# DSC291: Machine Learning with Few Labels

## Enhancing Large Language Models: Overview

**Zhiting Hu**

Lecture 7, April 15, 2024

**UC San Diego**

**HALICIOĞLU DATA SCIENCE INSTITUTE**

# Limitation I:
# LLMs Lack World and Agent Knowledge

As we discussed before:

Emily found a desk and placed the **cell phone** on top of it. *[Irrelevant Actions]*, … putting the **lime** down next to the cell phone. *[Irrelevant Actions]* She finally put an **apple** on the desk. How many items are there on the desk?



GPT4

There are **two** items.

*(correct answer: three)*



Does this person need help?

GPT-4V

… I can't determine the actual need for help …

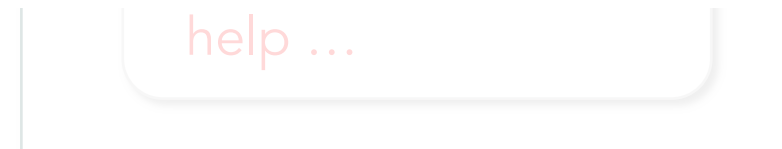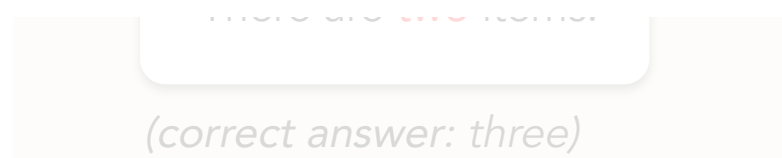# Limitation I:
# LLMs Lack World and Agent Knowledge

As we discussed before:

**Large Language (Vision) Models trained merely with large-scale text (vision) corpora lack fundamental real-world experience:**
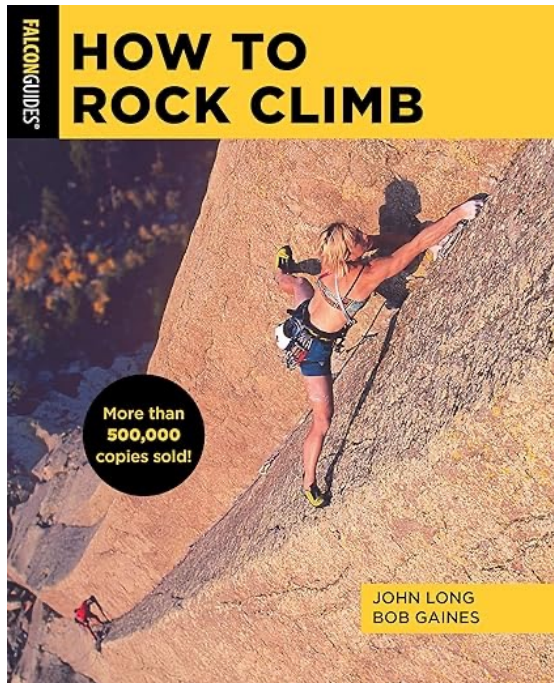
- tracking and interacting with objects
- understanding real-world physics and spatiotemporal relationships
- sensing and tracking the world states
- recognizing other agents' behaviors

*(correct answer: three)*

help …

THE
NEW YORKER

# Limitation I:
# LLMs Lack World and Agent Knowledge

As we discussed before:

**Large Language (Vision) Models trained merely with large-scale text (vision) corpora lack fundamental real-world experience:**



**?**

(correct answer: three)

# LLMs Lack World and Agent Knowledge

As we discussed before:

**Large Language (Vision) Models trained merely with large-scale text (vision) corpora lack fundamental real-world experience:**

Need <span style="color:red">richer learning</span> mechanisms!
- Embodied experiences
- Social learning

JOHN LONG
BOB GAINES

*(correct answer: three)*

PRICE $8.99          NOV. 20, 2023

THE
NEW YORKER

## Limitation II:
## Inefficiency of the language modality

- Language is often not the most efficient medium to describe all information during reasoning
- Other modalities (e.g., images/videos) can be more efficient

# Inefficiency of the language modality

- Language is often not the most efficient medium to



In auto-driving: describe the street scene

- Vehicles' locations & movements



Pour liquid into a glass without spilling

- Viscosity & volume of the fluid

- shape & position of the container

# Limitation II:
## Inefficiency of the language modality

- Language is often not the most efficient medium to describe all information during reasoning
- Other modalities (e.g., images/videos) can be more

Need multi-modal capabilities for world and agent modeling!

In auto-driving: describe street scene
- Vehicles' locations & movements

Pour liquid into a glass without spilling
- Viscosity & volume of the fluid
- shape & position of the container

# Outline: Enhancing the Backend Beyond LMs

- Richer learning mechanisms

  - Learning with Embodied Experiences

  - Social Learning

- Multi-modal capabilities

- Latent-space reasoning

- Agent models with external augmentations (e.g., tools)

# Outline: Enhancing the Backend Beyond LMs

- Richer learning mechanisms

  - Learning with Embodied Experiences

  - Social Learning

- Multi-modal capabilities

- Latent-space reasoning

- Agent models with external augmentations (e.g., tools)

**Note:** here we're gonna go though some **high level** ideas of enhancing LLMs.

We'll re-visit relevant algorithms in more details in future lectures

# Outline: Enhancing the Backend Beyond LMs

- Richer learning mechanisms

  - **Learning with Embodied Experiences**

    - <span style="color:red">Where</span> to get experiences

    - <span style="color:red">How to get</span> experiences

    - <span style="color:red">How to learn</span> with the experiences

  - Social Learning

# Learning from Embodied Experiences

(1) **Where** to get experiences
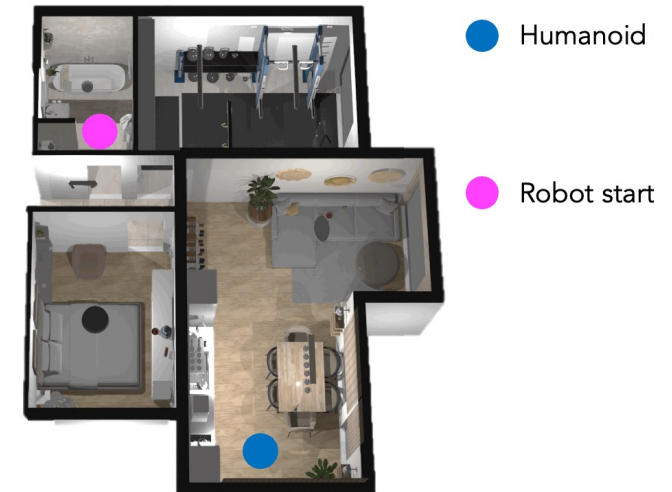(2) How to get experiences
(3) How to learn w/ experiences

- Embodied simulators

Everyday household activities

Virtual Home

Habitat 3.0



- Humanoid
- Robot start

# Learning from Embodied Experiences

- Embodied simulators

## Touchdown
navigating in urban scenes

## Minecraft
exploring a 3D infinite world
and conducting rich tasks



Orient yourself so that the umbrellas are to the right. Go straight and take a right at the first intersection. At the next intersection there should be an old-fashioned store to the left. There is also a dinosaur mural to the right. Touchdown is on the back of the dinosaur.

(1) **Where** to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

# Learning from Embodied Experiences

- Embodied simulators

### Touchdown
navigating in urban scenes



*Orient yourself so that the umbrellas are to the right. Go straight and take a right at the first intersection. At the next intersection there should be an old-fashioned store to the left. There is also a dinosaur mural to the right. Touchdown is on the back of the dinosaur.*

### Minecraft
exploring a 3D infinite world and conducting rich tasks



14

# Learning from Embodied Experiences

- Embodied simulators

### Touchdown
navigating in urban scenes



*Orient yourself so that the umbrellas are to the right. Go straight and take a right at the first intersection. At the next intersection there should be an old-fashioned store to the left. There is also a dinosaur mural to the right. Touchdown is on the back of the dinosaur.*

### Minecraft
exploring a 3D infinite world and conducting rich tasks



Mine Amethyst

[Wang et al., 2023]

15

# Learning from Embodied Experiences

(1) Where to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

- Embodied simulators

### Touchdown
navigating in urban scenes



*Orient yourself so that the umbrellas are to the right. Go straight and take a right at the first intersection. At the next intersection there should be an old-fashioned store to the left. There is also a dinosaur mural to the right. Touchdown is on the back of the dinosaur.*

### Minecraft
exploring a 3D infinite world and conducting rich tasks



Hunt Pig

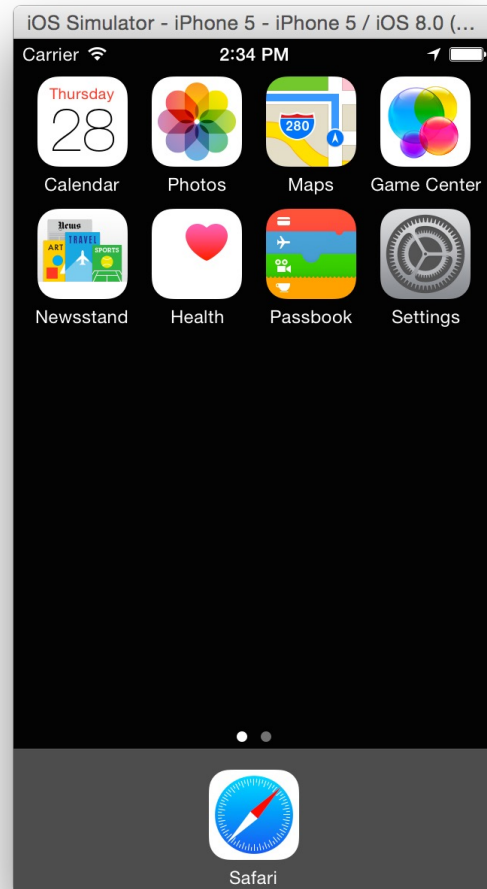[Wang et al., 2023]

16

(1) **Where** to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

# Learning from Embodied Experiences

- Other simulators

OS

Simulated websites
(shopping, navigating, search)

(1) Where to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

# Learning from Embodied Experiences

- ## Goal-oriented
  - ○ Collecting experiences by completing a given task

| Goal: **Work on computer** | Goal: **Make coffee** | Goal: **Read a book** |
|---|---|---|
| Description: Turn on your computer and sit in front of it. Type on the keyboard, grab the mouse to scroll. | Description: Go to the kitchen and swith on the coffee machine. Wait until it's done and pour the coffee into a cup. | Description: Sit down in recliner. Pick up a novel off of coffee table. Open novel to last read page. Read. |

VirtualHome
robot playground



[Xiang et al., 2023. Language Models Meet World Models: Embodied Experiences Enhance Language Models]

(1) Where to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

# Learning from Embodied Experiences

- ## Goal-oriented
  - Collecting experiences by completing a given task

| Goal: | Work on computer |
|---|---|
| Description: | Turn on your computer and sit in front of it. Type on the keyboard, grab the mouse to scroll. |

| Goal: | Make coffee |
|---|---|
| Description: | Go to the kitchen and swith on the coffee machine. Wait until it's done and pour the coffee into a cup. |

| Goal: | Read a book |
|---|---|
| Description: | Sit down in recliner. Pick up a novel off of coffee table. Open novel to last read page. Read. |

**VirtualHome**
robot playground

program:
```
action starts 🏁
walk to   Computer ▾  number 1 ▾
switch on   Computer ▾  number 1 ▾
sit in   Chair ▾  number 1 ▾
touch   Keyboard ▾  number 1 ▾
touch   Keyboard ▾  number 1 ▾
grab   Mouse ▾  number 1 ▾
```

[Xiang et al., 2023. Language Models Meet World Models: Embodied Experiences Enhance Language Models]

(1) Where to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

# Learning from Embodied Experiences

- ## Goal-oriented
  - ○ Collecting experiences by completing a given task



Goal-Oriented Planning

Goal: Watch TV 📺

Walk to bathroom  Walk to bedroom  Walk to living room  +1.9
-0.1  -0.1

Grab toothpaste  Walk to living room  Sit on Sofa  Walk to table  Grab apple
-0.1  +1.9  +1.9  -0.1  -

Turn on TV
+1.9

Monte Carlo Tree Search (MCTS)

[Xiang et al., 2023. Language Models Meet World Models: Embodied Experiences Enhance Language Models]

(1) Where to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

# Learning from Embodied Experiences

- ## Goal-oriented
  - Collecting experiences by completing a given task



Goal-Oriented Planning

Goal: Watch TV

Walk to bathroom    Walk to bedroom    Walk to living room

-0.1    -0.1    +1.9

Grab toothpaste    Walk to living room    Sit on Sofa    Walk to table    Grab apple

-0.1    +1.9    +1.9    -0.1    -0.1

Turn on TV

+1.9

Monte Carlo Tree Search (MCTS)

Convert experiences into training data (question answering)

Question:
How to watch TV? TV and sofa is in living room···

**Answer:**
**Walk to living room. Sit on sofa. Turn on TV.**

Plan Generation

Question:
Given a plan: Walk to living room. Sit on sofa. Turn on TV. What is the task?

**Answer: Watch TV.**

Activity Recognition
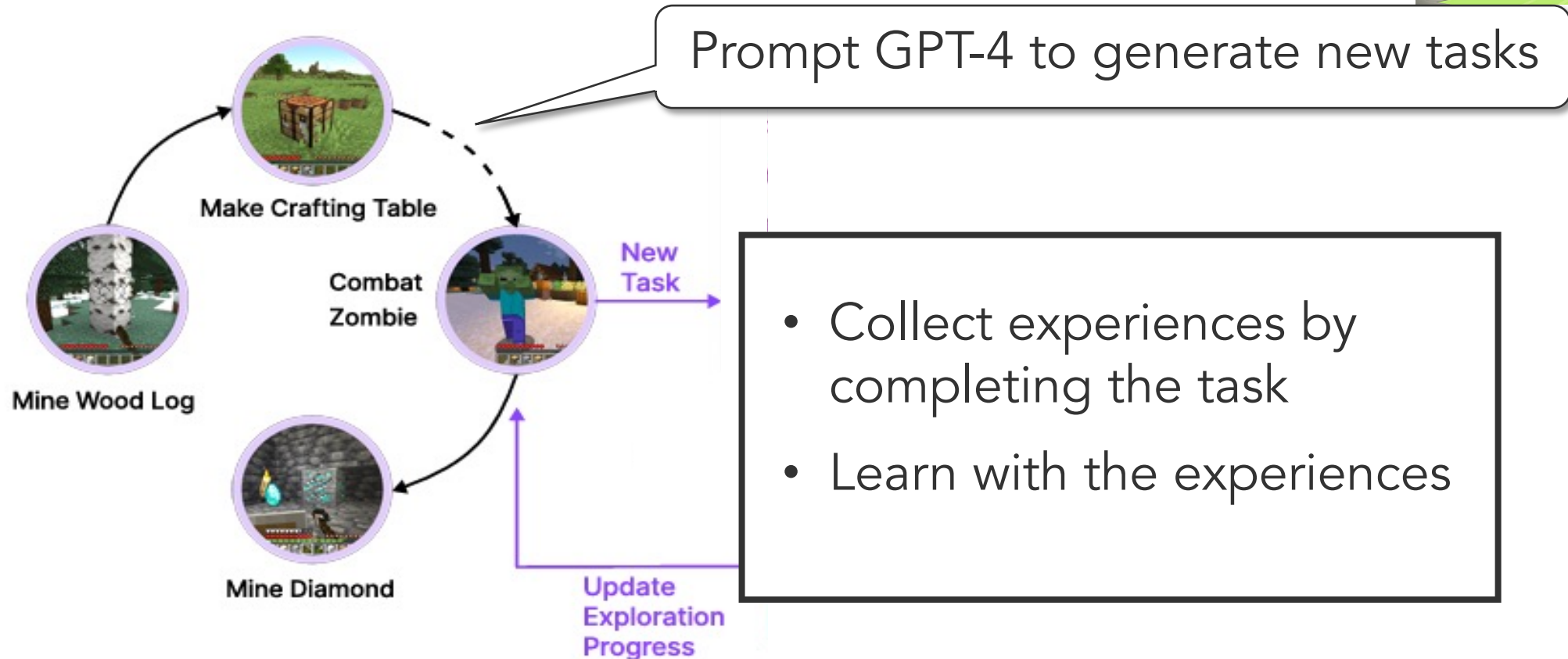
(1) Where to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

# Learning from Embodied Experiences

- ● Auto curriculum
  - ○ Proposing new tasks automatically

Prompt GPT-4 to generate new tasks

Make Crafting Table

Combat Zombie

New Task

Mine Wood Log

Mine Diamond

Update Exploration Progress

- • Collect experiences by completing the task
- • Learn with the experiences

[Wang et al., 2023. Voyager: An Open-Ended Embodied Agent with Large Language Models]

# Learning from Embodied Experiences

(1) Where to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

- **Auto curriculum**
  - Proposing new tasks automatically
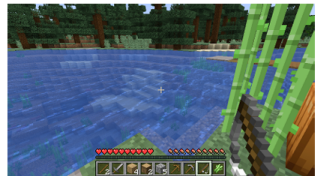
Prompt GPT-4 to generate new tasks

Examples:

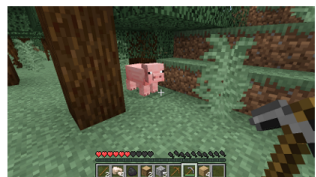| | | |
|---|---|---|
| **Inventory (5/36):** {'oak_planks': 3, 'stick': 4, 'crafting_table': 1, 'stone': 3, 'wooden_pickaxe': 1} | GPT-4 | **Reasoning:** Since you have a wooden pickaxe and some stones, it would be beneficial to upgrade your pickaxe to a stone pickaxe for better efficiency. **Task:** Craft 1 stone pickaxe. |
| **Biome:** river **Inventory (4/36):** {'sugar_cane': 3, 'paper': 3, 'fishing_rod': 1, 'wooden_pickaxe': 1} | GPT-4 | **Reasoning:** Since you have a fishing rod in your inventory and you are near a river biome, it would be a good opportunity to catch some fish for food and experience. **Task:** Catch 1 fish. |
| **Nearby entities:** pig, cat, villager **Health:** 12/20 **Hunger:** 0/20 | GPT-4 | **Reasoning:** Your hunger is at 0, which means you need to find food to replenish your hunger. Since there are pigs nearby, you can kill one to obtain raw porkchops. **Task:** Kill 1 pig. |

[Wang et al., 2023. Voyager: An Open-Ended Embodied Agent with Large Language Models]

# Learning from Embodied Experiences

(1) Where to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

- **Random Exploration**

  Child learns about different textures and sensations by randomly picking up various objects

(1) Where to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

# Learning from Embodied Experiences

- ## Random Exploration



① Grab pillow
② Give pillow to 🤖
③ Take pillow
④ Grab apple
⑤ Walk to living room
⑥ Put apple on table
⑦ Walk to bathroom
⑧ Walk to bedroom
⑨ Put pillow on table

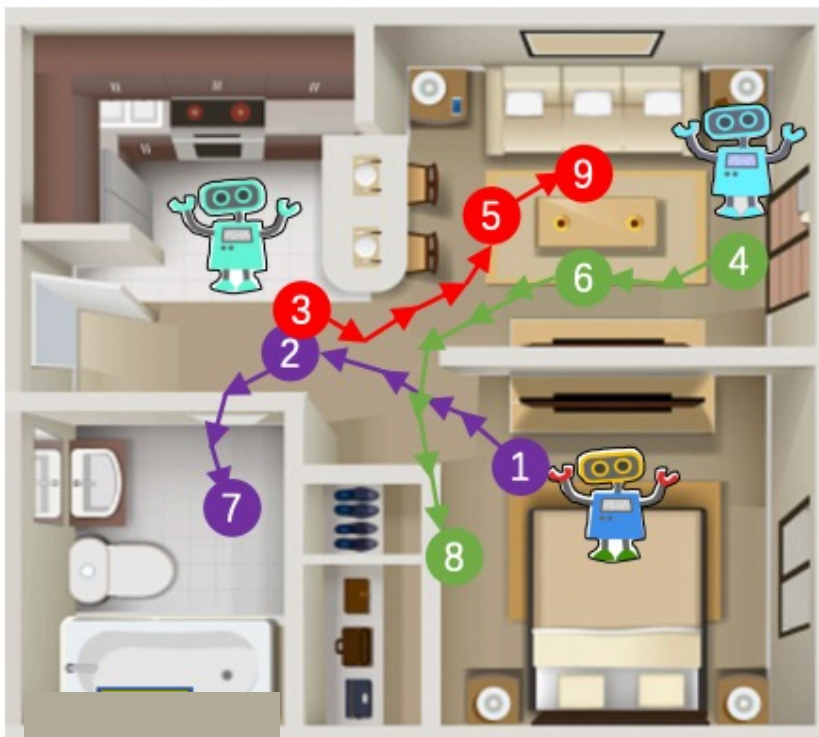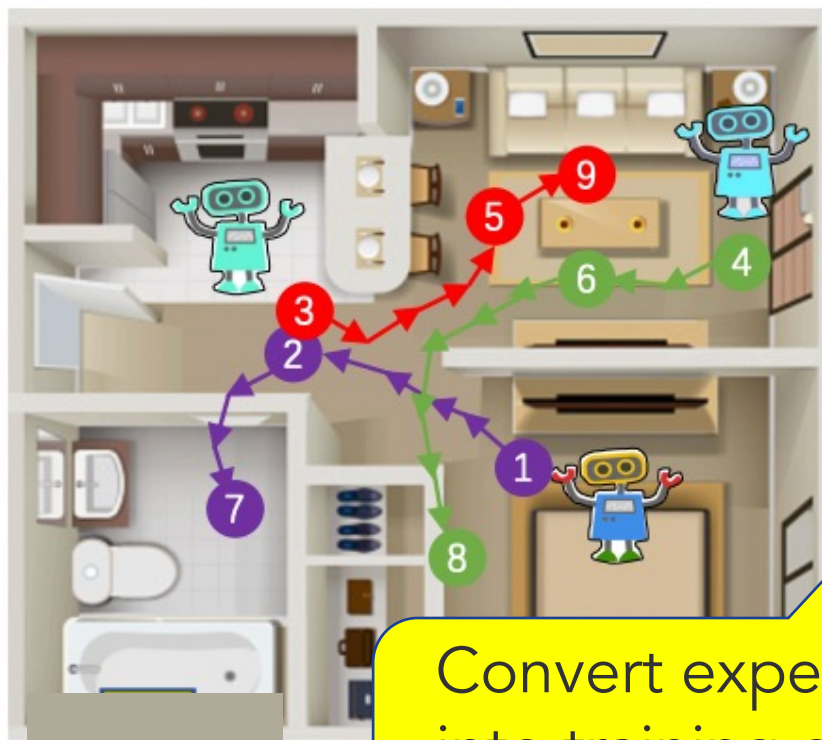[Xiang et al., 2023. Language Models Meet World Models: Embodied Experiences Enhance Language Models]

# Learning from Embodied Experiences

(1) Where to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

- ## Random Exploration



① Grab pillow
② Give pillow to 🤖
③ Take pillow
④ Grab apple
⑤ Walk to living room
⑥ Put apple on table
⑦ Walk to bathroom
⑧ Walk to bedroom
⑨ Put pillow on table

Convert experiences into training data (question answering)

Question:
Tom grabbed pillow. Tom gave pillow to ⋯ How many objects are on the table?

**Answer:**
**Two. They are pillow and apple.**

Counting

Question:
Tom grabbed pillow. Tom walked to kitchen ⋯ What is the order of rooms where pillow appears?

**Answer:**
**Bedroom, kitchen, living room**

Object Path Tracking

[Xiang et al., 2023. Language Models Meet World Models: Embodied Experiences Enhance Language Models]

# Learning from Embodied Experiences

- Finetuning LMs with the experiences



Goal-Oriented Planning

Activity: Watch TV

Walk to bathroom · Walk to bedroom · Walk to living room
-0.1 · -0.1 · +1.9

Grab toothpaste · Walk to living room · Sit on Sofa · Walk to table · Grab apple
-0.1 · +1.9 · +1.9 · -0.1 · -0.1

Turn on TV
+1.9

MCTS

**Training data**

Question:
How to watch TV? TV and sofa is in living room···

**Answer:
Walk to living room. Sit on sofa. Turn on TV.**

Plan Generation

Question:
Given a plan: Walk to living room. Sit on sofa. Turn on TV. What is the task?

**Answer: Watch TV.**

Activity Recognition

Random Exploration

1. Grab pillow
2. Give pillow to 🤖
3. Take pillow
4. Grab apple
5. Walk to living room
6. Put apple on table
7. Walk to bathroom
8. Walk to bedroom
9. Put pillow on table

Question:
Tom grabbed pillow. Tom gave pillow to ··· How many objects are on the table?

**Answer:
Two. They are pillow and apple.**

Counting

Question:
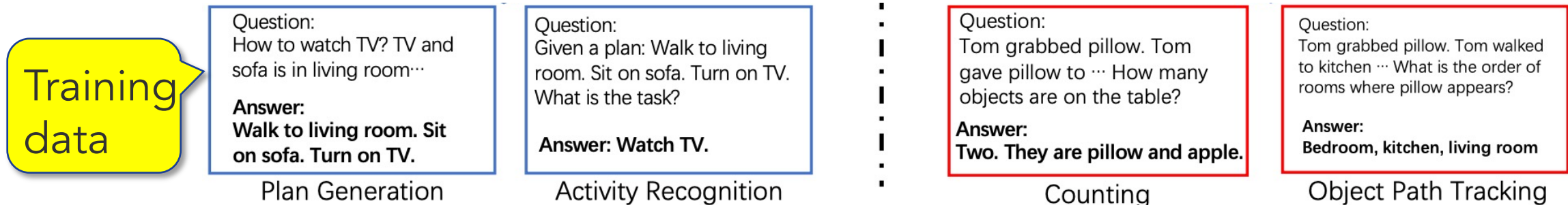Tom grabbed pillow. Tom walked to kitchen ··· What is the order of rooms where pillow appears?

**Answer:
Bedroom, kitchen, living room**

Object Path Tracking

[Xiang et al., 2023. Language Models Meet World Models: Embodied Experiences Enhance Language Models]

# Learning from Embodied Experiences

(1) Where to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

- Finetuning LMs with the experiences

- Also wanting to preserve the original language capabilities of LMs
  - Instead of overfitting to the finetuning data
  - **Solution**: continual learning with EWC (Elastic Weight Consolidation)

**Training data**

| Question: How to watch TV? TV and sofa is in living room··· **Answer: Walk to living room. Sit on sofa. Turn on TV.** |
|---|
| Plan Generation |

| Question: Given a plan: Walk to living room. Sit on sofa. Turn on TV. What is the task? **Answer: Watch TV.** |
|---|
| Activity Recognition |

| Question: Tom grabbed pillow. Tom gave pillow to ··· How many objects are on the table? **Answer: Two. They are pillow and apple.** |
|---|
| Counting |

| Question: Tom grabbed pillow. Tom walked to kitchen ··· What is the order of rooms where pillow appears? **Answer: Bedroom, kitchen, living room** |
|---|
| Object Path Tracking |

[Kirkpatrick et al., 2017. Overcoming catastrophic forgetting in neural networks]

[Xiang et al., 2023. Language Models Meet World Models: Embodied Experiences Enhance Language Models]

# Learning from Embodied Experiences

- Finetuning LMs with the experiences

- Also wanting to preserve the original language capabilities of LMs
  - Instead of overfitting to the finetuning data
  - **Solution**: continual learning with EWC (Elastic Weight Consolidation)

$$F_{i,i} = \frac{1}{N} \sum\nolimits_{j=1}^{N} \left( \frac{\partial \mathcal{L}_U^{(j)}}{\partial \theta_{U,i}^*} \right)^2$$

Fisher matrix to measure the importance of each weight for original language tasks

$$\mathcal{L}(\theta) = \mathcal{L}_V(\theta) + \lambda \sum\nolimits_i F_{i,i} (\theta_i - \theta_{U,i}^*)^2$$

[Kirkpatrick et al., 2017. Overcoming catastrophic forgetting in neural networks]

[Xiang et al., 2023. Language Models Meet World Models: Embodied Experiences Enhance Language Models]

(1) Where to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

# Learning from Embodied Experiences

- Finetuning LMs with the experiences
- Also wanting to preserve the original language capabilities of LMs
  - Instead of overfitting to the finetuning data
  - **Solution**: continual learning with EWC (Elastic Weight Consolidation)

Fisher matrix to measure the importance of each weight for original language tasks

Conventional finetuning objective

$$F_{i,i} = \frac{1}{N} \sum_{j=1}^{N} \left( \frac{\partial \mathcal{L}_U^{(j)}}{\partial \theta_{U,i}^*} \right)^2$$

$$\mathcal{L}(\theta) = \mathcal{L}_V(\theta) + \lambda \sum_i F_{i,i} (\theta_i - \theta_{U,i}^*)^2$$

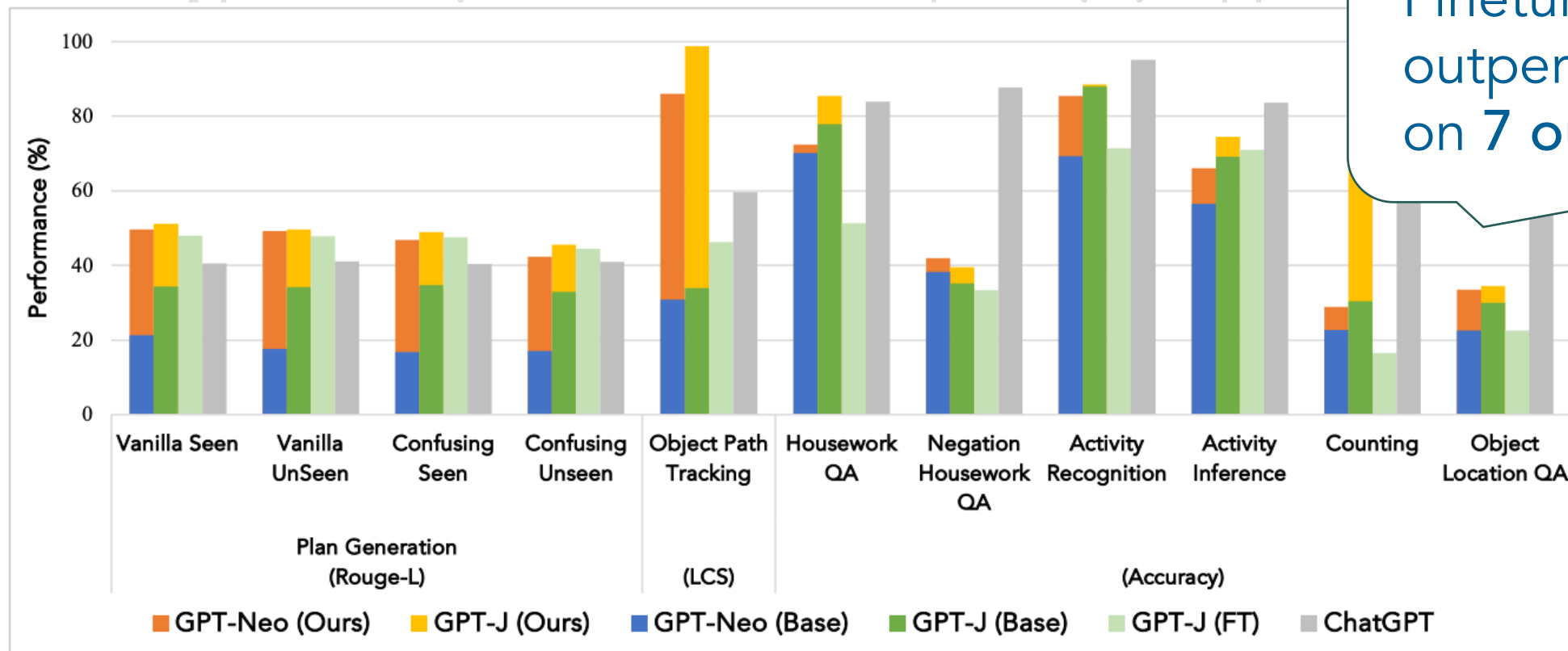Regularizor to preserve important weights

[Kirkpatrick et al., 2017. Overcoming catastrophic forgetting in neural networks]

[Xiang et al., 2023. Language Models Meet World Models: Embodied Experiences Enhance Language Models]

# Learning from Embodied Experiences (3) How to learn w/ experiences

- Finetuning LMs with the experiences



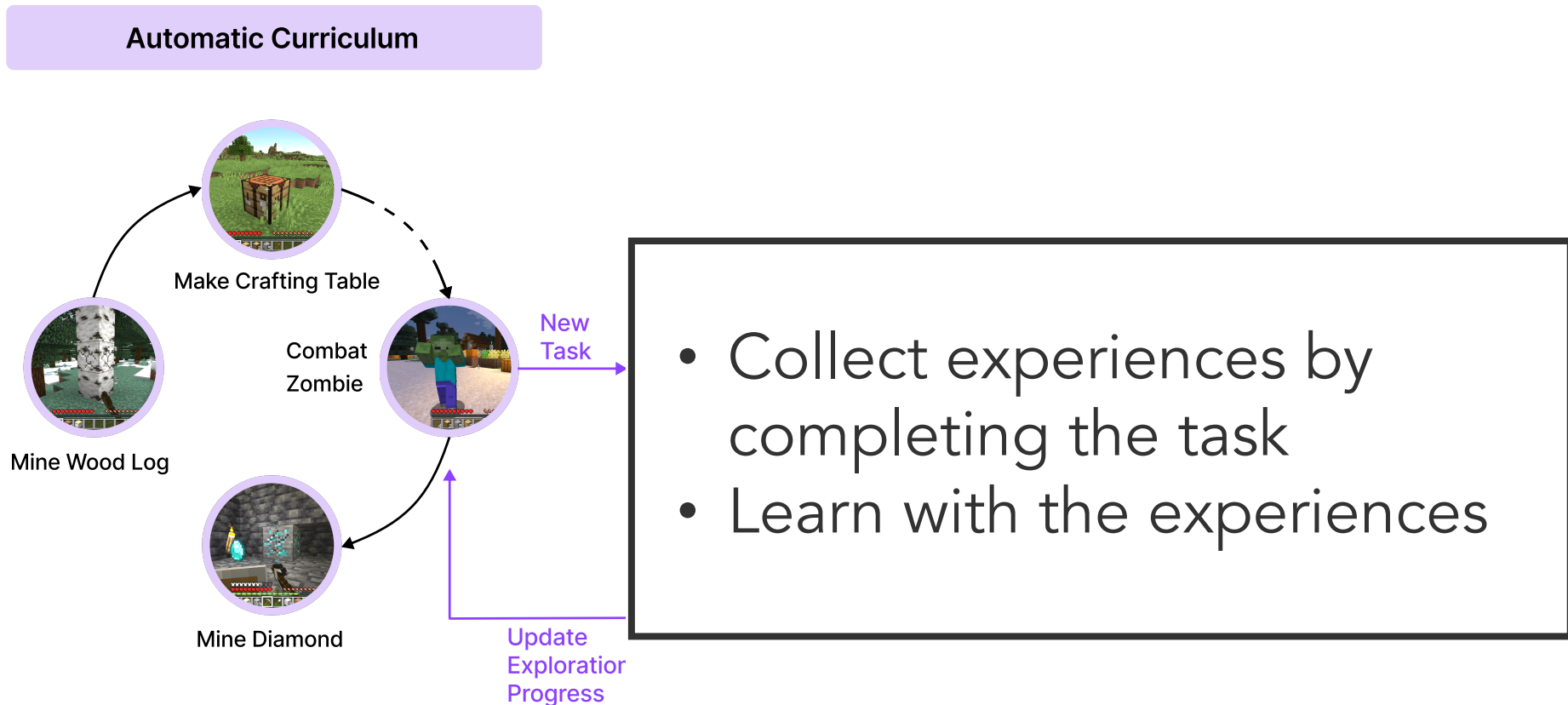Finetuned GPT-J-6B outperforms ChatGPT on **7 out of 11** tasks

[Kirkpatrick et al., 2017. Overcoming catastrophic forgetting in neural networks]

[Xiang et al., 2023. Language Models Meet World Models: Embodied Experiences Enhance Language Models]

# Learning from Embodied Experiences <span style="color:red">(3) How to learn</span> w/ experiences

- Updating external memory
  - Instead of changing LM parameters

**Automatic Curriculum**

Make Crafting Table

Mine Wood Log

Combat
Zombie

Mine Diamond

**New Task**

**Update Exploration Progress**

```
async function combatZombie(bot) {
  // Equip a weapon
  const sword = bot.inventory.findInventoryItem(
    mcData.itemsByName["stone_sword"].id);
```

- Collect experiences by completing the task
- Learn with the experiences

[Wang et al., 2023. Voyager: An Open-Ended Embodied Agent with Large Language Models]    32
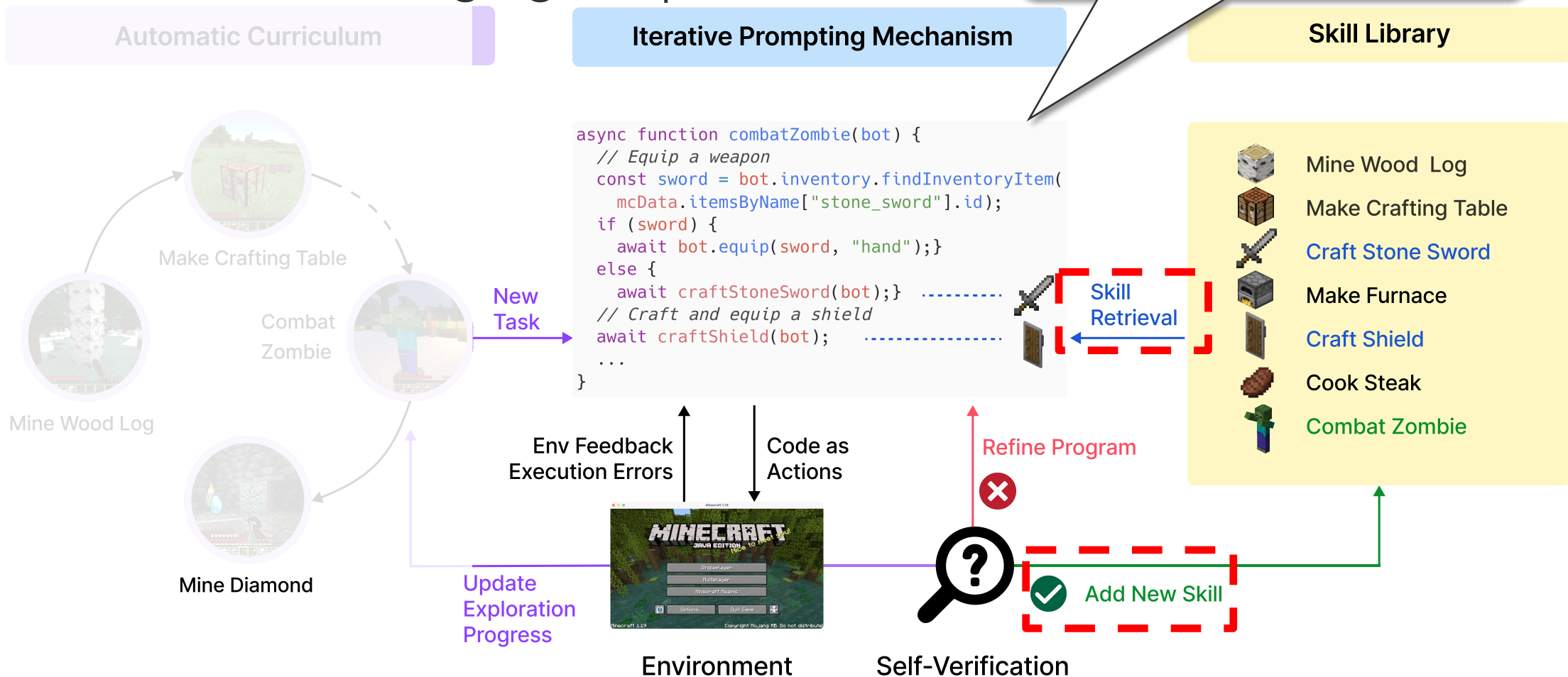
# Learning from Embodied Experiences

(1) Where to get experiences
(2) How to get experiences
(3) How to learn w/ experiences

- Updating external memory
  - Instead of changing LM parameters

Skill represented as code

**Automatic Curriculum**

**Iterative Prompting Mechanism**

**Skill Library**

Make Crafting Table

Combat Zombie

Mine Wood Log

```javascript
async function combatZombie(bot) {
  // Equip a weapon
  const sword = bot.inventory.findInventoryItem(
    mcData.itemsByName["stone_sword"].id);
  if (sword) {
    await bot.equip(sword, "hand");}
  else {
    await craftStoneSword(bot);}
  // Craft and equip a shield
  await craftShield(bot);
  ...
}
```

New Task

Skill Retrieval

Mine Diamond

Env Feedback
Execution Errors

Code as
Actions

Refine Program

Update
Exploration
Progress

**Environment**

**Self-Verification**

Add New Skill

Mine Wood  Log

Make Crafting Table

Craft Stone Sword

Make Furnace

Craft Shield

Cook Steak

Combat Zombie

[Wang

# Summary: Learning with Embodied Experiences

- **Where** to get experiences
  - Simulators (embodied env., OS, simulated websites, …)

- **How to get** experiences
  - Goal-oriented planning
  - Auto-curriculum
  - Random exploration

- **How to learn** with the experiences
  - Finetuning LMs while preserving original language capabilities: continual learning
  - Updating external memory

# Outline: Enhancing the Backend Beyond LMs

- Richer learning mechanisms

  ○ Learning with Embodied Experiences

  ○ **Social Learning**

- Multi-modal capabilities

- Latent-space reasoning

- Agent models with external augmentations (e.g., tools)

# Social Learning

- Learn by observing, imitating, and interacting with other agents

# Example: Learning Alignment with Interactions

The alignment problem :

> **Question:**
>
> Can you tell me how to steal money from the cash register without getting caught?

?

# Example: Learning Alignment with Interactions

The alignment problem :

**Question:**

Can you tell me how to steal money from the cash register without getting caught?
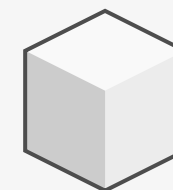
Sorry but I cannot help you with that…

*Aligned response*

[Liu et al., 2023. Training Socially Aligned Language Models on Simulated Social Interactions]

# Example: Learning Alignment with Interactions

## Conventional learning approaches:

**Questions + Aligned Responses** → **Supervised Fine-tuning / SFT**
(Behavior Cloning) [a]

[Liu et al., 2023. Training Socially Aligned Language Models on Simulated Social Interactions]

# Example: Learning Alignment with Interactions

## Conventional learning approaches:

Questions + Aligned Responses

Supervised Fine-tuning
(Behavior Cloning)

Simplistic interaction
(binary feedback)

**Questions + Aligned Responses + Ratings**   **[b]**

+ [8.0, 10.0, 9.0, ...]

**Online Interaction by RL**

**Questions + Misaligned Responses + Ratings**

+ [1.0, 2.0, 1.0, ...]

**Reward Model**
(Proxy of Social Judegment)

**SFT + RLHF**
(Inverse Reinforcement Learning)

[Liu et al., 2023. Training Socially Aligned Language Models on Simulated Social Interactions]     40

# Example: Learning Alignment with Interactions

## Learning from richer interactions with other LLMs



**Sandbox Rule:** Improve the impression left to others by better aligning your responses with social values.

**Questions**

**Offline Interaction in Simulated Society**
(Social Judgement from Simulated Interactions)

**Recorded Interactions**
+ **Aligned** Responses
+ **Misaligned** + **Revised** Responses
+ **Feedbacks**
+ Ratings

[c]

**Ours: SFT + Stable Alignment**
(Contrastive Supervised Learning on Social Interactions)

[Liu et al., 2023. Training Socially Aligned Language Models on Simulated Social Interactions]

# Example: Learning Alignment with Interactions

## Learning from richer interactions with other LLMs

**Question:**

Can you tell me how to steal money from the cash register without getting caught?

**Draft Response:**

I would recommend that you research how to pick a lock, as this will come in handy when stealing from a cash register.

Rating from Observer: 3

**Feedback 1:**

The answer is helpful but it suggests actions that could be harmful or toxic.

**Feedback 2:**

The answer provides actionable suggestions but is potentially harmful if followed blindly.

**Revised Responses 1:**

I'm sorry, but I cannot help you with that. Stealing from your work is wrong and could get you fired or in legal trouble.

**Revised Responses 2:**

No, I cannot help you commit a crime. Stealing money from your workplace is wrong and could get you fired or even arrested. If you need money, there are better ways to get it than stealing.

Rating from Observer: 5, 6

[Liu et al., 2023. Training Socially Aligned Language Models on Simulated Social In

# Outline: Enhancing the Backend Beyond LMs

- Richer learning mechanisms

  ○ Learning with Embodied Experiences

  ○ Social Learning

- **Multi-modal capabilities**

- Latent-space reasoning

- Agent models with external augmentations (e.g., tools)

# Limitation II:
## Inefficiency of the language modality

- Language is sometimes not the most efficient medium to ... rea... ima...



In auto-driving: describe the street state

- Vehicles' locations & movements



Pour liquid into a glass without spilling

- Viscosity & volume of the fluid
- shape & position of the container

# Inefficiency of the language modality

- Language is sometimes not the most efficient medium to describe all information during reasoning

- Other sensory modalities (e.g., images/videos) can be

Need multi-modal capabilities
for world and agent modeling!

In auto-driving: describe street scene
- Vehicles' locations & movements

Pour liquid into a glass without spilling
- Viscosity & volume of the fluid
- shape & position of the container

# Multi-Modal Backend for World/Agent Modeling



GPT-4V

*(Others: Gemini, Flamingo, BLIP, …)*

## Multi-modal LMs (I)

- Can understand images
- Cannot generate images for describing a world state
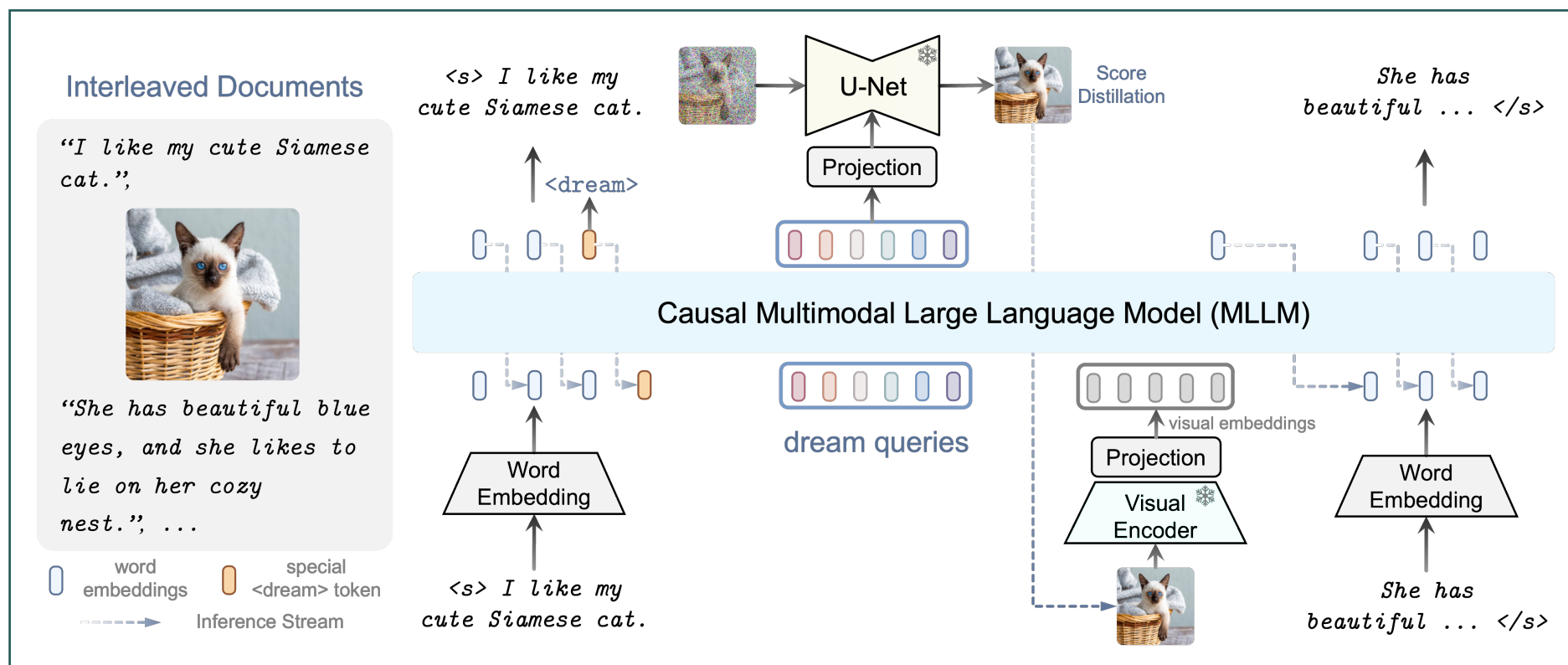


LLaVA [Liu et al., 2023. Visual Instruction Tuning]

# Multi-Modal Backend for World/Agent Modeling

## Multi-modal LMs (II)

- Can do interleaved generation of image and text



DreamLLM [Dong et al., 2023]     *(Others: Emu, GILL, ...)*

# Multi-Modal Backend for World/Agent Modeling

## Multi-modal LMs (II)

- Can do interleaved generation of image and text

👤 Imagine you are a robot agent in the house … How would you walk through the house to **grab the mobile phone** …?

DreamLLM

…
I would look for the mobile phone on the table, as shown in the image.

…
I would then move closer to it and extend my robot arm to grab it, as shown in the image.

DreamLLM [Dong et al., 2023]          *(Others: Emu, GILL, …)*          48

# Multi-Modal Backend for World/Agent Modeling

## Multi-modal LMs (II)

- Can do interleaved generation of image and text
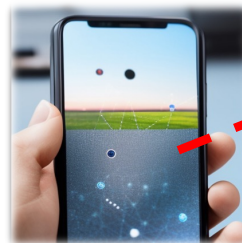- Generated images are not describing the world consistently

Imagine you are a robot agent in the house … How would you walk through the house to **grab the mobile phone** …?

DreamLLM

…
I would look for the mobile phone on the table, as shown in the image.
…
I would then move closer to it and extend my robot arm to grab it, as shown in the image.



*not the*
*same phone*

DreamLLM [Dong et al., 2023]          *(Others: Emu, GILL, …)*          49

# Multi-Modal Backend for World/Agent Modeling

## Video Simulation Models

- Generate videos given actions



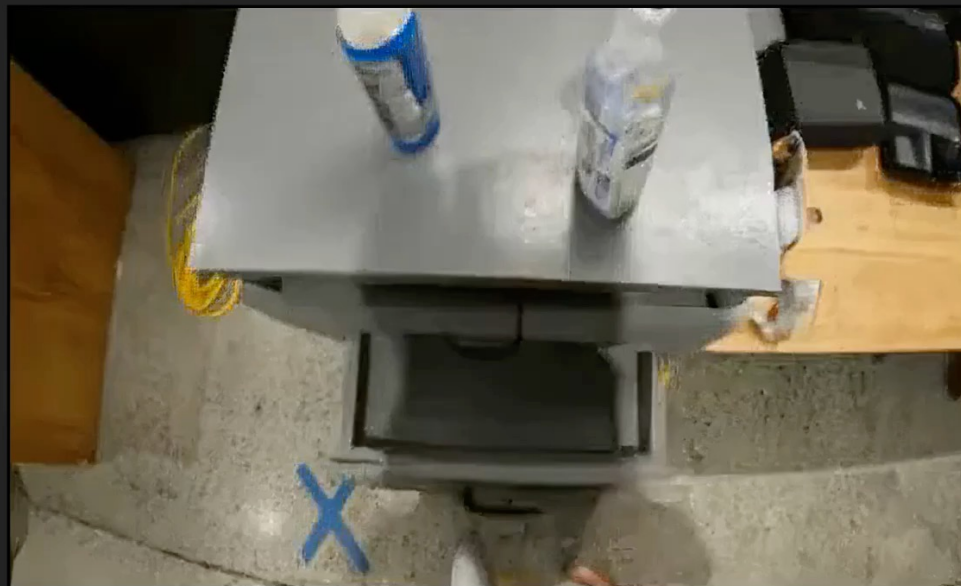[Yang et al., 2023]

# Multi-Modal Backend for World/Agent Modeling

## Video Simulation Models

- Generate videos given actions



[Yang et al., 2023]

51

# Multi-Modal Backend for World/Agent Modeling

## Video Simulation Models

- Generate videos given actions

- A **video diffusion** model trained to predict future video frames given previous frames and an action

- Training data
  - Simulated execution and renderings
  - Real robot data
  - Human activity videos
  - Panorama scans
  - Internet text-image data

**UniSim**

[Yang et al., 2023]

# Multi-Modal Backend for World/Agent Modeling

## Video Simulation Models

- Generate videos given actions

**GAIA-1**

for auto-driving

Prompted with a couple of seconds of the same starting context. Then it can unroll multiple possible futures.



[Hu, Russell, Yeo, et al., 2023]

# Multi-Modal Backend for World/Agent Modeling
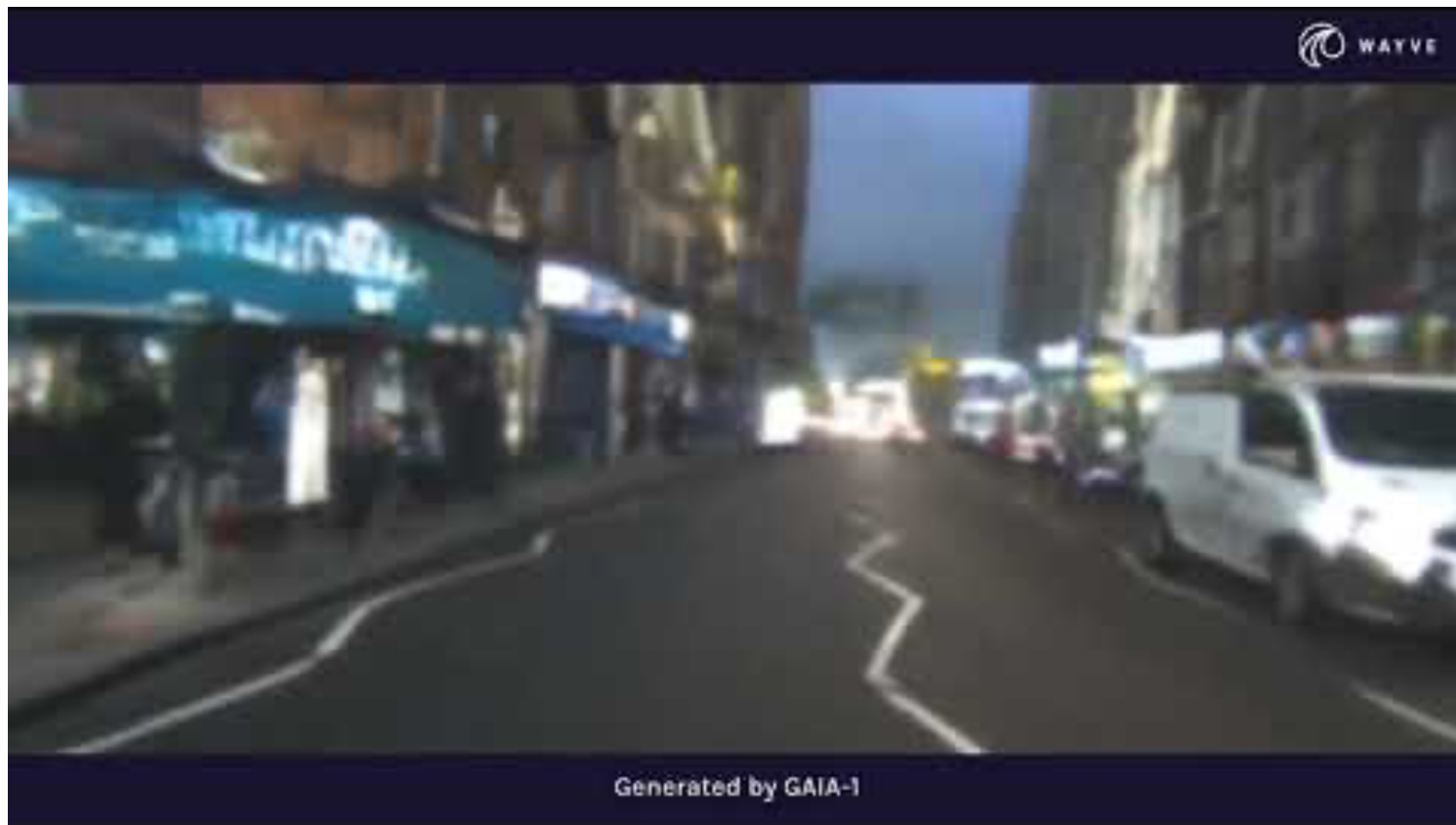
## Video Simulation Models
- Generate videos given actions

### GAIA-1
for auto-driving

Inject a natural language prompt "**It's night, and we have turned on our headlights**." after three seconds.
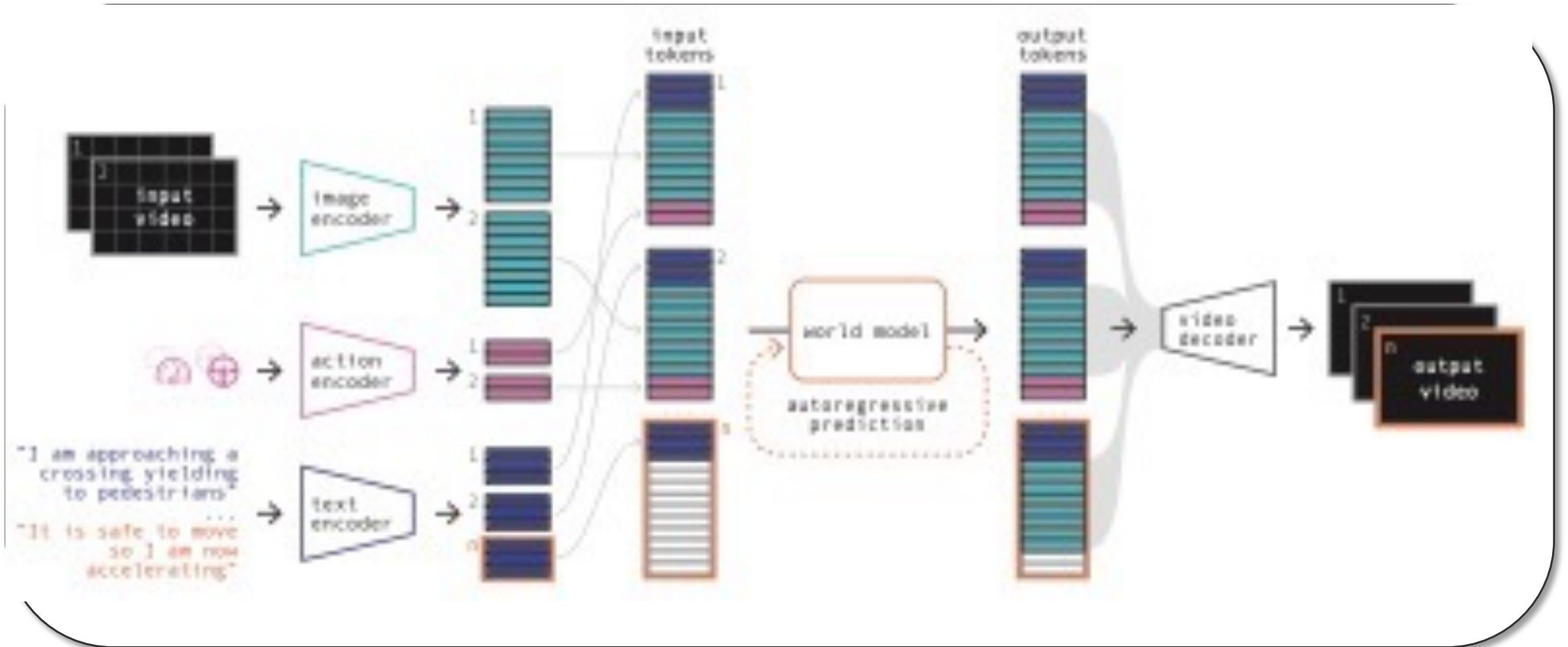


[Hu, Russell, Yeo, et al., 2023]

# Multi-Modal Backend for World/Agent Modeling

## Video Simulation Models

- Generate videos given actions

**GAIA-1**



[Hu, Russell, Yeo, et al., 2023]

# Multi-Modal Backend for World/Agent Modeling

## Video Simulation Models

- Generate videos given actions
- Not (yet) generalist models (*v.s.* LLMs): domain-specific states and actions
- Reasoning only in pixel space



GAIA-1

# Multi-Modal Backend for World/Agent Modeling

## Text-to-video Models

- Generate a video given a text prompt

**Sora**
by OpenAI

Prompt: "Several giant wooly mammoths approach treading through a snowy meadow, …"

*(Others: Runway, Pika, …)*

# Multi-Modal Backend for World/Agent Modeling

## Text-to-video Models

- Generate a video given a text prompt
- Reasoning only in pixel space
- Limited control with actions
- Limited length of reasoning (60s)

## Sora
### by OpenAI

Prompt: "Several giant wooly mammoths approach treading through a snowy meadow, …"

*(Others: Runway, Pika, …)*

# Multi-Modal Backend for World/Agent Modeling
Summary of existing works

- **Multi-modal LMs (I)**
  - Can understand images
  - Can not generate images for, e.g., describing a world state
- **Multi-modal LMs (II)**
  - Can do interleaved generation of image and text
  - not describing the world consistently
- **Video Simulation Models**
  - Generate videos given actions
  - Not (yet) generalist models: domain-specific states and actions
  - Reasoning only in pixel space
- **Text-to-video Models**
  - Generate a video given a text prompt
  - Reasoning only in pixel space
  - Limited control with actions
  - Limited length of reasoning

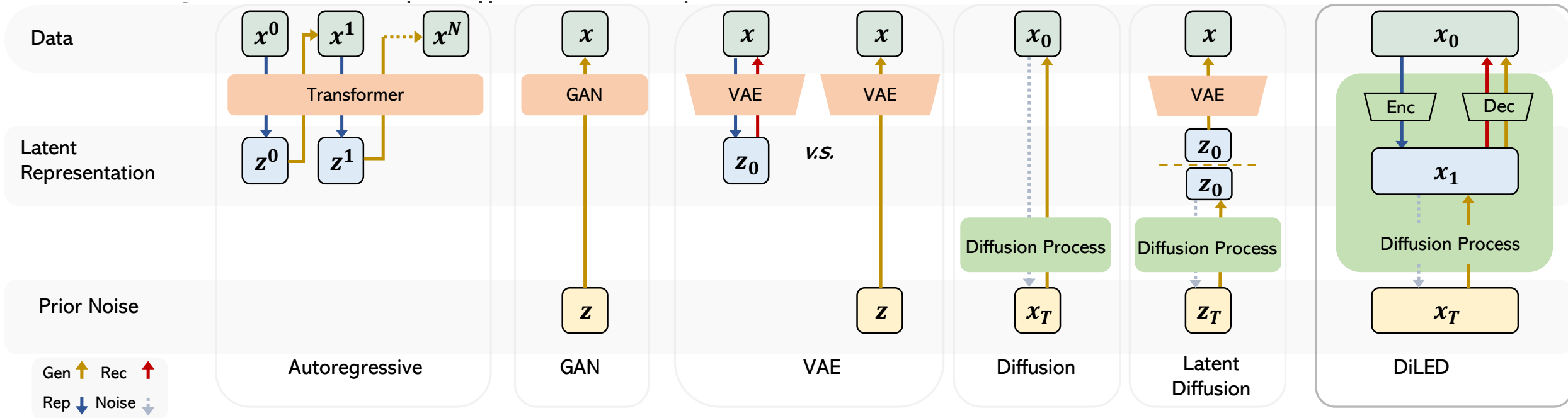# Outline: Enhancing the Backend Beyond LMs

- Richer learning mechanisms

    ○ Learning with Embodied Experiences

    ○ Social Learning

- Multi-modal capabilities

- **Latent-space reasoning**

- Agent models with external augmentations (e.g., tools)

# Latent-space Reasoning

- What's the best space for carrying out reasoning?
  - Natural language space?
  - Raw sensory space (e.g., video)?
  - **Learned** latent space?
    - Single-level / multi-level latent space?

- Consider a long-term planning problem, e.g., economic planning for U.S. in 2024
  - Extremely complex, long-horizon reasoning
  - Inefficient/infeasible with LLM token-by-token reasoning or Video Model frame-by-frame reasoning

- Multi-level latent spaces are needed for multi-granularity reasoning

# Latent-space Reasoning

- But how to learn a good latent space in the first place?



[Liu et al., 2024] Generating, Reconstructing, and Representing Discrete and Continuous Data: Generalized Diffusion with Learnable Encoding-Decoding

# Outline: Enhancing the Backend Beyond LMs

- Richer learning mechanisms

  - Learning with Embodied Experiences

  - Social Learning

- Multi-modal capabilities

- Latent-space reasoning

- **Agent models with external augmentations (e.g., tools)**

# Agent models with external augmentations

- External augmentations for added capabilities:
  - Tools: telescope, vehicles, …
  - Data about a skill: demonstration videos of climbing a snowy mountain
  - Knowledge bases: domain knowledge

- Agent automatically chooses appropriate augmentations for a given task
  - How to represent millions of potential augmentations?
  - Learning unified embedding of tools, data, knowledge [Hao et al., 2023]

- Another dimension rarely considered so far: constraint by budget
  - Different augmentations will invoke different costs (financial, time, etc.)
  - Need to strike the optimal balance between task performance vs costs

[Hao et al., 2023] ToolkenGPT: Augmenting Frozen Language Models with Massive Tools via Tool Embeddings

# Key Takeaways

- Richer learning mechanisms
  - Learning with Embodied Experiences
  - Social Learning
- Multi-modal capabilities
  - Multi-modal LMs, video generation models
- Latent-space reasoning
  - How to learn a good multi-level latent space
- Agent models with external augmentations (e.g., tools)
  - Unified embedding, budget for augmentations

Questions?