
Situated Instruction Following Under Ambiguous Human Intent

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Language is never spoken in a vacuum. It is expressed and comprehended within the
2 holistic backdrop of the speaker’s history, actions, and environment. Since humans
3 are used to communicating efficiently with situated language, the practicality of
4 robotic assistants hinge on their ability to understand and act upon implicit and
5 situated instructions. In traditional instruction following paradigms, the agent
6 acts alone in an empty house, leading to language use that is both simplified and
7 artificially “complete.” In contrast, we propose *situated instruction following*
8 (SIF), which embraces the inherent underspecification and ambiguity of real-world
9 communication with the physical presence of a human speaker. The meaning of
10 situated instructions naturally unfold through the past actions and the expected
11 future behaviors of the human involved. Specifically, within our settings we
12 have instructions that (1) are ambiguously specified, (2) have temporally evolving
13 intent, (3) can be interpreted more precisely with the agent’s dynamic actions. Our
14 experiments indicate that state-of-the-art Embodied Instruction Following (EIF)
15 models lack holistic understanding of situated human intention.

16 1 Introduction

17 Humans communicate efficiently by providing only the necessary information, relying on shared
18 context like history, actions, and environment. For example, the request "Can you bring me a cup?"
19 varies based on context—if said near a kitchen sink with gloves, it likely refers to a dirty cup, while
20 near a bathroom sink, it suggests a clean one. Although clarification is possible, humans often
21 interpret such requests accurately using contextual cues, showing our ability to derive nuanced,
22 situation-specific meanings from ambiguous language.

23 As robotic agents increasingly become integral to our daily lives, their effectiveness and utility
24 critically depend on their ability to comprehend and respond to situated language— natural language
25 spoken by humans. Without this capability, agents may prove more of a hindrance than a help, forcing
26 users to perform tasks themselves rather than entrusting them to an assistant. As discussed in the field
27 of agent alignment [Leike et al., 2018], it is often difficult for users to precisely define or articulate
28 ideal task specifications. Consequently, an agent that demands detailed explanations might render
29 manual task execution by humans more attractive.

30 Current instruction-following tasks prioritize accurate low-level instruction interpretation [Anderson
31 et al., 2018, Gu et al., 2022, Padmakumar et al., 2021, Shridhar et al., 2020] or use commonsense
32 to achieve underspecified goals like object navigation [Chaplot et al., 2020, Das et al., 2018]. In
33 contrast, our work SIF aims to generalize *Embodied* Instruction Following to *Situated* Instruction
34 Following, with instructions closer to the language naturally spoken by humans. Specifically, we
35 focus on three dimensions of situated reasoning:

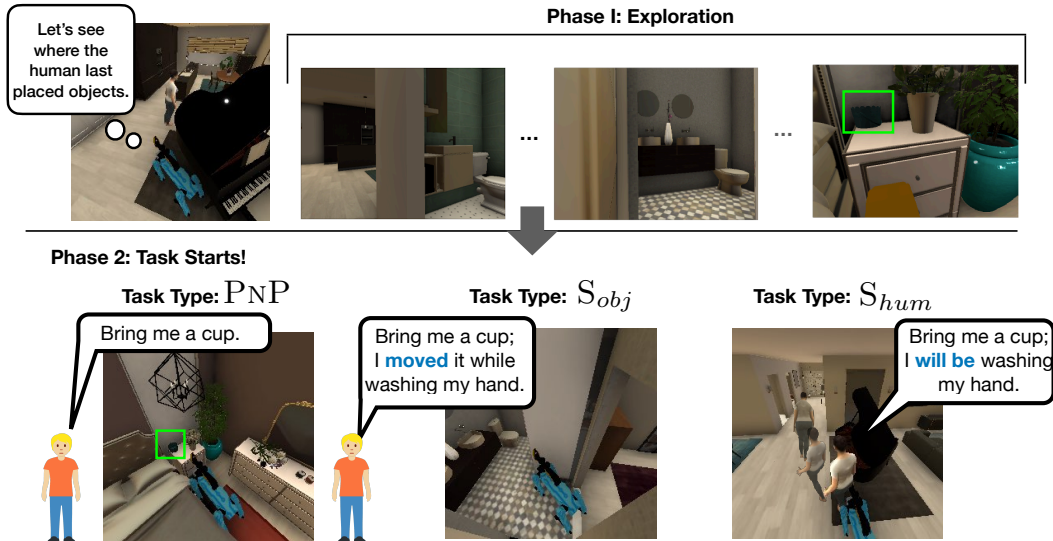


Figure 1: **Overview.** The tasks in SIF consist of two phases: an exploration phase (phase 1) and a task phase (phase 2). PNP represents a conventional static Pick-and-Place task used for comparison, wherein the environment remains unchanged after the exploration phase. S_{hum} and S_{obj} introduce two novel types of situated instruction following tasks. In these tasks, the *objects* and *human* subjects move during the task phase. Nuanced communication regarding these movements is provided, necessitating reasoning about ambiguous and temporally evolving human intent.

- 36 1. **Ambiguity:** As in the cup example above, there is ambiguity in the instruction given by the
37 speaker.
- 38 2. **Temporal:** A speaker’s actions change how their instruction should be interpreted (e.g.,
39 clarifying an underspecified reference).
- 40 3. **Dynamic:** When the environment changes, the agent needs to decide what actions will
41 reduce their uncertainty (e.g., following the human).

42 We implement our tasks in Habitat 3.0 [Puig et al., 2023], which includes simulated human agents.
43 To ensure fair comparison with prior work, we include both static (prior work) and dynamic (this
44 work) tasks (Fig. 1). The static task follows the classic pick-and-place paradigm where the agent is
45 instructed to Put [Obj] in/on [Recep]. We simplify the setup by allowing the agent to explore,
46 minimizing the role of mapping in our reasoning benchmark.

47 Our benchmark focuses on dynamic tasks where the agent must combine instruction understanding
48 with human movement. The dynamic tasks include S_{obj} (object moved by human) and S_{hum} (human
49 is the receptacle). In these, the agent receives goal instructions (e.g., “Bring me a mug” for S_{hum} or
50 “Put the mug in the bathroom” for S_{obj}) along with relocation hints. In S_{hum} , the human moves as
51 the task begins, signaling intent through both words and movement. The agent must efficiently follow
52 instructions, retrieve the object, and place it in the correct location (e.g., with the moving human in
53 S_{hum} tasks).

54 We specifically target evaluation of state-of-the-art Embodied Instruction Following (EIF) baselines.
55 We implement two such systems inspired by papers on related tasks. The first baseline, which we
56 refer to as REASONER, is a closed-loop system incorporating a semantic map, a prompt generator,
57 and a Large Language Model (LLM) planner. For the prompt generator, we integrated components
58 from Voyager[Wang et al., 2023], LLMPlanner[Song et al., 2023], and ReAct[Yao et al., 2022],
59 tailoring them to suit our dataset’s specific requirements. The second baseline, PROMPTER[Inoue and
60 Ohashi, 2022], was very successful at executing ALFRED [Shridhar et al., 2020] tasks despite being
61 open-loop. We see the desired result that our static scenarios match those from existing EIF datasets
62 [Inoue and Ohashi, 2022, Song et al., 2023], and these LLM based approaches perform very well in
63 tasks requiring common sense. However, their performance significantly declines when faced with
64 situations that require reasoning about the human’s behavior.

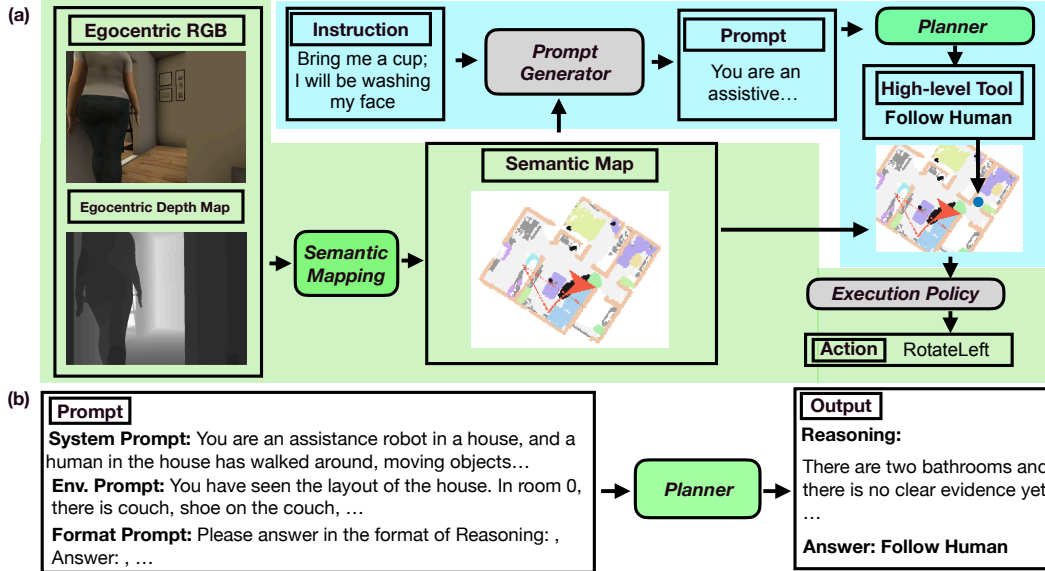


Figure 2: **REASONER**: (a) The semantic mapper is updated at every timestep, whereas the prompt generator and planner are activated either upon completion of the last high-level action or when a new decision is required. (b) The prompt consists of system prompt, environment prompt, format prompt.

65 2 Task

66 Our tasks (1) are structured into two distinct phases: (1) the exploration phase and (2) the task
 67 phase. During the exploration phase, the agent is allotted N steps to navigate around a static house
 68 environment where object assets are positioned. The value of N is determined to ensure the agent
 69 has sufficient steps to thoroughly scan the environment; specifically, $N = 1.5 \times$ (the number of
 70 steps required to achieve a complete map using frontier-based exploration techniques). Following
 71 the exploration phase, some objects are repositioned without the agent’s knowledge. As the task
 72 phase commences, the agent receives an instruction (e.g., “Bring me a cup,” “Put the cup in the
 73 sink”), accompanied by either direct or ambiguous information regarding which objects have been
 74 moved (e.g., “I took a cup with me. I’ll be getting ready for bed”). If the task involves delivering an
 75 object to a human, the human walks into the agent’s field of view as the task begins, simultaneously
 76 providing hints about their intended location (“I will be in the bathroom washing my face”). These
 77 elements, along with other strategic design decisions, ensure that the exploration phase effectively
 78 contextualizes the language directives, rendering tasks sufficiently solvable.

79 3 Baselines

80 Many recent state-of-the-art EIF agents are modular models with an LLM planner, connected to
 81 learned/engineered episodic memory, perception, and execution tools. We present a baseline within
 82 this high-performing family — REASONER, a closed-loop baseline that adapts FILM[Min et al.,
 83 2021] and the prompts of llm-planner[Song et al., 2023], and ReAct [Yao et al., 2022], and prompter
 84 [Inoue and Ohashi, 2022], an open-loop SOTA agent built for ALFRED [Shridhar et al., 2020].

85 **Semantic Mapper.** The semantic mapper creates a global representation for visual observation. As
 86 in previous work[Chaplot et al., 2020, Min et al., 2021], we process egocentric RGB and depth into
 87 an allocentric top-down map of obstacles and semantic categories using Detic[Zhou et al., 2022]. The
 88 semantic categories of interest are [ObjectCat], [Recep], and “human.” In contrast to previous
 89 works[Chaplot et al., 2020, Min et al., 2021], the most recent human and object positions are refreshed
 90 post new observations and pick/place actions, ensuring a dynamic and accurate representation of the
 91 environment.

92 **Text representation generator.** The semantic map and other contexts are converted into prompts.
 93 It is a concatenation of three components: the system prompt, environment prompt, and the format
 94 prompt:

Table 1: **SPL** performance of REASONER across splits. In each sectioned-row, the top row assumes oracle perception (semantic segmentation and manipulation); the bottom row assumes learned semantic segmentation and heuristic manipulation. To minimize the burden on API costs and time, we have limited LLM API calls for plan generation to 15 times.

Model		Val Seen			Val Unseen			Test Seen			Test Unseen		
Planning	Perception	PNP	S_{obj}	S_{hum}	PNP	S_{obj}	S_{hum}	PNP	S_{obj}	S_{hum}	PNP	S_{obj}	S_{hum}
Oracle	Oracle	98	100	95	100	100	100	98	93	98	95	100	98
	Learned	46	46	59	41	30	54	52	30	69	44	47	46
REASONER	Oracle	82	61	23	78	49	39	73	58	29	81	49	34
	Learned	21	8	12	24	11	12	29	2	15	18	14	15

- 95
- **System:** The system prompt outlines the agent’s role and encourages it to account for uncertainty. It is presented as “You are an assistive robot in a house, aiding a human. Your observations may be incomplete or wrong.”
- 96
- **Environment:** The environment prompt is a conversion of the episodic memory into text format, and contains information of the agent’s current state and previously completed/failed actions. It is given in the following sequence: (1) observation of P_e during exploration phase, based on the semantic map, (2) C , regarding object/ human movements, (3) the goal instruction I , (4) the high-level action executed by agents at timesteps and their observed consequences (success/fail), (5) the agent’s latest observation, based on the semantic map.
- 98
- **Format:** The format prompt explains action affordance and a format for chain of thought [Wei et al., 2022]. It also explains the desired effect of actions (e.g. “If you want to keep searching for object(s) or human that might exist (but you have not detected) in the current room, choose ‘Explore Room X ’ (Table 2).”)
- 99
- 100
- 101
- 102
- 103
- 104
- 105
- 106
- 107

108 **Execution Tools** Upon receiving the prompt, the planner is prompted to choose a high-level action (Tab. 2); then corresponding execution tools are called. A complete list of tools are listed in Table 2. When the execution is done, the tool sends this message, and the prompt generator creates a new prompt and the planner calls a new tool.

109

110

111

112 4 Results

113 Results from our experiments are presented in Table 1. This table notably shows the following facts about our dataset and baselines. First, the gap of model performance across PNP versus S_{hum} , S_{obj} shows that PNP can be solved with commonsense and mechanistic combination, and the rest two tasks cannot. The reasoning challenges of S_{obj} and S_{hum} are backed by the performance of REASONER with oracle perception/manipulation; it shows a stark contrast in PNP tasks ($\sim 80\%$) and S_{obj} , S_{hum} tasks ($\sim 45\%$).

114

115

116

117

118

119 Table 3 examines model performance on clear versus ambiguous tasks. Ambiguity in S_{obj} tasks emerges when multiple potential locations exist for an object, as indicated by communicative cues. For example, the statement “I am washing my face” becomes ambiguous when multiple bathrooms are available. Similarly, ambiguous S_{hum} tasks occur when the human could be in several different locations. In S_{hum} tasks, REASONER underperforms in clear tasks due to a tendency to conservatively judge that there is insufficient evidence of the human’s destination, even when only one plausible location exists. REASONER attempts some calibration but generally leans towards following the human. Qualitative analysis reveals that in ambiguous tasks, REASONER often disengages prematurely, assuming it has accumulated enough evidence.

120

121

122

123

124

125

126

127

128 5 Conclusion

129 We present Situated Instruction Following (SIF), a new dataset to evaluate situated and holistic understanding of language instructions. Our dataset reflects aspects of real-world instruction following: (1) ambiguous task specification, (2) evolving intent over time, and (3) dynamic interpretation influenced by agent action. We show that current state-of-the-art models struggle with this level of understanding, further highlighting the complexity and uniqueness of our dataset.

130

131

132

133

134 References

- 135 Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid,
136 Stephen Gould, and Anton Van Den Hengel. Vision-and-language navigation: Interpreting visually-
137 grounded navigation instructions in real environments. In *Proceedings of the IEEE conference on*
138 *computer vision and pattern recognition*, pages 3674–3683, 2018.
- 139 Devendra Singh Chaplot, Dhiraj Prakashchand Gandhi, Abhinav Gupta, and Russ R Salakhutdinov.
140 Object goal navigation using goal-oriented semantic exploration. *Advances in Neural Information*
141 *Processing Systems*, 33, 2020.
- 142 Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied
143 question answering. In *Proceedings of the IEEE conference on computer vision and pattern*
144 *recognition*, pages 1–10, 2018.
- 145 Jing Gu, Eliana Stefani, Qi Wu, Jesse Thomason, and Xin Eric Wang. Vision-and-language navigation:
146 A survey of tasks, methods, and future directions. *arXiv preprint arXiv:2203.12667*, 2022.
- 147 Yuki Inoue and Hiroki Ohashi. Prompter: Utilizing large language model prompting for a data
148 efficient embodied instruction following. *arXiv preprint arXiv:2211.03267*, 2022.
- 149 Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. Scalable agent
150 alignment via reward modeling: a research direction. *arXiv preprint arXiv:1811.07871*, 2018.
- 151 So Yeon Min, Devendra Singh Chaplot, Pradeep Ravikumar, Yonatan Bisk, and Ruslan Salakhutdinov.
152 Film: Following instructions in language with modular methods. *arXiv preprint arXiv:2110.07342*,
153 2021.
- 154 Aishwarya Padmakumar, Jesse Thomason, Ayush Shrivastava, Patrick Lange, Anjali Narayan-Chen,
155 Spandana Gella, Robinson Piramuthu, Gokhan Tur, and Dilek Hakkani-Tur. Teach: Task-driven
156 embodied agents that chat, 2021. URL <https://arxiv.org/abs/2110.00534>.
- 157 Xavier Puig, Eric Undersander, Andrew Szot, Mikael Dallaire Cote, Tsung-Yen Yang, Ruslan Partsey,
158 Ruta Desai, Alexander William Clegg, Michal Hlavac, So Yeon Min, et al. Habitat 3.0: A co-habitat
159 for humans, avatars and robots. *arXiv preprint arXiv:2310.13724*, 2023.
- 160 Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi,
161 Luke Zettlemoyer, and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions
162 for everyday tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
163 *recognition*, pages 10740–10749, 2020.
- 164 Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su.
165 Llm-planner: Few-shot grounded planning for embodied agents with large language models. In
166 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2998–3009,
167 2023.
- 168 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and
169 Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv*
170 *preprint arXiv:2305.16291*, 2023.
- 171 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
172 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
173 *Neural Information Processing Systems*, 35:24824–24837, 2022.
- 174 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
175 React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*,
176 2022.
- 177 Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-
178 thousand classes using image-level supervision. In *European Conference on Computer Vision*,
179 pages 350–368. Springer, 2022.

180 **A Execution Tools**

181 Execution tools for REASONER/PROMPTER and their working details/affordance are in Table 2.

Execution Tool	Description & Affordance
Navigation	
Go to Room X	FMM Planner navigates to a random point in Room X.
Explore Room X	FMM Planner navigates to a random point in Room X; then, agent turns 15 times to the right to look around.
Follow Human	The last observed position of the human is given as the goal, to the human-following wrapper (more explanation is Sec. ??) on top of FMM Planner.
Manipulation	
Grab Obj	The closest object within 2 meters of the grasper is grabbed, and agent’s grasper is closed.
Put Obj	Grasped object is put on the closest receptacle within 2 meters of the grasper is grabbed, and agent’s grasper is opened.
Give Obj to Human	The agent goes within 1 meter of the human and gives grasped object to human, if human is visible from current view.

Table 2: Execution tools for REASONER/PROMPTER and their working details/affordance.

182 **B Ablations Across Ambiguous/ Clear Tasks**

Table 3: **Ambiguous vs Clear tasks.** SPL and SR of REASONER and PROMPTER with G.T./learned vision and manipulation on Val seen & unseen combined.

Model	Metric	G.T. Vis. & Man.		Learned Vis. & Man.					
		S_{obj}	S_{hum}	S_{obj}	S_{hum}				
		Clear Amb.	Clear Amb.	Clear Amb.	Clear Amb.				
REASONER	SPL	62	52	13	42	9	11	3	17
	SR	76	71	14	67	15	14	6	26
PROMPTER	SPL	38	29	3	42	11	8	0	17
	SR	54	36	4	66	18	10	0	27