# Robotic Foundation Models
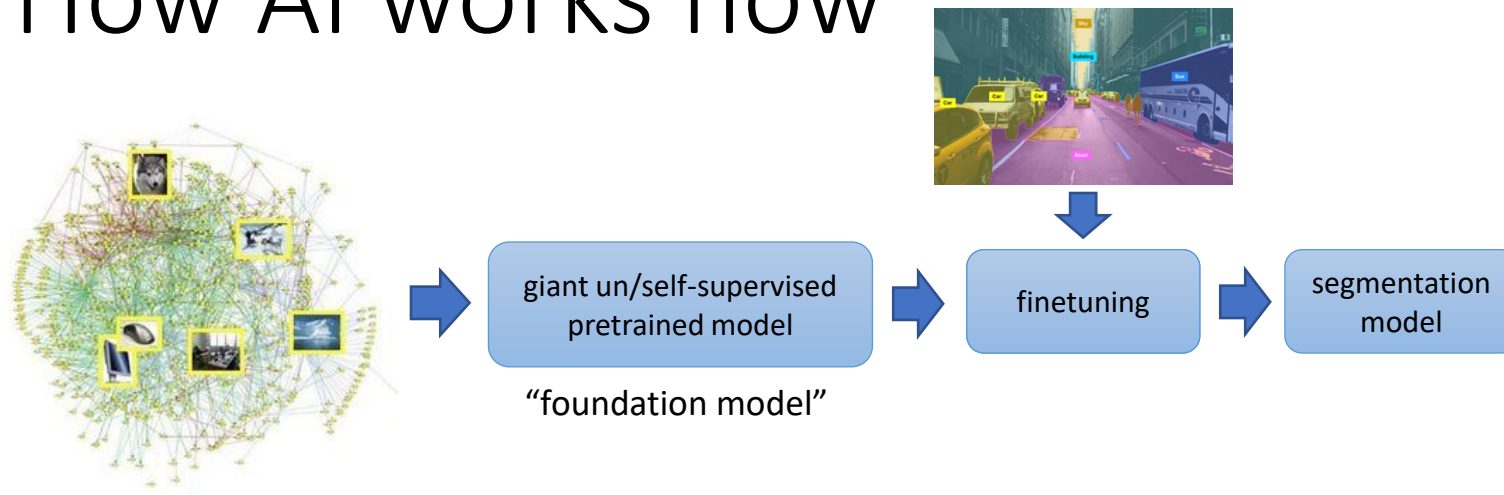
**Sergey Levine**

**UC Berkeley**

**Physical Intelligence**

# How AI used to work



segmentation model

classification model

captioning model

visual QA model

sentiment model

summarization model

**container ship**
container ship
lifeboat
amphibian
fireboat
drilling platform

**A group of people shopping at an outdoor market.**

What is the mustache made of?

bananas

"Horrible services. The room was dirty and unpleasant. Not worth the money."

NEGATIVE

Natural Language Processing

Natural language processing (NLP) is a subfield of linguistics, computer science, and artificial intelligence concerned with the interactions between computers and human language, in particular how to program computers to process and analyze large amounts of natural language data. The result is a computer capable of "understanding" the contents of documents, including the contextual nuances of the language within them. The technology can then accurately extract information and insights contained in the documents as well as categorize and organize the documents themselves.

Natural Language Processing

Natural language processing (NLP) is a subfield of linguistics, computer science, and artificial intelligence concerned with the interactions between computers and human language, in particular how to program computers to process and analyze large amounts of natural language data.

# How AI works now

giant un/self-supervised pretrained model

finetuning

segmentation model

"foundation model"

2

# How robotic learning works now



PR2 pancake model

WAM pancake model

UR10 box picking model

lipstick robot model

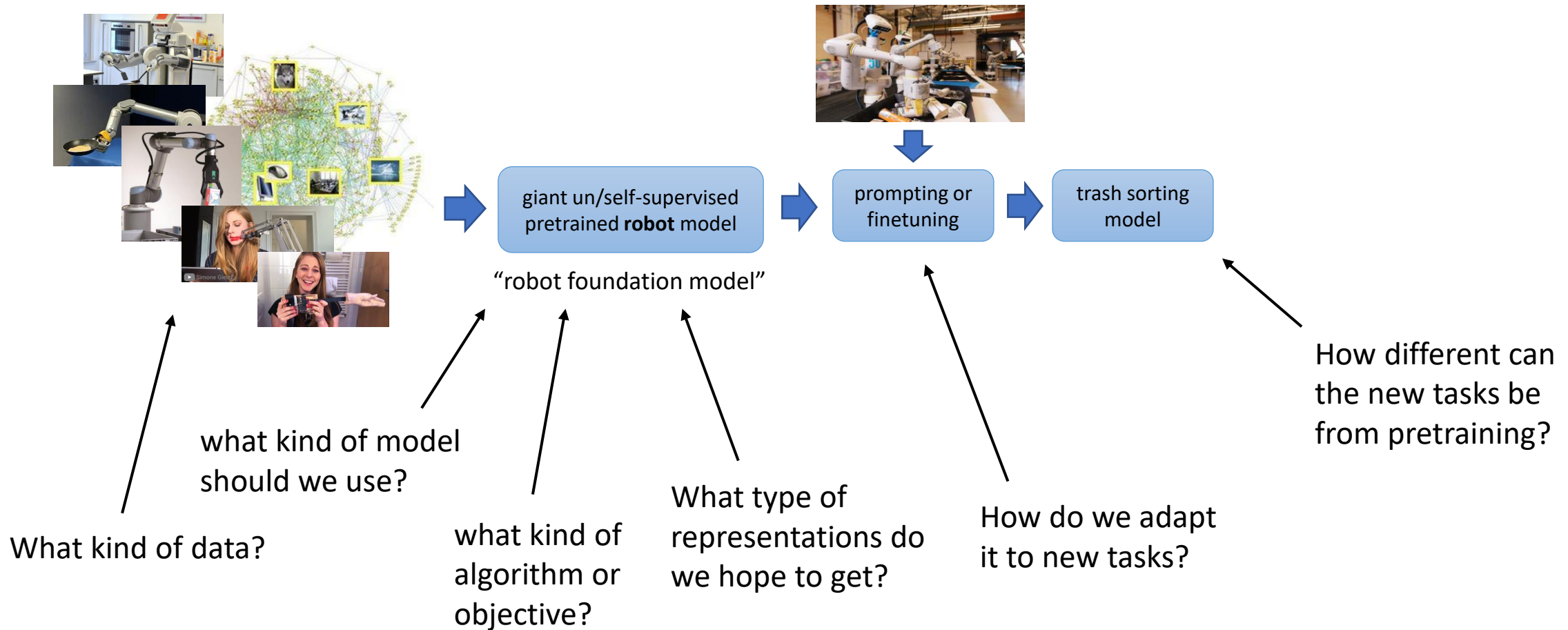mousetrap... hand... shake... model??

Image credit: Simone Giertz

# How robotic learning will work in the future



giant un/self-supervised pretrained **robot** model

"robot foundation model"

prompting or finetuning

trash sorting model

# What do we need to figure out?



giant un/self-supervised pretrained **robot** model

"robot foundation model"

prompting or finetuning

trash sorting model

How different can the new tasks be from pretraining?

what kind of model should we use?

What kind of data?

what kind of algorithm or objective?

What type of representations do we hope to get?

How do we adapt it to new tasks?

# How do we build robotic foundation models?



Robotic foundation models for navigation



Manipulation, VLAs, and open-source models



Taking cross-embodied learning to the limit

# How do we build robotic foundation models?

 Robotic foundation models for navigation

 Manipulation, VLAs, and open-source models

 Taking cross-embodied learning to the limit

# Robotic foundation models for navigation

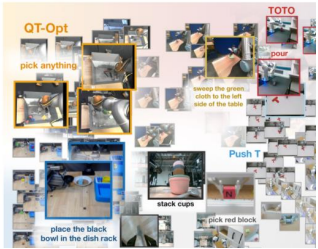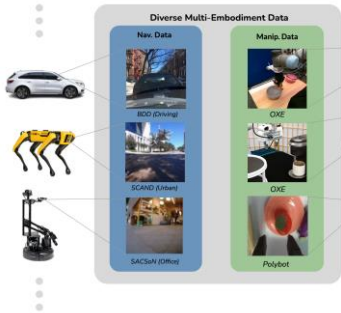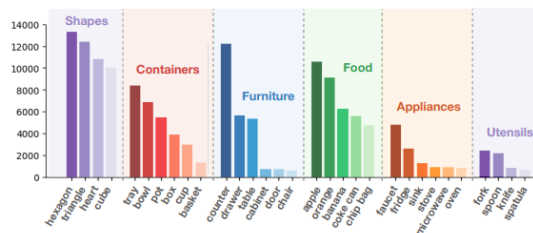| | Dataset | Platform | Speed | Amt. | Environment |
|---|---|---|---|---|---|
| 1 | GoStanford [26] | TurtleBot2 | 0.5m/s | 14h | office |
| 2 | RECON [32] | Jackal | 1m/s | 25h | off-road |
| 3 | CoryHall [35] | RC Car | 1.2m/s | 2h | hallways |
| 4 | Berkeley [33] | Jackal | 2m/s | 4h | suburban |
| 5 | SCAND-S [36] | Spot | 1.5m/s | 8h | sidewalks |
| 6 | SCAND-J [36] | Jackal | 2m/s | 1h | sidewalks |
| 7 | Seattle [37] | Warthog | 5m/s | 1h | off-road |
| 8 | TartanDrive [38] | ATV | 10m/s | 5h | off-road |
| | Ours | | | 60h | |



RC-Car
*(Kahn et al. 2018)*

TurtleBot
*(Hirose et al. 2019)*

Jackal
*(Shah et al. 2021, 2022)*

Spot
*(Karnan et al. 2022)*

Warthog
*(Shaban et al. 2021)*

ATV
*(Triest et al. 2022)*

Shah*, Sridhar*, Bhorkar, Hirose, Levine. **GNM: A General Navigation Model to Drive Any Robot**. 2022.

# Scaling it up with Transformers



## ViNT: Visual Navigation Transformer

https://general-navigation-models.github.io/

ViNT Foundation Model

Pre-trained ViNT

Shah, Sridhar, Dashora, Stachowicz, Black, Hirose, Levine. **ViNT: A Foundation Model for Visual Navigation.** 2023.

# Now make it go fast!



Stage 1: Offline Learning

Prior Dataset → Encoder → Critic
IQL Training

Stage 2: Online Learning

Encoder (Frozen) → Actor / Critic
RLPD Training

5 min

All Videos at 1x (Real-Time)

10

Stachowicz, Shah, Bhorkar, Kostrikov, Levine. **FastRLAP: A System for Learning High-Speed Driving via Deep RL and Autonomous Practicing.** 2023.

All Videos at 1x (Real-Time)

All Videos at 1x (Real-Time)

25 min

# How do we build robotic foundation models?


Robotic foundation models for navigation


Manipulation, VLAs, and open-source models
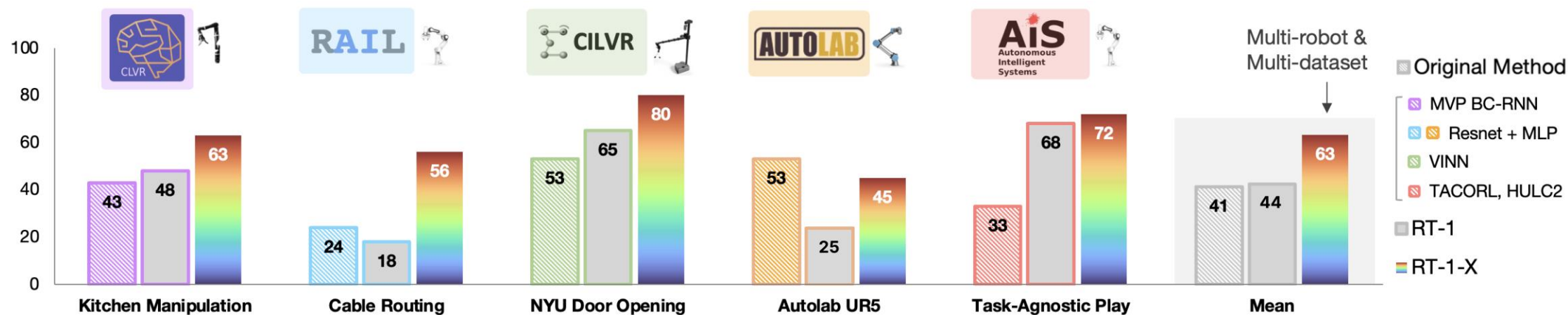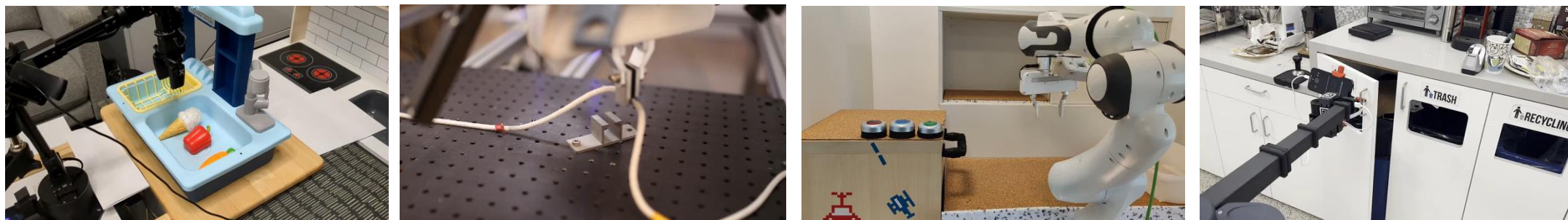

Taking cross-embodied learning to the limit

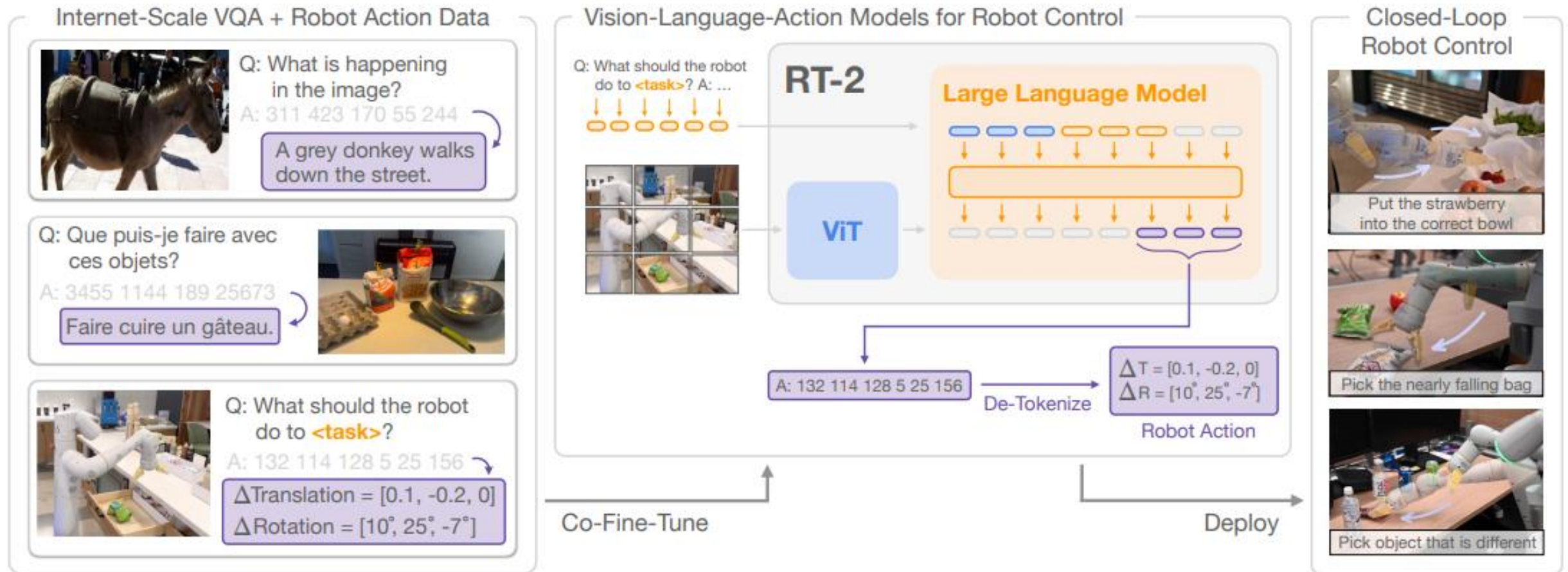# RT-X: Combining many datasets for cross-embodiment manipulation

**Visualizations:** https://dibyaghosh.com/rtx_viz/
By Dibya Ghosh

17

# RT-X: Combining many datasets for cross-embodiment manipulation





RT-X: https://robotics-transformer-x.github.io/
Open X-Embodiment Collaboration

# RT-2: A **vision-language-action (VLA)** model

# RT-2: A vision-language-action (VLA) model



put strawberry into the correct bowl

pick up the bag about to fall off the table

move apple to Denver Nuggets

move redbull can to H

move soccer ball to basketball

move banana to Germany

move coke can to Taylor Swift

move coke can to X

move bag to Google

Prompt:
Given <img> Instruction: Bring me a drink.
Prediction:
Plan: pick 7up can.
Action: 1 143 129 123 145 114 115 127

Prompt:
Given <img> Instruction: Move all the objects together.
Prediction:
Plan: move green can near green rice chip bag.
Action: 1 128 126 127 135 123 119 127

Prompt:
Given <img> Instruction: Pick the object that is different from all other objects
Prediction:
Plan: pick rxbar chocolate. Action: 1 128 129 125 131 125 128 127

Prompt:
Given <img> Instruction: Move the green objects together.
Prediction:
Plan: move green can near green rice chip bag.
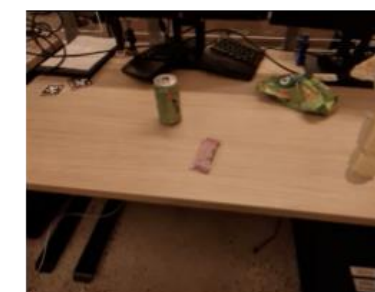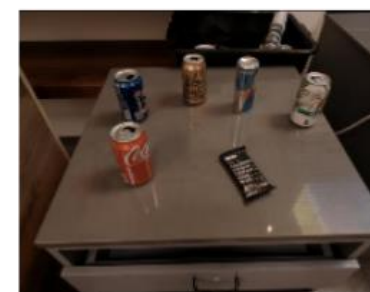Action: 1 130 129 121 131 127 128 127

Prompt:
Given <img> I need to hammer a nail, what object from the scene might be useful?
Prediction:
Rocks. Action: 1 129 138 122 132 135 106 127

# RT-2-X: Does cross embodiment training help VLAs?
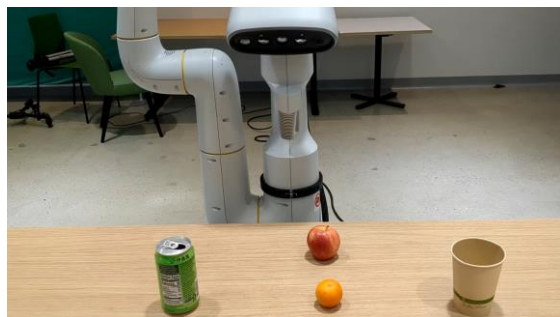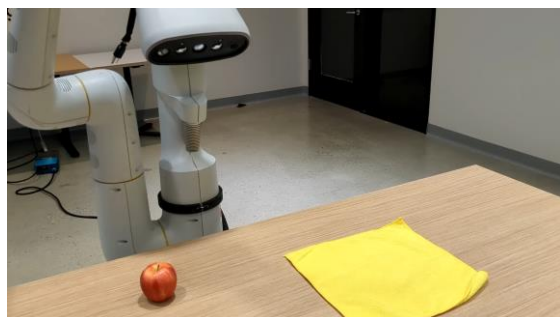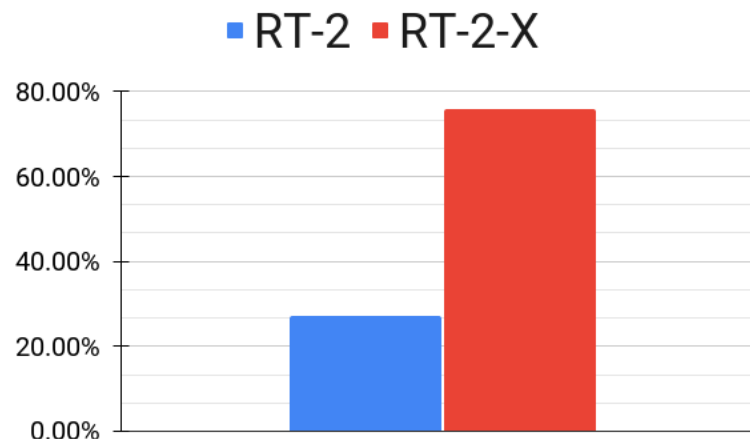


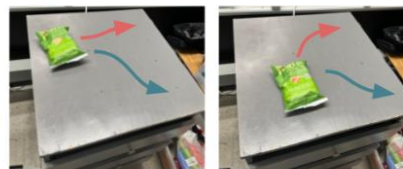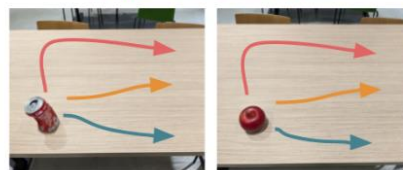"move apple between can & orange"



"move apple near cloth"



"move apple on cloth"



RT-2  RT-2-X



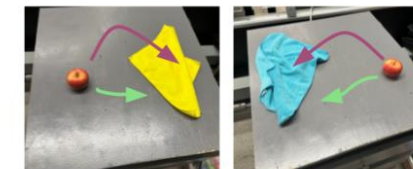(a) Absolute Motion
move the chip bag to the top / bottom right of the counter

move to top right / right / bottom right

(b) Object-Relative Motion
move apple between coke and cup / coke and sponge / cup and sponge
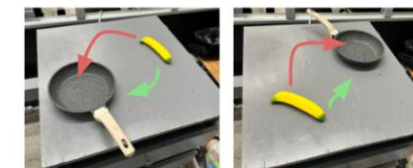
(c) Preposition Alters Behavior
put apple on cloth / move apple near cloth

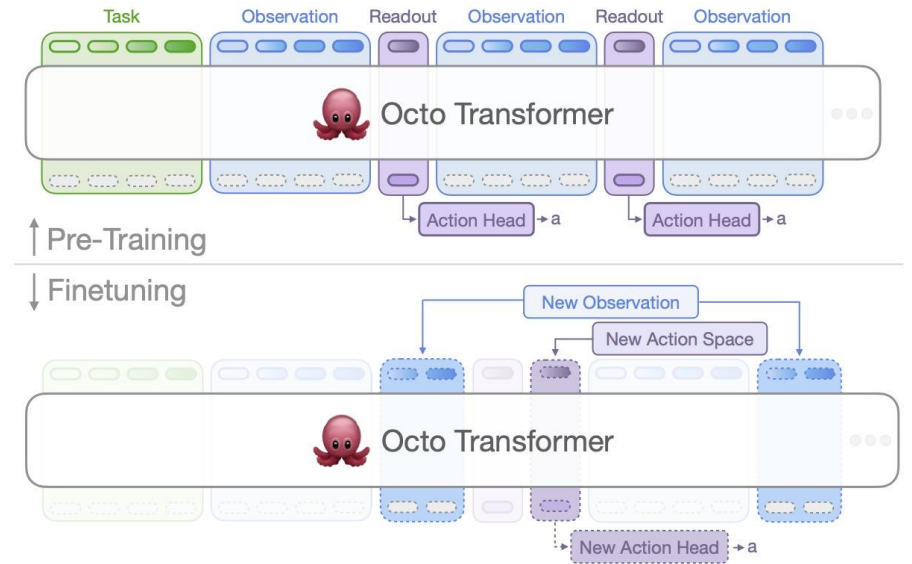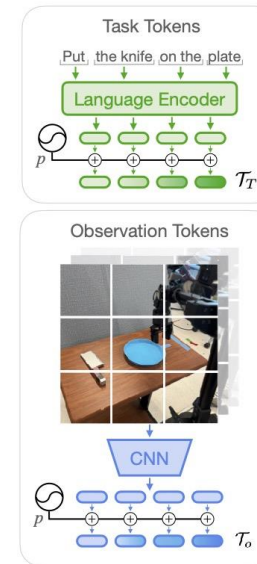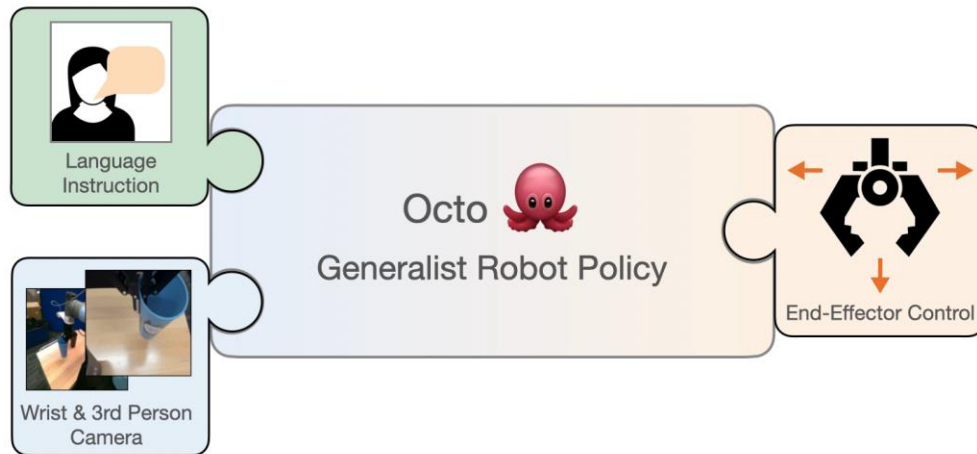put orange into the pot / move orange near pot

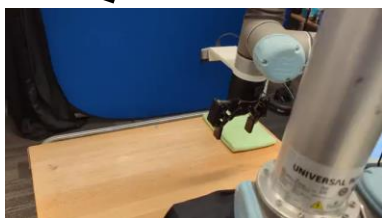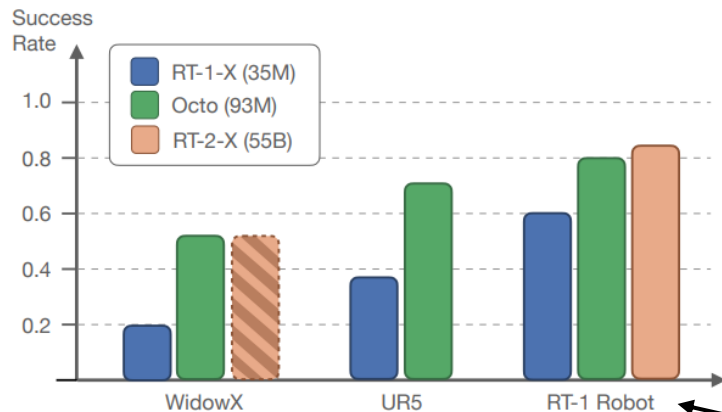put banana on top of the pan / move banana near pan
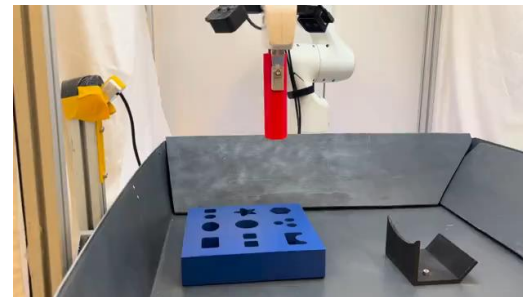
# Octo: an open-source robotic foundation model

# Octo: an open-source robotic foundation model

**zero-shot evaluation**



Berkeley Insertion

Stanford Coffee

CMU Baking

Berkeley Bimanual

**finetuning**

| | Berkeley Insertion* | Stanford Coffee | CMU Baking | Berkeley Pick-Up[†] | Berkeley Coke | Berkeley Bimanual[†] | Average |
|---|---|---|---|---|---|---|---|
| ResNet+Transformer Scratch | 10% | 45% | 25% | 0% | 20% | 20% | 20% |
| VC-1 [57] | 5% | 0% | 30% | 0% | 10% | 50% | 15% |
| Octo (Ours) | **70%** | **75%** | **50%** | **60%** | **100%** | **80%** | **72%** |

# OpenVLA: an open-source **vision-language-action** model

**OpenVLA:** https://openvla.github.io/
Kim*, Pertsch*, Karamcheti*, et al.

# How do we build robotic foundation models?



Robotic foundation models for navigation



Manipulation, VLAs, and open-source models



Taking cross-embodied learning to the limit

# How diverse can the data be?

# An "extreme" cross-embodiment recipe

Yang, Glossop, Bhorkar, Shah, Vuong, Finn, Sadigh, Levine. **Extreme Cross-Embodiment Learning for Manipulation and Navigation.** 2024.

# Why might this work?

Yang, Glossop, Bhorkar, Shah, Vuong, Finn, Sadigh, Levine. **Extreme Cross-Embodiment Learning for Manipulation and Navigation.** 2024.

# Some results

## Does navigation help with manipulation?





Legend: ■ Manip. ■ GNM + Manip. ■ GNM + Driving + Manip.

- Average: 51, 64, 71
- Two-Object Grasp: 70, 80, 80
- Cluttered Grasp: 65, 75, 80
- Toy Kitchen: 70, 65, 80
- Shelf Manipulation: 30, 50, 65
- Novel Cluttered Grasp: 20, 50, 50

Yang, Glossop, Bhorkar, Shah, Vuong, Finn, Sadigh, Levine. **Extreme Cross-Embodiment Learning for Manipulation and Navigation.** 2024.

# Some results

## Does manipulation help with navigation?





Yang, Glossop, Bhorkar, Shah, Vuong, Finn, Sadigh, Levine. **Extreme Cross-Embodiment Learning for Manipulation and Navigation.** 2024.

# Can we make it even more "extreme"?

Doshi, Walke, Mees, Dasari, Levine. **Scaling Cross-Embodied Learning: One Policy for Manipulation, Navigation, Locomotion and Aviation.** 2024.

# The CrossFormer architecture

Doshi, Walke, Mees, Dasari, Levine. **Scaling Cross-Embodied Learning: One Policy for Manipulation, Navigation, Locomotion and Aviation.** 2024.

# How diverse do the embodiments get?



- Robotic manipulation matches prior robotic foundation models (e.g., Octo)
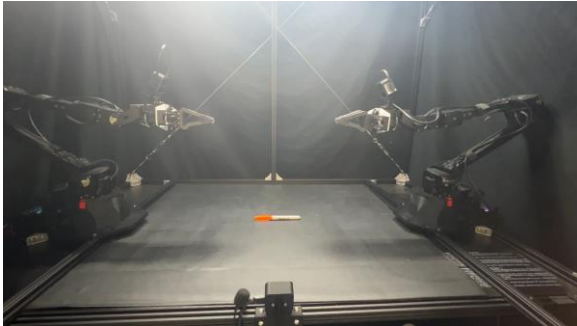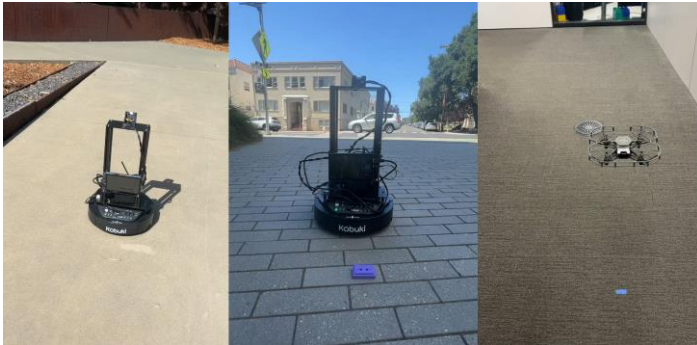- Can use **either** third person or wrist-mounted cameras



High-frequency bimanual manipulation (50 Hz) matches dedicated bimanual models
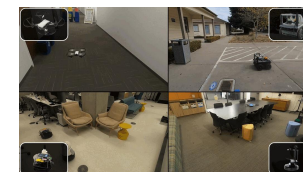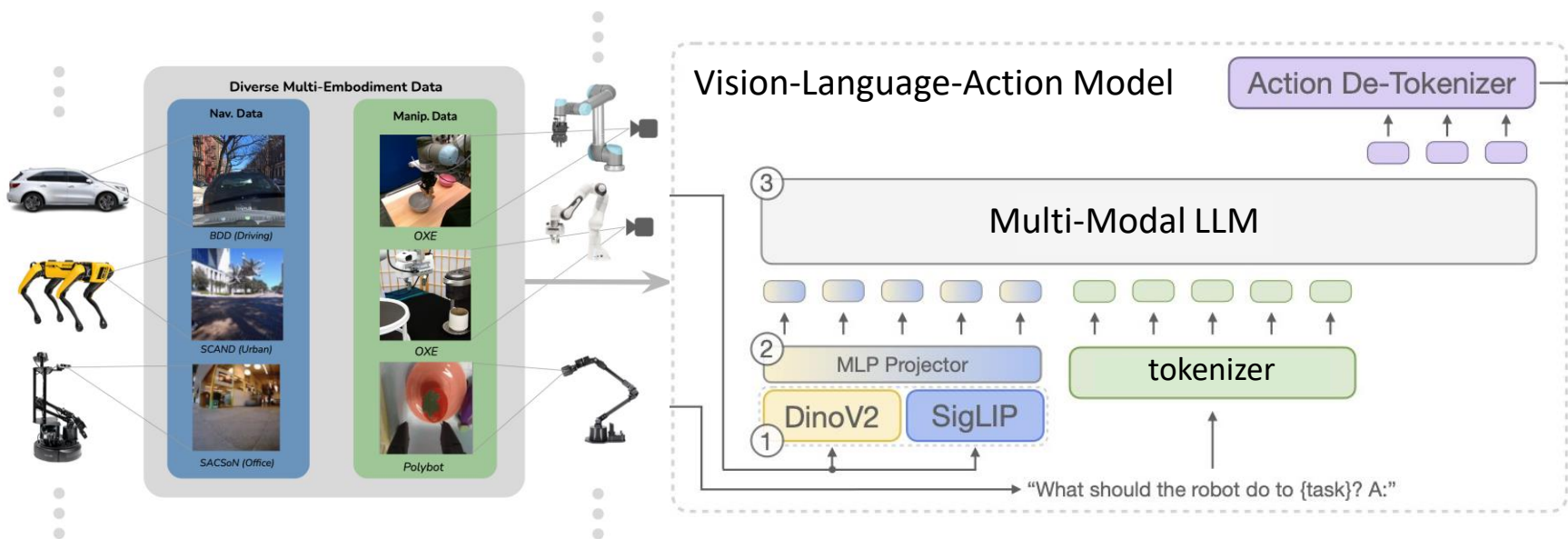


Integrates with topological graphs for long-horizon navigation (ground robots & quadcopters)



Same model performs low-level joint control for a quadruped

Doshi, Walke, Mees, Dasari, Levine. **Scaling Cross-Embodied Learning: One Policy for Manipulation, Navigation, Locomotion and Aviation.** 2024.

# Summary



Diverse Multi-Embodiment Data

Nav. Data — BDD (Driving), SCAND (Urban), SACSoN (Office)

Manip. Data — OXE, OXE, Polybot

Vision-Language-Action Model

Action De-Tokenizer

Multi-Modal LLM

MLP Projector

tokenizer

DinoV2   SigLIP

"What should the robot do to {task}? A:"

zero shot to new robots

few shot to new tasks

Turn Left @ Intersection   Continue Straight   Turn Right @ Intersection

Stage 1: Offline Learning — Prior Dataset — Encoder — Critic — IQL Training

Stage 2: Online Learning — Encoder (Frozen) — Actor — Critic — RLPD Training
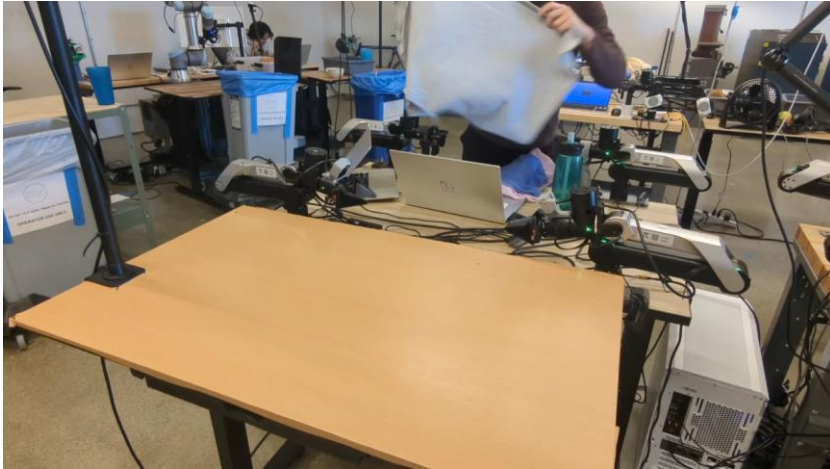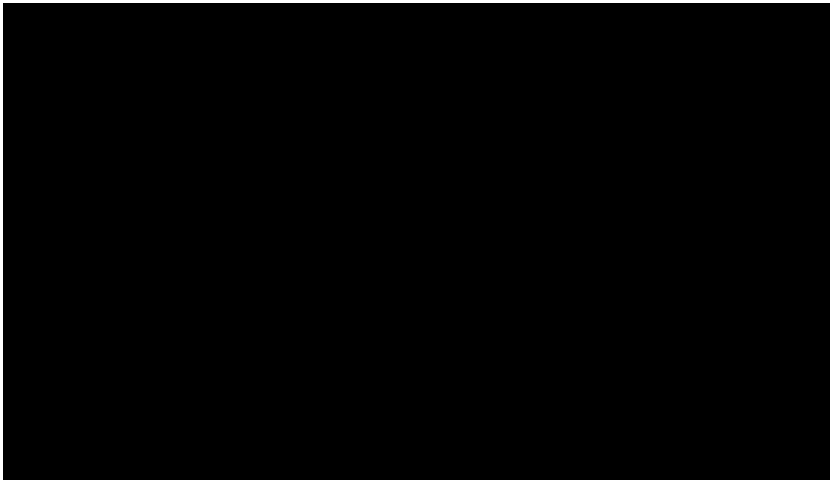
pretrain for super fast online RL

utilize for downstream instruction following

# $\pi$ Physical Intelligence



Can we scale up robotic foundation models to tackle the breadth of real-world tasks and robotic platforms?

http://physicalintelligence.company

**RAIL**
Robotic AI & Learning Lab