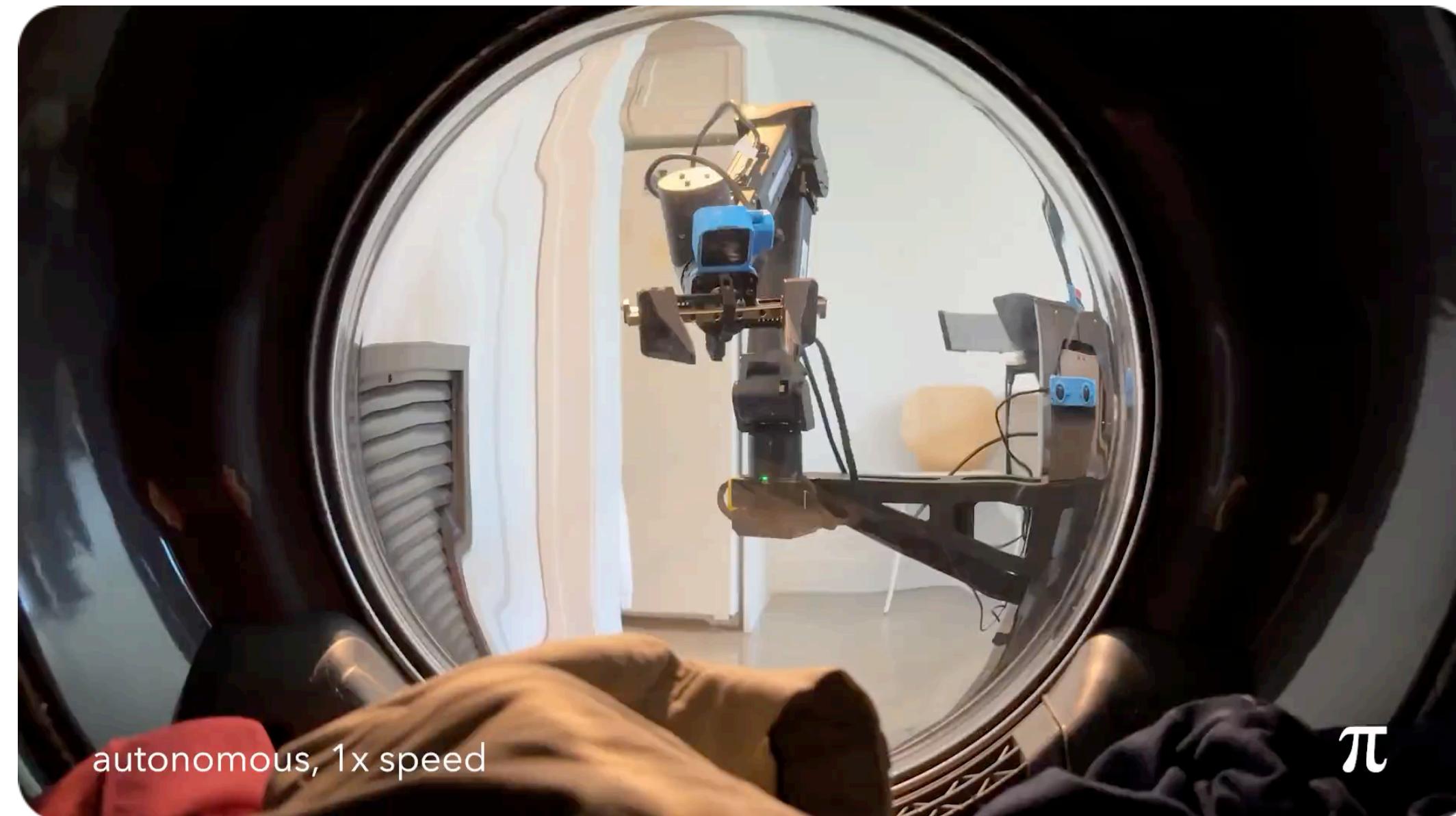




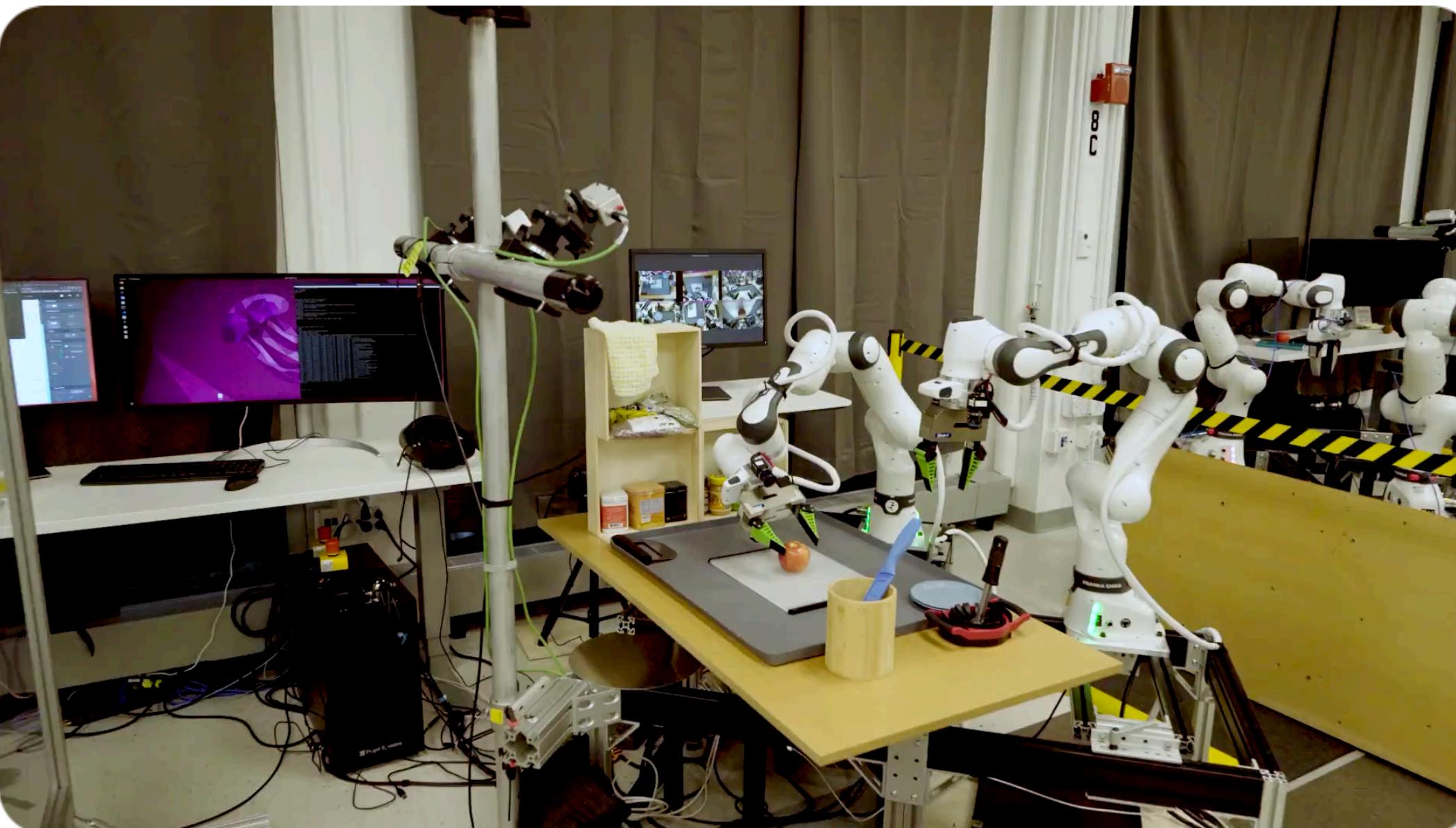
Perspectives on Designing Vision Language Action Models

Ankit Goyal

The Evolving State of Robotics



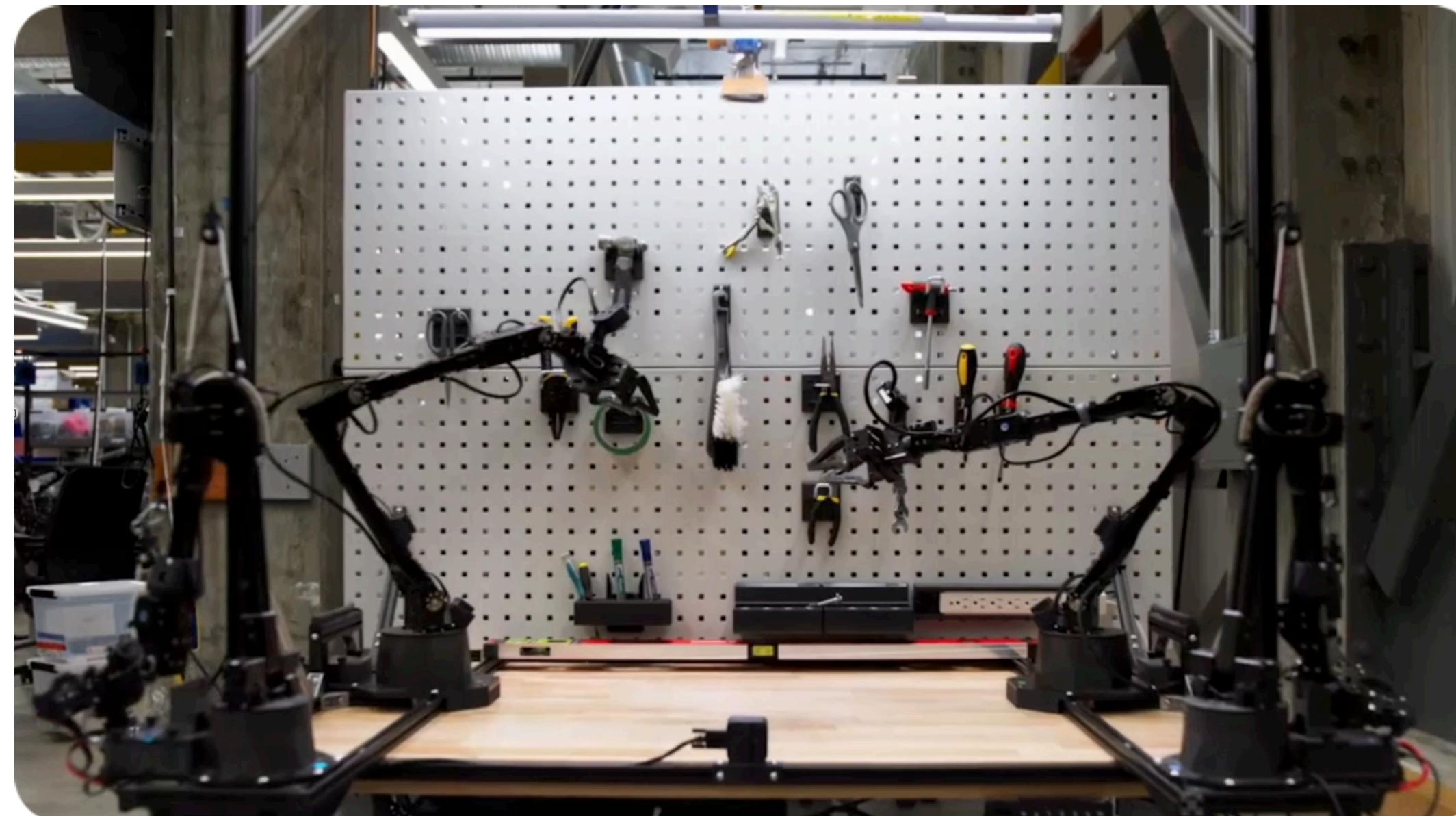
Physical Intelligence



Toyota Research Institute



Dyna Robotics



Google Robotics



Figure Robotics



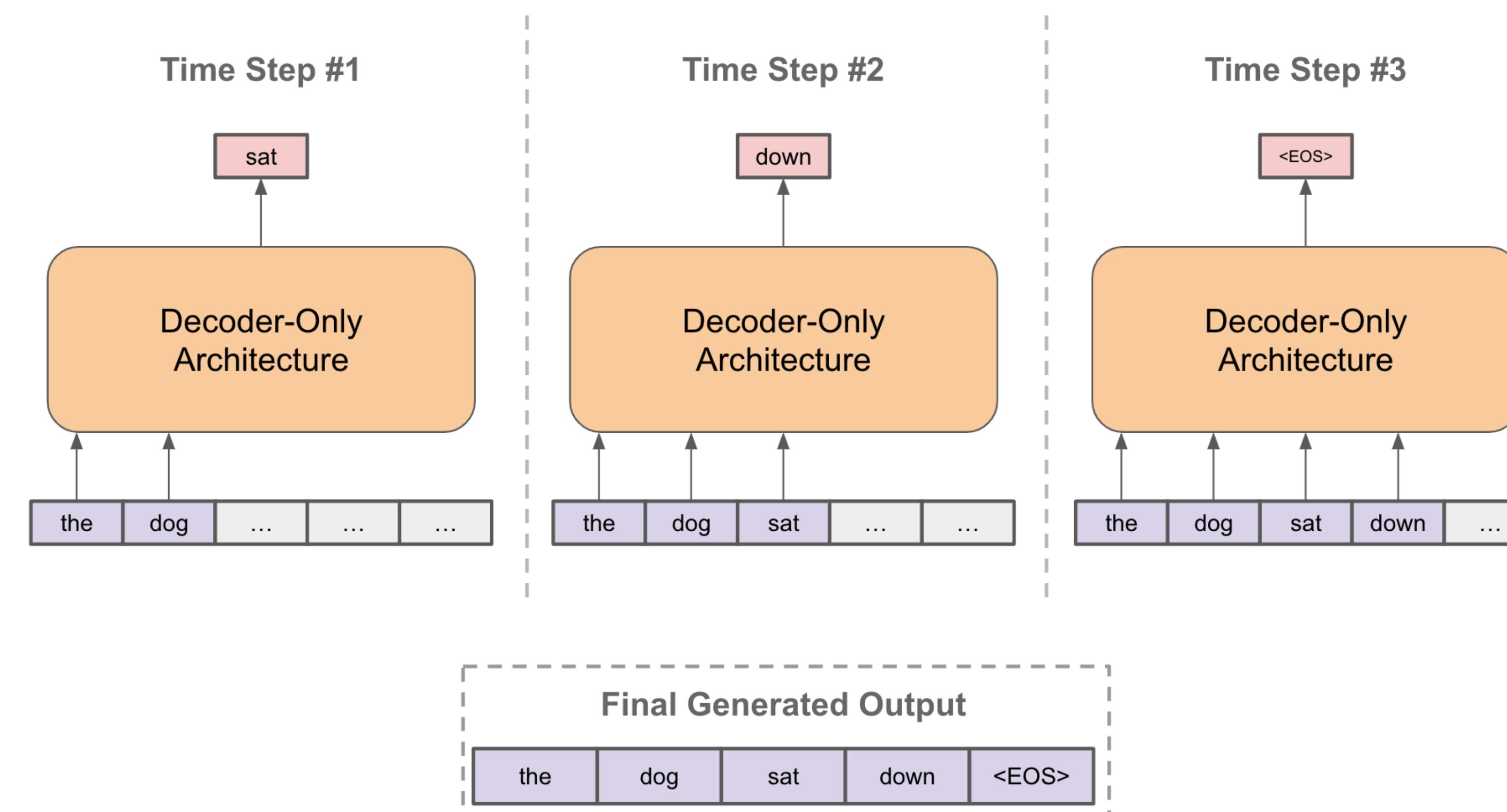
Nvidia GR00T

[Courtesy: Physical Intelligence, TRI, Dyna, Google Robotics, Figure Robotics, NVIDIA GR00T]

Background

LLMs

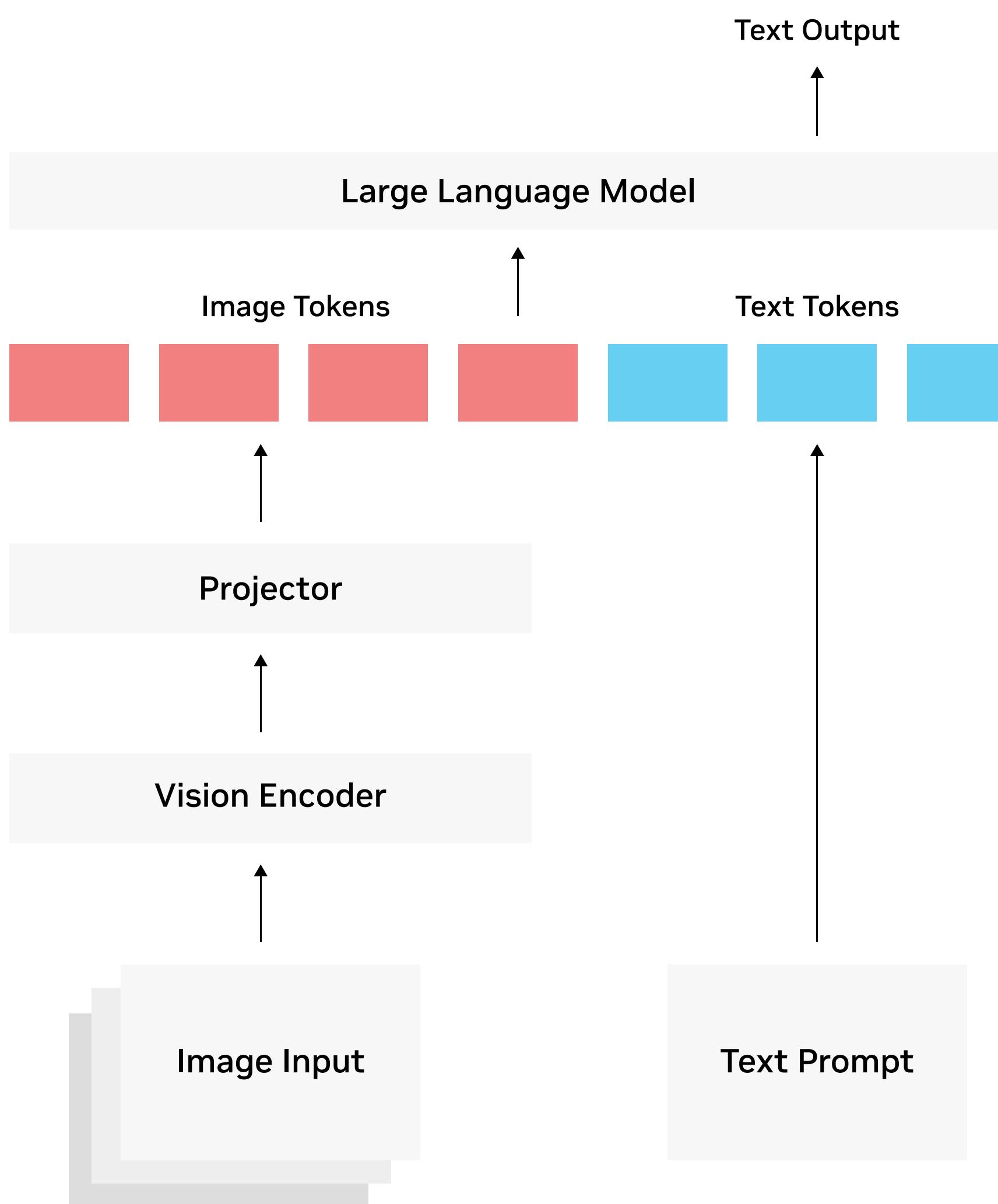
- LLMs (Large Language Models) — They predict one token (similar to sub-word) at a time
- Trained with large scale data and many other tricks!



Background

VLMs

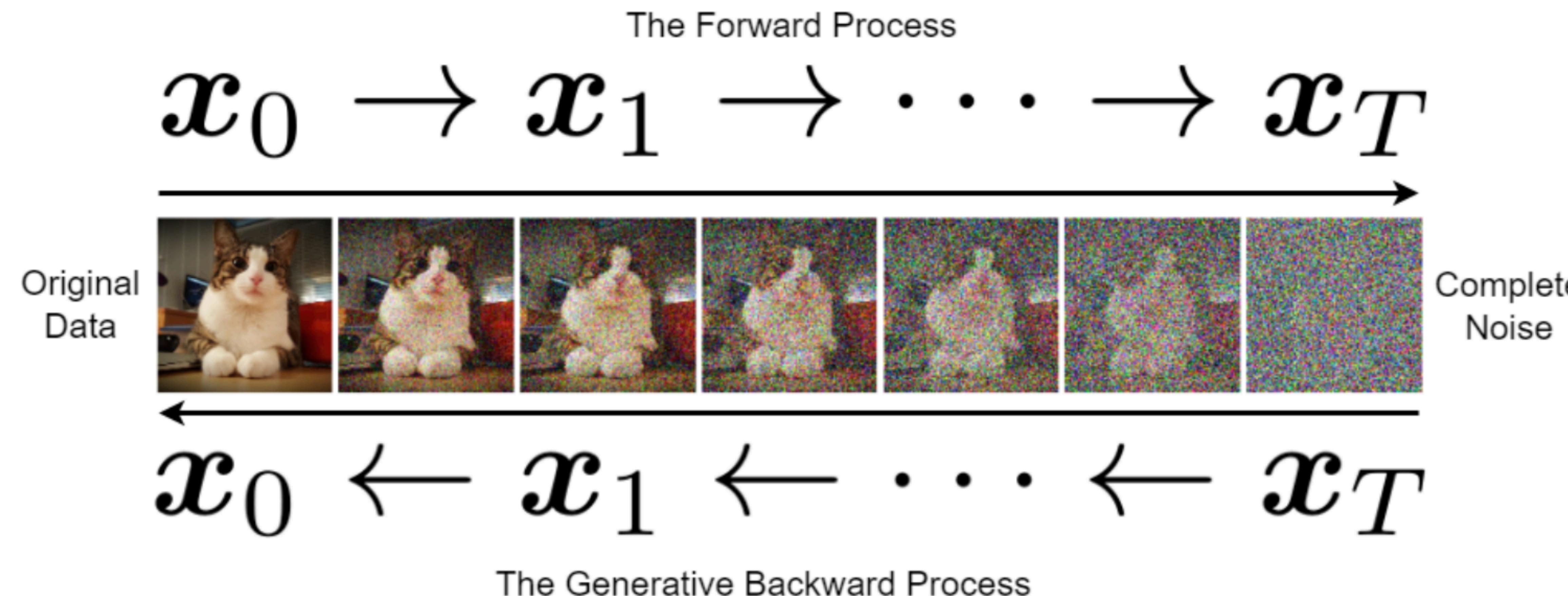
- We start with an LLM and then introduce image tokens to it
- Image tokens come from a pre-trained Vision Encoder
- Fine-tune the LLM on this joint task to create a VLM



Background

Diffusion Model

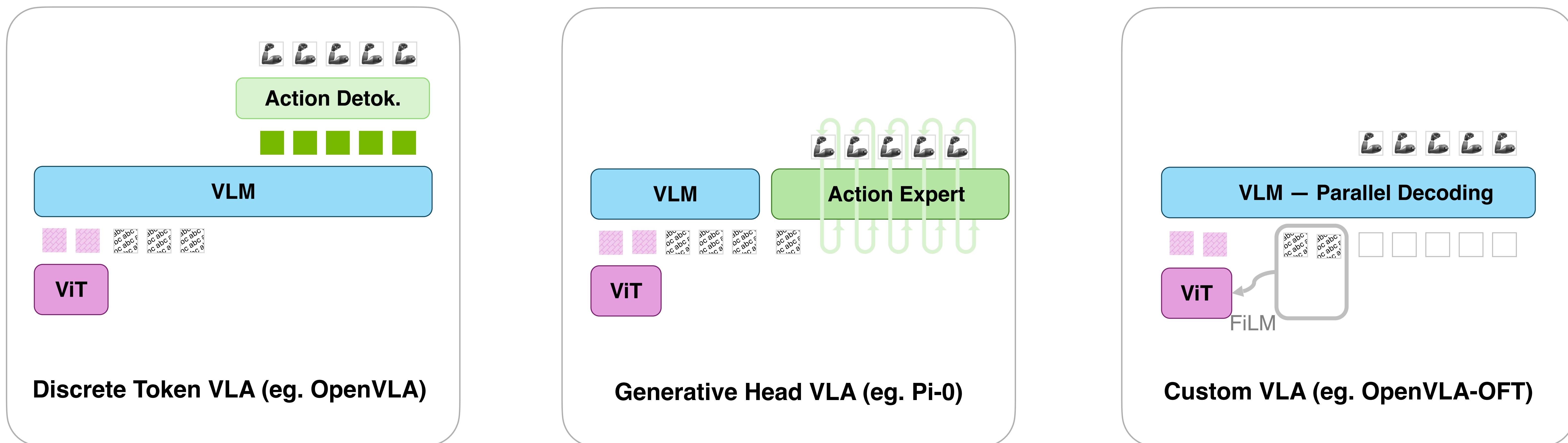
- A class of generative model — i.e. they can be used to generate
- They are trained via adding progressive adding noise to a sample and then asking a network to denoise it
- During generation, we start with noise and progressively denoise it



What are Vision Language Action Models ?

Family of VLAs

- LLMs power Language Understanding, VLMs bridge Vision and Language – VLAs extend this to Actions
- VLAs build on pretrained VLMs – Adapting them to reason about Actions
- VLAs in the existing literature can be classified into three categories



Text Tokens



Image Tokens



Empty / Query Tokens



Robot Actions

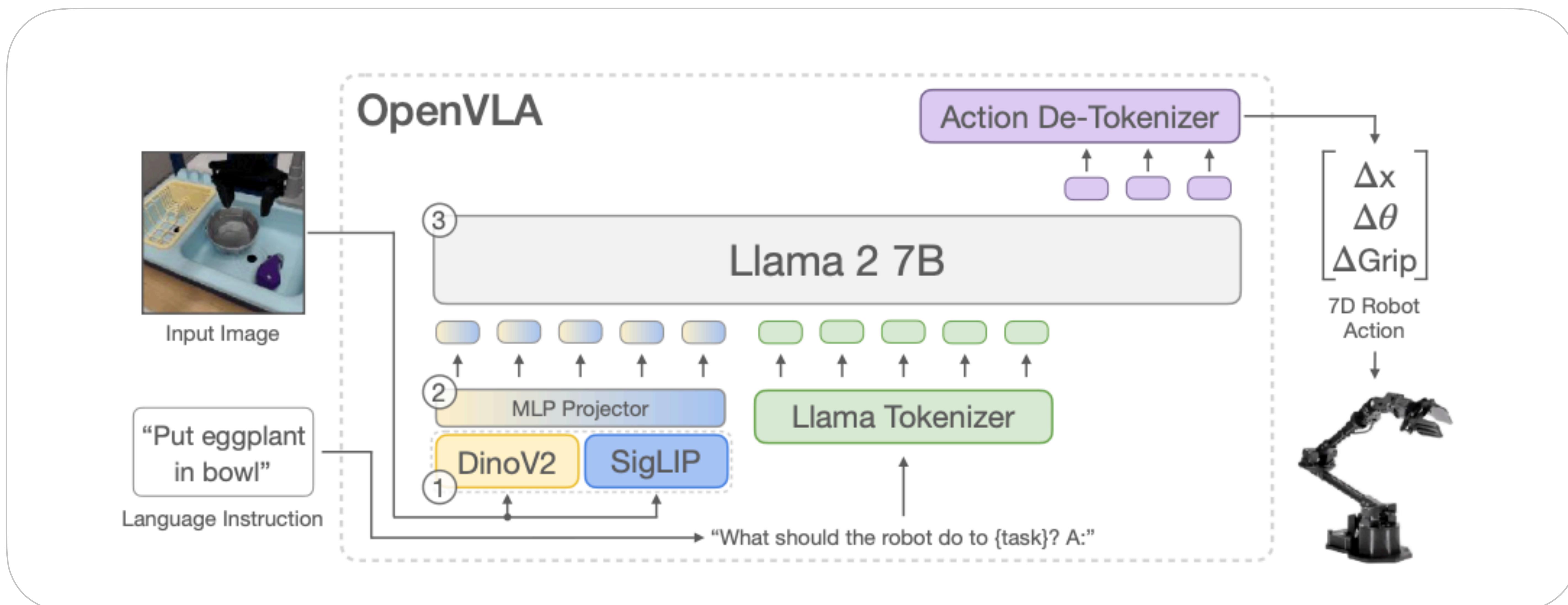


Action Tokens

I. Discrete Token VLAs

E.g. OpenVLA

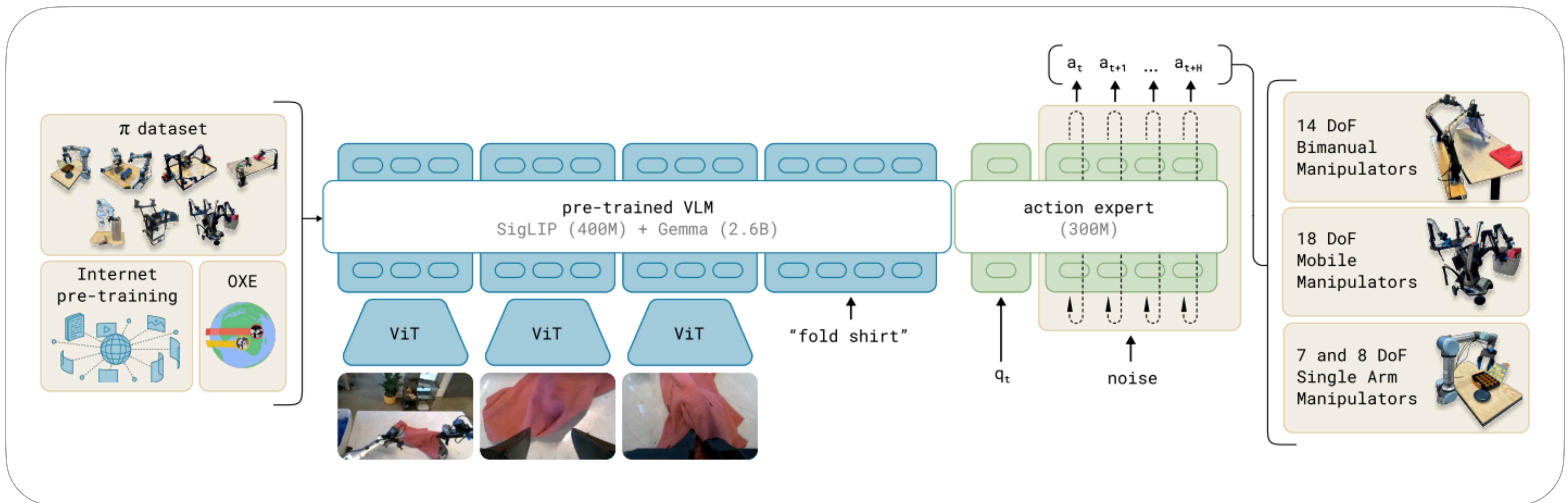
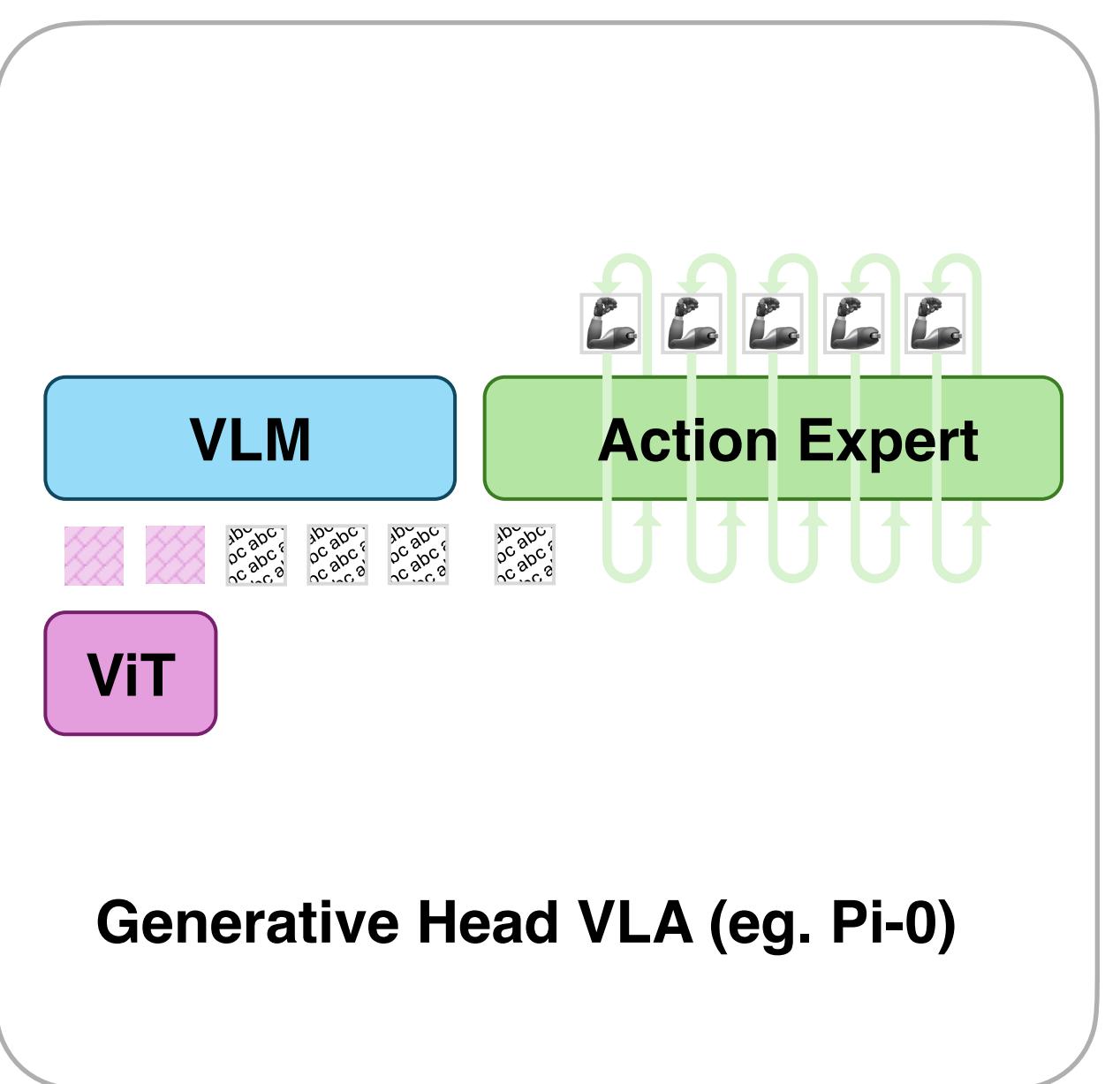
- Replace or modify the vocabulary to introduce Discrete Action Tokens
- Limited Action Resolution
- Compromised pretrained language representations



II. Generative Head VLAs

Eg. Pi-0

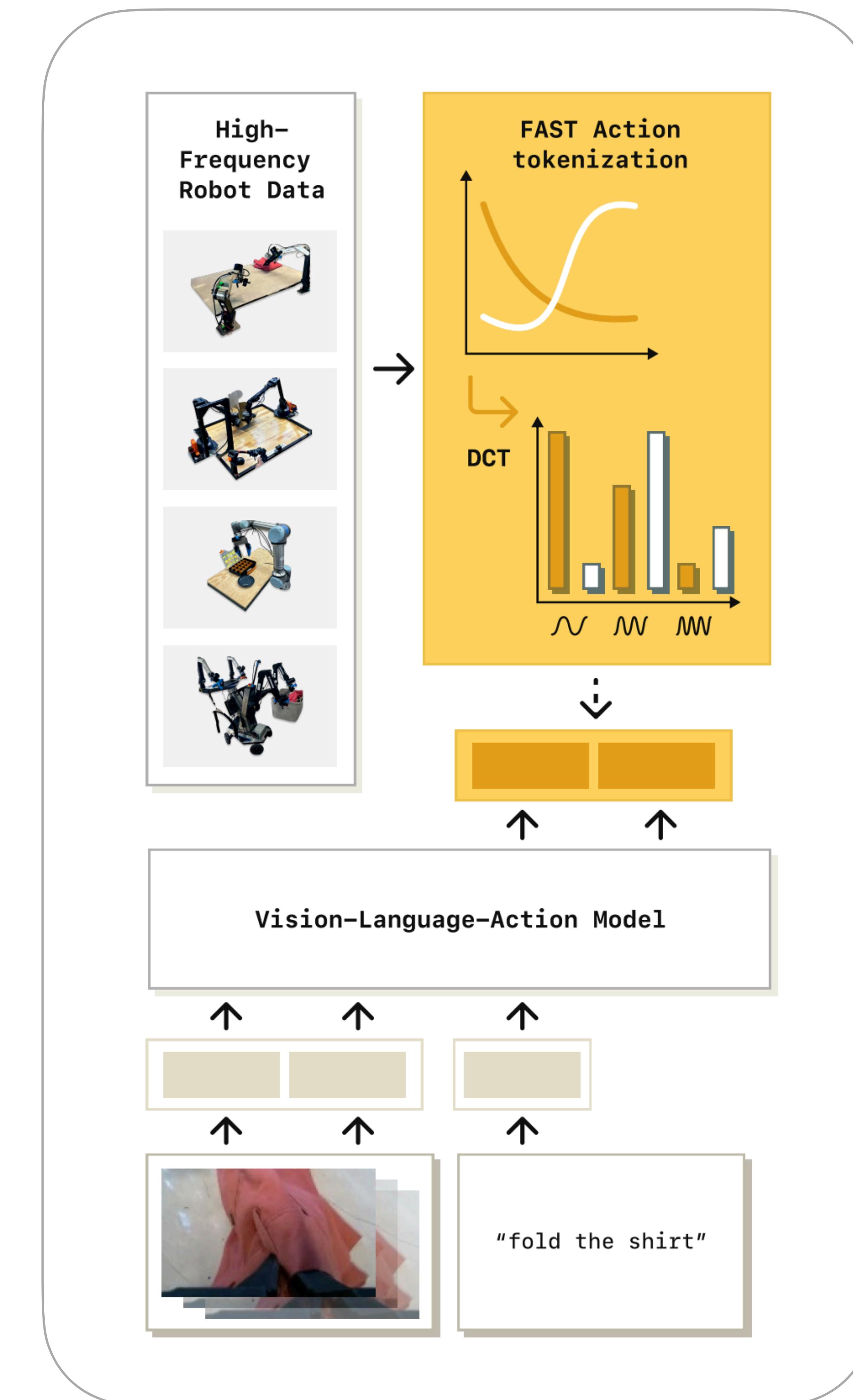
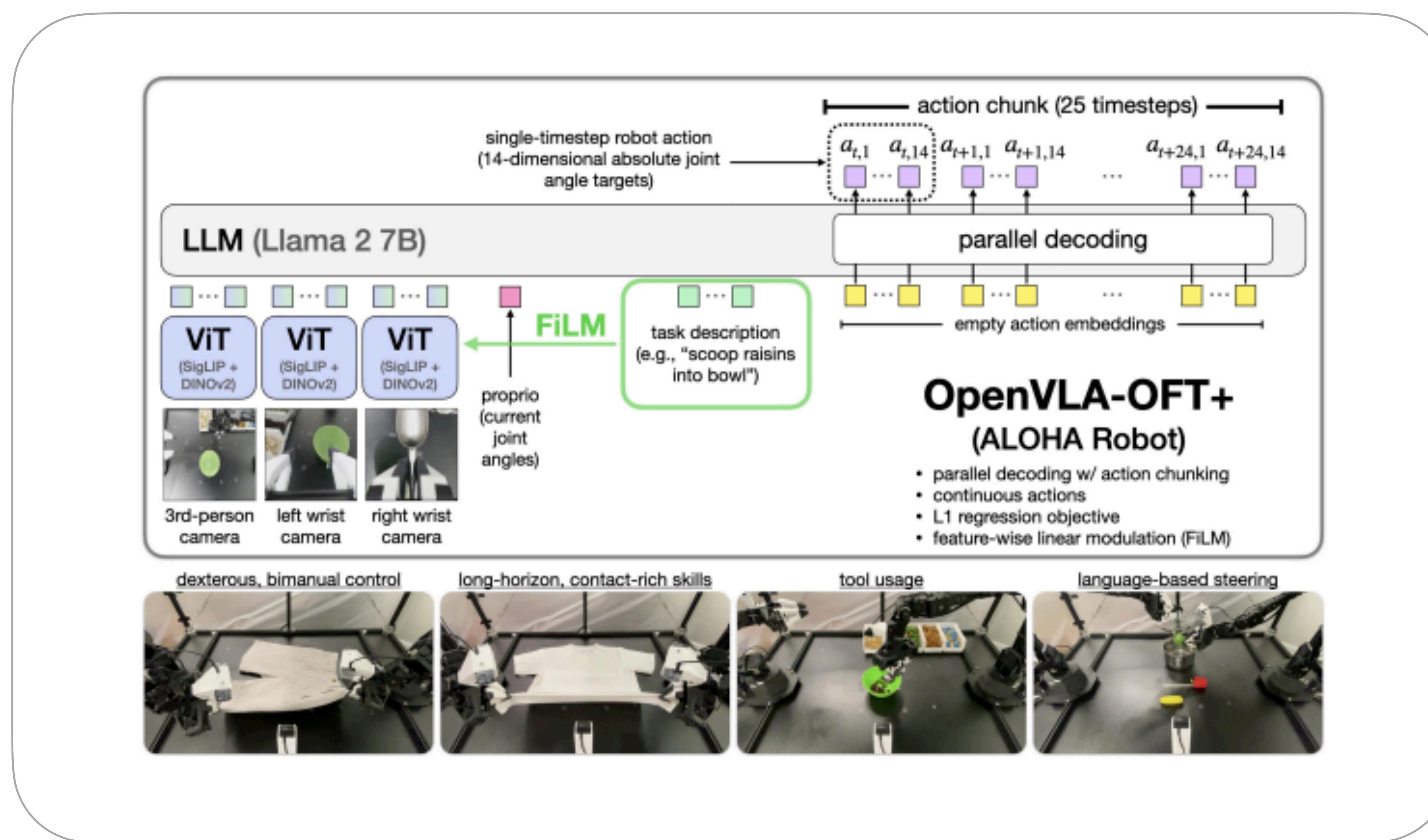
- VLM predicts latent vector, fed into a Generative Head (Diffusion/Flow Match.) to predict actions
- Non-pretrained components reduce model generalization
- Compromised pretrained language representations



III. Custom VLA Designs

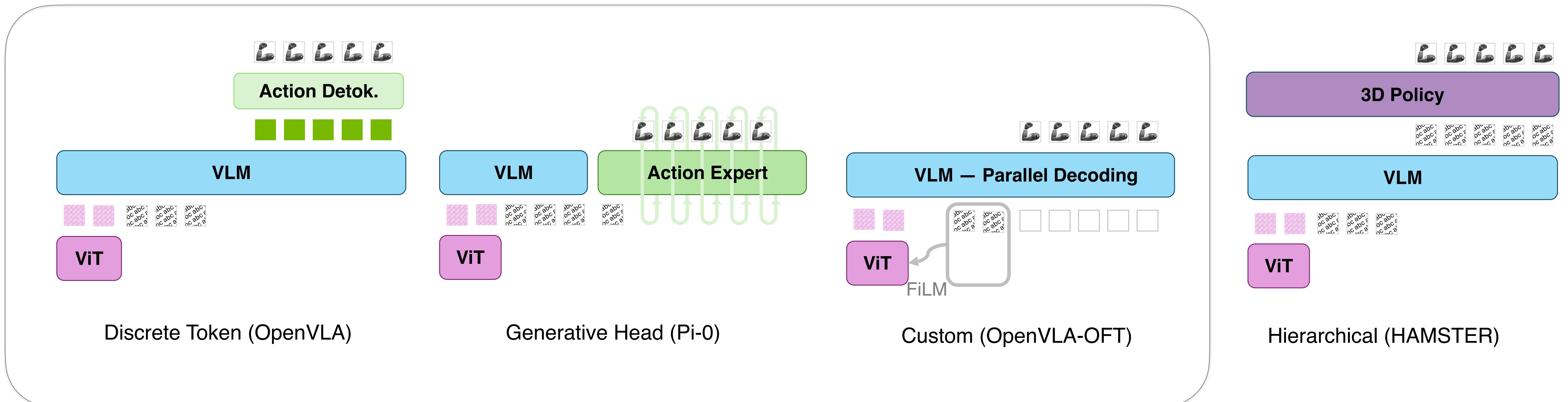
Eg. OpenVLA-OFT

- Not in the above two general categories
- Significant architectural changes
- Specialized Training Pipelines



Family of VLAs

- Hierarchical design instead of monolithic design
- Leverage VLMs for generalization and specialized policies (3D) for action prediction



HAMSTER: Hierarchical Action Models For Open-World Robot Manipulation

ICLR 2025



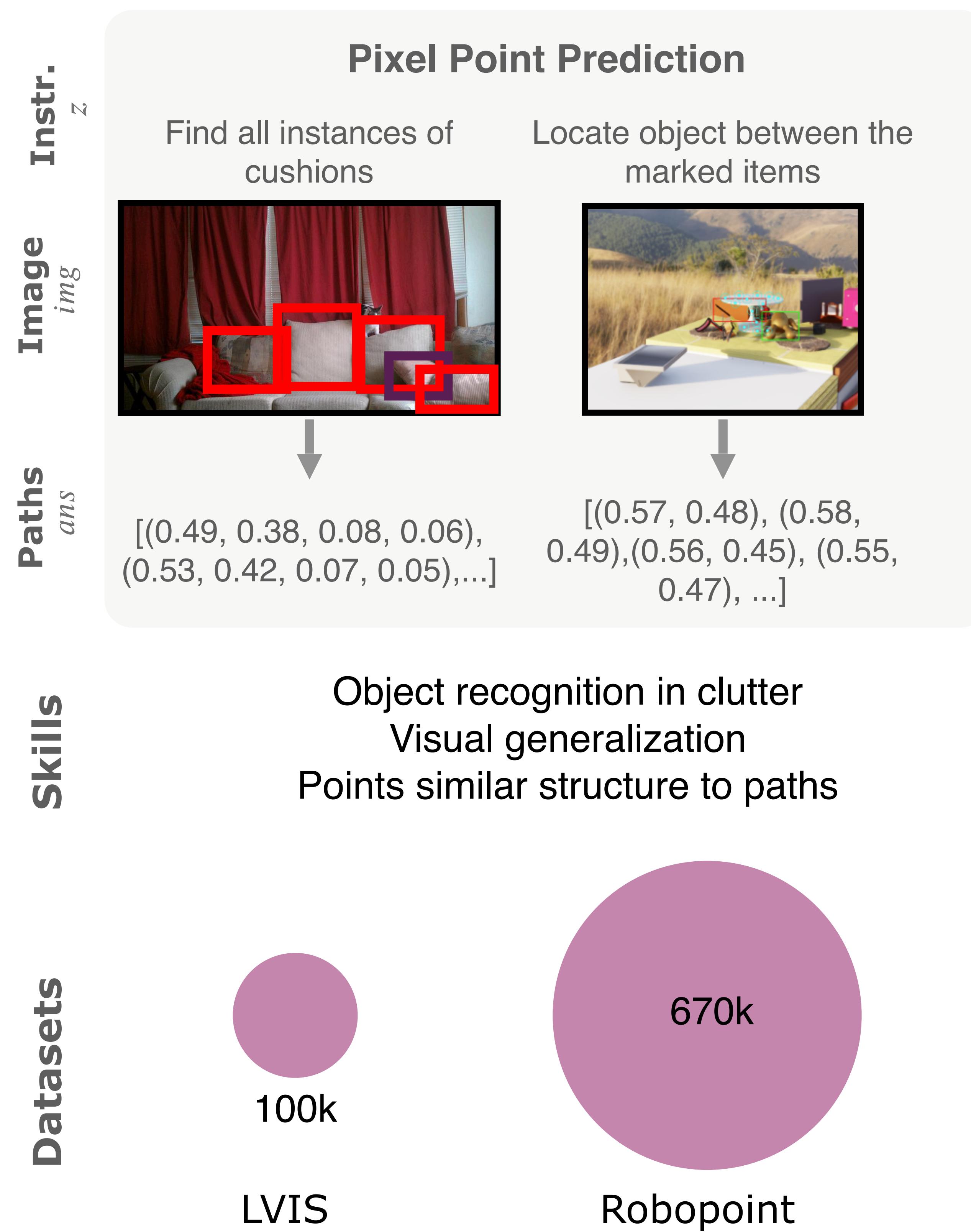
Put M on Mona Lisa



The HAMSTER VLM predicts points to denote the 2D paths

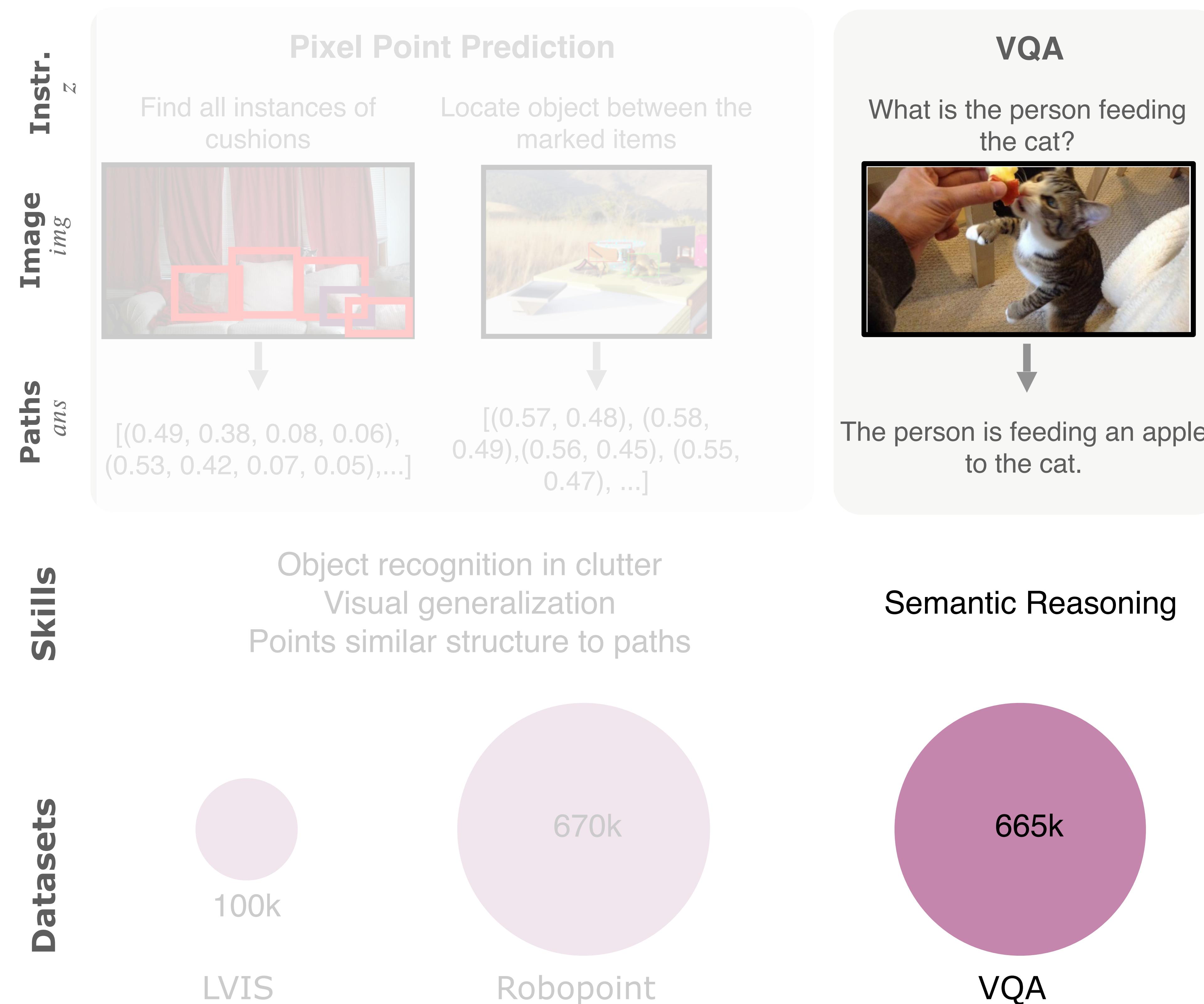
Hamster VLA Training Data

High-Level model is trained on a variety of data sources (Sim, Real, Point-QA)



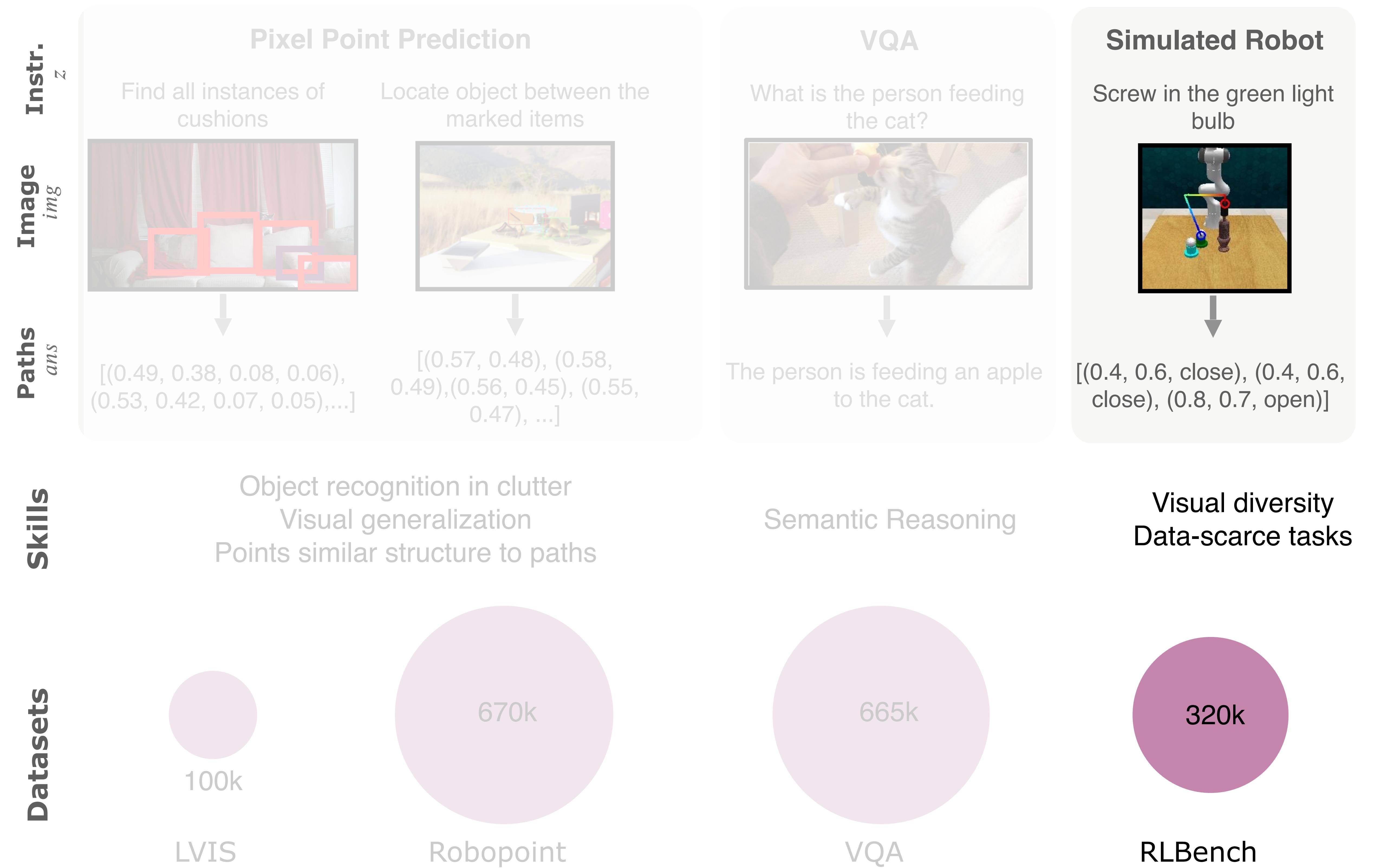
Hamster VLA Training Data

High-Level model is trained on a variety of data sources (Sim, Real, Point-QA)



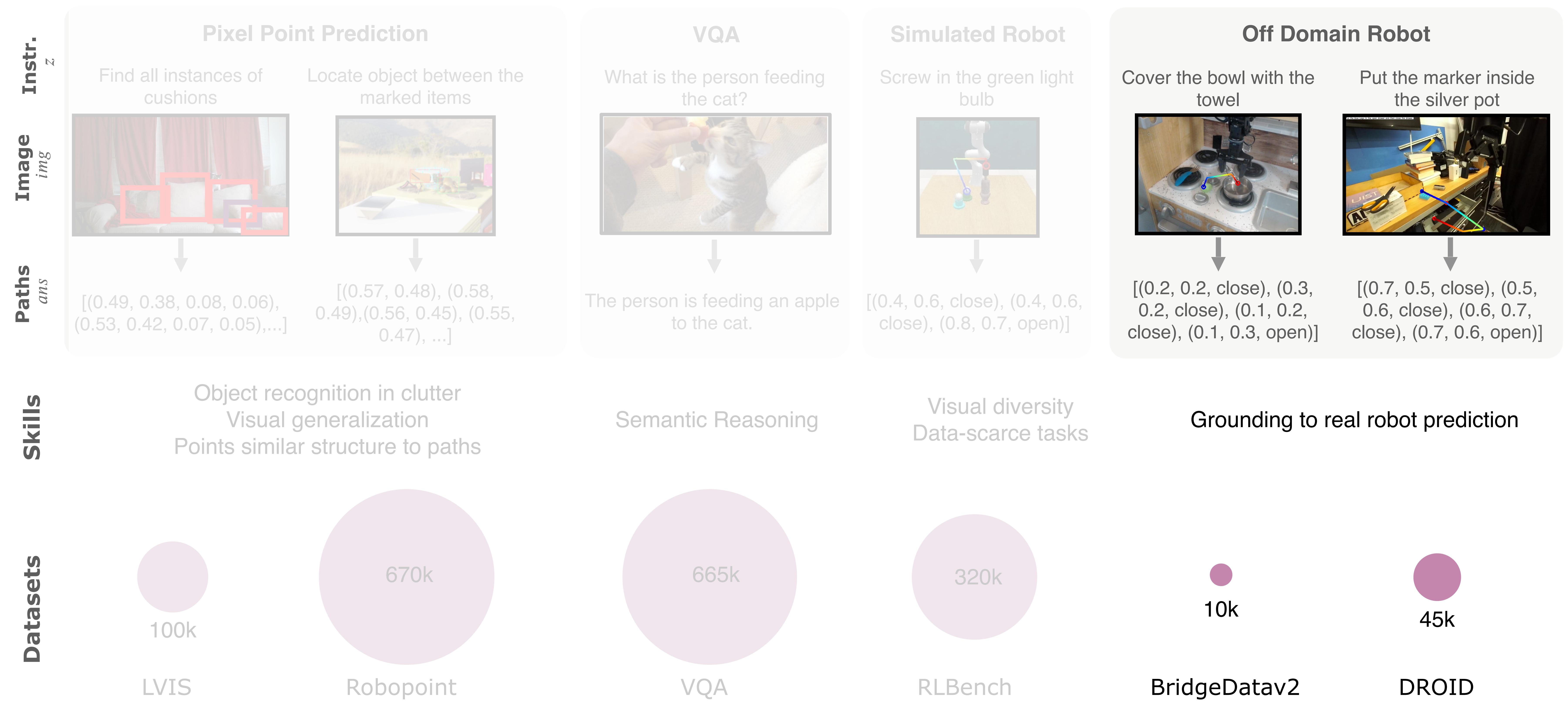
Hamster VLA Training Data

High-Level model is trained on a variety of data sources (Sim, Real, Point-QA)



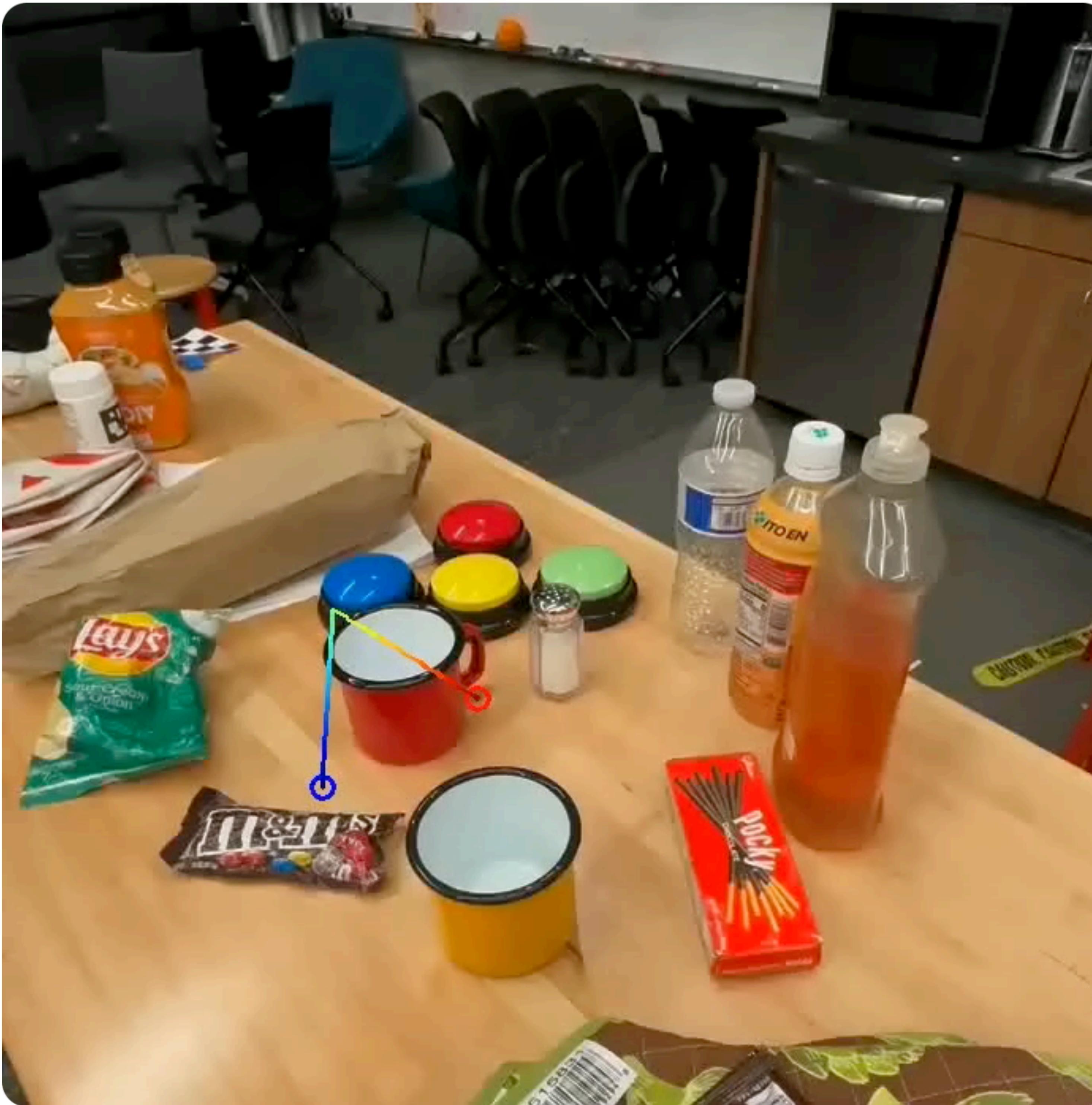
Hamster VLA Training Data

High-Level model is trained on a variety of data sources (Sim, Real, Point-QA)



Hamster VLA Results

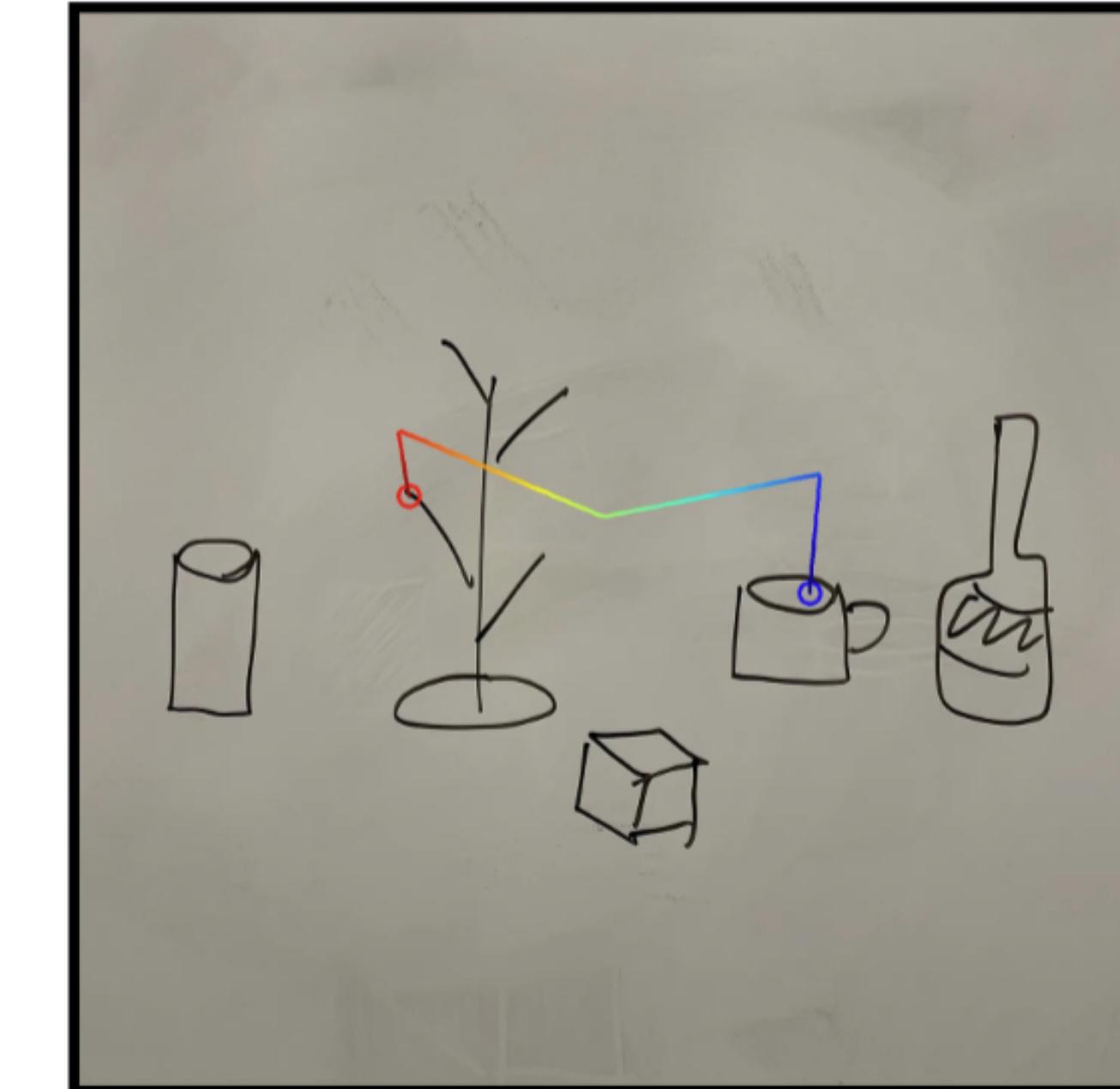
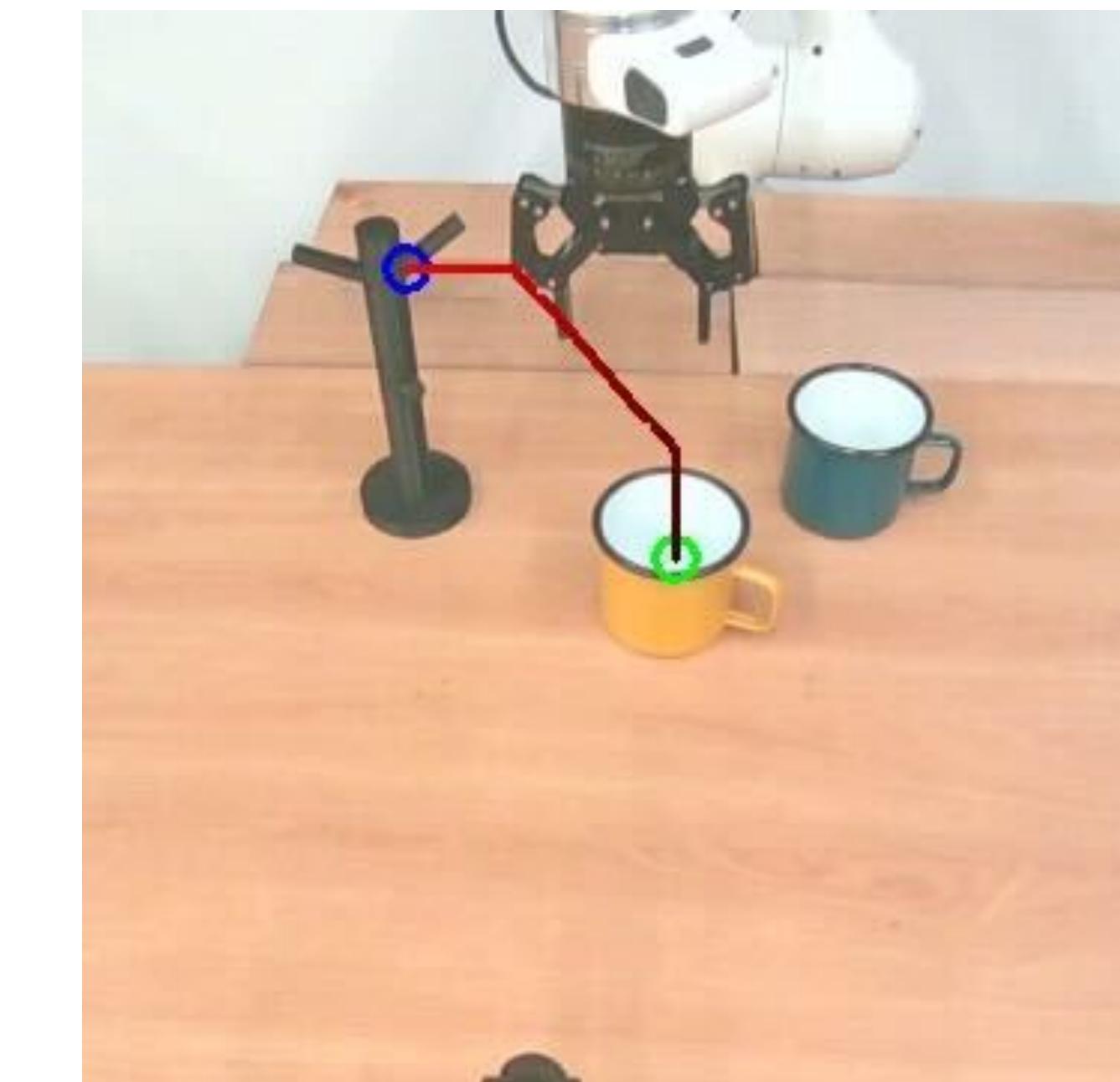
Robust To Novel Camera Position



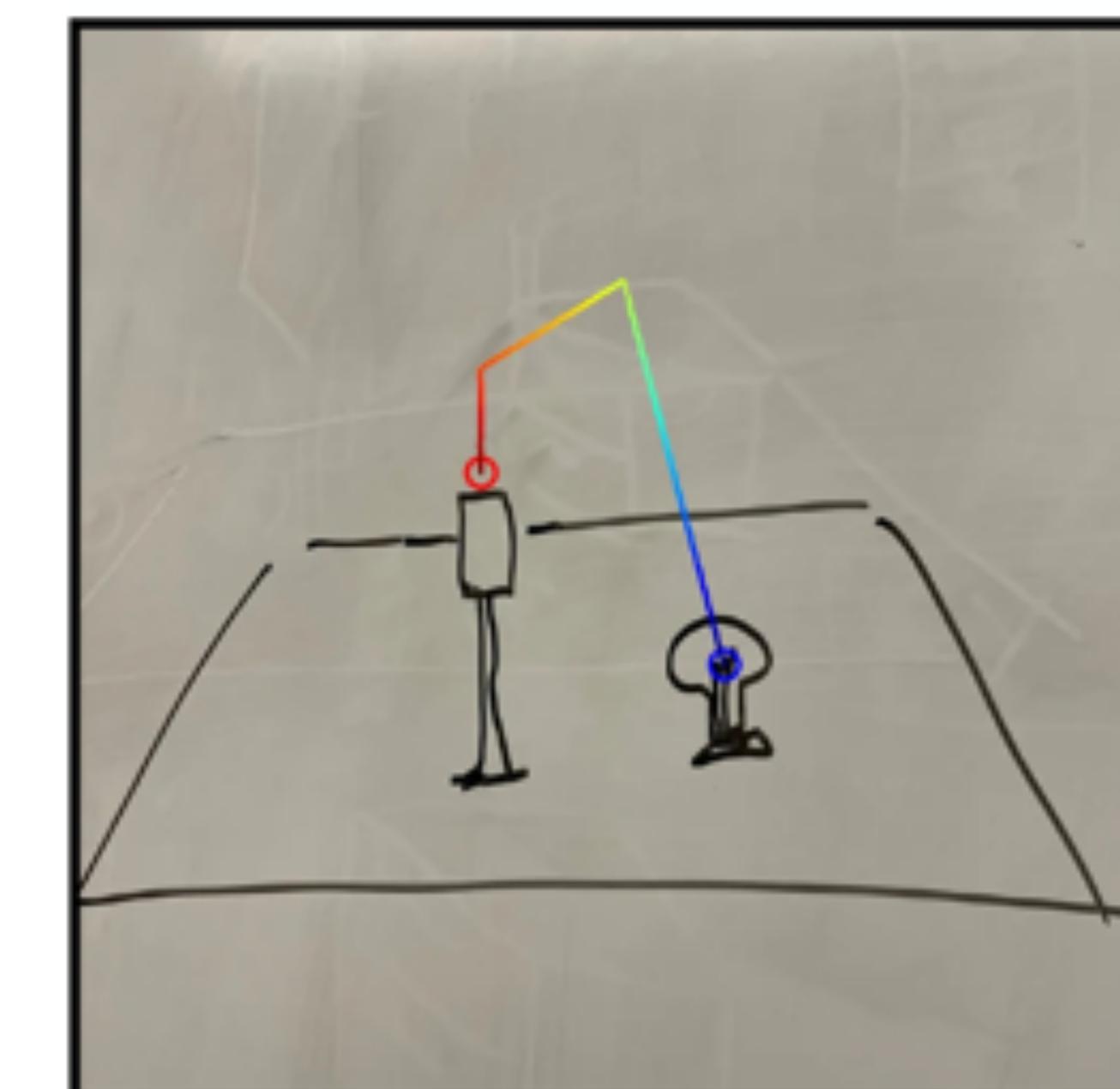
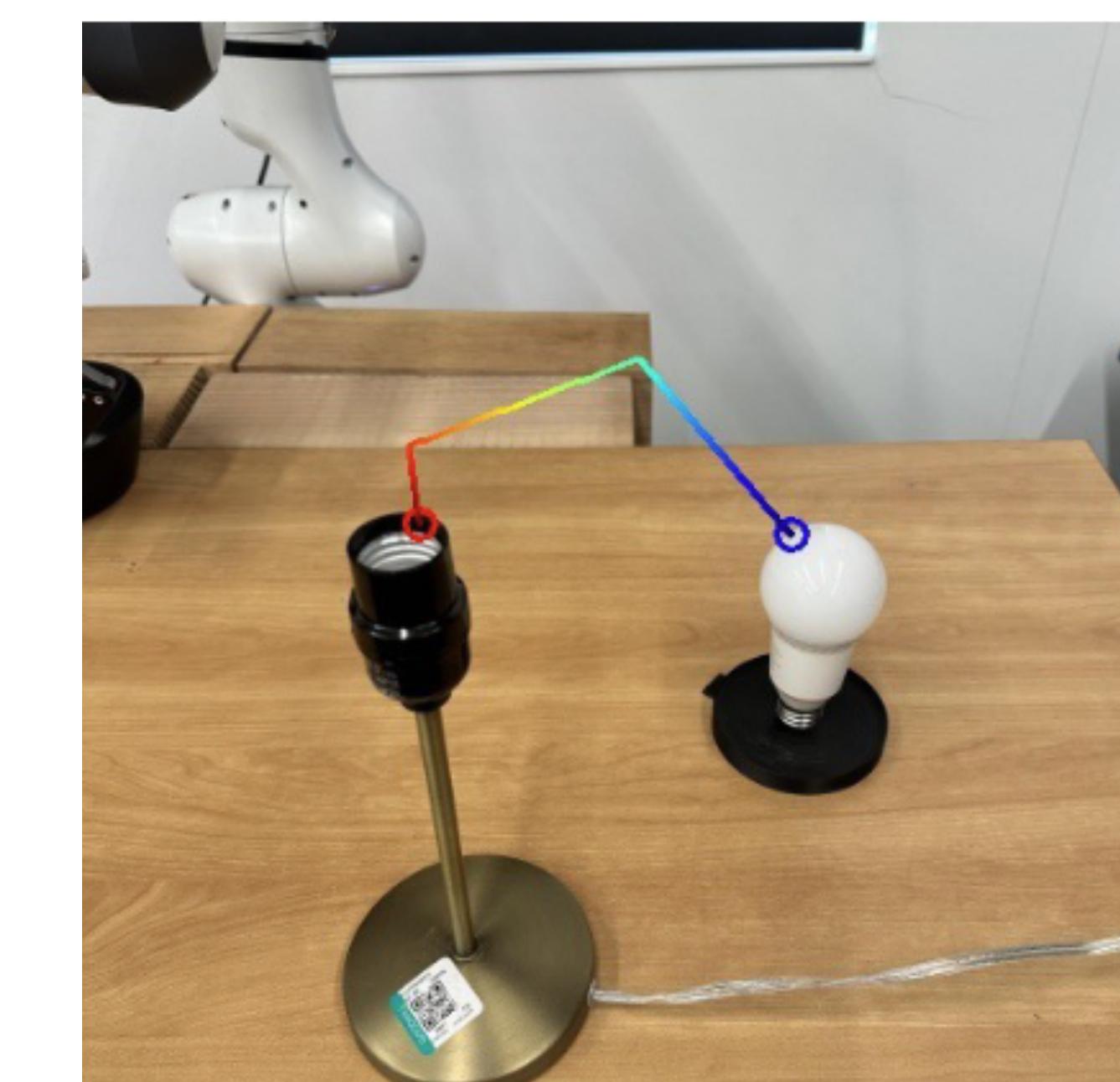
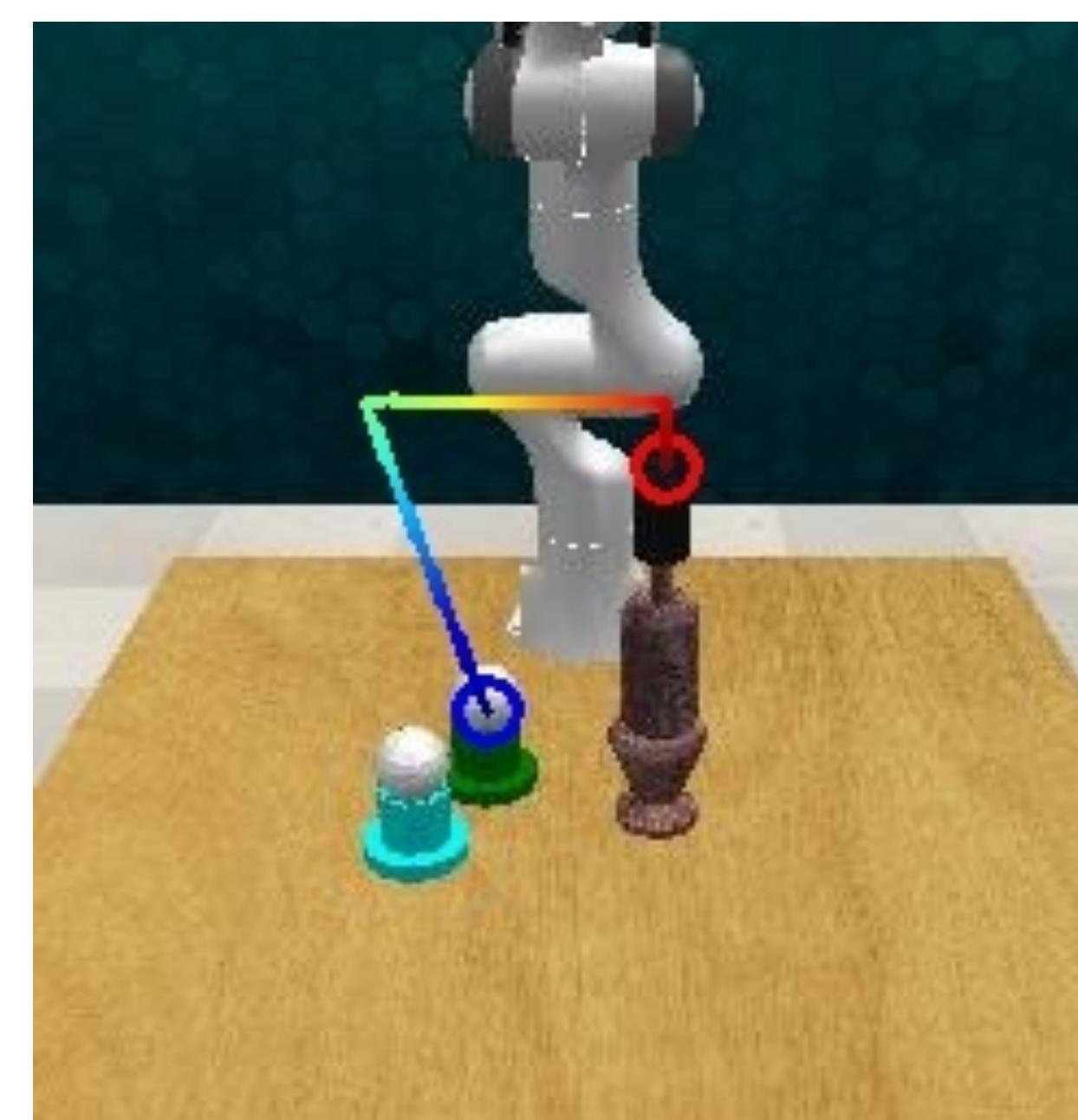
Instruction: Pick up the M&M chocolate and put it in the yellow mug

Hamster VLA Results

Sim to Real to Sketch



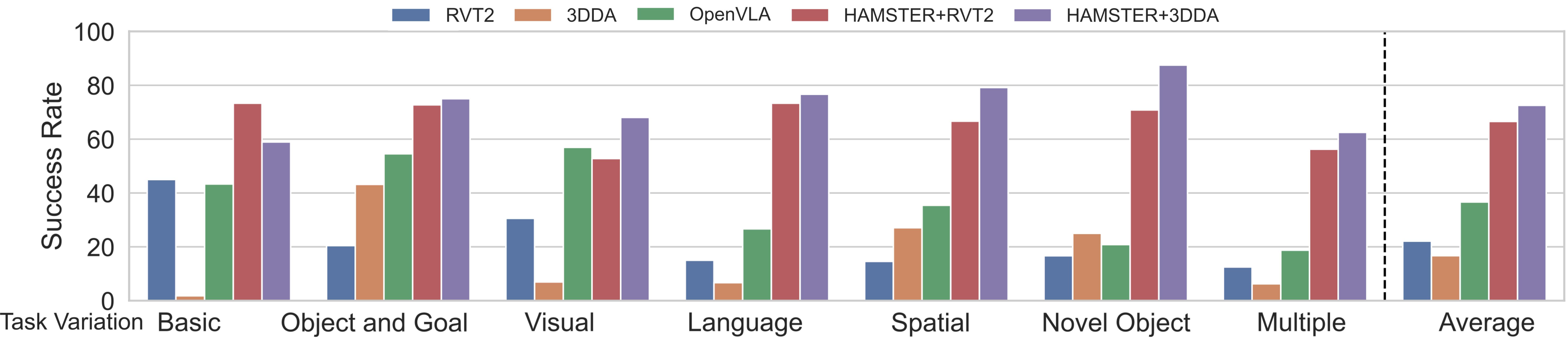
Instruction: Put the cup in the cup holder



Instruction: Screw in the light bulb

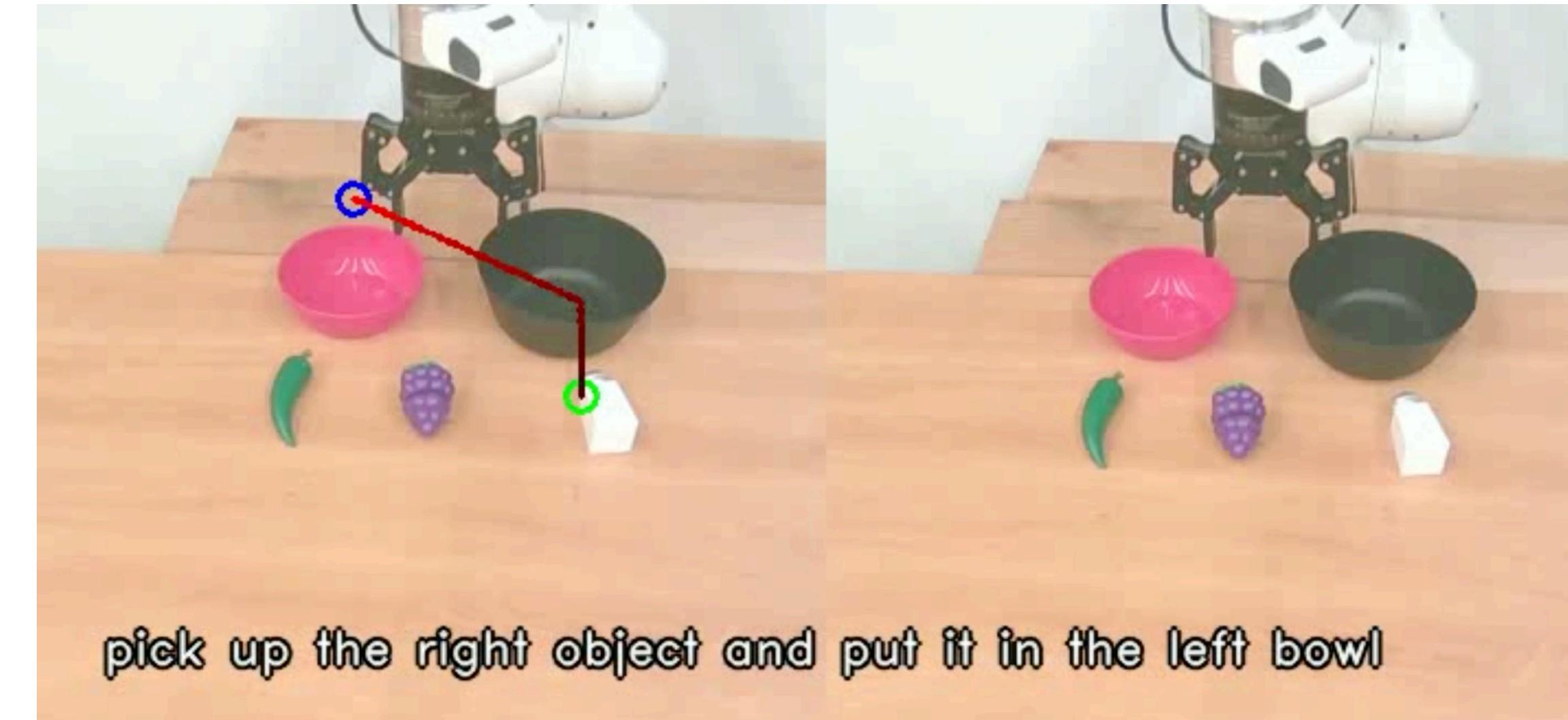
Hamster Results

Outperforms OpenVLA as well as 3D Policies

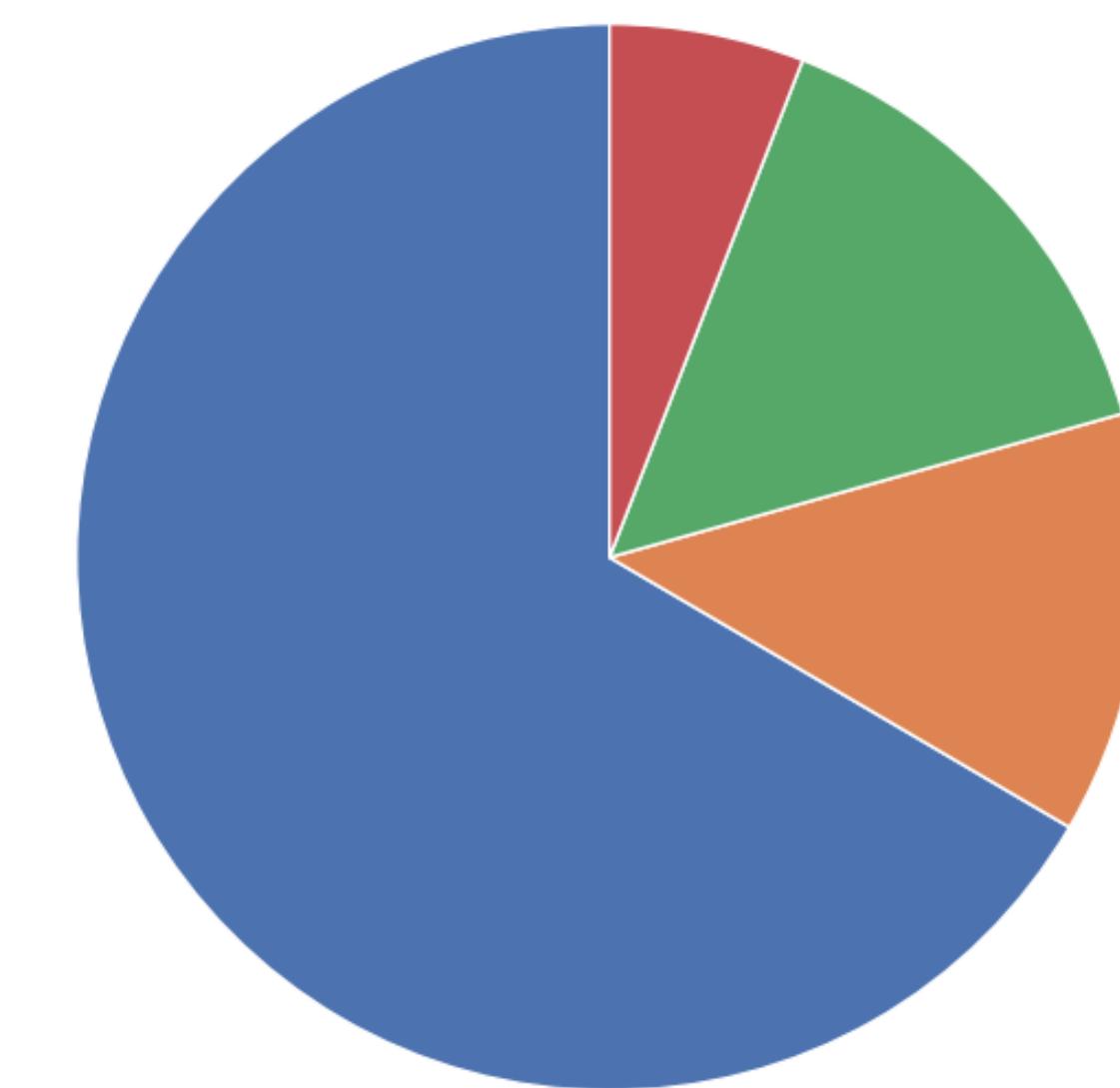


Failure Analysis

Low-Level Policy Struggles

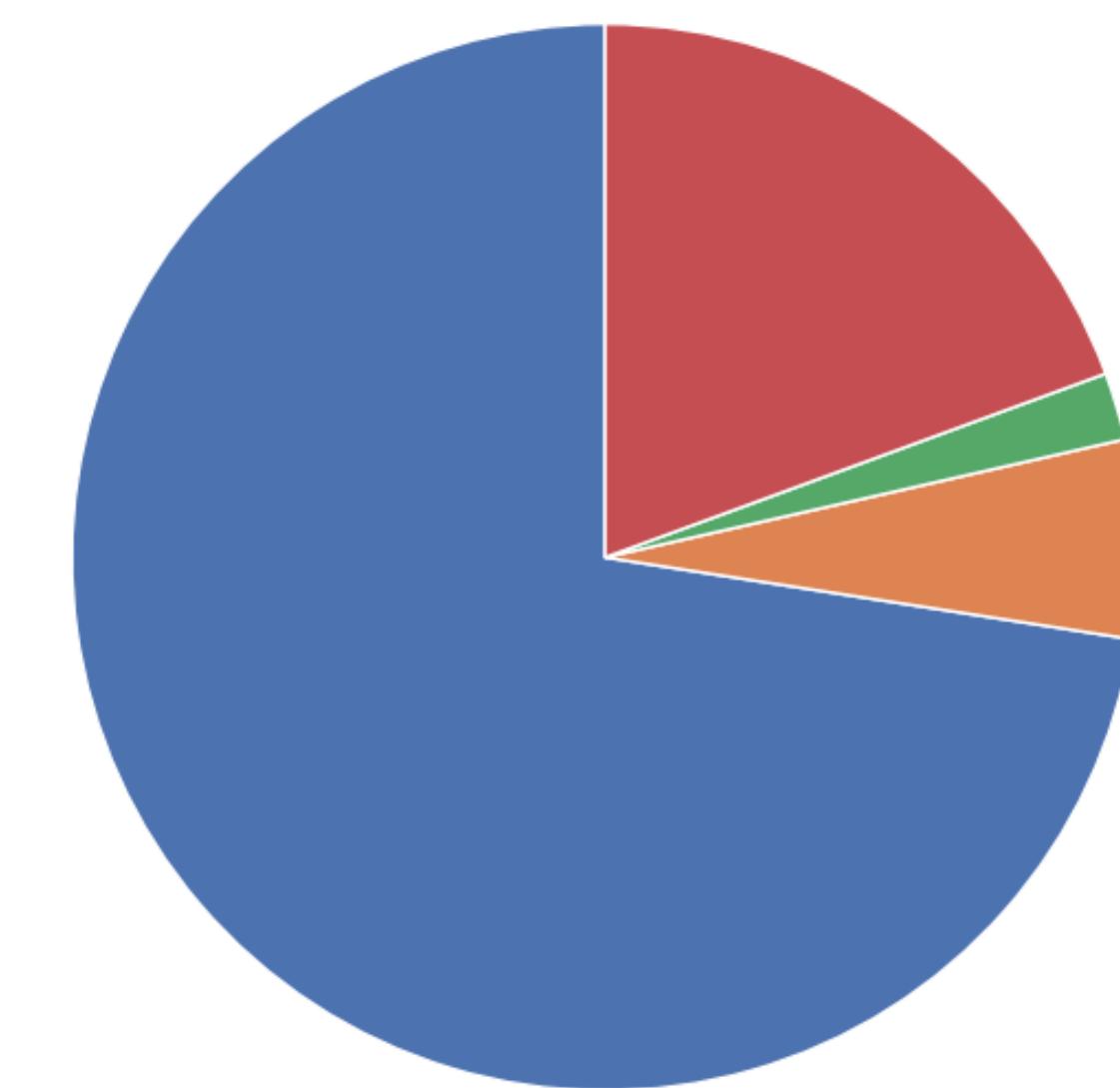


VLM Failure

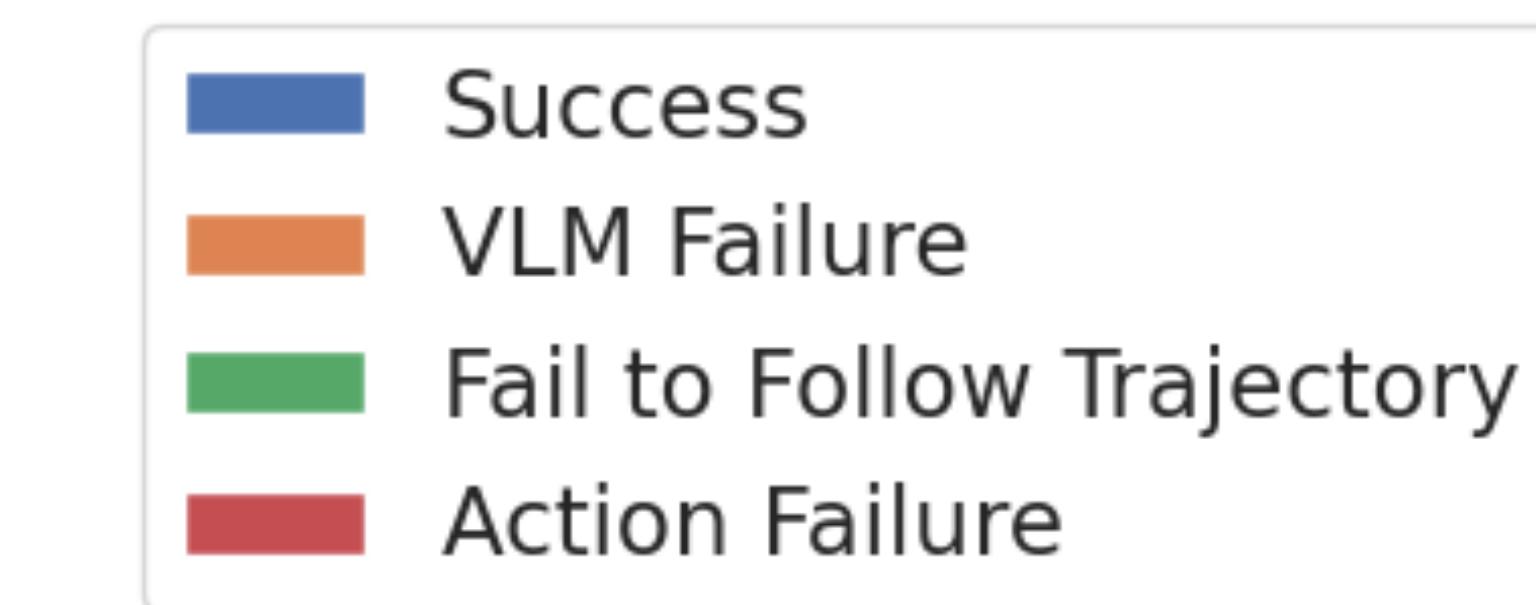


HAMSTER with RVT-2

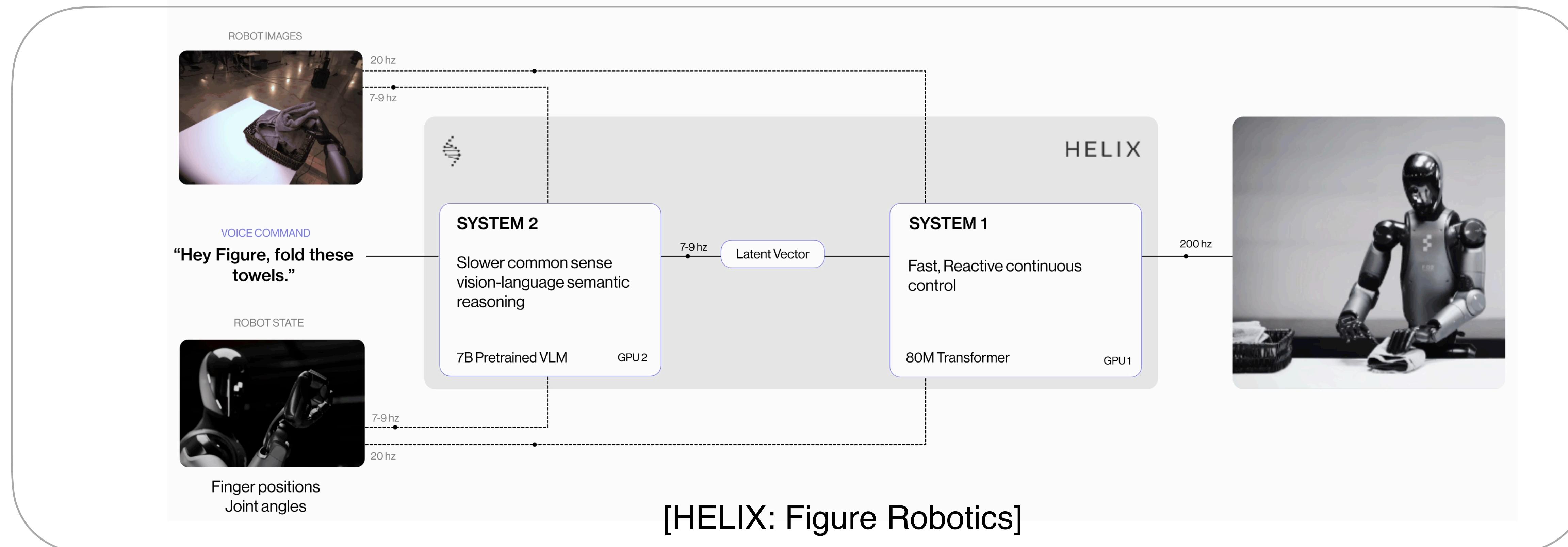
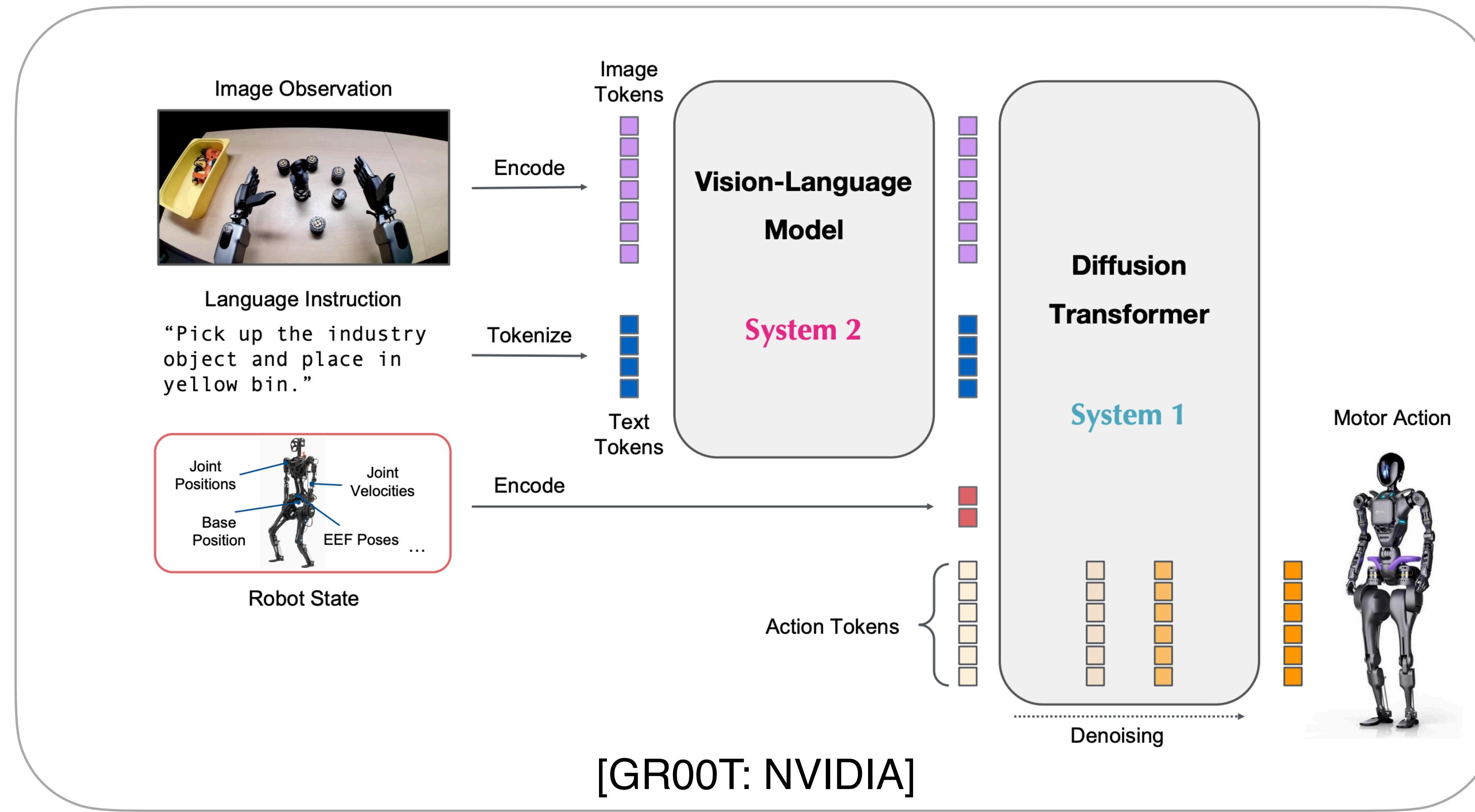
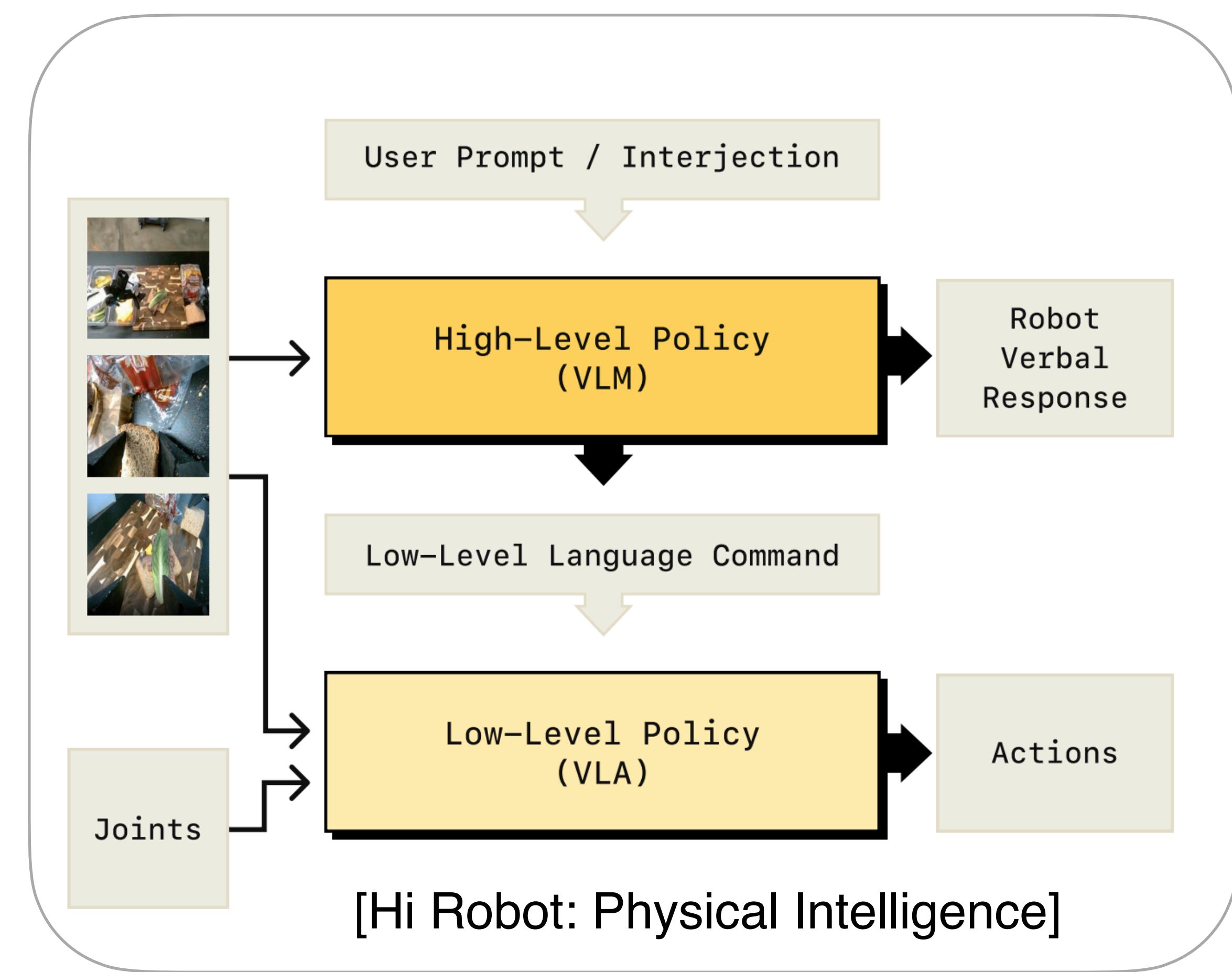
Fail to Follow Trajectory



HAMSTER with 3D-DA

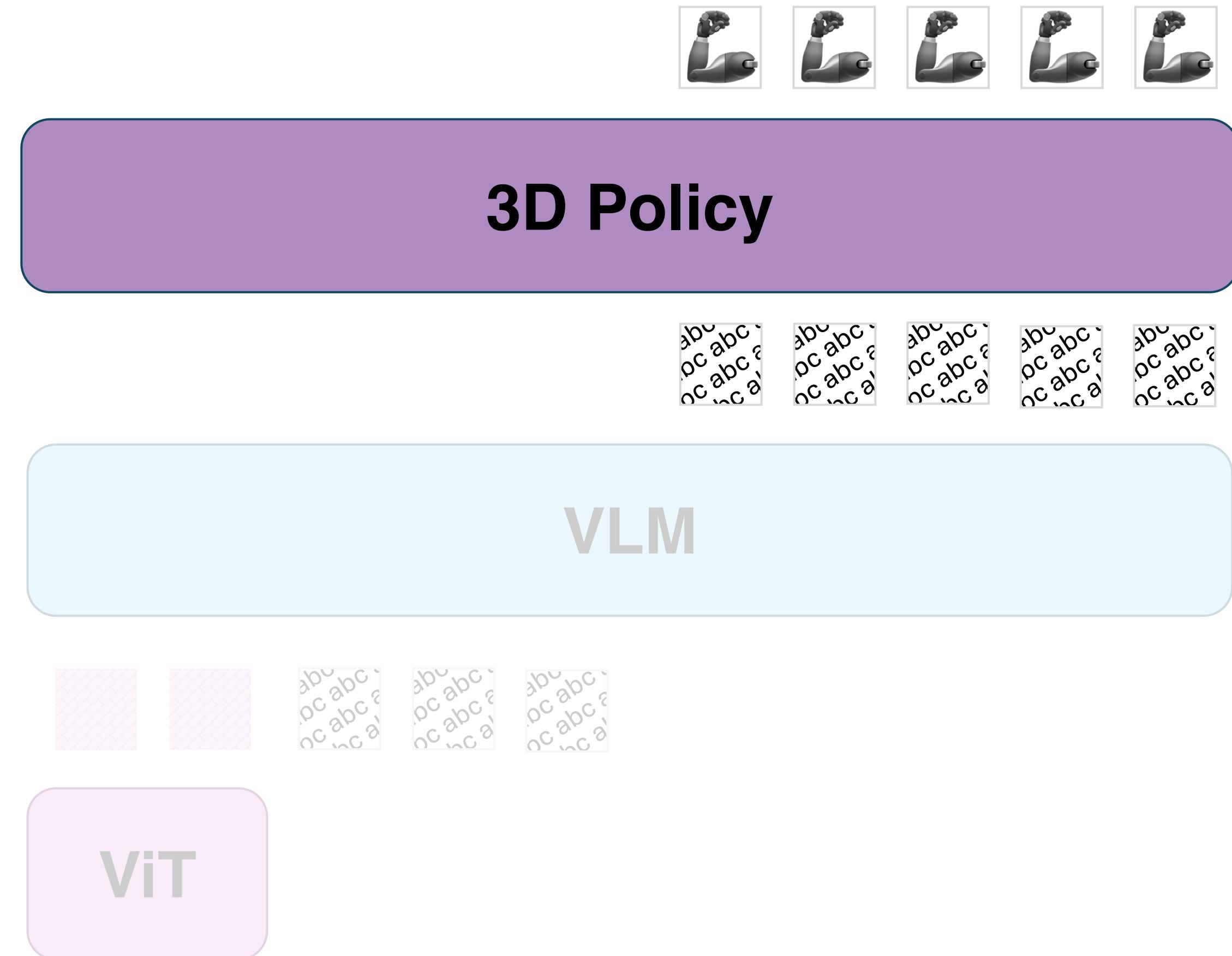


Hierarchical VLAs: Concurrent Works



3D Policies

Specialized and Efficient



- Take in Scene depth and / or camera calibration
- Specialized policies – Require very few demonstrations

RVT: Robotic View Transformer for 3D Object Manipulation

CoRL 2023



Ankit Goyal



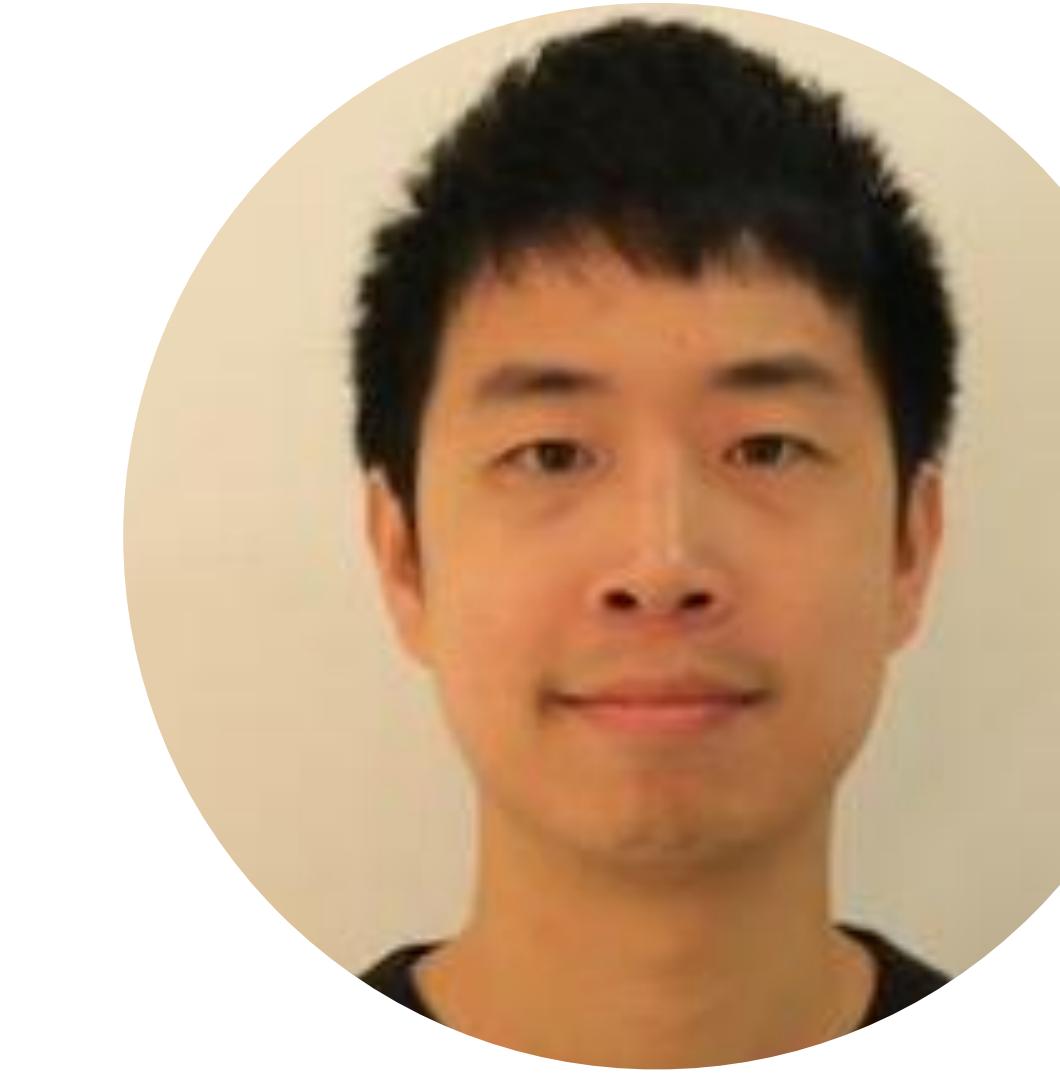
Jie Xu



Yijie Guo



Valts Blukis

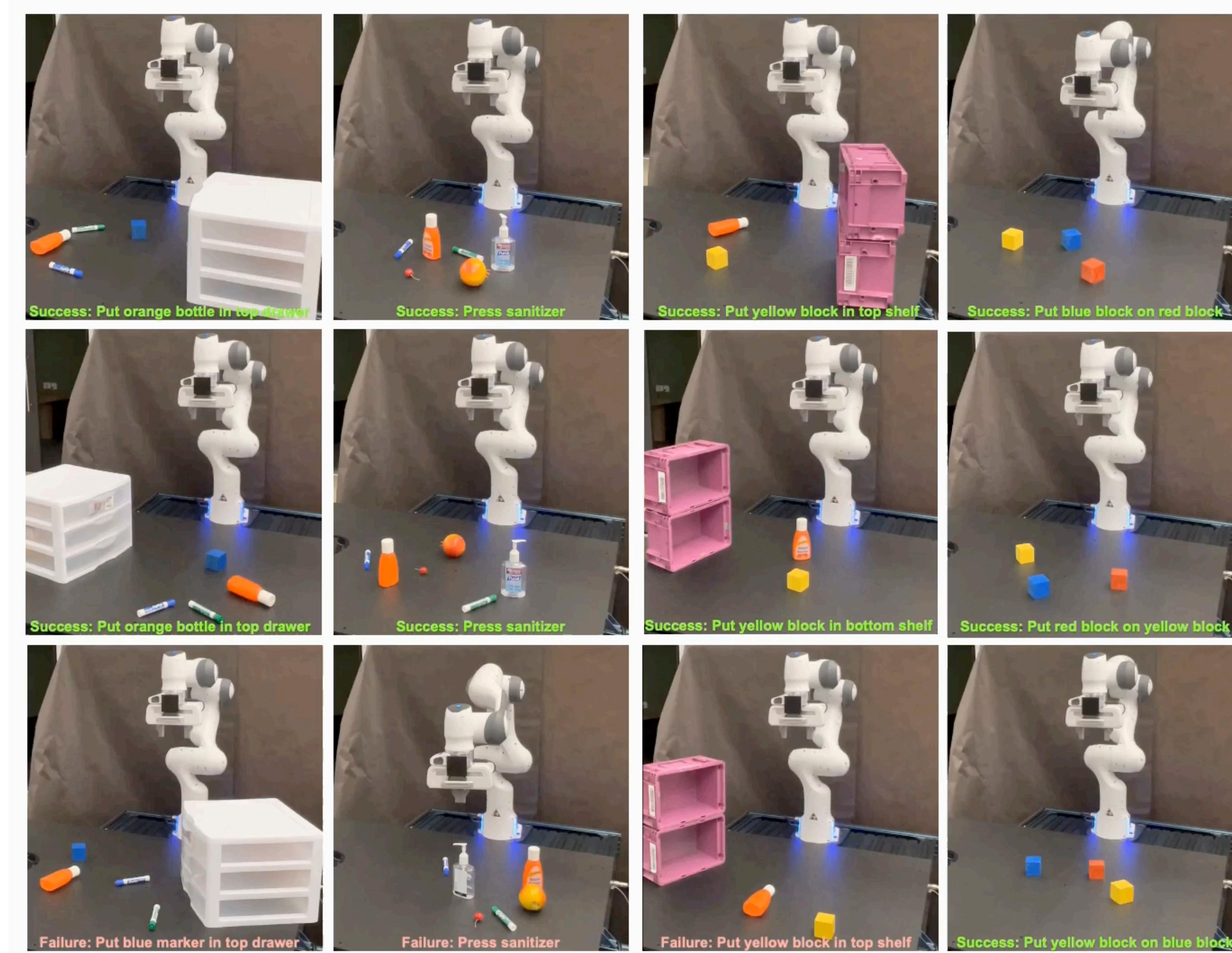


Yu-Wei Chao



Dieter Fox

RVT: Robotic View Transformer



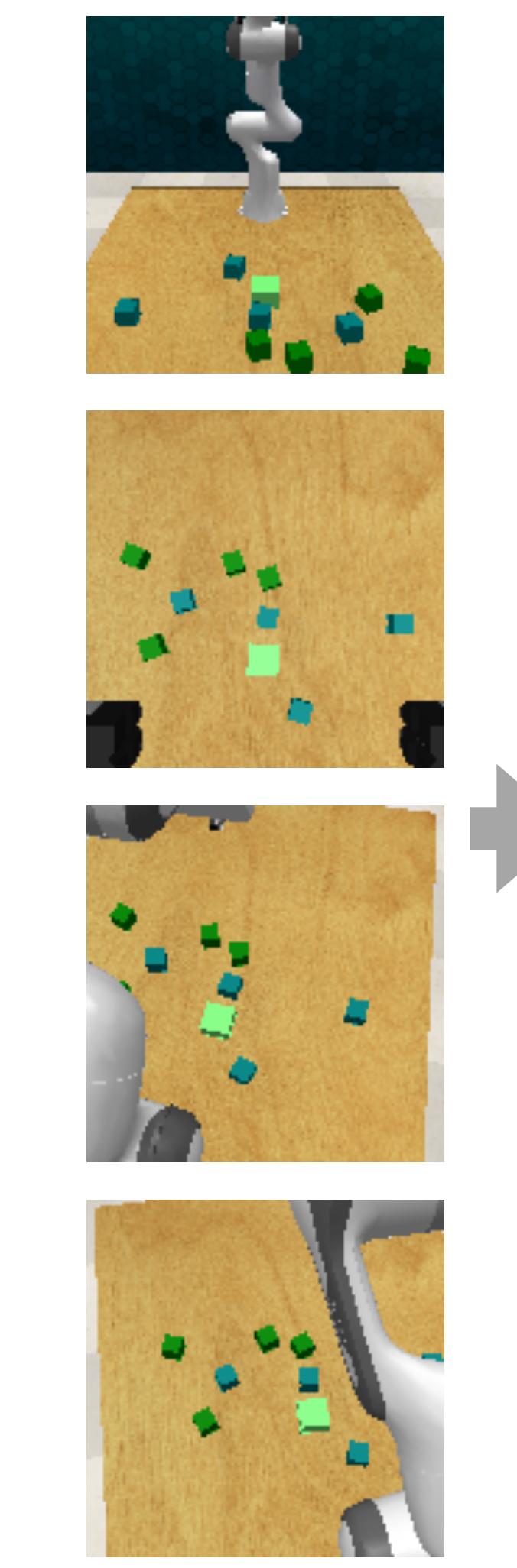
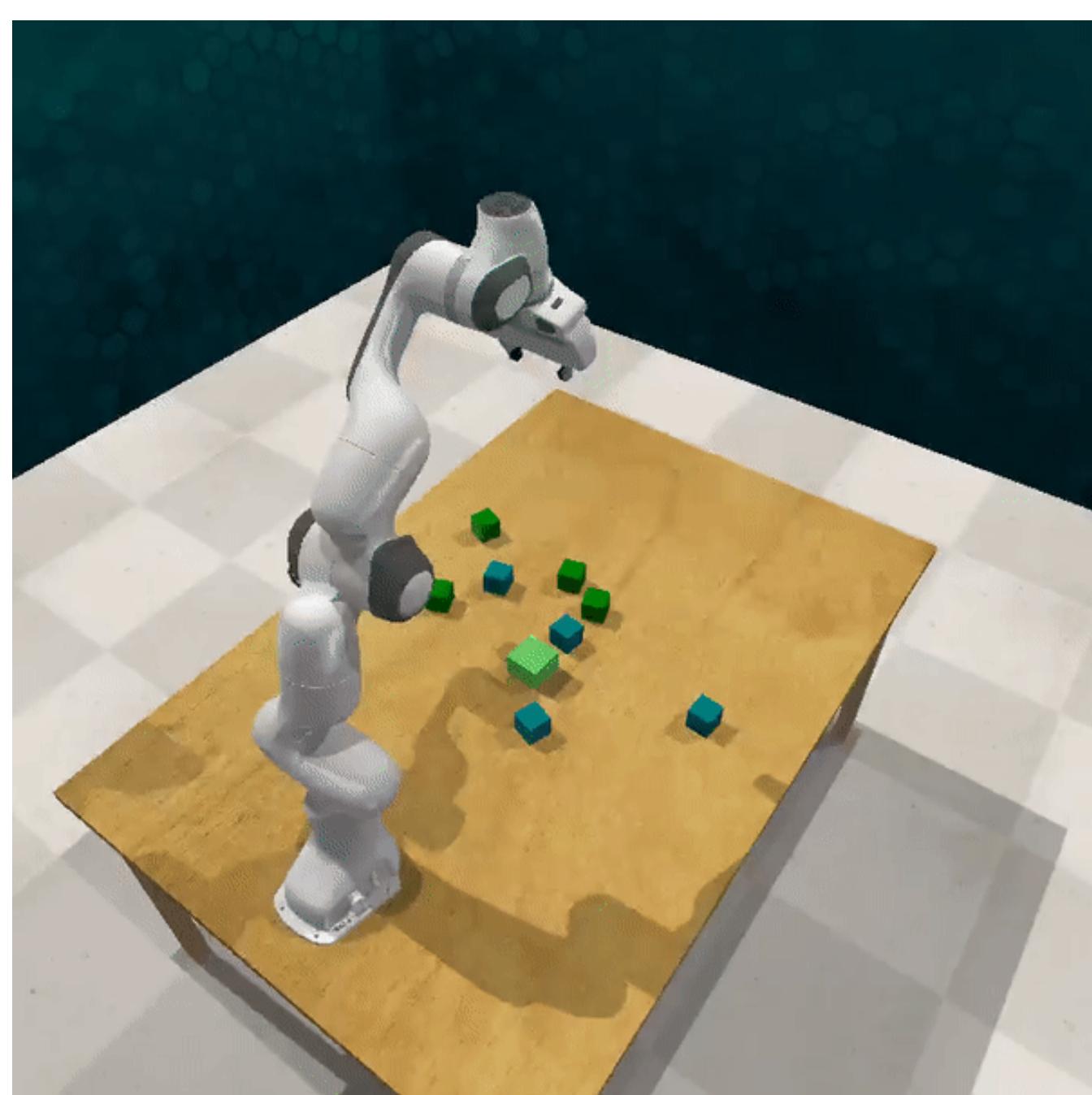
A single model achieves 90% success across tasks with just ~10 demos each.

RVT: Robotic View Transformer for 3D Manipulation

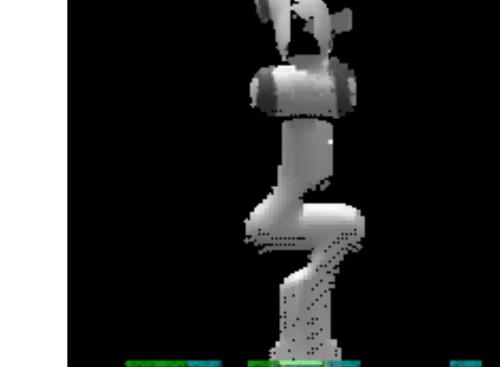
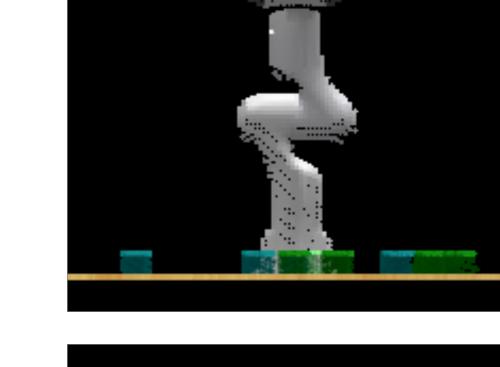
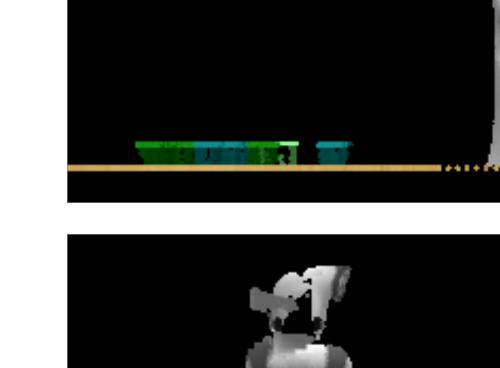
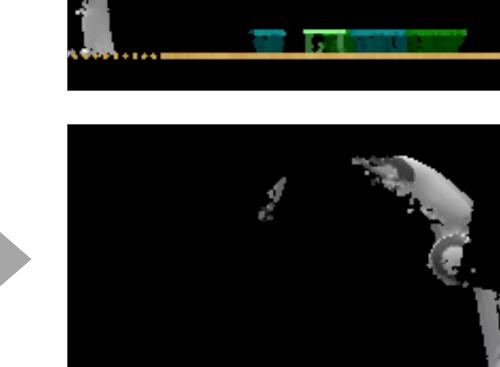
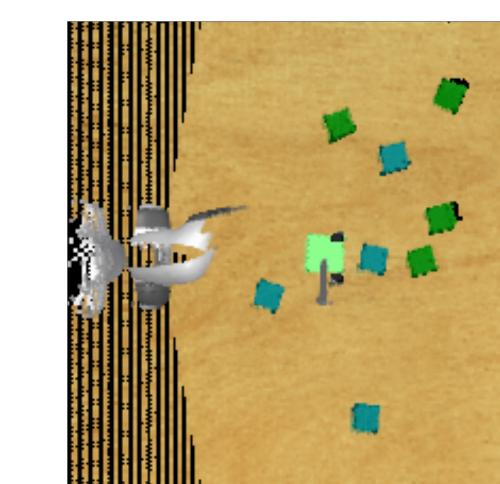
Pipeline Illustration

Instructions

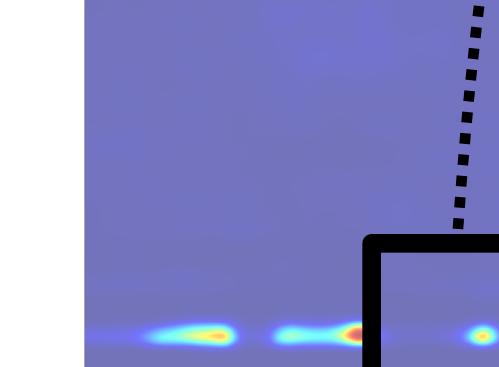
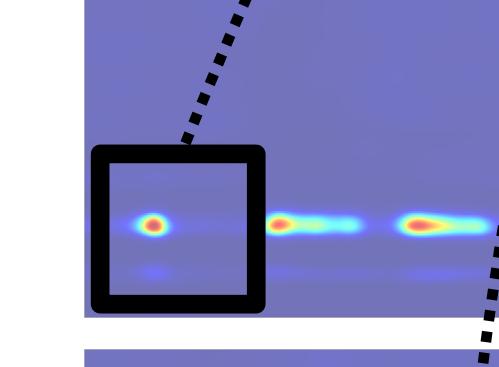
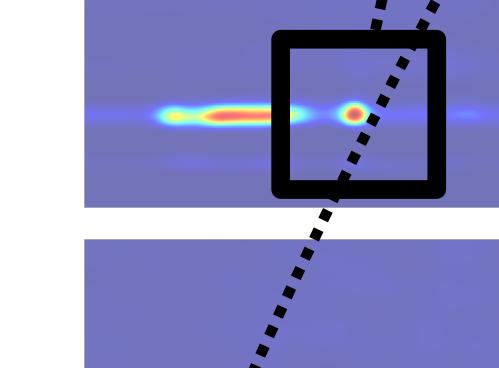
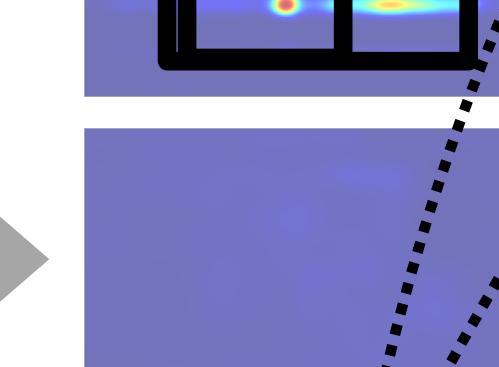
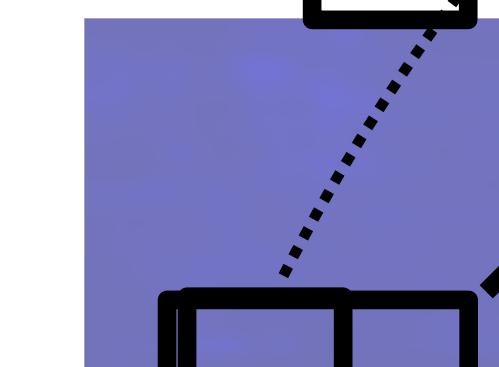
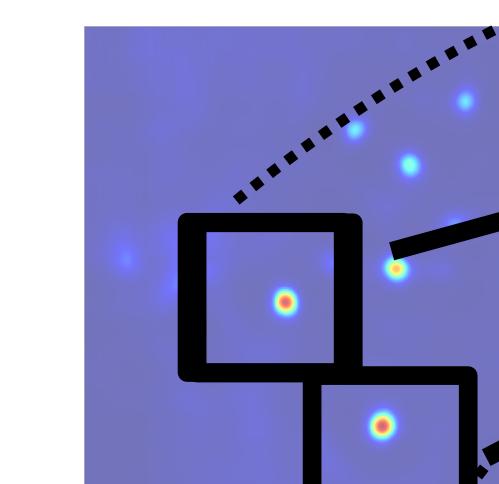
"Stack the teal blocks"



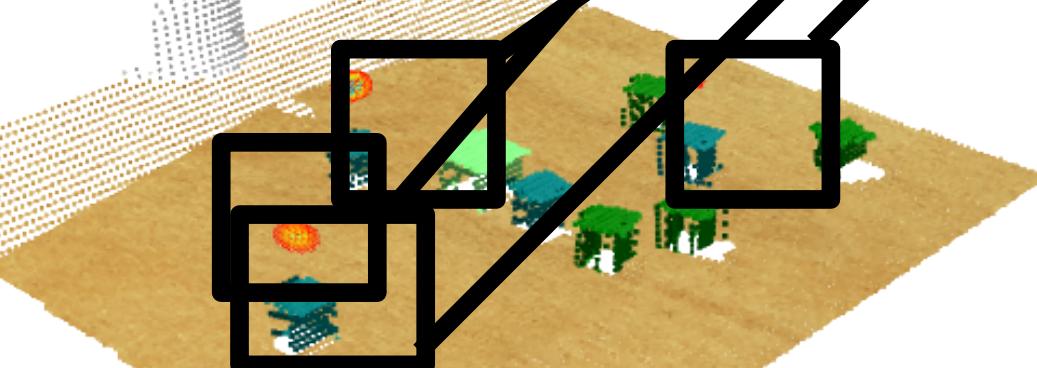
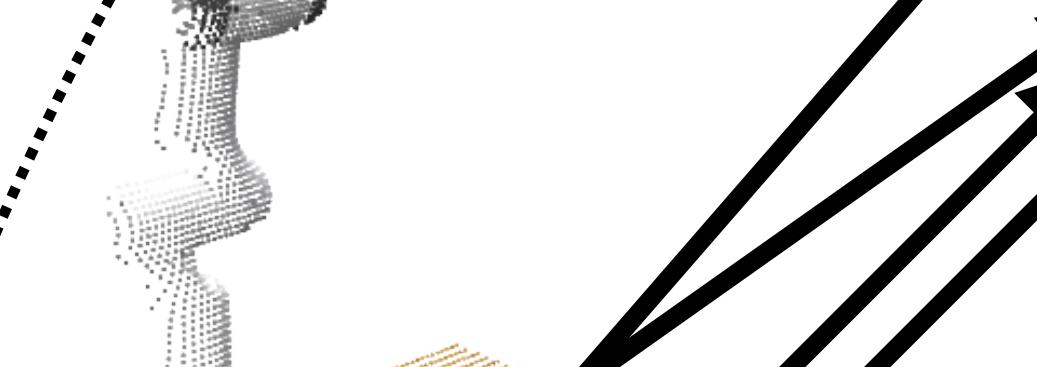
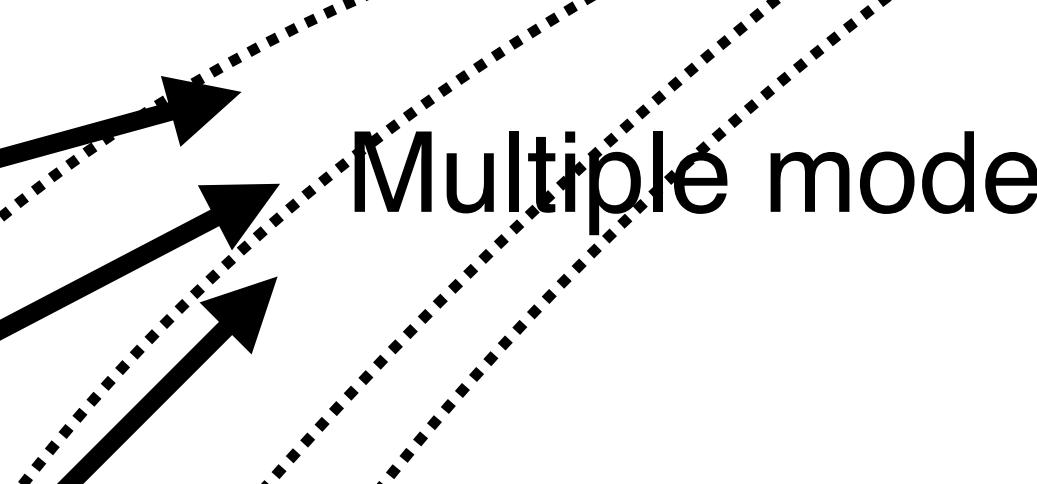
Reconstructed Point Cloud



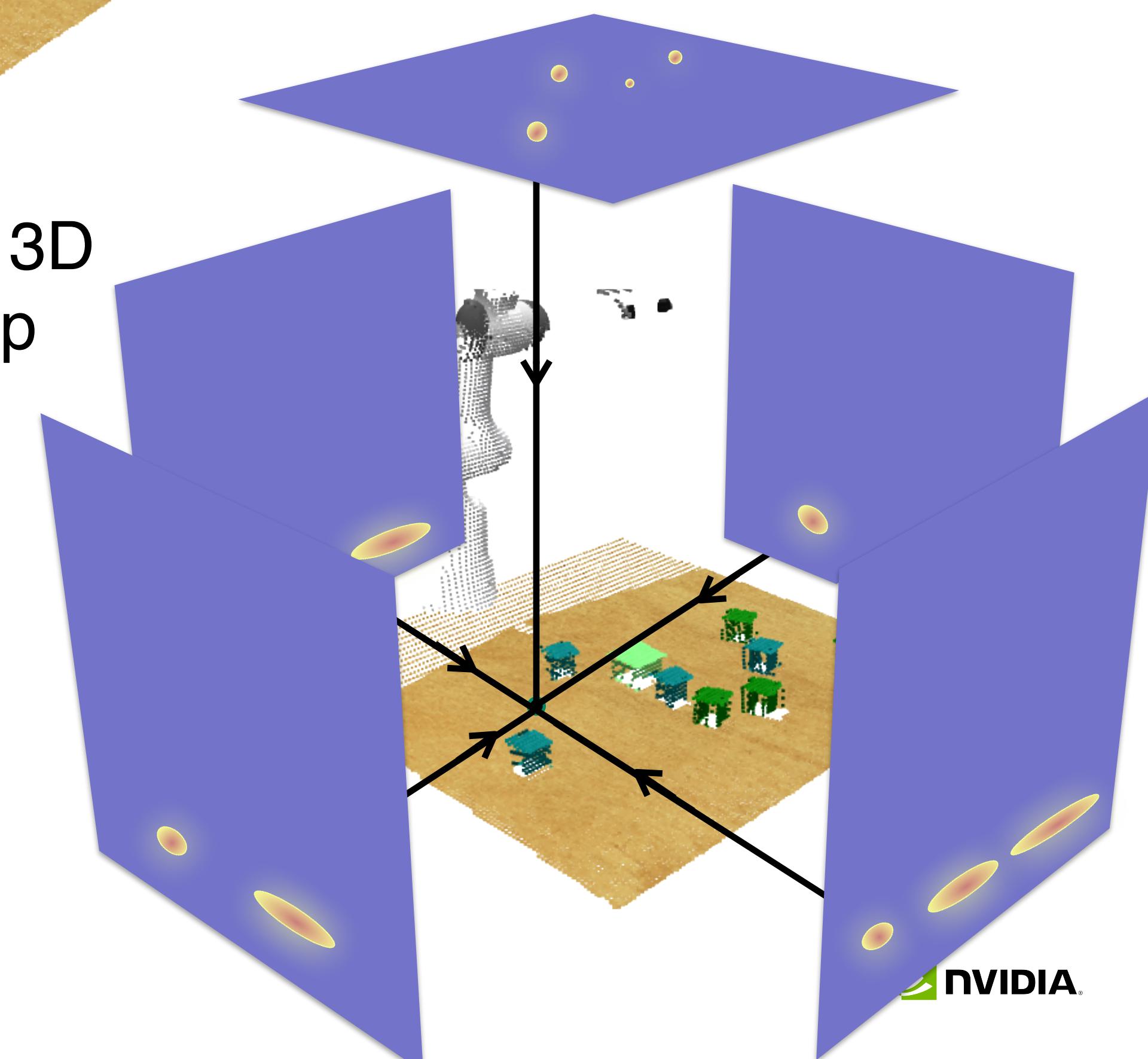
Virtual Images



Predicted 2D Heatmap



Predicted 3D Heatmap



Gripper Rotation
and State

Gripper Location

NVIDIA

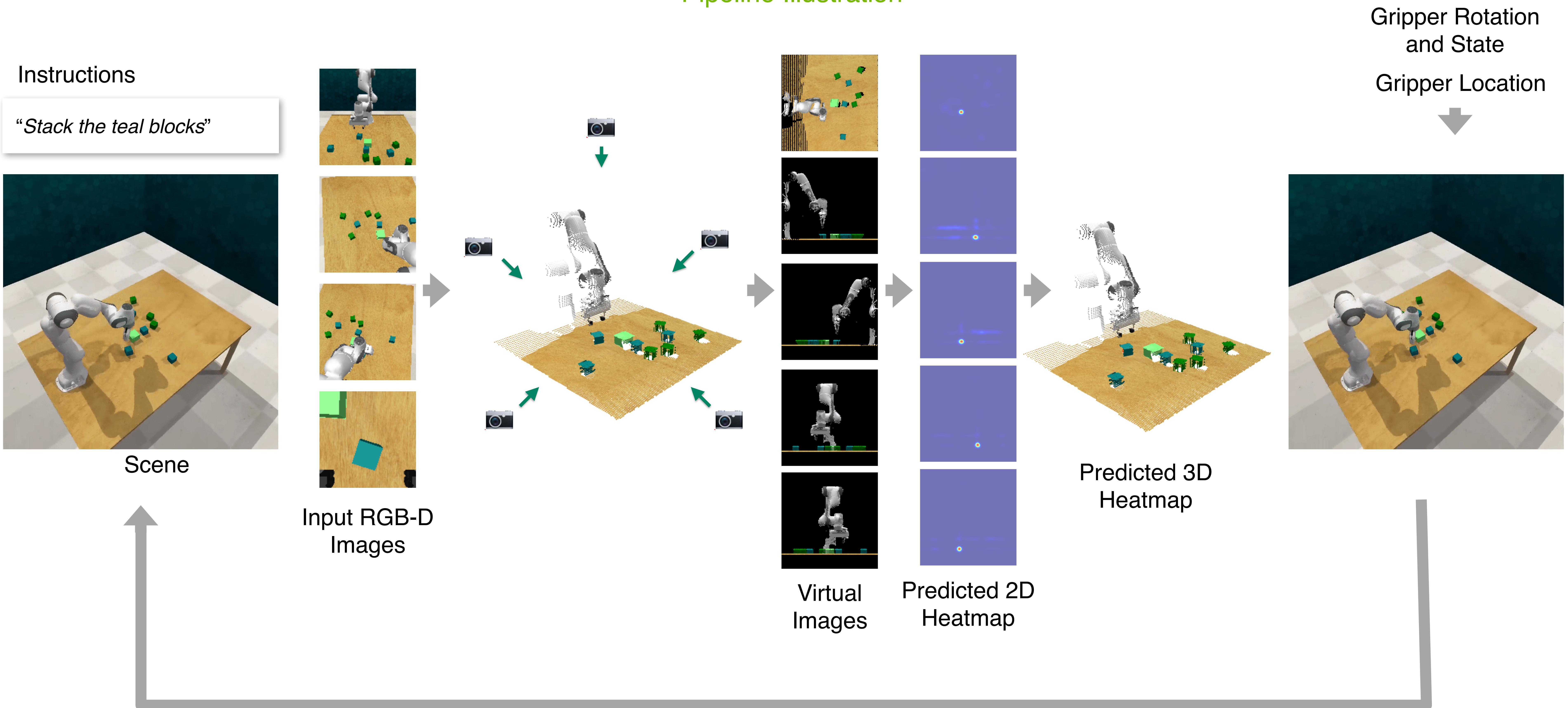
RVT: Robotic View Transformer for 3D Manipulation

Pipeline Illustration



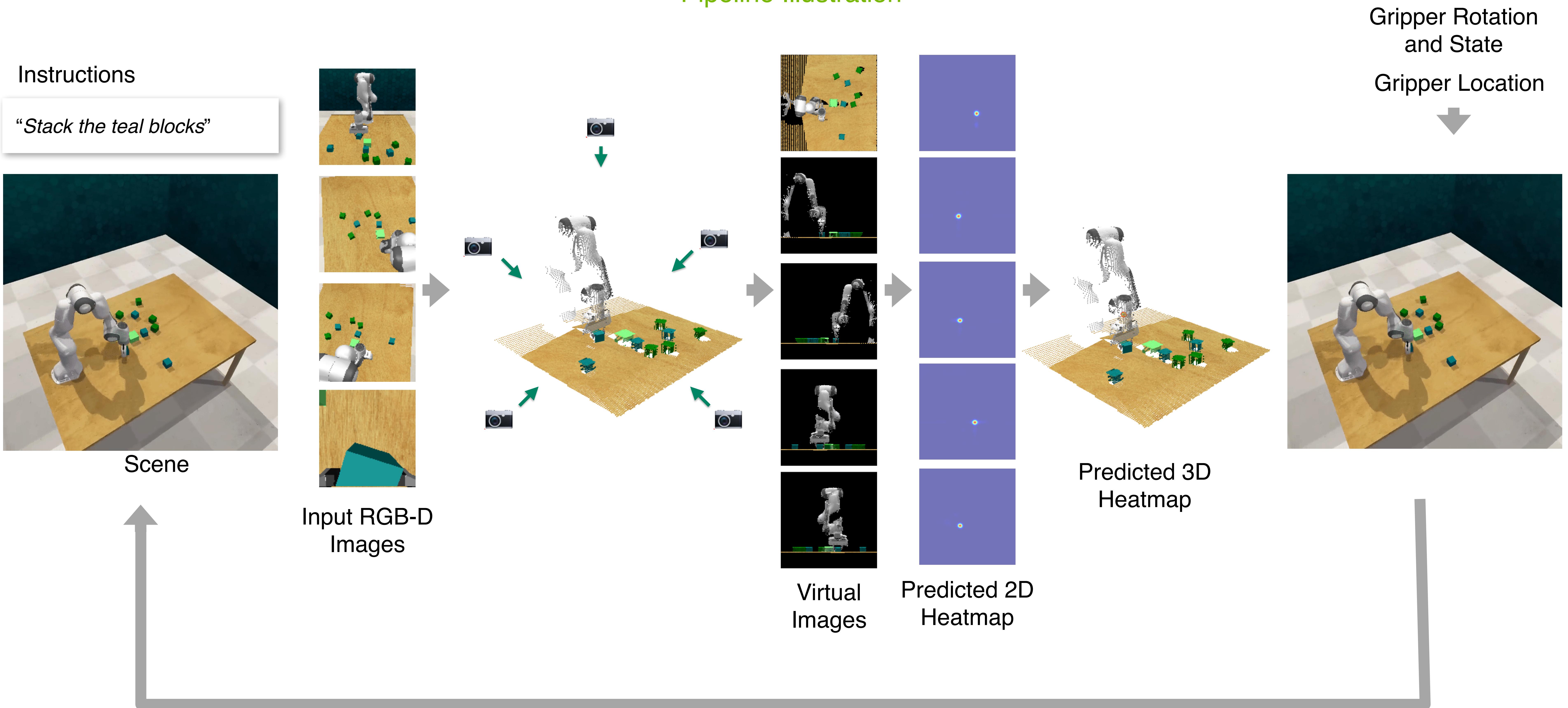
RVT: Robotic View Transformer for 3D Manipulation

Pipeline Illustration



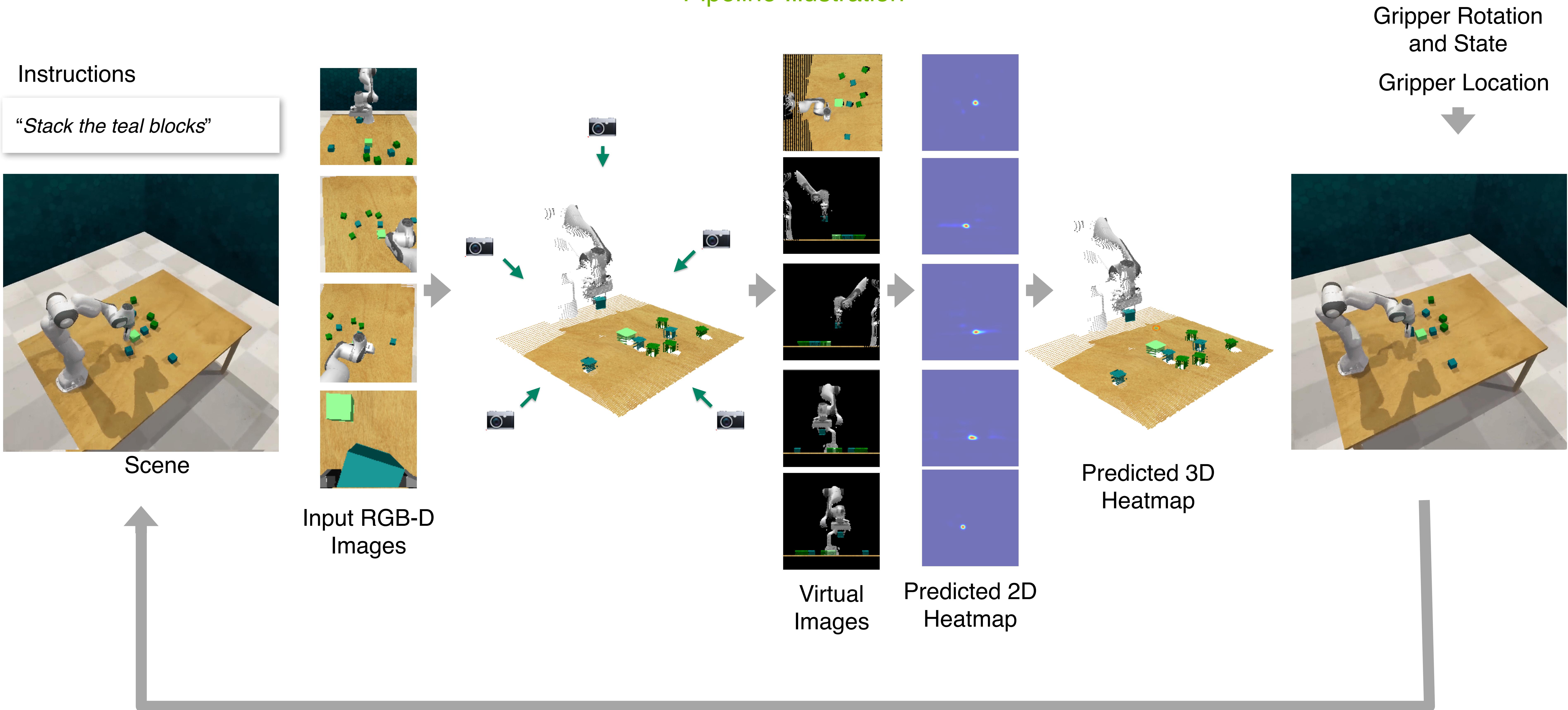
RVT: Robotic View Transformer for 3D Manipulation

Pipeline Illustration



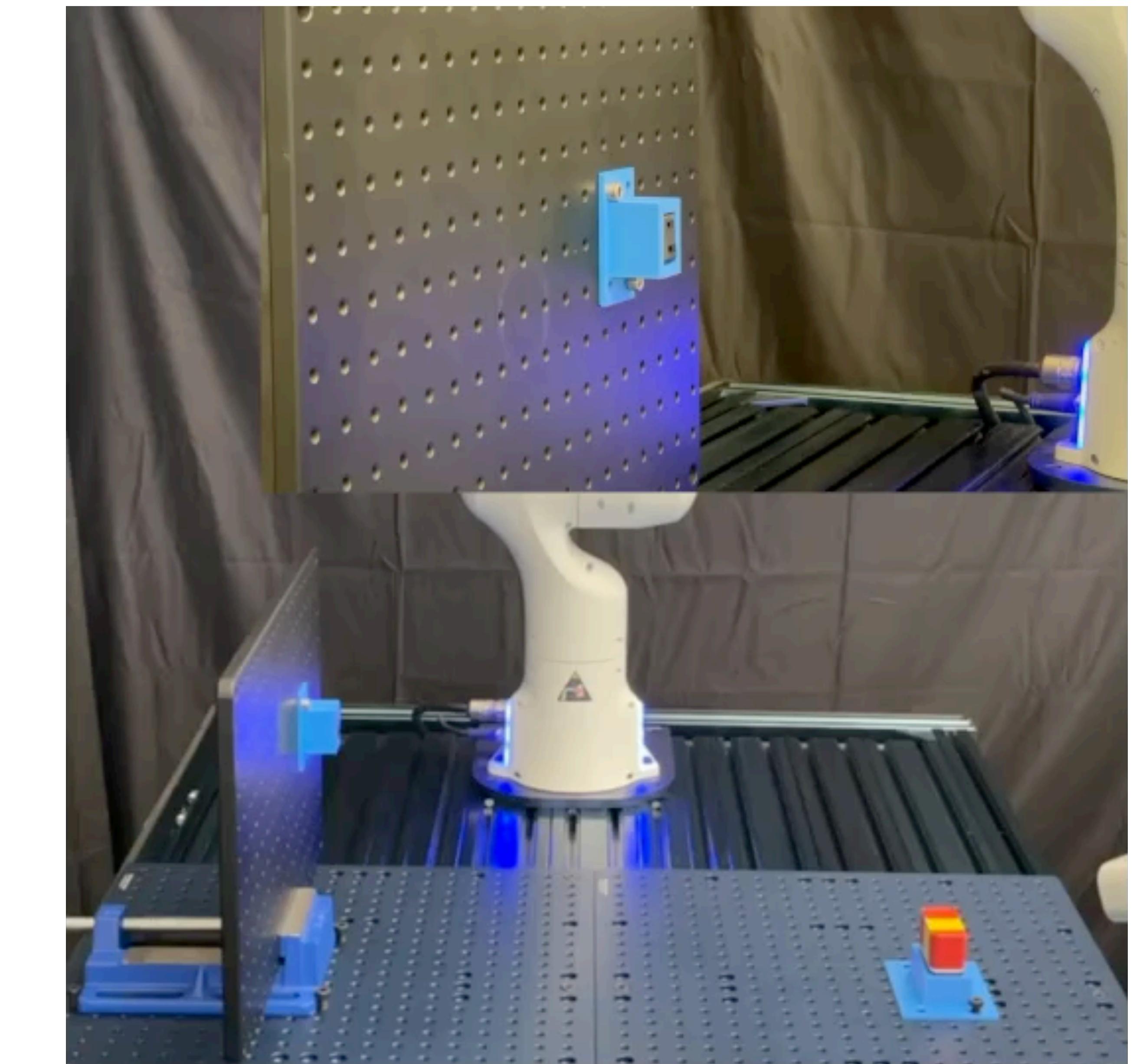
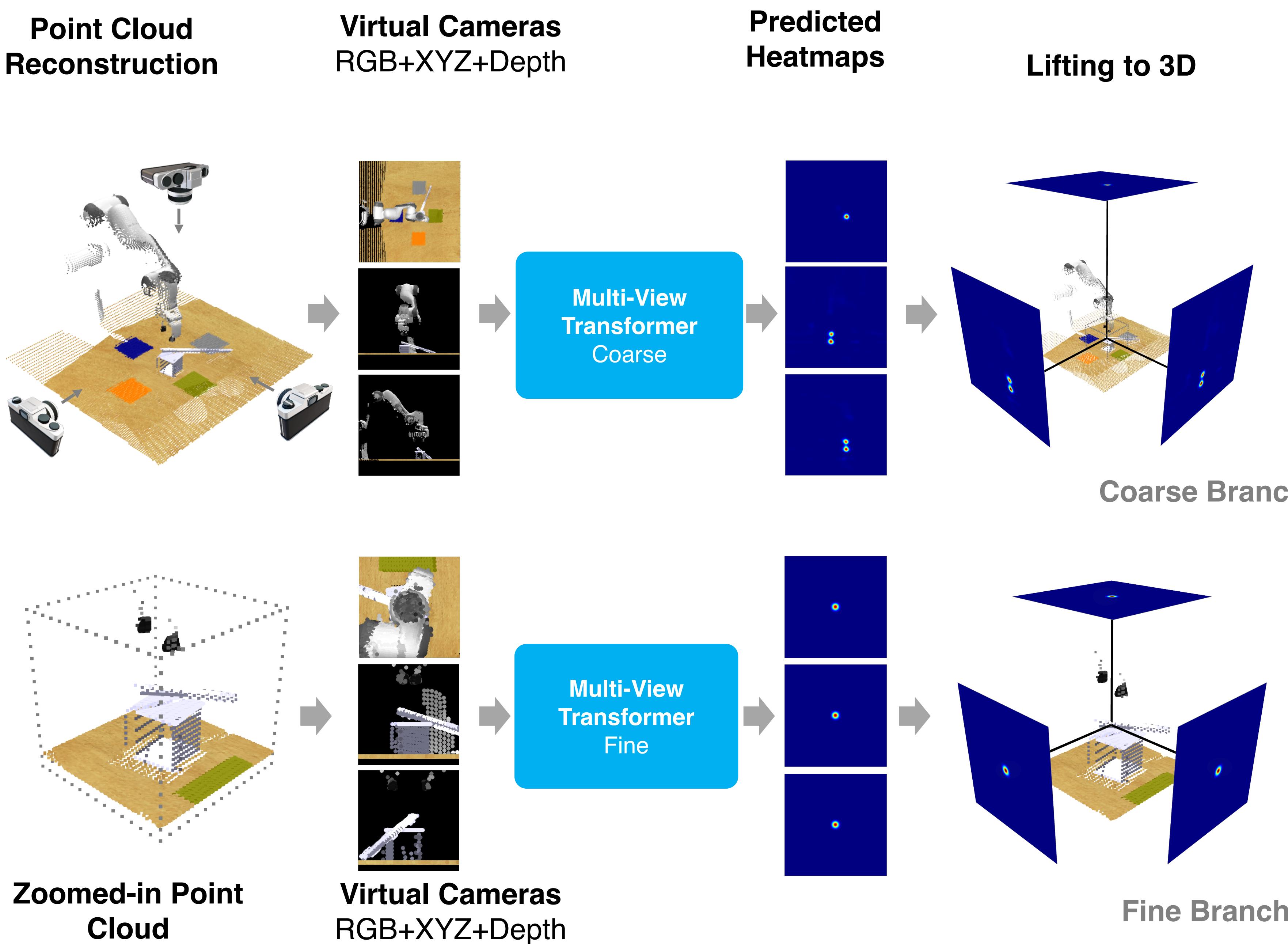
RVT: Robotic View Transformer for 3D Manipulation

Pipeline Illustration



RVT-2: Learning Precise Manipulation from Few Examples

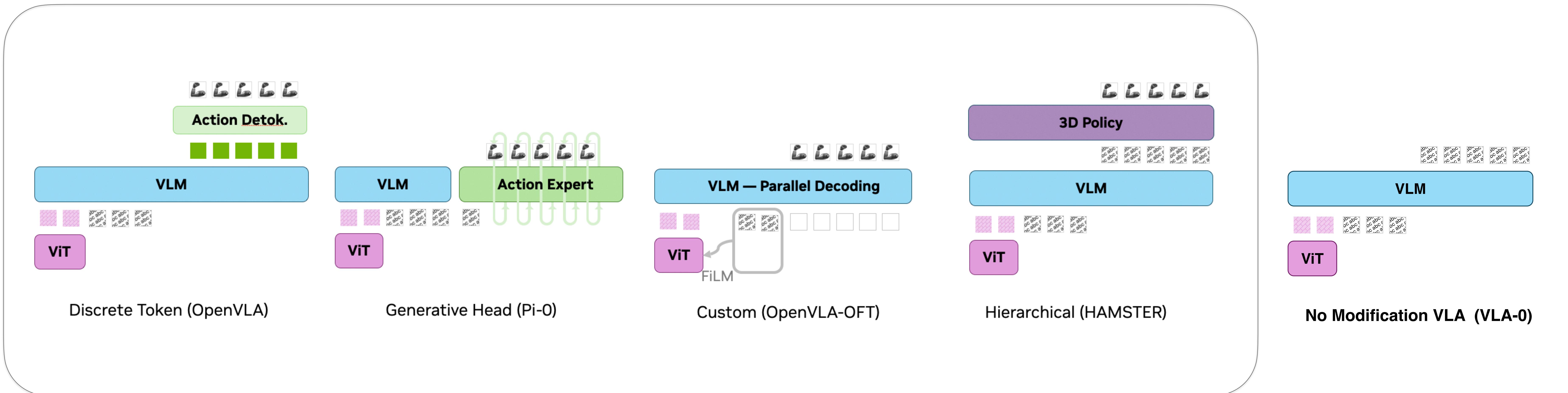
RSS 2024



Multiscale design allows for precise manipulation with just 10 demos

Family of VLAs

- How about the simplest variant?
- Predict Action as Text — No modification



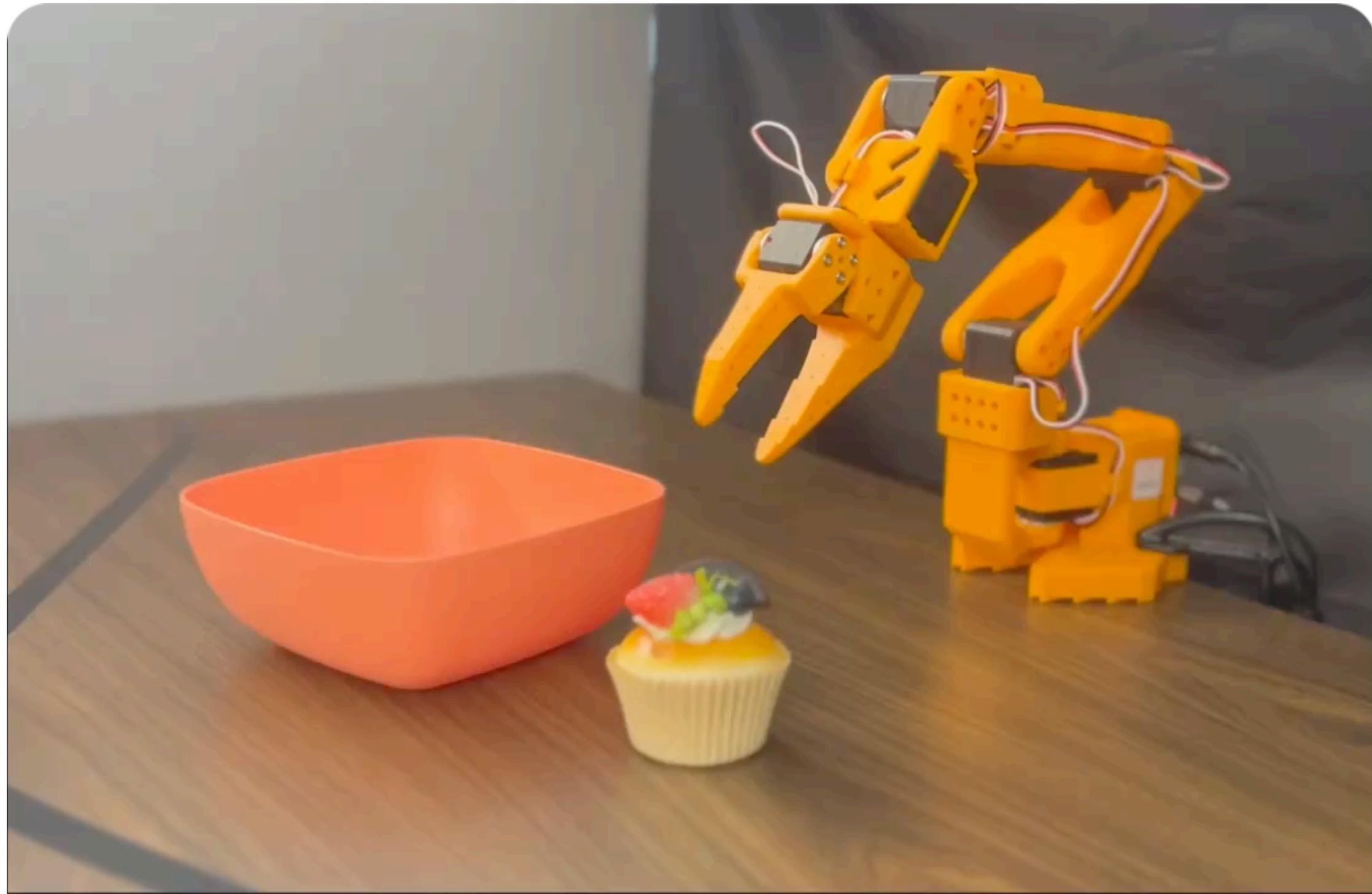
VLA-0: Building State-of-the-Art VLAs with Zero Modification



VLA-0

No Modification to the VLM

- VLA with no modification to the VLM
- Predict Action as Text
- No change to tokenization
- No new architectural component
- Training recipe is important - Check the paper for details



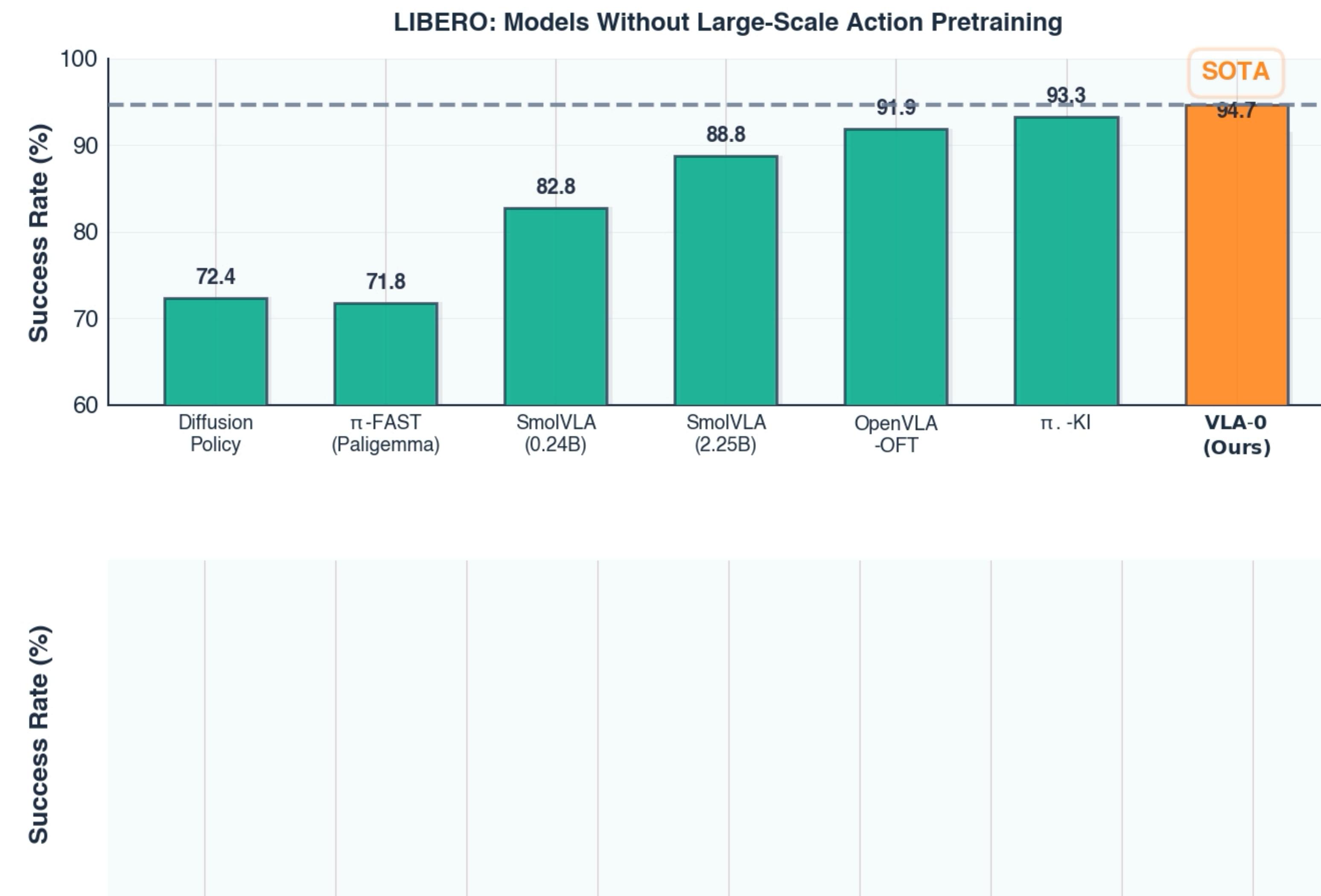
VLA-0

Results on Libero

- #1 among all non-pretrained architectures — Outperforming Pi-FAST, Pi-0.5-KI, SmoVLA, OpenVLA-OFT

Even more surprising:

- Without any action pretraining, outperforms leading pretrained models like Pi-0, Pi-0.5-K, GR00T-N1, MomlmoAct



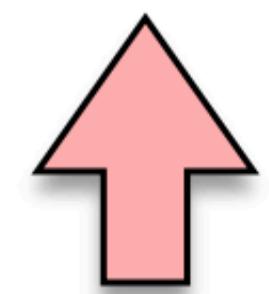
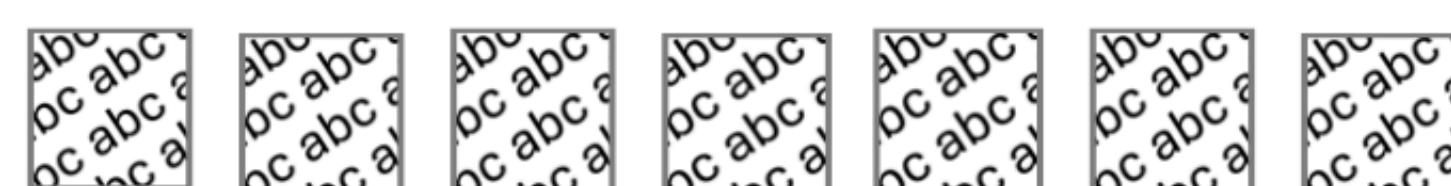
VLA-0

Output Action:

4 12 98 3 0 0 13 5 123 23 0 0 24 0 132 34 13 0 ...

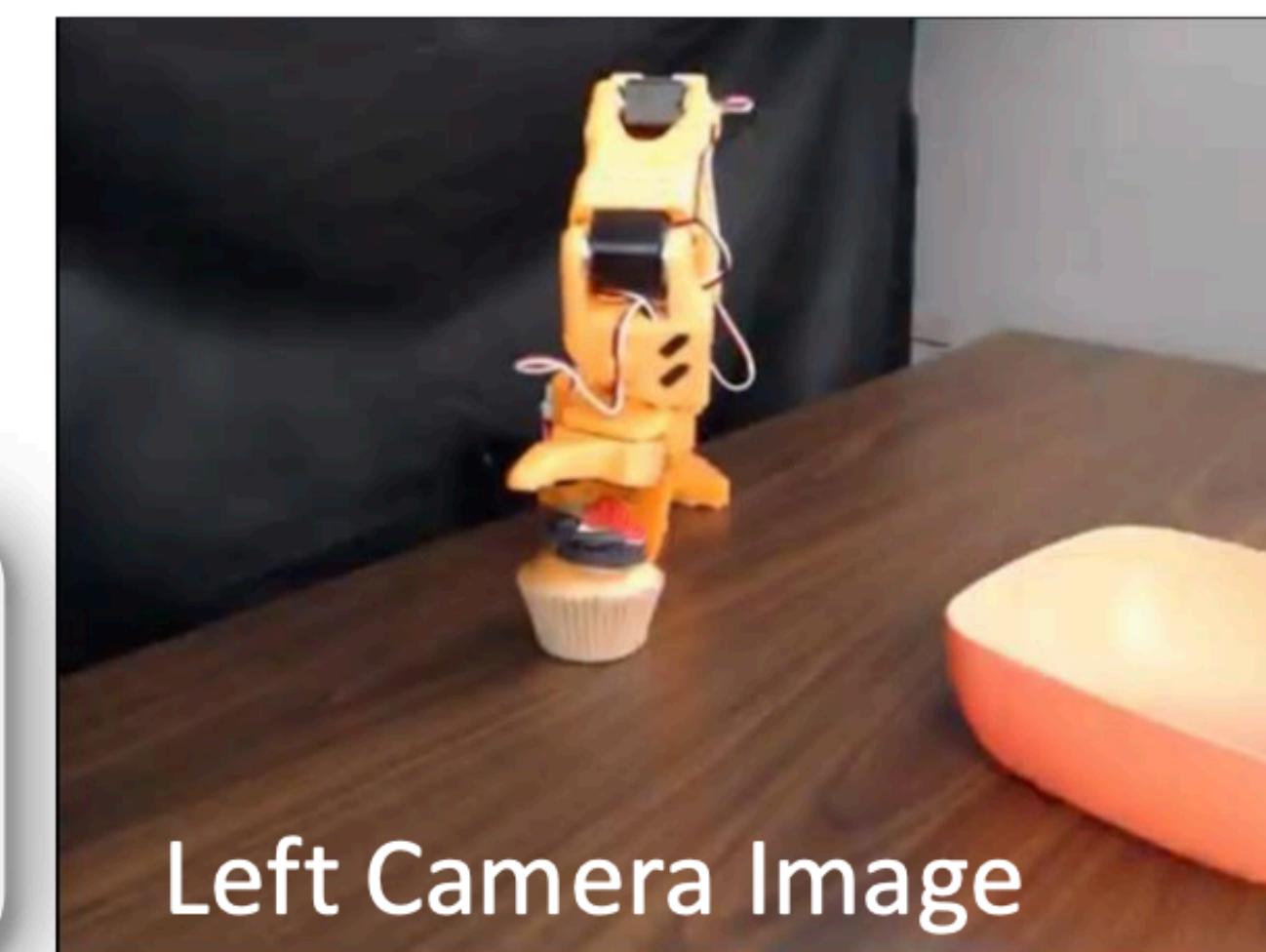


Vision Language Model

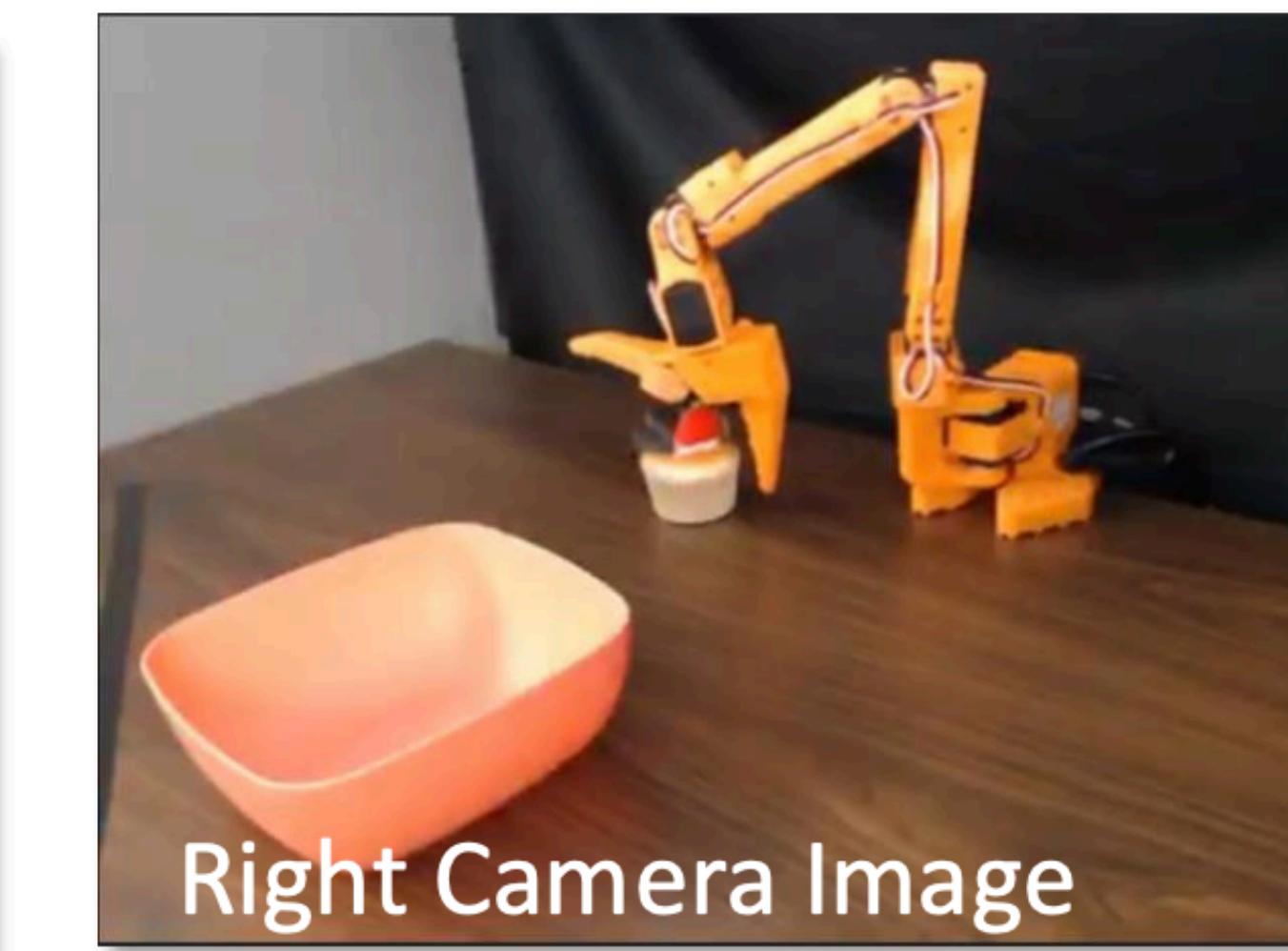


System Prompt

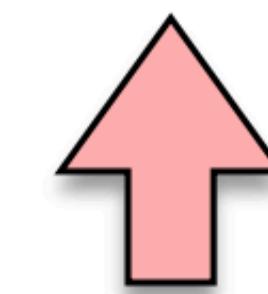
Analyze the input image and predict
robot actions for the next H
timesteps



Left Camera Image



Right Camera Image



User Prompt

Put the cupcake
in the bowl

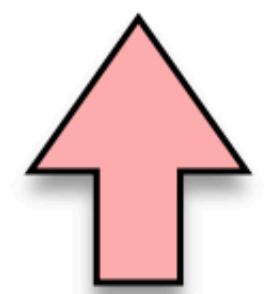
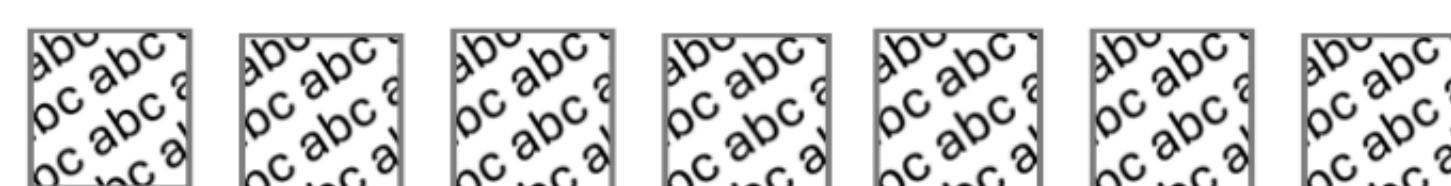
VLA-0

Output Action:

4 12 98 3 0 0 13 5 123 23 0 0 24 0 132 34 13 0 ...

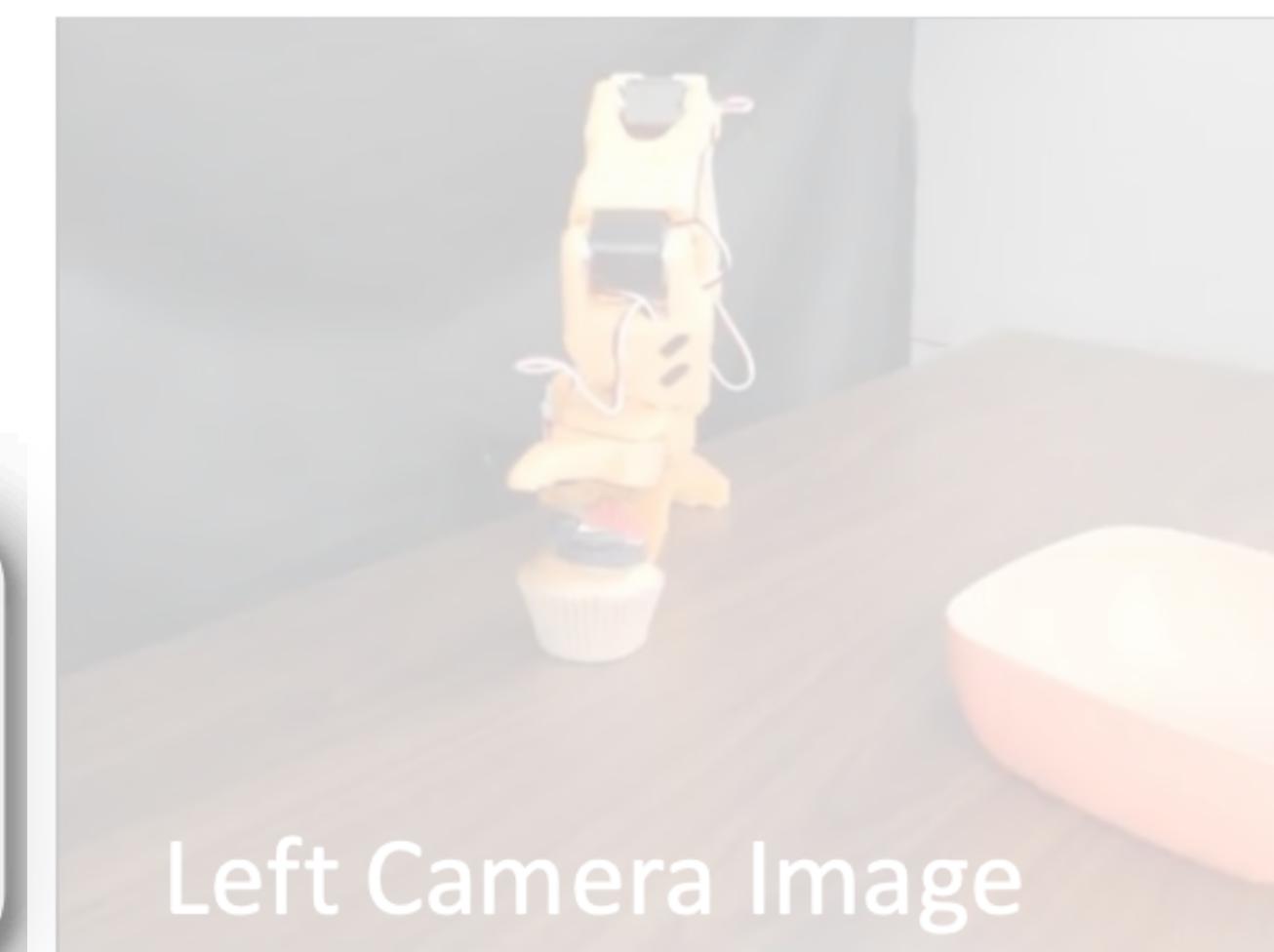


Vision Language Model



System Prompt

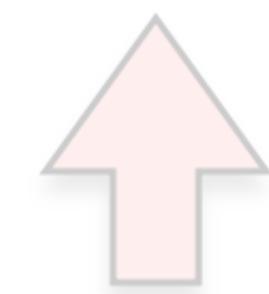
Analyze the input image and predict robot actions for the next H timesteps



Left Camera Image



Right Camera Image



User Prompt

Put the cupcake
in the bowl

VLA-0

Output Action:

4 12 98 3 0 0 13 5 123 23 0 0 24 0 132 34 13 0 ...

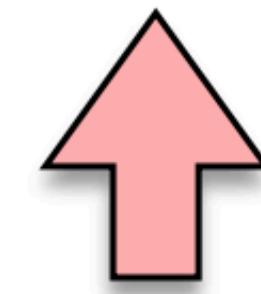
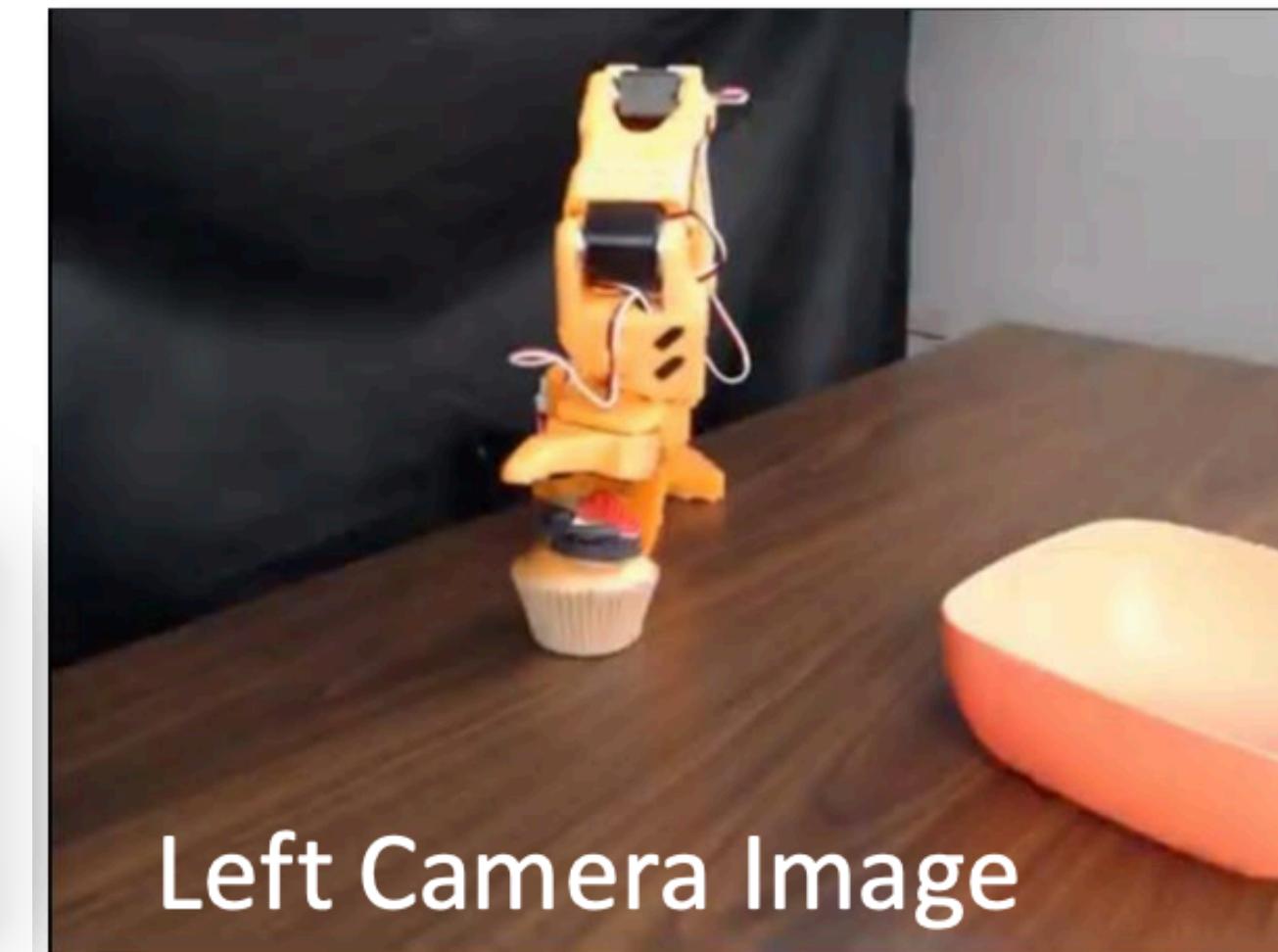


Vision Language Model



System Prompt

Analyze the input image and predict
robot actions for the next H
timesteps



User Prompt

Put the cupcake
in the bowl

VLA-0

Output Action:

4 12 98 3 0 0 13 5 123 23 0 0 24 0 132 34 13 0 ...

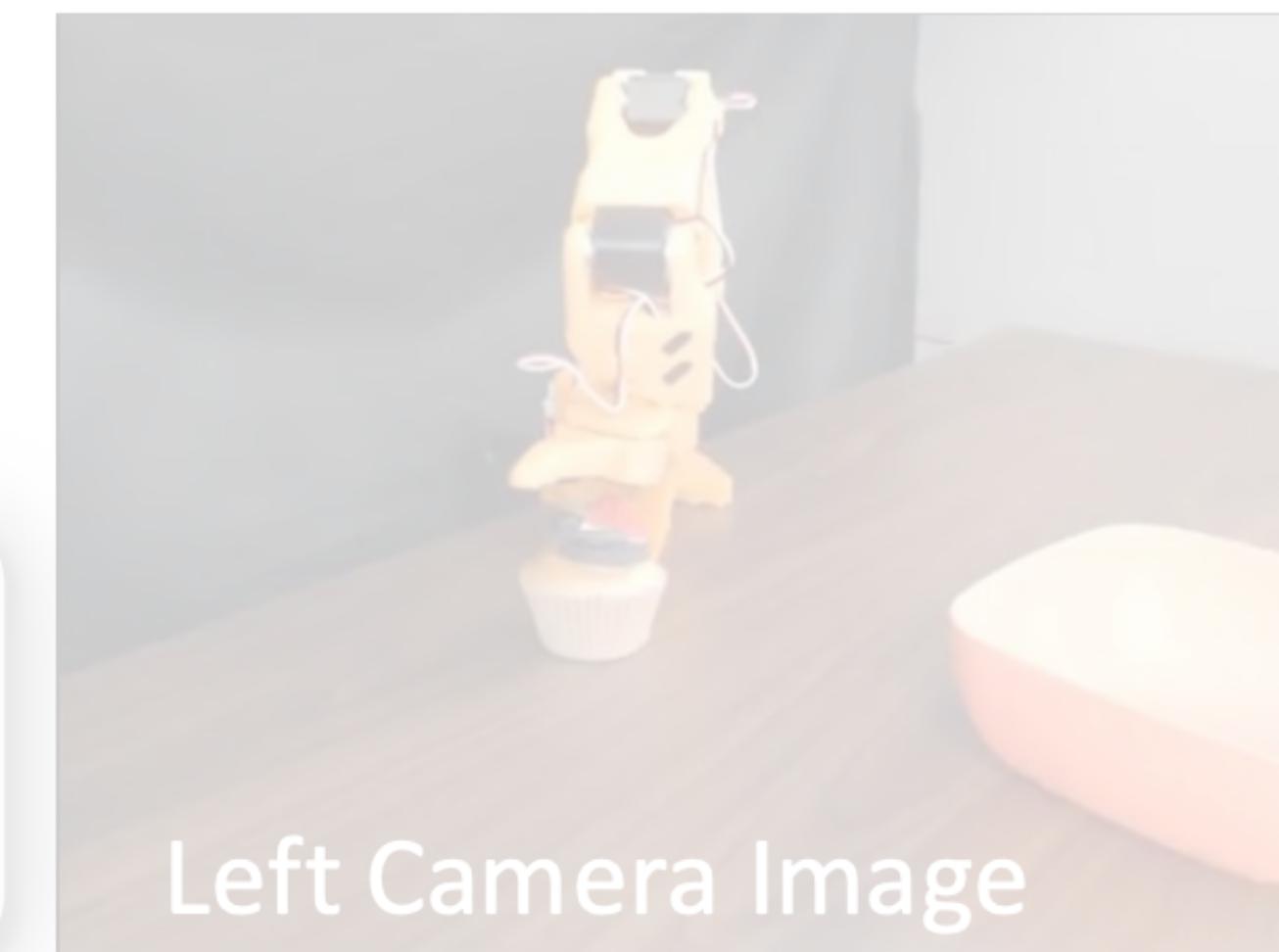


Vision Language Model

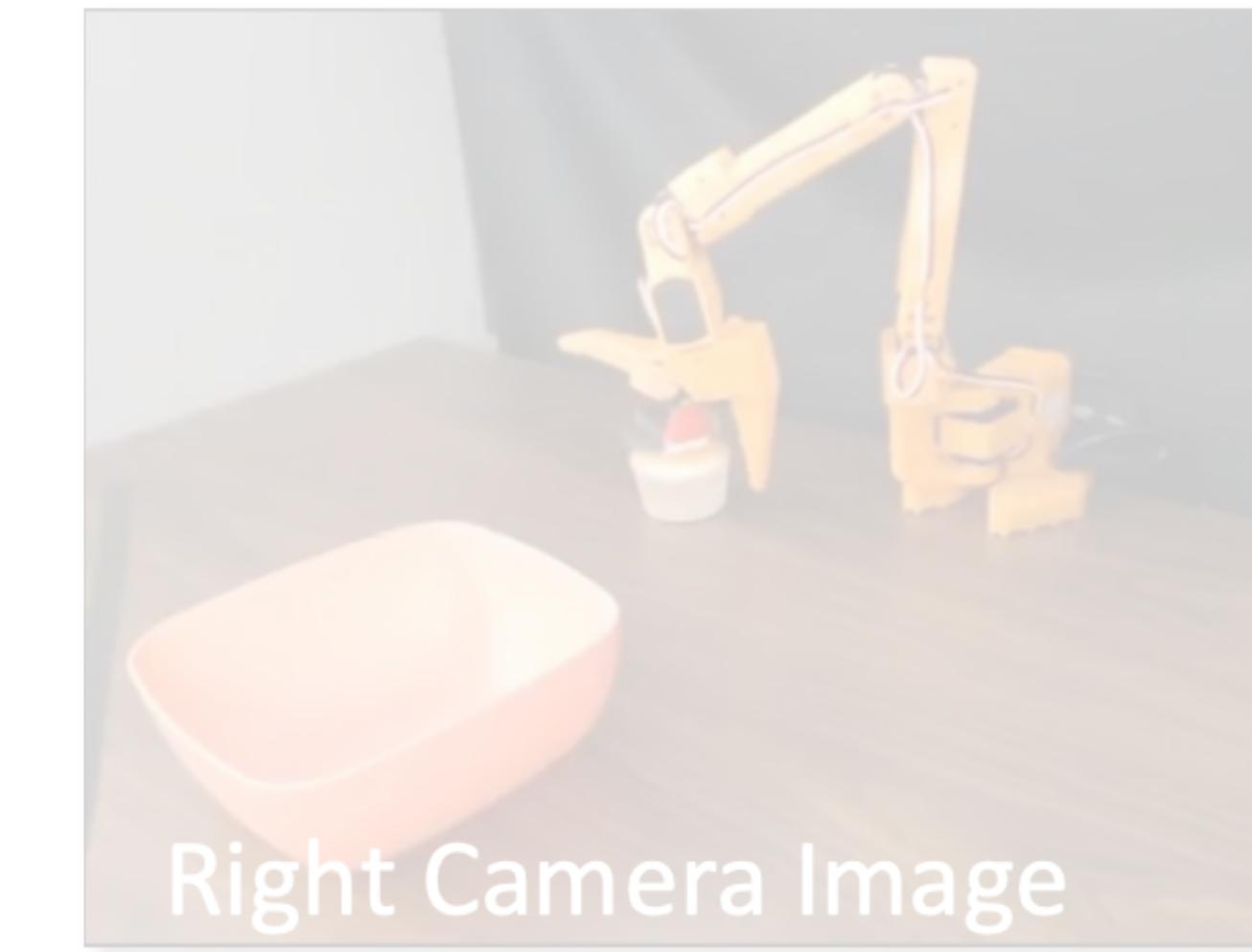


System Prompt

Analyze the input image and predict
robot actions for the next H
timesteps



Left Camera Image



Right Camera Image

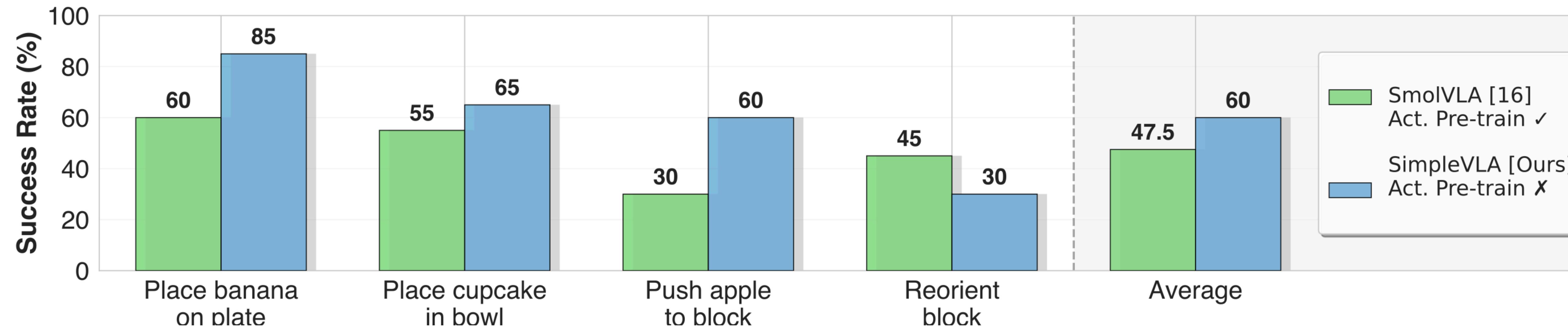


User Prompt

Put the cupcake
in the bowl

VLA-0

Results in Real



- On real world data, outperforms SmoVLA
- SmoVLA — Pretrained on large-scale SO-100 data + Finetuned on 100 demos per task
- VLA-0 trained from scratch with 100 demos per task

VLA-0

vla0.github.io



Fig. 1: Schematic representation of VLA-0. VLA-0 converts a VLM into a VLA by prompting the VLM to predict action as text. This strategy is surprisingly effective and achieves state-of-the-art results akin to alternatives.



Summary

Thank you! Questions?

- Hierarchical VLAs – Combine the strength of VLMs and 3D Policies
- 3D Policies – Efficient task specific learners
- Among monolithic VLAs – Simplest design (VLA-0) is surprisingly effective

