# Supplementary Materials for REALFRED: An Embodied Instruction Following Benchmark in Photo-Realistic Environments

Taewoong Kim<sup>1,\*</sup><sup>(b)</sup>, Cheolhong Min<sup>1,\*</sup><sup>(b)</sup>, Byeonghwi Kim<sup>1</sup><sup>(b)</sup>, Jinyeon Kim<sup>1,2</sup><sup>(b)</sup>, Wonje Jeung<sup>1</sup><sup>(b)</sup>, and Jonghyun Choi<sup>1,†</sup><sup>(c)</sup>

<sup>1</sup> Seoul National University <sup>2</sup> Yonsei University {twoongg.kim, cheolhong.min, byeonghwikim}@snu.ac.kr, jinyeonkim@yonsei.ac.kr, {specific0924, jonghyunchoi}@snu.ac.kr

Note: Green denotes reference to the main paper.

# A Benchmark Details

We provide all 114 types of objects in Fig. 1. The **bold** text denotes the uniquely introduced object in ours.

# A.1 Annotation Interface

In this section, we describe the overall process of acquiring language annotations. Fig. 2a illustrates the interface of the Mechanical Turk used to collect human annotations from Mechanical Turk workers. We provided workers with an expert demonstration video and divided the timeline segments that have the intended subgoal (e.q., 'pick up the pencil case,' 'qo to the sofa'). Workers were asked to fill each segment with their own words (e.g., 'Pick up the pencil case from the coffee table,' 'Walk around the table to get closer to the sofa'). Workers were paid 0.7 per annotation as following the previous work [61]. Moreover, we adopted a voting survey to filter out inappropriate annotations. Fig. 2b illustrates the interface of getting votes from workers. We conducted the voting with a minimum of 2 and up to 5 reviewers per annotation. Only annotations that received more than a majority of accepts in all cases were included in our set of annotations. For annotations that did not achieve a majority of accepts, we re-collected annotations and implemented a voting system to prevent the inclusion of low-quality annotations. We paid workers \$0.35 to compare 5 sets of annotations following [61].

# A.2 Vocabulary Distribution

We provide vocabulary statistics for the language instructions in the REAL-FRED benchmark in Fig. 4.

<sup>&</sup>lt;sup>\*</sup>Equal contribution. <sup>†</sup>Corresponding author.



**Fig. 1: Object list in ALFRED and REALFRED.** Listed in alphabetical order. The object classes in the REALFRED benchmark are a superset of those in ALFRED, with newly introduced objects highlighted in bold.

		Cho If all	ose one ins instruction Summary: Description	struction IT as any officent contents with other instructions a describe the same actions and well-tighted with video, choose "All good" Calcet to twile up par pass or to the design table Calcet to twile exect any other same action of the	
	• car/friz	[2]	Summary: Description	Part for glate along with the both on the table. 1. On one to the disease. 2. Part for the both one the both. 2. Part for each one the gates and the table. 4. Part to the both grates and the table. 4. Part to the table one the gates and the table. 4. Part type is an operator to the table. 4. Part type is an operator to the table. 4. Part type is an operator to the table.	
I. Re 2. R 3. R 4. R	levent Object: desser d	[3]	Summary: Description	Tenden the grine with here in the table.   So are named and gap over the distance.  So are the values of the distance.  So are the values	
5. R	elevant Object: coffeetable	Vide	90	+	
Goal:	n instruction that summarizes what the robot should accomplish by the end of the vides	Pick (1) (2) (3) (3) (4)	one amon aaya different saya different saya different I good. (All inst	g the options or wang contents with intern or wang contents with intern or wang contents with intern therefore descele the source actives and web-signed with video.) Submitt	
	(a) Annotation interface.		(b	) Voting survey interface.	

Fig. 2: Mechanical Turk interface.

## A.3 Scanned Indoor Houses

Among collected 150 scenes, we split scenes into 135 seen and 15 unseen environments. Then we further split into validation (both seen and unseen) and test (both seen and unseen) folds. The details are presented in Table 1. Note that validation and test unseen scenes are exclusive.

**Degree of photorealism.** We compare a degree of photorealism by measuring FID [27] and KID [4] scores following Ramakrishnan *et al.* [56]. We use rendederd

	Train	Vali	dation	Test		
		Seen	Unseen	Seen	Unseen	
# of scenes	135	131	6	135	9	

Table 1: Indoor house splits.

Table 2: Photorealism comparison.						
Environments	HM3D		Gibson		Gibson HQ	
	$\mathrm{FID}{\downarrow}$	$\mathrm{KID}{\times}10^{3}\downarrow$	$\mathrm{FID}{\downarrow}$	$\mathrm{KID}{\times}10^{3}\downarrow$	$\mathrm{FID}{\downarrow}$	$\mathrm{KID}{\times}10^{3}\downarrow$
TDW [20]	129.52	$89.13 {\pm} 1.52$	122.68	$80.38 {\pm} 1.36$	124.89	$79.05 \pm 1.49$
BEHAVIOR-1K [41]	113.59	$81.58 {\pm} 2.96$	99.52	$66.03 {\pm} 1.85$	107.53	$63.03 {\pm} 1.99$
REALFRED	81.06	$69.96{\pm}1.47$	83.25	$71.05{\pm}1.51$	101.60	$89.98{\pm}1.72$
Gibson [73]	43.79	$31.66 {\pm} 1.05$	_	_	46.67	$33.74{\pm}1.21$
HM3D [56]			43.79	$31.66 {\pm} 1.05$	26.33	$21.70{\pm}1.03$

RGB images from each synthetic environments [20, 41]. We compare the image quality with a set of RGB images rendered from previous dataset, HM3D [56] and Gibson [73]. To compare with previous scanned environments, we acquire a collection of real RGB images derived from high-resolution raw panoramas in Gibson. We designate this collection as Gibson HQ. Results are presented in Table 2. We observe that ours achieves the lowest (*i.e.*, best) FID and KID scores when compared to the previous environments [20, 41] that provide interaction with the environments including the objects. However, despite these promising metrics, our results reflect lower photorealism compared to the previous scanned environments [73, 56] which do not provide interactable environments. This may be due to manually added agent-interactable objects.

Manual removal of 3D objects from background. Designers use automated tools (*i.e.*, Blender) to extract objects from the background meshes easily and set their properties (e.g., labels, colliders, interactability, etc.).

In Fig. 3, we provide the qualitative examples of how the source data were fixed by data correction. Workers fill in the missing parts or smooth out the uneven surface by checking the overall alignment composition after mapping and alignment using automated tools such as Blender.

Qualitative examples. We show several houses used in the REALFRED benchmarks in Fig. 8 and 9.

### A.4 Examples of Expert Demonstration

Fig. 10-13 illustrate the examples of expert demonstrations for 7 task types. The agent has to solve the task in interactive environments by understanding the language instructions and planning the sequential and executable actions.



Fig. 3: Qualitative examples of the 3D mesh data correction. In each row, we provide source data and the corrected data.

# A.5 Diverse Episodes with More Objects

Each episode is generated based on the combination of task-relevant objects (e.g., put a 'knife' on the 'table.') and this indicates that more object classes can result in more object-diverse episodes (e.g., put a 'potato' in a 'fridge.'). We observe that our REALFRED provides more diverse episodes compared to the ALFRED benchmark [61] by enriching the number of object classes. Here, we denote an episode whose combination of task-relevant objects does not overlap with the others by a unique episode. We observe that our REALFRED benchmark provides the 4,649 unique episodes while ALFRED provides 2,522 ones in the combined train and valid splits. In addition, the ratio of the unique episodes among the total ones is 53.3% in REALFRED while it is 35.6% in ALFRED. This indicates that our REALFRED benchmark provides not only a larger number of episodes but also more diverse combinations of episodes.

## A.6 Qualitative Comparison of Indoor Houses

We provide a qualitative comparison of the indoor houses used in our REAL-FRED and ALFRED [61] in the attached video files (*video1.mp4*, *video2.mp4*, and *video3.mp4*). We provide the agent's egocentric view on the left-hand side of a video and a top-down view with a red circle denoting the agent's corresponding current location. While ALFRED's indoor house environments consist of single room types on a room scale, ours are provided on a house scale, featuring multiple rooms within a single house. This implies that agents developed with our environments are enabled to perform instruction-following tasks that require navigating through multiple rooms.



Fig. 4: Vocabulary statistics in collected human language instructions.

# **B** Details of State of the Art Models

We provide details of state-of-the-art models in imitation learning and spatial map reconstruction, respectively.

# **B.1** Imitation Learning

The Seq2Seq [61] model encodes the visual input with the frozen backbone visual encoder. The natural language goal and instructions are encoded with a bidirectional LSTM encoder to produce an embedding for each word. Alongside the previous action, embeddings are passed as input to an LSTM cell to produce the current hidden state. The action and corresponding mask are finally predicted using a hidden state. MOCA [62] exploits separate branches for action prediction and object localization to better address different semantic understanding. ABP [33] extends MOCA [62] by perceiving surrounding perception for a better understanding of environments with the enlarged field of view.

# **B.2** Spatial Map Reconstruction

HLSM [5] uses a hierarchical controller to bridge the gap between natural language instructions and agent executable actions. The high-level controller predicts the next subgoal given the instruction and the map, and then the low-level controller outputs a sequence of actions to achieve the subgoal. FILM [46] utilizes a pre-designed template as a high-level action sequence. It uses two submodules of BERT classifiers to predict the type of instruction and the arguments to fill in the template. Finally, it uses a deterministic algorithm [60] for obstacle-free path planning. LLM-Planner [63] leverages large language model to generate subgoal sequence with a few examples. To enhance LLMs planning accuracy, it updates plans that are physically grounded in the environment. CAPEAM [34] uses context-aware planning to plan a subgoal sequence and conduct the respective subgoal with the corresponding detailed planners. It also uses additional memory to prevent the interaction of inappropriate objects.

Table 3: Task and Goal-Condition Success Rate (valid split). Path-lengthweighted (PLW) metrics are given in parentheses for each value. We report mean and stndadrd deviation over multiple runs. <sup>†</sup>Authors' implementation as the code is not publicly available.

		Validation						
Learning	Model	Se	en	Unseen				
		Success Rate	Goal Condition	Success Rate	Goal Condition			
Imitation Learning	Seq2Seq [61] MOCA [62] ABP <sup>†</sup> [33]	$\begin{array}{c} 0.77\pm 0.06 \ (0.47\pm 0.06) \\ 12.64\pm 0.12 \ (8.35\pm 0.16) \\ 24.71\pm 0.05 \ (15.49\pm 0.34) \end{array}$	$\begin{array}{c} 6.93 \pm 0.06 \ (4.73 \pm 0.06) \\ 20.95 \pm 0.18 \ (13.43 \pm 0.16) \\ 33.80 \pm 0.14 \ (23.27 \pm 0.32) \end{array}$	$\begin{array}{c} 0.00 \pm 0.00 \ (0.00 \pm 0.00) \\ 1.44 \pm 0.05 \ (0.56 \pm 0.06) \\ 4.22 \pm 0.05 \ (1.70 \pm 0.08) \end{array}$	$\begin{array}{c} 4.03\pm 0.06 \ (2.50\pm 0.00) \\ 6.76\pm 0.04 \ (3.64\pm 0.06) \\ 11.71\pm 0.27 \ (5.42\pm 0.13) \end{array}$			
Spatial Map Reconst.	HLSM [5] FILM [46] LLM-Planner <sup>†</sup> [63] CAPEAM <sup>†</sup> [34]	$\begin{array}{c} 4.23\pm0.08\ (0.72\pm0.08)\\ 7.08\pm0.28\ (1.87\pm0.11)\\ 5.80\pm0.19\ (1.51\pm0.03)\\ 13.45\pm0.05\ (3.43\pm0.06)\end{array}$	$\begin{array}{c} 9.14\pm0.09~(2.67\pm0.06)\\ 11.93\pm0.23~(4.82\pm0.15)\\ 11.69\pm0.35~(4.76\pm0.07)\\ 18.16\pm0.27~(4.50\pm0.05) \end{array}$	$\begin{array}{l} 1.08\pm0.14 \ (0.19\pm0.03) \\ 4.44\pm0.17 \ (1.25\pm0.10) \\ 3.33\pm0.22 \ (0.96\pm0.05) \\ 4.92\pm0.22 \ (1.22\pm0.03) \end{array}$	$\begin{array}{c} 6.12\pm 0.23 \ (1.52\pm 0.02) \\ 9.26\pm 0.13 \ (3.84\pm 0.11) \\ 8.29\pm 0.19 \ (3.49\pm 0.09) \\ 9.47\pm 0.23 \ (1.79\pm 0.04) \end{array}$			

Table 4: Task and Goal-Condition Success Rate (test split). Path-lengthweighted (PLW) metrics are given in parentheses for each value. We report mean and studadrd deviation over multiple runs. <sup>†</sup>Authors' implementation as the code is not publicly available.

_	Model	Test						
Learning		Se	een	Unseen				
		Success Rate	Goal Condition	Success Rate	Goal Condition			
Imitation Learning	Seq2Seq [61] MOCA [62] ABP <sup>†</sup> [33]	$\begin{array}{c} 1.10\pm 0.00 \ (0.05\pm 0.01) \\ 14.11\pm 0.03 \ (9.20\pm 0.05) \\ 27.44\pm 0.40 \ (16.96\pm 0.16) \end{array}$	$\begin{array}{c} 6.60 \pm 0.00 \ (5.00 \pm 0.00) \\ 22.84 \pm 0.04 \ (16.42 \pm 0.04) \\ 35.81 \pm 0.23 \ (24.57 \pm 0.19) \end{array}$	$\begin{array}{l} 0.00\pm 0.00 \ (0.00\pm 0.00) \\ 0.62\pm 0.08 \ (0.35\pm 0.05) \\ 3.54\pm 0.23 \ (1.51\pm 0.08) \end{array}$	$\begin{array}{c} 3.50\pm 0.00 \ (2.80\pm 0.00) \\ 5.14\pm 0.08 \ (3.39\pm 0.06) \\ 10.57\pm 0.22 \ (5.59\pm 0.10) \end{array}$			
Spatial Map Reconst.	HLSM [5] FILM [46] LLM-Planner <sup>†</sup> [63] CAPEAM <sup>†</sup> [34]	$\begin{array}{c} 6.27\pm 0.04 \ (0.88\pm 0.10) \\ 8.79\pm 0.07 \ (2.36\pm 0.01) \\ 8.16\pm 0.20 \ (2.20\pm 0.06) \\ 15.61\pm 0.15 \ (3.68\pm 0.09) \end{array}$	$\begin{array}{c} 10.44\pm 0.13 \ (2.78\pm 0.10) \\ 13.03\pm 0.08 \ (5.58\pm 0.08) \\ 13.20\pm 0.13 \ (5.72\pm 0.06) \\ 20.22\pm 0.11 \ (5.39\pm 0.09) \end{array}$	$\begin{array}{l} 0.49\pm 0.16 \ (0.08\pm 0.03) \\ 2.15\pm 0.18 \ (0.56\pm 0.04) \\ 1.90\pm 0.13 \ (0.57\pm 0.04) \\ 2.87\pm 0.13 \ (0.84\pm 0.02) \end{array}$	$\begin{array}{c} 4.28 \pm 0.13 \ (1.37 \pm 0.16) \\ 6.56 \pm 0.15 \ (3.16 \pm 0.05) \\ 6.33 \pm 0.02 \ (3.09 \pm 0.04) \\ 7.36 \pm 0.07 \ (2.01 \pm 0.03) \end{array}$			

# C Extended Quantitative Results

We present experiment results with path-length-weighted success rate and goal condition (*i.e.*, PLWSR and PLWGC) over multiple runs in Table 3 and 4.

# D Map Reconstruction Strategy

We provide a more detailed analysis of the challenges in recognizing narrow passages (e.g., doors, aisles, etc.). Our observation is as follows: failure to recognize narrow navigable pathways leads an agent to *stuck* within the initial room.

To quantify this challenge, we establish a criterion for *leaving* a room by taking 1 step (*i.e.* 0.25 meter) further from the entrance of the room where the agent was initiated. Fig. 5 represents the top-down view of one of *valid unseen* scenes labeled for each space from room 1 to room 6. We observe that only 5.3% of the agents, from the total episodes conducted on the scene, left one room to another for further exploration. We also noticed that above mentioned *leaving* occurred within a short span, not exceeding 77 steps. This implies that if an agent initially overlooks narrow spaces, they will be mistaken for walls when viewed at an oblique angle.

We provide qualitative examples in Fig. 6, illustrating failure cases in reconstructing narrow passages. We observe that an agent that adopted a spatial map



Fig. 5: Top-down view of a scene with labeling. We select one of the largest scene in *valid unseen* fold and annotate the space from room 1 to room 6 based on the door gap for further analysis. The area inside the outer black line is a navigable area.

fails to reconstruct its surroundings, fails to recognize narrow doorways, gets *stuck* in a single room, and eventually, fails to complete a task.

We further explore the likelihood of an agent, under a random navigation policy, to be positioned where it can directly observe narrow spaces or door gaps to construct a semantic map accurately. We assess the likelihood within 100 steps (*i.e.* a slight buffer extended to 77 steps). In detail, we consider the agent's position, viewing direction, and the horizontal angle of the head to measure the likelihood. As a result, the agent has a 6.9% likelihood of aligning with and facing the exit passage to recognize it, similar to the empirical result of the agent leaving the room where it was initiated (*i.e.*, 5.3%). This may imply that not being positioned correctly to see doors shortly after its deployment increases the chance of getting stuck in a room. This consideration may not have been necessary in the ALFRED benchmark [61], which consists of single rooms. However, since REALFRED is composed of multiple rooms, 'Spatial Map Reconst.' baselines may achieve much lower performance compared to [61].

# E Qualitative examples for domain adaptation

We provide qualitative examples of real-to-sim domain adaptation in Fig. 7. 'Source' column denotes a visual frame from the REALFRED benchmark, 'Cy-



Fig. 6: Reconstructed spatial map and egocentric view. A reconstructed spatial map is presented on the left-hand side and an egocentric view of the agent is presented on the right-hand side on each figures. The green area denotes a predicted navigable area and the dark gray area denotes a predicted obstacle. Agent fails to recognize narrow navigable space (highlighted in red on the right-hand side of the each figure).



Fig. 7: Qualitative examples of the real-to-sim domain adaptation. In each column, we provide source image, domain adapted image for the source image, and an image from target domain.

cleGan' column denotes an adapted visual frame with CycleGan [79], 'UVCGANv2' column denotes an adapted visual frame with UVCGAN-v2 [70], and 'Target' column denotes an image from ALFRED target domain image. The advantage of real-to-sim domain adaptation is to make an agent feel at home, reducing the visual domain gap. We expect a domain adapted image to resemble some characteristics that are well represented in the target domain, where sim2real agent is trained.

We begin our examination with an example in the first row. We observe that the flooring, dominating the image frame, is adapted to resemble the flooring that frequently appears in the target domain (highlighted with  $\square$ ). We also notice that the image generated with CycleGan adapts the color of the wall to brown, while the image generated with UVCGAN-v2 adapts the wall to white tone (highlighted with  $\square$ ).

We now examine an example in the second row. We observe that the flooring, dominating the image frame as in the first example, is generated to resemble checkerboard tiles which are represented in target domain (highlighted with  $\Box$ ).



Fig. 8: Example of houses used in the REALFRED benchmark.



Fig. 9: Additional example of houses used in the REALFRED benchmark.



Annotation #1

Goal: Take the egg from the pot and put it on the coffee table

Instructions: Instructions: Make from the chair and turn right and get out from the bed room and turn left wake towards to living room and turn left and wake straight and turn left and get have a straight and turn left and wake straight and turn left and get have a straight and turn left. Wake head towards to the egg in the living room. Turn left and turn around to the living room. Again turn right and turn around in the coffeetable. Put the pot on the coffeetable in the living room.

### Annotation #2

Goal: Ready for cook to prepare well

Instructions: Turn left and right move to the shelf in the kitchen. After reach the kitchen go to the egg rack. In the pot rack below the rack take one egg. Take the egg from the pot and move to the kitchen. In that egg put the pot on the table. Take that bow and move to the coffeetable area. Then place the pot in the coffeetable in the hall area.

Annotation #3

Goal: Place the cooking vessel with egg on the coffeetable.

Instructions: Turn around and go near the shelf table. Take the egg from the table. Turn around and go near the pot. Put the egg into the pot on the table. Take the pot with egg from the table. Turn around and go to the coffeetable. Put the egg on the coffeetable.



### Annotation #1

 Goal:

 Put two ladles in the bathroom sink.

 Instructions:

 Turn left and head to the fridge.

 Pick up the ladle from the fridge.

 Turn left and set the passage on the left and reach the bathroom sink on the right.

 Put the ladle on the sink.

 Turn left and set the room and turn left to head right to the living room coffeetable.

 Pick up the ladle from the fiving room table.

 Turn left and set down the passage on the left and reach the bathroom sink on the right.

 Put the ladle in the sink.

### Annotation #2

Goal: Place two ladles in a bathroom sink.

Instructions: Go a bit straight and left through this room, approaching the fridge on the left. Remove the ladle from the fridge. Turn back towards the area you started, then turn left at the white table, towards the wall with a doorway on each side, entering the sink on the right. Place the ladle inside the sink. Exit this room the way you came in and go back into the larger room, turning right to go towards the dark coffeetable in front of the couch. Place the batter from the table. Return to the bathroom you placed the first spoon in the sink.

### Annotation #3

Goal: Obtaining spatula from various locations and arranging them in a table. Instructions: Turm around and head straight, then proceed to turn left towards the fridge. Pick up spatula from the desk. Take spatula on the table. The around a going straight, proceed to turn to the left towards the coffee Pick up spatula from the desk. Turn right and head straight, then hang left towards the sink. Place spatula on the table.

Fig. 10: Examples of expert demonstration and human annotation ('Heat & Place' on the left and 'Pick Two & Place' on the right). We provide examples of expert demonstration for tasks 'Heat & Place' and 'Pick Two & Place.' The black lines denote the expert's trajectories, and several egocentric views are presented alongside a top-down view of the scene.



### Annotation #1

Goal: Put the ButterKnife with the Bowl on the CounterTop. Instructions: Turn left to reach the shelf at the bottom of the bed. Pick up the butterknife from the shelf of the dresser. Turn around to head out of the bedroom and walk down the living room bowl on the left. Put the butterknife in the brown bowl. Pick up the bowler from the table. The table of the microwave. Put the bowl on the countertop.

### Annotation #2

Goal: Take the ButterKnife and Bowl to store the CounterTop.

Instructions: Turn left and left move to the shelf. After reach the bedroom to take butterknife from bedroom table. Take the knife and move to the bowl. After reaching the hall, put the butterknife into the bowl on the hall table. Take the bowl and move to reach the kitchen. After reaching the kitchen, place the bowl on the countertop. Place the bowl on the countertop near the oven.

### Annotation #3

Goal: Move ButterKnife and Bowl to CounterTop.

# Moto statustication of the state of the shelf. Walk into the bedroom towards the white shelf. Pick up the butterknife from off the dresser Walk with the stick towards the living area towards the bowl on the coffee table. Set the butterknife in the brown bowl. Pick up the brown bowl from the coffee table. Walk towards the countertop with the bowl. Set the butterknife the content of the microwave on the countertop.

Pick & Place



### Annotation #1

Goal: Take the lettuce to store the fridge

The the breact of some inter integritude of the int

### Annotation #2

Goal: Cut the lettuce, keep the butter knife on the sofa and keep the piece of lettuce inside the fridge.

Instructions: Stand up turn around walk towards the dining table. Pick up the butter knife from the dining table. Turn around walk towards the sofa. Cut the lettuce on the sofa. Keep the butter knife on the top of the sofa. Turn around move towards the sofa. Turn around move towards the sofa. Pin up this piece of the lettuce forgine. Keep the butter of the lettuce forgine.

### Annotation #3

Goal: Slice the lettuce and put it safely on the fridge

Instructions: Turn around your right and cross the living room and move head to the dining

Turn around your right and cross the wring iscanse table. Pick up the butter knife near the watch. Turn your left and move towards sofa on the living room. Turn your right towards the lettuce on the sofa. Out the lettuce into pieces using the butter knife. Put the knife on the sofa near the lettuce. Turn your right turn left towards the lettuce. Pick up one of the siled lettuce on the sofa. Turn your right and across the room and move towards the fridge. Put the lettuce into the top rack on the fridge near the bread.

Fig. 11: Examples of expert demonstration and human annotation ('Stack & Place' on the left and 'Pick & Place' on the right). We provide examples of expert demonstration for tasks 'Stack & Place' and 'Pick & Place.' The black lines denote the expert's trajectories, and several egocentric views are presented alongside a top-down view of the scene.



### Annotation #1

Goal: Place a slice of Banana on the CoffeeTable. Instructions: Turn right oward the chair and coffeetable in the corner, go across to the coffeetable and chair. Pick up the knife on the table. Turn right go to the first door on the left, go across to the banana. Slice the banana on the bed. Place the knife on the third shelf on the middle. Place the knife on the third shelf on the middle. Turn right at the doorway, go to the circle shaped table and go straight across to the fridge. Place the banana in the fridge and pick it back up again. Turn right and go across to the coffeetable.

### Annotation #2

Goal: Put the cooled slice of Banana on the CoffeeTable.

Instructions: Turn around to head out of the room and reach the coffeetable. Pick up the knife from the table. Turn around and take a left turn to head to the banana. Cut the banana on the bed with the knife. Put the knife on the shoff. Pick up the aligned of banana from the bed. Turn around to head to the kitchen on the right and reach the fridge. Head to the round coffeetable on the right. Put the slice of banana on the coffeetable.

### Annotation #3

Goal: Put the half Banana on the CoffeeTable. Instructions: Go near to the coffeetable. Take the knife and walk to the left. Go to the banana room. Knife to cut the banana room the bed. Cut the banana in half and take it to the shelf to exit the room. Go to reach the banana roupboard. Put the banana into the rdge. Walk hock and from the cupbrard. Half hock and from the cupbrard. Put the half banana on the coffeetable.



### Annotation #1

Goal: Examine a carry bag by the light of a tall lamp.

Instructions: Turn right go straight to the kitchen. Open to the cupboard under the sink. Take the carry bag from the cupboard turn around go straight walk near to the lamp.. Turn on the lamp

### Annotation #2

Goal: Switch on the floor lamp to view the cookie.

Instructions: Turn around in the living room and move towards the kitchen sink. Pick the cookie from the cupboard under the sink. Turn left and walk towards the floor lamp in the living room. Switch on the floor lamp.

### Annotation #3

Goal: Take the cookie in the sink turn on the floor lamp.

Instructions: Walk around the room. Take the cookie in the sink. With the cookie in hand to reach the floor lamp. Turn on the floor lamp.

Fig. 12: Examples of expert demonstration and human annotation ('Cool & Place' on the left and 'Examine in Light' on the right). We provide examples of expert demonstration for tasks 'Cool & Place' and 'Examine in Light.' The black lines denote the expert's trajectories, and several egocentric views are presented alongside a top-down view of the scene.



Fig. 13: Example of expert demonstration and human annotation ('Clean & Place'). We provide an example of an expert demonstration for task 'Clean & Place.' The black line denotes the expert's trajectory, and several egocentric views are presented alongside a top-down view of the scene.