

Vision-Language Navigation Policy Learning and Adaptation

Xin Wang, *Member, IEEE*, Qiuyuan Huang, *Member, IEEE*, Asli Celikyilmaz, *Member, IEEE*, Jianfeng Gao, *Fellow, IEEE*, Dinghan Shen, *Member, IEEE*, Yuan-Fang Wang, *Senior Member, IEEE*, William Yang Wang, *Member, IEEE*, and Lei Zhang, *Senior Member, IEEE*

(Invited Paper)

Abstract—Vision-language navigation (VLN) is the task of navigating an embodied agent to carry out natural language instructions inside real 3D environments. In this paper, we study how to address three critical challenges for this task: the cross-modal grounding, the ill-posed feedback, and the generalization problems. First, we propose a novel Reinforced Cross-Modal Matching (RCM) approach that enforces cross-modal grounding both locally and globally via reinforcement learning (RL). Particularly, a matching critic is used to provide an intrinsic reward to encourage global matching between instructions and trajectories, and a reasoning navigator is employed to perform cross-modal grounding in the local visual scene. Evaluation on a VLN benchmark dataset shows that our RCM model significantly outperforms baseline methods by 10% on Success Rate weighted by Path Length (SPL) and achieves the state-of-the-art performance. To improve the generalizability of the learned policy, we further introduce a Self-Supervised Imitation Learning (SIL) method to explore and adapt to unseen environments by imitating its own past, good decisions. We demonstrate that SIL can approximate a better and more efficient policy, which tremendously minimizes the success rate performance gap between seen and unseen environments (from 30.7% to 11.7%).

Index Terms—Vision-Language Navigation, Reinforcement Learning, Imitation Learning, Multimodal Machine Learning.

1 INTRODUCTION

VISION and language grounded embodied agents have received increased attention [1]–[3] due to their popularity in many intriguing real-world applications, e.g., in-home robots and personal assistants. Meanwhile, such an agent pushes forward visual and language grounding by putting itself in an active learning scenario through first-person vision. In particular, Vision-Language Navigation (VLN) [4] is the task of navigating an agent inside real environments by following natural language instructions. VLN requires a deep understanding of linguistic semantics, visual perception, and most importantly, the alignment of the two. The agent must reason about the vision-language dynamics in order to move towards the target that is inferred from the instructions.

VLN presents some unique challenges. First, reasoning over visual images and natural language instructions can be difficult. As we demonstrate in Figure 1, to reach a destination, the agent needs to ground an instruction in the local visual scene, represented as a sequence of words, as well as match the instruction to the visual trajectory in the global temporal space. Secondly, except for strictly following expert demonstrations, the feedback is rather coarse, since the “Success” feedback is provided only when the agent reaches a target position, completely ignoring whether the agent has followed the instructions (e.g., Path A in Figure 1) or followed a random path to reach the destination (e.g., Path C in Figure 1). Even a “good” trajectory that matches an instruction can

Instruction

Turn right and head towards the *kitchen*. Then turn left, pass a *table* and enter the *hallway*. Walk down the hallway and turn into the *entry* way to your right *without doors*. Stop in front of the *toilet*.

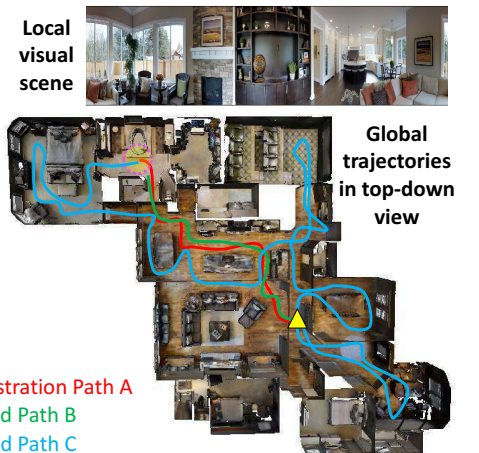


Fig. 1: Demonstration of the VLN task. The instruction, the local visual scene, and the global trajectories in a top-down view is shown. The agent does not have access to the top-down view. Path A is the demonstration path following the instruction. Path B and C are two different paths executed by the agent.

be considered unsuccessful if the agent stops marginally earlier than it should be (e.g., Path B in Figure 1). An ill-posed feedback can potentially deviate from the optimal policy learning. Thirdly, existing work suffers from the generalization problem, causing a huge performance gap between seen and unseen environments.

In this paper, we propose to combine the power of reinforcement learning (RL) and imitation learning (IL) to address the challenges above. First, we introduce a novel Reinforced Cross-Modal

- X. Wang, Y. Wang and W. Wang are with the Department of Computer Science, University of California, Santa Barbara, CA 93106, USA. E-mail: xwang@cs.ucsb.edu
- Q. Huang, A. Celikyilmaz, J. Gao and L. Zhang are with Microsoft Research, Redmond, USA.
- D. Shen is with Duke University, Durham, USA.

Matching (RCM) approach that enforces cross-modal grounding both locally and globally via RL. Specifically, we design a *reasoning navigator* that learns the cross-modal grounding in both the textual instruction and the local visual scene, so that the agent can infer which sub-instruction to focus on and where to look at. From the global perspective, we equip the agent with a *matching critic* that evaluates an executed path by the probability of reconstructing the original instruction from it, which we refer to as the cycle-reconstruction reward. Locally, the cycle-reconstruction reward provides a fine-grained intrinsic reward signal to encourage the agent to better understand the language input and penalize the trajectories that do not match the instructions. For instance, using the proposed reward, Path B is considered better than Path C (see Figure 1).

Being trained with the intrinsic reward from the matching critic and the extrinsic reward from the environment, the reasoning navigator learns to ground the natural language instruction on both local spatial visual scene and global temporal visual trajectory. Our RCM model significantly outperforms the existing methods and achieves new state-of-the-art performance on the Room-to-Room (R2R) dataset.

Our experimental results indicate a large performance gap between seen and unseen environments. To narrow the gap, we propose an effective solution to explore unseen environments with self-supervision. This technique is valuable because it can facilitate lifelong learning and adaptation to new environments. For example, domestic robots can explore a new home it arrives at and iteratively improve the navigation policy by learning from previous experience. Motivated by this fact, we introduce a Self-Supervised Imitation Learning (SIL) method in favor of exploration on unseen environments that do not have labeled data. The agent learns to imitate its own past, good experience. Specifically, in our framework, the navigator performs multiple roll-outs, of which good trajectories (determined by the matching critic) are stored in the replay buffer and later used for the navigator to imitate. In this way, the navigator can approximate its best behavior that leads to a better policy. To summarize, our contributions are mainly four-fold:

- We propose a novel Reinforced Cross-Modal Matching (RCM) framework that utilizes both extrinsic and intrinsic rewards for reinforcement learning, of which we introduce a cycle-reconstruction reward as the intrinsic reward to enforce the global matching between the language instruction and the agent’s trajectory.
- Our reasoning navigator learns the cross-modal contexts and makes decisions based on trajectory history, textual context, and visual context.
- Experiments show the effectiveness of the RCM model, which achieves the state-of-the-art performance on the R2R dataset.
- We introduce a new learning scenario for VLN, where exploring unseen environments prior to testing is allowed, and then propose a Self-Supervised Imitation Learning (SIL) method for exploration and adaptation, whose effectiveness and efficiency are validated on the R2R dataset.

2 RELATED WORK

2.1 Vision-and-Language Grounding

Recently, researchers in both computer vision and natural language processing are striving to bridge vision and natural language

towards a deeper understanding of the world [5]–[11], e.g., captioning an image or a video with natural language [12]–[18] or localizing desired objects within an image given a natural language description [19]–[22]. Moreover, visual question answering [23] and visual dialog [24] aim to generate one-turn or multi-turn response by grounding it on both visual and textual modalities. However, those tasks focus on passive visual perception in the sense that the visual inputs are usually fixed. In this work, we are particularly interested in solving the dynamic multi-modal grounding problem in both temporal and spatial spaces. Thus, we focus on the task of vision-language navigation (VLN) [4] which requires the agent to actively interact with the environment.

2.2 Embodied Navigation Agent

Navigation in 3D environments [25]–[28] is an essential capability of a mobile intelligent system that functions in the physical world. In the past two years, a plethora of tasks and evaluation protocols [1], [2], [4], [29], [30] have been proposed as summarized in [31]. VLN [4] focuses on language-grounded navigation in the real 3D environment. In order to solve the VLN task, Anderson et al. [4] set up an attention-based sequence-to-sequence baseline model. Then Wang et al. [32] introduced a hybrid approach that combines model-free and model-based reinforcement learning (RL) to improve the model’s generalizability. Lately, Fried et al. [33] proposed a speaker-follower model that adopts data augmentation, panoramic action space and modified beam search for VLN, establishing the current state-of-the-art performance on the Room-to-Room dataset. Extending prior work, we propose a Reinforced Cross-Modal Matching (RCM) approach to VLN. The RCM model is built upon [33] but differs in many significant aspects: (1) we combine a novel multi-reward RL with imitation learning for VLN while Speaker-Follower models [33] only uses supervised learning as in [4]. (2) Our reasoning navigator performs cross-modal grounding rather than the temporal attention mechanism on single-modality input. (3) Our matching critic is similar to Speaker in terms of the architecture design, but the former is used to provide the cycle-reconstruction intrinsic reward for both RL and SIL training while the latter is used to augment training data for supervised learning. Moreover, we introduce a self-supervised imitation learning method for exploration in order to explicitly address the generalization issue, which is a problem not well-studied in prior work. Concurrent to our work, [34]–[37] studies the VLN tasks from various aspects, and [38] introduces a variant of the VLN task to find objects by requesting language assistance when needed. Note that we are the first to propose to explore unseen environments for the VLN task.

2.3 Exploration

Much work has been done on improving exploration [39]–[43] because the trade-off between exploration and exploitation is one of the fundamental challenges in RL. The agent needs to exploit what it has learned to maximize reward and explore new territories for better policy search. Curiosity or uncertainty has been used as a signal for exploration [44]–[46]. Most recently, Oh et al. [47] proposed to exploit past good experience for better exploration in RL and theoretically justified its effectiveness. Our Self-Supervised Imitation Learning (SIL) method shares the same spirit. But instead of testing on games, we adapt SIL and validate its effectiveness and efficiency on the more practical task of VLN.

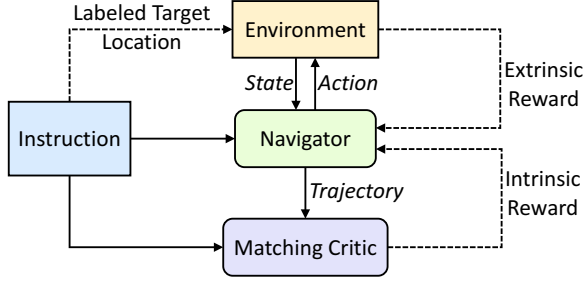


Fig. 2: Overview of our RCM framework.

3 REINFORCED CROSS-MODAL MATCHING

3.1 Overview

Here we consider an embodied agent that learns to navigate inside real indoor environments by following natural language instructions. The RCM framework mainly consists of two modules (see Figure 2): a **reasoning navigator** π_θ and a **matching critic** V_β . Given the initial state s_0 and the natural language instruction (a sequence of words) $\mathcal{X} = x_1, x_2, \dots, x_n$, the reasoning navigator learns to perform a sequence of actions $a_1, a_2, \dots, a_T \in \mathcal{A}$, which generates a trajectory τ , in order to arrive at the target location s_{target} indicated by the instruction \mathcal{X} . The navigator interacts with the environment and perceives new visual states as it executes actions. To promote the generalizability and reinforce the policy learning, we introduce two reward functions: an **extrinsic reward** that is provided by the environment and measures the success signal and the navigation error of each action, and an **intrinsic reward** that comes from our matching critic and measures the alignment between the language instruction \mathcal{X} and the navigator’s trajectory τ .

3.2 Model

Here we discuss the reasoning navigator and matching critic in details, both of which are end-to-end trainable.

3.2.1 Cross-Modal Reasoning Navigator

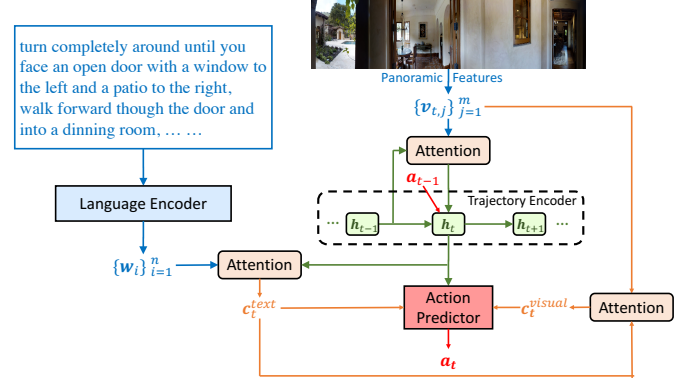
The navigator π_θ is a policy-based agent that maps the input instruction \mathcal{X} onto a sequence of actions $\{a_t\}_{t=1}^T$. At each time step t , the navigator receives a state s_t from the environment and needs to ground the textual instruction in the local visual scene. Thus, we design a cross-modal reasoning navigator that learns the trajectory history, the focus of the textual instruction, and the local visual attention in order, which forms a cross-modal reasoning path to encourage the local dynamics of both modalities at step t .

Figure 3 shows the unrolled version of the navigator at time step t . Similar to [33], we equip the navigator with a panoramic view, which is split into image patches of m different viewpoints, so the panoramic features that are extracted from the visual state s_t can be represented as $\{v_{t,j}\}_{j=1}^m$, where $v_{t,j}$ denotes the pre-trained CNN feature of the image patch at viewpoint j .

History Context: Once the navigator runs one step, the visual scene would change accordingly. The history of the trajectory $\tau_{1:t}$ till step t is encoded as a history context vector h_t by an attention-based trajectory encoder LSTM [48]:

$$h_t = LSTM([v_t, a_{t-1}], h_{t-1}) \quad (1)$$

where a_{t-1} is the action taken at previous step, and $v_t = \sum_j \alpha_{t,j} v_{t,j}$, the weighted sum of the panoramic features. $\alpha_{t,j}$

Fig. 3: Cross-modal reasoning navigator at step t .

is the attention weight of the visual feature $v_{t,j}$, representing its importance with respect to the previous history context h_{t-1} . Note that we adopt the dot-product attention [49] hereafter, which we denote as (taking the attention over visual features above for an example)

$$v_t = attention(h_{t-1}, \{v_{t,j}\}_{j=1}^m) \quad (2)$$

$$= \sum_j softmax(h_{t-1} W_h (v_{t,j} W_v)^T) v_{t,j} \quad (3)$$

where W_h and W_v are learnable projection matrices.

Visually Grounded Textual Context: Memorizing the past can enable the recognition of the current status and thus understanding which words or sub-instructions to focus on next. Hence, we further learn the textual context c_t^{text} conditioned on the history context h_t . We let a language encoder LSTM to encode the language instruction \mathcal{X} into a set of textual features $\{w_i\}_{i=1}^n$. Then at every time step, the textual context is computed as

$$c_t^{text} = attention(h_t, \{w_i\}_{i=1}^n) \quad (4)$$

Note that c_t^{text} weighs more on the words that are more relevant to the trajectory history and the current visual state.

Textually Grounded Visual Context: Knowing where to look at requires a dynamic understanding of the language instruction; so we compute the visual context c_t^{visual} based on the textual context c_t^{text} :

$$c_t^{visual} = attention(c_t^{text}, \{v_j\}_{j=1}^m) \quad (5)$$

Action Prediction: In the end, our action predictor considers the history context h_t , the textual context c_t^{text} , and the visual context c_t^{visual} , and decides which direction to go next based on them. It calculates the probability p_k of each navigable direction using a bilinear dot product as follows:

$$p_k = softmax([h_t, c_t^{text}, c_t^{visual}] W_c (u_k W_u)^T) \quad (6)$$

where u_k is the action embedding that represents the k -th navigable direction, which is obtained by concatenating an appearance feature vector (CNN feature vector extracted from the image patch around that view angle or direction) and a 4-dimensional orientation feature vector $[\sin\psi; \cos\psi; \sin\omega; \cos\omega]$, where ψ and ω are the heading and elevation angles respectively. The learning objectives for training the navigator are introduced in Section 3.3.

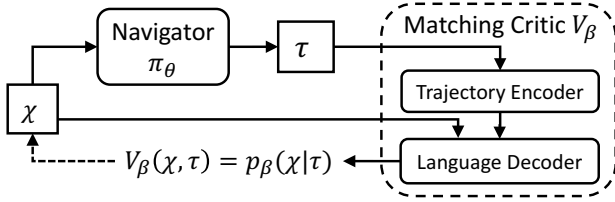


Fig. 4: Cross-modal matching critic that provides the cycle-reconstruction intrinsic reward.

3.2.2 Cross-Modal Matching Critic

In addition to the extrinsic reward signal from the environment, we also derive an intrinsic reward R_{intr} provided by the matching critic V_β to encourage the global matching between the language instruction \mathcal{X} and the navigator π_θ 's trajectory $\tau = \{ \langle s_1, a_1 \rangle, \langle s_2, a_2 \rangle, \dots, \langle s_T, a_T \rangle \}$:

$$R_{intr} = V_\beta(\mathcal{X}, \tau) = V_\beta(\mathcal{X}, \pi_\theta(\mathcal{X})) \quad (7)$$

One way to realize this goal is to measure the cycle-reconstruction reward $p(\hat{\mathcal{X}} = \mathcal{X} | \pi_\theta(\mathcal{X}))$, the probability of reconstructing the language instruction \mathcal{X} given the trajectory $\tau = \pi_\theta(\mathcal{X})$ executed by the navigator. The higher the probability is, the better the produced trajectory is aligned with the instruction.

Therefore as shown in Figure 4, we adopt an attention-based sequence-to-sequence language model as our matching critic V_β , which encodes the trajectory τ with a trajectory encoder and produces the probability distributions of generating each word of the instruction \mathcal{X} with a language decoder. Hence the intrinsic reward

$$R_{intr} = p_\beta(\mathcal{X} | \pi_\theta(\mathcal{X})) = p_\beta(\mathcal{X} | \tau) = \frac{1}{n} \sum_i^n \log p_\beta(x_i | \tau) \quad (8)$$

which is normalized by the instruction length n . In our experiments, the matching critic is pre-trained with human demonstrations (the ground-truth instruction-trajectory pairs $\langle \mathcal{X}^*, \tau^* \rangle$) via supervised learning.

3.3 Learning

In order to quickly approximate a relatively good policy, we use the demonstration actions to conduct supervised learning with maximum likelihood estimation (MLE). The training loss L_{sl} is defined as

$$L_{sl} = -\mathbb{E}[\log(\pi_\theta(a_t^* | s_t))] \quad (9)$$

where a_t^* is the demonstration action provided by the simulator. Warm starting the agent with supervised learning can ensure a relatively good policy on the seen environments. But it also limits the agent's generalizability to recover from erroneous actions in unseen environments, since it only clones the behaviors of expert demonstrations.

To learn a better and more generalizable policy, we then switch to reinforcement learning and introduce the extrinsic and intrinsic reward functions to refine the policy from different perspectives.

3.3.1 Extrinsic Reward

A common practice in RL is to directly optimize the evaluation metrics. Since the objective of the VLN task is to successfully reach the target location s_{target} , we consider two metrics for the reward design. The first metric is the relative navigation distance

similar to [32]. We denote the distance between s_t and s_{target} as $\mathcal{D}_{target}(s_t)$. Then the immediate reward $r(s_t, a_t)$ after taking action a_t at state s_t ($t < T$) becomes:

$$r(s_t, a_t) = \mathcal{D}_{target}(s_t) - \mathcal{D}_{target}(s_{t+1}), \quad t < T \quad (10)$$

This indicates the reduced distance to the target location after taking action a_t . Our second choice considers the "Success" as an additional criterion. If the agent reaches a point within a threshold measured by the distance d from the target (d is preset as 3m in the R2R dataset), then it is counted as "Success". Particularly, the immediate reward function at last step T is defined as

$$r(s_T, a_T) = \mathbb{1}(\mathcal{D}_{target}(s_T) \leq d) \quad (11)$$

where $\mathbb{1}()$ is an indicator function. To incorporate the influence of the action a_t on the future and account for the local greedy search, we use the discounted cumulative reward rather than the immediate reward to train the policy:

$$R_{extr}(s_t, a_t) = \underbrace{r(s_t, a_t)}_{\text{immediate reward}} + \underbrace{\sum_{t'=t+1}^T \gamma^{t'-t} r(s_{t'}, a_{t'})}_{\text{discounted future reward}} \quad (12)$$

where γ is the discounted factor (0.95 in our experiments).

3.3.2 Intrinsic Reward

As discussed in Section 3.2.2, we pre-train a matching critic to calculate the cycle-reconstruction intrinsic reward R_{intr} (see Equation 8), promoting the alignment between the language instruction \mathcal{X} and the trajectory τ . It encourages the agent to respect the instruction and penalizes the paths that deviate from what the instruction indicates.

With both the extrinsic and intrinsic reward functions, the RL loss can be written as

$$L_{rl} = -\mathbb{E}_{a_t \sim \pi_\theta} [A_t] \quad (13)$$

where the advantage function $A_t = R_{extr} + \delta R_{intr}$. δ is a hyper-parameter weighing the intrinsic reward. Based on REINFORCE algorithm [50], the gradient of non-differentiable, reward-based loss function can be derived as

$$\nabla_\theta L_{rl} = -A_t \nabla_\theta \log \pi_\theta(a_t | s_t) \quad (14)$$

4 SELF-SUPERVISED IMITATION LEARNING

The last section introduces the effective RCM method for generic vision-language navigation task, whose standard setting is to train the agent on seen environments and test it on unseen environments without exploration. In this section we discuss a different setting where the agent is allowed to explore unseen environments without ground-truth demonstrations. This is of practical benefit because it facilitates lifelong learning and adaption to new environments.

To this end, we propose a Self-Supervised Imitation Learning (SIL) method to imitate the agent's own past good decisions. As shown in Figure 5, given a natural language instruction \mathcal{X} without paired demonstrations and ground-truth target location, the navigator produces a set of possible trajectories and then stores the best trajectory $\hat{\tau}$ that is determined by matching critic V_β into a replay buffer, in formula,

$$\hat{\tau} = \arg \max_{\tau} V_\beta(\mathcal{X}, \tau) \quad (15)$$

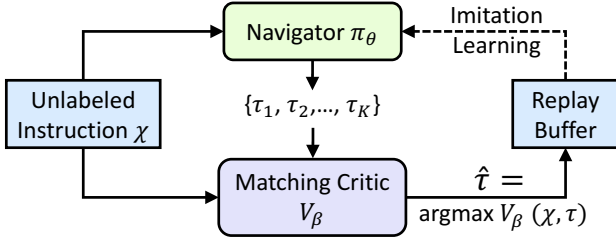


Fig. 5: SIL for exploration on unlabeled data.

The matching critic evaluates the trajectories with the cycle-reconstruction reward as introduced in Section 3.2.2. Then by exploiting the good trajectories in the replay buffer, the agent is indeed optimizing the following objective with self-supervision. The target location is unknown and thus there is no supervision from the environment.

$$L_{sil} = -R_{intr} \log \pi_{\theta}(a_t | s_t) \quad (16)$$

Note that L_{sil} can be viewed as the loss for policy gradient except that the off-policy Monte-Carlo return R_{intr} is used instead of on-policy return. L_{sil} can also be interpreted as the supervised learning loss with $\hat{\tau}$ as the ‘‘ground truths’’:

$$L_{sil} = -\mathbb{E}[\log(\pi_{\theta}(\hat{a}_t | s_t))] \quad (17)$$

where \hat{a}_t is the action stored in the replay buffer using Equation 15. Paired with a matching critic, the SIL method can be combined with various learning methods to approximate a better policy by imitating the previous best of itself.

5 EXPERIMENTAL SETUP

5.1 R2R Dataset

We evaluate our approaches on the Room-to-Room (R2R) dataset [4] for vision-language navigation in real 3D environments. The R2R dataset is built upon the Matterport3D dataset [51], which consists of 10,800 panoramic views constructed from 194,400 RGB-D images of 90 building-scale scenes (Many of the scenes can be viewed in the Matterport 3D spaces gallery2). The R2R dataset samples 7,189 paths capturing most of the visual diversity in the dataset and collects 21,567 navigation instructions with an average length of 29 words (each path is paired with 3 different instructions). The R2R dataset is split into training (14,025 instructions), seen validation (1,020), unseen validation (2,349), and test (4,173) sets. The seen validation set shares the same environments with the training set. While both the unseen validation and test sets contain distinct environments that do not appear in the other sets.

5.2 Testing Scenarios

The standard testing scenario of the VLN task is to train the agent in seen environments and then test it in previously unseen environments in a zero-shot fashion. There is no prior exploration on the test set. This setting is preferred and able to clearly measure the generalizability of the navigation policy, so we evaluate our RCM approach under the standard testing scenario.

Furthermore, exploration in unseen environments is certainly meaningful in practice, e.g., in-home robots are expected to explore and adapt to a new environment. So we introduce a lifelong

learning scenario where the agent is encouraged to learn from trials and errors on the unseen environments. In this case, how to effectively explore the unseen validation or test set where there are no expert demonstrations becomes an important task to study.

5.3 Evaluation Metrics

We report five evaluation metrics as used by the VLN Challenge: (1) Path Length (PL): the total length of the executed path; (2) Navigation Error (NE): the shortest-path distance between the agent’s final position and the target; (3) Oracle Success Rate (OSR): the success rate at the closest point to the goal that the agent has visited along the trajectory; (4) Success Rate (SR): the percentage of predicted end-locations within 3m of the target locations; (5) Success rate weighted by inverse Path Length (SPL): SPL trades-off Success Rate against Path Length. Among those metrics, SPL is the recommended primary measure of navigation performance [31], as it considers both effectiveness and efficiency. The other metrics are also reported as auxiliary measures.

5.4 Training Details

Following prior work [4], [32], [33], ResNet-152 CNN features [52] are extracted for all images without fine-tuning. The pretrained GloVe word embeddings [53] are used for initialization and then fine-tuned during training. All the hyper-parameters are tuned on the validation sets. We adopt the panoramic action space [33] where the action is to choose a navigable direction from the possible candidates. We set the maximal length of the action path as 10. The maximum length of the instruction is set as 80 and longer instructions are truncated.

We train the matching critic with a learning rate 1e-4 and then fix it during policy learning. Then we warm start the policy via supervised learning loss with a learning rate 1e-4, and then switch to RL training with a learning rate 1e-5. Self-supervised imitation learning can be performed to further improve the policy: during the first epoch of SIL, the loaded policy produces 10 trajectories, of which the one with the highest intrinsic reward is stored in the replay buffer; those saved trajectories are then utilized to fine-tune the policy for a fixed number of iterations (the learning rate is 1e-5). Early stopping is used for all the training and Adam optimizer [54] is used to optimize all the parameters. To avoid overfitting, we use an L2 weight decay of 0.0005 and a dropout ratio of 0.5. The discounted factor γ of our cumulative reward is 0.95. The weight σ of the intrinsic reward is set as 2.

5.5 Network Architecture

5.5.1 Reasoning Navigator

The language encoder consists of an LSTM with hidden size 512 and a word embedding layer of size 300. The inner dimensions of the three attention modules used to compute the history context, the textual context, and the visual context are 256, 512, and 256 respectively. The trajectory encoder is an LSTM with hidden size 512. The action embedding is a concatenation of the visual appearance feature vector of size 2048 and the orientation feature vector of size 128 (the 4-dimensional orientation feature $[\sin\psi; \cos\psi; \sin\omega; \cos\omega]$ are tiled 32 times as used in [33]). The action predictor is composed of three weight matrices: the projection dimensions of W_c and W_u are both 256, and then an output layer W_o together with a softmax layer are followed to obtain the probabilities over the possible navigable directions.

Test Set (VLN Challenge Leaderboard)					
Model	PL ↓	NE ↓	OSR ↑	SR ↑	SPL ↑
Low-level Visuomotor Action Space					
Random	9.89	9.79	18.3	13.2	12
seq2seq [4]	8.13	7.85	26.6	20.4	18
RPA [32]	9.15	7.53	32.5	25.3	23
High-level Panoramic Action Space					
Speaker-Follower [33]	14.82	6.62	44.0	35.0	28
+ beam search	<u>1257.38</u>	4.87	96.0	53.5	<u>1</u>
RCM	15.22	6.01	50.8	43.1	35
RCM + SIL (train)	11.97	6.12	49.5	43.0	38
RCM + SIL (unseen) ¹	9.48	4.21	66.8	60.5	59

TABLE 1: Comparison on the R2R test set [4]. Our RCM model significantly outperforms the prior work, especially on SPL (the primary metric for navigation tasks [31]). Moreover, using SIL to imitate itself on the training set can further improve its efficiency: the path length is shortened by 3.25m. Note that with beam search, the agent executes K trajectories at test time and chooses the most confident one as the ending point, which results in a super long path and is heavily penalized by SPL.

5.5.2 Matching Critic

The matching critic consists of an attention-based trajectory encoder with the same architecture as the one in the navigator, its own word embedding layer of size 300, and an attention-based language decoder. The language decoder is composed of an attention module (whose projection dimension is 512) over the encoded features, an LSTM of hidden size 512, and a multi-layer perceptron (Linear \rightarrow Tanh \rightarrow Linear \rightarrow SoftMax) that converts the hidden state into probabilities of all the words in the vocabulary.

6 EXPERIMENTS AND ANALYSIS

6.1 Results on the Test Set

6.1.1 Comparison with Prior Work

We mainly compare the performance of RCM to the previous baseline methods on the test set of the R2R dataset, which is held out as the VLN Challenge. The results are shown in Table 1, where we compare RCM to a set of baselines: (1) *Random*: randomly take a direction to move forward at each step until five steps. (2) *seq2seq*: the best-performing sequence-to-sequence model as reported in the original dataset paper [4], which is trained with the student-forcing method. (3) *RPA*: a reinforced planning-ahead model that combines model-free and model-based reinforcement learning for VLN [32]. (4) *Speaker-Follower*: a compositional Speaker-Follower method that combines data augmentation, panoramic action space, and beam search for VLN [33].

As can be seen in Table 1, RCM significantly outperforms the existing methods, improving the SPL score from 28% to 35%.² The improvement is consistently observed on the other metrics,

1. The results of using SIL to explore unseen environments are only used to validate its effectiveness for lifelong learning, which is not directly comparable to other models due to different learning scenarios.

2. Note that our RCM model also utilizes the panoramic action space and augmented data in [33] for a fair comparison.

e.g., the success rate is increased by 8.1%. Moreover, using SIL to imitate the RCM agent’s previous best behaviors on the training set can approximate a more efficient policy, whose average path length is reduced from 15.22m to 11.97m and which achieves the best result (38%) on SPL. Therefore, we submit the results of *RCM + SIL (train)* to the VLN Challenge, ranking first among prior work in terms of SPL. It is worth noticing that beam search is not practical in reality, because it needs to execute a very long trajectory before making the decision, which is punished heavily by the primary metric SPL. So we are mainly comparing the results without beam search.

6.1.2 Self-Supervised Imitation Learning

As mentioned above, for a standard VLN setting, we employ SIL on the training set to learn an efficient policy. For the lifelong learning scenario, we test the effectiveness of SIL on exploring unseen environments (the validation and test sets). It is noticeable in Table 1 that SIL indeed leads to a better policy even without knowing the target locations. SIL improves RCM by 17.5% on SR and 21% on SPL. Similarly, the agent also learns a more efficient policy that takes less number of steps (the average path length is reduced from 15.22m to 9.48m) but obtains a higher success rate. The key difference between SIL and beam search is that SIL optimizes the policy itself by play-and-imitate while beam search only makes a greedy selection of the rollouts of the existing policy. But we would like to point out that due to different learning scenarios, the results of *RCM + SIL (unseen)* cannot be directly compared with other methods following the standard settings of the VLN challenge.

6.2 Effect of Individual Components

We conduct an ablation study to illustrate the effect of each component on both seen and unseen validation sets in Table 2. Comparing Row 1 and Row 2, we observe the efficiency of the learned policy by imitating the best of itself on the training set. Then we start with the RCM model in Row 2, and successively remove the *intrinsic reward*, *extrinsic reward*, and *cross-modal reasoning* to demonstrate their importance.

Removing the intrinsic reward (Row 3), we notice that the success rate (SR) on unseen environments drops 1.9 points while it is almost fixed on seen environments (0.2 \uparrow). It evaluates the alignment between instructions and trajectories, serving as a complementary supervision besides of the feedback from the environment, therefore it works better for the unseen environments that require more supervision due to lack of exploration. This also indirectly validates the importance of exploration on unseen environments.

Furthermore, the results of Row 4 (the RCM model with only supervised learning) validate the superiority of reinforcement learning compared to purely supervised learning on the VLN task. Meanwhile, since eventually the results are evaluated based on the success rate (SR) and path length (PL), directly optimizing the extrinsic reward signals can guarantee the stability of the reinforcement learning and bring a big performance gain.

We then verify the strength of our cross-modal reasoning navigator by comparing it (Row 4) with an attention-based sequence-to-sequence model (Row 5) that utilizes the previous hidden state h_{t-1} to attend to both the visual and textual features at decoding time. Everything else is exactly the same except the cross-modal attention design. Evidently, our navigator improves upon the

#	Model	Seen Validation				Unseen Validation			
		PL ↓	NE ↓	OSR ↑	SR ↑	PL ↓	NE ↓	OSR ↑	SR ↑
0	Speaker-Follower (no beam search) [33]	-	3.36	73.8	66.4	-	6.62	45.0	35.5
1	RCM + SIL (train)	10.65	3.53	75.0	66.7	11.46	6.09	50.1	42.8
2	RCM	11.92	3.37	76.6	67.4	14.84	5.88	51.9	42.5
3	– intrinsic reward	12.08	3.25	77.2	67.6	15.00	6.02	50.5	40.6
4	– extrinsic reward = pure SL	11.99	3.22	76.7	66.9	14.83	6.29	46.5	37.7
5	– cross-modal reasoning	11.88	3.18	73.9	66.4	14.51	6.47	44.8	35.7
6	RCM + SIL (unseen)	10.13	2.78	79.7	73.0	9.12	4.17	69.31	61.3

TABLE 2: Ablation study on seen and unseen validation sets as reported in [55]. We report the performance of the speaker-follower model without beam search as the baseline. Row 1-5 shows the influence of each individual component by successively removing it from the final model. Row 6 illustrates the power of SIL on exploring unseen environments with self-supervision. Please see Section 6.2 for more detailed analysis.

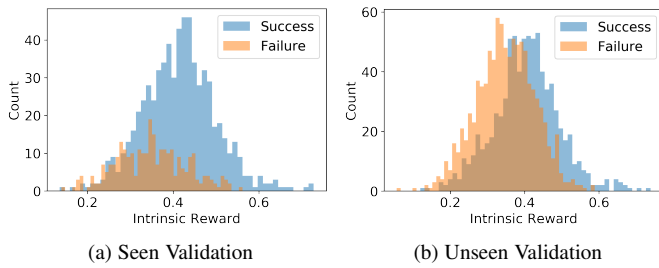


Fig. 6: Visualization of the intrinsic reward on seen and unseen validation sets.

baseline by considering history context, visually-conditioned textual context, and textually-conditioned visual context for decision making.

In the end, we demonstrate the effectiveness of the proposed SIL method for exploration in Row 6. Considerable performance boosts have been obtained on both seen and unseen environments, as the agent learns how to better execute the instructions from its own previous experience.

6.3 Generalizability

Another observation from the experiments (e.g., see Table 2) is that our RCM approach is much more generalizable to unseen environments compared with the baseline. The improvements on the seen and unseen validation sets are 0.3 and 7.1 points, respectively. So is the SIL method, which explicitly explores the unseen environments and tremendously reduces the success rate performance gap between seen and unseen environments from 30.7% (Row 5) to 11.7% (Row 6).

6.4 Visualizing Intrinsic Reward

In Figure 6, we plot the histogram distributions of the intrinsic rewards (produced by our submitted model) on both seen and unseen validation sets. On the one hand, the intrinsic reward is aligned with the success rate to some extent, because the successful examples are receiving higher averaged intrinsic rewards than the failed ones. On the other hand, the complementary intrinsic reward provides more fine-grained reward signals to reinforce multi-modal grounding and improve the navigation policy learning.

6.5 Case Study

For a more intuitive view of how our model works for the VLN task, we visualize two qualitative examples in Figure 7. Particularly, we choose two examples, both with high intrinsic rewards. In (a), the agent successfully reaches the target destination, with a comprehensive understanding of the natural language instruction. While in (b), the intrinsic reward is also high, which indicates most of the agent’s actions are good, but it is also noticeable that the agent fails to recognize *the laundry room* at the end of the trajectory, which shows the importance of more precise visual grounding in the navigation task.

6.6 ML + RL vs. MIXER

Inspired by the recent work by Tan et al. [58], instead of using MIXER [57] where the policy model is first warmed up with behavior cloning and then finetuned with reinforcement learning objective, we conduct additional experiments where we train the model with a weighted sum of both supervised learning and reinforcement learning objectives (ML + RL) as in [56]. The updated results are shown in Table 3. As you can see, switching from MIXER to ML + RL, the RCM model can achieve 48.1% Success Rate on the unseen validation set, improving by 5.6%. The experiments also demonstrate that it is not stable to only use the intrinsic reward to finetune the model and the mixed reward works the best.

6.7 Error Analysis

In this section, we further analyze the negative examples and showcase a few common errors in the vision-language navigation task. First, a common mistake comes from the misunderstanding of the natural language instruction. Figure 8 demonstrate such a qualitative example, where the agent successfully perceived the concepts “hallway”, “turn left”, and “mirror” etc., but misinterpreted the meaning of the whole instruction. It turned left earlier and mistakenly entered the bathroom instead of the bedroom at Step 3.

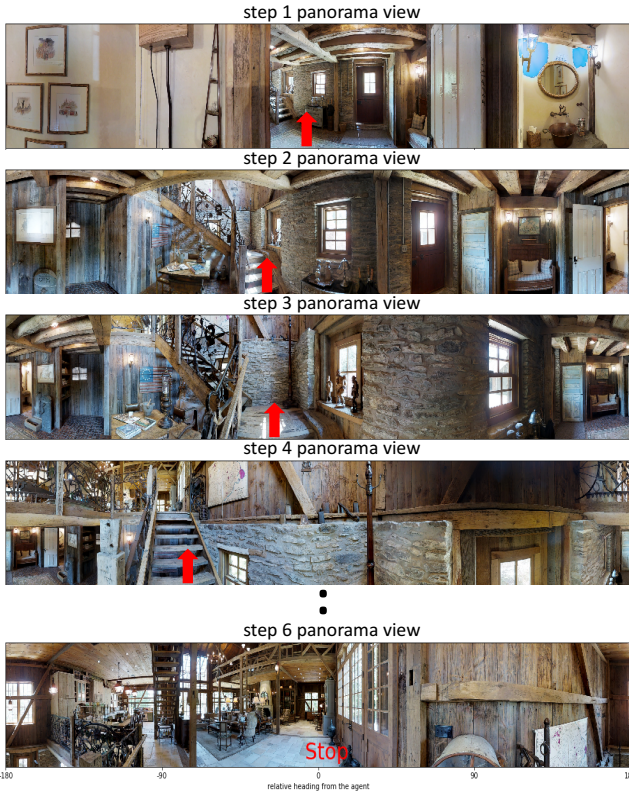
Secondly, failing to ground objects in the visual scene can usually result in an error. As shown in Figure 9 (a), the agent did not recognize the “mannequins” in the end (Step 5) and stopped at a wrong place even though it executed the instruction pretty well. Similar in Figure 9 (b), the agent failed to detect the “red ropes” at the beginning (Step 1) and thus took a wrong direction which also has the “red carpet”. Note that “mannequins”

Training Method	Seen Validation				Unseen Validation			
	PL ↓	NE ↓	OSR ↑	SR ↑	PL ↓	NE ↓	OSR ↑	SR ↑
MIXER	11.92	3.37	76.6	67.4	14.84	5.88	51.9	42.5
ML + RL	10.47	4.68	66.3	58.4	12.16	5.45	51.6	44.5
+ finetuning w Extrinsic Reward	11.41	4.32	64.4	56.6	14.69	5.35	54.7	46.7
+ finetuning w Intrinsic Reward	9.66	5.14	56.4	48.3	9.42	6.08	48.7	41.8
+ finetuning w Extrinsic & Intrinsic Rewards	11.92	4.21	66.3	57.4	14.62	5.32	56.8	48.1

TABLE 3: Updated results on seen and unseen validation sets with ML + RL objective [56] instead of MIXER [57].

Instruction: Exit the door and turn left towards the staircase. Walk all the way up the stairs, and stop at the top of the stairs.

Intrinsic Reward: 0.53 Result: Success (error = 0m)



(a) A successful case

Instruction: Turn right and go down the stairs. Turn left and go straight until you get to *the laundry room*. Wait there.

Intrinsic Reward: 0.54 Result: Failure (error = 5.5m)



(b) A failure case

Fig. 7: Qualitative examples from the unseen validation set.

is an out-of-vocabulary word in the training data; besides, both “mannequins” and “red ropes” do not belong to the 1000 classes of the ImageNet [59], so the visual features extracted from a pretrain ImageNet model [52] are not able to represent them.

In Figure 10, we illustrate a long negative trajectory which our agent produced by following a relatively complicated instruction. In this case, the agent match “the floor is in a circle pattern” with the visual scene, which seems to be another limitation of the current visual recognition systems. The above examples also suffer from the error accumulation issue as pointed out by Wang et al. [32], where one bad decision leads to a series of bad decisions during the navigation process. Therefore, an agent capable of being aware of and recovering from errors is desired for future study.

7 CONCLUSION

In this paper we present two novel approaches, RCM and SIL, which combine the strength of reinforcement learning and self-supervised imitation learning for the vision-language navigation task. Experiments illustrate the effectiveness and efficiency of our methods under both the standard testing scenario and the lifelong learning scenario. Moreover, our methods show strong generalizability in unseen environments. The proposed learning frameworks are modular and model-agnostic, which allow the components to be improved separately. We also believe that the idea of learning more fine-grained intrinsic rewards, in addition to the coarse external signals, is commonly applicable to various embodied agent tasks, and the idea SIL can be generally adopted to explore other unseen environments.

Instruction: Through hallway toward clock on the wall. Turn left at the mirror. Enter bedroom. Walk straight through the bedroom stopping just inside of walk-in closet.

Intrinsic Reward: 0.18 Result: Failure (error = 13.5m)



Fig. 8: Misunderstanding of the instruction.

ACKNOWLEDGMENT

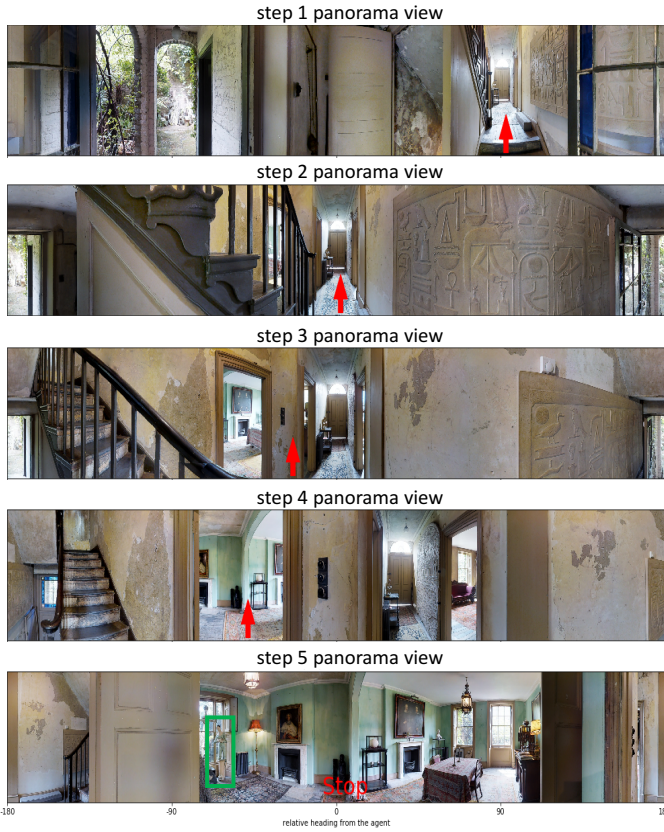
This work was partly performed when Xin Wang was interning at Microsoft Research. The authors would like to thank Peter Anderson and Pengchuan Zhang for their helpful discussions, and Ronghang Hu for his visualization code.

REFERENCES

- [1] M. Savva, A. X. Chang, A. Dosovitskiy, T. Funkhouser, and V. Koltun, "Minos: Multimodal indoor simulator for navigation in complex environments," *arXiv preprint arXiv:1712.03931*, 2017.
- [2] E. Kolve, R. Mottaghi, D. Gordon, Y. Zhu, A. Gupta, and A. Farhadi, "Ai2-thor: An interactive 3d environment for visual ai," *arXiv preprint arXiv:1712.05474*, 2017.
- [3] A. Das, S. Datta, G. Gkioxari, S. Lee, D. Parikh, and D. Batra, "Embodied Question Answering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [4] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sünderhauf, I. Reid, S. Gould, and A. van den Hengel, "Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 2, 2018.
- [5] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *International Conference on Machine Learning*, 2015, pp. 2048–2057.
- [6] X. Wang, W. Chen, Y.-F. Wang, and W. Y. Wang, "No metrics are perfect: Adversarial reward learning for visual storytelling," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2018.
- [7] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3128–3137.
- [8] X. Chen and C. Lawrence Zitnick, "Mind's eye: A recurrent visual representation for image caption generation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 2422–2431.
- [9] R. Hu, M. Rohrbach, and T. Darrell, "Segmentation from natural language expressions," in *European Conference on Computer Vision*. Springer, 2016, pp. 108–124.
- [10] M. Tapaswi, Y. Zhu, R. Stiefelwagen, A. Torralba, R. Urtasun, and S. Fidler, "Movieqa: Understanding stories in movies through question-answering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4631–4640.
- [11] Q. Huang, P. Zhang, D. Wu, and L. Zhang, "Turbo learning for captionbot and drawingbot," in *Advances in Neural Information Processing Systems (NIPS)*, 2018.
- [12] J. Donahue, L. A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell, "Long-term recurrent convolutional networks for visual recognition and description," in *CVPR*, 2015.
- [13] H. Fang, S. Gupta, F. Iandola, R. K. Srivastava, L. Deng, P. Dollár, J. Gao, X. He, M. Mitchell, J. C. Platt *et al.*, "From captions to visual concepts and back," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1473–1482.
- [14] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: A neural image caption generator," in *Computer Vision and Pattern Recognition (CVPR), 2015 IEEE Conference on*. IEEE, 2015, pp. 3156–3164.
- [15] X. Wang, W. Chen, J. Wu, Y.-F. Wang, and W. Y. Wang, "Video captioning via hierarchical reinforcement learning," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [16] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola, "Stacked attention networks for image question answering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 21–29.
- [17] H. Yu, J. Wang, Z. Huang, Y. Yang, and W. Xu, "Video paragraph captioning using hierarchical recurrent neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4584–4593.
- [18] X. Wang, Y.-F. Wang, and W. Y. Wang, "Watch, listen, and describe: Globally and locally aligned cross-modal attentions for video captioning," in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, 2018.
- [19] B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, J. Hockenmaier, and S. Lazebnik, "Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2641–2649.
- [20] R. Hu, H. Xu, M. Rohrbach, J. Feng, K. Saenko, and T. Darrell, "Natural language object retrieval," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4555–4564.
- [21] L. Yu, P. Poirson, S. Yang, A. C. Berg, and T. L. Berg, "Modeling context in referring expressions," in *European Conference on Computer Vision*. Springer, 2016, pp. 69–85.
- [22] D. Zhang, X. Dai, X. Wang, Y.-F. Wang, and L. S. Davis, "Man: Moment alignment network for natural language moment retrieval via iterative graph adjustment," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1247–1257.
- [23] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, "VQA: Visual Question Answering," in *International Conference on Computer Vision (ICCV)*, 2015.
- [24] A. Das, S. Kottur, K. Gupta, A. Singh, D. Yadav, J. M. Moura, D. Parikh, and D. Batra, "Visual Dialog," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [25] Y. Zhu, R. Mottaghi, E. Kolve, J. J. Lim, A. Gupta, L. Fei-Fei, and A. Farhadi, "Target-driven visual navigation in indoor scenes using deep reinforcement learning," in *Robotics and Automation (ICRA), 2017 IEEE International Conference on*. IEEE, 2017, pp. 3357–3364.
- [26] P. Mirowski, R. Pascanu, F. Viola, H. Soyer, A. J. Ballard, A. Banino, M. Denil, R. Goroshin, L. Sifre, K. Kavukcuoglu *et al.*, "Learning to navigate in complex environments," *arXiv preprint arXiv:1611.03673*, 2016.
- [27] A. Mousavian, A. Toshev, M. Fiser, J. Kosecka, and J. Davidson, "Visual representations for semantic target driven navigation," *arXiv preprint arXiv:1805.06066*, 2018.
- [28] S. Hemachandra, F. Duvallat, T. M. Howard, N. Roy, A. Stentz, and M. R. Walter, "Learning models for following natural language directions in unknown environments," in *2015 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2015, pp. 5608–5615.
- [29] S. Song, F. Yu, A. X. Chang, M. Savva, and T. Funkhouser, "Semantic scene completion from a single depth image," *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [30] F. Xia, A. R. Zamir, Z.-Y. He, A. Sax, J. Malik, and S. Savarese, "Gibson Env: real-world perception for embodied agents," in *Computer Vision*

Instruction: Go up the stairs to the right, turn left and go into the room on the left. Turn left and stop near the **mannequins**.

Intrinsic Reward: 0.51 Result: Failure (error = 3.1m)



(a)

Instruction: With the **red ropes** to your right, walk down the room on the red carpet past the display. Turn left when another red carpet meets the one you are on in a right angle. Stop on the carpet where these two directions of carpet meet.

Intrinsic Reward: 0.17 Result: Failure (error = 19.6m)



(b)

Fig. 9: Ground errors where objects were not recognized from the visual scene.

and Pattern Recognition (CVPR), 2018 IEEE Conference on. IEEE, 2018.

[31] P. Anderson, A. Chang, D. S. Chaplot, A. Dosovitskiy, S. Gupta, V. Koltun, J. Kosecka, J. Malik, R. Mottaghi, M. Savva *et al.*, “On evaluation of embodied navigation agents,” *arXiv preprint arXiv:1807.06757*, 2018.

[32] X. Wang, W. Xiong, H. Wang, and W. Y. Wang, “Look before you leap: Bridging model-free and model-based reinforcement learning for planned-ahead vision-and-language navigation,” in *The European Conference on Computer Vision (ECCV)*, September 2018.

[33] 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 *Advances in Neural Information Processing Systems (NIPS)*, 2018.

[34] J. Thomason, D. Gordon, and Y. Bisk, “Shifting the baseline: Single modality performance on visual navigation & QA,” in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 1977–1983.

[35] L. Ke, X. Li, Y. Bisk, A. Holtzman, Z. Gan, J. Liu, J. Gao, Y. Choi, and S. Srinivasa, “Tactical rewind: Self-correction via backtracking in vision-and-language navigation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6741–6749.

[36] C.-Y. Ma, J. Lu, Z. Wu, G. AlRegib, Z. Kira, R. Socher, and C. Xiong, “Self-monitoring navigation agent via auxiliary progress estimation,” *arXiv preprint arXiv:1901.03035*, 2019.

[37] C.-Y. Ma, Z. Wu, G. AlRegib, C. Xiong, and Z. Kira, “The regretful agent: Heuristic-aided navigation through progress estimation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern*

Recognition, 2019, pp. 6732–6740.

[38] K. Nguyen, D. Dey, C. Brockett, and B. Dolan, “Vision-based navigation with language-based assistance via imitation learning with indirect intervention,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 527–12 537.

[39] M. Bellemare, S. Srinivasan, G. Ostrovski, T. Schaul, D. Saxton, and R. Munos, “Unifying count-based exploration and intrinsic motivation,” in *Advances in Neural Information Processing Systems*, 2016, pp. 1471–1479.

[40] J. Gao, M. Galley, L. Li *et al.*, “Neural approaches to conversational ai,” *Foundations and Trends® in Information Retrieval*, vol. 13, no. 2-3, pp. 127–298, 2019.

[41] R. Houthoofd, X. Chen, Y. Duan, J. Schulman, F. De Turck, and P. Abbeel, “Vime: Variational information maximizing exploration,” in *Advances in Neural Information Processing Systems*, 2016, pp. 1109–1117.

[42] G. Ostrovski, M. G. Bellemare, A. van den Oord, and R. Munos, “Count-based exploration with neural density models,” in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 2721–2730.

[43] H. Tang, R. Houthoofd, D. Foote, A. Stooke, O. X. Chen, Y. Duan, J. Schulman, F. DeTurck, and P. Abbeel, “# exploration: A study of count-based exploration for deep reinforcement learning,” in *Advances in Neural Information Processing Systems*, 2017, pp. 2753–2762.

[44] J. Schmidhuber, “Adaptive confidence and adaptive curiosity,” in *Institut für Informatik, Technische Universität München, Arcisstr. 21, 800 München 2*. Citeseer, 1991.

[45] A. L. Strehl and M. L. Littman, “An analysis of model-based interval estimation for markov decision processes,” *Journal of Computer and System Sciences*, vol. 74, no. 8, pp. 1309–1331, 2008.

Instruction: Turn around and exit the room to the right of the TV. Once out turn left and walk to the end of the hallway and then turn right. Walk down the hallway past the piano and then stop when you enter the next doorway and the floor is in a circle pattern.

Intrinsic Reward: 0.39 Result: Failure (error = 5.4m)

step 1 panorama view



step 2 panorama view



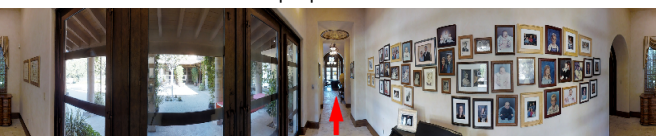
step 3 panorama view



step 4 panorama view



step 5 panorama view



step 6 panorama view



step 7 panorama view



step 8 panorama view



step 9 panorama view



180 -90 0 90 180
relative heading from the agent

Fig. 10: Failure of executing a relatively complicated instruction.

- [46] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell, "Curiosity-driven exploration by self-supervised prediction," in *International Conference on Machine Learning (ICML)*, vol. 2017, 2017.
- [47] J. Oh, Y. Guo, S. Singh, and H. Lee, "Self-imitation learning," in *Proceedings of the 35th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, J. Dy and A. Krause, Eds., vol. 80. Stockholmssan, Stockholm Sweden: PMLR, 10–15 Jul 2018, pp. 3878–3887.
- [48] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [49] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.
- [50] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine learning*, vol. 8, no. 3-4, pp. 229–256, 1992.
- [51] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Nießner, M. Savva, S. Song, A. Zeng, and Y. Zhang, "Matterport3d: Learning from rgb-d data in indoor environments," *arXiv preprint arXiv:1709.06158*, 2017.
- [52] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [53] J. Pennington, R. Socher, and C. Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 2014, pp. 1532–1543.
- [54] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [55] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6629–6638.
- [56] R. Paulus, C. Xiong, and R. Socher, "A deep reinforced model for abstractive summarization," *arXiv preprint arXiv:1705.04304*, 2017.
- [57] M. Ranzato, S. Chopra, M. Auli, and W. Zaremba, "Sequence level training with recurrent neural networks," *arXiv preprint arXiv:1511.06732*, 2015.
- [58] H. Tan, L. Yu, and M. Bansal, "Learning to navigate unseen environments: Back translation with environmental dropout," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019.
- [59] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

Xin Wang Xin Wang is a Ph.D. candidate at University of California, Santa Barbara. He received B.S. in computer science from Zhejiang University. His major research interests lie in compute vision, natural language processing, robotics, and machine learning, especially the intersection of vision and language. He organized the third workshop of Closing the Loop Between Vision and Language (CLVL) in ICCV 2019 and the first VATEX Challenge for Multilingual Video Captioning. He served as a session chair for the NLP session at AAI 2019 and frequently serves as the Program Committee Member for the top-tier CV/NLP/ML conferences such as CVPR, ICCV, ECCV, ACL, NAACL, EMNLP, AAI, NeurIPS etc. He is the recipient of the Best Student Paper Award at CVPR 2019.

Qiuyuan Huang Qiuyuan Huang is a senior researcher in the Deep Learning group at Microsoft Research, Redmond, WA. She obtained her Ph.D. degree from the University of Florida in 2017 and has worked as a postdoctoral researcher in the Deep Learning group at Microsoft Research for one year (2017-2018). Her current research interests are in the areas of deep learning and natural language processing in general, including neural-symbolic intelligence, self-supervised learning, reinforcement learning, generative adversarial networks, and multi-modal intelligence. She was a recipient of the Outstanding Young Researcher Award in Mathematics and Computer Science by Heidelberg Laureate Forum in 2015, the Best Poster Paper Award in ACM MobiCom 2016, the Rising Stars in EECSS by MIT in 2018, and the Best Student Paper Award in CVPR 2019.

Asli Celikyilmaz Asli Celikyilmaz is a Principal Researcher at Microsoft Research AI in Redmond and Adjunct Professor of the Paul G. Allen School of Computer Science & Engineering Department at the University of Washington. Her primary research interests are in the fields of Natural Language Processing, Machine Learning, Artificial Intelligence, with broader interests in Computer Vision. She is a co-recipient of the best student paper award at CVPR 2019, a recipient of the best paper award at Semantic Computing in 2010 and best paper award at NAFIPS in 2007. She has done her Post Doc research work in Computer Science at University of California, Berkeley, and received her Ph.D. and Ma.SC. in Information Science at University of Toronto, in Canada.

Jianfeng Gao Jianfeng Gao is Partner Research Manager at Microsoft Research AI, Redmond. IEEE fellow. He leads the development of AI systems for machine reading comprehension (MRC), question answering (QA), social bots, goal-oriented dialogue, and business applications. From 2014 to 2017, he was Partner Research Manager at Deep Learning Technology Center at Microsoft Research, Redmond, where he was leading the research on deep learning for text and image processing. From 2006 to 2014, he was Principal Researcher at Natural Language Processing Group at Microsoft Research, Redmond, where he worked on Web search, query understanding and reformulation, ads prediction, and statistical machine translation. From 2005 to 2006, he was a Research Lead in Natural Interactive Services Division at Microsoft, where he worked on Project X, an effort of developing natural user interface for Windows. From 2000 to 2005, he was Research Lead in Natural Language Computing Group at Microsoft Research Asia, where he and his colleagues developed the first Chinese speech recognition system released with Microsoft Office, the Chinese/Japanese Input Method Editors (IME) which were the leading products in the market, and the natural language platform for Microsoft Windows.

Dinghan Shen Dinghan Shen is Ph.D. student from the Electrical and Computer Engineering Department at Duke University, advised by Prof. Lawrence Carin. Prior to joining Duke, he received the bachelor degree from Peking University in 2015. He is a co-recipient of the Best Paper Honorable Mention award at ACL 2018.

Yuan-Fang Wang Yuan-Fang Wang is a Professor in the Department of Computer Science at the University of California at Santa Barbara. He received his bachelor's degree in electrical engineering from National Taiwan University and his master's and Ph.D. degrees in electrical and computer engineering from the University of Texas at Austin. Dr. Wang is also the founder of Visualsize Inc. that develops 3D computer vision products based on his federally-sponsored research and a co-founder of Proximex Corp. in physical security, which was recently acquired by ADT Security. His research activities center on computer vision, medical image analysis, computer graphics, bioinformatics, and digital image and video libraries. He has published over 120 peer-reviewed papers and holds two patents. He has received funding support from many federal agencies including NSF, NASA, DARPA, US Army, US Navy, and also from many private industries. He served as a consultant for a number of companies and was a faculty consultant at LG Electronics Research Center of America in summer 1998. Dr. Wang was the program co-chair of the 1998 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. He is an Associate Editor of IEEE Transactions on Pattern Analysis and Machine Intelligence from 1998 to 2002, the Pattern Recognition Journal from 2000 to 2007, and Journal of Ambient Intelligence and Smart Environment from 2007 to 2009.

William Wang William Wang is an Assistant Professor in the Department of Computer Science at the University of California, Santa Barbara. He directs UCSB's Natural Language Processing Group, and he is also the Director of UCSB's Responsible Machine Learning Center. He received his PhD from School of Computer Science, Carnegie Mellon University. He has broad interests in machine learning approaches to language and vision. He has published more than 80 papers at leading NLP/AI/CV/ML conferences and journals, and received best paper awards (or nominations) at ASRU 2013, CIKM 2013, EMNLP 2015, and CVPR 2019, a DARPA Young Faculty Award (Class of 2018), a Google Faculty Research Award (2018), two IBM Faculty Awards in 2017 and 2018, a Facebook Research Award in 2018, an Adobe Research Award in 2018. He frequently serves as an Area Chair for NAACL, ACL, EMNLP, and AAAI.

Lei Zhang Lei Zhang (Senior Member, IEEE) received his Ph.D. degree in computer science from Tsinghua University in 2001. He is a principal researcher and research manager in Microsoft, working on computer vision and machine learning. Prior to his current post, he was a senior researcher at Microsoft Research Asia. He is interested in image understanding, visual pattern recognition, and machine learning, and holds 50 U.S. patents in these fields. He has also served as program area chairs or committee members for many related conferences.