

Abstract

In this paper, we introduce a large-scale EXCALIBUR benchmark to encourage and evaluate embodied agents' ability to explore. We also collected human performance with Virtual Reality headsets through an immersive integration of our environment. In the supplementary material, we provide the following items that shed further insight into these contributions:

- A An immersive interface to EXCALIBUR, that enables human annotators to transfer their knowledge from the real world to the virtual environment.*
- B Objects in EXCALIBUR have a large diversity in terms of categories, colors, materials, sizes, and weights. Within each category, there are different handcrafted realistic object models.*
- C A large diversity of procedurally generated houses.*
- D Question filtering process that automatically chooses the most challenging questions for all scenes.*
- E Handcrafted templates that are used for generating questions.*
- F An explanation of Eq. 1 and how the constant is found based on human performance.*
- G Detail description of the baseline model we used.*
- H Discussion on detailed evaluation of agents and humans performance.*

A. Immersive Human Performance Collection for EXCALIBUR

In the supplementary material package, we provide a sample of a human episode (`human_traj.mp4`). Here, we mainly discuss the technical details of designing the VR interface for AI2THOR.

A.1. Design choices of VR interface

The most straightforward reason to use VR for data collection is that controlling the arm and grasping with keyboard input can be non-intuitive, while VR makes it possible to directly apply the sensorimotor experience. The agent's arm movement replicates human arm movement in all spatial dimensions, developing hinge joints for elbows, ball-and-socket joints for shoulders, and synovial joints for the wrist. These simulations are done using the Oculus hand controllers, facilitating a smooth interface for human performance. A user can pick up and drop objects in the VR, using the Oculus Grip Button, wherein the intensity with which the grip button is pressed determines the force with which the arm will exert to pick up an object in VR. Since different objects in the room have different masses, we require different amounts of force to pick up different objects. This will map a human's intuition of how heavy an object is to that of the agent's arm, making it a nearly exact emulation of the human mind coupled with actual physics.

This realistic interface also provides a great challenge for comparing human performance with agents. With OpenXR, we can use the exactly same scenes in both VR and agents' experiments. However, the VR interface has a continuous action space (except for opening drawers, fridges and closets), while agents' action space is discrete. To make the comparison as fair as possible, we also count humans' actions as discrete time steps. All participants completed the data collection while sitting and try to only controller to move and rotate. The grasp force is discretized: whenever the force changes more than 0.05 kg, we count it as a timestep. The same applies to joint movements. This underestimates humans' performance, but we still find the huge gap between humans and best model as shown in Tab. 2.

A.2. Data collection procedure

We conducted the whole human performance collection in the following steps:

1. 7 college students with no prior experience of using VR headsets are recruited as participants;
2. To train each participant, we use a hand-crafted ARCHITECTHOR scene, which is not in the test set. After the participants wear the headset, we ask them to freely explore the houses for as long as they want, and to try out all of the actions on the controller. We look at a mirrored screen and give instructions whenever a collision happens. A training session ends when the participant lifts and drops every pickupable object and correctly answers 20 questions consecutively. The training sessions last about one hour per participant.
3. Each participant does 4 scenes with replay and 5 scenes without replay or 5 scenes with replay and 4 scenes without replay. The difference between with and without replay is whether the participant can look at the replay of their own

recording of Phase I and Phase III in Phase II and Phase IV. Before a participant enters the scene, we would tell them whether this episode is with replay or without. When a participant is exploring in the first Phase, we look at the live number of time steps they consumed. Whenever it reaches 2,500 steps, they are forced to stop. In practice, most of the participants spend less than 2,500 steps before deciding to stop themselves. One episode typically lasts about one hour per participant.

B. Full object list

See Table 4 for a listing of objects in AI2-THOR and their default properties (*e.g.* mass, volume, \parallel).

C. House distribution

Alg 1 shows the procedure to randomize houses, agent spawn locations and object physical properties.

Algorithm 1 House randomization procedure

Require: ProcTHOR houses \mathcal{S}

```

for  $S \in \mathcal{S}$  do
   $\mathcal{P} \sim$  reachable locations in  $S$ 
  while Agent collides with objects at location  $\mathcal{P}$  do
     $\mathcal{P} \sim$  reachable locations in  $S$ 
  end while
  for object in  $\mathcal{H}$  do
    Sample scale factor  $\alpha \sim [0.8, 1]$ , and mass factor  $\beta \sim [0.5, 1.5]$ 
    object.size =  $\alpha \times$  object.size
    object.mass =  $\beta \times$  object.mass
  end for
end for

```

D. Question filtering

Alg. 2 shows the procedure to filter questions based on the distributions of answers. For each question, we calculate the answer distribution among all houses. We choose the questions with the highest

Algorithm 2 Question filtering

Require: Template set \mathcal{T} , Houses \mathcal{S} , number of questions generated per scene N , maximum number of questions per template M

```

 $Q \leftarrow \emptyset$ 
for  $T \in \mathcal{T}$  do
   $Q' \leftarrow$  Dict[question  $\rightarrow$  Counter[answer]]
  for  $S \in \mathcal{S}$  do
    for  $q \in$  Generate( $T, S$ ):  $N$  do
      if  $\exists$  valid  $a \leftarrow$  Answer( $S, q$ ) then
         $Q'[q][a] \leftarrow Q'[q][a] + 1$ 
      end if
    end for
  end for
   $Q' \leftarrow$  Rank $_{q \in Q'}$  Ent  $\leftarrow$  Shannon-Entropy  $Q'[q]$ 
   $Q \leftarrow Q \cup Q'[: M]$ 
end for

```

Index	Template
1	What number of [C2] [M2] [S2]s are <R> the [C1] [M1] [S1]?
2	Are there any [C2] [M2] [S2]s <R> the [C1] [M1] [S1]?
3	What color is the [M2] [S2] [that is] <R> the [C1] [M1] [S1]?
4	What is the material of the [C2] [S2] [that is] <R> the [C1] [M1] [S1]?
5	What is the [C2] [M2] object [that is] <R> the [C1] [M1] [S1]?
6	Is there anything else that has the same color as the [C1] [M1] [S1]?
7	Is there anything else that has the same material as the [C1] [M1] [S1]?
8	Do the [M1] [S1] and the [M2] [S2] have the same color?
9	Do the [C1] [S1] and the [C2] [S2] have the same material?
10	Do the [M1] [S1] <R> the [C3] [M3] [S3] and the [M2] [S2] have the same color?
11	Do the [C1] [S1] <R> the [C3] [M3] [S3] and the [C2] [S2] have the same material?

Table 5. Eleven different question families. C_i , M_i , and S_i ($i = 1, 2, 3$) are one of the colors, materials, and object types listed in Tab. 4. [] denotes that word in the parenthesis can be omitted (S_i ($i = 1, 2, 3$) can be changed to “object”). R denotes one of the relations specified in Tab. 6

Index	Relation / Property	Explanation
1	Color	the most salient one or two colors of the object
2	Material	the most salient one to three materials of the object
3	CONTAINEDBY	the relationship describing an object is located in a receptacle
4	ADJACENTTO	two objects are within 0.5 meters and no objects are on the line connecting the two
5	ONTOPOF	One object’s bounding box collides with another’s on the bottom and top respectively
6	HEAVIERTHAN	an object’s mass is higher than the other’s by at least 0.05 kg
7	LIGHTERTHAN	an object’s mass is lower than the other’s by at least 0.05 kg
8	LARGERTHAN	an object’s longest, middle and shortest dimensions are all longer than the other’s longest, middle and shortest dimensions respectively
9	SMALLERTHAN	an object’s longest, middle and shortest dimensions are all shorter than the other’s longest, middle and shortest dimensions respectively
10	LONGERTHAN	for a few object types, including a baseball bat, candle, ladle, and plungers, if one object’s longest dimension is longer than the other’s
11	SHORTERTHAN	for a few object types, including a baseball bat, candle, ladle, and plunger, if one object’s longest dimension is shorter than the other’s

Table 6. Relations and properties in EXCALIBUR.

E. Question families

F. Exploration score

We consider a normal distribution prior to the exploration score, i.e. $\text{ExQA} \sim \mathcal{N}(\mu, \sigma)$, and maximize the likelihood of the scores of human annotators

$$k^* = \arg \max_k \max_{\mu, \sigma} \mathbb{E}_{\text{ExQA}} \left[-\frac{1}{2} \left(\frac{\text{ExQA} - \mu}{\sigma} \right)^2 - \log \sigma \right]. \quad (5)$$

G. Baseline Model

Our baseline model is a GRU [13] f_ϕ^{GRU} with parameters ϕ , which at each step t takes observation o_t , action at last time step a_{t-1} , and questions q as input, and outputs belief:

$$h_t = f_\phi^{\text{GRU}}(h_{t-1}, f_\psi^{\text{obs}}([\text{MLP}_1(\text{CLIP}(o_t)); \text{MLP}_2(\text{action-emb}(a_t)); \text{MLP}_3(\text{T5}^{\text{encoder}}(q))])), \quad (6)$$

where f_{ψ}^{obs} and $\text{MLP}_i, i = 1, 2, 3$ are all three-layer MLPs, CLIP [48] encodes the observation frame at time step t , action-emb is an embedding layer, and the encoder of T5 [49] encodes a text string composed of questions \mathcal{Q} after Phase II and outputs zero vectors in Phase II, since the questions are not available.

The belief h_t is used for two different purposes: (1) predict actions

$$a_t \sim \text{softmax}(f_{\xi}^{\text{action}}(h_t)), \quad (7)$$

where a_t is an action in the space shown in Fig. 2, and f_{ξ}^{action} is a three-layer MLP converting belief to action logits; (2) predict answers to questions

$$\text{answer}_i = \text{beam-search}(p_{\text{T5}}(a \mid [f_{\theta}^{\text{prefix}}(h_t), f^{\text{emb}}(q_i)])), \quad (8)$$

where $f_{\theta}^{\text{prefix}}$ is an MLP generating soft prompts for the frozen T5 similar to [60], and $f^{\text{emb}}(q)$ is the embedding layer of question.

H. Further Analysis of Results

To better understand the different failing patterns of humans and reinforcement learning agents, we study the following four metrics as well

1. TTI – number of time steps till touching an object for the first time,
2. Area – the ratio of regions covered by the agents,
3. % Objects Seen – the ratio of objects that are seen,
4. and % Objects Interacted – the ratio of objects that are touched.

Where do humans go wrong? Human failures mainly result from forgetting, and model failures mainly result from a lack of exploration. Through comparing humans w/o replay and w/ replay we can confirm the first part, and through comparing human area, we can clearly see that when exposed to an immersive setup, an average user covers more area, hence performing better. Also, when we compare Phase III to Phase I, the performance definitely improves as the humans now know what exactly to observe. Another reason why humans lag in the first phase is because humans tend to remember objects only when they have a pre-defined goal in mind. A general exploration, doesn't trigger the humans to learn every thing they see. Which also explains why Phase III performance is better than Phase I.

Relation of % Objects seen/interacted to Performance: We see an increase in the number of Objects Seen and Objects Interacted for human performance, which is again easier to do so in the VR environment. When the humans are able to interact with a greater number of items, they firstly make better idea of the things placed in the setup. Also, interaction stimulates the memory, making them remember which objects they interacted on each of the rooms. Overall, effect of these help them get a better spatial mapping and they tend to perform better. If we closely observe the results, better performance has a direct relation to the percentage of objects seen or interacted.

TTI v/s Accuracy: An interesting observation is that for human performance, the TTI in the third phase increases compared to the first phase as the humans know which part of the house to explore and which object to interact with. Due to this, human tends to avoid interacting with all objects and checks only the relevant ones. Also, TTI for humans after watching the replay is higher than without replay case because after watching the replay, re-interaction is required only for a few objects, which explains the increase in TTI. Furthermore, humans intelligently explore those parts of the house which they missed during Phase I exploration, thereby reducing the area covered in the re-entering phase.

	Phase I Exploration				Phase III Reentering			
	TTI	Area	% Obj Seen	% Obj Inter	TTI	Area	% Obj Seen	% Obj Inter
QA reward	303.5	11.2	23.4	10.4	202.3	21.2	15.1	24.9
Novelty reward	571.2	31.4	33.5	5.6	172.1	23.2	23.0	3.2
Novelty+QA	436.1	32.8	37.9	9.9	153.2	23.6	20.5	25.8
Human w/o replay	53.7	59.3	96.7	83.8	65.7	28.5	83.1	54.0
Human w/ replay	68.5	70.7	98.5	81.2	154.3	26.6	75.4	39.4

Table 7. Analysis on exploration and reentering behavior of agents and humans.