

---

# A<sup>2</sup>Nav : Action-Aware Zero-Shot Robot Navigation Using Vision-Language Ability of Foundation Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We tackle the challenging task of zero-shot vision-and-language navigation (ZS-  
2 VLN), where an agent learns to follow complex path instructions without annotated  
3 data. We introduce A<sup>2</sup>Nav, an action-aware ZS-VLN method leveraging founda-  
4 tion models like GPT and CLIP. Our approach includes an instruction parser and an  
5 action-aware navigation policy. The parser breaks down complex instructions into  
6 action-aware sub-tasks, which are executed using the learned action-specific naviga-  
7 tion policy. Extensive experiments show A<sup>2</sup>Nav achieves promising ZS-VLN  
8 performance and even surpasses some supervised learning methods on R2R-Habitat  
9 and RxR-Habitat datasets.

## 10 1 Introduction

11 In vision-and-language navigation (VLN) tasks, an agent  
12 is required to navigate in a novel environment according to  
13 language navigation instructions. Current dominant meth-  
14 ods (6; 14; 21; 11) attempt to learn VLN ability in a su-  
15 pervised learning manner. However, creating high-quality  
16 labeled data requires a significant amount of human effort,  
17 which can be time-consuming and expensive. Addition-  
18 ally, the labeled data may not cover all possible scenarios,  
19 making it challenging for the model to generalize to new,  
20 unseen environments. To address these challenges, exploit-  
21 ing the knowledge from large foundation models (3; 18; 8)  
22 for learning navigation ability without requiring down-  
23 stream task annotated data is a potential solution. We call  
24 it **zero-shot navigation ability**.

25 Recently, researchers have made some attempts (10; 16;  
26 1; 25) at solving object navigation tasks in a zero-shot  
27 manner. They use a foundation vision-and-language model  
28 (VLM) (18) to localize the object (10) or use it to encode  
29 the object goal features (16), enabling the agent to navigate  
30 to any object goal described by natural language. Although  
31 some progress has been made, existing methods fail to take into account the varied action demands  
32 (e.g., “*proceed beyond*”, “*depart from*”) in navigation instructions. This may lead the agent to the

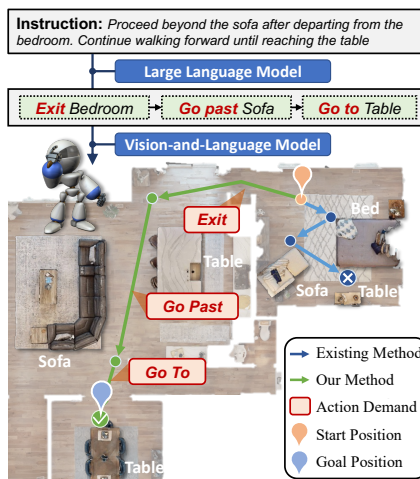


Figure 1: Existing zero-shot VLN overlooks the action demands.

33 wrong destination. For the example in Figure 1, the agent is expected to “*exist the bedroom*”, but the  
34 existing methods only take the landmark “*bedroom*” into consideration. The agent mistakenly goes  
35 into the bedroom, which is in the opposite direction of the path described by the instruction.

36 To solve this problem, the agent must correctly figure out the expected action demand associated with  
37 each landmark and accurately execute them. In this paper, we propose an action-aware navigation  
38 method, named A<sup>2</sup>Nav, for the zero-shot VLN task. Our method consists of two components: an  
39 instruction parser empowered by LLMs for figuring out landmarks and associated action demands;  
40 and an action-aware navigation policy empowered by CLIP for executing these action demands  
41 sequentially for navigation. Extensive experiments demonstrate that our A<sup>2</sup>Nav achieves promising  
42 performance on zero-shot VLN task, getting 22.6% and 16.8% success rates on R2R-Habitat and  
43 RxR-Habitat, respectively.

44 Our main contributions are as follows: 1) Instead of treating vision-and-language navigation as  
45 a sequence of object navigation tasks, we take into account the instruction action demands and  
46 decompose the instruction into a sequence of action-specific object navigation sub-tasks, where  
47 the agent is expected to not only localize the landmarks but also navigate to different goal position  
48 according to the associated action demand. 2) To address the problem that existing zero-shot  
49 navigators cannot satisfy different action demands, we identify and summarize five fundamental  
50 action demands and learn a unique navigator for executing each one without requiring manual  
51 path-instruction annotation, leading to more accurate and explainable navigation results.

## 52 2 Related Works

53 **Zero-Shot Object Navigation.** Since the navigation instruction is often described by several  
54 landmarks, it can be decomposed into sequential object navigation tasks. The object navigation  
55 task has been explored by previous literature (10; 16; 17; 5; 19; 1; 4; 24). Among these methods,  
56 we notice that some trails that design the object navigation agent in a zero-shot manner show great  
57 potential. Gadre *et al.* (10) design a heuristic algorithm to navigate to an object using the open-world  
58 object recognition ability of the foundation vision-and-language model (*i.e.*, CLIP (18)). Some works  
59 like ZER (1) and ZSON (16) learn an image navigation agent first, and then map the image goal  
60 representation into object text goal embedding space, and thus transfer to the object navigation task.

61 **Zero-Shot Vision-and-Language Navigation.** Based on the previous success on VLN and zero-shot  
62 object navigation, we aim to tackle the VLN task in the zero-shot manner, releasing the agent from  
63 expensive manual-labeled path-instruction training data. This problem has not been fully exploited  
64 yet. Pioneering works (22; 9; 7) have already verified the effectiveness of foundation models (LLM (3)  
65 and VLM (18)) in this scenario. These methods leverage GPT-3 (3) to extract navigation landmarks  
66 from the instruction and then initialize a heuristic object navigator using CLIP (18) to find out  
67 the landmark from visual observation and to navigate to the front of the landmark. Concurrent  
68 work (26) leverages a GPT model for inferring navigation actions on a discrete navigation graph. The  
69 performance in continuous environments has not been well explored. Our proposed A<sup>2</sup>Nav solved  
70 these issues using a learnable action-aware object navigator.

## 71 3 Action-Aware Zero-Shot VLN

72 We consider a practical but challenging problem zero-shot VLN, where the agent is expected to  
73 complete the VLN task without requiring path-instruction annotation. we leverage a large language  
74 model as an instruction parser for parsing all landmarks and their associated action demands. The  
75 instruction is then decomposed into a sequence of action-specific object navigation sub-tasks, in  
76 which the agent is required to localize the landmark and navigate based on the specific action demands  
77 associated with that landmark. For executing each sub-tasks sequentially, an action-aware navigation  
78 policy comprising five action-specific navigators is learned in a zero-shot manner. The general scheme  
79 is shown in Figure 3.

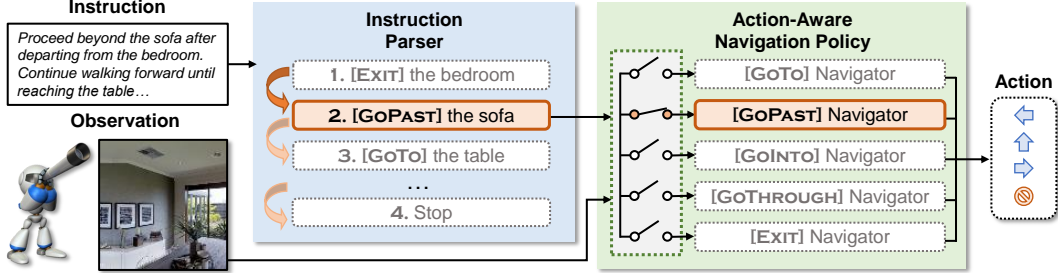


Figure 3: General scheme of A<sup>2</sup>Nav for zero-shot VLN task.

### 80 3.1 Instruction Parser

81 **Action-Specific Object Navigation Sub-Task.** The instruction parser aims to transfer a complicated instruction  
 82 into several sequential executable action-specific object navigation sub-tasks. Each sub-task contains a landmark  
 83 and an associated action demand, such as “*departing from the bedroom*”. The sub-task can be represented by a template  
 84 “(ACTION, LANDMARK)”. We explicitly summarize 5 basic sub-tasks shown in Figure 2, including “(GoTO,  
 85 OBJECT)”, “(GoPAST, OBJECT)”, and “(GoINTO, REGION)”, “(GoTHROUGH, REGION)”, “(EXIT, REGION)”.  
 86  
 87  
 88  
 89  
 90

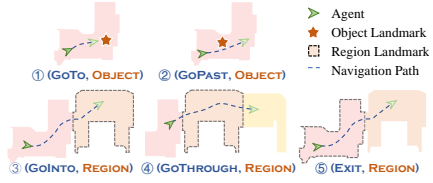


Figure 2: Illustration of sub-task types.

91 **Decomposing Instruction into Sub-Tasks.** We use the few-shot learning ability of GPT-3 LLM (3)  
 92 for decomposing an instruction into a sequence of sub-tasks described above. The prompt contains  
 93 several correct instruction parsing examples. As the predicted sub-tasks from GPT-3 are in the  
 94 free-form language, we need to map each prediction to the predefined sub-tasks. In most cases, the  
 95 “ACTION” predictions made by GPT-3 accurately match one of the “ACTION” in predefined sub-tasks,  
 96 and thus we can directly map it to this sub-task type. In cases where a prediction does not match any,  
 97 we follow (13) to perform mapping through semantic translation. Specifically, we use BERT (8) to  
 98 encode the predicted “ACTION” and the “ACTION” in all predefined sub-tasks. Then we compute the  
 99 cosine similarity between them and consider the predefined sub-task with the highest score as the  
 100 predicted sub-task.

### 101 3.2 Action-Aware Navigation Policy

102 With the sub-task sequence parsed by LLMs, we learn an action-aware navigation policy to execute them sequentially.  
 103 The policy consists of five action-specific navigators, each of which is responsible for a specific sub-task type.  
 104 We follow ZSON (16) to transform the question of learning such a navigator into learning an image-goal navigator  
 105 on a freely collected action-specific image-path dataset. Each sample in this dataset contains an image sampled  
 106 from the environment and a navigation path that is compatible with the action demands. Specifically, for training  
 107 **GOPAST** navigator that expects the agent to go to the object and keep going forward past the object,  
 108 we capture the goal image in the middle of the path. For the **GOINTO** action demand that expects the agent to go  
 109 cross a doorway into the target region, we sample the path that crosses over two regions and sample the goal image at the end of this path. For the **GO**  
 110 **THROUGH** action demand that expects the agent to go from one side to the other side of a region, we randomly sample the path  
 111 that starts near one entrance and ends near the other one of a region. The goal image is captured in  
 112 the middle of the path. For the **EXIT** action demand, the path is sampled in the same way as the  
 113 **GOINTO** action demand, while the goal image is captured at the beginning of the path. We fine-tune  
 114 the trained ZSON model on these datasets using the same learning pipeline as ZSON, which is shown  
 115 in Figure 4. For the **GoTO** action demands, we directly utilize a trained ZSON model as a navigator.

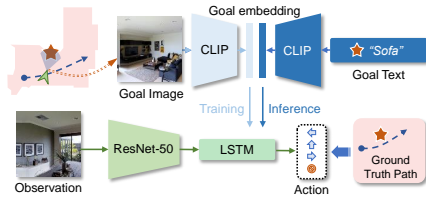


Figure 4: Paradigm of navigators.

116  
 117  
 118  
 119  
 120  
 121

	Method	Extra Info.	R2R-Habitat		RxR-Habitat		CSR
			SR	SPL	SR	SPL	
Supervised	Seq2Seq (14)	Depth	25.0%	22.0%	-	-	-
	LAW (21)	Depth	35.0%	31.0%	10.0%	9.0%	28.6%
	WS-MGMap (6)	Depth	38.9%	34.3%	15.0%	12.1%	38.6%
Zero-Shot	Random	-	0.0%	0.0%	6.0%	6.0%	0.0%
	CLIP-Nav (9)	Panoramic	5.6%	2.9%	9.8%	3.2%	57.4%
	Seq CLIP-Nav (9)	Panoramic	7.1%	3.7%	9.1%	3.3%	77.8%
	Cow (10)	Depth	7.8%	5.8%	7.9%	6.1%	<b>98.3%</b>
	ZSON (16)	-	19.3%	9.3%	14.2%	4.8%	73.6%
	<b>A<sup>2</sup>Nav (Ours)</b>	-	<b>22.6%</b>	<b>11.1%</b>	<b>16.8%</b>	<b>6.3%</b>	74.3%

Table 1: Comparisons with zero-shot and supervised methods on VLN datasets.

## 122 4 Experiments

123 We compare our A<sup>2</sup>Nav with existing zero-shot and supervised-learning navigation methods on  
 124 R2R-Habitat and RxR-Habitat datasets. We introduce the experimental setup, agent configurations,  
 125 and baselines in Appendix.

126 **Comparisons with Zero-Shot Methods.** In Table 1, our A<sup>2</sup>Nav outperforms other zero-shot methods.  
 127 On R2R-Habitat, it surpasses CLIP-Nav, Seq CLIP-Nav, CoW, and ZSON by 17.0%, 15.5%, 14.8%,  
 128 and 3.3% in success rate, respectively. On RxR-Habitat, it outperforms them by 7.0%, 7.7%, 8.9%,  
 129 and 2.6%, respectively.

130 **Comparisons with Supervised Learning Methods.** We  
 131 compare our zero-shot A<sup>2</sup>Nav with three supervised VLN  
 132 methods: vanilla Seq2Seq (15), LAW (21), and WS-  
 133 MGMap (6). In Table 1, our zero-shot A<sup>2</sup>Nav achieves  
 134 comparable performance compared with the vanilla  
 135 Seq2Seq on R2R-Habitat and outperforms all super-  
 136 vised learning methods on RxR-Habitat, indicating that  
 137 A<sup>2</sup>Nav is more effective at generalizing to different  
 138 datasets and can adapt more easily to varying environ-  
 139 ments. We also compare A<sup>2</sup>Nav with supervised methods  
 140 trained on limited data. In Figure 5, A<sup>2</sup>Nav outperforms  
 141 the SOTA method (*i.e.*, WS-MGMap) if less than 50%  
 142 training episodes are available for it.

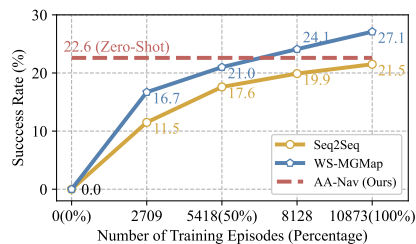


Figure 5: Comparison with the supervised learning methods that are trained on partial training data.

143 **Effectiveness of Action-Aware Navigation Policy.** To verify the effectiveness of each navigator,  
 144 we create multiple navigation policy variants that progressively include 5 navigators in an order of  
 145 GoTo, GoPAST, GoINTO, GoTHROUGH, and EXIT. By default, the sub-task is executed by the  
 146 GoTo navigator if its corresponding navigator is not included. In R2R-Habitat, the SR of policy  
 147 with 1~5 navigators are 19.3%, 20.9%, 21.5%, 22.3%, and 22.6%, respectively, demonstrating the  
 148 importance of each navigator.

## 149 5 Conclusion

150 In this paper, we take into account the instruction action demands and decompose the VLN task into  
 151 a sequence of action-specific object navigation sub-tasks. To execute these sub-tasks, we further  
 152 propose an action-aware navigation policy that learns different navigation abilities without requiring  
 153 any manual path-instruction annotation. The proposed A<sup>2</sup>Nav achieves the best zero-shot VLN  
 154 performance on two benchmark datasets (*i.e.*, R2R-Habitat and RxR-Habitat) and outperforms the  
 155 state-of-the-art supervised learning methods on RxR-Habitat. Furthermore, our A<sup>2</sup>Nav is able to  
 156 more accurately follow navigation instructions that contain specific action demands, demonstrating  
 157 its potential for the scenario that needs human-robot communication and interaction.

## References

- [1] Ziad Al-Halah, Santhosh K. Ramakrishnan, and Kristen Grauman. Zero experience required: Plug & play modular transfer learning for semantic visual navigation. In *CVPR*, pages 17010–17020, 2022.
- [2] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian D. Reid, Stephen Gould, and Anton van den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *CVPR*, pages 3674–3683, 2018.
- [3] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *NeurIPS*, 2020.
- [4] Devendra Singh Chaplot, Dhiraj Gandhi, Abhinav Gupta, and Ruslan Salakhutdinov. Object goal navigation using goal-oriented semantic exploration. In *NeurIPS*, 2020.
- [5] Peihao Chen, Dongyu Ji, Kunyang Lin, Weiwen Hu, Wenbing Huang, Thomas H. Li, Mingkui Tan, and Chuang Gan. Learning active camera for multi-object navigation. *NeurIPS*, 2022.
- [6] 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. In *NeurIPS*, 2022.
- [7] Zhenfang Chen, Qinhong Zhou, Yikang Shen, Yining Hong, Hao Zhang, and Chuang Gan. See, think, confirm: Interactive prompting between vision and language models for knowledge-based visual reasoning. *arXiv*, 2023.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, pages 4171–4186, 2019.
- [9] Vishnu Sashank Dorbala, Gunnar A. Sigurdsson, Robinson Piramuthu, Jesse Thomason, and Gaurav S. Sukhatme. Clip-nav: Using CLIP for zero-shot vision-and-language navigation. *CoRR*, abs/2211.16649, 2022.
- [10] Samir Yitzhak Gadre, Mitchell Wortsman, Gabriel Ilharco, Ludwig Schmidt, and Shuran Song. CLIP on wheels: Zero-shot object navigation as object localization and exploration. *CoRR*, abs/2203.10421, 2022.
- [11] Georgios Georgakis, Karl Schmeckpeper, Karan Wanchoo, Soham Dan, Eleni Miltsakaki, Dan Roth, and Kostas Daniilidis. Cross-modal map learning for vision and language navigation. In *CVPR*, pages 15439–15449, 2022.
- [12] Yicong Hong, Cristian Rodriguez Opazo, Qi Wu, and Stephen Gould. Sub-instruction aware vision-and-language navigation. In *EMNLP*, pages 3360–3376, 2020.
- [13] Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *ICML*, pages 9118–9147, 2022.
- [14] Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In *ECCV*, pages 104–120, 2020.
- [15] Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In *ECCV*, pages 104–120, 2020.
- [16] Arjun Majumdar, Gunjan Aggarwal, Bhavika Devnani, Judy Hoffman, and Dhruv Batra. ZSON: zero-shot object-goal navigation using multimodal goal embeddings. In *NeurIPS*, 2022.
- [17] So Yeon Min, Yao-Hung Hubert Tsai, Wei Ding, Ali Farhadi, Ruslan Salakhutdinov, Yonatan Bisk, and Jian Zhang. Object goal navigation with end-to-end self-supervision. *CoRR*, abs/2212.05923, 2022.
- [18] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya

- 213 Sutskever. Learning transferable visual models from natural language supervision. In *ICML*,  
214 pages 8748–8763, 2021.
- 215 [19] Santhosh Kumar Ramakrishnan, Devendra Singh Chaplot, Ziad Al-Halah, Jitendra Malik, and  
216 Kristen Grauman. PONI: potential functions for objectgoal navigation with interaction-free  
217 learning. In *CVPR*, pages 18868–18878, 2022.
- 218 [20] Santhosh Kumar Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alexan-  
219 der Clegg, John Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X. Chang,  
220 Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-matterport 3d dataset (HM3D): 1000 large-  
221 scale 3d environments for embodied AI. In *NeurIPS Datasets and Benchmarks*, 2021.
- 222 [21] Sonia Raychaudhuri, Saim Wani, Shivansh Patel, Unnat Jain, and Angel X. Chang. Language-  
223 aligned waypoint (LAW) supervision for vision-and-language navigation in continuous environ-  
224 ments. In *EMNLP*, pages 4018–4028, 2021.
- 225 [22] Dhruv Shah, Blazej Osinski, Brian Ichter, and Sergey Levine. Lm-nav: Robotic navigation with  
226 large pre-trained models of language, vision, and action. volume abs/2207.04429, 2022.
- 227 [23] Erik Wijmans, Abhishek Kadian, Ari Morcos, Stefan Lee, Irfan Essa, Devi Parikh, Manolis  
228 Savva, and Dhruv Batra. DD-PPO: learning near-perfect pointgoal navigators from 2.5 billion  
229 frames. In *ICLR*, 2020.
- 230 [24] Karmesh Yadav, Ram Ramrakhya, Arjun Majumdar, Vincent-Pierre Berges, Sachit Kuhar,  
231 Dhruv Batra, Alexei Baevski, and Oleksandr Maksymets. Offline visual representation learning  
232 for embodied navigation. *CoRR*, abs/2204.13226, 2022.
- 233 [25] Hongxin Zhang, Weihua Du, Jiaming Shan, Qinhong Zhou, Yilun Du, Joshua B Tenenbaum,  
234 Tianmin Shu, and Chuang Gan. Building cooperative embodied agents modularly with large  
235 language models. *arXiv preprint arXiv:2307.02485*, 2023.
- 236 [26] Gengze Zhou, Yicong Hong, and Qi Wu. Navgpt: Explicit reasoning in vision-and-language  
237 navigation with large language models. *arXiv preprint arXiv:2305.16986*, 2023.

# APPENDIX

238 In the supplementary, we provide more implementation details and visualization results of our method.  
239 We organize the supplementary as follows.

- 240 • In Section A, we present the implementation details, datasets, metrics, and baselines.
- 241 • In Section B, we present more details on action-specific image-path dataset collection.
- 242 • In Section C, we present more details on action-specific navigator training and inference.
- 243 • In Section D, we present more details on prompt design for the instruction parser.
- 244 • In Section E, we present more details on zero-shot navigation baselines.
- 245 • In Section F, we present more visualization results.

## 246 A Experimental Setup

247 **Evaluation Datasets and Metrics.** We conduct experiments on the validation unseen split of three  
248 datasets, namely R2R-Habitat (15), RxR-Habitat (15), and Fine-Grained R2R (FG-R2R) (12). These  
249 three datasets contain 1,839, 1,079, and 1,839 validation episodes on 11 scenes in Matterport3D,  
250 respectively. RxR-Habitat contains instructions in three languages, and we only use the English  
251 split in our experiments. FG-R2R is an extension of R2R (2), where instructions are chunked into  
252 several sub-instructions and each sub-instruction is labeled with a corresponding sub-path, resulting  
253 in 6,687 sub-instruction-sub-path pairs. Following the existing works (14; 21; 6), we evaluate  
254 VLN performance using Success Rate (**SR**) and Success weighted by inverse Path Length (**SPL**).  
255 Besides, to evaluate the generalization ability among datasets, we follow Dorbala *et al.* (9) to propose  
256 Consistency on SR (**CSR**) for computing the relative change in SR among datasets. Specifically,  
257  $CSR = 1 - \frac{|SR_a - SR_b|}{\max\{SR_a, SR_b\}} \times 100\%$ , where  $SR_a$  and  $SR_b$  are success rates for different datasets.

258 **Agent Configurations.** Following ZSON (16), the agent has a height of 1.25m, with a radius of  
259 0.1m. It is equipped with one  $128 \times 128$  RGB sensor with  $90^\circ$  horizontal field of view. The agent  
260 can execute four low-level actions, namely STOP indicating the end of an episode, FORWARD that  
261 moves itself forward by 0.25 meters and TURNLEFT and TURNRIGHT that turn itself by  $30^\circ$ . The  
262  $m_s$  is empirically set to 100 and 50 for R2R-Habitat and RxR-Habitat, respectively. The  $m_e$  is set to  
263 500 for all three datasets following exiting works (6; 21).

264 **Baselines.** We decompose the instruction into an object navigation sub-task sequence using GPT-3  
265 and execute these sub-tasks sequentially using four zero-shot object navigation methods.

- 266 • **CLIP-Nav** (9) is designed for navigating among discrete navigable nodes. The agent uses CLIP to  
267 determine which adjacent node has the highest possibility of containing landmarks and then moves  
268 to this node. We adapt it to continuous environments using a waypoint navigation algorithm. More  
269 details can be found in Appendix.
- 270 • **Seq CLIP-Nav** (9) is an extended version of CLIP-Nav with an additional backtracking mechanism,  
271 which allows the agent to go back to the previous location if it cannot find the landmarks for several  
272 steps.
- 273 • **CoW** (10) uses CLIP gradient for object localization and a path-planning algorithm for action  
274 determination.
- 275 • **ZSON** (16) uses the CLIP for encoding both image and landmark text to the same semantic feature  
276 space and then trains an image navigator for object navigation.

## 277 B More Details on Action-Specific Image-Path Dataset Collection

278 For learning a navigator for executing each action demand, we need to collect an action-specific  
279 image-path dataset for fine-tuning a trained ZSON model. In Section 3.3.2 of the paper, we have

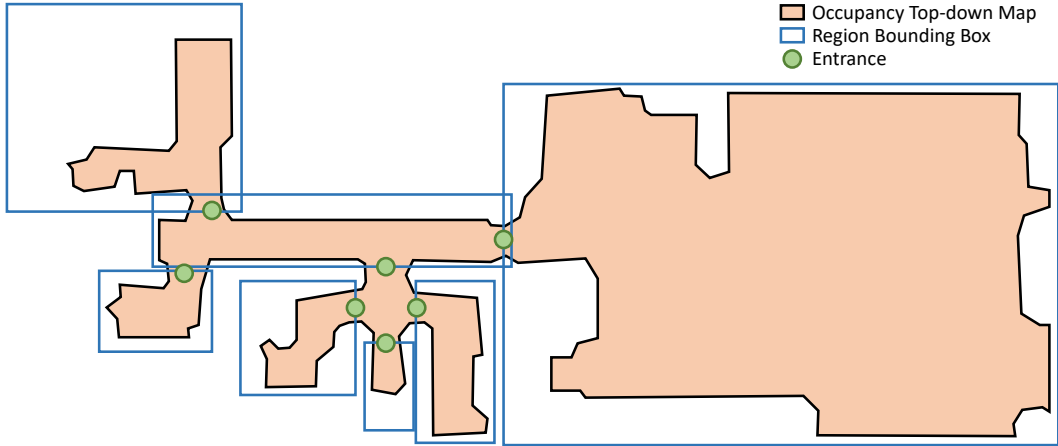


Figure 6: Obtaining entrance positions from the intersections between regions and top-down map.

introduced the basic principle for collecting the episodes (*i.e.*, the path and the corresponding goal image) in the dataset. In this section, we present more data collection details.

- **GOPAST Dataset.** We randomly sample two points whose geometric distance is longer than 1.5m, and consider the shortest navigation path as the ground truth path. The goal image is sampled in the middle of the path facing the direction of the agent’s advancement. We introduce some randomness to the angle by jittering it by  $\pm 45^\circ$ .
- **GOINTO Dataset.** We randomly choose a region from the scene. The start point is sampled near the entrance of this region (the geometric distance is less than 1.5m). The goal point is sampled randomly inside this region. The goal image is taken in a random direction at the goal point.
- **GOTROUGH Dataset.** We randomly select a region with two different entrances and sample a random point near each entrance respectively to form a path. The geometric distance from the start or end point to the entrance is less than 1.5m. Goal image is taken in the middle of the path and faces the direction of the agent’s advancement.
- **EXIT Dataset.** The ground truth path of the EXIT action is similar to the GOINTO action besides switching the position of start point and goal point.

We utilize room region bounding box annotations for obtaining the entrance position of regions. Specifically, we consider the intersection between the room region bounding box and the occupancy top-down map as the entrance of a region. An example is shown in Figure 6. Since collecting GOINTO, GOTROUGH, and EXIT datasets requires the entrance position, we collect these three datasets from 131 scenes in HM3D (20) dataset that have region bounding box annotations. For the GOPAST dataset, we collect from all 800 scenes in the train split of HM3D. We sample 9,000 image-path pairs from each scene.

## C More Details on Action-Specific Navigator Training and Inference

We fine-tune a trained ZSON model on the action-specific dataset for learning an action-specific navigator. We use the ZSON model that is trained on agent configuration A described by Majumdar *et al.* (16), *i.e.*, the agent has a height of 1.25m, with a radius of 0.1m and is equipped with one  $128 \times 128$  RGB sensor with  $90^\circ$  horizontal field of view. We fine-tune this model using a reinforcement learning algorithm (*i.e.*, DD-PPO (23)) for 100M steps with the same navigation reward as ZSON. This reward encourages the agent to go to the end of the path in an episode and look toward the goal image:

$$r_t = r_{\text{success}} + r_{\text{angle-success}} - \Delta_{\text{dtg}} - \Delta_{\text{atg}} + r_{\text{slack}} \quad (1)$$



310 where  $r_{\text{success}} = 5$  if STOP is predicted when the agent is within 1m of the goal position,  
 311  $r_{\text{angle-success}} = 5$  if the agent is within 1m of the goal position and within  $25^\circ$  of the goal ori-  
 312 entation (and 0 otherwise). Besides,  $\Delta_{\text{dtg}}$  is the change in the agent’s distance-to-goal, and  $\Delta_{\text{atg}}$  is  
 313 the change in the agent’s angle-to-goal.  $\Delta_{\text{atg}}$  is set to 0 if the agent is not within a circle of 1m radius  
 314 from the goal position. We also use a slack reward  $r_{\text{slack}} = -0.01$  to encourage the agent to reach  
 315 the goal as soon as possible.

Q: You are a robot walking in a house. You should parse the navigation instruction into several subtasks and then execute them one by one.

Each subtask consists of an {action} and a {landmark}. Action should be chosen from {turn left}, {turn right}, {go to}, {go past}, {go into}, {go through}, and {exit}. Landmark should be {a specific object} or {a region}. Here is the definition of each action:  
 {go to} means go to the front of {a specific object};  
 {go past} means go to {a specific object} and then go pass it;  
 {go into} means go into {a region};  
 {go through} means walk along {a region} or through {a region};  
 {exit} means find a door and go out of {a region}.

Now help me parse the following instruction: "Exit the bedroom and turn left. Walk straight passing the gray couch and stop near the rug. "

A: Let's think step by step:  
 Subtask 1: {exit} the {bedroom};  
 Subtask 2: {turn left};  
 Subtask 3: {go past} the {gray couch};  
 Subtask 4: {go to} the {rug}.

Figure 7: An example of prompt design 1: a brief description of sub-task.

Q: Parse the instruction using the following subtasks: 1. {go to} [landmark], 2. {go past} [landmark], 3. {turn left}, 4. {turn right}, 5. {go through} [region], 6. {go into} [region], 7. {exit} [region], 8. {stop}.

Here are several examples.  
 (1) Instruction: "Walk into the hallway and turn left. Walk to the left of the railing and across the hall past the plant. Stop to the left of the stairs.", and the subtasks should be:  
 1. {go into} {the hallway}.  
 2. {turn left}.  
 3. {go to} {the left of the railing}.  
 4. {go past} {the plant}.  
 5. {go to} {the left of the stairs}.  
 6. {stop}.  
 .....(with another 4 examples)

Now help me parse the following instruction: "Exit the bedroom and turn left. Walk straight passing the gray couch and stop near the rug. "

A: Let's think step by step:  
 1. {exit} {the bedroom}.  
 2. {turn left}.  
 3. {go past} {the gray couch}.  
 4. {go to} {the rug}.  
 5. {stop}.

Figure 8: An example of prompt design 2: a collection of parsing examples. This prompt design performs the best.

316 After fine-tuning, we use the trained navigator for executing a sub-task. Specifically, we feed the  
 317 landmark (*i.e.*, text description of an object or a region) to the CLIP for extracting goal embedding.  
 318 The navigator take the current RGB observation and the goal embedding as input for predicting a  
 319 low-level action for this sub-task.

## 320 D More Details on Prompt Design for the Instruction Parser

321 We have tried three prompt designs for parsing instructions using the large language model GPT-3.

- 322 • **Prompt Design 1:** a brief description of each sub-task definition.
- 323 • **Prompt Design 2:** a collection of instruction parsing examples

Q: You are a robot walking in a house. You should parse the navigation instruction into several subtasks and then execute them one by one.

Each subtask consists of an [action] and a {landmark}. Action should be chosen from [turn left], [turn right], [go to], [go past], [go into], [go through], and [exit]. Landmark should be {a specific object} or {a region}. Here is the definition of each action:  
 [go to] means go to the front of {a specific object};  
 [go past] means go to {a specific object} and then go pass it;  
 [go into] means go into {a region};  
 [go through] means walk along {a region} or through {a region};  
 [exit] means find a door and go out of {a region}.

Here are several examples.  
 (1) Instruction: "Walk into the hallway and turn left. Walk to the left of the railing and across the hall past the plant. Stop to the left of the stairs.", and the subtasks should be:  
 1. [go into] {the hallway}.  
 2. [turn left].  
 3. [go to] {the left of the railing}.  
 4. [go past] {the plant}.  
 5. [go to] {the left of the stairs}.  
 6. [stop].  
 .....(with another 4 examples)

Now help me parse the following instruction: " Exit the bedroom and turn left. Walk straight passing the gray couch and stop near the rug. "

A: Let's think step by step :  
 1. [exit] {the bedroom}.  
 2. [turn left].  
 3. [go past] {the gray couch}.  
 4. [go to] {the rug}.  
 5. [stop].

Figure 9: An example of prompt design 3: a combination of both sub-task definition description and examples.

324 • **Prompt Design 3:** a combination of both brief description and examples

325 Experimental results in Table 4 in the paper show that the second prompt design performs the best.  
 326 We show the examples of these prompt designs in Figures 7, 8 and 9, respectively. We mark the  
 327 GPT-3 output in brown color.

## 328 E More Details on Zero-Shot Navigation Baselines

329 We decompose the instruction into an object navigation sub-task sequence using GPT-3 and execute  
 330 these sub-tasks sequentially using four zero-shot object navigation methods.

- 331 • **CLIP-Nav (9)** is designed for navigating among discrete navigable nodes. The agent uses CLIP  
 332 to determine which adjacent node has the highest possibility of containing landmarks and then  
 333 moves to this node. To adapt it to continuous environments, we capture 4 RGB images uniformly  
 334 in different directions and use CLIP (18) to select one image that has the highest possibility  
 335 of containing landmarks. Then, we randomly set a waypoint in that direction and use a path-  
 336 planing algorithm to plan low-level actions for navigating to the waypoint. If the CLIP softmax  
 337 score is higher than the threshold of 0.8, we switch to the next object navigation sub-task. For  
 338 implementation convenience, we use the “*shortest\_path\_follower*” API in the Habitat simulator for  
 339 path planning, which assumes the complete occupancy top-down map is available.
- 340 • **Seq CLIP-Nav (9)** is an extended version of CLIP-Nav with an additional backtracking mechanism,  
 341 which allows the agent to go back to the previous location if it cannot find the landmarks for several  
 342 steps. In our implementation, we directly set the agent to the position 15 step before for performing  
 343 backtracking.
- 344 • **CoW (10)** uses CLIP gradient for object localization and a path-planning algorithm for action  
 345 determination. For implementation convenience, we use the “*shortest\_path\_follower*” API in the  
 346 Habitat simulator for path planning, which assumes the complete occupancy top-down map is  
 347 available. Even using the oracle occupancy information, our A<sup>2</sup>Nav still performs better than this  
 348 baseline.

- 349 • **ZSON (16)** uses the CLIP for encoding both image and landmark text to the same semantic feature  
350 space and then trains an image navigator for object navigation. We use the model trained on the  
351 HM3D dataset using the config A setting.

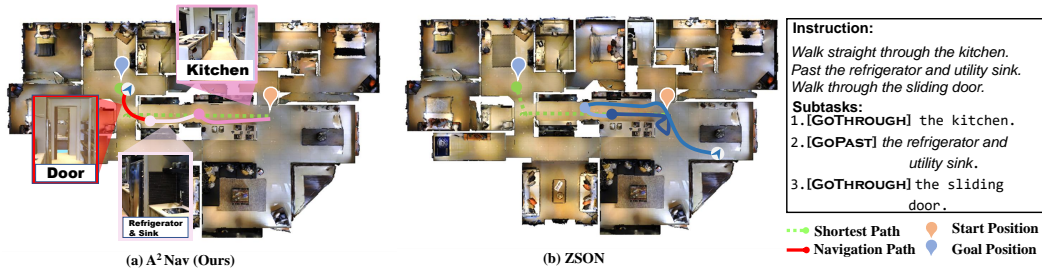


Figure 10: Visualization of the navigation path. Our method successfully goes through the kitchen, while the baseline fails to do it.



Figure 11: Visualization of the navigation path. Our method successfully exits the bedroom and goes past the stair, while the baseline is stuck in the bedroom.



Figure 12: Visualization of the navigation path. Our method successfully goes through the kitchen and finds the refrigerator, while the baseline fails to do it.

## 352 F More Visualization Results

353 In this section, we provide more visualization examples for comparing the method between ZSON (16)  
354 and ours. In Figure 10, the instruction requires the agent to go across the kitchen and exit

355 this area through the door. Our A<sup>2</sup>Nav successfully follows the instruction because of the learned  
356 “GOTROUGH” ability, which leads the agent to completely go through the area. However, ZSON (16)  
357 just goes to the kitchen area of a refrigerator, which directly causes the task to fail. In Figure 11 the  
358 instruction requires the agent to go past the stair which is outside the bedroom. Our A<sup>2</sup>Nav success-  
359 fully exits the bedroom, navigates past the stair and then stop at the correct doorway. In contrast, the  
360 ZSON model fails to exit the bedroom and finally stop at the doorway of the bedroom incorrectly. In  
361 Figure 12, the instruction requires the agent to get out of the kitchen and stop near the refrigerator.  
362 Our A<sup>2</sup>Nav successfully walks across the kitchen and goes by the hallway, finally finding the target.  
363 ZSON (16) tries to go to the area which suggests the higher confidence of the kitchen, which is the  
364 opposite direction of the shortest path to the target. All examples demonstrate the effectiveness of our  
365 action-aware agent.