# Generative (Mental) World Explorer

**Jieneng Chen**
PhD candidate, Siebel Scholar Class 2025
Dept. Computer Science, Johns Hopkins University
04/04/2025

GenEx

# Mental imagery

- Simple test: **Close your eyes and visualize an apple.**

- How vivid?



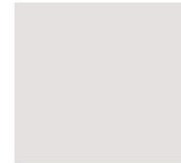1 - When you close your eyes you see an apple

2 - You see an apple, but it isn't sharp with highlights
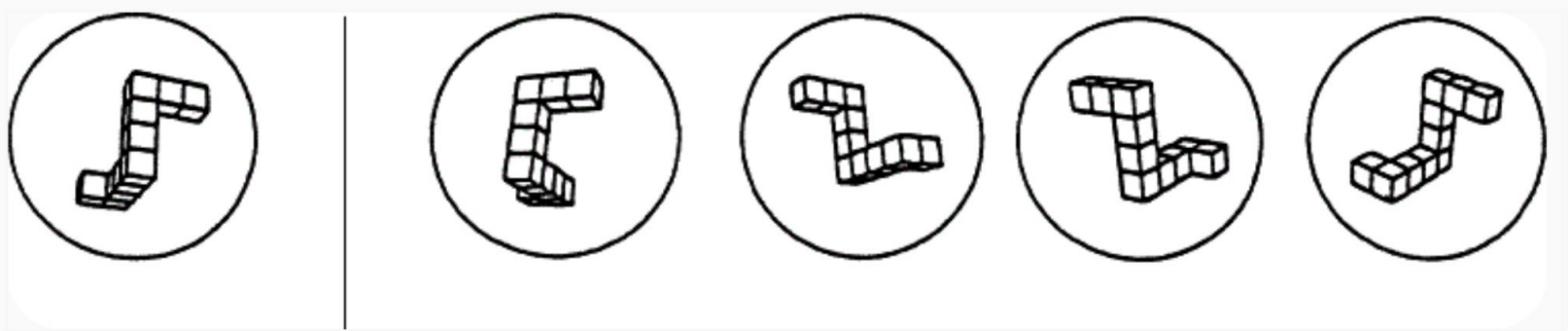
3 - You see an apple, but it looks like a solid color
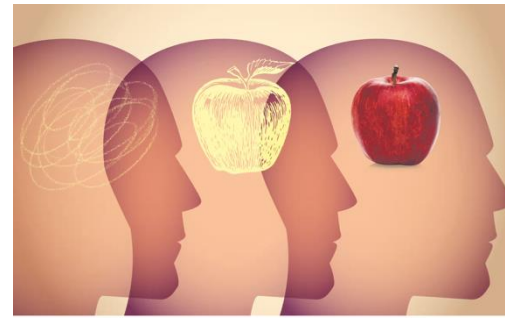
4 - You see an apple, but it looks grey or black

5 - You can't see anything. You only know you are seeing an apple.

# Mental rotation



Which image is the same as the original image, aside from its orientation?

# Mental imagery

The experience of "seeing" (or otherwise sensing)

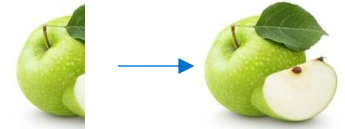in the **mind's eye** without direct external sensory input.

# Mental imagery

Why do different people visualize apples differently in their mind's eye?

# Gestalt Psychology

- **Holism: the whole is greater than the sum of its parts.**
  - Structure & Part-whole relationship;
  - How the world emerges from the integration of its parts.
- **Law of past experience**:
  - Our perception of a apple is not solely derived from its shape, color, or size as sensory inputs; it also incorporates our past experiences and impressions of flowers.
  - Together, these elements form our holistic perception of the apple.

# Generative Prior ⬅ Gestalt Psychology

- Holism: the whole is greater than the sum of its parts.

- Law of past experience

- The era of **generative priors**:

    learn the visual commonsense (e.g., **holism**)
    from huge amount of data (**past experience**),
    encoding high-level structural regularities, as parameters (**deep learning**)

# Generative Priors → Mental World Models

**Novel view from observed view**
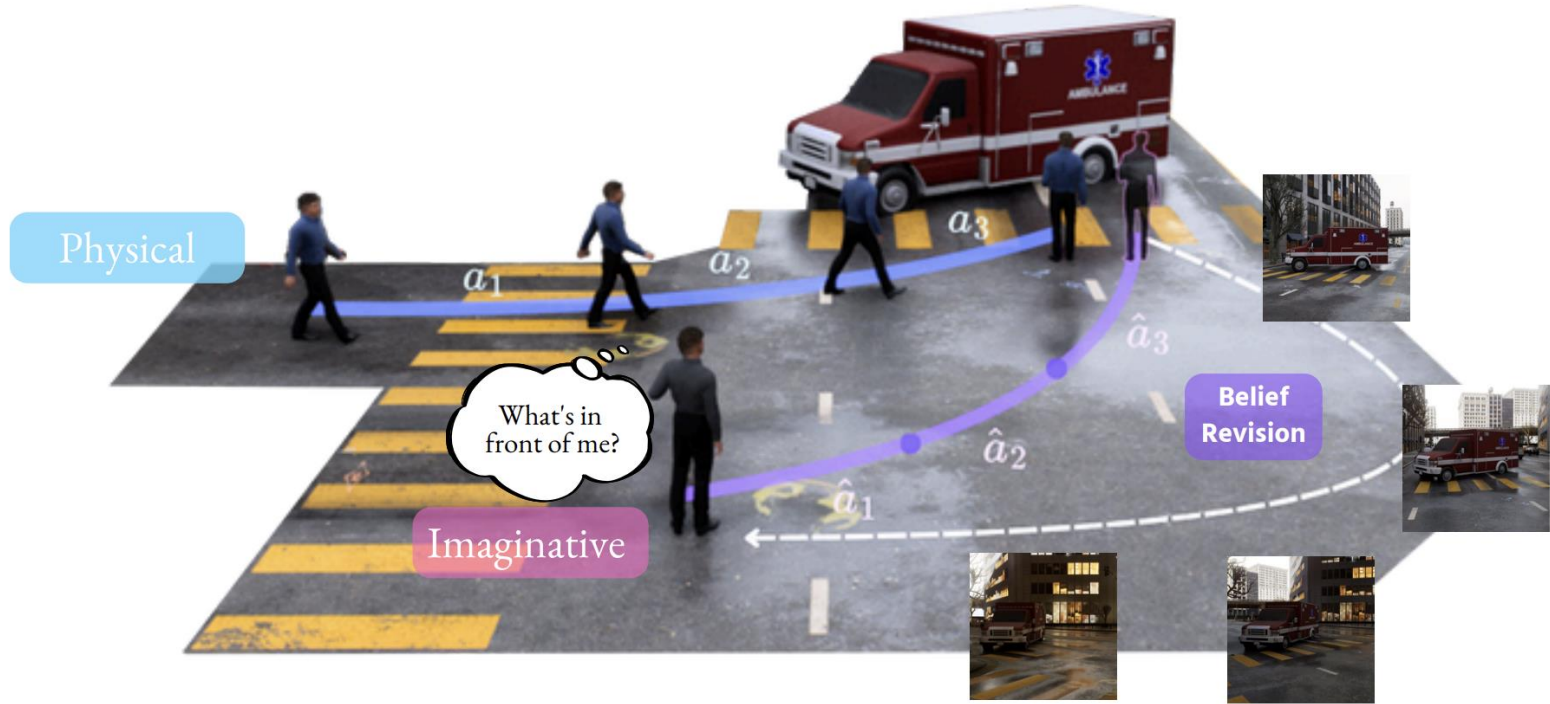
Mental models

**Part-to-Whole Relationship**

Mental models

# Why Mental Imagery/Models Matter?

- Planning with partial observation is challenging.

- Humans can imagine unseen parts of the world through a mental exploration and revise their beliefs with imagined observations.

- Such updated beliefs can allow them to make more informed decisions, without necessitating the physical exploration of the world at all times.

# Mental Exploration

# Mental Exploration Enhances Decision Making



Youtube:
https://www.youtube.com/watch?v=cf4apIcnPtU&ab_channel=CenterforLanguage%26SpeechProcessing%28CLSP%29%2CJHU
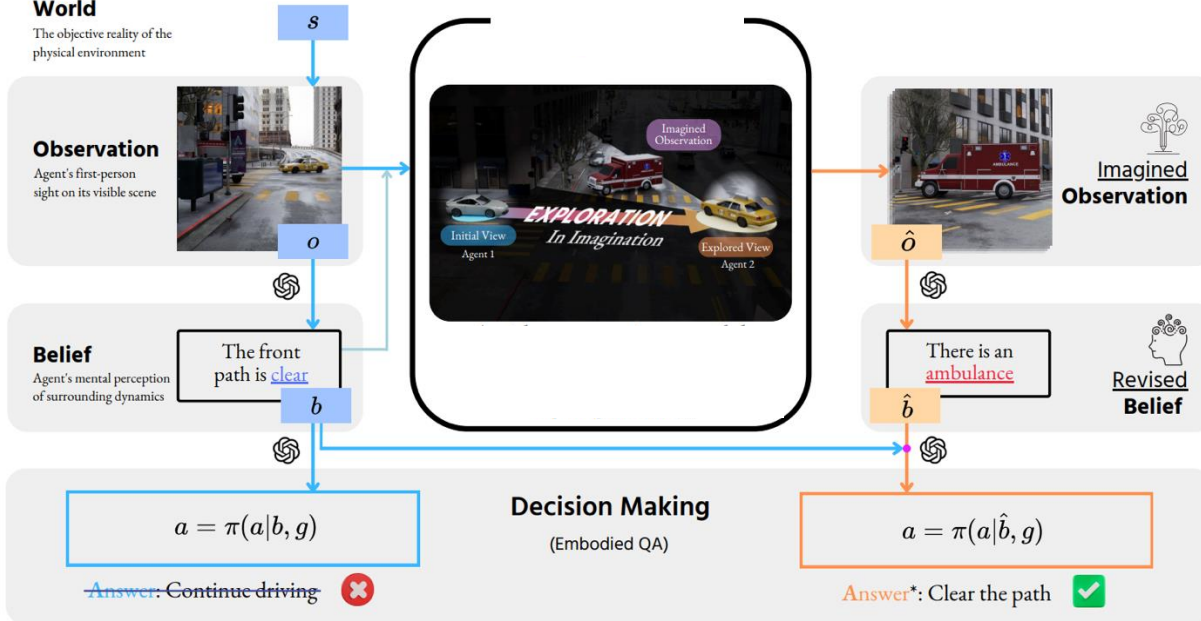
# Mental Exploration Enhances Decision Making



**Question:**
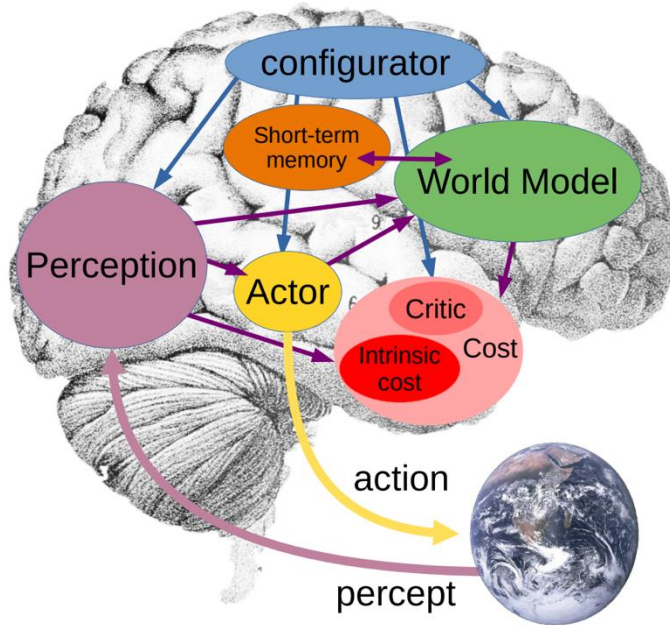Given my observation, what should I do now to cross the street? I can see the taxi ahead suddenly stops.

**World**
The objective reality of the physical environment

$s$

**Observation**
Agent's first-person sight on its visible scene

$o$

**Belief**
Agent's mental perception of surrounding dynamics

The front path is clear

$b$

EXPLORATION In Imagination
Initial View Agent 1
Imagined Observation
Explored View Agent 2

Imagined **Observation**

$\hat{o}$

There is an ambulance

Revised **Belief**

$\hat{b}$

$a = \pi(a|b, g)$

**Decision Making**
(Embodied QA)

$a = \pi(a|\hat{b}, g)$

Answer: Continue driving ❌

Answer*: Clear the path ✅

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# World models: Computational Counterpart of Mental Models
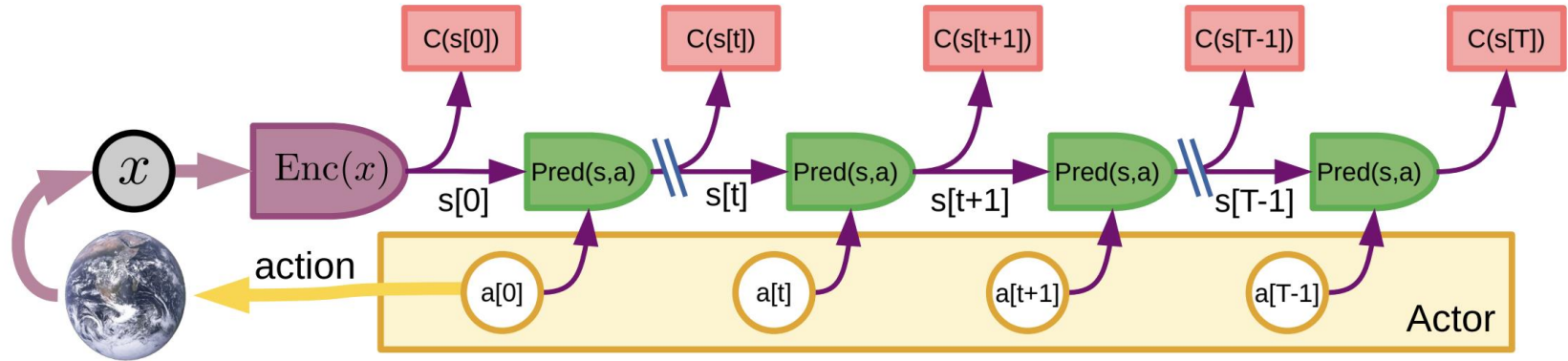
- **Definition**: multi-sensory neural networks that offer a <mark>**predictive distribution over "changes" in the world.**</mark>
  - $p(s_t | s_{t-1}, a_t)$
- **Functionality**: mimic human understanding and interaction by predicting **future world states** (e.g., the existence, properties and location of the objects in a scene) to help agents make informed decisions.

# World Models



The world model module predicts possible future world states as a function of imagined actions sequences proposed by the actor

# World Models



The world model recursively predicts an estimate of the world state sequence using $s[t+1] = \text{Pred}(s[t], a[t])$

# World Models Summary

- multi-sensory neural network that offer a <mark>**predictive distribution over "changes" in the world.**</mark>
  - $p(s_t|s_{t-1}, a_t)$

# Engineering Mental World Models

- Develop **generative models** grounded in physical world.
- The models are capable of predicting world dynamics conditioned on actions.

**Inference**

- **Gather** imagined observation from (interactively) imaginative exploration.
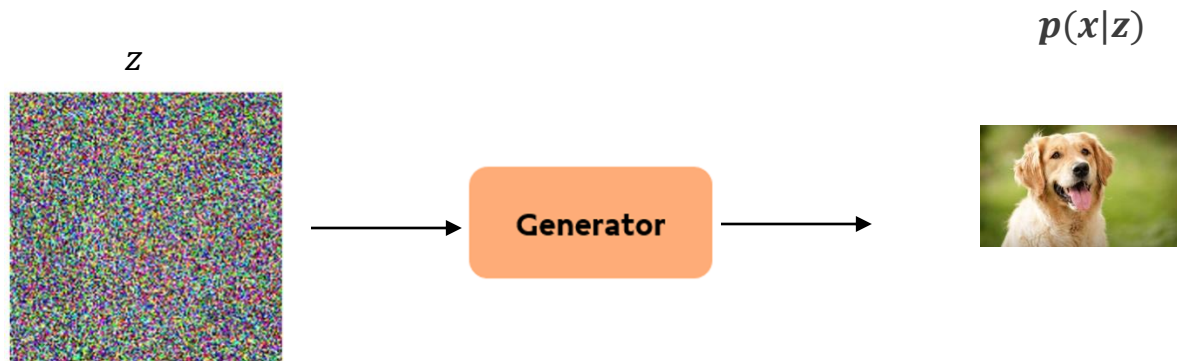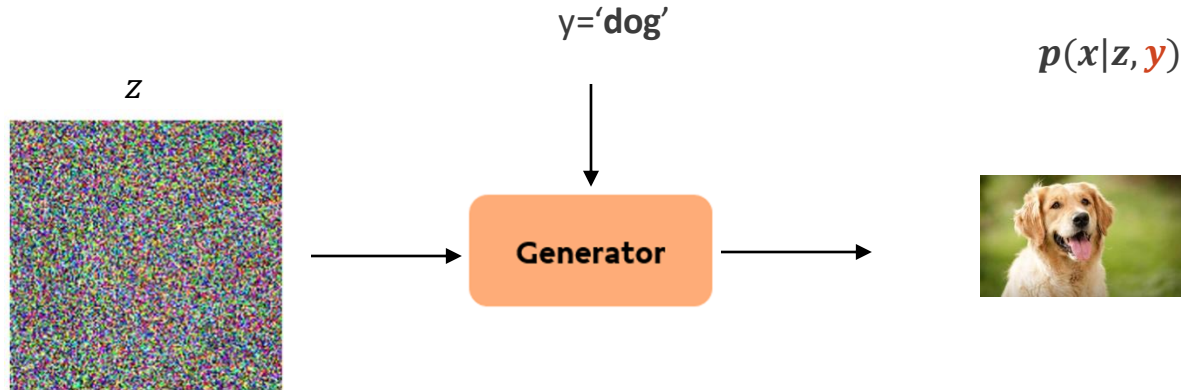- **Planning** with imagined observation.

# Generative Models

# Generative Models

- $z$ : a random variable sample from normal distribution

- $x$ : a predicted data, with the learnt distribution $P(x|z)$

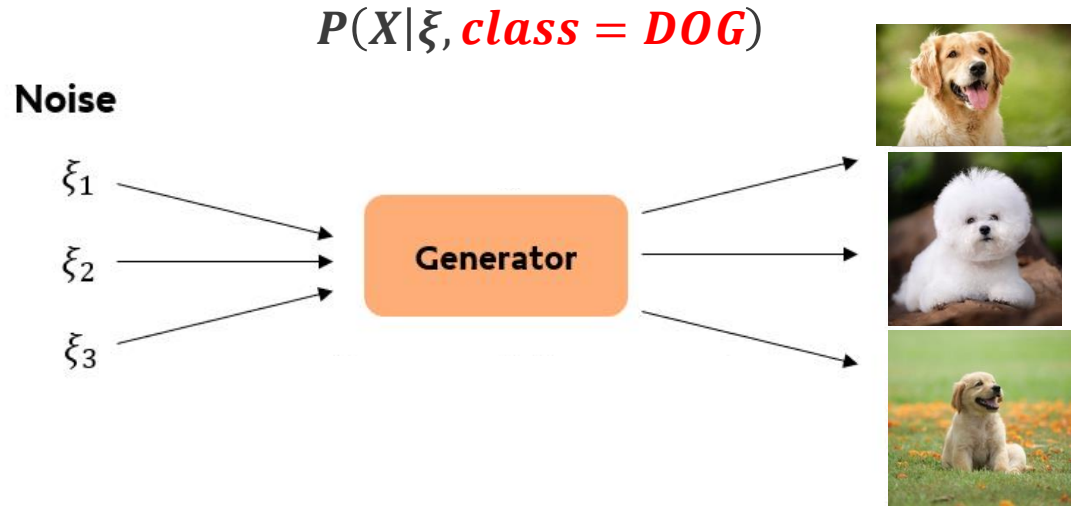- Ideally, we expect the output X is a real image without corruption.

$$p(x|z)$$

$z$

Generator

# Conditional Generative Models

- $z$ : a random variable sample from normal distribution

- $x$ : a predicted data, with the conditional distribution $P(x|z, y)$

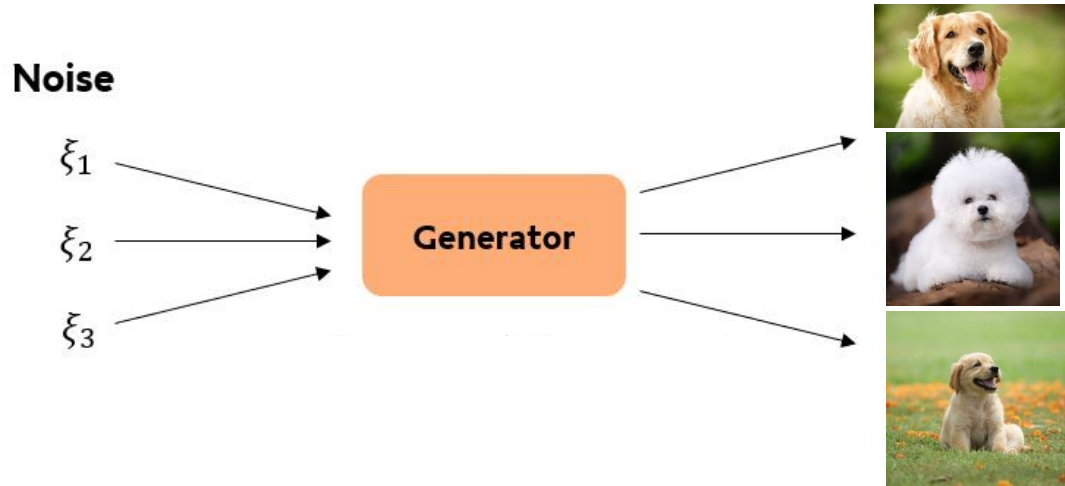- Ideally, we expect the output X is a real image without corruption.

y='**dog**'

$$z$$

$$p(x|z, y)$$

Generator

# Class-Conditioned Generative Models

$$P(X|\xi, \textcolor{red}{\textbf{class = DOG}})$$

**Noise**

$\xi_1$

$\xi_2$

$\xi_3$

Generator

# Class-Conditioned Generative Models

- Class condition can have the same effect of text condition.

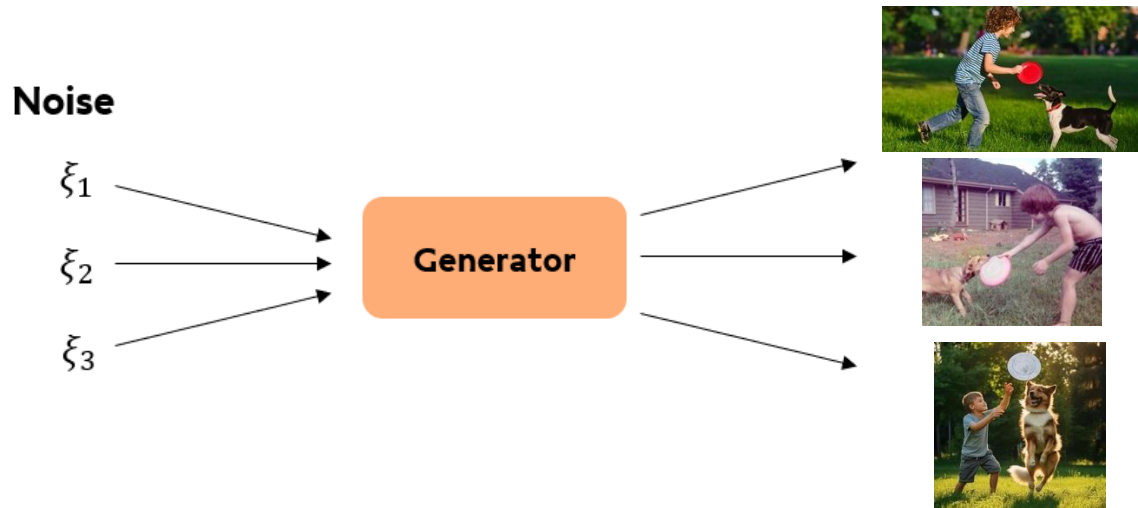- Essentially, the class label of 'DOG' has the same meaning of the text 'a photo of dog'.

$$P(X|\xi, class = DOG) = P(X|\xi, {'A\ photo\ of\ dog'})$$

# Text-Conditioned Generative Models

- progressed from class-conditioned to text-conditioned approaches.

$$P(X|\xi,' \textbf{\textit{a boy is playing frisbee with a dog}}')$$
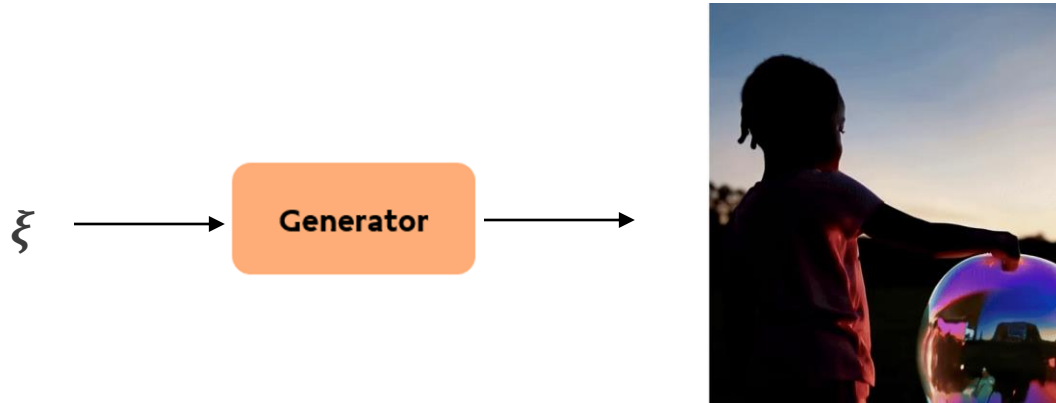
**Noise**

$\xi_1$

$\xi_2$ → Generator →

$\xi_3$

# Text-Conditioned Video Generative Models

- The prediction is not limited to image!
- Video includes dynamics, and thus generating video is harder.

OpenAI
Sora

$$P(X|\xi, 'A\ kid\ throws\ a\ bubble\ into\ air')$$



$\xi \longrightarrow$ **Generator** $\longrightarrow$

# Limitation of Sora-like Models

- **3D consistency?**

- **Physical commonsense?**

- **Interaction?**

# Ground Generative Models in the Physical world

- **Data collected** from 3D physical world (rather than Youtube video)

- **Action as condition** (rather than text)
  - ○ 'the agent is moving two meters forward'
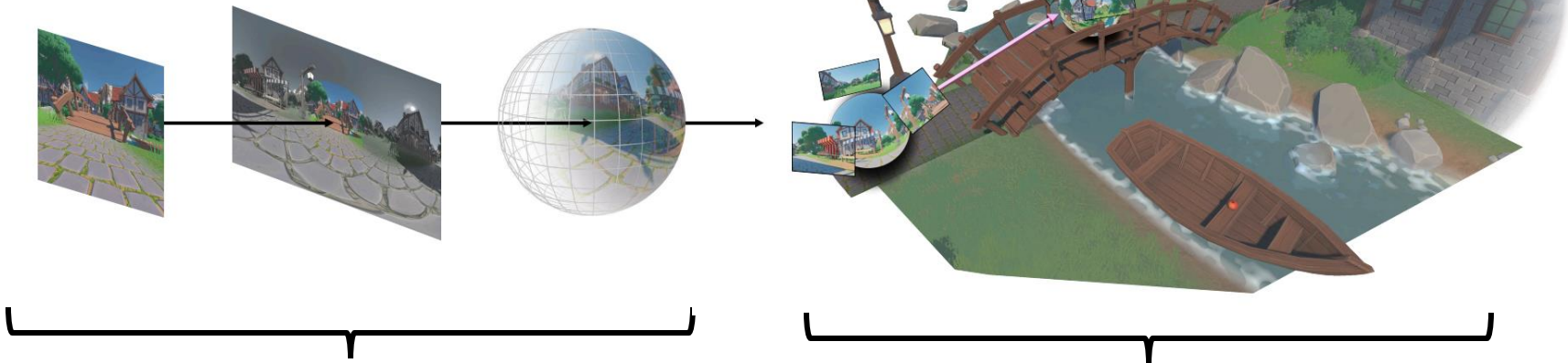
- **Predict the world dynamics**

# Generative World Explorer

Singel Image Input

- **World initialization** (§2.2): Given the initial image $i_0$ and a language description $l_0$, the anchor 360° world view $x_0$ is sampled from:
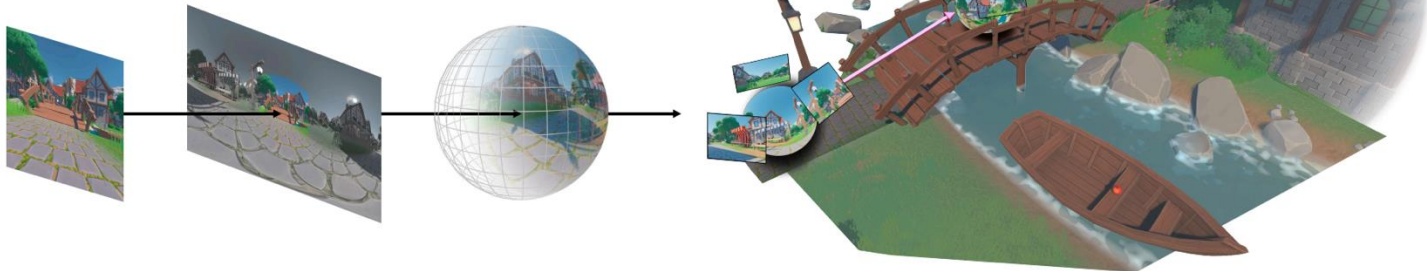
$$x_0 \sim p_{\theta_1}(x \mid i_0, l_0),$$

- **World transition** (§2.3): Given the chosen action $a_t$, the next world view $\mathbf{x}_t$ is sampled from:

$$\mathbf{x}_t = (x_t^0, x_t^1, \ldots, x_t^S) \sim p_{\theta_2}(\mathbf{x} \mid x_{t-1}^S, a_t),$$

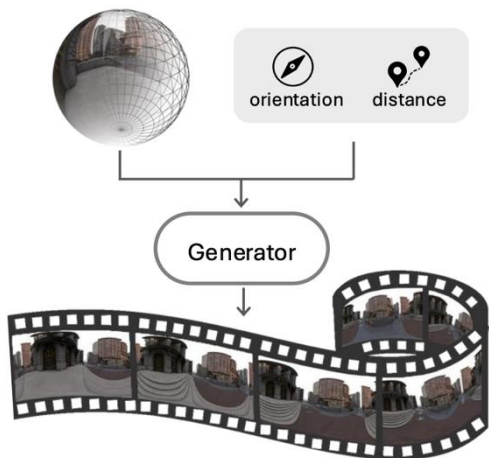where $\theta_2$ is a 360° panoramic video generator, $t = 1, \ldots, T$, and $x_0^S := x_0$.

# World Exploration

Singel Image Input
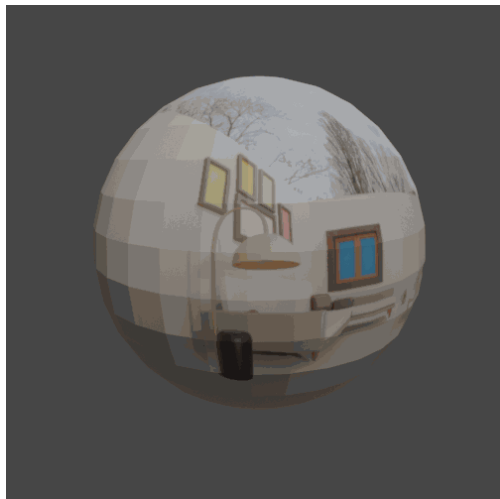


Panoramic World     Action Sampling

orientation    distance

Generator

- **World transition** (§2.3): Given the chosen action $a_t$, the next world view $\mathbf{x}_t$ is sampled from:

$$\mathbf{x}_t = (x_t^0, x_t^1, \ldots, x_t^S) \sim p_{\theta_2}(\mathbf{x} \mid x_{t-1}^S, a_t),$$

where $\theta_2$ is a 360° panoramic video generator, $t = 1, \ldots, T$, and $x_0^S := x_0$.
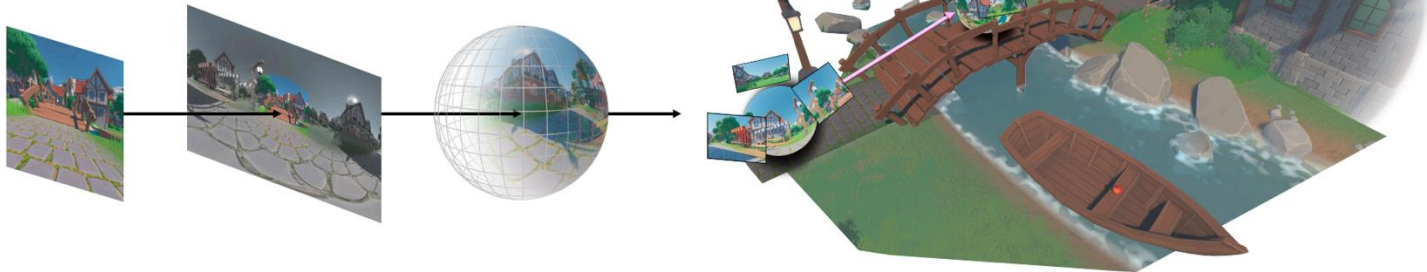
# Exploration at Any Direction

# Forward Exploration (using diffusion model)

## World Exploration

Singel Image Input

## Action Control

- **World transition** (§2.3): Given the chosen action $a_t$, the next world view $\mathbf{x}_t$ is sampled from:

$$\mathbf{x}_t = (x_t^0, x_t^1, \ldots, x_t^S) \sim p_{\theta_2}(\mathbf{x} \mid x_{t-1}^S, a_t),$$

where $\theta_2$ is a 360° panoramic video generator, $t = 1, \ldots, T$, and $x_0^S := x_0$.

# Train on the Data from 3D Synthetic Engines

# Test and Explore in Diverse Scenes

Step into the picture. Imagine the world within.  **Explore**

# Generating Bird's-Eye Worlds

Initialized
Panorama
→
GenEx Upward
Exploration
→
Bird's-Eye
World

# 3D Consistency



Baseline image-to-3D models | GenEx

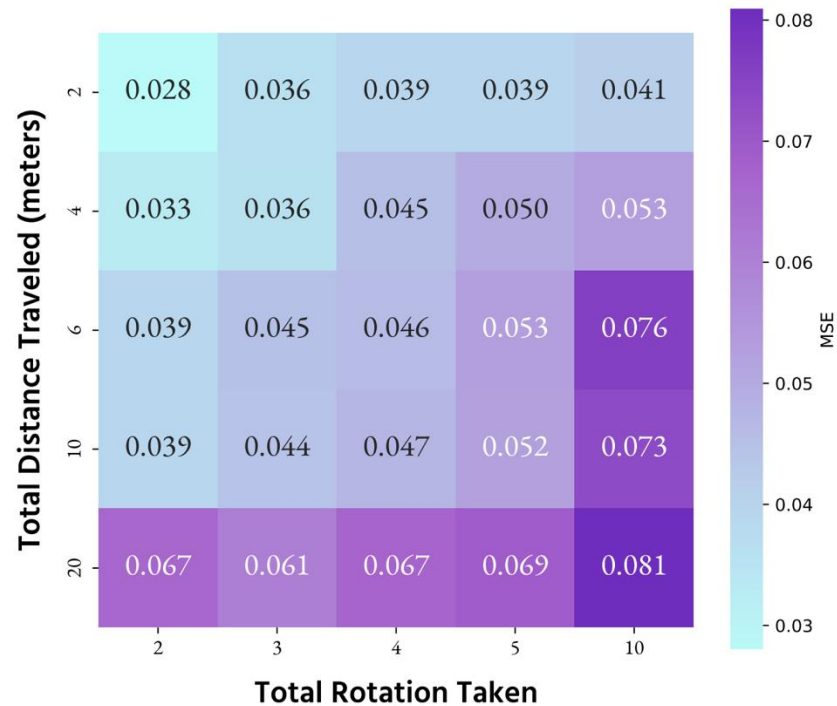**Panorama Input** | **2D Input** | TripoSR | SV3d | **Stable Zero123** | **Ours** | **Ground Truth**
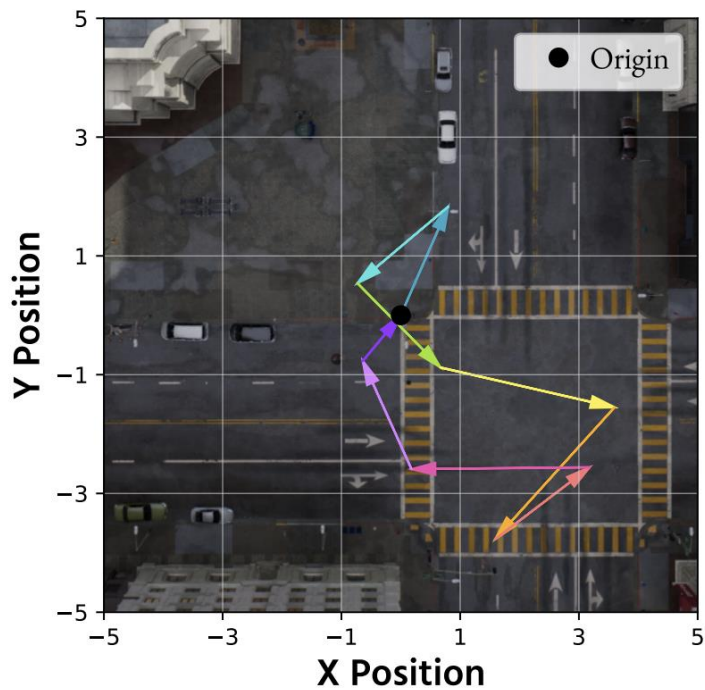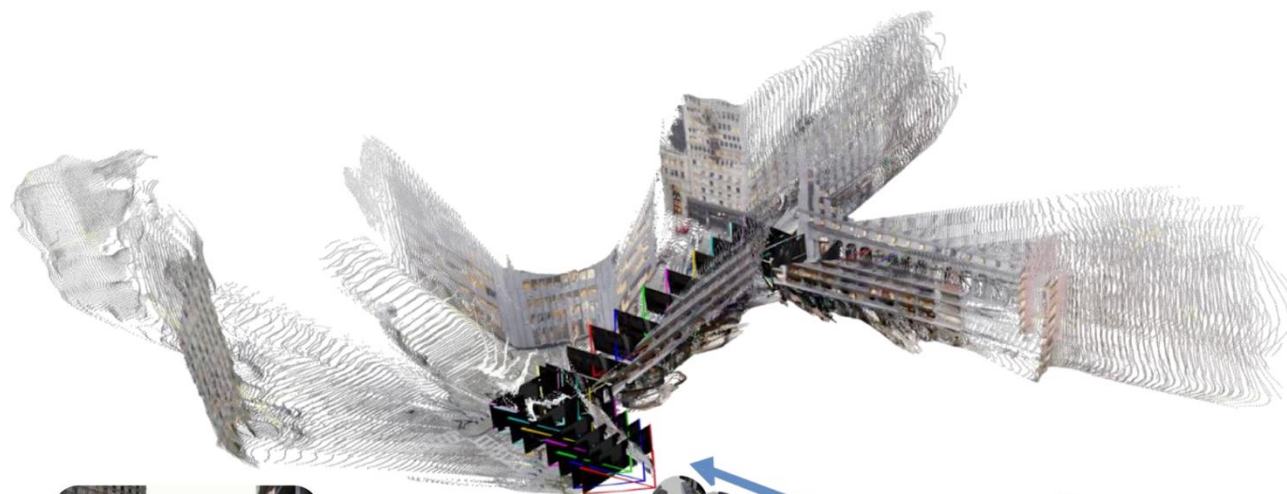
# Loop Consistency when Navigating in the city

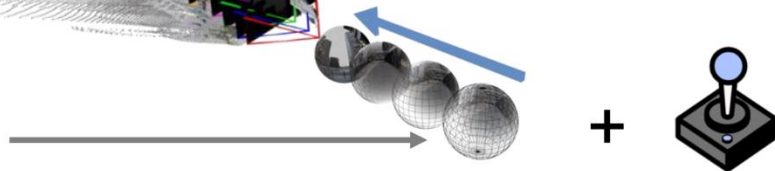# Loop Consistency

# Active 3D mapping through exploration



Single Image

Active 3D Mapping Through Exploration
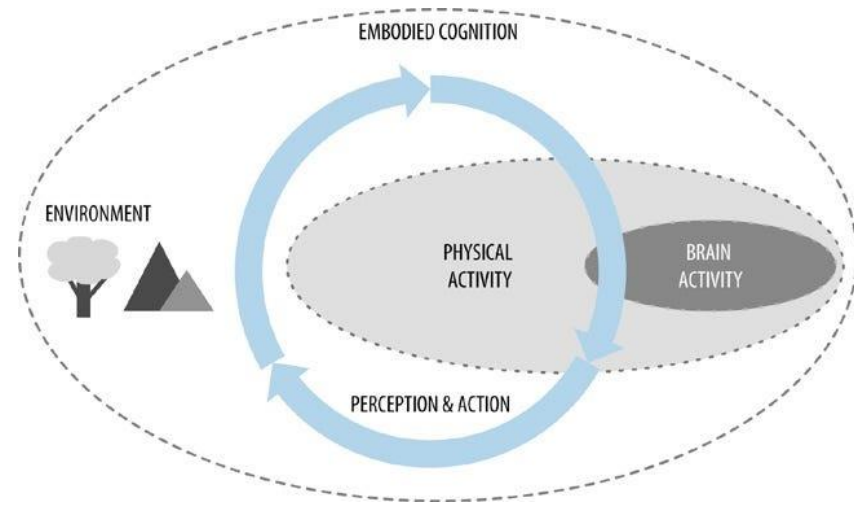
# Active 3D mapping through exploration
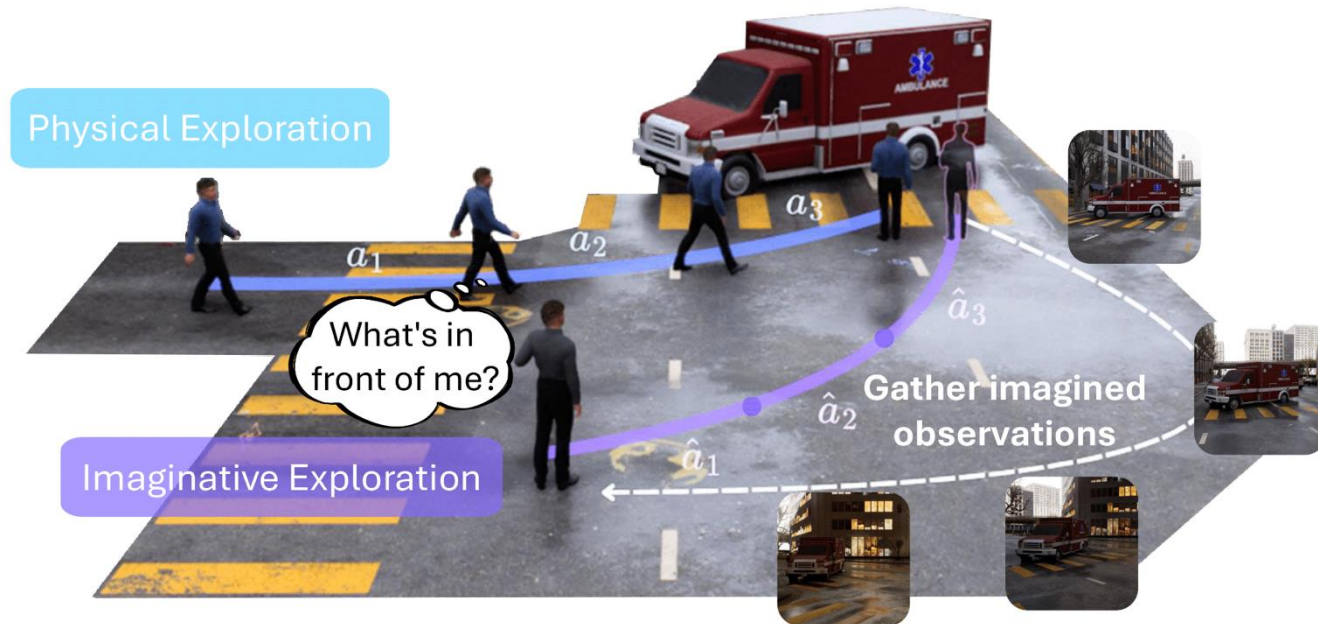
# Connecting GenEx to Embodied AI

# Embodied AI



- Definition of embodied AI:
  - The **embodiment hypothesis**, also known as embodied cognition, is the idea that intelligence is a result of how an agent **interacts** with its environment.

- Connect GenEx to Embodied AI:
  - Predict the change of environment after interaction (agent exploration).

# Replacing Physical Exploration

# Exploration Policy

- The exploration action is decided by a policy:

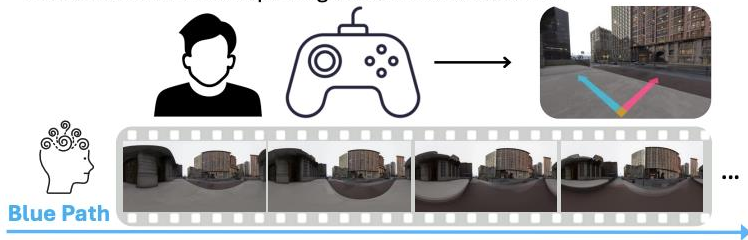$$a_t = \arg\max_a \pi_{explore}(a|x^S_{t-1}, \mathcal{I})$$

- $\mathcal{I}$ is the instruction that specifies the exploration mode to be either human interaction or assisted by a GPT
- $x^S_{t-1}$ denotes the latest explored view from the previous step $t-1$.
- action $a_t = (\alpha_t, d_t)$ defines how the agent rotates its field of view with the rotation angle $\alpha t$ and moves forward with $dt$ distance

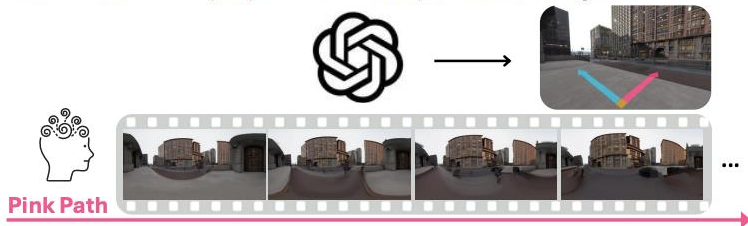# Embodied Exploration: Three Modes



(a) Interactive Exploration

Humans control the exploring direction and distance

Blue Path

(b) GPT-Assisted Free Exploration

Instruction: "Freely explore to observe your surroundings"

Pink Path

(c) Goal-Driven Navigation

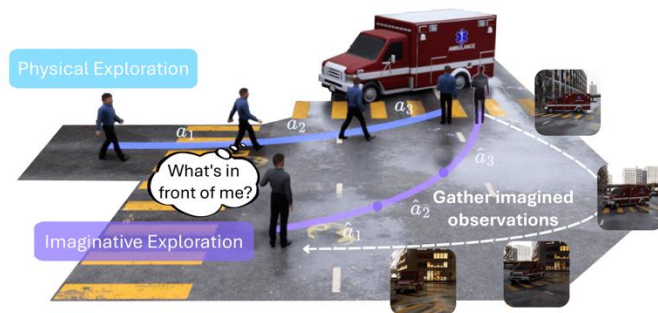Instruction: "Plan to move to the position of the blue car, then turn back."

Turn Right — Move Forward — Turn Left

Turn Back — Move Forward

# **Imagination–Augmented** Policy

**Require:** • Initial observation $i_0$ and world initialization description $l_0$

- A goal $g$ to answer embodied questions. *E.g,* "Danger ahead—stop or go ahead?"
- A navigation instruction $\mathcal{I}$. *E.g,* "Navigate to the unseen parts of the environment."
- GenEx $p(\mathbf{x}_{0:T}|i_0, l_0, \mathcal{I})$ defined in § 2.1 and Algorithm 1.
- An embodied policy $\pi_{\theta_3}(A|o, g)$ conditioned on observation variable $o$ and goal $g$.

1: **Gather imagined observations** with GenEx:

$$\mathbf{x}_{0:T} \sim p(\mathbf{x}_{0:T} \mid i_0, l_0, \mathcal{I})$$

2: **Select an action with imagined observations** to maximize the policy:

$$A = \arg\max_A \pi_\theta(A \mid i_0, \mathbf{x}_{0:T}, g)$$

# **Multi–Agent Imagination–Augmented** Policy

- Step 1 : Gather imagined observations by exploring the position to agent-k

$$\mathbf{x}_{0:T}^{(k)} \sim p(\mathbf{x}_{0:T} \mid i_0, l_0, \mathcal{I}_k)$$

- Step 2: Repeat Step 1 a total of $K$ times, then imaginatively explore the resulting positions of all $K$ agents in our generated explorable world
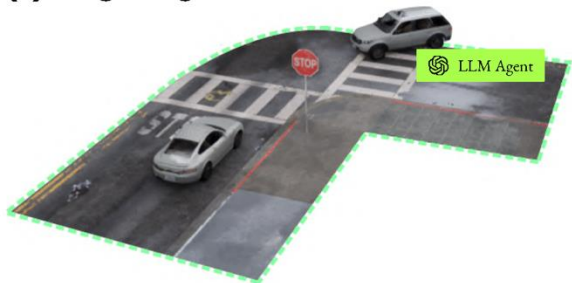
$$\{\mathbf{x}_{1:T}^{(k)}\}_{k=1}^{K} = (\mathbf{x}_{1:T}^{(1)}, \mathbf{x}_{1:T}^{(2)}, \ldots, \mathbf{x}_{1:T}^{(K)})$$

- Step 3: Select an embodied action $A$ with imagined observations to maximize the policy

$$A = \arg\max_{A} \pi_{\theta_3}(A \mid i_0, \{\mathbf{x}_{1:T}^{(k)}\}_{k=1}^{K}, g)$$
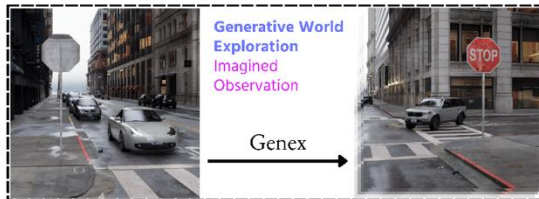
# Embodied Decision Making



**(a) Single-Agent**

**Observation**

I'm turning left at an intersection with no traffic lights. A silver car is slowly moving ahead, and I'm unsure if it will stop. Should I wait?

**Generative World Exploration**
Imagined Observation

Genex

I should stop to avoid a potential collision, as the car might not stop.

The car sees a stop sign and will stop, so I should move to avoid blockage

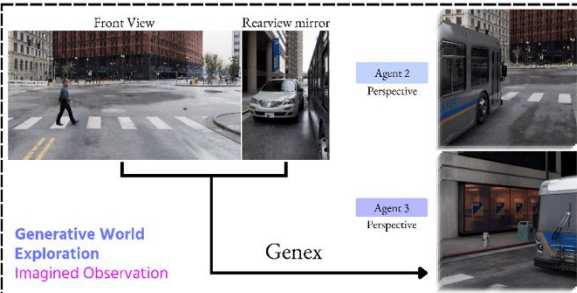**Egocentric Single-View Decision:** Stop in place ❌

**Decision with Imagination:** Continue driving ✅

**Observation**

I'm waiting at the light to move forward, where the right turn is allowed. The front path is clear. A car is driving fast and about to turn right, and a pedestrian is crossing. What should I do?

Front View   Rearview mirror

Agent 2 Perspective

Agent 3 Perspective

**Generative World Exploration**
Imagined Observation

Genex

I want to drive forward, but the light is red, so I should wait in place.

I'm blocking the view between the car and pedestrian, and they might collide.
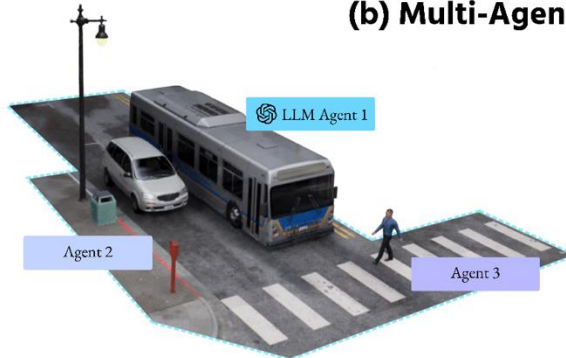
**Egocentric Single-View Decision:** Stop in place ❌

**Decision with Imagination:** Warn both parties ✅

**(b) Multi-Agent**

LLM Agent 1

Agent 2

Agent 3

# Multi-Agent Imagination-Augmented Policy

Enrich real observation with imaginative observation

| Method | Acc. (%) | Confidence (%) | Logic Acc. (%) |
|---|---|---|---|
| Random | 25.00 | 25.00 | - |
| Human Text-only | 21.21 | 11.56 | 13.50 |
| Human with Image | 55.24 | 58.67 | 46.49 |
| Human with **GenEx** | **77.41** | **71.54** | **72.73** |
| Unimodal Gemini-1.5 | 26.04 | 24.37 | 5.56 |
| Unimodal GPT-4o | 25.88 | 26.99 | 5.00 |
| Multimodal Gemini-1.5 | 11.54 | 15.35 | 0.0 |
| Multimodal GPT-4o | 21.88 | 21.16 | 6.25 |
| **GPT4-o with GenEx** | **94.87** | **69.21** | **72.11** |

- Augment human decision making

- Augment GPT decision making

# Imagination–Augmented Policy

Enrich real observation with imaginative observation

| Method | Acc. (%) | Confidence (%) | Logic Acc. (%) |
|---|---|---|---|
| Random | 25.00 | 25.00 | - |
| Human Text-only | 44.82 | 52.19 | 46.82 |
| Human with Image | 91.50 | 80.22 | 70.93 |
| Human with **GenEx** | **94.00** | **90.77** | **86.19** |
| Unimodal Gemini-1.5 | 30.56 | 29.46 | 13.89 |
| Unimodal GPT-4o | 27.71 | 26.38 | 20.22 |
| Multimodal Gemini-1.5 | 46.73 | 36.70 | 0.0 |
| Multimodal GPT-4o | 46.10 | 44.10 | 12.51 |
| **GPT4-o with GenEx** | **85.22** | **77.68** | **83.88** |

- Augment human decision making

- Augment GPT decision making

# Thank you! Question?



Acknowledgement: