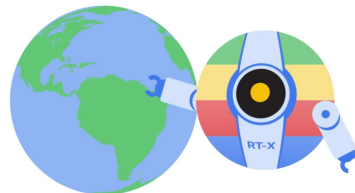
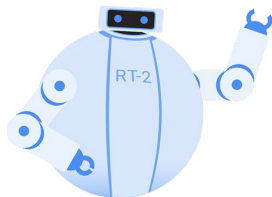
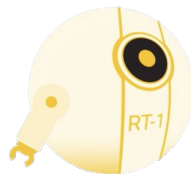


# What's **Missing** for Robotics-First Foundation Models?

Ted Xiao

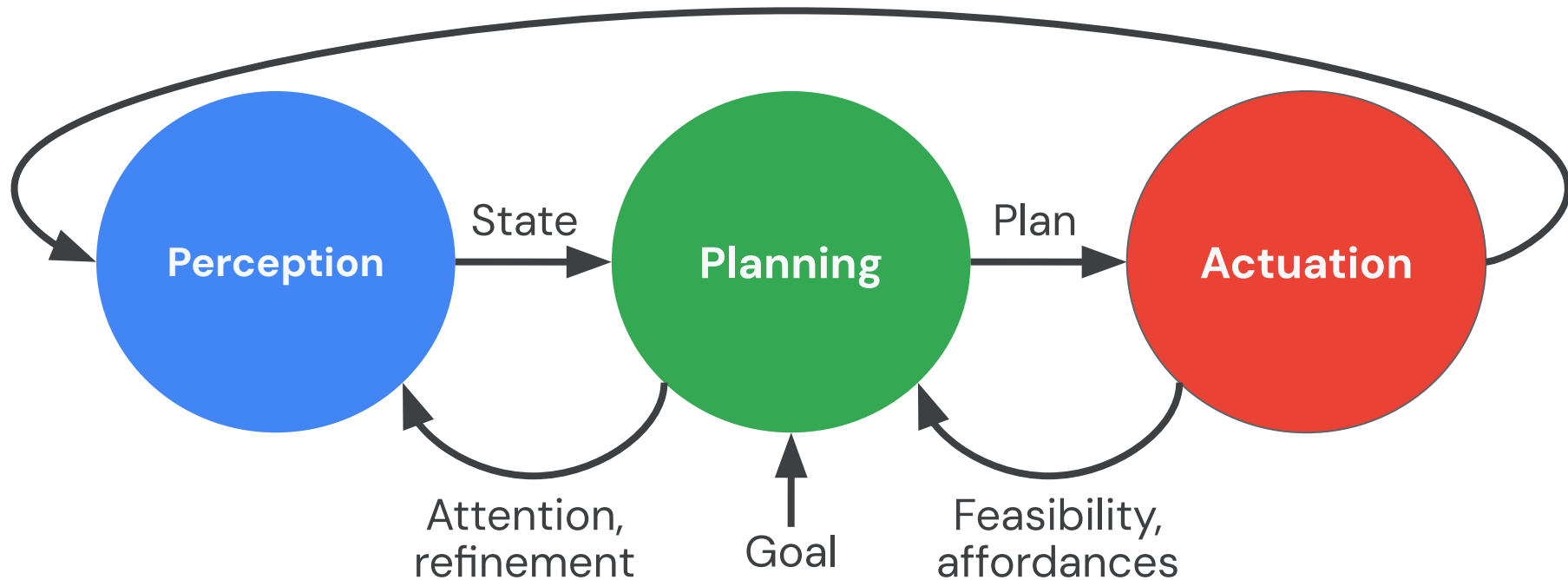




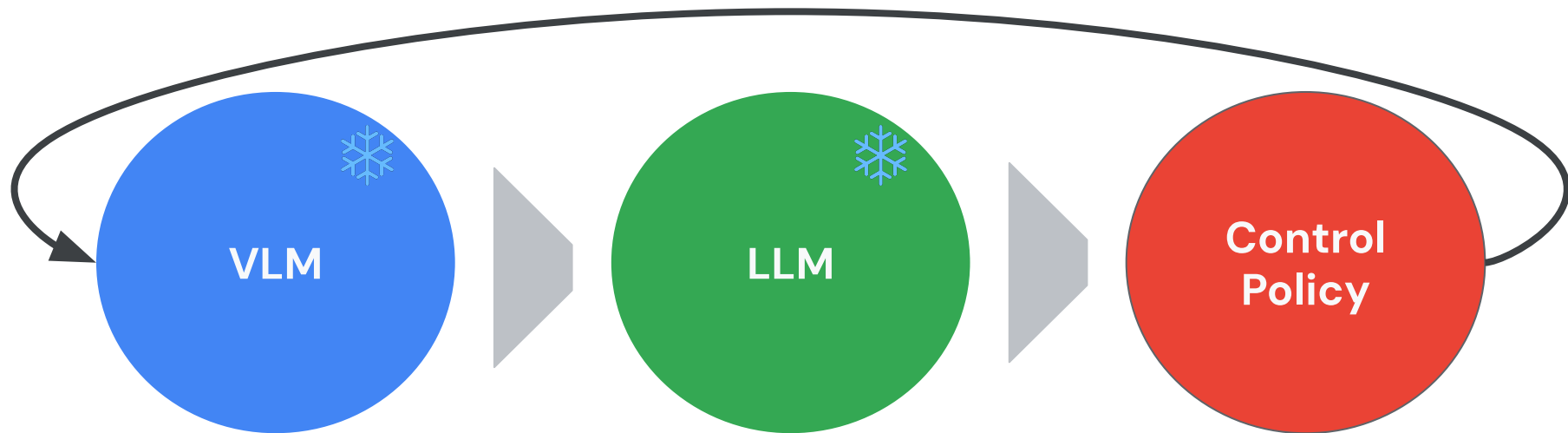
# Agenda

- 01** Why Robot Foundation Models?
- 02 Piece #1: Positive Transfer from Scaling
- 03 Piece #2: Steerability
- 04 Piece #3: Scalable Evaluation
- 05 Horizons

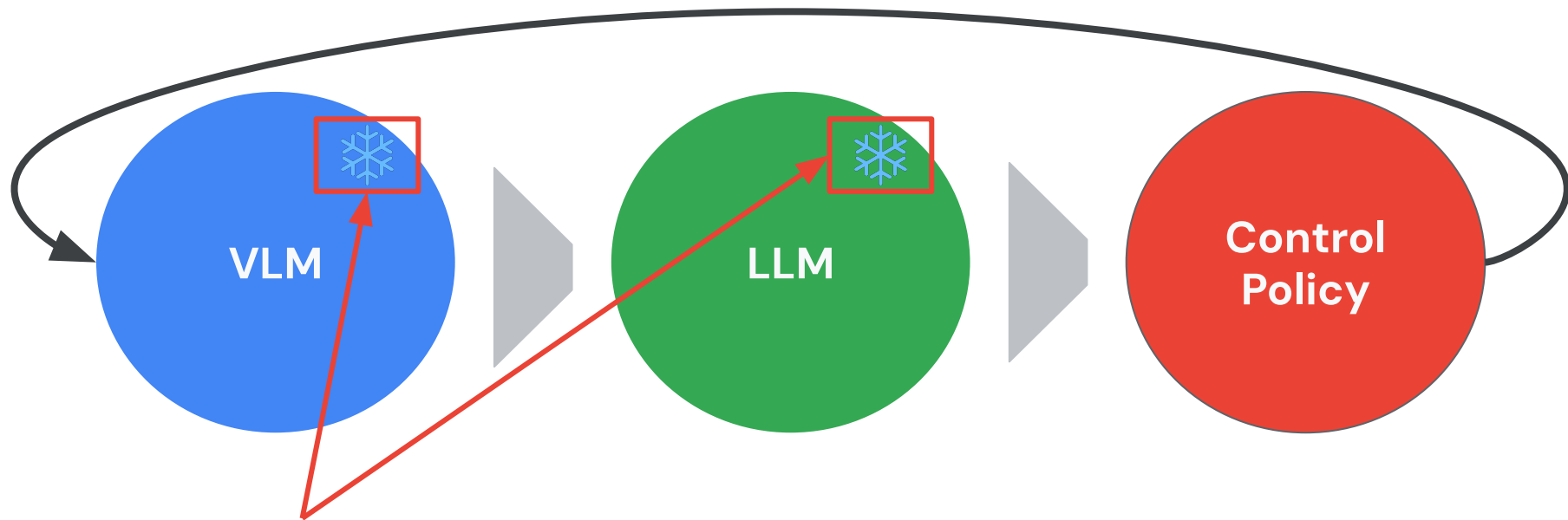
# The Robotics Information Flow



# Foundation Models as Experts

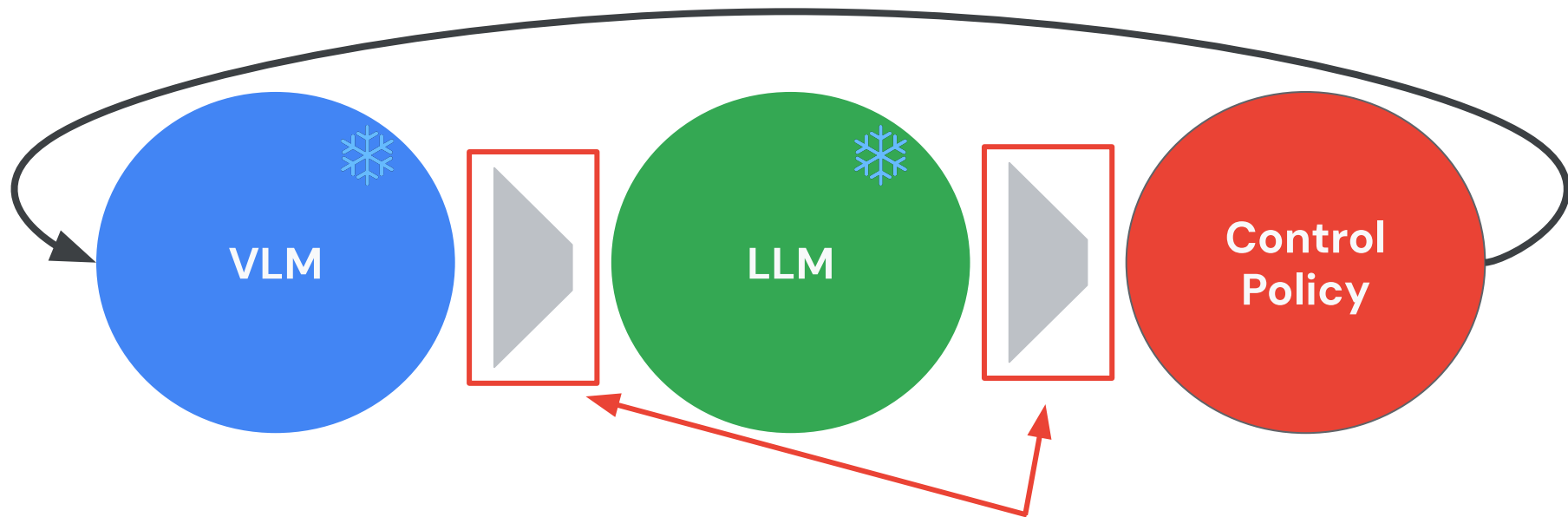


# Foundation Models as Experts



**Issue #1: Not optimized for robotics**

# Foundation Models as Experts



**Issue #1: Not optimized for robotics**

**Issue #2: Narrow communication bandwidth between "intelligence modules"**

# Foundation Model-fication of Robotics?

Sentiment  
Classification

Translation

Summarization

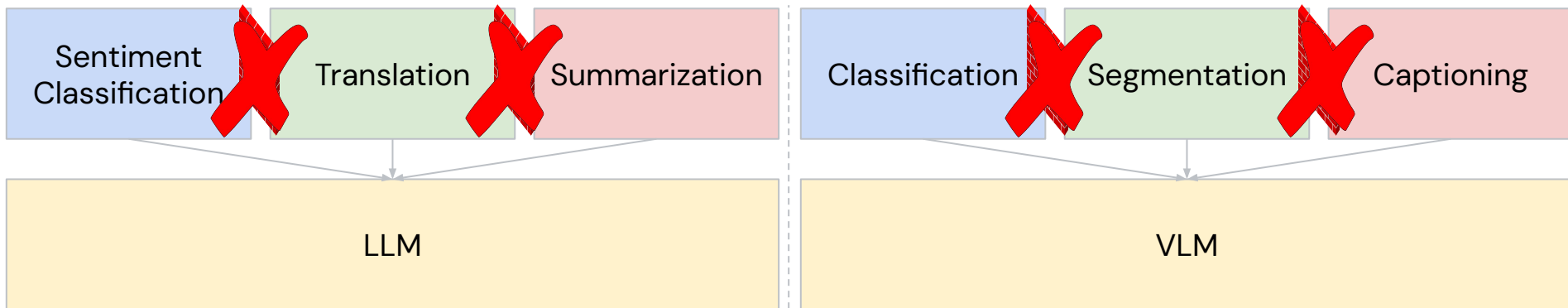
Classification

Segmentation

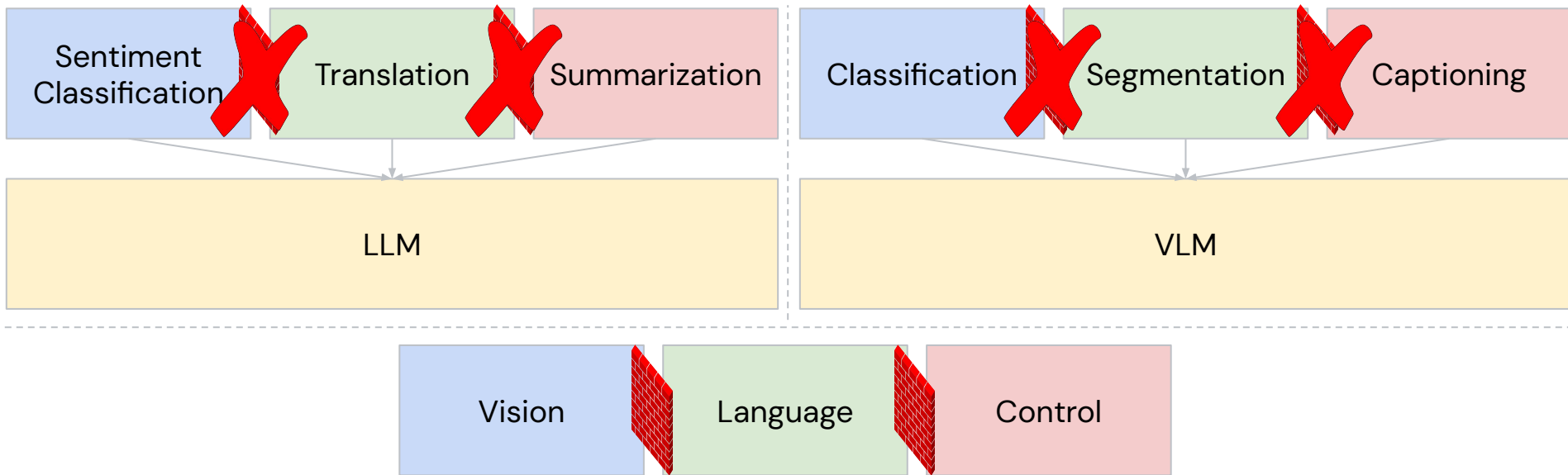
Captioning



# Foundation Model-fication of Robotics?



# Foundation Model-fication of Robotics?



# Foundation Model-fication of Robotics?

Sentiment  
Classification

Translation

Summarization

Classification

Segmentation

Captioning

LLM

VLM

Vision

Language

Control

**Robotics-first Foundation Model**

# Missing Foundation Model Pieces

**Non-robotics Foundation  
Models**

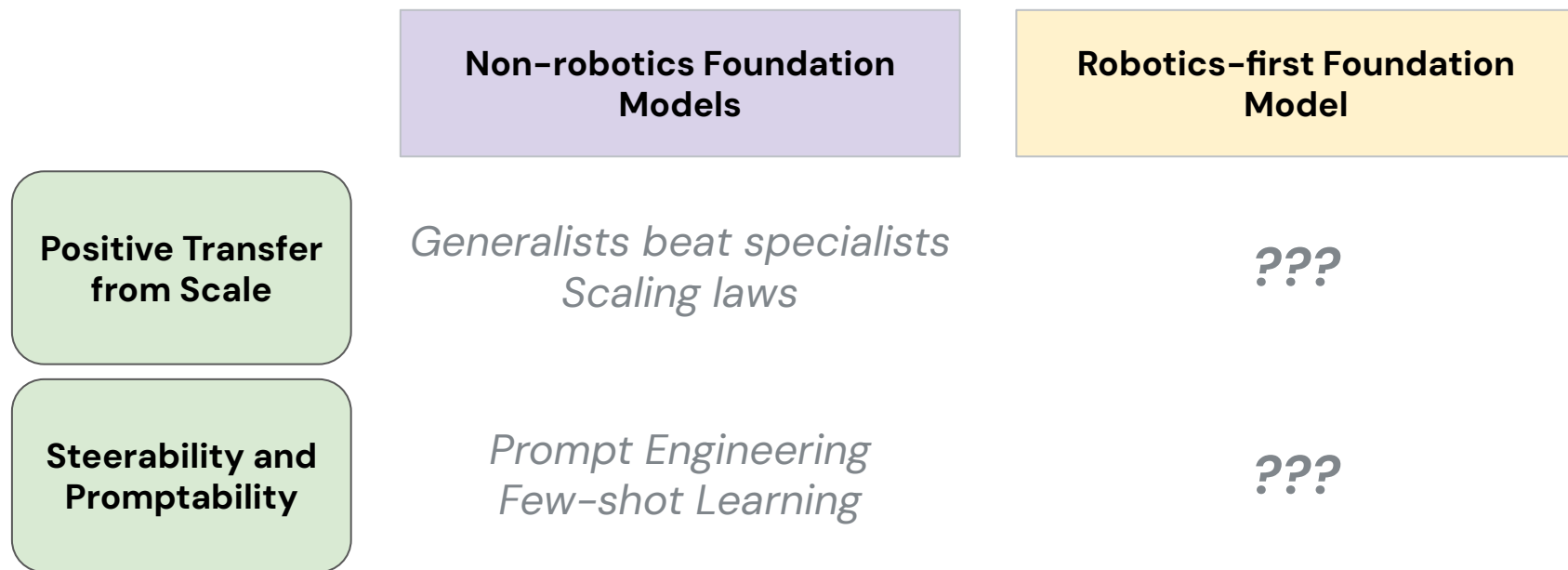
**Robotics-first Foundation  
Model**

**Positive Transfer  
from Scale**

*Generalists beat specialists*  
*Scaling laws*

???

# Missing Foundation Model Pieces



# Missing Foundation Model Pieces

	Non-robotics Foundation Models	Robotics-first Foundation Model
Positive Transfer from Scale	<i>Generalists beat specialists</i> <i>Scaling laws</i>	???
Steerability and Promptability	<i>Prompt Engineering</i> <i>Few-shot Learning</i>	???
Scalable Evaluations	<i>Realistic Evals</i> <i>Predictive Benchmarks</i>	???

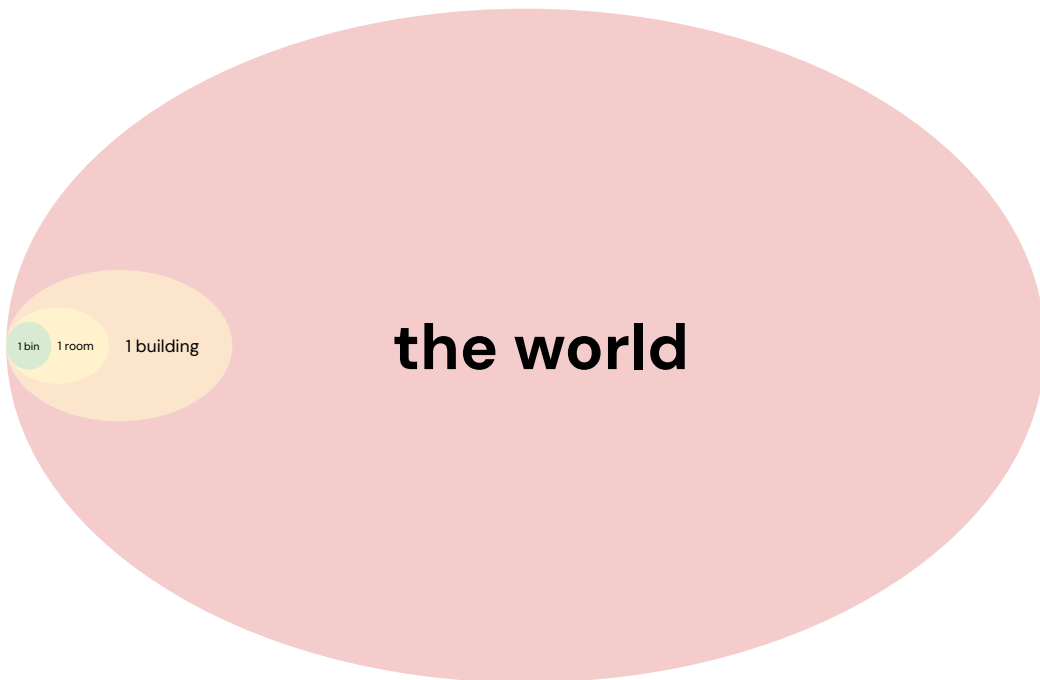
# Missing Foundation Model Pieces

*Claim: These missing properties are necessary for robotics to operate in the real world*

**Positive Transfer  
from Scale**

**Steerability and  
Promptability**

**Scalable  
Evaluations**



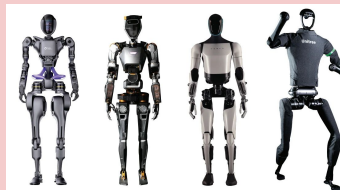
# Missing Foundation Model Pieces

*Claim: These missing properties are necessary for robotics to operate in the real world*

Positive Transfer  
from Scale

Steerability and  
Promptability

Scalable  
Evaluations



1 bin 1 room 1 building

the world





# Missing Foundation Model Pieces

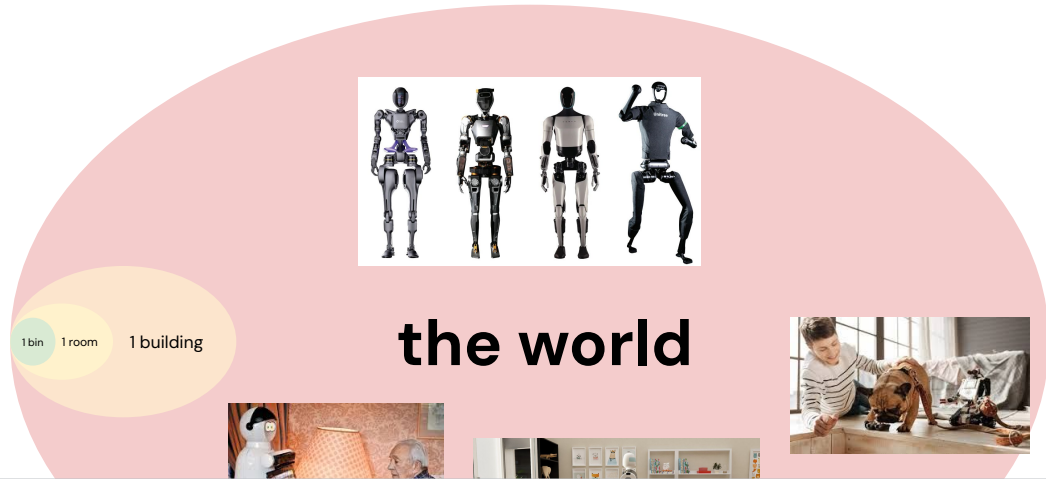
*Claim: These missing properties are necessary for robotics to operate in the real world*

Positive Transfer  
from Scale

Steerability and  
Promptability



**2024 level SoTA technology is not sufficient for general robotics.  
At least one or two paradigm shifts (algorithms and data) required**

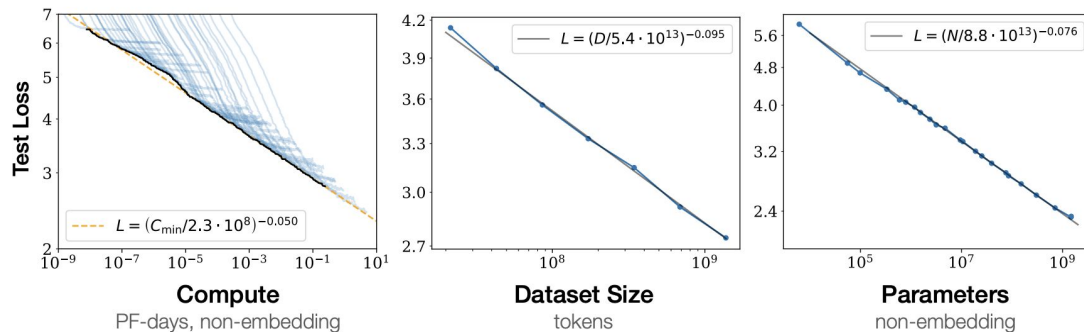


# Agenda

- 01 Why Robot Foundation Models?
- 02 **Piece #1: Positive Transfer from Scaling**
- 03 Piece #2: Steerability
- 04 Piece #3: Scalable Evaluation
- 05 Horizons

# Lessons from Foundation Modeling: Data Scaling

- **Data scaling** a key ingredient in LLMs and VLMs
- ...but the internet already exists. No equivalent for robot data yet!



Source: Kaplan et al. 2020

# Lessons from Foundation Modeling: Data Scaling

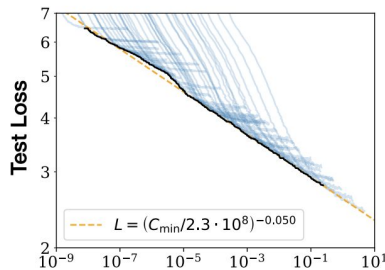
- **Data scaling** a key ingredient in LLMs and VLMs
- ...but the internet already exists. No equivalent for robot data yet!

#1

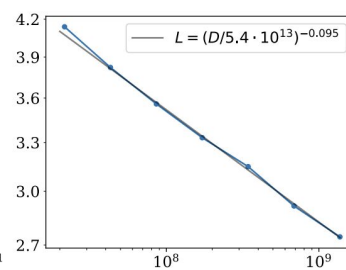
Merge robot data with internet data?

#2

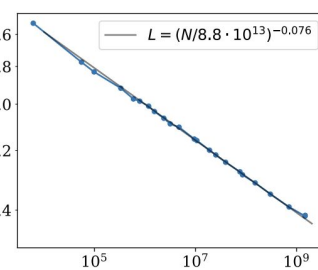
Merge all kinds of robot data?



Compute  
PF-days, non-embedding



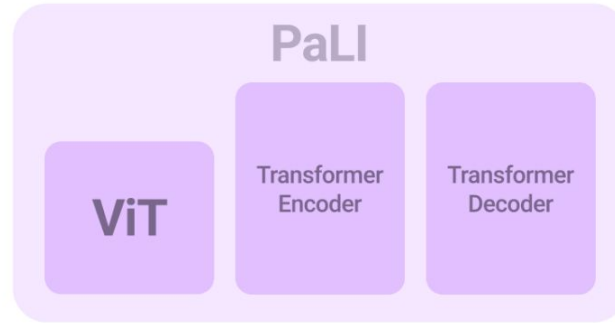
Dataset Size  
tokens



Parameters  
non-embedding

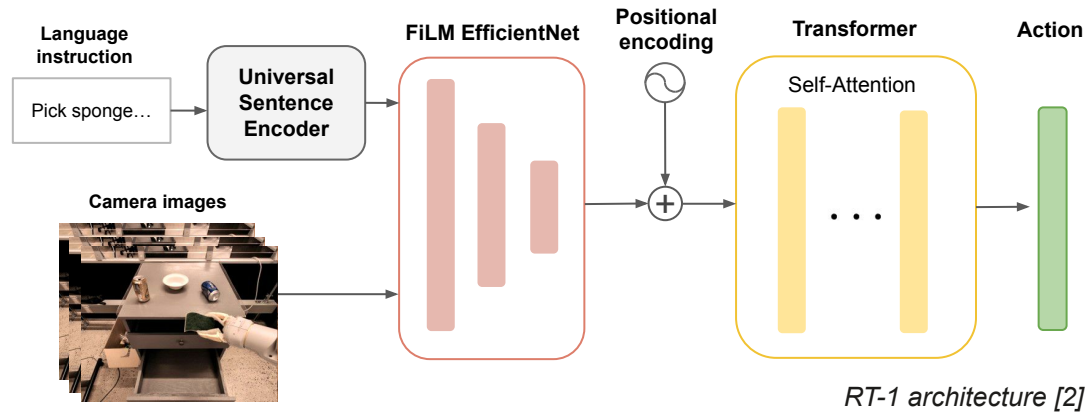
Source: Kaplan et al. 2020

# Vision-Language Models

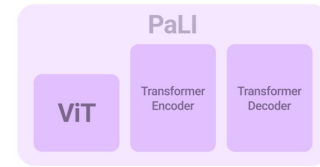


- VLMs encompass both **visual** and **semantic** understanding of the world

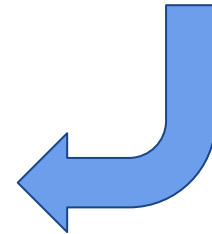
# VLMs as Robot Policies



RT-1 architecture [2]



PaLI architecture [1]

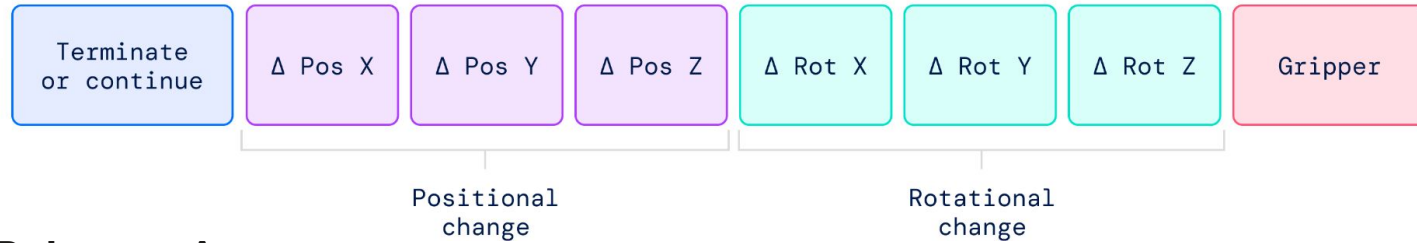


- **RT-1:** image + text  $\rightarrow$  **discretized actions**
- Similar to a Visual-Language Model (VLM) with different **output tokens**
- Use large pre-trained VLMs directly as the **policy!**
- How do we **deal with actions** when using pre-trained VLMs?

[1] PaLI: A Jointly-Scaled Multilingual Language-Image Model. Chen et al. 2022.

[2] RT-1: Robotics Transformer for Real-World Control at Scale, Robotics at Google and Everyday Robots, 2022.

# Representing Actions in VLMs



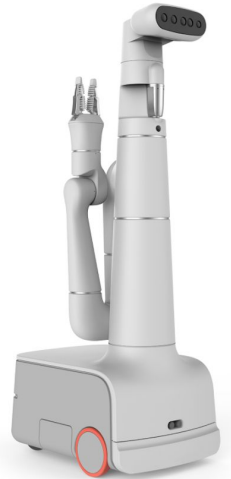
- **Robot actions:**

- Moving the robot arm and gripper
- Discretized into 256 bins

- **Actions in VLMs**

- Convert to a string of numbers
- Example: "1 127 115 218 101 56 90 255"
- Alternatives:
  - *Float numbers* – more tokens needed
  - *Extra-IDs, least used language tokens*
  - *Human language (left, right etc.)* – can't be directly executed on a robot

→ **Vision-Language-Action (VLA) model!**



# Training data and underlying models

## Models

- PaLI-X (5B, 55B)
- PaLM-E (12B)

## Data

- Pretraining: Web-data
- Robot data
  - RT-1 data
  - 13 robots
  - 17 months
  - 130k demos

### Internet-Scale VQA + Robot Action Data



Q: What is happening in the image?

A grey donkey walks down the street.



Q: Que puis-je faire avec ces objets?

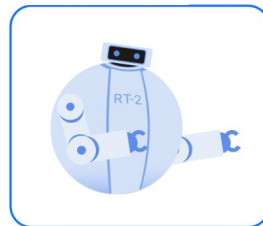
Faire cuire un gâteau.



Q: What should the robot do to <task>?

$\Delta$  Translation =  $[0.1, -0.2, 0]$   
 $\Delta$  Rotation =  $[10^\circ, 25^\circ, -7^\circ]$

Co-Fine-Tune





# Results: Emergent skills



*put strawberry into the correct bowl*



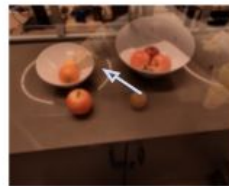
*pick up the bag about to fall off the table*



*move apple to Denver Nuggets*



*pick robot*



*place orange in the matching bowl*



*move redbull can to H*



*move soccer ball to basketball*



*move banana to Germany*



*move cup to the wine bottle*



*pick animal with different color*



*move coke can to Taylor Swift*



*move coke can to X*



*move bag to Google*

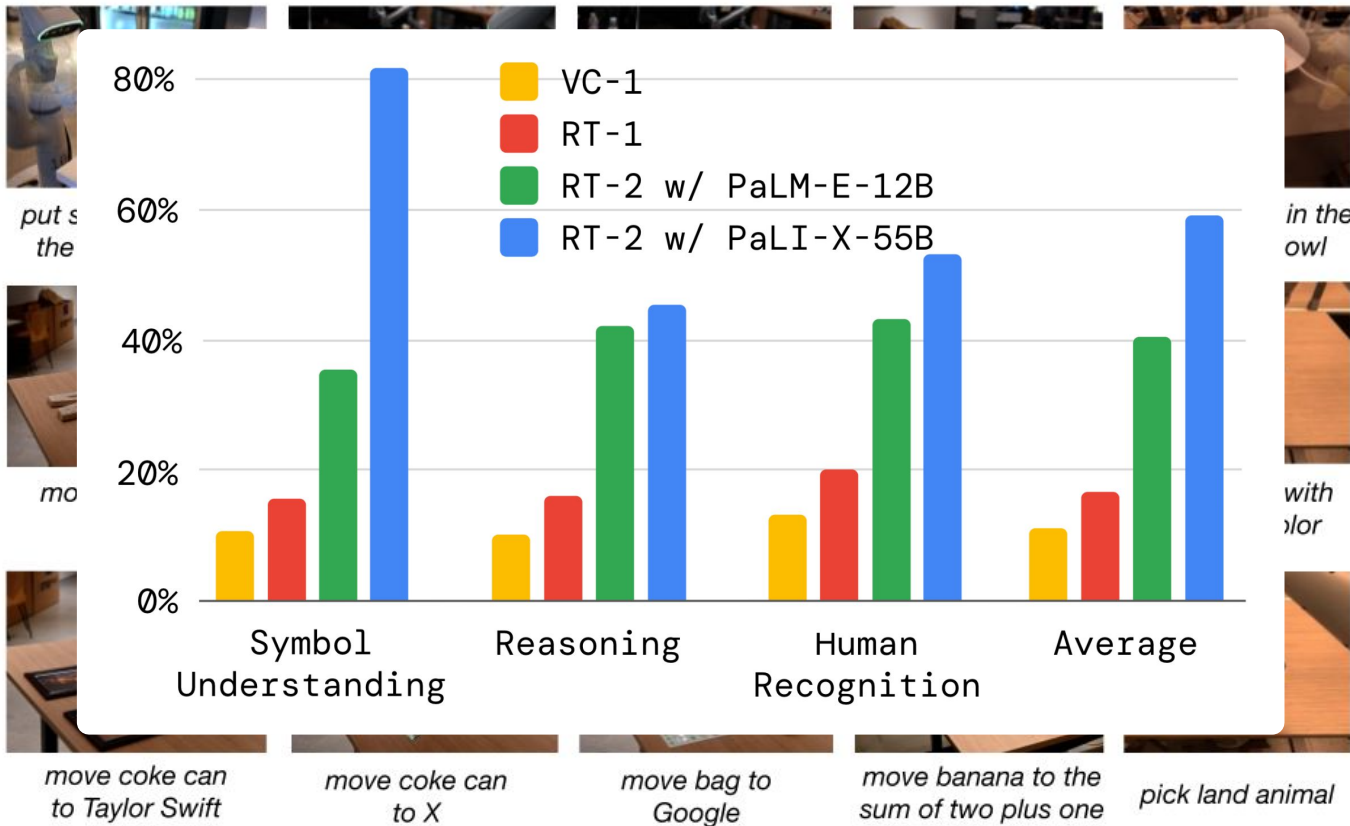


*move banana to the sum of two plus one*



*pick land animal*

# Results: Emergent skills



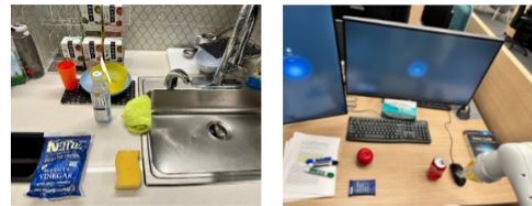
# Results: Quantitative evals



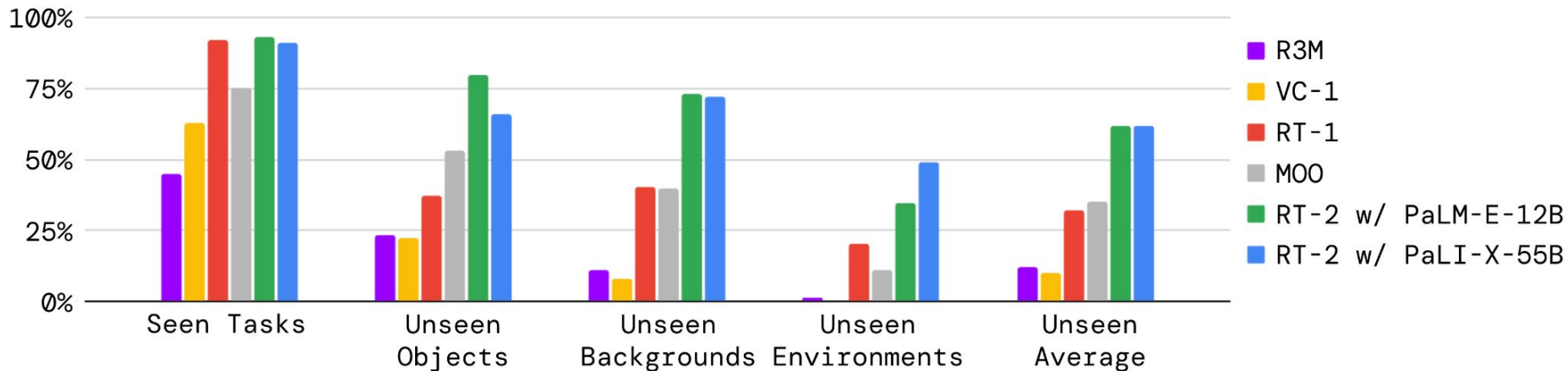
(a) Unseen Objects



(b) Unseen Backgrounds



(c) Unseen Environments



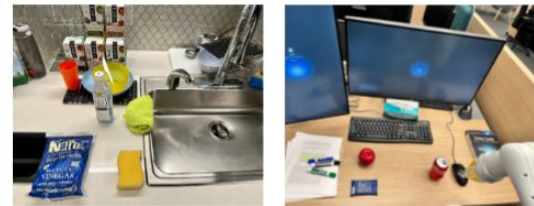
# Results: Quantitative evals



(a) Unseen Objects



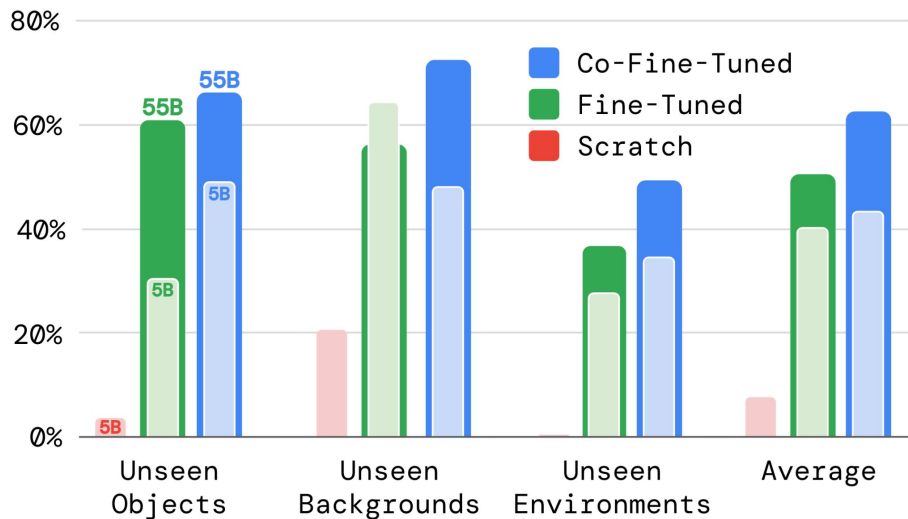
(b) Unseen Backgrounds



(c) Unseen Environments

## RT2 w/ PaLI-X-55B ablations

- Co-Fine-Tuning with VQA data
- Fine-Tuning on robot data only
- Training on robot data from scratch



# Positive Transfer: Chain-of-Thought with RT-2-PaLM-E

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



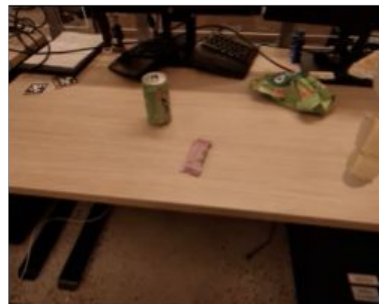
Prompt:

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?

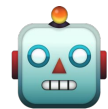
Prediction:

Rocks. Action: 1 129 138  
122 132 135 106 127





# The Open X-Embodiment Dataset



**1M+** Real Robot Episodes



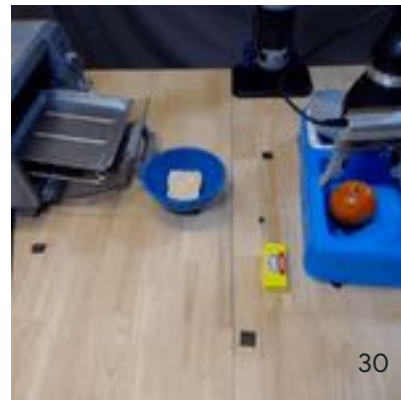
**22** Robot Embodiments



**34** Research Labs

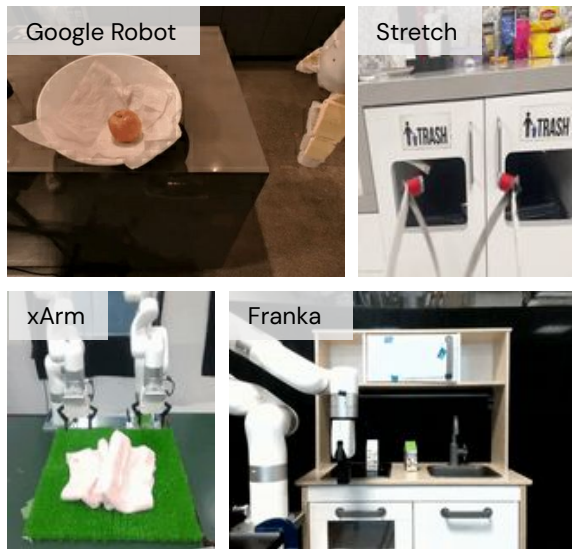


**300+** Scenes

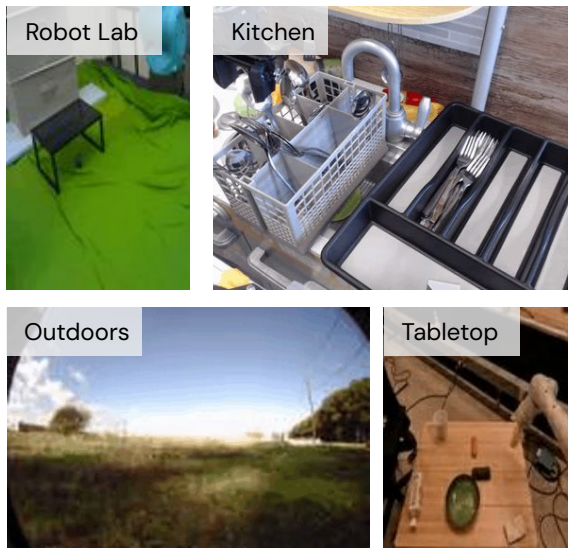


# The Open X-Embodiment Dataset

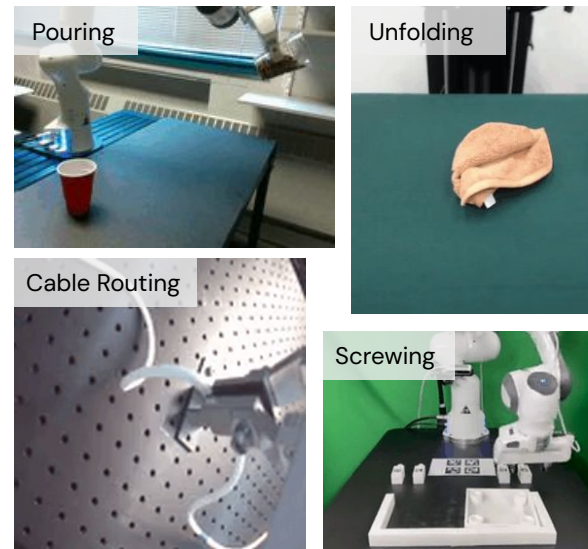
## Many Embodiments



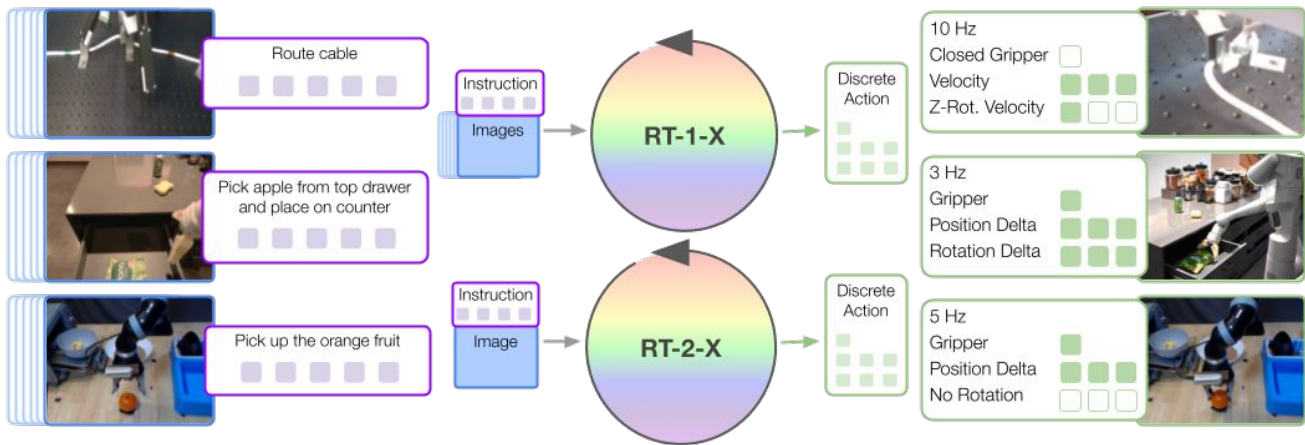
## Many Scenes



## Many Skills



# Model Architectures



Inputs: RGB images and text instructions

Outputs: discretized end-effector actions

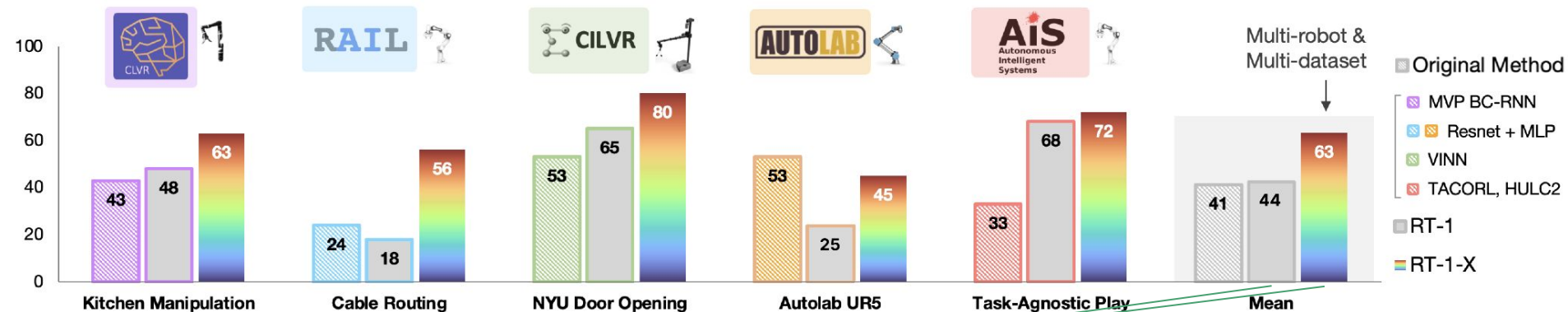
Just RT-1 and RT-2 trained on X-Embodiment datasets

Velocity, delta position, absolute position

Different evaluations run at different frequencies



# Results: Signs of Positive Transfer

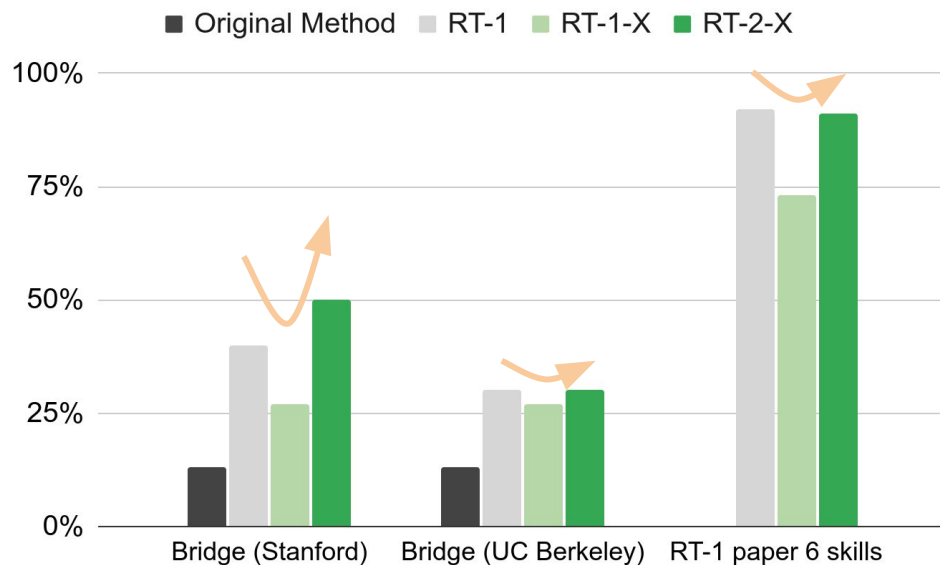


Generalist (RT-1-X) vs. Specialists (RT-1, Baselines)

- Training on data from **all robots** outperforms training on data from the particular evaluation robot

50% improvement

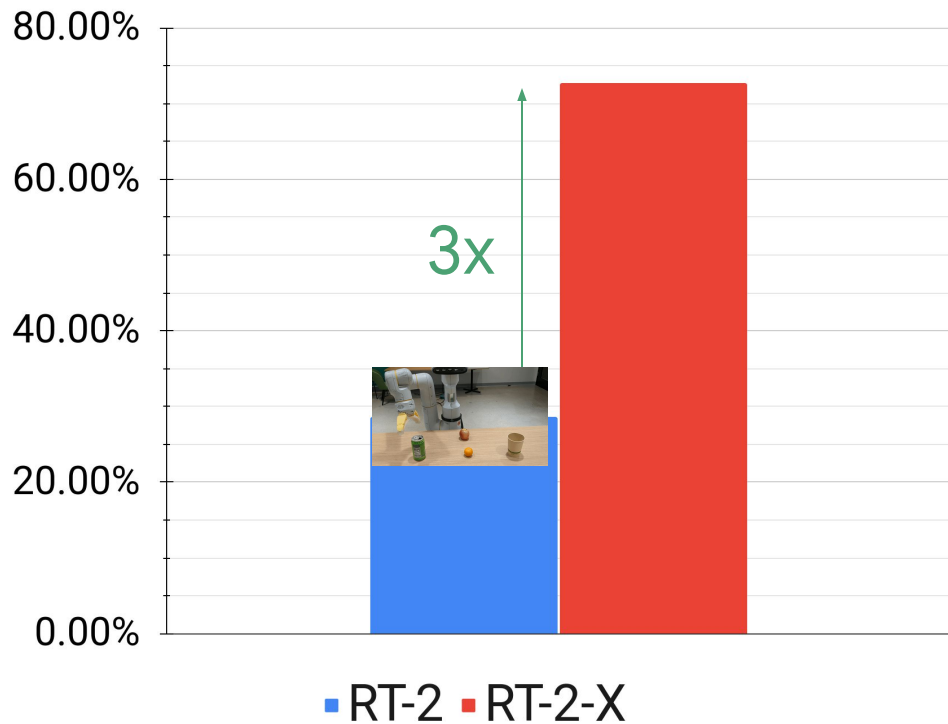
# Results: Small Models Underfit



RT-1-X underfits for large datasets

RT-2-X recovers performance

# Is Web-scale Data Sufficient?



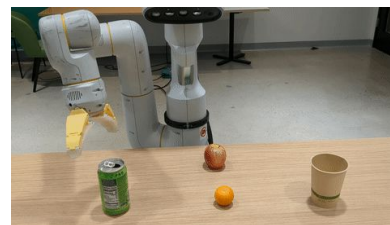
RT-2-X outperforms RT-2 by 3x  
in emergent skill evaluations



put apple on cloth



move apple near cloth



move apple between  
cup and apple

# Data Scaling and Positive Transfer Recap

Real-world robot  
demonstration dataset

*[RT-1]*

Co-train on robot data  
alongside internet data

*[RT-2]*

Add robot data from  
different embodiments

*[RT-X]*

**Increasing data interoperability by treating robot actions  
as just another data modality**

# Data Scaling and Positive Transfer Recap

Real-world robot demonstration dataset

*[RT-1]*

Co-train on robot data alongside internet data

*[RT-2]*

Add robot data from different embodiments

*[RT-X]*

Better **in-distribution** performance

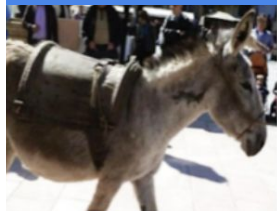
Generalization to internet **semantics**

Generalization to **spatial** concepts

# ...But Many Open Challenges!

VLA overfit to robotics data distributions

VQA Image



Robot Image



Robot Image



VQA Prompt

Q: What is happening in the image?

Action Prompt

Q: What action should the robot take to pick coke can?

VQA Prompt

Q: What is happening in the image?

A: A grey donkey walks down the street.



A: 1 130 129 121  
131 127 128 127



A: 1 127 127 127  
127 127 127 127



# ...But Many Open Challenges!

VLA overfit to robotics data distributions

VQA Image



Robot Image



Robot Image



VQA Prompt

Q: What is happening in the image?

Action Prompt

Q: What action should the robot take to pick coke can?

VQA Prompt

Q: What is happening in the image?

A: A grey donkey walks down the street. ✓

A: 1 130 129 121  
131 127 128 127 ✓

A: 1 127 127 127  
127 127 127 127 ✗

Reasoning mixes unpredictably with low-level robot action control

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

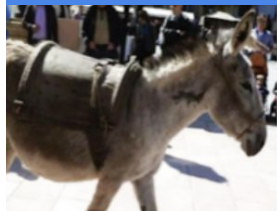


*If coke can added to scene, planning works but actions break!*

# ...But Many Open Challenges!

VLA overfit to robotics data distributions

VQA Image



Robot Image



Robot Image



VQA Prompt

Q: What is happening in the image?

Action Prompt

Q: What action should the robot take to pick coke can?

VQA Prompt

Q: What is happening in the image?

A: A grey donkey walks down the street. ✓

A: 1 130 129 121  
131 127 128 127 ✓

A: 1 127 127 127  
127 127 127 127 ✗

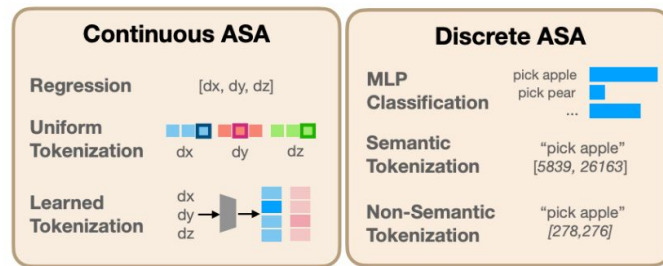
Reasoning mixes unpredictably with low-level robot action control

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



*If coke can added to scene, planning works but actions break!*

Action representations and tokenization decision choices are underexplored





# Agenda

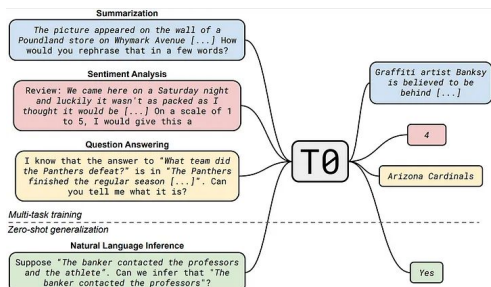
- 01 Why Robot Foundation Models?
- 02 Piece #1: Positive Transfer from Scaling
- 03 Piece #2: Steerability**
- 04 Piece #3: Scalable Evaluation
- 05 Horizons




We convey intent to robot policies  
via very constrained interfaces...

...but LLM reasoning is enabled by  
large context bandwidths.

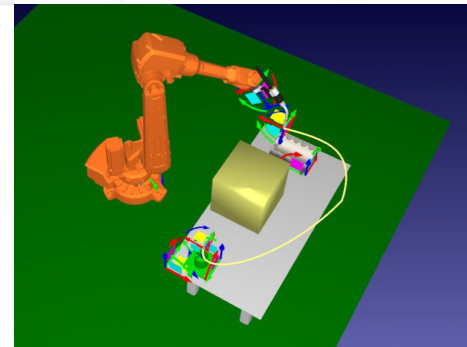
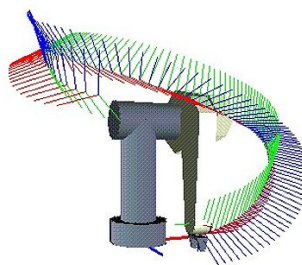
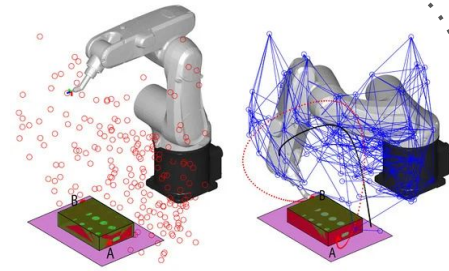
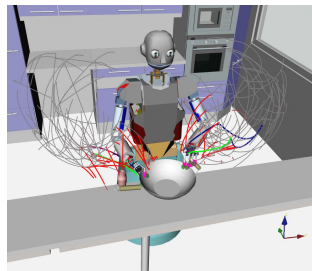
***Where is my promptable generalist robot??***

# Strengths and Limitations of Language



	<b>Q:</b> What does the man who sits have trouble doing?	<b>MC Answers:</b> (a) Riding (b) Breathing (c) Walking (d) Magic <b>Direct Answers:</b> Walking, Walking, ...	<b>Rationale:</b> The vehicle being used is for people who cannot use their legs properly and need it for assistance in being mobile.
	<b>Q:</b> What could block the washer's door?	<b>MC Answers:</b> (a) Stool (b) Stove (c) Window (d) Sink <b>Direct Answers:</b> Stove, Oven, ...	<b>Rationale:</b> The washer door is right in front of the range preventing it from opening.
	<b>Q:</b> How many people will dine at this table?	<b>MC Answers:</b> (a) Two (b) One (c) None (d) Five <b>Direct Answers:</b> One, One person, ...	<b>Rationale:</b> There is only one cup of water and main dish at this table.

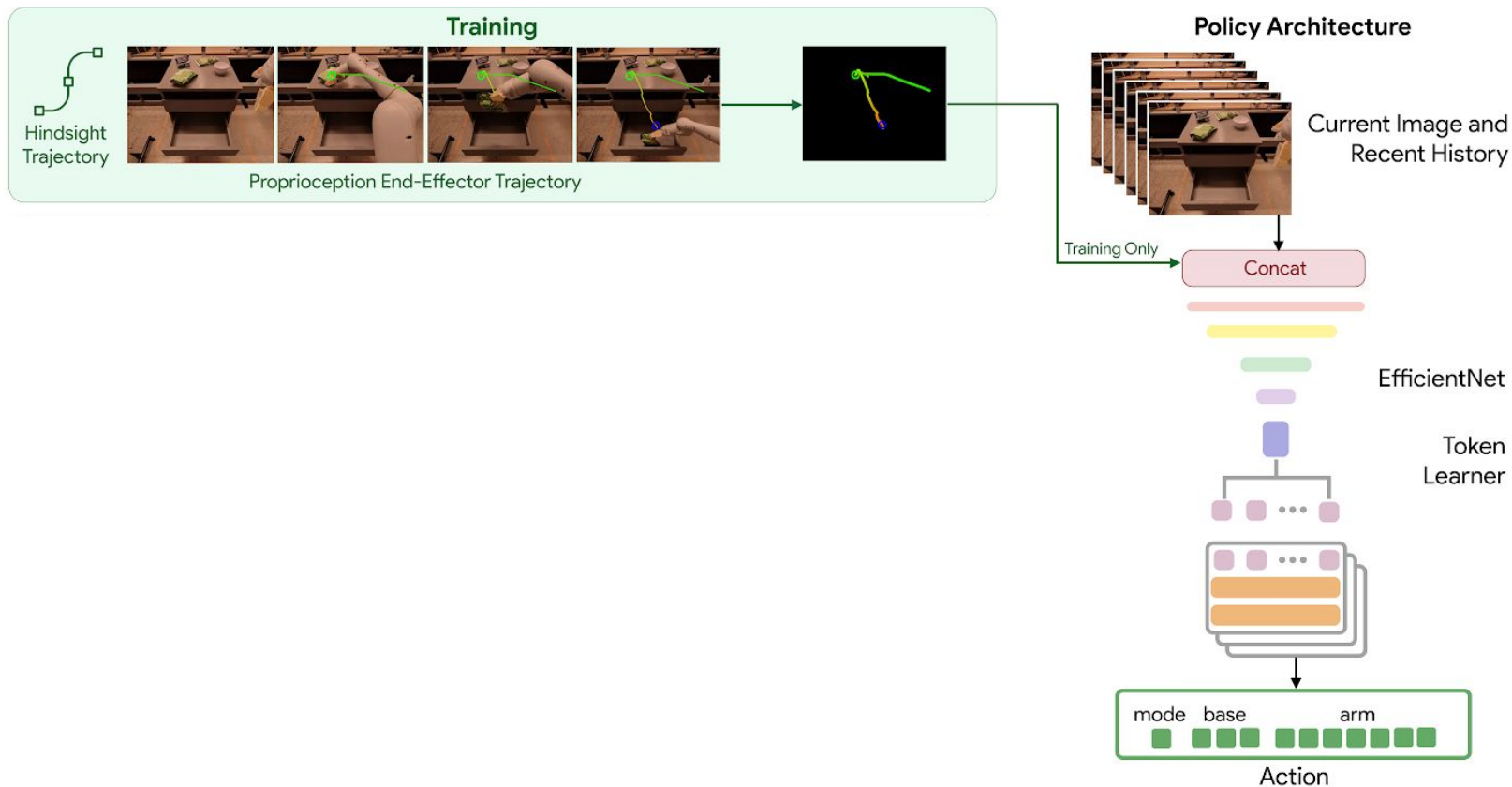
High-level Language Knowledge



Low-level Robotics Knowledge

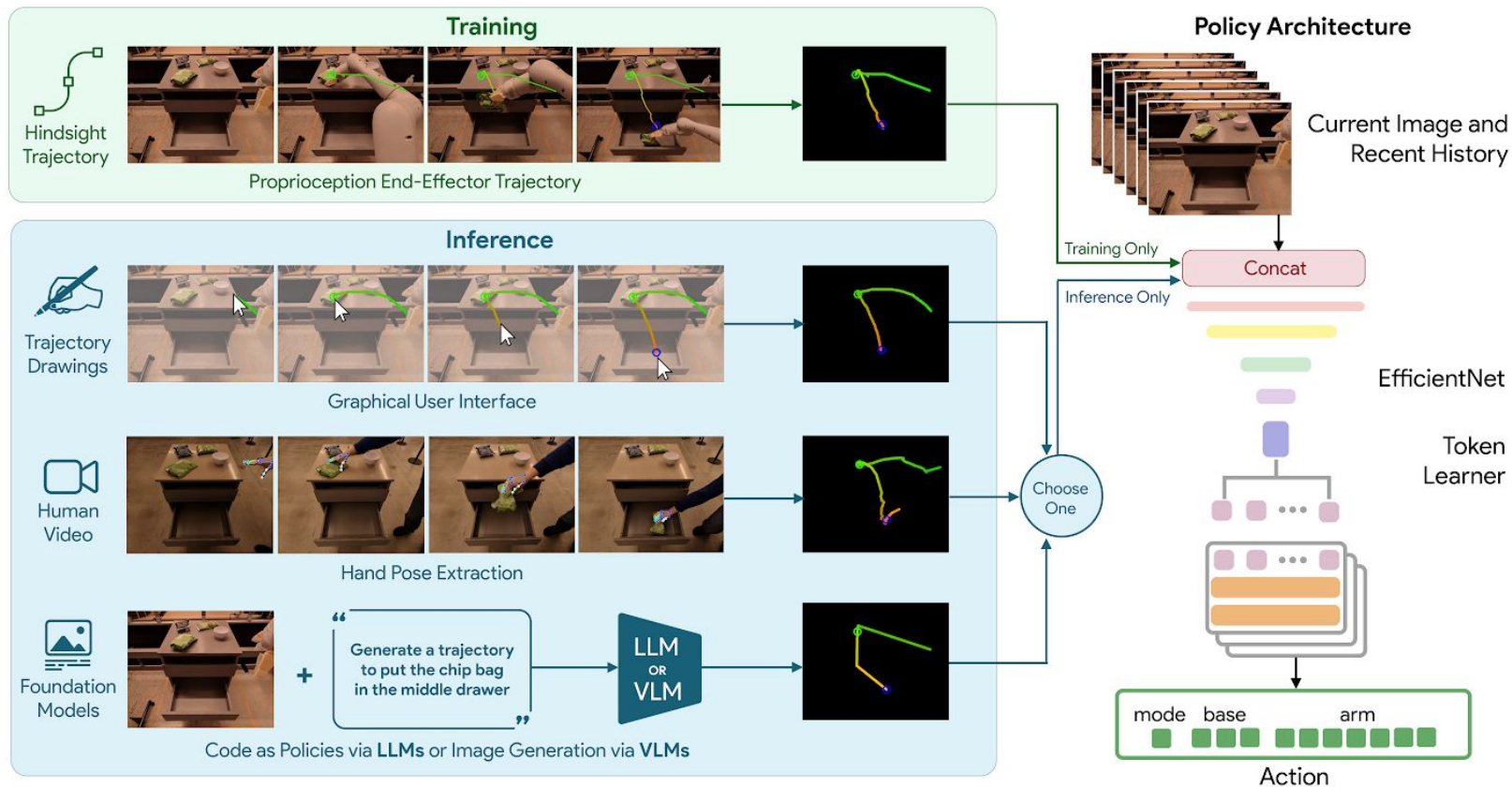
# Motion-centric Representations: Hindsight Trajectories

## RT-Trajectory

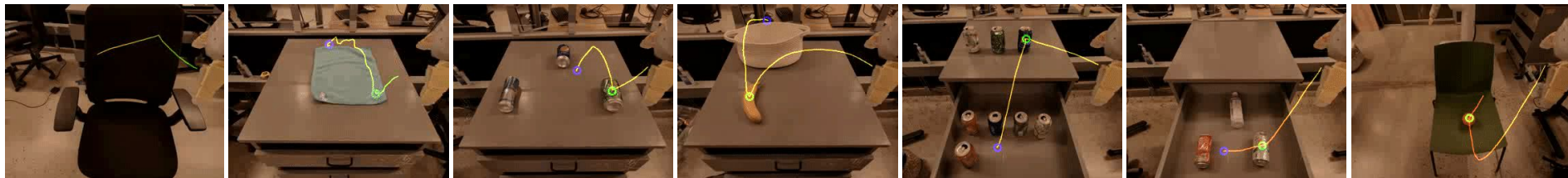
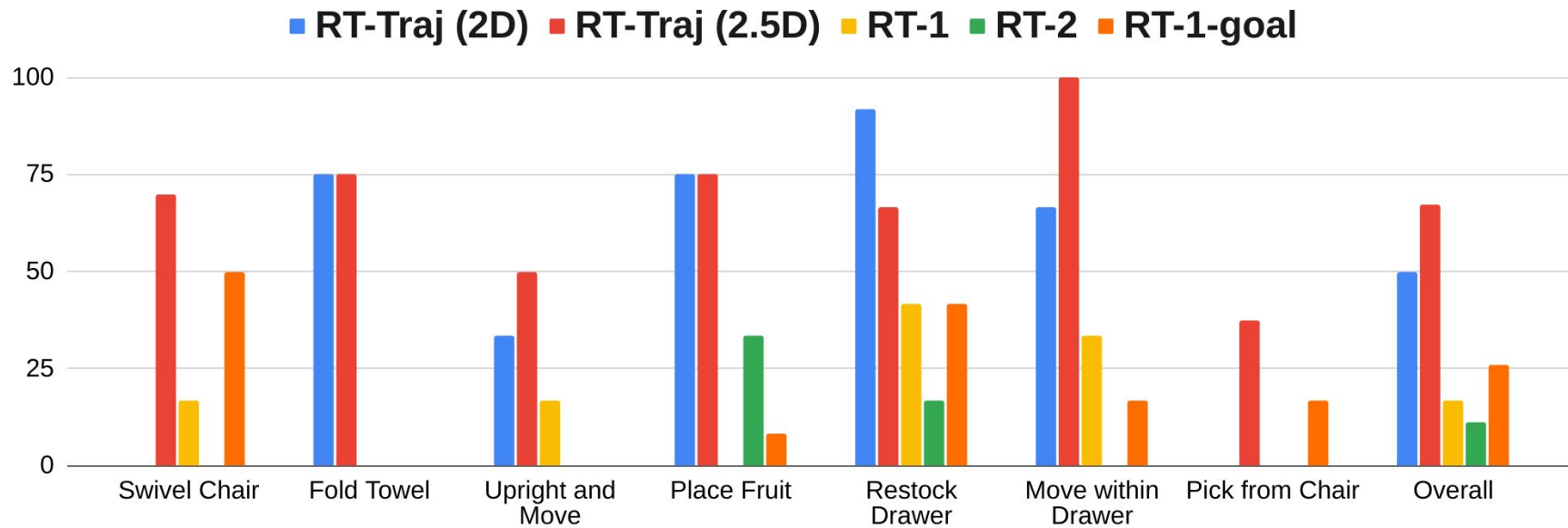


# Motion-centric Representations: Hindsight Trajectories

## RT-Trajectory



# Results: Quantitative Evaluations



# Results: Prompt Engineering via Trajectories

## Ego-centric trajectory representations enable broad generalization:

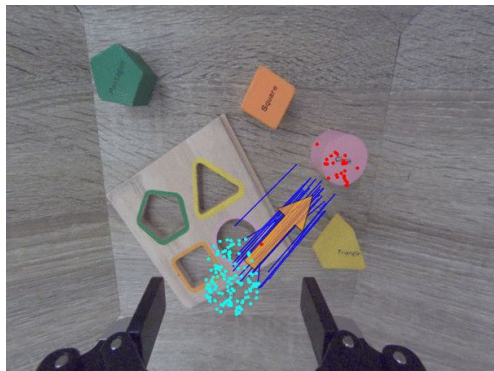
- Novel motions (new heights, new shapes, new curvatures)
- Visual distribution shifts (new furniture, new rooms, new objects, new lighting)
- Behavior modulation within skills (specify exactly *how* to accomplish the task)





# Concurrent Work: Tracks, Flow, Motion

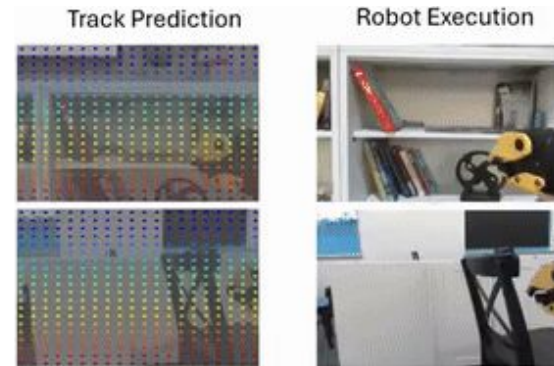
Motions and trajectories are a powerful representation which capture the unique properties of robotics: actions, dynamics, physics, change



RoboTAP



Any-point Trajectory Modeling



Track2Act

[4] *RoboTAP: Tracking Arbitrary Points for Few-Shot Visual Imitation*, Vecerik et al., 2023.

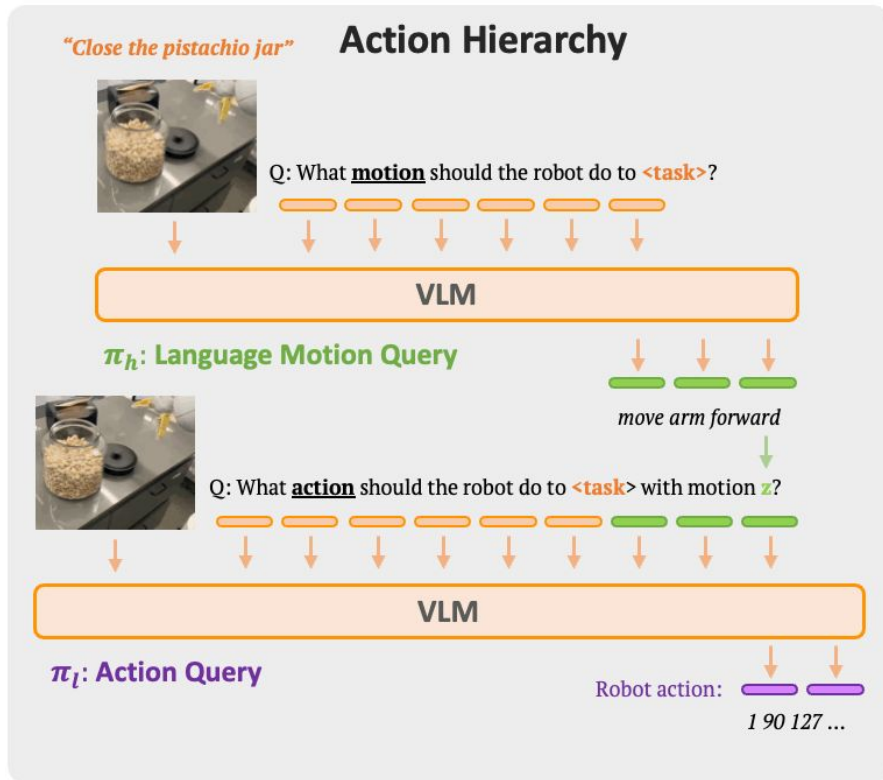
[5] *Any-point Trajectory Modeling for Policy Learning*, Wen et al., 2024.

[6] *Track2Act: Predicting Point Tracks from Internet Videos enables Diverse Zero-shot Robot Manipulation*, Bharadhwaj et al. 2024.



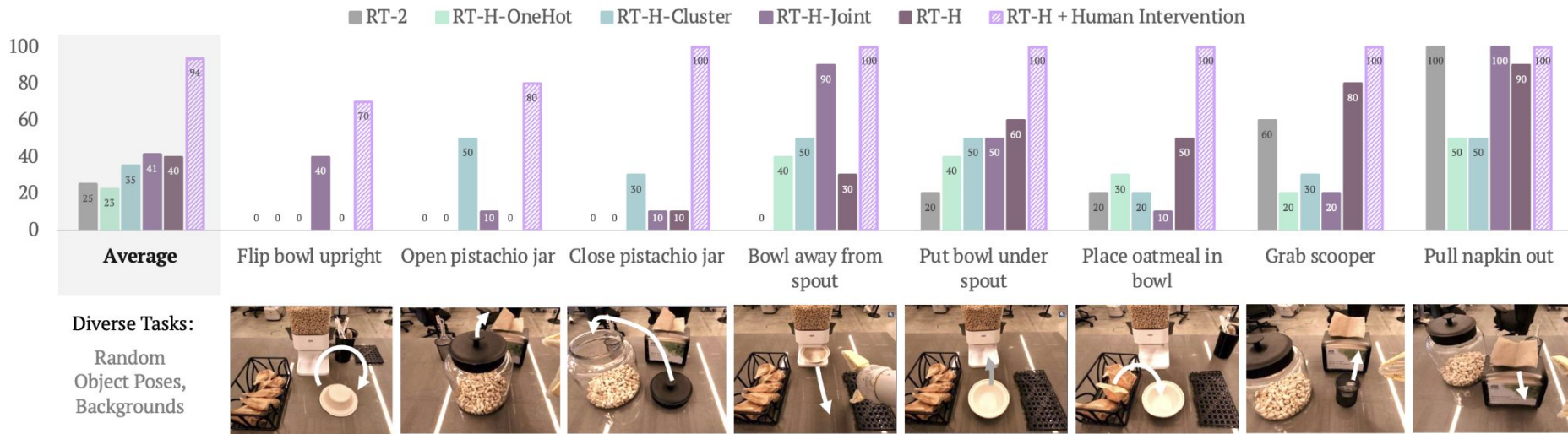
# Is language enough, if it's *hierarchical* and *granular*?

## RT-Hierarchy



- Idea: predict granular language motions before predicting low-level robot actions
  - "move arm forward", "rotate arm clockwise", "close gripper"
- Can be viewed as chain-of-thought / planning for language-based skills

# Results: RT-H Outperforms RT-2



*No other policy class (RT-1, RT-2) was able to learn from challenging new data*

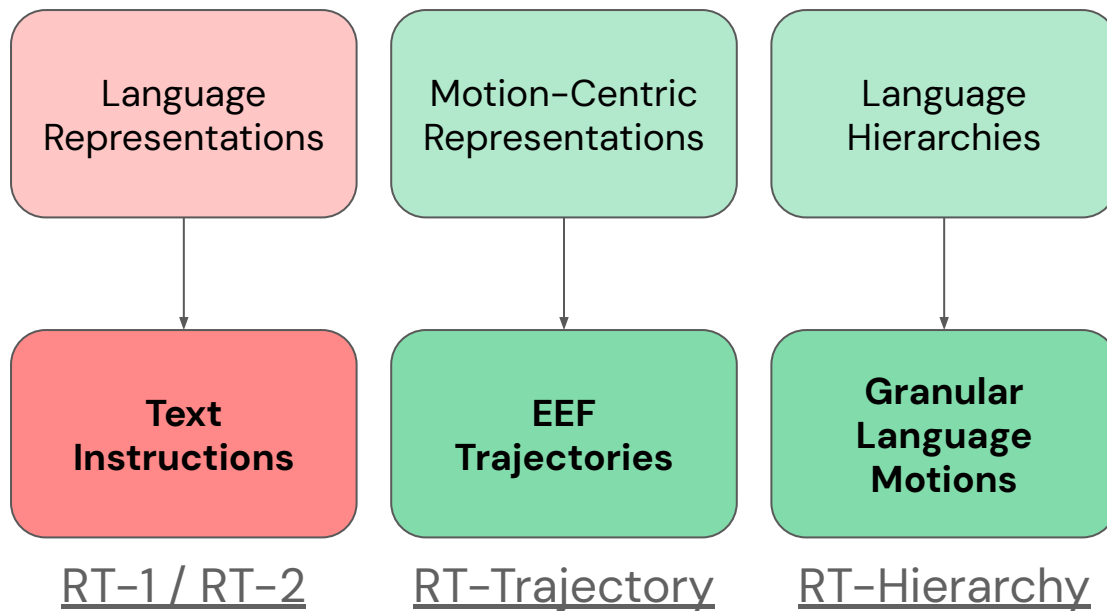
# Results: Language Interventions

**Task: “Close the pistachio jar”**

**Action Hierarchies Improve Performance and Enable Intervention**

RT-H bottleneck often was language motion prediction rather than low-level action prediction: language motions easier to collect interventions for!

# Steerability Recap



We have proofs of concept for  
promptable robots...

...but do we have enough robot data  
to support these algorithms?

We have proofs of concept for promptable robots...

...but do we have enough robot data to support these algorithms?



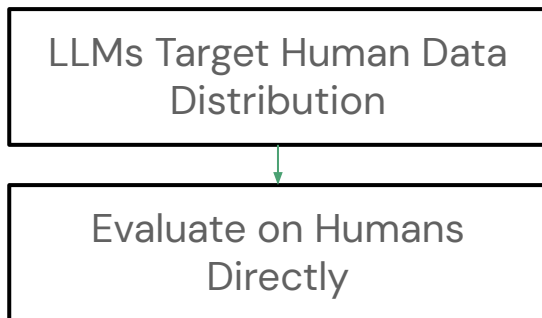
***Robot data is not guaranteed to be a bottleneck because we don't yet know what kind of robot data we need***

# Agenda

- 01 Why Robot Foundation Models?
- 02 Piece #1: Positive Transfer from Scaling
- 03 Piece #2: Steerability
- 04 Piece #3: Scalable Evaluation**
- 05 Horizons

# AI has an Evaluation Problem

- All roads lead to generalist models, but generalist models that can "do anything" need to be evaluated on "everything"!
- How do you scalably evaluate a broad set of capabilities?



**LMSYS Chatbot Arena Leaderboard**

MSYS Chatbot Arena is a crowdsourced open platform for LLM evals. We've collected over 300,000 human preference comparisons to rank LLMs with the [GPT-4o](#) and display the model ratings in this table. You can view more details in our [help page](#).

Model: [GPT-4o](#)

Test Arena ID: Test-Arena-072-236. Last updated: April 13, 2024.

NOTE: View leaderboard for different categories (e.g., coding, long form evals).

Users receive leaderboard titles and points in this [contest](#). You can contribute your vote @ [#\\_ChatbotArena](#)

Rank	Model	Score	Score E1s	95% CI	Notes	Organization	Licensed	Roundtrip Cost/1k
1	<a href="#">GPT-4o (2024-06-09)</a>	1240	+4.5	1376	OpenAI	Proprietary	2023/12	
2	<a href="#">Claude 3.5 Sonnet</a>	1205	+3.4	95381	Anthropic	Proprietary	2023/9	
3	<a href="#">GPT-4o (2024-04-15)</a>	1204	+3.3	65109	OpenAI	Proprietary	2023/9	
4	<a href="#">GPT-4o (2024-04-15)</a>	1200	+3.4	59923	OpenAI	Proprietary	2023/12	
5	<a href="#">DeepSeek-V3</a>	1200	+5.5	13468	Google	Proprietary	On-line	
6	<a href="#">Claude 3.5 Sonnet</a>	1203	+3.3	62856	Anthropic	Proprietary	2023/9	
7	<a href="#">Gemini 1.5 Pro</a>	1193	+4.4	24397	Google	CC-BY-NC-SA	2024/3	
8	<a href="#">GPT-4o (2024-04-15)</a>	1189	+4.4	62920	OpenAI	Proprietary	2024/9	

## HumanEval: Hand-Written Evaluation Set

This is an evaluation harness for the HumanEval problem solving dataset described in [Language Models Trained on Code](#).





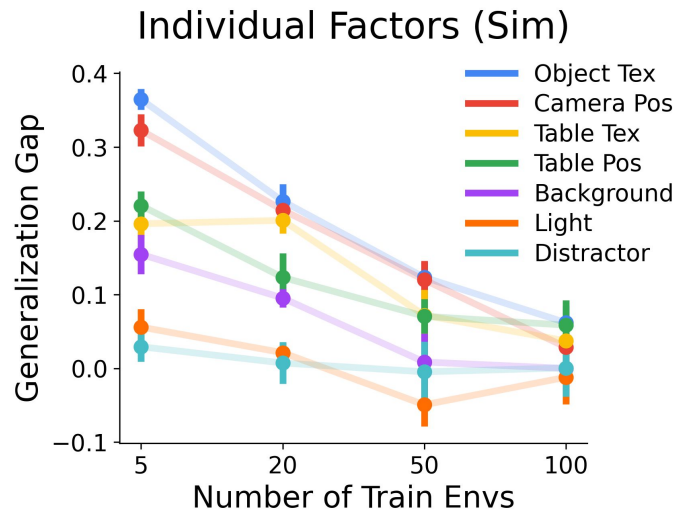
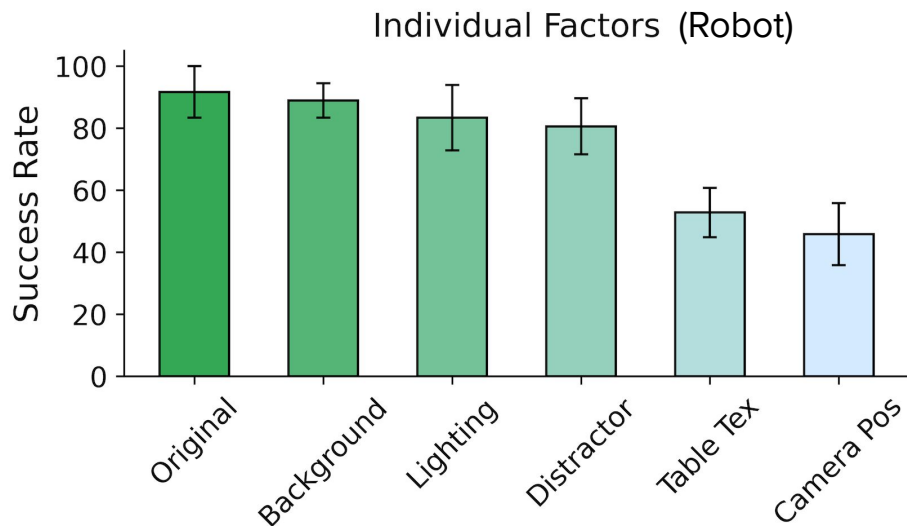
# Measuring Axes of Generalization

Can we *systematically* measure policy generalization?



**Evaluation Metrics:** success rate, generalization gap (train – test success rate)

# Impact of Individual Factors



**“Easier” factors:** background, lighting, distractor

**“Harder” factors:** table position, table texture, camera position, object texture

# Real-to-Sim Evaluation for Real-world Robot Policies

## Real Robot Evaluation

(Train on real, eval on real)

- ⊖ Slow
- ⊖ Expensive
- ⊖ Not Reproducible

REAL

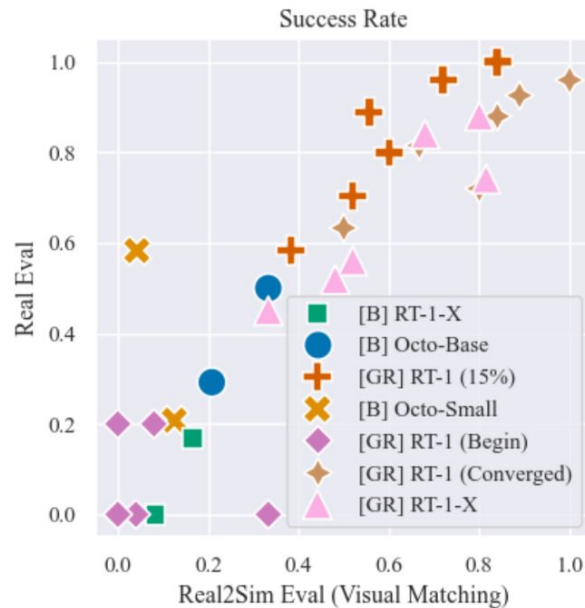
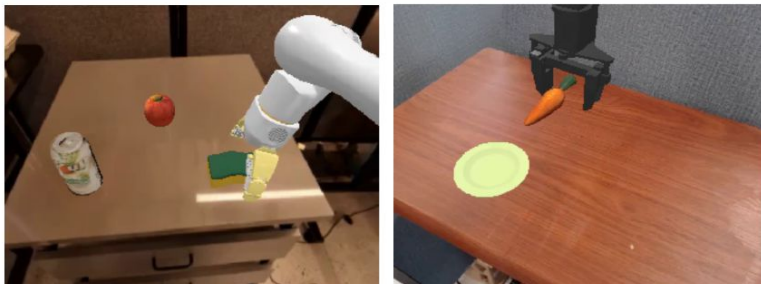


## Real-to-Sim Evaluation

(Train on real, eval in sim)

- + Cheap
- + Scalable
- + Fully Reproducible

SIM



Key Insight: A simulation "good enough" for useful evaluation signal may be much easier to build than a full digital clone for training

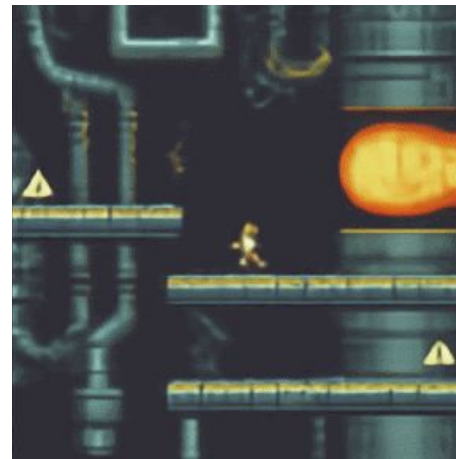
# World Models for Evaluation



PRISM-1



UniSim



Genie

[4] PRISM-1, Wayve, 2024

[5] UniSim: Learning Interactive Real-World Simulators, Yang et al., 2024

[6] Genie: Generative Interactive Environments, Bruce et al., 2024



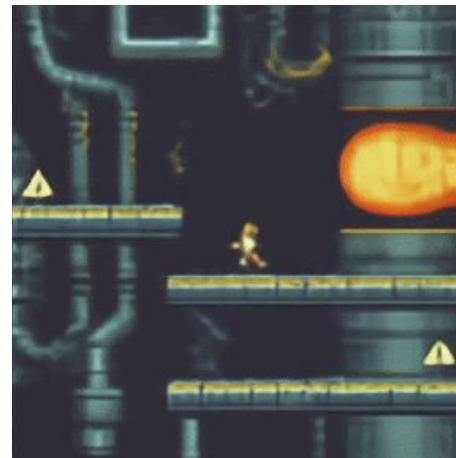
# World Models for Evaluation



PRISM-1



UniSim



Genie



***Real world evaluations will always be the gold standard.  
Scaled evaluations will be solved by unit economics and products.***

[4] PRISM-1, Wayve, 2024

[5] UniSim: Learning Interactive Real-World Simulators, Yang et al., 2024

[6] Genie: Generative Interactive Environments, Bruce et al., 2024

# Agenda

- 01 Why Robot Foundation Models?
- 02 Piece #1: Positive Transfer from Scaling
- 03 Piece #2: Steerability
- 04 Piece #3: Scalable Evaluation
- 05 Horizons**

Missing  
Piece

Positive Transfer  
from Scale

Bleeding  
Edge

VLA Models and  
X-Embodiment

Progress

**6/10**

Horizon

Overfitting and  
little understood  
robot  
post-training



Missing  
Piece

Positive Transfer  
from Scale

Steerability and  
Promptability

Bleeding  
Edge

VLA Models and  
X-Embodiment

Going Beyond  
Language

Progress

**6/10**

**4/10**

Horizon

Overfitting and  
little understood  
robot  
post-training

Robotics-specific  
data is sparse  
with low coverage

Missing  
Piece

Positive Transfer  
from Scale

Steerability and  
Promptability

Scalable  
Evaluations

Bleeding  
Edge

VLA Models and  
X-Embodiment

Going Beyond  
Language

Generalization  
and Simulation

Progress

**6/10**

**4/10**

**3/10**

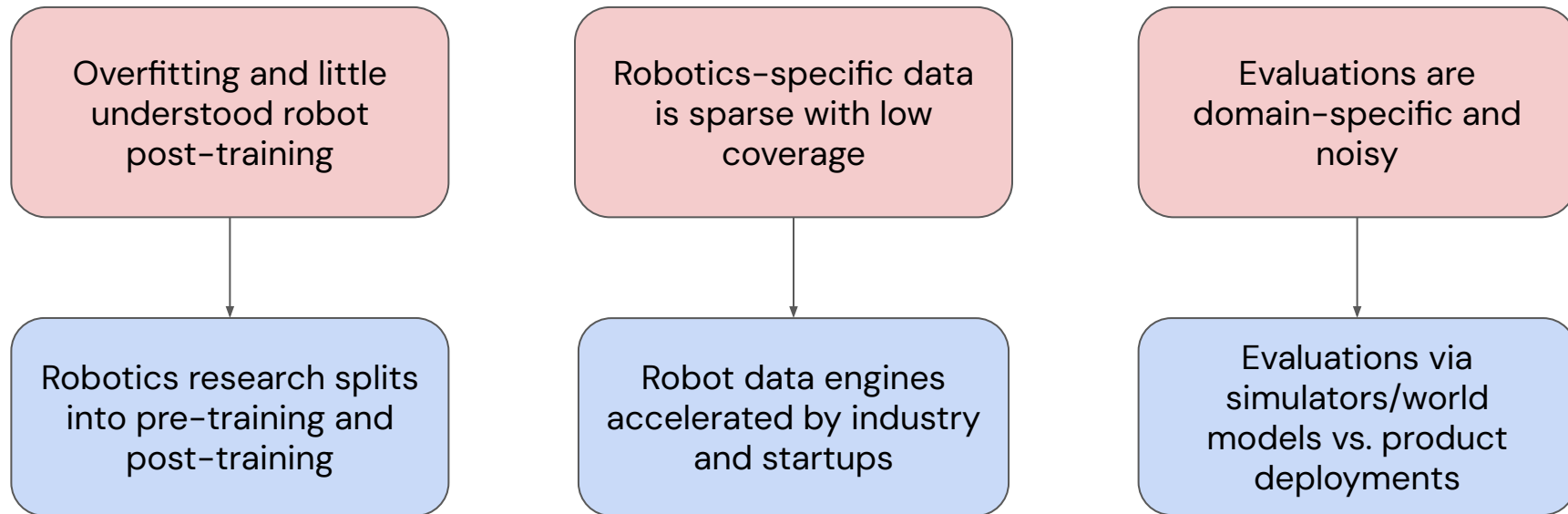
Horizon

Overfitting and  
little understood  
robot  
post-training

Robotics-specific  
data is sparse  
with low coverage

Evaluations are  
domain-specific  
and noisy

# Predictions



# Thank you!

tedxiao@google.com



Google DeepMind