

REVERIE: Remote Embodied Visual Referring Expression in Real Indoor Environments

Yuankai Qi¹, Qi Wu^{*1}, Peter Anderson², Xin Wang³, William Yang Wang³, Chunhua Shen¹, and Anton van den Hengel¹

¹Australia Centre for Robotic Vision, The University of Adelaide

²Georgia Institute of Technology

³University of California, Santa Barbara

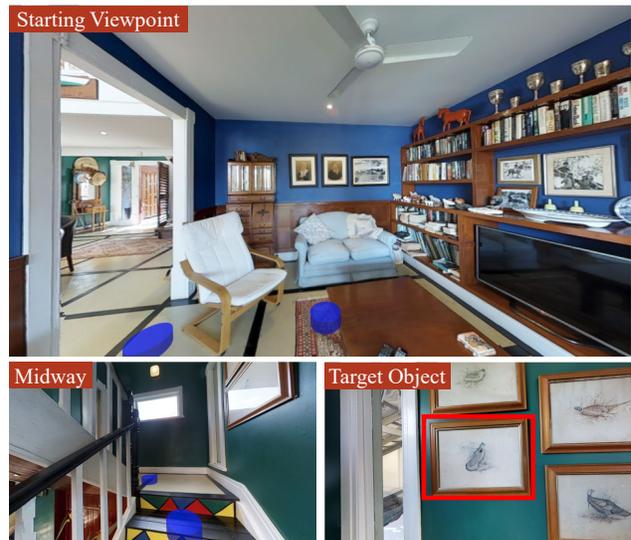
Abstract

One of the long-term challenges of robotics is to enable robots to interact with humans in the visual world via natural language, as humans are visual animals that communicate through language. Overcoming this challenge requires the ability to perform a wide variety of complex tasks in response to multifarious instructions from humans. In the hope that it might drive progress towards more flexible and powerful human interactions with robots, we propose a dataset of varied and complex robot tasks, described in natural language, in terms of objects visible in a large set of real images. Given an instruction, success requires navigating through a previously-unseen environment to identify an object. This represents a practical challenge, but one that closely reflects one of the core visual problems in robotics. Several state-of-the-art vision-and-language navigation, and referring-expression models are tested to verify the difficulty of this new task, but none of them show promising results because there are many fundamental differences between our task and previous ones. A novel Interactive Navigator-Pointer model is also proposed that provides a strong baseline on the task. The proposed model especially achieves the best performance on the unseen test split, but still leaves substantial room for improvement compared to the human performance.

1. Introduction

You can ask a 10-year-old child to bring you a cushion, and there is a good chance that they will succeed (even in an unfamiliar environment), while the probability that a robot will achieve the same task is significantly lower. Children

*corresponding author; qi.wu01@adelaide.edu.au



Instruction: Bring me the bottom picture that is next to the top of stairs on level one.

Figure 1. REVERIE task: an agent is given a natural language instruction referring to a remote object (here in the red bounding box) in a photo-realistic 3D environment. The agent must navigate to an appropriate location and identify the object from multiple distracting candidates. The blue discs indicate nearby navigable viewpoints provided by the simulator.

have a wealth of knowledge learned from similar environments that they can easily apply to such tasks in an unfamiliar environment, including the facts that cushions generally inhabit couches, that couches inhabit lounge rooms, and that lounge rooms are often connected to the rest of a building through hallways. Children are also able to interpret natural language instructions and associate them with the visual world. However, the fact that robots currently lack these capabilities dramatically limits their domain of application.

-
1. Fold the towel in the bathroom with the fishing theme.
 2. Enter the bedroom with the letter E over the bed and turn the light switch off.
 3. Go to the blue family room and bring the framed picture of a person on a horse at the top left corner above the TV.
 4. Push in the bar chair, in the kitchen, by the oven.
 5. Windex the mirror above the sink, in the bedroom with the large, stone fireplace.
 6. Could you please dust the light above the toilet in the bathroom that is near the entry way?
 7. At the top of the stairs, the first set of potted flowers in front of the stairs need to be dusted off.
 8. To the right at the end of the hall, where the large blue table foot stool is, there is a mirror that needs to be wiped.
 9. Go to the hallway area where there are three pictures side by side and get me the one on the right.
 10. There is a bottle in the office alcove next to the piano. It is on the shelf above the sink on the extreme right. Please bring it here.
-

Table 1. Indicative instruction examples from the REVERIE dataset illustrating various interesting linguistic phenomena such as dangling modifiers (e.g. 1), spatial relations (e.g. 3), imperatives (e.g. 9), co-references (e.g. 10), etc. Note that the agent in our task is required to identify the referent object, but is not required to complete any manipulation tasks (such as folding the towel).

Therefore, to equip robots with such abilities and to advance real-world vision-and-language research, we introduce a new problem, which we refer to as *Remote Embodied Visual referring Expression in Real Indoor Environments* — *REVERIE*. An example of the REVERIE task is illustrated in Fig. 1. A robot spawns at a starting location and is given a natural language instruction that refers to a remote target object at another location within the same building. To carry out the task, the agent is required to navigate closer to the object and return a bounding box encompassing the target object specified by the instruction. It demands the robot to infer the probable location of the object using knowledge of the environments, and explicitly identify the object according to the language instruction.

In distinction to other embodied tasks such as Vision-and-Language Navigation (VLN) [1] and Embodied Question Answering (EQA) [7], REVERIE evaluates the success based on explicit object grounding rather than the point navigation in VLN or the question answering in EQA. This more clearly reflects the necessity of robots’ capability of natural language understanding, visual navigation, and object grounding. More importantly, the concise instructions in REVERIE represent more practical tasks that humans would ask a robot to perform (see Tab. 1). Those high-level instructions are fundamentally different from the fine-grained visuomotor instructions in VLN, and would empower high-level reasoning and real-world applications. Moreover, compared to the task of Referring Expression (RefExp) [14, 22] that selects the desired object from a single image, REVERIE is far more challenging in the sense that the target object is not visible in the initial view and needs to be discovered by actively navigating in the environment. Hence, in REVERIE, there are at least an order of magnitude more object candidates to choose from than RefExp.

We build the REVERIE dataset upon the Matterport3D Simulator [1, 4], which provides panoramas of all the navigable locations and the connectivity graph in a building. To provide object-level information of the environ-

ments, we have extended the simulator to incorporate object annotations, including labels and bounding boxes from Chang *et al.* [4]. The extended simulator can project bounding boxes onto images of different viewpoints and angles, thus able to accommodate evaluation on every possible location. The REVERIE dataset comprises 10,567 panoramas within 90 buildings containing 4,140 target objects, and 21,702 crowd-sourced instructions with an average length of 18 words. Tab. 1 demonstrates sample instructions from the dataset, which illustrate various linguistic phenomena, including spatial relations, multiple long and dangling modifiers, and coreferences, etc.

We investigate the difficulty of the REVERIE task by directly combining state-of-the-art (SoTA) navigation methods and referring expression methods, and none of them shows promising results. We then propose an Interactive Navigator-Pointer model serving as a strong baseline for the REVERIE task. We also provide the human performance of the REVERIE task on the test set to quantify the machine-human gap.

In summary, our main contributions are:

1. A new embodied vision-and-language problem, Remote Embodied Visual referring Expressions in Real 3D Indoor Environments (REVERIE), where given a natural language instruction that represents a practical task to perform, an agent must navigate and identify a remote object in real indoor environments.
2. The first benchmark dataset for the REVERIE task, which contains large-scale human-annotated instructions and extends the Matterport3D Simulator [1] with additional object annotations.
3. A novel interactive navigator-pointer model that provides strong baselines for the REVERIE dataset under several evaluation metrics.

2. Related Work

Referring Expression Comprehension. The referring expression comprehension task requires an agent to localise

Dataset	Language Context				Visual Context			Goal
	Human	Main Content	Unamb	Guidance Level	BBox	Real-world	Temporal	
EQA [7], IQA [10]	✗	QA-pair	✓	–	✗	✗	Dynamic	QA
MARCO [21], DRIF [3]	✓	Nav-Instruction	✓	Detailed	✗	✗	Dynamic	Navigation
R2R [1]	✓	Nav-Instruction	✓	Detailed	✗	✓	Dynamic	Navigation
TouchDown [5]	✓	Nav-Instruction	✓	Detailed	✗	✓	Dynamic	Navigation
VLNA [24], HANNA[23]	✗	Nav-Dialog	✗	High	✗	✓	Dynamic	Find Object
TtW [8]	✓	Nav-Dialog	✓	High	✗	✓	Dynamic	Navigation
CVDN [25]	✓	Nav-Dialog	✗	High	✗	✓	Dynamic	Find Room
ReferCOCO [30]	✓	RefExp	✓	–	✓	✓	Static	Localise Object
REVERIE	✓	Remote RefExp	✓	High	✓	✓	Dynamic	Localise Remote Object

Table 2. Compared to existing datasets involving embodied vision and language tasks. Symbol instruction: ‘QA’: ‘Question-Answer’, ‘Unamb’: ‘Unambiguous’, ‘BBox’: ‘Bounding Box’, ‘Dynamic’/‘Static’: visual context temporally changed or not.

an object in an image given a natural language expression. Recent work casts this task as looking for the object that can generate its paired expressions [13, 18, 30] or jointly embedding the image and expression for matching estimation [6, 12, 16, 29]. Yu *et al.* [30] propose to compute the appearance difference of the same category objects to enhance the visual features for expression generation. Instead of treating each expression as a unit, [29] learns to decompose an expression into appearance, location, and object relationship three components.

Different from referring expression, REVERIE introduces three new challenges: i) The referred object is not visible in the initial scene and only can be accessed after navigating to a closed location. ii) In contrast to previous referring expression tasks that select the target object from a single image, object candidates in REVERIE come from panoramas of all the possible viewpoints. iii) The objects in referring expression are normally captured from the front view, while in our setting, the visual appearances of objects may vary largely due to different observation angles and viewpoints.

Vision-and-Language Navigation. Vision-and-language navigation (VLN) is the task where an agent is to navigate to a goal location in a 3D simulator given detailed natural language instructions such as “Turn right and go through the kitchen. Walk past the couches on the right and into the hallway on the left. Go straight until you get to a room that is to the left of the pictures of children on the wall. Turn left and go into the bathroom. Wait near the sink.” [1]. A range of VLN methods [9, 15, 19, 20, 27, 28] have been proposed to address this VLN task.

Although the proposed REVERIE task also requires an agent to navigate to a goal location, it differs from existing VLN tasks in two important aspects: i) The challenge is much more closely related to the overarching objective of enabling natural language robot tasking because the goal is to localise a target object specified in an instruction, not just a location. This removes the artificial constraint that the instruction is restricted to solely to navigation, and reflects the reality of the fact that most objects can be seen from multiple viewpoints. ii) Our collected navigation in-

structions are semantic-level commands which better reflect the way humans communicate. They are thus closer to ‘the cold tap in the first bedroom on level two’ rather than step by step navigation instructions such as ‘go to the top of the stairs then turn left and walk along the hallway and stop at the first bedroom on your right’.

The most closely related challenge to that proposed here is that addressed in [23, 24, 25] whereby an agent must identify an object by requesting and interpreting natural language assistance. The instructions are of the form ‘Find a mug’, and the assumption is that there is an oracle following the agent around the environment willing to provide natural language assistance. The question is then whether the agent can effectively exploit the assistance provided by the omniscient oracle. REVERIE, in contrast, evaluates whether the agent can carry out a natural-language instruction alone. Another closely related work is TOUCHDOWN [5], that requires an agent to find a location in an urban outdoor environment on the basis of detailed navigation instructions.

Embodied Question Answering. Embodied question answering (EQA) [7] requires an agent to answer a question about an object or a room in a synthetic environment. Gordon *et al.* [10] introduce an interactive version of the EQA task, where the agent may need to interact with the environment/objects to correctly answer questions. Our REVERIE task differs from previous works that only output a simple answer or a series of actions, as we ask the agent to output a bounding box around a target object. This is a more challenging but realistic setting because if we want a robot to carry out a task that relates to an object, we need its precise location. Tab. 2 displays the difference between our task and other related embodied vision-language tasks.

3. The REVERIE Dataset

We now describe the REVERIE task and dataset, including the task definition, evaluation metrics, simulator, data collection policy, and analysis of the collected instructions.

3.1. The REVERIE Task

As shown in Fig. 1, our REVERIE task requires an intelligent agent to correctly localise a remote target object (can not be observed at starting location) specified by a concise high-level natural language instruction (see samples in Tab. 1). It is worth noting that our instructions are much closer to practical scenarios in daily life than the detailed instructions in VLN [1], because the latter ones are so complex and long that are unrealistic for human to command robots. Since the target object is in a different room from the starting one, the agent needs first to navigate to the goal location.

Formally, at the beginning of each episode, the agent is given as input a high-level natural language instruction $\mathcal{X} = \langle w_1, w_2, \dots, w_L \rangle$, where L is the length of the instruction and w_i is a single word token. Following the common practice in VLN, the agent has access to surrounding panoramic images $\mathcal{V}_0 = \{v_{0,k}, k \in 1, \dots, 36\}$ and navigable viewpoints at the current location, where $v_{0,k}$ is determined by the agent’s states comprising a tuple of 3D position, heading and elevation $s_{0,k} = \langle p_0, \phi_{0,k}, \theta_{0,k} \rangle$ (3 elevation and 12 heading angles are used). Then the agent needs to make a sequence of actions $\langle a_0, \dots, a_T \rangle$ to reach the goal location, where each action is choosing one of the navigable viewpoints or choosing the current viewpoint which means to stop. The action can also be a ‘detecting’ action that outputs the target object bounding-box referred by the instruction. It is worth noting that the agent can attempt to localise the target at any step, which is totally up to algorithm design. But we only allow the agent output once in each episode, which means the agent only can guess the answer once in a single run. If the agent ‘thinks’ it has localised the target object and decides to output it, it is required to output a bounding box or choose from several candidates provided by the simulator. A bounding box is denoted as $\langle b_x, b_y, b_w, b_h \rangle$, where b_x and b_y are the coordinate of the left-top point, b_w and b_h denote the width and height of the bounding box, respectively. The episode ends after the agent outputs the target bounding box.

3.2. Evaluation Metrics

The performance of a model is mainly measured by REVERIE success rate, which is the number of successful tasks over the total number of tasks. A task is considered successful if it selects the correct bounding box of the target object from a set of candidates (or the IoU between the predicted bounding box and the ground-truth bounding box ≥ 0.5 , when candidate objects bounding boxes are not given). Because the target object can be observed at different viewpoints or camera views, we treat it as a success as long as the agent can identify the target within 3 meters, regardless of from different viewpoints or views. We also measure the navigation performance with four kinds

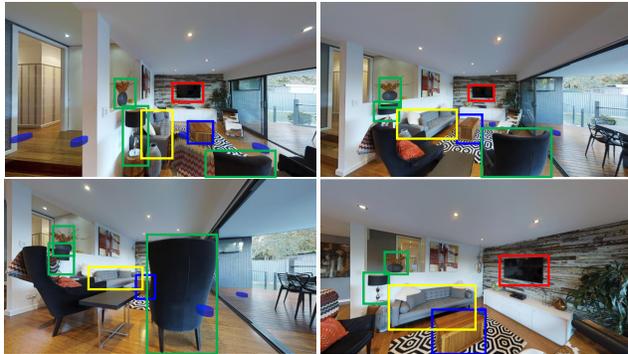


Figure 2. Object bounding boxes (BBox) in our simulator. The BBox size and aspect ratio of the same object may change after the agent moves to another viewpoint or changes its camera view.

of metrics, including success rate, oracle success rate, success rate weighted by path length (SPL), and path length (in meters) [1]. Please note that in our task, a navigation is considered successful only when the agent stops at a location within 3 meters from the target object. More details of evaluation metrics can be found in supplementary materials.

3.3. The REVERIE Simulator

Our simulator is based on the Matterport3D Simulator [1], a large-scale interactive environment constructed from the Matterport3D dataset [4]. In the simulator, an embodied agent is able to virtually ‘move’ throughout each building by iteratively selecting adjacent nodes from the graph of panoramic viewpoints and adjusting the camera pose at each viewpoint. At each viewpoint, it returns a rendered colour image that captures the current view, as shown in Fig. 1.

Adding Object-level Annotations. Object bounding boxes are needed in our proposed task, which are either provided as object hypotheses or used to assess the agent’s ability to localise the object that is referred to by a natural expression. The main challenge of adding the object bounding boxes into the simulator is that we need to handle the changes in visibility and coordinate of 2D bounding boxes as the camera moves or rotates.

To address these issues, we calculate the overlap between bounding boxes and object depth in each view. If a bounding box is fully covered by another one and it has a larger depth, we treat it as an occluded case. Specifically, for each building the Matterport3D dataset provides all the objects appearing in it with centre point position $\mathbf{c} = \langle c_x, c_y, c_z \rangle$, three axis directions $\mathbf{d}_i = \langle d_i^x, d_i^y, d_i^z \rangle, i \in \{1, 2, 3\}$, and three radii r_i , one for each axis direction. To correctly render objects in the web simulator, we first calculate the eight vertexes using \mathbf{c}, \mathbf{d}_i and r_i . Then these vertexes are projected into the camera space by the camera pose provided by Matterport3D dataset. Both C++ and web simulators will be released with the code. Fig. 2 presents an exam-

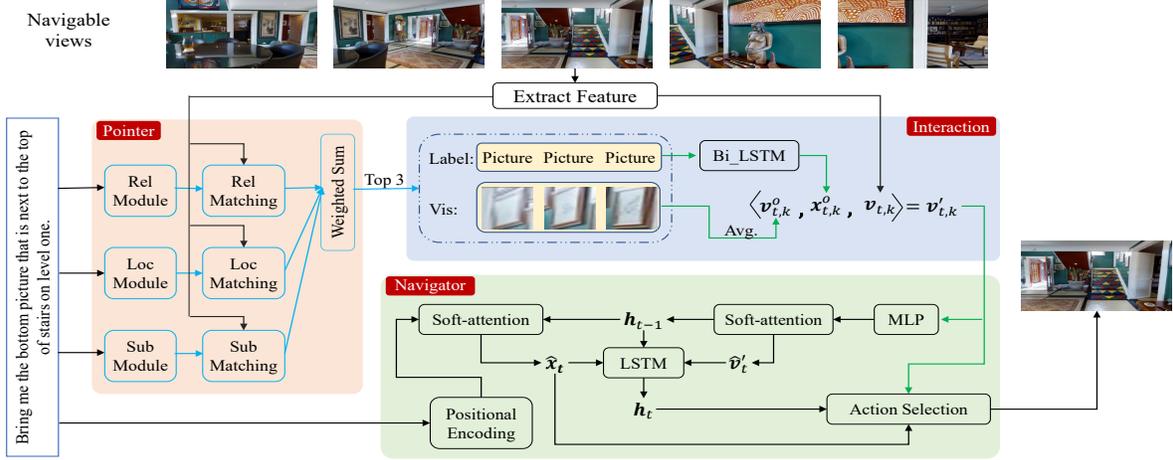


Figure 5. Our Interactive Navigator-Pointer Model

responses for sending the referring expression comprehension information obtained from the pointer to the navigator to guild it to make more accurate action prediction.

4.1. The Navigator Module

The backbone of our navigator module is a ‘short’ version of FAST [15], which uses a sequence-to-sequence LSTM architecture with an attention mechanism and a backtracking mechanism to increase the action accuracy. Specifically, let $\mathbf{X} \in \mathbb{R}^{L \times 512}$ denote instruction features obtained from \mathcal{X} by an LSTM, and $\mathbf{V}' = [v'_{t,1}; \dots; v'_{t,K}] \in \mathbb{R}^{K \times 4736}$ denote updated visual features obtained by our interactive module (Sec. 4.3) for panoramic images \mathcal{V}_t at step t . FAST-short learns the local logit l_t signal, which is calculated by a visual and textual co-grounding model adopted from [20]. First, grounded text $\hat{\mathbf{x}}_t = \alpha_t^\top \mathbf{X}$ and grounded visual $\hat{\mathbf{v}}'_t = \beta_t^\top \mathbf{V}'$ are learned by

$$\alpha_t = \text{softmax}(PE(\mathbf{X})(\mathbf{W}_x \mathbf{h}_{t-1})) \quad (1)$$

$$\beta_t = \text{softmax}(g(\mathbf{V}')(\mathbf{W}_v \mathbf{h}_{t-1})) \quad (2)$$

where $\alpha_t \in \mathbb{R}^{L \times 1}$ is textual attention weight, $\beta_t \in \mathbb{R}^{K \times 1}$ is visual attention weight, \mathbf{W}_x and \mathbf{W}_v are learnable parameters, $PE(\cdot)$ is the positional encoding [26] that captures the relative position between each word within an instruction, $g(\cdot)$ is a one-layer Multi-Layer Perceptron (MLP), $\mathbf{h}_{t-1} \in \mathbb{R}^{512 \times 1}$ is previous encoder context. The new context is updated by an LSTM taking as input the newly grounded text and visual features as well as previous selected action

$$(\mathbf{h}_t, \mathbf{c}_t) = \text{LSTM}([\hat{\mathbf{x}}_t, \hat{\mathbf{v}}'_t, \mathbf{a}_{t-1}], (\mathbf{h}_{t-1}, \mathbf{c}_{t-1})). \quad (3)$$

Then the logit l_t can be computed via an inner-product between each candidate’s encoded context and instruction by

$$l_{t,k} = (\mathbf{W}_a [\mathbf{h}_t, \hat{\mathbf{x}}_t])^\top g(v'_{t,k}) \quad (4)$$

where \mathbf{W}_a is a learnable parameter matrix.

Based on logit l_t , FAST-short maintains one candidate queue and one ending queue. All navigable viewpoints (including the current viewpoint) at the current location are pushed into the candidate queue, but only the viewpoint with the largest accumulated logit $\sum_{\tau=0}^t l_\tau$ is popped out as the selected next step. Each passed viewpoint is pushed into the ending queue. One episode ends if the current viewpoint is selected or the candidate queue is empty or the maximum step is reached. Finally, the viewpoint with the largest accumulated logits is chosen as the actual stop location.

4.2. The Pointer Module

We use MAttNet [29] as our pointer because of its good generalisation ability. It decomposes an expression into three modular components related to subject appearance, location, and relationship to other objects via the attention mechanism $q^m = \sum_{j=1}^L a_{m,j} e_j$, where $m \in \{\text{subj}, \text{loc}, \text{rel}\}$, e_j is the embedding of each word in the expression/instruction \mathcal{X} . $a_{m,j}$ is the attention on each word for each module.

Then three kinds of matching scores $S(o_i | q^m)$ are computed for each object o_i conditioned on each modular phrase embedding q^m . Specifically, $S(o_i | q^{\text{subj}}) = F(\tilde{v}_i^{\text{subj}}, q^{\text{subj}})$, $S(o_i | q^{\text{loc}}) = F(\tilde{l}_i^{\text{loc}}, q^{\text{loc}})$, and $S(o_i | q^{\text{rel}}) = \max_{j \neq i} F(\tilde{v}_{ij}^{\text{rel}}, q^{\text{rel}})$, where $F(\cdot)$ is a two-layer MLP, $\tilde{v}_i^{\text{subj}}$ is a ‘in-box’ attended feature for each object using a 14×14 grid. \tilde{l}_i^{loc} is the location representation of object o_i obtained by a fully-connected layer taking as input the relative position offset and area ratio to its up to five surrounding objects of the same category. $\tilde{v}_{ij}^{\text{rel}}$ is the visual representation of the surrounding object o_j regardless of categories.

The final matching score of object o_i and the instruction \mathcal{X} is a weighted sum:

$$S = \sum S(o_i|q^m)w_m \quad (5)$$

where $w_m = \text{softmax}(W_m^L[h_0, h_L] + b_m)$.

4.3. The Interaction Module

Intuitively, we want the Navigator and Pointer to interact with each other so that both navigation and referring expression accuracy can be improved. For example, the navigator can use the visual grounding information to decide when and where to stop, and the pointer accuracy can be improved if the navigator can reach the correct target location. To this end, we propose an interaction module that can plug the pointer’s output into the navigator. Specifically, we first perform referring expression comprehension using the above pointer module to select the top-3 matching objects in each candidate view. Then we use a trainable bi-direction LSTM to encode the category labels of these selected objects $\mathcal{X}_O = \{\text{Label}_{i \in \text{top}3}\}$

$$\mathbf{x}_{t,k}^o = \text{bi-LSTM}(\mathcal{X}_O) \quad (6)$$

as the textual representation for the k -th candidate viewpoint. In addition, the averaged output of ResNet FC7 layer of these object regions is used as the visual representation $\mathbf{v}_{t,k}^o$. Finally, we update the candidate viewpoint feature by concatenation

$$\mathbf{v}'_{t,k} = [\mathbf{v}_{t,k}, \mathbf{x}_{t,k}^o, \mathbf{v}_{t,k}^o] \quad (7)$$

which is send to the navigator (see Equ. 2 and 4). The pointer in such an interaction serves as hard attention for each candidate viewpoint, which highlights the most target-related objects for the navigator to take into account.

4.4. Loss Functions

Our final loss consists of two parts, the navigation loss L_{nav} and referring expression loss L_{exp} . The L_{nav} is a cross-entropy loss for action selection and a mean squared error loss for progress monitor:

$$L_{nav} = -\lambda_1 \sum_{t=1}^T y_t^a \log(l_{t,k}) - (1 - \lambda_1) \sum_{t=1}^T (y_t^{pm} - p_t^{pm})^2 \quad (8)$$

where y_t^a is the ground truth action at step t , $\lambda_1 = 0.5$ is the weight balancing the two loss, $y_t^{pm} \in [0, 1]$ is the normalised distance in units of length from the current viewpoint to the goal, p_t^{pm} is the predicted progress.

The referring expression loss L_{exp} is a ranking loss:

$$L_{exp} = \sum_i [\lambda_2 \max(0, \delta + S(o_i|r_j) - S(o_i|r_i)) + \lambda_3 \max(0, \delta + S(o_k|r_i) - S(o_i|r_i))] \quad (9)$$

where $\lambda_2 = 1.0$, $\lambda_3 = 1.0$, (o_i, r_i) is a positive (object, expression) pair, (o_i, r_j) and (o_k, r_i) are negative (object,

expression) pairs, δ is the distance margin between positive and negative pairs. All of the losses are summarised together:

$$L = L_{nav} + \lambda_4 L_{exp} \quad (10)$$

to train our Interactive Navigator-Pointer model. We set λ_4 to 1.0 by default.

5. Experiments

In this section, we first present the training details of the interactive navigator–pointer model. Then, we provide extensive evaluation and analysis.

5.1. Implementation Details

The simulator image resolution is set to 640×480 pixels with a vertical field of view of 60 degrees. For each instruction in the train split, images and object bounding boxes at the goal viewpoint (for the views where the target object is visible) are organised following the format as in MAttNet for pointer training. With the trained pointer, assistant object information is provided for the navigator as described in Section 4.3 to train the navigator. The code and dataset will be released.

5.2. REVERIE Experimental Results

We first evaluate several baseline models and state-of-the-art (SoTA) navigation models, combined with the MAttNet, *i.e.*, the pointer module. After the navigation models decide to stop, the pointer module is used to predict the target object. In addition, we also test human performance using an interactive web interface (see details in the supplementary).

Below is a brief introduction of the evaluated baseline and state-of-the-art models. There are four baseline models, which are:

- **Random** exploits the characteristics of the dataset by randomly choosing a path with random steps (maximum 10) and then randomly choose an object as the predicted target.
- **Shortest** always follows the shortest path to the goal.
- **R2R-TF and R2R-SF** [1] are the first batch of navigation baselines, which trains a basic LSTM [11] with attention mechanism [2]. The difference between R2R-TF and R2R-SF is that R2R-TF is trained with the ground truth action at each step (Teacher-Forcing, TF) while R2R-SF adopts an action sampled from the predicted probability over its action space (Student-Forcing, SF).

The evaluated four SoTA navigation models are:

- **SelfMonitor** [20] uses a visual-textual co-grounding module to highlight the instruction for the next action and a progress monitor to reflect the progress.

Methods	Val Seen					Val UnSeen					Test (Unseen)				
	Navigation Acc.				REVERIE Succ.	Navigation Acc.				REVERIE Succ.	Navigation Acc.				REVERIE Succ.
	Succ.	OSucc.	SPL	Length		Succ.	OSucc.	SPL	Length		Succ.	OSucc.	SPL	Length	
Random	2.74	8.92	1.91	11.99	1.97	1.76	11.93	1.01	10.76	0.96	2.30	8.88	1.44	10.34	1.18
Shortest	100	100	100	10.46	68.45	100	100	100	9.47	56.63	100	100	100	9.39	48.98
R2R-TF [1]	7.38	10.75	6.40	11.19	4.22	3.21	4.94	2.80	11.22	2.02	3.94	6.40	3.30	10.07	2.32
R2R-SF [1]	29.59	35.70	24.01	12.88	18.97	4.20	8.07	2.84	11.07	2.16	3.99	6.88	3.09	10.89	2.00
RCM [27]	23.33	29.44	21.82	10.70	16.23	9.29	14.23	6.97	11.98	4.89	7.84	11.68	6.67	10.60	3.67
SelfMonitor [20]	41.25	43.29	39.61	7.54	30.07	8.15	11.28	6.44	9.07	4.54	5.80	8.39	4.53	9.23	3.10
FAST-Short [15]	45.12	49.68	40.18	13.22	31.41	10.08	20.48	6.17	29.70	6.24	14.18	23.36	8.74	30.69	7.07
FAST-Lan-Only	8.36	23.61	3.67	49.43	5.97	9.37	29.76	3.65	45.03	5.00	8.15	28.45	2.88	46.19	4.34
Ours	50.53	55.17	45.50	16.35	31.97	14.40	28.20	7.19	45.28	7.84	19.88	30.63	11.61	39.05	11.28
Human	-	-	-	-	-	-	-	-	-	-	81.51	86.83	53.66	21.18	77.84

Table 3. REVERIE success rate achieved by combining state-of-the-art navigation methods with the RefExp method MAttNet [29].

- **RCM [27]** employs reinforcement learning to encourage global matching between instructions and trajectories, and performs cross-model grounding.
- **FAST-Short [15]** introduces backtracking into Self-Monitor.
- **FAST-Lan-Only** employs above FAST-Short model but we only input the language instruction without any visual input. This model is used to check whether our task/dataset has a bias on language input.

Results. The detailed experimental results are presented in Tab. 3, of which the first four rows are results for baselines, the following four rows are for SoTA methods, and the last two rows are for our model and human performance.

According to the baseline section in Tab. 3, the Random model only achieves a REVERIE success around 1%, which indicates the REVERIE task has a huge solution space. The sequence-to-sequence baselines R2R-TF and R2R-SF [1] achieve good results on the Val-Seen split but decrease a lot on the unseen splits. Student-Forcing is generally better than Teacher-Forcing. The Shortest model achieves the perfect performance because the ground-truth path to the goal is directly given.

In the second part, the best REVERIE success rate is achieved by the combination of SoTA navigation (FAST) and referring expression (MAttNet) models. However, the REVERIE success rate is only 7.07% on the test split, falling far behind human performance 77.84%. The navigation-only accuracy of these SoTA navigation models indicates the challenge of our navigation task. Nearly 30% drops on the unseen splits are observed compared to the performance on previous R2R [1]. For example, the navigation SPL score of FAST-Short [15] on Val-UnSeen split drops from 43% on the R2R dataset to 6.17% on REVERIE.

To test whether our dataset has strong language bias, *i.e.*, whether a language-only model can achieve good performance, we implement a FAST-Lan-Only model with only instructions as its input. We observe a big drop on both seen and unseen splits, which suggests jointly considering language and visual information is necessary to our task.

Overall, these SoTA model results show that a simple combination of SoTA navigation and referring expression methods would not necessarily lead to the best performance

	Val Seen	Val UnSeen	Test
Baseline	30.69	18.63	16.18
MAttNet [29]	68.45	56.63	48.98
CM-Erase [17]	65.21	54.02	45.25
Human	-	-	90.76

Table 4. Referring expression comprehension success rate (%) at the ground truth goal viewpoint of our REVERIE dataset.

as failures from either the navigator or the pointer would decrease the overall success. In this paper, we make the first attempt to enable the navigator and pointer to work interactively as described in Sec. 4.3. The results in Tab. 3 show that our interactive model achieves consistently better results than non-interactive ones. The FAST-Short can be treated as our ablated model that without our proposed interaction module. Our final model achieves a gain of 4.2% on the test split.

Referring Expression-Only. We also report the Referring Expression-Only performance. In this setting, the agent is placed at the ground-truth target location, and then referring expression comprehension models are tested.

We test the SoTA models such as MAttNet [29] and CM-Erase [17] as well as a simple CNN-RNN baseline model with triplet ranking loss. Tab. 4 presents the results with human performance. It shows that the SoTA models achieve around 50% accuracy on the test split¹ which are far more better than the results when jointly considering the navigation and referring expression shown in Tab. 3. Even though, there is still a 40% gap to human performance, suggesting that our proposed REVERIE task is challenging.

6. Conclusion

Enable human-robots collaboration is a long-term goal. In this paper, we make a step further towards this goal by proposing a Remote Embodied Visual referring Expressions in Real Indoor Environments (REVERIE) task and dataset. The REVERIE is the first one to evaluate the capability of an agent to follow high-level natural languages instructions

¹These SoTA models achieve 80% accuracy on ReferCOCO [30], a golden benchmark for referring expression.

to navigate and identify the target object in previously unseen real images rendered buildings. We investigate several baselines and an interactive Navigator-Pointer agent model, of which the performance consistently demonstrate the significant necessity of further researches in this field.

We reach three main conclusions: First, REVERIE is interesting because existing vision and language methods can be easily plugged in. Second, the challenge of understanding and executing high-level instructions is significant. Finally, the combination of instruction navigation and referring expression comprehension is a challenging task due to the large gap to human performance.

7. Acknowledgments

We thank Sam Bahrami and Phil Roberts for their great help in the building of the REVERIE dataset.

References

- [1] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian D. Reid, Stephen Gould, and Anton van den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *CVPR*, pages 3674–3683, 2018. [2](#), [3](#), [4](#), [5](#), [7](#), [8](#)
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *ICLR*, 2015. [7](#)
- [3] Valts Blukis, Dipendra Kumar Misra, Ross A. Knepper, and Yoav Artzi. Mapping navigation instructions to continuous control actions with position-visitation prediction. In *2nd Annual Conference on Robot Learning, CoRL 2018, Zürich, Switzerland, 29-31 October 2018, Proceedings*, pages 505–518, 2018. [3](#)
- [4] Angel X. Chang, Angela Dai, Thomas A. Funkhouser, Maciej Halber, Matthias Nießner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from RGB-D data in indoor environments. In *2017 International Conference on 3D Vision, 3DV 2017, Qingdao, China, October 10-12, 2017*, pages 667–676, 2017. [2](#), [4](#)
- [5] Howard Chen, Alane Suhr, Dipendra Misra, Noah Snavely, and Yoav Artzi. Touchdown: Natural language navigation and spatial reasoning in visual street environments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12538–12547, 2019. [3](#)
- [6] Kan Chen, Rama Kovvuri, and Ram Nevatia. Query-guided regression network with context policy for phrase grounding. In *ICCV*, pages 824–832, 2017. [3](#)
- [7] Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied question answering. In *CVPR*, pages 1–10, 2018. [2](#), [3](#)
- [8] Harm de Vries, Kurt Shuster, Dhruv Batra, Devi Parikh, Jason Weston, and Douwe Kiela. Talk the walk: Navigating new york city through grounded dialogue. *CoRR*, abs/1807.03367, 2018. [3](#)
- [9] Daniel Fried, Ronghang Hu, Volkan Cirik, Anna Rohrbach, Jacob Andreas, Louis-Philippe Morency, Taylor Berg-Kirkpatrick, Kate Saenko, Dan Klein, and Trevor Darrell. Speaker-follower models for vision-and-language navigation. In *NeurIPS*, pages 3318–3329, 2018. [3](#)
- [10] Daniel Gordon, Aniruddha Kembhavi, Mohammad Rastegari, Joseph Redmon, Dieter Fox, and Ali Farhadi. IQA: visual question answering in interactive environments. In *CVPR*, pages 4089–4098, 2018. [3](#)
- [11] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997. [7](#)
- [12] Ronghang Hu, Marcus Rohrbach, Jacob Andreas, Trevor Darrell, and Kate Saenko. Modeling relationships in referential expressions with compositional modular networks. In *CVPR*, pages 4418–4427, 2017. [3](#)
- [13] Ronghang Hu, Huazhe Xu, Marcus Rohrbach, Jiashi Feng, Kate Saenko, and Trevor Darrell. Natural language object retrieval. In *CVPR*, pages 4555–4564, 2016. [3](#)
- [14] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 787–798, 2014. [2](#)
- [15] Liyiming Ke, Xiujun Li, Yonatan Bisk, Ari Holtzman, Zhe Gan, Jingjing Liu, Jianfeng Gao, Yejin Choi, and Siddhartha S. Srinivasa. Tactical rewind: Self-correction via backtracking in vision-and-language navigation. In *CVPR*, pages 6741–6749, 2019. [3](#), [6](#), [8](#)
- [16] Jingyu Liu, Liang Wang, and Ming-Hsuan Yang. Referring expression generation and comprehension via attributes. In *ICCV*, pages 4866–4874, 2017. [3](#)
- [17] Xihui Liu, Zihao Wang, Jing Shao, Xiaogang Wang, and Hongsheng Li. Improving referring expression grounding with cross-modal attention-guided erasing. In *CVPR*, pages 1950–1959, 2019. [8](#)
- [18] Ruotian Luo and Gregory Shakhnarovich. Comprehension-guided referring expressions. In *CVPR*, pages 3125–3134, 2017. [3](#)
- [19] Chih-Yao Ma, Zuxuan Wu, Ghassan AlRegib, Caiming Xiong, and Zsolt Kira. The regretful agent: Heuristic-aided navigation through progress estimation. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 6732–6740, 2019. [3](#)
- [20] Chih-Yao Ma, Jiasen Lu, Zuxuan Wu, Ghassan AlRegib, Zsolt Kira, Richard Socher, and Caiming Xiong. Self-monitoring navigation agent via auxiliary progress estimation. In *ICLR*, 2019. [3](#), [6](#), [7](#), [8](#)
- [21] Matt MacMahon, Brian Stankiewicz, and Benjamin Kuipers. Walk the talk: Connecting language, knowledge, and action in route instructions. In *Proceedings, The Twenty-First National Conference on Artificial Intelligence and the Eighteenth Innovative Applications of Artificial Intelligence Conference, July 16-20, 2006, Boston, Massachusetts, USA*, pages 1475–1482, 2006. [3](#)
- [22] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan Yuille, and Kevin Murphy. Generation

- and comprehension of unambiguous object descriptions. In *CVPR*, 2016. 2
- [23] Khanh Nguyen and Hal Daumé III. Help, anna! visual navigation with natural multimodal assistance via retrospective curiosity-encouraging imitation learning. *arXiv preprint arXiv:1909.01871*, 2019. 3
- [24] Khanh Nguyen, Debadepta Dey, Chris Brockett, and Bill 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*, pages 12527–12537, 2019. 3
- [25] Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer. Vision-and-dialog navigation. *CoRR*, abs/1907.04957, 2019. 3
- [26] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017. 6
- [27] Xin Wang, Qiuyuan Huang, Asli Çelikyilmaz, Jianfeng Gao, Dinghan Shen, Yuan-Fang Wang, William Yang Wang, and Lei Zhang. Reinforced cross-modal matching and self-supervised imitation learning for vision-language navigation. *CoRR*, abs/1811.10092, 2018. 3, 8
- [28] Xin Wang, Wenhan Xiong, Hongmin Wang, and William Yang Wang. Look before you leap: Bridging model-free and model-based reinforcement learning for planned-ahead vision-and-language navigation. In *ECCV*, pages 38–55, 2018. 3
- [29] Licheng Yu, Zhe Lin, Xiaohui Shen, Jimei Yang, Xin Lu, Mohit Bansal, and Tamara L. Berg. Mattnet: Modular attention network for referring expression comprehension. In *CVPR*, pages 1307–1315, 2018. 3, 6, 8
- [30] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C. Berg, and Tamara L. Berg. Modeling context in referring expressions. In *ECCV*, pages 69–85, 2016. 3, 5, 8

Supplements

In this supplementary material, we provide detailed explanation of evaluation metrics, examples of collected data, data collecting tools, how human test are performed and visualisation of several REVERIE results.

1. Evaluation Metrics

- **Navigation Success:** a navigation is considered successful only if the target object can be observed at the stop viewpoint. Please note that to encourage the agent to approach closer to the target object, we set the objects visible if they are within 3 meters away from the current location.
- **Navigation Oracle Success:** a navigation is considered oracle successful if the target object can be observed at one of its passed viewpoints.
- **Navigation SPL:** it is the navigation success weighted by the length of navigation path, which is

$$\frac{1}{N} \sum_{i=1}^N S_i \frac{\ell_i}{\max(\ell_i, p_i)} \quad (1)$$

where N is the number of tasks, $S_i \in \{0, 1\}$ is a binary indicator of success of task i , ℓ_i is the shortest length between the starting viewpoint and the goal viewpoint of task i , and p_i is the path length of an agent for task i .

- **Navigation Length:** trajectory length in meters.
- **REVERIE Success:** a task is considered REVERIE successful if the output bounding box has an IoU (intersection over union) ≥ 0.5 with the ground truth.

2. Typical Samples of The REVERIE Task

In Figure 1, we present several typical samples of the proposed REVERIE task. It shows the diversity in object category, goal region, path instruction, and target object referring expression.

3. Data Collecting Tools

To collect data for the REVERIE task, we develop a WebGL based data collecting tool as shown in Figure 2 and Figure 3. To facilitate the workers, we provide real-time updated reference information in the web page according to the location of the agent, including the current level/total level, the current region, and the number of regions in the build having the same function. At the goal location, in addition to highlighting the target object with a red 3D rectangle, we also provide the label of the target object and

the number of objects falling in the same category with the target object. Text and video instructions are provided for workers to understand how to make high quality annotations as shown in Figure 2.

4. Human Performance Test

To obtain the machine-human performance gap, we develop a WebGL based tool as shown in Figure 4 to test human performance. In the tool, we show an instruction about a remote object to the worker. Then the worker needs to navigate to the goal location and select one object as the target object from a range of object candidates. The worker can look around and go forward/backward by dragging or clicking.

5. Visualisation of REVERIE Results

In Figure 5, we provide the visualisation of several REVERIE results obtained by the typical state-of-the-art method, FAST-short, and the typical baseline method, R2R-SF.



Enter the bathroom with the red and black walls and turn on the sink.

Go down the hall to the bathroom of the bedroom with the large three section window and turn on the sink.

Go to the bathroom with black walls and clean out the sink.



Go to the office and clean the picture above the yellow stapler.

Go to the office room in the first level and bring me the picture to the right of the lamp.

Go to the office and clean the black and white picture of a child.



Go to the bedroom next to the bathroom on the second level and open the window on the left.

Move to the bedroom with the picture of a soup can and open the window on the far left.

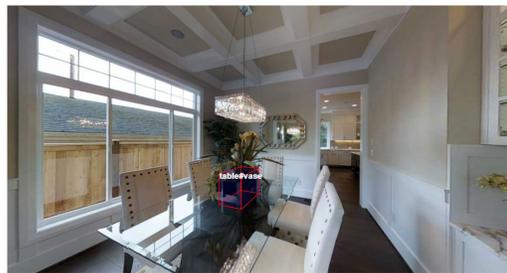
Go to second level bed room with a Campbell's Tomato Soup picture and clean the window nearest to this picture.



Go to the bathroom with the window and the two towels and take the towel that isn't a hand towel.

Go to the bathroom of the beige bedroom with an E on the wall and fold the towel on the right.

Go to the bathroom with a frilly white shower curtain and grab the towel directly across from the toilet.



Bring me the vase on the dining room table.

Go into the dining room and take the vase off the table.

Go to the dining room and take the vase.



Up the stairs in the room with the black painting on the wall take the left candle off the table.

Light the candle furthest from the windows in the lounge with green walls

Move to the lounge on level 2 with the egg sculpture and light the candle furthest from the windows

Figure 1. Several typical samples of the collected dataset, which involves various object category, goal region, path instruction, and object referring expression.

Instructions: Command a Smart Robot (Click to collapse)

Give a command to a smart robot to **find and interact with an object** in an indoor environment.

Please read the following instructions as they have been updated, hopefully to make things a bit clearer. Each individual hit should take less than one minute once you know how they're done

Instructions:

1. You need to first watch an animation (**press the Play/Replay button**), in which the robot automatically moves through a path from a **start location** to a **goal location** in a building.
 - The goal location is indicated by a **red** cylinder marker. **Green** cylinder indicates the start location. **Blue** cylinders indicate intermediate positions on the path.
2. At the goal location, a target object is marked in a **red bounding box**. Use the **Left/Right/Up/Down** arrow direction keys to look around and find it.
3. Now you need to think of a command and fill in the box below. It should be a command you could ask another human to get to the room and interact with the target object.
 - **Examples**
 - Open the left window in the kitchen.
 - Go to the living room on level 2 and bring me a pillow.
 - **About the navigation part in a command**
 - Focus on describing the goal location and **not the path** itself (e.g. Walk to the Kitchen).
 - Include information about the floor number if the path moves between floors (e.g. Go to the kitchen on level 2).
 - Specify the room if there are **multiple of the same room** on that level (e.g. Go to the bedroom with the yellow walls and a black couch). You can see information about the number of rooms and level in the reference information section below the viewer.
 - **About the interaction part in a command**
 - It should be the one that you could ask another human or a smart robot to interact with the highlighted target object (e.g. Open the cupboard under the sink).
 - It might help to imagine a future scenario where you are commanding a robot butler in your home (pick up, open, place, turn on, bring me, give etc.).
 - Specify the target object if there are **multiple of similar objects** in the goal location.
4. Finally press the "submit" button to complete the hit

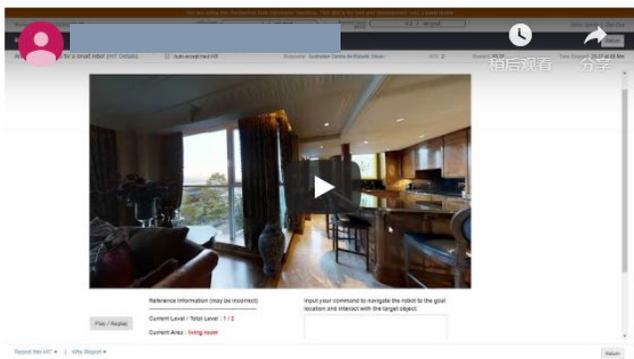
Controls:

1. Click the **'Play / Replay'** button at any time to watch a 15-20 second animated fly-through from the start to the goal.
2. **Left/Right/Up/Down keys on your keyboard control the camera direction** to look around.

Notes for hard situations:

- There are many rooms and repeats of the same objects, e.g. a building may have 3 bedrooms and 2 laptops **so try to mention nearby objects or landmarks in your commands** that could help to accurately locate the goal location and the target object. Your commands should be sufficient for another human to find the exact object in the correct room.
- **Do not use detailed navigation instructions. Some turkers once helped us on a similar task, where detailed navigations are required. But this is a different task. So please read the whole instruction carefully and watch the demo.**
- **The action should be directly operated on the target object itself not other items. For example, if a table is highlighted, then**
 - **Good interaction case: clean the table**
 - **Bad interaction case: put a plate on the table / remove the pot from the table**
- **Do not ask the robot to answer Wh-Questions or Wh-Question clause. For example, if a table is highlighted, then**
 - **Allowed: Could you please help to clean the table?**
 - **Not allowed: What is on the table? / Tell me what is one the table. / What color is the table? / Tell me the color of the table.**
- If the target object is not able to be seen, is obscured by another object, is not bound correctly by the red bounding box, or the object is not worthy to be found and interacted with, **mark the checkbox** and leave the command box blank.
- Do not mention the green/red/blue cylinder markers in your commands as the robot can not see them.
- The reference information labels for rooms and objects are **not always correct**, if you can think of better words to describe them you should use those instead
- Please use full sentences with punctuation (,) and correct spelling.
- The robot understands language and recognizes objects about as well as a typical person. However, you should assume that the robot is visiting this building for the first time.

Before you start, **please watch this short training video** (watch on youtube for fullscreen). It contains guided examples that will help you complete these tasks efficiently.



Note: This task is not suitable for devices with small screens or touch screen devices. Recommended browsers are Chrome, Firefox and Safari (not Internet Explorer).

These tasks relate to academic research conducted by the [redacted]. We estimate that on average each HIT to take around 1-1.5 minutes to complete. Please send your queries and feedback to [redacted] (should respond pretty quickly!). We will be continually releasing more HITs for this task.

Figure 2. Data collecting interface part I: instructions for AMT workers.



Reference Information (may be incorrect)

Play / Replay

Current Level / Total Level : 1 / 1

Current Area : **kitchen**

Number of areas having the same type as the current area:

1 kitchen in current level

Current Target Object Label: **phone**

Number of objects having the same type as the target object: **1**

Input your command to navigate the robot to the goal location and interact with the target object:

(1) Rich diversity in actions/verbs and sentence styles are encouraged.

***** If you want to avoid mistakes, read the instruction carefully.**

Email me if you are unsure about your commands.

The target object is not visible or is not suitable for a robot to find or interact with.

Submit

Figure 3. Data collecting interface Part II: assistant information and user input field.

Instructions: Navigate through a building from a command (Click to collapse)

Below the image you will see a navigation command that describes a goal location, and a command about a target object (both in red). **Your task is (1) to follow the direction, by exploring through the building to find the goal location, and (2) input the ID of the goal object.** Typically the goal location will be between 5 and 15 meters away, usually in a different room. We will award a **\$0.20 bonus** for every HIT finds the correct object. We will reject HITs from workers that are consistently far below average in performance.

Further requirements:

- You **should not explore the building unnecessarily**. Please go directly to the goal whenever possible.
- Please **do not submit until you are as close as possible** to where you think the goal is.
- You will be **assessed on the distance from your final location to the goal and object selection** (but it doesn't matter which direction you are facing).

Mouse Controls:

1. **Left-click and drag the image** to look around.
2. **Right-click on a blue cylinder** to move to that position (note: sometimes the blue cylinders are close to your feet, so you may need to look down).
3. Use the scroll wheel to zoom in and out.

Note: This task is not suitable for devices with small screens or touch screen devices. Recommended browsers are Chrome, Firefox and Safari (not Internet Explorer). Please do not submit HITs if you experience difficulties with the interface.

These tasks relate to academic research conducted by the [redacted]. We estimate that on average each HIT to take around 1 minute to complete. Please send your queries and feedback to [redacted].



Tips: Left-click and drag the panoramic image to start. Don't submit until you reach the goal and input the object ID. Must read full instructions at top. Different object number will be shown when you stand at different blue cylinders, move closer to the goal if the object is not highlighted.

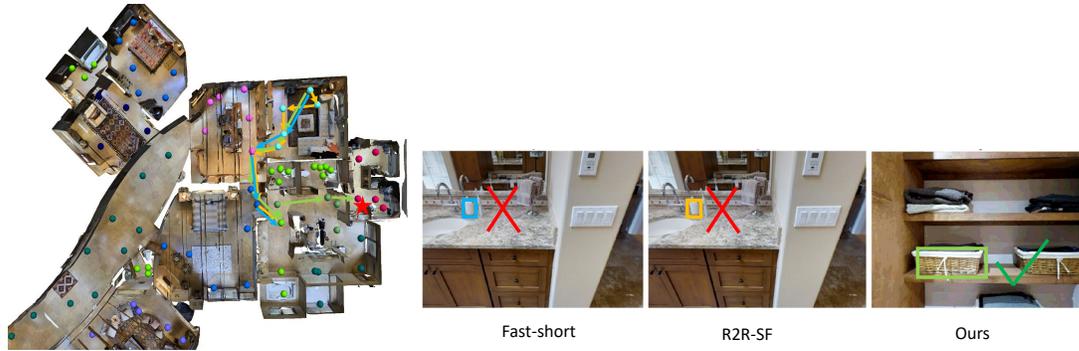
Instructions for target object (You need to deduce the object from the instruction. e.g. "Fluff the grey pillow" the object would be a grey pillow):

Clean the armchair furthest away from the front entrance in the living room.

Input the object ID (number) of the object from the image view above:

Submit

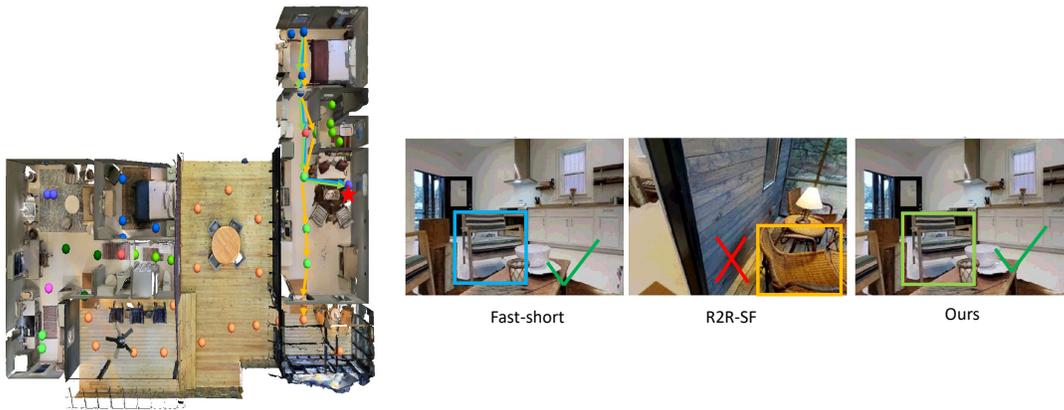
Figure 4. Human test interface. Workers need first to navigate to the goal location by clicking or dragging mouse, and then identify the target object. Only objects within 3 meters from the current location are highlighted. Different colors are used to facilitate workers to distinguish one object from others.



Go to the closet and bring me the pile of clothes on the left side of the shelf second from the top.



Go to the bathroom with a frilly white shower curtain and clean the toilet.



Go to the living room facing the kitchen and pull out the chair that is closer to the kitchen.



In the bathroom with the long mirror and towels and brush resting on the tub, wipe the sink out.



Figure 5. Visualisation of several REVERIE results, including trajectories and referring expression grounding of three typical methods. Colorized dots denote reachable locations, and different colors mark locations belonging to different regions according to the Matterport dataset.