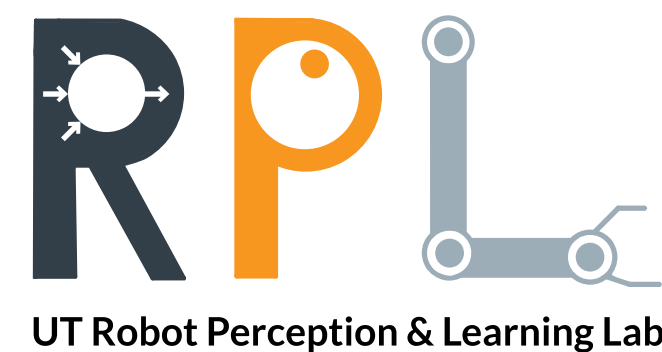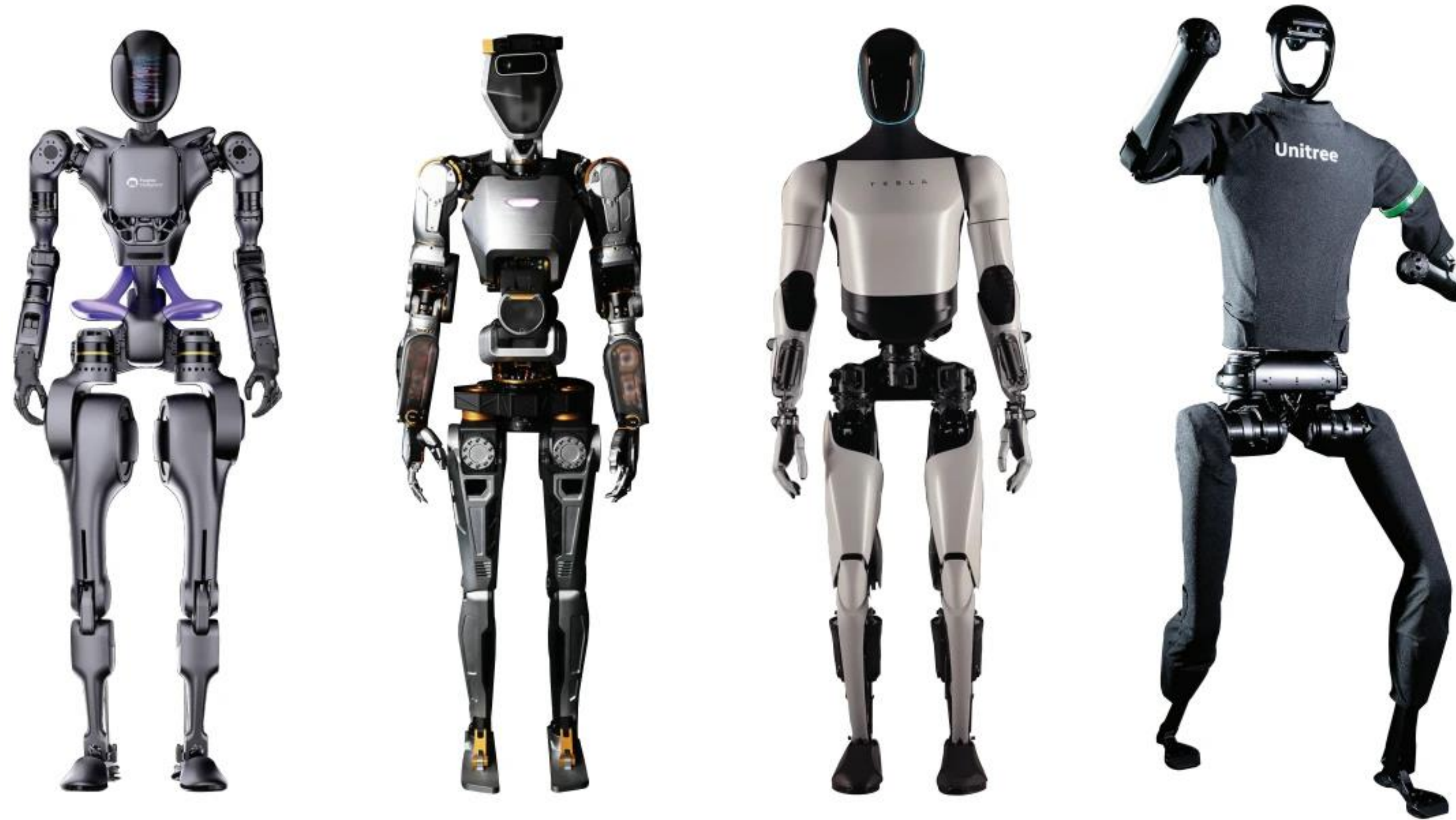# **Data Pyramid** and **Data Flywheel**
# for Robotic Foundation Models
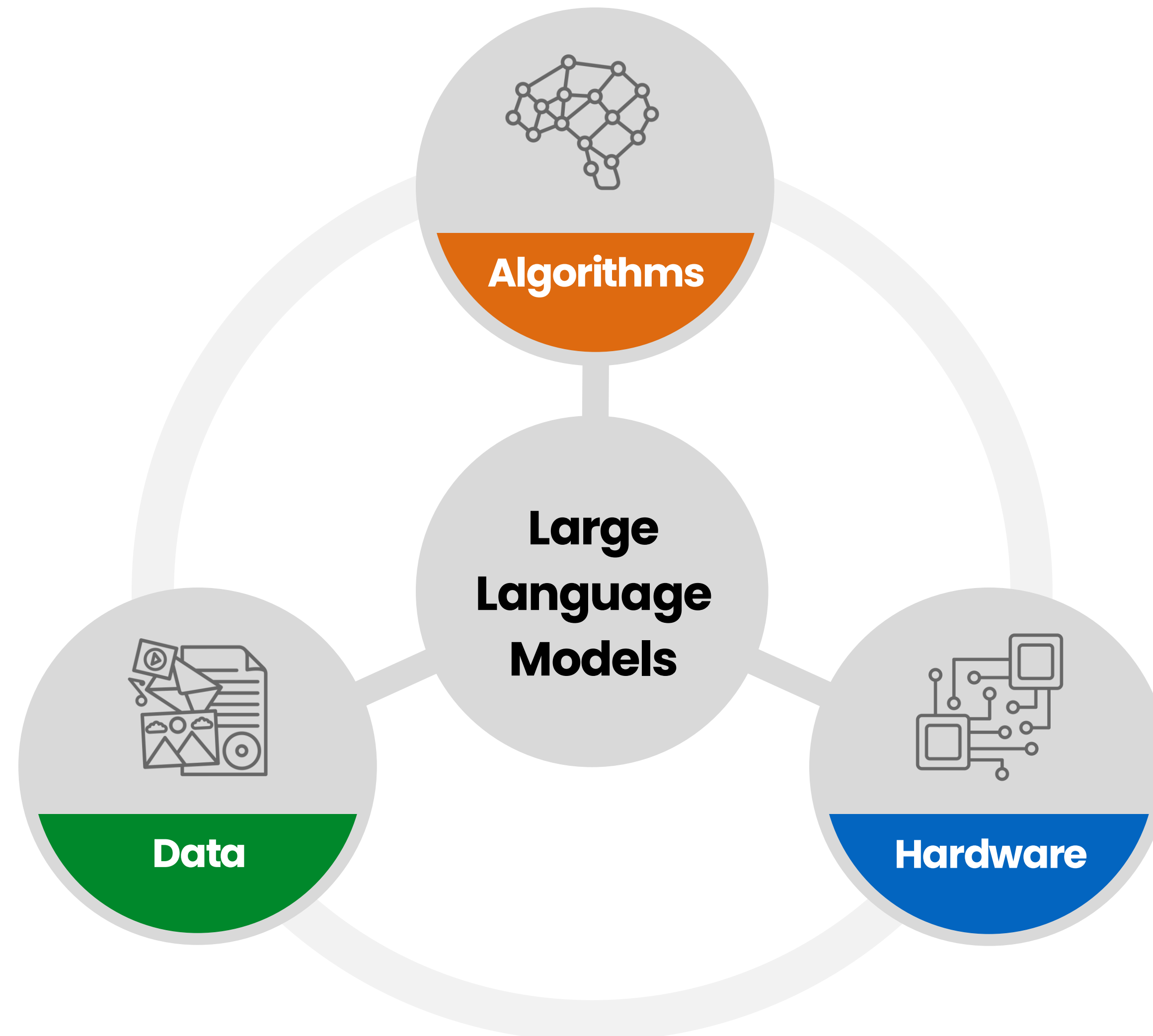
## Yuke Zhu

UT Austin / NVIDIA

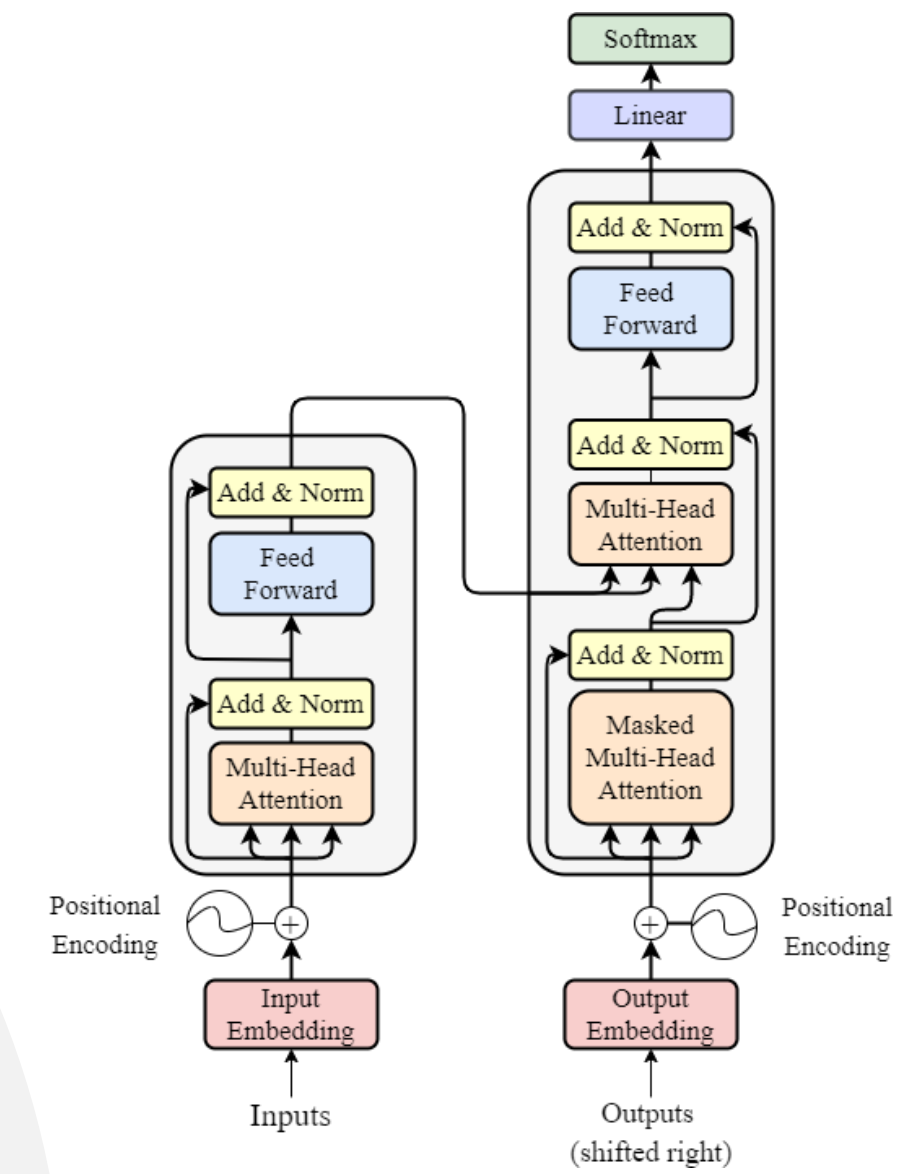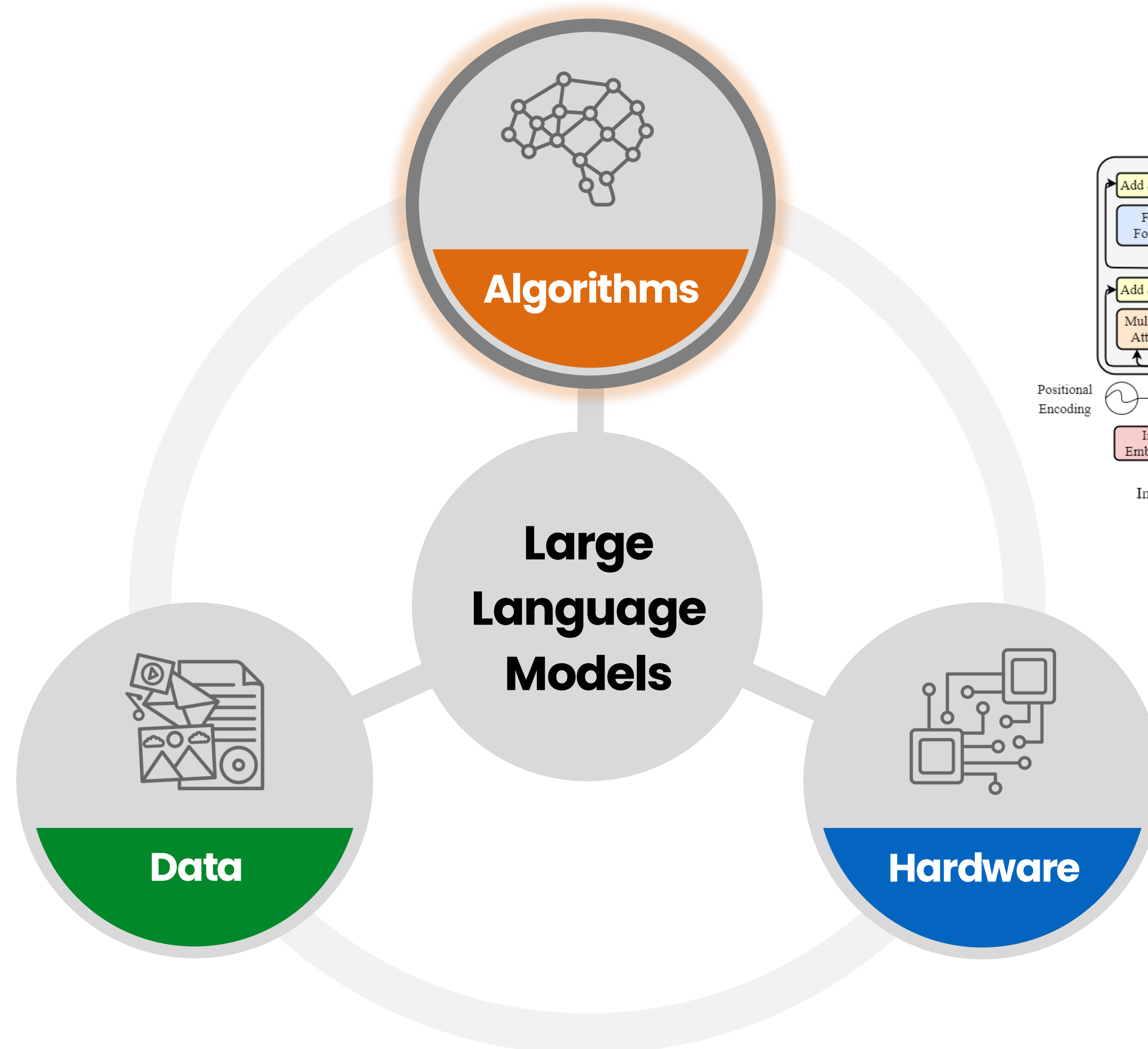# Building Robotic Foundation Models



## One "AI Brain" for All (Humanoid) Robots

# Recipe for Building Large Language Models

# Recipe for Building Large Language Models
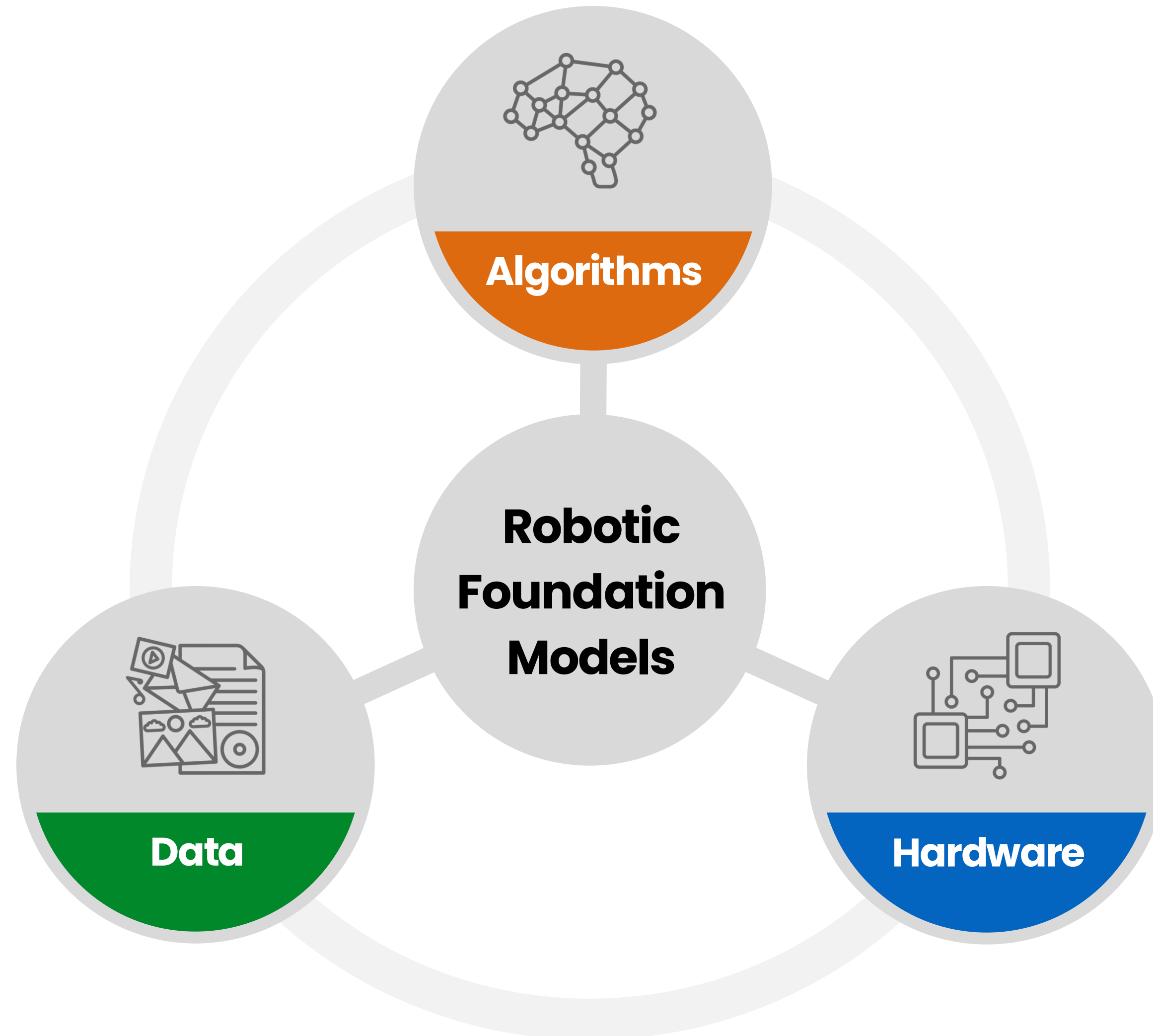
# Recipe for Building Large Language Models

# Recipe for Building Large Language Models

# Recipe for Building Robotic Foundation Models

# Recipe for Building Robotic Foundation Models

**Scalable Algorithms**

Powerful robot learning models

that scale with data and compute

**Algorithms**

**Robotic Foundation Models**

**Data**

**Hardware**

**Data Engine**

New mechanisms to produce

massive training data

**Human-like Embodiment**

Humanoid robot platform for
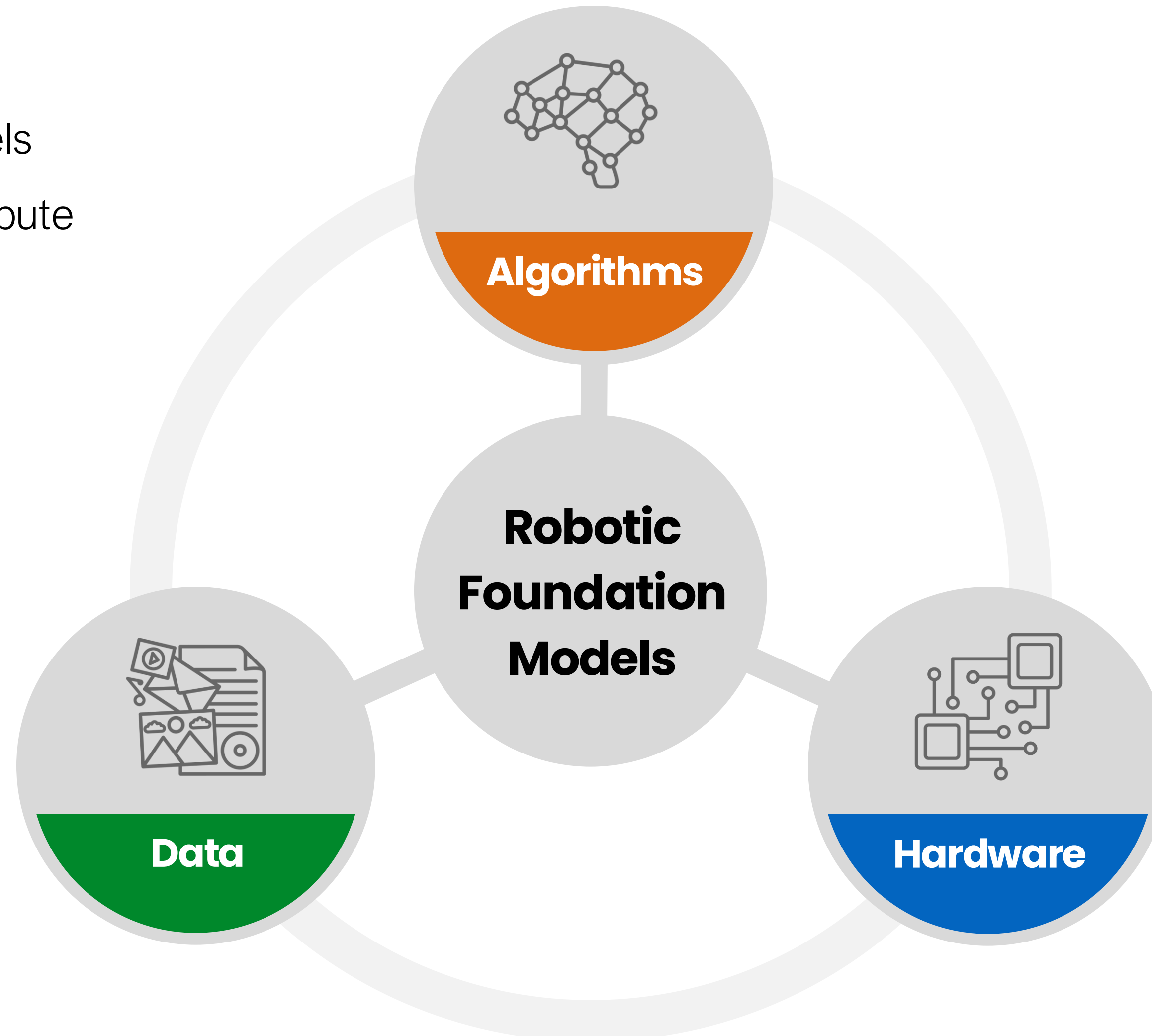
broad applications

# Recipe for Building Robotic Foundation Models

**Scalable Algorithms**

Powerful robot learning models
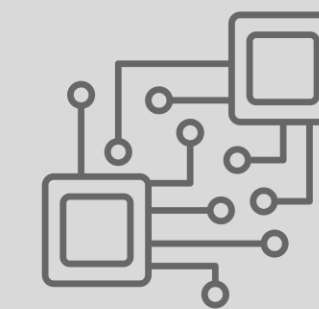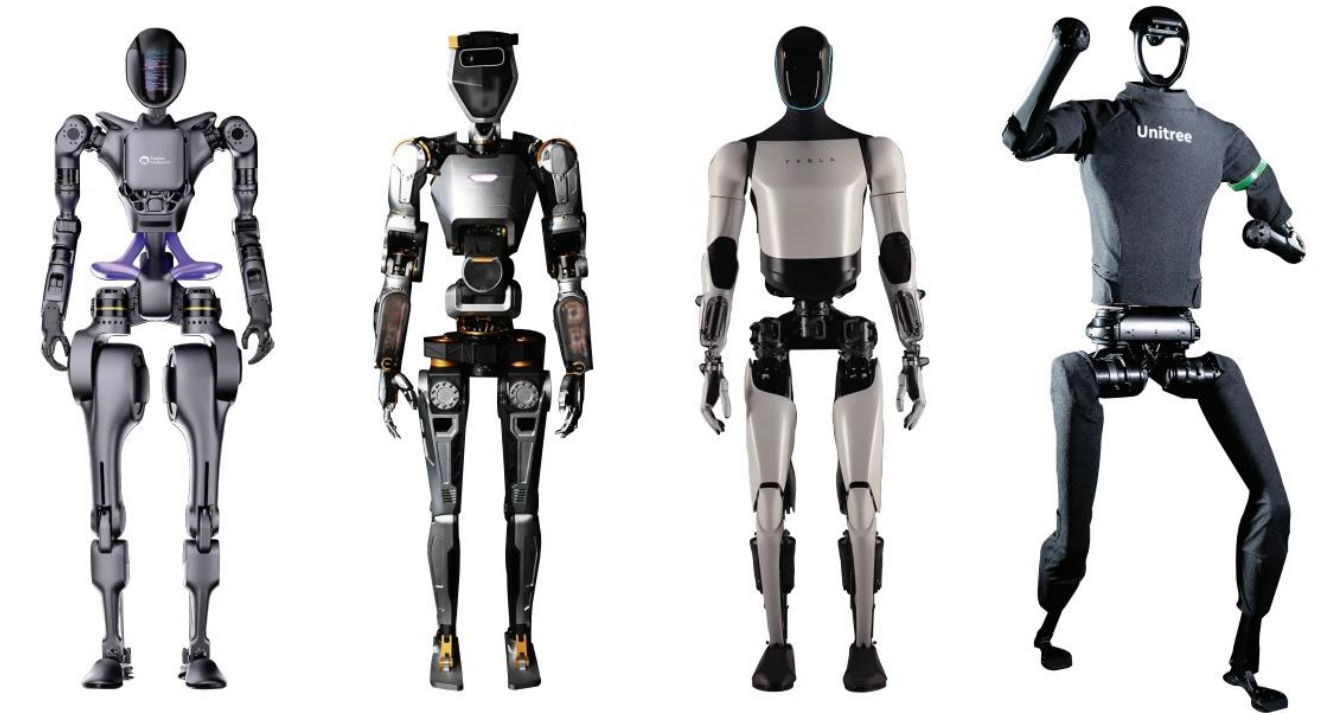
that scale with data and compute



**Data Engine**

New mechanisms to produce

massive training data

**Human-like Embodiment**

Humanoid robot platform for

broad applications

# Why Humanoids?

❖ **Versatility:** General-purpose robot autonomy needs a versatile body.

❖ **Costs:** Hardware becomes cheaper and more robust to democratize transformative research.

❖ **Safety:** Humanoid robots can be more predictable and safer for human-robot interaction.

❖ **Data:** Their similar physique unlocks Internet-scale, human-centered data sources.

❖ …

Research Principle #1:

**First Generalist, then Better Specialist**



semantic parser    sentiment analyzer

text summarizer    information extractor

Large Language Models
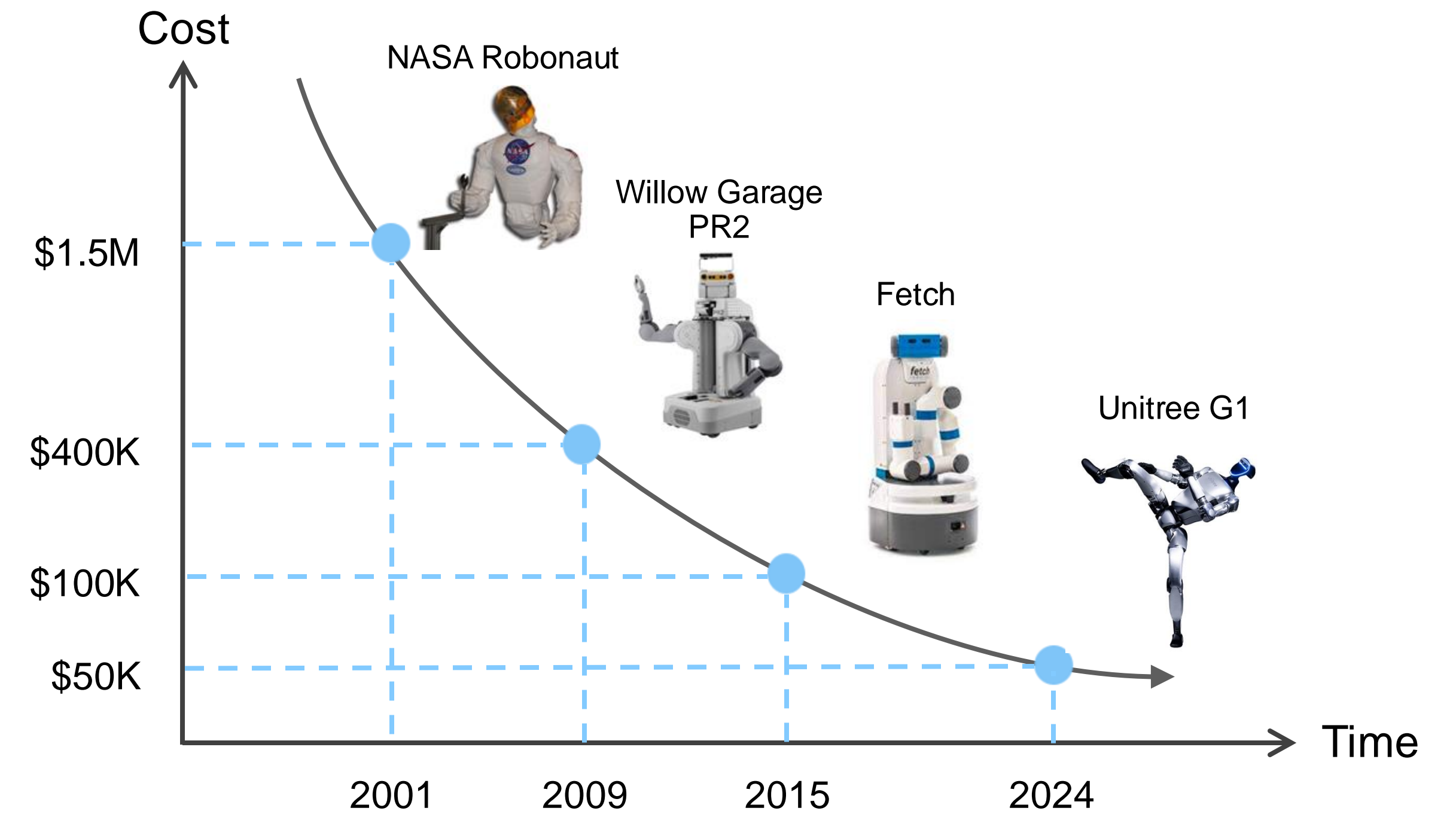(a.k.a. "Foundation Model")

creative writer

coding assistant

travel planner
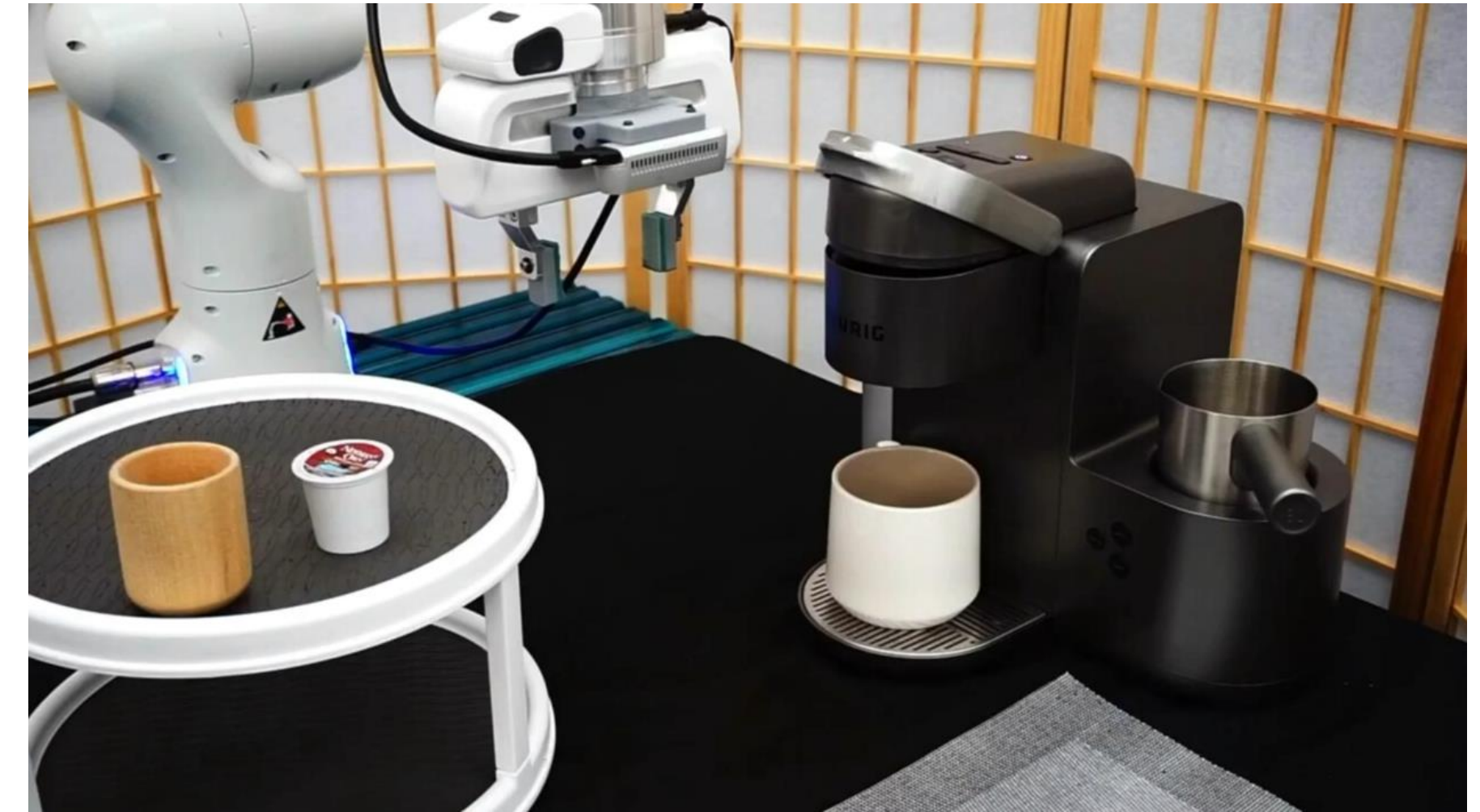
[Credit: Yuke Zhu CoRL'23 Keynote]

# Why Humanoids?

❖ **Versatility:** General-purpose robot autonomy needs a versatile body.

❖ **Costs:** Hardware becomes cheaper and more robust to democratize transformative research.

❖ Safety: Humanoid robots can be more predictable and safer for human-robot interaction.

❖ Data: Their similar physique unlocks Internet-scale, human-centered data sources.

❖ …



Cost

NASA Robonaut

Willow Garage PR2

Fetch

Unitree G1

$1.5M
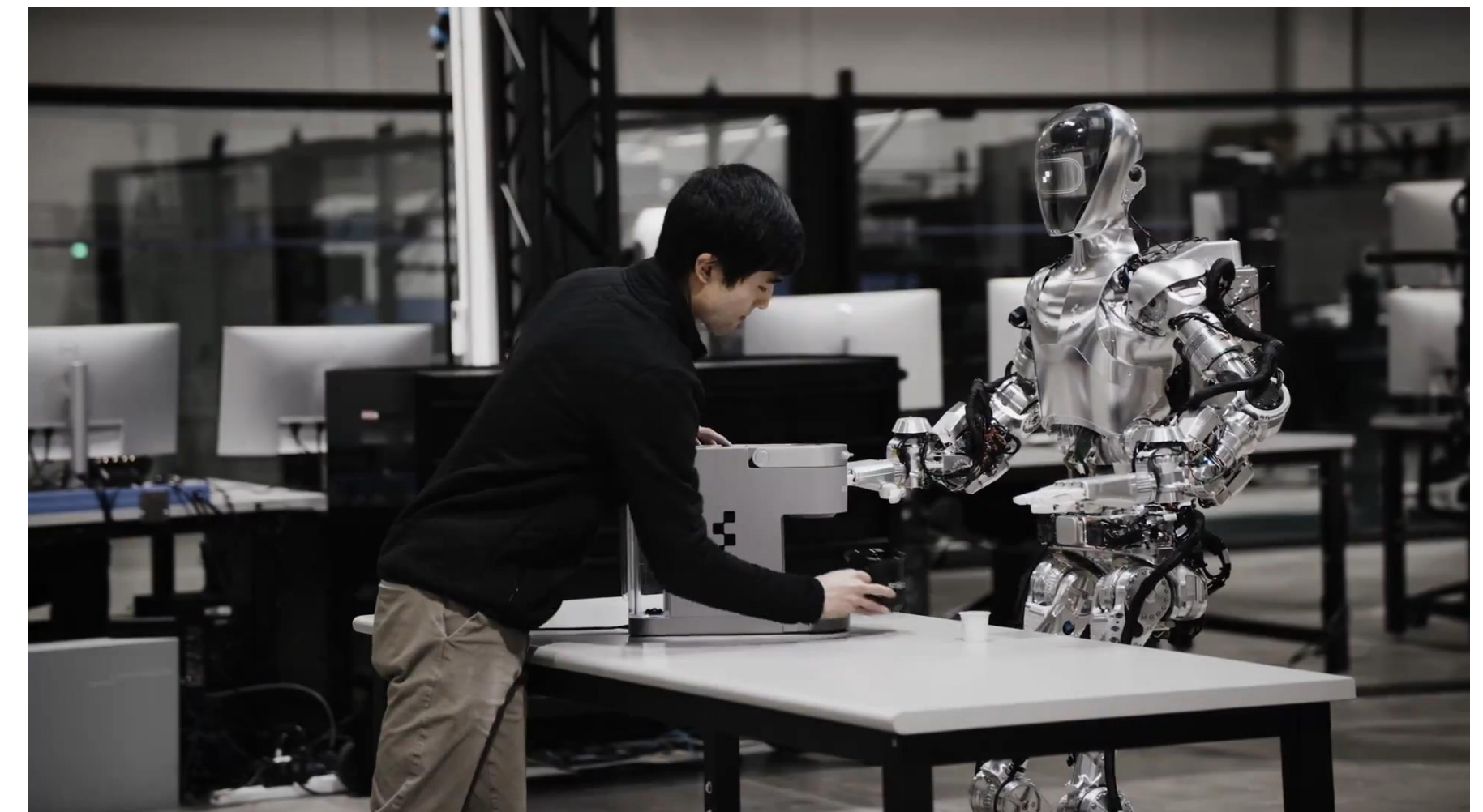
$400K

$100K

$50K

2001    2009    2015    2024

Time

[Credit: Chad Jenkins]

# Why Humanoids?

❖ **Versatility:** General-purpose robot autonomy needs a versatile body.

❖ **Costs:** Hardware becomes cheaper and more robust to democratize transformative research.

❖ **Safety:** Humanoid robots can be more predictable and safer for human-robot interaction.

❖ **Data:** Their similar physique unlocks Internet-scale, human-centered data sources.

❖ ...



[VIOLA, Zhu et al. CoRL 2022]



[Credit: Figure AI 2024]

# Why Humanoids?

❖ **Versatility:** General-purpose robot autonomy needs a versatile body.

❖ **Costs:** Hardware becomes cheaper and more robust to democratize transformative research.

❖ **Safety:** Humanoid robots can be more predictable and safer for human-robot interaction.

❖ **Data:** Their similar physique unlocks Internet-scale, human-centered data sources.

❖ …



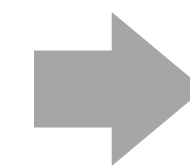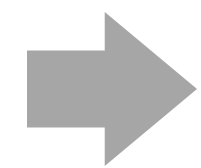Ego4D dataset, CVPR 2022

AMASS dataset, ICCV 2019

# Why Humanoids?

❖ **Versatility:** General-purpose robot autonomy needs a versatile body.

❖ **Costs:** Hardware becomes cheaper and more robust to democratize transformative research.

❖ **Safety:** Humanoid robots can be more predictable and safer for human-robot interaction.

❖ **Data:** Their similar physique unlocks Internet-scale, human-centered data sources.

❖ …



Ego4D dataset, CVPR 2022

AMASS dataset, ICCV 2019

**Note:** humanoid robotics is still incredibly hard (!) — huge challenges in mechanical designs, dynamics & control, sensor technologies, compute and power, AI algorithm designs…
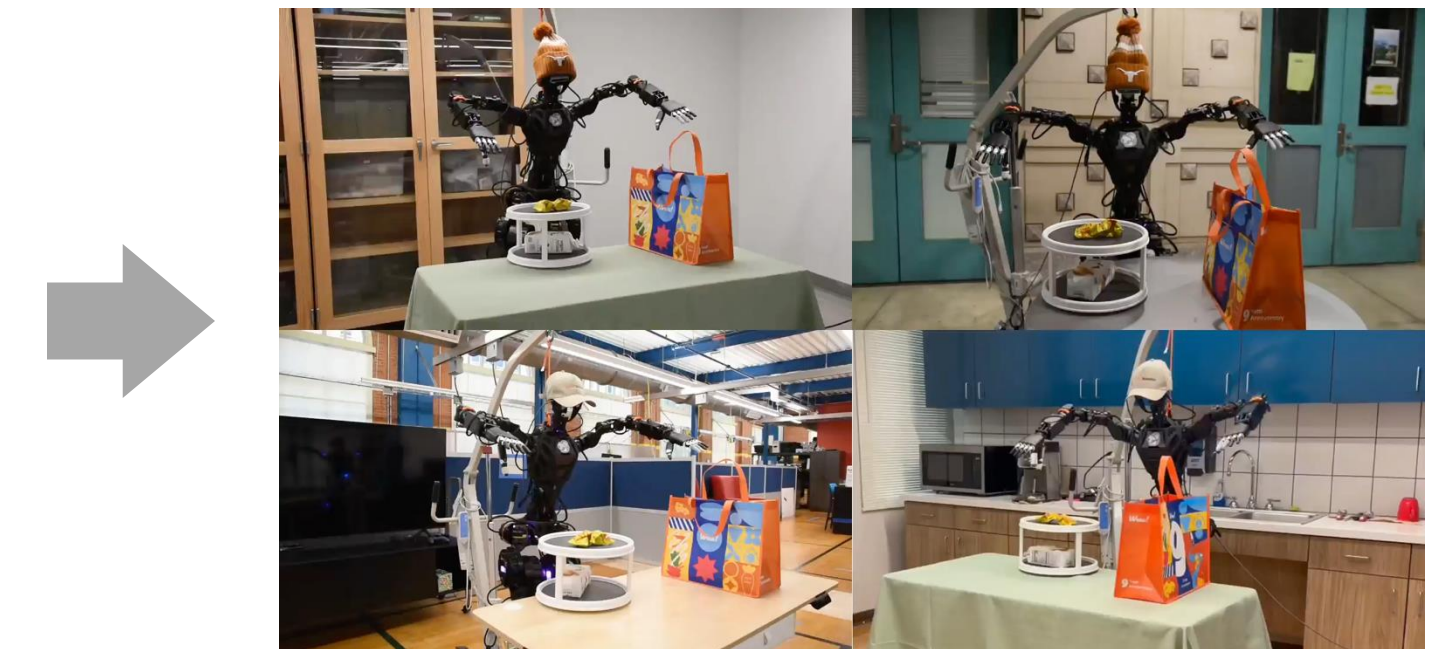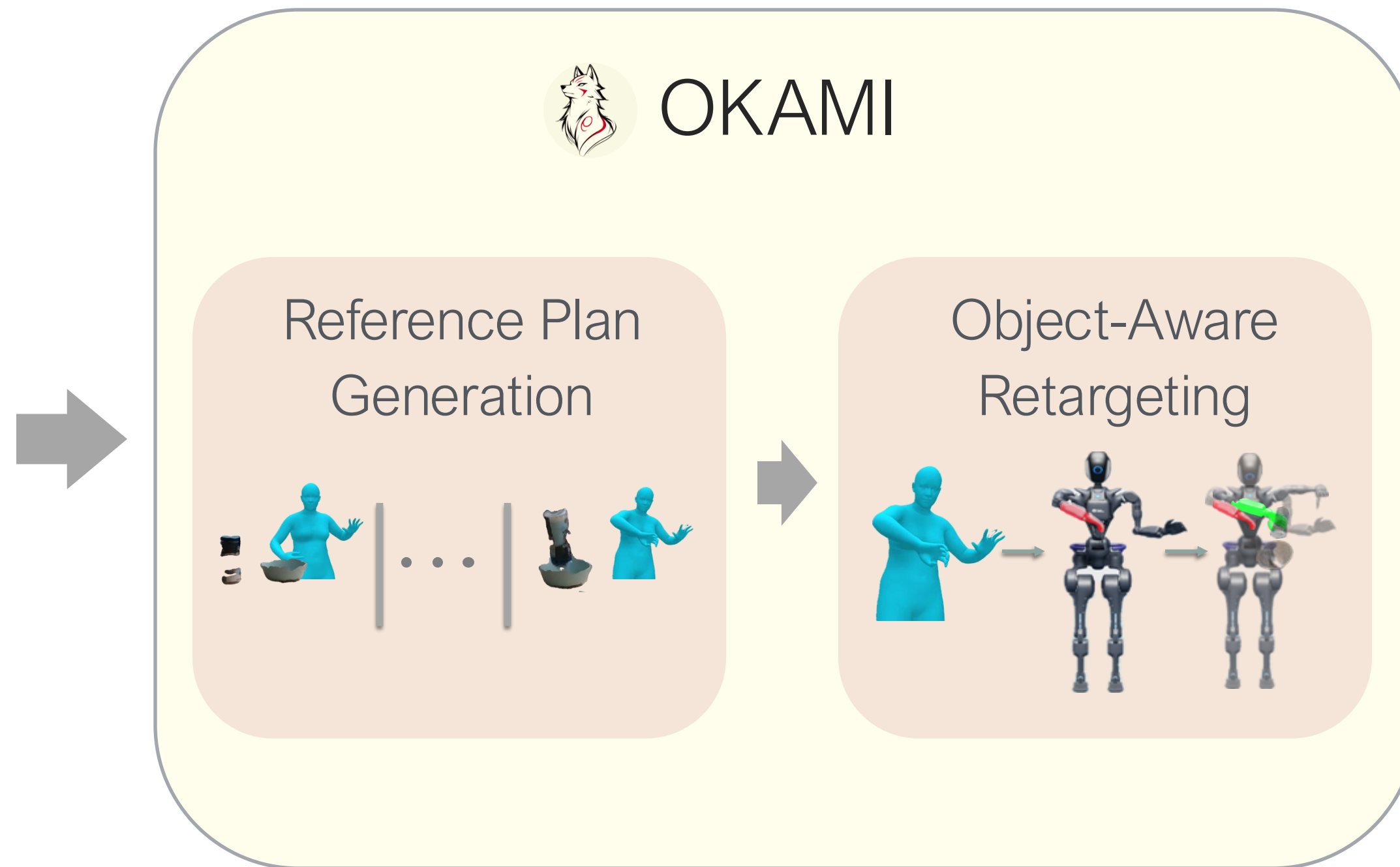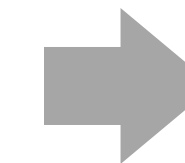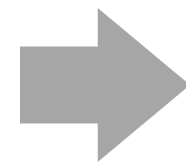
# Learning from Human Videos



single video
demonstration

OKAMI

trajectory rollouts
in diverse scenes

# Learning from Human Videos



single video
demonstration

OKAMI

Reference Plan
Generation

Object-Aware
Retargeting

trajectory rollouts
in diverse scenes

"OKAMI: Teaching Humanoid Robots Manipulation Skills through Single Video Imitation." Li et al. CoRL 2024

# Learning from Human Videos



single video
demonstration

OKAMI

REC

trajectory rollouts
in diverse scenes

# Learning from Human Videos



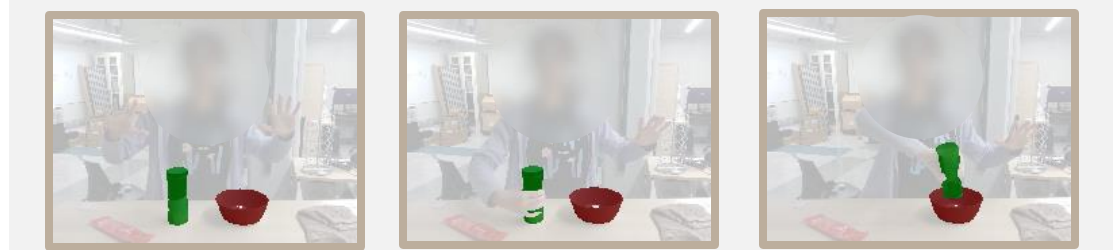OKAMI

Reference Plan Generation
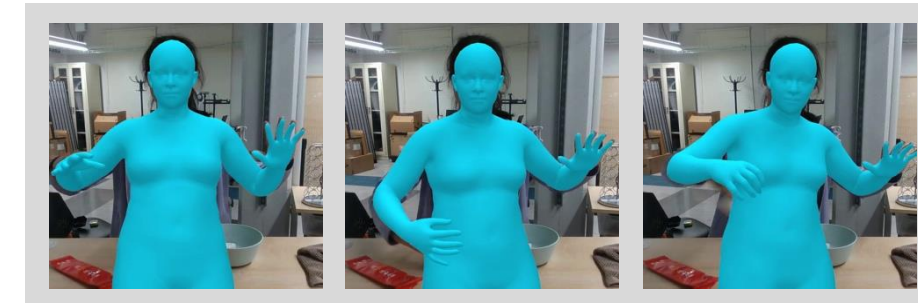
Demonstration Video

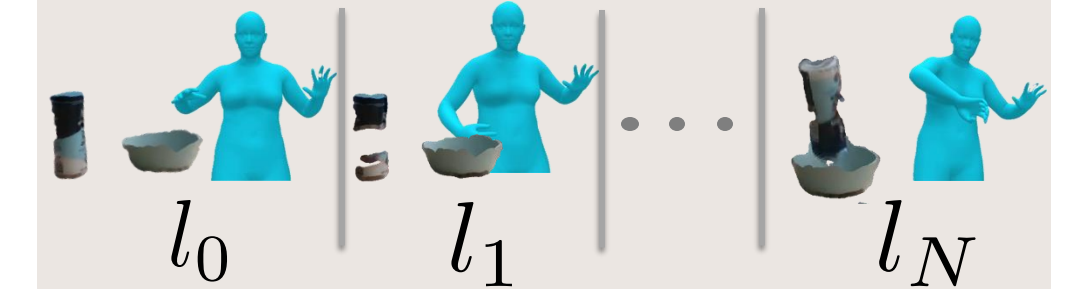GPT-4V

"bottle"
"bowl"

Track objects across the video

Identify keyframes through changepoint detection
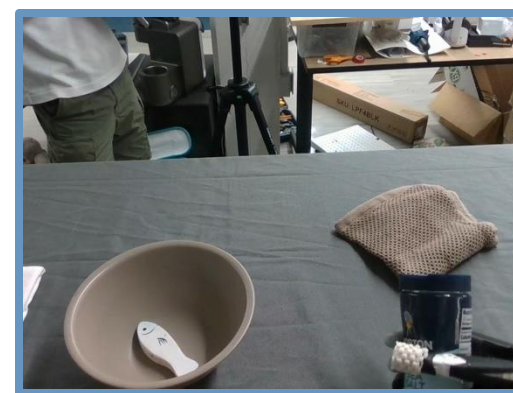
Human Reconstruction Model

SMPL-H trajectory

Reference Plan

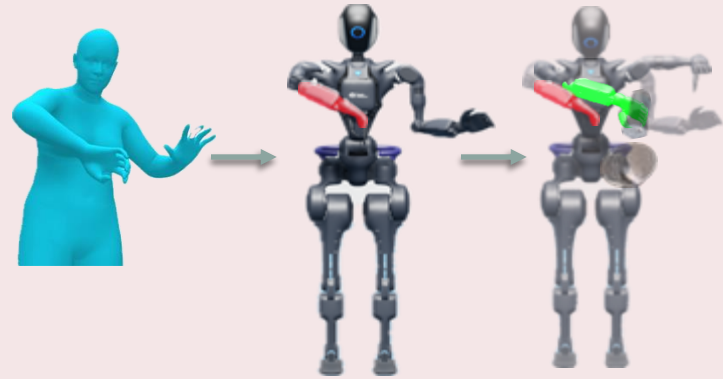$l_0$    $l_1$    $\cdots$    $l_N$

"OKAMI: Teaching Humanoid Robots Manipulation Skills through Single Video Imitation." Li et al. CoRL 2024

# Learning from Human Videos



OKAMI

Object-Aware Retargeting

Reference Plan

$$l_0 \quad l_1 \quad \cdots \quad l_N$$

Robot Observation

Localize relevant objects at test time

Estimate transformation between point clouds
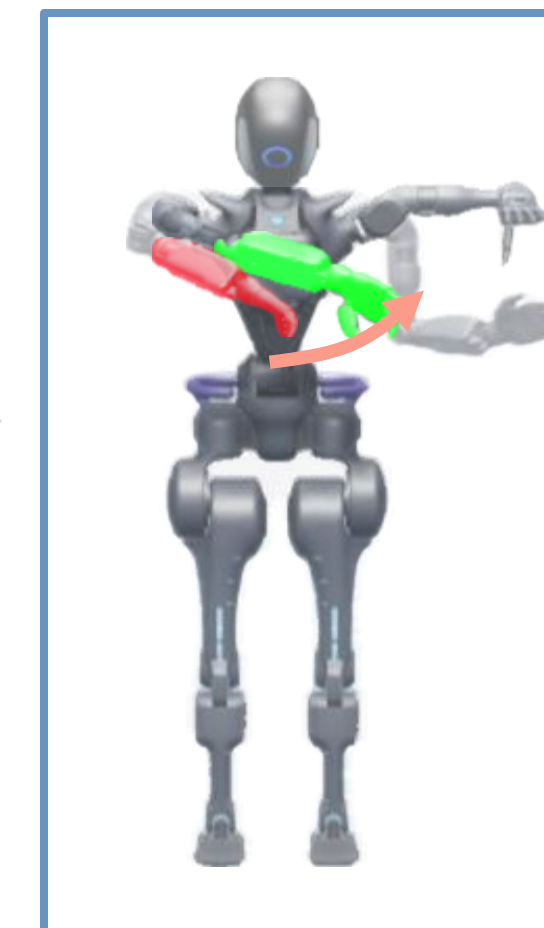
Target and reference objects

SMPL-H trajectory segment

Retarget motions Using SMPL-H

Warped motions

Robot execution

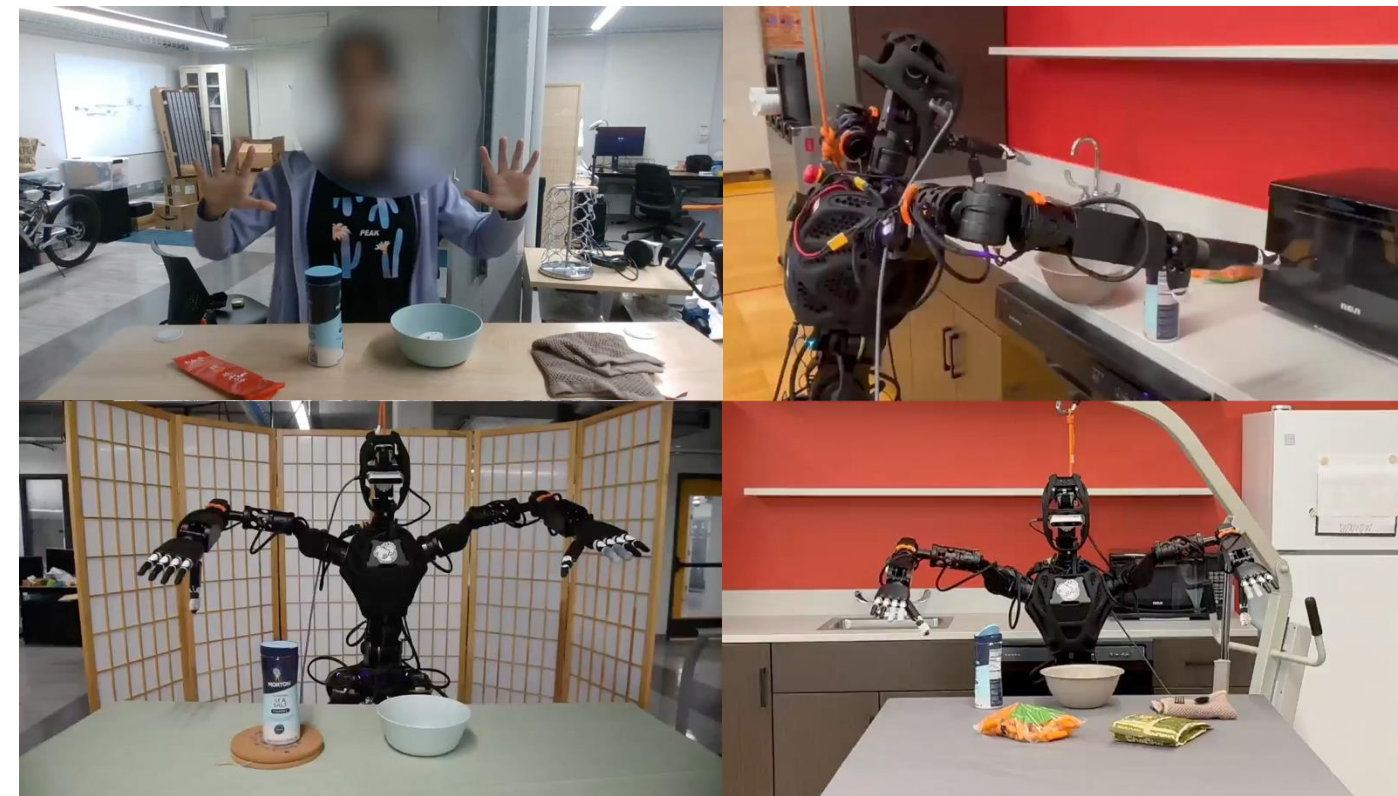"OKAMI: Teaching Humanoid Robots Manipulation Skills through Single Video Imitation." Li et al. CoRL 2024

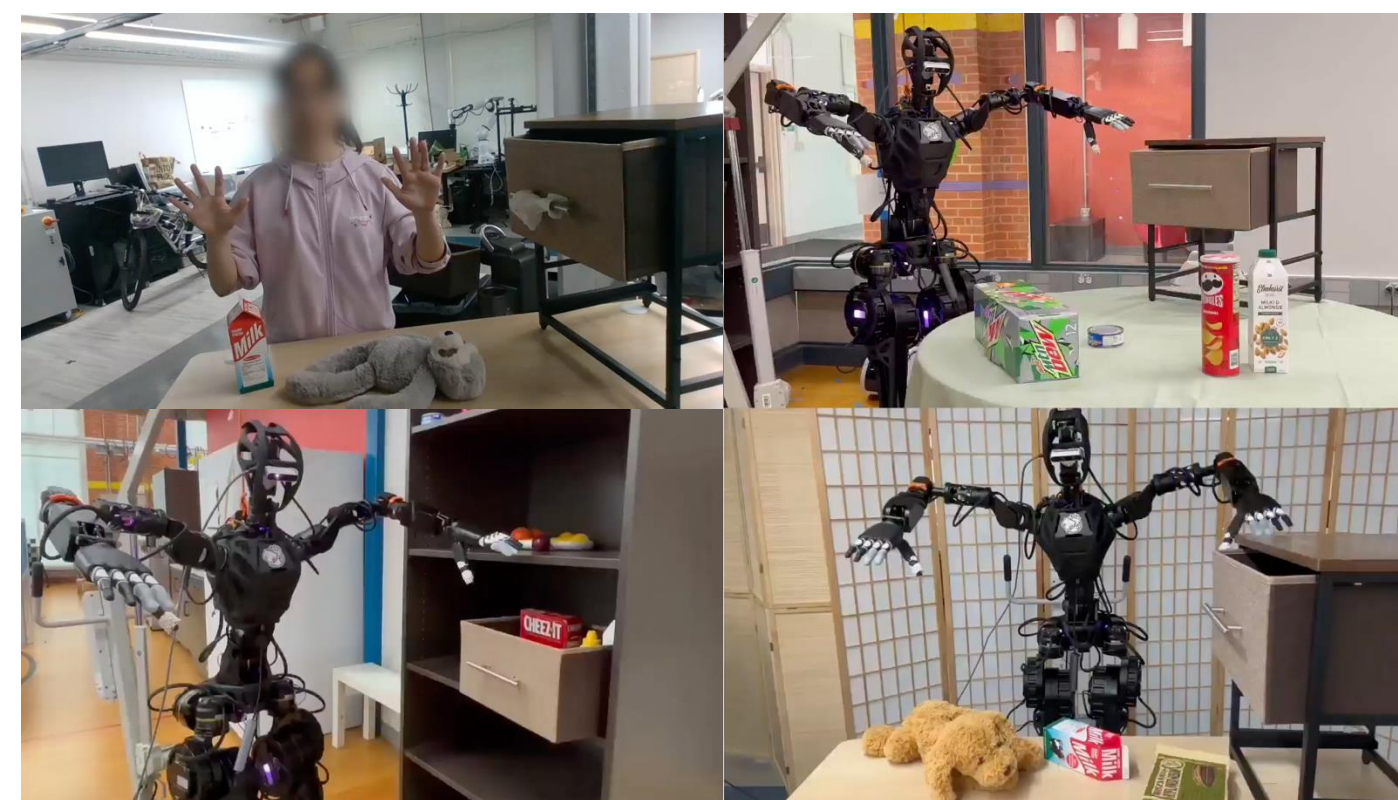# Learning from Human Videos

bagging (58.3%)

sprinkling salt (58.3%)

putting toy in basket (66.7%)



placing snacks on plate (75.0%)

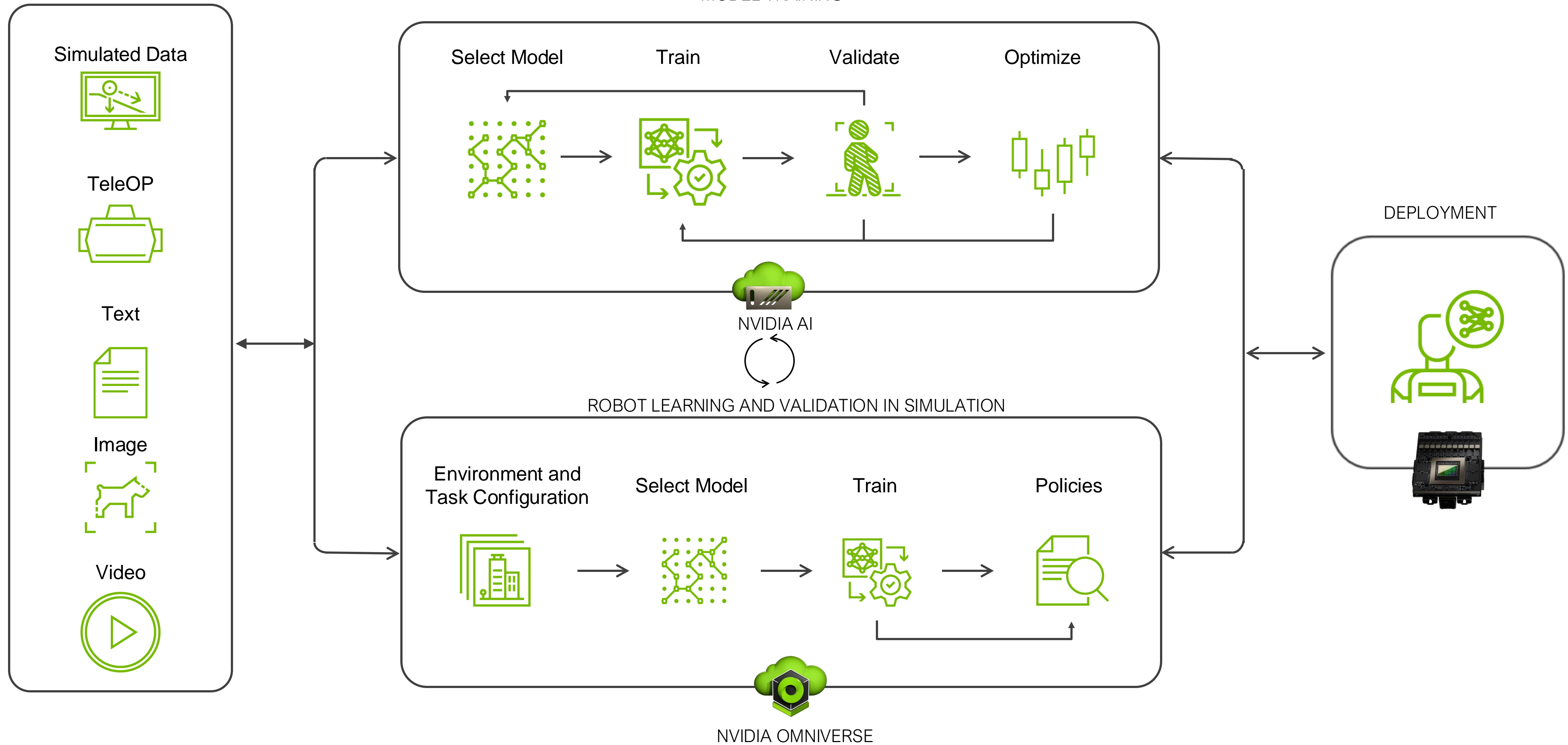closing the drawer (75.0%)

closing the laptop (83.3%)

"OKAMI: Teaching Humanoid Robots Manipulation Skills through Single Video Imitation." Li et al. CoRL 2024

# Project GR00T



DATA PROCESSING AND GENERATION

Simulated Data

TeleOP

Text

Image

Video

MODEL TRAINING

Select Model → Train → Validate → Optimize

NVIDIA AI

ROBOT LEARNING AND VALIDATION IN SIMULATION

Environment and Task Configuration → Select Model → Train → Policies
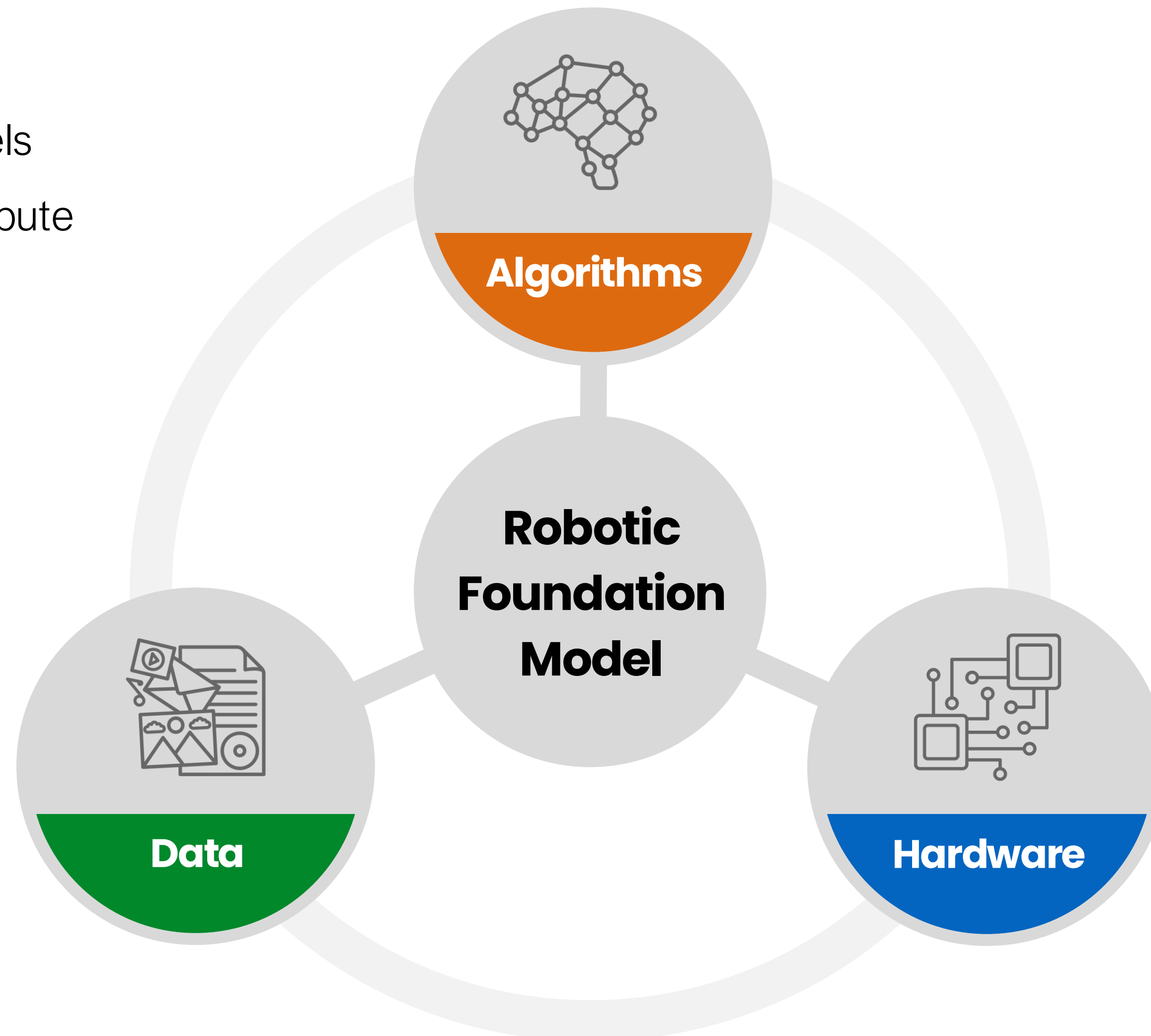
NVIDIA OMNIVERSE

DEPLOYMENT

# Recipe for Building Robotic Foundation Models

**Scalable Algorithms**

Powerful robot learning models

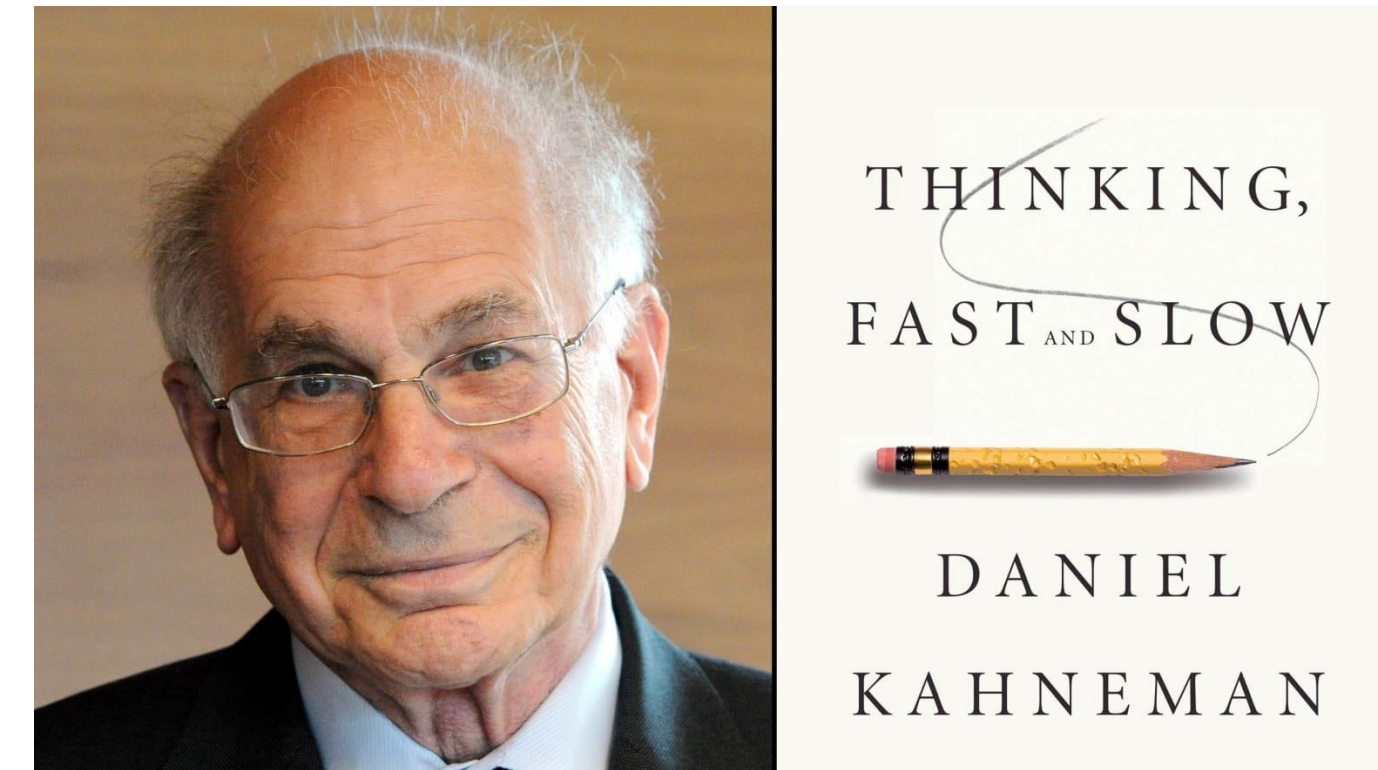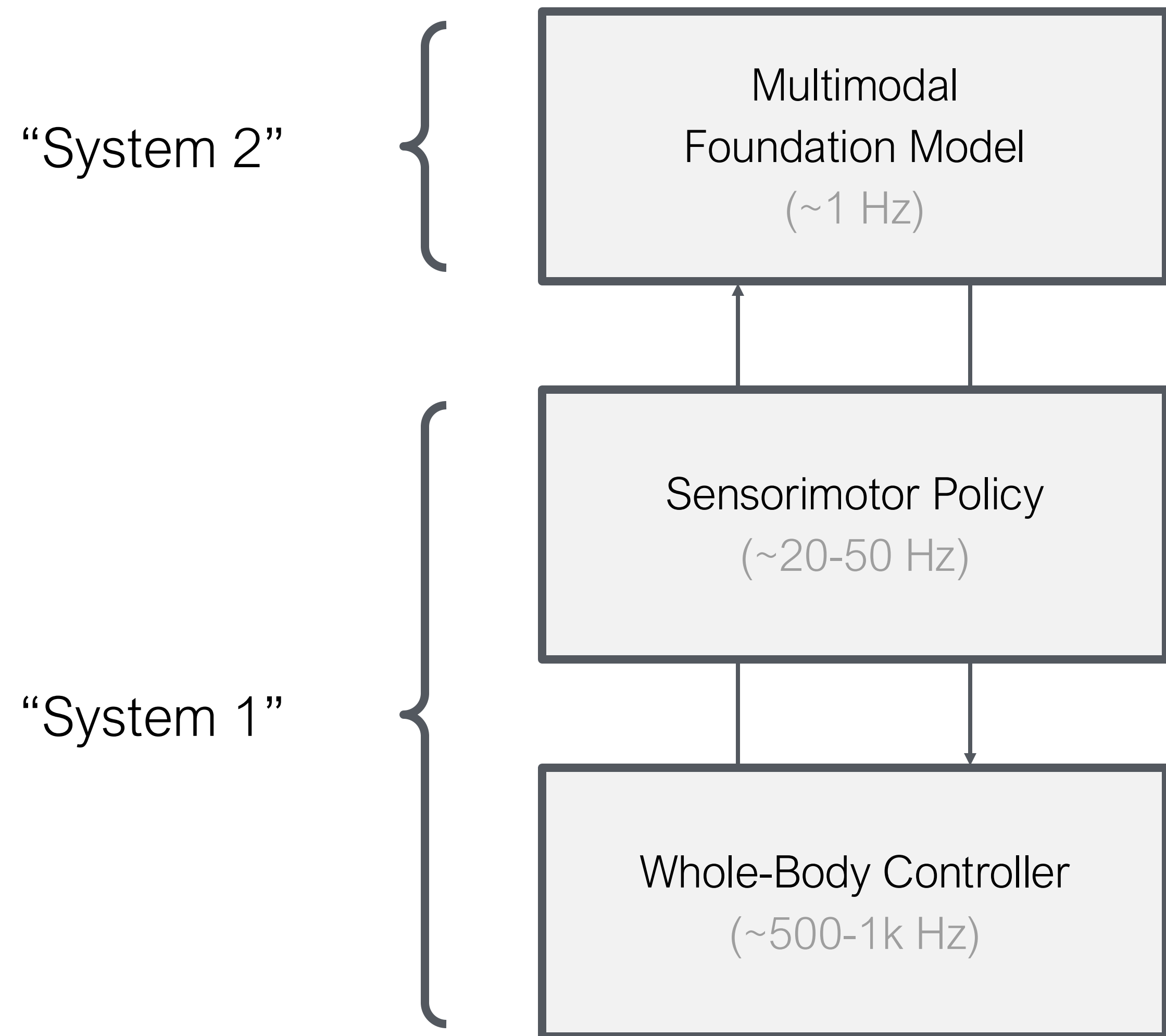that scale with data and compute



**Data Engine**

New mechanisms to produce

massive training data

**Human-like Embodiment**

Humanoid robot platform for

broad applications

# Hierarchical Autonomy Stack: System 1-System 2

# The Data Pyramid for Generalist Robots



Web Data

- Massive scale and ever-growing
- Multimodal and unstructured
- Human-centered data

# The Data Pyramid for Generalist Robots

The "Cambrian explosion" of Vision-Language Models

INTRODUCING
Lightweight and multimodal Llama models

ON-DEVICE 3B
ON-DEVICE 1B
MULTIMODAL 90B
MULTIMODAL 11B

Gemini

✳ Claude 3

ChatGPT-4o

Pixtral 12B
MISTRAL AI_

Qwen2

LONG VILA

[Xue et al. 2024]

Web Data

YouTube

reddit

Common Crawl

WIKIPEDIA
The Free Encyclopedia

- Massive scale and ever-growing
- Multimodal and unstructured
- Human-centered data

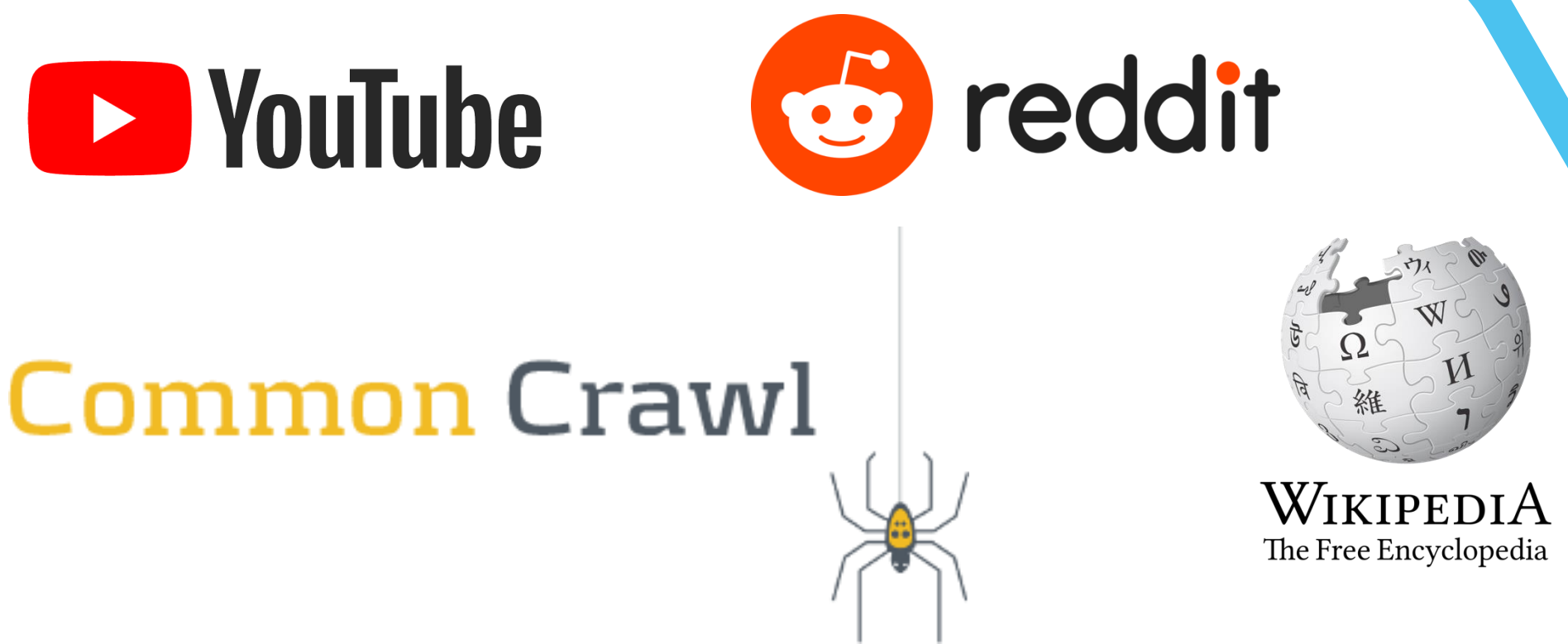📍 2501 Wichita St.    📍 2500 Speedway    📍 100 E 24th St.

# The Data Pyramid for Generalist Robots



**Synthetic Data**

- Unlimited simulated data (in theory)
- Content creation challenge, reality gap, computational burden

**Web Data**

- Massive scale and ever-growing
- Multimodal and unstructured
- Human-centered data

# The Data Pyramid for Generalist Robots



**Real-World Data**

- Small scale and expensive to collect
- Ease of use for imitation learning, direct transfer

**Synthetic Data**

- Unlimited simulated data (in theory)
- Content creation challenge, reality gap, computational burden

**Web Data**

- Massive scale and ever-growing
- Multimodal and unstructured
- Human-centered data
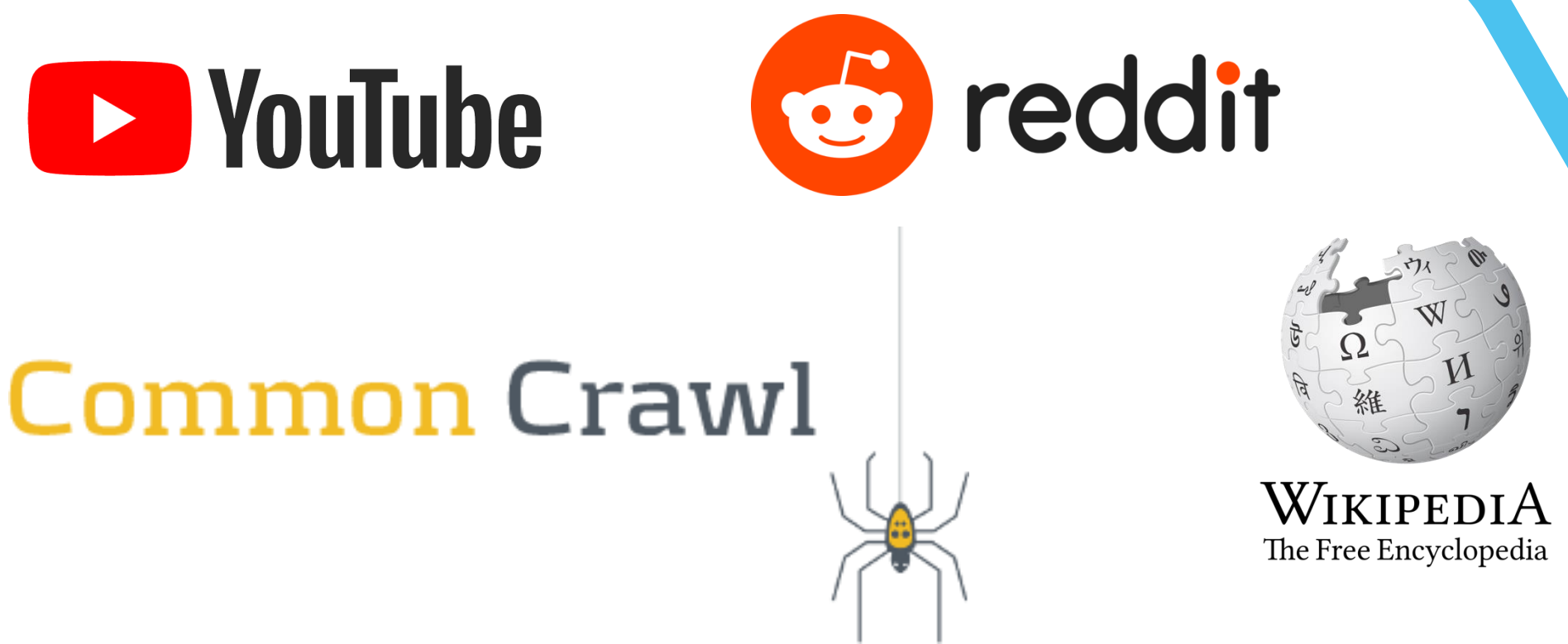
# The Data Pyramid for Generalist Robots



Real-World Data

Synthetic Data

Web Data

Research Principle #2:

**Learning Across the Data Pyramid**

# The Data Pyramid for Generalist Robots



**Real-World Data**

**Synthetic Data**

**Web Data**

▶ YouTube

reddit

Common Crawl

WIKIPEDIA
The Free Encyclopedia

real-time teleoperation

Data grows **linearly** with respect to time, money, human efforts, …

# The Data Pyramid for Generalist Robots



Real-World Data

Synthetic Data

Web Data

YouTube  reddit

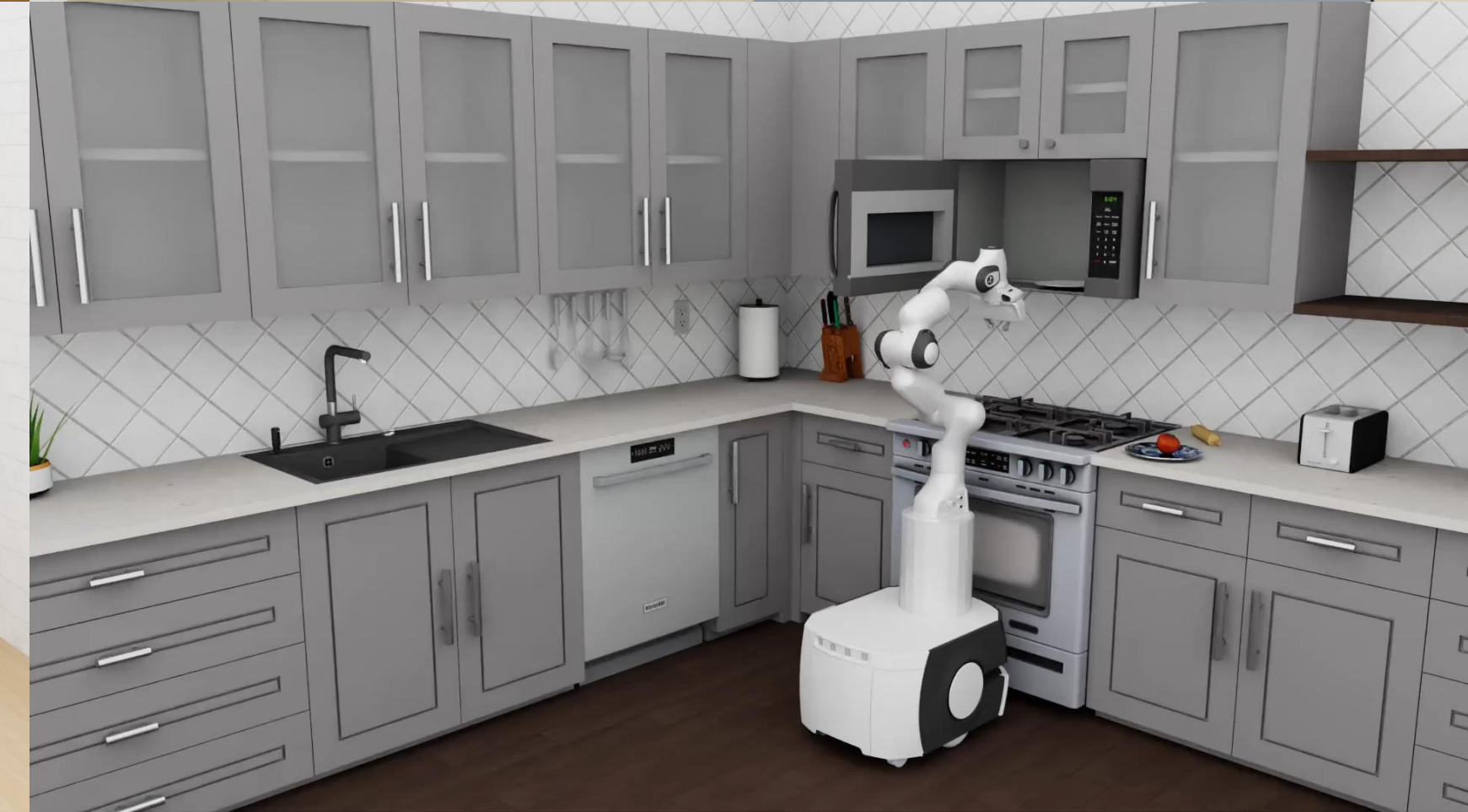Common Crawl  WIKIPEDIA The Free Encyclopedia

real-time teleoperation (Tesla)

Data grows **linearly** with respect to time, money, human efforts, …

# The Data Pyramid for Generalist Robots



Real-World Data

Synthetic Data

Web Data

synthetic data generation

Data grows **exponentially** with automated generation in simulation.

# RoboCasa

Large-Scale Simulation of Everyday Tasks for Generalist Robots

Creating diverse object assets with text-to-3D models

Interactable Furniture and Appliances
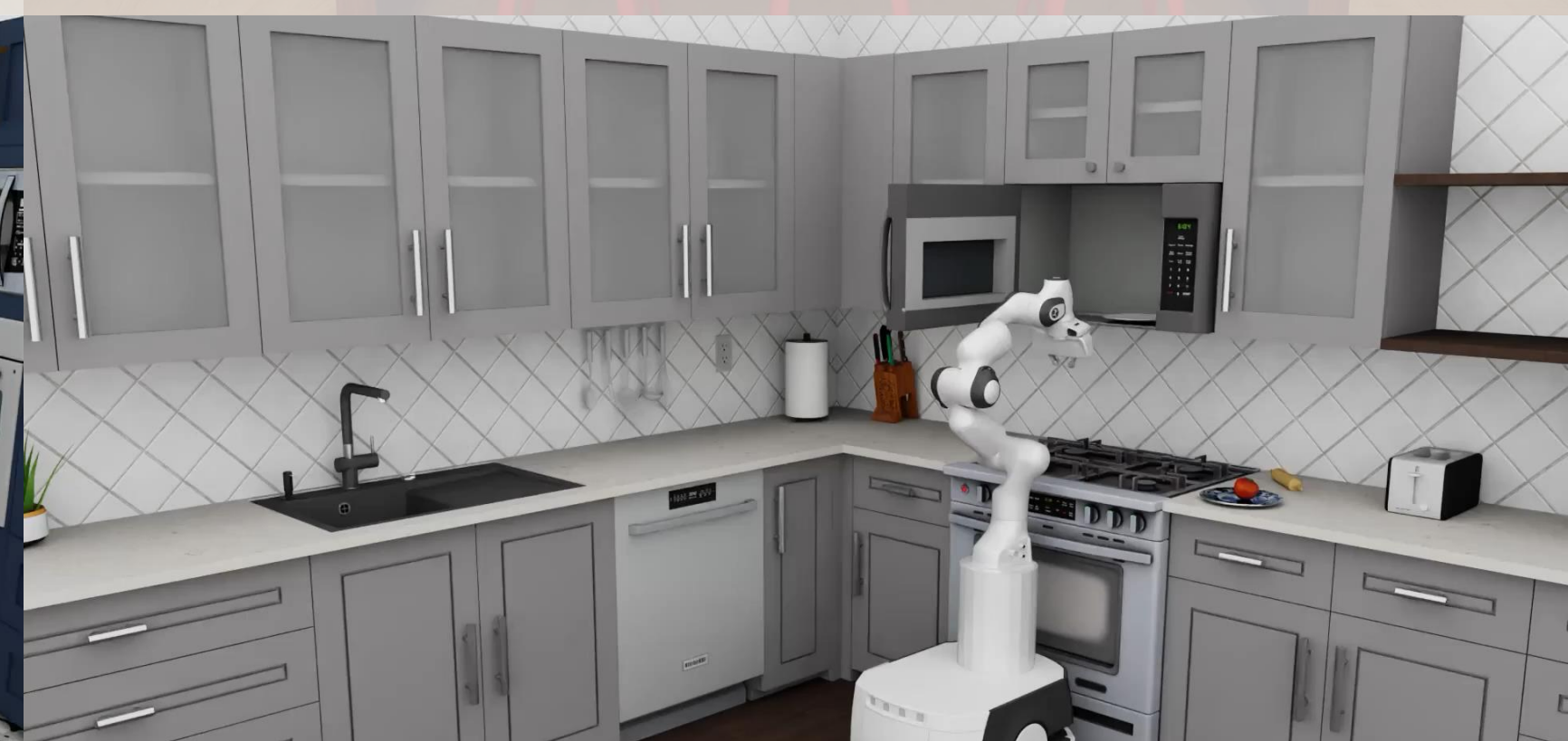
Farmhouse

Modern

Industrial

Rustic

Scandinavian

Traditional

Traditional

Coastal

Transitional

# RoboCasa: Generative Robotic Simulation

## Diverse tasks generated with LLM guidance

**Activity Prompting**
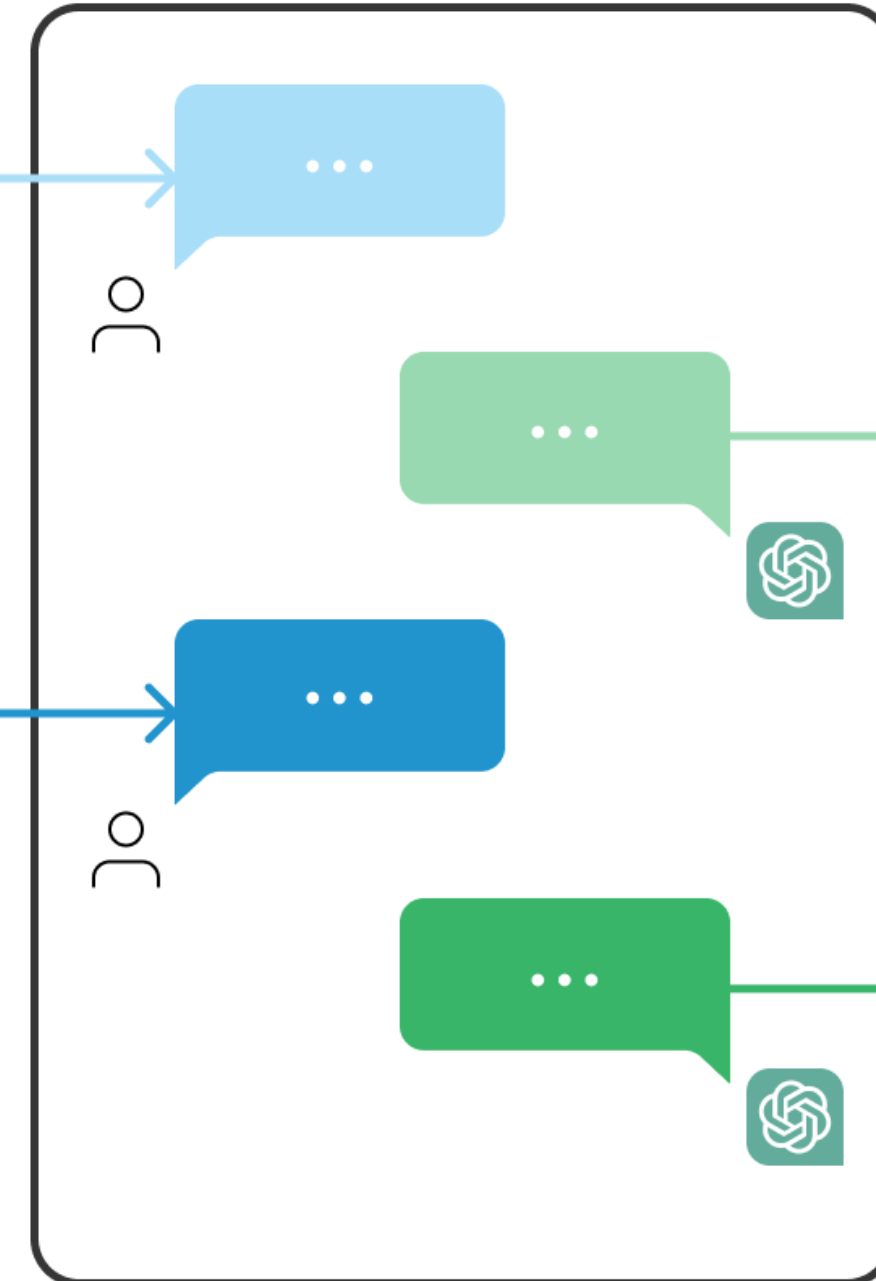
Can you give me 30 simple everyday kitchen activities?

**Task Prompting**

Your goal is to come up with 15 unique tasks that a robot can complete that all fall under {ACTIVITY FROM GPT}.

Available objects and skills:
…

Example tasks:
…

**GPT-4**

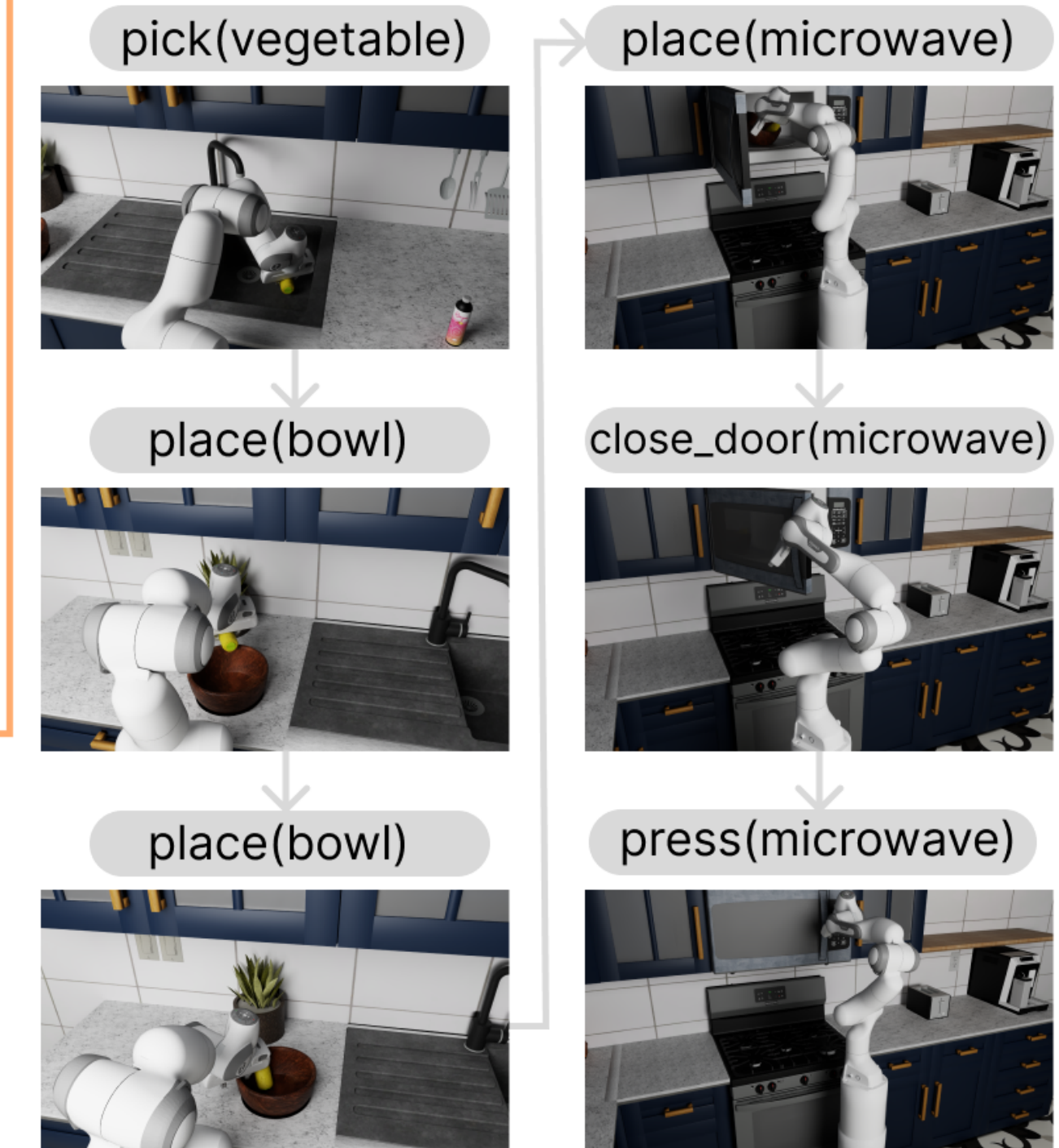**List of activities**
1. Chopping Food
2. Frying
3. Serving Food …

**Task:** Prepare Microwave Steaming
**Goal:** Put a bowl of vegetables inside the microwave to steam them there.
**Objects:** bowl, vegetables
**Fixtures:** sink, microwave
**Skills** (6):
   1. pick(vegetable)
   2. place(bowl)
   3. pick(bowl)
   4. place(microwave)
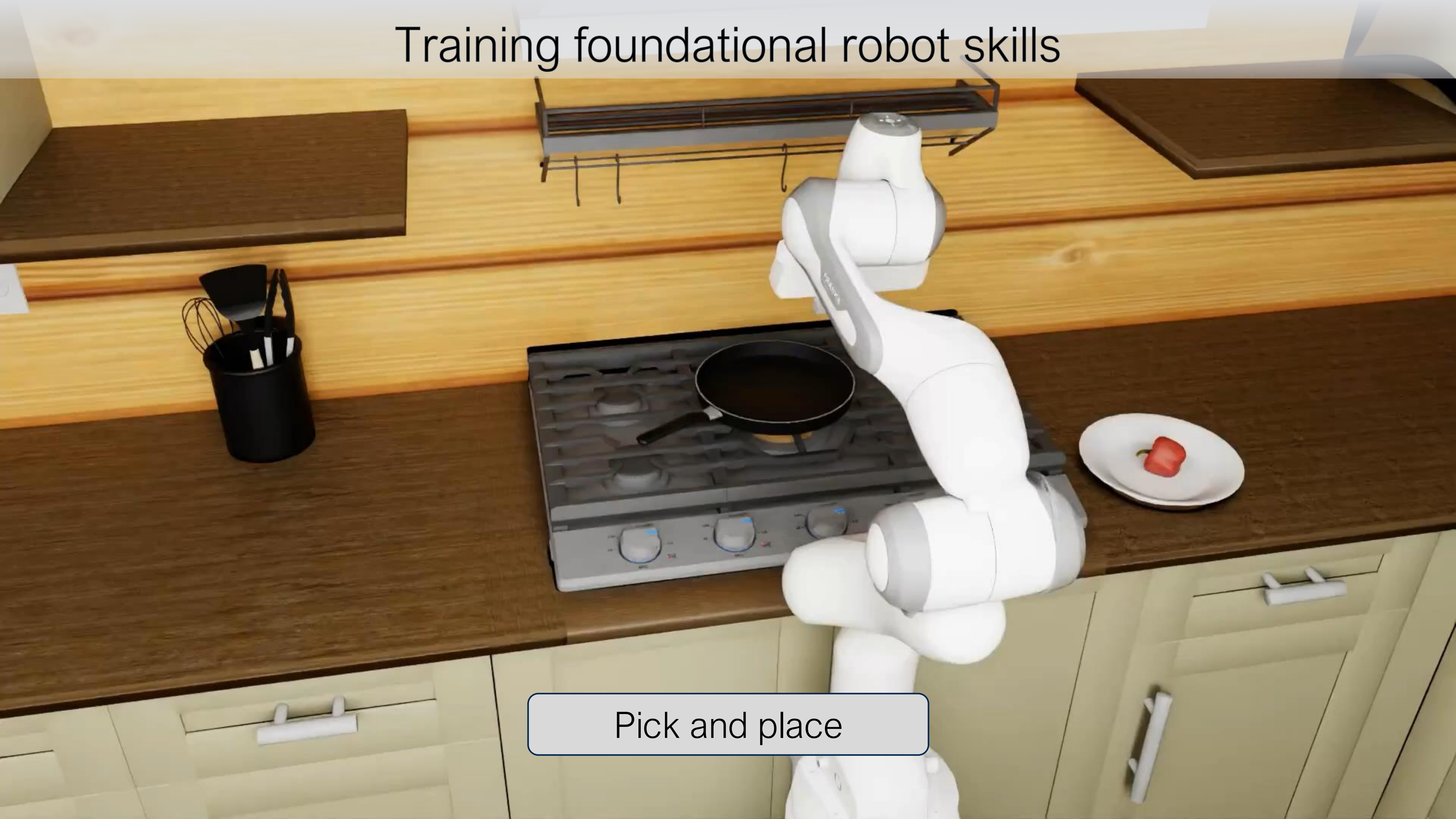   5. close_door(microwave)
   6. press(microwave)

**Task Generation Process**



pick(vegetable) → place(microwave)

place(bowl) → close_door(microwave)
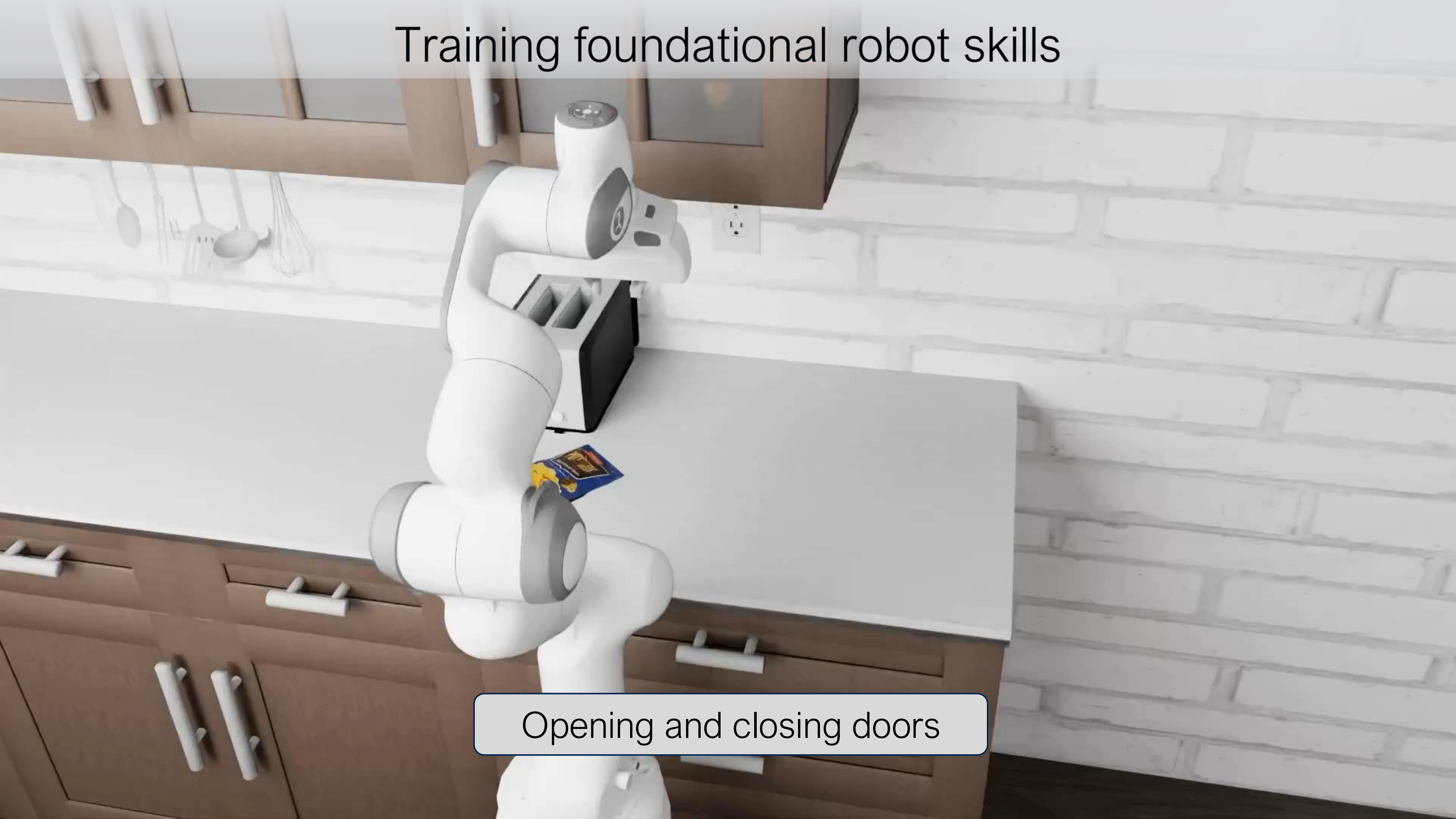
place(bowl) → press(microwave)

Cross-embodiment support
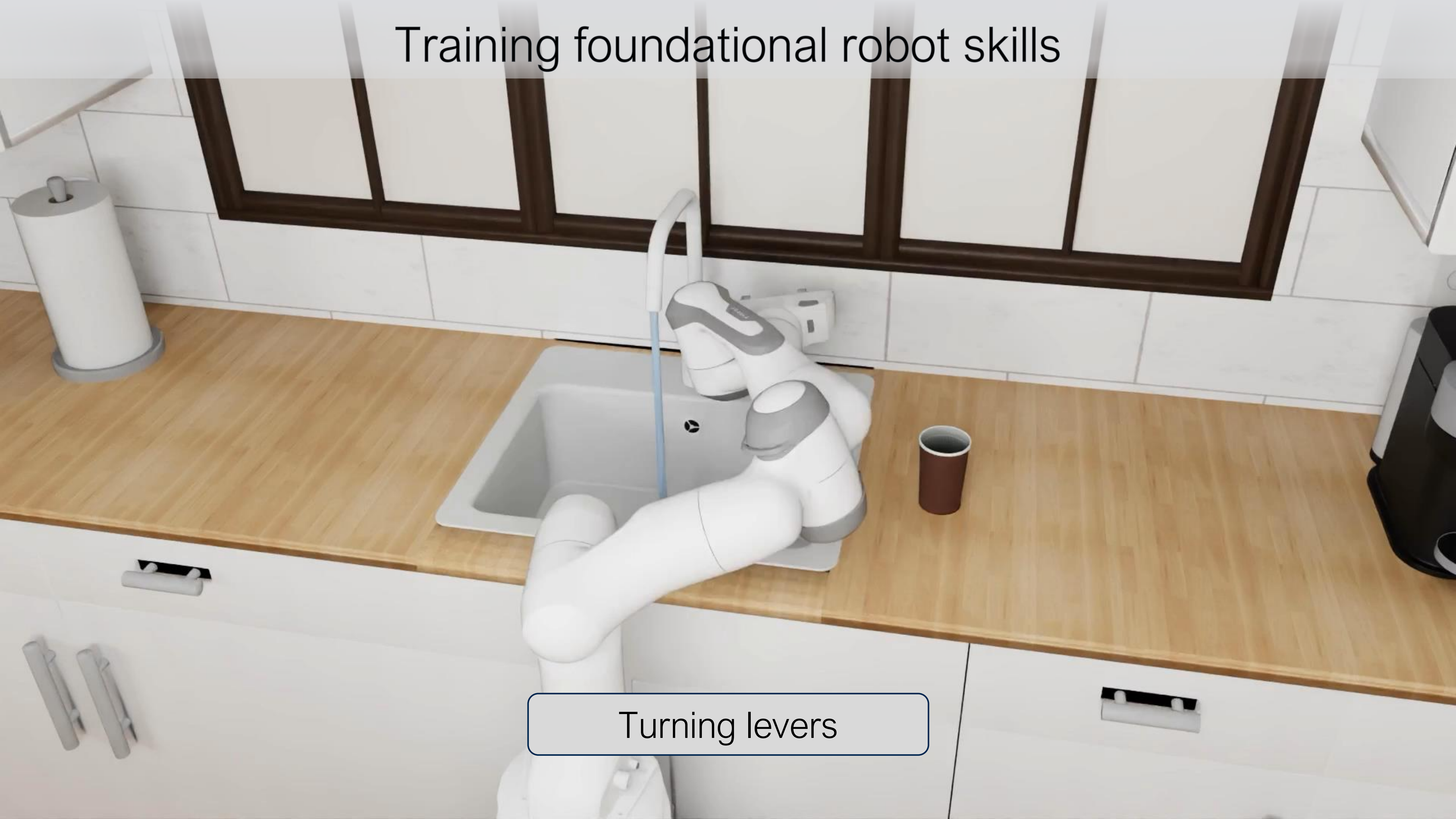
Training foundational robot skills

Pick and place

Training foundational robot skills

Opening and closing doors

Training foundational robot skills
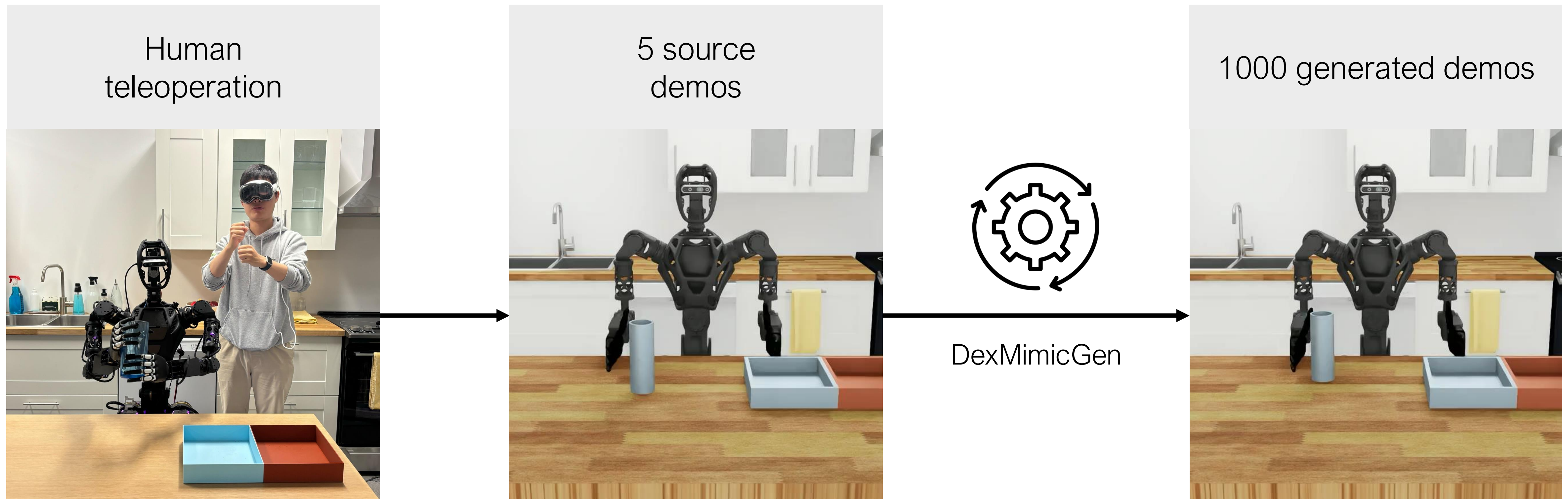
Turning levers

Training foundational robot skills
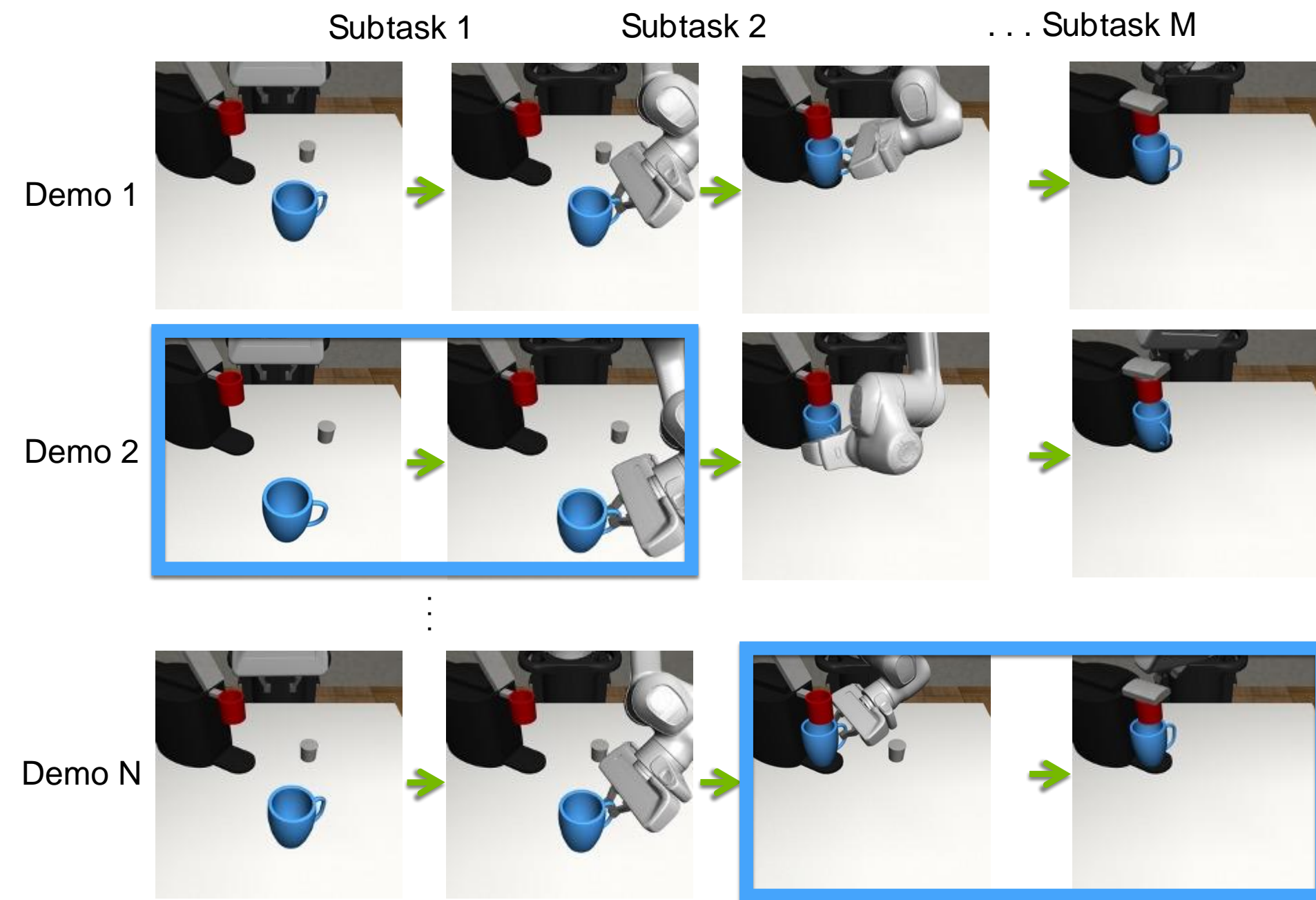
Twisting knobs

Training foundational robot skills

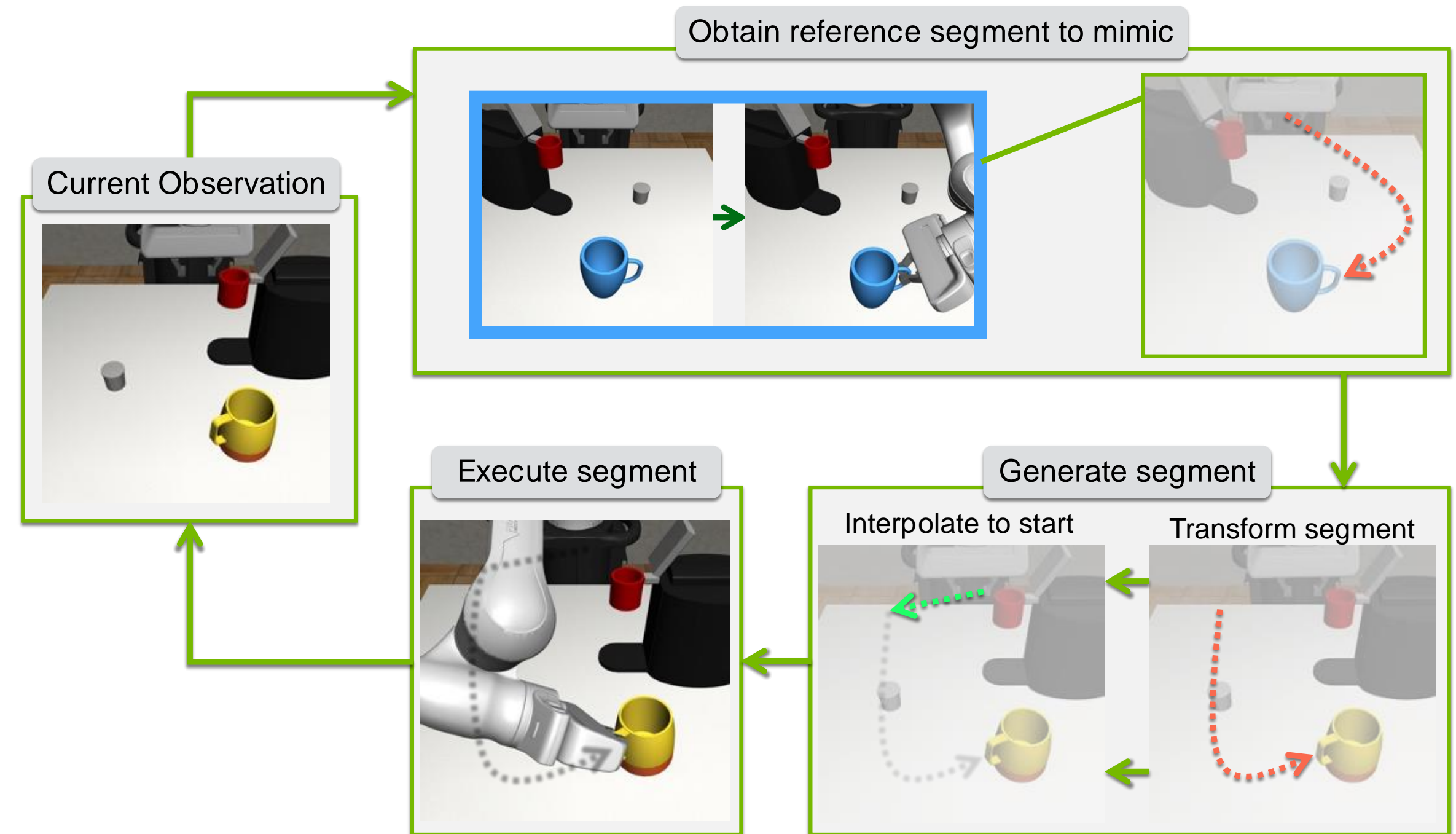Pressing buttons

# DexMimicGen: Automated Data Generation System



Human teleoperation

5 source demos

DexMimicGen

1000 generated demos

"DexMimicGen: Automated Data Generation for Bimanual Dexterous Manipulation via Imitation Learning." Jiang*, Xie*, Lin*, et al. 2024

# DexMimicGen: Automated Data Generation System

## Parse source demonstrations into segments

Subtask 1    Subtask 2    . . . Subtask M

Demo 1

Demo 2

Demo N

Source demos are split into object-centric pieces

## Pipeline for generating new trajectories

Obtain reference segment to mimic

Current Observation

Execute segment    Generate segment

Interpolate to start    Transform segment

Source demo pieces are transformed and replayed in the new scene one by one

"MimicGen: A Data Generation System for Scalable Robot Learning using Human Demonstrations." Mandlekar et al. CoRL 2023

# MimicGen: Data Generation Example



Source dataset trajectory

Generated trajectory

Interpolation segment

# MimicGen: Data Generation Example



Execute transformed segment

Source dataset trajectory

Generated trajectory

# MimicGen: Data Generation Example



Mug grasp is consistent!

Source dataset trajectory

Generated trajectory

# DexMimicGen: Automated Data Generation System
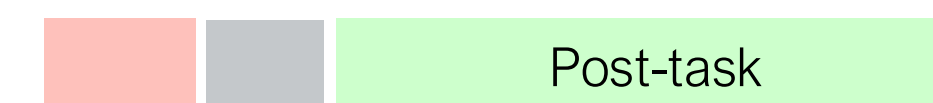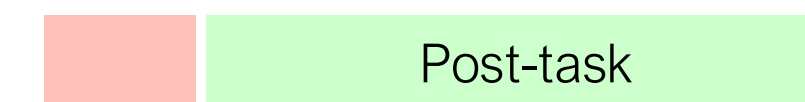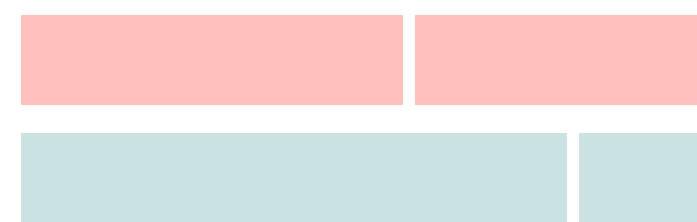


Parallel Subtasks      Coordination Subtasks      Sequential Subtasks

MimicGen

DexMimicGen

Coordination

Coordination

Coordination

Coordination

Post-task

Pre-task

Post-task

Pre-task

Parallel subtask (right)     Parallel subtask (left)     Coordination subtask     Sequential subtask (pre-task)     Sequential subtask (post-task)

# DexMimicGen: Automated Data Generation System

## Source demo segmentation

**Parallel:** pick cube
**Ref object:** blue cube

**Parallel:** place cube
**Ref object:** tray

**Coordination:** lift tray
**Ref object:** tray

Right Arm

Full Demo

Left Arm

**Parallel:** pick cube
**Ref object:** green cube

**Parallel:** place cube
**Ref object:** tray

**Coordination:** lift tray
**Ref object:** tray

## New trajectory generation and execution

Reference Subtask

Reference Trajectory

Current Observation

Object-Centric Trajectory Transformation

Executed Trajectory

Generated Trajectory

DexMimicGen generates data for a large range of tasks.

Contact-rich tasks

DexMimicGen generates data for a large range of tasks.

Long-horizon tasks

DexMimicGen can be used to train real-world visuomotor policy.

Human teleoperation

Real source demo

Real2Sim

DexMG

Generated demos

Sim source demo

Transfer real demo to sim using digital twin to ensure the sim demos are valid in real

DexMimicGen can be used to train real-world visuomotor policy.

Generated demo (sim)

Sim2Real

$\pi_\theta$

Generated demo (real)

Real-world visuomotor policy

Transfer only **successful** generated demos from sim to real to train a visuomotor policy

DexMimicGen can be used to train real-world visuomotor policy.

Real-world visuomotor policy rollouts (10X)

# DexMimicGen: Automated Data Generation System

## Multi-task imitation learning evaluation with RoboCasa simulation tasks



65% relative gain over human data

Legend: Human-50, Generated-100, Generated-300, Generated-3000

"RoboCasa: Large-Scale Simulation of Everyday Tasks for Generalist Robots." Nasiriany et al. RSS 2024

# DexMimicGen: Automated Data Generation System



Training on 50 real-robot demonstrations: 13.6%

Co-training with real (50) + sim (45k) datasets: 24.4%

"RoboCasa: Large-Scale Simulation of Everyday Tasks for Generalist Robots." Nasiriany et al. RSS 2024

# Recipe for Building Robotic Foundation Models



**Scalable Algorithms**

Powerful robot learning models that scale with data and compute

**Algorithms**

**Robotic Foundation Models**

**Data**

**Hardware**

**Data Engine**

New mechanisms to produce massive training data

**Human-like Embodiment**

Humanoid robot platform for broad applications

# Three-Phase Training for Robotic Foundation Models



| | | | |
|---|---|---|---|
| Dataset size | x00 billions to 1.x trillion tokens | ~xk to x0k (prompt, response) | ~x0k prompt |
| Example of models | GPT-3, LLaMA, Falcon, BLOOM | Dolly-v2, Falcon-instruct | Claude, GPT-4, ChatGPT |
| | (a)Pre-training | (b) Instruction fine-tuning | (c) Reinforcement learning from human feedback |

[Source: RBC Borealis]

Training process of **LLMs** (ChatGPT, Claude, etc.)

Real-Robot Data

Synthetic Data

Web Data

**(a) pre-training** on data pyramid

**(b) fine-tuning** on domain-specific data

Deploy Model

Improve Model

**(c) alignment** during deployment

Training process of **Robotic Foundation Models**

# The Robot Learning Data Flywheel

# The Robot Learning Data Flywheel

# The Robot Learning Data Flywheel

# The Robot Learning Data Flywheel

Research Principle #3:

**Data Flywheel through Trustworthy and Safe Deployment**

# Robot Learning on the Job: Building the Data Flywheel

## The Sirius Framework for Human-Robot Teaming

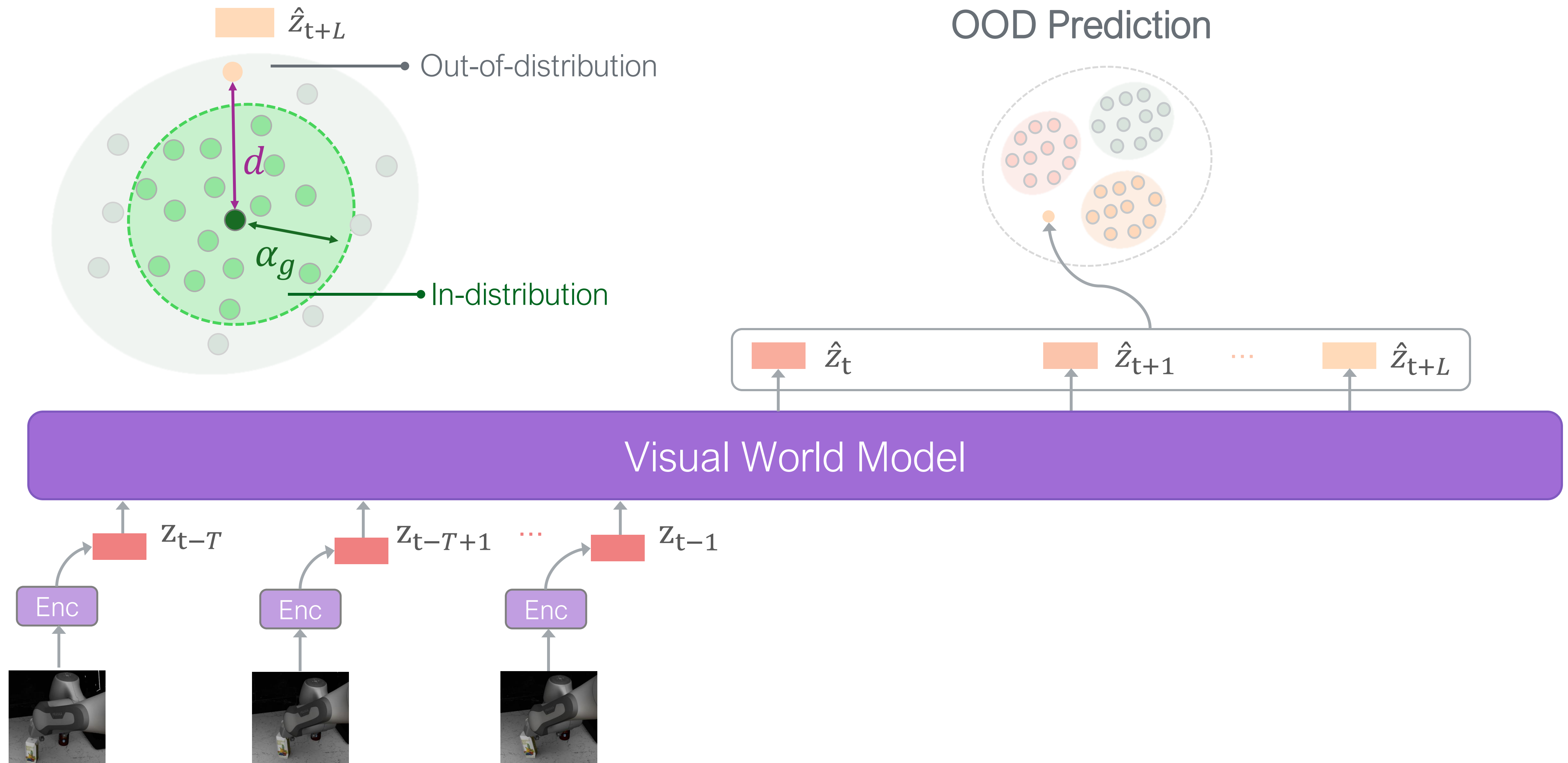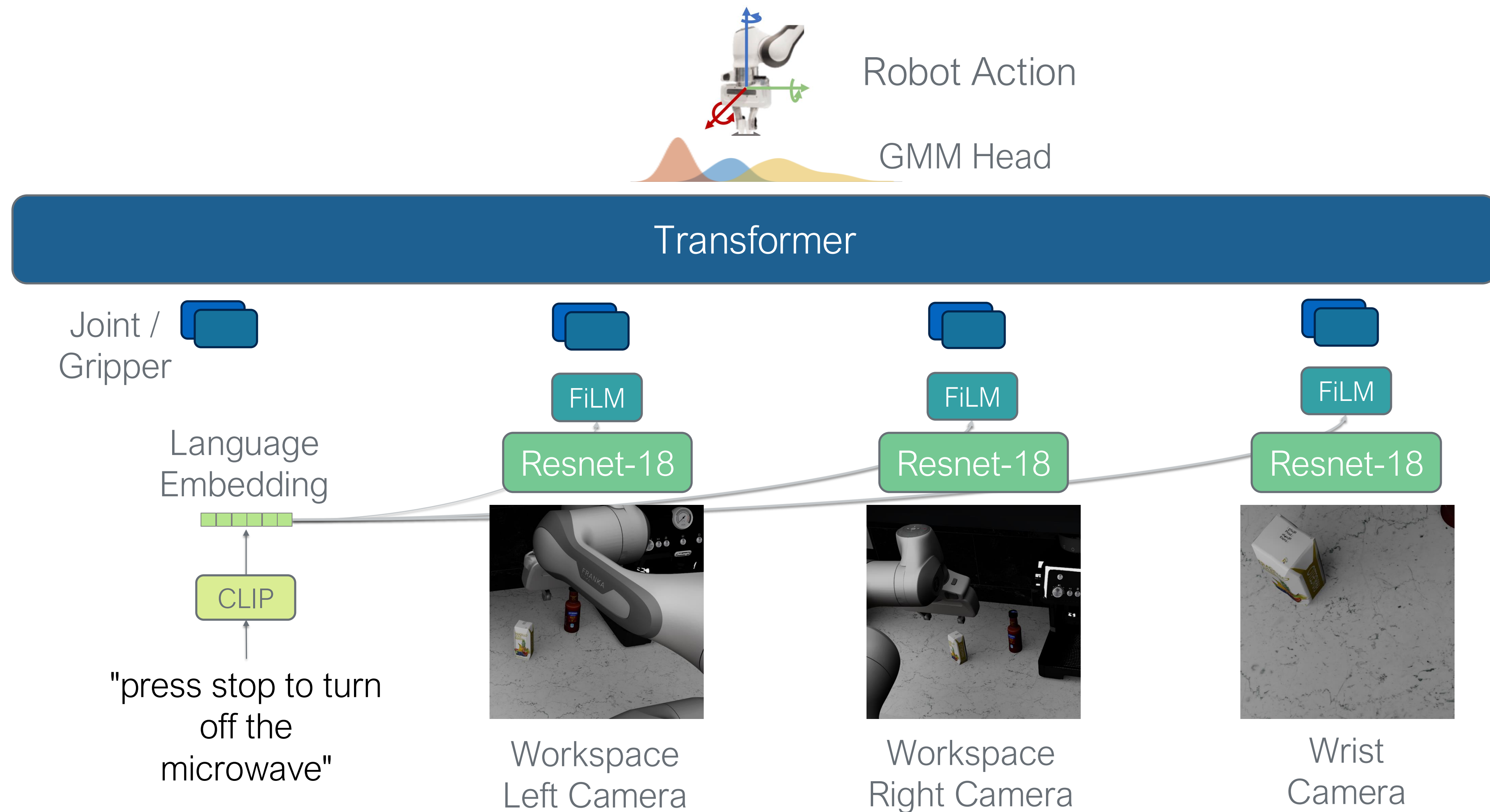"Robot Learning on the Job: Human-in-the-Loop Autonomy and Learning During Deployment." Liu et al. RSS 2023

# Robot Learning on the Job: Building the Data Flywheel

Robot
Deployment



**Model Update**

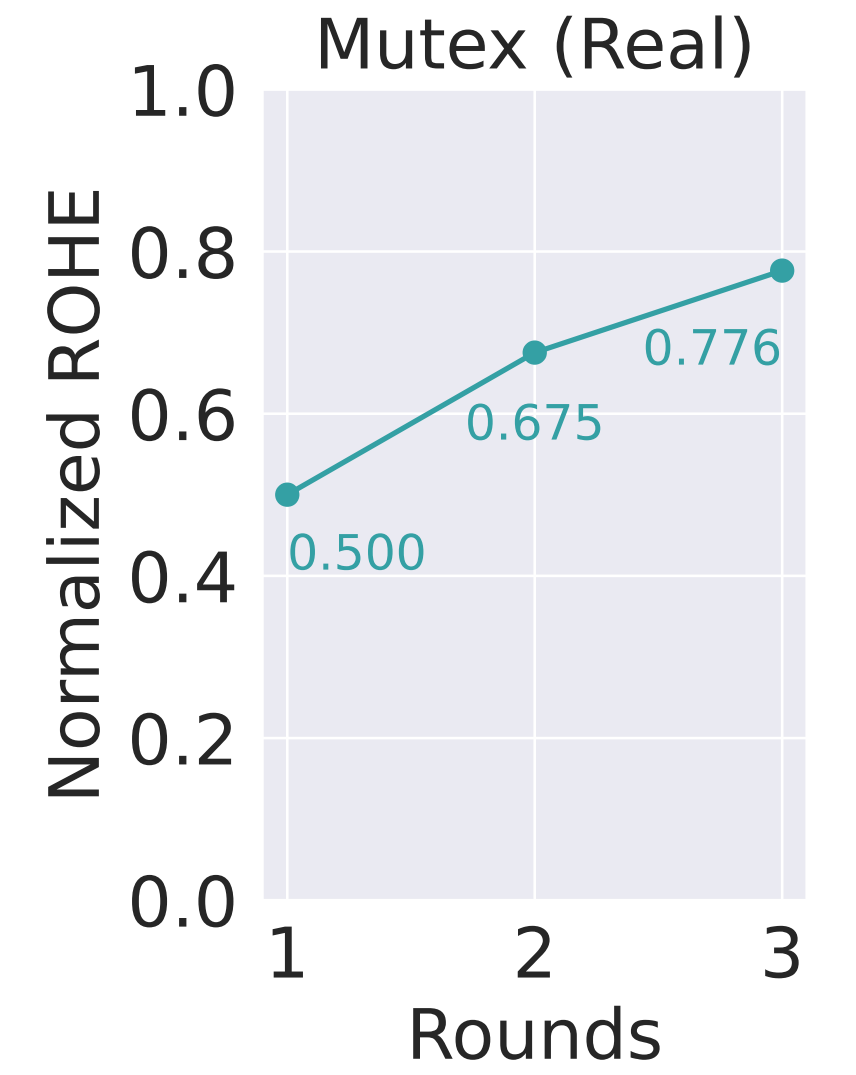"Robot Learning on the Job: Human-in-the-Loop Autonomy and Learning During Deployment." Liu et al. RSS 2023

# Robot Learning on the Job: Building the Data Flywheel



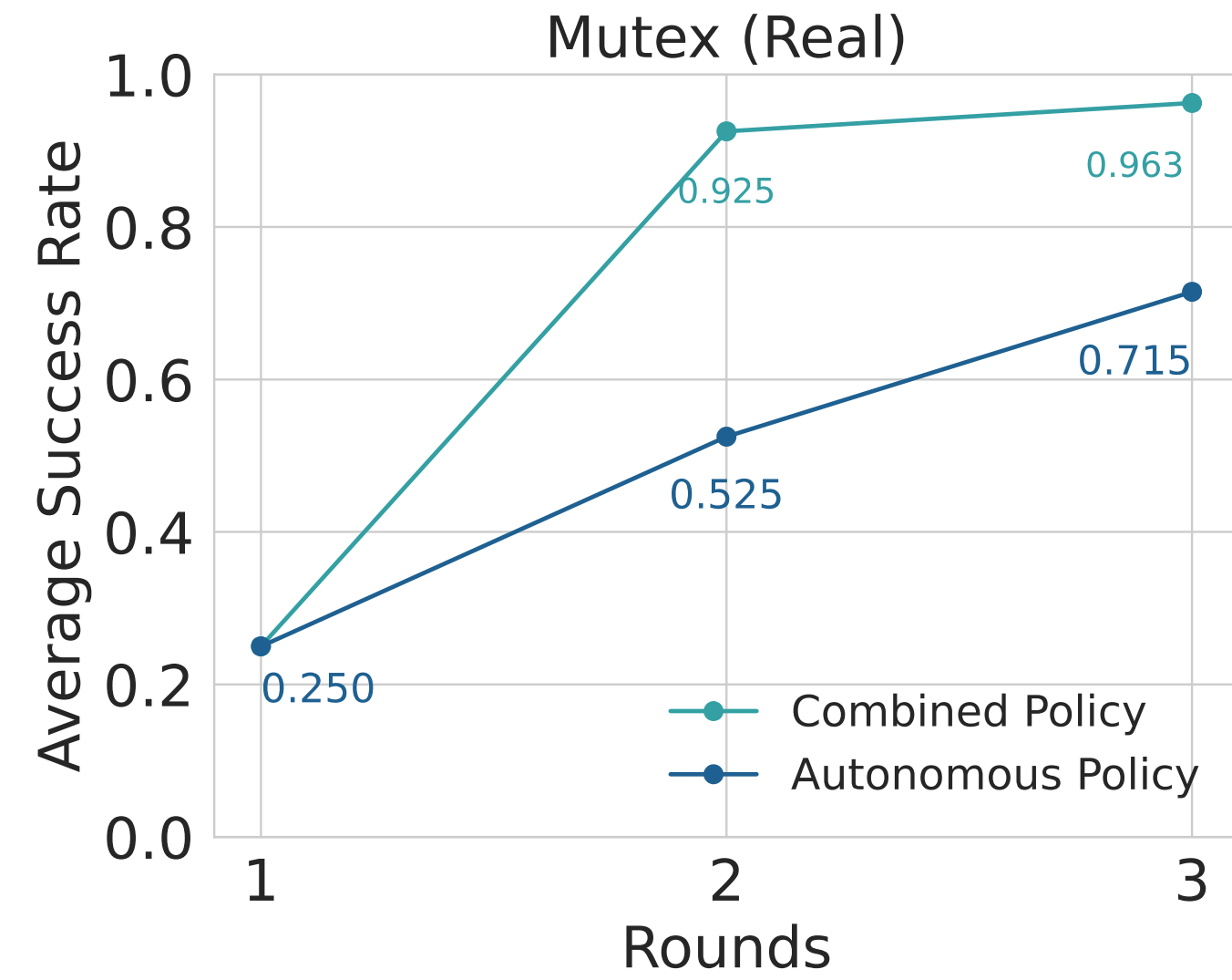"Robot Learning on the Job: Human-in-the-Loop Autonomy and Learning During Deployment." Liu et al. RSS 2023

2x

Robot takes
control again

Human provides
intervention

"Model-Based Runtime Monitoring with Interactive Imitation Learning." Liu et al. ICRA 2024

# Robot Learning on the Job: Building the Data Flywheel

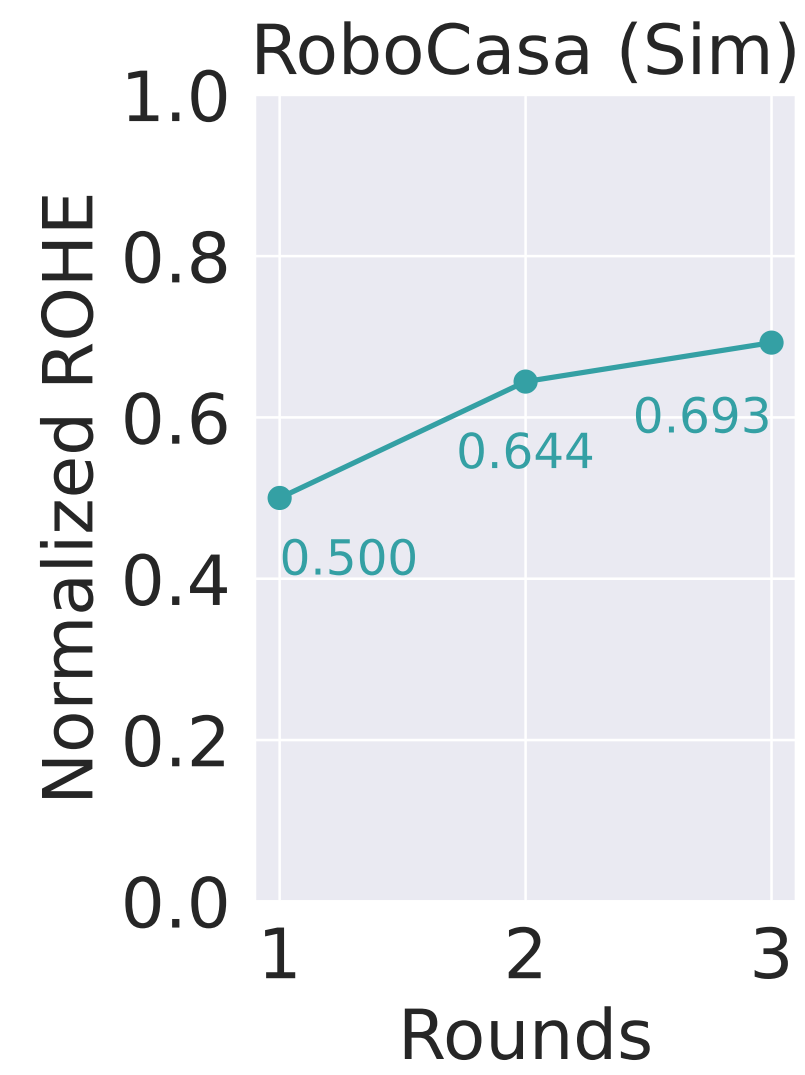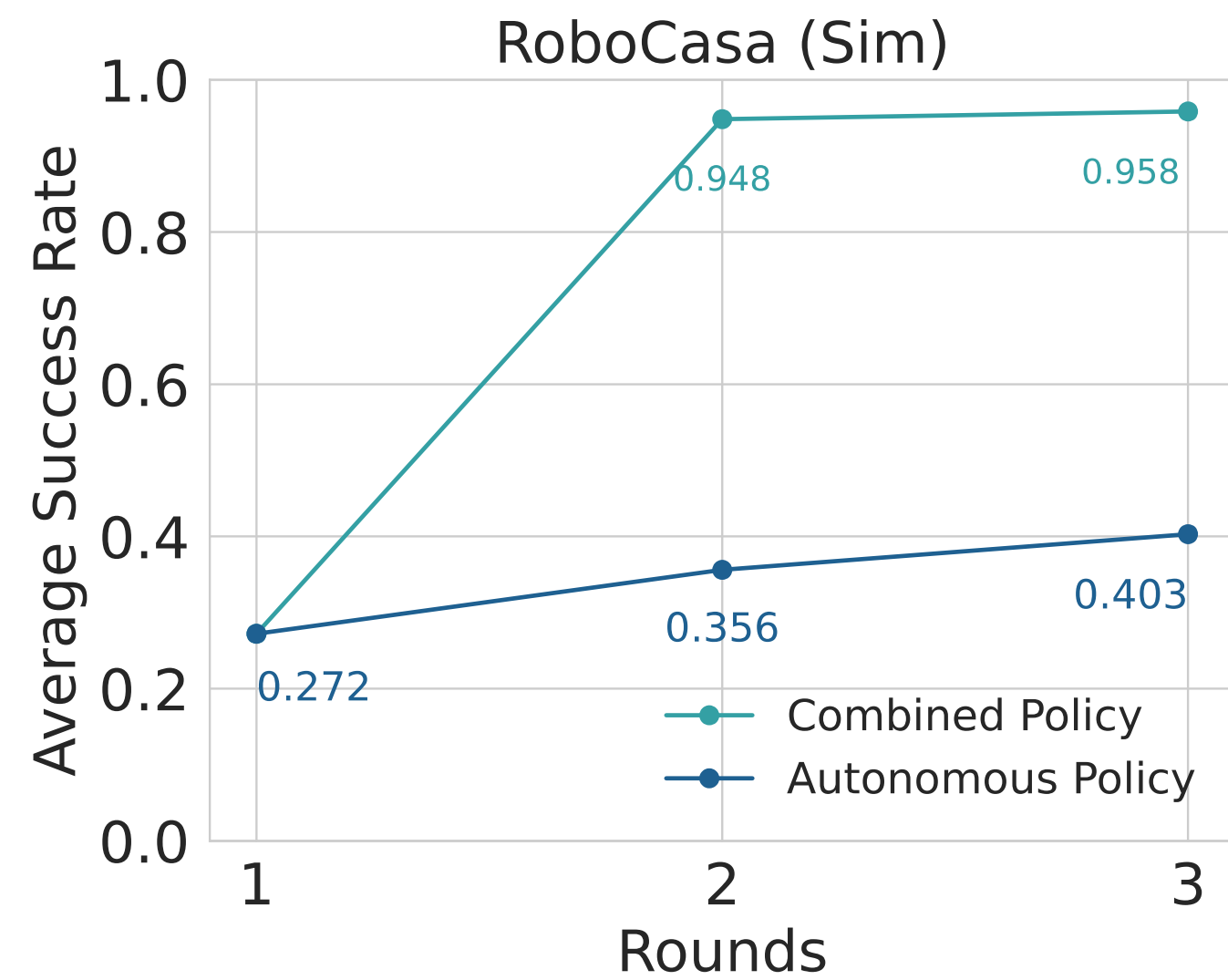### Round 1 Deployment



† Green masks indicate human intervention.

### Intervention Distribution



### Round 3 Deployment



### Intervention Distribution

# Robot Learning on the Job: Building the Data Flywheel

**Human**

**Robot Fleet**

Runtime Monitoring

Robot Fleet Deployment

**Anomaly Predictor**

**Visual World Model**

**Memory Storage**

**Next Policy**

**Model Update**

"Multi-Task Interactive Robot Fleet Learning with Visual World Models." Liu et al. CoRL 2024

# Robot Learning on the Job: Building the Data Flywheel
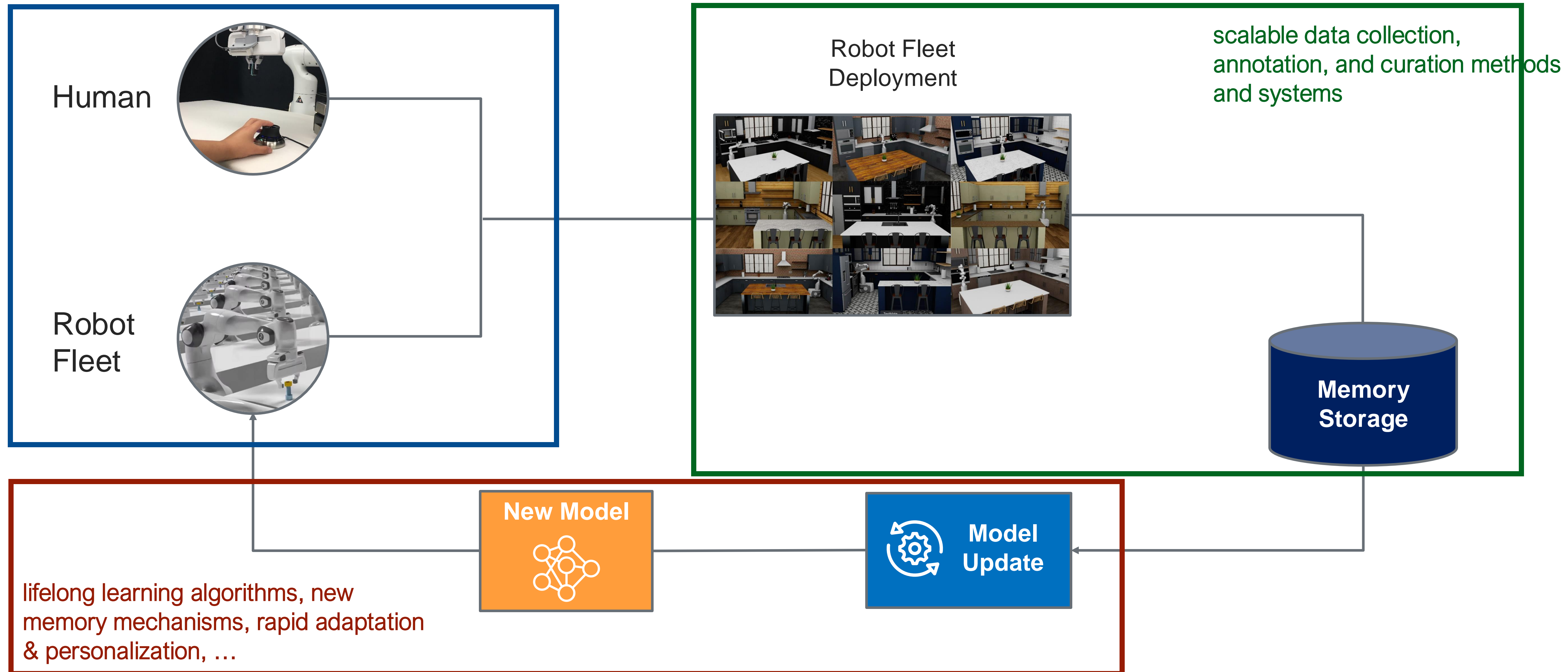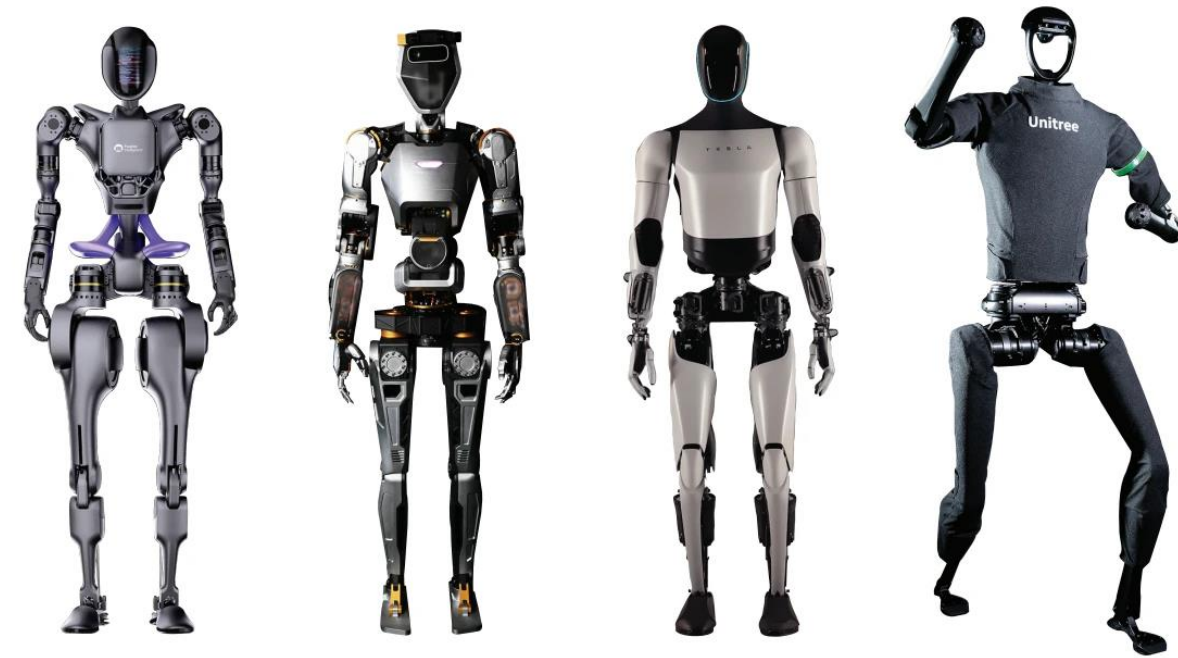
# Robot Learning on the Job: Building the Data Flywheel



"Multi-Task Interactive Robot Fleet Learning with Visual World Models." Liu et al. CoRL 2024

# Robot Learning on the Job: Building the Data Flywheel



OOD Prediction

Out-of-distribution

In-distribution

$\hat{z}_{t+L}$

$d$

$\alpha_g$

$\hat{z}_t$  $\hat{z}_{t+1}$  $\cdots$  $\hat{z}_{t+L}$

Visual World Model

$z_{t-T}$  $z_{t-T+1}$  $\cdots$  $z_{t-1}$

Enc  Enc  Enc

"Multi-Task Interactive Robot Fleet Learning with Visual World Models." Liu et al. CoRL 2024

# Robot Learning on the Job: Building the Data Flywheel



"Multi-Task Interactive Robot Fleet Learning with Visual World Models." Liu et al. CoRL 2024

# Robot Learning on the Job: Building the Data Flywheel



Human efforts reduce over time as policy performance continually improves.

"Multi-Task Interactive Robot Fleet Learning with Visual World Models." Liu et al. CoRL 2024

# Robot Learning on the Job: Building the Data Flywheel

# Robot Learning on the Job: Building the Data Flywheel

# Turn the Data Flywheel, Flip Data Pyramid Upside Down



Real-World Data

Synthetic Data

Web Data

Real-World Data
(through widespread deployments)

Synthetic Data
(turbocharged by generative AI)

Web Data
(growing but dwarfed by the other two)

The Present

The Future

# Talk Summary



Research Principle #1: **First Generalist, then Better Specialist**

[OKAMI, CoRL 2024; NVIDIA Project GR00T]

Real-World Data

Synthetic Data

Web Data

Research Principle #2: **Learning Across the Data Pyramid**

[MimicGen, CoRL 2023; RoboCasa, RSS 2024; BUMBLE, arXiv 2024; DexMimicGen, arXiv 2024]

Increased Deployments

More Training Data

More Capable Robots

Better Learning

Research Principle #3: **Data Flywheel through Trustworthy and Safe Deployment**

[Sirius, RSS 2023; Sirius-RM, ICRA 2024; Sirius-Fleet, CoRL 2024]

Papers can be found at https://yukezhu.me/

# Acknowledgement