



**HAL**  
open science

# Think Global, Act Local: Dual-scale Graph Transformer for Vision-and-Language Navigation

Shizhe Chen, Pierre-Louis Guhur, Makarand Tapaswi, Cordelia Schmid, Ivan Laptev

► **To cite this version:**

Shizhe Chen, Pierre-Louis Guhur, Makarand Tapaswi, Cordelia Schmid, Ivan Laptev. Think Global, Act Local: Dual-scale Graph Transformer for Vision-and-Language Navigation. CVPR 2022 - IEEE/CVF Conference on Computer Vision and Pattern Recognition, Jun 2022, New Orleans, United States. hal-03696868

**HAL Id: hal-03696868**

**<https://inria.hal.science/hal-03696868v1>**

Submitted on 16 Jun 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Think Global, Act Local: Dual-scale Graph Transformer for Vision-and-Language Navigation

Shizhe Chen<sup>†</sup>, Pierre-Louis Guhur<sup>†</sup>, Makarand Tapaswi<sup>‡</sup>, Cordelia Schmid<sup>†</sup> and Ivan Laptev<sup>†</sup>

<sup>†</sup>Inria, École normale supérieure, CNRS, PSL Research University <sup>‡</sup>IIT Hyderabad

[https://cshizhe.github.io/projects/vln\\_duet.html](https://cshizhe.github.io/projects/vln_duet.html)

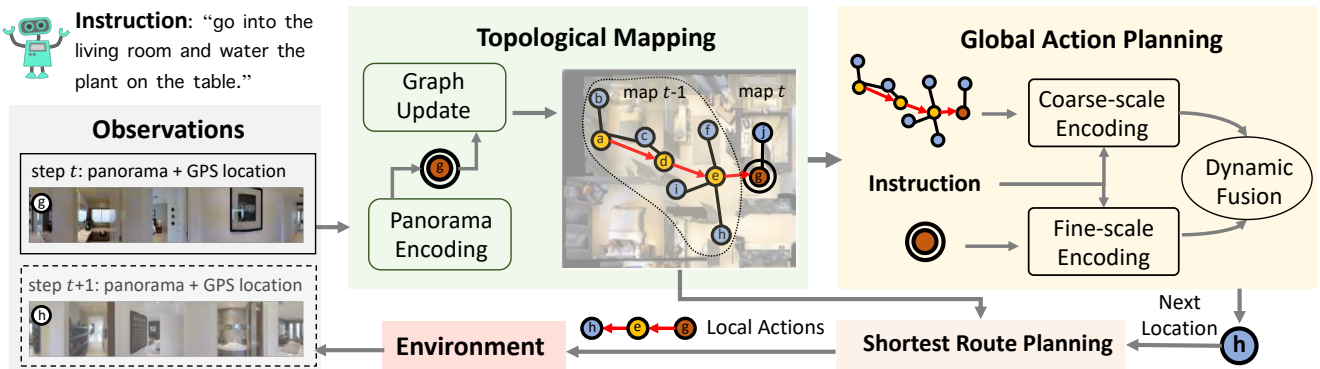


Figure 1. An agent is required to navigate in unseen environments to reach target locations according to language instructions. It only obtains local observations of the environment and is allowed to make local actions, *i.e.*, moving to neighboring locations. In this work, we propose to build topological maps on-the-fly to enable long-term action planning. The map contains visited nodes ● and navigable nodes ○ that can be reached from the previously visited nodes. Our method predicts global actions, *i.e.*, all navigable nodes in the map, and trades off complexity by combining a coarse-scale graph encoding with a fine-scale encoding ◎ of observations at the current node ●.

## Abstract

Following language instructions to navigate in unseen environments is a challenging problem for autonomous embodied agents. The agent not only needs to ground languages in visual scenes, but also should explore the environment to reach its target. In this work, we propose a **dual-scale graph transformer (DUET)** for joint long-term action planning and fine-grained cross-modal understanding. We build a topological map on-the-fly to enable efficient exploration in global action space. To balance the complexity of large action space reasoning and fine-grained language grounding, we dynamically combine a fine-scale encoding over local observations and a coarse-scale encoding on a global map via graph transformers. The proposed approach, DUET, significantly outperforms state-of-the-art methods on goal-oriented vision-and-language navigation (VLN) benchmarks REVERIE and SOON. It also improves the success rate on the fine-grained VLN benchmark R2R.

## 1. Introduction

Autonomous navigation is an essential ability for intelligent embodied agents. Given the convenience of natu-

ral language for human-machine interaction, autonomous agents should also be able to understand and act according to human instructions. Towards this goal, Vision-and-Language Navigation (VLN) [1] is a challenging problem that has attracted a lot of recent research [2–9]. VLN requires an agent to follow language instructions and to navigate in unseen environments to reach a target location. Initial approaches to VLN [2–4] use fine-grained instructions providing step-by-step navigation guidance such as “Walk out of the bedroom. Turn right and walk down the hallway. At the end of the hallway turn left. Walk in front of the couch and stop”. This fine-grained VLN task enables grounding of detailed instructions but is less practical due to the need of step-by-step guidance. A more convenient interaction with agents can be achieved by goal-oriented instructions [7, 8] such as “Go into the living room and water the plant on the table”. This task, however, is more challenging as it requires both the grounding of rooms and objects as well as the efficient exploration of environments to reach the target.

In order to efficiently explore new areas, or correct previous decisions, an agent should keep track of already executed instructions and visited locations in its memory. Many existing VLN approaches [2, 10–14] implement

memory using recurrent architectures, *e.g.* LSTM, and condense navigation history in a fixed-size vector. Arguably, such an implicit memory mechanism can be inefficient to store and utilize previous experience with a rich space-time structure. A few recent approaches [15, 16] propose to explicitly store previous observations and actions, and to model long-range dependencies for action prediction via transformers [17]. However, these models only allow for local actions, *i.e.*, moving to neighboring locations. As a result, an agent has to run its navigation model  $N$  times to backtrack  $N$  steps, which increases instability and compute.

A potential solution is to build a map [18] that explicitly keeps track of all visited and navigable locations observed so far. The map allows an agent to make efficient long-term navigation plans. For example, the agent is able to select a long-term goal from all navigable locations in the map, and then uses the map to calculate a shortest path to the goal. Topological maps have been explored by previous VLN works [8, 19, 20]. These methods, however, still fall short in two aspects. Firstly, they rely on recurrent architectures to track the navigation state as shown in the middle of Figure 2, which can greatly hinder the long-term reasoning ability for exploration. Secondly, each node in topological maps is typically represented by condensed visual features. Such coarse representations reduce complexity but may lack details to ground fine-grained object and scene descriptions in instructions.

Our approach addresses both of these shortcomings, the first one based on a transformer architecture and the second one with a dual-scale action planning approach. We propose a **Dual-scale graph Transformer (DUET)** with topological maps. As illustrated in Figure 1, our model consists of two modules: topological mapping and global action planning. In topological mapping, we construct a topological map over time by adding newly observed locations to the map and updating visual representations of nodes. Then at each step, the global action planning module predicts a next location in the map or a stop action. To balance fine-grained language grounding and reasoning over large graphs, we propose to dynamically fuse action predictions from dual scales: a fine-scale representation of the current location and a coarse-scale representation of the map. In particular, we use transformers to capture cross-modal vision-and-language relations, and improve the map encoding by introducing the knowledge of graph topology into transformers. We pretrain the model with behavior cloning and auxiliary tasks, and propose a pseudo interactive demonstrator to further improve policy learning. DUET significantly outperforms state-of-the-art methods on goal-oriented VLN benchmarks REVERIE and SOON. It also improves success rate on fine-grained VLN benchmark R2R. In summary, the contributions of our work are three-fold:

- We propose a dual-scale graph transformer (DUET)

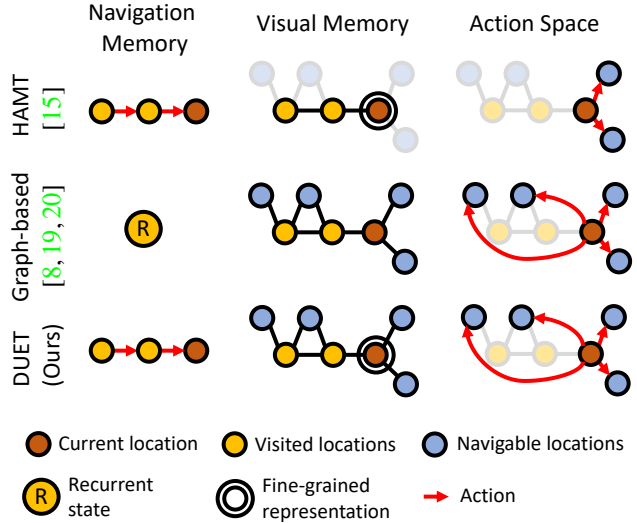


Figure 2. Method comparison. HAMT [15] stores navigation and visual memories to capture long-range dependency in action prediction, but is limited to a local action space. Graph-based approaches [8, 19, 20] use topological maps to support a global action space, but suffer from a recurrent navigation memory and a coarse-scale visual representation. Our DUET model overcomes previous limitations with a dual-scale encoding over the map.

with topological maps for VLN. It combines coarse-scale map encoding and fine-scale encoding of the current location for efficient planning of global actions.

- We employ graph transformers to encode the topological map and to learn cross-modal relations with the instruction, so that action prediction can rely on a long-range navigation memory.
- DUET achieves state of the art on goal-oriented VLN benchmarks, with more than 20% improvement on success rate (SR) on the challenging REVERIE and SOON datasets. It also generalizes to fine-grained VLN task, *i.e.*, increasing SR on R2R dataset by 4%.

## 2. Related work

**Vision-and-language navigation (VLN).** Navigation tasks involving instruction following [2–6, 9, 21–23] have become increasingly popular. Initial VLN methods mainly adopt recurrent neural networks with cross-modal attention [2, 10, 13, 24, 25]. More recently, transformer-based architectures have been shown successful in VLN tasks [26], notably by leveraging pre-trained architectures. For example, PRESS [27] adopts BERT [28] for instruction encoding. Different variants of ViLBERT are used in [29, 30] to measure compatibility between instructions and visual paths, but cannot be used for sequential action prediction. Recurrent VLN-BERT [14] addresses the limitation by injecting a recurrent unit in transformer architectures for action prediction. Instead of relying on one recurrent state,

E.T. [16] and HAMT [15] directly use transformers to capture long-range dependency to all past observations and actions (see first row in Figure 2).

**Maps for navigation.** The work on visual navigation has a long tradition of using SLAM [31] to construct metric maps [32] of the environment, using non-parametric methods [33], neural networks [34,35], or a mixture of both [36]. Anderson *et al.* [37] employ such metric maps for VLN tasks. However, it is challenging and requires accurate determination to construct metric map in real-time navigation. Therefore, several works [38,39] propose to represent the map as topological structures for pre-exploring environments [40], or for back-tracking to other locations, trading-off navigation accuracy with the path length [10,24]. A few recent VLN works [8,19,20] used topological maps to support global action planning, but they suffer from using recurrent architectures for state tracking and also lack a fine-scale representation for language grounding as shown in Figure 2. We address the above limitations via a dual-scale graph transformer with topological maps.

**Training algorithms for sequential prediction.** Behavior cloning is the most widely used training algorithm for sequential prediction. Nevertheless, it suffers from distribution shifts between training and testing. To address the limitation, different training algorithms have been proposed such as scheduled sampling [41], DAgger [42], reinforcement learning (RL) [43]. Most VLN works [13,14] combine behavior cloning and A3C RL [44]. Wang *et al.* [45] propose to learn rewards via soft expert distillation. Due to the difficulty of using RL in tasks with sparse rewards, we instead use an interactive demonstrator to mimic an expert and provide supervision in sequential training.

### 3. Method

**Problem formulation.** In the standard VLN setup for discrete environments [2,7,8], the environment is an undirected graph  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ , where  $\mathcal{V} = \{V_i\}_{i=1}^K$  denotes  $K$  navigable nodes, and  $\mathcal{E}$  denotes connectivity edges. An agent is equipped with an RGB camera and a GPS sensor, and is initialized at a starting node in a previously unseen environment. The goal of the agent is to interpret natural language instructions and to traverse the graph to the target location and find the object specified by the instruction.  $\mathcal{W} = \{w_i\}_{i=1}^L$  are word embeddings of the instruction with  $L$  words. At each time step  $t$ , the agent receives a panoramic view and position coordinates of its current node  $V_t$ . The panorama is split into  $n$  images  $\mathcal{R}_t = \{r_i\}_{i=1}^n$ , each represented by an image feature vector  $r_i$  and a unique orientation. To enable fine-grained visual perception,  $m$  object features  $\mathcal{O}_t = \{o_i\}_{i=1}^m$  are extracted in the panorama using annotated object bounding boxes or automatic object detectors [46]. In addition, the agent is aware of a few naviga-

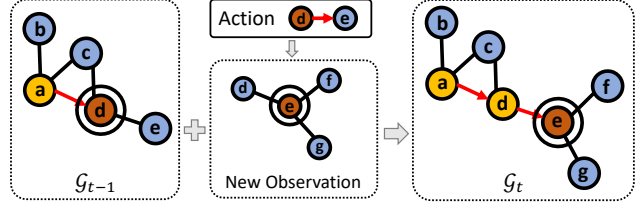


Figure 3. Illustration of graph updating at time step  $t$ . Given a new action  $d \rightarrow e$ , an agent receives new observations at node  $e$ . It then adds new nodes and updates node representations.

ble views corresponding to its neighboring nodes  $\mathcal{N}(V_t)$  as well as their coordinates. The navigable views of  $\mathcal{N}(V_t)$  are a subset of  $\mathcal{R}_t$ . The possible local action space  $\mathcal{A}_t$  at step  $t$  contains navigating to  $V_i \in \mathcal{N}(V_t)$  and stopping at  $V_i$ . After the agent decides to stop at a location, it needs to predict the location of the target object in the panorama.

Exploration and language grounding are two essential abilities for VLN agents. However, existing works either only allow for local actions  $\mathcal{A}_t$  [13–15] which hinders long-range action planning, or lack object representations  $\mathcal{O}_t$  [8,19,20] which might be insufficient for fine-grained grounding. Our work addresses both issues with a dual-scale representation and global action planning.

**Overview.** As illustrated in Figure 1, our model consists of two learnable modules, namely topological mapping and global action planning. The topological mapping module gradually constructs a topological map over time. The global action planning module then performs dual-scale reasoning based on coarse-scale global observations and fine-scale local observations. In the following, we introduce topological mapping in Sec. 3.1 and global action planning in Sec. 3.2. We end this section by presenting our approach to train our model and use it for inference in Sec. 3.3.

#### 3.1. Topological Mapping

The environment graph  $\mathcal{G}$  is initially unknown to the agent, hence, our model gradually builds its own map using observations along the path. Let  $\mathcal{G}_t = \{\mathcal{V}_t, \mathcal{E}_t\}$  with  $K_t$  nodes,  $\mathcal{G}_t \subset \mathcal{G}$  be the map of the environment observed after  $t$  navigation steps. There are three types of nodes in  $\mathcal{V}_t$  (see Figure 1): (i) visited nodes  $\bullet$  (yellow circle); (ii) navigable nodes  $\circ$  (blue circle); and (iii) the current node  $\bullet$  (orange circle). The agent has access to panoramic views for visited nodes and the current node. Navigable nodes are unexplored and are only partially observed from already visited locations, hence, they have different visual representations. At each step  $t$ , we add the current node  $V_t$  and its neighboring unvisited nodes  $\mathcal{N}(V_t)$  to  $\mathcal{V}_{t-1}$ , and update  $\mathcal{E}_{t-1}$  accordingly as illustrated in Figure 3. Given the new observation at  $V_t$ , we also update visual representations of the current node and navigable nodes as follows.

**Visual representations for nodes.** At time step  $t$ , the agent

receives image features  $\mathcal{R}_t$  and object features  $\mathcal{O}_t$  of node  $V_t$ . We use a multi-layer transformer [17] to model spatial relations among images and objects. The core of the transformer is the self-attention block:

$$[\mathcal{R}'_t, \mathcal{O}'_t] = \text{SelfAttn}([\mathcal{R}_t, \mathcal{O}_t]), \quad (1)$$

$$\text{SelfAttn}(X) = \text{Softmax}\left(\frac{XW_q(XW_k)^T}{\sqrt{d}}\right)XW_v, \quad (2)$$

where  $W_* \in \mathbb{R}^{d \times d}$  are parameters and biases are omitted. For ease of notation, we still use  $\mathcal{R}_t, \mathcal{O}_t$  in the following instead of  $\mathcal{R}'_t, \mathcal{O}'_t$  to denote the encoded embeddings.

Then we update visual representation of the current node  $\bullet$  by average pooling of  $\mathcal{R}_t$  and  $\mathcal{O}_t$ . As the agent also partially observes  $\mathcal{N}(V_t)$  at  $V_t$ , we accumulate visual representations of these navigable nodes  $\bullet$  based on the corresponding view embedding in  $\mathcal{R}_t$ . If a navigable node has been seen from multiple locations, we average all the partial view embeddings as its visual representation. We use  $v_i$  to denote the pooled visual representation for each node  $V_i$ . Such a coarse-scale representation enables efficient reasoning over large graphs, but may not provide sufficient information for fine-grained language grounding especially for objects. Therefore, we keep  $\mathcal{R}_t, \mathcal{O}_t$  as a fine-grained visual representation  $\bullet$  for the current node  $V_t$  to support detailed reasoning at a fine-scale.

## 3.2. Global Action Planning

Figure 4 illustrates the global action planning module. The coarse-scale encoder makes predictions over all previously visited nodes, but uses a coarse-scale visual representation. The fine-scale encoder instead predicts local actions given fine-grained visual representations of the current location. The dynamic fusion of both encoders combines predictions of global and local actions.

### 3.2.1 Text Encoder

To each word embedding in  $\mathcal{W}$  is added a positional embedding [28] corresponding to the position of the word in the sentence and a type embedding for text [47]. All word tokens are then fed into a multi-layer transformer to obtain contextual word representations, denoted here as  $\hat{\mathcal{W}} = \{\hat{w}_1, \dots, \hat{w}_L\}$ .

### 3.2.2 Coarse-scale Cross-modal Encoder

The module takes the coarse-scale map  $\mathcal{G}_t$  and encoded instruction  $\hat{\mathcal{W}}$  to make navigation predictions over a global action space ( $\cup_{i=1}^t \mathcal{A}_i$ ).

**Node embedding.** To the node visual feature  $v_i$  is added a location encoding and a navigation step encoding. The location encoding embeds the location of a node in the map in an egocentric view, which is the orientation and distance relative to the current node. The navigation step encoding embeds the latest visited time step for visited nodes and 0

for unexplored nodes. In this way, visited nodes are encoded with a different navigation history to improve alignment with the instruction. We add a ‘stop’ node  $v_0$  in the graph to denote a stop action and connect it with all other nodes.

**Graph-aware cross-modal encoding.** The encoded node and word embeddings are fed into a multi-layer graph-aware cross-modal transformer. Each transformer layer consists of a cross-attention layer [47] to model relations between nodes and instructions, and a graph-aware self-attention layer to encode environment layout. The standard attention in Eq. (2) only considers visual similarity among nodes, and thus it might overlook nearby nodes which are more relevant than distant nodes. To address the problem, we propose the graph-aware self-attention (GASA) which further takes into account the structure of the graph to compute attention as follows:

$$\text{GASA}(X) = \text{Softmax}\left(\frac{XW_q(XW_k)^T}{\sqrt{d}} + M\right)XW_v, \quad (3)$$

$$M = EW_e + b_e, \quad (4)$$

where  $X$  denotes node representations,  $E$  is the pair-wise distance matrix obtained from  $\mathcal{E}_t$ , and  $W_e, b_e$  are two learnable parameters. We stack  $N$  layers in the encoder and denote the output embedding of node  $V_i$  as  $\hat{v}_i$ .

**Global action prediction.** We predict a navigation score for each node  $V_i$  in  $\mathcal{G}_t$  as below:

$$s_i^c = \text{FFN}(\hat{v}_i), \quad (5)$$

where FFN denotes a two-layer feed-forward network. To be noted,  $s_0^c$  is the stop score. In most VLN tasks, it is not necessary for an agent to revisit a node, and thus we mask the score for visited nodes if not specially mentioned.

### 3.2.3 Fine-scale Cross-modal Encoder

This part attends to the current location  $V_t$  in the map to enable fine-scale cross-modal reasoning. The input is the instruction  $\hat{\mathcal{W}}_t$  and fine-grained visual representations  $\{\mathcal{R}_t, \mathcal{O}_t\}$  of the current node. The module predicts navigation actions in a local action space ( $\mathcal{A}_t$ ), and grounds the object at the final time step.

**Visual Embedding.** We add two types of location embeddings to  $\mathcal{R}_t, \mathcal{O}_t$ . The first type is the current location in the map relative to the start node. This embedding helps understand absolute locations in instruction such as “go to the living room in first floor”. Then for  $V_i \in \mathcal{N}(V_t)$ , we add a second location embedding, the relative position of each neighboring node to the current node. It helps the encoder to realize egocentric directions such as “turn right”. A special ‘stop’ token  $r_0$  is added for stop action.

**Fine-grained cross-modal reasoning.** We concatenate  $[r_0; \mathcal{R}_t; \mathcal{O}_t]$  as visual tokens and exploit a standard multi-

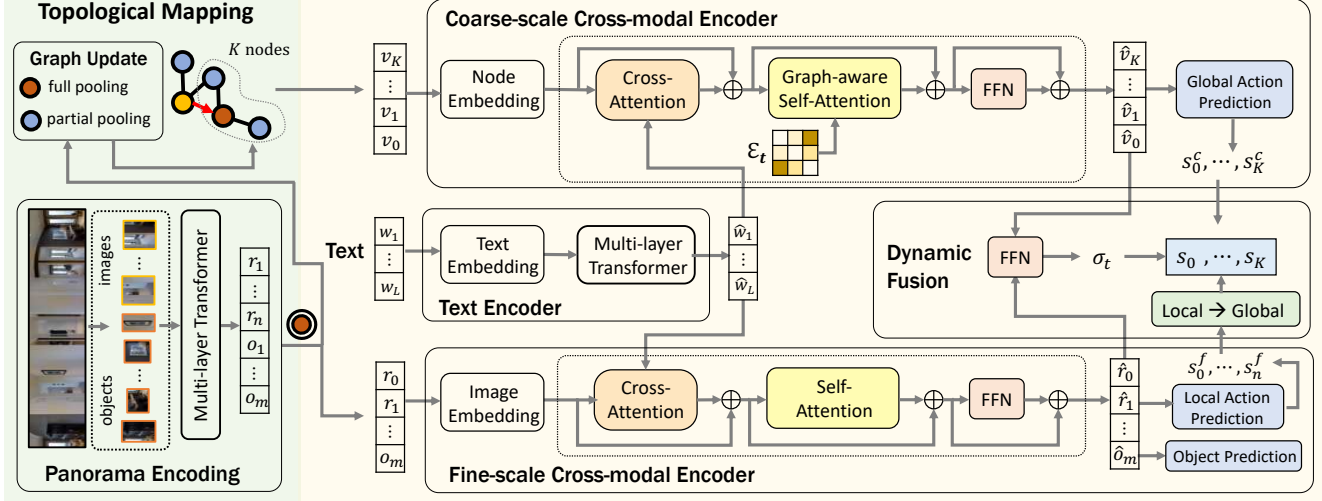


Figure 4. DUET consists of topological mapping (left) and global action planning (right). The mapping module outputs a graph with  $K$  node features  $\{v_i\}_{i=1}^K$ , and the current panorama encoding with image features  $\{r_i\}_{i=1}^n$  and object features  $\{o_i\}_{i=1}^m$ . Node feature  $v_0$  and image feature  $r_0$  are used to indicate the ‘stop’ action. The global action planning uses transformers for coarse- and fine-scale cross-modal encoding and fuses the two scales to obtain a global action score  $s_i$  for each node.

layer cross-modal transformer [47] to model vision and language relations. The output embeddings of visual tokens are represented as  $\hat{r}_0, \hat{\mathcal{R}}_t, \hat{\mathcal{O}}_t$  respectively.

**Local action prediction and object grounding.** We predict a navigation score  $s_i^f$  in local action space  $\mathcal{A}_t$  similar to Eq. (5). Moreover, as the goal-oriented VLN task requires object grounding, we further use a FFN to generate object scores based on  $\hat{\mathcal{O}}_t$ .

### 3.2.4 Dynamic Fusion

We propose to dynamically fuse coarse- and fine-scale action predictions for better global action prediction. However, the fine-scale encoder predicts actions in a local action space which does not match with the coarse-scale encoder. Therefore, we first convert local action scores  $s_i^f \in \{\text{stop}, \mathcal{N}(V_t)\}$  into the global action space. In order to navigate to other unexplored nodes that are not connected with the current node, the agent needs to backtrack through its neighboring visited nodes. Therefore, we sum over scores of visited nodes in  $\mathcal{N}(V_t)$  as an overall backtrack score  $s_{\text{back}}$ . We keep the values for  $s_i^f \in \{\text{stop}, \mathcal{N}(V_t)\}$  and use the constant  $s_{\text{back}}$  for the others. Hence, the converted global action scores are:

$$s_i^{f'} = \begin{cases} s_{\text{back}}, & \text{if } V_i \in \mathcal{V}_t - \mathcal{N}(V_t), \\ s_i^f, & \text{otherwise.} \end{cases} \quad (6)$$

At each step, we concatenate  $\hat{v}_0$  from coarse-scale encoder and  $\hat{r}_0$  from fine-scale encoder to predict a scalar for fusion:

$$\sigma_t = \text{Sigmoid}(\text{FFN}([\hat{v}_0; \hat{r}_0])). \quad (7)$$

The final navigation score for  $V_i$  is:

$$s_i = \sigma_t s_i^c + (1 - \sigma_t) s_i^{f'}. \quad (8)$$

## 3.3. Training and Inference

**Pretraining.** As shown in [15, 16, 26], it is beneficial to pre-train transformer-based VLN models with auxiliary tasks as initialization. Therefore, we first pretrain our model based on off-line expert demonstrations with behavior cloning and other common vision-and-language proxy tasks. We use masked language modeling (MLM) [28], masked region classification (MRC) [48], single-step action prediction (SAP) [15] and object grounding (OG) [49] if object annotations are available. The SAP and OG loss in behavior cloning given a demonstration path  $\mathcal{P}^*$  is as follows:

$$L_{\text{SAP}} = \sum_{t=1}^T -\log p(a_t^* | \mathcal{W}, \mathcal{P}_{<t}^*) \quad (9)$$

$$L_{\text{OG}} = -\log p(o^* | \mathcal{W}, \mathcal{P}_T) \quad (10)$$

where  $a_t^*$  is the expert action of a partial demonstration path  $\mathcal{P}_{<t}^*$ , and  $o^*$  is the groundtruth object at the last location  $\mathcal{P}_T$ . More details are presented in the supplementary material.

**Policy learning via an interactive demonstrator.** Behavior cloning suffers from distribution shifts between training and testing. Therefore, we propose to further train the policy with the supervision from a pseudo interactive demonstrator (PID)  $\pi^*$  similar to the DAgger algorithm [42]. During training we have access to the environment graph  $\mathcal{G}$ , hence  $\pi^*$  can utilize  $\mathcal{G}$  to select the next target node, i.e., a navigable node with the overall shortest distance from the current node and to the final destination. In each iteration, we use the current policy to sample a trajectory  $\mathcal{P}$  and use  $\pi^*$  to obtain pseudo supervision:

$$L_{\text{PID}} = \sum_{t=1}^T -\log p(a_t^{\pi^*} | \mathcal{W}, \mathcal{P}_{<t}) \quad (11)$$

where  $a_t^*$  is our pseudo target at step  $t$ . We combine the original expert demonstrations with our pseudo demonstrations in policy learning with a balance factor  $\lambda$ :

$$L = \lambda L_{SAP} + L_{PID} + L_{OG}. \quad (12)$$

**Inference.** At each time step during testing, we update the topological map as introduced in Sec. 3.1 and then predict a global action as explained in Sec. 3.2. If it is a navigation action, the shortest route planning module employs the Floyd algorithm to obtain a shortest path from the current node to the predicted node given the map, otherwise the agent stops at the current location. The agent is forced to stop if it exceeds the maximum action steps. In such case, it will return to a node with maximum stop probability as its final prediction. At the stopped location, the agent selects an object with maximum object prediction score.

## 4. Experiments

### 4.1. Datasets

We focus our evaluation on goal-oriented VLN benchmarks REVERIE [7] and SOON [8], which require fine-grained object grounding and advanced exploration capabilities to find a remote object. We also evaluate our model on the widely used VLN benchmark R2R [2], which has step-by-step instructions and no object localization.

**REVERIE** contains high-level instructions mainly describing target locations and objects. Instructions contain 21 words on average. Given predefined object bounding boxes provided for each panorama, the agent should select the correct object bounding box at the end of the navigation path. The length of expert paths ranges from 4 to 7 steps.

**SOON** also provides instructions describing target rooms and objects. The average length of instructions is 47 words. SOON does not provide object boxes and requires the agent to predict object center locations in the panorama. Hence, we use an automatic object detector [46] to obtain candidate object boxes. The length of expert paths ranges from 2 to 21 steps with 9.5 steps on average.

**R2R** contains step-by-step navigation instructions. The average length of instructions is 32 words. The average length of expert paths is 6 steps.

Examples from REVERIE and R2R are illustrated in Figure 5. Further details are in the supplementary material.

### 4.2. Evaluation Metrics

**Navigation metrics.** We use standard metrics [1] to measure navigation performance, i.e., Trajectory Length (TL): average path length in meters; Navigation Error (NE): average distance in meters between agent’s final location and the target; Success Rate (SR): the ratio of paths with NE

Table 1. Comparison of different scales and dual-scale fusion strategy on REVERIE val unseen split.

scale	fusion	OSR↑	SR↑	$\frac{SR}{OSR}$ ↑	SPL↑	RGS↑	RG SPL↑
fine	-	30.96	28.86	<b>93.22</b>	23.57	20.39	16.64
coarse	-	46.44	36.52	78.64	25.98	-	-
multi	average	<b>51.86</b>	45.81	88.33	31.94	<b>32.49</b>	22.78
	dynamic	51.07	<b>46.98</b>	91.40	<b>33.73</b>	32.15	<b>23.03</b>

less than 3 meters; Oracle SR (OSR): SR given oracle stop policy; and SR penalized by Path Length (SPL).

**Object grounding metrics.** To evaluate both the navigation and object grounding, we follow [7] and adopt Remote Grounding Success (RGS): the proportion of successfully executed instructions. We also use RGS penalized by Path Length (RG SPL). All the metrics are the higher the better except for TL and NE.

### 4.3. Implementation Details

**Features.** For images, we adopt ViT-B/16 [50] pretrained on ImageNet to extract features. For objects, we use the same ViT on the REVERIE dataset as it provides bounding boxes, while we use the BUTD object detector [46] on the SOON dataset. The orientation feature [11] contains  $\sin(\cdot)$  and  $\cos(\cdot)$  values for heading and elevation angles.

**Model architecture.** We use 9, 2, 4 and 4 transformer layers in the text encoder, panorama encoder, coarse-scale cross-modal encoder and fine-scale cross-modal encoder, respectively. Other hyper-parameters are set the same as in LXMERT [47], e.g., the hidden layer size is 768. We utilize the pretrained LXMERT for initialization.

**Training details.** On the REVERIE dataset, we first pre-train DUET with the batch size of 32 for 100k iterations using 2 Nvidia Tesla P100 GPUs. We automatically generate synthetic instructions to augment the dataset [10]. Then we use Eq. (12) to fine-tune the policy with the batch size of 8 for 20k iterations on a single Tesla P100. The best epoch is selected by SPL on val unseen split. More details are provided in supplementary material.

### 4.4. Ablation Study

We ablated our approach on the REVERIE dataset. All results in this section are reported on the val unseen split.

**1) Coarse-scale vs. fine-scale encoders.** We first evaluate coarse-scale and fine-scale encoders separately for the REVERIE navigation task in the upper part of Table 1. As the coarse-scale encoder is not fed with object representations, it is unable to select target objects for the REVERIE task. However, it outperforms the fine-scale version except for  $\frac{SR}{OSR}$ , for which the fine-scale encoder achieves much higher performance. This ratio estimates the performance of the stop action (the OSR is the success rate under oracle

Table 2. Ablation of graph-aware self-attention (GASA) for graph encoding on REVERIE val unseen split.

Fusion	GASA	OSR↑	SR↑	SPL↑	RGS↑	RGSPL↑
average	×	49.22	44.50	30.90	29.88	20.73
	✓	<b>51.86</b>	45.81	31.94	32.49	22.78
dynamic	×	49.25	45.24	32.88	29.91	21.57
	✓	51.07	<b>46.98</b>	<b>33.73</b>	<b>32.15</b>	<b>23.03</b>

stop policy) and indicates that fine-grained visual representations are essential to determine the target location specified in the instruction. However, the fine-scale encoder obtains a low OSR score, suggesting it lacks exploration due to a limited action space. The coarse-scale encoder instead benefits from the constructed map and is able to efficiently explore more areas with high OSR and SPL metrics.

**2) Dual-scale fusion strategy.** As the fine- and coarse-scale encoders are complementary, we compare different approaches to fuse the two encoders in the bottom part of Table 1. Both fusion methods outperform the fine-scale and coarse-scale encoder by a large margin. Our proposed dynamic fusion achieves more efficient exploration compared to the average fusion with 1.79% improvement on SPL.

**3) Graph-aware self-attention.** Table 2 ablates models with or without graph topology encoded in the transformer as in Eq. (3). It shows that the awareness of the graph structures is more beneficial to improve the SPL score, which emphasizes navigating to the target with shorter distance.

**4) Training losses.** In Table 3, we compare different training losses for DUET. The first row only uses  $L_{SAP}$  in behavior cloning. As it is not trained for object grounding, we can ignore RGS and RGSPL metrics. The second row adds the object supervision in training. It also improves navigation performance, which suggests that additional cross-modal supervisions such as association between words and objects can be beneficial to VLN tasks. In the third row, we add common auxiliary proxy tasks MLM and MRC in training, which are more helpful for object grounding. As instructions in REVERIE mainly describe the final target, these two losses are more relevant to object grounding. We further fine-tune the model with reinforcement learning (RL) [14, 15] or our PID in the last two rows to address distribution shift issue in behavior cloning. Both RL and PID achieve significant improvement and PID outperforms RL.

**5) Data augmentation with synthetic instructions.** We evaluate contributions of augmenting training data with synthetic instructions. The upper block of Table 4 presents results of pretraining with or without the augmented data. We can see that the synthetic data is beneficial in the pretraining stage and improves SPL and RGSPL by 1.63% and 1.76% respectively. Based on the initialization of the model in row 2, we use PID to further improve the policy. The results are

Table 3. Ablation of training losses on REVERIE val unseen split.

Pretrain		Finetune			OSR↑	SR↑	SPL↑	RGS↑	RGSPL↑
SAP	OG	Aux	RL	PID					
✓	×	×	×	×	38.45	35.30	24.55	-	-
✓	✓	×	×	×	40.24	37.80	26.40	23.89	16.36
✓	✓	✓	×	×	37.63	36.81	27.19	25.05	18.40
✓	✓	✓	✓	×	47.51	42.35	32.97	29.91	<b>23.53</b>
✓	✓	✓	×	✓	<b>51.07</b>	<b>46.98</b>	<b>33.73</b>	<b>32.15</b>	23.03

Table 4. Ablation of augmented speaker data in training on REVERIE val unseen split.

PID	Aug	OSR↑	SR↑	SPL↑	RGS↑	RGSPL↑
×	×	37.29	34.56	25.56	23.00	16.64
	✓	37.63	36.81	27.19	25.05	18.40
✓	×	51.07	<b>46.98</b>	<b>33.73</b>	<b>32.15</b>	<b>23.03</b>
	✓	<b>52.09</b>	46.58	32.72	31.75	22.18

Table 5. Comparison with the state of the art on SOON dataset.

Split	Methods	TL	OSR↑	SR↑	SPL↑	RGSPL↑
Val	GBE [8]	28.96	28.54	19.52	13.34	1.16
	DUET (Ours)	36.20	<b>50.91</b>	<b>36.28</b>	<b>22.58</b>	<b>3.75</b>
Test	GBE [8]	27.88	21.45	12.90	9.23	0.45
	DUET (Ours)	41.83	<b>43.00</b>	<b>33.44</b>	<b>21.42</b>	<b>4.17</b>

shown in the bottom block of Table 4. The synthetic data however does not bring improvements to the performance. We hypothesize that auxiliary proxy tasks in pretraining help to take advantage from the noisy synthetic data, but the policy learning still requires cleaner data.

#### 4.5. Comparison with State of the Art

**REVERIE.** Table 6 compares our final model with state-of-the-art models on the REVERIE dataset. Our model significantly beats the state of the arts on all evaluation metrics on the three splits. For example, on the val unseen split, our model outperforms the previous best model HAMT [15] by 14.03% on SR, 3.53% on SPL and 5.75% on RGSPL. Our model also generalizes better on the test unseen split, where we improve over HAMT by 22.11% on SR, 9.39% on SPL and 8.98% on RGSPL. This clearly demonstrates the effectiveness of our dual-scale action planning model with topological maps. Note that none of the previous methods has employed a map for navigation on this dataset.

**SOON.** Table 5 presents the results on the SOON dataset. Our model also achieves significant better performance than the previous graph-based approach GBE [8], with 20.54% gains on SR and 12.19% on SPL on test unseen split. The results, however, are much lower than those on REVERIE. This is because SOON contains fewer and more challenging



Table 6. Comparison with the state-of-the-art methods on REVERIE dataset.

Methods	Val Seen						Val Unseen						Test Unseen					
	Navigation			Grounding			Navigation			Grounding			Navigation			Grounding		
	TL	OSR $\uparrow$	SR $\uparrow$	SPL $\uparrow$	RGS $\uparrow$	RGSPL $\uparrow$	TL	OSR $\uparrow$	SR $\uparrow$	SPL $\uparrow$	RGS $\uparrow$	RGSPL $\uparrow$	TL	OSR $\uparrow$	SR $\uparrow$	SPL $\uparrow$	RGS $\uparrow$	RGSPL $\uparrow$
Human	-	-	-	-	-	-	-	-	-	-	-	-	21.18	86.83	81.51	53.66	77.84	51.44
Seq2Seq [2]	12.88	35.70	29.59	24.01	18.97	14.96	11.07	8.07	4.20	2.84	2.16	1.63	10.89	6.88	3.99	3.09	2.00	1.58
RCM [12]	10.70	29.44	23.33	21.82	16.23	15.36	11.98	14.23	9.29	6.97	4.89	3.89	10.60	11.68	7.84	6.67	3.67	3.14
SMNA [11]	7.54	43.29	41.25	39.61	30.07	28.98	9.07	11.28	8.15	6.44	4.54	3.61	9.23	8.39	5.80	4.53	3.10	2.39
FAST-MATTN [7]	16.35	55.17	50.53	45.50	31.97	29.66	45.28	28.20	14.40	7.19	7.84	4.67	39.05	30.63	19.88	11.61	11.28	6.08
SIA [49]	13.61	65.85	61.91	57.08	45.96	42.65	41.53	44.67	31.53	16.28	22.41	11.56	48.61	44.56	30.80	14.85	19.02	9.20
RecBERT [14]	13.44	53.90	51.79	47.96	38.23	35.61	16.78	35.02	30.67	24.90	18.77	15.27	15.86	32.91	29.61	23.99	16.50	13.51
Airbert [30]	15.16	48.98	47.01	42.34	32.75	30.01	18.71	34.51	27.89	21.88	18.23	14.18	17.91	34.20	30.28	23.61	16.83	13.28
HAMT [15]	12.79	47.65	43.29	40.19	27.20	25.18	14.08	36.84	32.95	30.20	18.92	17.28	13.62	33.41	30.40	26.67	14.88	13.08
DUET (Ours)	13.86	<b>73.86</b>	<b>71.75</b>	<b>63.94</b>	<b>57.41</b>	<b>51.14</b>	22.11	<b>51.07</b>	<b>46.98</b>	<b>33.73</b>	<b>32.15</b>	<b>23.03</b>	21.30	<b>56.91</b>	<b>52.51</b>	<b>36.06</b>	<b>31.88</b>	<b>22.06</b>

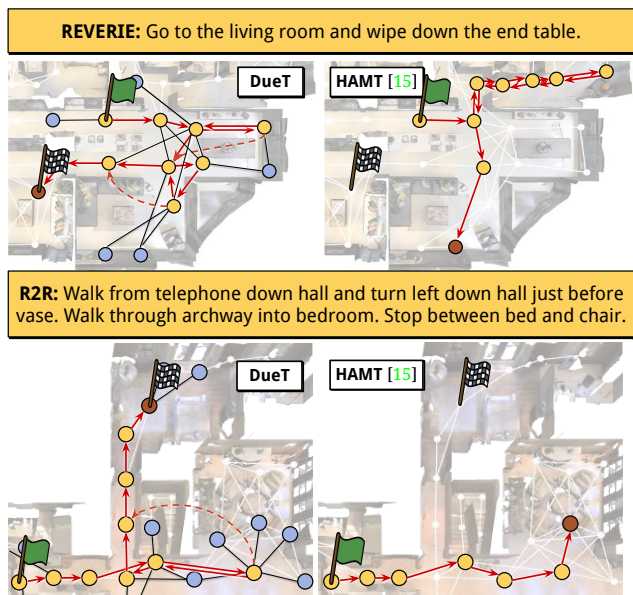


Figure 5. Predicted trajectories of DUET and the state-of-the-art HAMT [15]. The green and checkered flags denote start and target locations respectively. The dashed lines denote global actions. DUET is able to make more efficient explorations and correct its previous decisions, while HAMT is limited by its local actions.

training data (see supplementary material for analysis).

**R2R.** As shown in Table 7, DUET beats state-of-the-art approaches on success rate (SR) by 6% and 4% on val unseen and test unseen split respectively. However, it achieves comparable performances on SPL. This can be explained by the fact that for map-based approaches backtracking is encouraged which makes the trajectory length longer. We further compare a coarse-scale DUET for fair comparison with previous graph-based approaches [8, 19, 20] which do not use a fine-scale encoder. Even without using the fine-scale representation, DUET still outperform them by a margin, showing the effectiveness of our graph transformer. It also demonstrates DUET is able to backtrack more efficiently. Figure 5 visualizes some qualitative examples.

Table 7. Comparison with the state of the art on R2R dataset. Methods are grouped according to the used memories: ‘Rec’ for recurrent state, ‘Seq’ for sequence and ‘Map’ for topological map.

Mem	Methods	Val Unseen				Test Unseen			
		TL $\downarrow$	NE $\downarrow$	SR $\uparrow$	SPL $\uparrow$	TL $\downarrow$	NE $\downarrow$	SR $\uparrow$	SPL $\uparrow$
Rec	Seq2Seq [2]	8.39	7.81	22	-	8.13	7.85	20	18
	SF [10]	-	6.62	35	-	14.82	6.62	35	28
	PRESS [27]	10.36	5.28	49	45	10.77	5.49	49	45
	EnvDrop [13]	10.70	5.22	52	48	11.66	5.23	51	47
	AuxRN [51]	-	5.28	55	50	-	5.15	55	51
	PREVALENT [26]	10.19	4.71	58	53	10.51	5.30	54	51
	RelGraph [52]	9.99	4.73	57	53	10.29	4.75	55	52
RecBERT [14]	12.01	3.93	63	57	12.35	4.09	63	57	
Seq	HAMT [15]	11.87	3.65	65	59	12.65	4.11	63	58
	HAMT-e2e [15]	11.46	<b>2.29</b>	66	<b>61</b>	12.27	3.93	65	<b>60</b>
Map	EGP [19]	-	4.83	56	44	-	5.34	53	42
	GBE [8]	-	5.20	54	43	-	5.18	53	43
	SSM [20]	20.7	4.32	62	45	20.4	4.57	61	46
	DUET-coarse	12.96	3.67	68	59	13.08	3.93	67	58
	DUET (Ours)	13.94	3.31	<b>72</b>	60	14.73	<b>3.65</b>	<b>69</b>	59

## 5. Conclusion

We propose DUET (dual-scale graph transformer) for vision-and-language navigation (VLN) based on online constructed topological maps. It uses graph transformers to reason over a coarse-scale map representation for long-term action planning and a fine-scale local representation for fine-grained language grounding. The two scales are dynamically combined in the navigation policy. DUET achieves state-of-the-art performance on VLN benchmarks REVERIE, SOON and R2R.

**Acknowledgement.** This work was granted access to the HPC resources of IDRIS under the allocation 101002 made by GENCI. This work is funded in part by the French government under management of Agence Nationale de la Recherche as part of the ‘‘Investissements d’avenir’’ program, reference ANR19-P3IA-0001 (PRAIRIE 3IA Institute) and by Louis Vuitton ENS Chair on Artificial Intelligence.

## References

- [1] Peter Anderson, Angel Chang, Devendra Singh Chaplot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, et al. On evaluation of embodied navigation agents. *arXiv preprint arXiv:1807.06757*, 2018. 1, 6
- [2] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian 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. 1, 2, 3, 6, 8
- [3] Howard Chen, Alane Suhr, Dipendra Misra, Noah Snaveley, and Yoav Artzi. Touchdown: Natural language navigation and spatial reasoning in visual street environments. In *CVPR*, pages 12538–12547, 2019. 1, 2
- [4] Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, and Jason Baldridge. Room-across-room: Multilingual vision-and-language navigation with dense spatiotemporal grounding. In *EMNLP*, pages 4392–4412, 2020. 1, 2
- [5] Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. Beyond the nav-graph: Vision-and-language navigation in continuous environments. In *ECCV*, pages 104–120. Springer, 2020. 1, 2
- [6] Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *CVPR*, pages 10740–10749, 2020. 1, 2
- [7] Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang, William Yang Wang, Chunhua Shen, and Anton van den Hengel. Reverie: Remote embodied visual referring expression in real indoor environments. In *CVPR*, pages 9982–9991, 2020. 1, 3, 6, 8
- [8] Fengda Zhu, Xiwen Liang, Yi Zhu, Qizhi Yu, Xiaojun Chang, and Xiaodan Liang. Soon: Scenario oriented object navigation with graph-based exploration. In *CVPR*, pages 12689–12699, 2021. 1, 2, 3, 6, 7, 8
- [9] Muhammad Zubair Irshad, Chih-Yao Ma, and Zsolt Kira. Hierarchical cross-modal agent for robotics vision-and-language navigation. In *ICRA*, pages 13238–13246, 2021. 1, 2
- [10] 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. 1, 2, 3, 6, 8
- [11] Chih-Yao Ma, Jiasen Lu, Zuxuan Wu, Ghassan Al-Regib, Zsolt Kira, Richard Socher, and Caiming Xiong. Self-monitoring navigation agent via auxiliary progress estimation. In *ICLR*, 2019. 1, 6, 8
- [12] Xin Wang, Qiuyuan Huang, Asli Celikyilmaz, 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. In *CVPR*, pages 6629–6638, 2019. 1, 8
- [13] Hao Tan, Licheng Yu, and Mohit Bansal. Learning to navigate unseen environments: Back translation with environmental dropout. In *NAACL*, pages 2610–2621, 2019. 1, 2, 3, 8
- [14] Yicong Hong, Qi Wu, Yuankai Qi, Cristian Rodriguez-Opazo, and Stephen Gould. Vln BERT: A recurrent vision-and-language BERT for navigation. In *CVPR*, pages 1643–1653, 2021. 1, 2, 3, 7, 8
- [15] Shizhe Chen, Pierre-Louis Guhur, Cordelia Schmid, and Ivan Laptev. History aware multimodal transformer for vision-and-language navigation. In *NeurIPS*, 2021. 2, 3, 5, 7, 8
- [16] Alexander Pashevich, Cordelia Schmid, and Chen Sun. Episodic transformer for vision-and-language navigation. In *ICCV*, pages 15942–15952, 2021. 2, 3, 5
- [17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017. 2, 4
- [18] Devendra Singh Chaplot, Ruslan Salakhutdinov, Abhinav Gupta, and Saurabh Gupta. Neural topological SLAM for visual navigation. In *CVPR*, pages 12875–12884, 2020. 2
- [19] Zhiwei Deng, Karthik Narasimhan, and Olga Russakovsky. Evolving graphical planner: Contextual global planning for vision-and-language navigation. In *NeurIPS*, volume 33, 2020. 2, 3, 8
- [20] Hanqing Wang, Wenguan Wang, Wei Liang, Caiming Xiong, and Jianbing Shen. Structured scene memory for vision-language navigation. In *CVPR*, pages 8455–8464, 2021. 2, 3, 8

- [21] Vihan Jain, Gabriel Magalhaes, Alexander Ku, Ashish Vaswani, Eugene Ie, and Jason Baldridge. Stay on the path: Instruction fidelity in vision-and-language navigation. In *ACL*, pages 1862–1872, 2019. 2
- [22] Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied question answering. In *CVPR*, pages 1–10, 2018. 2
- [23] Licheng Yu, Xinlei Chen, Georgia Gkioxari, Mohit Bansal, Tamara L Berg, and Dhruv Batra. Multi-target embodied question answering. In *CVPR*, pages 6309–6318, 2019. 2
- [24] Chih-Yao Ma, Zuxuan Wu, Ghassan AlRegib, Caiming Xiong, and Zsolt Kira. The regretful agent: Heuristic-aided navigation through progress estimation. In *CVPR*, pages 6732–6740, 2019. 2, 3
- [25] Arun Balajee Vasudevan, Dengxin Dai, and Luc Van Gool. Talk2nav: Long-range vision-and-language navigation with dual attention and spatial memory. *International Journal of Computer Vision*, 129(1):246–266, 2021. 2
- [26] Weituo Hao, Chunyuan Li, Xiujun Li, Lawrence Carin, and Jianfeng Gao. Towards learning a generic agent for vision-and-language navigation via pre-training. In *CVPR*, pages 13137–13146, 2020. 2, 5, 8
- [27] Xiujun Li, Chunyuan Li, Qiaolin Xia, Yonatan Bisk, Asli Celikyilmaz, Jianfeng Gao, Noah A Smith, and Yejin Choi. Robust navigation with language pretraining and stochastic sampling. In *EMNLP*, pages 1494–1499, 2019. 2, 8
- [28] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, pages 4171–4186, 2019. 2, 4, 5
- [29] Arjun Majumdar, Ayush Shrivastava, Stefan Lee, Peter Anderson, Devi Parikh, and Dhruv Batra. Improving vision-and-language navigation with image-text pairs from the web. In *ECCV*, pages 259–274. Springer, 2020. 2
- [30] Pierre-Louis Guhur, Makarand Tapaswi, Shizhe Chen, Ivan Laptev, and Cordelia Schmid. Airbert: In-domain pretraining for vision-and-language navigation. In *ICCV*, pages 1634–1643, 2021. 2, 8
- [31] Sebastian Thrun. Probabilistic robotics. *Communications of the ACM*, 45(3):52–57, 2002. 3
- [32] Sebastian Thrun. Learning metric-topological maps for indoor mobile robot navigation. *Artificial Intelligence*, 99(1):21–71, 1998. 3
- [33] Albert S Huang, Abraham Bachrach, Peter Henry, Michael Krainin, Daniel Maturana, Dieter Fox, and Nicholas Roy. Visual odometry and mapping for autonomous flight using an rgb-d camera. In *Robotics Research*, pages 235–252. Springer, 2017. 3
- [34] Jingwei Zhang, Lei Tai, Ming Liu, Joschka Boedecker, and Wolfram Burgard. Neural SLAM: Learning to explore with external memory. *arXiv preprint arXiv:1706.09520*, 2017. 3
- [35] Saurabh Gupta, James Davidson, Sergey Levine, Rahul Sukthankar, and Jitendra Malik. Cognitive mapping and planning for visual navigation. In *CVPR*, pages 2616–2625, 2017. 3
- [36] Devendra Singh Chaplot, Dhiraj Gandhi, Saurabh Gupta, Abhinav Gupta, and Ruslan Salakhutdinov. Learning to explore using active neural SLAM. In *ICLR*, 2020. 3
- [37] Peter Anderson, Ayush Shrivastava, Devi Parikh, Dhruv Batra, and Stefan Lee. Chasing ghosts: Instruction following as bayesian state tracking. *NeurIPS*, 32:371–381, 2019. 3
- [38] Nikolay Savinov, Alexey Dosovitskiy, and Vladlen Koltun. Semi-parametric topological memory for navigation. *ICLR*, 2018. 3
- [39] Kuan Fang, Alexander Toshev, Li Fei-Fei, and Silvio Savarese. Scene memory transformer for embodied agents in long-horizon tasks. In *CVPR*, pages 538–547, 2019. 3
- [40] Kevin Chen, Junshen K Chen, Jo Chuang, Marynel Vázquez, and Silvio Savarese. Topological planning with transformers for vision-and-language navigation. In *CVPR*, pages 11276–11286, 2021. 3
- [41] Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In *NeurIPS*, volume 28, 2015. 3
- [42] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *AISTATS*, pages 627–635. JMLR Workshop and Conference Proceedings, 2011. 3, 5
- [43] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018. 3

- [44] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *ICML*, pages 1928–1937. PMLR, 2016. 3
- [45] Hu Wang, Qi Wu, and Chunhua Shen. Soft expert reward learning for vision-and-language navigation. In *ECCV*, pages 126–141. Springer, 2020. 3
- [46] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *CVPR*, pages 6077–6086, 2018. 3, 6
- [47] Hao Tan and Mohit Bansal. LXMERT: Learning cross-modality encoder representations from transformers. In *EMNLP*, pages 5103–5114, 2019. 4, 5, 6
- [48] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *NeurIPS*, volume 32, 2019. 5
- [49] Xiangru Lin, Guanbin Li, and Yizhou Yu. Scene-intuitive agent for remote embodied visual grounding. In *CVPR*, pages 7036–7045, 2021. 5, 8
- [50] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth  $16 \times 16$  words: Transformers for image recognition at scale. *ICLR*, 2020. 6
- [51] Fengda Zhu, Yi Zhu, Xiaojun Chang, and Xiaodan Liang. Vision-language navigation with self-supervised auxiliary reasoning tasks. In *CVPR*, pages 10012–10022, 2020. 8
- [52] Yicong Hong, Cristian Rodriguez, Yuankai Qi, Qi Wu, and Stephen Gould. Language and visual entity relationship graph for agent navigation. *NeurIPS*, 33:7685–7696, 2020. 8
- [53] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, pages 2048–2057. PMLR, 2015.
- [54] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *EMNLP*, pages 1532–1543, 2014.
- [55] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niebner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. In *3DV*, pages 667–676. IEEE, 2017.
- [56] Federico Landi, Lorenzo Baraldi, Marcella Cornia, Massimiliano Corsini, and Rita Cucchiara. Perceive, transform, and act: Multi-modal attention networks for vision-and-language navigation. *arXiv preprint arXiv:1911.12377*, 2019.
- [57] Gabriel Ilharco, Vihan Jain, Alexander Ku, Eugene Ie, and Jason Baldridge. General evaluation for instruction conditioned navigation using dynamic time warping. In *NeurIPS Workshop*, 2019.