

Semantic Map-based Generation of Navigation Instructions

Chengzu Li^{1,2,*}, Chao Zhang², Simone Teufel¹,
Rama Sanand Doddipatla², Svetlana Stoyanchev²

¹University of Cambridge, ²Toshiba Europe Limited

{cl917, sht25}@cam.ac.uk

{chao.zhang, rama.doddipatla, svetlana.stoyanchev}@toshiba.eu

Abstract

We are interested in the generation of navigation instructions, either in their own right or as training material for robotic navigation task. In this paper, we propose a new approach to navigation instruction generation by framing the problem as an image captioning task using semantic maps as visual input. Conventional approaches employ a sequence of panorama images to generate navigation instructions. Semantic maps abstract away from visual details and fuse the information in multiple panorama images into a single top-down representation, thereby reducing computational complexity to process the input. We present a benchmark dataset for instruction generation using semantic maps, propose an initial model and ask human subjects to manually assess the quality of generated instructions. Our initial investigations show promise in using semantic maps for instruction generation instead of a sequence of panorama images, but there is vast scope for improvement. We release the code for data preparation and model training at <https://github.com/chengzu-li/VLGen>.

Keywords: semantic map, navigation instruction generation, Room2Room

1. Introduction

Vision and Language Navigation (VLN) is a task that involves an agent navigating in a physical environment in response to natural language instructions (Wu et al., 2021). The data annotation for the VLN task is time-consuming and costly to scale up, and the development of models that address the task is severely limited by the availability of training data (Gu et al., 2022). Navigation instruction generation (VL-GEN) is the reverse of the VLN task in that it generates natural language instructions for a path in the virtual (or physical) environment, which is helpful for interactions with users and explainability. Previous work has also demonstrated the effectiveness of VL-GEN in improving the performance of VLN systems such as the Speaker-Follower model (Fried et al., 2018) and Env Drop (Tan et al., 2019). This paper explores the VL-GEN task of generating navigation instruction framing it as an image captioning task.

VL-GEN requires the model to generate language instruction in the context of the physical environment, grounding objects references and action instructions to the given space. Previous studies use photo-realistic RGB panoramic images as the visual input; they frame VL-GEN as the end-to-end task of generating text from a sequence of photo-realistic RGB images (Fried et al., 2018; Tan et al., 2019; Wang et al., 2022d). While Zhao et al. (2021) report that the overall quality of instructions generated with end-to-end models is only slightly better than that of template-based generation, the application of object grounding to the panoramic images

achieves a better result (Wang et al., 2022d).

The existing approach to this task has two shortcomings. From the perspective of representation, using panoramic images is resource-intensive as it requires processing of multiple image inputs corresponding to different points on the path. Second, panoramic images contain many details that are irrelevant for the task. The model has to learn to interpret the environments from RGB panoramas, such as object recognition, and generate instructions at the same time. As it is natural for humans to understand navigation instructions from a top-down map (as in Google Maps) (Paz-Argaman et al., 2024), we propose to separate the VL-GEN task into two steps: 1) environment interpretation, which is addressed by semantic SLAM in physical robotic systems (Chaplot et al., 2020), and 2) spatial reasoning. In this paper, we focus on the second step and explore the feasibility of using top-down semantic map for VL-GEN.

Our research question is whether it is feasible to use the top-down semantic map (a single RGB image) as our main source of information. We also explore which other data sources, in addition to the semantic map, can further improve performance. To address this question, we formalize the VL-GEN task as image captioning with the input of a semantic map with the path (see Figure 1). We extract the images of top-down maps from the Habitat simulator (Savva et al., 2019) based on Room-to-Room dataset (Anderson et al., 2018) and VLN-CE (Krantz et al., 2020). Our key contributions and findings include the following:

- We extend the R2R dataset with semantic

*Work done on internship at Toshiba Europe Limited.

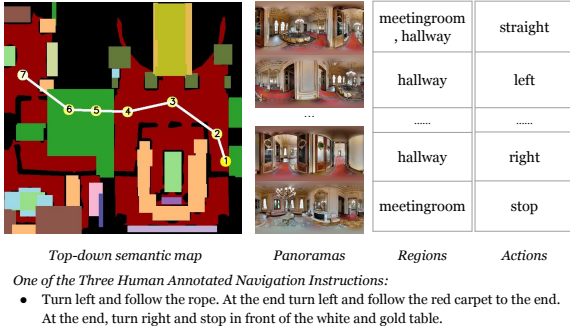


Figure 1: An example navigation scenario from our new dataset for instruction generation, with the navigation path overlaid on the semantic map.

split	size	Avg. # points	Avg. # regions	Avg. # objects
train	10623	5.95	3.26	22.64
val seen	768	6.07	3.3	22.36
val unseen	1839	5.87	3.11	22.13

Table 1: Statistics of extracted semantic maps. Avg. # region: average number of distinct regions along the path. Avg. # object: average number of object types in the semantic map.

maps, providing a new benchmark dataset and a baseline that demonstrates the feasibility of using semantic maps for VL-GEN task.

- We demonstrate experimentally with both automatic and human evaluations that including additional information (namely, region, action, and prompt) leads to more accurate and robust navigation instructions than using only semantic maps.
- We also conduct an intrinsic human evaluation of the quality of the generated instructions with fine-grained error analysis.

2. Task Definition and Data

A semantic map M_s is a top-down view of the scene s , which contains a path $P = \{p_1, \dots, p_K\}$, represented as a sequence of points connected by a line, and a set of N objects $O = \{o_1, \dots, o_N\}$.

In light of the success of image captioning models (Li et al., 2022; Wang et al., 2022b), we frame the VL-GEN task as image captioning task. Given a semantic map M_s , the task is to generate a natural language description D_P that describes the path P shown. Our task description replaces the photo-realistic RGB images used previously, with a semantic map. The processing of RGB images is resource-intensive, while our task definition has the advantage of abstracting away from the object recognition task, concentrating on the instruction generation task instead.

We also experiment with providing the model with additional features of the navigation path beyond the semantic maps alone, including actions, names of regions, and panoramic images. There is a fixed set of action types (LEFT, RIGHT, STRAIGHT, STOP), which are determined heuristically from the path shape at each navigation point. For each navigation point, we use the name of its associated region (e.g., hallway, meeting room). We do not think that panoramic images constitute ideal input to the system, but it is possible that they may provide additional visual information not shown in the map. Therefore, we also conduct experiments with panoramic images as part of the input information to the model.

We extract semantic maps, region and action information from the Habitat (Savva et al., 2019; Krantz et al., 2020) simulation environment. In a deployed robot, it may be obtained with a semantic SLAM component (Chaplot et al., 2020). Each object type on the map is represented in a unique color. We adopt the navigation paths and human annotations from the R2R dataset (Anderson et al., 2018). Panoramic images in RGB are obtained from the Matterport3D simulator (Chang et al., 2017) at each discrete navigation point. An example of the new dataset derived from R2R, including a semantic map with a path, language instruction, panorama images, actions, and region names, is shown in Figure 1.

Statistics about the semantic maps are presented in Table 1. The data splits we use are inherited from the original R2R dataset. The difference between seen validation set and the unseen validation set in R2R is whether the room environment is included in the train set.¹

3. Method

Motivated by the success of the multimodal pre-trained models, we construct a multimodal text generation model using BLIP² (Li et al., 2022). Figure 2 illustrates the architecture of the proposed model with modules that process different inputs; these will be described in Section 3.1. In Section 3.2, we describe the augmentations applied to the BLIP model in our experiments.

3.1. Model Input

Top-down semantic map (TD) The semantic map forms the main input used in all experiments.

¹Further details on the dataset are presented in Appendix A.1.

²The implementation is based on the Huggingface transformers library (Wolf et al., 2019): [Salesforce/blip-image-captioning-base](https://github.com/huggingface/transformers)

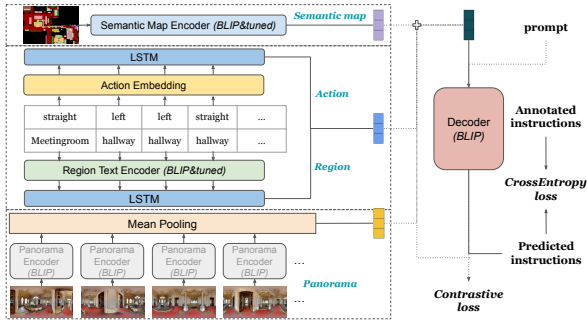


Figure 2: Illustration of the overall model architecture. Text input is encoded with pretrained BLIP text encoder and LSTM, and image input is encoded with the pretrained BLIP encoder. Modules shown in the same color share the weights. The weights of the panorama encoder are fixed.

It is encoded by the image encoder in the BLIP model. We first resize the image by nearest sampling to 384×384 and then feed it to the vision transformer with patch size 16.

Regions (Reg) and actions (Act) Region names and actions are frequently mentioned in human navigation instructions. To give the model information about the relevant region names, we represent them as a sequence of strings for each navigation point. We use a text encoder from the pre-trained BLIP model to represent the region names. The region embedding for each point is obtained by applying a mean pooling operation to the word embeddings. For actions, we apply an embedding layer to the discrete action values and get action embeddings in the same dimension as the region embedding. We add the region and the action embeddings together at each point and use a 3-layer LSTM model to embed the sequential information along the navigation path.

Panoramic images (Pano) Based on our analysis, visual object properties such as color and shape are mentioned in more than 25% of human instructions. As semantic maps only include object types but not the properties of visual objects, we augment the model input with panoramic images. This might enable the model to learn the visual properties mentioned in the instructions. We initialize the image encoder based on the pre-trained image encoder in BLIP model. We freeze its parameters during training because the model is pre-trained on photo-realistic images, which we believe endows the model with capabilities of recognizing panoramic images in our case. In order to increase the flexibility of the visual embedding, we apply an additional MLP with two linear layers on top of the panoramic vision encoder. Following the

methods in the video captioning task (Tang et al., 2021; Luo et al., 2022), we treat the panoramas as discrete frames and use the mean average of all panoramic embeddings to represent the panorama information of the navigation path.

Finally, the embedded input representations are added together to form the input to the decoder that outputs natural language instructions.

3.2. Model Augmentation

Multimodal alignment with contrastive loss

Contrastive learning is an effective method used in self-supervised learning for visual representation learning (Radford et al., 2021; Li et al., 2022) and multimodal pre-training in BLIP (Li et al., 2022). We investigate the effectiveness of introducing contrastive training for navigation instruction generation task as an auxiliary loss. We define the positive examples $P^+(C_{gt}, I_{gt})$ as pairs of the combined input embedding and the instruction embedding. The negative examples $P^-(C_{gt}, I_{rnd})$ consist of the pairs of the input embedding and the embedding of a randomly sampled instruction. Following CLIP (Radford et al., 2021), we multiply the multimodal input matrix E_{input} and textual instruction matrix E_{text} to obtain the predicted compatible matrix C_{pred} between inputs and labels and then compute the CrossEntropy loss on C_{pred} with the ground-truth correspondence C_{gt} .

Augmentation and grounding with prompt

The prompting of LLMs has demonstrated its effectiveness across various domains in previous works (Li and Liang, 2021; Liu et al., 2021; Tang et al., 2022; Keicher et al., 2022; Song et al., 2022). We generate the prompt from a template, which describes the nearby objects and regions, such as *Starting from the dark yellow point near sofa cushion in the living room region*. We tune the model with prompting and feed the prompt template to the decoder during inference. We argue that prompting can benefit the generation task in two ways. First, it can help visual-language grounding because the prompting template describes nearby landmarks and regions. Second, at inference time, the instructions that are generated are conditioned on the prompt template in an auto-regressive way, resulting in more controllable generation in VL-GEN task.

4. Experiments

We perform two evaluations over experiments: an automatic evaluation according to performance on the task (extrinsic) and a human evaluation of the quality of the instructions (intrinsic). These evaluations can tell us about the influence of region,

Input	P	C	SPICE		Human Score
			seen	unseen	unseen
TD (baseline)	-	-	20.50	16.19	3.42 (5)
	✓	-	20.79	15.77	-
	✓	✓	21.78*	17.10	-
TD+Reg+Act	-	-	21.00	17.00	4.20 (3)
	✓	-	21.86*	17.84**	4.29 (2)
	✓	✓	19.96	17.09	3.98 (4)
TD+Reg+Act+Pano	-	-	19.87	17.44*	4.36* (1)
	✓	-	22.14**	17.79**	-
	✓	✓	20.36	17.08	-

Table 2: Automatic (SPICE) and human evaluation results with inputs of different modalities in seen and unseen environments, where P is short for prompt and C is short for contrastive loss. ** and * indicate statistically significant difference with the baseline ($p \leq 0.01$) and ($p \leq 0.05$).

actions, prompting, and contrastive loss on the quality of the instructions both quantitatively and qualitatively.

4.1. Experimental setup

We train the model using the train split of the R2R dataset and evaluate it both on validation seen and unseen sets. We use the BLIP-*base* model for experiments. We setup the baselines with different combinations of the input: 1) top-down semantic map (TD) 2) + regions (Reg) and actions (Act); 3) + panoramic images (Pano). We also experiment with contrastive loss and prompting, making 9 system variants for experiments in total.

In the intrinsic human evaluation, we use a Latin Square design of size 5. We therefore compare only a subset of the above system variants with different combinations of input (*TD*, *TD+Reg+Act* and *TD+Reg+Act+Pano*), and prompting and contrastive loss on *TD+Reg+Act*.

4.2. Human Participants and Procedure

For the human experiment, we recruit 5 evaluators who have never contributed to or been involved in the project before under the consent from the Ethics Committee. The evaluation workload for each participant is designed to be within 30 minutes for them to concentrate on the task. We also provide two specific illustration examples about the evaluation task for the human participants. The evaluation materials consist of 15 navigation paths in the unseen environments, randomly sampled. The experiment is performed online using an evaluation interface. The participants are shown the semantic map with the path as well as panorama images. They are asked to assign a score from 0 (worst) to 10 (best) based on the quality of the

instruction candidates generated by different systems.

4.3. Automatic Evaluation Metrics

In the automatic evaluation, we compare the performance of 9 system variants based on an automatic metric SPICE (Semantic Propositional Image Caption Evaluation) (Anderson et al., 2016), following Zhao et al. (2021). SPICE is a metric used to evaluate the quality of image captions, focusing on the semantic content of captions. It identifies semantic propositions within the parse trees and compares the semantic propositions from the generated caption with those from the reference captions.

When comparing different systems, we use the two-sided permutation test to see if the arithmetic means of the two systems' performances are equal. If the p-value is larger than 0.05, we consider the performance of the two systems to be not significantly different.

4.4. Evaluation Results

Table 2 shows the SPICE and human evaluation scores in seen and unseen environments. As expected, the models perform better in seen than in unseen setting by 3.88 in SPICE score on average across all 9 systems. For both settings, we observe that using region and action information with the prompt improves the model's performance with $p \leq 0.05$, while contrastive learning does not seem to help. Adding panoramic images tends to improve the performance, but not significantly ($p \geq 0.1$). When comparing with previous methods in SPICE score, our systems (17.84/22.14) perform on par or even achieve higher SPICE scores than Speaker Fol. (Fried et al., 2018) (17.0/18.7) and EnvDrop (Tan et al., 2019) (18.1/20.2) on unseen/seen settings.

In the results for the human evaluation, shown in Table 2, we observe that using the semantic map as the only input results in the lowest average score across all systems (3.42). This repeats the observations from the automatic evaluation. Using regions, actions, and panoramas achieves the highest rating (4.36) which is significantly better than the baseline ($p=0.05$), followed by using regions, actions, and prompts (4.29). However, incorporating *Pano* (4.36) alongside *TD+Reg+Act* (4.20) does not show a noteworthy difference.

In addition to the results above, we were also curious about the degree to which our automatic results in SPICE correlate with the human judgments. We measure a Kendall τ correlation between SPICE and human evaluation results of 0.6

Input Information	P	C	Incorrect	Hallucination	Redundancy	Linguistic
TD	-	-	15	10	0	0
TD + Reg + Act	-	-	15	10	0	1
TD + Reg + Act	✓	-	12	6	1	2
TD + Reg + Act	✓	✓	12	6	1	2
TD + Reg + Act + Pano	-	-	11	6	0	0

Table 3: Error analysis on randomly selected predictions from the systems in unseen environments, where P is short for prompt and C is short for contrastive loss.

and conclude that this is satisfactory, justifying the use of SPICE for automatic evaluation.³

Our findings indicate that incorporating more information in different modalities tends to improve the performance for the generation task. Our semantic map abstracts information in a way that is useful for current systems, although it consists of only a single image. Most of our system variants that do not use panorama images performs on-par with the existing LSTM-based end-to-end approaches that use only panoramic images. However, the absolute performance of all models is still low, indicating that there is much room for improvement.

4.5. Error Analysis

Further to human evaluation score, we manually analyze the quality of the instructions generated by the same 5 system variants according to the following four aspects:

- **Incorrectness:** Does the prediction contain incorrect information?
- **Hallucination:** Does the prediction contain a description not corresponding to the input?
- **Redundancy:** Does the prediction contain redundant expressions and information?
- **Linguistic problems:** Is the generated instruction grammatically wrong or not fluent?

For each experimental setting, we randomly select 15 examples. The counts for each error type are given in Table 3. We can see that the systems that do not use prompting or panorama images contain errors in all cases. Most of these errors are caused by hallucinations. Analyzing hallucinations further, we find that the action descriptions are most prone to hallucinations, such as when left and right are confused with each other. When regions and actions are used as input, the number of hallucinations in action descriptions goes down, but remains high in regions.

³We also computed BLEU and ROUGE scores, however they show lower correlation with the human-assigned scores, which are omitted here.

Apart from changing the input information, when we train the model with prompting, the resulting instructions are less likely to include hallucinations in terms of actions and objects. Yet after introducing the contrastive loss, it causes redundancy and linguistic problems in the predictions. The language quality problems mainly consist of spelling mistakes in objects and regions, and punctuation errors when introducing the prompt and contrastive loss for training. This may be because the contrastive loss influences the CrossEntropy loss and thus interferes with the language generation task.

5. Conclusion

Our longer-term goal is to build mobile robots with spatial awareness and reasoning capabilities which can follow natural language instructions and express their intentions in natural language. We propose to use semantic maps as the intermediate representation for spatial reasoning as it is a human-interpretable and light-weight approach that encodes information necessary for the navigation in a single abstract image.

In this work, we create the dataset with top-down semantic maps for R2R corpus and reframe instruction generation task as image captioning, using abstract top-down semantic map as main input. We set a baseline for the instruction generation from semantic map input. Our experimental results show that using the top-down semantic map performs on-par with the end-to-end methods that use sequence of panorama images as input.

Limitations

The current approach to the semantic map representation is missing some of the information required to generate or interpret instructions. For example, room names, such as *bathroom*, *bedroom*, or *sitting room*, are naturally used in indoor navigation instructions. However, the current single-layer semantic map representation does not encode the information about such region names. To address this in our current approach, we provide region names for each navigation point as a separate textual input. The limitation of this approach is that it only includes the region names for the navigation points. For example, an instruction ‘*Stop in front of the bathroom*’, the *bathroom* will not be included in the input because the navigation point is outside of the bathroom region. In future work, we plan to introduce a multi-layered semantic map where, in addition to encoding objects, a separate layer encodes information about regions.

Another limitation is that current semantic map encoding does not encode object properties, such as color, material, or shape. According to our analy-

sis, object properties are mentioned in one-third of the instructions, but these would not be captured by the map. To address this limitation, in future work, we will encode the object properties in the semantic map.

Data and Code availability

We release the code for data preparation, model training and inference, and evaluation at <https://github.com/chengzu-li/VLGen>, along with the prompt templates and hyper-parameter settings for experiments. We also release the the top-down semantic maps extracted from Habitat environment extending the existing R2R dataset, which can be obtained upon request following the guideline at <https://github.com/chengzu-li/VLGen>.

Bibliographical References

- Dong An, Yuankai Qi, Yangguang Li, Yan Huang, Liang Wang, Tieniu Tan, and Jing Shao. 2023. Bevbort: Multimodal map pre-training for language-guided navigation. *Proceedings of the IEEE/CVF International Conference on Computer Vision*.
- Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. 2016. SPICE: Semantic propositional image caption evaluation. In *ECCV*.
- Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton van den Hengel. 2018. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Satanjeev Banerjee and Alon Lavie. 2005. **ME-TEOR: An automatic metric for MT evaluation with improved correlation with human judgments**. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. 2017. Matterport3d: Learning from rgb-d data in indoor environments. *arXiv preprint arXiv:1709.06158*.
- Devendra Singh Chaplot, Dhiraj Gandhi, Saurabh Gupta, Abhinav Gupta, and Ruslan Salakhutdinov. 2020. Learning to explore using active neural slam. In *International Conference on Learning Representations (ICLR)*.
- Howard Chen, Alane Suhr, Dipendra Misra, Noah Snaveley, 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.
- Ta-Chung Chi, Minmin Shen, Mihail Eric, Seokhwan Kim, and Dilek Hakkani-tur. 2020. Just ask: An interactive learning framework for vision and language navigation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 2459–2466.
- Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. 2021. Unifying vision-and-language tasks via text generation. In *International Conference on Machine Learning*, pages 1931–1942. PMLR.
- Chenhe Dong, Yinghui Li, Haifan Gong, Miaoxin Chen, Junxin Li, Ying Shen, and Min Yang. 2022. A survey of natural language generation. *ACM Computing Surveys*, 55(8):1–38.
- Vishnu Sashank Dorbala, Gunnar Sigurdsson, Robinson Piramuthu, Jesse Thomason, and Gaurav S Sukhatme. 2022. Clip-nav: Using clip for zero-shot vision-and-language navigation. *arXiv preprint arXiv:2211.16649*.
- Weixi Feng, Tsu-Jui Fu, Yujie Lu, and William Yang Wang. 2022. Uln: Towards underspecified vision-and-language navigation. *arXiv preprint arXiv:2210.10020*.
- 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. *Advances in Neural Information Processing Systems*, 31.
- Albert Gatt and Emiel Krahmer. 2017. Survey of the state of the art in natural language generation: Core tasks, applications and evaluation. *arXiv preprint arXiv:1703.09902*.
- Jing Gu, Eliana Stefani, Qi Wu, Jesse Thomason, and Xin Wang. 2022. **Vision-and-language navigation: A survey of tasks, methods, and**

- future directions. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7606–7623, Dublin, Ireland. Association for Computational Linguistics.
- David Hall, Ben Talbot, Suman Raj Bista, Haoyang Zhang, Rohan Smith, Feras Dayoub, and Niko Sünderhauf. 2020. [The robotic vision scene understanding challenge](#).
- Xu Han, Weilin Zhao, Ning Ding, Zhiyuan Liu, and Maosong Sun. 2022. Ptr: Prompt tuning with rules for text classification. *AI Open*, 3:182–192.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Karl Moritz Hermann, Mateusz Malinowski, Piotr Mirowski, Andras Banki-Horvath, Keith Anderson, and Raia Hadsell. 2020. Learning to follow directions in street view. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11773–11781.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Roman Le Bras, and Yejin Choi. 2021. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*.
- Yicong Hong, Cristian Rodriguez-Opazo, Qi Wu, and Stephen Gould. 2020a. Sub-instruction aware vision-and-language navigation. *arXiv preprint arXiv:2004.02707*.
- Yicong Hong, Qi Wu, Yuankai Qi, Cristian Rodriguez-Opazo, and Stephen Gould. 2020b. A recurrent vision-and-language bert for navigation. *arXiv preprint arXiv:2011.13922*.
- Yicong Hong, Qi Wu, Yuankai Qi, Cristian Rodriguez-Opazo, and Stephen Gould. 2021. A recurrent vision-and-language bert for navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1643–1653.
- Nikolai Ilinykh, Yasmeen Emampoor, and Simon Dobnik. 2022. [Look and answer the question: On the role of vision in embodied question answering](#). In *Proceedings of the 15th International Conference on Natural Language Generation*, pages 236–245, Waterville, Maine, USA and virtual meeting. Association for Computational Linguistics.
- Vihan Jain, Gabriel Magalhaes, Alexander Ku, Ashish Vaswani, Eugene Ie, and Jason Baldridge. 2019. [Stay on the path: Instruction fidelity in vision-and-language navigation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1862–1872, Florence, Italy. Association for Computational Linguistics.
- Aishwarya Kamath, Peter Anderson, Su Wang, Jing Yu Koh, Alexander Ku, Austin Waters, Yinfei Yang, Jason Baldridge, and Zarana Parekh. 2022. A new path: Scaling vision-and-language navigation with synthetic instructions and imitation learning. *arXiv preprint arXiv:2210.03112*.
- Matthias Keicher, Kamilia Mullakaeva, Tobias Czempel, Kristina Mach, Ashkan Khakzar, and Nassir Navab. 2022. Few-shot structured radiology report generation using natural language prompts. *arXiv preprint arXiv:2203.15723*.
- Maurice G Kendall. 1938. A new measure of rank correlation. *Biometrika*, 30(1/2):81–93.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Thomas Kollar, Jayant Krishnamurthy, and Grant P Strimel. 2013. Toward interactive grounded language acquisition. In *Robotics: Science and systems*, volume 1, pages 721–732.
- Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. 2020. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVIII 16*, pages 104–120. Springer.
- Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, and Jason Baldridge. 2020. Room-across-room: Multilingual vision-and-language navigation with dense spatiotemporal grounding. *arXiv preprint arXiv:2010.07954*.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*.
- Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*.

- Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. 2020. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXX 16*, pages 121–137. Springer.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Junyang Lin, An Yang, Yichang Zhang, Jie Liu, Jingren Zhou, and Hongxia Yang. 2020. Interbert: Vision-and-language interaction for multi-modal pretraining. *arXiv preprint arXiv:2003.13198*.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35.
- Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2021. P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks. *arXiv preprint arXiv:2110.07602*.
- Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. 2022. Clip4clip: An empirical study of clip for end to end video clip retrieval and captioning. *Neuro-computing*, 508:293–304.
- 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 *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2667–2678, Brussels, Belgium. Association for Computational Linguistics.
- OpenAI. 2023. Gpt-4 technical report. *ArXiv*, abs/2303.08774.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [BLEU: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Seth Pate, Wei Xu, Ziyi Yang, Maxwell Love, Siddarth Ganguri, and Lawson L. S. Wong. 2021. [Natural language for human-robot collaboration: Problems beyond language grounding](#). *ArXiv*, abs/2110.04441.
- Tzuf Paz-Argaman, John Palowitch, Sayali Kulkarni, Jason Baldridge, and Reut Tsarfaty. 2024. [Where do we go from here? multi-scale allocentric relational inference from natural spatial descriptions](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1026–1040, St. Julian’s, Malta. Association for Computational Linguistics.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.
- Laria Reynolds and Kyle McDonell. 2021. Prompt programming for large language models: Beyond the few-shot paradigm. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–7.
- Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. 2019. Habitat: A platform for embodied ai research. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9339–9347.
- Sheng Shen, Liunian Harold Li, Hao Tan, Mohit Bansal, Anna Rohrbach, Kai-Wei Chang, Zhewei Yao, and Kurt Keutzer. 2021. How much can clip benefit vision-and-language tasks? *arXiv preprint arXiv:2107.06383*.
- Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. 2020. [AL-FRED: A Benchmark for Interpreting Grounded Instructions for Everyday Tasks](#). In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Xuemeng Song, Liqiang Jing, Dengtian Lin, Zhongzhou Zhao, Haiqing Chen, and Liqiang Nie. 2022. [V2P: Vision-to-prompt based multi-modal product summary generation](#). In *Proceedings of the 45th International ACM SIGIR Con-*

- ference on Research and Development in Information Retrieval, SIGIR '22, page 992–1001, New York, NY, USA. Association for Computing Machinery.
- Xiuchao Sui, Shaohua Li, Hong Yang, Hongyuan Zhu, and Yan Wu. 2023. [Language models can do zero-shot visual referring expression comprehension](#). In *International Conference on Learning Representations (ICLR)*.
- Hao Tan and Mohit Bansal. 2019. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*.
- Hao Tan, Licheng Yu, and Mohit Bansal. 2019. [Learning to navigate unseen environments: Back translation with environmental dropout](#). In *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)*, pages 2610–2621, Minneapolis, Minnesota. Association for Computational Linguistics.
- Mingkang Tang, Zhanyu Wang, Zhenhua Liu, Fengyun Rao, Dian Li, and Xiu Li. 2021. Clip4caption: Clip for video caption. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 4858–4862.
- Tianyi Tang, Junyi Li, Wayne Xin Zhao, and Ji-Rong Wen. 2022. Context-tuning: Learning contextualized prompts for natural language generation. *arXiv preprint arXiv:2201.08670*.
- Yoad Tewel, Yoav Shalev, Idan Schwartz, and Lior Wolf. 2021. Zero-shot image-to-text generation for visual-semantic arithmetic. *arXiv preprint arXiv:2111.14447*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575.
- Hanqing Wang, Wei Liang, Jianbing Shen, Luc Van Gool, and Wenguan Wang. 2022a. Counterfactual cycle-consistent learning for instruction following and generation in vision-language navigation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15471–15481.
- Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022b. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*, pages 23318–23340. PMLR.
- Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022c. Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. *arXiv preprint arXiv:2202.03052*.
- Su Wang, Ceslee Montgomery, Jordi Orbay, Vighnesh Birodkar, Aleksandra Faust, Izzeddin Gur, Natasha Jaques, Austin Waters, Jason Baldrige, and Peter Anderson. 2022d. Less is more: Generating grounded navigation instructions from landmarks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15428–15438.
- Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. 2021. Simvlm: Simple visual language model pre-training with weak supervision. *arXiv preprint arXiv:2108.10904*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.
- Wansen Wu, Tao Chang, and Xinmeng Li. 2021. Visual-and-language navigation: A survey and taxonomy. *arXiv preprint arXiv:2108.11544*.
- Teng Xue, Weiming Wang, Jin Ma, Wenhai Liu, Zhenyu Pan, and Mingshuo Han. 2020. [Progress and prospects of multimodal fusion methods in physical human–robot interaction: A review](#). *IEEE Sensors Journal*, 20(18):10355–10370.
- Jianing Yang, Xuweiyi Chen, Shengyi Qian, Nikhil Madaan, Madhavan Iyengar, David F. Fouhey, and Joyce Chai. 2023. [LLM-Grounder: Open-vocabulary 3d visual grounding with large language model as an agent](#). *ArXiv*, abs/2309.12311.
- Tian Yun, Chen Sun, and Ellie Pavlick. 2021. Does vision-and-language pretraining improve lexical grounding? *arXiv preprint arXiv:2109.10246*.
- Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. 2021. Vinvl: Revisiting visual representations in vision-language models.

In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5579–5588.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Ming Zhao, Peter Anderson, Vihan Jain, Su Wang, Alexander Ku, Jason Baldridge, and Eugene Ie. 2021. On the evaluation of vision-and-language navigation instructions. *arXiv preprint arXiv:2101.10504*.

Luwei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason Corso, and Jianfeng Gao. 2020. Unified vision-language pre-training for image captioning and vqa. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 13041–13049.

Shuyan Zhou, Uri Alon, Frank F Xu, Zhengbao Jlang, and Graham Neubig. 2022. Doccoder: Generating code by retrieving and reading docs. *arXiv preprint arXiv:2207.05987*.

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.

Zhuofan Zong, Guanglu Song, and Yu Liu. 2022. [DETRs with collaborative hybrid assignments training](#). *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6725–6735.

A. Data Extraction

A.1. Conditions of Data Extraction

This section describes how the semantic map is extracted from the Habitat environment.

Objects Objects on a semantic map are represented by the bounding box with a unique color assigned to each object type. We use the (X, Y) coordinates of the object’s bounding box in Matterplot3D to represent them in the 2D semantic map.

There are 40 different object types labeled in the simulation environment. We filter out the following object types from the simulator because they are seldom mentioned in the instructions but take up a large area in the semantic map:

```
['misc', 'ceiling', 'curtain',
'objects', 'floor', 'wall',
'void']
```

For the buildings with multiple floors, we extract a semantic map for each floor. Given the 3D coordinates of the object’s center (x_i, y_i, h_i) and the size of the object’s bounding box is (w_x, w_y, w_h) , we use the agent’s vertical position h_{agent} to filter the objects for a given floor by including all objects that satisfy one of the following:

$$h_i - \frac{1}{2}w_h \leq h_{agent} \leq h_i + \frac{1}{2}w_h$$

$$|h_i - h_{agent}| \leq 1.6$$

Regions For each navigation point, we determine the corresponding region by calculating whether the agent’s current position is within the area of the region. The region’s area is defined by the coordinates of the center (x_c, y_c) and the sizes in width and length (w_x, w_y) as a rectangle. We define that if the agent’s location (l_x, l_y) satisfies the following requirement, the region would be added into the information of this navigation point.

$$x_c - \frac{1}{2}w_x \leq l_x \leq x_c + \frac{1}{2}w_x$$

$$y_c - \frac{1}{2}w_y \leq l_y \leq y_c + \frac{1}{2}w_y$$

Actions The actions are a closed set of LEFT, RIGHT, STRAIGHT, STOP. They are determined based on the coordination of the navigation points. We calculate the differences in angles between the previous navigation point and the current position and define that if the differences stay within 20 degrees, the agent is heading straight. Otherwise, the agent makes a turn, with the corresponding direction depending on whether the difference is positive or negative.

A.2. More Information about the Extracted Data

To evaluate the quality of the dataset, we randomly sample 30 examples from the training set and look into the instructions regarding the way they describe the path. we find that the instructions mention 2.2 landmark objects and 1.3 regions on average. 20 out of 30 sampled instructions can be inferred simply from the top-down view. This finding justifies the use of top-down semantic maps for navigation instruction generation to some degree. For the other 10 instructions, 7 out of 10 can not be inferred only from the top-down view due to the descriptive expressions about the environment. The descriptions can be obtained from the panoramic images as supportive information (such as *going upstairs* or *downstairs*) and requires the interactions between different types of input. The other 3 instructions miss the region annotations from the

simulator in R2R. These observations indicate that the new task we propose has the problems of weak supervision and requires the model to connect different types of inputs with each other.

A.3. Alignment between Colors and Objects

Below shows the mapping between RGB pixel values and the object types (textual names) in the semantic map.

```
[31, 119, 180], "void",
[174, 199, 232], "wall",
[255, 127, 14], "floor",
[255, 187, 120], "chair",
[44, 160, 44], "door",
[152, 223, 138], "table",
[214, 39, 40], "picture",
[255, 152, 150], "cabinet",
[148, 103, 189], "cushion",
[197, 176, 213], "window",
[140, 86, 75], "sofa",
[196, 156, 148], "bed",
[227, 119, 194], "curtain",
[247, 182, 210], "chest_of_drawers",
[127, 127, 127], "plant",
[199, 199, 199], "sink",
[188, 189, 34], "stairs",
[219, 219, 141], "ceiling",
[23, 190, 207], "toilet",
[158, 218, 229], "stool",
[57, 59, 121], "towel",
[82, 84, 163], "mirror",
[107, 110, 207], "tv_monitor",
[156, 158, 222], "shower",
[99, 121, 57], "column",
[140, 162, 82], "bathtub",
[181, 207, 107], "counter",
[206, 219, 156], "fireplace",
[140, 109, 49], "lighting",
[189, 158, 57], "beam",
[231, 186, 82], "railing",
[231, 203, 148], "shelving",
[132, 60, 57], "blinds",
[173, 73, 74], "gym_equipment",
[214, 97, 107], "seating",
[231, 150, 156], "board_panel",
[123, 65, 115], "furniture",
[165, 81, 148], "appliances",
[206, 109, 189], "clothes",
[222, 158, 214], "objects",
[255, 255, 102], "[POINT]",
[255, 255, 0], "[START]",
[255, 255, 204], "[END]",
[255, 255, 255], "[LINE]",
[0, 0, 0], "[NONNAVIGABLE]",
[150, 0, 0], "[NAVIGABLE]"
```

B. Experimental Setup

B.1. Hyperparameters

We train our model for a maximum of 25 epochs using an initial learning rate of $5e-5$ with linear lr scheduler. The batch size is set to 32 for training and 64 for validation. When balancing the contrastive loss and CrossEntropy loss, we assign a weight of 0.1 to the contrastive loss to make the model more focused on the generation task.

B.2. Data Preprocessing

We mainly adopt the original BLIP processor for image and text inputs, but make a few modifications to the top-down semantic map. Because the maps are in different sizes for different rooms, we first pad them to the size of 1024×1024 with black pixels and apply a masking strategy to the top-down map by only selecting the nearby regions of the path. We keep the receptive field of the top-down view to a certain value (default to 40 pixels) within the path area. Then, in order to avoid introducing new pixel values when resizing, we resize the masked image with the nearest resampling interpolation strategy to 386×386 , following the default setting of BLIP.

C. Prompt Design

```
prompt      : Starting from the dark
              yellow point [objects]
              [regions], [instruction]
```

For example:

```
[objects]   : near sofa cushion
[regions]   : in the living room region
[instruction]: exit the living room, turn
              left, wait at the bottom
              of the stairs.
```

D. Experiment Results

D.1. Significance Test Results on SPICE Scores

We first define the indices for all of our system variants from 1 to 9 following the order of systems in Table 2. Table 4 shows the full significance test on SPICE scores in unseen environments.

D.2. Human Evaluations

For the evaluation page, it shows the top-down semantic map, the panoramic images and the region information as well. The page provides 5 generated instructions to describe the navigation path from 5 different generator models. The evaluator is supposed to give the quality scores for these 5 instructions based on the guidance in the instruction documentation. Figure 3 shows a screenshot

	2(15.77)	3(17.10)	4(17.00)	5(17.84)	6(17.09)	7(17.44)	8(17.79)	9(17.08)
1(16.19)	0.4421	0.7891	0.1554	0.0030	0.1050	0.0269	0.0055	0.1115
2(15.77)		0.2992	0.0310	0.0002	0.0175	0.0033	0.0004	0.0194
3(17.10)			0.2432	0.0072	0.1726	0.0505	0.0113	0.1825
4(17.00)				0.1459	0.8774	0.4536	0.1825	0.8890
5(17.84)					0.1822	0.4809	0.9318	0.1830
6(17.09)						0.5423	0.2256	0.9913
7(17.44)							0.5463	0.5401
8(17.79)								0.2271

Table 4: Two-sided permutation test p-values on SPICE in validation unseen environments. The row names and column names are the system indices for different systems, with the SPICE values in the parenthesis brackets. The numbers in bold are the p-values below 0.05.

of the interface for human evaluation.

Significance test based on human evaluation

Table 5 presents the two-sided permutation test results based on the human evaluation.

	4(4.20)	5(4.29)	6(3.98)	7(4.36)
1(3.42)	0.10	0.06	0.27	0.05
4(4.20)		0.88	0.68	0.77
5(4.29)			0.55	0.92
6(3.98)				0.47

Table 5: Two-sided permutation test results between systems based on the human evaluation. The row indices and column indices are the system indices following Table 4 and their quality scores in the parenthesis.

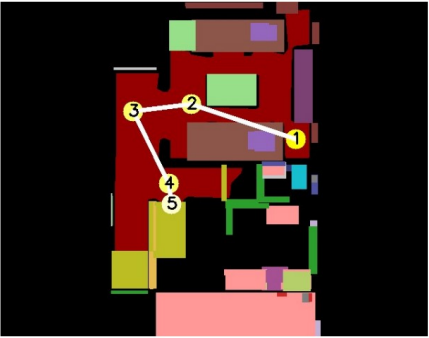
Ranking Evaluator for Robot Navigation Instruction

select an index

0 - +

Top-down Images

Top-down images ^




Path Image, starting from point 1 to point 5



Near Objects Labels

See panoramic view ^



Evaluation

Regions

	Point-1	Point-2	Point-3	Point-4	Point-5
0	familyroom/lounge	familyroom/lounge, hallway	hallway	hallway	hallway

Candidates

go past the statue and into the doorway on the right. walk straight, turn right, and wait in the hallway.

Quality (0 being worst and 10 being best)

0
0
10

exit the room, go to the patio door and take a left, and go to the stairs.

Quality (0 being worst and 10 being best)

0
0
10

walk out of the bathroom and turn right. walk past the stairs and turn left. walk down the stairs and stop in front of the bathroom.

Quality (0 being worst and 10 being best)

0
0
10

turn left. turn left at the stairs. go past the curved entryway and go into the bathroom. wait near the sink.

Quality (0 being worst and 10 being best)

0
0
10

exit the room. turn left and go down the hallway. go past the large glass vase and into the room. wait near the black chair.

Quality (0 being worst and 10 being best)

0
0
10

Figure 3: Screenshot of the evaluation interface for human evaluation.