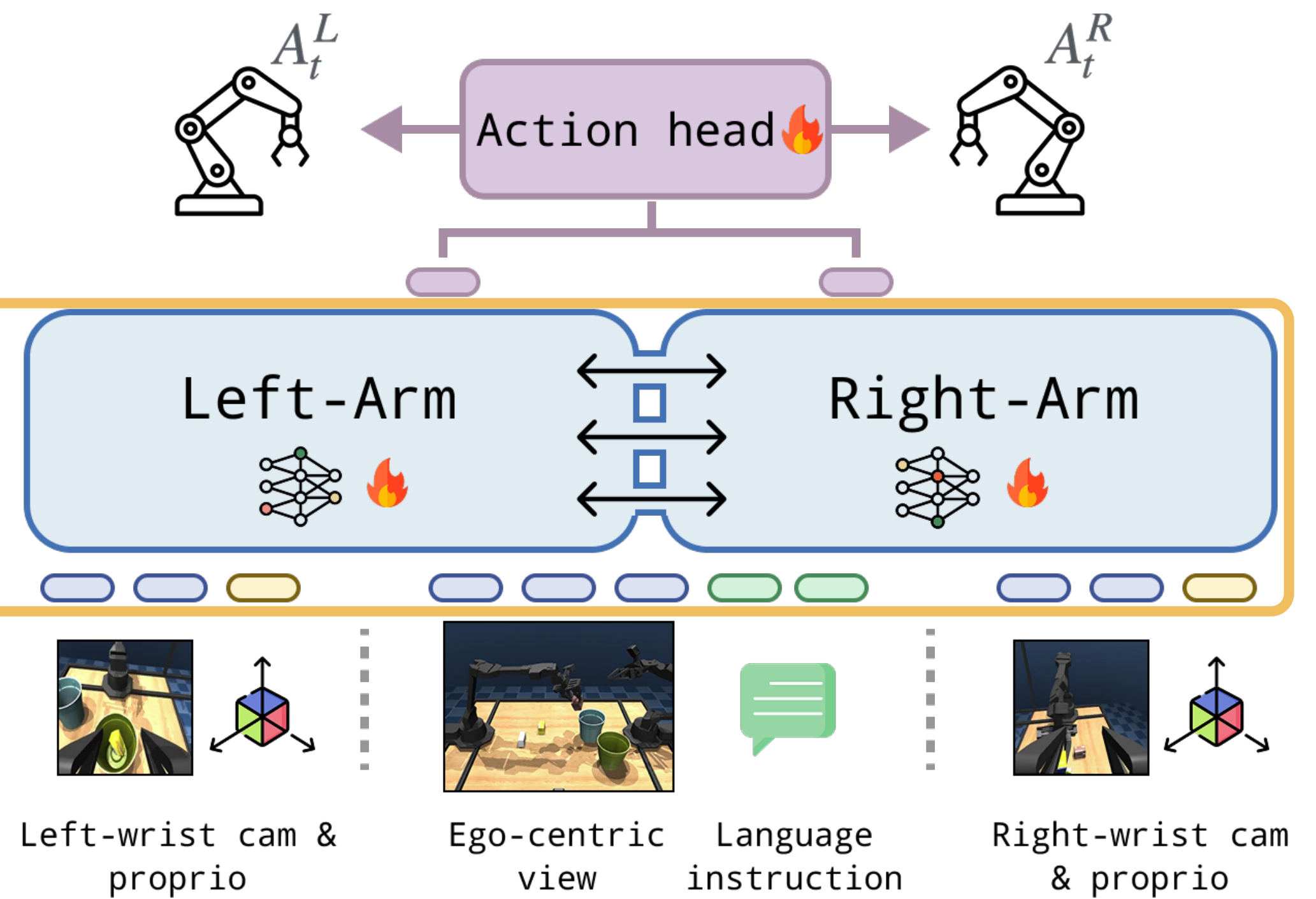
TwinVLA : Data-Efficient Bimanual VLA



1. Single-Arm VLA pretraining

TwinVLA-1.3B
Joint attention & MoE integration

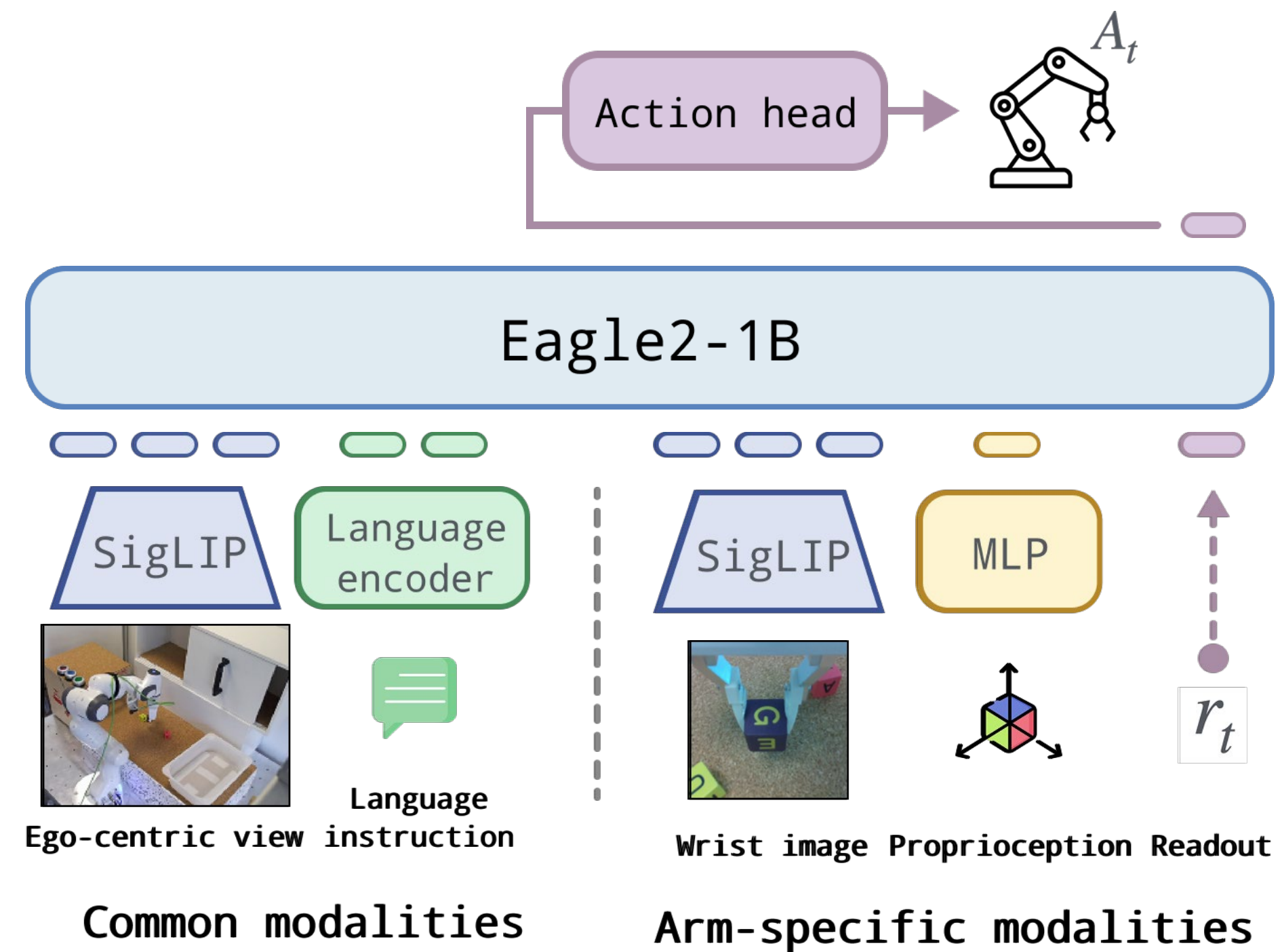
Selective duplication



2. Cross-arm fusion using Joint attention

Cross-arm fusion using Joint Attention

SingleVLA
Pretrained on single arm data



$$o_t^{\text{twin}} = (I_{\text{ego}}, l)_t$$

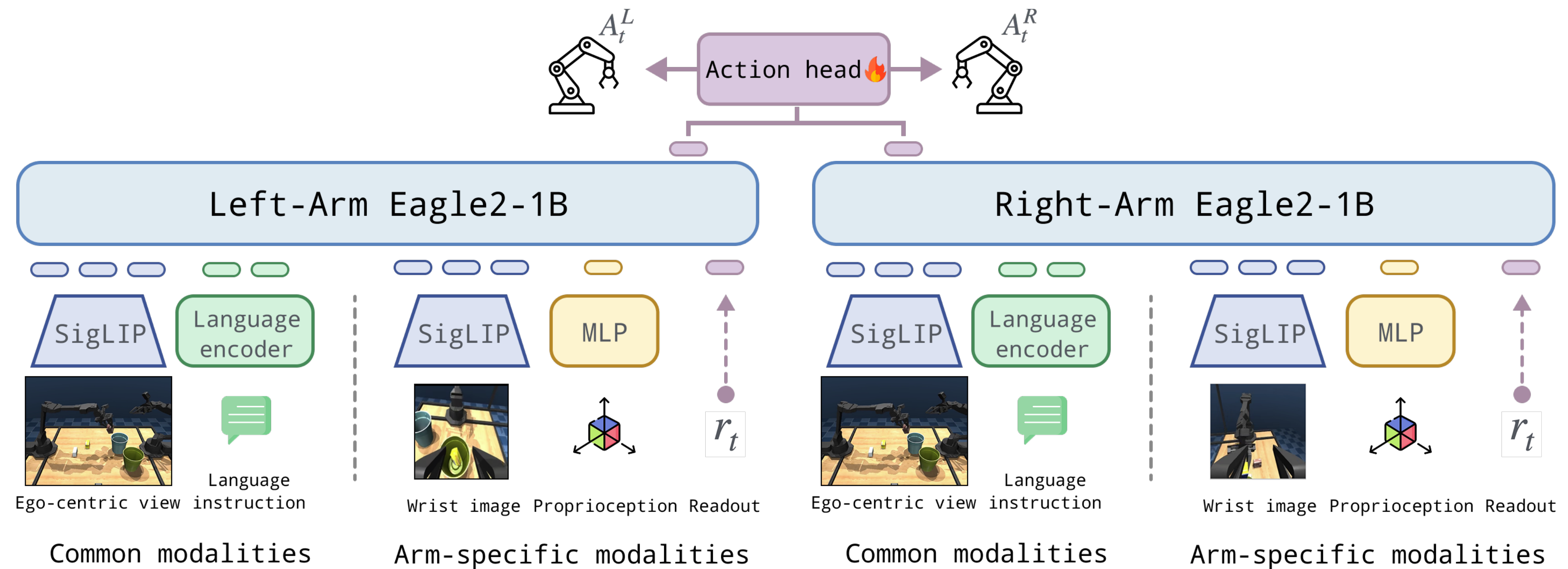
$$(I_{\text{wrist}}^L, d^L)_t$$

$$(I_{\text{ego}}, l)_t$$

$$(I_{\text{wrist}}^R, d^R)_t$$

Cross-arm fusion using Joint Attention

SingleVLA
Pretrained on single arm data



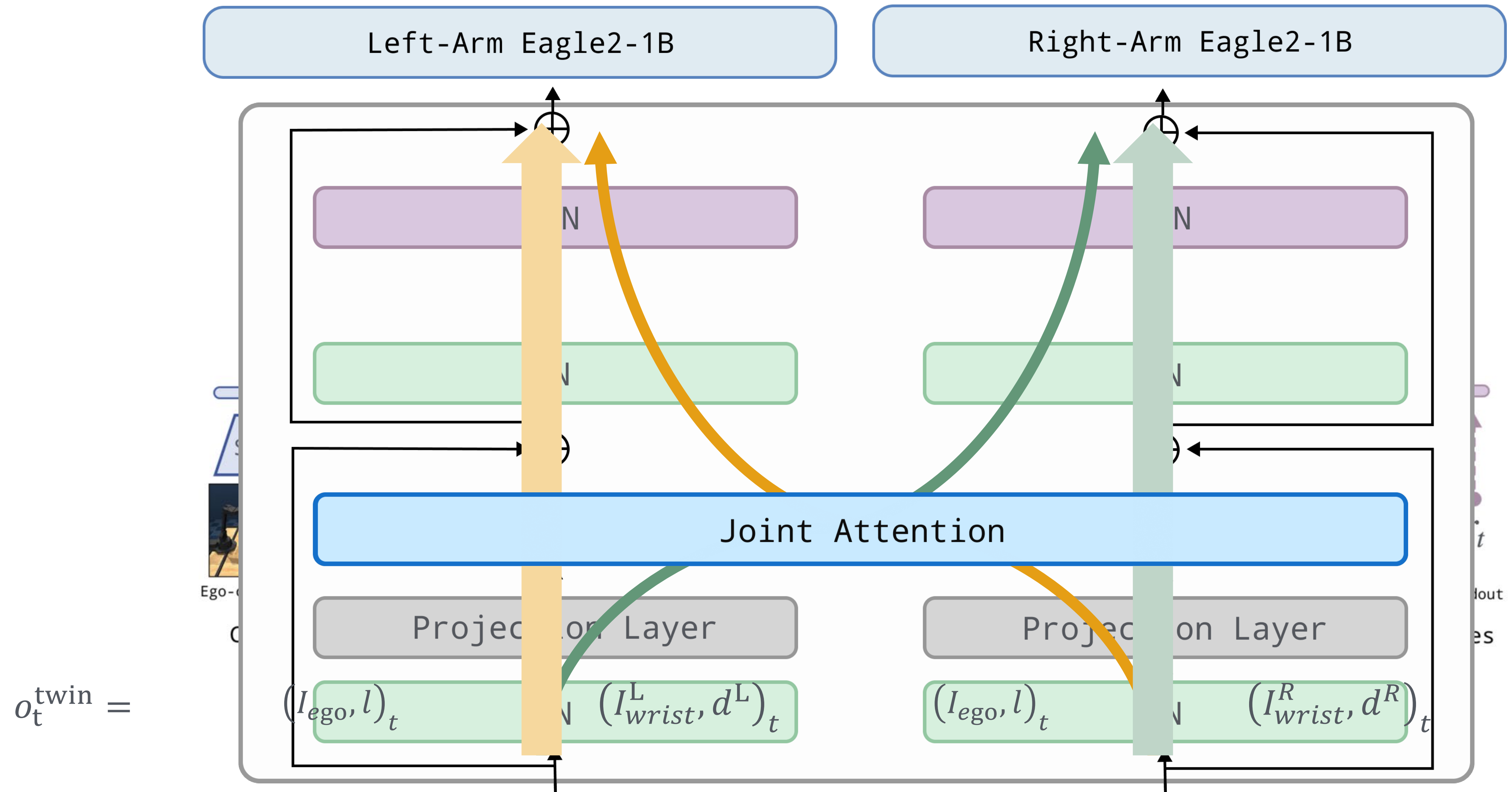
$$o_t^{\text{twin}} = (I_{\text{ego}}, l)_t$$

$$(I_{\text{wrist}}^L, d^L)_t$$

$$(I_{\text{ego}}, l)_t$$

$$(I_{\text{wrist}}^R, d^R)_t$$

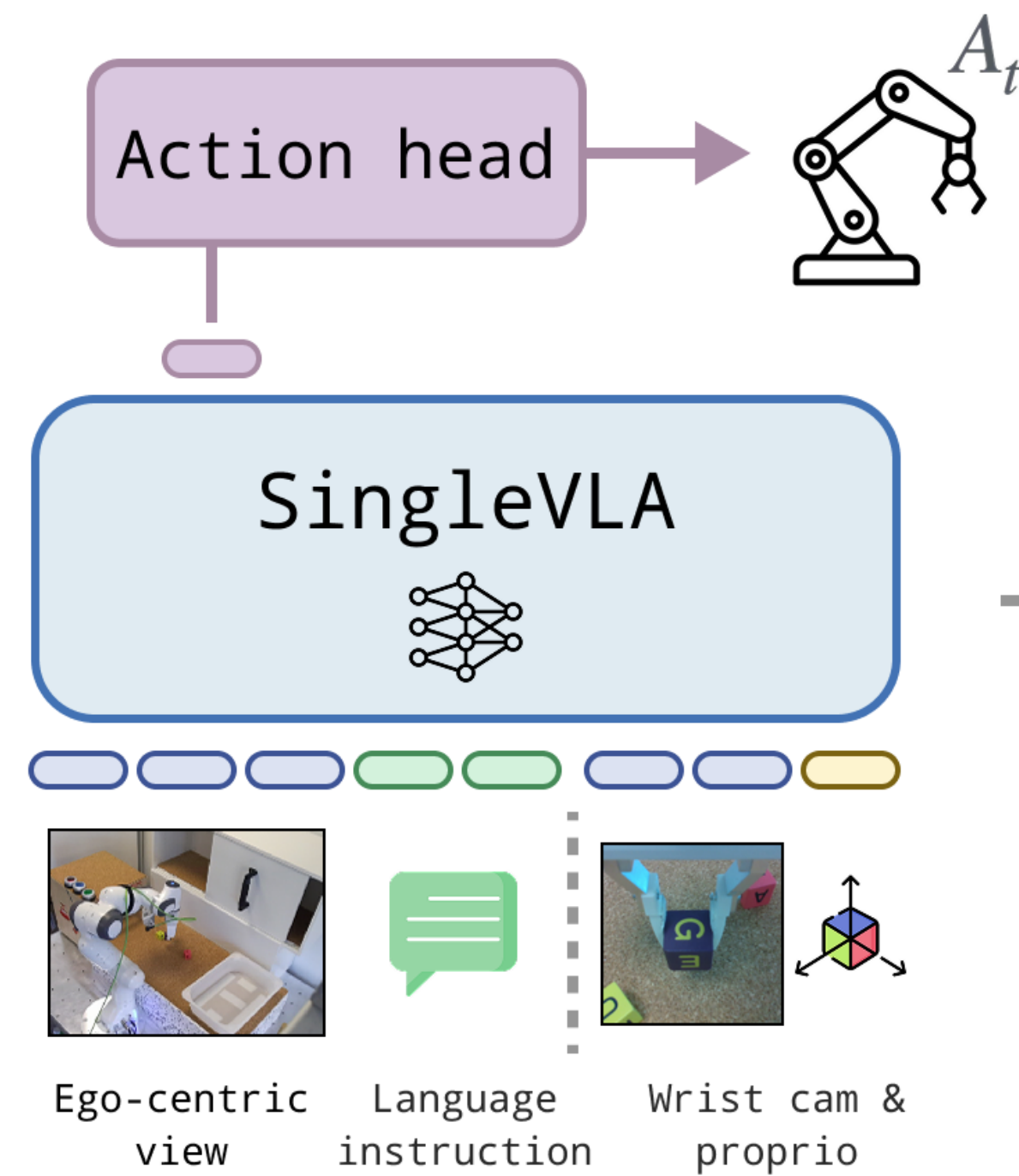
Cross-arm fusion using Joint Attention



TwinVLA : Data-Efficient Twin Single-Arm VLA

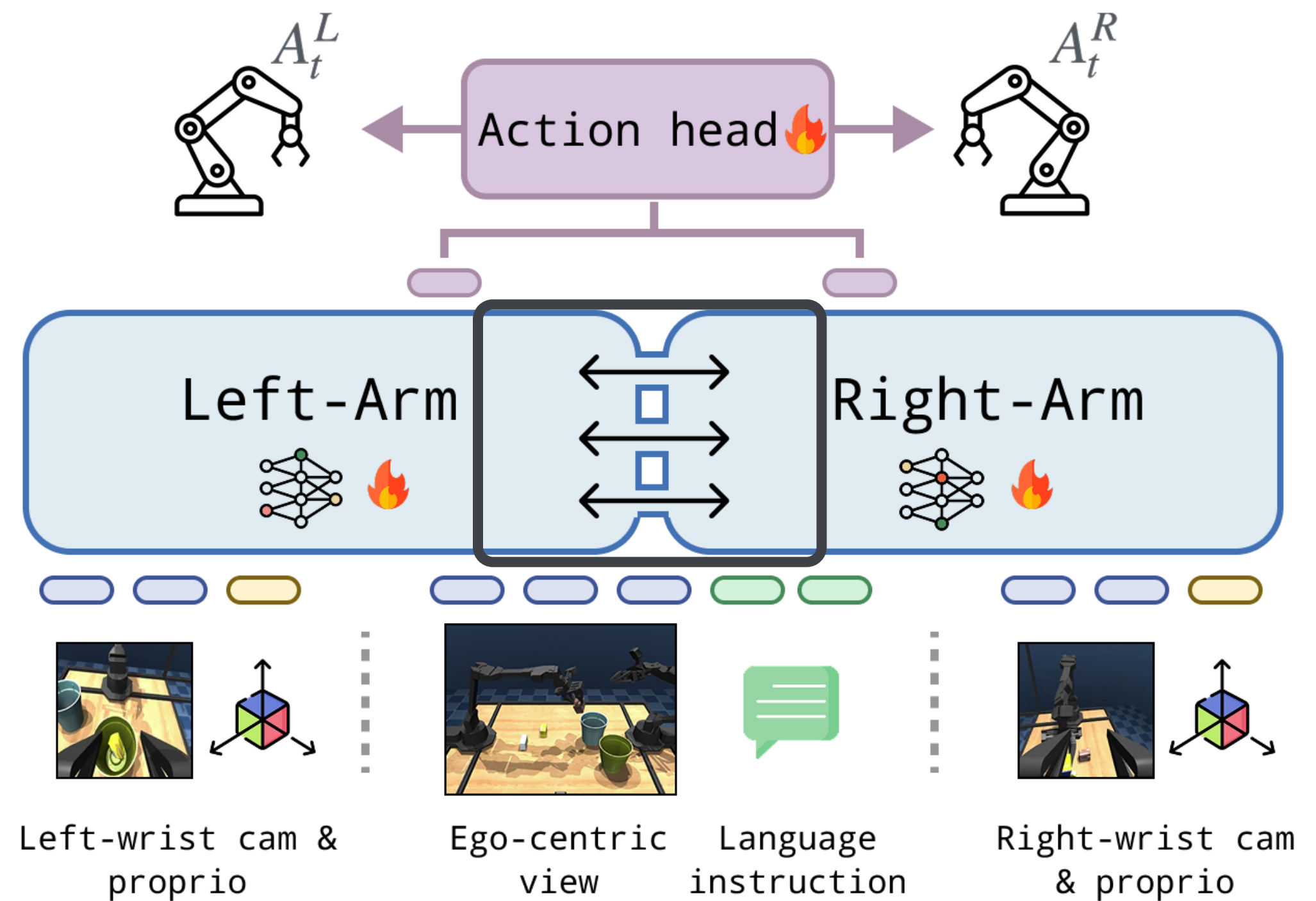
SingleVLA

Pretrained on single arm data



TwinVLA-1.3B

Joint attention & MoE integration

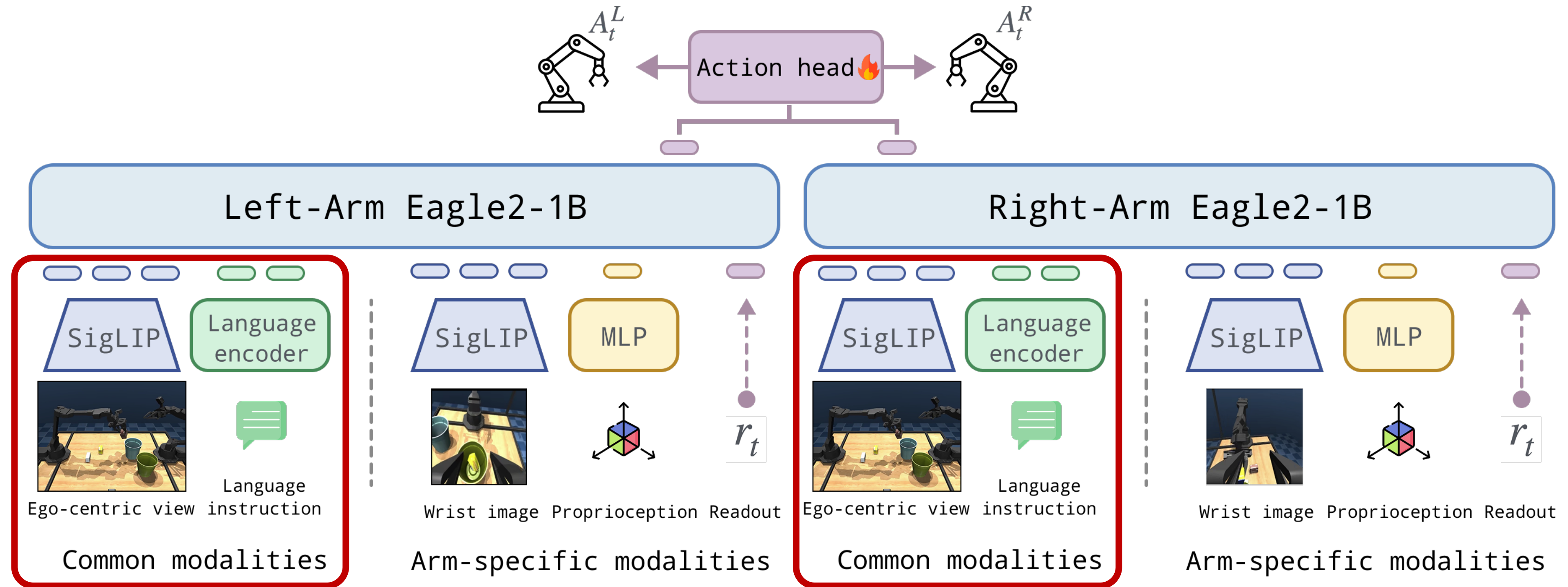


1. Single-Arm VLA pretraining

2. Cross-arm fusion using Joint attention

3. MoE for compute efficiency

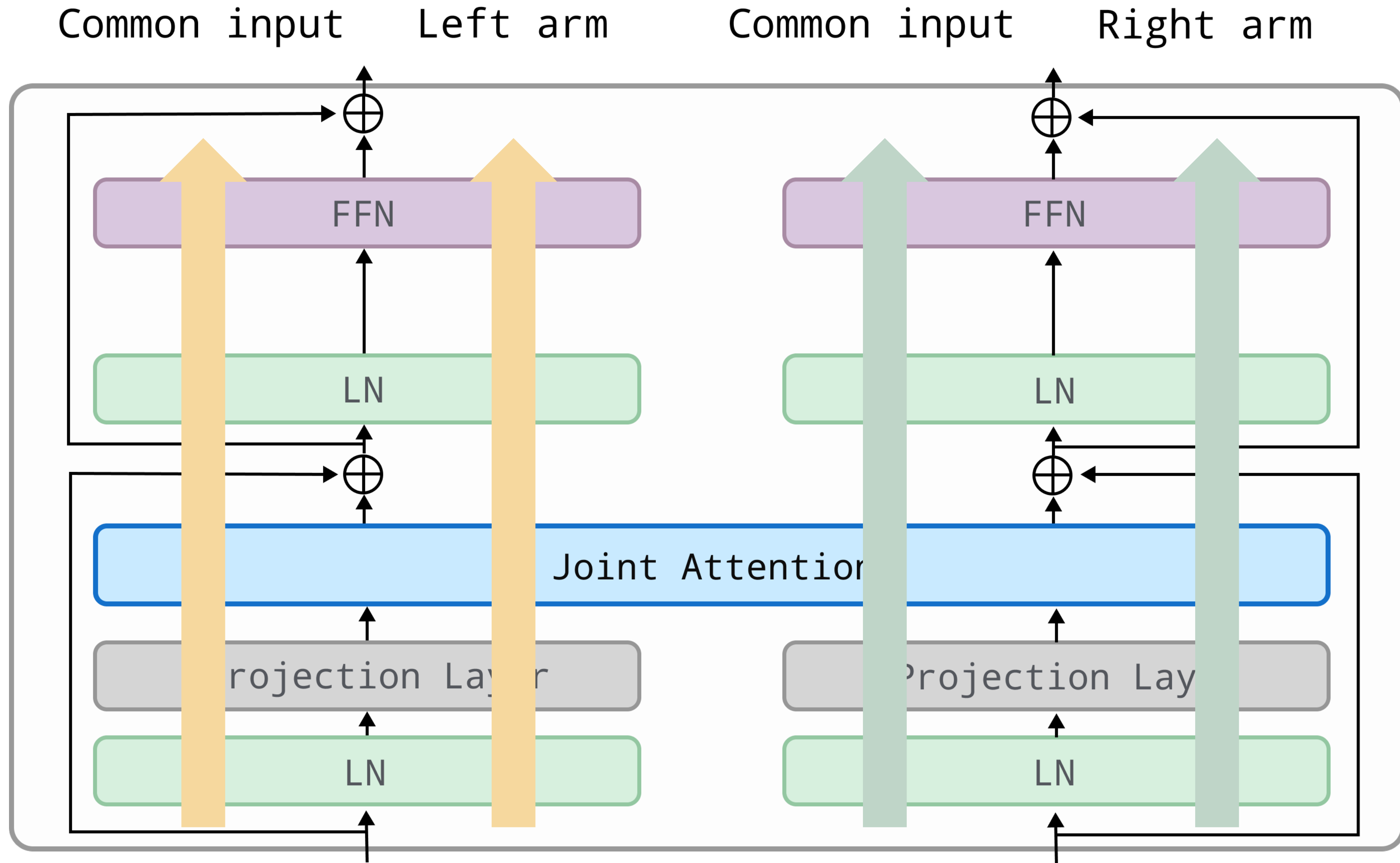
MoE for compute efficiency



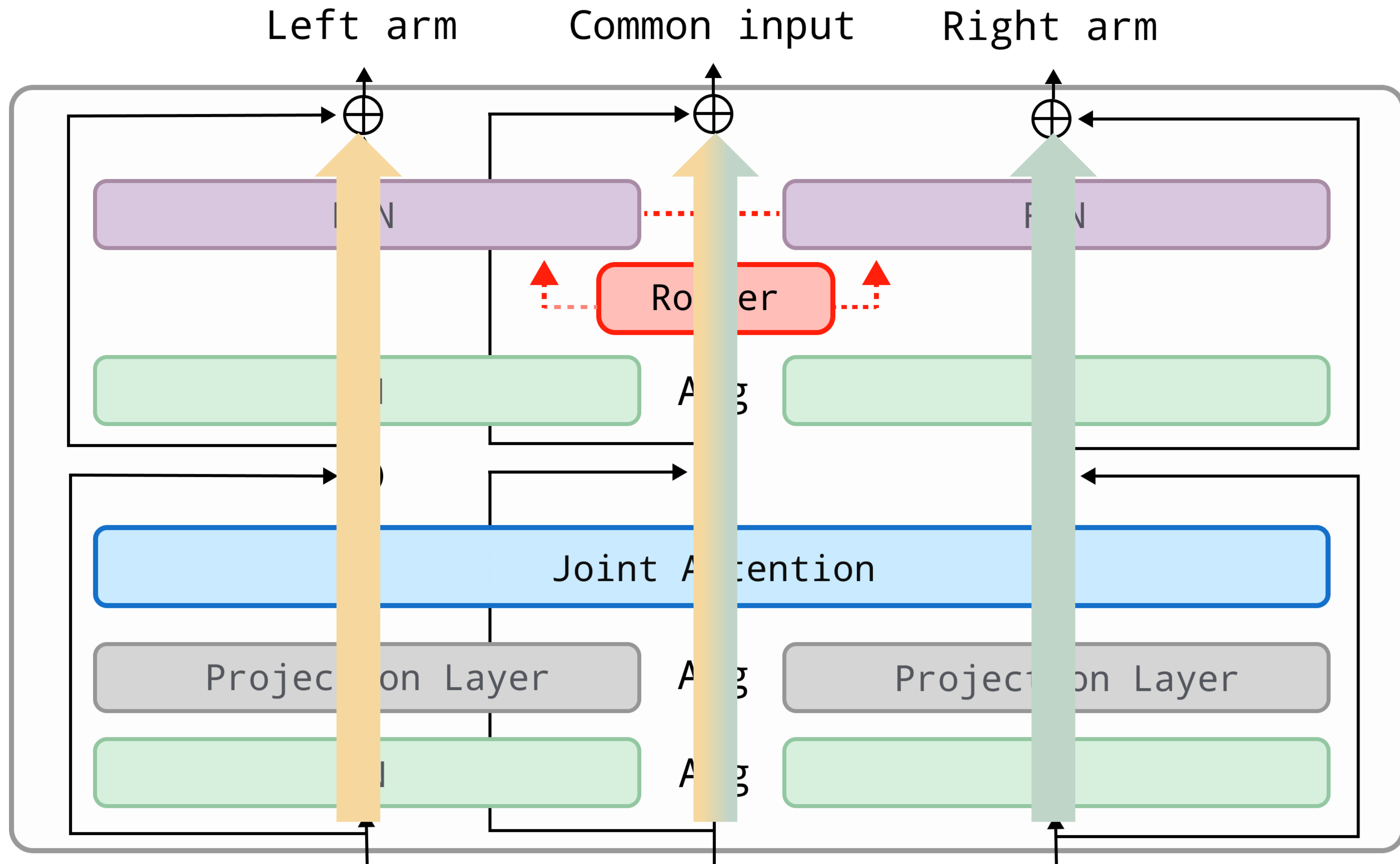
Increased input length

Increased VRAM usage

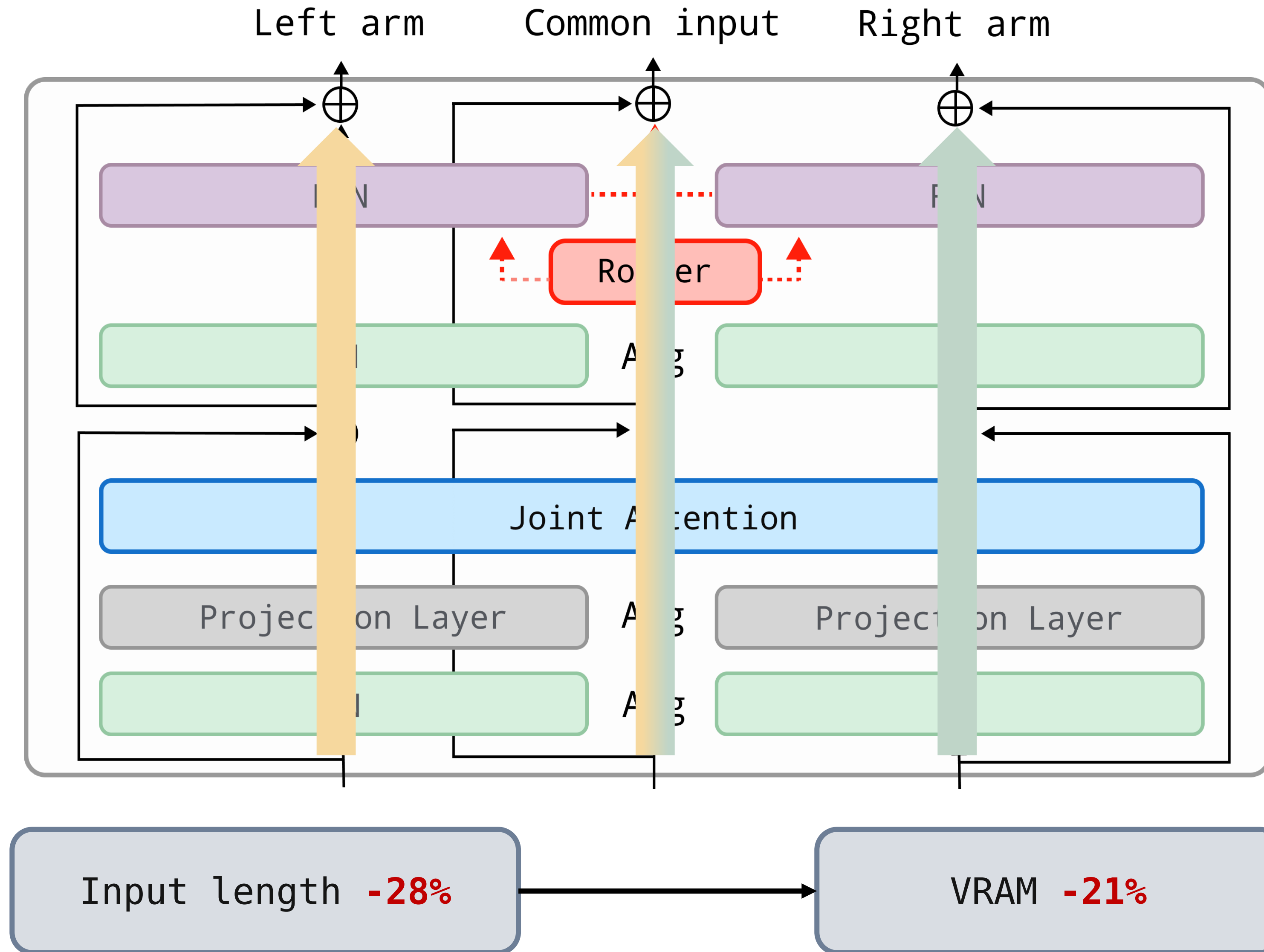
MoE for compute efficiency



MoE for compute efficiency



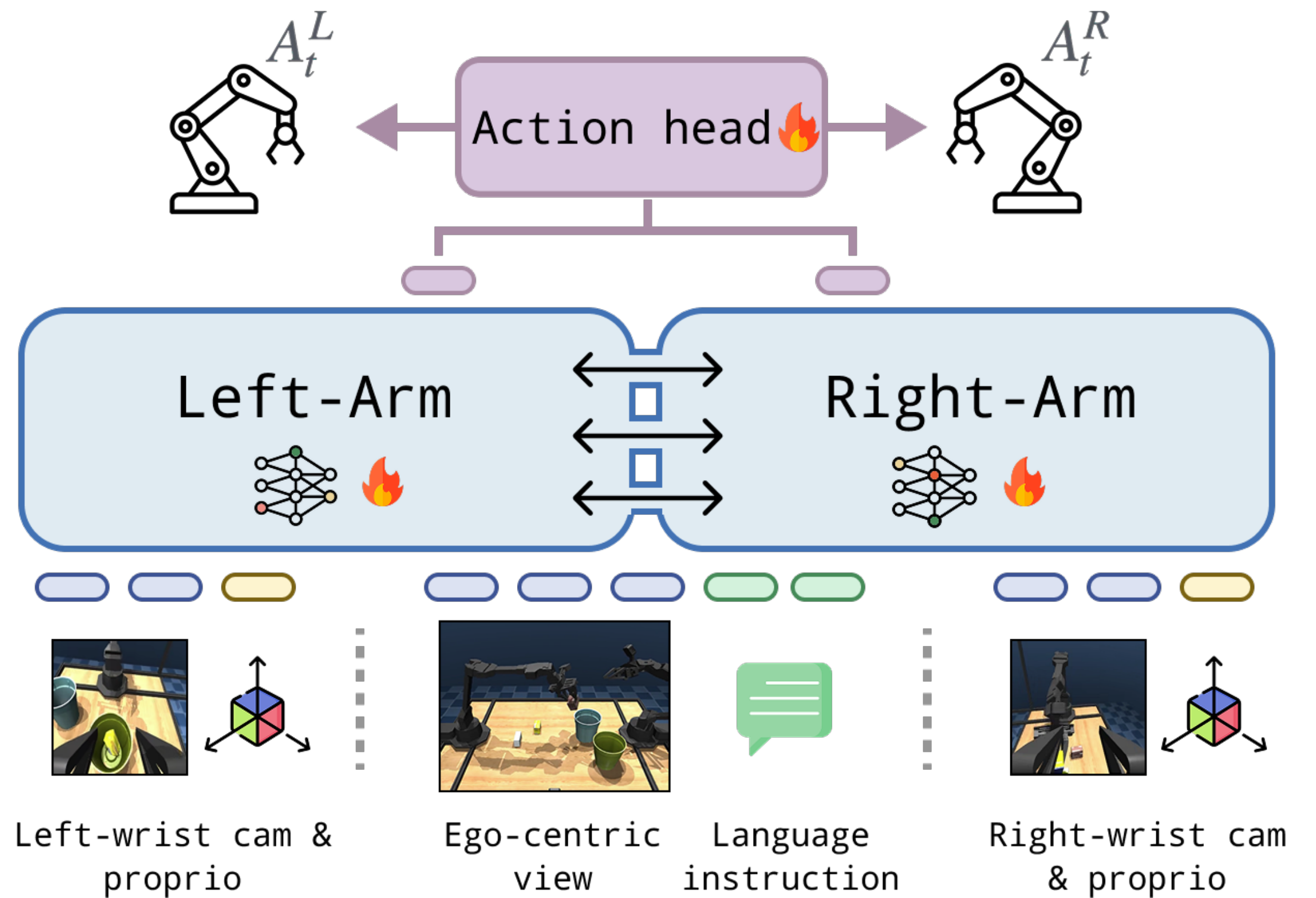
MoE for compute efficiency



MoE for compute efficiency

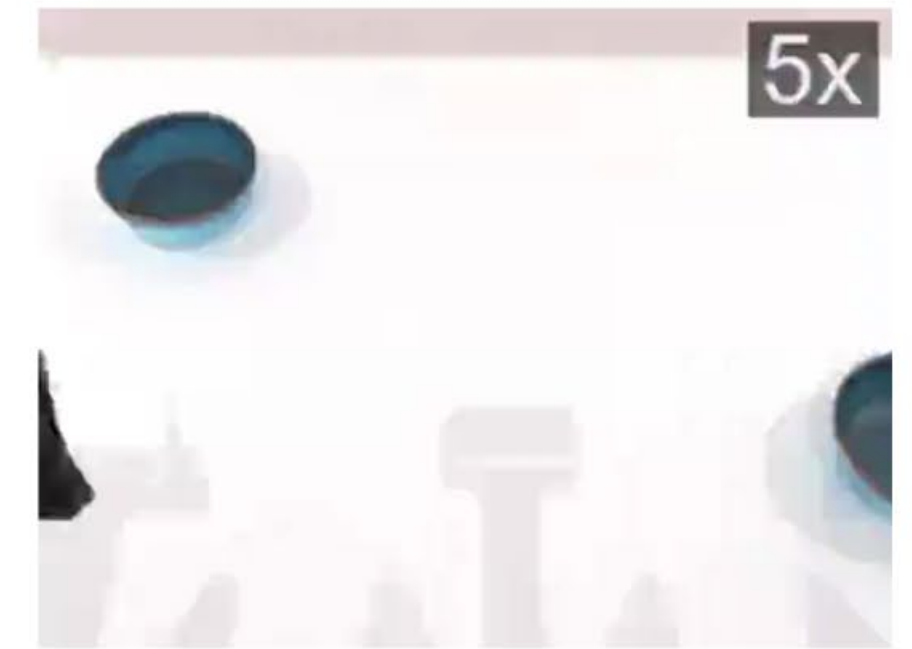
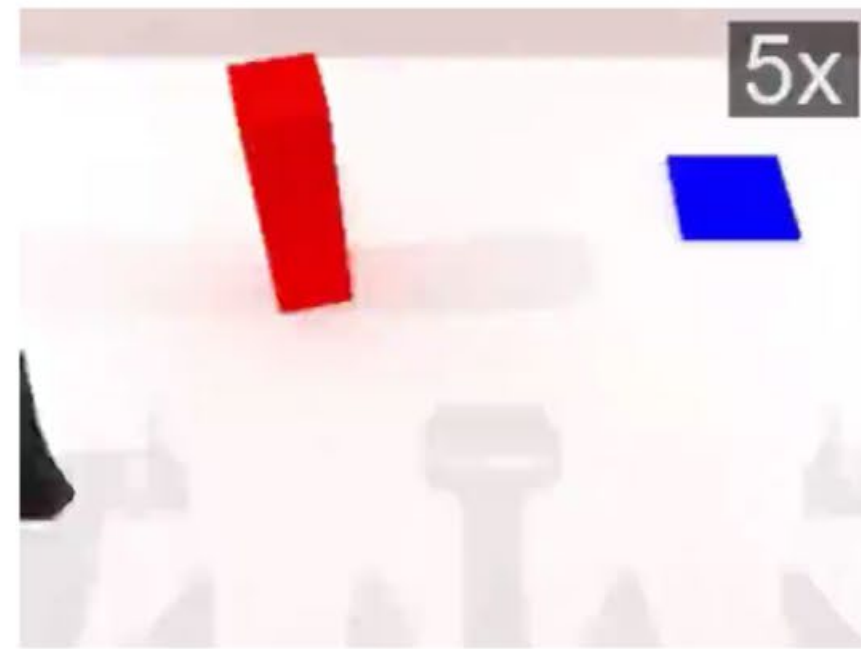
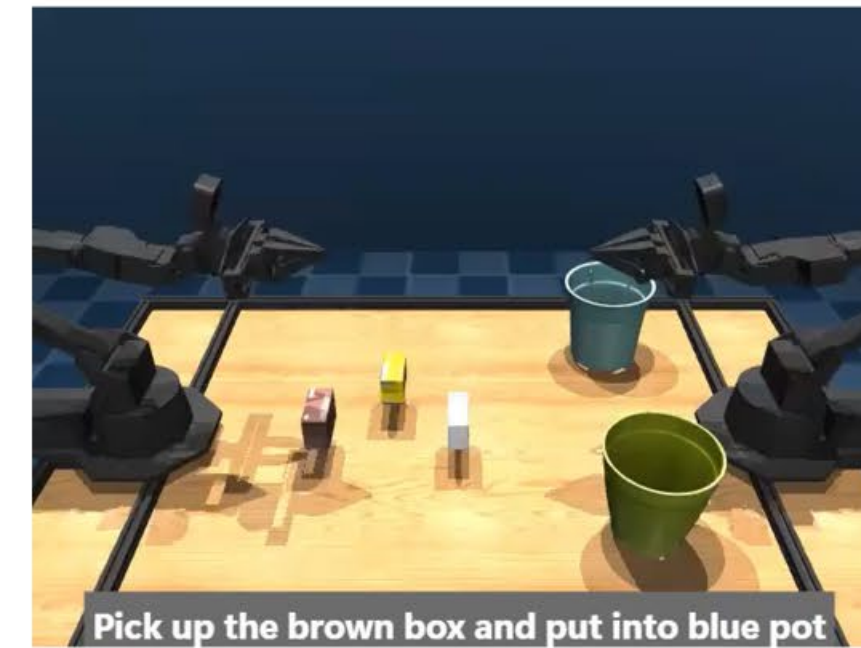
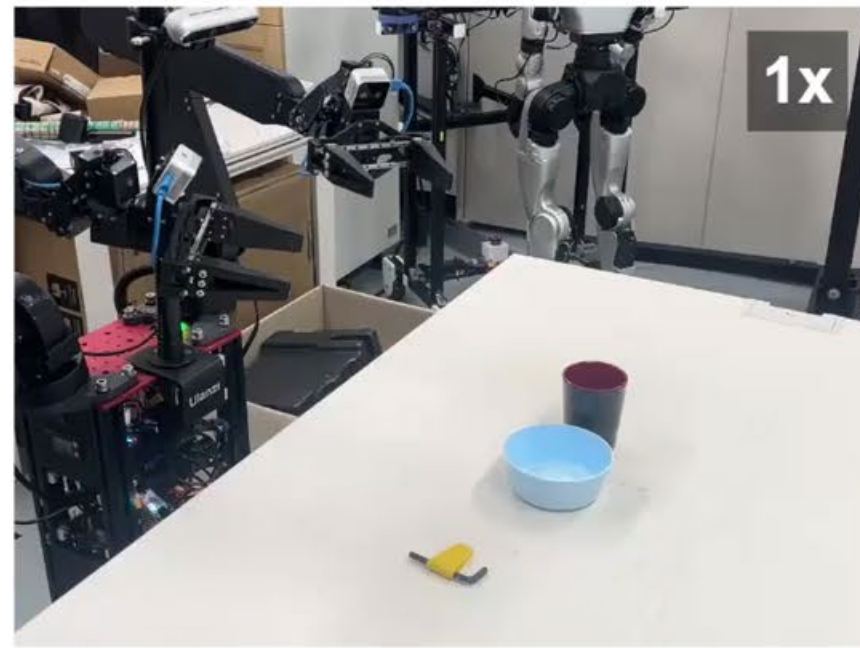
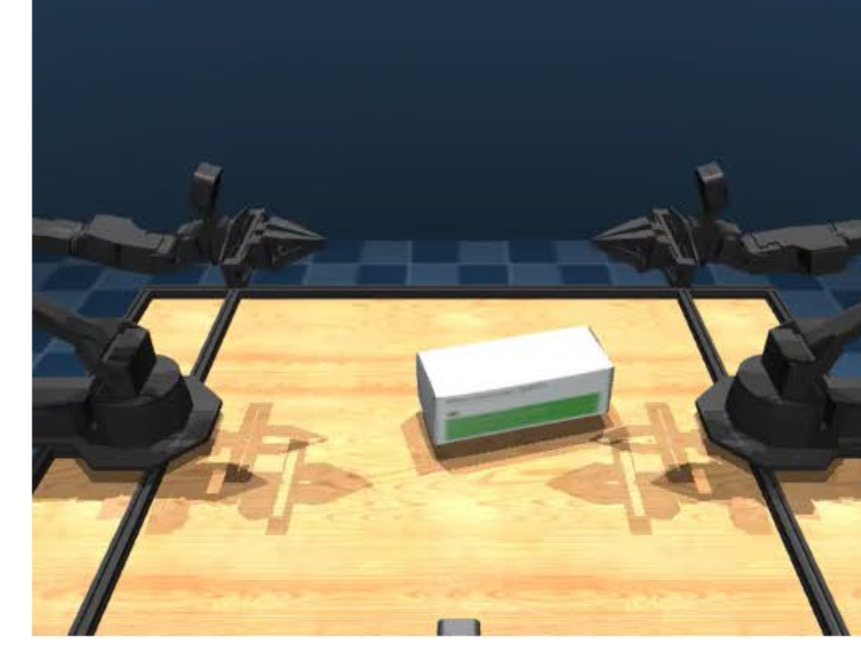
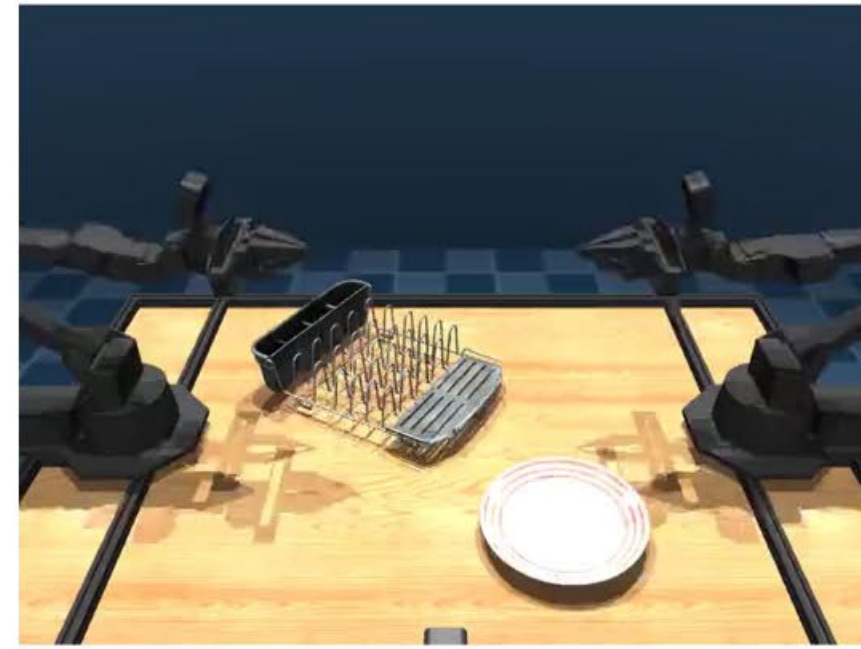
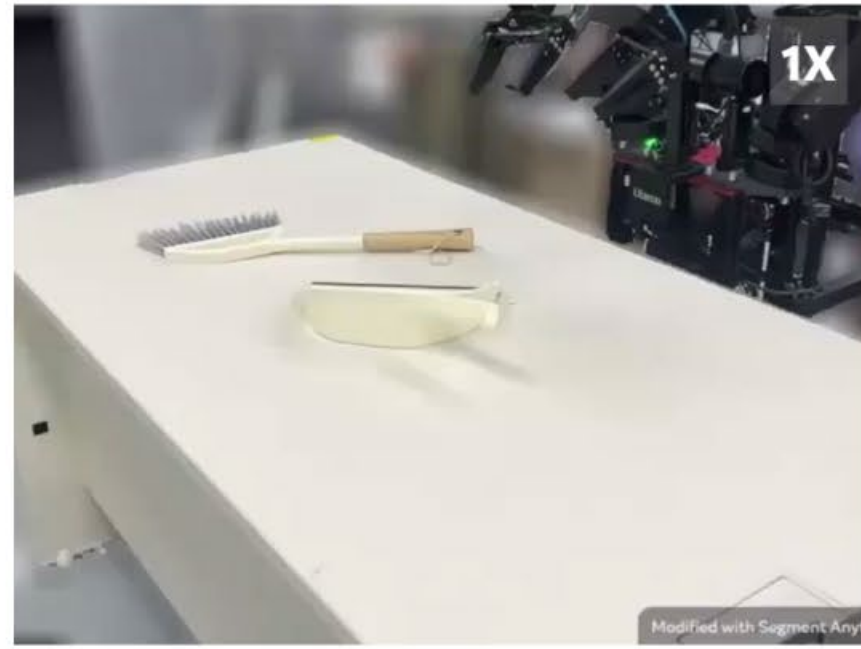
TwinVLA-1.3B

Joint attention & MoE integration



$$o_t^{\text{twin}} = \left(I_{\text{wrist}}^L, d^L \right)_t \quad \left(I_{\text{ego}}, l \right)_t \quad \left(I_{\text{wrist}}^R, d^R \right)_t$$

Experiments



Real World (Anubis Bimanual Robot)

Simulation (Tabletop-Sim & RoboTwin2.0)

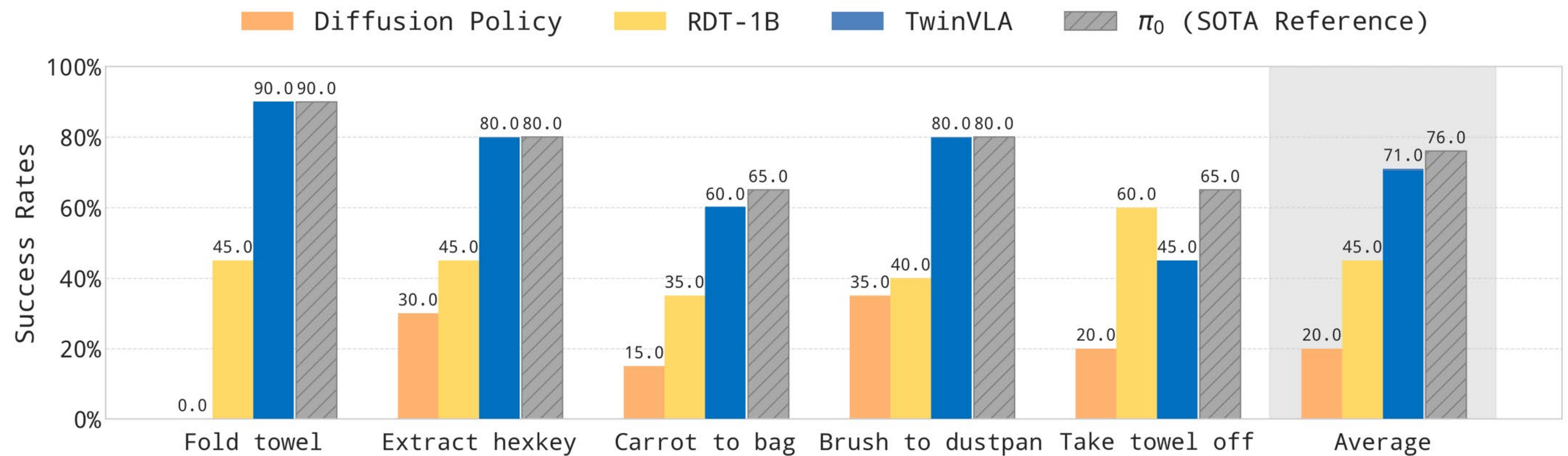
Experiments – Real World



Long-Horizon

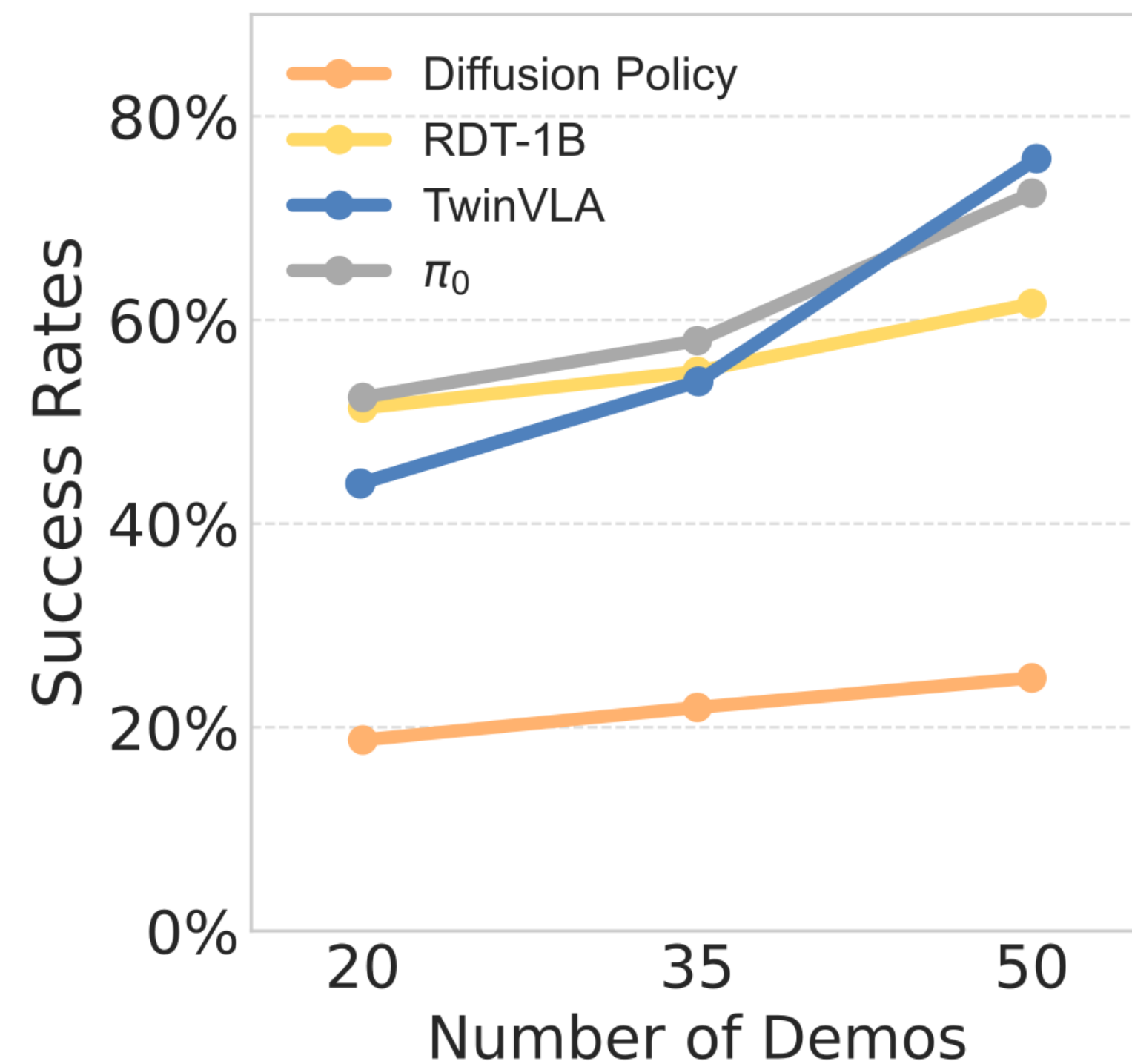


Arm Coordination

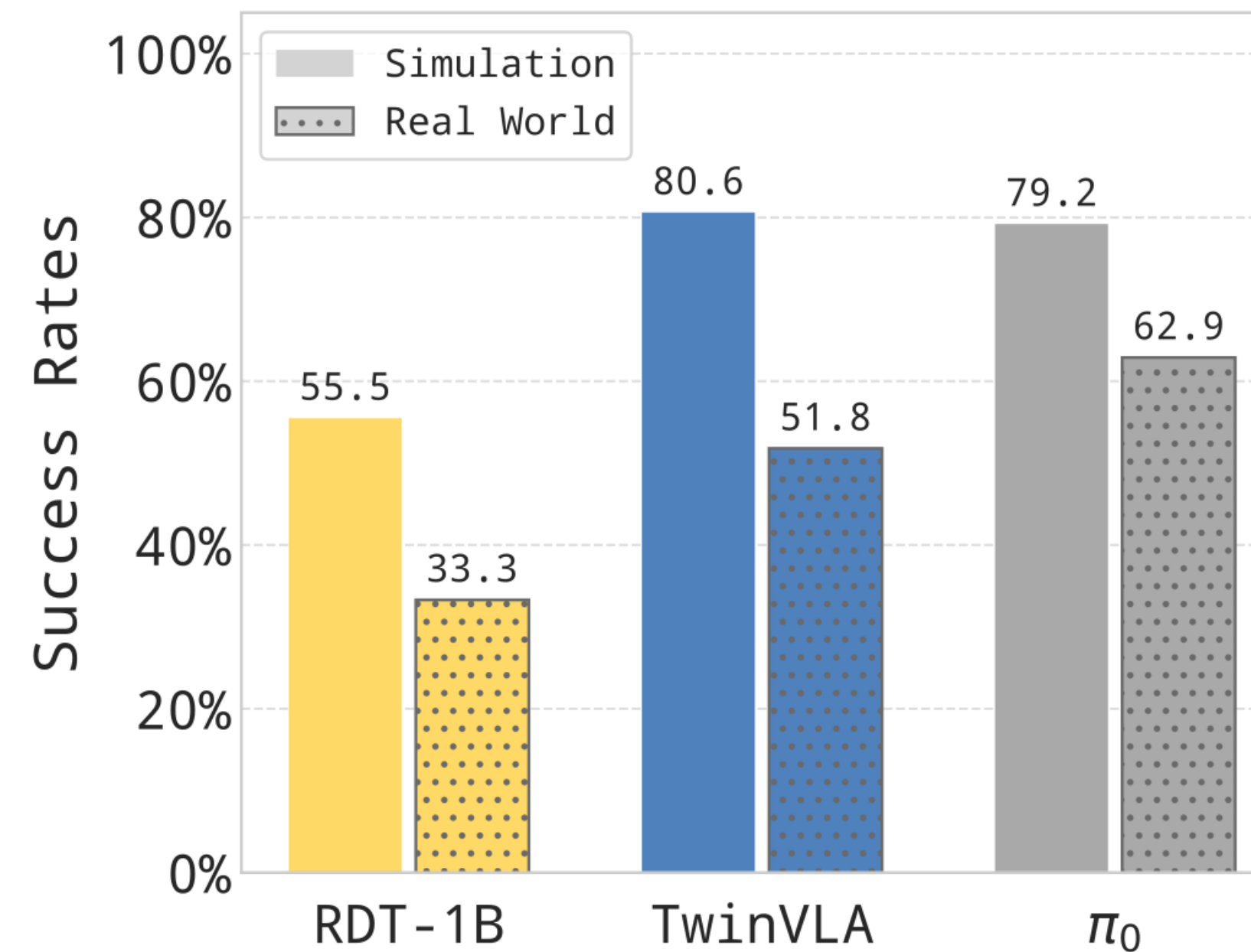


TwinVLA outperforms RDT, while comparable to π_0

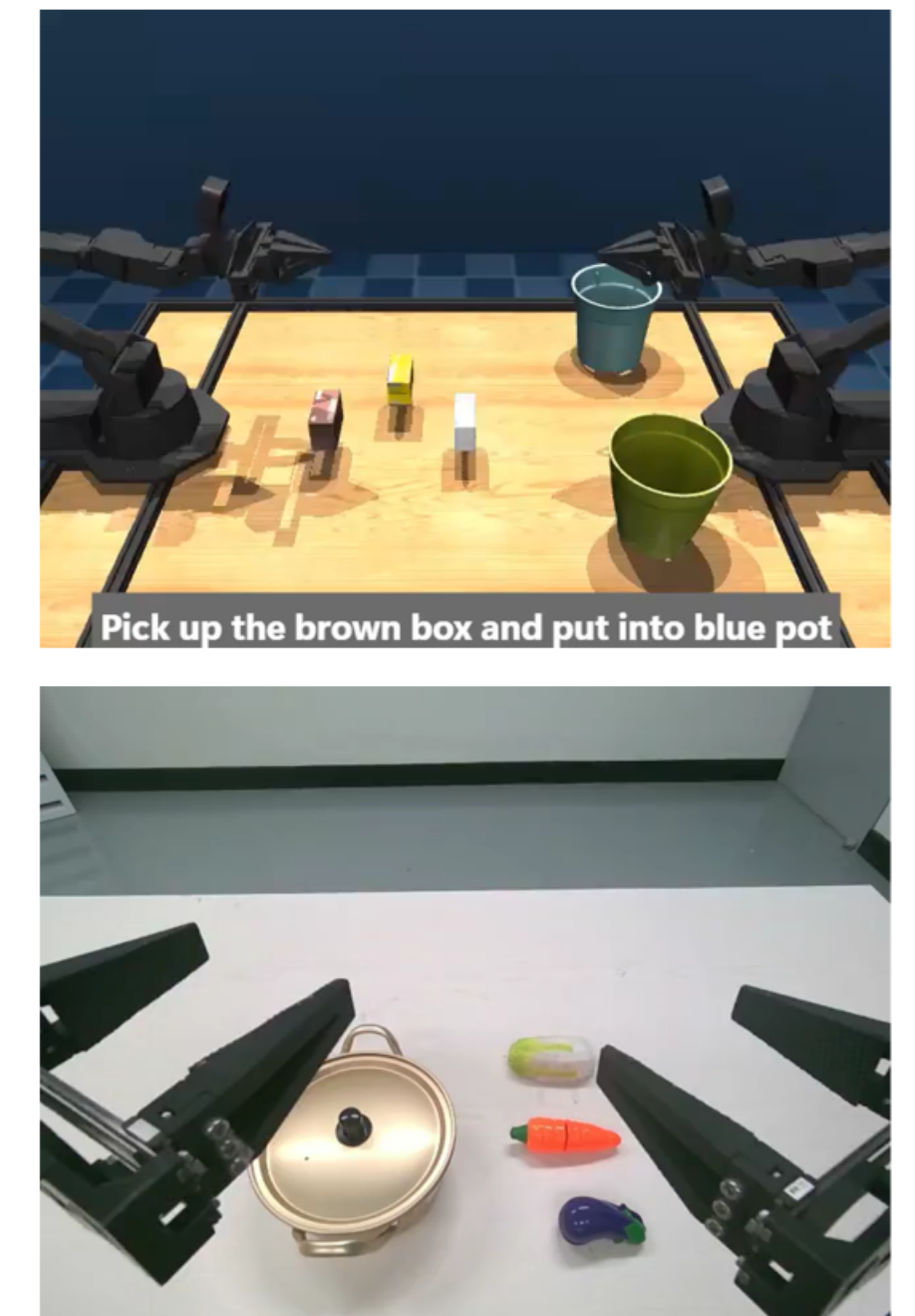
Experiments – Data efficiency & Language following



TwinVLA exhibits exceptional
data efficiency



Robust performance on
language-following tasks



TwinVLA : Data-Efficient Bimanual VLA



Summary

1. Zero Bimanual Data for Pretraining:

We can make a bimanual VLA only from single-arm data

2. Modular Architecture:

Single-arm VLAs + Joint Attention + MoE

3. Outperforms Monolithic Baselines