

Figure 5: Task gap between ALFRED and R2R. We highlight notable distinctions between the navigation tasks in ALFRED and R2R, encompassing variations in appearance, step size, and instruction complexity. See appendix B for more details.

## A IMPLEMENTATIONS DETAILS

We used the LLaMA-7B model [Touvron et al. \(2023a\)](#) and the LLaMA2-7B model [Touvron et al. \(2023b\)](#) for our method, fine-tuning it on 72 V100-32GB GPUs with a batch size of 144. The training tokens had a maximum length of 1024, while during inference, the maximum length was set to 2048. The AdamW optimizer [Loshchilov & Hutter \(2017\)](#) with a learning rate of  $2 \times 10^{-5}$  and weight decay of 0 was employed for optimization. The WarmupDecayLR learning rate scheduler was used for learning rate scheduling. For image captioning in both the R2R and ALFRED tasks, BLIP [Li et al. \(2022a\)](#) was utilized. Deformable DETR [Zhu et al. \(2020\)](#) was used for object detection in the R2R dataset, with suppression of outdoor object categories. We used the ground-truth object detection results provided in ALFRED when we generated the instruction-following pairs in § 4.2. When prompting GPT-4 API, we set the temperature as 1 and top\_p as 1. The cost of collecting the generated trajectories by prompting GPT-4 API [OpenAI \(2023\)](#) was around \$500. In the few-shot learning experiments in § 4.1 and § 4.2, we set  $\rho = 0$ . While when fine-tuning with the full train set in § 5, we set  $\rho = 0.2$ . We pretrain on 128K ALFRED instruction-following pairs whose format is given in § 3.2. We augment the observations in ALFRED to 12 views and randomly mask a variable number of views to mimic the irregular number of candidates in R2R.

## B DIFFERENCES BETWEEN ALFRED AND R2R.

There are significant differences between ALFRED and R2R which makes straightforward sim2real transfer challenging. These differences include:

**Visual appearance.** ALFRED uses images rendered from the synthetic AI2THOR environment, while R2R, based on the Matterport3D, incorporates images captured from real indoor environments. These image sources differ in texture, occlusion, illumination, and other visual aspects.

**Step size.** There is a difference in step sizes between the two tasks (see the right part of fig. 5). ALFRED uses a step size of 0.25 meters, while R2R has larger and more variable step sizes. To bridge this gap, we consolidate four consecutive MoveAhead steps into a single step along the ALFRED trajectory.

**Action type.** A complete ALFRED trajectory includes not only navigation actions but also interaction actions, where the interaction actions are combined with a target object to change the state of the surrounding environment. In order to filter the interaction actions in ALFRED, we divide each ALFRED trajectory into multiple sub-trajectories and keep the sub-trajectories that are labeled with the GotoLocation tag.

Table 5: Performance of the Multi-task Model on R2R. We demonstrate the multi-task capability of the LM agent. For single-task models, each model is trained within the task data. We trained the multi-task model with data from both R2R and ALFRED tasks.

Models	R2R Seen		R2R Unseen	
	SR $\uparrow$	SPL $\uparrow$	SR $\uparrow$	SPL $\uparrow$
Single-Task	55.0	51.0	43.2	37.9
Multi-Task	<b>55.9</b>	<b>51.7</b>	<b>45.6</b>	<b>40.0</b>

Table 6: Performance of the Multi-task Model on ALFRED. ST: Single-Task. MT: Multi-Task.

	ALFRED Seen		ALFRED Unseen	
	Task $\uparrow$	GC $\uparrow$	Task $\uparrow$	GC $\uparrow$
ST	0.0 (0.0)	6.0 (4.7)	0.5 (0.1)	9.5(7.8)
MT	0.0 (0.0)	<b>6.4 (5.0)</b>	<b>0.6 (0.2)</b>	<b>9.8 (7.8)</b>

**Instruction complexity.** Due to trajectory splitting, ALFRED’s navigation trajectories and instructions appear simpler and shorter compared to R2R’s instructions. R2R instructions involve guiding the agent between rooms, whereas ALFRED trajectories mainly keep the agent within a single room.

**Action space.** In ALFRED, the agent is limited to rotating left/right by  $90^\circ$  and moving forward, while in R2R, the agent can move in any combination of 12 candidate heading directions and 3 elevation directions. The number of available movement directions is irregular. This difference in action space makes R2R trajectories more human-like. To address this, we introduce randomness by adding or reducing a heading offset of  $\pm 30^\circ$  to the agent’s direction at each step in ALFRED, allowing rotations of  $30^\circ$  or  $60^\circ$  in addition to  $90^\circ$ .

## C MULTI-TASK PERFORMANCE

One of the advantages of our approach is its inherent suitability for multitasking. Similar to LLMs use instruction to handle multiple language tasks concurrently, we consolidate task information and inputs into instructions. To validate the multitasking capability of our method, we extend its application to the ALFRED task.

**Metrics on ALFRED.** We evaluate our model on ALFRED using two metrics: *Task Success* (Task) and *Goal-Condition Success* (GC). Task Success measures the ratio of trajectories where object positions and state changes accurately match all task goal conditions at the end. GC assesses the ratio of completed goal conditions in each action sequence. Task Success is only considered successful when GC is also 1. On average, each ALFRED task has 2.55 goal conditions. We also calculate the *Path Length Weighted Metrics* (PLW) for both Task and GC, which normalize the metrics based on the actual action sequence length.

**Results of the Multi-Task Model.** In ALFRED task, we set  $\rho = 0$  as the expert policy in ALFRED is suboptimal. To save training time and balance the data amount between R2R and ALFRED, we utilize only 50% of the training dataset, resulting in a dataset for ALFRED with 386K data pairs. For R2R task training, we maintain  $\rho = 0.2$  and run each demonstration trajectory twice, resulting in a training set size of 235K for R2R. Consequently, the merged dataset for the multitask model contains a total of 621K instruction-following data pairs. We select VLN Bert [Hong et al. \(2021\)](#) as the baseline for the R2R task and Seq2seq model [Shridhar et al. \(2020\)](#) for the ALFRED task. Given the substantial differences between the R2R task and the ALFRED task (§ 4.2), our method is, to the best of our knowledge, the first model that simultaneously addresses these two tasks. In table 5 and table 6, we find that the multitask model exhibits superior performance compared to the single-task models. These results underscore the capability of our method to effectively handle multiple highly diverse tasks.

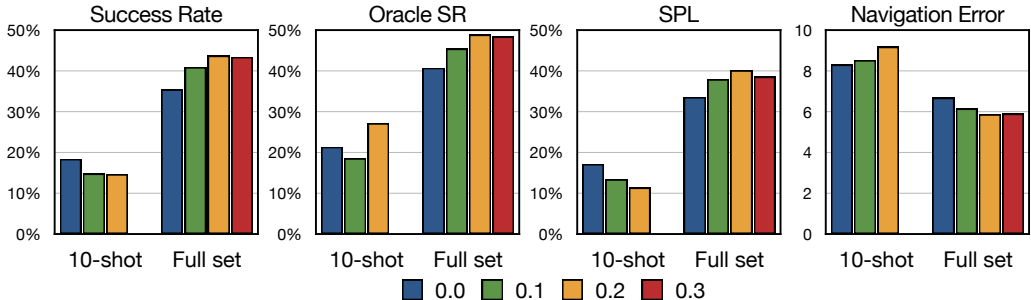


Figure 6: Investigating the Impact of the Randomness Factor  $\rho$  on Model Performance. This image caption depicts an ablation study exploring the influence of the randomness factor  $\rho$  on our model’s performance in both few-shot learning and full-set training scenarios. We test  $\rho$  with values of 0.0, 0.1, 0.2, and 0.3.

### D IMPACT OF THE RANDOMNESS FACTOR

We conduct the ablation study to investigate the impact of the randomness factor  $\rho$  on the model’s performance in both few-shot learning and full-set fine-tuning scenarios. Interestingly, we observe different behaviors of the model with varying  $\rho$  values in these scenarios. Fig. 6 illustrates our findings. In the 10-shot scenario, increasing  $\rho$  negatively affected the model’s performance. However, in the full-set scenario, the model achieved optimal performance at  $\rho = 0.2$ . We propose a metaphorical explanation: for effective few-shot learning, the language model behaves like an infant, relying on highly accurate demonstrations, while for large-scale imitation learning, the language model behaves like an adult, benefitting from occasional detours (introduction of randomness factor  $\rho$ ) to enhance robustness.

### E EXAMPLE OF THE FULL TRAJECTORY

We plot an example of a full text-based trajectory in R2R as we mentioned in § 3.2 as bellow:

You are a navigation agent who must navigate according to instructions given only descriptions of your current position via natural language. The natural language description is sometimes incorrect.

### Instruction:  
Go across the kitchen passed the pizza oven into the next room. Stop just inside the bedroom.

### Trajectory:  
Step 1:  
  
To your 60 degree left is,  
a kitchen with a stove, sink, and refrigerator  
Details: oven, bowl, dining table, microwave, refrigerator

To your 30 degree left is,  
a kitchen with white cabinets and a white stove top oven  
Details: oven, bowl, microwave, potted plant, chair

To your 60 degree right is,  
a living room filled with furniture and a fire place  
Details: couch, chair

To your back is,

a room with a bed and a painting on the wall

You chose:

a living room filled with furniture and a fire place

Step 2:

To your 150 degree left is,

a kitchen with a stove top oven sitting next to a counter  
Details: oven, dining table, microwave, refrigerator, chair

To your back is,

a kitchen with a wooden table and white appliances  
Details: oven, bowl, dining table, microwave, refrigerator, apple, couch, chair

To your 30 degree right is,

a living room with a couch a chair and a window  
Details: couch

To your 120 degree right is,

a dining room table with a bowl of fruit on it  
Details: chair, bowl, dining table

To your 150 degree right is,

a bowl of fruit sits on a wooden table  
Details: couch, chair, bowl, dining table

You chose:

a living room with a couch a chair and a window

Step 3:

To your back is,

a kitchen with a table, chairs, and stairs  
Details: oven, dining table, refrigerator, potted plant, chair

To your 150 degree right is,

a room with a table, chairs, and stairs  
Details: oven, chair, refrigerator, dining table

To your straight ahead and 30 degree down is,

a view of a hallway from the top of stairs  
Details: refrigerator

To your 90 degree right and 30 degree up is,

a room with a staircase and a mirror on the wall  
Details: toilet

You chose:

a view of a hallway from the top of stairs

Step 4:

To your back and 60 degree up is,

a living room filled with furniture and a ceiling fan  
Details: oven, potted plant, refrigerator

To your 90 degree left and 30 degree up is,

a living room with a couch and a table

To your straight ahead and 30 degree up is,  
a bedroom with a bed and a mirror on the wall  
Details: bed

You chose:  
a bedroom with a bed and a mirror on the wall

Step 5:

To your back is,  
a hallway leading to a kitchen and living room  
Details: refrigerator, potted plant

To your 30 degree left is,  
a room with a wooden door and a mirror

To your straight ahead is,  
a bedroom with a bed, dresser, mirror and a ceiling fan  
Details: potted plant, bed

To your 30 degree right is,  
a bedroom with a bed and a ceiling fan  
Details: potted plant, bed

To your 60 degree right is,  
a bedroom with a bed, dresser and mirror  
Details: potted plant, bed

You chose:  
stop

## F COMPLETE PROMPT TEMPLATE OF GENERATING TRAJECTORIES FOR GPT-4

We list our complete templates for prompting GPT-4 to generate synthetic instructions (Phase I) and synthetic trajectories to fulfill the instruction (Phase II).

**Phase I:** The template of phase I is listed as follows:

I am going to give you example instructions written by humans to train a deep learning-based navigation agent acting inside a home. These example instructions are intended to be completed by the navigation agent in 5-7 steps.

- {real\_instruction\_1}
- {real\_instruction\_2}
- {real\_instruction\_3}

Your goal is to write 10 more instructions like the above that can be used to train a navigation agent. Since the navigation agent will be navigating in different home environments, your instructions should also be diverse and cover a wide range of home environments and rooms. You should make sure that the instruction can be completed by an agent in 5 to 7 steps.

**Phase II:** The template of phase II is listed as follows:

Here is an example of a large language model acting as a blind navigation agent in an indoor environment through text descriptions. The agent is given an instruction at the start and must follow the instruction. At each time step, the agent is given descriptions of its field of view via the following template:

To your [VIEW] is [CAPTION]

- [VIEW] consists of the agent’s visible field of view (e.g., 30 degrees right, 120 degrees left, etc.)
- [CAPTION] is the text description of that view obtained from an image captioning model

#Example 1

### Instruction: {real\_instruction\_example}

### Trajectory: {real\_trajectory\_example}

Now I will give you another instruction. Please generate a trajectory of 5-7 steps that would complete the instruction.

#Example 2

### Instruction: {synthetic\_instruction}

## G PROMPTS OF ZERO-SHOT AND FEW-SHOT NAVIGATION FOR GPT-4

Here we attach the the task description  $D$  in the prompt template for prompting GPT-4 to navigate in the R2R evaluation dataset.

Zero-shot:

You are a navigation agent who must navigate according to instructions given only descriptions of your current position via natural language. The natural language description is sometimes incorrect.

At each step, you will be given several directions and captions for each direction. You must choose one direction by printing only the [caption\_of\_the\_direction] or choose "Stop" if you think the goal is reached.

For example:

Input:

To your [direction\_1] is, [caption of the direction\_1].

.....

To your [direction\_N] is, [caption of the direction\_N].

You choose:

Output: [caption of the direction\_3]

Hint: You should use the information inside the instructions, history steps, and current observations to make the decision.

Few-shot:

You are a navigation agent who must navigate according to instructions given only descriptions of your current position via natural language. The natural language description is sometimes incorrect.

At each step, you will be given several directions and captions for each direction. You must choose one direction by printing only the [caption\_of\_the\_direction] or choose "Stop" if you think the goal is reached.

For example:

Input:

To your [direction\_1] is, [caption of the direction\_1].

.....

To your [direction\_N] is, [caption of the direction\_N].

You choose:

Output: [caption of the direction\_3]

And here is an example trajectory:

### Instruction:

Go down the stairs. Turn right and go down the hallway. Turn right and stand near the fireplace.

### Trajectory:

Step 1:

To your straight ahead is,  
an ornate doorway leading to another room

To your 60 degree right is,  
a red carpeted staircase leading to a chandelier

To your 120 degree right is,  
a room with a red carpet and a large mirror

To your back and 30 degree down is,  
a room with a red carpet and two windows

To your 120 degree left is,  
a room with a red carpet and gold trim

You chose:  
a room with a red carpet and gold trim

Step 2:

To your 150 degree right is,  
a very ornate staircase in a house with red and white striped chairs

To your back is,  
a red carpeted hallway leading to a staircase

To your 150 degree left is,  
a hallway with a red carpet and a chandelier

To your 120 degree left is,  
a room with a red carpet and a chandelier

To your 90 degree left is,  
a room with a chandelier and two windows

To your 60 degree left is,  
a room with a red carpet and a large mirror

To your 30 degree right is,  
a hallway with a red carpet and wooden doors

You chose:  
a hallway with a red carpet and wooden doors

Step 3:

To your back is,  
a hallway with a red carpet and a chandelier

To your straight ahead is,  
a hallway with a red carpet and a gold ceiling  
a hallway with a red carpet and a gold ceiling

You chose:  
a hallway with a red carpet and a gold ceiling

Step 4:

To your 90 degree right is,  
a living room with a chandelier and a fireplace

To your 120 degree right is,  
a room with a fireplace and a chandelier in it

To your back is,  
a hallway with a red carpet and gold trim

To your 90 degree left is,  
a room with a chandelier and a table in it

To your 30 degree right is,  
a living room with a chandelier and a couch

You chose:  
a living room with a chandelier and a fireplace

Step 5:

To your straight ahead is,  
a living room filled with furniture and a fire place

To your 30 degree left is,  
a living room filled with furniture and a chandelier

To your 150 degree right is,  
a room with a chandelier and a table in it

To your back is,  
a hallway with a chandelier and a mirror on the wall

To your 90 degree left is,



a living room filled with furniture and a chandelier

You chose:

Stop

Hint: You should use the information inside the instructions, history steps, and current observations to make the decision.

Now let's start!

## H MORE EXAMPLES OF GENERATED TRAJECTORIES

We plot more examples of the generated trajectories from GPT-4 below:

Example 1:

### Instruction

Begin in the basement, walk towards the pool table. Pass the pool table and ascend the staircase. At the top of the stairs, turn right and walk through the hallway. Turn left into the den, walk through and stop by the window. ### Trajectory:

Step 1:

To your straight ahead is,  
a basement with a pool table, a bar, and a seating area

To your 90 degree left is,  
a wall with various framed photos and a door

To your 150 degree right is,  
a storage area filled with boxes and shelves

You chose:  
a basement with a pool table, a bar, and a seating area

Step 2:

To your back is,  
a basement entrance with a staircase leading down

To your 120 degree left is,  
a basement with a bar, television, and several chairs

To your straight ahead is,  
a basement with a pool table, a dartboard, and a bar

To your 90 degree right is,  
a staircase leading upwards with a wooden railing

You chose:  
a basement with a pool table, a dartboard, and a bar

Step 3:

To your back is,  
a basement with a seating area and framed photos on the wall

To your straight ahead is,  
a pool table surrounded by chairs and a dartboard on the wall

To your 150 degree left is,  
a basement with a bar, television, and several chairs

To your 90 degree right is,  
a staircase leading upwards with a wooden railing

You chose:  
a staircase leading upwards with a wooden railing

Step 4:

To your back is,  
a basement with a pool table and a dartboard on the wall

To your 90 degree left is,  
a hallway with various framed photos on the walls

To your straight ahead is,  
a den with a large bookshelf, couch, and coffee table

You chose:  
a hallway with various framed photos on the walls

Step 5:

To your back is,  
a staircase leading downwards into the basement

To your straight ahead is,  
a long hallway with several doors leading to other rooms

To your 120 degree right is,  
a den with a comfortable seating area and a bookshelf

You chose:  
a long hallway with several doors leading to other rooms

Step 6:

To your back is,  
a hallway with various framed photos on the walls

To your 90 degree left is,  
a den filled with cozy furniture, bookshelves, and a television

To your 90 degree right is,  
a bedroom with a large bed, dresser, and a window

You chose:  
a den filled with cozy furniture, bookshelves, and a television

Step 7:

To your back is,  
a hallway leading to other parts of the house

To your 90 degree left is,  
a den with several bookshelves, a fireplace, and a large rug

To your straight ahead is,  
a window overlooking the garden outside

You chose:  
Stop

Example 2:

### Instruction

From the garage entrance, walk into the kitchen. Turn left and walk past the stove. Continue straight into the family room, then walk towards the sliding glass doors. Step outside and stop near the barbecue grill.

### Trajectory:

Step 1:

To your straight ahead is,  
a clean garage with a parked car and a workbench

To your 90 degree left is,  
a laundry room with a washer and a dryer

To your back is,  
a driveway leading up to a garage door

To your 120 degree left is,  
a garage entrance with an open door leading inside the house

You chose:  
a garage entrance with an open door leading inside the house

Step 2:

To your 30 degree left is,  
a kitchen with a stove, a refrigerator, and a microwave oven

To your 90 degree right is,  
a small mudroom with a shoe rack and a coat hook

To your back is,  
a view of a garage with a car and a workbench

You chose:  
a kitchen with a stove, a refrigerator, and a microwave oven

Step 3:

To your straight ahead is,  
a large kitchen with a stove and a range hood

To your 30 degree right is,  
a kitchen with a center island and a sink on the countertop

To your 150 degree left is,

a kitchen with a refrigerator and a microwave oven

To your 120 degree right is,  
a tiled floor leading to a hallway with a picture on the wall

You chose:  
a large kitchen with a stove and a range hood

Step 4:

To your straight ahead is,  
a cozy family room with a TV and a fireplace

To your 90 degree right is,  
a kitchen with a sink and a dishwasher

To your back is,  
a kitchen with a countertop, a cooktop, and an oven

You chose:  
a cozy family room with a TV and a fireplace

Step 5:

To your 30 degree left is,  
a sliding glass door leading to the backyard

To your 90 degree right is,  
a comfortable sofa with a coffee table in front of a TV

To your back is,  
a large kitchen with a stove and a range hood

You chose:  
a sliding glass door leading to the backyard

Step 6:

To your straight ahead is,  
a well-maintained backyard with a barbecue grill and a dining table

To your 90 degree left is,  
a sliding glass door leading inside the family room

To your back is,  
a cozy family room with a TV and a fireplace

You chose:  
a well-maintained backyard with a barbecue grill and a dining table

Step 7:

To your 90 degree left is,  
a barbecue grill on a patio in the backyard

To your straight ahead is,  
a backyard with a dining table and chairs

To your 120 degree right is,

a backyard with a garden and a swing set

To your back is,  
a sliding glass door leading inside the family room

You chose:  
Stop

Example 3:

### Instruction

Stand in front of the fireplace in the living room. Walk diagonally across the room, passing the sofa. Proceed through the door into the hallway. Turn left and continue toward the guest room. Stop just before entering and wait.

### Trajectory:

Step 1:

To your straight ahead is,  
a living room with a fireplace and two chairs

To your 30 degree left is,  
a living room with a large window and a dark brown couch

To your 90 degree right is,  
a living room with a flat screen tv and a view of the backyard

You chose:  
a living room with a fireplace and two chairs

Step 2:

To your 60 degree left is,  
a living room with a fireplace and a floor lamp

To your 30 degree right is,  
a living room with a dark brown couch and a coffee table

To your straight ahead is,  
a living room with a white rug in the middle

You chose:  
a living room with a dark brown couch and a coffee table

Step 3:

To your back is,  
a living room with a fireplace and a christmas tree

To your 150 degree left is,  
a living room with two chairs and a painting on the wall

To your straight ahead is,  
a wooden door leading to a hallway

You chose:  
a wooden door leading to a hallway

Step 4:

To your 90 degree left is,  
a hallway with a view of a staircase and a table

To your straight ahead is,  
a hallway with a painting on the wall and an open door

You chose:  
a hallway with a painting on the wall and an open door

Step 5:

To your back is,  
a hallway with a wooden floor and a closed door

To your 120 degree left is,  
a guest bedroom with a neatly made bed and a dresser

To your 30 degree right is,  
a hallway with white walls and floor-to-ceiling mirrors

You chose:  
Stop just before entering the guest bedroom