

# Vision-and-Language Navigation: A Survey of Tasks, Methods, and Future Directions

Anonymous ACL submission

## Abstract

A long-term goal of AI research is to build intelligent agents that can communicate with humans in natural language, perceive the environment, and perform real-world tasks. Vision-and-Language Navigation (VLN) is a fundamental and interdisciplinary research topic towards this goal, and receives increasing attention from the natural language processing, computer vision, and machine learning communities. In this paper, we review contemporary studies in the emerging field of VLN, covering tasks, evaluation metrics, methods, etc. Through structured analysis of current progress and challenges, we also highlight the limitations of current VLN and opportunities for future work. This paper serves as a thorough reference for the VLN research community.

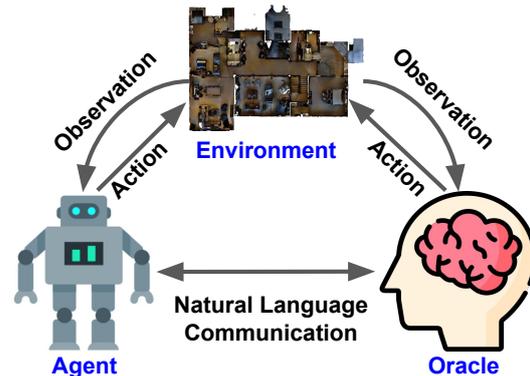


Figure 1: Relation between environment, agent, and oracle (human) in VLN. The agent and oracle discuss the task in natural language. Both the oracle and agent observe and interact with the navigable environment to accomplish a task (although current benchmarks usually do not involve human interaction).

## 1 Introduction

Humans communicate with each other using natural language to issue tasks and request help. A robot that can understand human language and navigate intelligently would significantly benefit human society, both personally and professionally. Such a robot can be spoken to in natural language, and would autonomously execute tasks such as household chores indoors, repetitive delivery work outdoors, or work in hazardous conditions following human commands (bridge inspection; fire-fighting). Scientifically, developing such a robot explores how an artificial agent interprets natural language from humans, perceives its visual environment, and utilizes that information to execute a sequence of actions to complete a task successfully.

Vision-and-Language Navigation (VLN) (Anderson et al., 2018b; Chen et al., 2019; Thomason et al., 2019b) is an emerging research field that aims to build such an embodied agent that can communicate with humans in natural language and navigate in real 3D environments. VLN extends visual navigation in both simulated (Zhu et al., 2017;

Mirowski, 2019) and real environments (Mirowski et al., 2018) with natural language communication. As illustrated in Figure 1, VLN is a task that involves the oracle (frequently a human), the robot agent, and the environment. The agent and the oracle communicate in natural language. The agent may ask for guidance and the oracle could respond. The agent navigates and interacts with the environment to complete the task according to the instructions received and the environment observed. Meanwhile, the oracle observes the environment and agent status, and may interact with the environment to help the agent.

Since the development and release of Room-to-Room (R2R) (Anderson et al., 2018b), many VLN datasets have been introduced. Regarding the degree of communication, researchers create benchmarks where the agent is required to passively understand one instruction before navigation, to benchmarks where agents converse with the oracle in free-form dialog. Regarding the task objective, the requirements for the agent range from strictly following the route described in the initial

064 instruction, to actively exploring the environment  
065 and interacting with objects.

066 Many challenges exist in VLN tasks. First, VLN  
067 faces a complex environment and requires effective  
068 understanding and alignment of information from  
069 different modalities. Second, VLN agents require a  
070 reasoning strategy for the navigation process. Data  
071 scarcity is also an obstacle. Lastly, the general-  
072 ization of a model trained in seen environments  
073 to unseen environments is also essential. We cat-  
074 egorize the solutions according to the respective  
075 challenges. (1) *Representation learning* methods  
076 help understand information from different modal-  
077 ities. (2) *Action strategy learning* aims to make  
078 reasonable decisions based on gathered informa-  
079 tion. (3) *Data-centric learning* methods effectively  
080 utilize the data and address data challenges such  
081 as data scarcity. (4) *Prior exploration* helps the  
082 model familiarize itself with the test environment,  
083 improving its ability to generalize.

084 We make three primary contributions. (1) We  
085 systematically categorize current VLN benchmarks  
086 from *communication complexity* and *task objective*  
087 perspectives, with each category focusing on a dif-  
088 ferent type of VLN task. (2) We hierarchically  
089 classify current solutions and the papers within the  
090 scope. (3) We discuss potential opportunities and  
091 identify future directions.

## 092 2 Tasks and Datasets

093  
094 The ability for an agent to interpret natural lan-  
095 guage instructions (and in some instances, request  
096 feedback during navigation) is what makes VLN  
097 unique from visual navigation (Bonin-Font et al.,  
098 2008). In Table 2, we mainly categorize current  
099 datasets on two axes, *Communication Complexity*  
100 and *Task Objective*.

101 **Communication Complexity** defines the level  
102 at which the agent may converse with the oracle,  
103 and we differentiate three levels: In the first level,  
104 the agent is only required to understand an *Initial*  
105 *Instruction* before navigation starts. In the second  
106 level, the agent sends a signal for help whenever it  
107 is unsure, utilizing the *Guidance* from the oracle.  
108 In the third level, the agent with *Dialog* ability asks  
109 questions in the form of natural language during the  
110 navigation and understands further oracle guidance.

111 **Task Objective** defines *how* the agent attains  
112 its goal. In the first objective type, *Fine-grained*  
113 *Navigation*, the agent can find the target accord-

114 ing to a detailed step-by-step route description. In  
115 the second type, *Course-grained Navigation*, the  
116 agent is required to find a distant target goal with a  
117 coarse navigation description, requiring the agent  
118 to reason a path in an unseen environment. Tasks  
119 in the previous two types only require the agent  
120 to navigate to complete the mission. In the third  
121 type, *Navigation and Object Interaction*, besides  
122 reasoning a path, the agent also needs to interact  
123 with objects in the environment to achieve the goal  
124 since the object might be hidden or need to change  
125 physical states.<sup>1</sup>

### 126 2.1 Initial Instruction

127 Currently, in the setting of many benchmarks, the  
128 agent is given a natural language instruction for  
129 the whole navigation process, such as “Go upstairs  
130 and pass the table in the living room. Turn left and  
131 go through the door in the middle.”

132 **Fine-grained Navigation** An agent needs to  
133 strictly follow the natural language instruction to  
134 reach the target goal. Anderson et al. (2018b) cre-  
135 ate R2R dataset and build the Matterport3D simu-  
136 lator. An embodied agent in R2R moves through a  
137 house in the simulator traversing edges on a navi-  
138 gation graph, jumping to adjacent nodes containing  
139 panoramic views. R2R is extended to create other  
140 VLN benchmarks. Room-for-Room joins paths  
141 in R2R to longer trajectories (Jain et al., 2019).  
142 Yan et al. (2020) collect XL-R2R to extend R2R  
143 with Chinese instructions. RxR (Ku et al., 2020)  
144 contains instructions from English, Hindi, Telegu.  
145 The dataset has more samples and the instructions  
146 in it are time-aligned to the virtual poses of the  
147 instruction. The English split of RxR is further  
148 extended to build Landmark-RxR (He et al., 2021)  
149 by incorporating landmark information.

150 In most current datasets, agents traverse a navi-  
151 gation graph at predefined viewpoints. To facil-  
152 itate transfer learning to real robots, VLN tasks  
153 should provide a continuous action space and a  
154 freely navigable environment. To this end, Krantz  
155 et al. (2020) reconstruct the navigation graph based  
156 R2R trajectories in continuous environments and  
157 create VLNCE. Irshad et al. (2021) propose Robo-  
158 VLN task where the agent operates in a continuous  
159 action space over long-horizon trajectories.

160 Outdoor environments are usually more com-

<sup>1</sup>Navigation and Object Interaction includes both fine-grained and coarse-grained instructions, which ideally should be split further. But given that there are only few datasets in this category, we keep the current categorization in Table 2.

Obj Comms	Fine-grained Navigation	Coarse-grained Navigation	Navigation + Object Interaction
Initial Instruction	Room-to-Room (Anderson et al., 2018b), Room-for-Room (Jain et al., 2019), Room-Across-Room (Ku et al., 2020), XL-R2R (Yan et al., 2020), Landmark-RxR (He et al., 2021), VLNCE (Krantz et al., 2020), TOUCHDOWN (Chen et al., 2019), StreetLearn (Mirowski et al., 2019), StreetNav (Hermann et al., 2020), Talk2Nav (Vasudevan et al., 2021), LANI (Misra et al., 2018)	RoomNav (Wu et al., 2018), REVERIE (Qi et al., 2020b), SOON (Zhu et al., 2021a)	ALFRED (Shridhar et al., 2020)
Guidance	Just Ask (Chi et al., 2020)	VNLA (Nguyen et al., 2019), HANNA (Nguyen and Daumé III, 2019)	None
Dialog	None	CVDN (Thomason et al., 2019b), RobotSlang (Banerjee et al., 2020), Talk the Walk (de Vries et al., 2018), CEREALBAR (Suhr et al., 2019)	TEACh (Padmakumar et al., 2021), Minecraft Collaborative Building (Narayan-Chen et al., 2019)

Table 1: Vision-and-Language Navigation datasets organized by **Communication Complexity** versus **Task Objective**. Please refer to Appendix for more details about the datasets and the commonly used underlying simulators.

plex and contain more objects than indoor environments. In TOUCHDOWN (Chen et al., 2019), an agent follows instructions to navigate a streetview rendered simulation of New York City to find a hidden object. Most photo-realistic outdoor VLN datasets including TOUCHDOWN (Chen et al., 2019), StreetLearn (Mirowski et al., 2019; Mehta et al., 2020), StreetNav (Hermann et al., 2020), and Talk2Nav (Vasudevan et al., 2021) are proposed based on Google Street View.

Research is exploring the use of natural language to guide drones. LANI (Misra et al., 2018) is a 3D synthetic navigation environment, where an agent navigates between landmarks following natural language instructions. Current datasets on drone navigation usually fall in a synthetic environment such as Unity3D (Blukis et al., 2018, 2019).

**Coarse-grained Navigation** In real life, detailed information about the route may not be available since it may be unknown to the human instructor (oracle). Usually, instructions are more concise and contain merely information of the target goal.

RoomNav (Wu et al., 2018) requires agent navigate according to instruction “go to X”, where X is a predefined room or object. The instructions in REVERIE (Qi et al., 2020b) are annotated by humans, and thus more complicated and diverse. The agent navigates through the rooms and differentiates the object against multiple competing candidates. In SOON (Zhu et al., 2021a), an agent receives a long complex coarse-to-fine instruction which gradually narrows down the search scope.

**Navigation+Object Interaction** For some tasks, the target object might be hidden (e.g., the spoon in a drawer), or need to change status (e.g., a sliced apple is requested but only a whole apple is available). In these scenarios, it is necessary to interact with the objects to accomplish the task (e.g., opening the drawer or cutting the apple). Based on indoor scenes in AI2-THOR (Kolve et al., 2017), Shridhar et al. (2020) propose the ALFRED dataset, where agents provided with both coarse-grained and fine-grained instructions complete household tasks in an interactive visual environment.

## 2.2 Guidance

Agents in Guidance VLN tasks may receive further natural language guidance from the oracle during navigation. For example, if the agent is unsure of the next step (e.g., entering the kitchen), it can send a [help] signal, and the oracle would assist by responding “go left”.

**Fine-grained Navigation** The initial fine-grained navigation instruction may still be ambiguous in a complex environment. Guidance from the oracle could clarify possible confusion. Chi et al. (2020) introduce Just Ask—a task where an agent could ask oracle for help during navigation.

**Coarse-grained Navigation** With only a coarse-grained instruction, the agent tends to be more confused and spends more time exploring. Further guidance resolves this ambiguity. VNLA (Nguyen et al., 2019) and HANNA (Nguyen and Daumé III, 2019) both train an agent to navigate indoors to

find objects. The agent could request help from the oracle, which responds by providing a subtask which helps the agent make progress. While oracle in VNLA uses predefined script to respond, the oracle in HANNA uses a neural network to generate natural language responses.

**Navigation+Object Interaction** While VLN is still in its youth, there are no VLN datasets in support of Guidance and Object Interaction.

### 2.3 Dialog

It is human-friendly to use natural language to request help (Banerjee et al., 2020; Thomason et al., 2019b). For example, when agent is not sure about what fruit the human wants, it could ask “What fruit do you want, the banana in the refrigerator or the apple on the table?”, and the human response would provide clear navigation direction.

**Fine-grained Navigation** No datasets are in the scope of this category. Currently, route-detailed instruction with possible guidance could help the agent achieve relatively good performance in most simulated environments. We expect datasets to be developed for this category for complex environments especially with rich dynamics where dialog is necessary to clear confusions.

**Coarse-grained Navigation** CVDN (Thomason et al., 2019b) is a dataset of human-human dialogues. Besides interpreting a natural language instruction and deciding on the following action, the VLN agent also needs to ask questions in natural language for guidance. The oracle, with knowledge of the best next steps, needs to understand and correctly answer said questions. CEREALBAR (Suhr et al., 2019) is a collaborative task between a leader and a follower. Both agents move in a virtual game environment to collect valid sets of cards.

Dialog is important in complex outdoor environments. de Vries et al. (2018) introduce the Talk the Walk dataset, where the guide has knowledge from a map and guides the tourist to a destination, but does not know the tourist’s location; while the tourist navigates a 2D grid via discrete actions.

**Navigation+Object Interaction** Minecraft Collaborative Building (Narayan-Chen et al., 2019) studies how an agent places blocks into a building by communicating with the oracle. TEACH (Padmakumar et al., 2021) is a dataset that studies object interaction and navigation with free-form dialog. The follower converses with the commander and interacts with the environment to complete various

house tasks such as making coffee.

## 3 Evaluation

**Goal-oriented Metrics** mainly consider the agent’s proximity to the goal. The most intuitive is *Success Rate (SR)*, which measures how frequently an agent completes the task within a certain distance of the goal. *Goal Progress* (Thomason et al., 2019b) measures the reduction in remaining distance to the target goal. *Path Length (PL)* measures the total length of the navigation path. *Shortest-Path Distance (SPD)* measures the mean distance between the agent’s final location and the goal. Since a longer path length is undesirable (increases duration and wear-and-tear on actual robots), *Success weighted by Path Length (SPL)* (Anderson et al., 2018a) balances both Success Rate and Path Length. Similarly, *Success weighted by Edit Distance (SED)* (Chen et al., 2019) compares the expert’s actions/trajectory to the agent’s actions/trajectory, also balancing SR and PL. *Oracle Navigation Error (ONE)* takes the shortest distance from any node in the path rather than just the last node, and *Oracle Success Rate (OSR)* measures whether any node in the path is within a threshold from the target location.

**Path-fidelity Metrics** evaluate to what extent an agent follows the desired path. The fidelity between the instruction and the path is important when evaluating an agent’s performance. *Coverage weighted by LS (CLS)* (Jain et al., 2019) is the product of the *Path Coverage (PC)* and *Length Score (LS)* with respect to the reference path. It measures how closely an agent’s trajectory follows the reference path. *Normalized Dynamic Time Warping (nDTW)* (Ilharco et al., 2019) softly penalizes deviations from the reference path to calculate the match between two paths. *Success weighted by normalized Dynamic Time Warping (SDTW)* (Ilharco et al., 2019) further constrains nDTW to only successful episodes to capture both success and fidelity.

## 4 VLN Methods

As shown in Figure 2, we categorize existing methods into *Representation Learning*, *Action Strategy Learning*, *Data-centric Learning*, and *Prior Exploration*. Representation learning methods help agent understand relations between these modalities since VLN involves multiple modalities, including vision, language, and action. Moreover, VLN is a complex reasoning task where mission re-

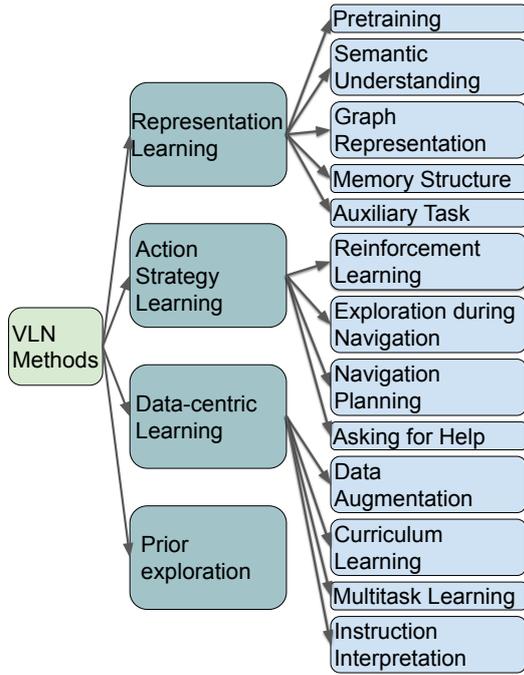


Figure 2: Categories of VLN methods. Methods may not be mutually exclusive to an individual category.

323 results depend on the accumulating steps, and better  
 324 action strategies help the decision-making process.  
 325 Additionally, VLN tasks face challenges within  
 326 their training data. One severe problem is scarcity.  
 327 Collecting training data for VLN is expensive and  
 328 time-consuming, and the existing VLN datasets are  
 329 relatively small with respect to the complexity of  
 330 VLN tasks. Therefore, data-centric methods help  
 331 to utilize the existing data and create more training  
 332 data. Prior exploration helps adapt agents to  
 333 previously unseen environments, improving their  
 334 ability to generalize, decreasing the performance  
 335 gap between seen versus unseen environments.

## 336 4.1 Representation Learning

337 Representation learning helps the agent understand  
 338 how the words in the instruction relate to the per-  
 339 ceived features in the environment.

### 340 4.1.1 Pretraining

341 **Vision or Language** Using a pretrained model to  
 342 initialize a vision or text encoder provides agents  
 343 with single-modality knowledge. pretrained vision  
 344 models may use a ResNet (He et al., 2016)  
 345 or Vision Transformers (Dosovitskiy et al., 2020).  
 346 Other navigation tasks (Wijmans et al., 2019b)  
 347 may also provide visual initialization (Krantz et al.,  
 348 2020). Large pretrained language models such as  
 349 BERT (Devlin et al., 2019) and GPT (Radford et al.,  
 350 2019) can encode language and improve instruction

351 understanding (Li et al., 2019), which can be fur-  
 352 ther pretrained with VLN instructions (Pashevich  
 353 et al., 2021) before fine-tuning in VLN task.

354 **Vision and Language** Vision-and-language pre-  
 355 trained models provide good joint representation  
 356 for text and vision. A common practice is to ini-  
 357 tialize the VLN agent with a pretrained model such  
 358 as ViLBERT (Lu et al., 2019). The agent may be  
 359 further trained with VLN-specific features such as  
 360 objects and rooms (Qi et al., 2021).

361 **VLN** Downstream tasks benefit from being closely  
 362 related to the pretraining task. Researchers also  
 363 explored pretraining on the VLN domain directly.  
 364 VLN-BERT (Majumdar et al., 2020) pretrains nav-  
 365 igation models to measure the compatibility be-  
 366 tween paths and instructions. PREVALENT (Hao  
 367 et al., 2020) is trained from scratch on image-text-  
 368 action triplets to learn textual representations in  
 369 VLN tasks. The [CLS] token in BERT-based pre-  
 370 training models could be leveraged in a recurrent  
 371 fashion to represent history state (Hong et al., 2021;  
 372 Moudgil et al., 2021). Airbert (Guhur et al., 2021)  
 373 achieve good performance on few-shot setting after  
 374 pretraining on a large-scale in-domain dataset.

### 375 4.1.2 Semantic Understanding

376 Semantic understanding of VLN tasks incorporates  
 377 knowledge about important features in VLN. In  
 378 addition to the raw features, high-level semantic  
 379 representations also improve performance.

380 **Intra-Modality** Visual or textual modalities can  
 381 be decomposed into many features, which matter  
 382 differently in VLN. The overall visual features ex-  
 383 tracted by a neural model may actually hurt the per-  
 384 formance in some cases (Thomason et al., 2019a;  
 385 Hu et al., 2019; Zhang et al., 2020b). Therefore, it  
 386 is important to find the feature(s) that best improve  
 387 performance. High-level features such as visual  
 388 appearance, route structure, and detected objects  
 389 outperform the low level visual features extracted  
 390 by CNN (Hu et al., 2019). Different types of tokens  
 391 within the instruction also function differently (Zhu  
 392 et al., 2021b). Extracting these tokens and encod-  
 393 ing the object tokens and directions tokens are cru-  
 394 cial (Qi et al., 2020a; Zhu et al., 2021b).

395 **Inter-Modality** Semantic connections between  
 396 different modalities: actions, scenes, observed ob-  
 397 jects, direction clues, and objects mentioned in in-  
 398 structions can be extracted and then softly aligned  
 399 with attention mechanism (Qi et al., 2020a; Gao  
 400 et al., 2021). The soft alignment also highlights rel-  
 401 evant parts of the instruction with respect to the cur-

rent step (Landi et al., 2019; Zhang et al., 2020a).

### 4.1.3 Graph Representation

Graph has been widely applied to model relationships. Building graph to incorporate structured information from instruction and observation provides explicit semantic relation to guide the navigation. The graph neural network may encode the relation between text and vision to better interpret the context information (Hong et al., 2020a). The graph could record the location information during the navigation, which can be used to predict the most likely trajectory (Anderson et al., 2019a) or probability distribution over action space (Deng et al., 2020). When connected with prior exploration, an overview graph about the navigable environment (Chen et al., 2021a) can be built to improve navigation interpretation.

### 4.1.4 Memory-augmented Model

Information accumulates as the agent navigates, which is not efficient to utilize directly. Memory structure helps the agent effectively leverage the navigation history. Some solutions leverage memory modules such as LSTMs or recurrently utilize informative states (Hong et al., 2021), which can be relatively easily implemented, but may struggle to remember features at the beginning of the path as path length increases. Another solution is to build a separate memory model to store the relevant information (Zhu et al., 2020c; Lin et al., 2021; Nguyen and Daumé III, 2019). Notably, by hierarchically encoding a single view, a panorama, and then all panoramas in history, HAMT (Chen et al., 2021b) successfully utilized the full navigation history for decision-making.

### 4.1.5 Auxiliary Tasks

Auxiliary tasks help the agent better understand the environment and its own status without extra labels. From the machine learning perspective, an auxiliary task is usually achieved in the form of an additional loss function. The auxiliary task could, for example, explain its previous actions, or predict information about future decisions (Zhu et al., 2020a). Auxiliary tasks could also involve the current mission such as current task accomplishment, and vision instruction alignment (Ma et al., 2019a; Zhu et al., 2020a). Notably, auxiliary tasks are effective when adapting pretrained representations into the VLN domain (Huang et al., 2019).

## 4.2 Action Strategy Learning

With a variety of possible action sequences, action strategy learning provides a variety of methods to help the agent decide on those best actions.

### 4.2.1 Reinforcement Learning

VLN is a sequential decision-making problem and can naturally be modeled as a Markov decision process. So Reinforcement Learning (RL) methods are proposed to learn better policy for VLN tasks. A critical challenge for RL methods is that VLN agents only receive the success signal at the end of the episode, so it is difficult to know which actions to attribute success to, and which to penalize. To address the ill-posed feedback issue, Wang et al. (2019) propose RCM model to enforce cross-modal grounding both locally and globally, with goal-oriented extrinsic reward and instruction-fidelity intrinsic reward. He et al. (2021) propose to utilize the local alignment between the instruction and critical landmarks as the reward. Evaluation metrics such as CLS (Jain et al., 2019) or nDTW (Ilharco et al., 2019) can also provide informative reward signal (Landi et al., 2020).

To model rich dynamics in the environment, Wang et al. (2018) leverage model-based reinforcement learning to predict the next state and improve the generalization in unseen environment. Zhang et al. (2020a) find recursively alternating the learning schemes of imitation and reinforcement learning improve the performance.

### 4.2.2 Exploration during Navigation

Exploring and gathering environmental information while navigating provides a better understanding of the state space. Student-forcing is a frequently used strategy, where the agent keeps navigating based on sampled actions and is supervised by the shortest-path action (Anderson et al., 2018b).

There is a tradeoff between exploration versus exploitation: with more exploration, the agent sees better performance at the cost of a longer path and longer duration, so the model needs to determine when and how deep to explore (Wang et al., 2020a). After having gathered the local information, the agent needs to decide which step to choose, or whether to backtrack (Ke et al., 2019). Notably, Koh et al. (2021) designed Pathdreamer, a visual world model to synthesize visual observation future viewpoints without actually looking ahead.

499	<b>4.2.3 Navigation Planning</b>	548
500	Planing future navigation steps leads to a better	549
501	action strategy. From the visual side, predicting	550
502	the waypoints (Krantz et al., 2021), next state and	551
503	reward (Wang et al., 2018), generate future obser-	552
504	vation (Koh et al., 2021) or incorporating neigh-	553
505	bor views (An et al., 2021) has proven effective.	554
506	The natural language instruction also contains land-	
507	marks and direction clues to plan detailed steps.	
508	Anderson et al. (2019b) predict the forthcoming	
509	events based on the instruction, which is used to	
510	predict actions with a semantic spatial map.	
511	<b>4.2.4 Asking for Help</b>	
512	An intelligent agent asks for help when uncertain	
513	about the next action. Action probabilities or a	
514	separately trained model (Chi et al., 2020; Zhu	
515	et al., 2021c; Nguyen et al., 2021) can be leveraged	
516	to decide whether to ask for help. Using natu-	
517	ral language to converse with the oracle covers a	
518	wider problem scope than sending a signal. Both	
519	rule-based methods (Padmakumar et al., 2021) and	
520	neural-based methods (Roman et al., 2020; Nguyen	
521	et al., 2021) have been developed to build naviga-	
522	tion agents with dialog ability. Meanwhile, for	
523	tasks (Thomason et al., 2019b; Padmakumar et al.,	
524	2021) that do not provide an oracle agent to answer	
525	question in natural language, researchers also needs	
526	to build a rule-based (Padmakumar et al., 2021) or	
527	neural-based (Roman et al., 2020) oracle.	
528	<b>4.3 Data-centric Learning</b>	
529	Compared with previously discussed works that	
530	focus on building a better VLN agent structure,	
531	data-centric methods most effectively utilize the	
532	existing data, or create synthetic data.	
533	<b>4.3.1 Data Augmentation</b>	
534	<b>Trajectory-Instruction Augmentation</b> Aug-	
535	mented path-instruction pairs could be used in VLN	
536	directly. Currently the common practice is to train	
537	a speaker module to generate instructions given a	
538	navigation path (Fried et al., 2018). This generated	
539	data could have varying quality (Zhao et al., 2021).	
540	Therefore an alignment scorer (Huang et al., 2019)	
541	or adversarial discriminator (Fu et al., 2020) can	
542	select high-quality pairs for augmentation.	
543	<b>Environment Augmentation</b> Generating more en-	
544	vironment data not only helps generate more trajec-	
545	tories, but also alleviates the problem of overfitting	
546	in seen environments. Randomly masking the same	
547	visual feature across different viewpoints (Tan et al.,	
	2019) or simply splitting the house scenes and re-	548
	mixing them (Liu et al., 2021) could create new	549
	environments, which could further be used to gener-	550
	ate more trajectory-instruction pairs (Fried et al.,	551
	2018). Training data may also be augmented by	552
	replacing some visual features with counterfactual	553
	ones (Parvaneh et al., 2020).	554
	<b>4.3.2 Curriculum Learning</b>	555
	Curriculum learning (Bengio et al., 2009) gradually	556
	increases the task’s difficulty during the training	557
	process. The instruction length could be a metric	558
	for task difficulty. BabyWalk (Zhu et al., 2020b)	559
	keep increasing training samples’ instruction length	560
	during the training process. Attributes from the	561
	trajectory may also be used to rank task difficulty.	562
	Zhang et al. (2021) rearrange the R2R dataset using	563
	the number of rooms each path traverses. They	564
	found curriculum learning helps smooth the loss	565
	landscape and find a better local optima.	566
	<b>4.3.3 Multitask Learning</b>	567
	Different VLN tasks can learn from each other by	568
	cross-task knowledge transfer. Wang et al. (2020c)	569
	propose an environment-agnostic multitask naviga-	570
	tion model for both VLN and Navigation from Di-	571
	alog History tasks (Thomason et al., 2019b). Chap-	572
	lot et al. (2020) propose an attention module to train	573
	a multitask navigation agent to follow instructions	574
	and answer questions (Wijmans et al., 2019a).	575
	<b>4.3.4 Instruction Interpretation</b>	576
	A trajectory instruction phrased multiple times in	577
	different ways may help the agent better under-	578
	stand its objective. LEO (Xia et al., 2020) leverages	579
	and encodes all the instructions with a shared set	580
	of parameters to enhance the textual understand-	581
	ing. Shorter, and more concise instructions pro-	582
	vide clearer guidance for the agent compared to	583
	longer, semantically entangled instructions, thus	584
	Hong et al. (2020b) breaks long instructions into	585
	shorter ones, allowing the agent to track progress	586
	and focus on each atomic instruction individually.	587
	<b>4.4 Prior Exploration</b>	588
	Good performance in seen environments often can-	589
	not generalize to unseen environments (Parvaneh	590
	et al., 2020; Tan et al., 2019). Prior exploration	591
	methods allow the agent to observe and adapt to	592
	unseen environments <sup>2</sup> , bridging the performance	593

<sup>2</sup>Thus prior exploration methods are not directly comparable with other VLN methods.

594 gap between seen and unseen environments.

595 Wang et al. (2019) introduce a self-supervised  
596 imitation learning to learn from the agent’s own  
597 past, good behaviors. The best navigation path  
598 determined to align the instruction the best by a  
599 matching critic will be used to update the agent.  
600 Tan et al. (2019) leverage the testing environments  
601 to sample and augment paths for adaptation. Fu  
602 et al. (2020) propose environment-based prior ex-  
603 ploration, where the agent can only explore a partic-  
604 ular environment where it is deployed. When con-  
605 necting with graph, prior exploration may construct  
606 a map or overview about the unseen environment,  
607 providing explicit guidance for navigation (Chen  
608 et al., 2021a; Zhou et al., 2021).

## 609 5 Conclusion and Future Directions

610 In this paper, we discuss the importance of VLN  
611 agents as a part of society, how their tasks vary as  
612 a function of communication level versus task ob-  
613 jective, and how different agents may be evaluated.  
614 We broadly review VLN methodologies and cate-  
615 gorize them. This paper only discusses these issues  
616 broadly at an introductory level. In reviewing these  
617 papers, we can see the immense progress that has  
618 already been made, as well as directions that this  
619 research topic can be expanded on.

620 Current methods usually do not explicitly uti-  
621 lize external knowledge such as objects and house  
622 descriptions in Wikipedia. Incorporating knowl-  
623 edge also improves the interpretability and trust of  
624 embodied AI. Moreover, currently several naviga-  
625 tion agents learn which direction to move and with  
626 what to interact, but there is a last-mile problem  
627 of VLN—how to interact with objects. Anderson  
628 et al. (2018b) asked whether a robot could learn to  
629 “Bring me a spoon”; new research may ask how a  
630 robot can learn to “Pick up a spoon”. The environ-  
631 ments also lack diversity: most interior terrestrial  
632 VLN data consists of American houses, but never  
633 warehouses or hospitals: the places where these  
634 agents may be of most use.

635 Below we detail additional future directions:

636 **Collaborative VLN** Current VLN benchmarks  
637 and methods predominantly focus on tasks where  
638 only one agent navigates, yet complicated real-  
639 world scenarios may require several robots collabo-  
640 rating. Multi-agent VLN tasks require development  
641 in swarm intelligence, information communication,  
642 and performance evaluation. VLN studies the rela-  
643 tionship between the human and the environment

644 in Figure 1, yet here humans are oracles simply  
645 observing (but not acting on) the environment. Col-  
646 laboration between humans and robots is crucial for  
647 them to work together as teams (e.g., as personal  
648 assistants or helping in construction). Future work  
649 may target at collaborative VLN between multiple  
650 agents or between human and agents.

651 **Simulation to Reality** There is a performance loss  
652 when transferred to real-life robot navigation (An-  
653 derson et al., 2020). Real robots function in contin-  
654 uous space, but most simulators only allow agents  
655 to “hop” through a pre-defined navigation graph  
656 which is unrealistic for three reasons (Krantz et al.,  
657 2020). Navigation graphs assume: (1) perfect  
658 localization—in the real world is a noisy estimate;  
659 (2) oracle navigation—real robots cannot “teleport”  
660 to a new node; (3) known topology—in reality an  
661 agent may not have access to a preset list of naviga-  
662 ble nodes. Continuous implementations of realistic  
663 environments may contain patches of the images,  
664 be blurred, or have parallax errors, making them  
665 unrealistic. A simulation that is based on both  
666 a 3D model and realistic imagery could improve  
667 the match between virtual sensors (in simulation)  
668 and real sensors. Lastly, most simulators assume a  
669 static environment only changed by the agent. This  
670 does not account for other dynamics such as people  
671 walking or objects moving, nor does it account for  
672 lighting conditions through the day. VLN environ-  
673 ments with probabilistic transition function may  
674 also narrow the gap between simulation and reality.

675 **Ethics & Privacy** During both training and in-  
676 ference, VLN agents may observe and store sen-  
677 sitive information that can get leaked or misused.  
678 Effective navigation with privacy protection is cru-  
679 cially important. Relevant areas such as federated  
680 learning (Konečný et al., 2016) or differential pri-  
681 vacy (Dwork et al., 2006) could also be studied in  
682 VLN domain to preserve the privacy of training  
683 and inference environments.

684 **Multicultural VLN** VLN lacks diversity in  
685 3D environments: most outdoor VLN datasets  
686 use Google Street View recorded in major Amer-  
687 ican cities, but lacks data in developing countries.  
688 Agents trained on American data face potential  
689 generalization problems in other city or housing  
690 layouts. Future work should explore more diverse  
691 environments across multiple cultures and regions.  
692 Multilingual VLN datasets (Ku et al., 2020; Yan  
693 et al., 2020) could be good resources to study multi-  
694 cultural differences from the linguistic perspective.

## References

- 696 Dong An, Yuankai Qi, Yan Huang, Qi Wu, Liang Wang,  
697 and Tieniu Tan. 2021. Neighbor-view enhanced  
698 model for vision and language navigation. *arXiv*  
699 *preprint arXiv:2107.07201*.
- 700 Peter Anderson, Angel Chang, Devendra Singh Chap-  
701 lot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen  
702 Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mot-  
703 taghi, Manolis Savva, et al. 2018a. On evalua-  
704 tion of embodied navigation agents. *arXiv preprint*  
705 *arXiv:1807.06757*.
- 706 Peter Anderson, Ayush Shrivastava, Devi Parikh, Dhruv  
707 Batra, and Stefan Lee. 2019a. Chasing ghosts: In-  
708 struction following as bayesian state tracking. In  
709 *Advances in Neural Information Processing Systems*  
710 *(NeurIPS)*.
- 711 Peter Anderson, Ayush Shrivastava, Devi Parikh, Dhruv  
712 Batra, and Stefan Lee. 2019b. Chasing ghosts: In-  
713 struction following as bayesian state tracking. *Ad-*  
714 *vances in Neural Information Processing Systems*,  
715 32:371–381.
- 716 Peter Anderson, Ayush Shrivastava, Joanne Truong, Ar-  
717 jun Majumdar, Devi Parikh, Dhruv Batra, and Ste-  
718 fan Lee. 2020. Sim-to-real transfer for vision-and-  
719 language navigation. In *Conference on Robot Learn-*  
720 *ing (CoRL)*.
- 721 Peter Anderson, Qi Wu, Damien Teney, Jake Bruce,  
722 Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen  
723 Gould, and Anton van den Hengel. 2018b. Vision-  
724 and-language navigation: Interpreting visually-  
725 grounded navigation instructions in real environ-  
726 ments. In *Proceedings of the IEEE Conference on*  
727 *Computer Vision and Pattern Recognition (CVPR)*.
- 728 Shurjo Banerjee, Jesse Thomason, and Jason J. Corso.  
729 2020. [The RobotSlang Benchmark: Dialog-guided](#)  
730 [robot localization and navigation](#). In *Conference on*  
731 *Robot Learning (CoRL)*.
- 732 Yoshua Bengio, Jérôme Louradour, Ronan Collobert,  
733 and Jason Weston. 2009. Curriculum learning. In  
734 *Proceedings of the 26th annual international confer-*  
735 *ence on machine learning*, pages 41–48.
- 736 Valts Blukis, Nataly Brukhim, Andrew Bennett, Ross A.  
737 Knepper, and Yoav Artzi. 2018. Following high-level  
738 navigation instructions on a simulated quadcopter  
739 with imitation learning. In *Robotics: Science and*  
740 *Systems (RSS)*.
- 741 Valts Blukis, Yannick Terme, Eyvind Niklasson, Ross A.  
742 Knepper, and Yoav Artzi. 2019. Learning to map  
743 natural language instructions to physical quadcopter  
744 control using simulated flight. In *Conference on*  
745 *Robot Learning (CoRL)*.
- 746 Francisco Bonin-Font, Alberto Ortiz, and Gabriel Oliver.  
747 2008. Visual navigation for mobile robots: A survey.  
748 *Journal of intelligent and robotic systems*, 53(3):263–  
749 296.
- Angel Chang, Angela Dai, Thomas Funkhouser, Maciej  
Halber, Matthias Niessner, Manolis Savva, Shuran  
Song, Andy Zeng, and Yinda Zhang. 2017. Matter-  
port3D: Learning from RGB-D data in indoor envi-  
ronments. *International Conference on 3D Vision*  
(3DV). 750 751 752 753 754 755
- Devendra Singh Chaplot, Lisa Lee, Ruslan Salakhutdi-  
nov, Devi Parikh, and Dhruv Batra. 2020. Embodied  
multimodal multitask learning. In *Proceedings of*  
*the Twenty-Ninth International Joint Conference on*  
*Artificial Intelligence, IJCAI-20*. International Joint  
Conferences on Artificial Intelligence Organization. 756 757 758 759 760 761
- Howard Chen, Alane Suhr, Dipendra Misra, Noah  
Snively, and Yoav Artzi. 2019. [Touchdown: Natural](#)  
[language navigation and spatial reasoning in visual](#)  
[street environments](#). In *2019 IEEE/CVF Conference*  
*on Computer Vision and Pattern Recognition (CVPR)*,  
pages 12530–12539. 762 763 764 765 766 767
- Kevin Chen, Junshen K Chen, Jo Chuang, Marynel  
Vázquez, and Silvio Savarese. 2021a. Topological  
planning with transformers for vision-and-language  
navigation. In *Proceedings of the IEEE/CVF Confer-*  
*ence on Computer Vision and Pattern Recognition*,  
pages 11276–11286. 768 769 770 771 772 773
- Shizhe Chen, Pierre-Louis Guhur, Cordelia Schmid, and  
Ivan Laptev. 2021b. History aware multimodal trans-  
former for vision-and-language navigation. *arXiv*  
*preprint arXiv:2110.13309*. 774 775 776 777
- Ta-Chung Chi, Minmin Shen, Mihail Eric, Seokhwan  
Kim, and Dilek Hakkani-tur. 2020. Just ask: An in-  
teractive learning framework for vision and language  
navigation. In *Proceedings of the AAAI Conference*  
*on Artificial Intelligence*, volume 34, pages 2459–  
2466. 778 779 780 781 782 783
- Harm de Vries, Kurt Shuster, Dhruv Batra, Devi Parikh,  
Jason Weston, and Douwe Kiela. 2018. [Talk the](#)  
[walk: Navigating new york city through grounded](#)  
[dialogue](#). 784 785 786 787
- Zhiwei Deng, Karthik Narasimhan, and Olga Rus-  
sakovsky. 2020. Evolving graphical planner: Con-  
textual global planning for vision-and-language navi-  
gation. *Advances in Neural Information Processing*  
*Systems*, 2020-December. 788 789 790 791 792
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
Kristina Toutanova. 2019. Bert: Pre-training of deep  
bidirectional transformers for language understand-  
ing. In *NAACL-HLT (1)*. 793 794 795 796
- Alexey Dosovitskiy, Lucas Beyer, Alexander  
Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
Thomas Unterthiner, Mostafa Dehghani, Matthias  
Minderer, Georg Heigold, Sylvain Gelly, et al. 2020.  
An image is worth 16x16 words: Transformers  
for image recognition at scale. In *International*  
*Conference on Learning Representations*. 797 798 799 800 801 802 803

804	Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. 2006. Calibrating noise to sensitivity in private data analysis. In <i>Theory of cryptography conference</i> , pages 265–284. Springer.	<i>Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 3360–3376, Online. Association for Computational Linguistics.	859 860 861
808	Daniel Fried, Ronghang Hu, Volkan Cirik, Anna Rohrbach, Jacob Andreas, Louis-Philippe Morency, Taylor Berg-Kirkpatrick, Kate Saenko, Dan Klein, and Trevor Darrell. 2018. Speaker-follower models for vision-and-language navigation. In <i>Neural Information Processing Systems (NeurIPS)</i> .	Yicong Hong, Qi Wu, Yuankai Qi, Cristian Rodriguez-Opazo, and Stephen Gould. 2021. Vln bert: A recurrent vision-and-language bert for navigation. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 1643–1653.	862 863 864 865 866 867
814	Tsu-Jui Fu, Xin Eric Wang, Matthew Peterson, Scott Grafton, Miguel Eckstein, and William Yang Wang. 2020. Counterfactual vision-and-language navigation via adversarial path sampler. In <i>European Conference on Computer Vision (ECCV)</i> .	Ronghang Hu, Daniel Fried, Anna Rohrbach, Dan Klein, Trevor Darrell, and Kate Saenko. 2019. <a href="#">Are you looking? grounding to multiple modalities in vision-and-language navigation</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 6551–6557, Florence, Italy. Association for Computational Linguistics.	868 869 870 871 872 873 874
819	Chen Gao, Jinyu Chen, Si Liu, Luting Wang, Qiong Zhang, and Qi Wu. 2021. Room-and-object aware knowledge reasoning for remote embodied referring expression. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 3064–3073.	Haoshuo Huang, Vihan Jain, Harsh Mehta, Alexander Ku, Gabriel Magalhaes, Jason Baldrige, and Eugene Ie. 2019. Transferable representation learning in vision-and-language navigation. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</i> .	875 876 877 878 879 880
825	Pierre-Louis Guhur, Makarand Tapaswi, Shizhe Chen, Ivan Laptev, and Cordelia Schmid. 2021. Airbert: In-domain pretraining for vision-and-language navigation. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</i> , pages 1634–1643.	Gabriel Ilharco, Vihan Jain, Alexander Ku, Eugene Ie, and Jason Baldrige. 2019. General evaluation for instruction conditioned navigation using dynamic time warping. <i>arXiv preprint arXiv:1907.05446</i> .	881 882 883 884
831	Weituo Hao, Chunyuan Li, Xiujun Li, Lawrence Carin, and Jianfeng Gao. 2020. Towards learning a generic agent for vision-and-language navigation via pre-training. <i>Conference on Computer Vision and Pattern Recognition (CVPR)</i> .	Muhammad Zubair Irshad, Chih-Yao Ma, and Zsolt Kira. 2021. Hierarchical cross-modal agent for robotics vision-and-language navigation. <i>arXiv preprint arXiv:2104.10674</i> .	885 886 887 888
836	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 770–778.	Vihan Jain, Gabriel Magalhaes, Alexander Ku, Ashish Vaswani, Eugene Ie, and Jason Baldrige. 2019. <a href="#">Stay on the path: Instruction fidelity in vision-and-language navigation</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1862–1872, Florence, Italy. Association for Computational Linguistics.	889 890 891 892 893 894 895
841	Keji He, Yan Huang, Qi Wu, Jianhua Yang, Dong An, Shuanglin Sima, and Liang Wang. 2021. Landmark-rrr: Solving vision-and-language navigation with fine-grained alignment supervision. In <i>NeurIPS</i> .	Liyiming Ke, Xiujun Li, Yonatan Bisk, Ari Holtzman, Zhe Gan, Jingjing Liu, Jianfeng Gao, Yejin Choi, and Siddhartha Srinivasa. 2019. Tactical rewind: Self-correction via backtracking in vision-and-language navigation. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> .	896 897 898 899 900 901 902
845	Karl Moritz Hermann, Mateusz Malinowski, Piotr Mirowski, Andras Banki-Horvath, Keith Anderson, and Raia Hadsell. 2020. Learning to follow directions in street view. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 34, pages 11773–11781.	Jing Yu Koh, Honglak Lee, Yinfei Yang, Jason Baldrige, and Peter Anderson. 2021. Pathdreamer: A world model for indoor navigation. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</i> , pages 14738–14748.	903 904 905 906 907
851	Yicong Hong, Cristian Rodriguez, Yuankai Qi, Qi Wu, and Stephen Gould. 2020a. Language and visual entity relationship graph for agent navigation. <i>Advances in Neural Information Processing Systems</i> , 33:7685–7696.	Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. 2017. AI2-THOR: An Interactive 3D Environment for Visual AI. <i>arXiv</i> .	908 909 910 911 912
856	Yicong Hong, Cristian Rodriguez, Qi Wu, and Stephen Gould. 2020b. <a href="#">Sub-instruction aware vision-and-language navigation</a> . In <i>Proceedings of the 2020</i>		

913	Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. 2016. Federated learning: Strategies for improving communication efficiency. <i>arXiv preprint arXiv:1610.05492</i> .	968	Chih-Yao Ma, Zuxuan Wu, Ghassan AlRegib, Caiming Xiong, and Zsolt Kira. 2019b. The regretful agent: Heuristic-aided navigation through progress estimation. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> .	969
914		970		971
915		971		972
916		972		973
917		973		974
918	Jacob Krantz, Aaron Gokaslan, Dhruv Batra, Stefan Lee, and Oleksandr Maksymets. 2021. Waypoint models for instruction-guided navigation in continuous environments. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</i> , pages 15162–15171.	974	Arjun Majumdar, Ayush Shrivastava, Stefan Lee, Peter Anderson, Devi Parikh, and Dhruv Batra. 2020. Improving vision-and-language navigation with image-text pairs from the web. In <i>Proceedings of the European Conference on Computer Vision (ECCV)</i> .	975
919		975		976
920		976		977
921		977		978
922		978	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. 2019. Habitat: A Platform for Embodied AI Research. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</i> .	979
923		979		980
924	Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. 2020. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In <i>Computer Vision – ECCV 2020</i> , pages 104–120, Cham. Springer International Publishing.	980		981
925		981		982
926		982		983
927		983		984
928		984		985
929	Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, and Jason Baldrige. 2020. Room-Across-Room: Multilingual vision-and-language navigation with dense spatiotemporal grounding. In <i>Conference on Empirical Methods for Natural Language Processing (EMNLP)</i> .	985	Harsh Mehta, Yoav Artzi, Jason Baldrige, Eugene Ie, and Piotr Mirowski. 2020. <a href="#">Retouchdown: Releasing touchdown on StreetLearn as a public resource for language grounding tasks in street view</a> . In <i>Proceedings of the Third International Workshop on Spatial Language Understanding</i> , pages 56–62, Online. Association for Computational Linguistics.	986
930		986		987
931		987		988
932		988		989
933		989		990
934		990		991
935	Federico Landi, Lorenzo Baraldi, Marcella Cornia, Massimiliano Corsini, and Rita Cucchiara. 2020. <a href="#">Perceive, transform, and act: Multi-modal attention networks for vision-and-language navigation</a> .	991		992
936		992		993
937		993		994
938		994		995
939	Federico Landi, Lorenzo Baraldi, Massimiliano Corsini, and Rita Cucchiara. 2019. Embodied vision-and-language navigation with dynamic convolutional filters. In <i>Proceedings of the British Machine Vision Conference</i> .	995	Piotr Mirowski. 2019. <a href="#">Learning to navigate</a> . In <i>1st International Workshop on Multimodal Understanding and Learning for Embodied Applications</i> , MULEA '19, page 25, New York, NY, USA. Association for Computing Machinery.	996
940		996		997
941		997		998
942		998		999
943		999		1000
944	Xiujun Li, Chunyuan Li, Qiaolin Xia, Yonatan Bisk, Asli Celikyilmaz, Jianfeng Gao, Noah Smith, and Yejin Choi. 2019. Robust navigation with language pretraining and stochastic sampling.	1000	Piotr Mirowski, Andras Banki-Horvath, Keith Anderson, Denis Teplyashin, Karl Moritz Hermann, Mateusz Malinowski, Matthew Koichi Grimes, Karen Simonyan, Koray Kavukcuoglu, Andrew Zisserman, et al. 2019. The streetlearn environment and dataset. <i>arXiv preprint arXiv:1903.01292</i> .	1001
945		1001		1002
946		1002		1003
947		1003		1004
948		1004		1005
949	Xiangru Lin, Guanbin Li, and Yizhou Yu. 2021. Scene-intuitive agent for remote embodied visual grounding. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 7036–7045.	1005	Piotr Mirowski, Matthew Koichi Grimes, Mateusz Malinowski, Karl Moritz Hermann, Keith Anderson, Denis Teplyashin, Karen Simonyan, Koray Kavukcuoglu, Andrew Zisserman, and Raia Hadsell. 2018. Learning to navigate in cities without a map. In <i>Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS'18</i> , page 2424–2435, Red Hook, NY, USA. Curran Associates Inc.	1006
950		1006		1007
951		1007		1008
952		1008		1009
953	Chong Liu, Fengda Zhu, Xiaojun Chang, Xiaodan Liang, Zongyuan Ge, and Yi-Dong Shen. 2021. Vision-language navigation with random environmental mixup. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</i> , pages 1644–1654.	1009		1010
954		1010		1011
955		1011		1012
956		1012		1013
957		1013		1014
958		1014		1015
959	Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In <i>Advances in Neural Information Processing Systems</i> , volume 32. Curran Associates, Inc.	1015	Dipendra Misra, Andrew Bennett, Valts Blukis, Eyvind Niklasson, Max Shatkhin, and Yoav Artzi. 2018. Mapping instructions to actions in 3d environments with visual goal prediction. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2667–2678.	1016
960		1016		1017
961		1017		1018
962		1018		1019
963		1019		1020
964	Chih-Yao Ma, Jiasen Lu, Zuxuan Wu, Ghassan AlRegib, Zsolt Kira, Richard Socher, and Caiming Xiong. 2019a. Self-monitoring navigation agent via auxiliary progress estimation.	1020	Abhinav Moudgil, Arjun Majumdar, Harsh Agrawal, Stefan Lee, and Dhruv Batra. 2021. Soat: A scene- and object-aware transformer for vision-and-language navigation. In <i>NeurIPS</i> .	1021
965		1021		1022
966		1022		1023
967		1023		

1024	Minecraft. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , Florence, Italy. Association for Computational Linguistics.	
1025		
1026		
1027		
1028	Khanh Nguyen, Yonatan Bisk, and Hal Daumé III au2.	
1029	2021. Learning when and what to ask: a hierarchical reinforcement learning framework.	
1030		
1031	Khanh Nguyen and Hal Daumé III. 2019. <a href="#">Help, anna! visual navigation with natural multimodal assistance via retrospective curiosity-encouraging imitation learning</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 684–695, Hong Kong, China. Association for Computational Linguistics.	
1032		
1033		
1034		
1035		
1036		
1037		
1038		
1039		
1040	Khanh Nguyen, Debadepta Dey, Chris Brockett, and Bill Dolan. 2019. Vision-based navigation with language-based assistance via imitation learning with indirect intervention. In <i>The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> .	
1041		
1042		
1043		
1044		
1045	Aishwarya Padmakumar, Jesse Thomason, Ayush Shrivastava, Patrick Lange, Anjali Narayan-Chen, Spandana Gella, Robinson Piramithu, Gokhan Tur, and Dilek Hakkani-Tur. 2021. <a href="#">Teach: Task-driven embodied agents that chat</a> .	
1046		
1047		
1048		
1049		
1050	Amin Parvaneh, Ehsan Abbasnejad, Damien Teney, Qinfeng Shi, and Anton van den Hengel. 2020. Counterfactual vision-and-language navigation: Unravelling the unseen. In <i>NeurIPS</i> .	
1051		
1052		
1053		
1054	Alexander Pashevich, Cordelia Schmid, and Chen Sun. 2021. Episodic transformer for vision-and-language navigation. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</i> , pages 15942–15952.	
1055		
1056		
1057		
1058		
1059	Yuankai Qi, Zizheng Pan, Yicong Hong, Ming-Hsuan Yang, Anton van den Hengel, and Qi Wu. 2021. The road to know-where: An object-and-room informed sequential bert for indoor vision-language navigation. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)</i> , pages 1655–1664.	
1060		
1061		
1062		
1063		
1064		
1065		
1066	Yuankai Qi, Zizheng Pan, Shengping Zhang, Anton van den Hengel, and Qi Wu. 2020a. Object-and-action aware model for visual language navigation. In <i>Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part X 16</i> , pages 303–317. Springer.	
1067		
1068		
1069		
1070		
1071		
1072	Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang, William Yang Wang, Chunhua Shen, and Anton van den Hengel. 2020b. Reverie: Remote embodied visual referring expression in real indoor environments. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> .	
1073		
1074		
1075		
1076		
1077		
1078	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.	
1079		
1080		
	Homero Roman Roman, Yonatan Bisk, Jesse Thomason, Asli Celikyilmaz, and Jianfeng Gao. 2020. <a href="#">RMM: A recursive mental model for dialog navigation</a> . In <i>Findings of Empirical Methods in Natural Language Processing (EMNLP Findings)</i> .	1081
		1082
		1083
		1084
		1085
	Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. 2020. ALFRED: A Benchmark for Interpreting Grounded Instructions for Everyday Tasks. In <i>The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> .	1086
		1087
		1088
		1089
		1090
		1091
	Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva, and Thomas Funkhouser. 2017. Semantic scene completion from a single depth image. <i>CVPR</i> .	1092
		1093
		1094
		1095
	Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J. Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, Anton Clarkson, Mingfei Yan, Brian Budge, Yajie Yan, Xiqiang Pan, June Yon, Yuyang Zou, Kimberly Leon, Nigel Carter, Jesus Briales, Tyler Gillingham, Elias Mueggler, Luis Pesqueira, Manolis Savva, Dhruv Batra, Hauke M. Strasdat, Renzo De Nardi, Michael Giesele, Steven Lovegrove, and Richard Newcombe. 2019. The Replica dataset: A digital replica of indoor spaces. <i>arXiv preprint arXiv:1906.05797</i> .	1096
		1097
		1098
		1099
		1100
		1101
		1102
		1103
		1104
		1105
		1106
	Alane Suhr, Claudia Yan, Jack Schluger, Stanley Yu, Hadi Khader, Marwa Mouallem, Iris Zhang, and Yoav Artzi. 2019. <a href="#">Executing instructions in situated collaborative interactions</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2119–2130, Hong Kong, China. Association for Computational Linguistics.	1107
		1108
		1109
		1110
		1111
		1112
		1113
		1114
		1115
	Q. Sun, Y. Zhuang, Z. Chen, Y. Fu, and X. Xue. 2021. <a href="#">Depth-guided adain and shift attention network for vision-and-language navigation</a> . In <i>2021 IEEE International Conference on Multimedia and Expo (ICME)</i> , pages 1–6, Los Alamitos, CA, USA. IEEE Computer Society.	1116
		1117
		1118
		1119
		1120
		1121
	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. 2021. Habitat 2.0: Training home assistants to rearrange their habitat. <i>arXiv preprint arXiv:2106.14405</i> .	1122
		1123
		1124
		1125
		1126
		1127
		1128
		1129
		1130
	Hao Tan, Licheng Yu, and Mohit Bansal. 2019. <a href="#">Learning to navigate unseen environments: Back translation with environmental dropout</a> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 2610–2621, Minneapolis, Minnesota. Association for Computational Linguistics.	1131
		1132
		1133
		1134
		1135
		1136
		1137
		1138

1139	Jesse Thomason, Daniel Gordon, and Yonatan Bisk.	Erik Wijmans, Abhishek Kadian, Ari Morcos, Stefan	1194
1140	2019a. <a href="#">Shifting the baseline: Single modality performance on visual navigation &amp; QA</a> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 1977–1983, Minneapolis, Minnesota. Association for Computational Linguistics.	Lee, Irfan Essa, Devi Parikh, Manolis Savva, and Dhruv Batra. 2019b. <a href="#">Dd-ppo: Learning near-perfect pointgoal navigators from 2.5 billion frames</a> . In <i>International Conference on Learning Representations</i> .	1195
1141			1196
1142			1197
1143			1198
1144		Yi Wu, Yuxin Wu, Georgia Gkioxari, and Yuandong Tian. 2018. <a href="#">Building generalizable agents with a realistic and rich 3d environment</a> .	1199
1145			1200
1146			1201
1147	Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer. 2019b. <a href="#">Vision-and-dialog navigation</a> . In <i>Conference on Robot Learning (CoRL)</i> .	Fei Xia, Amir R. Zamir, Zhiyang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. 2018. <a href="#">Gibson env: Real-world perception for embodied agents</a> . In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> .	1202
1148			1203
1149			1204
1150	Arun Balajee Vasudevan, Dengxin Dai, and Luc Van Gool. 2021. <a href="#">Talk2nav: Long-range vision-and-language navigation with dual attention and spatial memory</a> . <i>International Journal of Computer Vision</i> , 129(1):246–266.	Qiaolin Xia, Xiujun Li, Chunyuan Li, Yonatan Bisk, Zhifang Sui, Jianfeng Gao, Yejin Choi, and Noah A. Smith. 2020. <a href="#">Multi-view learning for vision-and-language navigation</a> .	1205
1151			1206
1152			1207
1153			1208
1154			1209
1155			1210
1156	Hanqing Wang, Wenguan Wang, Wei Liang, Caiming Xiong, and Jianbing Shen. 2021. <a href="#">Structured scene memory for vision-language navigation</a> . In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 8455–8464.	An Yan, Xin Eric Wang, Jiangtao Feng, Lei Li, and William Yang Wang. 2020. <a href="#">Cross-lingual vision-language navigation</a> .	1211
1157			1212
1158			1213
1159			1214
1160			1215
1161			1216
1162	Hanqing Wang, Wenguan Wang, Tianmin Shu, Wei Liang, and Jianbing Shen. 2020a. <a href="#">Active visual information gathering for vision-language navigation</a> .	Jiwen Zhang, Zhongyu Wei, Jianqing Fan, and Jiajie Peng. 2021. <a href="#">Curriculum learning for vision-and-language navigation</a> . In <i>NeurIPS</i> .	1217
1163			1218
1164			1219
1165			1220
1166			1221
1167			1222
1168			1223
1169			1224
1170			1225
1171			1226
1172			1227
1173			1228
1174			1229
1175			1230
1176			1231
1177			1232
1178			1233
1179			1234
1180			1235
1181			1236
1182			1237
1183			1238
1184			1239
1185			1240
1186			1241
1187			1242
1188			1243
1189			1244
1190			1245
1191			1246
1192			1247
1193			1248
			1249
			1249

1250 Wang Zhu, Hexiang Hu, Jiacheng Chen, Zhiwei Deng,  
1251 Vihan Jain, Eugene Ie, and Fei Sha. 2020b. Baby-  
1252 Walk: Going farther in vision-and-language naviga-  
1253 tion by taking baby steps. In *Proceedings of the*  
1254 *58th Annual Meeting of the Association for Computa-*  
1255 *tational Linguistics*, pages 2539–2556. Association  
1256 for Computational Linguistics.

1257 Wanrong Zhu, Yuankai Qi, Pradyumna Narayana, Ka-  
1258 zuo Sone, Sugato Basu, Xin Eric Wang, Qi Wu,  
1259 Miguel Eckstein, and William Yang Wang. 2021b.  
1260 [Diagnosing vision-and-language navigation: What](#)  
1261 [really matters.](#)

1262 Yi Zhu, Yue Weng, Fengda Zhu, Xiaodan Liang,  
1263 Qixiang Ye, Yutong Lu, and Jianbin Jiao. 2021c.  
1264 Self-motivated communication agent for real-world  
1265 vision-dialog navigation. In *Proceedings of the*  
1266 *IEEE/CVF International Conference on Computer*  
1267 *Vision (ICCV)*, pages 1594–1603.

1268 Yi Zhu, Fengda Zhu, Zhaohuan Zhan, Bingqian Lin,  
1269 Jianbin Jiao, Xiaojun Chang, and Xiaodan Liang.  
1270 2020c. Vision-dialog navigation by exploring cross-  
1271 modal memory. In *Proceedings of the IEEE/CVF*  
1272 *Conference on Computer Vision and Pattern Recog-*  
1273 *niton*, pages 10730–10739.

1274 Yuke Zhu, Roozbeh Mottaghi, Eric Kolve, Joseph J Lim,  
1275 Abhinav Gupta, Li Fei-Fei, and Ali Farhadi. 2017.  
1276 Target-driven visual navigation in indoor scenes us-  
1277 ing deep reinforcement learning. In *2017 IEEE in-*  
1278 *ternational conference on robotics and automation*  
1279 *(ICRA)*, pages 3357–3364. IEEE.

## A Dataset Details

Here in Table 2, we introduce more information about the datasets. Compared with the number of the datasets, the simulators are limited. More specifically, most indoor datasets are based on Matterport3D and most outdoor datasets are based on Google Street View. Also, more datasets are about indoor environments rather than outdoor environments. Outdoor environments are usually more complex and contain more objects compared with indoor environments.

## B Simulator

The virtual features of the dataset are deeply connected with the simulator in which datasets are built. Here we summarize simulators frequently used during the VLN dataset creation process.

House3D (Wu et al., 2018) is a realistic virtual 3D environment built based on the SUNCG (Song et al., 2017) dataset. An agent in the environment has access to first-person view RGB images, together with semantic/instance masks and depth information.

Matterport3D (Anderson et al., 2018b) simulator is a large-scale visual reinforcement learning simulation environment for research on embodied AI based on the Matterport3D dataset (Chang et al., 2017). Matterport3D contains various indoor scenes, including houses, apartments, hotels, offices, and churches. An agent can navigate between viewpoints along a pre-defined graph. Most indoors VLN datasets such as R2R and its variants are based on the Matterport3D simulator.

Habitat (Manolis Savva\* et al., 2019; Szot et al., 2021) is a 3D simulation platform for training embodied AI in 3D physics-enabled scenarios. Compared with other simulation environments, Habitat 2.0 (Szot et al., 2021) shows strength in system response speed. Habitat has the following datasets built-in: Matterport3D (Chang et al., 2017), Gibson (Xia et al., 2018), and Replica (Straub et al., 2019). AI2-THOR (Kolve et al., 2017) is a near photo-realistic 3D indoor simulation environment, where agents could navigate and interact with objects. Based on the object interaction function, it helps to build a dataset that requires object interaction, such as ALFRED (Shridhar et al., 2020).

Gibson (Xia et al., 2018) is a real-world perception interactive environment with complex semantics. Each viewpoint has a set of RGB panoramas with global camera poses and reconstructed 3D

meshes. Matterport3D dataset (Chang et al., 2017) is also integrated into the Gibson simulator.

House3D (Wu et al., 2018) converts SUNCG’s static environment into a virtual environment, where the agent can navigate with physical constraints (e.g. it cannot pass through walls or objects).

LANI (Misra et al., 2018) is a 3D simulator built in Unity3D platform. The environment in LANI is a fenced, square, grass field containing randomly placed landmarks. An agent needs to navigate between landmarks following the natural language instruction. Drone navigation tasks (Blukis et al., 2018, 2019) are also built based on LANI.

Currently, most datasets and simulators focus on indoors navigable scenes partly because of the difficulty of building an outdoor photo-realistic 3D simulator out of the increased complexity. Google Street View<sup>3</sup>, an online API that is integrated with Google Maps, is composed of billions of realistic street-level panoramas. It has been frequently used to create outdoor VLN tasks since the development of TOUCHDOWN (Chen et al., 2019).

## C Room-to-Room Leaderboard

Room-to-Room (R2R) (Anderson et al., 2018b) is the benchmark used most frequently for evaluating different methods. Here we collect all the reported performance metrics in the corresponding papers and the official R2R leaderboard<sup>4</sup>. Since beam search explores more routes, and since prior exploration has additional observations in the test environment, their performance can not be directly compared with other methods.

<sup>3</sup><https://developers.google.com/maps/documentation/streetview/overview>

<sup>4</sup><https://eval.ai/web/challenges/challenge-page/97/leaderboard/270>

Name	Simulator	Language-Active	Environment
Room-to-Room (Anderson et al., 2018b)	Matterport3D	✗	Indoor
Room-for-Room (Jain et al., 2019)	Matterport3D	✗	Indoor
Room-Across-Room (Ku et al., 2020)	Matterport3D	✗	Indoor
Landmark-RxR (He et al., 2021)	Matterport3D	✗	Indoor
XL-R2R (Yan et al., 2020)	Matterport3D	✗	Indoor
VLNCE (Krantz et al., 2020)	Habitat	✗	Indoor
StreetLearn (Mirowski et al., 2019)	Google Street View	✗	Outdoor
StreetNav (Hermann et al., 2020)	Google Street View	✗	Outdoor
TOUCHDOWN (Chen et al., 2019)	Google Street View	✗	Outdoor
Talk2Nav (Vasudevan et al., 2021)	Google Street View	✗	Outdoor
LANI (Misra et al., 2018)	-	✗	Outdoor
RoomNav (Wu et al., 2018)	House3D	✗	Indoor
REVERIE (Qi et al., 2020b)	Matterport3D	✗	Indoor
SOON (Zhu et al., 2021a)	Matterport3D	✗	Indoor
ALFRED (Shridhar et al., 2020)	AI2-THOR	✗	Indoor
VNLA (Nguyen et al., 2019)	Matterport3D	✓	Indoor
HANNA (Nguyen and Daumé III, 2019)	Matterport3D	✓	Indoor
CEREALBAR (Suhr et al., 2019)	-	✓	Indoor
Just Ask (Chi et al., 2020)	Matterport3D	✓	Indoor
CVDN (Thomason et al., 2019b)	Matterport3D	✓	Indoor
RobotSlang (Banerjee et al., 2020)	-	✓	Indoor
Talk the Walk (de Vries et al., 2018)	-	✓	Outdoor
MC Collab (Narayan-Chen et al., 2019)	Minecraft	✓	Outdoor
TEACh (Padmakumar et al., 2021)	AI2-THOR	✓	Indoor

Table 2: Vision-and-Language Navigation datasets. Language-Active means the agent needs to use natural language to request help, including both Guidance datasets and Dialog datasets in Table 1.

Simulator	Photo-realistic	3D
House3D (Wu et al., 2018)	✓	✓
Matterport3D (Chang et al., 2017)	✓	✓
Habitat (Manolis Savva* et al., 2019)	✓	✓
AI2-THOR (Kolve et al., 2017)	✗	✓
Gibson (Xia et al., 2018)	✓	✓
LANI (Misra et al., 2018)	✗	✓
*Google Street View	✓	✓

Table 3: Common simulators used to build VLN datasets. \*Google Street View is online API, providing similar functionality as a simulator for building VLN datasets.

Leader-Board (Test Unseen)	Single Run					Prior Exploration					Beam Search				
	TL↓	NE↓	OSR↑	SR↑	SPL↑	TL↓	NE↓	OSR↑	SR↑	SPL↑	TL↓	NE↓	OSR↑	SR↑	SPL↑
Models															
Random	9.89	9.79	0.18	0.13	0.12	-	-	-	-	-	-	-	-	-	-
Human	11.85	1.61	0.90	0.86	0.76	-	-	-	-	-	-	-	-	-	-
Seq-to-Seq (Anderson et al., 2018b)	8.13	20.4	0.27	0.20	0.18	-	-	-	-	-	-	-	-	-	-
RPA (Wang et al., 2018)	9.15	7.53	0.32	0.25	0.23	-	-	-	-	-	-	-	-	-	-
Speaker-Follower (Fried et al., 2018)	14.82	6.62	0.44	0.35	0.28	-	-	-	-	-	1257.38	4.87	0.96	0.54	0.01
Chasing Ghosts (Anderson et al., 2019a)	10.03	7.83	0.42	0.33	0.30	-	-	-	-	-	-	-	-	-	-
Self-Monitoring (Ma et al., 2019a)	18.04	5.67	0.59	0.48	0.35	-	-	-	-	-	373.1	4.48	0.97	0.61	0.02
RCM (Wang et al., 2019)	11.97	6.12	0.50	0.43	0.38	9.48	4.21	0.67	0.60	0.59	357.6	4.03	0.96	0.63	0.02
Regretful Agent (Ma et al., 2019b)	13.69	5.69	0.56	0.48	0.40	-	-	-	-	-	-	-	-	-	-
FAST (Ke et al., 2019)	22.08	5.14	0.64	0.54	0.41	-	-	-	-	-	196.5	4.29	0.90	0.61	0.03
ALTR (Huang et al., 2019)	10.27	5.49	0.56	0.48	0.45	-	-	-	-	-	-	-	-	-	-
EnvDrop (Tan et al., 2019)	11.66	5.23	0.59	0.51	0.47	9.79	3.97	0.70	0.64	0.61	686.8	3.26	0.99	0.69	0.01
PRESS (Li et al., 2019)	10.52	4.53	0.63	0.57	0.53	-	-	-	-	-	-	-	-	-	-
PTA (Landi et al., 2020)	10.17	6.17	0.47	0.40	0.36	-	-	-	-	-	-	-	-	-	-
EGP (Deng et al., 2020)	-	5.34	0.61	0.53	0.42	-	-	-	-	-	-	-	-	-	-
SERL (Wang et al., 2020b)	12.13	5.63	0.61	0.53	0.49	-	-	-	-	-	690.61	3.21	0.99	0.70	0.01
QAAM (Qi et al., 2020a)	10.40	-	0.61	0.53	0.50	-	-	-	-	-	-	-	-	-	-
CMG-AAL (Zhang et al., 2020a)	12.07	3.41	0.76	0.67	0.60	-	-	-	-	-	-	-	-	-	-
AuxRN (Zhu et al., 2020a)	-	5.15	0.62	0.55	0.51	10.43	3.69	0.75	0.68	0.65	40.85	3.24	0.81	0.71	0.21
RelGraph (Hong et al., 2020a)	10.29	4.75	0.61	0.55	0.52	-	-	-	-	-	-	-	-	-	-
PRRVALENT (Hao et al., 2020)	10.51	5.30	0.61	0.54	0.51	-	-	-	-	-	-	-	-	-	-
Active Exploration (Wang et al., 2020a)	21.03	4.34	0.71	0.60	0.43	9.85	3.30	0.77	0.70	0.68	176.2	3.07	0.94	0.70	0.05
VLN-BERT (Majumdar et al., 2020)	-	-	-	-	-	-	-	-	-	-	686.62	3.09	0.99	0.73	0.01
DASA (Sun et al., 2021)	10.06	5.11	-	0.54	0.52	-	-	-	-	-	-	-	-	-	-
ORIST (Qi et al., 2021)	11.31	5.10	-	0.57	0.52	-	-	-	-	-	-	-	-	-	-
NvEM (An et al., 2021)	12.98	4.37	0.66	0.58	0.54	-	-	-	-	-	-	-	-	-	-
SSM (Wang et al., 2021)	20.39	4.57	0.70	0.61	0.46	-	-	-	-	-	-	-	-	-	-
Recurrent VLN BERT (Hong et al., 2021)	12.35	4.09	0.70	0.63	0.57	-	-	-	-	-	-	-	-	-	-
SOAT (Moudgil et al., 2021)	12.26	-	4.49	58	53	-	-	-	-	-	-	-	-	-	-
REM (Liu et al., 2021)	13.11	3.87	0.72	0.65	0.59	-	-	-	-	-	-	-	-	-	-
HAMT(Chen et al., 2021b)	12.27	3.93	0.72	0.65	0.60	-	-	-	-	-	-	-	-	-	-
Spatial Route Prior (Zhou et al., 2021)	-	-	-	-	-	-	-	-	-	-	625.27	3.55	0.99	0.74	0.01
Airbert (Guhur et al., 2021)	-	-	-	-	-	-	-	-	-	-	686.54	2.58	0.99	0.78	0.01

Table 4: Leaderboard of Room-to-Room benchmark as of November, 2021.