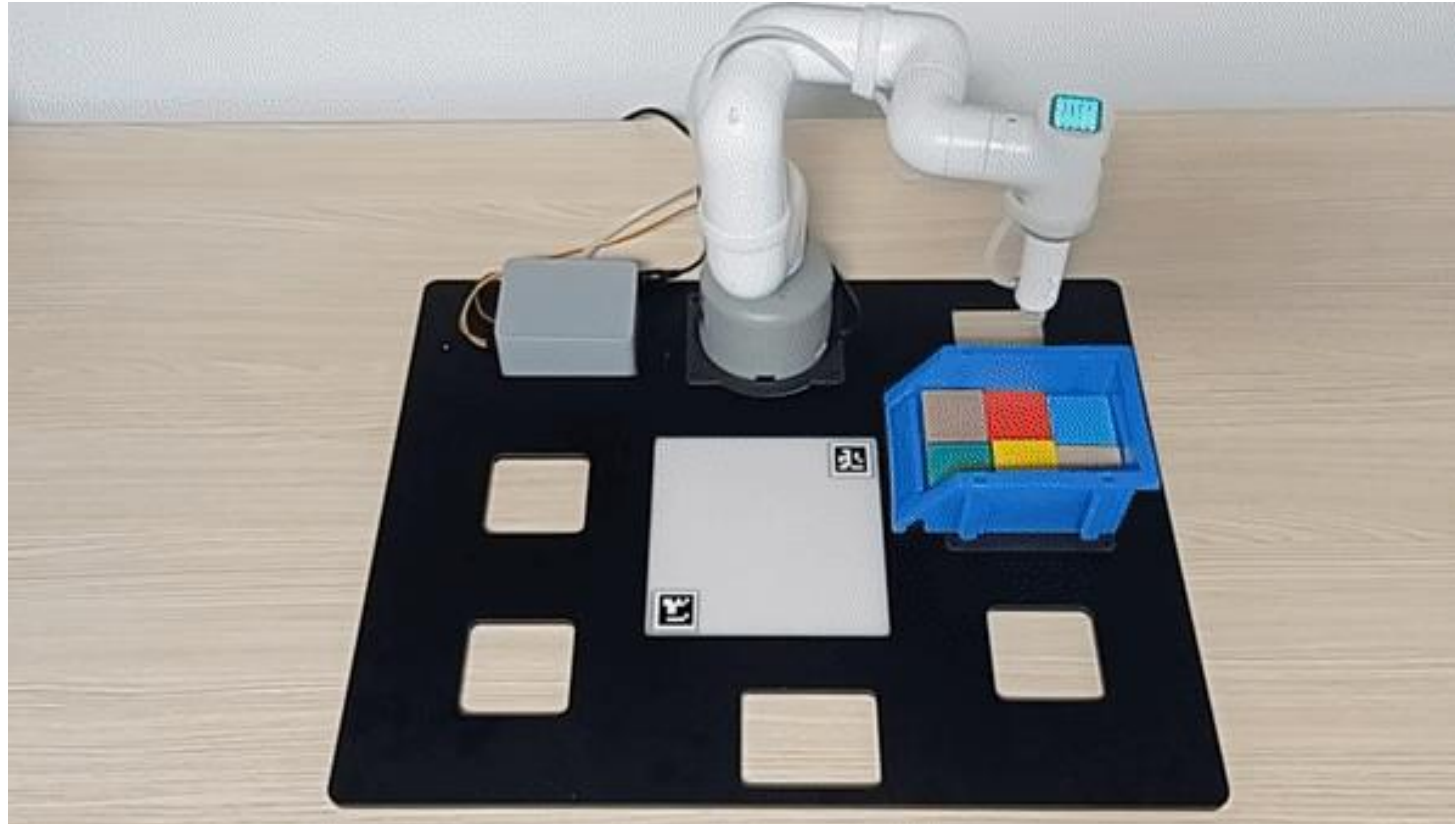


Foundational Models for Robotics: Removing the Engineer from the Loop



Rogério Bonatti
April 27, 2023

Speaker **Bio**: Rogerio Bonatti



Senior Researcher

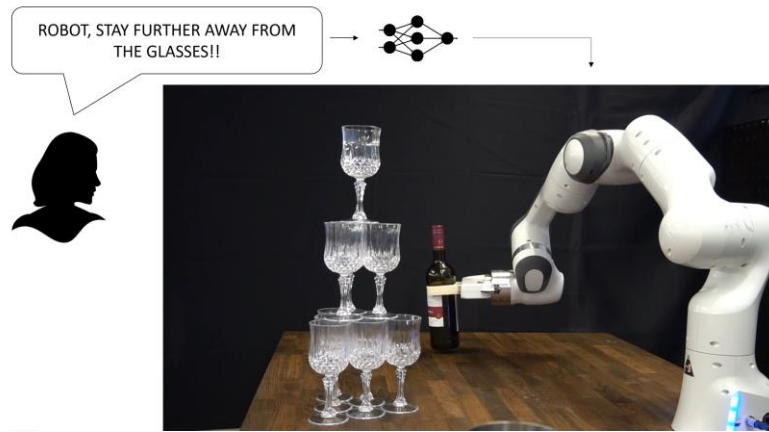
PhD in Robotics,
School of Computer Science

Research Internships,
FAIR and MSR

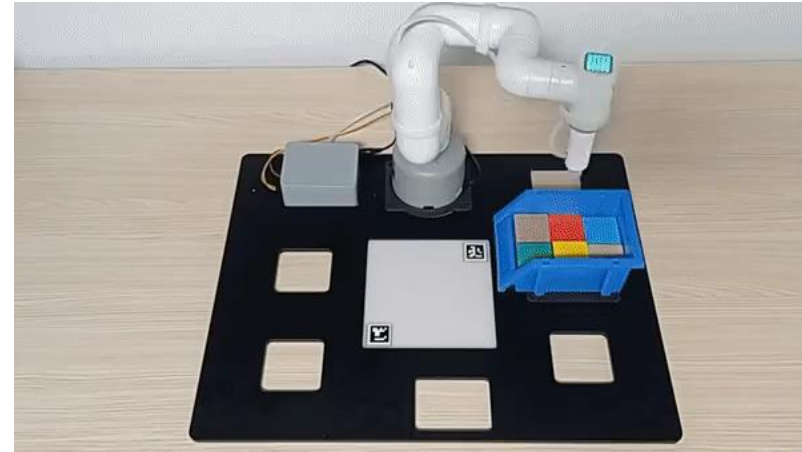


Research projects:

Multi-modal robot control (vision + language):



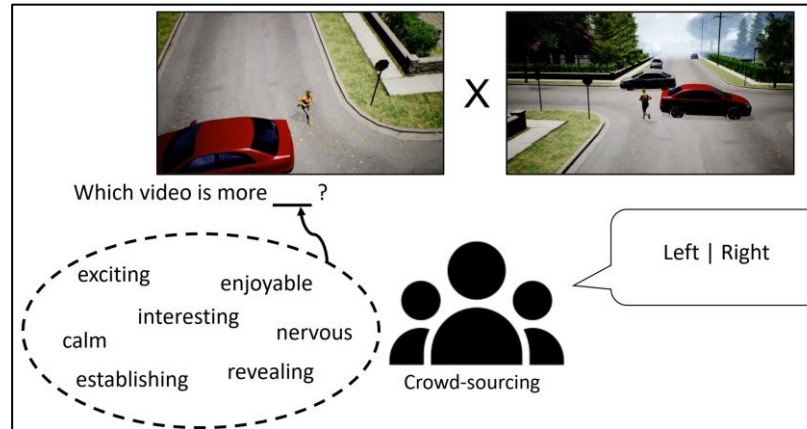
Generative models for decision-making:



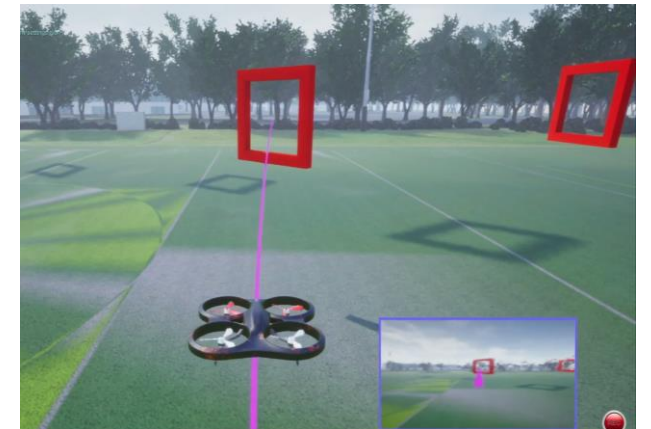
Autonomous drone cinematography:



Learning emotional camera control:

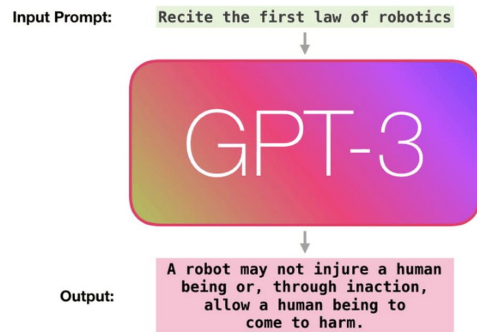


Visuo-motor representation learning:



Large models are revolutionizing generative AI

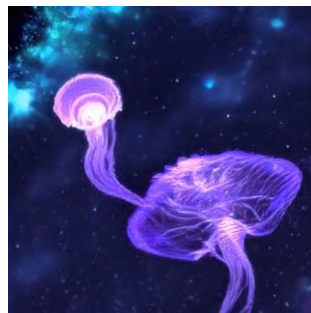
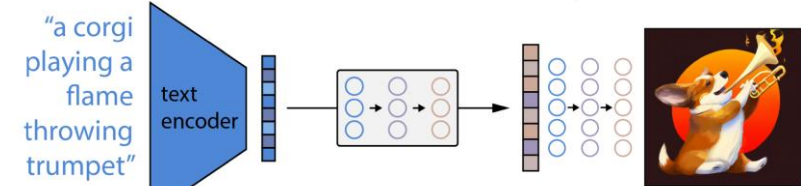
Text generation
[GPT-X, LLaMA]



Conversation
[ChatGPT]



Image generation
[DALL-E, Stable Diffusion]



Text-to-video
[NVIDIA Picasso]



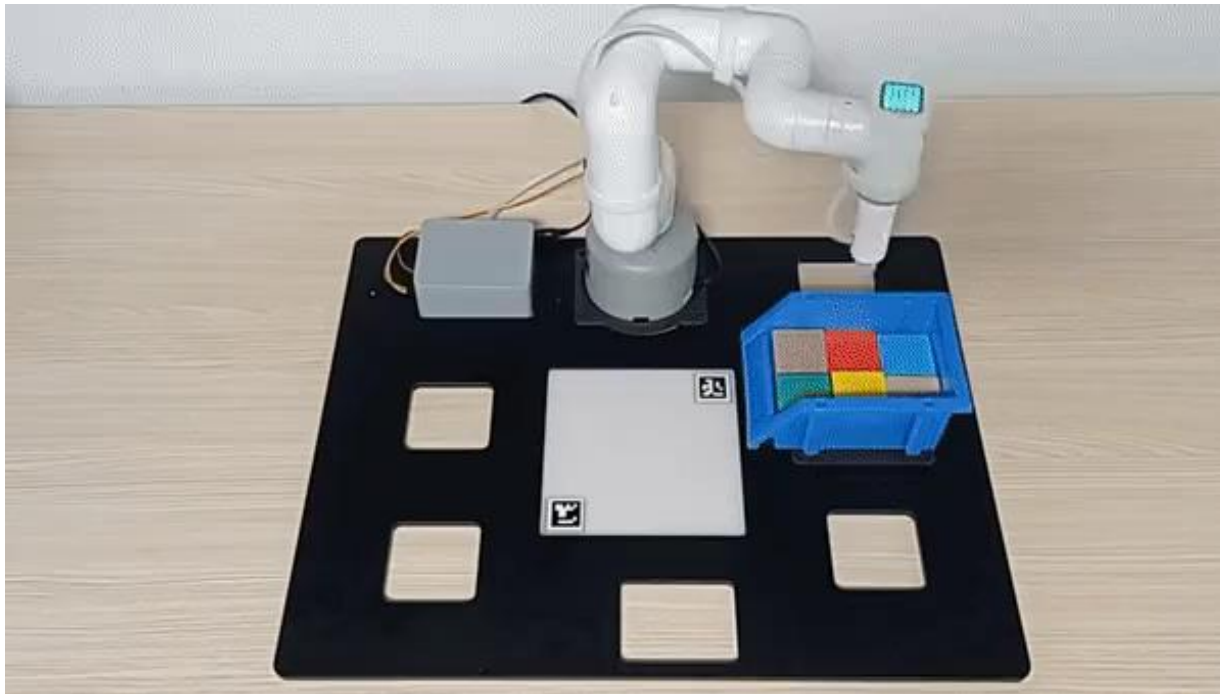
Text-to-3D
[Point-e, NVIDIA Picasso]



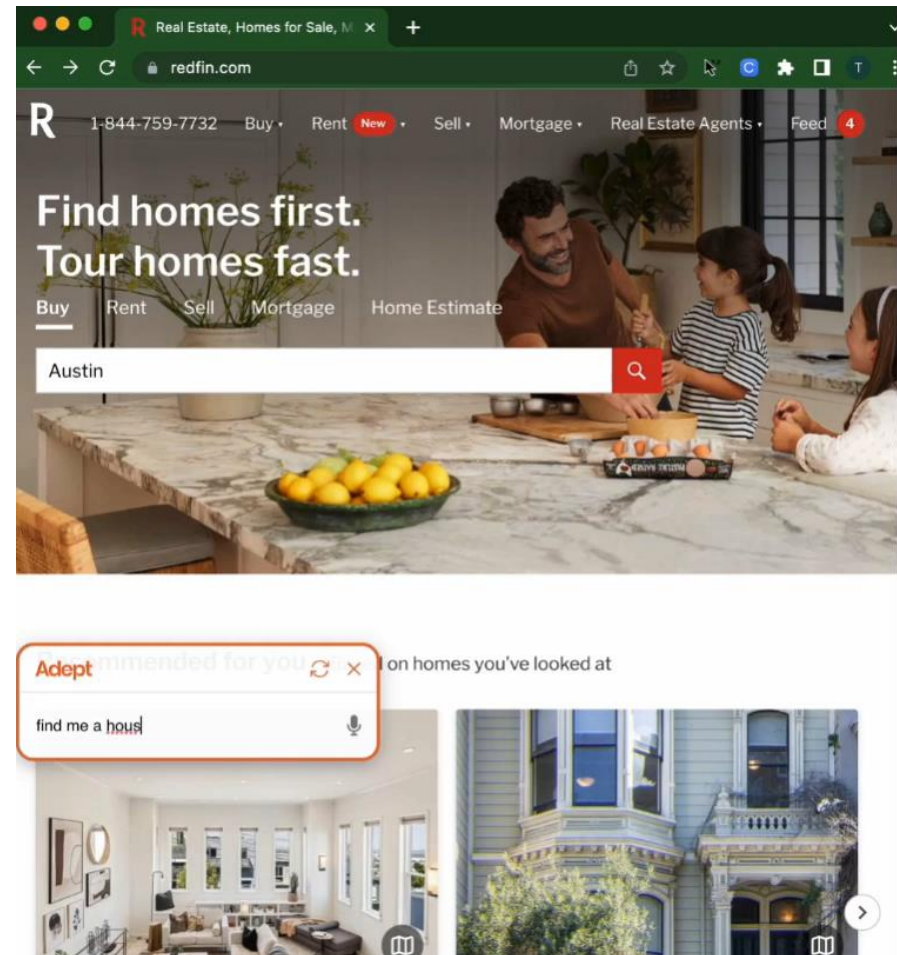
Image editing
[InstructPix2Pix, Prompt-to-Prompt]

But generative AI is just starting to revolutionize **decision-making**:

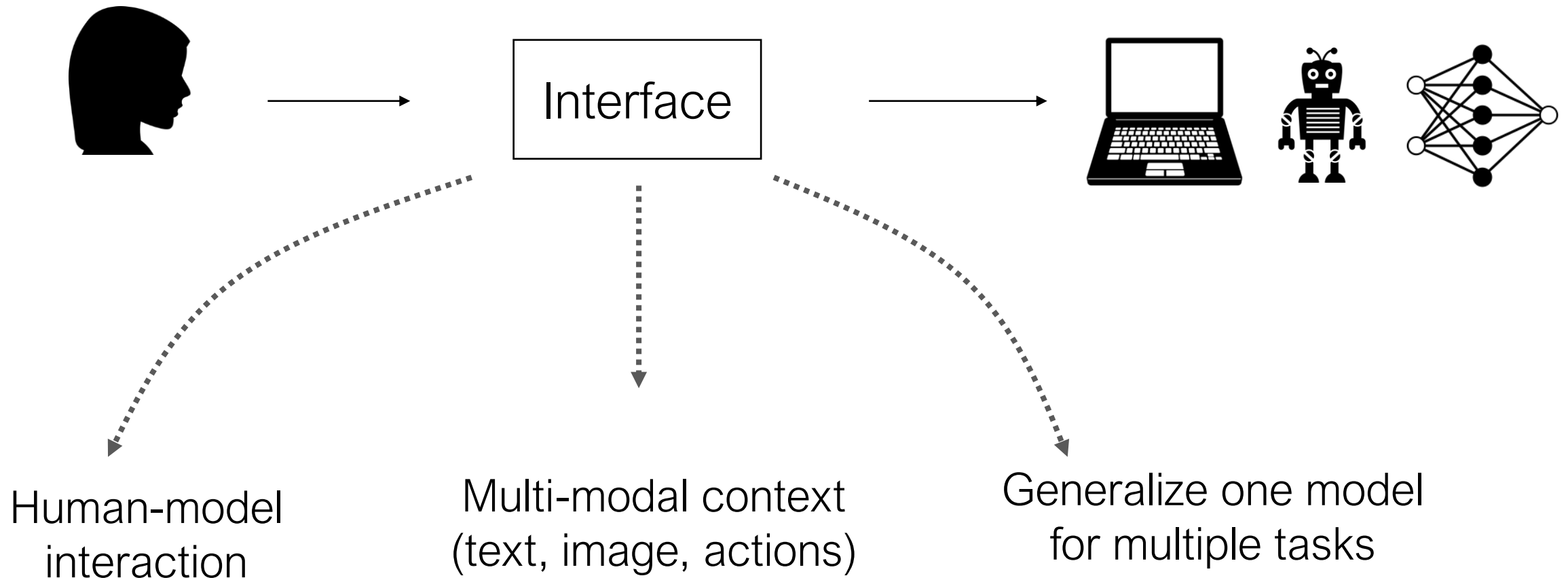
Human-robot interfaces:
[ChatGPT for Robotics]



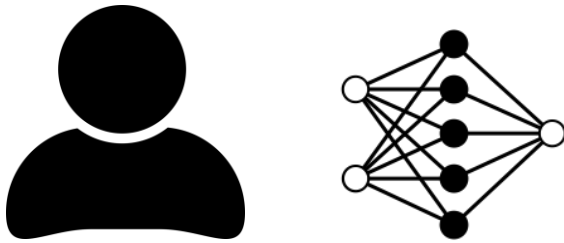
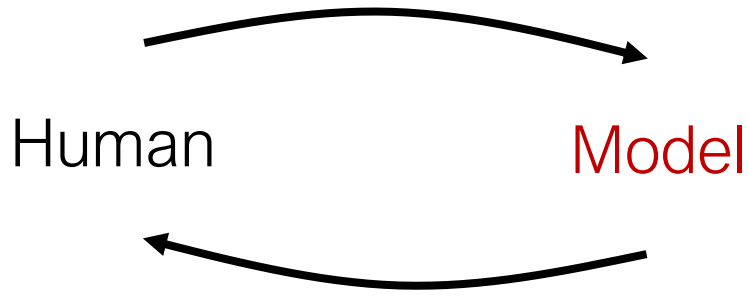
Human-computer interfaces
[video from Adept AI]



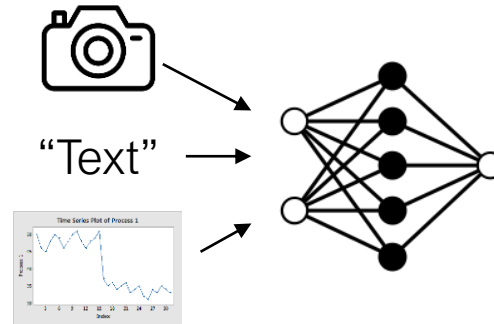
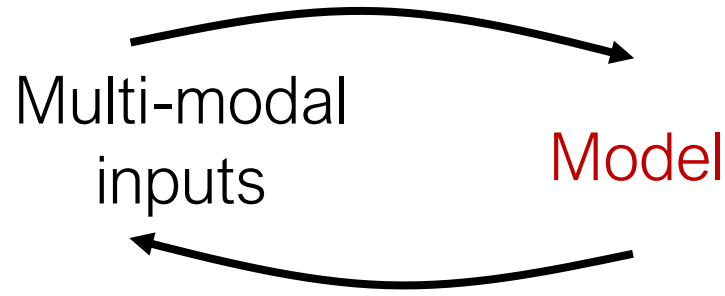
Goal: use large models to **empower any user** with the capabilities of generative AI



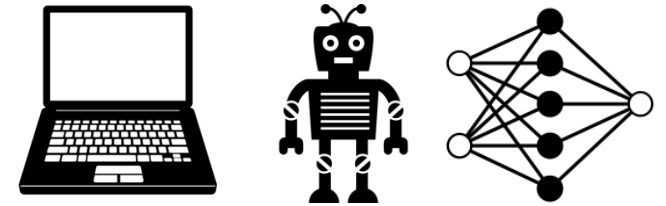
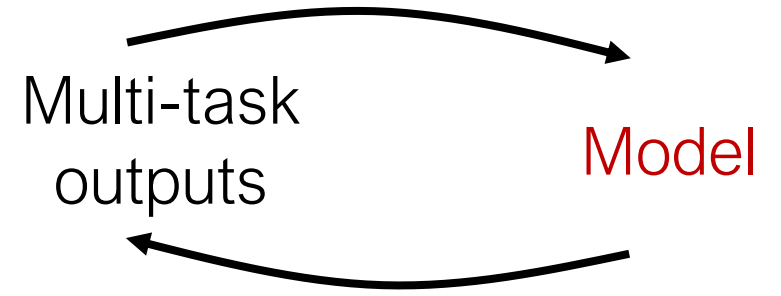
Goal: use large models to **empower any user** with the capabilities of generative AI



Close the loop between the human user and the model



Fusing multi-modal information in I/O

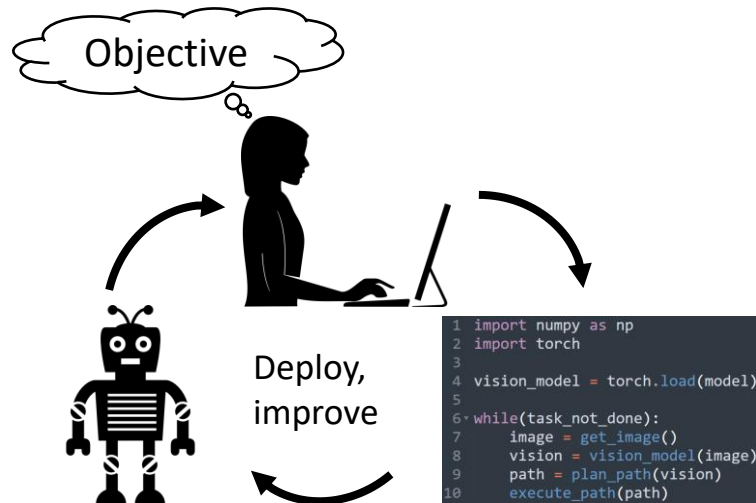


Generalizing a single model to multiple tasks

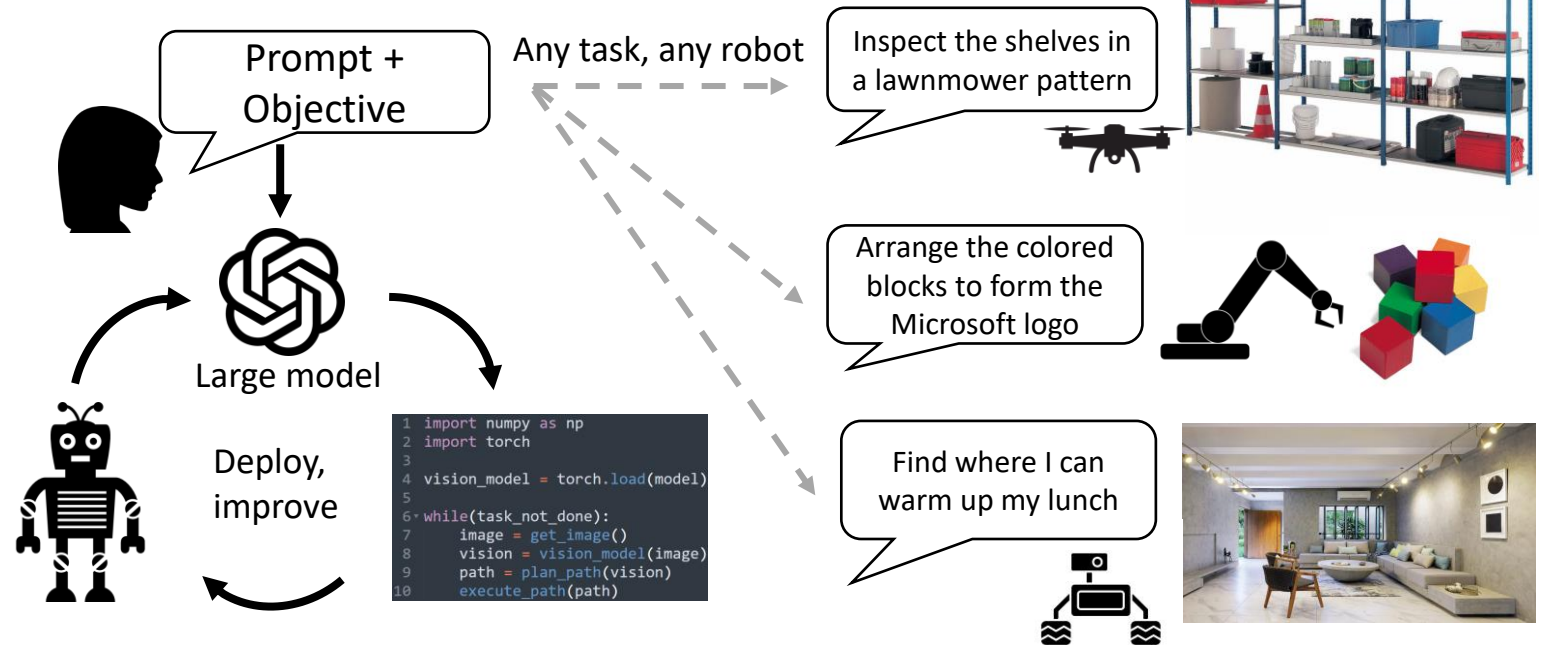
What is the future of human-computer **interfaces**?

Remove the **engineer** from the loop. Bring in the **user**.

Interfaces today: engineer *in the loop*



Goal with ChatGPT: user *on the loop*



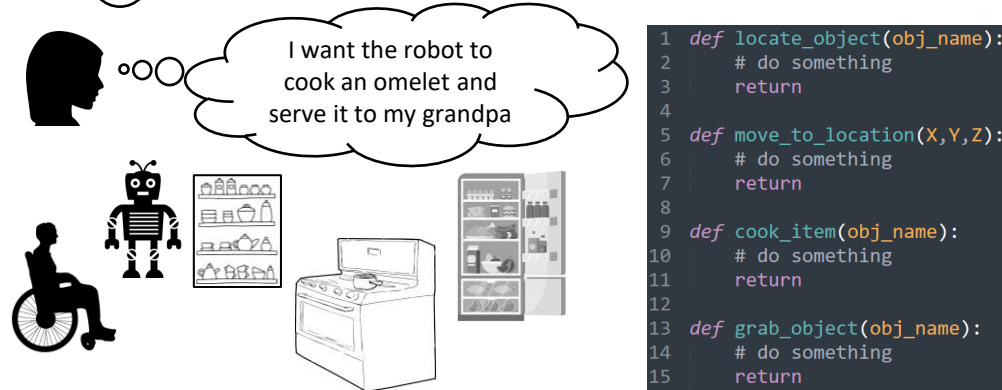
ChatGPT for Robotics: Deployment example

Autonomous flight experiments



Design principles for effective ChatGPT usage:

① Define a task-relevant robot API library*



I want the robot to cook an omelet and serve it to my grandpa

```
1 def locate_object(obj_name):
2     # do something
3     return
4
5 def move_to_location(X,Y,Z):
6     # do something
7     return
8
9 def cook_item(obj_name):
10    # do something
11    return
12
13 def grab_object(obj_name):
14    # do something
15    return
```

*APIs should be easily implementable on the robot and have descriptive text names for the LLM. They can be chained together to form more complex functions.

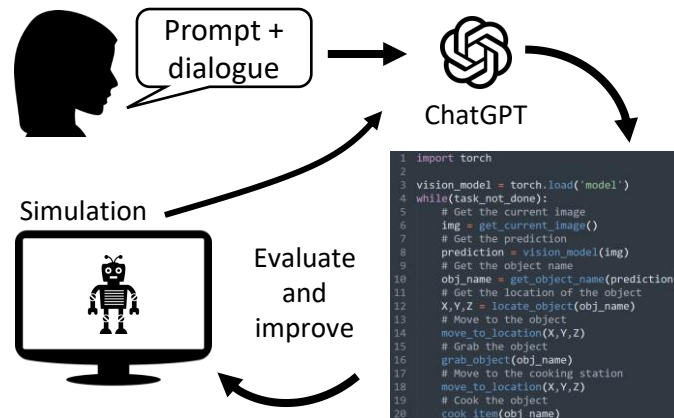
② Build prompt following engineering principles

Consider you are a home assistant robot. Your goal is to prepare an omelet for an elderly person. You are equipped with functions:

locate_object(obj_name): returns a X,Y,Z tuple representing the location of the desired object defined by string "obj_name";
move_to_location(X,Y,Z): moves the robot's hands to a specific X,Y,Z location in space. Returns nothing;
cook_item(obj_name): cooks a particular item defined by "obj_name". Returns nothing;
grab_object(obj_name): picks a particular object defined by "obj_name". Returns nothing;

Output python code with the sequence of steps that achieves your objective.

③ User on the loop: iterate on solution quality and safety

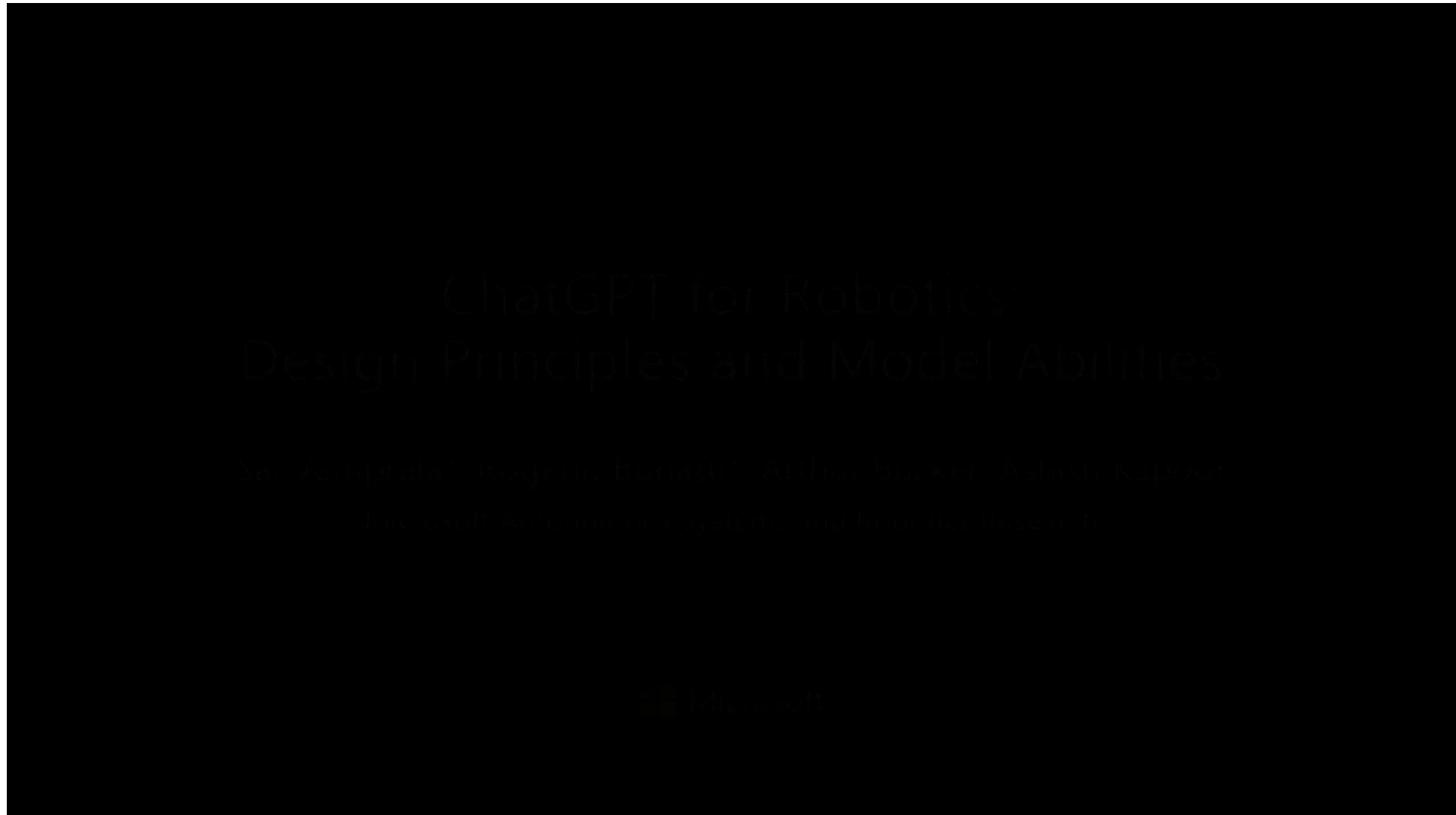


④ Execute!




Deployment with the **user in the loop**




When a **conversation** is needed for complex tasks






Multi-modal emerging properties


Hallucination over colors, shapes and geometry


SA Hi, what does the Microsoft logo look like? 

 The Microsoft logo is a colored rectangular with four squares arranged in a grid pattern. The colors used in the logo are blue, green, orange, and yellow. The logo has gone through several iterations over the years, but this is the most recent and recognizable design.  

SA Can you generate an SVG with those colors?

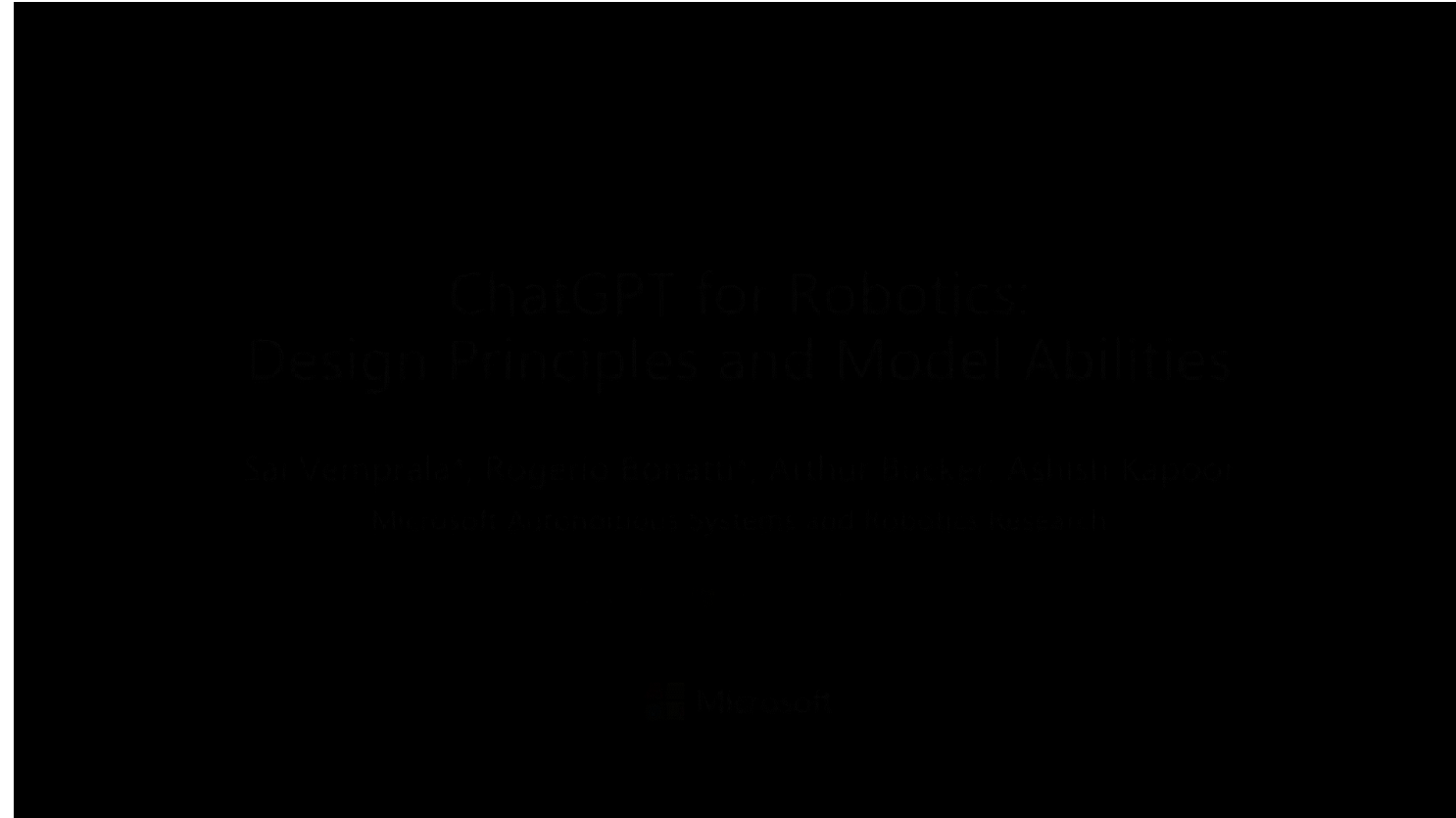
 Yes, I can generate an SVG file with the Microsoft logo colors. Here is an example:  

```
svg  Copy code
<svg width="100" height="100" viewBox="0 0 100 100" xmlns="http://www.w3.org/2000/svg">
  <rect x="0" y="0" width="50" height="50" fill="#F25022"/>
  <rect x="50" y="0" width="50" height="50" fill="#7FBA00"/>
  <rect x="0" y="50" width="50" height="50" fill="#00A4EF"/>
  <rect x="50" y="50" width="50" height="50" fill="#FFB900"/>
</svg>
```



Rendered SVG

This SVG file has a width and height of 100 pixels and uses the Microsoft logo colors: #F25022 for the red square, #7FBA00 for the green square, #00A4EF for the blue square, and #FFB900 for the yellow square.



Community support: prompt library and simulator

PromptCraft-Robotics

The PromptCraft-Robotics repository serves as a community for people to test and share interesting prompting examples for large language models (LLMs) within the robotics domain. We also provide a sample [robotics simulator](#) (built on Microsoft AirSim) with ChatGPT integration for users to get started.

We currently focus on OpenAI's [ChatGPT](#), but we also welcome examples from other LLMs (for example open-sourced models or others with API access such as [GPT-3](#) and [Codex](#)).

Users can contribute to this repository by submitting interesting prompt examples to the [Discussions](#) section of this repository. A prompt can be submitted within different robotics categories such as [Manipulation](#), [Home Robotics](#), [Physical Reasoning](#), among many others. Once submitted, the prompt will be reviewed by the community (upvote your favorites!) and added to the repository by a team of admins if it is deemed interesting and useful. We encourage users to submit prompts that are interesting, fun, or useful. We also encourage users to submit prompts that are not necessarily "correct" or "optimal" but are interesting nonetheless.

We encourage prompt submissions formatted as markdown, so that they can be easily transferred to the main repository. Please specify which LLM you used, and if possible provide other visuals of the model in action such as videos and pictures.

ChatGPT Prompting Guides & Examples

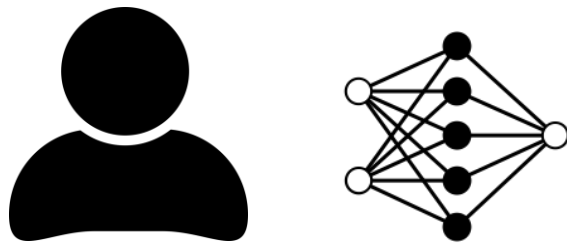
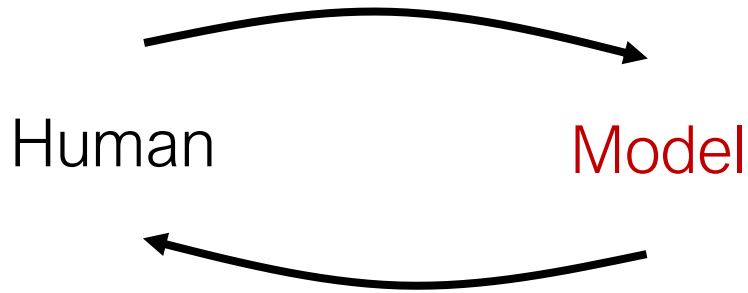
The list below contains links to the different robotics categories and their corresponding prompt examples. We welcome contributions to this repository to add more robotics categories and examples. Please submit prompt examples to the [Discussions](#) page, or submit a pull request with your category and examples.

- Embodied agent
 - [ChatGPT - Habitat, closed loop object navigation 1](#)
 - [ChatGPT - Habitat, closed loop object navigation 2](#)
 - [ChatGPT - AirSim, object navigation using RGBD](#)
- Aerial robotics
 - [ChatGPT - Real robot: Tello deployment | Video Link](#)
 - [ChatGPT - AirSim turbine Inspection | Video Link](#)
 - [ChatGPT - AirSim solar panel Inspection](#)
 - [ChatGPT - AirSim obstacle avoidance | Video Link](#)
- Manipulation
 - [ChatGPT - Real robot: Picking, stacking, and building the MSFT logo | Video Link](#)
 - [ChatGPT - Manipulation tasks](#)
- Spatial-temporal reasoning
 - [ChatGPT - Visual servoing with basketball](#)

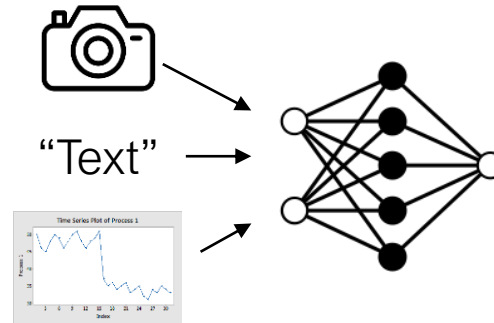
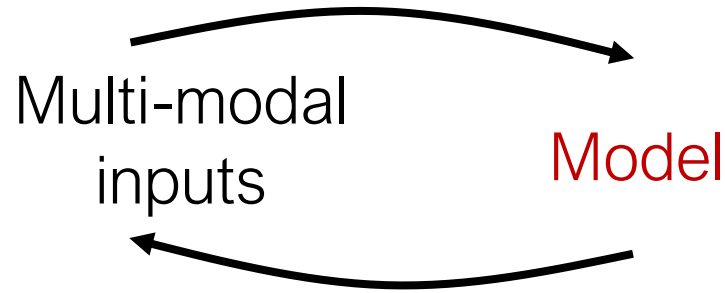
Open-sourced AirSim simulator with ChatGPT API:



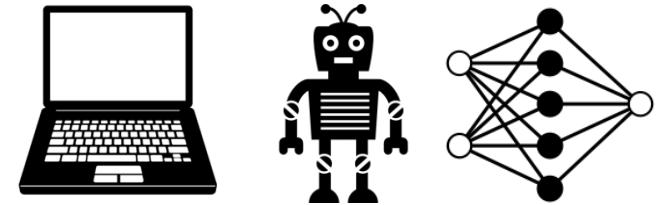
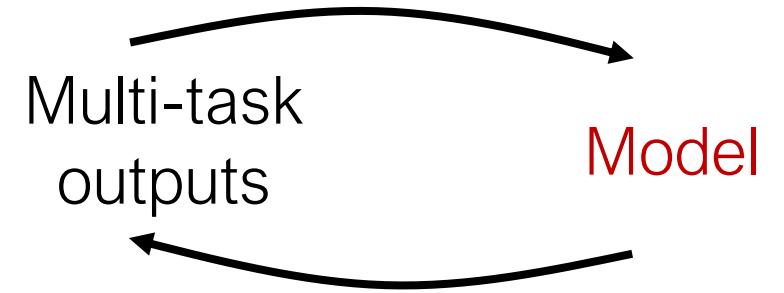
Goal: use large models to **empower any user** with the capabilities of generative AI



Close the loop between the human user and the model

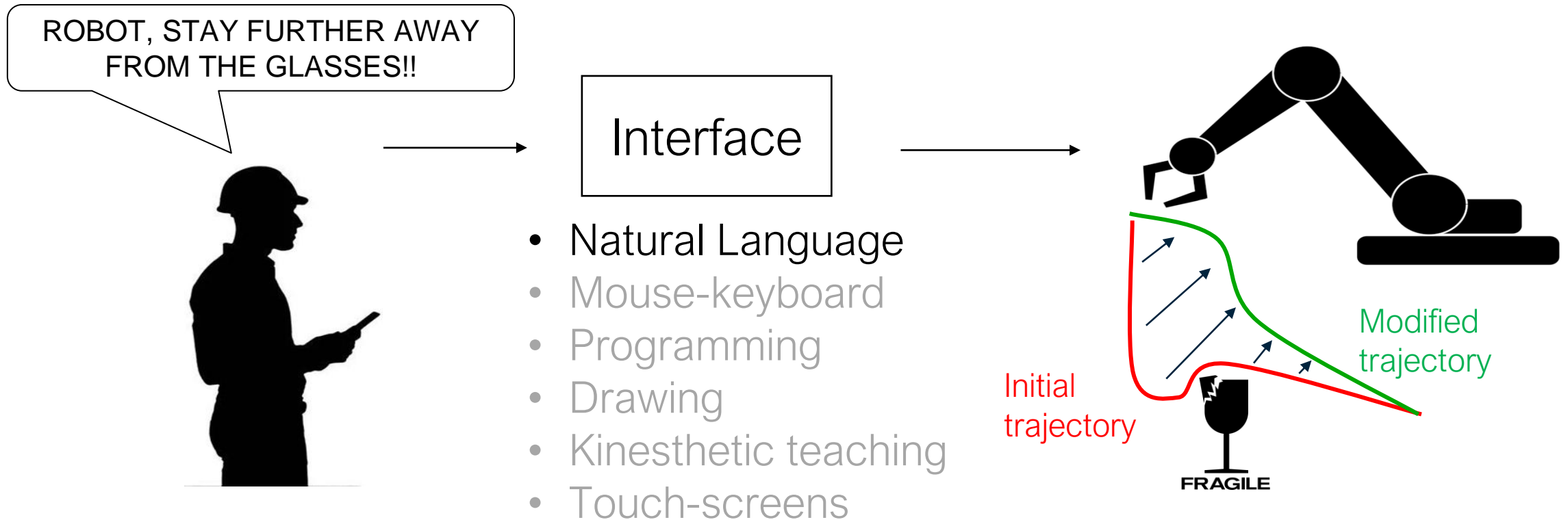


Fusing multi-modal information in I/O



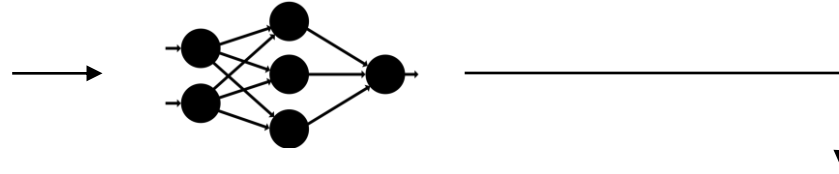
Generalizing a single model to multiple tasks

Goal: reshape robot trajectories based on user **language inputs**




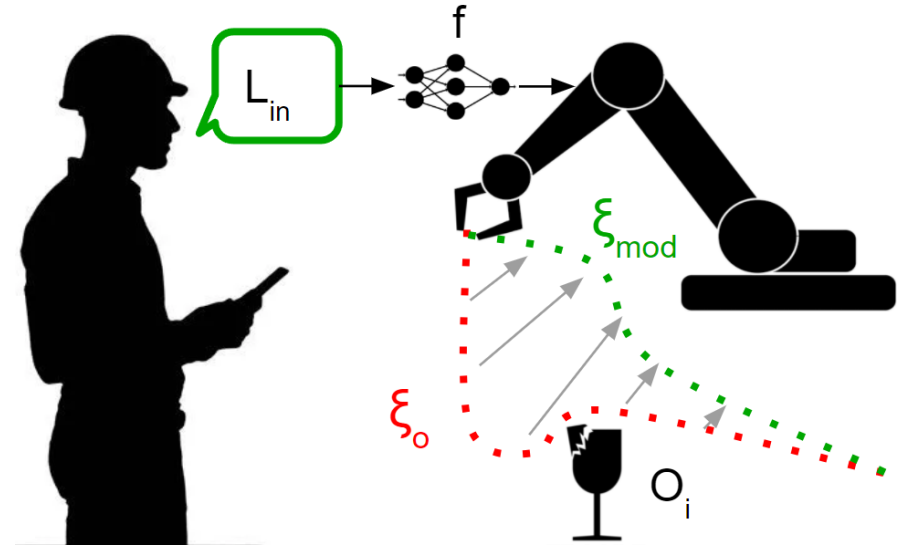


ROBOT, STAY FURTHER AWAY FROM
THE GLASSES!!

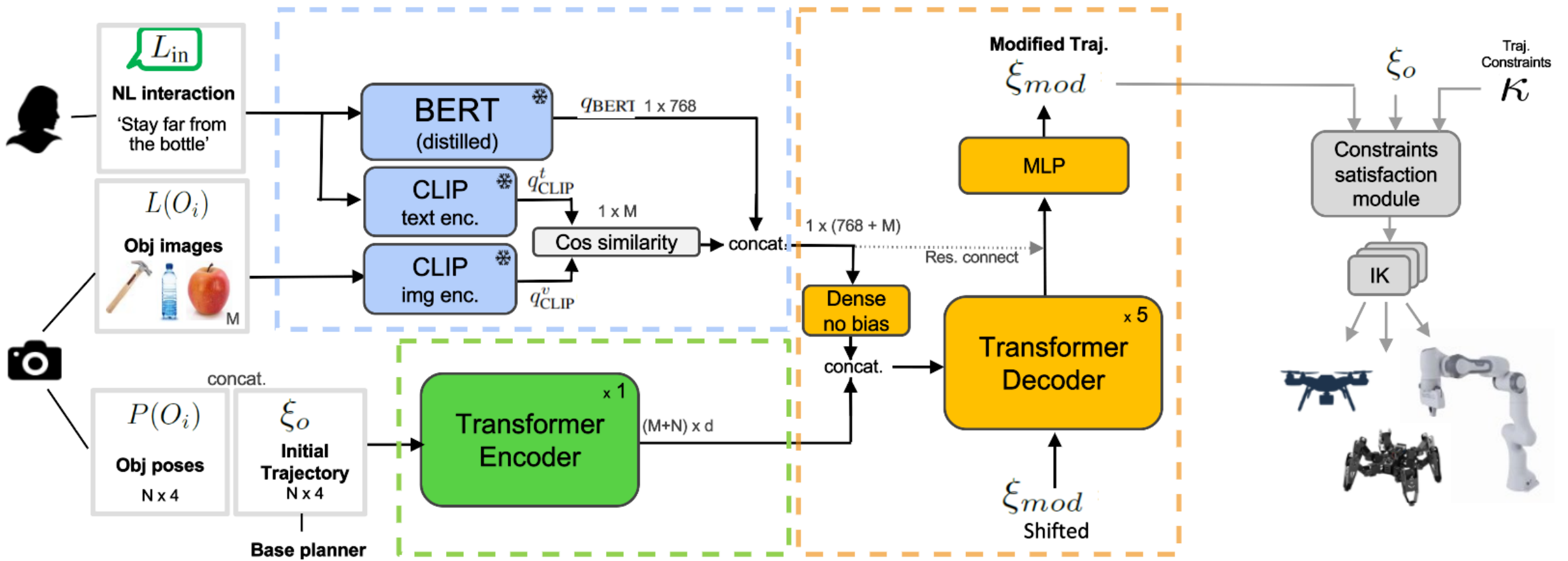


Mathematical definitions

- Original trajectory: $\xi_o : [-1, 1] \rightarrow \mathbb{R}^4 \quad \{(x_1, y_1, z_1, v_1), \dots, (x_N, y_N, z_N, v_N)\}$
- Objects in the scene: $\mathcal{O} = \{O_1, \dots, O_M\}$
 - Position $\vec{P}(O_i) \in \mathbb{R}^3$
 - Image label: $I(O_i)$ 
- Language input: L_{in}
- Learning objective: $\xi_{mod} = f(\xi_o, L_{in}, \mathcal{O})$

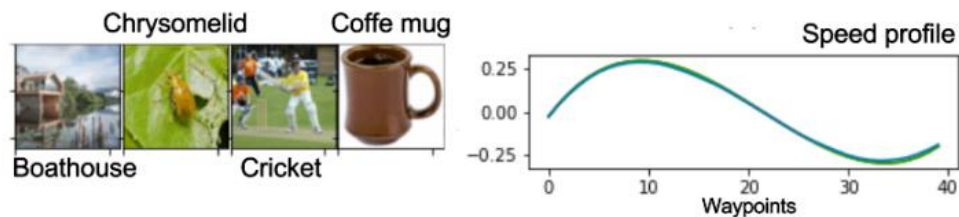
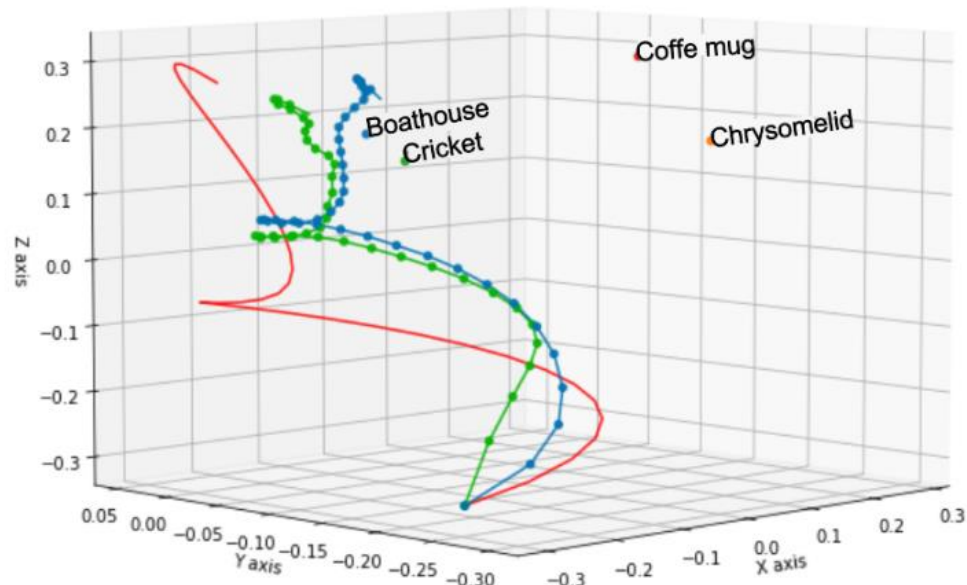


System overview

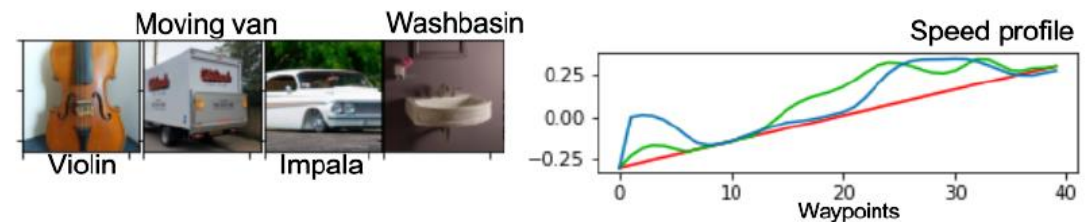
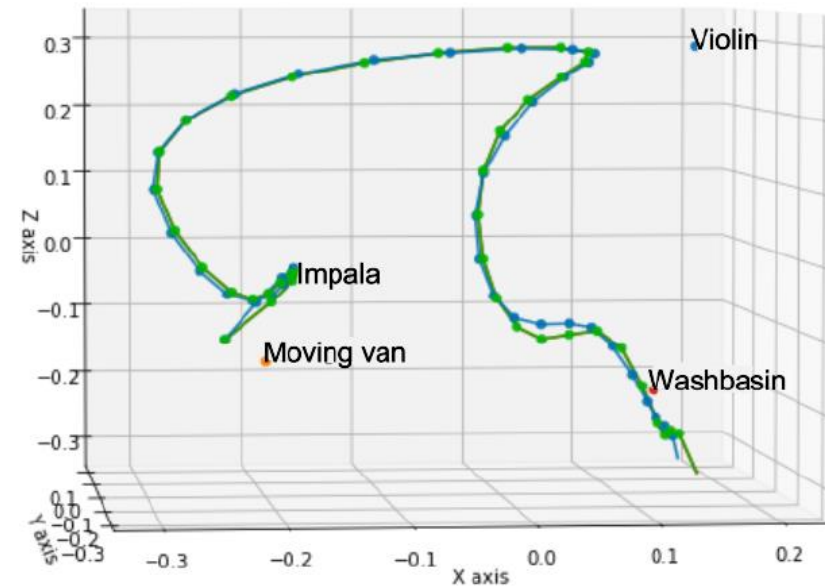


It is **expensive** to collect real-world trajectory correction examples. We use **procedural** data generation.

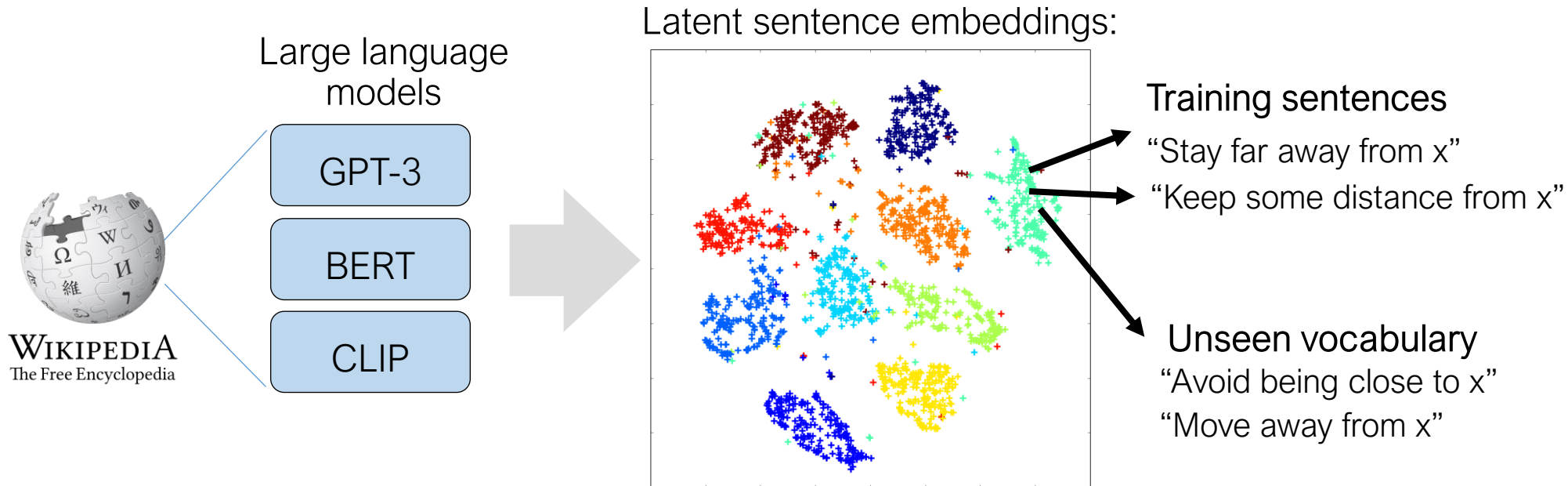
"Drive much closer to the boathouse"



"walk faster when passing next to the violin"



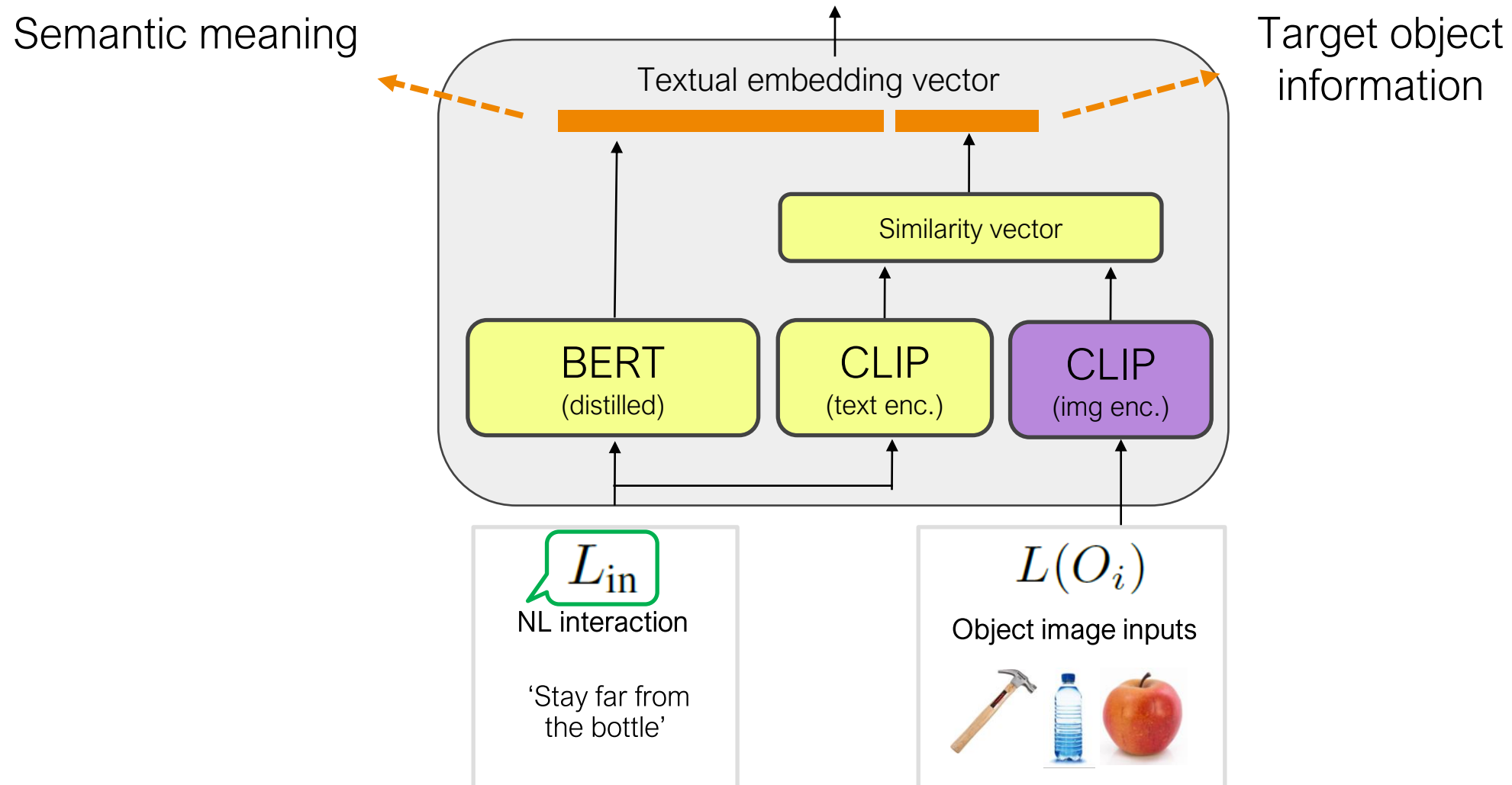
Key idea: use frozen LLM embeddings to compensate for low synthetic dataset variations

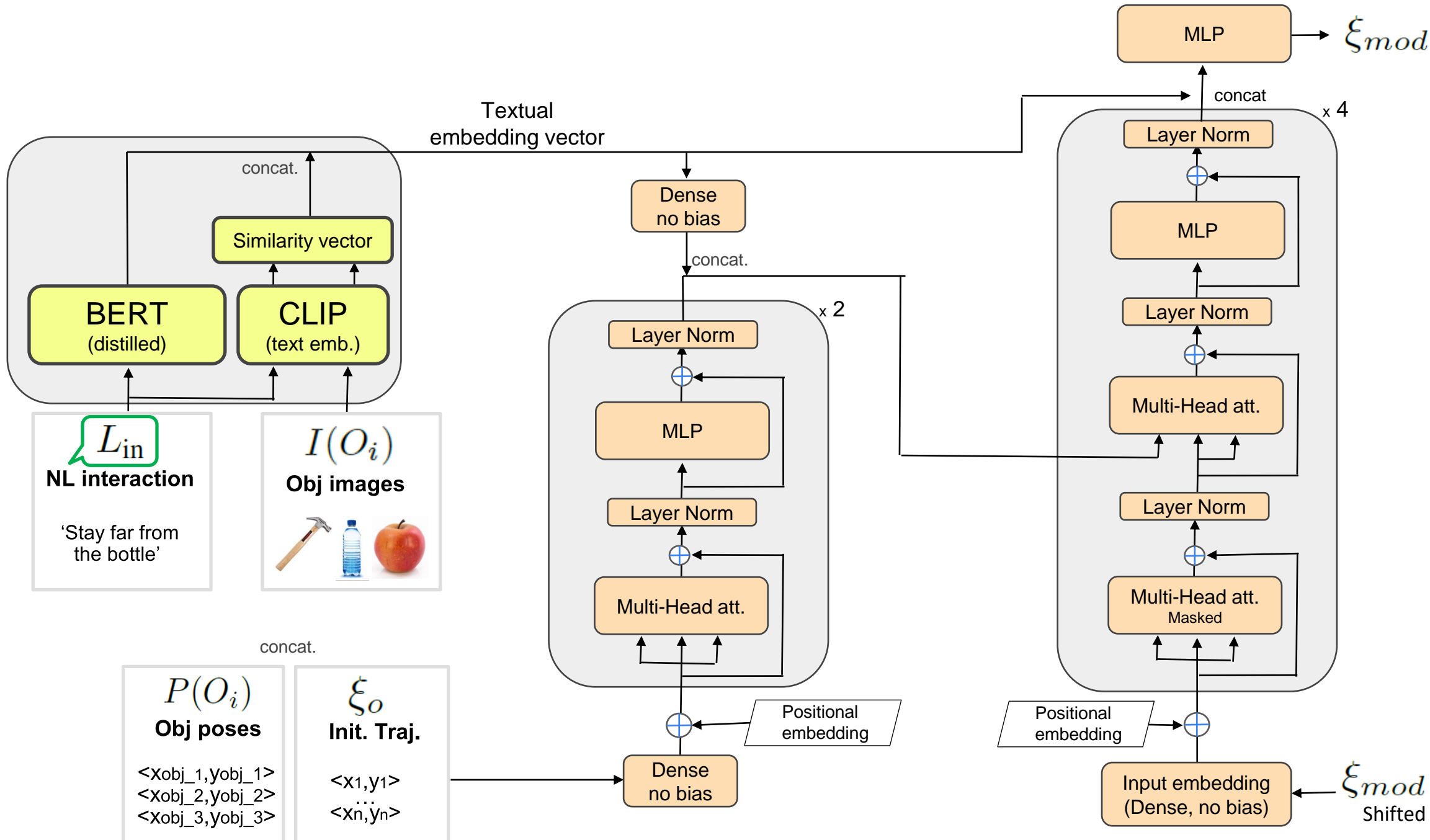


Advantages:

- Richer representations
- Less training data required
- More robust to vocabulary variations

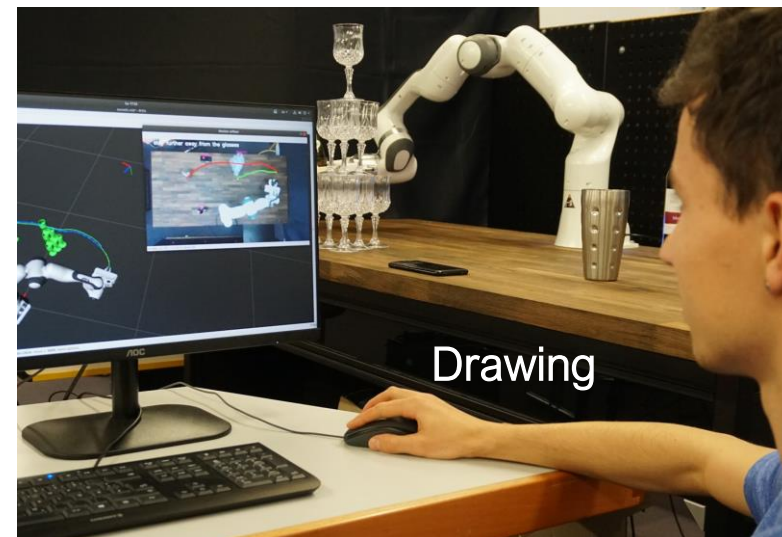
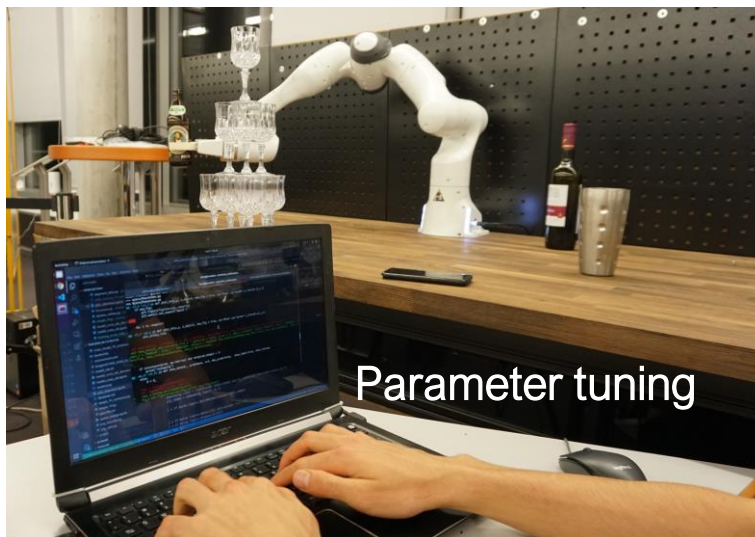
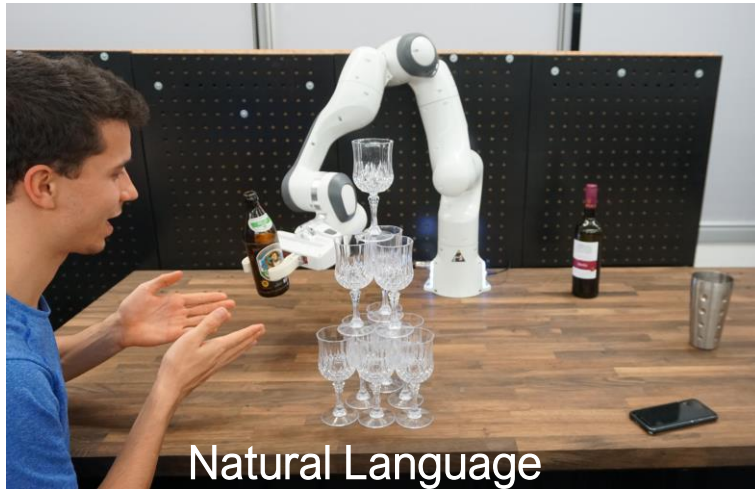
Semantic sentence embedding



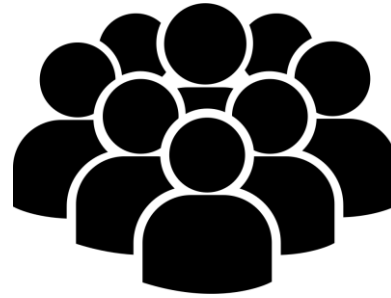


LaTTe: Language Trajectory TransformEr

User study: comparison between interfaces



Results: multi-model model is **effective, precise and intuitive**

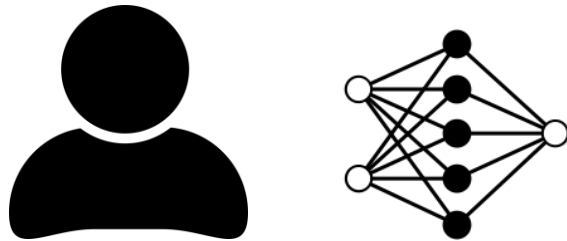
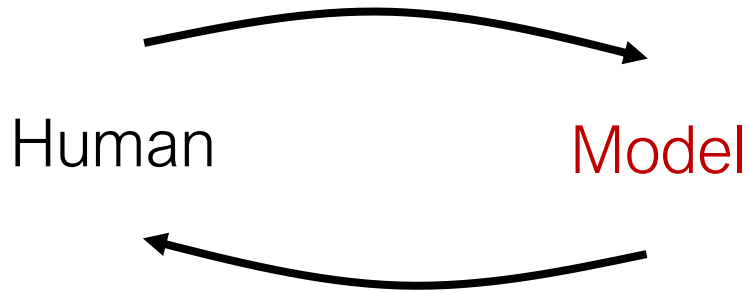


Interface	Satisfied	Easy to use	Safety	Natural	Predictable	Av. Iterations	Success rate (%)	Avr. Time (s)
Natural Language (ours)	90	92	92	98	72	1.33	100	81
Kinesthetic Teaching	90	88	88	78	96	1.78	56	139
Drawing	88	74	100	80	88	1.89	65	120
Parameter tuning	62	58	88	62	48	4.00	92	284

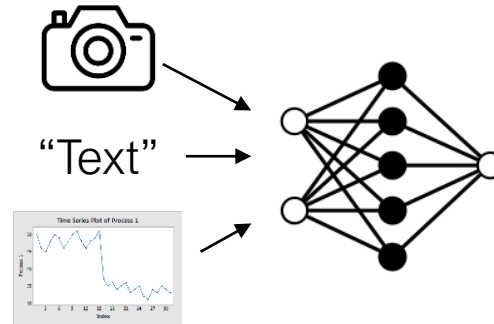
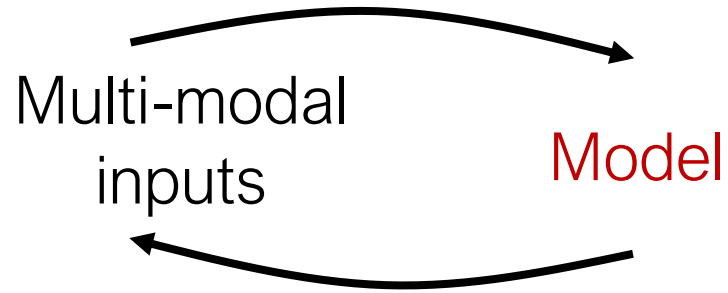
User ratings collected in the user study

User study metadata

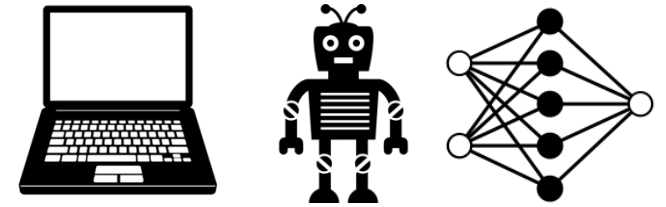
Goal: use large models to **empower any user** with the capabilities of generative AI



Close the loop between the human user and the model



Fusing multi-modal information in I/O



Generalizing a single model to multiple tasks

General Pre-Training for Decision-Making

Language models

Input:

Recite the first law of robotics

GPT-3

Output:

Sequence of words

A robot may not injure a human being or, through inaction allow a human to be harmed...

Autonomous systems models

Input:

Images/
Depth



LiDAR



Language:
'Go fast!'

PACT

Perception-Action
Causal Transformer

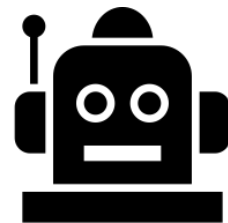
Output:

Sequence of actions

$\theta_0 = 10^\circ, V_0 = 0.5m/s$
 $\theta_1 = 15^\circ, V_1 = 0.1m/s$
 $\theta_2 = 18^\circ, V_2 = -0.7m/s$

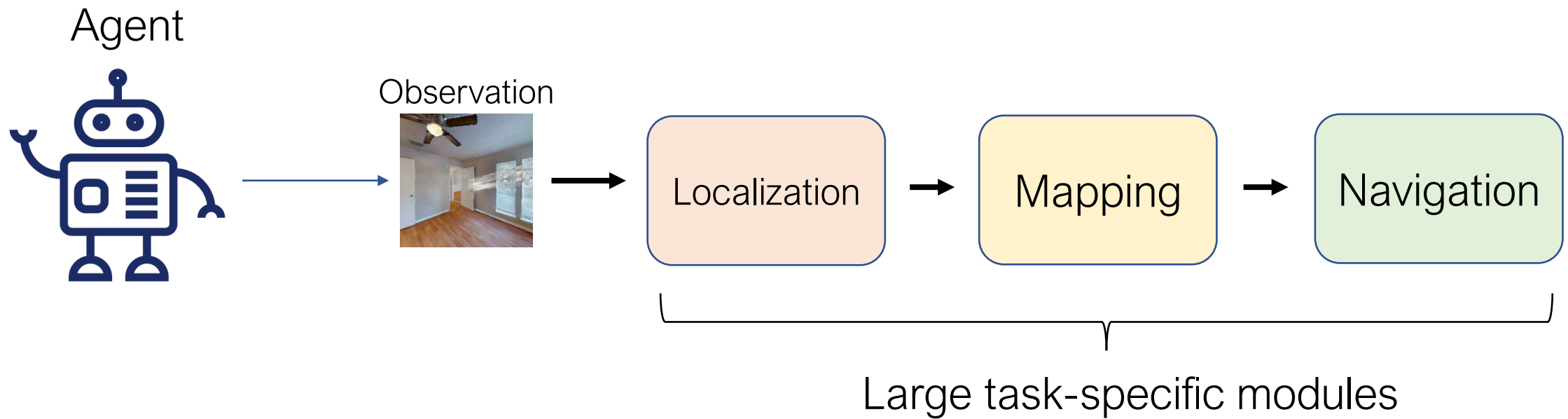
New sensor inputs

Environment

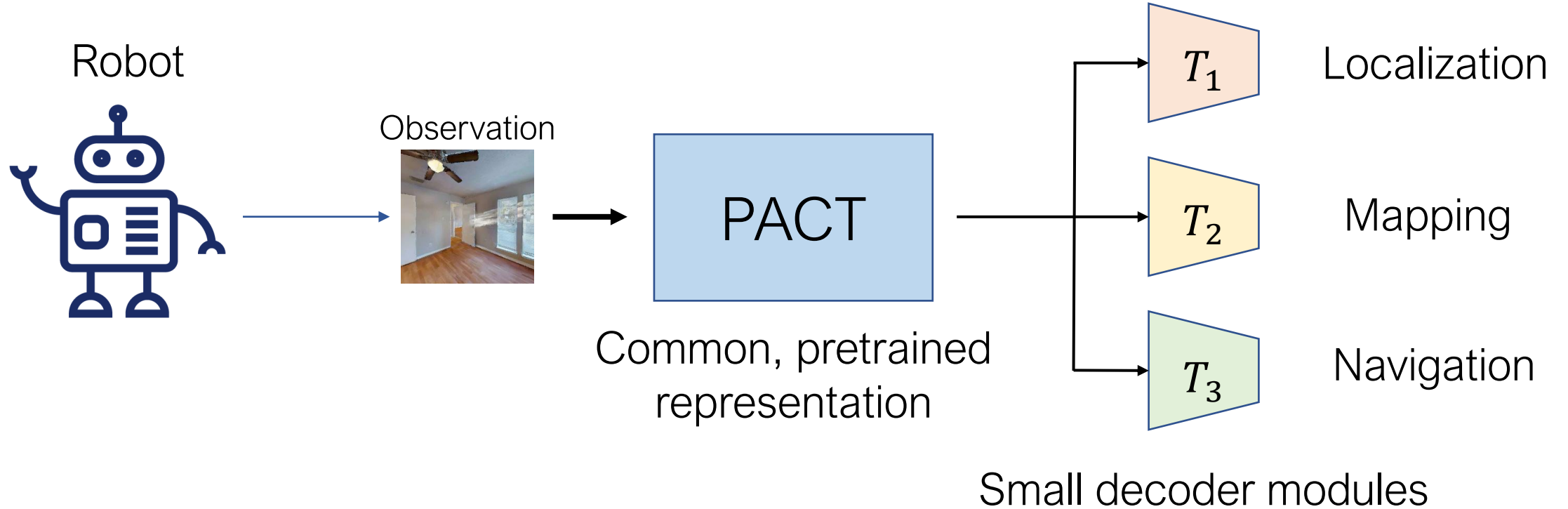


Actuators

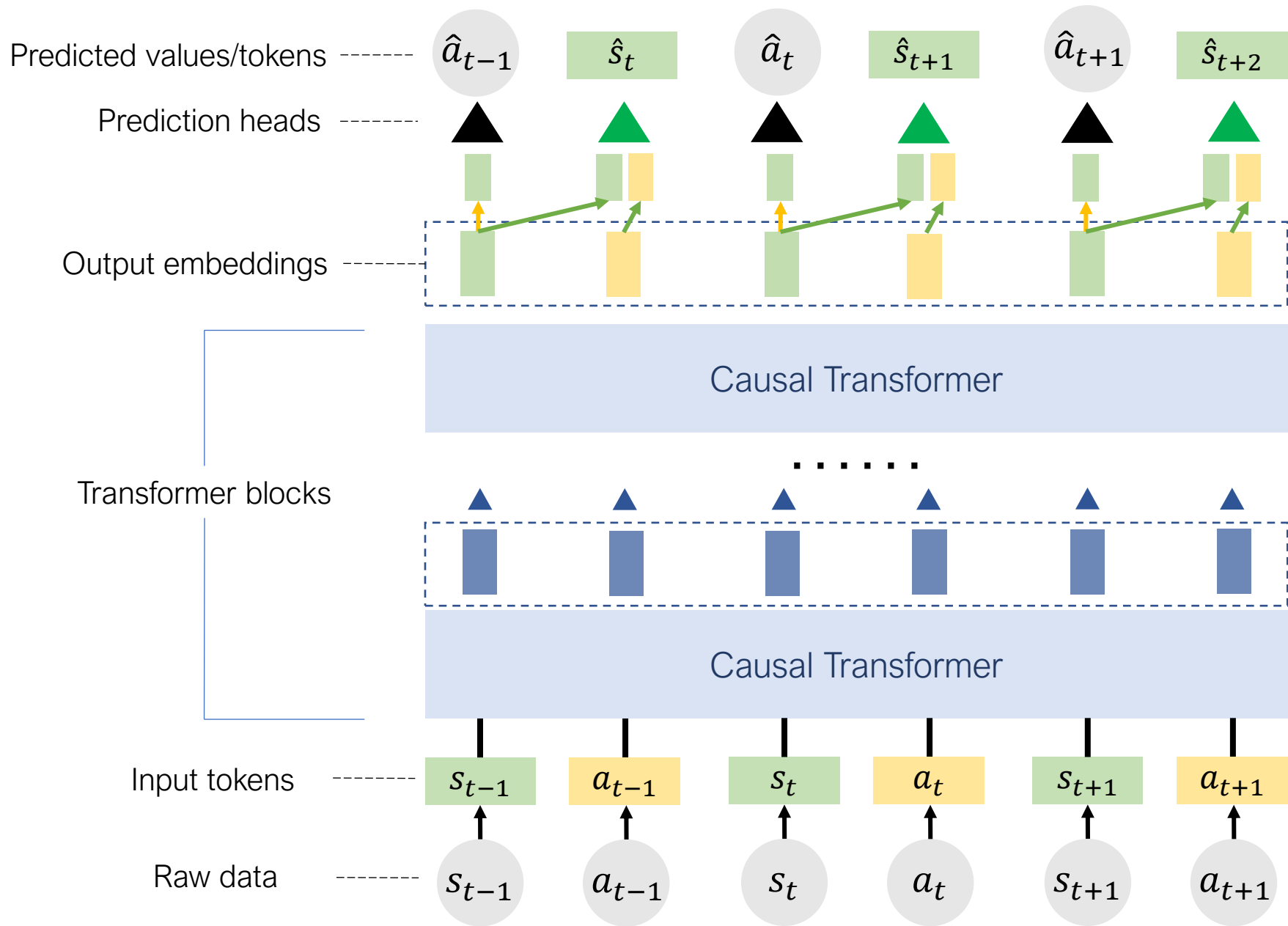
Traditional Autonomous Systems Pipeline



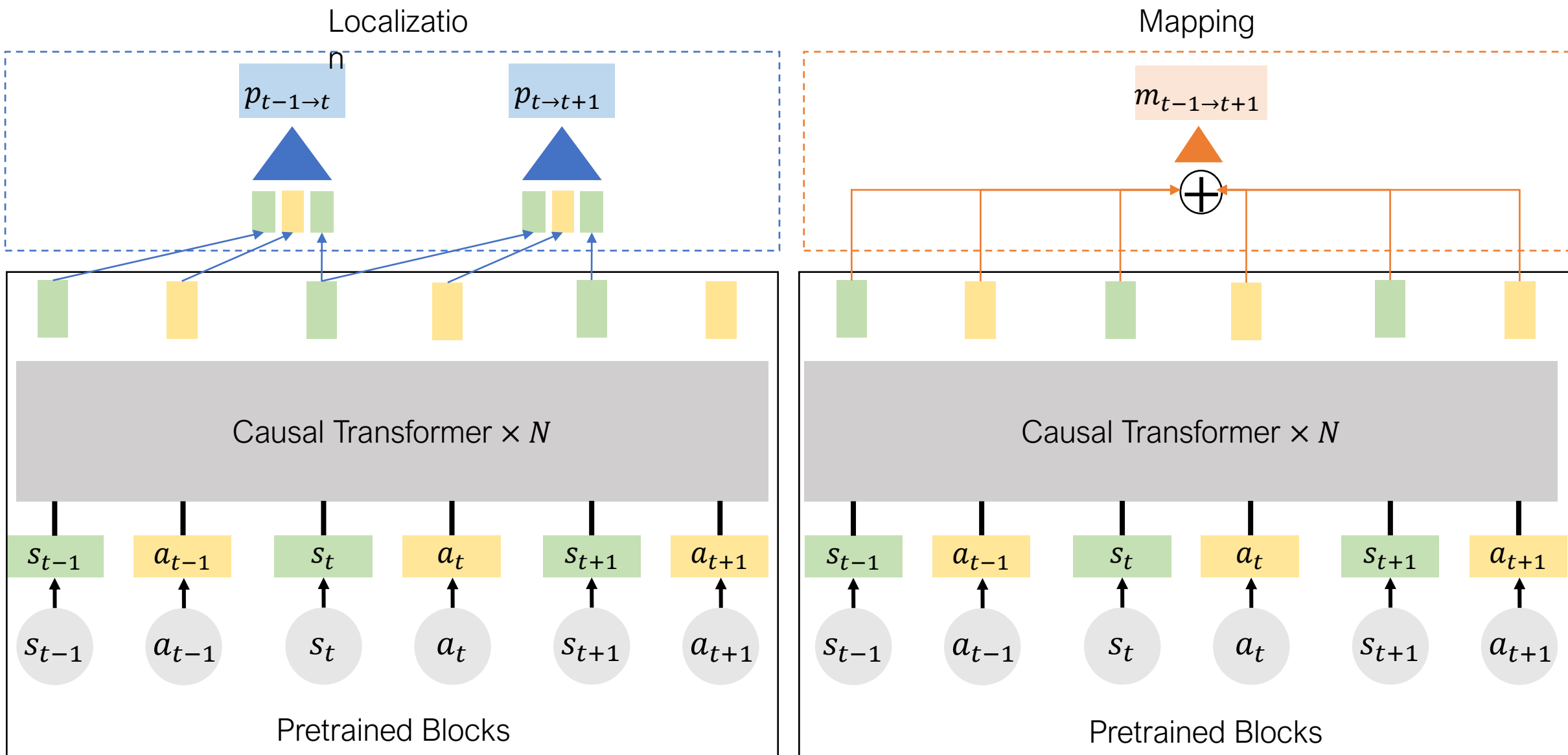
Foundational models for decision-making: Deployment in multiple downstream tasks



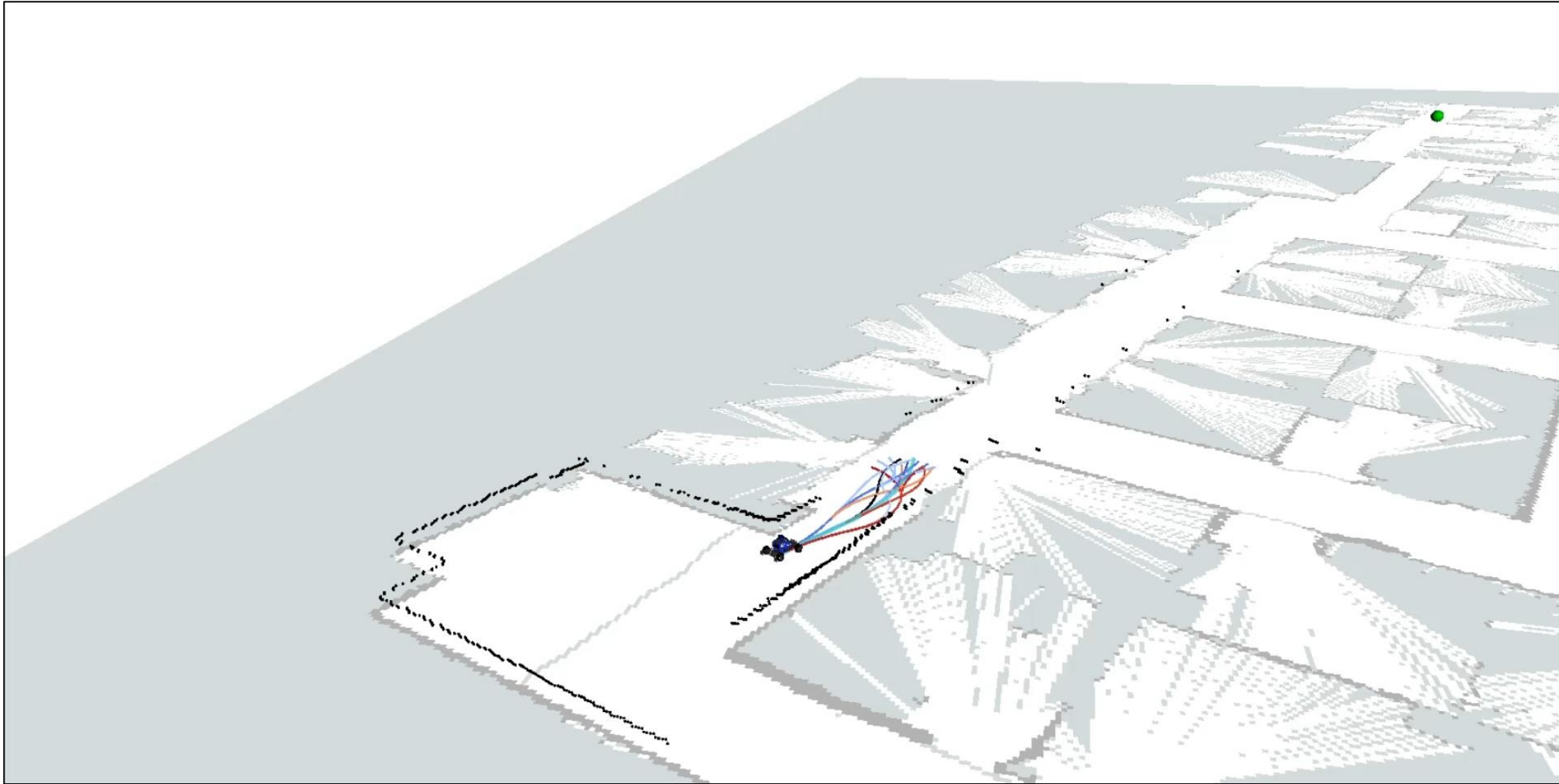
Pre-training network:



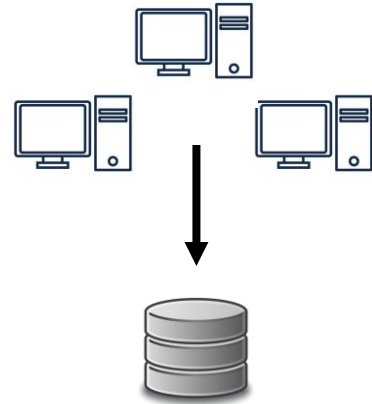
Downstream task networks:



Pre-training data collection for MuSHR: millions of perception-action pairs

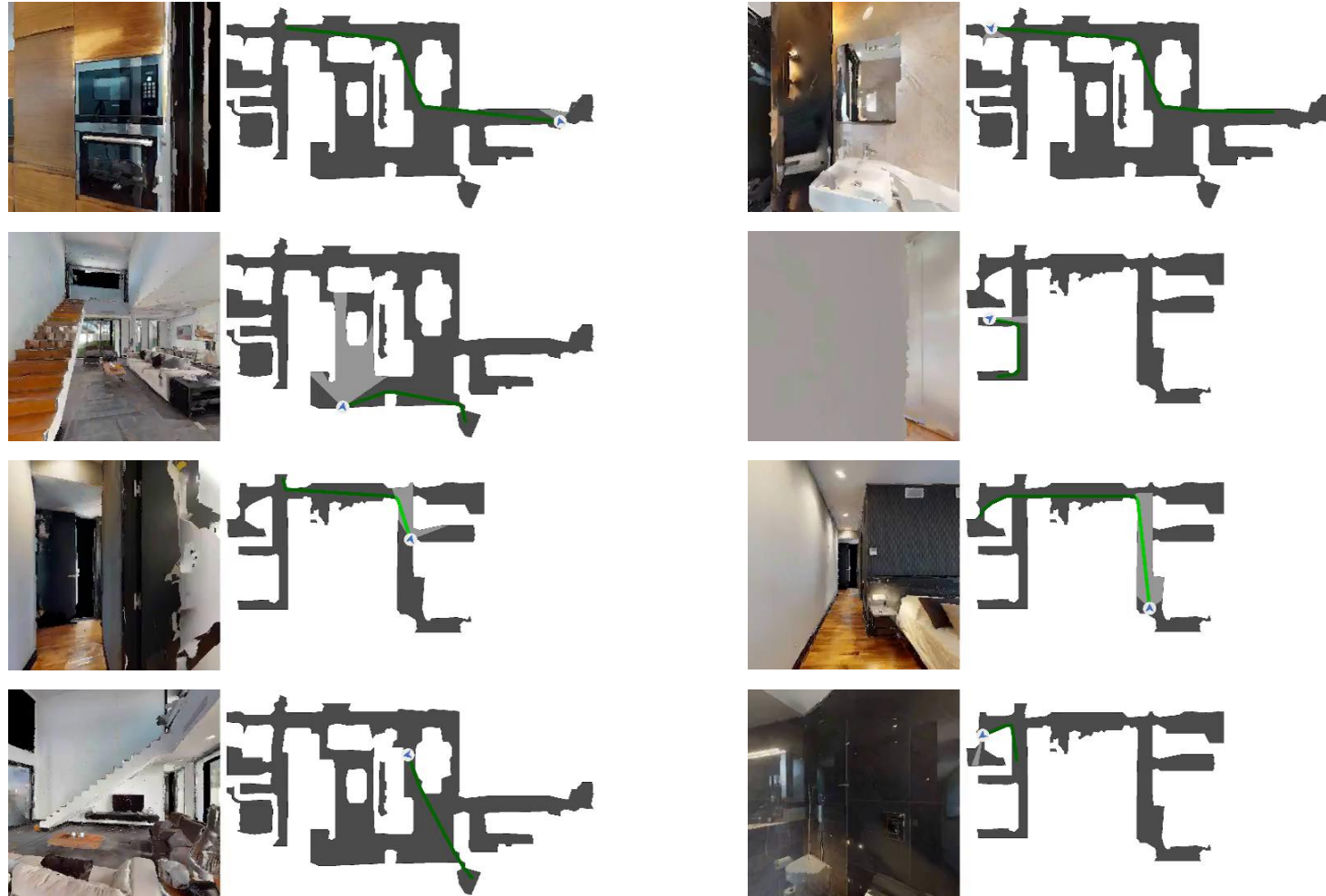


Parallel data
collection

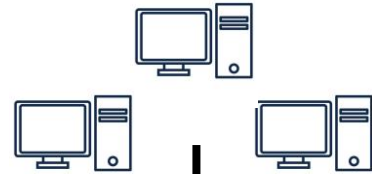


Perception-
Action Pairs

Pre-training data collection for Habitat: millions of perception-action pairs

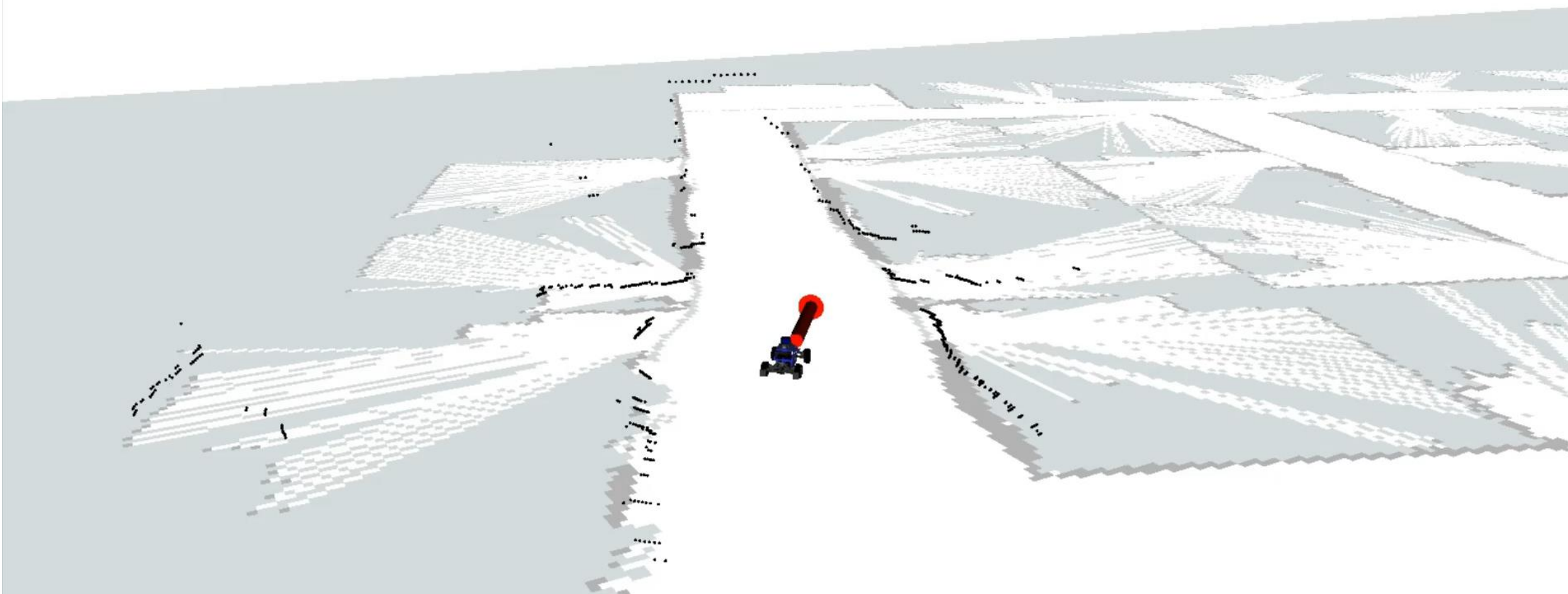


Parallel data
collection



Perception
-Action
Pairs

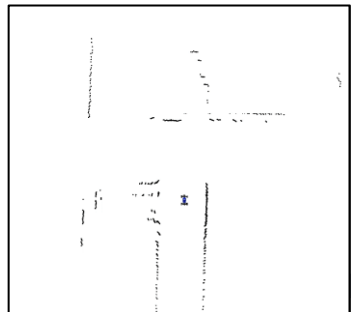
What does the pre-trained model learn?



PACT applied towards multiple downstream tasks

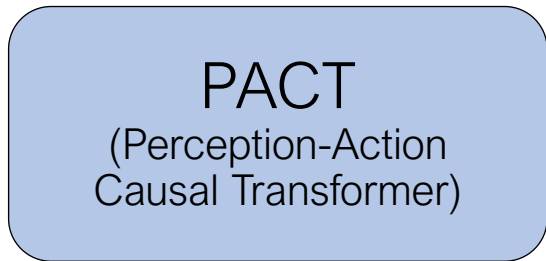
MuSHR
vehicle:

Perception
input: LiDAR

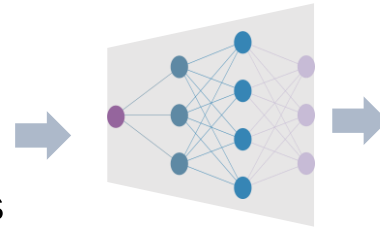
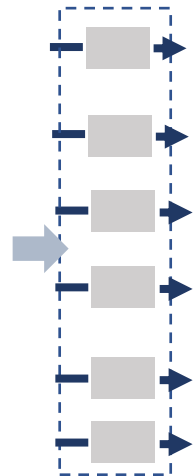


Action input:
wheel
angles

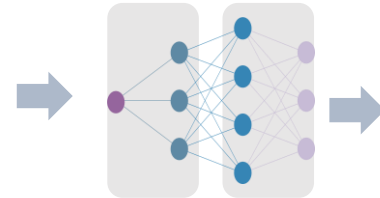
Common Representation



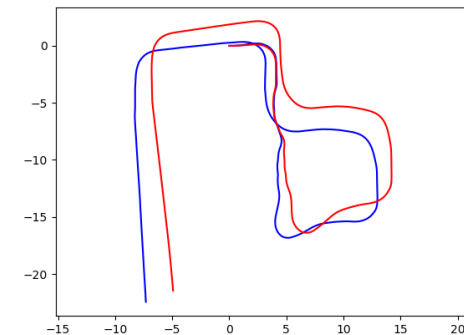
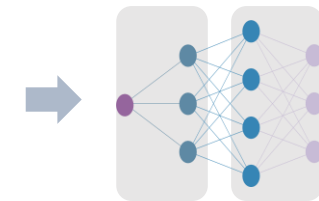
Embeddings



Local
Mapping



Real-World
Navigation

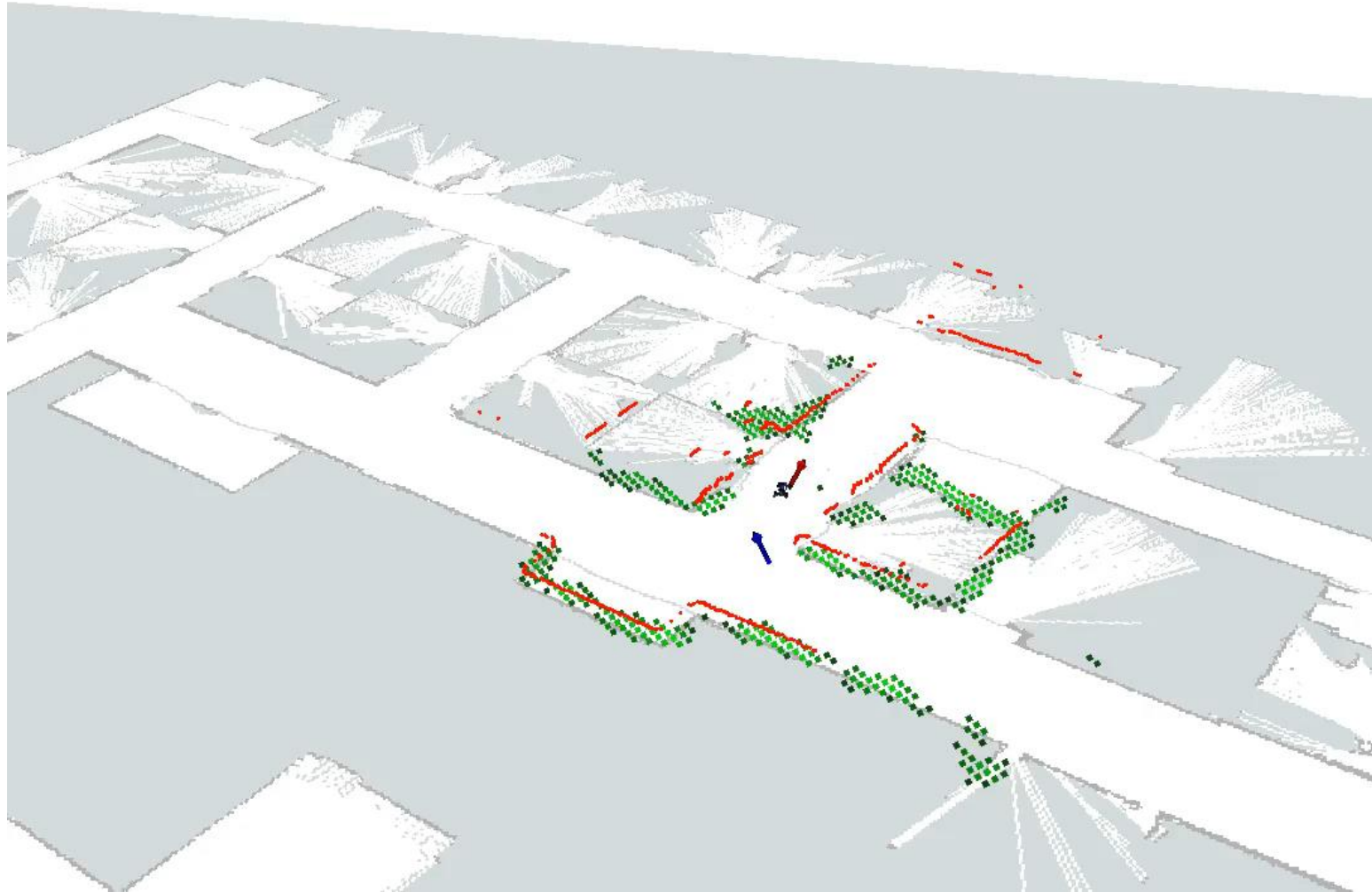


Localization

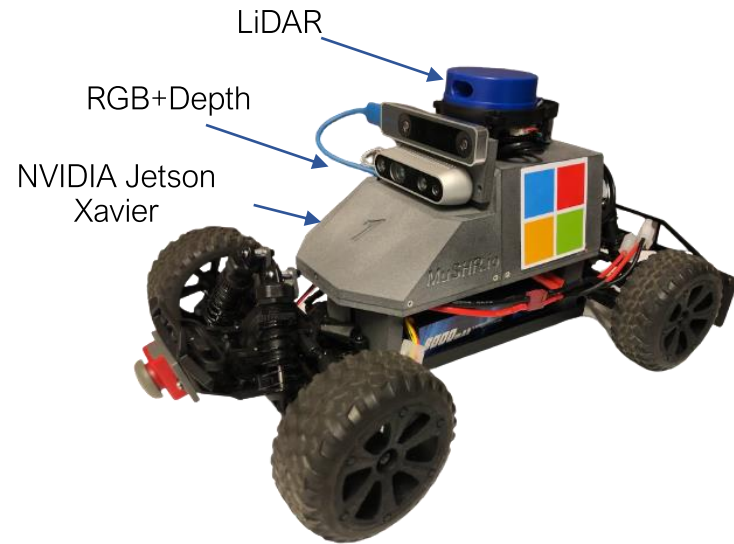
— GT pose
— Prediction

Simultaneous navigation, localization and mapping model deployment in simulation

→ Reference frame for localization (zero) → Estimated vehicle pose (integration over time) ■ Decoded local map

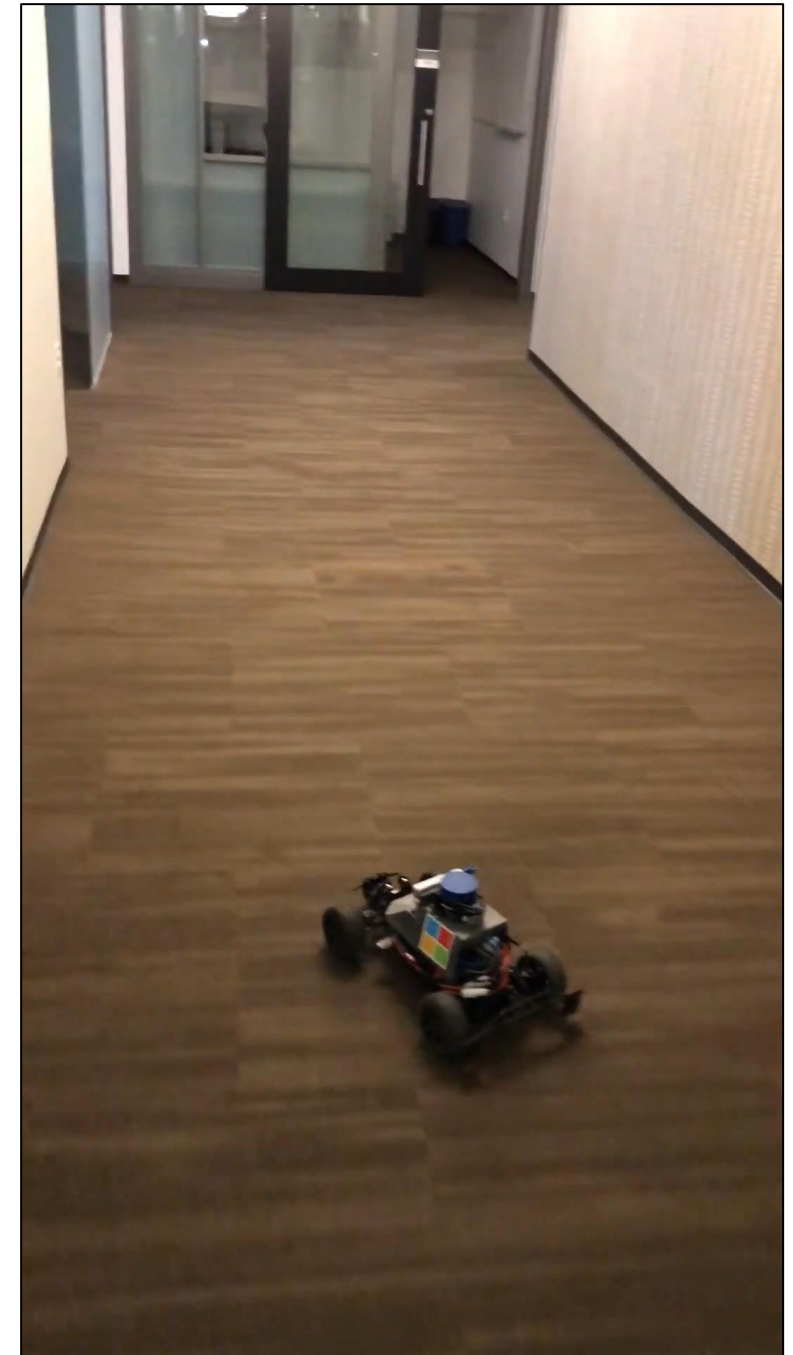


MuSHR deployment **in the wild**:

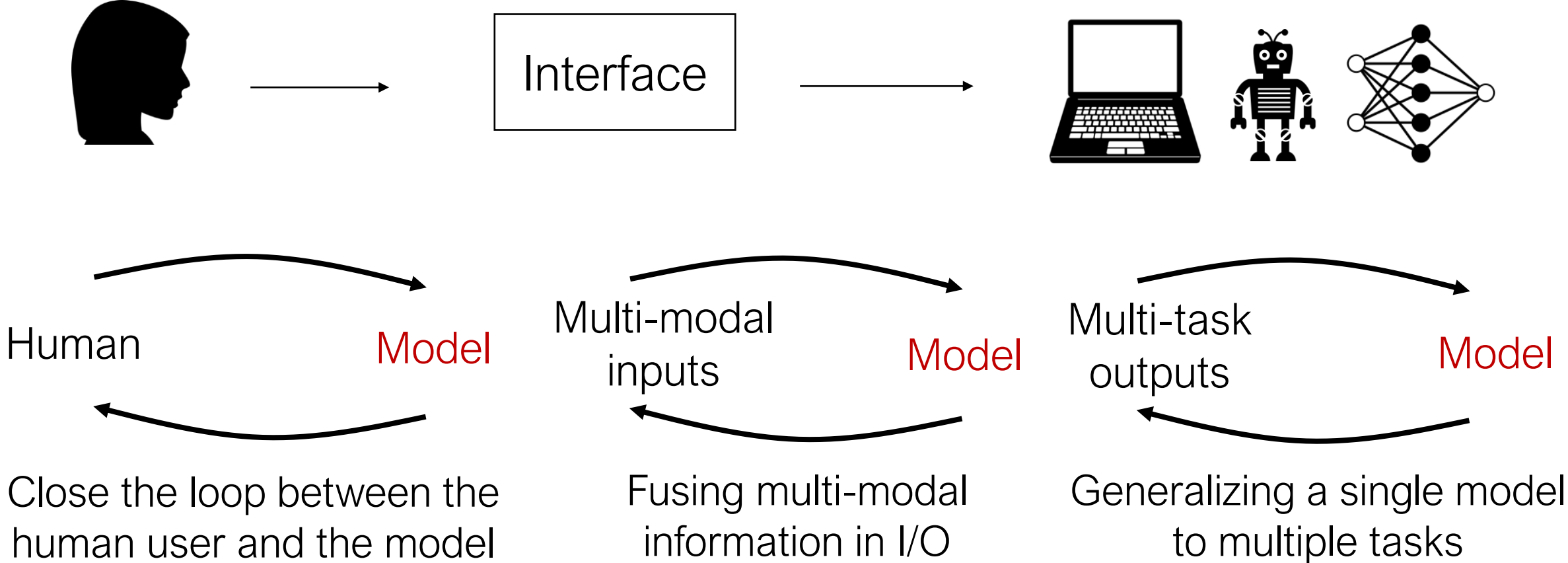


Pre-trained model overcomes sim-to-real gap:

- Model dynamics
- Actuator and processing delays
- Sensor noise and imperfect LiDAR returns (e.g. glass surfaces)



Use large models to **empower any user** with the capabilities of generative AI



Thanks to collaborators:

- Microsoft:

- Ashish Kapoor
- Sai Vemprala
- Shuang Ma
- Shuhang Chen
- Felipe Frujeri



- TUM, Germany:

- Arthur Bucker
- Luis Figueredo
- Sami Haddadin



Q&A + discussions