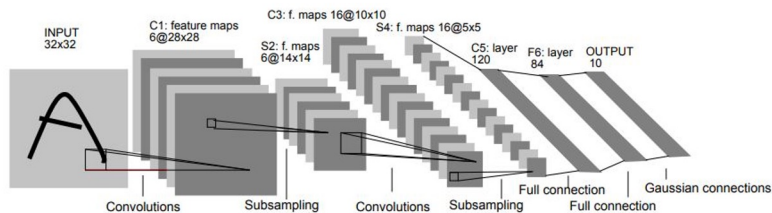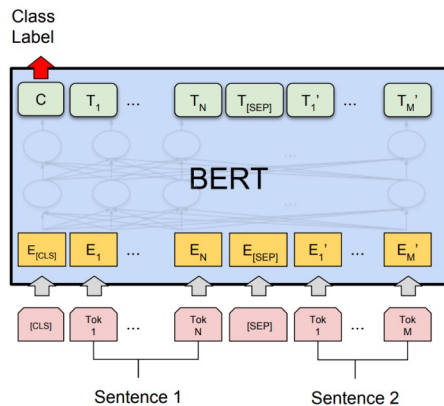# Generalist Embodied AI in an Open World

Xiaojian Ma
Machine Learning @ BIGAI
11/24/2023

# ML is stepping into a new era



1990s            2010s            2020s

- More data, O(10k) -> O(10M) -> O(1T);

- More parameters, O(1M) -> O(1b) -> O(100b);

- More computation, GFLOPS -> TFLOPS
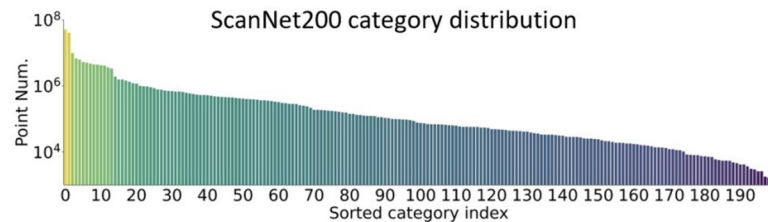
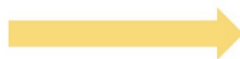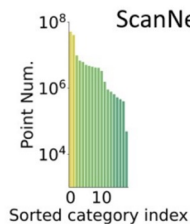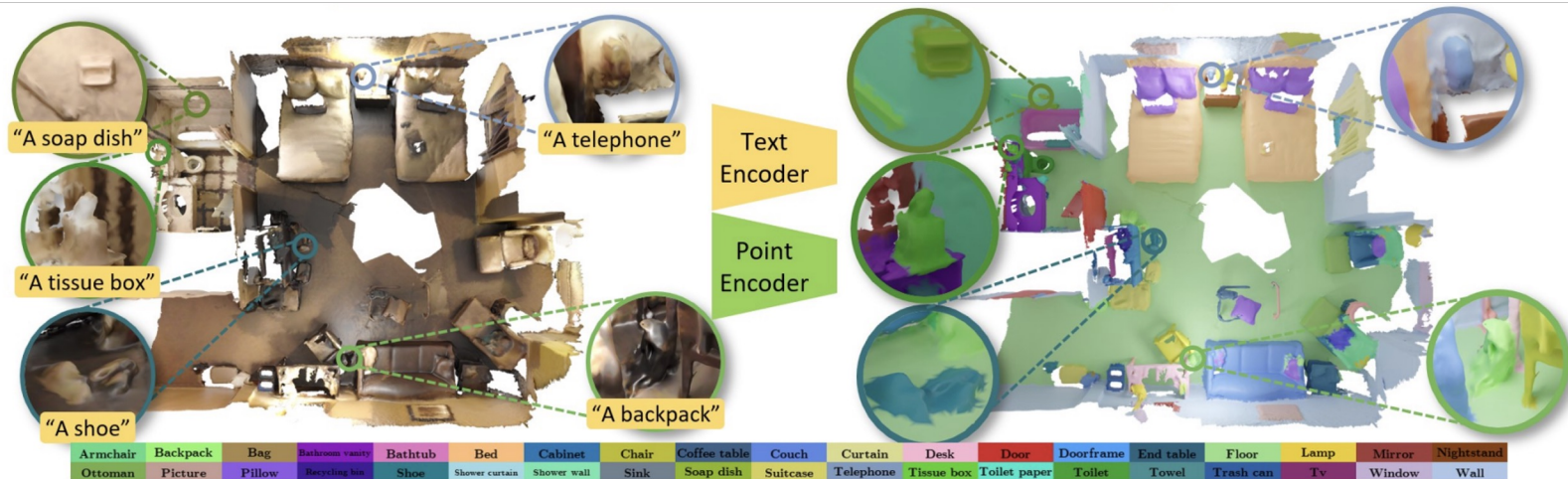# ML is stepping into a new era



1990s



2010s



2020s

- Complex domains and semantics

- Close world (vocabulary) -> open world (vocabulary)
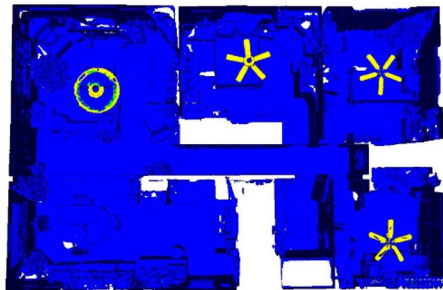
# …and so is embodied AI



ScanNet-200: Language-Grounded Indoor 3D Semantic Segmentation in the Wild

# ...and so is embodied AI



Input 3D Point Cloud

"fan" - Object

"made of metal" - Material

"kitchen" – Room Type

Zero-shot Semantic Segmentation

"anything soft" - Property

"where to sit" - Affordance

"work" - Activity

OpenScene: 3D Scene Understanding with Open Vocabularies

# ...and so is embodied AI



Top-down visualization

Task: Find the gingerbread house

CoWs on Pasture:
Baselines and Benchmarks for Language-Driven Zero-Shot Object Navigation

# A paradigm shift for embodied AI



NeurIPS 2023 HomeRobot: Open Vocabulary Mobile Manipulation (OVMM) Challenge



contrived; limited tasks;
static; close world…

realistic; massive tasks;
dynamic open world…

# Embodied AI

Machine Learning @ BIGAI

# Generalist Embodied AI in an Open World

LEO: An Embodied Generalist Agent in 3D World

CraftJarvis: Multi-task Embodied Agents in an Open World

stay tuned…

LEO: An Embodied Generalist
Agent in 3D World

CraftJarvis: Multi-task Embodied
Agents in an Open World

stay tuned…

?

# LEO: An Embodied Generalist Agent

An Embodied Generalist Agent in 3D World,
arXiv preprint 2023

3D-VisTA: Pre-trained Transformer for 3D Vision and Text Alignment,
ICCV 2023

SQA3D: Situated Question Answering in 3D Scenes,
ICLR 2023

embodied-generalist.github.io

# Embodied Generalist Agent

**Capabilities**: *Perception*, *Grounding*, *Reasoning*, *Planning*, *Acting*

## Tasks

| | | | |
|---|---|---|---|
| **3D Object Captioning** | **3D Question Answering** | **3D Dialogue** | **Embodied Navigation** |
| **Scene Captioning** | **Embodied Reasoning** | **Task Planning** | **Robotic Manipulation** |

## 3D World

# Single-modal agent

**System Message**

You are an assistant situated in a 3D scene.

Tokenizer

↓

**Instruction**

User: Please set up a home office.

Tokenizer

↓

## Large Language Model

1. Choose a room…
2. Plan the layout…
3. Create a workspace…

**Without scene awareness:**
ambiguous, hallucination

---

**Human**

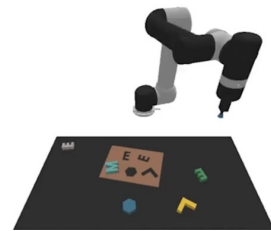Task: Throw away paper
Step 1: Walk to home office
Step 2: Walk to desk
Step 3: Find desk
Step 4: Turn to desk
Step 5: Find chair
Step 6: Sit on chair
Step 7: Find check
Step 8: Grab check
Step 9: Squeeze check
Step 10: Stand up
Step 11: Walk to trashcan
Step 12: Put check on trashcan

**GPT-2 1.5B**

Task: Brush teeth
Step 1: Go to bathroom

**GPT-3 175B**

Task: Throw away paper
Step 1: Walk to home office
Step 2: Walk to wastebasket
Step 3: Find wastebasket
Step 4: Grab wastebasket
Step 5: Walk to desk
Step 6: Drop paper in wastebasket

**Codex 12B**

Task: Brush teeth
Step 1: Walk to bathroom
Step 2: Walk to sink
Step 3: Find toothbrush
Step 4: Pick up toothbrush
Step 5: Put toothbrush in mouth
Step 6: Move brush around mouth for two minutes
Step 7: Spit out toothpaste and brush into sink
Step 8: Turn on water in sink and rinse brush for one minute
Step 9: Turn off water in sink and return brush to cupboard

**Translated GPT-3 175B**

Task: Throw away paper
Step 1: Walk to home office
Step 2: Walk to table
Step 3: Find table
Step 4: Turn to table
Step 5: Find paper
Step 6: Grab paper
Step 7: Walk to trashcan
Step 8: Open trashcan
Step 9: Put paper on trashcan
Step 10: Close trashcan

**Translated Codex 12B**

Task: Brush teeth
Step 1: Walk to bathroom
Step 2: Open door
Step 3: Walk to sink
Step 4: Put pot on sink
Step 5: Put brush on toothbrush
Step 6: Turn to toothpaste
Step 7: Put toothpaste on toothbrush
Step 8: Put teeth on toothbrush

Language Models as Zero-Shot Planners: Extracting Actionable Knowledge for Embodied Agents

## scene to text



### System Message
You are an assistant situated in a 3D scene.

Tokenizer
↓

### Scene Caption
There is a white table, wooden …

Tokenizer
↓

### Instruction
User: Please set up a home office.

Tokenizer
↓

## Large Language Model

Tedious text, intractable to embed complex 3D information

---

**Robot Planning & Interaction**

Human
Can you bring me the drink from the table?

Robot
Action: "go to table"

Robot
Do you want water or coke?

Human
Coke please.

Robot
Action: "pick up the coke"

Robot
Action: "pick up the coke"

Robot
Action: "bring it to you"

**Grounded Closed-Loop Feedback**

Robot
Scene Descriptor      Success Detector
Human

Scene Descriptor
I see: coke, water, chocolate bar.

Success Detector
Action was not successful.

Success Detector
Action was successful.

Inner Monologue: Embodied Reasoning through Planning with Language Models

| System Message | 3D Scene | Instruction |
|---|---|---|
| You are an assistant situated in a 3D scene. |  | User: Please set up a home office. |
| Tokenizer | 3D Encoder | Tokenizer |
| ↓ | ↓ | ↓ |

**Large Language Model**

1. Place the desk in the desired position in the room…
2. Position the chair next to the desk, to the right of it.
3. Set up the shelf to the left of the desk…
4. Place the lamp on the desk…
5. Arrange the showcase to the right of the desk.
6. Place the plants on the shelf…
7. Hang the curtains on the wall behind the desk…

Scene-aware agent with capacity of perceiving (3D) scenes

# Scene representation



2D branch

2D Encoder

Ego-centric features

d

h

w

3D branch

3D Encoder

Object-centric point clouds

Object-centric features

...

...

# Embodied Generalist Agent in 3D World

| **System Message** | **Egocentric Image** | **3D Scene** | **Instruction** |
|---|---|---|---|
| You are an assistant situated in a 3D scene. | | | User: Please describe the toy house over ... |
| ❄ Tokenizer | 🔥 2D Encoder | 🔥 3D Encoder | ❄ Tokenizer |
| ↓ | ↓ | ↓ | ↓ |

## Large Language Model

LoRA🔥$\delta$

↓

| **Text Response** | | **Action Response** | |
|---|---|---|---|
| There is a sofa next to the TV. | It's a kitchen for cooking. | $P = [0.1, -0.2, 0]$ <br> $R = [0, 0, 0, 1]$ | "Turn right" |

← De-tokenize

**Unified task sequence**

$$\underbrace{\texttt{You are...}}_{\text{system message}} \underbrace{s_{\text{2D}}^{(1)},...,s_{\text{2D}}^{(M)}}_{\substack{\text{2D image tokens} \\ \text{(optional)}}} \underbrace{s_{\text{3D}}^{(1)},...,s_{\text{3D}}^{(N)}}_{\substack{\text{object-centric} \\ \text{3D tokens}}} \underbrace{\texttt{USER:...\ ASSISTANT:}}_{\text{instruction}} \underbrace{s_{\text{res}}^{(1)},...s_{\text{res}}^{(T)}}_{\text{response}}.$$

**Auto-regressive objective**

$$\mathcal{L}(\theta, \mathcal{B}) = -\sum_{b=1}^{|\mathcal{B}|} \sum_{t=1}^{T} \log p_\theta(s_{\text{res}}^{(b,t)} | s_{\text{res}}^{(b,<t)}, s_{\text{prefix}}^{(b,1)}, ..., s_{\text{prefix}}^{(b,L)})$$
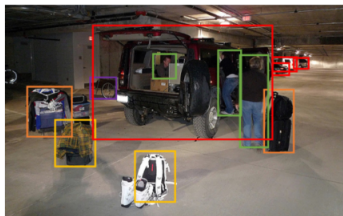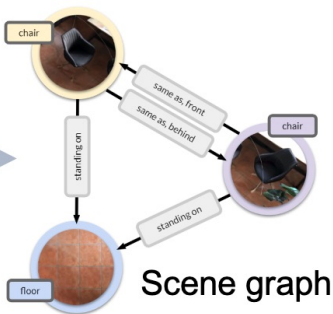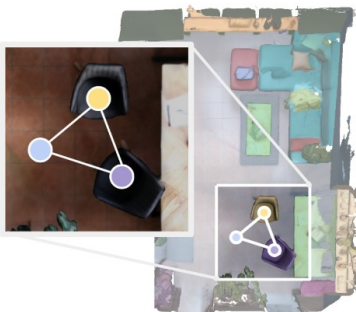
18

Captioning

Planning

Reasoning

Navigation

Dialogue

Manipulation

An Embodied Generalist
Agent in 3D World

LEO

Machine Learning @ BIGAI

# Scene-Graph-based Prompting

## Messages

**1. System Message**
You are an AI visual assistant in a 3D scene…

**2. Demonstrations**
**Scene Graph Context:**
{'sofa-1': {'attributes': {'color': 'red'}, 'relations': ['to the right of chair-2', 'in front of table-3']}, 'chair-2': {'attributes': {'color': 'brown'}, 'relations': []}}
**Human-labeled Responses**

**3. Query (a new scene)**
**Scene Graph Context**

## Responses

1) Object scene caption
2) Scene caption
3) Dialogue (O-CoT)
4) QA (O-CoT)
5) Planning

Potential Issues:
1) Wrong Answers
2) Negative Responses
3) Responses with IDs
4) Others

# Refinement Procedures

| Raw Responses | Refined Responses |
| --- | --- |
| 1. Wrong Answers<br>There are two curtains in the room. | 1. Answer Refinement<br>There are 3 curtains in the room. |
| 2. Negative Responses<br>Question: Where is the football table?<br>Thought: football table-17<br>Answer: unknown | 2. Removing Negative Responses<br>The negative responses will be removed. |
| 3. Responses with IDs<br>The kitchen cabinet-9 and kitchen counter-10 are parts of the kitchen. | 3. GPT-based rewriting<br>The kitchen features a cabinet and a counter. |
| 4. Others<br>Unnecessary contents (Thought, etc.) | 4. Others<br>Remove unnecessary contents. |

**Dialogue (O-CoT):**
Dialogue Context: high level task: organize the bedroom. low level task: check some objects.
Human: I need help organizing my bedroom. Can you assist me? Thought: no specific objects
Assistant: Of course! What specifically do you need help with?
Human: Can you check if there are any clean clothes hanging on the wall? Thought: clothes-13
Assistant: I'm sorry, but the clothes hanging on the wall are messy.
Human: Thanks for your help! Thought: no specific object
Assistant: You're welcome! If you need any further assistance, feel free to ask.

**Scene Caption:**
In this room, there is a wooden floor that is clean and flat. A tall wardrobe stand on the right side of a desk, close to a basket. The wardrobe is in front of a chair, a bed, and behind a bedside table. Clothes are hanging on a white wall, to the right of a closed curtain. The bed is wide, clean, and covered with blue blanket. The room has a comfortable and organized setting with functional furniture.

**Object Scene Caption:**
The showcase is supported by the wall and positioned behind, close to, and to the left of the shelf. The wall is flat, architectural, and dark in color. The shelf is made of wood and has a low size.

**QA (O-CoT):**
Question: Where is the printer located? Thought: printer-8
Answer: standing on the desk
Question: How many blankets are on the bed? Thought: blanket-16, blanket-17 Answer: 2
Question: What is the type of the room? Thought:wardrobe-2, desk-7, chair-11, bed-15 Answer: bedroom

**Planning:**
High-Level Task: Organize and tidy up the bedroom.
Low-Level Actions:
1. Clean the floor by sweeping to remove any dirt.
2. Make the bed by arranging the blanket and pillows.
3. Place any loose items or belongings into the basket.
4. Arrange items on the shelves and showcase in a tidy way.

**3DVL**

| | Scan2Cap (val) | | | | | ScanQA (val) | | | | | SQA3D |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | C | B-4 | M | R | Sim | C | B-4 | M | R | EM@1 | EM@1 |
| *Task-specific models* | | | | | | | | | | | |
| Scan2Cap (GPT-3) (Chen et al., 2021) | 35.2 | 22.4 | 21.4 | 43.5 | - | - | - | - | - | - | 41.0[†] |
| 3DJCG (Cai et al., 2022) | 47.7 | 31.5 | 24.3 | 51.8 | - | - | - | - | - | - | - |
| Vote2Cap-DETR (Chen et al., 2023) | 61.8 | 34.5 | 26.2 | 54.4 | - | - | - | - | - | - | - |
| ScanRefer+MCAN (Chen et al., 2020) | - | - | - | - | - | 55.4 | 7.9 | 11.5 | 30.0 | 18.6 | - |
| ClipBERT (Lei et al., 2021) | - | - | - | - | - | - | - | - | - | - | 43.3 |
| ScanQA (Azuma et al., 2022) | - | - | - | - | - | 64.9 | 10.1 | 13.1 | 33.3 | 21.1 | 47.2 |
| *Task-specific fine-tuned* | | | | | | | | | | | |
| 3D-VisTA (Zhu et al., 2023c) | 66.9 | 34.0 | 27.1 | 54.3 | 53.8 | 69.6 | 10.4 | 13.9 | 35.7 | 22.4 | 48.5 |
| 3D-LLM (FlanT5) (Hong et al., 2023) | - | - | - | - | - | 69.4 | **12.0** | 14.5 | 35.7 | 20.5 | - |
| LEO | **68.4** | **36.9** | **27.7** | **57.8** | **54.7** | **80.0** | 11.5 | **16.2** | **39.3** | **36.6** | **53.7** |

**CLIPort Manipulation**

| | separating-piles | | packing-google-objects-seq | | put-blocks-in-bowls | |
|---|---|---|---|---|---|---|
| | seen | unseen | seen | unseen | seen | unseen |
| Transporter | 48.4 | 52.3 | 46.3 | 37.3 | 64.7 | 18.7 |
| CLIP-only | 90.2 | 71.0 | 95.8 | 57.8 | 97.7 | 44.5 |
| RN50-BERT | 46.5 | 44.9 | 94.0 | 56.1 | 91.8 | 23.8 |
| CLIPort (single) | 98.0 | **75.2** | **96.2** | 71.9 | **100** | 25.0 |
| CLIPort (multi) | 89.0 | 62.8 | 84.4 | 70.3 | **100** | **45.8** |
| LEO | **98.8** | **75.2** | 76.6 | **79.8** | 86.2 | 35.2 |

**ObjNav Navigation**

| | MP3D-val | | HM3D-val | |
|---|---|---|---|---|
| | S(↑) | L(↑) | S(↑) | L(↑) |
| H.w. (shortest) | 4.4 | 2.2 | - | - |
| H.w. (70k demo) | **35.4** | 10.2 | - | - |
| VC-1 (ViT-B) | - | - | **57.1** | **31.4** |
| LEO | 23.1 | **15.2** | 23.1[†] | 19.1[†] |

# Related research



**Mobile Manipulation**

Human: Bring me the rice chips from the drawer. Robot: 1. Go to the drawers, 2. Open top drawer. I see **<img>**. 3. Pick the green rice chip bag from the drawer and place it on the counter.

**Visual Q&A, Captioning …**

Given **<img>**. Q: What's in the image? Answer in emojis.
A: 🍎🥝🍌🍓🍎🍒

**PaLM-E: An Embodied Multimodal Language Model**

Given **<emb>** ... **<img>** Q: How to grasp blue block? A: First, grasp yellow block

**Large Language Model (PaLM)**

**Control**    A: First, grasp yellow block and ...

Describe the following **<img>**:
A dog jumping over a hurdle at a dog show.

**Language Only Tasks**

**Task and Motion Planning**

Given **<emb>** Q: How to grasp blue block? A: First grasp yellow block and place it on the table, then grasp the blue block.

**Tabletop Manipulation**

Given **<img>** Task: Sort colors into corners. Step 1. Push the green star to the bottom left. Step 2. Push the green circle to the green star.

PaLM-E: a comprehensive generalist exceling in multi-modal reasoning and planning

**Internet-Scale VQA + Robot Action Data**

Q: What is happening in the image?
A: 311 423 170 55 244
A grey donkey walks down the street.

Q: Que puis-je faire avec ces objets?
A: 3455 1144 189 25673
Faire cuire un gâteau.

Q: What should the robot do to **<task>**?
A: 132 114 128 5 25 156
ΔTranslation = [0.1, -0.2, 0]
ΔRotation = [10°, 25°, -7°]

**Vision-Language-Action Models for Robot Control**

Q: What should the robot do to **<task>**? A: ...
**RT-2**    **Large Language Model**

ViT

A: 132 114 128 5 25 156    De-Tokenize    ΔT = [0.1, -0.2, 0]
ΔR = [10°, 25°, -7°]
**Robot Action**

**Closed-Loop Robot Control**

Put the strawberry into the correct bowl
Pick the nearly falling bag
Pick object that is different

Co-Fine-Tune    Deploy

RT-2: bridging the gap between vision, language and action

LEO: An Embodied Generalist Agent in 3D World

CraftJarvis: Multi-task Embodied Agents in an Open World

CraftJarvis

stay tuned…

# CraftJarvis: Embodied Agents in an Open World

Open-World Multi-Task Control Through Goal-Aware Representation Learning and Adaptive Horizon Prediction,
CVPR 2023

Describe, Explain, Plan and Select: Interactive Planning with Large Language Models Enables Open-World Multi-Task Agents,
Best paper award, ICML '23 TEACH Workshop
NeurIPS 2023

JARVIS-1: Open-world Multi-task Agents with Memory-Augmented Multimodal Language Models,
arXiv 2023

[craftjarvis-jarvis1.github.io](craftjarvis-jarvis1.github.io)

# Minecraft: embodied AI in an open world



Today's embodied AI
- Restrictive objectives
- Very few tasks
- Limited knowledge

Embodied AI in an open world
- Open-ended objectives
- Massively multitask
- Web-scale knowledge

# Challenges in open world environments



Planning success plummet in open worlds due to new challenges

Challenge #1: Complex Sub-task Dependency

Manipulation in Tabletop environment

Mine diamond in Minecraft environment

## Challenge #1: long-horizon planning

Open worlds have highly abundant object types with complex dependency and relation. As a result, ground-truth plans typically involve a long sequence of sub-goals with strict dependencies.

=> Planning Success Rate will drops significantly on long-horizon tasks.

# Challenges in open world environments



## Challenge #2: state-aware planning

When dealing with a task that can be completed by executing multiple possible sequences of sub-goals, the planner should be able to select the best route base on the current state of the agent.

=> the complex and diverse state distribution of open-world environments makes state-awareness hard to achieve.

# Challenges in open world environments



Task: Shave
Step 1: Grab razor
Step 2: Switch on razor
Step 3: Put razor on face
**Prompt**

Task: Apply lotion

**Frozen**
**Pre-Trained Causal LLM**

Step 1: Squeeze out a glob of lotion

Zero-Shot Planning via Causal LLM

Step 1: Squeeze out a glob of lotion

**Frozen**
**Pre-Trained Masked LLM**

Step 1: Pour lotion into right hand

Translation to Admissible Action

Task: Shave
Step 1: Grab razor
Step 2: Wash razor
Step 3: Switch on razor
**Prompt**

Task: Apply lotion
Step 1: Pour lotion into right hand
Step 2:

**Frozen**
**Pre-Trained Causal LLM**

Step-By-Step
Autoregressive Generation

Instruction Relevance with LLMs | Combined | Skill Affordances with Value Functions

Prompt Examples

How would you put an apple on the table?

*I would: 1. _____*

LLM

| | | |
|---|---|---|
| -6 | **Find an apple** | 0.6 |
| -30 | Find a coke | 0.6 |
| -30 | Find a sponge | 0.6 |
| -4 | Pick up the apple | 0.2 |
| -30 | Pick up the coke | 0.2 |
| ... | ... | ... |
| -5 | Place the apple | 0.1 |
| -30 | Place the coke | 0.1 |
| -10 | Go to the table | 0.8 |
| -20 | Go to the counter | 0.8 |

Value Functions

*I would: 1. **Find an apple**, 2. ___*

LLM    VF

LLM for planning in close worlds

Language Models as Zero-Shot Planners: Extracting Actionable Knowledge for Embodied Agents
Do As I Can, Not As I Say: Grounding Language in Robotic Affordance

# CraftJarvis: an embodied agents in Minecraft

```python
def craft_wooden_axe(initial_inventory={}):
    # step 1: mine 3 logs
    mine(obj = {"log":3}, tool = None)
    # step 2: craft 12 planks from 3 logs
    craft(obj = {"planks":12}, materials = {"log":3},
            tool = None)
    # step 3: craft 4 sticks from 2 planks
    craft(obj = {"stick":4}, materials = {"planks"
            :2}, tool = None)
    # step 4: craft 1 crafting_table from 4 planks
    craft(obj = {"crafting_table":1}, materials = {"
            planks":4}, tool = None)
    # step 5: craft 1 wooden_axe from 3 planks and 2
            sticks on crafting table
    craft(obj = {"wooden_axe":1}, {"planks": 3, "
            stick": 2}, "crafting_table")
    return "wooden_axe"
```

Lots of errors!

# Self-correction

**<task>**: Obtain a diamond 💎 in *Minecraft* step-by-step?; **<obs>**:

**original <plan>**:

**Self-check**: When simulating on the goal ⛏, I find ∕ are not enough (lack of 2 ∕). So I need craft more ∕ from 🪵. More 🪵 require more 🪵. So I need to mine more 🪵.

**refined <plan>**:

**multi-modal <feedback>**: I **failed** on 💎. My current state is: ⛏ is broken; I still have 🪵 ∕ 🪨 🪨 in the inventory. My position is …

**Self-explain**: Because mining 💎 needs ⛏, which I do not have in the inventory. Crafting ⛏ needs 🔩. So I need to smelt 🪵 into 🔩 first.

**new <plan> by re-planning**: 🔩 → ⛏ → 💎

# Embodied RAG (retrieval-augmented generation)

Query generation via reasoning

*User:* My current task is 🔲 , but I have never accomplished this task before. What related tasks might be helpful for me to complete 🔲 ?

*Assistant:* 

**initial query (text)**     query generation via reasoning     *reasoning stops*

Enchanting Table  →  ✗  →  Obsidian  →  ✗  →  Diamond Pickaxe  →  ✗  →  Diamond  →  ✓

Diamond  →  ✓  →  Leather  →  ✓  →  Iron Pickaxe  →  ✓

Book  →  ✗  →  Paper  →  ✓  →  Diamond axe  →  ✗

✗ not in memory
✓ in memory
→ reasoning

**Multi-Modal Memory**

```
[
  {
    "task":[entity]  🪓 wooden pickaxe
    "plan":[language]  3  12  1  4  1
    "state":[image]
  },
  {
    "task":[entity]  ⛏ stone pickaxe
    "plan":[language]  ...  1  1  3  1
    "state":[image]
  },
  ...
]
```

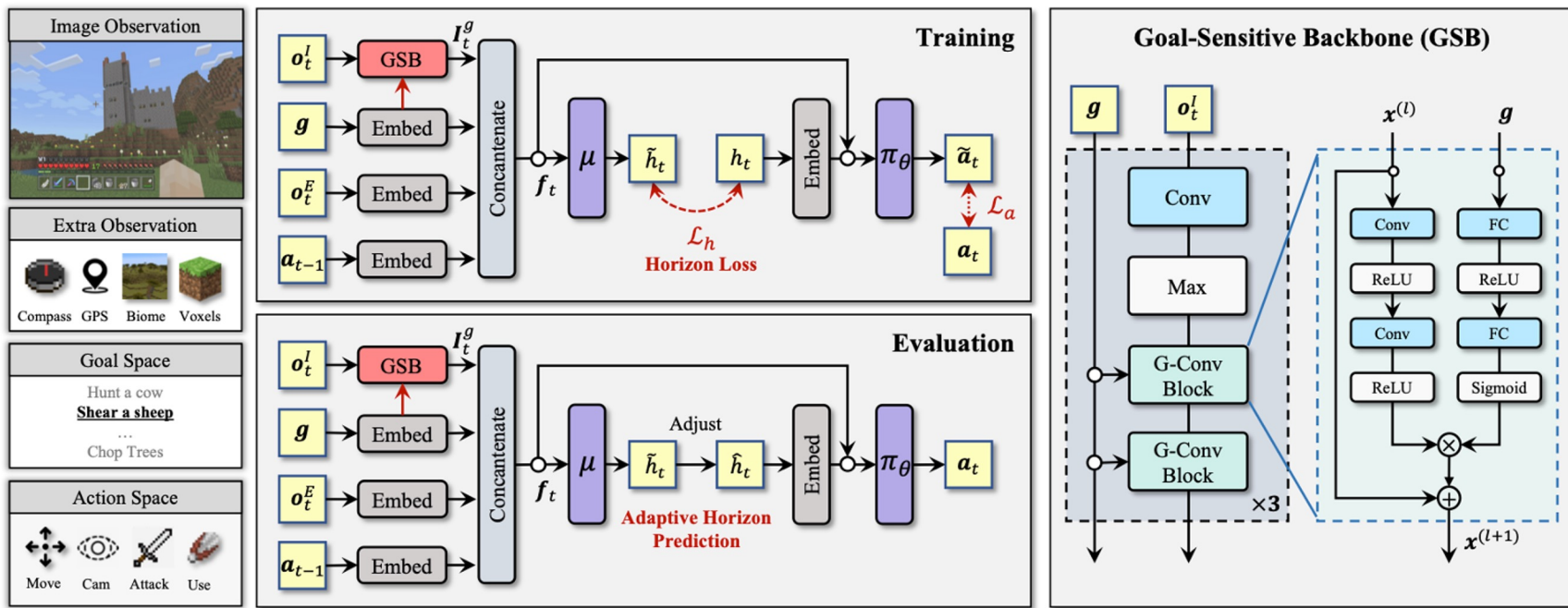query gen     retrieve

**final query (text):** 💎 Diamond  🟫 Leather  ⬜ Paper  🪓 Iron Pickaxe  **+ final query (obs):**     Query
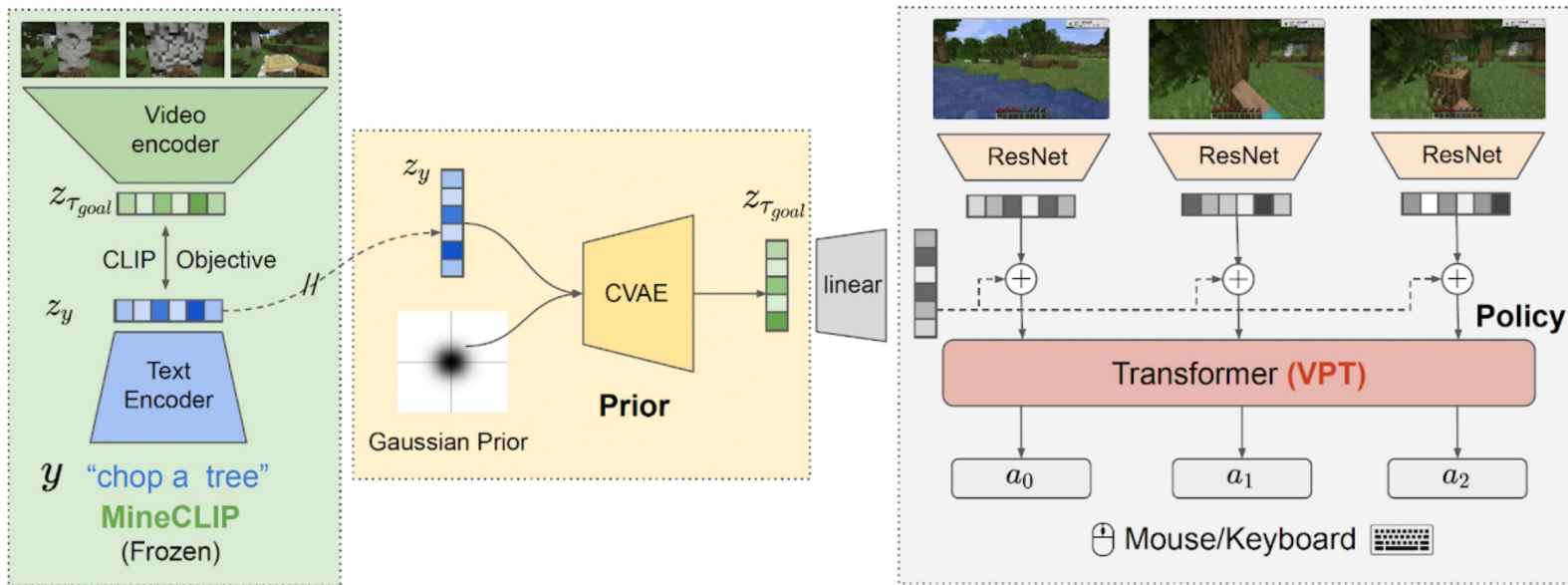
Machine Learning @ BIGAI

# Open world embodied control: goal-aware representation learning and horizon prediction

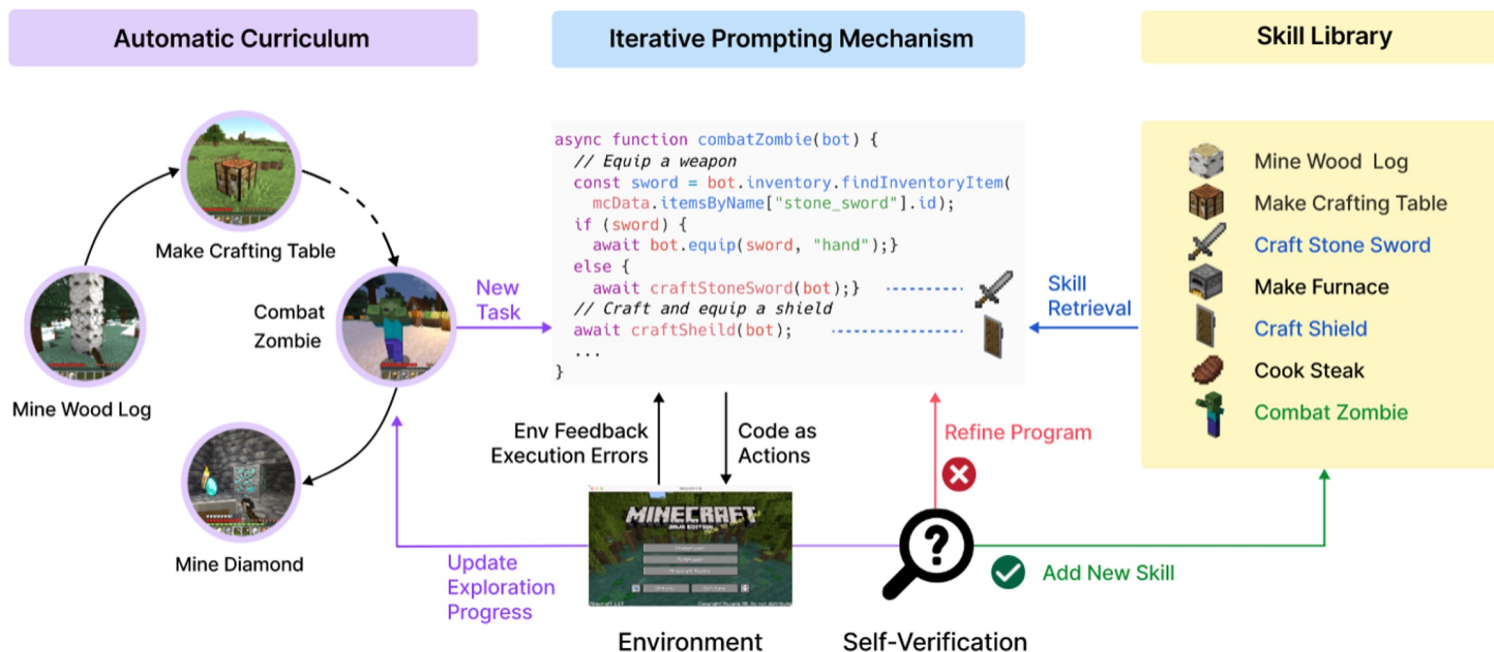# Open world embodied control: pretraining and alignment



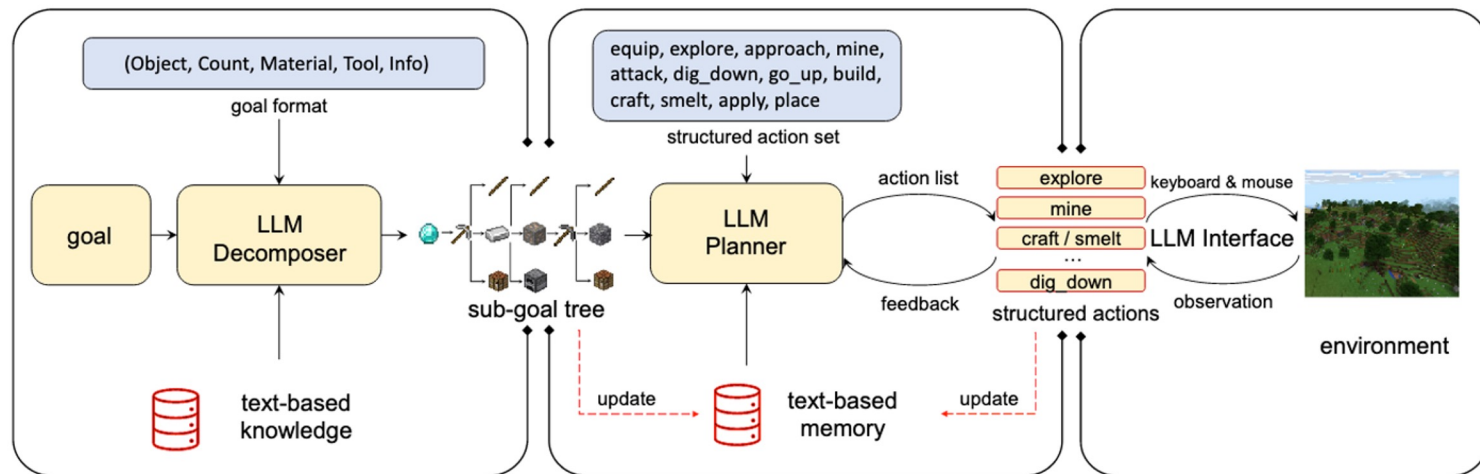STEVE-1: A Generative Model for Text-to-Behavior in Minecraft

# Some follow-up research projects built upon CraftJarvis
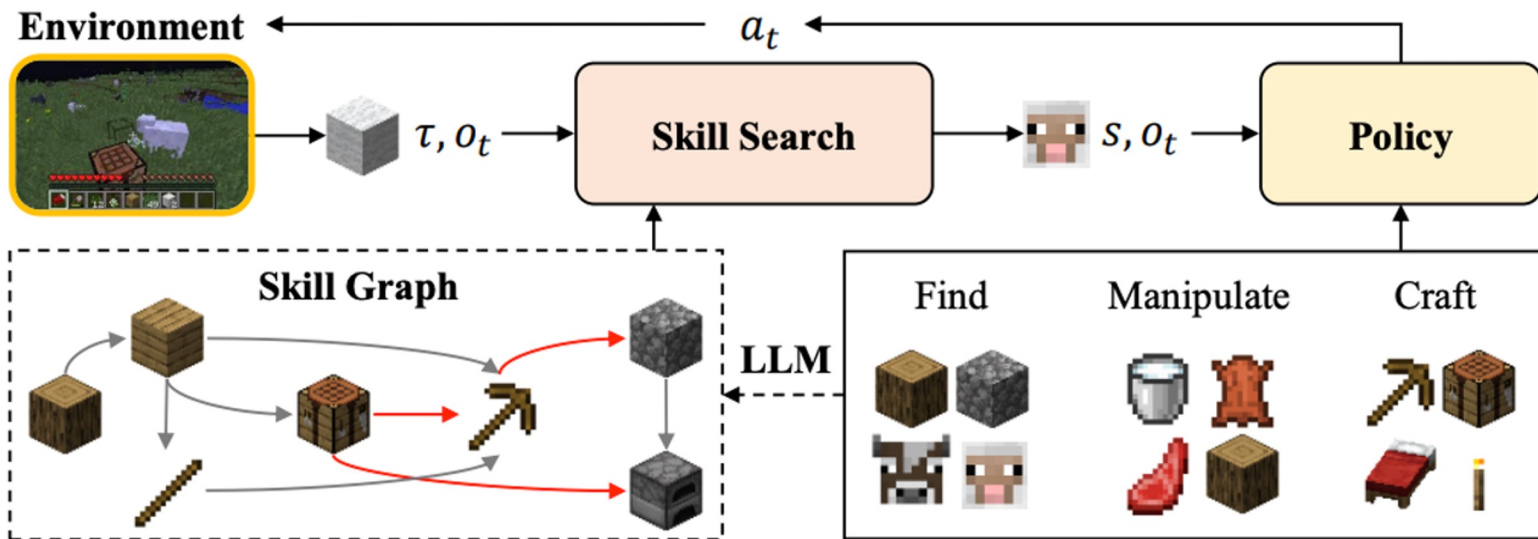
# Voyager: GPT-4 based language agent



Voyager: An Open-Ended Embodied Agent with Large Language Models

# GITM: language agent with structured knowledge library

Ghost in the Minecraft: Generally Capable Agents for Open-World Environments via Large Language Models with Text-based Knowledge and Memory
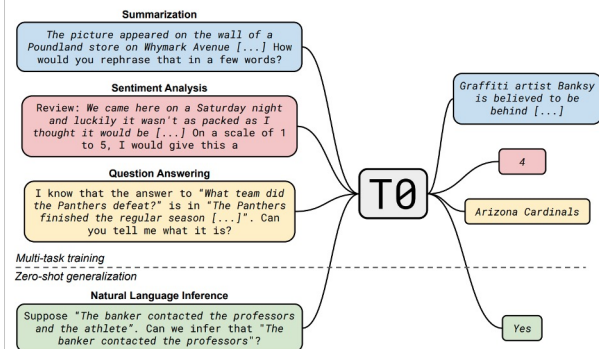
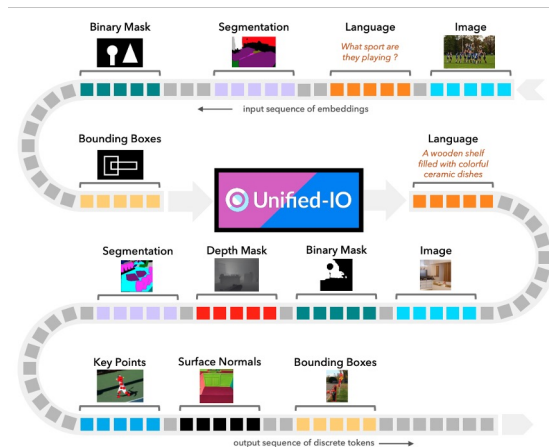# Plan4MC: language model + RL skills

Plan4MC: Skill Reinforcement Learning and Planning for Open-World Minecraft Tasks

# What's next?

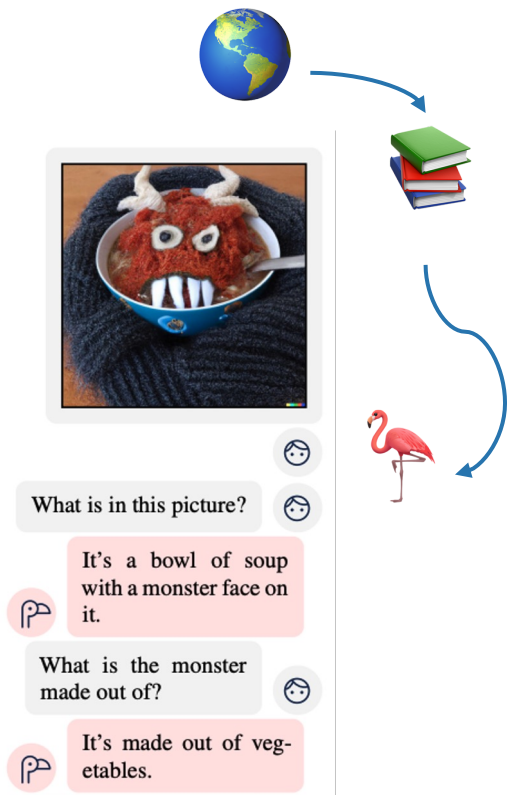# From unified models to unified agents



unified text models

unified multimodal models

unified agents

# From unified models to unified agents



What is in this picture?

It's a bowl of soup with a monster face on it.

What is the monster made out of?

It's made out of vegetables.

The following facts seems to be true:

1. Learning from massive web-scale data 📚
2. Large scale architecture O(10B) 🐘
3. Multi-tasking 🐙
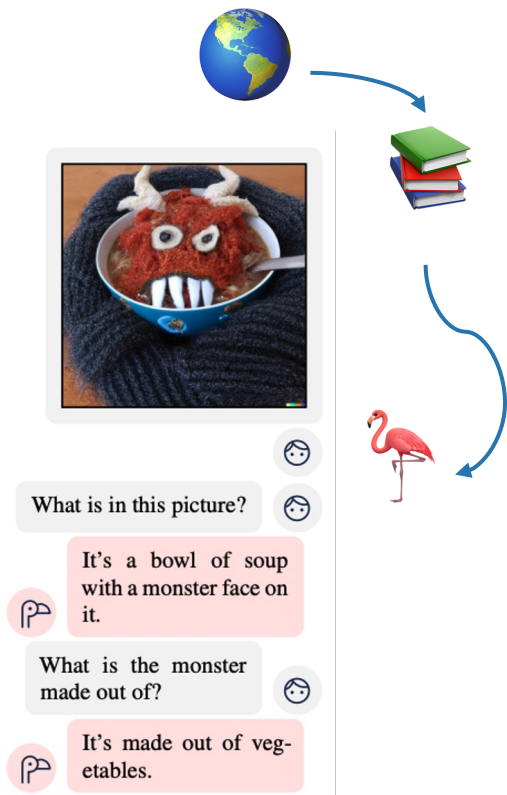4. (optional) Multimodal understanding 📖👀👂

# From unified models to unified agents



The following facts seems to be true:

1. Learning from massive web-scale data 📚
2. Large scale architecture O(10B) 🐘
3. Multi-tasking 🐙
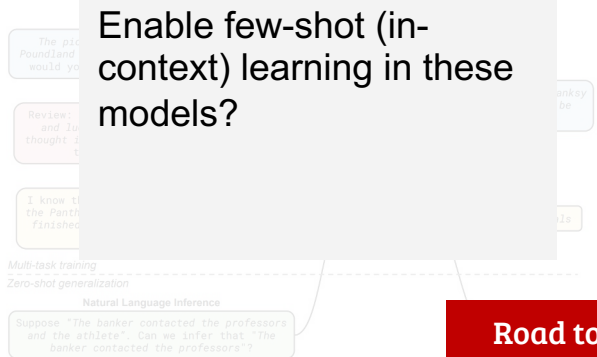4. (optional) Multimodal understanding 📖👀👂

### Definition

A language-piloted, large-scale agent that can fulfill arbitrary goals from multimodal input in embodied environments.

# From unified models to unified agents

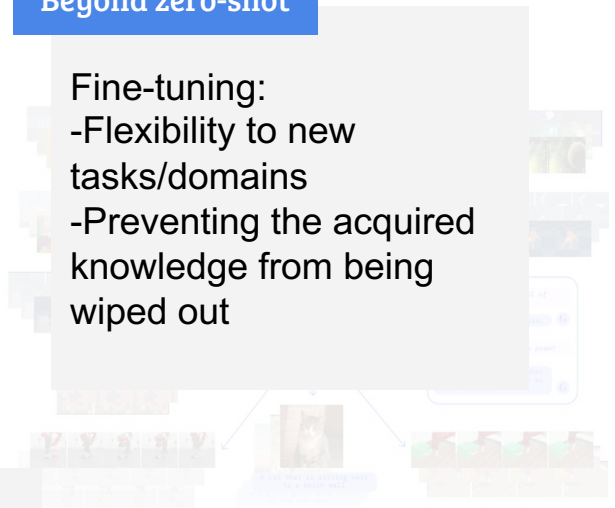Enable few-shot (in-context) learning in these models?

**More modalities**

Unified models for other modalities (3D, egocentric videos, proprioception, high-res structured input, etc)?
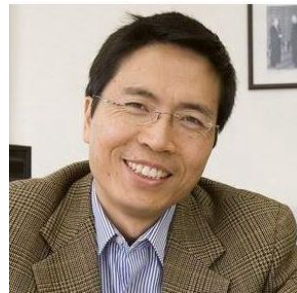
**Beyond zero-shot**

Fine-tuning:
-Flexibility to new tasks/domains
-Preventing the acquired knowledge from being wiped out

**Road to agents**

We need better learning algorithms for:
-episodic memory & situation awareness
-learning from interactions

unified text models
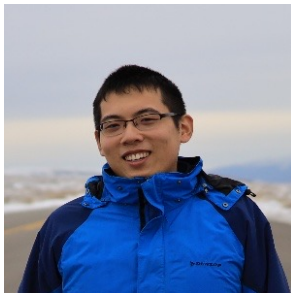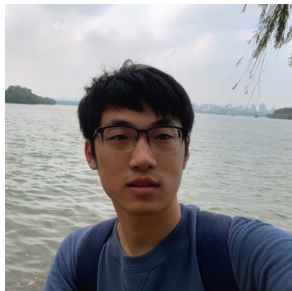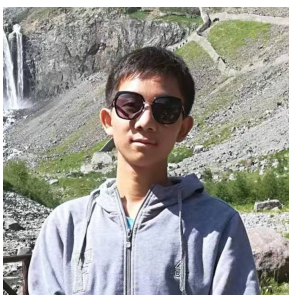
unified agents

*Multi-task training*
*Zero-shot generalization*

Natural Language Inference

Suppose "The banker contacted the professors and the athlete". Can we infer that "The banker contacted the professors"?

Surface Normals    Bounding Boxes

# Thank you