Controllable Navigation Instruction Generation with Chain of Thought Prompting

Xianghao Kong¹[®]^{*}, Jinyu Chen¹[®]^{*}, Wenguan Wang²[®][⊠], Hang Su³[®], Xiaolin Hu³[®], Yi Yang²[®], and Si Liu¹[®][⊠]

 ¹ School of Artificial Intelligence, Beihang University
 ² College of Computer Science and Technology, Zhejiang University
 ³ Department of Computer Science and Technology, Tsinghua University https://github.com/refkxh/C-Instructor

Abstract. Instruction generation is a vital and multidisciplinary research area with broad applications. Existing instruction generation models are limited to generating instructions in a single style from a particular dataset, and the style and content of generated instructions cannot be controlled. Moreover, most existing instruction generation methods also disregard the spatial modeling of the navigation environment. Leveraging the capabilities of Large Language Models (LLMs), we propose C-INSTRUCTOR, which utilizes the chain-of-thought-style prompt for stylecontrollable and content-controllable instruction generation. Firstly, we propose a Chain of Thought with Landmarks (CoTL) mechanism, which guides the LLM to identify key landmarks and then generate complete instructions. CoTL renders generated instructions more accessible to follow and offers greater controllability over the manipulation of landmark objects. Furthermore, we present a Spatial Topology Modeling Task to facilitate the understanding of the spatial structure of the environment. Finally, we introduce a Style-Mixed Training policy, harnessing the prior knowledge of LLMs to enable style control for instruction generation based on different prompts within a single model instance. Extensive experiments demonstrate that instructions generated by C-INSTRUCTOR outperform those generated by previous methods in text metrics, navigation guidance evaluation, and user studies.

Keywords: Instruction generation · Vision-and-language navigation

1 Introduction

Developing an agent capable of communicating with humans in natural language and accomplishing specific tasks in its environment is a crucial goal for researchers in the field of embodied AI. Such an agent needs two key abilities: the first one is to execute specific tasks based on human instructions, and the second one is to provide interactive feedback and guidance to humans based on environmental information. Regarding the first ability, one of the most typical

^{*} Equal contribution. \bowtie Corresponding author.



Fig. 1: C-INSTRUCTOR possesses the ability to control the linguistic style of generated instructions, and the ability to manipulate landmarks within the instructions (§1).

tasks is vision-and-language navigation (VLN) [5], which has garnered extensive research interest [26, 31, 44, 45, 48, 51, 54, 73] and developed fast in recent years [2, 3, 9, 17-19, 23, 34, 39, 40, 49, 59-61, 64, 67, 71].

Regarding the implementation of the second capability, *i.e.*, machine feedback, one of its prominent facets, instruction generation, has been a long-standing area of multidisciplinary research dating back to the 1960s [43]. The instruction generation model can be used for describing a path explored by a robot to a human in human-robot collaboration tasks. In practical scenarios, it can be applied to intelligent guidance for the visually impaired [26], foster human-machine trust [63], and provide guidance in hazardous scenarios, etc. An instruction generation model fulfilling the prerequisites of human-machine collaboration must possess the following two capabilities [31,48], *i.e.*, executability and controllability. For executability, instructions are supposed to exhibit high linguistic quality and provide clear guidance at navigational junctions. For controllability, control over instruction generation in style and content is also of essential importance to improve communication efficiency. For example, when the instruction recipient is acquainted with the environment, it is more efficient to generate instructions with higher levels of abstraction. Additionally, the guidance provided in the instructions may need adjustments based on the landmarks that the instruction recipient focuses on in the environment.

To enhance the executability and controllability of instruction generation models, we propose a <u>C</u>ontrollable Navigation <u>Instructor</u> (C-INSTRUCTOR), which possesses the ability to generate easily executable instructions with high linguistic quality, as well as the capability to controllably generate instructions in various linguistic styles with different landmarks (Fig. 1). C-INSTRUCTOR primarily encompasses the following four technological contributions: **First**, to enhance the linguistic quality of instruction generation and handle different styles of instructions neatly, we propose an adapter structure that effectively incorporates path information into the GPT-based Large Language Model (LLM) [20]. **Second**, to

3

improve the executability of generated instructions, we present a training strategy involving a Chain-of-Thought with Landmarks (CoTL) mechanism and a Spatial Topology Modeling Task (STMT). CoTL employs a step-by-step thinking [66] approach to guide the model to identify crucial landmarks before generating complete instructions; STMT incorporates spatial connectivity prediction as an auxiliary task in training to facilitate the understanding of the topological structure of the environment. Third, in order to generate instructions in various styles with a single model instance, we introduce a Style-Mixed Training (SMT) policy, in which different styles of instructions are jointly learned. Distinct instruction styles are trained using prompts as differentiation, enabling control over the style of generated instructions. Fourth, the collaboration between CoTL and SMT enhances the capabilities of crucial navigation waypoints localization and spatial direction guiding, thus improving the executability of the generated instructions. Benefiting from SMT and CoTL, C-INSTRUCTOR allows control over the generation style of instructions and attention to specific objectives while maintaining high linguistic quality of generated instructions.

In our experiments, C-INSTRUCTOR significantly outperforms previous instruction generation methods [16,52,58,63] across different linguistic metrics on four indoor/outdoor benchmarks [5,26,31,48]. In addition, it proves to be an effective means of data augmentation for VLN training over previous speaker models [16,52,58,63]. Moreover, instructions generated by C-INSTRUCTOR demonstrate enhanced navigation guidance capabilities in both instruction following model experiments and human evaluations.

2 Related Work

Navigation Instruction Generation. The study of generating linguistic instruction for navigation can date back to Lynch's work [43] in the 1960s. Early efforts [1,65] investigated the human cognitive mechanism for describing routes. They found that navigation direction is associated with the cognitive map [32] and influenced by various factors including cultural background [56] and genders [27]. This area has long been overlooked by the computer vision academia and is simply viewed as a data augmentation tool for VLN. However, it holds significant practical relevance, e.g., establishing human-machine trust [63] and facilitating blind navigation [26]. Fried et al [16] first proposed a LSTM-based instruction generation model to augment training samples and re-weight the route choice of the navigator. There are three primary aspects for the advancement of instruction generation: elevated linguistic quality, finer-grained directives, and longer, more intricate instructions. In order to enhance the quality of instructions, some methods introduce supplementary information like external knowledge [68] and landmark information [62, 70], build instruction template [70] and utilize larger language models [62]. [22,24,29,70,74] generate fine-grained alignment between language and navigation paths. To build more intricate instructions, [28, 38, 74] cross-connect paths to generate longer instruction-trajectory pairs. Methods like [15, 58, 63] also consider instruction generation and follow-

ing as dual tasks, and employ joint-optimization or cycle-consistent learning to promote navigation performance and instruction generation quality.

Previous deep-learning-based methods [16, 52, 58, 63] can only generate navigation instructions in a single style with limited linguistic quality and no controllability. By leveraging LLMs, C-INSTRUCTOR notably enhances the linguistic quality of instructions. Moreover, C-INSTRUCTOR provides style and content controllability in a single model instance via SMT and CoTL respectively.

Parameter-Efficient Fine-Tuning. The pre-training and fine-tuning paradigm has demonstrated remarkable efficacy in VLN and various other tasks. However, as model parameters grow exponentially and downstream task data remain limited, full-scale fine-tuning fails to yield robust performance on downstream tasks due to overfitting and catastrophic forgetting. The approach known as Parameter-efficient Fine-tuning (PEFT), involving the selective freezing of a significant portion of the model's parameters while training only a small subset, has met success in numerous domains. PEFT has proven highly effective in adapting pre-trained models like CLIP [50], BERT [12], and GPT [8,55] to downstream tasks. There are three main types of PEFT methods, namely prefix finetuning, reparameterization, and adapter. Prefix finetuning methods like [33, 36, 41, 72] feed learnable prompts into the model to learn task-specific knowledge. The methods [25, 30] use reparameterization to reduce the amount of trainable parameters. Approaches employing adapters [20, 69] adeptly accommodate inputs from diverse modalities and various downstream tasks by incorporating additional layers into the pre-trained network.

Understanding the spatial topology of the navigation environment is essential for the instruction generator to guide the instruction follower. Based on adapter PEFT methods [20, 69], C-INSTRUCTOR introduces a trajectory encoder to incorporate spatial information into the LLM. Moreover, C-INSTRUCTOR includes STMT to facilitate the understanding of spatial connectivity of the environment.

3 Methodology

3.1 Task Formulation

The instruction generation model is required to generate the instruction $X = \{x_1, x_2, ..., x_S\}$ with S words that provides guidance for following the given path $R = \{r_1, r_2, ..., r_T\}$ with T steps. At a given time step t, r_t is composed of the panoramic observation o_t and action a_t . The objective of model parameters $\boldsymbol{\theta}$ is to maximize the likelihood of the target instruction X^* :

$$\boldsymbol{\theta}^* = \operatorname*{arg\,max}_{\boldsymbol{\theta}} \log p(X^* | R, \boldsymbol{\theta}). \tag{1}$$

3.2 Overall Framework

To leverage the linguistic capabilities of LLMs, we employ an adapter-based [20] approach in C-INSTRUCTOR to embed actions and visual observations. The



(a) The overall framework of C-INSTRUCTOR (§3.2) including Trajectory Encoder and LLM Adapter.

(b) Details of STMT (§3.3). In STMT, C-INSTRUCTOR selects backtracking action that leads back to previous viewpoint.

Fig. 2: Overall framework of C-INSTRUCTOR (§3.2) and details of STMT (§3.3).

adapter consists of two components: the Trajectory Encoder and the LLM Adapter. The overall structure is shown in Fig. 2a.

Trajectory Encoder. The trajectory encoder encodes the viewpoint and action information for each step along the path into visual features. In the Matterport3D Simulator [46], a panoramic observation o_t at time step t is partitioned into K=36 subview images $\{v_{t,k}\}_{k=1}^{K}$, where the action a_t is represented using the index of the subview image corresponding to the motion direction. First, we extract visual features for each subview image using the CLIP [50] visual encoder followed by a linear projection layer with Layer Normalization [6]:

$$I_{t,k} = \texttt{layer_norm}(\texttt{linear}(f_{CLIP}(v_{t,k}))), \tag{2}$$

where $I_{t,k} \in \mathbb{R}^{1 \times D_I}$, $v_{t,k} \in \mathbb{R}^{224 \times 224 \times 3}$. To distinguish the spatial and temporal relation of each view, we add a spatial positional encoding pos_k^v and a history encoding pos_t^h to $I_{t,k}$. To represent action information, we introduced a special token pos^a for the action view a_t and another token pos^o for non-action views:

$$\hat{I}_{t,k} = \begin{cases} I_{t,k} + pos_k^v + pos_t^h + pos^a, & \text{if } k = a_t \\ I_{t,k} + pos_k^v + pos_t^h + pos^o, & \text{otherwise.} \end{cases}$$
(3)

Subsequently, we concatenate M aggregator tokens $p_{1:M}^v$ with $\hat{I}_{t,1:K}$ along the length dimension and then feed them into several ViT [13] blocks to aggregate global features for step t:

$$[\overline{\boldsymbol{p}}_{t,1:M}^{v};\overline{\boldsymbol{I}}_{t,1:K}] = f_{ViT}([\boldsymbol{p}_{1:M}^{v};\hat{\boldsymbol{I}}_{t,1:K}]), \qquad (4)$$

where $\mathbf{p}_{1:M}^{v} \in \mathbb{R}^{M \times D_{p}}$; $\overline{\mathbf{p}}_{t,1:M}^{v}$ is the trajectory feature representation at step t. **LLM Adapter.** We introduce the trajectory features into LLM via layer-wise adapting. We utilize $\mathtt{adapter}_{l}(\cdot, \cdot)$ to integrate the trajectory features $\overline{\mathbf{p}}_{t,1:M}^{v}$ into $\mathbf{x}_{l,1:S}$, which is the output of l-th LLM transformer block:

$$\widetilde{\boldsymbol{x}}_{l,1:S} = \text{adapter}_{l}(\overline{\boldsymbol{p}}_{t,1:M}^{v}, \boldsymbol{x}_{l,1:S}).$$
(5)

Here $\tilde{\boldsymbol{x}}_{l,1:S}$ replaces the $\boldsymbol{x}_{l,1:S}$ in the subsequent LLM blocks. Next, we will detail the structure of $\operatorname{adapter}_{l}(\cdot, \cdot)$. We add the trajectory features $\overline{\boldsymbol{p}}_{t,1:M}^{v}$ with the *l*th-layer's adapter query $\boldsymbol{q}_{l,1:M}$ and map them to the textual space through a linear layer linear_l(·):

$$\widetilde{p}_{l,t,1:M} = \operatorname{linear}_{l}(\overline{p}_{t,1:M}^{v} + q_{l,1:M}).$$
(6)

Next, we concatenate the $\{\widetilde{\boldsymbol{p}}_{l,t,1:M}\}_{t=1}^{T}$ in the order of t:

$$\boldsymbol{\rho}_{l,1:V} = \operatorname{concat}(\{\widetilde{\boldsymbol{p}}_{l,t,1:M}\}_{t=1}^T), \quad V = T \times M.$$
(7)

To preserve the natural language capabilities of the LLM, we use zero-initialized attention [69] to get $\tilde{x}_{l,1:S}$:

$$\widetilde{\boldsymbol{x}}_{l,1:S} = \texttt{zero_attn}([\boldsymbol{\rho}_{l,1:V}; \boldsymbol{x}_{l,1:S}]). \tag{8}$$

Based on this model structure, we design STMT (§3.3) to improve the model's spatial awareness, and CoTL (§3.4) to enhance the model's perception of landmarks. Finally, through SMT (§3.5), we achieve style-controlled instruction generation. In subsequent sections, we utilize [R; W] to denote the model's input, where R represents the path input, and W stands for the language input.

3.3 Spatial Topology Modeling Task (STMT)

Understanding the spatial relationships between different viewpoints is fundamental for generating navigation instructions. LLMs and visual encoders are typically trained on data from the Internet with few embodied-type data. Consequently, they possess limited spatial cognition abilities. Therefore, we introduce STMT as an auxiliary task to enhance the model's spatial perception capability.

In STMT, the model predicts actions between adjacent viewpoints along a trajectory. As the actions along the navigation path are already represented through location encoding, we make the model predict how to return to the previous location from the current viewpoint, as shown in Fig.2b. Given a trajectory $\{r_1, r_2, ..., r_t\}$, the model needs to predict a_t^p in order to transit from r_t back to r_{t-1} . We use prompt_a to distinguish this task and introduce a new special token \boldsymbol{x}_0^a for predicting a_t^p . The model input is:

$$[r_1, r_2, \dots, r_t; \texttt{prompt}_a, \boldsymbol{x}_0^a]. \tag{9}$$

We denote the output corresponding to \mathbf{x}_0^a at the *l*-th LLM block as $\mathbf{x}_l^a \in \mathbb{R}^{1 \times D_p}$. We then aggregate the visual features at step *t* through an attention layer:

$$\widetilde{\boldsymbol{x}}_{l}^{a} = \texttt{cross_attn}(\boldsymbol{x}_{l}^{a}, \boldsymbol{I}_{t,1:36}). \tag{10}$$

 $\tilde{\boldsymbol{x}}_{l}^{a}$ replaces \boldsymbol{x}_{l}^{a} as the input for the following layers. To mitigate the impact on the primary model and enhance training stability, the aggregation operation only starts from the output of L_{s} -th LLM block. We replace the original word prediction layer with an attention mechanism to predict a_{t}^{p} :

$$\boldsymbol{A}_{t} = \texttt{softmax}(\boldsymbol{x}_{L}^{a}\boldsymbol{W}\boldsymbol{I}_{t,1:36}^{\dagger}), \tag{11}$$

where $\boldsymbol{W} \in \mathbb{R}^{D_p \times D_I}$ is a learnable projection matrix, \boldsymbol{x}_L^a is the output of the LLM and \boldsymbol{A}_t is the predicted distribution. We apply cross entropy loss over \boldsymbol{A}_t :

$$\mathcal{L}_a = \operatorname{cross_entropy}(a_t^p, \boldsymbol{A}_t).$$
(12)

During the training process, \mathcal{L}_a is jointly optimized with the auto-regressive loss for instruction generation.

3.4 Chain of Thought with Landmarks (CoTL)

Distinguished from image or video captioning, navigation instructions encompass more than just visual descriptions. An easily executable navigation instruction usually includes several *landmarks* for directional guidance at crucial turning points. Besides, according to research in human cognitive psychology [43], it has been observed that humans, when providing path guidance, tend to first identify key navigation points within their cognitive maps before structuring their language. Therefore, the ability to determine landmarks is crucial for instruction generation. CoT [66] has been validated as an effective means of guiding the reasoning process of LLMs. Consequently, we introduce CoTL to direct the model to utilize critical landmarks in the navigation trajectory to generate instructions. Landmark Selection. For the provided annotation pairs of instructions and paths in the training set, we initially extract nouns from the instructions as linguistic landmarks $\Lambda_x = \{\lambda_n^x\}_{n=1}^{N_x}$. Since valuable landmarks may not be fully specified in the annotated instructions, we supplement the landmark set by considering the visual characteristics of the path, as shown in Fig. 3. We select visual landmarks from two perspectives, *i.e.*, the temporal perspective and the spatial perspective. From the temporal perspective, we identify crucial viewpoints along the trajectory, where landmarks are more essential for guidance. Specifically, when the trajectory leads into a new scene, e.g., transitioning from a corridor to a room, the navigator often requires a landmark for guidance. We compute the feature difference of panoramic views along a trajectory to locate these viewpoints. For a given path, we construct a sequence comprising the mean-pooled features of panoramic views $\{I_t^*\}_{t=1}^T$. We then compute the temporal importance score δ_t^{τ} via cosine distance between I_t^* and I_{t+1}^* :

$$\delta_t^{\tau} = 1 - \frac{\boldsymbol{I}_t^* \cdot \boldsymbol{I}_{t+1}^*}{\||\boldsymbol{I}_t^*\|| \cdot \||\boldsymbol{I}_{t+1}^*\|}, \quad \boldsymbol{I}_t^* = \frac{1}{K} \sum_{k=1}^K \boldsymbol{I}_{t,k},$$
(13)

where δ_t^{τ} indicates the temporal importance of landmarks appearing at time step t. From the spatial perspective, we need to identify the most distinctive object to serve as a landmark. Distinctive objects are primarily the ones that appear in the action view and not in any other candidate views. At time step t, we first extract all objects appearing in v_{t,a_t} as the candidate landmark set $\{\lambda_{t,n}^*\}_{n=1}^{N_t}$. Then, we assign distinctive scores according to the occurrence of landmarks in other candidate views. For example, the landmark $\lambda_{t,n}^*$ that also appears in candidate views $\{c_1, c_2, c_3\}$ is assigned the spatial importance score $\delta_{t,n}^a$:

$$\delta^a_{t,n} = 1 - d^a_{t,c_1} - d^a_{t,c_2} - d^a_{t,c_3},\tag{14}$$



Fig. 3: Details of Landmark Selection (left) and CoT Inference (right) in CoTL (§3.4). In Spatial Selection, candidate views are partitioned in **blue boxes**, and only objects that are distinct in action view are selected as landmarks (marked with a green tick \checkmark). In Temporal Selection, the action that leads to a new scene is treated as a significant viewpoint (marked in red box).

where d_{t,c_i}^a is the cosine distance between view a_t and view c_i . The final score for landmark $\lambda_{t,n}^*$ is:

$$\delta_{t,n} = \delta^a_{t,n} \cdot \delta^\tau_t. \tag{15}$$

We select landmarks with $\delta_{t,n} \geq \beta$ from all $\lambda^*_{t,n}$ in the trajectory to build the visual landmark set $\Lambda_v = \{\lambda^v_n\}_{n=1}^{N_v}$. Finally, the full landmark set of trajectory R can be constructed as:

$$A = A_x \cup A_v. \tag{16}$$

CoT Training and Inference. To enable the model to comprehensively identify landmarks, we utilize extracted landmarks Λ to construct training data. For a trajectory R, its corresponding data item consists of:

$$[R; \mathtt{prompt}_{\lambda}, \Lambda], \tag{17}$$

where $prompt_{\lambda}$ is the prompt for landmark generation. During training, only the Λ part is supervised.

To equip the model with the ability to generate instructions according to given landmarks, the training data for instruction generation corresponding to a path R can be constructed as:

$$[R; \texttt{prompt}_w, \Lambda_x, X], \tag{18}$$

where only the X part is supervised during training. We establish a strong correspondence between landmarks and instructions in this phase by using only Λ_x as the landmark input. This helps ensure the generation of diverse instructions by modifying landmarks.

Accordingly, the instruction generation process of the model (Fig. 3) is divided into two stages. Firstly, given a trajectory R, the model is guided by prompt_{λ} to predict landmarks M. Then, using the generated M and guided by prompt_w, the complete instruction is generated. There are two key advantages of this CoT paradigm. Firstly, it can highlight the landmarks within the path during training, enhancing the feasibility of instructions and reducing the risk of semantic errors in instruction generation. Secondly, by modifying the landmarks predicted in the first step, it allows for controlled alterations in the model's focus on landmarks in the trajectory. Further details of the prompts are discussed in the supplementary.

3.5 Style-Mixed Training (SMT)

In application, a model that can only generate step-by-step instructions is less practical. When the instruction follower is familiar with the environment, finegrained instructions lead to reduced communication efficiency. Additionally, due to the extensive amount of labor required for annotating navigation instructions, the data available is limited, especially for instructions with specified styles. This results in LLMs being susceptible to overfitting, makes it challenging to achieve accurate cross-modal mapping, and leads to suboptimal instruction generation performance when the model is trained with single-style instructions.

To mitigate the issues above, we mix datasets with instructions in different linguistic styles for training. We devise descriptions that encapsulate diverse styles into prompts to enable the LLM to generate in different styles. By employing SMT, not only is the quality of instruction generation enhanced, but we also enable a single LLM instance to adaptively generate different styles of instructions for the same path R by switching between different prompts.

4 Experiments

4.1 Datasets and Evaluation Metrics

Datasets. We evaluate the instruction generation performance on three indoor navigation datasets [5,31,48] and one outdoor navigation dataset [26]:

- R2R [5]: It has four splits with step-by-step instructions, *i.e.*, train (61 scenes, 14,039 instructions), val seen (61 scenes, 1,021 instructions), val unseen (11 scenes, 2,349 instructions), and test unseen (18 scenes, 4,173 instructions). As test unseen is reserved for benchmarking instruction followers, we report the performance of instruction generation on val splits.
- REVERIE [48]: It contains high-level descriptions of target destinations and objects. It has three open-access splits, *i.e.*, train (61 scenes, 10, 466 instructions), val seen (61 scenes, 1, 371 instructions), and val unseen (10 scenes, 3, 753 instructions). We report the performance on two val splits.
- RxR [31]: It is a multilingual indoor navigation dataset with longer trajectories and more fine-grained aligned instructions. we specifically utilize the English instructions for comparison with previous methods. It has three publicly available splits, and we report the performance on two val splits.
- UrbanWalk [26]: It is an outdoor navigation dataset with 26,808 imageinstruction pairs simulated by CARLA [14]. We follow the setting in [68].

Table 1: Comparison to state-of-the-art methods (§4.3) on R2R [5].

Mathada	R2R val seen					R2R val unseen						
Methods	SPICE↑	BLEU-1 \uparrow	BLEU-4↑	$\mathtt{CIDEr}\uparrow$	$\texttt{Meteor}\uparrow$	$\texttt{Rouge} \uparrow$	$SPICE \uparrow$	BLEU-1↑	BLEU-4 \uparrow	$\mathtt{CIDEr}\uparrow$	$\texttt{Meteor}\uparrow$	Rouge \uparrow
BT-speaker [16] [NeurIPS2018]	0.173	0.670	0.236	0.373	0.209	0.443	0.113	0.600	0.149	0.113	0.167	0.376
EDrop-speaker [52] [NAACL2019]	0.168	0.660	0.228	0.362	0.208	0.447	0.117	0.590	0.157	0.160	0.174	0.389
CCC-speaker [58] [CVPR2022]	0.194	0.698	0.265	0.449	0.218	0.467	0.108	0.591	0.139	0.120	0.164	0.375
Lana [63] [CVPR2023]	0.170	0.657	0.215	0.265	0.205	0.433	0.174	0.667	0.236	0.295	0.213	0.448
C-Instructor w/o SMT	0.230	0.732	0.270	0.511	0.237	0.475	0.217	0.715	0.263	0.453	0.234	0.470
C-Instructor (Ours)	0.233	0.726	0.276	0.529	0.247	0.480	0.212	0.713	0.266	0.447	0.239	0.473

The val unseen splits in R2R [5], REVERIE [48], and RxR [31] contain trajectories whose corresponding scenes are not included in train splits, and thus are good testbeds for generalizability [11,15,68,70]. Consequently, we focus on those splits to better validate the generalizability of C-INSTRUCTOR.

Evaluation Metrics. We evaluate the linguistic quality of generated instructions with widely-used automatic text similarity metrics, including BLEU [47], SPICE [4], CIDEr [57], Meteor [7], and Rouge [37]. For each navigation path, all corresponding ground-truth instructions are used as references.

4.2 Implementation Details

Detailed Architecture. We use the multimodal LLaMA-Adapter [20] with 32 layers and 7B parameters as the LLM. We adopt CLIP-ViT-L-14 [50] and 8 ViT [13] blocks in the Trajectory Encoder. The score threshold β for landmark selection in §3.4 is set to 0.25, and L_s in §3.3 is set to 30.

Training. We only finetune the last 2 layers of LLM while fixing the other 30 layers. The CLIP [50] visual encoder is also fixed. We first pre-train C-INSTRUCTOR on PREVALENT [21] for 240K iterations with a batch size of 16, and then finetune C-INSTRUCTOR on multiple datasets jointly for 120K iterations with batch size 4. We use the AdamW [42] optimizer with base learning rate 1.0×10^{-4} . Four NVIDIA A100 80GB GPUs are used for training.

Inference. We set the generation temperature to 1.0 for RxR [31], and 0.1 for all other datasets. All other hyperparameters remain the same as [20].

4.3 Comparison to State-of-the-Art Methods

We compare C-INSTRUCTOR with four existing instruction generation models. For a fair comparison, we report the performance of C-INSTRUCTOR without SMT in addition to the performance of the full model. We employ the Penn Treebank tokenizer [53] to compute the linguistic metrics.

R2R [5]. The results on R2R are summarized in Tab. 1. C-INSTRUCTOR outperforms previous methods under all metrics on both val splits. In terms of SPICE, C-INSTRUCTOR demonstrates a superiority of 3.9% in absolute terms and 20.1% in relative terms on val seen as well as 3.8% in absolute terms and 21.8% in relative terms on val unseen compared to the previous best. This verifies that C-INSTRUCTOR exhibits good performance in generating fine-grained directives.

Table 2	2: Comparison	to state-of-the-art	methods (§4.3) on REVERIE	[48]
Table 2	2. Comparison	10 Statt-01-0110-art	memous (gr.o) OII ICL V LICIL	110

Mathada	REVERIE val seen						REVERIE val unseen					
Methods	SPICE↑	BLEU-1↑	BLEU-4↑	$CIDEr\uparrow$	Meteor↑	$\texttt{Rouge} \uparrow$	$\text{SPICE} \uparrow$	BLEU-1 \uparrow	BLEU-4↑	$\mathtt{CIDEr}\uparrow$	$\texttt{Meteor}\uparrow$	Rouge \uparrow
BT-speaker [16] [NeurIPS2018]	0.121	0.693	0.347	0.269	0.223	0.602	0.103	0.664	0.302	0.190	0.200	0.569
EDrop-speaker [52] [NAACL2019]	0.138	0.641	0.360	0.523	0.277	0.597	0.114	0.648	0.319	0.333	0.233	0.546
CCC-speaker [58] [CVPR2022]	0.144	0.727	0.408	0.502	0.272	0.589	0.115	0.681	0.357	0.334	0.232	0.548
Lana [63] [CVPR2023]	0.137	0.707	0.404	0.627	0.282	0.619	0.107	0.696	0.345	0.327	0.239	0.582
C-Instructor w/o SMT	0.184	0.785	0.480	0.844	0.319	0.649	0.139	0.739	0.369	0.464	0.259	0.577
C-Instructor (Ours)	0.182	0.775	0.459	0.805	0.311	0.647	0.141	0.754	0.419	0.545	0.267	0.591

Table 3: Comparison to state-of-the-art methods (§4.3) on RxR [31].

Mathada		Rx	R val se	en		RxR val unseen				
Methods	BLEU-1↑	BLEU-4 \uparrow	$\mathtt{CIDEr}\uparrow$	$\texttt{Meteor}\uparrow$	$\texttt{Rouge}\uparrow$	BLEU-1 \uparrow	BLEU-4 \uparrow	$\mathtt{CIDEr}\uparrow$	$\texttt{Meteor}\uparrow$	$\texttt{Rouge}\uparrow$
BT-speaker [16] [NeurIPS2018]	0.514	0.188	0.026	0.204	0.365	0.566	0.211	0.024	0.208	0.372
EDrop-speaker [52] [NAACL2019]	0.595	0.197	0.047	0.213	0.378	0.568	0.184	0.038	0.205	0.370
CCC-speaker [58] [CVPR2022]	0.526	0.194	0.024	0.185	0.355	0.518	0.187	0.026	0.184	0.353
Lana [63] [CVPR2023]	0.342	0.123	0.040	0.128	0.275	0.319	0.115	0.043	0.124	0.273
C-Instructor w/o SMT	0.683	0.233	0.081	0.243	0.381	0.667	0.224	0.080	0.236	0.379
C-Instructor (Ours)	0.685	0.234	0.082	0.238	0.382	0.678	0.233	0.077	0.239	0.382

REVERIE [48]. As depicted in Tab. 2, C-INSTRUCTOR also attains state-ofthe-art performance in generating high-level trajectory descriptions. It exhibits a relative improvement of 26.4% on val seen and 22.6% on val unseen in terms of SPICE, which is more pronounced compared to R2R [5].

RxR [31]. As shown in Tab. 3, C-INSTRUCTOR significantly outperforms existing instruction generation algorithms in all metrics. This suggests that C-INSTRUCTOR possesses the capability to manage visual contexts of extended trajectory and generate more intricate instructions.

UrbanWalk [26]. As shown in Tab. 4, C-INSTRUCTOR also achieves the best performance under all metrics on outdoor scenes. This indicates that our C-INSTRUCTOR possesses strong generalization capability and universality.

4.4 Diagnostic Experiment

To thoroughly study the effectiveness of C-INSTRUCTOR, we compare the full model with several ablative designs. We test the ablative models on REVERIE [48] and R2R [5] val unseen. The results are summarized in Tab. 5.

Mathada		UrbanWalk								
Methods	,	SPICE \uparrow	BLEU-1 \uparrow	BLEU-4 \uparrow	Meteor \uparrow	Rouge \uparrow				
BT-speaker [16]	[NeurIPS2018]	0.524	0.649	0.408	0.350	0.620				
EDrop-speaker [52]	[NAACL2019]	0.531	0.689	0.435	0.358	0.634				
ASSISTER [26]	[ECCV2022]	0.451	0.576	0.164	0.319	0.557				
Kefa-speaker [68]	[Arxiv2023]	0.566	0.711	0.450	0.378	0.655				
C-INSTRUCTOR	(Ours)	0.645	0.771	0.534	0.461	0.781				

Table 4: Comparison to state-of-the-art methods (§4.3) on UrbanWalk [26].

Table 5: Ablation study (§4.4) on REVERIE [48] val unseen and R2R [5] val unseen.

Methods		REVER	unseen	R2R val unseen						
memous	BLEU-1 \uparrow	BLEU-4 \uparrow	$\mathtt{CIDEr}\uparrow$	$\texttt{Meteor}\uparrow$	$\texttt{Rouge}\uparrow$	BLEU-1↑	BLEU-4 \uparrow	$\mathtt{CIDEr}\uparrow$	$\texttt{Meteor}\uparrow$	$\texttt{Rouge} \uparrow$
Vanilla LLM	0.399	0.131	0.432	0.156	0.400	0.307	0.059	0.292	0.139	0.303
Baseline	0.648	0.308	0.347	0.248	0.547	0.676	0.232	0.356	0.225	0.449
Baseline + SMT	0.679	0.344	0.397	0.254	0.562	0.685	0.254	0.407	0.233	0.466
Baseline $+$ SMT $+$ STMT	0.737	0.402	0.490	0.258	0.590	0.689	0.262	0.445	0.228	0.479
Baseline $+$ SMT $+$ STMT $+$ CoTL	0.754	0.419	0.545	0.267	0.591	0.713	0.266	0.447	0.239	0.473

Vanilla LLM. We assess the performance of vanilla LLM by captioning views along the trajectory using BLIP [35] and feeding those captions with devised prompts into pre-trained LLaMA [20] to generate navigation instructions. The performance of this vanilla method fine-tuned on REVERIE [48] and R2R [5] respectively (#1) remains largely inferior to the baseline in §3.2 (#2), which still significantly lags behind our full method (#5). This underscores the inherent information loss through captioning as well as the effectiveness of our design.

SMT. To train a model with instructions from diverse domains yields performance benefits. In comparison to #2, the model trained using SMT (#3) exhibits an improvement in SPICE on the REVERIE val unseen from 0.127 to 0.129. It concurrently achieves a performance improvement on the R2R val unseen. This suggests that enhancing linguistic diversity will foster the quality of instructions generated by C-INSTRUCTOR.

STMT. The model trained with STMT (#4) demonstrates a notable impact on generating highly abstract instructions. It lifts BLEU-4 from 0.344 to 0.402 and CIDEr from 0.397 to 0.490 on the REVERIE val unseen. This highlights the significance of understanding the environment layout.

CoTL. Compared to #4, the model with CoTL (#5) significantly improves the semantic consistency with the ground truth instruction. The improvement on REVERIE is more significant: SPICE increases from 0.129 to 0.141. This suggests that incorporating CoTL enhances the alignment between generated instructions and the visual environment, especially for high-level instructions.

4.5 Instruction Quality Analysis

Evaluating the quality of instructions solely based on text similarity metrics is insufficient as those metrics do not thoroughly assess the semantic alignment between instructions and trajectories. Thus, we further analyze the semantic quality of instructions generated by C-INSTRUCTOR from three aspects through the following experiments:

Path Guiding Proficiency. The success rate (SR) of navigators with instructions from different instruction generators can be used as an index for the quality of instructions. We regenerate instructions for the paths in REVERIE [48] val unseen and employ two navigators (HAMT [10] and DUET [11]) to assess SR and SPL (SR weighted by Path Length) when guided by regenerated instructions. As depicted in Tab. 6b, SR and SPL of instructions provided by C-INSTRUCTOR

Data Source	REVERIE val unseen						
Data Source	$\mathtt{SR}\uparrow$	$\mathtt{SPL}\uparrow$	$\mathtt{RGS}\uparrow$	$\mathtt{RGSPL}\uparrow$			
Original [48]	32.95	30.20	18.92	17.28			
+BT-speaker [16]	31.84	28.37	17.35	15.14			
+EDrop-speaker [52]	30.45	27.18	18.60	16.24			
+CCC-speaker [58]	29.65	26.20	16.33	14.58			
+Lana [63]	33.05	29.76	19.14	17.20			
+C-Instructor (Ours)	34.25	31.25	19.99	18.08			

Table 6: Instruction quality analysis based on performance of navigation models (§4.5).

(a) Performance of HAMT [10] using different instruction generator for data augmentation on REVERIE [48] val unseen (§4.5). Training with instructions generated by C-INSTRUCTOR yields the most significant improvement.

Follower HAMT [10] Instruction Generator DUET [11] SR↑ SPL ↑ $SR\uparrow$ $SPL\uparrow$ 33.73 32.95 Human annotation [48] 30.20 46.98BT-speaker [16] 24.8521.7430.4721.46EDrop-speaker [52] 26.1923.5527.8917.00CCC-speaker [58] 23.2920.6929.7419.55Lana [63] 26.8424.3831.39 20.44C-INSTRUCTOR (Ours) 31.35 29.27 43.34 30.13

(b) Performance of HAMT [10] and DUET [11] in following instructions generated on REVERIE [48] val unseen (§4.5). SR and SPL are provided as metrics to evaluate the pathguiding proficiency of different instruction generation models.

significantly exceeds that of those generated by prior models and remarkably aligns with the navigation accuracy of human-annotated instructions.

Data Augmentation. The enhancement of navigation accuracy of instruction followers via data augmentation can also serve as an indicator for the improved quality of instruction generation. Hence, we leverage 17,533 instructions generated by various instruction generation models on randomly sampled paths along with the original train split of REVERIE [48] to train HAMT [10]. As shown in Tab. 6a, the model utilizing data generated by C-INSTRUCTOR exhibits an increase in the accuracy of navigation including SR, SPL, RGS (Remote Grounding Success rate), and RGSPL (RGS weighted by Path Length). RGS and RGSPL measure the success rate of the agent's finding the target object indicated in the given instruction and are used as navigator performance metrics on REVERIE [48]. In contrast, employing other models for data augmentation results in an unintended performance drop for the navigator. This indicates that C-INSTRUCTOR, when utilized as a means of data augmentation, exhibits superior efficacy in generating instructions with high-level abstraction. User Study. To provide a more comprehensive evaluation of the semantic quality of generated instructions, we conduct a series of human evaluations. Specifically, 15 college students are individually tasked with scoring from 0 to 5 according to the semantic alignment between the given instructions and the corresponding trajectories. The instructions provided are generated by C-INSTRUCTOR, Lana [63], CCC [58], BT-Speaker [16], and EnvDrop-Speaker [52] from a total of 100 paths. The paths are sampled from the val unseen split of REVERIE [48]. C-INSTRUCTOR garners a higher average score, *i.e.*, 3.50, vs Lana 2.26, CCC 2.14, BT-Speaker 2.10 and EnvDrop-Speaker 2.10.

4.6 Qualitative Results

We visualize an example of indoor navigation trajectory and corresponding instruction generation results in Fig. 4. As seen, C-INSTRUCTOR can identify critical landmarks in the path and generate high-quality instructions accordingly in specified styles. Moreover, we can control the focus of C-INSTRUCTOR by mod-



Fig. 4: Visualizations of navigation trajectory and instruction generation results on R2R [5] and REVERIE [48] (§4.6).



Fig. 5: Visualizations of path and generated instruction on UrbanWalk [26] (§4.6).

ifying landmarks. Fig. 5 displays a result on UrbanWalk [26]. We can observe that C-INSTRUCTOR can also provide practical instructions in outdoor scenes.

5 Conclusion and Discussion

In this work, we propose C-INSTRUCTOR, which generates style-controllable and content-controllable instructions with high linguistic quality. It uses an adapterbased structure to leverage the language capability of LLMs and distinct style prompts in SMT to achieve style control. To enhance the executability of generated instructions, we adopt CoTL to help identify crucial landmarks and provide content controllability. We also devise STMT to enhance the model's understanding of the environment's spatial topology. The instructions generated by C-INSTRUCTOR not only achieve high scores in text metrics but also demonstrate strong competence in guiding navigators, further validating the strong correspondence between generated instructions and given trajectories. We expect that C-INSTRUCTOR can greatly enhance agent-human communication and significantly contribute to the development of versatile embodied agents.

Acknowledgments. This research is supported in part by the National Science and Technology Major Project (No. 2022ZD0115502 and 2023ZD0121300), the National Natural Science Foundation of China (No. 62122010, U23B2010, and 62372405), Zhejiang Provincial Natural Science Foundation of China (No. LDT23F02022F02), Beijing Natural Science Foundation (No. L231011), the Fundamental Research Funds for the Central Universities (No. 226-2024-00058), and Beihang World TOP University Cooperation Program.

References

- 1. Allen, G.L.: From knowledge to words to wayfinding: Issues in the production and comprehension of route directions. In: International Conference on Spatial Information Theory (1997)
- An, D., Qi, Y., Li, Y., Huang, Y., Wang, L., Tan, T., Shao, J.: Bevbert: Multimodal map pre-training for language-guided navigation. In: ICCV (2023)
- An, D., Wang, H., Wang, W., Wang, Z., Huang, Y., He, K., Wang, L.: Etpnav: Evolving topological planning for vision-language navigation in continuous environments. IEEE TPAMI (2024)
- 4. Anderson, P., Fernando, B., Johnson, M., Gould, S.: Spice: Semantic propositional image caption evaluation. In: ECCV (2016)
- Anderson, P., Wu, Q., Teney, D., Bruce, J., Johnson, M., Sünderhauf, N., Reid, I., Gould, S., Van Den Hengel, A.: Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In: CVPR (2018)
- Ba, J.L., Kiros, J.R., Hinton, G.E.: Layer normalization. arXiv preprint arXiv:1607.06450 (2016)
- 7. Banerjee, S., Lavie, A.: Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In: ACL Workshop (2005)
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. NeurIPS (2020)
- Chen, J., Gao, C., Meng, E., Zhang, Q., Liu, S.: Reinforced structured stateevolution for vision-language navigation. In: CVPR (2022)
- 10. Chen, S., Guhur, P.L., Schmid, C., Laptev, I.: History aware multimodal transformer for vision-and-language navigation. NeurIPS (2021)
- 11. Chen, S., Guhur, P.L., Tapaswi, M., Schmid, C., Laptev, I.: Think global, act local: Dual-scale graph transformer for vision-and-language navigation. In: CVPR (2022)
- Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
- 14. Dosovitskiy, A., Ros, G., Codevilla, F., Lopez, A., Koltun, V.: Carla: An open urban driving simulator. In: CoRL (2017)
- Dou, Z.Y., Peng, N.: Foam: A follower-aware speaker model for vision-and-language navigation. In: NAACL (2022)
- Fried, D., Hu, R., Cirik, V., Rohrbach, A., Andreas, J., Morency, L.P., Berg-Kirkpatrick, T., Saenko, K., Klein, D., Darrell, T.: Speaker-follower models for vision-and-language navigation. NeurIPS (2018)
- 17. Gao, C., Chen, J., Liu, S., Wang, L., Zhang, Q., Wu, Q.: Room-and-object aware knowledge reasoning for remote embodied referring expression. In: CVPR (2021)
- Gao, C., Liu, S., Chen, J., Wang, L., Wu, Q., Li, B., Tian, Q.: Room-object entity prompting and reasoning for embodied referring expression. IEEE TPAMI (2023)
- Gao, C., Peng, X., Yan, M., Wang, H., Yang, L., Ren, H., Li, H., Liu, S.: Adaptive zone-aware hierarchical planner for vision-language navigation. In: CVPR (2023)
- Gao, P., Han, J., Zhang, R., Lin, Z., Geng, S., Zhou, A., Zhang, W., Lu, P., He, C., Yue, X., et al.: Llama-adapter v2: Parameter-efficient visual instruction model. arXiv preprint arXiv:2304.15010 (2023)

- 16 X. Kong et al.
- Hao, W., Li, C., Li, X., Carin, L., Gao, J.: Towards learning a generic agent for vision-and-language navigation via pre-training. In: CVPR (2020)
- He, K., Huang, Y., Wu, Q., Yang, J., An, D., Sima, S., Wang, L.: Landmark-rxr: Solving vision-and-language navigation with fine-grained alignment supervision. NeurIPS (2021)
- 23. He, K., Si, C., Lu, Z., Huang, Y., Wang, L., Wang, X.: Frequency-enhanced data augmentation for vision-and-language navigation. NeurIPS (2024)
- Hong, Y., Rodriguez, C., Wu, Q., Gould, S.: Sub-instruction aware vision-andlanguage navigation. In: EMNLP (2020)
- Hu, E.J., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W., et al.: Lora: Low-rank adaptation of large language models. In: ICLR (2021)
- 26. Huang, Z., Shangguan, Z., Zhang, J., Bar, G., Boyd, M., Ohn-Bar, E.: Assister: Assistive navigation via conditional instruction generation. In: ECCV (2022)
- 27. Hund, A.M., Minarik, J.L.: Getting from here to there: Spatial anxiety, wayfinding strategies, direction type, and wayfinding efficiency. Spatial cognition and computation (2006)
- Jain, V., Magalhaes, G., Ku, A., Vaswani, A., Ie, E., Baldridge, J.: Stay on the path: Instruction fidelity in vision-and-language navigation. In: ACL (Jan 2019)
- Kamath, A., Anderson, P., Wang, S., Koh, J., Ku, A., Waters, A., Yang, Y., Baldridge, J., Parekh, Z.: A new path: Scaling vision-and-language navigation with synthetic instructions and imitation learning. In: CVPR (2023)
- 30. Karimi Mahabadi, R., Henderson, J., Ruder, S.: Compacter: Efficient low-rank hypercomplex adapter layers. NeurIPS (2021)
- Ku, A., Anderson, P., Patel, R., Ie, E., Baldridge, J.: Room-across-room: Multilingual vision-and-language navigation with dense spatiotemporal grounding. In: EMNLP (2020)
- 32. Kuipers, B.: Modeling spatial knowledge. Cognitive science (1978)
- 33. Lester, B., Al-Rfou, R., Constant, N.: The power of scale for parameter-efficient prompt tuning. In: EMNLP (2021)
- 34. Li, J., Bansal, M.: Panogen: Text-conditioned panoramic environment generation for vision-and-language navigation. NeurIPS (2024)
- 35. Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In: ICML. PMLR (2022)
- Li, X.L., Liang, P.: Prefix-tuning: Optimizing continuous prompts for generation. In: ACL-IJCNLP (2021)
- 37. Lin, C.Y.: Rouge: A package for automatic evaluation of summaries. In: Text summarization branches out (2004)
- Liu, C., Zhu, F., Chang, X., Liang, X., Ge, Z., Shen, Y.D.: Vision-language navigation with random environmental mixup. In: ICCV (2021)
- Liu, R., Wang, W., Yang, Y.: Volumetric environment representation for visionlanguage navigation. In: CVPR (2024)
- 40. Liu, R., Wang, X., Wang, W., Yang, Y.: Bird's-eye-view scene graph for visionlanguage navigation. In: ICCV (2023)
- Liu, X., Zheng, Y., Du, Z., Ding, M., Qian, Y., Yang, Z., Tang, J.: Gpt understands, too. AI Open (2023)
- 42. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101 (2017)
- 43. Lynch, K.: The image of the city (1964)
- 44. Mehta, H., Artzi, Y., Baldridge, J., Ie, E., Mirowski, P.: Retouchdown: Adding touchdown to streetlearn as a shareable resource for language grounding tasks in street view. arXiv preprint arXiv:2001.03671 (2020)

- Nguyen, K., Daumé III, H.: Help, anna! visual navigation with natural multimodal assistance via retrospective curiosity-encouraging imitation learning. In: EMNLP-IJCNLP (2019)
- 46. Nguyen, K., Dey, D., Brockett, C., Dolan, B.: Vision-based navigation with language-based assistance via imitation learning with indirect intervention. In: CVPR (2018)
- 47. Papineni, K., Roukos, S., Ward, T., Zhu, W.J.: Bleu: a method for automatic evaluation of machine translation. In: ACL (2002)
- Qi, Y., Wu, Q., Anderson, P., Wang, X., Wang, W.Y., Shen, C., Hengel, A.v.d.: Reverie: Remote embodied visual referring expression in real indoor environments. In: CVPR (2020)
- 49. Qiao, Y., Qi, Y., Yu, Z., Liu, J., Wu, Q.: March in chat: Interactive prompting for remote embodied referring expression. In: ICCV (2023)
- Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: ICML (2021)
- Shridhar, M., Thomason, J., Gordon, D., Bisk, Y., Han, W., Mottaghi, R., Zettlemoyer, L., Fox, D.: Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In: CVPR (2020)
- 52. Tan, H., Yu, L., Bansal, M.: Learning to navigate unseen environments: Back translation with environmental dropout. In: NAACL (2019)
- Taylor, A., Marcus, M., Santorini, B.: The penn treebank: an overview. Treebanks: Building and using parsed corpora pp. 5–22 (2003)
- Thomason, J., Murray, M., Cakmak, M., Zettlemoyer, L.: Vision-and-dialog navigation. In: CoRL (2020)
- 55. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al.: Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023)
- 56. Vanetti, E.J., Allen, G.L.: Communicating environmental knowledge: The impact of verbal and spatial abilities on the production and comprehension of route directions. Environ Behav (1988)
- Vedantam, R., Lawrence Zitnick, C., Parikh, D.: Cider: Consensus-based image description evaluation. In: CVPR (2015)
- Wang, H., Liang, W., Shen, J., Van Gool, L., Wang, W.: Counterfactual cycleconsistent learning for instruction following and generation in vision-language navigation. In: CVPR (2022)
- Wang, H., Liang, W., Van Gool, L., Wang, W.: Dreamwalker: Mental planning for continuous vision-language navigation. In: ICCV (2023)
- Wang, H., Wang, W., Liang, W., Xiong, C., Shen, J.: Structured scene memory for vision-language navigation. In: CVPR (2021)
- Wang, H., Wang, W., Shu, T., Liang, W., Shen, J.: Active visual information gathering for vision-language navigation. In: ECCV (2020)
- Wang, S., Montgomery, C., Orbay, J., Birodkar, V., Faust, A., Gur, I., Jaques, N., Waters, A., Baldridge, J., Anderson, P.: Less is more: Generating grounded navigation instructions from landmarks. In: CVPR (2022)
- Wang, X., Wang, W., Shao, J., Yang, Y.: Lana: A language-capable navigator for instruction following and generation. In: CVPR (2023)
- 64. Wang, Z., Li, X., Yang, J., Liu, Y., Jiang, S.: Gridmm: Grid memory map for vision-and-language navigation. In: ICCV (2023)

- 18 X. Kong et al.
- Ward, S.L., Newcombe, N., Overton, W.F.: Turn left at the church, or three miles north a study of direction giving and sex differences. Environment and Behavior (1986)
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q.V., Zhou, D., et al.: Chain-of-thought prompting elicits reasoning in large language models. NeurIPS (2022)
- 67. Yang, Z., Chen, G., Li, X., Wang, W., Yang, Y.: Doraemongpt: Toward understanding dynamic scenes with large language models (exemplified as a video agent). In: Forty-first International Conference on Machine Learning (2024)
- Zeng, H., Wang, X., Wang, W., Yang, Y.: Kefa: A knowledge enhanced and fine-grained aligned speaker for navigation instruction generation. arXiv preprint arXiv:2307.13368 (2023)
- Zhang, R., Han, J., Zhou, A., Hu, X., Yan, S., Lu, P., Li, H., Gao, P., Qiao, Y.: Llama-adapter: Efficient fine-tuning of language models with zero-init attention. arXiv preprint arXiv:2303.16199 (2023)
- 70. Zhang, Y., Kordjamshidi, P.: Vln-trans, translator for the vision and language navigation agent. In: ACL (2023)
- Zhao, Y., Chen, J., Gao, C., Wang, W., Yang, L., Ren, H., Xia, H., Liu, S.: Targetdriven structured transformer planner for vision-language navigation. In: ACM MM (2022)
- 72. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Learning to prompt for vision-language models. IJCV (2022)
- 73. Zhu, F., Liang, X., Zhu, Y., Yu, Q., Chang, X., Liang, X.: Soon: Scenario oriented object navigation with graph-based exploration. In: CVPR (2021)
- 74. Zhu, W., Hu, H., Chen, J., Deng, Z., Jain, V., Ie, E., Sha, F.: Babywalk: Going farther in vision-and-language navigation by taking baby steps. In: ACL (2020)