

RILA: Reflective and Imaginative Language Agent for Zero-Shot Semantic Audio-Visual Navigation

Anonymous CVPR submission

Paper ID 987

Abstract

We leverage Large Language Models (LLM) for zero-shot Semantic Audio Visual Navigation (SAVN). Existing methods utilize extensive training demonstrations for reinforcement learning, yet achieve relatively low success rates and lack generalizability. The intermittent nature of auditory signals further poses additional obstacles to inferring the goal information. To address this challenge, we present the **Reflective and Imaginative Language Agent (RILA)**. By employing multi-modal models to process sensory data, we instruct an LLM-based planner to actively explore the environment. During the exploration, our agent adaptively evaluates and dismisses inaccurate perceptual descriptions. Additionally, we introduce an auxiliary LLM-based assistant to enhance global environmental comprehension by mapping room layouts and providing strategic insights. Through comprehensive experiments and analysis, we show that our method outperforms relevant baselines without training demonstrations from the environment and complementary semantic information¹.

1. Introduction

Intelligent agents are anticipated to navigate intricate environments, leveraging both auditory and visual stimuli [29, 36]. Considering a scenario that a vase falls and breaks, a robot must swiftly pinpoint a target within a room, relying primarily on transient auditory cues. This need underpins our focus on the Semantic Audio-Visual Navigation (SAVN) task [9]. In SAVN, the target object within the scene emits intermittent sounds, which the agent must use, in conjunction with visual information, to find the object. In addition to the ambiguous goal information conveyed through sporadic sounds, intricate room layouts and complex navigation trajectories also present significant challenges [42], rendering the SAVN task notably difficult. Previous research [9] concentrated on the end-to-end train-

¹<https://rila-savn.github.io/RILA>

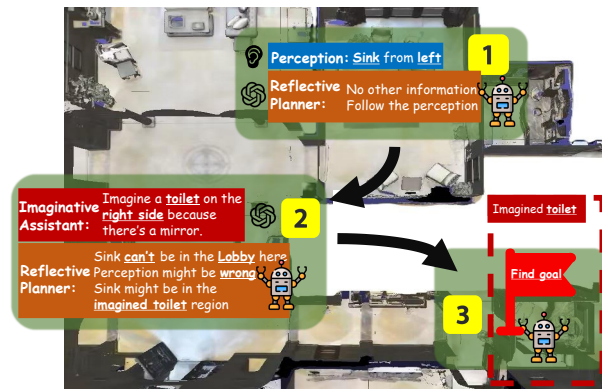


Figure 1. An illustration of our agent’s strategy for semantic audio-visual navigation. The **Reflective Planner** initiates navigation by relying on perceptual information for exploration. When exploration leads to an incorrect region, it subsequently discounts the perceptual descriptions, redirecting its focus. Throughout this process, the **Imaginative Assistant** persistently contributes spatial insights and suggestions, thereby assisting in reasoning.

ing of reinforcement learning models, yielding inadequate performance despite the use of extensive training trajectories. Recent approaches enhance performance by integrating auxiliary modules [42] or employing oracle instructions [29, 36], which may not be feasible in real-world applications.

Large language models (LLMs) [33, 34] have shown remarkable progress [28, 45]. Beyond the promising performance on natural language tasks [35, 38], the integration of LLMs into embodied robotics applications has also resulted in substantial improvements [2, 13, 14, 46, 48]. Recent methods [53, 54] equip LLMs with multi-modal models [26, 27] that provide perception and feedback from the environment, either explicitly [47, 49] or implicitly [18, 22], in vision-and-language navigation tasks [4]. However, these applications also fail on SAVN due to their reliance on precise perception information and explicit goal descriptions. Consequently, realizing zero-shot SAVN, as anticipated for

053 intelligent agents, remains a formidable challenge.

054 Therefore, we propose our *Reflective and Imaginative*
055 *Language Agent* (RILA), leveraging the inherent common-
056 sense reasoning capabilities of LLMs to perform zero-shot
057 SAVN. Practically, we design distinct perception models
058 that process audio and visual signals, which further guide
059 a frozen LLM in strategic planning. Through active ex-
060 ploration of the environment, our agent adaptively identi-
061 fies and deprioritizes misleading goal descriptions. Fur-
062 thermore, we introduce an LLM-based imaginative assis-
063 tant, which extracts room layouts and provides high-level
064 guidance. Incorporating this assistant enables our agent
065 to achieve comprehensive environmental understanding and
066 navigate toward the target object in a zero-shot manner.
067 Fig. 1 provides an illustration of our agent’s navigation.

068 To validate our approach, we conduct experiments
069 within the SoundSpaces framework. Experimental results
070 show that our method surpasses relevant baselines without
071 reliance on training demonstrations or complementary mod-
072 ules. Notably, our agent exhibits a success rate exceeding
073 60% when paired with oracle perceptions, highlighting the
074 strong planning capability of LLMs. Additionally, we con-
075 duct a thorough analysis of the bottleneck of the current task
076 configuration. We summarize our contributions as follows:

- 077 • We propose RILA for zero-shot SAVN, exploiting the
078 commonsense reasoning capabilities of LLMs to navigate
079 effectively without precise goal descriptions.
- 080 • We introduce an imaginative assistant, designed to deduce
081 the environment’s room layout and provide comprehen-
082 sive suggestions, thereby enhancing the navigation.
- 083 • Experiments substantiate that RILA surpasses previous
084 baselines, which require training, in a zero-shot manner.
085 We also conduct a thorough analysis of the SAVN task.

086 2. Related Work

087 2.1. Semantic Audio Visual Navigation

088 Semantic audio-visual navigation is defined in Habitat [31,
089 39] with the SoundSpaces dataset [8, 11]. Previous re-
090 search [7, 50] extract features from RGB-D images and
091 two-channel spectrograms using pre-trained encoders sep-
092 arately [3, 8], and then train an end-to-end policy network
093 by reinforcement learning to predict the next action. How-
094 ever, these methods lack generalizability, failing in unsu-
095 pervised scenes [42] despite necessitating extensive training
096 demonstrations. Recent methods [29, 36] query for human
097 instructions during the navigation. K-SAVEN [42] further
098 constructs a knowledge graph to provide spatial comprehen-
099 sion. Instead of training on massive demonstrations, our
100 method exploits the commonsense reasoning capabilities of
101 LLMs to perform solve the task in a zero-shot manner.

2.2. Navigation with Large Language Models

LLMs have recently demonstrated impressive reasoning
abilities across a range of tasks [19, 37], including embod-
ied tasks [15]. Recent studies [41, 52] investigate visual-
language navigation with LLMs. For instance, ESC [54]
employs LLMs to deduce relationships between objects,
thereby aiding navigation. [12, 40], on the other hand, uti-
lize visual foundation models to convert perceptions into
natural language instructions. However, the application of
LLMs in SAVN remains underexplored, especially since
prior methods often rely on ground-truth goal descriptions.
In contrast, RILA reflectively navigates toward the target,
handling potentially misleading goal descriptions.

2.3. Layout Complementary

Spatial understanding, particularly regarding room lay-
out, is crucial for comprehending complex environments.
LGD[25] employs a room-type codebook to conceptual-
ize room layouts from image clips. Text2Room [21], con-
versely, creates entire rooms guided by textual instructions.
Recent LayoutGPT [16] taps into the visual planning ca-
pabilities of LLMs to produce plausible layouts for visual
generation. In our work, RILA utilizes LLMs to progres-
sively deduce the room layout and type, thereby achieving
a global understanding of the environment.

3. Method

In this work, we consider solving the Semantic Audio Vi-
sual Navigation (SAVN) task [9] in a zero-shot manner,
challenging agents to locate the sounding object within an
intricate and unseen environment. Notably, the audio sig-
nals here are sporadic and often absent, posing a signifi-
cant challenge to the agent’s decision-making process. In-
stead of training on trajectories from the simulated environ-
ment or incorporating additional semantic information, we
leverage the intrinsic commonsense reasoning capabilities
of LLMs for navigation planning.

3.1. Overview

In this section, we provide an overview of our RILA frame-
work, illustrated in Fig. 2. RILA consists of three parts:
the perception module, the Imaginative Assistant, and the
Reflective Planner, which we will introduce separately.

The perception module transforms the sensory data into
natural language descriptions. Visual perceptions o_i^v are di-
rectly processed via a pre-trained visual-language model,
which discerns and catalogs the observed objects, thereby
facilitating the construction of a semantic top-down map.
We develop distinct modules for auditory perceptions o_i^a to
pinpoint the goal location and identify pertinent semantic
cues, given the intermittent nature. Both perceptions are

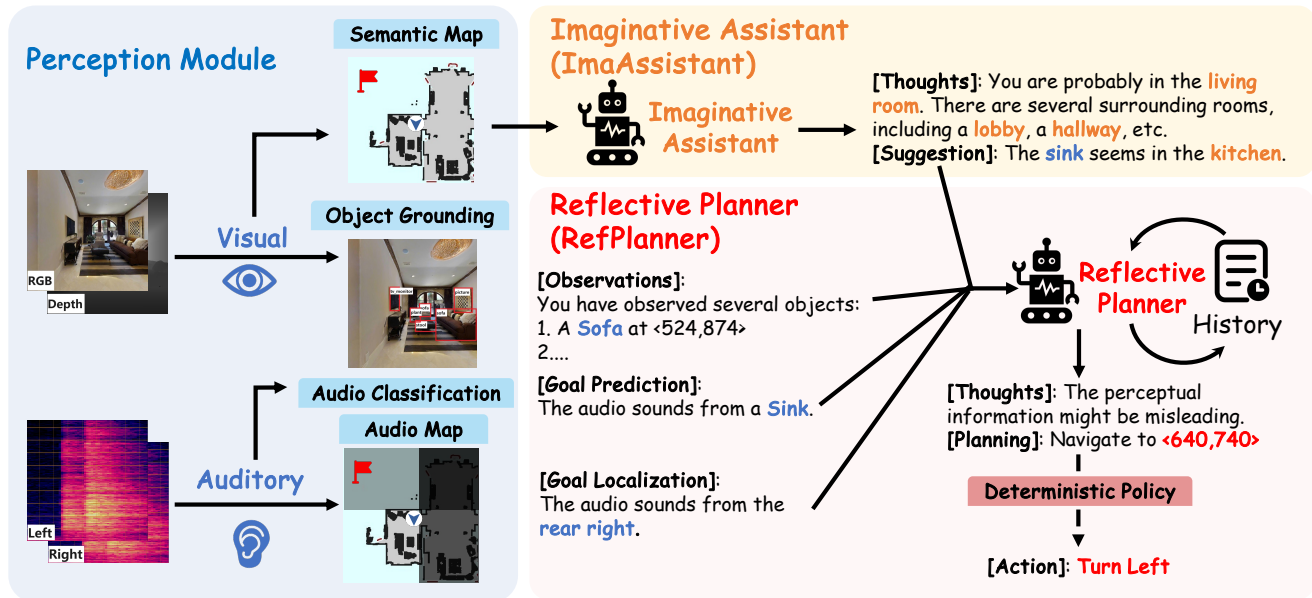


Figure 2. The architecture of our agent comprises three primary components. Firstly, the perception module transforms sensory inputs into text-based descriptions. Secondly, the Imaginative Assistant analyzes regional information and offers strategic guidance from a global perspective. Lastly, by integrating the two components, our Reflective Planner assesses perceptual data and navigates toward the target.

150 then synthesized into a text-based format for planning. A
151 detailed description is illustrated in Section 3.2.

152 Extending beyond individual objects, we integrate an
153 LLM-based Imaginative Assistant (ImaAssistant) to de-
154 deduce room layouts, thereby enriching the spatial compre-
155 hension of intricate environments. ImaAssistant then uti-
156 lizes the layout of both explored and partially observed ar-
157 eas to provide strategic planning guidance, aiding in navi-
158 gation. A thorough explanation is provided in Section 3.3.
159 By amalgamating insights from the perception module and
160 ImaAssistant, our Reflective Planner (RefPlanner) lever-
161 ages inherent commonsense reasoning abilities to explore
162 the environment and identify misleading auditory descrip-
163 tions, circumventing the need for exact sound localiza-
164 tion. Detailed explanations are shown in Section 3.4.

165 3.2. Audio Visual Perception

166 Following [42], we use the same pre-trained audio classifi-
167 cation model M_{obj}^a to infer the target object. Considering the
168 transient nature of audio signals, which presents a consid-
169 erable obstacle in precise identification, we employ a pro-
170 gressive strategy. Upon an audio signal o_t^a at time step t ,
171 we make a prediction $M_{\text{obj}}^a(\{o_{1,\dots,t}^a\})$ by amalgamating the
172 current audio with the accumulative history, facilitating a
173 refined accuracy. The object \hat{g}_t with the highest cumulative
174 prediction score at time t is thus designated as the current

goal object:

$$175 \hat{g}_t = \underset{g}{\operatorname{argmax}} \left(\sum_{i=1}^t \mathbb{1}_{M_{\text{obj}}^a(\{o_{1,\dots,i}^a\})=g} \right), \quad (1) \quad 176$$

177 where $\mathbb{1}$ denotes the indicator function. Guided by the pre-
178 diction \hat{g}_t , we aim to further localize it, thereby improving
179 the distinction of the target from analogous entities in the
180 environment. Nonetheless, the complex reverberation of the
181 simulation poses a significant challenge for localization, as
182 evidenced by an error margin of about 8 meters [9].

183 Therefore, we partition the localization into independ-
184 ent estimations of distance and direction. To quantify
185 sound distance, we collected 10,000 unheard auditory sam-
186 ples from the training environment to delineate the simu-
187 lation’s dimensional attributes. A pre-trained ResNet-18
188 model fine-tuned on this dataset demonstrates commend-
189 able accuracy in estimating distances. Predicting direction,
190 however, is substantially more arduous.

191 Instead of ascertaining the precise angle, we shift to
192 identify the binary directionality, greatly simplified by the
193 dual-sensor configuration. Nonetheless, techniques such
194 as Interaural Time Difference (ITD) [6, 20] and fine-tuned
195 models fall short of the task, which is further discussed in
196 Section 5. Consequently, we employ weighted predictions
197 based on the Root Mean Square (RMS) intensity of audi-
198 tory signals from the dual channels, denoted by R_l^t and R_r^t .
199 Practically, we consider the audio source to be from the side
200 with the larger RMS intensity. For each point p and time t ,

201 the confidence C_p^t is calculated as:

$$202 \quad C_p^t = \sum_{i=1}^t w_i^a \cdot \mathbb{1}_{RMS}(p, o_i^a), \quad (2)$$

203 where $\mathbb{1}_{RMS}(p, o)$ is an indicator function which is equal
204 to 1 if p is located, with respect to the agent, in the side cor-
205 responding to the larger RMS intensity given observation
206 o , and the weight w_t^a is calculated as $w_t^a = \frac{|R_i^t - R_r^t|}{\max(R_i^t, R_r^t)}$.
207 Through iteratively accumulating the weighted predictions,
208 we construct an audio map that facilitates an approximate
209 localization of the goal.

210 To transform visual signals into linguistic representa-
211 tions, we employ the pre-trained GroundingDINO for both
212 delineating bounding boxes and identifying the objects
213 within the RGB observation, thereby furnishing a rudimen-
214 tary environmental understanding. Besides, we separately
215 prompt to detect the predicted goal object in case the tar-
216 get is missed. Simultaneously, a semantic top-down map is
217 constructed from the Depth observations, with the map seg-
218 mented into distinct regions demarcated by detected walls,
219 enabling the assistant to provide a region-level comprehen-
220 sion. A more detailed illustration of our perception modules
221 is further provided in Appendix.

222 3.3. Imaginative Assistant

223 Given the restricted information from the perception mod-
224 ule, the planning relies mainly on discrete objects. How-
225 ever, a global environmental understanding substantially
226 benefits planning, especially for distant goals requiring
227 multi-room navigation. To address this, we integrate an
228 auxiliary LLM-based Imaginative Assistant (ImaAsssis-
229 tant), offering strategic suggestions to bolster navigation.

230 In practice, ImaAssistant infers room layouts. By par-
231 titioning the semantic map into regions using the detected
232 walls, we instruct ImaAssistant to determine closed room
233 types from observed objects. Yet, as a comprehensive ex-
234 ploration of a room rarely occurs, partially observed rooms
235 are more frequently encountered. Therefore, we utilize the
236 spatial imagination capabilities of LLMs to conceptualize
237 the layout of these rooms, subsequently directing it to de-
238 duce room types by interior objects and adjacent rooms. We
239 present below simplified versions of the prompts.

/* Task Description */

Please infer the room type and precise layout of the
provided interested region.

/* Room Layouts */

Observed Rooms: living room, etc.

Partially Observed Room: wall₁, wall₂, etc.

Internal Objects: chair₁, chair₂, table, etc.

240

241 Through iterative deduction of both observed and par-
242 tially observed rooms, RILA attains a comprehensive un-
243 derstanding of the environment, which yields additional in-
244 sights beyond the scope of individual objects. To augment
245 the planning, ImaAssistant is further instructed to provide
246 strategic navigation advice. Rather than specific waypoints,
247 ImaAssistant reasons about the potential goal locations,
248 considering spatial layouts and semantic attributes. These
249 insights enable ImaAssistant to make suggestions that as-
250 sist in selecting waypoints more effectively. A simplified
251 version of the prompt template is presented below.

/* Task Description */

Given the room layout, infer where the **Counter** is.

Give your advice about which room to explore.

/* Information */

Current room: living room

Surrounding rooms: kitchen, hallway, etc.

252

253 3.4. Reflective Planner

254 By incorporating layouts and suggestions from ImaAsssis-
255 tant, our LLM-based Reflective Planner (RefPlanner) har-
256 nesses the inherent commonsense reasoning capabilities in
257 planning based on perceptions. At each time step t , au-
258 dio and visual perceptions are formatted as *Goal Descrip-*
259 *tion* and *Observation*, respectively. Additionally, a *Task De-*
260 *scription* is articulated at the outset. A simplified template
261 for the perception prompt is as follows:

/* Task Description */

You are performing a navigation task.

/* Goal Description */

Navigate to the object that sounds like a **Counter**.

/* Observations */

You have observed the following objects.

262

263 With a natural language synopsis of the environment and
264 the designated navigational objective, we commission Ref-
265 Planner to strategize high-level planning. Rather than spec-
266 ifying actions outright, we implement a heuristic method,
267 frontier-based exploration (FBE), which discerns the junc-
268 tures between explored and uncharted territories as poten-
269 tial waypoints for environmental reconnaissance. Instead of
270 determining specific action, RefPlanner is directed to rea-
271 son and select an exploration frontier based on current per-
272 ceptions in a zero-shot manner. The navigation history of
273 perceptions and reasonings is also provided. Practically, we
274 implement a deterministic policy for decomposing the way-
275 point into action sequences. Utilizing a connected graph de-
276 rived from the semantic top-down map, we apply Dijkstra's
277 algorithm to determine the shortest path to the waypoint.

278 Moreover, as outlined in Section 3.2, the perception de-

279 scriptions, particularly the goal location, are often ambigu-
 280 ous and may lead to misconceptions, while an intelligent
 281 agent is anticipated to actively interact with the environment
 282 to make judgments about uncertain perceptions. Therefore,
 283 along with the localization confidence of the frontier from
 284 the perception module, we hint to RefPlanner about the po-
 285 tential inaccuracy, which empowers it to explore the en-
 286 vironment adaptively and reflect the reliability of percep-
 287 tion, thus enhancing its proficiency in locating the target
 288 object. The layouts and suggestions from ImaAssistant are
 289 included as well. We present below a simplified version of
 290 the template used for the navigation prompt.

```

/* Agent Position */
You are at  $\langle x, y \rangle$ 
/* Hint */
The perceptual confidence is not always accurate.
/* Frontier Candidates */
Frontier 1:  $\langle x, y \rangle$  in the living room
Perceptual confidence:  $c$ 
Surrounding objects: chair1, chair2, table, etc.
/* Suggestions */
The goal object may be in the kitchen.
  
```

291
 292 As shown in Fig. 1, RefPlanner adaptively selects appro-
 293 priate waypoints from a global perspective. When RefPlan-
 294 ner fails to find the target after exploring an area based on
 295 perceptions, it identifies perceptual inaccuracies and navi-
 296 gates using object characteristics. The full prompt scheme
 297 and a detailed example of the navigation are provided in
 298 Appendix. In practice, we implement all LLMs using the
 299 March 2023 version of gpt-3.5-turbo, leveraging the
 300 OpenAI LLM API service² with a temperature of 0.0.

301 4. Experiments

302 4.1. Experimental Setup

303 4.1.1 Datasets

304 We use SoundSpaces [8, 11] from Habitat [31, 39] envi-
 305 ronment to simulate navigation in 3D environments. We
 306 adopt the Matterport3D (MP3D) dataset for its ground-truth
 307 region layout labels and object labels. In particular, we
 308 evaluate our RILA on 1,000 test episodes within 10 unseen
 309 scenes with unheard sounds from 21 goal objects.

310 4.1.2 Baselines

311 We compare our model with several baselines:

- 312 • **AudioGoal** [8] uses a GRU state encoder to acquire the
 313 following action with an end-to-end RL policy network.

²<https://platform.openai.com/docs/models>

- **AV-WAN** [10] designs a waypoint predictor and leverages 314
 a local path planner to navigate to the waypoint. 315
 - **SAVi** [9] incorporates a goal descriptor network to predict 316
 both the classification and location of the sounding object. 317
 - **AVLEN** [36] adopts a hierarchical RL policy with goal 318
 predictor and memory unit, and queries oracle instruc- 319
 tions from humans if necessary. 320
 - **K-SAVEN** [42] proposes an end-to-end policy network 321
 with a knowledge graph constructed on the training data, 322
 presenting the relationship between regions and objects. 323
- In addition, we incorporate two zero-shot methods based 324
 on foundation models to facilitate a more comprehensive 325
 comparison. The ground truth goal object is provided here. 326
- **ImageBind-LLM** [18] is a novel multi-modality 327
 model that aggregates ImageBind [17] and LLaMA- 328
 Adapter [51] and we use the perfect stop strategy. 329
 - **ESC** [54] leverages LLMs and Probabilistic Soft Logic 330
 (PSL) [5] to choose a frontier for a visual-language navi- 331
 gation task. We provide our audio goal description. 332

4.1.3 Metrics 333

Following previous work [7, 36, 42], we report agent per- 334
 formance with the following metrics: Success Rate (SR), 335
 Success Rate weighted by Path Length (SPL), Success Rate 336
 weighted by Number of Actions (SNA), and Success When 337
 Silent (SWS), all in percentage (%). We also report the av- 338
 erage Distance To Goal (DTG) in meters at episode end. 339

4.1.4 Implementation Details 340

Consistent with previous studies, the agent is provided with 341
 RGB and depth images at a resolution of 256×256 . It 342
 also receives two-channel audio clips in the form of $65 \times$ 343
 26 spectrograms. The action space includes *MoveForward*, 344
TurnRight, *TurnLeft*, and *Stop*, with a movement step set at 345
 1 meter. Additionally, the agent obtains its GPS location 346
 at each time step. Detailed implementation details are pro- 347
 vided in Appendix. 348

4.2. Experimental Results 349

The comparative results are presented in Table 1. We derive 350
 the results of major baselines from their respective papers. 351
 For ESC and ImageBind-LLM, we incorporate ground-truth 352
 audio descriptions for the SAVN task. Implementation de- 353
 tails are provided in the Supplementary Material. Accord- 354
 ing to Table 1, our agent surpasses baselines that utilize 355
 end-to-end reinforcement learning training, such as SAVi, 356
 in a zero-shot manner. Even when juxtaposed with base- 357
 lines that utilize additional information, RILA achieves a 358
 higher success rate. Besides, we notice that Imagebind- 359
 LLM fails on the SAVN task, despite incorporating ground- 360
 truth audio descriptions, reflecting the limited performance 361

	Method	SR (%) \uparrow	SPL (%) \uparrow	SNA (%) \uparrow	DTG (m) \downarrow	SWS (%) \uparrow
Supervised	AudioGoal [8]	16.5	15.5	10.4	12.8	5.6
	AV-WAN [10]	17.2	13.2	12.7	11.0	6.9
	SAVi [9]	24.8	17.2	13.2	9.9	14.7
	AVLEN [36]	26.2	17.6	14.2	9.2	15.8
	K-SAVEN [42]	34.4	23.4	21.7	6.6	14.3
Unsupervised	Imagebind-LLM [†] [18] + Audio*	2.4	1.5	1.1	22.6	1.4
	ESC [54] + Audio*	23.6	8.0	4.8	17.7	14.2
	Ours w/o Assistant	31.4	9.6	6.8	12.2	15.3
	Ours	35.4	11.8	8.7	11.4	20.4

Table 1. Comparison with relevant baselines on SoundSpaces **Matterport3D** test dataset. AVLEN incorporates extra oracle instructions. \dagger denotes the perfect stop strategy and *Audio** indicates that the ground-truth audio description is provided. In contrast, our method requires no training trajectories or additional semantic information.

Method	SR (%) \uparrow	SPL (%) \uparrow	SWS (%) \uparrow
Random [†]	19.8	11.8	16.2
Nearest [†]	9.8	22.6	6.4
Llama-2 7B	39.4	22.2	35.4
Ours	60.8	39.6	56.6

Table 2. Ablation study on RefPlanner by replacing it with heuristic frontier selection methods and replacing the ChatGPT with Llama-2. \dagger indicates using oracle stop.

of open-source multi-modality foundation models on complex embodied tasks. Notably, our approach significantly outperforms previous works in terms of SWS, with over 40% improvement over K-SAVEN. This underscores the exceptional efficacy of our method in scenarios involving long distances and intermittent sounds, thereby highlighting the potential of harnessing the commonsense reasoning abilities of LLMs for navigation in physical environments.

We observe a relatively lower SPL of our method, attributed to the fact that RILA requires holistic exploration of the environment to ascertain the target object due to the absence of end-to-end training. Additionally, given the vague nature of the goal descriptions, RILA adopts a more cautious strategy for navigation, often traversing longer distances before reaching the objective. For better illustration, we provide two cases of snapshots of the navigation process using RILA in Fig. 3. As demonstrated in the left case study, RILA initially explores the living room, guided by erroneous perceptual cues. Upon realizing the absence of the goal object, RILA shifts its navigation toward the bathroom, utilizing object characteristics to locate the toilet. This process highlights RILA’s ability to effectively reflect on potentially misleading goal descriptions, a factor that inevitably results in a lower SPL. We posit that enhancing audio lo-

calization, perhaps through the well-established Neural Radiance Fields (NeRF) [32], could further improve the SPL. Moreover, as depicted in the right case of Fig. 3, when the RefPlanner encounters unexplored areas, the ImaAssistant supplies conjectural room layouts. The spatial insight directs the RefPlanner to explore the kitchen instead of the dining room in search of the sink, underscoring the ImaAssistant’s utility. Overall, RILA demonstrates the capacity to adaptively navigate complex environments.

4.3. Ablation Study

Ablation on ImaAssistant. As shown in Table 1, the integration of ImaAssistant markedly improves performance, underscoring the impact of strategic guidance. We also observed considerable advancements in SWS, demonstrating the crucial role of comprehensive layout understanding for long-distance navigation in intricate settings.

Ablation on RefPlanner. We replace our frontier selection RefPlanner with two heuristic frontier-based exploration methods, namely Random which selects a frontier randomly, and Nearest which selects the nearest frontier. We also compare the ability of GPT-3.5 and Llama-2 for frontier selection by replacing GPT-3.5 in RefPlanner with Llama-2. To eliminate the effect from perception, we use ground-truth perceptions (*i.e.*, acoustic object, audio map) in these experiments. In the two heuristic approaches, we automatically execute the *Stop* action when the distance to the goal is less than 1m. As illustrated in Table 2, despite access to ground-truth perceptions, these heuristic methods exhibit poor performance. Notably, Llama-2 7B [43] also struggles to locate the goal object, indicating the lack of spatial reasoning ability of Llama-2 for navigation tasks.

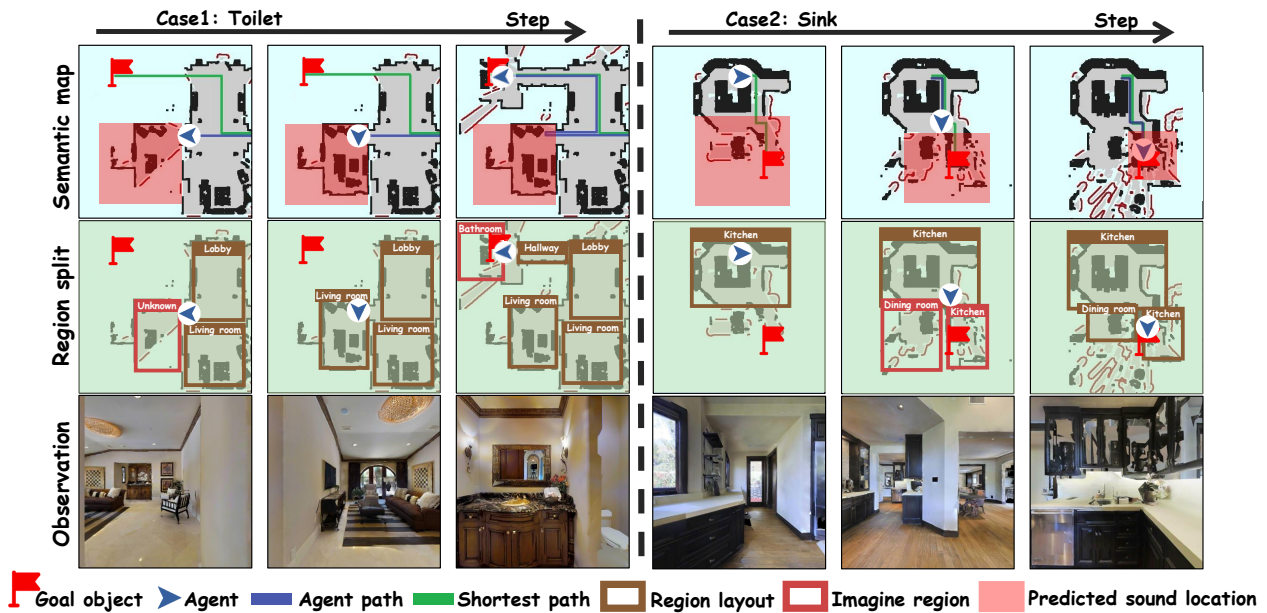


Figure 3. Visualization of two navigation trajectories, including region layouts and egocentric observations. The left case demonstrates how **RefPlanner** reflects on a misleading perception, whereas the right case illustrates that **ImaAssistant** makes imagination and suggestions, guiding **RefPlanner** in waypoint selection based on semantic relevance.

Perception	Accuracy (\uparrow)
Object Recognition	83.9%
Audio Classification	93.0%
Audio Distance	83.8%
Audio Direction	73.7%

Table 3. Accuracy results of different perception modules. Object recognition accuracy represents the probability that the detected item is correctly classified. Audio distance prediction is deemed accurate within a 4-meter error range.

Ablation on Perception Module. Furthermore, we conducted a comprehensive evaluation of the perception modules across 500 episodes from 10 scenes. Results are shown in Table 3. GroundingDINO achieves an 85.0% recall rate on object recognition, indicating only a 15.0% error rate in goal object identification. For all recognized objects, the accuracy also reaches a notable 83.9%. Similarly, the audio classifier distinguishes among 21 classes with an accuracy rate of up to 93.0%. By progressively refining the prediction, RILA made correct predictions in almost all episodes. The accuracy of audio distance prediction is also commendable, reaching 83.8% within a 4-meter margin of error, and has an average distance error of 2.8 meters. Conversely, the accuracy of binary judgments on audio direction is limited to 73.7%, indicating a significant likelihood of error accumulation over steps. To investigate whether the direction

judgment is impacted by complex reverberations in intricate environments, we further separately evaluate episodes based on whether the goal distance is less or more than 15 meters. Notably, accuracy reached 85.6% for shorter distances, in stark contrast to only 59.5% for longer distances. These findings underscore the difficulty of making binary direction determinations in SAVN, particularly over extended distances.

In conclusion, each component of RILA demonstrates competitive performance, with the exception of direction classification, which tends to be less reliable. To delve deeper into the capabilities of RILA, we present a comprehensive analysis in Section 5.

5. Analysis and Discussion

In this section, we focus on the following research questions: (i) Are LLMs adequate for completing complex navigation tasks? (ii) Does the sensory data provided by the SoundSpaces simulation offer clarity and sufficiency for effective navigation? (iii) Are there any inappropriate scenario settings within the current task configuration?

LLMs excel in intricate language-based navigation with inherent commonsense reasoning capabilities. By integrating ground-truth perceptual information, we investigate the navigational planning capabilities of LLMs. Rather than specifying precise goal locations, we provide only a rough

Method	SR \uparrow	SPL \uparrow	DTG \downarrow
Ours	30.2	9.0	11.8
+ GT Audio Semantic	30.2	11.2	11.6
+ GT Audio Localization	52.4	24.6	6.4
+ GT Visual Perception	62.0	39.2	4.8

Table 4. Comparison of incorporating different ground-truth perceptions on the validation dataset. Experiments in each row include the ground-truth information from all previous rows.

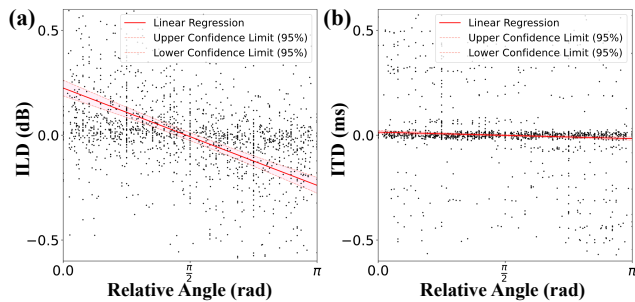


Figure 4. ILT and ITD of the sampled data points. We present the linear regression and the corresponding confidence intervals.

458 area. According to the results in Table 4, our agent achieves
 459 a success rate exceeding 60% with a DTG under 5 on the
 460 validation dataset. Failures typically arise from encounter-
 461 ing similar objects in the target area or due to the inherent
 462 limitations of FBE in long-distance navigation. These find-
 463 ings further confirm the adequacy of LLMs’ planning abili-
 464 ties for navigational tasks.

465 Besides, we observe that providing only ground-truth au-
 466 ditory data yields commendable performance. Conversely,
 467 the success rate markedly decreases in the absence of pre-
 468 cise audio location information, consistent with the experi-
 469 mental results of the perception modules. Although RILA
 470 can effectively utilize potentially imprecise perceptual de-
 471 scription, it remains vulnerable to misdirection caused by
 472 similar objects, thereby constraining the overall perfor-
 473 mance. These observations suggest that the current bottle-
 474 neck in the SAVN task lies in sound source localization.

475 **The auditory sensory data is inadequate for precise lo-**
 476 **calization.** To further investigate the audio localization,
 477 we sampled 4,000 dual-channel audio data points from the
 478 environment and computed two metrics: Interaural Level
 479 Difference (ILD) [44] and ITD. These metrics, crucial for
 480 sound source identification in dual-channel audio [1, 30],
 481 measure differences in sound intensity and arrival time, re-
 482 spectively. The results are depicted in Fig. 3, where the x-
 483 axis represents the sound source angle relative to the agent.
 484 Ideally, these metrics should display a pronounced negative
 485 correlation with the angle [23]. Our analysis reveals that
 486 while ILD demonstrates some negative correlation, serv-

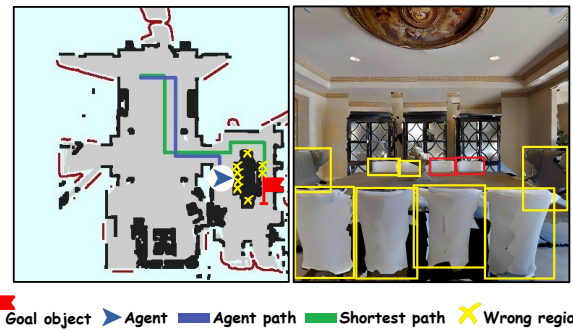


Figure 5. An example of an episode where the goal object is indistinguishable. In this case, the target is far from the agent and surrounded by similar, incorrect items.

ing as the basis for our direction classification, ITD does
 487 not effectively indicate the sound’s relative direction. This
 488 underlines the constraints of the current audio input con-
 489 figuration [24], complicating precise localization based on
 490 auditory inputs. Detailed analysis is provided in Appendix.
 491

Some cases could be further improved. Even in the ab-
 492 sence of precise localization, semantic cues are expected to
 493 guide the agent to the target. However, our observations re-
 494 veal situations where both audio localization is imprecise
 495 and semantic information fails to sufficiently differentiate
 496 between objects. For instance, as illustrated in Figure 5,
 497 the sounding object is distant from the agent, surrounded
 498 by numerous similar items, such as eight chairs in this case.
 499 In SAVN, where sounds are intermittent, the agent must se-
 500 mantically discern the correct stopping point. In this ex-
 501 ample, only two positions would lead to success. Lacking
 502 adequate reasoning cues, the agent resorts to random selec-
 503 tion, leading to failure without exact goal location details.
 504 We postulate that these episodes could be improved by in-
 505 troducing distinct visual differences in target objects, such
 506 as overturning chairs, thus providing definitive cues for the
 507 agent to accurately identify the target.
 508

6. Conclusion

509
 510 In this work, we propose RLIA, a reflective and imaginative
 511 agent for zero-shot semantic audio-visual navigation. By
 512 utilizing distinct models for sensory data processing, RILA
 513 guides an LLM-based reflective planner in active environ-
 514 mental exploration. Throughout this exploration process,
 515 RILA reflectively assesses and disregards erroneous sen-
 516 sory perceptions, especially the goal descriptions. Besides,
 517 we integrate an LLM-based auxiliary imaginative assistant,
 518 designed to generate room layouts and offer strategic guid-
 519 ance, thereby improving global understanding of the envi-
 520 ronment. Comprehensive experimental results demonstrate
 521 the efficacy of RILA.

References

- 522
- 523 [1] Neil Aaronson and William Hartmann. Testing, correcting,
524 and extending the Woodworth model for interaural time differ-
525 ence. *The Journal of the Acoustical Society of America*,
526 135, 2014. 8
- 527 [2] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Cheb-
528 otar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu,
529 Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog,
530 Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Ir-
531 pan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally
532 Jesmonth, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov,
533 Yuheng Kuang, Kuang-Huei Lee, Sergey Levine, Yao Lu,
534 Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao,
535 Kanishka Rao, Jarek Rettinghouse, Diego Reyes, Pierre Ser-
536 manet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vin-
537 cent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu,
538 Mengyuan Yan, and Andy Zeng. Do as I can and not as I say:
539 Grounding language in robotic affordances. *arXiv preprint*
540 *arXiv:2204.01691*, 2022. 1
- 541 [3] Ziad Al-Halah, Santhosh K. Ramakrishnan, and Kristen
542 Grauman. Zero experience required: Plug & play modu-
543 lar transfer learning for semantic visual navigation. *arXiv*
544 *preprint arXiv:2202.02440*, 2022. 2
- 545 [4] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark
546 Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and
547 Anton Van Den Hengel. Vision-and-language navigation:
548 Interpreting visually-grounded navigation instructions in real
549 environments. In *CVPR*, pages 3674–3683, 2018. 1
- 550 [5] Stephen H. Bach, Matthias Broecheler, Bert Huang, and Lise
551 Getoor. Hinge-loss markov random fields and probabilistic
552 soft logic. *arXiv preprint arXiv:1505.04406*, 2017. 5
- 553 [6] Joshua G. W. Bernstein, Olga A. Stakhovskaya, Ger-
554 ald I. Schuchman, Kenneth Kragh Jensen, and Matthew J.
555 Goupell. Interaural time-difference discrimination as a mea-
556 sure of place of stimulation for cochlear-implant users with
557 single-sided deafness. *Trends in Hearing*, 22, 2018. 3
- 558 [7] Changan Chen, Ziad Al-Halah, and Kristen Grauman.
559 Semantic audio-visual navigation. *arXiv preprint*
560 *arXiv:2012.11583*, 2020. 2, 5
- 561 [8] Changan Chen, Unnat Jain, Carl Schissler, Sebastia Vi-
562 cenc Amengual Gari, Ziad Al-Halah, Vamsi Krishna Ithapu,
563 Philip Robinson, and Kristen Grauman. SoundSpaces:
564 Audio-visual navigation in 3D environments. In *ECCV*,
565 pages 17–36. Springer, 2020. 2, 5, 6
- 566 [9] Changan Chen, Ziad Al-Halah, and Kristen Grauman. Se-
567 mantic audio-visual navigation. In *CVPR*, pages 15516–
568 15525, 2021. 1, 2, 3, 5, 6
- 569 [10] Changan Chen, Sagnik Majumder, Ziad Al-Halah, Ruohan
570 Gao, Santhosh Kumar Ramakrishnan, and Kristen Grauman.
571 Learning to set waypoints for audio-visual navigation. In
572 *ICLR*, 2021. 5, 6
- 573 [11] Changan Chen, Carl Schissler, Sanchit Garg, Philip
574 Kobernik, Alexander Clegg, Paul Calamia, Dhruv Batra,
575 Philip W Robinson, and Kristen Grauman. SoundSpaces 2.0:
576 A simulation platform for visual-acoustic learning. *NeurIPS*
577 *Datasets and Benchmarks Track*, 2022. 2, 5
- [12] Peihao Chen, Dongyu Ji, Kunyang Lin, Runhao Zeng,
Thomas H. Li, Mingkui Tan, and Chuang Gan. Weakly-
supervised multi-granularity map learning for vision-and-
language navigation. *arXiv preprint arXiv:2210.07506*,
2022. 2
- [13] Peihao Chen, Xinyu Sun, Hongyan Zhi, Runhao Zeng,
Thomas H Li, Gaowen Liu, Mingkui Tan, and Chuang Gan.
 α^2 nav: Action-aware zero-shot robot navigation by exploit-
ing vision-and-language ability of foundation models. *arXiv*
preprint arXiv:2308.07997, 2023. 1
- [14] Matt Deitke, Winson Han, Alvaro Herrasti, Aniruddha
Kembhavi, Eric Kolve, Roozbeh Mottaghi, Jordi Salvador,
Dustin Schwenk, Eli VanderBilt, Matthew Wallingford, Luca
Weihs, Mark Yatskar, and Ali Farhadi. RoboTHOR: An open
simulation-to-real embodied AI platform. In *CVPR*, 2020. 1
- [15] Vishnu Sashank Dorbala, James F. Mullen Jr. au2, and Di-
nesh Manocha. Can an embodied agent find your "cat-
shaped mug"? LLM-guided exploration for zero-shot object
navigation. *arXiv preprint arXiv:2303.03480*, 2023. 2
- [16] Weixi Feng, Wanrong Zhu, Tsu jui Fu, Varun Jampani, Ar-
jun Akula, Xuehai He, Sugato Basu, Xin Eric Wang, and
William Yang Wang. LayoutGPT: Compositional visual
planning and generation with large language models. *arXiv*
preprint arXiv:2305.15393, 2023. 2
- [17] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat
Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan
Misra. ImageBind: One embedding space to bind them all,
2023. 5
- [18] Jiaming Han, Renrui Zhang, Wenqi Shao, Peng Gao, Peng
Xu, Han Xiao, Kaipeng Zhang, Chris Liu, Song Wen, Ziyu
Guo, Xudong Lu, Shuai Ren, Yafei Wen, Xiaoxin Chen,
Xiangyu Yue, Hongsheng Li, and Yu Qiao. ImageBind-
LLM: Multi-modality instruction tuning. *arXiv preprint*
arXiv:2309.03905, 2023. 1, 5, 6
- [19] Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen
Wang, Daisy Zhe Wang, and Zhiting Hu. Reasoning with
language model is planning with world model. *arXiv*
preprint arXiv:2305.14992, 2023. 2
- [20] David B. Hawkins, Lamar L. Young, and Cheryl Parker. An
investigation of the interaural time difference threshold for
speech. *Perception & Psychophysics*, 24:168–170, 1978. 3
- [21] Lukas Höllein, Ang Cao, Andrew Owens, Justin Johnson,
and Matthias Nießner. Text2Room: Extracting textured 3D
meshes from 2D text-to-image models. In *ICCV*, 2023. 2
- [22] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng,
Yilun Du, Zhenfang Chen, and Chuang Gan. 3D-LLM:
Injecting the 3D world into large language models. *arXiv*
preprint arXiv:2307.12981, 2023. 1
- [23] Maike Klingel, Norbert Kopčo, and Bernhard Laback.
Reweighting of binaural localization cues induced by later-
alization training. *JARO: Journal of the Association for Re-*
search in Otolaryngology, 22:551 – 566, 2020. 8
- [24] Christine Köppl and Catherine Emily Carr. Maps of interau-
ral time difference in the chicken’s brainstem nucleus lami-
naris. *Biological Cybernetics*, 98:541–559, 2008. 8
- [25] Mingxiao Li, Zehao Wang, Tinne Tuytelaars, and Marie-
Francine Moens. Layout-aware dreamer for embod-

- ied referring expression grounding. *arXiv preprint arXiv:2212.00171*, 2022. 2
- [26] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 1
- [27] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding DINO: Marrying DINO with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 1
- [28] Xiao Liu, Hanyu Lai, Hao Yu, Yifan Xu, Aohan Zeng, Zhengxiao Du, Peng Zhang, Yuxiao Dong, and Jie Tang. WebGLM: Towards an efficient web-enhanced question answering system with human preferences. *arXiv preprint arXiv:2306.07906*, 2023. 1
- [29] Xiulong Liu, Sudipta Paul, Moitrey Chatterjee, and Anoop Cherian. Active sparse conversations for improved audio-visual embodied navigation. *arXiv preprint arXiv:2306.04047*, 2023. 1, 2
- [30] Louise Loisel, Michael Dorman, William Yost, Sarah Natale, and René Gifford. Using ILD or ITD cues for sound source localization and speech understanding in a complex listening environment by listeners with bilateral and with hearing-preservation cochlear-implants. *Journal of Speech Language and Hearing Research*, 59:1, 2016. 8
- [31] Manolis Savva*, Abhishek Kadian*, Oleksandr Maksymets*, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A platform for embodied AI research. In *ICCV*, 2019. 2, 5
- [32] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 6
- [33] OpenAI. Introducing ChatGPT, 2022. (Accessed on Jun 18, 2023). 1
- [34] OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1
- [35] Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel Bowman. BBQ: A hand-built bias benchmark for question answering. In *Findings of the Association for Computational Linguistics (ACL)*, pages 2086–2105, Dublin, Ireland, 2022. Association for Computational Linguistics. 1
- [36] Sudipta Paul, Amit K Roy-Chowdhury, and Anoop Cherian. AVLEN: Audio-visual-language embodied navigation in 3d environments. In *NeurIPS*, 2022. 1, 2, 5, 6
- [37] Dhruv Shah, Błażej Osiniński, Sergey Levine, et al. LM-Nav: Robotic navigation with large pre-trained models of language, vision, and action. In *Conference on Robot Learning*, pages 492–504. PMLR, 2023. 2
- [38] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*, 2023. 1
- [39] Andrew Szot, Alex Clegg, Eric Undersander, Erik Wijmans, Yili Zhao, John Turner, Noah Maestre, Mustafa Mukadam, Devendra Chaplot, Oleksandr Maksymets, Aaron Gokaslan, Vladimir Vondrus, Sameer Dharur, Franziska Meier, Wojciech Galuba, Angel Chang, Zsolt Kira, Vladlen Koltun, Jitendra Malik, Manolis Savva, and Dhruv Batra. Habitat 2.0: Training home assistants to rearrange their habitat. In *NeurIPS*, 2021. 2, 5
- [40] Andrew Szot, Max Schwarzer, Bogdan Mazouze, Harsh Agrawal, Walter Talbott, Katherine Metcalf, Natalie Mackraz, Devon Hjelm, and Alexander Toshev. Large language models as generalizable policies for embodied tasks. *arXiv preprint arXiv:2310.17722*, 2023. 2
- [41] Yujin Tang, Wenhao Yu, Jie Tan, Heiga Zen, Aleksandra Faust, and Tatsuya Harada. SayTap: Language to quadrupedal locomotion, 2023. 2
- [42] Gyan Tatiya, Jonathan Francis, Luca Bondi, Ingrid Navarro, Eric Nyberg, Jivko Sinapov, and Jean Oh. Knowledge-driven scene priors for semantic audio-visual embodied navigation. *arXiv preprint arXiv:2212.11345*, 2022. 1, 2, 3, 5, 6
- [43] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 6
- [44] Balemir Uragun and Ramesh Rajan. The discrimination of interaural level difference sensitivity functions: development of a taxonomic data template for modelling. *BMC Neuroscience*, 14:114 – 114, 2013. 8
- [45] Zekun Wang, Ge Zhang, Kexin Yang, Ning Shi, Wangchunshu Zhou, Shaochun Hao, Guangzheng Xiong, Yizhi Li, Mong Yuan Sim, Xiuying Chen, Qingqing Zhu, Zhenzhu Yang, Adam Nik, Qi Liu, Chenghua Lin, Shi Wang, Ruibo Liu, Wenhua Chen, Ke Xu, Dayiheng Liu, Yike Guo, and Jie Fu. Interactive natural language processing. *arXiv preprint arXiv:2305.13246*, 2023. 1
- [46] Justin Wasserman, Karmesh Yadav, Girish Chowdhary, Abhinav Gupta, and Unnat Jain. Last-mile embodied visual navigation. In *Conference on Robot Learning*, 2022. 1
- [47] Canwen Xu, Yichong Xu, Shuohang Wang, Yang Liu, Chengguang Zhu, and Julian McAuley. Small models are valuable plug-ins for large language models. *arXiv preprint arXiv:2305.08848*, 2023. 1
- [48] Karmesh Yadav, Santhosh Kumar Ramakrishnan, John Turner, Aaron Gokaslan, Oleksandr Maksymets, Rishabh Jain, Ram Ramrakhya, Angel X Chang, Alexander Clegg, Manolis Savva, Eric Undersander, Devendra Singh Chaplot, and Dhruv Batra. Habitat challenge 2022. <https://aihabitat.org/challenge/2022/>, 2022. 1
- [49] Jianing Yang, Xuweiyi Chen, Shengyi Qian, Nikhil Madaan, Madhavan Iyengar, David F Fouhey, and Joyce Chai. LLM-Grounder: Open-vocabulary 3D visual grounding with large language model as an agent. *arXiv preprint arXiv:2309.12311*, 2023. 1
- [50] Abdelrahman Younes, Daniel Honerkamp, Tim Welschhold, and Abhinav Valada. Catch me if you hear me: Audio-visual navigation in complex unmapped

750 environments with moving sounds. *arXiv preprint*
751 *arXiv:2111.14843*, 2023. 2

752 [51] Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Aojun
753 Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, and
754 Yu Qiao. LLaMA-Adapter: Efficient fine-tuning of language
755 models with zero-init attention, 2023. 5

756 [52] Zhen Zhang, Anran Lin, Chun Wai Wong, Xiangyu Chu, Qi
757 Dou, and KW Au. Interactive navigation in environments
758 with traversable obstacles using large language and vision-
759 language models. *arXiv preprint arXiv:2310.08873*, 2023.
760 2

761 [53] Gengze Zhou, Yicong Hong, and Qi Wu. NavGPT: Explicit
762 reasoning in vision-and-language navigation with large lan-
763 guage models. *arXiv preprint arXiv:2305.16986*, 2023. 1

764 [54] Kaiwen Zhou, Kaizhi Zheng, Connor Pryor, Yilin Shen,
765 Hongxia Jin, Lise Getoor, and Xin Eric Wang. ESC: Ex-
766 ploration with soft commonsense constraints for zero-shot
767 object navigation. *arXiv preprint arXiv:2301.13166*, 2023.
768 1, 2, 5, 6