# Graph-based Cluttered Scene Generation and Interactive Exploration using Deep Reinforcement Learning

K. Niranjan Kumar
Georgia Institute of Technology

Irfan Essa
Georgia Institute of Technology

Sehoon Ha
Georgia Institute of Technology

*Abstract*— We introduce a novel method to teach a robotic agent to interactively explore cluttered yet structured scenes, such as kitchen pantries and grocery shelves, by leveraging the physical plausibility of the scene. We propose a novel learning framework to train an effective scene exploration policy to discover hidden objects with minimal interactions. First, we define a novel scene grammar to represent structured clutter. Then we train a Graph Neural Network (GNN) based *Scene Generation* agent using deep reinforcement learning (deep RL), to manipulate this *Scene Grammar* to create a diverse set of stable scenes, each containing multiple hidden objects. Given such cluttered scenes, we then train a *Scene Exploration* agent, using deep RL, to uncover hidden objects by interactively rearranging the scene. We show that our learned agents hide and discover significantly more objects than the baselines. We present quantitative results that prove the generalization capabilities of our agents. We also demonstrate sim-to-real transfer by successfully deploying the learned policy on a real UR10 robot to explore real-world cluttered scenes. The supplemental video can be found at: https://www.youtube.com/watch?v=T2Jo7wwaXss.

## I. INTRODUCTION

A long-standing goal of robotics is to build household robots that help with everyday chores like cooking or cleaning. Such tasks involve challenging real-world scenarios, which require a clear understanding of the environment. Computer Vision (CV) is a useful sensory modality to understand such scenes. Nevertheless, it only provides limited information when the scene is cluttered with objects, with many objects partially or completely occluded. For instance, a typical kitchen has dozens of products, all stacked away in small pantry cupboards and fridges. To estimate the underlying state of such a scene, a robotic agent might rely on interactions via manipulators, to augment its visual perception. However, developing an efficient agent that fully understands a scene, with minimal interaction, is not a straightforward problem, especially in scenes rampant with heavy visual occlusions and dense clutter.

The key for developing such a robotic agent is to *connect* visual observations to physics priors. Let us consider Figure 1a, that shows a cluttered scene of kitchen objects. While we cannot see the entire scene from the top, we can strategically explore it by, for instance, looking under the largest object as it is the most likely to be occluding another object. However, an intelligent agent should utilize additional clues of physical plausibility. For example, it can realize that some configurations (Figure 1b) are not physically plausible and that there is likely to be an object underneath (Figure 1c), even though these objects are not visible. Indeed, humans
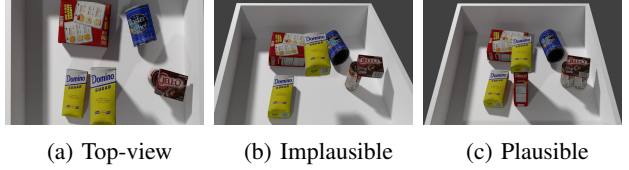


| (a) Top-view | (b) Implausible | (c) Plausible |

Fig. 1: We show how the agent can leverage clues about physical plausibility to search for objects better. When encountering the scene in Figure 1a, an agent equipped with a static camera cannot completely observe it. However, from just the top-view, a knowledge about physics priors can hint that scene (b) is impossible. The agent can then plan smarter interactions to reveal the true state of the scene (c).

have this intuitive physics engine built into our cognition [1] that helps us notice violations of physical laws, thereby aiding us in exploration. However, building this intuition in a machine, by explicitly defining a set of rules is tedious and impractical. We instead want the agent to *learn* these priors implicitly from a large number of experiences.

In this work, we propose a novel framework for learning an effective strategy to explore cluttered environments by (1) training *Scene Generation* and *Scene Exploration* agents, (2) that utilize a shared graphical representation of the scene, built upon a *Scene Grammar*. First, we define the Scene Grammar to represent the space of possible structured clutter scenes compactly. In this space of possible scenes afforded by the grammar, a majority would either be physically unstable or would not contain any hidden objects. To reliably generate stable scenes with multiple hidden objects, we train a Scene Generation Agent using deep reinforcement learning (deep RL). Then, we train an efficient Scene Exploration agent that can explore these challenging scenes generated with minimal number of interactions. For efficient training, we adopt the "learning by cheating" framework [2] that first trains a teacher agent with privileged information and then clones the learned behaviors to a student agent. We model our agents as Message Passing Neural Networks (MPNNs) [3] and train them with RL to maximize their respective objectives.

We demonstrate that our framework trains an effective Scene Exploration agent that explores complex scenes generated by the Scene Generation agent. We show quantitative comparisons that indicate significant improvements over baselines and provide evidence of sim-to-real transfer, by testing our method trained in simulation on real world clutter.

The contributions of our work are as follows:

1) A scene grammar to generate scene graphs of various cluttered multi-object scenes.
2) A learning framework to train an effective scene exploration agent that discovers hidden objects, along with a scene generation agent that specializes in hiding objects in clutter.
3) Demonstration of sim-to-real transfer of our learned exploration policy.

## II. RELATED WORK

Our work falls within the broad domain of interactive perception [4], which involves leveraging interaction to infer otherwise hidden properties of an environment, such as masses of objects [5], [6], state of a cluttered scene [7], [8] or kinematic structure of articulated objects [9].

*Object Search in Robotics:* Building robots that help in domestic households is a long-standing goal of robotics. It is important for such robots to search and retrieve objects in cluttered scenes. However, even when all objects are completely visible, retrieving is quite challenging if many objects are in close proximity. Therefore, a body of work in object singulation focuses on moving objects in clutter such that they are spread out and do not touch each other [10], [11], [12]. Once the objects are singulated, planning of feasible grasps becomes much easier as grasping does not depend on the clutter surrounding it.

Other approaches deal with occluded objects that may be partially or completely invisible [13], [14], [15], [16], [17]. Often in such heavily cluttered scenes, such as a bin [18], [13], [7], [14] or a shelf [19], [20], the agent has to move other objects around to search for the target. While these approaches learn to search for a given target object, a more generalized version of this problem is to identify and locate *all* objects present in the scene. Estimating the *state* of such structured cluttered scenes has been explored by prior works [21], [22] by using passive and active perception techniques. Another approach to tackle the problem is to get the scene to a state where all objects are clearly visible [8]. However, this may not always be feasible since additional space might not be available in some scenarios.

*Scene Generation:* The goal of scene generation is to create synthetic counterparts of realistic scenes that share resemblance at a structural and semantic level. A natural choice of representation for scenes is a *scene graph*, which typically represents objects as nodes and uses the edges to represent relational attributes between objects [23], [24], [25]. A body of work in scene generation [26], [27], [28] focuses on creating 3D scenes by first creating scene graphs and then instantiating the represented scene using a renderer. Formal systems [29] offer a promising approach to generate scene graphs. A formal system is defined by a set of non-terminal nodes, terminal nodes and production rules that define transformation from one sequence to another. Production rules can then be defined to add, remove or modify nodes and edges in the graph. Scene grammar was utilized by prior work [27], [30] to generate simulated scenes that mimicked the structure

and semantics of a real-world dataset. In contrast, we take a self-supervised approach to scene generation that is grounded in physical plausibility, by building scenes with cluttered but stable object arrangements in a *physics simulator*. In the work of Sui *et al.* [21], scene graphs are used along with a particle filter to represent multiple hypothesis of the underlying object-object relationships in the scene. They however assume that all the objects in the scene are visible and identifiable. In our case, the scenes frequently contain heavily occluded objects that even a state-of-the-art object detector will not be able to identify. Instead, we take an RL-based approach and learn manipulation strategies that can efficiently uncover (dis-occlude) these objects.

*Object Rearrangement:* Object rearrangement in the context of robotic manipulation has been investigated by prior works [31], [32], [33]. In a typical object rearrangement scenario, the agent is given a goal state (either explicitly or implicitly, such as organizing a pile of grocery items in a pantry). Similarly, in our work we train an agent to rearrange with the objective of revealing hidden objects. However, in contrast to prior approaches in object search [18], [14], we *do not require an auxiliary bin* to discard occluding objects. This provides our agent two advantages: *(i)* in real world environments, the robot may not have access to additional space where it can discard objects, giving our approach a competitive edge; and, *(ii)* at the end of its sequence of interaction, our robot knows the location and identity of every object present in the scene and can retrieve multiple objects without any additional exploration.

## III. GRAPH-BASED SCENE GENERATION AND EXPLORATION

The core idea of our method is to utilize scene graphs as a unified representation that generate and rearrange cluttered scenes with many objects. Our method consists of two agents: *(i)* a Scene Generation Agent that exploits a scene grammar and constructs scene graphs to represent difficult, physically stable and cluttered multi-object scenes; and, *(ii)* a Scene Exploration Agent that starts with scenes generated by the Scene Generation Agent and outputs a sequence of efficient manipulation strategies that uncovers hidden objects and reveals the underlying scene graph. We represent both these agents as MPNNs and train them with RL to maximize their respective objectives.

### A. Scene Grammar

In our approach, both the scene generation and scene exploration agents need to handle various scenes with multiple objects, thereby requiring the algorithm to efficiently represent and synthesize different configurations. To model the space of possible scenes an agent might encounter, we propose a *scene grammar* defined by a set of production rules that represent transformations from one scene graph to another. These rules are sequentially applied to add or remove objects to a given scene graph, thereby letting us control the structure and composition of a scene. Our scene grammar is recursive and contains production rules that
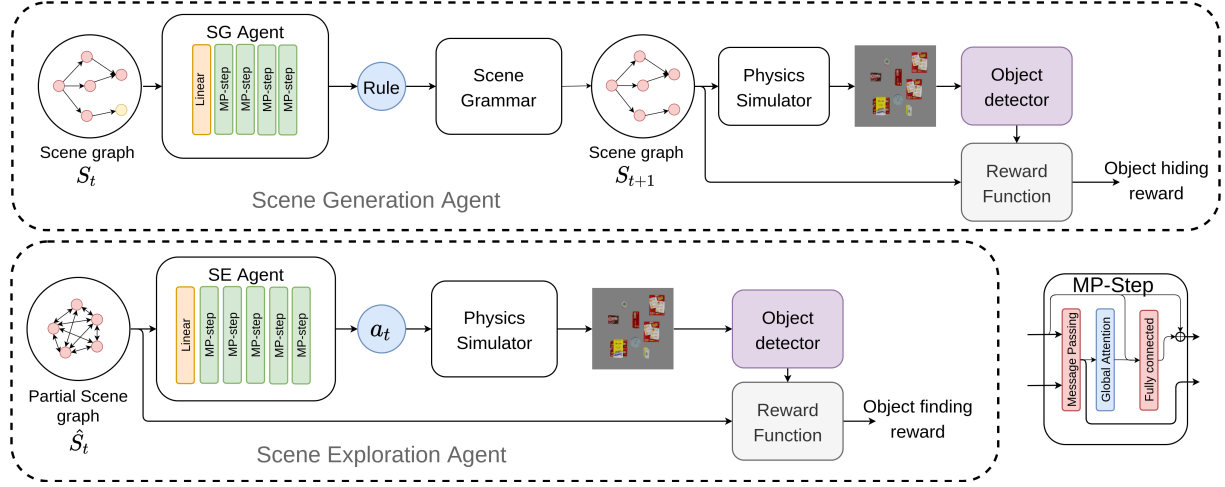
Fig. 2: Overview of our method. We train two agents with conflicting objectives while ensuring that the scene remains stable. The Scene Generation agent (SG) tries to hide more objects in the scene while the scene exploration agent (SE) tries to find them. Each MP-step consists of a message passing layer followed by Global Attention [34] with a skip connection.

recursively add structures to the graph to mimic stacking objects on top of and/or next to each other. We define the scene grammar $\mathcal{G}$ as follows:

$$\mathcal{G} = (N_{NT}, N_T, E, R, S) \tag{1}$$

where $N_{NT}$ are the non-terminal nodes, $N_T$ correspond to the terminal nodes, and $S$ is the start graph. $E$ represents the set of edges available in the grammar. Terminal attributes represent *real* aspects of the scene such as object type or orientation, while the non-terminal attributes are intermediate entities that help with the construction of the graph. For a graph to be simulated as a scene, all its attributes must be terminal. $R$ represents the production rules that are applied to transform the graph.

Each production rule is represented with two graphs, $G_L$ and $G_R$, that define the transformation it can perform. To apply a rule, we first search (using BFS) for a subgraph in the input graph that is isomorphic to $G_L$ and then replace it with $G_R$. Both $G_L$ and $G_R$ can contain terminal and/or non-terminal attributes. In our grammar we use a non-terminal node called *object_slot* as a placeholder that can later be substituted by object nodes. This design choice disentangles the structure of the graph from the actual objects it will be instantiated with, leading to a reduction in the number of grammar rules required.

Consider the example in Figure 3. We start out with a scene graph that contains two objects, a tomato soup can and a sugar box. To add new objects near the existing stack, we first apply the *drop_object* rule that creates a new *object_slot* attached to the tray. Next, we apply *stack_object* to stack another *object_slot* on-top of this node. Note that the resulting scene graph has non-terminal nodes (i.e. two *object_slot* nodes), and cannot be simulated as a scene. We then apply rules that substitute the *object_slot's* with actual objects that can be simulated in a scene. First we apply an *insert_meta_pbox_1* rule that inserts two pudding boxes, right next to each other, upon which more objects can be stacked.

Finally, we apply an *insert_cracker_box* rule that substitutes the topmost *object_slot* with a cracker box. This results in a scene graph that just contains terminal nodes that can be simulated as shown. For more details about implementation, see our project page [1].

### B. Scene Generation Agent

In the previous section, we defined a scene grammar that can represent arrangements of objects in a scene. But not all scene graphs that can be generated, would actually be stable in the real world. In addition, we want to generate scenes that have many hidden objects that escape detection from a vision-based perception system. This raises the question: *how can we procedurally generate difficult, stable and cluttered scenes from the set of production rules?*

We model the process of generating a scene graph using our scene grammar as a Markov Decision Process (MDP) defined by the 5-tuple; $\mathcal{M} : \langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma \rangle$ where $\mathcal{S}$, $\mathcal{A}$ and $\mathcal{R}$ are the state space, action space and reward function respectively. $\mathcal{T}$ is a transition function that determines the next state given the current state and action. $\gamma \in (0, 1)$ is the discount factor. We define how these terms relate to the scene generation procedure below.

**States**: We define the state at time-step $t$ to be the corresponding scene graph at $t$. This graph can contain both terminal and non-terminal nodes. Each node feature is defined to be a combination of the class the object belongs to, a boolean flag that denotes whether the node has been simulated and rendered in the environment and the position of the object, if it has been rendered in the scene. The edges in this graph represent the orientation of the child node.

**Actions**: We define the action space to be the set of all production rules $R$ (we use 30 rules in our implementation).

**Reward**: The reward function for our agent consists of two components. A stability penalty $r_s$ that penalizes the agent if
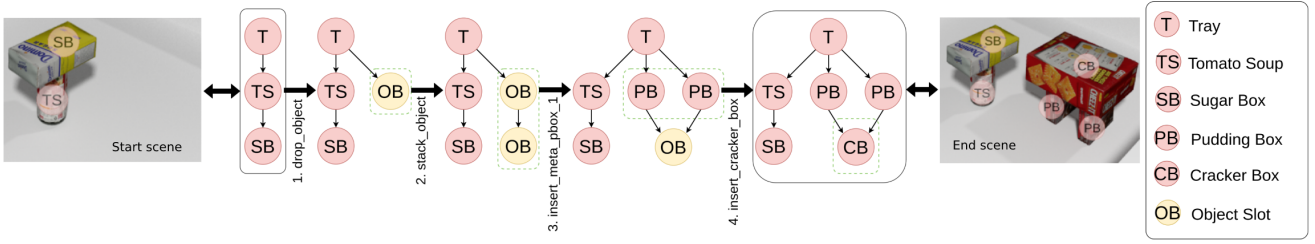
Fig. 3: We show an example to demonstrate our proposed scene grammar by applying a sequence of rules to a scene graph. The agent starts with a scene that contains a sugar box stacked on top of a tomato soup can on a tray. The agent then applies a series of rules to add more objects, namely two pudding boxes and a cracker box, to the scene.

the constructed scene is physically unstable. An uncertainty reward $r_u$ that is proportional to the number of objects the agent is able to hide in the scene.

$$\mathcal{R}(s) = w_s r_s + w_u r_u \qquad (2)$$

We use constant weights $w_s = -1$ and $w_u = 0.5$ for all our experiments. The transition function $\mathcal{T}$ is determined by possible transitions induced by the scene grammar.

In addition to maximizing the above objective, we also want to generate high entropy policies to ensure that that we create a *distribution* of cluttered scenes instead of just one that the Scene Exploration Agent can easily over-fit to. To this end, we exploit a maximum entropy RL based approach to this problem. We learn a policy with a MPNN architecture using Soft-Q learning [35] to output a categorical probability distribution over production rules available at a given state.

*C. Scene Exploration Agent*

The next component of our approach involves training an agent to discover hidden objects and estimate the underlying state of the scene. We formulate this problem as a Partially Observed Markov Decision Process (PoMDP) where the agent does not know the true underlying state of the scene $s$ and only has access to a partial observation $o$.

Consider a scene with $N$ objects (unknown during test time) represented by the set $\mathbf{S} : \{0, 1, \ldots N\}$, where 0 represents the ground and $\{1, \ldots N\}$ represents objects present in the tray. Since the scene can have occluded objects, we define a set of currently visible objects $\mathbf{V}_t \subseteq \mathbf{S}$, which contains all the objects detectable by the agent at time $t$, and a set of historically identified objects $\mathbf{K}_t \subseteq \mathbf{S}$, which contains all the objects the agent has seen *at least* once. At time $t$, the agent updates $\mathbf{K_t}$ using $\mathbf{K_t} = \mathbf{K_{t-1}} \cup \mathbf{V_t}$. The agent interacts with the scene until $\mathbf{K}_t \cup \{0\} = \mathbf{S}$. Given this premise, we define a PoMDP for the structured clutter search problem with the 7-tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{O}, \Omega, \mathcal{R}, \gamma \rangle$:

**States** ($\mathcal{S}$): The underlying state of the scene containing $N$ objects, at a given timestep $t$, is a scene graph $S_t$ with objects represented as nodes and object-object relationships represented as edges. The node for each object $i \in [1, N]$ contains the following: Pose $p_i$, the class ID, a boolean flag that denotes if the object has ever been detected by the robot.

**Actions** ($\mathcal{A}$): A tuple $(x, y) \mid x \in \mathbf{V_t}, y \in \mathbf{V_t} \cup \{0\}$ where object $x$ represents the object that is to be picked up and $y$ represents the object on top of which $x$ will be placed.

**Observations** ($\Omega$): A partial scene graph $\hat{S}_t$ that contains objects seen by the agent, i.e. objects in $\mathbf{K}_t$. Since, object-object relationships are unknown during test time, we define $\hat{S}_t$ as a fully connected graph with pose $p_i$ and class ID as the node features.

**Reward** ($\mathcal{R}$): The agent gets a reward of $r_d$ for every newly detected object and a penalty of $r_s$ if the scene becomes unstable, same as Equation 2.

$\mathcal{T}, \mathcal{O}$ denote the transition function $p(s_{t+1}|s_t, a_t)$ and the observation model $p(o_t|s_t)$ respectively, with $\gamma$ being the discount factor. With this setup, we train an agent to discover all the objects present in the scene by rearrangement. Such an agent should select actions that maximize the expected discounted sum of rewards defined by:

$$J = \mathbb{E}_{p(\tau)} \left[ \sum_{t=1}^{T} \gamma^{t-1} r_t \right] \qquad (3)$$

where $p(\tau)$ is a distribution over trajectories $\tau$ imposed by the policy. However, for a PoMDP we only have access to $o_t$ and not the underlying state $s_t$ during test time.

**Training with Privileged Learning:** In our experience, it is not straightforward to solve the given PoMDP using a naïve end-to-end learning approach. Inspired by the work of [2], we propose a solution that uses a privileged agent to first solve the problem and then distill the network into a student agent that can mimic the actions taken by the former.

Consider an alternate problem: *how should an agent with complete knowledge of the scene rearrange objects such that every object gets revealed during the episode?* The agent in this case gets access to the complete scene graph $S_t$ (instead of the partial graph $\hat{S}_t$) as input. Since it has access to the underlying state, we can formulate our problem as an MDP and solve it using traditional RL methods. We use a MPNN architecture to model our policy and train it with an actor-critic framework to maximize Equation 3. We call this agent the privileged Scene Exploration agent $\pi^*(a|s)$. We then collect a dataset $\mathcal{D}$ consisting of tuples $\langle \hat{S}, a^* \rangle$ by running the privileged policy. Subsequently, we use Behaviour Cloning to train the student agent $\pi(a|s)$ by minimizing the following objective.

$$L = - \sum_{a^*, \hat{S}_t \in \mathcal{D}} \log \pi(a^*|\hat{S}_t) \qquad (4)$$

where $a^*$ is the action taken by $\pi^*$ for state $S_t$.

Fig. 4: Examples of scenes generated by the (a) Scene Generation Agent trained on 7 objects (b) Scene generation Agent trained on 14 objects.

| Num. Nodes | Num. Objects Hidden | | | Stability Rate | | |
|---|---|---|---|---|---|---|
| | RG | SG | SG* | RG | SG | SG* |
| 10 | 1.13 | 2.12 | **4.62** | 0.84 | **0.99** | 0.97 |
| 15 | 1.64 | 2.41 | **5.88** | 0.64 | **0.99** | 0.85 |
| 20 | 2.18 | 2.46 | **6.67** | 0.55 | 0.93 | **0.97** |
| 25 | 2.51 | 3.01 | **7.45** | 0.41 | 0.81 | **0.88** |

TABLE I: Performance of the Scene Generation Agent (SG) with graph size = 15 nodes and Scene Generation Agent* (SG*) with graph size = 25 nodes, compared to Random Generation Agent (RG). Stability rate is defined as the fraction of simulated scenes that are physically stable.

This forces the student policy to discover patterns in the arrangement of *visible* objects, namely priors about physical plausibility, giving clues about where to find hidden objects. Once trained, the student policy operates with incomplete information about the scene and takes actions that have a high likelihood of uncovering hidden objects in the scene.

## IV. Experiments

We evaluate our method in simulation using iGibson [36] and Pybullet [37] and on real world clutter using a UR10 robot. We utilize objects from the YCB [38] dataset to create multi-object structured clutter that an interactive agent can explore. We train Scaled-YOLOv4 [39] network and use it as the object detector. To train the object detector, we first create a training dataset of simulated scenes by dropping 20 objects (randomly sampled from the set of objects defined previously) into a bin and generate 100k annotated images. We then augment this dataset with 10k images of scenes generated with the scene grammar. In each case, we capture an image of the simulated scene from the overhead RGBD camera and record the bounding box annotations for all objects in the scene. Our simulated scenes consist of seven objects from the YCB dataset: Cracker box, Pudding box, Master chef can, Tomato soup can, Sugar box, Gelatin box and Large marker. These objects vary both in shape and size, presenting abundant opportunities for objects to be hidden efficiently. We further demonstrate the scalability of our method, by running experiments with an extended version of this set that has 7 additional kitchen items (2×cereal, granola, pasta, lasagna, chips boxes and a pringles can). The Scene Generation Agent strategically creates stable scenes with cluttered arrangements of these objects, which the Scene Exploration Agent explores through rearrangement.

**Scene Generation agent:** The Scene Generation Agent takes a scene graph as input and outputs a distribution over

| Num. Nodes | Success Rate | | | | Num objects found | | | |
|---|---|---|---|---|---|---|---|---|
| | RE | LF | SE$^r$ | SE* | RE | LF | SE$^r$ | SE* |
| 10 | 0.05 | 0.88 | **0.92** | 0.91 | 0.76 | 2.93 | 3.06 | **3.13** |
| 15 | 0.03 | 0.58 | **0.84** | 0.82 | 1.15 | 3.33 | **4.40** | 4.27 |
| 20 | 0.00 | 0.32 | 0.36 | **0.64** | 0.93 | 3.64 | 4.24 | **4.71** |
| 25* | 0.00 | 0.10 | 0.52 | **0.91** | 0.90 | 2.70 | 4.87 | **5.82** |

TABLE II: Performance of Scene Exploration Agent (SE*) trained with scenes generated by SG*, compared to the Random Exploration Agent (RE), Largest First Agent (LF), and an agent (SE$^r$) trained with scenes generated by RG. SE* shows higher success rates, particularly for larger scenes.

| Object set | Objects hidden/found |
|---|---|
| 7 objects | 7.45/5.82 |
| 14 objects | 6.03/5.33 |

TABLE III: Scalability analysis: we test our approach on two object sets, containing 7 and 14 objects. Note that in each case, we train the Scene Generation agent to generate scene graphs with 25 nodes using the available object types.

the set of production rules $R$. Before feeding the scene graphs into a network, we first embed each node using a fully connected (FC) layer followed by a non-linearity (Leaky-ReLU), resulting in a vector of 32 dimensions. We use a similar network to embed the edges into a vector of size 3. The architecture of our policy is similar to the one used by [40] for BlockWorld, but with four message passing layers. We pass the output of the final layer through a FC layer to generate Q-values for all the rules in our graph grammar. We train the agent for a total of $640k$ steps starting with a learning rate of $10^{-3}$ and halving it every $128k$ steps.

We compare the performance of our trained Scene Generation (SG) agent to a Random Scene Generation policy (RG) that samples uniformly from the set of feasible production rules. To evaluate the performance of our scene generators, we define two metrics: *(i)* the number of objects they successfully hide; and, *(ii)* the stability rate, which indicates the probability of generating physically stable scenes. We train two agents SG (trained to generate graphs with 15 nodes) and SG* (trained to generate graphs with 25 nodes) and compare their performance to RG (Table I).

We find that both our agents outperform the baseline random agent for the graph sizes they were trained on. SG* can generate scene graphs with 25 nodes that contain an average of 7.45 objects ($\approx 5$ more than RE), with a stability rate of 88% (43% more than RE). Interestingly, we also observe that our agents can generalize to graph sizes they were not trained to generate. For example, even though SG has never been trained to generate graphs with $10, 20$ and 25 nodes, during test time, we can adjust the graph size to reliably generate stable cluttered scenes of different sizes. To test the scalability of our approach we train another Scene Generation agent on the extended version of the object set (containing 14 objects) as described earlier. We augment our grammar with an additional set of 19 rules (49 in total) to support these additional objects and structures. This agent hides 6.03 (stability rate = 70%) objects per scene when

Fig. 5: Example actions performed by the Scene Exploration Agent (a) In simulation (b) Real world. In each case, the agent rearranges the scene and discovers hidden objects (5 in the simulated scene and 2 in the real world clutter.)

compared to a random agent, that only hides an average of $2.44$ (stability rate $9\%$) objects. We observe that in spite of doubling the object set size, the performance of our method is not degraded (Table III). A few examples of generated scenes are shown in Figure 4.

**Scene Exploration agent:** We use architecture identical to [40] for both Privileged and Student Exploration agents. To evaluate the performance, we define the following metrics: *(i)* success rate (SR), which denotes the fraction of scenes where all objects were found; and, *(ii)* the average number of objects found per scene (OF). We first generate a dataset of $3k$ scene graphs each containing $25$ nodes using our Scene Generation agent (SG\*). We randomly sample a scene graph from this dataset to initialize the scene, but randomize the object locations on the tray.

We compare our Scene Exploration Agent (SE\*) to three baselines: *(i)* a Random Exploration Agent (RE), that randomly rearranges objects; *(ii)* the Largest First agent (LF) that moves the largest objects in the scene to look under them, as described in Danielczuk *et al.* [18]; and, *(iii)* the Scene Exploration Agent (SE$^r$) which is trained similar to SE, but on stable scenes generated by the random generator (RG) instead of the learned SG. We observe that the baselines could work for simpler scenes, but as complexity increases, our method provides a significant margin of improvement. Largest First agent (LF) performs well for graphs with $10$ or $15$ nodes, but fails on larger scenes. While SE$^r$ and SE\* perform comparably for small scenes, the importance of the Scene Generation Agent becomes evident with larger scenes. SE\* has a Success Rate of $91\%$ which is a $75\%$ relative improvement over SE$^r$(SR= $52\%$). To test the scalability of our method, we train another Scene Exploration agent on scenes generated by the $14$ objects version of the Scene Generation agent. This agent is able to find all hidden objects in $64\%$ of the scenes outperforming LF ($56\%$) and RE ($3\%$). As shown in Table III, our Exploration agent is able to reliably find objects despite doubling the object set size. An example exploration sequence is shown in Figure 5a.

**Real Robot Experiments:** We show evidence of sim-to-real transfer by testing the Scene Exploration agent (SE) in the real world. Our hardware experiments consist of a UR10 robot fitted with an e-pick Vacuum gripper and an Intel 435 Realsense camera mounted over the workspace to get a top-down view of the scene (Figure 6). We first process the input image captured by the camera using PoseRBPF [41] to detect the objects in the scene and get their respective poses. We then feed the partial graph $\hat{S}$ to SE and get the
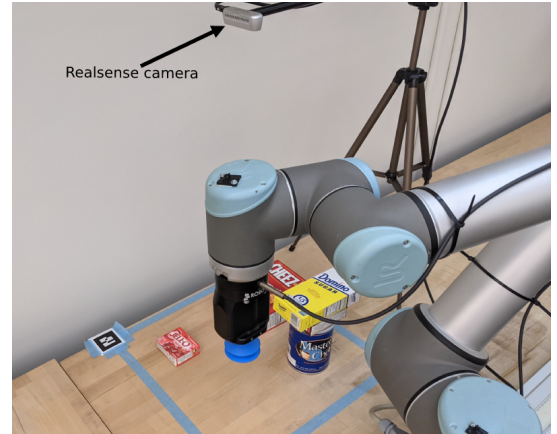


Fig. 6: We use a UR10 robot to execute the manipulation actions predicted by our Scene Exploration Agent (SE\*) and discover hidden objects. Note that the Realsense camera takes a top-down image of the scene and cannot objects hidden under the sugar box.

manipulation action that can reveal new objects. We repeat this process until all objects have been detected. We show the robot successfully solving 3 different cluttered scenes (Figure 5b) in the supplementary video.

## V. CONCLUSION

We investigate generating and interactively exploring structured clutter with multiple hidden objects. We present a graph-based framework to represent, generate and explore complex structured clutter scenes, while ensuring physical stability. We introduce two agents, a Scene Generation agent that skillfully hides objects and a Scene Exploration agent that rearranges the scene to find hidden objects, and demonstrated their effectiveness against multiple baselines. We showed that our agents generalize to scene sizes not encountered during training and provided additional experimental verification proving the scalability of our method. Finally, we test our Scene Exploration Agent with real world clutter, demonstrating successful sim-to-real transfer. For future work, we intend to test the scalability of our approach to a wider range of objects, in other cluttered scenarios [42].

## ACKNOWLEDGMENT

## References

[1] E. S. Spelke and K. D. Kinzler, "Core knowledge," *Developmental Science*, vol. 10, no. 1, pp. 89–96, jan 2007.

[2] D. Chen, B. Zhou, V. Koltun, and P. Krähenbühl, "Learning by cheating," in *Conference on Robot Learning*. PMLR, 2020, pp. 66–75.

[3] J. Zhou, G. Cui, S. Hu, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li, and M. Sun, "Graph neural networks: A review of methods and applications," *AI Open*, vol. 1, pp. 57–81, 2020.

[4] J. Bohg, K. Hausman, B. Sankaran, O. Brock, D. Kragic, S. Schaal, and G. S. Sukhatme, "Interactive perception: Leveraging action in perception and perception in action," *IEEE Transactions on Robotics*, vol. 33, no. 6, pp. 1273–1291, 2017.

[5] Z. Xu, J. Wu, A. Zeng, J. B. Tenenbaum, and S. Song, "Densephysnet: Learning dense physical object representations via multi-step dynamic interactions," *arXiv preprint arXiv:1906.03853*, 2019.

[6] K. N. Kumar, I. Essa, S. Ha, and C. K. Liu, "Estimating mass distribution of articulated objects using non-prehensile manipulation," *arXiv preprint arXiv:1907.03964*, 2019.

[7] T. Novkovic, R. Pautrat, F. Furrer, M. Breyer, R. Siegwart, and J. Nieto, "Object Finding in Cluttered Scenes Using Interactive Perception," in *Proceedings - IEEE ICRA*. Institute of Electrical and Electronics Engineers Inc., may 2020, pp. 8338–8344.

[8] J. Poon, Y. Cui, J. Ooga, A. Ogawa, and T. Matsubara, "Probabilistic active filtering for object search in clutter," in *Proceedings - IEEE ICRA*, vol. 2019-May. Institute of Electrical and Electronics Engineers Inc., may 2019, pp. 7256–7261.

[9] D. Katz and O. Brock, "Interactive segmentation of articulated objects in 3d," in *Workshop on mobile manipulation at ICRA*, vol. 2011, 2011.

[10] I. Sarantopoulos, M. Kiatos, Z. Doulgeri, and S. Malassiotis, "Split Deep Q-Learning for Robust Object Singulation," in *Proceedings - IEEE ICRA*. IEEE, sep 2020, pp. 6225–6231.

[11] L. Chang, J. R. Smith, and D. Fox, "Interactive singulation of objects from a pile," in *Proceedings - IEEE ICRA*. Institute of Electrical and Electronics Engineers Inc., 2012, pp. 3875–3882.

[12] A. Eitel, N. Hauff, and W. Burgard, "Learning to Singulate Objects Using a Push Proposal Network," jul 2020, pp. 405–419.

[13] A. Kurenkov, J. Taglic, R. Kulkarni, M. Dominguez-Kuhne, A. Garg, R. Martín-Martín, and S. Savarese, "Visuomotor Mechanical Search: Learning to Retrieve Target Objects in Clutter," aug 2020. [Online]. Available: http://arxiv.org/abs/2008.06073

[14] M. Danielczuk, A. Kurenkov, A. Balakrishna, M. Matl, D. Wang, R. Martin-Martin, A. Garg, S. Savarese, and K. Goldberg, "Mechanical search: Multi-step retrieval of a target object occluded by clutter," in *Proceedings - IEEE ICRA*, vol. 2019-May. Institute of Electrical and Electronics Engineers Inc., may 2019, pp. 1614–1621.

[15] A. Zeng, S. Song, S. Welker, J. Lee, A. Rodriguez, and T. Funkhouser, "Learning synergies between pushing and grasping with self-supervised deep reinforcement learning," in *IROS*. IEEE, 2018, pp. 4238–4245.

[16] Y. Deng, X. Guo, Y. Wei, K. Lu, B. Fang, D. Guo, H. Liu, and F. Sun, "Deep reinforcement learning for robotic pushing and picking in cluttered environment," in *2019 IROS*. IEEE, 2019, pp. 619–626.

[17] Y. Yang, H. Liang, and C. Choi, "A Deep Learning Approach to Grasping the Invisible," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 2232–2239, apr 2020.

[18] M. Danielczuk, A. Angelova, V. Vanhoucke, and K. Goldberg, "X-Ray: Mechanical Search for an Occluded Object by Minimizing Support of Learned Occupancy Distributions," apr 2020.

[19] J. K. Li, D. Hsu, and W. S. Lee, "Act to See and See to Act: POMDP planning for objects search in clutter," in *2016 IEEE/RSJ IROS*. IEEE, oct 2016.

[20] H. Huang, M. Dominguez-Kuhne, J. Ichnowski, V. Satish, M. Danielczuk, K. Sanders, A. Lee, A. Angelova, V. Vanhoucke, and K. Goldberg, "Mechanical search on shelves using lateral access x-ray," *arXiv preprint arXiv:2011.11696*, 2020.

[21] Z. Sui, L. Xiang, O. C. Jenkins, and K. Desingh, "Goal-directed robot manipulation through axiomatic scene estimation," *IJRR*, vol. 36, no. 1, pp. 86–104, 2017.

[22] Z. Sui, H. Chang, N. Xu, and O. C. Jenkins, "Geofusion: geometric consistency informed scene estimation in dense clutter," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 5913–5920, 2020.

[23] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, *et al.*, "Visual genome: Connecting language and vision using crowdsourced dense image annotations," *International journal of computer vision*, vol. 123, no. 1, pp. 32–73, 2017.

[24] T. Liu, S. Chaudhuri, V. G. Kim, Q. Huang, N. J. Mitra, and T. Funkhouser, "Creating consistent scene graphs using a probabilistic grammar," *ACM Transactions on Graphics (TOG)*, vol. 33, no. 6, pp. 1–12, 2014.

[25] K. Wang, Y.-A. Lin, B. Weissmann, M. Savva, A. X. Chang, and D. Ritchie, "Planit: Planning and instantiating indoor scenes with relation graph and spatial prior networks," *ACM Transactions on Graphics (TOG)*, vol. 38, no. 4, pp. 1–15, 2019.

[26] Y. Zhou, Z. While, and E. Kalogerakis, "Scenegraphnet: Neural message passing for 3d indoor scene augmentation," in *Proceedings of the IEEE/CVF ICCV*, 2019, pp. 7384–7392.

[27] J. Devaranjan, A. Kar, and S. Fidler, "Meta-sim2: Unsupervised learning of scene structure for synthetic data generation," in *European Conference on Computer Vision*. Springer, 2020, pp. 715–733.

[28] S. Qi, Y. Zhu, S. Huang, C. Jiang, and S.-C. Zhu, "Human-centric indoor scene synthesis using stochastic grammar," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5899–5908.

[29] N. Chomsky, *Syntactic structures*. De Gruyter Mouton, 2009.

[30] P. Purkait, C. Zach, and I. Reid, "Sg-vae: Scene grammar variational autoencoder to generate new indoor scenes," in *ECCV*. Springer, 2020, pp. 155–171.

[31] 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.

[32] A. Zeng, P. Florence, J. Tompson, S. Welker, J. Chien, M. Attarian, T. Armstrong, I. Krasin, D. Duong, V. Sindhwani, *et al.*, "Transporter networks: Rearranging the visual world for robotic manipulation," *arXiv preprint arXiv:2010.14406*, 2020.

[33] M. Danielczuk, A. Mousavian, C. Eppner, and D. Fox, "Object rearrangement using learned implicit collision functions," *arXiv preprint arXiv:2011.10726*, 2020.

[34] Y. Li, D. Tarlow, M. Brockschmidt, and R. Zemel, "Gated graph sequence neural networks," *arXiv preprint arXiv:1511.05493*, 2015.

[35] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine, "Reinforcement learning with deep energy-based policies," in *ICML*. PMLR, 2017, pp. 1352–1361.

[36] B. Shen, F. Xia, C. Li, R. Martín-Martín, L. Fan, G. Wang, S. Buch, C. D'Arpino, S. Srivastava, L. P. Tchapmi, *et al.*, "igibson, a simulation environment for interactive tasks in large realisticscenes," *arXiv preprint arXiv:2012.02924*, 2020.

[37] E. Coumans and Y. Bai, "PyBullet, a Python module for physics simulation for games, robotics and machine learning," \url{http://pybullet.org}.

[38] B. Calli, A. Singh, A. Walsman, S. Srinivasa, P. Abbeel, and A. M. Dollar, "The ycb object and model set: Towards common benchmarks for manipulation research," in *2015 ICAR*. IEEE, 2015, pp. 510–517.

[39] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao, "Scaled-YOLOv4: Scaling cross stage partial network," *arXiv preprint arXiv:2011.08036*, 2020.

[40] J. Janisch, T. Pevný, and V. Lisý, "Symbolic relational deep reinforcement learning based on graph neural networks," *arXiv preprint arXiv:2009.12462*, 2020.

[41] X. Deng, A. Mousavian, Y. Xiang, F. Xia, T. Bretl, and D. Fox, "Poserbpf: A rao–blackwellized particle filter for 6-d object pose tracking," *IEEE Transactions on Robotics*, 2021.

[42] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijmans, 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 ICCV*, 2019, pp. 9339–9347.