

NavCoT: Boosting LLM-Based Vision-and-Language Navigation via Learning Disentangled Reasoning

Bingqian Lin*, Yunshuang Nie*, Ziming Wei, Jiaqi Chen, Shikui Ma, Jianhua Han, Hang Xu, Xiaojun Chang, Xiaodan Liang†

Abstract—Vision-and-Language Navigation (VLN), as a crucial research problem of Embodied AI, requires an embodied agent to navigate through complex 3D environments following natural language instructions. Recent research has highlighted the promising capacity of large language models (LLMs) in VLN by improving navigational reasoning accuracy and interpretability. However, their predominant use in an offline manner usually suffers from substantial domain gap between the VLN task and the LLM training corpus. This paper introduces a novel strategy called **Navigational Chain-of-Thought (NavCoT)**, where we fulfill parameter-efficient in-domain training to enable self-guided navigational decision, leading to a significant mitigation of the domain gap in a cost-effective manner. Specifically, at each timestep, the LLM is prompted to forecast the navigational chain-of-thought by: 1) acting as a world model to imagine the next observation according to the instruction, 2) selecting the candidate observation that best aligns with the imagination, and 3) determining the action based on the reasoning from the prior steps. In this way, the action prediction can be effectively simplified benefiting from the disentangled reasoning. Through constructing formalized labels for training, the LLM can learn to generate desired and reasonable chain-of-thought outputs for improving the action decision. Experimental results across various training settings and popular VLN benchmarks (e.g., Room-to-Room (R2R), Room-across-Room (RxR), Room-for-Room (R4R)) show the significant superiority of NavCoT over the direct action prediction variants. Through simple parameter-efficient finetuning, our NavCoT outperforms a recent GPT4-based approach with $\sim 7\%$ relative improvement on the R2R dataset. We believe that NavCoT will help unlock more task-adaptive and scalable LLM-based embodied agents, which are helpful for developing real-world robotics applications. Code is available at <https://github.com/expectorlin/NavCoT>.

Index Terms—Vision-and-language navigation, large language models, disentangled reasoning

A INTRODUCTION

IN Vision-and-Language Navigation (VLN) [1], [2], [3], [4], [5], an embodied agent is required to reach the target position following a language instruction. As a representative Embodied AI task, VLN has attracted increasing attention in recent years for its practicality and flexibility. It imposes great challenges on the embodied agent since successful navigation requires complex reasoning ability, e.g., long-term planning for following different sub-instructions and monitoring the navigation progress.

With the rapid development of the large language models (LLMs) [6], [7], [8], emerging works have attempted to introduce LLMs to solve Embodied AI tasks due to their rich real-world commonsense and powerful reasoning ability [9],

[10], [11]. These works have revealed the great potential of LLMs for assisting the embodied task completion. To enable LLMs to interact with the physical world, some works have introduced the off-the-shelf vision-to-text system [12], [13] to transform the visual information into a linguistic representation. Then, the LLM can reason the action according to the textual representation of the surrounding observation. A few recent VLN works also introduce the LLM as the navigation backbone to study how LLM can improve navigation action decisions [14], [15]. However, they tend to utilize some high-cost LLMs such as GPT-4 [16], which suffers from poor scalability and a large domain gap with the VLN tasks. Moreover, LLMs are required to make action decisions straightforwardly without guidance about how to filter noisy textual-represented visual information.

In this paper, we propose **Navigational Chain-of-Thought (NavCoT)**, where we conduct parameter-efficient in-domain training to enable LLMs to perform self-guided navigational reasoning for facilitating action decisions. Inspired by the world model theory [17], [18], when humans interact with the world, we tend to build a mental model that summarizes the surroundings we have seen before and helps us to predict the future. Then, we can make action decisions sequentially to complete different tasks based on this mental model. Therefore, we adapt the above process to the Chain-of-Thought (CoT) [19] reasoning mechanism in a trainable manner. The resulting strategy, termed *nav-*

• *These two authors contribute equally to this work.

• †Xiaodan Liang is the corresponding author.

• Bingqian Lin, Yunshuang Nie, Ziming Wei and Xiaodan Liang are with Shenzhen Campus of Sun Yat-sen University, Shenzhen, China.

E-mail: {linbq6@mail2.sysu.edu.cn, nieysh@mail2.sysu.edu.cn, weizm3@mail2.sysu.edu.cn, liangxd9@mail.sysu.edu.cn}

• Jiaqi Chen is with the University of Hong Kong.

E-mail: jqchen@cs.hku.hk.

• Shikui Ma is with Dataa Robotics company.

E-mail: maskey.bj@gmail.com.

• Jianhua Han and Hang Xu are with Huawei Noah's Ark Lab.

E-mail: hanjianhua4@huawei.com, chromexbjxh@gmail.com.

• Xiaojun Chang is with ReLER, AAIL, University of Technology Sydney.

E-mail: Xiaojun.chang@uts.edu.au.

igational chain-of-thought, transforms the LLM into both a world model and a navigation reasoning agent, i.e., the LLM learns to imagine the future surroundings, filter the confusing observations based on the imagination, and then make final action decisions at each navigation timestep with the customized chain-of-thought labels. As shown in Fig. 1, through NavCoT, the LLM is able to filter redundant visual information based on the imagination and therefore significantly simplify the action decision.

We explore various training strategies, including pre-training and finetuning, full data and low-resource settings, to study how in-domain data can contribute to the performance improvement of LLM-based VLN comprehensively. To facilitate training, we create formalized ground-truths to constrain the LLM to generate navigational chain-of-thoughts with a unified format. For encouraging better instruction following, we constrain the imagination labels at each navigation timestep to be one of the mentioned objects/scenes in a given instruction, which essentially enables the construction of a task-oriented world model and significantly simplifies the training in the meanwhile. We adopt parameter-efficient finetuning, which can be supported by a single NVIDIA V100 GPU for two recently proposed language models (LLaMA-Adapter [20] and LLaMA 2 [8]), for improving the scalability and efficiency.

We conduct experiments on various popular VLN benchmarks, including R2R [1], R4R [4], RxR [5], and REVERIE [2]. Experimental results show that NavCoT significantly outperforms both the direct action prediction and zero-shot inference variants, demonstrating the effectiveness of the navigational chain-of-thoughts generation in a trainable manner. Through simple parameter-efficient finetuning, NavCoT surpasses a recent GPT4-based VLN model [14] by ~ 7 points in both SR and SPL on R2R.

To summarize, the main contributions of this paper are:

- We introduce NavCoT, where we repurpose the LLM to be both a world model and a navigation reasoning agent in a trainable manner to simplify the action decision process and improve interpretability.
- We adopt parameter-efficient in-domain training for adapting LLMs to the VLN task in a cost-effective way, making a solid step towards developing scalable LLM-based VLN methods.
- Experimental results show the superiority of NavCoT over high-cost LLM-based approaches and direct action prediction variants on multiple VLN datasets. Through explicit reasoning generation, NavCoT also exhibits much better explainability than traditional cross-modal based VLN models.

B RELATED WORK

B.1 Vision-Language Navigation

VLN has received great attention and many works have been proposed in the past few years. Early approaches mainly focus on exploring data augmentation techniques [21], [22], [23], [24] and useful model architectures [25], [26], [27], [28] to alleviate data scarcity and improve cross-modal alignment. To further improve the generalization to unseen environments, pretraining-based

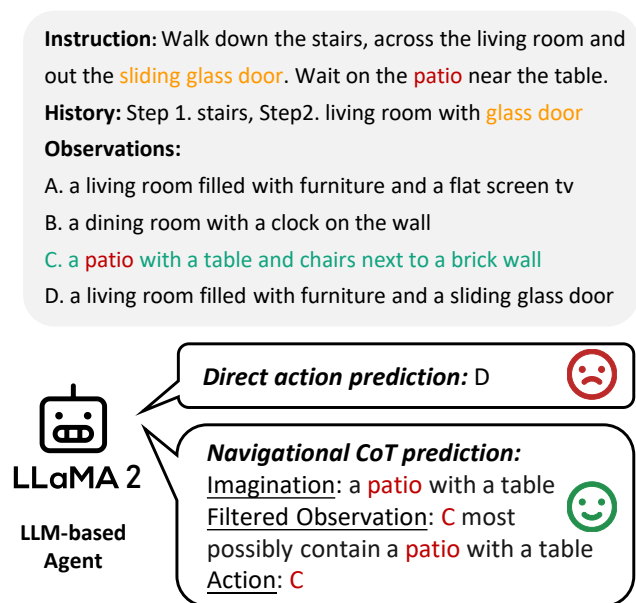


Fig. 1: Comparison between direct action decision and our NavCoT. According to the instruction (finding *patio* after the *sliding glass door*) and history (*glass door*), NavCoT successfully predicts the future imagination *patio*, selects the observation C that best matches the imagination and determines the correct action.

approaches [29], [30], [31], [32], [33], [34], [35] have been widely developed in the VLN field. However, adapting to realistic application scenarios that require rich common-sense knowledge is still quite challenging for existing VLN agents. Furthermore, previous methods usually lack enough interpretability in action decisions. To address the above issues, a few recent works introduce LLMs with great knowledge storage as the navigation backbone in a zero-shot manner [14], [15]. Nevertheless, the severe domain gap and the dependency on high-cost LLMs significantly harm navigation performance and scalability.

This paper proposes a new LLM-based VLN approach called NavCoT, where we conduct parameter-efficient in-domain training for teaching LLMs to perform self-guided navigational reasoning for action decisions. Through NavCoT, the navigation performance and interpretability of the action prediction can be greatly enhanced cost-effectively.

B.2 LLMs for Embodied AI

Introducing large language models (LLMs) into Embodied AI tasks has gained widespread interest recently. Benefiting from training on the ultra-large-scale corpus, LLMs exhibit the brilliant ability of planning, reasoning, and reflection to assist embodied task completion [9], [10], [36], [37], [38], [39]. SayCan [9] combines LLMs with affordance functions to produce feasible plans for completing household tasks. Inner Monologue [10] makes further improvements on [9] by injecting feedbacks from the environment. Although these approaches enable LLMs to interact with specific environments in different tasks, their offline use of LLMs inevitably brings noise. Recent methods have employed in-domain training to better adapt LLMs to embodied tasks [40], [41], [42]. For example, EmbodiedGPT [40] crafts a large-scale

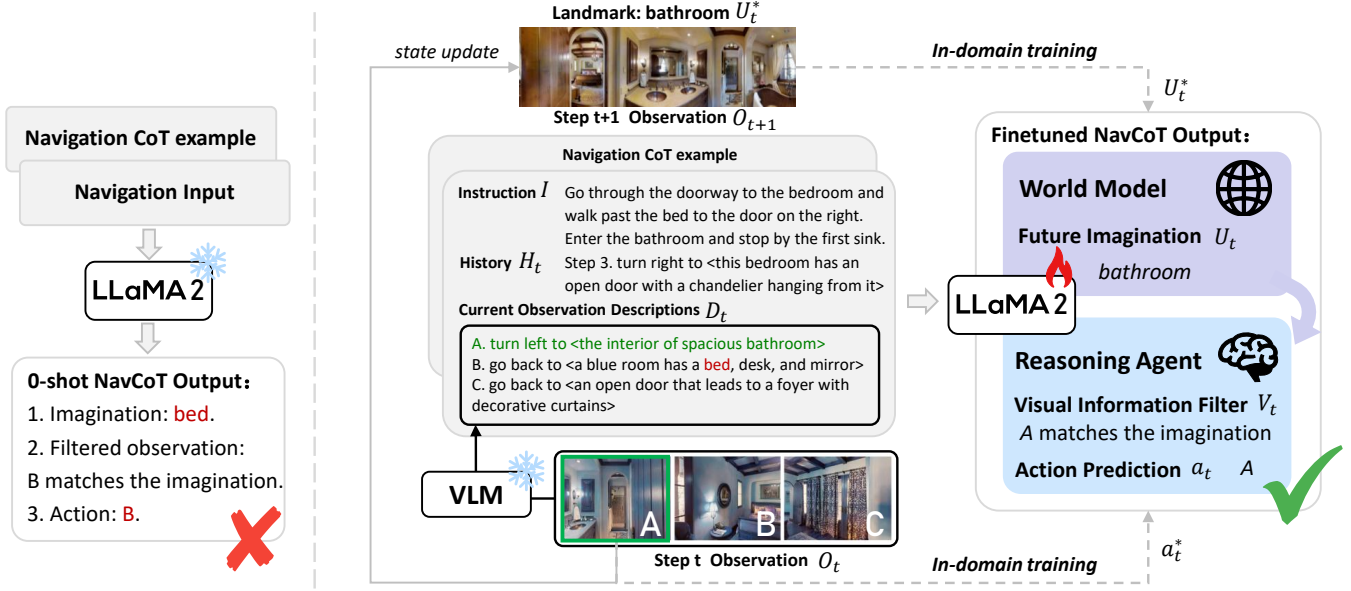


Fig. 2: Overview of NavCoT. At timestep t , we employ a VLM to translate the observation information into textual description. Then, the LLM is prompted with the example and the textual represented navigation input to produce the navigational chain-of-thought. We conduct in-domain training to enable the LLM to learn to generate reasonable navigational reasoning for action decisions.

embodied planning dataset and adapts the LLM to it by introducing an additional embodied-former. In our work, we adopt a parameter-efficient training scheme and utilize the available VLN training data to construct the navigational chain-of-thought labels, which effectively adapts the LLM to VLN tasks at a much lower cost.

B.3 Chain-of-Thought Prompting

Chain-of-Thought (CoT) prompting, firstly proposed in [19], is a powerful in-context learning technique to elicit multi-step reasoning abilities of LLMs. By elaborating intermediate reasoning steps to form the CoT rather than generating the answer only in the prompt, LLMs can learn to generate the output accordingly for the specific task and therefore improve the reasoning accuracy. Following [19], different works improve standard CoTs through self-consistency [43], least-to-most prompting [44], bootstrapping [45], tree-of-thought prompting [46], [47], etc. However, most of them prompt LLMs to produce CoTs in an offline and unconstrained manner.

In this work, we introduce the theory of the world model into the CoT mechanism in a trainable way and constrain the LLM to produce CoT outputs with the unified format by collecting formalized ground-truths. As a result, the LLM can learn to produce self-guided navigational reasoning, and the training process can be greatly simplified.

C PRELIMINARIES

C.1 Problem Setup

VLN requires an agent to follow a language instruction I to navigate from a start viewpoint to the target viewpoint. At timestep t , the agent receives a panoramic observation O_t containing K single-view observations $O_{t,k}$, i.e., $O_t = \{O_{t,k}\}_{k=1}^K$. There are N navigable views among K

views. The navigable views and a stop token $[stop]$ form the action space, from which the agent chooses one as the action prediction a_t . Actions before step t are viewed as the navigation history, which is denoted as $H_t = \{a_0, \dots, a_{t-1}\}$. A navigation trajectory is successful when the agent stops within 3m of the target viewpoint.

C.2 Large Language Models (LLMs)

Applying LLMs to non-linguistic embodied tasks has received more and more attention recently. We roughly divide these methods into two categories: the first one is to employ closed-source LLMs such as GPT-4 [16] in an offline way [9], [14], [37], which may suffer from poor scalability and severe domain gap. The second one is to introduce smaller open-source LLMs which can be deployed and trained locally [40]. Our method lies in the latter, and we adopt two open-source LLMs, LLaMA-Adapter [20] and LLaMA 2 [8], as the navigation backbones.

LLaMA-Adapter [20] is a lightweight adaption method that finetunes LLaMA 1 [7] efficiently with less time and parameters. The core idea is to introduce learnable adaption prompts and a zero-initialized attention mechanism with zero gating. LLaMA-Adapter can generate comparable responses to Alpaca 7B [48] trained with fully fine-tuned parameters. LLaMA 2 [8] is an updated version of LLaMA 1 [7]. It is trained on 2 trillion tokens and with twice the context length of the LLaMA 1. LLaMA 2 contains variants with different parameters scales, such as 7B, 13B, and 70B. We adopt a bias tuning strategy [49] for implementing parameter-efficient fine-tuning on LLaMA 2 7B.

D METHOD

The overview of NavCoT is presented in Fig. 2. At each timestep t , a vision-to-text system is first employed to

convert the surrounding observation into linguistic representations (Sec. D.1). Then, the LLM is prompted with the in-context example and textual represented navigation input to produce the navigational chain-of-thought (Sec. D.2). We adopt various training settings such as imitation-learning based finetuning and task-decomposed pretraining by collecting customized labels (Sec. D.3) to enable the LLM to generate reasonable and constrained outputs for improving the accuracy of action decisions (Sec. D.4).

D.1 Vision-to-Text System

The observation $O_{t,n}$ at each timestep t contains an RGB image $B_{t,n}$ and the direction information $A_{t,n} = \{\psi_{t,n}, \theta_{t,n}\}$, where $\psi_{t,n}$ and $\theta_{t,n}$ represent the heading and elevation, respectively. Therefore, the vision-to-text system translates both the vision information in $B_{t,n}$ and the direction information $A_{t,n}$ into textual descriptions and feeds them into the LLM for action decisions. We use an image captioning model BLIP [13], denoted as F_v , to translate the vision information in $B_{t,n}$ to a caption $D_{t,n}^v$:

$$D_{t,n}^v = F_v(B_{t,n}). \quad (1)$$

We map the direction information $A_{t,n}$ into the textual represented direction space containing six basic directions such as “turn left” and “go up”, following the direction mapping rules in VLN [1]. Denote the mapped direction information as $D_{t,n}^a$, the final textual description $D_{t,n}$ for each observation $O_{t,n}$ is obtained by:

$$D_{t,n} = \text{cat}(D_{t,n}^a, D_{t,n}^v), \quad (2)$$

where $\text{cat}(\cdot)$ denotes the string concatenation. For convenience, we add the alphabetical represented label for each observation to convert it into an action option, as in Fig. 2.

D.2 Navigational Chain-of-Thought Prompt

Since different tasks require distinguishing reasoning ability, proper design of the intermediate reasoning steps is crucial in designing the CoT prompt, which may greatly affect the performance of LLM’s predictions. In this work, we aim to empower the LLM to generate two significant intermediate reasoning steps for guiding the navigation action predictions, inspired by the world model theory [17], [18].

Concretely, the first reasoning step Future Imagination (FI) is to imagine the next observation according to the instruction and the navigation history. With FI, the LLM can monitor the progress of the navigation to guide itself for the subsequent action prediction. The second reasoning step, Visual Information Filter (VIF), aims to select the candidate observation that best aligns with the imagination generated in FI. Then the LLM generates the action prediction in the final Action Prediction (AP) reasoning step.

As shown in Fig. 2, at each timestep t , the LLM receives the prompt consisting of a chain-of-thought reasoning example and a query navigation input. The reasoning example serves as a reference to guide the LLM to generate the desired format of reasoning based on the given navigation input. The navigation input at timestep t consists of the instruction I , the textual described observation D_t obtained in Sec. D.1, i.e., $D_t = \{D_{t,n}\}_{n=1}^N$ (N is the number of navigable

views), and the navigation history H_t . With the navigation input, we prompt the LLM to generate the constrained reasoning format for the FI, VIF and AP steps.

Future Imagination (FI). In FI, we want LLM to generate the imagination about the next observation, which can be an object or a scene. Denote the imagination generated by LLM as U_t , the desired output format for FI is:

Imagination: U_t .

Visual Information Filter (VIF). After generating the imagination in FI, we introduce a further reasoning step, VIF, to force LLM to explicitly select the observation that best matches the imagination from the redundant observation information. Through this explicit visual information filter procedure, the LLM can better learn to connect the imagination and action decision. Denote the option of observation that LLM predicts to align the imagination U_t best as V_t . The desired output format for VIF is:

Filtered observation: V_t matches the imagination.

Action Prediction (AP). By summarizing the reasoning in FI and VIF, the LLM can make the final action prediction. Denote the option of action that LLM predicts as a_t , we define the output format for AP as:

Action: a_t .

We give an in-context example to the LLM as follows:

Input: Instruction: Walk towards the mirror and walk through the open door.
 Observation: [A. stop, B. go forward to <a bedroom with a bed>, C. turn right to <an open door leading to a hallway>].
 History: Step 1. go forward to <a wall with a mirror>.
 Output: Imagination: open door. Filtered observation: C matches the imagination.
 Action: C.

Through this example, the LLM can know the desired reasoning format and principle, e.g., the navigation history indicates “mirror” and the resulting imagination is “open door”. Based on the given example, we ask the LLM to generate the desired reasoning by giving the following prompt:

Input: Instruction: $\{I\}$ Observation: $\{D_t\}$
 History: $\{H_t\}$
 Output:

D.3 Ground-Truth Collection

As shown in Fig. 3, due to the uncertainty of the LLM’s output and the complexity of the VLN task, it is hard for the LLM to generate the multi-step reasoning accurately for action decisions in a zero-shot manner. To address this issue, we collect the ground-truth of the navigational chain-of-thought based on the available VLN data for implementing in-domain training to improve the action decision.

We firstly collect the the ground-truth imagination U_t^* for the reasoning task FI. Ideally, the ground-truth imagination should be consistent with the object/scene appearing in the following observation, which corresponds to the ground-truth action. Moreover, to enable better cross-modal alignment between the observation and the instruction for action decisions, the imagination is desired to be one of the mentioned objects/scenes in the given instruction. To this end,

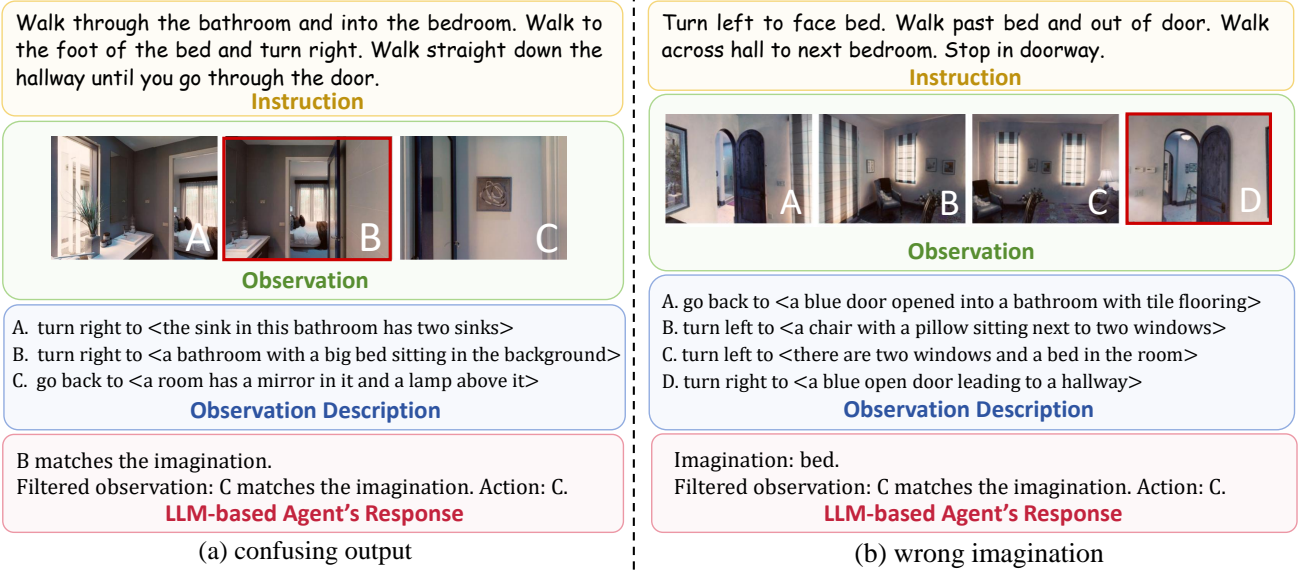


Fig. 3: Failure cases of LLM output in the zero-shot manner. The ground-truth actions are denoted by red boxes.

we resort to an LLM to extract the mentioned objects/scenes in the instruction and a cross-modal large model CLIP [50] to collect the ground-truth imagination U_t^* at different timesteps. Concretely, for a given instruction I , we provide the prompt with the customized task example to ask the LLM to extract the mentioned objects/scenes from I . Denote the extracted landmark list from I as $U^{la} = \{U_k^{la}\}_{k=1}^M$, where M is the number of objects/scenes mentioned in the instruction I . For the ground-truth observation B_t^* at each timestep t , we calculate the similarity between B_t^* and each landmark U_k^{la} in the list U^{la} , and take the one with the highest similarity as the ground-truth imagination label U_t^* :

$$U_t^* = \operatorname{argmax}_{U_k^{la}} \operatorname{Sim}(F_{\text{CLIP}}^t(U_k^{la}), F_{\text{CLIP}}^v(B_t^*)), \quad (3)$$

where F_{CLIP}^t and F_{CLIP}^v represent the text encoder and the image encoder of CLIP, respectively.

Since the reasoning task VIF aims to find aligned observation with the instruction for action decision, we set the ground-truth label of filtered observation to be consistent with the option of ground-truth action a_t^* . As a result, the ground truth of the navigational chain-of-thought, denoted as CoT_t^* , for instruction I at timestep t is defined by:

Imagination: U_t^* . Filtered observation: a_t^* matches the imagination. Action: a_t^* .

D.4 In-domain Training

We conduct two kinds of in-domain training schemes in NavCoT, i.e., pretraining and finetuning, which are illustrated in the following.

Pretraining. In NavCoT, we set each of the three navigational reasoning tasks defined in Sec. D.2 as a pretraining task and create corresponding instruction-following dataset. The pretraining objective \mathcal{L}_p is defined as follows:

$$\mathcal{L}_{\text{FI}} = -U^* \log(p_{\text{LLM}}(U|I, H, D)), \quad (4)$$

$$\mathcal{L}_{\text{VIF}} = -V^* \log(p_{\text{LLM}}(V|I, H, D)), \quad (5)$$

$$\mathcal{L}_{\text{AP}} = -a^* \log(p_{\text{LLM}}(a|I, H, D)), \quad (6)$$

$$\mathcal{L}_p = \mathcal{L}_{\text{FI}} + \mathcal{L}_{\text{VIF}} + \mathcal{L}_{\text{AP}}, \quad (7)$$

where D denotes the textual observations of a single navigation step and H is the history before the step. U , V , and a are the output of FI, VIF, AP, respectively. U^* , V^* , and a^* represent the ground-truth output extracted from CoT_t^* for FI, VIF, and AP, respectively.

Finetuning. Besides pretraining, we also conduct finetuning with the imitation learning strategy [1] to further adapt LLM to generate sequential action decisions for navigation. To this end, we reformulate the expert trajectory from original VLN dataset using CoT_t^* to construct instruction-following dataset. At each timestep t , we train the LLM to generate the complete navigational chain-of-thought CoT_t for action decisions. The finetuning objective \mathcal{L}_f is:

$$\mathcal{L}_f = - \sum_t \text{CoT}_t^* \log(p_{\text{LLM}}(\text{CoT}_t|I, H_t, D_t)), \quad (8)$$

After in-domain training, we prompt the trained LLMs to generate the navigational chain-of-thought for action decision based on the prompt and the in-context example described in Sec. D.2.

E EXPERIMENTS

E.1 Experimental Setup

Datasets. We evaluate NavCoT on four public VLN benchmarks: R2R [1], RxR [5], REVERIE [2], and R4R [4]. R2R is built on 90 real-world indoor environments containing 7189 trajectories, each corresponding to three fine-grained instructions. RxR contains much more complex instructions and trajectories than R2R. Since CLIP [51] is pretrained on English language data, we use the English subset of RxR (both en-IN and en-US) for verification, which includes 26464, 2939, 4551 path-instruction pairs for Training, Val Seen, and Val Unseen, respectively. REVERIE replaces the fine-grained instructions in R2R with high-level instructions. R4R concatenates two adjacent tail-to-head trajectories in R2R, forming longer instructions and trajectories.

Setting	Method	Val Seen					Val Unseen				
		TL	NE ↓	OSR ↑	SR ↑	SPL ↑	TL	NE ↓	OSR ↑	SR ↑	SPL ↑
cross-modal backbone	Seq2Seq [1]	11.33	6.01	53	39	-	8.39	7.81	28	21	-
	Speaker Follower [22]	-	3.36	74	66	-	-	6.62	45	36	-
	HAMT [30]	11.15	2.52	-	76	72	11.46	2.29	-	66	61
	DUET [31]	12.32	2.28	86	79	73	13.94	3.31	81	72	60
	BEVBert [34]	13.56	2.17	88	81	74	14.55	2.81	84	75	64
	ScaleVLN [35]	13.24	2.12	87	81	75	14.09	2.09	88	81	70
language only backbone	NavGPT [14]	-	-	-	-	-	11.45	6.46	42	34	29
	NavCoT+LLaMA-Adapter (ours)	-	-	-	-	-	9.19	8.20	28.91	21.97	19.99
	NavCoT+LLaMA 2 (ours)	-	-	-	-	-	9.83	6.67	43.98	36.40	33.17
	NavCoT+LLaMA 2 (ours)*	10.08	6.46	48.38	41.33	38.43	9.95	6.26	48.11	40.23	36.64

TABLE 1: Comparison with SOTAs on R2R. * denotes adding randomly chosen 12000 samples from the R2R augmentation dataset [22]. The best results for cross-modal and language-only backbones are denoted by bold and blue fonts, respectively.

Evaluation Metrics. The following standard metrics are used for evaluation on R2R [1] and REVERIE [2]: 1) Trajectory Length (TL): the average length of the agent’s navigated path, 2) Navigation Error (NE): the average distance between the agent’s destination and the target viewpoint, 3) Success Rate (SR): the ratio of success, where the agent stops within three meters of the target point, 4) Success rate weighted by Path Length (SPL) [1]: success rate normalized by the ratio between the length of the shortest path and the predicted path, 5) Oracle Success Rate (OSR): the ratio of containing a viewpoint along the path where the target position is visible. Three evaluation metrics related to the instruction following are added for R4R [4] and RxR [5], i.e., the Coverage weighted by Length Score (CLS) [4], the normalized Dynamic Time Warping (nDTW) [52], and the Success weighted by nDTW (SDTW) [52].

Implementation Details. We train LLaMA-Adapter [20] and LLaMA 2 [8] of size 7B with 1.2M and 1.6M trainable parameters, respectively. To speed up the training, we use 4 V100 GPUs with a batch size of 8. The total training time lasts ~ 10 h on 4 V100 GPUs. The inference is conducted on a single V100 GPU. We use the AdamW optimizer with the learning rate of 0.001 and the weight decay of 0.02. For fast evaluation on various ablation experiments, we randomly choose 90 instruction-trajectory pairs from 8 scans in the val unseen split as the Val Unseen Subset. This subset serves as an efficient testbed to indicate the performance gap among different methods.

For collecting the the ground-truth imagination U_t^* , we utilize a powerful LLM tool, ChatGPT [53], to extract the landmarks from a given instruction. An in-context example for the landmark extraction is as follows:

Instruction: Walk along the rug past the statue on the wooden table.

Landmarks:

- 1.rug;
- 2.statue;
- 3.wooden table.

The maximum length of tokens in training and inference are set to 400 and 512, respectively. During inference, the temperature of LLM is set as 0. For parameter-efficient in-domain training of LLaMA 2 [8], we initialize the model weights using the pretrained LLaMA 2 weights with bias tuning¹. The training epoch for NavCoT based on both LLaMA-Adapter [20] and LLaMA 2 [8] is set as 2.

1. https://llama2-accessory.readthedocs.io/en/latest/finetune/sg_peft.html#bias-norm-tuning-of-llama2-7b-on-alpaca

E.2 Quantitative Results

E.2.1 Comparison with Existing Approaches

Table 1 presents the comparison between NavCoT and methods with different backbones on R2R [1]. From Table 1, we can see that NavCoT with LLaMA 2 significantly outperforms a recent GPT4-based approach NavGPT [14] (e.g., 4.17% performance improvement on SPL), demonstrating that our method can effectively boost the performance of smaller LLMs (i.e., LLaMA 2) with a low training cost to outperform high-cost LLMs. We can also find that NavCoT with LLaMA-Adapter as backbone still shows a large performance gap with NavGPT, revealing the significant challenge of adapting LLMs with relatively small model capacity to the VLN task. Moreover, we can find that by adding a small amount of augmentation data (similar scale to that of R2R training data), the performance can be further enhanced largely, demonstrating the potential of improving our NavCoT via simple data augmentation.

Discussion. From the comparison results we can find that LLM-based agents like NavGPT [14] (using GPT-4) encounter a significant performance drop in VLN, demonstrating that original ability of LLM is far from solving such complex tasks. Our NavCoT proposes a promising and low-cost way to adapt smaller LLMs to VLN, which largely outperforms NavGPT. Note that although the language only backbone still shows a significant discrepancy with cross-modal backbone on different datasets, it provides much better interpretability through performing explicit reasoning than the latter, which is crucial in developing realistic robotic applications. In the future, it is promising to introduce our NavCoT into large vision-language models (VLMs) to further boost the navigation performance while enhancing interpretability.

E.2.2 Ablation Study

We conduct extensive ablation experiments on R2R [1] to analyze the effect of each component in NavCoT, including training settings, reasoning tasks, navigation histories, *etc.* **Training settings.** Fig. 4 presents a comprehensive comparison among various training settings. From Fig. 4, we can draw three crucial conclusions: 1) NavCoT surpasses the DAP variant in all training settings, highlighting the power of explicit disentangled reasoning; 2) The great superiority of training-based settings over zero-shot ones shows the effectiveness of our parameter-efficient in-domain training strategy; 3) In-domain pretraining and finetuning can encourage the LLM to select actions that are relevant to the

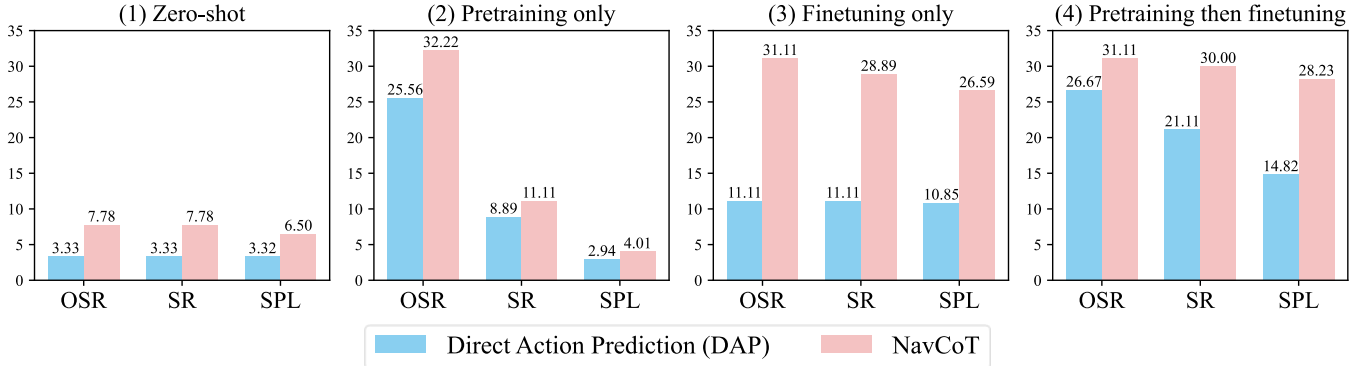


Fig. 4: Comparison of NavCoT with the Direct Action Prediction (DAP) variant under different training settings. In DAP, we directly prompt LLM to generate the action prediction.

Method	pretraining only			pretraining & finetuning		
	OSR	SR	SPL	OSR	SR	SPL
AP	25.56	8.89	2.94	26.67	21.11	14.82
VIF	28.89	10.00	3.43	30.00	21.11	19.08
AP+FI	23.22	8.89	4.20	27.78	24.44	23.78
AP+VIF	33.33	7.78	2.15	24.44	23.33	22.33
AP+FI+VIF (ours)	<u>32.22</u>	11.11	<u>4.01</u>	31.11	30.00	28.23

TABLE 2: Ablation results for different proxy tasks on R2R Val Unseen subset.

Setting	Val Unseen				
	TL	NE ↓	OSR ↑	SR ↑	SPL ↑
Example 1	10.34	6.91	44.10	34.53	31.00
Example 2	10.09	6.76	44.40	37.04	33.65
Example 3	10.34	6.74	44.83	36.14	32.50
Original Example	9.83	6.67	43.98	36.40	33.17

TABLE 3: Ablation results for example sensitivity. The backbone is LLaMA 2 [8].

target position, reflected in significant improvement of OSR. However, the finetuning contributes much higher to the enhancement of SR and SPL than pretraining by benefiting from the learning of sequential action decisions.

Reasoning tasks. Table 2 shows how our three reasoning tasks, i.e., Future Imagination (FI), Visual Information Filter (VIF), and Action Prediction (AP) impact on the navigation performance under both pretraining only and pretraining & finetuning settings. From Table 2 we can find that the complete combination of three reasoning tasks, i.e., AP+FI+VIF (ours) achieves the comprehensively best performance in both two training settings, showing the effectiveness of our proposed navigational chain of thoughts in improving the action decisions. The results of AP *vs.* AP+FI and AP *vs.* AP+VIF demonstrate the effects of the reasoning tasks FI and VIF, respectively. Moreover, through the results of AP+FI *vs.* AP+FI+VIF (ours) and AP+VIF *vs.* AP+FI+VIF (ours), we can also find that both FI and VIF are indispensable in the navigational chain of thoughts.

Example sensitivity. We test NavCoT on different in-context examples for in-domain training to examine the example sensitivity of NavCoT. Concretely, we define three more examples, using different manually written instructions and observation descriptions. From the results in Table 3, we can see that NavCoT encounters a small range of performance changes when given different in-context examples, showing

Setting	Method	Val Unseen Subset		
		NE ↓	SR ↑	SPL ↑
Finetuning	No	9.67	13.33	5.25
	All	7.45	21.11	14.55
	Last	7.16	28.89	26.59
Pretraining then Finetuning	No	7.92	22.22	8.21
	All	7.13	16.67	11.26
	Last	6.74	30.00	28.23

TABLE 4: Ablation results for navigation history. The backbone is LLaMA-Adapter [20].

Method	Val Unseen Subset			
	NE ↓	OSR ↑	SR ↑	SPL ↑
▷ Vision-to-text system:				
Coarse Direction & Object	7.23	33.33	26.67	24.20
Fine Direction & Caption	8.78	41.11	17.78	12.26
Coarse Direction & Caption (ours)	5.38	58.89	53.33	48.69
▷ Chain-of-Thought prompts:				
original CoT (GPT-4 [16])*	7.65	47.78	24.44	12.33
original CoT (LLaMA2)*	9.39	3.33	3.33	1.85
NavCoT (LLaMA2)*	9.02	26.67	16.67	5.79
NavCoT (ours)	5.38	58.89	53.33	48.69
▷ Imagination formats:				
Free Object	8.35	7.78	7.78	7.32
Direction & Object	5.99	53.33	44.44	41.12
Object (ours)	5.38	58.89	53.33	48.69
▷ Backbones:				
GPT-4*	7.80	22.22	15.56	8.70
LLaMA-Adapter*	9.64	7.78	7.78	6.50
LLaMA-Adapter	7.16	31.11	28.89	26.59
LLaMA 2*	9.02	26.67	16.67	5.79
LLaMA 2 (ours)	5.38	58.89	53.33	48.69
▷ In-context examples :				
No-example training	7.97	54.44	20.00	13.62
NavCoT (ours)	5.38	58.89	53.33	48.69

TABLE 5: Ablation results for vision-to-text systems, chain-of-thought prompts, imagination formats, backbones, and in-context examples.* denotes the zero-shot manner.

low example sensitivity of NavCoT.

Navigation history. We also conduct experiments to verify how navigation history influences navigational reasoning, as shown in Table 4. In our implementation, “No” means removing history information in the input. “All” and “Last” represent preserving all previous steps and only the last step as history, respectively. From Table 4, we can find that “No” performs the worst among all settings not surprisingly, since the navigation history is essential for monitoring

the progress of navigation. Moreover, the “Last” setting performs best among all settings. A possible reason is that the “All” setting may bring much more noise and also exceed the maximum sequence length, therefore leading to performance degradation.

Other ablations. Table 5 shows the ablation results for vision-to-text systems, chain-of-thought prompts, imagination formats, backbones, and in-context examples. In Table 5: 1) Coarse Direction means the textual direction description does not contain the concrete angle information (see Sec. D.1), and Fine Direction means adding the angle information into the textual direction description; 2) For Coarse Direction & Object setting, we use the GroundingDino [54] model combined with the RAM [55] model to extract the object information in the observation; 3) For original CoT, we directly prompt the LLM to generate thought and action like [14]; 4) For Free Object setting, we extract the object of ground-truth action (observation) rather than using the object mentioned in the instruction as the imagination label.

We can draw the following conclusions from Table 5: 1) **vision-to-text system:** By comparing Coarse Direction & Object and Coarse Direction & Caption (ours), we can find that captions are more suitable for our NavCoT, probably because captions contain spatial information for facilitating the alignment to the instruction, while pure object detections may have more noise. Through the comparison between Fine Direction and Coarse Direction, we can also find that fine-grained direction may bring unnecessary noise for confusing the navigation decision. This is reasonable since instructions generally contain coarse direction descriptions such as “go up” and “turn left”. 2) **Chain-of-Thought prompts:** We can see that NavCoT outperforms the original CoT, suggesting that our design of combining the world model into CoT can better activate reasonable navigational reasoning. Benefiting from training with formalized labels, NavCoT outperforms both the zero-shot variant and the original CoT for GPT-4 [16] largely (e.g., ~ 43 and ~ 36 points in SPL, respectively). Moreover, note that our NavCoT also shows great superiority than the original CoT for GPT-4 [16] in the computational efficiency, i.e., the inference time of one-step decision for GPT-4 and NavCoT are ~ 9.8 s and ~ 0.5 s, respectively. 3) **Imagination formats:** the great superiority of Object (ours) over Free Object reveals the reasonability and the effectiveness of our NavCoT by constraining the imagination to be mentioned objects in the instruction. This strategy essentially fulfills a more task-oriented world model for better instruction following and modality alignment. Moreover, asking the LLM to imagine both the direction and object brings performance degradation, probably because the imagined direction information is usually inaccurate and harms the alignment with the instruction, and therefore confuses the LLM for action decisions. 4) **Backbones:** we can find that NavCoT+LLaMA-Adapter is largely inferior to its alternative NavCoT+LLaMA2. This is not surprising since LLaMA-Adapter is based on LLaMA 1 [7] while LLaMA 2 is trained on 40% more data than the previous version. Moreover, LLaMA-Adapter has 0.4M fewer trainable parameters than the bias tuning setting [49]. However, NavCoT shows consistent improvement on both backbones under our in-domain training strategies. Interestingly, we can find that LLaMA 2 outperforms GPT-4

Method	Val Unseen				
	SR	SPL	CLS	nDTW	SDTW
EnvDrop [21]	38.5	34	54	51	32
HAMT* [30]	38.26	36.23	58.45	53.08	32.81
Direct Action Prediction (DAP)	20.46	18.68	37.19	33.64	16.02
NavCoT (ours)	24.52	22.58	45.06	38.94	19.63

TABLE 6: Comparison results on RxR English subset. * denotes our re-implementation results with the imitation learning strategy [30]. The best results for cross-modal and language-only backbones are denoted by bold and blue fonts, respectively.

Method	Val Unseen			
	TL	SR	OSR	SPL
Seq2Seq [2]	-	4.2	9.07	2.84
RCM [25]	11.98	9.29	14.23	6.97
FAST-MATTN [2]	45.28	14.40	28.20	7.19
HAMT* [30]	9.24	23.80	26.44	22.42
Direct Action Prediction (DAP)	16.30	3.12	7.36	1.74
NavCoT (ours)	12.36	9.20	14.20	7.18

TABLE 7: Comparison results on REVERIE. * denotes our re-implementation results with the imitation learning strategy [30]. The best results for cross-modal and language-only backbones are denoted by bold and blue fonts, respectively.

in the zero-shot NavCoT setting, probably because when facing the prompt with less information and constrained output format, the advantage of GPT-4 is not obvious over the smaller language model LLaMA 2. However, such kind of prompt is more beneficial for implementing in-domain training. 5) **In-context examples:** the comparison between No-example training and NavCoT shows that adding the example is still necessary under the training setting and the in-context example can effectively facilitate the training.

E.2.3 Generalization on Other Datasets

To further verify the generalization of NavCoT on different VLN datasets with much longer instructions (RxR [5]) and high-level instructions (REVERIE [2]), we conduct experiments on RxR and REVERIE, and the results are given in Table 6 and Table 7, respectively. Since there are no reported results of language-only backbone on RxR and REVERIE, we develop a naive baseline, which is a variant of NavCoT that makes direct action predictions (DAP) rather than generating navigational chain-of-thoughts. We also present the results of a strong pretrained cross-modal backbone HAMT [30] with the same training strategy as ours. From Table 6 we can find that although NavCoT shows the performance gap with the cross-modal methods, it significantly outperforms DAP especially in the CLS and nDTW metrics, showing that our method not only improves the navigation accuracy but also enables better instruction following, which is crucial for long-horizon navigation. Table 7 shows that different methods encounter significant performance drop on REVERIE compared to R2R, revealing the challenge of navigating under high-level instructions with limited information. However, NavCoT is still superior over DAP in a large margin, showing the effectiveness of the proposed method.

Method	R4R				R2R			
	OSR \uparrow	SR \uparrow	CLS \uparrow	NDTW \uparrow	NE \downarrow	OSR \uparrow	SR \uparrow	SPL \uparrow
HAMT* [30]	53.80	22.20	54.29	39.49	4.74	63.47	57.43	54.8
Direct Action Prediction (DAP)	25.80	12.80	19.73	16.88	7.89	23.50	20.52	19.96
NavCoT (ours)	42.20	13.00	45.99	31.79	7.08	37.46	31.33	29.08

TABLE 8: Low-resource experimental results on Val Unseen on R4R and R2R. * denotes our re-implementation results with the imitation learning strategy [30]. In R4R, CLS and NDTW are the main metrics.

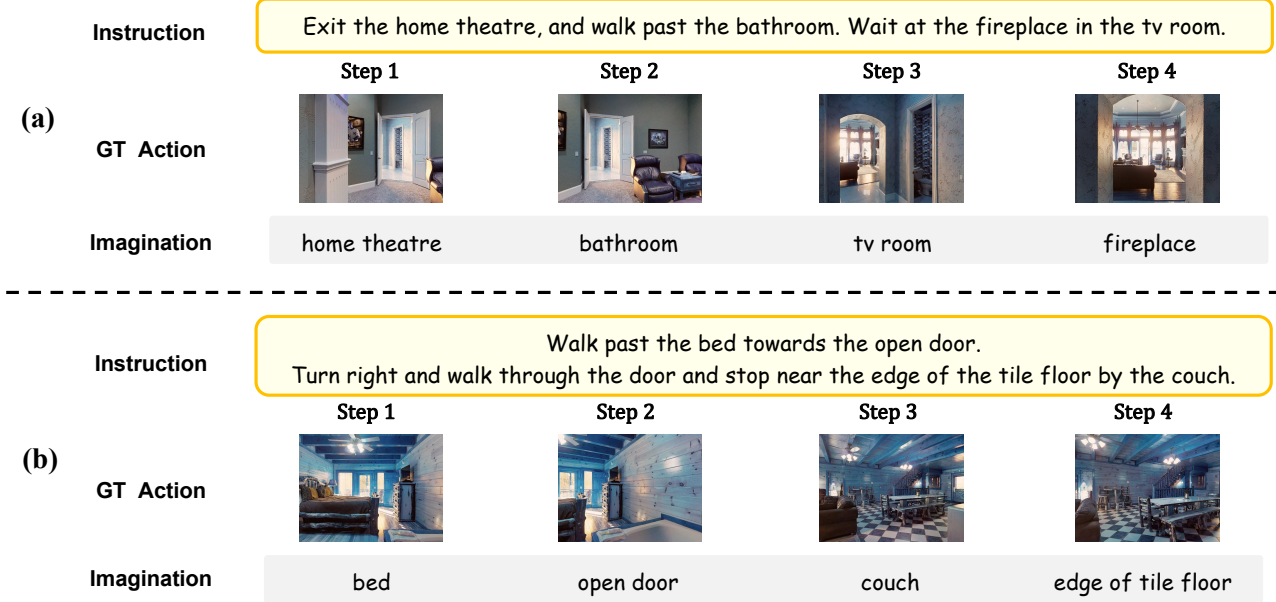


Fig. 5: Visualization examples of Imagination ground-truth (GT). We do not show the imagination GT for the final step which is “stop”.

E.2.4 Low-Resource Experiments

We further conduct the low-resource experiments on two popular VLN benchmarks R2R [1] and R4R [4] to validate the ability of NavCoT when facing a small amount of training data. R4R contains much longer instructions and trajectories than R2R. For R4R, we randomly choose 5000 trajectories from 50 scans in Train for training and 500 trajectories from 11 scans in Val Unseen for validation. For R2R, we randomly extract 3000 trajectories from 61 scans in Train for training. We compare NavCoT with Direct Action Prediction (DAP) and a strong cross-modal backbone method HAMT [30].

Table 8 shows that under the low-resource setting with the same training strategy, the performance gap between NavCoT and HAMT [30] is largely reduced compared to that in Table 1. For example, the performance gap of CLS between HAMT [30] and NavCoT is reduced to ~ 8 . Moreover, NavCoT outperforms DAP in a large margin on both R2R and R4R, showing the good generalization ability of NavCoT under low-resource settings. Note that the results on R4R show that NavCoT surpasses DAP significantly in OSR, CLS, and NDTW, demonstrating that beyond correct action decisions, NavCoT also shows promising instruction following ability for long trajectories.

E.3 Visualization

Quality of Imagination ground-truths (GTs). Fig. 5 gives some visualization examples of imagination GT. In Fig. 5, we can observe that the application of LLM and CLIP in

the ground-truth collection process effectively ensures the quality of the imagination GT even if when the landmark is rarely mentioned in the corpus (Step 4 in Fig. 5(b)) and the landmark only occupies a small region in the observation (Step 2 in Fig. 5(b) and Step 4 in Fig. 5(a)). We also conduct a quantitative human evaluation by randomly extracting 200 trajectories on R2R. The accuracy of CLIP reaches $\sim 73\%$, which is a relatively high value. Note that since the decisions are made based on comprehensive information, a small proportion of noisy imagination labels are tolerable.

Action decision visualizations. Fig. 6 gives the action decision visualization of NavCoT, where we can find that NavCoT generates reasonable navigational reasoning for guiding itself to make action decisions. For instance, in Step 2, from the history and observations, NavCoT infers that its position is possibly the *hall* and the option C may contain the desired *table*, then outputs the imagination that leads to the correct action. Notably, due to the domain gap between NavCoT and VLM, the observations inevitably do not always contain the imagination (Step 1, 3 and 4). However, with in-domain training, NavCoT shows the emerging ability to recognize potential connection between observations and imaginations. For example, in Step 4, the imagination *office* does not directly match any observation, whereas NavCoT correctly chooses the option C containing a *wooden desk* which suggests a scene of office.

Generalization to special cases. We present some visualization examples in Fig. 7 to show how NavCoT generalizes to typical special cases in VLN, e.g., there are little landmarks

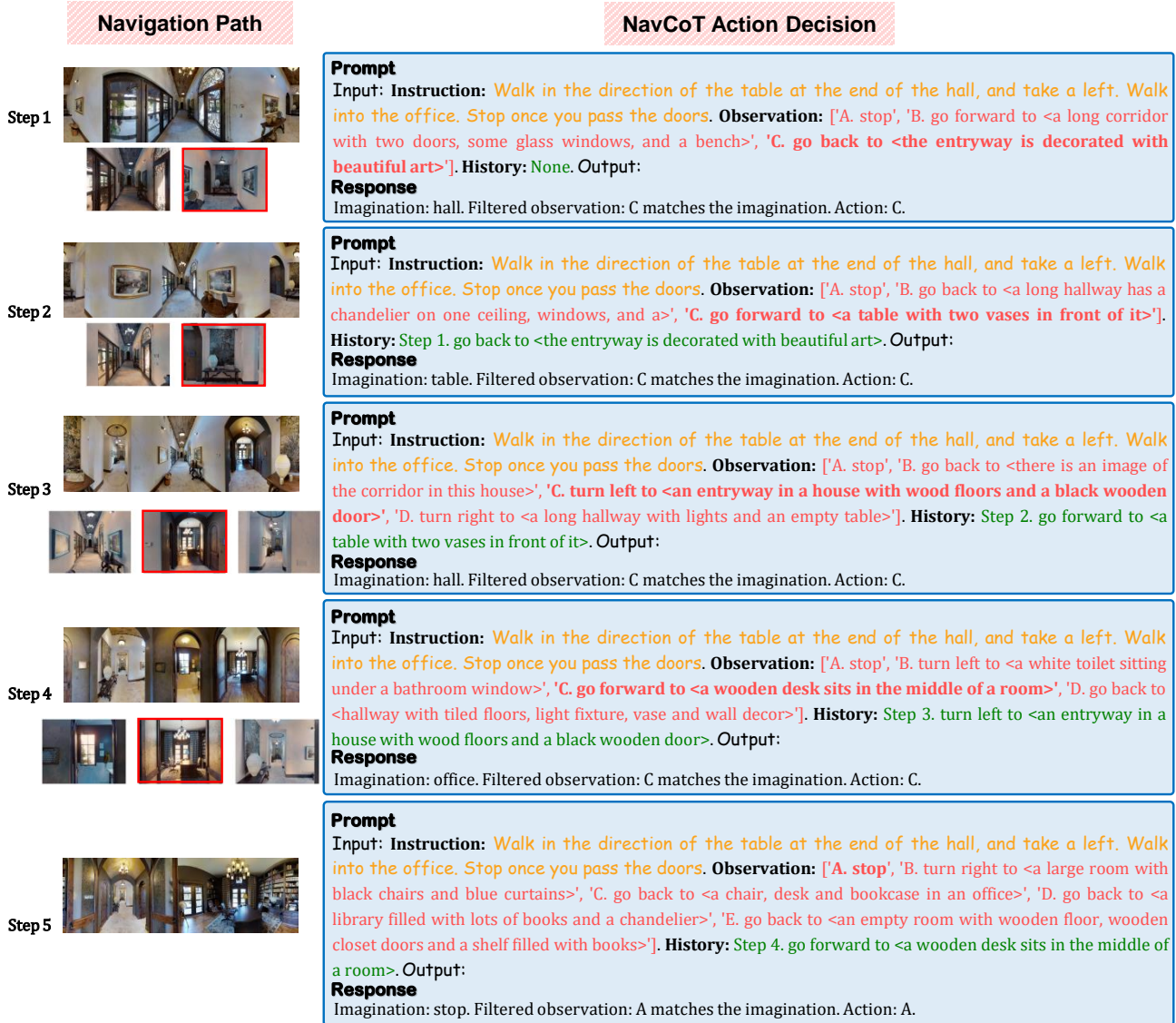


Fig. 6: Action decision visualization of NavCoT. The ground-truth actions are annotated by red boxes.

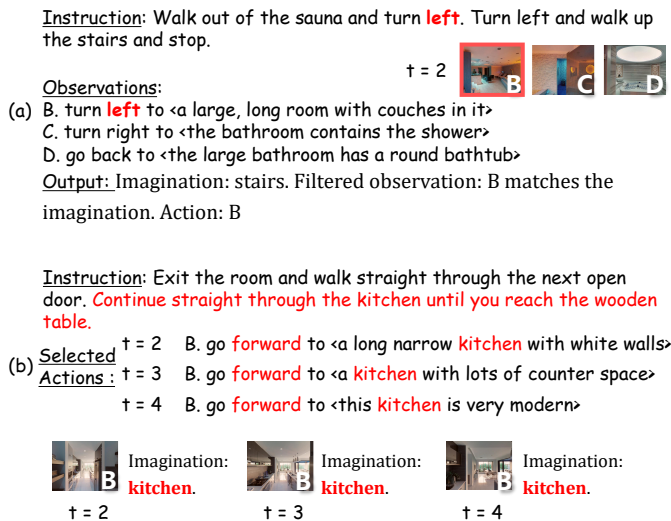


Fig. 7: Generalization to special cases: (a) Little landmarks: There are little landmarks mentioned in the instruction. (b) Same landmark: multiple continuous navigation steps point to the same landmark.

mentioned in the instruction or multiple continuous navigation steps pointing to the same landmark. From Fig. 7, we can find that NavCoT is capable of making correct navigation decisions under these cases. For example, in Fig. 7(a), even if the instruction contain little landmarks and limit the effect of the imagination, through the direction information provided by our vision-to-text system, NavCoT can still perform direction-level cross-modal alignment between the observation description and the instruction. Moreover, since the training data naturally contain many cases where continuous steps point to the same landmark, NavCoT can learn to generate reasonable imagination after training, as shown in Fig. 7(b).

Failure cases. We further analyze some failure cases of NavCoT in Fig. 8. We find two typical types of failures: 1) Indistinguishable observation descriptions. As shown in Fig. 8(a), since the observations in Options C and D in Step 2 are visually similar, their textual descriptions are also indistinguishable. Such cases also impose a big challenge on the cross-modal navigation agent. 2) Loss of global observation information. In Fig. 8(b), instruction mentions *the first door on the left*, which can not be easily inferred from



Fig. 8: Failure cases. The ground-truth (GT) actions and the actions chosen by NavCoT are denoted by green and red boxes, respectively.

independent candidate observation descriptions. The introduction of panoramic observation information may further help improve the global scene understanding.

F CONCLUSION

This work introduces NavCoT, which fulfills parameter-efficient in-domain training for enabling LLMs to perform self-guided navigational decisions. Experimental results show the great superiority of NavCoT over a recent high-cost LLM-based VLN approach and direct action prediction variants. We believe that our method makes a solid step towards developing scalable LLM-based VLN approaches and provides a meaningful reference in designing trainable navigational reasoning generation strategies for improving both the accuracy and interpretability of action decision. Constrained by the detail information loss during the vision-to-text transformation, the LLM may fail to make accurate decisions in some cases. Future direction includes introducing our NavCoT into powerful large vision-language models to further improve the navigation performance.

REFERENCES

[1] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sunderhauf, I. Reid, S. Gould, and A. van den Hengel, "Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments," in *CVPR*, 2018.

[2] Y. Qi, Q. Wu, P. Anderson, X. Wang, W. Y. Wang, C. Shen, and A. van den Hengel, "Reverie: Remote embodied visual referring expression in real indoor environments," in *CVPR*, 2020.

[3] H. Chen, A. Suhr, D. K. Misra, N. Snively, and Y. Artzi, "Touch-down: Natural language navigation and spatial reasoning in visual street environments," in *CVPR*, 2019.

[4] V. Jain, G. Magalhaes, A. Ku, A. Vaswani, E. Ie, and J. Baldridge, "Stay on the path: Instruction fidelity in vision-and-language navigation," in *ACL*, 2019.

[5] A. Ku, P. Anderson, R. Patel, E. Ie, and J. Baldridge, "Room-across-room: Multilingual vision-and-language navigation with dense spatiotemporal grounding," in *EMNLP*, 2020.

[6] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, and et al., "Language models are few-shot learners," in *NeurIPS*, 2020.

[7] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar et al., "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.

[8] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale et al., "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint arXiv:2307.09288*, 2023.

[9] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan et al., "Do as i can and not as i say: Grounding language in robotic affordances," *arXiv preprint arXiv:2204.01691*, 2022.

[10] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. R. Florence, A. Zeng, J. Tompson, I. Mordatch, Y. Chebotar, P. Sermanet, N. Brown, T. Jackson, L. Luu, S. Levine, K. Hausman, and B. Ichter, "Inner monologue: Embodied reasoning through planning with language models," *arXiv preprint arXiv:2207.05608*, 2022.

[11] D. Driess, F. Xia, M. S. M. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu et al., "Palm-

- e: An embodied multimodal language model," *arXiv preprint arXiv:2303.03378*, 2023.
- [12] J. Li, D. Li, S. Savarese, and S. Hoi, "BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models," in *ICML*, 2023.
- [13] J. Li, D. Li, C. Xiong, and S. Hoi, "Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation," in *ICML*, 2022.
- [14] G. Zhou, Y. Hong, and Q. Wu, "Navgpt: Explicit reasoning in vision-and-language navigation with large language models," in *AAAI*, 2024.
- [15] Y. Long, X. Li, W. Cai, and H. Dong, "Discuss before moving: Visual language navigation via multi-expert discussions," *arXiv preprint arXiv:2309.11382*, 2023.
- [16] O. OpenAI, "Gpt-4 technical report," Mar 2023.
- [17] P. N. Johnson-Laird, *Mental models: Towards a cognitive science of language, inference, and consciousness*. Harvard University Press, 1983, no. 6.
- [18] —, "Mental models and human reasoning," *Proceedings of the National Academy of Sciences*, vol. 107, no. 43, pp. 18 243–18 250, 2010.
- [19] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," in *NeurIPS*, 2022.
- [20] R. Zhang, J. Han, A. Zhou, X. Hu, S. Yan, P. Lu, H. Li, P. Gao, and Y. Qiao, "Llama-adapter: Efficient fine-tuning of language models with zero-init attention," *arXiv preprint arXiv:2303.16199*, 2023.
- [21] H. Tan, L. Yu, and M. Bansal, "Learning to navigate unseen environments: Back translation with environmental dropout," in *NAACL-HLT*, 2019.
- [22] D. Fried, R. Hu, V. Cirik, A. Rohrbach, J. Andreas, L.-P. Morency, T. Berg-Kirkpatrick, K. Saenko, D. Klein, and T. Darrell, "Speaker-follower models for vision-and-language navigation," in *NeurIPS*, 2018.
- [23] C. Liu, F. Zhu, X. Chang, X. Liang, Z. Ge, and Y.-D. Shen, "Vision-language navigation with random environmental mixup," in *ICCV*, 2021.
- [24] T.-J. Fu, X. E. Wang, M. F. Peterson, S. T. Grafton, M. P. Eckstein, and W. Y. Wang, "Counterfactual vision-and-language navigation via adversarial path sampling," in *ECCV*, 2020.
- [25] X. Wang, Q. Huang, A. Celikyilmaz, J. Gao, D. Shen, Y.-F. Wang, W. Y. Wang, and L. Zhang, "Reinforced cross-modal matching and self-supervised imitation learning for vision-language navigation," in *CVPR*, 2019.
- [26] C.-Y. Ma, Jiasen lu, Z. Wu, G. AlRegib, Z. Kira, richard socher, and C. Xiong, "Self-monitoring navigation agent via auxiliary progress estimation," in *ICLR*, 2019.
- [27] Z. Deng, K. Narasimhan, and O. Russakovsky, "Evolving graphical planner: Contextual global planning for vision-and-language navigation," in *NeurIPS*, 2020.
- [28] Y. Qi, Z. Pan, S. Zhang, A. van den Hengel, and Q. Wu, "Object-and-action aware model for visual language navigation," in *ECCV*, 2020.
- [29] Y. Hong, Q. Wu, Y. Qi, C. Rodriguez-Opazo, and S. Gould, "Vln bert: A recurrent vision-and-language bert for navigation," in *CVPR*, 2021.
- [30] S. Chen, P.-L. Guhur, C. Schmid, and I. Laptev, "History aware multimodal transformer for vision-and-language navigation," in *NeurIPS*, 2021.
- [31] S. Chen, P.-L. Guhur, M. Tapaswi, C. Schmid, and I. Laptev, "Think global, act local: Dual-scale graph transformer for vision-and-language navigation," in *CVPR*, 2022.
- [32] Y. Qiao, Y. Qi, Y. Hong, Z. Yu, P. Wang, and Q. Wu, "Hop: History-and-order aware pre-training for vision-and-language navigation," in *CVPR*, 2022.
- [33] P.-L. Guhur, M. Tapaswi, S. Chen, I. Laptev, and C. Schmid, "Airbert: In-domain pretraining for vision-and-language navigation," in *ICCV*, 2021.
- [34] D. An, Y. Qi, Y. Li, Y. Huang, L. Wang, T. Tan, and J. Shao, "Bevbert: Topo-metric map pre-training for language-guided navigation," in *ICCV*, 2023.
- [35] Z. Wang, J. Li, Y. Hong, Y. Wang, Q. Wu, M. Bansal, S. Gould, H. Tan, and Y. Qiao, "Scaling data generation in vision-and-language navigation," in *ICCV*, 2023.
- [36] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafraan, K. Narasimhan, and Y. Cao, "React: Synergizing reasoning and acting in language models," in *ICLR*, 2023.
- [37] D. Shah, B. Osinski, B. Ichter, and S. Levine, "Lm-nav: Robotic navigation with large pre-trained models of language, vision, and action," in *CoRL*, 2022.
- [38] R. Schumann, W. Zhu, W. Feng, T.-J. Fu, S. Riezler, and W. Y. Wang, "Velma: Verbalization embodiment of llm agents for vision and language navigation in street view," *arXiv preprint arXiv:2307.06082*, 2023.
- [39] G. Wang, Y. Xie, Y. Jiang, A. Mandlkar, C. Xiao, Y. Zhu, L. Fan, and A. Anandkumar, "Voyager: An open-ended embodied agent with large language models," *arXiv preprint arXiv:2305.16291*, 2023.
- [40] Y. Mu, Q. Zhang, M. Hu, W. Wang, M. Ding, J. Jin, B. Wang, J. Dai, Y. Qiao, and P. Luo, "Embodiedgpt: Vision-language pre-training via embodied chain of thought," *arXiv preprint arXiv:2305.15021*, 2023.
- [41] J. Yang, Y. Dong, S. Liu, B. Li, Z. Wang, C. Jiang, H. Tan, J. Kang, Y. Zhang, K. Zhou *et al.*, "Octopus: Embodied vision-language programmer from environmental feedback," *arXiv preprint arXiv:2310.08588*, 2023.
- [42] B. Zitkovich, T. Yu, S. Xu, P. Xu, T. Xiao, F. Xia, J. Wu, P. Wohlhart, S. Welker, A. Wahid *et al.*, "Rt-2: Vision-language-action models transfer web knowledge to robotic control," in *CoRL*, 2023.
- [43] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou, "Self-consistency improves chain of thought reasoning in language models," in *ICLR*, 2023.
- [44] D. Zhou, N. Schärli, L. Hou, J. Wei, N. Scales, X. Wang, D. Schuurmans, C. Cui, O. Bousquet, Q. Le *et al.*, "Least-to-most prompting enables complex reasoning in large language models," in *ICLR*, 2023.
- [45] E. Zelikman, J. Mu, N. D. Goodman, and Y. T. Wu, "Star: Self-taught reasoner bootstrapping reasoning with reasoning," in *NeurIPS*, 2022.
- [46] S. Yao, D. Yu, J. Zhao, I. Shafraan, T. L. Griffiths, Y. Cao, and K. Narasimhan, "Tree of thoughts: Deliberate problem solving with large language models," *arXiv preprint arXiv:2305.10601*, 2023.
- [47] J. Long, "Large language model guided tree-of-thought," *arXiv preprint arXiv:2305.08291*, 2023.
- [48] R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, and T. B. Hashimoto, "Stanford alpaca: An instruction-following llama model," https://github.com/tatsu-lab/stanford_alpaca, 2023.
- [49] P. Gao, J. Han, R. Zhang, Z. Lin, S. Geng, A. Zhou, W. Zhang, P. Lu, C. He, X. Yue *et al.*, "Llama-adapter v2: Parameter-efficient visual instruction model," *arXiv preprint arXiv:2304.15010*, 2023.
- [50] J. Yu, X. Jiang, Z. Qin, W. Zhang, Y. Hu, and Q. Wu, "Learning dual encoding model for adaptive visual understanding in visual dialogue," *IEEE TIP*, 2021.
- [51] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *ICML*, 2021.
- [52] G. Ilharco, V. Jain, A. Ku, E. Ie, and J. Baldridge, "General evaluation for instruction conditioned navigation using dynamic time warping," *ViGIL@NeurIPS*, 2019.
- [53] OpenAI, "Introducing chatgpt," <https://openai.com/blog/chatgpt>, 2022.
- [54] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu *et al.*, "Grounding dino: Marrying dino with grounded pre-training for open-set object detection," *arXiv preprint arXiv:2303.05499*, 2023.
- [55] Y. Zhang, X. Huang, J. Ma, Z. Li, Z. Luo, Y. Xie, Y. Qin, T. Luo, Y. Li, S. Liu *et al.*, "Recognize anything: A strong image tagging model," *arXiv preprint arXiv:2306.03514*, 2023.



Bingqian Lin received the B.E. and the M.E. degree in Computer Science from University of Electronic Science and Technology of China and Xiamen University, in 2016 and 2019, respectively. She is currently working toward the D.Eng in the school of intelligent systems engineering of Sun Yat-sen University. Her research interests include vision-and-language understanding and embodied AI.



Yunshuang Nie received the B.E. degree in Sun Yat-sen University, Shenzhen, China, in 2023. She is currently working toward the M.E. in the school of intelligent systems engineering of Sun Yat-sen University. Her current research interest is vision-and-language understanding.



Hang Xu is currently a senior researcher in Huawei Noah Ark Lab. He received his BSc in Fudan University and Ph.D in Hong Kong University in Statistics. His research Interest includes multi-modality learning, machine Learning, object detection, and AutoML. He has published 70+ papers in Top AI conferences: NeurIPS, CVPR, ICCV, AAAI and some statistics journals, e.g. CSDA, Statistical Computing.



Ziming Wei is currently an undergraduate in the school of intelligent systems engineering of Sun Yat-sen University. His current research interests include vision-and-language understanding, multimodality and embodied agent.



Jiaqi Chen received the B.E. and M.E. degrees from Xidian University and Sun Yat-sen University, in 2019 and 2021, respectively. He is currently working toward the Ph.D. degree in the Computer Science Department at The University of Hong Kong. His research interests include multi-modal reasoning and embodied AI.



Xiaojun Chang (Senior Member, IEEE) is a Professor at Australian Artificial Intelligence Institute, University of Technology Sydney. He is also an Honorary Professor at the School of Computing Technologies, RMIT University. Before joining UTS, he was an Associate Professor at the School of Computing Technologies, RMIT University, Australia. After graduation, he subsequently worked as a Postdoc Research Fellow at School of Computer Science, Carnegie Mellon University, Lecturer and Senior Lecturer in the Faculty of Information Technology, Monash University, Australia. He has spent most of his time working on exploring multiple signals (visual, acoustic, textual) for automatic content analysis in unconstrained or surveillance videos. He has achieved top performances in various international competitions, such as TRECVID MED, TRECVID SIN, and TRECVID AVS.



Shikui Ma received the M.S. degree in Northern Jiaotong University in 2003. He has been serving as an assistant vice president of Dataa Robotics company (<https://www.dataarobotics.com/en>) since 2015, where he is leading HARIX and AI R&D team to consistently enhance their HARIX intelligent system, especially significantly improve the smart vision and motion capabilities of their robots via real-time digital twin, multi-modal fusion perception and advanced cognition, and deep reinforcement learning technologies. He has approximately 19 years of experience in communications technology industry, expertized in large-scale carrier-grade applications, especially competitive in Robotics, AI, Cloud, OSS and IPTV fields.



Xiaodan Liang is currently an Associate Professor at Sun Yat-sen University. She was a postdoc researcher in the machine learning department at Carnegie Mellon University, working with Prof. Eric Xing, from 2016 to 2018. She received her PhD degree from Sun Yat-sen University in 2016, advised by Liang Lin. She has published several cutting-edge projects on human-related analysis, including human parsing, pedestrian detection and instance segmentation, 2D/3D human pose estimation, and activity recognition.



Jianhua Han received the Bachelor Degree in 2016 and Master Degree in 2019 from Shanghai Jiao Tong University, China. He is currently a researcher with the Noahs Ark Laboratory, Huawei Technologies Co., Ltd. His research interests lie primarily in deep learning and computer vision.