

Task-Driven Graph Attention for Hierarchical Relational Object Navigation

Michael Lingelbach¹, Chengshu Li¹, Minjune Hwang¹, Andrey Kurenkov¹, Alan Lou¹, Roberto Martín-Martín², Ruohan Zhang¹, Li Fei-Fei¹, Jiajun Wu¹

Abstract—Embodied AI agents in large scenes often need to navigate to find objects. In this work, we study a naturally emerging variant of the object navigation task, hierarchical relational object navigation (HRON), where the goal is to find objects specified by logical predicates organized in a hierarchical structure—objects related to furniture and then to rooms—such as finding an apple on top of a table in the kitchen. Solving such a task requires an efficient representation to reason about object relations and correlate the relations in the environment and in the task goal. HRON in large scenes (e.g. homes) is particularly challenging due to its partial observability and long horizon, which invites solutions that can compactly store the past information while effectively exploring the scene. We demonstrate experimentally that scene graphs are the best-suited representation compared to conventional representations such as images or 2D maps. We propose a solution that uses scene graphs as part of its input and integrates graph neural networks as its backbone, with an integrated task-driven attention mechanism, and demonstrate its better scalability and learning efficiency than state-of-the-art baselines.

I. INTRODUCTION

Searching for objects in large scenes is a challenging component of many embodied AI activities such as rearrangement tasks [1–3] and household activities [4, 5]. When searching for an object, the embodied AI agent needs to navigate across different rooms and uses what it observes along the way to make optimal navigation decisions until it finds the target. Given the complexity of natural scenes with multiple rooms and objects, a challenge in object navigation is to devise a scalable solution that efficiently represents and exploits known information for future decisions.

Most prior object search work focused on the version of the problem defined as finding *any* instance of the target object category (e.g., “find *any* pair of shoes”) – namely the **object navigation** problem [6–8]. An alternative definition requires finding a specific *instance* of the target object category, e.g., “the old shoes”. This task has been called **instance object navigation** [9] or ION. ION is oftentimes a more natural problem definition, as realistic downstream tasks usually require a specific object instance (“the red book”) rather than any instance (“any book”). In this work, we focus on a new instantiation of ION that introduces additional *hierarchical relational* constraints in the definition (object-furniture, furniture-room), such as “find the shoes under the bed in the bedroom”, or “find the mug on the table in the kitchen”. We call this problem **hierarchical relational object navigation** or **HRON**.

An ideal solution to an object navigation task should combine an appropriate input representation with an optimal mechanism to extract the necessary information to guide navigation; these can be task-dependent. While egocentric RGB-D images, segmented images, point clouds, and their integration into 2D or 3D (semantic) maps have been successfully used as inputs to object navigation [10–12], solutions that use them as input perform poorly for HRON (see Sec. V). This is because they cannot effectively represent relational information nor scale to large, multi-room natural scenes. In contrast, scene graphs [13–16] – graphs where nodes are objects or rooms, and edges are pairwise relations between them – provide a compact scene representation that captures the critical information to guide a HRON solution. Therefore, our proposed HRON solution uses a scene graph built during exploration as the input representation.

The information encoded in graph structures can be extracted and leveraged efficiently using graph neural networks (GNNs) [15, 17–22]. However, GNNs’ success may be limited if they are applied naively to large graphs with irrelevant nodes and edges like the ones representing realistic scenes with hundreds of objects. Hence, we propose to integrate task-conditioned attention into GNNs to focus on the task-relevant elements of the graph in order to better aggregate their features and solve the current task.

In summary, our contributions are threefold:

- We introduce the hierarchical relational object navigation (HRON) task. HRON requires more sophisticated reasoning about object and room relations than object navigation and instance object navigation.
- We propose a novel solution to HRON based on a scene graph representation, that combines graph neural networks and task-driven attention for better scalability and learning efficiency for HRON in large scenes. Through experimental evaluation, we show that a reinforcement learning (RL) agent with our proposed architecture outperforms prior work with better performance and sample efficiency.
- We introduce concrete instantiations of HRON in three tasks of increasing complexity and realism. We provide a symbolic implementation for the first task and a physically-grounded implementation for the other two tasks in iGibson 2.0 [23] – a 3D simulator that provides photorealistic rendering and physics simulations in large household scenes. These environments will be publicly accessible for future research.

¹ Department of Computer Science, Stanford University, CA, USA,
² University of Texas at Austin, TX, USA, mjlbach@stanford.edu

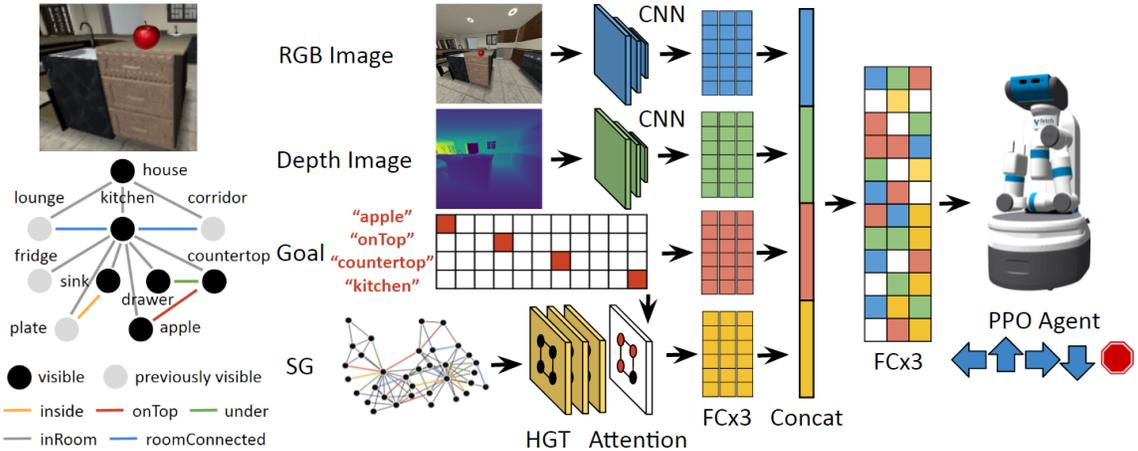


Fig. 1: Overview of our proposed model. We incrementally build a Scene Graph (SG) from RGB-D images where each node represents an object or a room, and each edge represents a physical relation between them. The SG is processed by a Heterogeneous Graph Transformer (HGT) with a task-driven attention mechanism to extract task-relevant information. The model also receives the current RGB-D image and a one-hot encoded goal description, which are processed by their respective learned encoders. All the features are then concatenated together before being fed to more fully connected (FC) layers, and eventually to a PPO agent that outputs navigation actions.

II. RELATED WORK

a) Object navigation: Embodied search for objects with partial observability is a long-studied problem in embodied AI and robotics [10, 21, 22, 24–39, 39–41]. Over the past several years many works have focused on learning-based approaches, since such approaches require less prior knowledge in novel environments [42]. However, they are typically limited to searching in single rooms only containing one instance of the target object [8]. Recently, ION (Instance-level Object Navigation) [9] and SOON (Scenario Oriented Object Navigation) [22] expanded the problem definition to locating an instance of an object with specific attributes and relations. However, ION still focuses on small single-room scenes, and SOON focuses on discrete selection over pre-defined waypoints trained with privileged information rather than embodied navigation. In contrast, our focus is on long-horizon learned embodied navigation in a large, multi-room (hierarchical) scene. We therefore choose to focus on defining object instances strictly in terms of relational constraints, since these constraints are naturally suited to directing the agent towards certain areas in larger scenes. A comparison between different problem setups can be seen in Table I.

b) Scene graphs in object navigation: Although learning a navigation policy directly from sensor input is viable [12], additional inductive biases and knowledge representations improve efficiency in larger scenes. Unlike other representations such as semantic maps, scene graphs scale with the number of objects rather than the size of the scene, making them suitable as a knowledge and memory representation for object navigation in large-scale scenes [46–48]. Scene graphs explicitly and compactly store information about objects’ geometry, placement, semantics, and relationships, which makes them ideal for reasoning about object relations. Furthermore, graphs naturally encode hierarchical relations and have been utilized in several prior works [15, 40, 49, 50].

However, while prior works focused on using scene graphs for the standard object navigation task with homogeneous graph edges, our focus is on studying how to best leverage scene graphs for HRON problems with directed heterogeneous edge types. Relational information in the edges provides useful information for targeted exploration of the scene. For example, when searching for an apple on a table in the kitchen, the presence of an "in room" edge between a particular table instance and a kitchen instance indicates this table should be prioritized during exploration.

c) Attention in object navigation: The importance of attention in human visual search has been long recognized [51–53]. Developing an attention mechanism for embodied AI agents is an active research topic. Such attention can be in the form of saliency map on egocentric RGB images [54], or weights on 2D maps [55, 56]. Attention can potentially help leverage large scene graphs in visual search, but an appropriate attention mechanism still needs to be explored since there are many forms of attention in graph neural networks developed for different purposes [57–59]. Our work builds upon existing research and proposes a framework for incorporating task-driven attention into a scene graph representation for the challenging HRON problem.

III. METHOD

Our solution for HRON comprises four elements (Fig. 1). First, at the core is a scene graph representation of the environment, incrementally constructed from current RGB-D images during task execution, and used as part of the input to our model (Sec. III-A). Second, a graph neural network and task-driven attention module are used to summarize the information of the scene graph into a single graph feature (Sec. III-B). Third, visual RGB-D inputs and the goal description, a tuple of one-hot encodings, are converted to feature vectors through learned layers, and all the input

	ON [8]	ION [9]	SOON [22]	HRON (Ours)
Goal (Example)	cup	blue, plastic cup near toaster	cylindrical, metallic and tall lamp which is set in the bright living room...The living room is on the first floor, next to the dining room and next to the kitchen	apple on top of counter in the kitchen
Environment	Habitat [43] single-room	AI2-Thor [44] single-room	Matterport3D [45] multi-room	iGibson 2.0 [23] multi-room
Embodiment Task horizon	navigation commands vary	navigation commands ~45 steps	node choice in a pre-defined graph ~10 steps	navigation commands ~100 steps

TABLE I: A comparison between different instances of object navigation (ON) tasks. HRON focuses on long-horizon learned embodied navigation in a large, multi-room (hierarchical) scene.

vectors are fused into a single vector (Sec. III-C). Finally, this fused vector is the learned representation used by a reinforcement learning agent that trains through interaction with the environment to find objects specified by hierarchical relational constraints (Sec. III-D). We then describe each component in detail.

A. Scene graph representation

A scene graph is comprised of a set of nodes and a set of directed edges, where each node represents a physical entity (e.g. objects, rooms). The node feature includes the attributes and the states of the physical entity (e.g. semantic class, 3D pose, size). All the nodes have their poses defined in the local coordinate frame of the agent, such that the agent knows its own position in the scene graph. The directed edges represent physical relations between the entities, e.g. “roomConnected” (room-room relation), “onTop” (object-object relation), and “inRoom” (object-room relation).

To build the scene graph, our method extracts information from the RGB-D images using a perfect object detector provided by the iGibson simulator at each step of the navigation task and accumulates it in an incremental fashion. In simulation, we emulate the execution of a scene graph building method such as the ones presented in prior work [14, 60, 60]. At any given step of the search, objects and rooms that are within the agent’s field of view but not yet in the scene graph are added to it as new nodes, and those that already exist will have their node features updated. The edges that connect to these nodes will be detected and updated as well. This procedure replicates a realistic incremental graph-building process in an embodied AI navigation agent.

B. Graph neural network architecture

a) Heterogeneous graph transformer (HGT): To compute a per node embedding, the graph is passed through three heterogeneous graph transformer layers [61] with ReLU activations. The HGT convolutions use distinct edge-based matrices for each edge type when computing attention, allowing the model to learn representations conditioned on different edge types, rather than connectivity alone.

b) Graph attention pooling: The final pooling layer applies a weighted mean pooling over all the node embeddings with a task-driven attention mechanism to create a single vector summarizing the scene graph. In order to effectively aggregate task-relevant nodes, the task-driven attention mechanism assigns a weight of 1 to all nodes for which the semantic

category matches any semantic category found in the current episodic goal description, and 0 otherwise. .

C. Multimodal feature fusion

To fuse features of global/history info (scene graphs), local info (current RGB-D images), and goal info, we process the RGB-D images with three consecutive convolutional layers, flattening, and three fully-connected (FC) layers, and process the one-hot encoded goal description with three FC layers. We then concatenate the embeddings from all three branches (including the aforementioned scene graph branch) and pass the concatenated feature through an additional stack of three FC layers for feature fusion.

D. Policy training

The model is trained end-to-end using a PPO [62] implementation adapted from RLlib [63], across 8 parallel environments for approximately 1.5 million environment steps of experience. The fused feature vector mentioned above is passed through separate, dedicated 3-layer FC for policy and value networks, respectively. We found that it was important to have sufficient network depth for these two networks in order for the policy to leverage the graph features.

IV. EXPERIMENTAL EVALUATION

We design our experiments to answer these two questions:

- **Q1:** Does the scene graph representation help the agent learn faster and perform better in the HRON tasks?
- **Q2:** Does the attention mechanism facilitate learning in large, populated scenes?

A. Experimental setup

We design three tasks that illustrate different aspects of relational object reasoning with increasing complexity and realism (see Fig. 2):

a) Relational object choice: The agent is presented with a 2D environment with two circles and rectangles: one circle above and the other below their respective rectangle. The agent is given a goal description of “circle above rectangle” or “circle below rectangle” and must choose one of two possible actions (`left/right`) to select the side of the environment that satisfies the goal description. The agent is given a reward of 1 if it selects the correct action, and a reward of 0 otherwise. The episode terminates after one step (bandit problem). To study the effect of high scene complexity, we add a random number (up to 75) of triangles as distractors.



Fig. 2: Three concrete tasks that illustrate relational object reasoning: relational object choice, directed object navigation, and exploratory object navigation. In relational object choice (left), given a relational object goal such as “the circle above the rectangle”, the agent should output a binary choice (left/right) that corresponds to the side of the environment that satisfies the goal. In directed object navigation and exploratory object navigation (middle, right), given a relational object goal such as “the apple on top of the table in the living room”, the agent should output a sequence of discrete navigation actions (forward/backwards/left/right/stop) to find and get close enough to the target object.

The other two tasks (directed and exploratory object navigation) are implemented in iGibson 2.0 [23], which provides photorealistic rendering and accurate physics simulation, with an onboard RGB-D camera. All assets (scenes and objects) from iGibson 2.0 [23] can be freely used within the iGibson simulator. Both the assets and the simulator are publicly available from their website under MIT license.

b) Directed object navigation: The agent is presented with a symmetrically arranged room, where an object (“bowl”, “gym shoe”, “apple”) is spawned according to a relational state (“onTop”, “inside”, “under”) with respect to an associated piece of furniture (“shelf”, “table”) on each half of the room, differing only in terms of relational state. The agent is initialized at the same starting position on each episode. The observation space includes 1) RGB-D images from onboard sensors, and 2) tokenized, one-hot encoded goal descriptions (object-relation-furniture). The agent outputs one of the five discrete navigation actions: *forward* (0.2m), *backward* (0.2m), *left* (turn 30 degrees), *right* (30 degrees), and *stop*, which will be physically executed. The agent achieves success if it navigates within a fixed distance of the goal object ($d = 1$ m). The episode terminates if the agent runs out of time, with a maximum episode length of 500 timesteps, or equivalently, 50 simulated seconds, or if the agent approaches the incorrect object ($d = 1$ m). The agent is given a reward of 10 if it achieves success, a reward of -5 if it approaches the incorrect object, and a reward of 0 otherwise. The agent is also provided with a geodesic-distance-based reward that encourages the agent to approach the goal object.

c) Exploratory object navigation: The agent is sampled randomly in one of the rooms in the realistic Wainscott_0_int iGibson scene [65] populated with furniture. The goal object category, relational state, associated furniture category, and room category of the target object are randomly selected at the beginning of each episode, e.g. “apple on top of the table in the living room”. An instance of an object model matching the object category is sampled to fulfill the given relational constraint in a physically stable manner (e.g., an apple is placed on top of an instance model of a

table in the living room). The observation space, action space, reward function, and termination conditions are identical to those of directed object navigation with two exceptions: 1) the goal description contains hierarchical relational constraints: object-relation-furniture and furniture-inRoom-room, 2) the episode does not terminate, and the agent doesn’t receive a negative reward when the agent approaches incorrect objects.

B. Baselines and ablation studies

We compare our method against two state-of-the-art baselines, Graph Convolution Network (GCN) by **Kipf et al.** [18] and the visual navigation method using scene priors by **Yang et al.** [21], both from the object navigation literature. Since the code for these works is not publicly available, we tried our best to replicate their approaches.

- **Kipf et al. [18]:** Graph Convolutional Networks (GCNs) are a simpler form of graph neural network compared to HGT, as it does not handle heterogeneous edges and involves no attention. Both [21] and [66] use GCNs to process their graphs. We implement this baseline by replacing our HGT with GCNs, and otherwise keeping all other aspects the same as our method.
- **Yang et al. [21]:** We replicate the approach used in [21]. Different from our method, this baseline represents the goal as a fastText [67] vector for the target object category, and creates node features from the fastText vectors for the nodes’ category and the ResNet-50 [68] softmax encoding of the current RGB image.

To showcase the effectiveness of each component of our model, we also perform extensive ablation studies.

- **RGB-D only:** The model only receives RGB-D images and one-hot encoded goal description (no scene graph).
- **RGB-D + MM:** Instead of incrementally building a scene graph, the model uses the accumulated RGB-D images to build a 2D, top-down, semantic metric map. Similar to SemExp [10], the model projects a semantically-segmented point cloud (extracted from depth images with ground-truth semantic class information from the

simulator) unto a 2D, top-down map as an image, where each pixel represents the physical space of $0.23\text{ m} \times 0.23\text{ m}$. If two points are projected onto the same 2D grid, the point with a higher z-value takes priority. The scene graph branch in the original model is replaced with a metric map branch that is identical to the RGB and Depth branches (conv layers, flattening, and MLP layers).

- **RGB-D + SG**: The attention mechanism is removed.
- **RGB-D + SG + TD ATTN**: This is our main model as described in Sec. III.

We also experiment with a set of ablations that provide additional information to the agent at the beginning of the episode. Specifically, we assume that the agent performs a pre-mapping procedure of the scene and stores all the furniture in a metric map or a scene graph. In other words, instead of incrementally building metric maps or scene graphs from scratch for each episode, the agent is provided with a pre-mapped scene representation that does not contain the target object as it is yet to be discovered. We argue that this is an alternative realistic setup that appears naturally when the agent pre-maps or searches consecutive for objects in the same scene. We call these variants of the models **PM: RGB-D + MM** and **PM: RGB-D + SG**. We run the full set of baselines and ablations for the exploratory object navigation task and a subset of them for the other two tasks when appropriate.

V. RESULTS

To answer the **Q1** and **Q2** raised in the previous section, our main finding is that scene graph representation does help embodied AI agents to perform better in tasks that require relational object reasoning. Moreover, in large, populated scenes, where scene graph size grows to around 100 nodes, the task-driven attention mechanism is essential for task performance since it significantly helps the aggregation of task-relevant information across nodes.

For quantitative results, we report the Success Rate (SR), whether or not the agent successfully approaches the target object within the time limit, and the Success weighted by Path Length (SPL), a ratio of the agent path length to the optimal path length conditioned on task success [69], averaged over 20 episodes across three random seeds.

a) Relational object choice: From Table II, we observed that when there are no distractors, both variants of our models quickly learn to output the correct action that matches the goal description. However, when there are a large number of distractors, our model without attention **SG** has a significant drop in performance. Our model with task-driven attention **SG + TD ATTN** can salvage most of the performance loss and achieve near-perfect success.

b) Directed object navigation: From Table III, we observed that the **RGB-D only** ablation is not able to solve the task, and achieves only chance-level performance. Similar to the previous task, when there are no distractors, both variants of our scene graph model **RGB-D + SG** and **RGB-D + SG + TD ATTN** are able to quickly learn to ground the goal description into the scene graph and choose a sequence of navigation actions that bring the robot close to the correct

TABLE II: Relational object choice: Success Rate

Distractors	Model	SR
No	SG	0.997±0.003
	SG + TD ATTN	0.991±0.009
Yes	SG	0.659±0.002
	SG + TD ATTN	0.993±0.007

TABLE III: Directed object navigation: Success Rate and Success weighted by Path Length

Distractors	Model	SR↑	SPL↑
No	RGB-D only	0.450±0.141	0.197 ± 0.039
	RGB-D + SG	0.967±0.024	0.949 ± 0.001
	RGB-D + SG + TD ATTN	0.983±0.024	0.958 ± 0.019
Yes	RGB-D only	0.417±0.103	0.173 ± 0.105
	RGB-D + SG	0.425±0.025	0.252 ± 0.046
	RGB-D + SG + TD ATTN	0.900±0.041	0.864 ± 0.039

TABLE IV: Exploratory object navigation: Success Rate and Success weighted by Path Length

Model	SR↑	SPL↑
Kipf et al. [18]	0.423±0.090	0.221±0.068
Yang et al. [21]	0.095±0.009	0.0302±0.0013
RGB-D only	0.586±0.021	0.309±0.061
RGB-D + MM	0.554±0.025	0.273±0.025
RGB-D + SG	0.458±0.135	0.183±0.099
RGB-D + SG + TD ATTN	0.879±0.048	0.577±0.041
PM: RGB-D + MM	0.405±0.023	0.404±0.022
PM: RGB-D + SG	0.634±0.108	0.402±0.037
PM: RGB-D + SG + TD ATTN	0.921±0.005	0.738±0.045

goal object. With distractors, on the other hand, task-driven attention is still critical for the effective aggregation of task-relevant information.

c) Exploratory object navigation: From the leftmost plot of Fig. 3 and Table IV, we observe that all models, including the **RGB-D only** model and the baseline models, are able to achieve some level of success. Naively introducing metric maps or scene graphs as external memory, i.e., **RGB-D + MM** and **RGB-D + SG** does not boost search performance by much. We hypothesize that simply accumulating information from the past RGB-D images introduces too much noise into the training process - after all, only a very small subset of the nodes are task-relevant for the given goal description.

The **Kipf et al. [18]** baseline performs comparably to **RGB-D + SG**, possibly because without attention, any processing of scene graphs does not yield significant improvements. Similarly, we believe that the **Yang et al. [21]** baseline is not able to learn at all due to its high dimensional node representation, which only serves to make learning harder.

As shown by the green curve in Fig. 3, the task-driven attention mechanism has a dramatic impact on the efficiency via which our method learns to extract task-relevant features from the graph. Our best model is able to converge to an 88% success rate with only 1 million environment steps, significantly outperforming baselines and ablations.

We also analyze the effectiveness of pre-mapping the scene. Comparing the dotted lines and the solid lines in the middle plot of Fig. 3, we can see that pre-mapping the

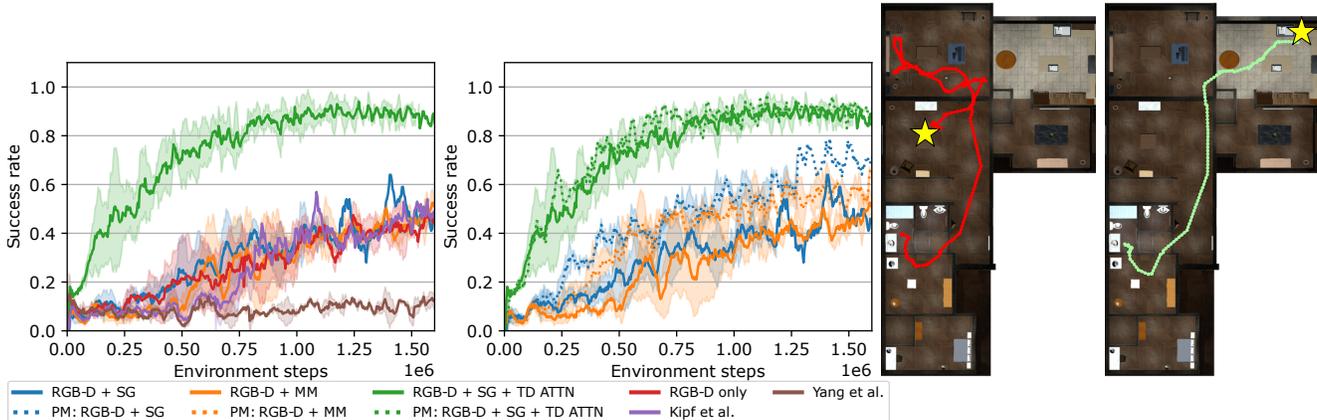


Fig. 3: Quantitative and qualitative results for the exploratory object navigation task. The left two plots depict Success Rate (SR) versus environment steps over training, without (left) and with (middle) pre-mapping, respectively. We observed that **RGB-D + SG** and **RGB-D + MM** slightly underperform **RGB-D only**. Adding the task-driven attention mechanism (**RGB-D + SG + TD ATTN**), however, results in a significant performance jump, doubling the success rate. The middle plot indicates that pre-mapping with furniture information helps guide more efficient exploration, and the scene graph models (**PM: RGB-D + SG** and **PM: RGB-D + SG + TD ATTN**) are better at extracting prior information than their metric map counterparts (**PM: RGB-D + MM**). Finally, the rightmost column showcases two example agent trajectories from the **RGB-D + SG + TD ATTN** model: the agent is able to leverage relational information stored in the scene graph to explore efficiently (right) and backtrack from its past mistakes after entering the wrong room (left).

scene and providing the agent with furniture information leads to a sizable performance boost across all three variants of our models. The improvement over the incrementally constructed scene graph indicates that our models **RGB-D + SG + TD ATTN** and **RGB-D + SG** can leverage prior information injected into the initial graph effectively. They also outperform the metric map counterpart (**RGB-D + MM**) with pre-mapping.

In the rightmost column of Fig. 3, we visualize two example agent trajectories of the exploratory object navigation task using the **RGB-D + SG + TD ATTN** model, where the star represents the location of the goal object. Thanks to the room-object and object-object relational information stored in the scene graph, our model can efficiently explore the scene to reach the goal object with the near-shortest path, or backtrack from past mistakes after entering the wrong room.

VI. DISCUSSION AND LIMITATIONS

In this work, we show the benefits of the scene graph as a representation for Hierarchical Relational Object Navigation (HRON) tasks that require reasoning about object relations. In large, populated scenes, having a task-driven attention mechanism is essential in aggregating task-relevant information and achieving a high success rate.

One advantage of explicit memory models such as scene graphs over latent representation such as weights in LSTMs [70] is the ability to store and retrieve prior scene knowledge such as the room configuration and furniture placement, which is invariant to the task goal. Providing scene priors by populating the scene graph with furniture information allows the scene graph-based model with pre-mapping (**PM: RGB-D + SG + TD ATTN**) to dramatically outperform other baselines.

The scene graph representation also has several desirable properties over the metric map representation in the context of relational object search. The graph complexity scales linearly with the number of objects in the scene, rather than the size of the physical space. Given a fixed memory, a metric map must make a tradeoff between representing small objects at high resolution, or large spaces at low resolution. Furthermore, due to the 2D top-down projection, the metric map has limitations when representing occlusion or containment.

The main challenge, addressed by task-driven attention, is transforming the object-centric scene graph into the scene graph embedding provided to the policy network. Naive global pooling, as seen in the results, mixes in information from non-relevant nodes. Our task-driven attention enriches this embedding with information from task-relevant nodes.

Our method, however, is not without limitations. The scene graphs in our model are constructed using privileged information from the simulator (a perfect 3D object detector), instead of from raw visual input, which has been accomplished in previous work [71]. Moreover, our problem setup is within the domain of embodied navigation without manipulation. Although the robot might collide with objects in the scene during navigation, the objects (and hence the scene graphs) stay largely static. How to leverage scene graphs and GNNs to solve mobile manipulation problems remains an active research area that is beyond the scope of this work.

ACKNOWLEDGMENT

This work is in part supported by ONR MURI N00014-22-1-2740, ONR MURI N00014-21-1-2801, NSF #2120095, AFOSR YIP FA9550-23-1-0127, Stanford Institute for Human-Centered AI (HAI), Amazon, Analog Devices, Bosch, JPMC, Meta, and Salesforce.

REFERENCES

- [1] D. Batra, A. X. Chang, S. Chernova, A. J. Davison, J. Deng, V. Koltun, S. Levine, J. Malik, I. Mordatch, R. Mottaghi, *et al.*, “Rearrangement: A challenge for embodied ai,” *arXiv preprint arXiv:2011.01975*, 2020.
- [2] C. Gan, S. Zhou, J. Schwartz, S. Alter, A. Bhandwaldar, D. Gutfreund, D. L. Yamins, J. J. DiCarlo, J. McDermott, A. Torralba, *et al.*, “The threedworld transport challenge: A visually guided task-and-motion planning benchmark for physically realistic embodied ai,” *arXiv preprint arXiv:2103.14025*, 2021.
- [3] L. Weihs, M. Deitke, A. Kembhavi, and R. Mottaghi, “Visual room rearrangement,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 5922–5931.
- [4] S. Srivastava, C. Li, M. Lingelbach, R. Martín-Martín, F. Xia, K. E. Vainio, Z. Lian, C. Gokmen, S. Buch, K. Liu, *et al.*, “Behavior: Benchmark for everyday household activities in virtual, interactive, and ecological environments,” in *Conference on Robot Learning*. PMLR, 2022, pp. 477–490.
- [5] T. Wisspeintner, T. Van Der Zant, L. Iocchi, and S. Schiffer, “Robocup@home: Scientific competition and benchmarking for domestic service robots,” *Interaction Studies*, vol. 10, no. 3, pp. 392–426, 2009.
- [6] Y. Zhu, X. Mottaghi, E. Kolve, J. J. Lim, A. Gupta, L. Fei-Fei, and A. Farhadi, “Target-driven visual navigation in indoor scenes using deep reinforcement learning,” in *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 3357–3364.
- [7] P. Anderson, A. Chang, D. S. Chaplot, A. Dosovitskiy, S. Gupta, V. Koltun, J. Kosecka, J. Malik, R. Mottaghi, M. Savva, *et al.*, “On evaluation of embodied navigation agents,” *arXiv preprint arXiv:1807.06757*, 2018.
- [8] D. Batra, A. Gokaslan, A. Kembhavi, O. Maksymets, R. Mottaghi, M. Savva, A. Toshev, and E. Wijnmans, “Objectnav revisited: On evaluation of embodied agents navigating to objects,” *arXiv preprint arXiv:2006.13171*, 2020.
- [9] W. Li, X. Song, Y. Bai, S. Zhang, and S. Jiang, “Ion: Instance-level object navigation,” in *Proceedings of the ACM International Conference on Multimedia*, 2021, pp. 4343–4352.
- [10] D. S. Chaplot, D. P. Gandhi, A. Gupta, and R. R. Salakhutdinov, “Object goal navigation using goal-oriented semantic exploration,” in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 4247–4258.
- [11] A. Wahid, A. Stone, K. Chen, B. Ichter, and A. Toshev, “Learning object-conditioned exploration using distributed soft actor critic,” *arXiv preprint arXiv:2007.14545*, 2020.
- [12] J. Ye, D. Batra, A. Das, and E. Wijnmans, “Auxiliary tasks and exploration enable objectnav,” *arXiv preprint arXiv:2104.04112*, 2021.
- [13] J. Johnson, R. Krishna, M. Stark, L.-J. Li, D. Shamma, M. Bernstein, and L. Fei-Fei, “Image retrieval using scene graphs,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3668–3678.
- [14] I. Armeni, Z.-Y. He, J. Gwak, A. R. Zamir, M. Fischer, J. Malik, and S. Savarese, “3d scene graph: A structure for unified semantics, 3d space, and camera,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 5664–5673.
- [15] Z. Ravichandran, L. Peng, N. Hughes, J. D. Griffith, and L. Carlone, “Hierarchical representations and explicit memory: Learning effective navigation policies on 3d scene graphs using graph neural networks,” in *IEEE International Conference on Robotics and Automation (ICRA)*, 2022.
- [16] N. Hughes, Y. Chang, and L. Carlone, “Hydra: A real-time spatial perception engine for 3d scene graph construction and optimization,” *arXiv preprint arXiv:2201.13360*, 2022.
- [17] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, “The graph neural network model,” *IEEE Transactions on Neural Networks*, vol. 20, no. 1, pp. 61–80, 2008.
- [18] T. N. Kipf and M. Welling, “Semi-supervised classification with graph convolutional networks,” in *International Conference on Learning Representations*, 2017.
- [19] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip, “A comprehensive survey on graph neural networks,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 1, pp. 4–24, 2020.
- [20] K. Chen, J. P. de Vicente, G. Sepulveda, F. Xia, A. Soto, M. Vazquez, and S. Savarese, “A behavioral approach to visual navigation with graph localization networks,” in *Proceedings of Robotics: Science and Systems*, 2019.
- [21] W. Yang, X. Wang, A. Farhadi, A. Gupta, and R. Mottaghi, “Visual semantic navigation using scene priors,” *arXiv preprint arXiv:1810.06543*, 2018.
- [22] F. Zhu, X. Liang, Y. Zhu, Q. Yu, X. Chang, and X. Liang, “Soon: scenario oriented object navigation with graph-based exploration,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 12 689–12 699.
- [23] C. Li, F. Xia, R. Martín-Martín, M. Lingelbach, S. Srivastava, B. Shen, K. E. Vainio, C. Gokmen, G. Dharan, T. Jain, *et al.*, “igibson 2.0: Object-centric simulation for robot learning of everyday household tasks,” in *5th Annual Conference on Robot Learning*, 2021.
- [24] Y. J. Tejwani, “Robot vision,” *IEEE International Symposium on Circuits and Systems*, pp. 1536–1539 vol.3, 1989.
- [25] L. E. Wixson and D. H. Ballard, “Using intermediate objects to improve the efficiency of visual search,” *International Journal of Computer Vision*, vol. 12, no. 2, pp. 209–230, 1994.
- [26] R. A. Baezayates, J. C. Culberson, and G. J. Rawlins, “Searching in the plane,” *Information and Computation*, vol. 106, no. 2, pp. 234–252, 1993.
- [27] P.-E. Forssén, D. Meger, K. Lai, S. Helmer, J. J. Little, and D. G. Lowe, “Informed visual search: Combining attention and object recognition,” in *2008 IEEE International Conference on Robotics and Automation*. IEEE, 2008, pp. 935–942.
- [28] S. Ekvall, D. Kragic, and P. Jensfelt, “Object detection and mapping for service robot tasks,” *Robotica*, vol. 25, no. 2, pp. 175–187, 2007.
- [29] S. P. Fekete, R. Klein, and A. Nüchter, “Online searching with an autonomous robot,” *Computational Geometry*, vol. 34, no. 2, pp. 102–115, 2006.
- [30] K. Sjö, D. G. López, C. Paul, P. Jensfelt, and D. Kragic, “Object search and localization for an indoor mobile robot,” *Journal of Computing and Information Technology*, vol. 17, no. 1, pp. 67–80, 2009.
- [31] D. Mishkin, A. Dosovitskiy, and V. Koltun, “Benchmarking classic and learned navigation in complex 3d environments,” *arXiv preprint arXiv:1901.10915*, 2019.
- [32] A. Mousavian, A. Toshev, M. Fišer, J. Koščeká, A. Wahid, and J. Davidson, “Visual representations for semantic target driven navigation,” in *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 8846–8852.
- [33] M. Wortsman, K. Ehsani, M. Rastegari, A. Farhadi, and R. Mottaghi, “Learning to learn how to learn: Self-adaptive visual navigation using meta-learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6750–6759.
- [34] H. Du, X. Yu, and L. Zheng, “Learning object relation graph and tentative policy for visual navigation,” in *European Conference on Computer Vision*. Springer, 2020, pp. 19–34.
- [35] J. Li, X. Wang, S. Tang, H. Shi, F. Wu, Y. Zhuang, and W. Y. Wang, “Unsupervised reinforcement learning of transferable meta-skills for embodied navigation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 123–12 132.
- [36] A. Staroverov, D. A. Yudin, I. Belkin, V. Adeshkin, Y. K. Solomentsev, and A. I. Panov, “Real-time object navigation with deep neural networks and hierarchical reinforcement learning,” *IEEE Access*, vol. 8, pp. 195 608–195 621, 2020.
- [37] S. Wani, S. Patel, U. Jain, A. Chang, and M. Savva, “Multion: Benchmarking semantic map memory using multi-object navigation,” in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 9700–9712.
- [38] H. Du, X. Yu, and L. Zheng, “Vtnet: Visual transformer network for object goal navigation,” in *International Conference on Learning Representations*, 2020.
- [39] S. Zhang, X. Song, Y. Bai, W. Li, Y. Chu, and S. Jiang, “Hierarchical object-to-zone graph for object navigation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 15 130–15 140.
- [40] A. Pal, Y. Qiu, and H. Christensen, “Learning hierarchical relationships for object-goal navigation,” in *Conference on Robot Learning*. PMLR, 2021, pp. 517–528.
- [41] F. Zhu, Y. Zhu, V. Lee, X. Liang, and X. Chang, “Deep learning for embodied vision navigation: A survey,” *arXiv preprint arXiv:2108.04097*, 2021.
- [42] X. Ye and Y. Yang, “From seeing to moving: A survey on learning for visual indoor navigation (vin),” *arXiv preprint arXiv:2002.11310*, 2020.
- [43] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijnmans, B. Jain, J. Straub, J. Liu, V. Koltun, J. Malik, *et al.*, “Habitat: A platform for

- embodied ai research,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9339–9347.
- [44] E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Herrasti, D. Gordon, Y. Zhu, A. Gupta, and A. Farhadi, “Ai2-thor: An interactive 3d environment for visual ai,” *arXiv preprint arXiv:1712.05474*, 2017.
- [45] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang, “Matterport3d: Learning from rgb-d data in indoor environments,” *arXiv preprint arXiv:1709.06158*, 2017.
- [46] S. Amiri, K. Chandan, and S. Zhang, “Reasoning with scene graphs for robot planning under partial observability,” *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 5560–5567, 2022.
- [47] G. Kumar, N. S. Shankar, H. Didwania, R. D. Roychoudhury, B. Bhowmick, and K. M. Krishna, “Gcexp: Goal-conditioned exploration for object goal navigation,” in *IEEE International Conference on Robot & Human Interactive Communication (RO-MAN)*. IEEE, 2021, pp. 123–130.
- [48] I. B. d. A. Santos and R. A. Romero, “A deep reinforcement learning approach with visual semantic navigation with memory for mobile robots in indoor home context,” *Journal of Intelligent & Robotic Systems*, vol. 104, no. 3, pp. 1–21, 2022.
- [49] Y. Zhu, J. Tremblay, S. Birchfield, and Y. Zhu, “Hierarchical planning for long-horizon manipulation with geometric and symbolic scene graphs,” in *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 6541–6548.
- [50] A. Kurenkov, R. Martín-Martín, J. Ichnowski, K. Goldberg, and S. Savarese, “Semantic and geometric modeling with neural message passing in 3d scene graphs for hierarchical mechanical search,” in *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 11 227–11 233.
- [51] G. Sperling and M. J. Melchner, “The attention operating characteristic: Examples from visual search,” *Science*, vol. 202, no. 4365, pp. 315–318, 1978.
- [52] M. J. Bravo and K. Nakayama, “The role of attention in different visual-search tasks,” *Perception & Psychophysics*, vol. 51, no. 5, pp. 465–472, 1992.
- [53] J. M. Wolfe and T. S. Horowitz, “Five factors that guide attention in visual search,” *Nature Human Behaviour*, vol. 1, no. 3, pp. 1–8, 2017.
- [54] B. Mayo, T. Hazan, and A. Tal, “Visual navigation with spatial attention,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 16 898–16 907.
- [55] Z. Seymour, K. Thopalli, N. Mithun, H.-P. Chiu, S. Samarasekera, and R. Kumar, “Maast: Map attention with semantic transformers for efficient visual navigation,” in *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 13 223–13 230.
- [56] S. K. Ramakrishnan, D. Jayaraman, and K. Grauman, “An exploration of embodied visual exploration,” *International Journal of Computer Vision*, vol. 129, no. 5, pp. 1616–1649, 2021.
- [57] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, “Graph attention networks,” *arXiv preprint arXiv:1710.10903*, 2017.
- [58] S. Yun, M. Jeong, R. Kim, J. Kang, and H. J. Kim, “Graph transformer networks,” in *Advances in Neural Information Processing Systems*, 2019.
- [59] D. Kim and A. Oh, “How to find your friendly neighborhood: Graph attention design with self-supervision,” *arXiv preprint arXiv:2204.04879*, 2022.
- [60] A. Rosinol, A. Gupta, M. Abate, J. Shi, and L. Carlone, “3d dynamic scene graphs: Actionable spatial perception with places, objects, and humans,” *arXiv preprint arXiv:2002.06289*, 2020.
- [61] Z. Hu, Y. Dong, K. Wang, and Y. Sun, “Heterogeneous graph transformer,” in *Proceedings of The Web Conference*, 2020.
- [62] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 2017.
- [63] E. Liang, R. Liaw, R. Nishihara, P. Moritz, R. Fox, K. Goldberg, J. Gonzalez, M. Jordan, and I. Stoica, “Rllib: Abstractions for distributed reinforcement learning,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 3053–3062.
- [64] M. Wise, M. Ferguson, D. King, E. Diehr, and D. Dymesich, “Fetch and freight: Standard platforms for service robot applications,” in *Workshop on Autonomous Mobile Service Robots*, 2016.
- [65] B. Shen, F. Xia, C. Li, R. Martín-Martín, L. Fan, G. Wang, C. Pérez-D’Arpino, S. Buch, S. Srivastava, L. Tchaptmi, *et al.*, “igibson 1.0: A simulation environment for interactive tasks in large realistic scenes,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 7520–7527.
- [66] Z. Seymour, N. C. Mithun, H.-P. Chiu, S. Samarasekera, and R. Kumar, “Graphmapper: Efficient visual navigation by scene graph generation,” in *International Conference on Pattern Recognition*, 2022.
- [67] A. Joulin, E. Grave, P. Bojanowski, and T. Mikolov, “Bag of tricks for efficient text classification,” *arXiv preprint arXiv:1607.01759*, 2016.
- [68] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [69] P. Anderson, A. X. Chang, D. S. Chaplot, A. Dosovitskiy, S. Gupta, V. Koltun, J. Kosecka, J. Malik, R. Mottaghi, M. Savva, and A. R. Zamir, “On evaluation of embodied navigation agents,” *CoRR*, vol. abs/1807.06757, 2018. [Online]. Available: <http://arxiv.org/abs/1807.06757>
- [70] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [71] S.-C. Wu, J. Wald, K. Tateno, N. Navab, and F. Tombari, “Scenegrph-fusion: Incremental 3d scene graph prediction from rgb-d sequences,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 7515–7525.