# SG-Bot: Object Rearrangement via Coarse-to-Fine Robotic Imagination on Scene Graphs

Guangyao Zhai[1], Xiaoni Cai[1], Dianye Huang[1], Yan Di[1,†]
Fabian Manhardt[2], Federico Tombari[1,2], Nassir Navab[1] and Benjamin Busam[1]
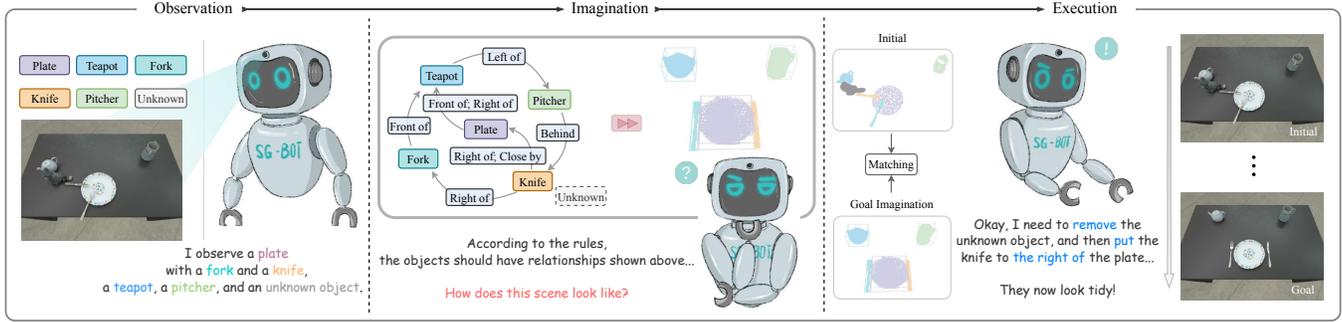
https://sites.google.com/view/sg-bot

Fig. 1. **SG-Bot workflow.** SG-Bot designs a three-fold procedure for robotic rearrangement. The cluttered initial scene is first perceived and processed into individual object nodes in *Observation*. Then, these objects are transitioned into the *Imagination* phase, where scene graph representation is adopted to facilitate the coarse-to-fine goal scene imagination, fusing all available priors and user commands. Finally, robotic action policies are generated in *Execution* by matching the initial and goal scenes. SG-Bot is lightweight, real-time, and controllable.

*Abstract*—**Object rearrangement is pivotal in robotic-environment interactions, representing a significant capability in embodied AI. In this paper, we present SG-Bot, a novel rearrangement framework that utilizes a coarse-to-fine scheme with a scene graph as the scene representation. Unlike previous methods that rely on either known goal priors or zero-shot large models, SG-Bot exemplifies lightweight, real-time, and user-controllable characteristics, seamlessly blending the consideration of commonsense knowledge with automatic generation capabilities. SG-Bot employs a three-fold procedure–observation, imagination, and execution–to adeptly address the task. Initially, objects are discerned and extracted from a cluttered scene during the observation. These objects are first coarsely organized and depicted within a scene graph, guided by either commonsense or user-defined criteria. Then, this scene graph subsequently informs a generative model, which forms a fine-grained goal scene considering the shape information from the initial scene and object semantics. Finally, for execution, the initial and envisioned goal scenes are matched to formulate robotic action policies. Experimental results demonstrate that SG-Bot outperforms competitors by a large margin.**

## I. INTRODUCTION

Object rearrangement is an essential but challenging task in robot-environment interaction, marking a crucial capability in embodied AI [1]. This interactive ability attains its zenith of automation by synergizing vision [2], [3], [4], [5], textual insights from sources [6], [7], [8], and strategic motion planning [9], [10]. Together, these elements culminate in a sophisticated physical embodiment for robots.

Robotic rearrangement refers to the process wherein a robotic agent, starting from an initial configuration within a scene, re-positions objects according to specific rules or instructions. The purpose is to achieve desired goal states, relying solely on sensory data and onboard perceptions. Recently proposed vision-based solutions to this task can be categorized into three approaches: **utilizing known geometric and semantic goal states**, **sequential object pose estimation**, and **zero-shot rearrangement with large models**. Typically, for goal-guided methods [11], [12], the quality of such priors significantly affects the accuracy of the rearrangement. When the goal state is unavailable, such methods become inapplicable for real-world use. Moreover, for pose estimation based approaches [13], while the sequential design aligns well with robotic manipulations, it can be affected by cumulative errors in autoregressive predictions. The last type of methods [14], [15], [16], [17], [18] tap into commonsense knowledge stored in zero-shot models. They necessitate either intricate post-filter procedures or prompt template designs, which tend to overlook scene-specific contextual cues and result in diverse undesired outcomes.

Orthogonal to the above methodologies, we explore a novel rearrangement routine embodied as *SG-Bot*, using goal imagination on scene graphs and goal-guided object matching as shown in Fig. 1. SG-Bot stacks three stages for the task, which are *observation*, *imagination*, and *execution*. Specifically, in the first stage, it processes initial scenes to extract objects by semantic instance segmentation. The imagination stage follows a coarse-to-fine solution, where objects are firstly treated as semantic nodes in a constructed goal scene graph. This graph is either directed by commonsense reasoning or user-defined rules, serving as coarse goal states. For a finer generation, the goal scene graph can already be decoded to an actual scene using a scene generative model,

arXiv:2309.12188v2 [cs.RO] 24 Mar 2024

Graph-to-3D [19]. However, inherited from the features of generative models, Graph-to-3D can produce diverse generation results inconsistent with the observation, potentially affecting the precision of subsequent object matching. We control the generation process by enriching the graph with shape priors to make a shape-aware graph, equipping the initial shape knowledge. Next, SG-Bot performs finer goal scene imagination conditioned on this graph, ensuring that the imagined shapes are coherent with the initial observation. Finally, in the execution stage, the imagined objects serve as anchors to guide the object matching by point cloud registration during the scene transformation. At each transformation step, we check occupancy between objects in the current observation and the imagination for safe rearrangement. The uniqueness of SG-Bot manifests in three aspects: First, SG-Bot does not need known goal priors but can self-generate goal scenes exclusively for the initial scenes, compared to the goal-required methods, *e.g.*, [11], [12]. Second, SG-Bot decouples the transformation policy using per-object matching to decrease the risk of error accumulation, compared to autoregressive methods, *e.g.*, [13]. Third, the concrete goal states and the closed-loop rearrangement strategy guarantee the rearrangement performance, compared to the loose-coupled zero-shot methods, *e.g.*, [16].

Our contributions are summarized as:

- We present *SG-Bot*, a new paradigm for the object rearrangement. The goal states are coarse-to-fine generated on the rules represented as scene graphs, with which goal-guided matching defines our motion policies.
- Ambiguous goal scene generation is alleviated by extracting shape priors from the initial observation. This leads to improved rearrangement performance.
- Experimental results in simulation show that SG-Bot can achieve competitive performance with state-of-the-art methods. Moreover, the rearrangement performance remains consistent in real-world scenarios.

## II. RELATED WORK

### A. Scene Graph

Scene graphs offer a rich symbolic and semantic representation of scenes [20], [21]. They can reason about objects and their relationships more explicitly than language [22]. This compact relationship description can be obtained through spatial grounding [23], [24], predicted from images [25], [26], [27], or even a GUI [28]. Scene graphs have applications in numerous computer vision areas such as 2D image generation [29], [22], manipulation [25], caption generation [30], camera localization [31], and 3D scene synthesis [32], [19], [33]. Recent robotics manipulation research also leverages scene graphs in planning [34], [35], [36]. In the context of this work, scene graphs serve to generate scenes, acting as anchors that guide the rearrangement.

### B. Object Rearrangement

The task necessitates that an embodied agent transition from initial states to goal states, adhering to specific rules based on perception and planning [1], as indicated by earlier

works [37], [38], [39], [40], [41]. By leveraging the development of visual perception [42], [43], [44], [45], [46], robotic grasping [47], [48], [49], motion planning [50], [51], [52], and research platforms [53], [54], [55], [56], [57], [58], a number of related methods have emerged. Solutions for this task fall into two categories. First, the goal states are given to the embodied agent, subsequently solving the problem by object matching, for example, using optical flow [11] or feature cosine similarity [12]. However, deriving such configurations can be challenging in real-world scenarios. Secondly, the goal states can be generated conditioned on the initial states. These goal states can be implicitly represented, such as by gradient fields [59], scene distributions [60], or sequential reasoning on the observation [13]. Alternatively, goals can be explicit in various formats, such as images [14] on prompts, bounding boxes [61] or poses [62] on descriptions, and direct language instructions [15], [63], [64], leveraging recent off-the-shelf large language models [65], [6], [66]. More powerful models even treat the initial-goal transformation as an end-to-end problem [67], [17], building on the large resource consumption. In this work, we generate the goal in a two-stage fashion, where coarse relationships are symbolized as a scene graph and finer concrete goals as the imagined scene given by the scene graph.

## III. PRELIMINARY

**Scene Graph.** The scene graph we use is semantic scene graph [20], denoted as $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$, which serves as a structured representation of a visual scene. In such representation, $\mathcal{V} = \{v_i \mid i = 1, \ldots, N\}$ refers to the set of object nodes, while $\mathcal{E} = \{e_{i \rightarrow j} \mid i, j = 1, \ldots, N, i \neq j\}$ represents the set of directed edges connecting each pair of nodes $v_i \rightarrow v_j$. As structured in the left of Fig. 3.b, each node $v_i$ can encompass various extensible attributes, *e.g.*, object category information $o_i \in O$, with $O$ containing all categories. As same as the node representation, each edge $e_{i \rightarrow j}$ is associated with a class label $\gamma_{i \rightarrow j} \in \Gamma$. In this paper, $\Gamma$ contains all pre-defined edge types, *i.e.*, {left/right, front/behind, standing on, close by}.

## IV. SG-BOT: OVERVIEW

### A. Problem Definition

From an initial layout state $\mathcal{S}_0$, the embodied agent is tasked with a sequential transformation of objects towards a desired goal state $\mathcal{S}^*$. This transformation is achieved by utilizing sequential motion policies $\mathcal{P}$, guided by sensory observations.

### B. Inference workflow

**Observation.** Given an RGB-D image capturing the initial object layout state $\mathcal{S}_0$, as shown in Fig. 2.a, SG-Bot first extracts all target objects as nodes $\mathcal{V}(O)$ via an arbitrary object detector, *e.g.*, MaskRCNN [68].

**Imagination.** The extracted object nodes are constructed as a scene graph $\mathcal{G}$ according to commonsense or user-defined rules, as shown in Fig. 2.b and explained in Sec. V-B. Next, we evolve $\mathcal{G}$ to a latent shape-aware scene graph $\mathcal{G}_z^\beta$ with
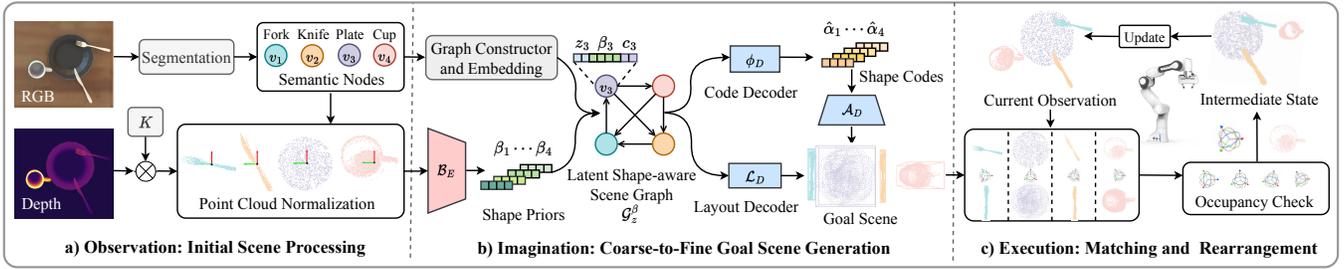
**Fig. 2. SG-Bot pipeline. a)** SG-Bot segments the input RGB image via MaskRCNN [68] to obtain individual object nodes $v_i$. Then, the corresponding point cloud of $v_i$ is obtained via back-projecting the depth map with camera intrinsics $K$. **b)** Coarse: the graph constructor connects each pair of nodes according to commonsense or user-defined rules, yielding scene graph $\mathcal{G}$. Fine: $\mathcal{G}$ is embedded and enhanced to $\mathcal{G}_z^\beta$ by combining estimated shape priors $\beta^*$ extracted from the normalized point clouds using the trained encoder $\mathcal{B}_E$ and latent code $z$ sampled from the learned layout-shape distribution. $\mathcal{G}_z^\beta$ then informs $\Phi_D$ and $\mathcal{L}_D$ of Graph-to-3D [19] to generate shape codes $\alpha^*$ and the scene layout respectively. $\alpha^*$ are decoded as shapes via $\mathcal{A}_D$, which are then populated in the layouts to form the goal scene. **c)** SG-Bot matches the initial and envisioned goal using point cloud registration and performs an occupancy check to determine the final movement in each step, as illustrated in V-E. The robot iteratively executes the action, transforming scenes into intermediate states and updating the observation until it reaches the goal state.

shape priors $\beta$ from the initial scene and learned layout-shape distribution $Z$ mentioned in Sec. V-C. Finally, SG-Bot imagines a goal scene $\mathcal{S}^*$ conditioned on $\mathcal{G}_z^\beta$ via the shape decoder $\Phi_D$ and layout decoder $\mathcal{L}_D$ of a scene generative model Graph-to-3D [19], where $\mathcal{S}^*$ comprises of dense point cloud and corresponding bounding box for each object.

**Execution.** Each target object in $\mathcal{S}_0$ is first extracted and represented as the back-projected point cloud from the depth map. Then, as shown in Fig. 2.c and explained in Sec. V-E, these objects are matched with the corresponding dense point clouds in $\mathcal{S}^*$ through iterative registration, *e.g.*, ICP [69], [70]. Based on the outcomes of this registration process, SG-Bot generates per-object manipulation policies $\mathcal{P}_t$ filtered and refined by object occupancy checking at each action step $t$. SG-Bot continues to iteratively reposition objects in $\mathcal{S}_0$ towards $\mathcal{S}^*$ until all objects are effectively rearranged.

## V. SG-BOT: METHODOLOGY

### A. Object Extraction

Given a cluttered scene $\mathcal{S}_0$ as the initial state, SG-Bot first performs semantic instance segmentation to segment all target objects, as shown in Fig. 2.a. Specifically, we adopt MaskRCNN to jointly predict the object masks and category labels. Then, each object is represented as the back-projected point cloud from the depth map. These objects, denoted as $\mathcal{V}(O) = \{v_i(o_i) \mid i = 1, \ldots, N\}$, are further collected and processed in the following *Imagination* module. This module aims to generate the desired goal scene by treating these objects as individual scene graph nodes.

After obtaining target objects $\mathcal{V}(O)$, we follow a coarse-to-fine scheme to generate the desired goal scene, which is leveraged to guide the object action.

### B. Coarse Stage: Goal Scene Graph Construction

SG-Bot establishes a goal scene graph $\mathcal{G} = \{\mathcal{V}(O), \mathcal{E}(\Gamma)\}$ via determining the edge type $\gamma_{i \to j} \in \Gamma$ for each edge in $\mathcal{E}(\Gamma)$, as shown in Fig. 2.b. In this paper, two modes are supported to define edges between nodes:

**Commonsense mode.** Following the recent trend of knowledge representation with graphs [71], we represent common

human knowledge in the form of edge attributes $\Gamma$ within a scene graph. For instance, for the scene containing a plate, the fork and knife are typically placed to the `left` `and` `right of` the plate. Additionally, the spoon needs to be placed `in front of` the plate if it exists. For the case without a plate, the spoon tends to be placed `close` `by` the bowl or cup. Moreover, other objects need to be placed `in front of` the plate, bowl, and cup, etc. Any unusual objects that appear on the table will be identified as obstacles and subsequently removed, which makes the final $M$ nodes from $N$ elements, $M \leq N$. Similar rules are naturally introduced based on the category of the object and commonsense. One way to achieve this is to use LLM to choose the optimal relationship according to the provided $\Gamma$.

**User-defined mode.** In contrast to the uncontrollable *Commonsense mode*, we demonstrate that one of the main advantages of introducing the scene graph representation is that it enables the controllable *User-defined mode*. Users can manipulate the scene graph from a long-term perspective, e.g., using a GUI, by directly editing the edges and nodes in $\mathcal{G}$ to interact with the edge database $\Gamma$ and nodes.

### C. Fine Stage: Graph to Scene Generation

SG-Bot stacks the architecture of Graph-to-3D [19] to generate a plausible goal scene. Graph-to-3D conditions on the latent shape-aware scene graph denoted as $\mathcal{G}_z^\beta$, which evolves from $\mathcal{G}$ and ensures the coherent shape transformation from the initial scene to the goal scene.

**Shape auto-encoders.** For this purpose, we first train two shape auto-encoder entities $\mathcal{A}, \mathcal{B}$ of AtlasNet [72] for different usages, as shown in Fig. 3.a. We train $\mathcal{A}(\mathcal{A}_E, \mathcal{A}_D)$ with full points under canonical view, whose encoder $\mathcal{A}_E$ offers shape codes $\alpha$ for training Graph-to-3D after. $\mathcal{B}(\mathcal{B}_E, \mathcal{B}_D)$ is trained with normalized object points under camera view in initial scenes to have initial shape priors $\beta$. The encoder $\mathcal{B}_E$ of $\mathcal{B}$ is preserved to produce $\beta$ during the training of Graph-to-3D and the final SG-Bot workflow. The training process of $\mathcal{A}, \mathcal{B}$ aligns with the original AtlasNet.

**Scene generative model.** After obtaining $\alpha$ and $\beta$, the training of Graph-to-3D starts with embedding $\mathcal{G}$ shown

in Fig. 3.b. The category information $c_i \in \mathcal{C}^{node}$ for $i$-th node is obtained by passing its textual information $o_i$ through node embedding layers $\mathcal{M}_O$, while $c_{i \to j} \in \mathcal{C}^{edge}$ is obtained by edge embedding layers $\mathcal{M}_\Gamma$ with $\gamma_{i \to j}$. Based on $\mathcal{G} \mapsto \mathcal{G} = \{\mathcal{V}(\mathcal{C}^{node}), \mathcal{E}(\mathcal{C}^{edge})\}$, Graph-to-3D, a subsequent dual-branch GCN architecture, is trained by modeling the layout-shape joint distribution $Z$ of goal scenes. As shown in Fig. 3.c, in training, the shape branch $\Phi(\Phi_E, \Phi_D)$ requires the graph to be augmented with ground truth shape codes $\alpha$ in goal scenes as input, whose output $\hat{\alpha}$ is supervised by the same shape codes. In the meantime, the layout branch $\mathcal{L}(\mathcal{L}_E, \mathcal{L}_D)$ takes the scene graph with ground truth bounding boxes $B = \{b_i \mid i = 1, .., M\}$ as input and the supervision labels. The two branches interact with each other in the bottleneck to model a latent graph $\mathcal{G}_z$, which shares the same idea of the concept of the latent code in the VAE [73]. $\mathcal{G}_z = \{\mathcal{V}(z, \mathcal{C}^{node}), \mathcal{E}(\mathcal{C}^{edge})\}$, consisting of $\mathcal{G}$ with sampled $z$ code from the modeled $Z$. More details can be found in [19]. Here, we change $\mathcal{G}_z$ as $\mathcal{G}_z^\beta$ by offering each node its shape prior $\beta$ extracted from its counterpart in the initial scene, i.e., $\mathcal{G}_z^\beta = \{\mathcal{V}(z, \beta, \mathcal{C}^{node}), \mathcal{E}(\mathcal{C}^{edge})\}$, to make $\hat{\alpha}$ and $\hat{b}$ aware of initial shapes.

**Controllable scene imagination.** After training, we subsequently engage in the process of generating the desired goal scene $\mathcal{S}^*$ conditioned on $\mathcal{G}_z^\beta$, shown in Fig. 2.b. This is accomplished through combination of code decoder $\Phi_D$, shape decoder $\mathcal{A}_D$, and layout decoder $\mathcal{L}_D$:

$$S = \mathcal{A}_D(\hat{\alpha}), \quad \hat{\alpha} = \Phi_D(\mathcal{G}_z^\beta), \quad \hat{\alpha} = \{\hat{\alpha}_i \mid i = 1, ..., M\}, \quad \text{(1a)}$$

$$\hat{B} = \mathcal{L}_D(\mathcal{G}_z^\beta), \quad \hat{B} = \{\hat{b}_i \mid i = 1, ..., M\}, \quad \text{(1b)}$$

where $\hat{\alpha}$ denotes the set of estimated shape codes, and $S$ is the set of normalized shapes decoded from $\hat{\alpha}$. $\hat{B}$ denotes the layout of object bounding boxes in the desired scene $\mathcal{S}^*$. $S$ then is transformed and populated into $\hat{B}$ to synthesize $\mathcal{S}^*$.

### D. Advantages of Coarse-to-Fine Scheme

SG-Bot features three key advantages: First, in the coarse stage, it utilizes a scene graph as an intermediary form of the target scene. This graph allows for multiple relationships between any two objects and enhances natural and intuitive human-computer interaction. Users can intuitively perceive the spatial distribution of objects within the scene through a 2D graphical scene graph, enabling direct editing through a GUI. Second, leveraging the scene graph as an intermediate representation allows for the seamless integration of commonsense knowledge, enabling automated scene rearrangement. Third, in the fine stage, we introduce the generative model to supplement missing fine-grained details, such as object shapes and poses, in the scene graph representation. This guides the robot in performing precise operations.

### E. Goal-Guided Object Matching and Manipulation

After obtaining $\mathcal{S}^*$, SG-Bot performs object matching by point cloud registration and rearranges objects after occupancy check in each round, as shown in Fig. 2.c, transferring $\mathcal{S}_0$ to $\mathcal{S}^*$. We illustrate the process with the first round:
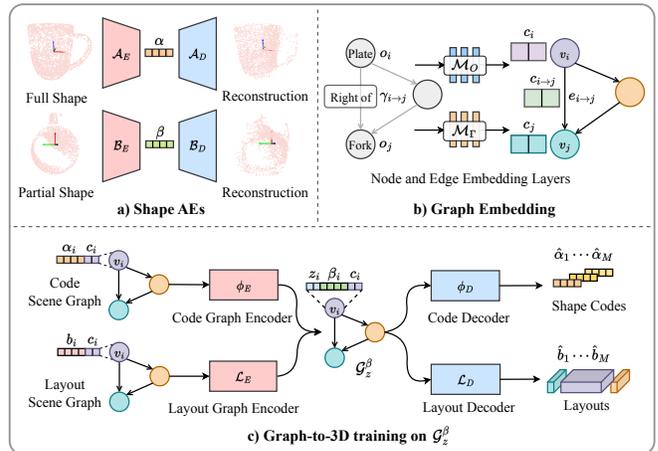


Fig. 3. **Modular Training. a)** $\mathcal{A}_E, \mathcal{A}_D$ are trained using full shapes in the canonical view to have the shape code $\alpha$, while $\mathcal{B}_E, \mathcal{B}_D$ are trained on partial shapes in the initial scenes under the camera view to have the shape priors $\beta$. $\mathcal{A}_D$ and $\mathcal{B}_E$ are retained during inference. **b)** A scene graph with textual information is processed through embedding layers $\mathcal{M}_O, \mathcal{M}_\Gamma$ to have implicit class features $c_i, c_{i \to j}$ on each node and edge. **c)** For training Graph-to-3D on goal scenes, the processed scene graph is first concatenated with $\alpha$ and a bounding box parameters $B$ on the shape box branch $\Phi(\Phi_E, \Phi_D)$ and layout branch $\mathcal{L}(\mathcal{L}_E, \mathcal{L}_D)$ respectively. $\Phi$ and $\mathcal{L}$ jointly model the layout-shape distribution $Z$ [19]. This model incorporates $\beta$ from initial scenes to create $\mathcal{G}_z^\beta$, subsequently estimating $\hat{\alpha}$ and $\hat{B}$. Modules in **b)** and **c)** are jointly trained, with $\mathcal{M}_O, \mathcal{M}_\Gamma, \Phi_D$ and $\mathcal{L}_D$ used during inference.

**Object matching.** SG-Bot compares $\mathcal{S}^*$ with the initial scene $\mathcal{S}_0$ to calculate the necessary transformation $\mathbf{T} = [\mathbf{R}|\mathbf{t}]$ for each object, where $\mathbf{R} \in \mathbb{R}^{3 \times 3}$ and $\mathbf{t} \in \mathbb{R}^3$ represent rotation and translation respectively. Therefore, in this module, the objective can be defined as,

$$[\mathbf{R}^*, \mathbf{t}^*] = \underset{\mathbf{R}, \mathbf{t}}{\arg\min} \sum_{i=1}^{N_P} (\min_{q \in Q} ||\mathbf{R}p_i + \mathbf{t} - q||^2) + I_{SO(3)}(\mathbf{R}), \quad \text{(2)}$$

where $\mathbf{R}^*$ and $\mathbf{t}^*$ represent the optimal rotation and translation parameters we aim to find. $p_i$ denotes one of the $N_P$ points in object $P$ of initial scene $\mathcal{S}_0$. After transforming $p_i$ from $\mathcal{S}_0$ to the goal scene $\mathcal{S}^*$ with $\mathbf{R}, \mathbf{t}$, its corresponding nearest point in $\mathcal{S}^*$ is denoted as $q$ inside object $Q$. $I_{SO(3)}(\mathbf{R})$ enforces $\mathbf{R}$ should lie in the special orthogonal group $SO(3)$ [74]. Since the generated objects in the goal scene are dense and complete, we observe that vanilla ICP can effectively solve the problem in Eq. 2 when provided with a well-suited initialization.

Given an object $P$ from the initial scene $\mathcal{S}_0$, its goal location is indicated by the generated object $Q$ in $\mathcal{S}^*$. We initialize the pose $\mathbf{T}$ by first centralizing each point cloud and then uniformly generating candidate rotations. We represent rotation using angles around the $x, y,$ and $z$ axes, dividing the interval of each axis's rotation angle $[-\pi, \pi]$ into $n$ segments, resulting in a total of $n^3$ candidate rotations, where $n = 5$ in the implementation. Finally, we apply ICP to estimate $\mathbf{R}^*, \mathbf{t}^*$, where $\mathbf{t}$ is initialized with $\mathbf{0}$ vector, while $\mathbf{R}$ is initialized with each candidate rotation. This will result in $n$ outcomes from ICP. We select the solution that minimizes Eq. 2 as the final result.

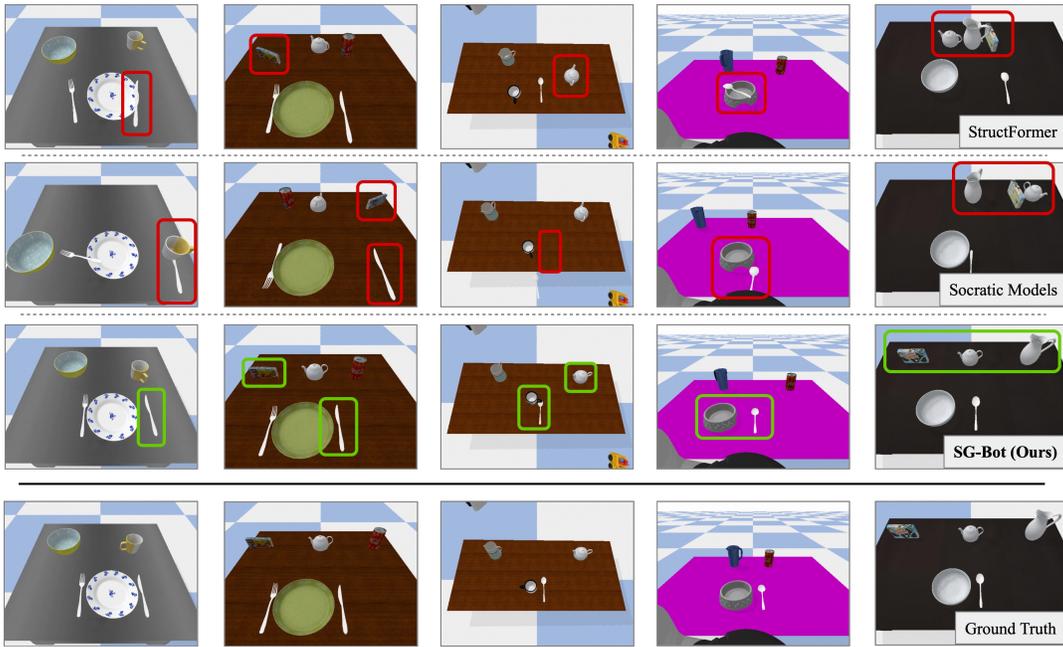**Object manipulation.** To determine the final robot action,

Fig. 4. **Visualization results in simulation.** We compare SG-Bot with state-of-the-art methods StructFormer [13] and Socratic Models [16]. We highlight the superiority of SG-Bot via rectangles.

TABLE I

PERFORMANCE EVALUATION ON THREE ASPECTS – ERRORS ($rad, cm$), SUCCESS RATE (%) AND FIDELITY.

| Method | Rearrangement Errors ($\downarrow$) | | | | Success Rate ($\uparrow$) | | Scene Fidelity ($\downarrow$) | |
|---|---|---|---|---|---|---|---|---|
| | $R_e$ | $t_e$ | $R_f$ | $t_f$ | IoU$_{0.25}$ | IoU$_{0.50}$ | FID | FID-CLIP |
| StructFormer [13] | **0.28** | 10.58 | 0.18 | 11.17 | 28.03 | 14.01 | 91.46 | 6.32 |
| Socratic Models [16] | – | 12.09 | – | 13.36 | 43.71 | **36.58** | 86.46 | 6.96 |
| **SG-Bot (Ours)** | 0.38 | **4.49** | **0.09** | **4.61** | **53.92** | 34.20 | **58.29** | **3.91** |

we select an object $P$ from $\mathcal{S}_0$ and check for occupancy: We measure the point-wise $L2$ distance between its counterpart $Q$ in $\mathcal{S}^*$, and all objects in $\mathcal{S}_0$. If the shortest distance $d$ is smaller than a set threshold $\sigma$, it implies a potential collision. We then bypass moving $P$ and evaluate the next object. This continues until an object with $d > \sigma$ is found, which is then moved to the target pose by its **T**.

The rearrangement ends in this manner when all objects are in their ideal poses.

## VI. EXPERIMENT

### A. Implementation Details

**Dataset.** We collect a synthetic dataset containing 1,042 realistic initial-goal RGB-D scene pairs with scene graph labels. First, we mix the meshes in Google Scanned Objects [75] and HouseCat6D [76] as the object database. Then, we randomly place objects on the tables to render the initial scenes into images using NVISII [77]. The goal scenes are set up using the rules mentioned in Sec. V-B. Then, we construct scene graph labels by comparing the spatial relations of the objects following [23], [33]. We define six types of relations as the edge class database $\Gamma$, including spatial, proximity, and support information, representing the *User-defined mode*.

**Trainval setup.** We use 952 scenes as the training split and 90 scenes as the validation (test) split. All modules in our

pipeline are trained on a single NVIDIA 3090 GPU. We adopt the Adam optimizer with an initial learning rate of 1e-4 to train each module. $\mathcal{A}$ is trained for 500 epochs on the meshes in the training split. $\mathcal{B}$ is trained for 5 epochs in terms of all partial points of each object in the training split. $\mathcal{M}_O, \mathcal{M}_\Gamma, \Phi, \mathcal{L}$ are jointly trained for 600 epochs.

### B. Evaluation Protocols

**Baselines.** We reproduce two methods representing different routines on the dataset for the comparison: *First,* StructFormer [13], a transformer-based method that autoregressively transforms objects to the goal state based on the current observation and previous states, is fully trained on our dataset. *Second,* Socratic Models [16], a LLM-based method that connects an object detection module [2], GPT [6], and a motion planning method CLIPort [78] in a series, where we use text-davinci-002 for LLM and train CLIPort solely using our dataset. All training and evaluation procedures use the same trainval splits as our method. More details about the reproduction can be found on our project website.

**Metrics.** *First,* for evaluating the rearrangement accuracy, we report the errors of estimated rotation $R_e$ and translation $t_e$ comparing final positions with ground truth following [13]. We also report the errors of final poses ($R_f, t_f$), as the final
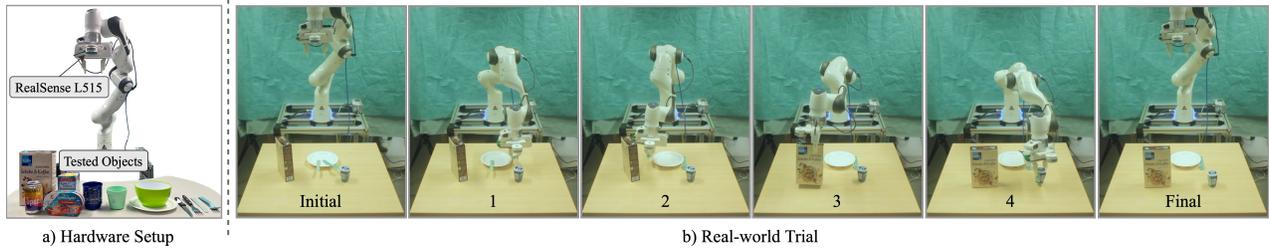
Fig. 5. **Real-world experiment.** a) We tested unseen cross-category objects with a physical manipulator. b) Action decomposition of one trial during the rearrangement.
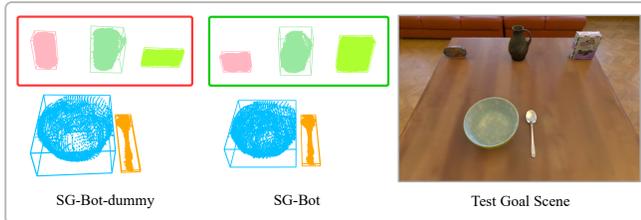


Fig. 6. **Functional shape priors.** Without shape priors, SG-Bot-dummy generates inconsistent shapes **(left)**. SG-Bot controls the generated shapes close to the ground truth (right) with the help of initial shape priors **(middle)**.

states of rearrangement are slightly different from the predicted ones because of the table-object physical interaction. *Second,* for the rearrangement success rate, we calculate the IoU between the bounding boxes of rearranged and ground truth objects. If IoU $> \sigma$, it counts as a success, $\sigma = 0.25, 0.50$. Note that this is a strict metric, as objects tend to be tiny, where even a small misalignment can cause failure. *Additionally,* inspired by some research on indoor scene synthesis [79], [80], [33], we believe that measuring the fidelity of the rearranged scene is critical for evaluating global performance. For this, we render rearranged scenes of all methods and ground truth scenes under a specific viewpoint, and then we employ the commonly adopted Fréchet Inception Distance (FID) [81] and recent FID-CLIP [82].

### C. Simulation Experiments

We import meshes with their initial poses to a PyBullet environment [41] to evaluate each method. In the simulation, we leverage ground truth instance masks and remove the effect of the robotic low-level control.

**Quantitative results.** As shown in Table I, our method surpasses the previous approaches on most metrics by a large margin. SG-Bot obtains lower rearrangement errors on the final states and yields competitive success rates, indicating that SG-Bot shows more accurate object-level rearrangement. For instance, SG-Bot decreases 50.0% on $R_{\mathrm{f}}$ and 58.7% on $t_{\mathrm{f}}$ compared with StructFormer [13]. When using $\mathrm{IoU}_{0.25}$, SG-Bot increases 10.21% on success rate compared with Socratic Models [16]. On the scene-level comparison, SG-Bot shows more fidelity in rearranged scenes than other methods, modeling a more similar scene distribution to ground truth supported by lower FID and FID-CLIP.

**Qualitative results.** We show several qualitative comparisons of rearranged scenes in Fig. 4, where our method shows clear advantages against others. For example, in the first

TABLE II
ABLATION − ERRORS ($rad, cm$), SUCCESS RATE (%) AND FIDELITY.

| Method | Errors (↓) | | Success Rate (↑) | | Scene Fidelity (↓) | |
|---|---|---|---|---|---|---|
| | $R_{\mathrm{f}}$ | $t_{\mathrm{f}}$ | $\mathrm{IoU}_{0.25}$ | $\mathrm{IoU}_{0.50}$ | FID | FID-CLIP |
| SG-Bot-dummy | 0.09 | 4.86 | 46.32 | 27.08 | 64.28 | 4.20 |
| SG-Bot | 0.09 | **4.61** | **53.92** | **34.20** | **58.29** | **3.91** |

scene, the rearranged knife collides with the plate or the cup in StructFormer and Socratic Models, which is better placed with our method. In the last scene, our method can separate objects at a sensible distance while others make them unevenly distributed.

**Ablation study.** We ablate the shape priors, resulting in *SG-Bot-dummy*, a framework only taking the original latent scene graph $\mathcal{G}_z$. As shown in Fig. 6, SG-Bot powered by $\mathcal{G}_z^{\beta}$ has more controllable ability than SG-Bot-dummy, generating more consistent shapes to the objects in the scenes. We also report quantitative comparisons in Table II.

### D. Real-world Experiments

We test SG-Bot in real-world scenarios using a 7-DoF Franka Panda robot with a parallel-jaw gripper as the end-effector. The sensor mounted on the gripper base is a RealSense L515 RGB-D camera. The framework is run on an NVIDIA 3080 laptop GPU. Different from the strategy in the simulation, we use Contact-GraspNet [48] to generate appropriate grasps on each masked object and rearrange them by reasoning the relative pose and executing the best grasp with Moveit! [83]. We show an example work stream in Fig. 5 out of 5 rounds where we test with unseen objects. More trials can be found on the project website. Our method can still maintain the rearrangement performance consistent with the one in the simulation.

## VII. CONCLUSIONS

In this paper, we present a novel robotic rearrangement framework, *SG-Bot*, which follows a three-phase procedure: observation, imagination, and execution to handle this task. With its unique coarse-to-fine design, SG-Bot embraces the synergy of commonsense priors and dynamic generation capabilities, all within a lightweight, real-time, and customizable pipeline. Extensive experiments on both simulation and real-world datasets demonstrate the superiority of SG-Bot. Future work will explore deformable point cloud matching for enhanced accuracy or accelerated point cloud alignment [84].

## 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," 2020. [Online]. Available: https://arxiv.org/abs/2011.01975

[2] X. Gu, T.-Y. Lin, W. Kuo, and Y. Cui, "Open-vocabulary object detection via vision and language knowledge distillation," in *ICLR*, 2022.

[3] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *ECCV*, 2020.

[4] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, "Segment anything," in *ICCV*, 2023.

[5] Y. Di, F. Manhardt, G. Wang, X. Ji, N. Navab, and F. Tombari, "So-pose: Exploiting self-occlusion for direct 6d pose estimation," in *ICCV*, 2021.

[6] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," in *NeurIPS*, 2020.

[7] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann *et al.*, "Palm: Scaling language modeling with pathways," 2022. [Online]. Available: https://arxiv.org/abs/2204.02311

[8] Y. Jiang, A. Gupta, Z. Zhang, G. Wang, Y. Dou, Y. Chen, L. Fei-Fei, A. Anandkumar, Y. Zhu, and L. Fan, "Vima: General robot manipulation with multimodal prompts," in *ICML*, 2023.

[9] C. Li, R. Zhang, J. Wong, C. Gokmen, S. Srivastava, R. Martín-Martín, C. Wang, G. Levine, M. Lingelbach, J. Sun *et al.*, "Behavior-1k: A benchmark for embodied ai with 1,000 everyday activities and realistic simulation," in *CoRL*, 2023.

[10] Y. Ding, X. Zhang, C. Paxton, and S. Zhang, "Task and motion planning with large language models for object rearrangement," 2023. [Online]. Available: https://arxiv.org/abs/2303.06247

[11] A. Goyal, A. Mousavian, C. Paxton, Y.-W. Chao, B. Okorn, J. Deng, and D. Fox, "Ifor: Iterative flow minimization for robotic object rearrangement," in *CVPR*, 2022.

[12] W. Goodwin, S. Vaze, I. Havoutis, and I. Posner, "Semantically grounded object matching for robust robotic scene rearrangement," in *ICRA*, 2022.

[13] W. Liu, C. Paxton, T. Hermans, and D. Fox, "Structformer: Learning spatial structure for language-guided semantic rearrangement of novel objects," in *ICRA*, 2022.

[14] I. Kapelyukh, V. Vosylius, and E. Johns, "Dall-e-bot: Introducing web-scale diffusion models to robotics," *RA-L*, vol. 8, no. 7, pp. 3956–3963, 2023.

[15] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan, K. Hausman *et al.*, "Do as i can, not as i say: Grounding language in robotic affordances," in *CoRL*, 2022.

[16] A. Zeng, M. Attarian, K. M. Choromanski, A. Wong, S. Welker, F. Tombari, A. Purohit, M. S. Ryoo, V. Sindhwani, J. Lee *et al.*, "Socratic models: Composing zero-shot multimodal reasoning with language," in *ICLR*, 2023.

[17] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn *et al.*, "Rt-2: Vision-language-action models transfer web knowledge to robotic control," in *CoRL*, 2023.

[18] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu *et al.*, "Palm-e: An embodied multimodal language model," in *ICML*, 2023.

[19] H. Dhamo, F. Manhardt, N. Navab, and F. Tombari, "Graph-to-3d: End-to-end generation and manipulation of 3d scenes using scene graphs," in *ICCV*, 2021.

[20] X. Chang, P. Ren, P. Xu, Z. Li, X. Chen, and A. Hauptmann, "A comprehensive survey of scene graphs: Generation and application," *T-PAMI*, vol. 45, no. 1, pp. 1–26, 2021.

[21] J. Johnson, R. Krishna, M. Stark, L.-J. Li, D. Shamma, M. Bernstein, and L. Fei-Fei, "Image retrieval using scene graphs," in *CVPR*, 2015.

[22] J. Johnson, A. Gupta, and L. Fei-Fei, "Image generation from scene graphs," in *CVPR*, 2018.

[23] J. Wald, H. Dhamo, N. Navab, and F. Tombari, "Learning 3d semantic scene graphs from 3d indoor reconstructions," in *CVPR*, 2020.

[24] 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 *ICCV*, 2019.

[25] H. Dhamo, A. Farshad, I. Laina, N. Navab, G. D. Hager, F. Tombari, and C. Rupprecht, "Semantic image manipulation using scene graphs," in *CVPR*, 2020.

[26] S.-C. Wu, J. Wald, K. Tateno, N. Navab, and F. Tombari, "Scenegraph-fusion: Incremental 3d scene graph prediction from rgb-d sequences," in *CVPR*, 2021.

[27] D. Xu, Y. Zhu, C. B. Choy, and L. Fei-Fei, "Scene graph generation by iterative message passing," in *CVPR*, 2017.

[28] M. Fessenden, "Scenegraph," https://github.com/mfessenden/SceneGraph, 2017.

[29] L. Yang, Z. Huang, Y. Song, S. Hong, G. Li, W. Zhang, B. Cui, B. Ghanem, and M.-H. Yang, "Diffusion-based scene graph to image generation with masked contrastive pre-training," 2022. [Online]. Available: https://arxiv.org/abs/2211.11138

[30] 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," *IJCV*, vol. 123, pp. 32–73, 2017.

[31] A. Rosinol, A. Violette, M. Abate, N. Hughes, Y. Chang, J. Shi, A. Gupta, and L. Carlone, "Kimera: From slam to spatial perception with 3d dynamic scene graphs," *IJRR*, vol. 40, no. 12-14, pp. 1510–1546, 2021.

[32] A. Luo, Z. Zhang, J. Wu, and J. B. Tenenbaum, "End-to-end optimization of scene layout," in *CVPR*, 2020.

[33] G. Zhai, E. P. Örnek, S.-C. Wu, Y. Di, F. Tombari, N. Navab, and B. Busam, "Commonscenes: Generating commonsense 3d indoor scenes with scene graphs," in *NeurIPS*, 2023.

[34] B. Tang and G. S. Sukhatme, "Selective object rearrangement in clutter," in *CoRL*, 2023.

[35] Y. Zhu, J. Tremblay, S. Birchfield, and Y. Zhu, "Hierarchical planning for long-horizon manipulation with geometric and symbolic scene graphs," in *ICRA*, 2021.

[36] K. Rana, J. Haviland, S. Garg, J. Abou-Chakra, I. Reid, and N. Suenderhauf, "Sayplan: Grounding large language models using 3d scene graphs for scalable task planning," *CoRL*, 2023.

[37] A. Cosgun, T. Hermans, V. Emeli, and M. Stilman, "Push planning for object placement on cluttered table surfaces," in *IROS*, 2011.

[38] J. E. King, M. Cognetti, and S. S. Srinivasa, "Rearrangement planning using object-centric and robot-centric action spaces," in *ICRA*, 2016.

[39] J. E. King, V. Ranganeni, and S. S. Srinivasa, "Unobservable monte carlo planning for nonprehensile rearrangement tasks," in *ICRA*, 2017.

[40] J. Lee, Y. Cho, C. Nam, J. Park, and C. Kim, "Efficient obstacle rearrangement for object manipulation tasks in cluttered environments," in *ICRA*, 2019.

[41] E. Coumans and Y. Bai, "Pybullet, a python module for physics simulation for games, robotics and machine learning," http://pybullet.org, 2016–2021.

[42] C. Wang, D. Xu, Y. Zhu, R. Martín-Martín, C. Lu, L. Fei-Fei, and S. Savarese, "Densefusion: 6d object pose estimation by iterative dense fusion," in *CVPR*, 2019.

[43] S. Peng, Y. Liu, Q. Huang, X. Zhou, and H. Bao, "Pvnet: Pixel-wise voting network for 6dof pose estimation," in *CVPR*, 2019.

[44] F. Manhardt, D. M. Arroyo, C. Rupprecht, B. Busam, T. Birdal, N. Navab, and F. Tombari, "Explaining the ambiguity of object detection and 6d pose from visual data," in *ICCV*, 2019.

[45] R. Zhang, Y. Di, F. Manhardt, F. Tombari, and X. Ji, "Ssp-pose: Symmetry-aware shape prior deformation for direct category-level object pose estimation," in *IROS*, 2022.

[46] Y. Di, R. Zhang, Z. Lou, F. Manhardt, X. Ji, N. Navab, and F. Tombari, "Gpv-pose: Category-level object pose estimation via geometry-guided point-wise voting," in *CVPR*, 2022.

[47] G. Zhai, Y. Zheng, Z. Xu, X. Kong, Y. Liu, B. Busam, Y. Ren, N. Navab, and Z. Zhang, "Da$^2$ dataset: Toward dexterity-aware dual-arm grasping," *RA-L*, vol. 7, no. 4, pp. 8941–8948, 2022.

[48] M. Sundermeyer, A. Mousavian, R. Triebel, and D. Fox, "Contact-graspnet: Efficient 6-dof grasp generation in cluttered scenes," in *ICRA*, 2021.

[49] G. Zhai, D. Huang, S.-C. Wu, H. Jung, Y. Di, F. Manhardt, F. Tombari, N. Navab, and B. Busam, "Monograspnet: 6-dof grasping with a single rgb image," in *ICRA*, 2023.

[50] R. Wang, K. Gao, D. Nakhimovich, J. Yu, and K. E. Bekris, "Uniform object rearrangement: From complete monotone primitives to efficient non-monotone informed search," in *ICRA*, 2021.

[51] S. H. Cheong, B. Y. Cho, J. Lee, C. Kim, and C. Nam, "Where to relocate?: Object rearrangement inside cluttered and confined environments for robotic manipulation," in *ICRA*, 2020.

[52] K. Gao, D. Lau, B. Huang, K. E. Bekris, and J. Yu, "Fast high-quality tabletop rearrangement in bounded workspace," in *ICRA*, 2022.

[53] 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 *ICCV*, 2019.

[54] A. Szot, A. Clegg, E. Undersander, E. Wijmans, Y. Zhao, J. Turner, N. Maestre, M. Mukadam, D. S. Chaplot, O. Maksymets *et al.*, "Habitat 2.0: Training home assistants to rearrange their habitat," in *NeurIPS*, 2021.

[55] E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Herrasti, M. Deitke, K. Ehsani, D. Gordon, Y. Zhu *et al.*, "Ai2-thor: An interactive 3d environment for visual ai," 2017. [Online]. Available: https://arxiv.org/abs/1712.05474

[56] S. James, Z. Ma, D. R. Arrojo, and A. J. Davison, "Rlbench: The robot learning benchmark & learning environment," *RA-L*, vol. 5, no. 2, pp. 3019–3026, 2020.

[57] F. Xiang, Y. Qin, K. Mo, Y. Xia, H. Zhu, F. Liu, M. Liu, H. Jiang, Y. Yuan, H. Wang *et al.*, "Sapien: A simulated part-based interactive environment," in *CVPR*, 2020.

[58] B. Shen, F. Xia, C. Li, R. Martín-Martín, L. Fan, G. Wang, C. Pérez-D'Arpino, S. Buch, S. Srivastava, L. Tchapmi *et al.*, "igibson 1.0: A simulation environment for interactive tasks in large realistic scenes," in *IROS*, 2021.

[59] M. Wu, F. Zhong, Y. Xia, and H. Dong, "Targf: Learning target gradient field to rearrange objects without explicit goal specification," in *NeurIPS*, 2022.

[60] Q. A. Wei, S. Ding, J. J. Park, R. Sajnani, A. Poulenard, S. Sridhar, and L. Guibas, "Lego-net: Learning regular rearrangements of objects in rooms," in *CVPR*, 2023.

[61] N. Gkanatsios, A. Jain, Z. Xian, Y. Zhang, C. Atkeson, and K. Fragki-adaki, "Energy-based models as zero-shot planners for compositional scene rearrangement," in *RSS*, 2023.

[62] I. Kapelyukh, Y. Ren, I. Alzugaray, and E. Johns, "Dream2real: Zero-shot 3d object rearrangement with vision-language models," in *ICRA*, 2024.

[63] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu, J. Tremblay, D. Fox, J. Thomason, and A. Garg, "Progprompt: Generating situated robot task plans using large language models," in *ICRA*, 2023.

[64] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng, "Code as policies: Language model programs for embodied control," in *ICRA*, 2023.

[65] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, "Hierarchical text-conditional image generation with clip latents," 2022. [Online]. Available: https://arxiv.org/abs/2204.06125

[66] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. d. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman *et al.*, "Evaluating large language models trained on code," 2021. [Online]. Available: https://arxiv.org/abs/2107.03374

[67] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu *et al.*, "Rt-1: Robotics transformer for real-world control at scale," 2022. [Online]. Available: https://arxiv.org/abs/2212.06817

[68] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *ICCV*, 2017.

[69] P. J. Besl and N. D. McKay, "Method for registration of 3-d shapes," in *Sensor fusion IV: control paradigms and data structures*, vol. 1611. Spie, 1992, pp. 586–606.

[70] Z. Zhang, "Iterative point matching for registration of free-form curves and surfaces," *IJCV*, vol. 13, no. 2, pp. 119–152, 1994.

[71] X. Chen, S. Jia, and Y. Xiang, "A review: Knowledge reasoning over knowledge graph," *Expert Systems with Applications*, vol. 141, p. 112948, 2020.

[72] T. Groueix, M. Fisher, V. G. Kim, B. C. Russell, and M. Aubry, "A papier-mâché approach to learning 3d surface generation," in *CVPR*, 2018.

[73] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," in *ICLR*, 2014.

[74] J. Zhang, Y. Yao, and B. Deng, "Fast and robust iterative closest point," *T-PAMI*, vol. 44, no. 7, pp. 3450–3466, 2021.

[75] L. Downs, A. Francis, N. Koenig, B. Kinman, R. Hickman, K. Reymann, T. B. McHugh, and V. Vanhoucke, "Google scanned objects: A high-quality dataset of 3d scanned household items," in *ICRA*, 2022.

[76] H. Jung, G. Zhai, S.-C. Wu, P. Ruhkamp, H. Schieber, P. Wang, G. Rizzoli, H. Zhao, S. D. Meier, D. Roth, N. Navab *et al.*, "Housecat6d–a large-scale multi-modal category level 6d object perception dataset with household objects in realistic scenarios," 2022. [Online]. Available: https://arxiv.org/abs/2212.10428

[77] N. Morrical, J. Tremblay, S. Birchfield, and I. Wald, "NVISII: Nvidia scene imaging interface," 2020, https://github.com/owl-project/NVISII/.

[78] M. Shridhar, L. Manuelli, and D. Fox, "Cliport: What and where pathways for robotic manipulation," in *CoRL*, 2022.

[79] D. Ritchie, K. Wang, and Y.-a. Lin, "Fast and flexible indoor scene synthesis via deep convolutional generative models," in *CVPR*, 2019.

[80] D. Paschalidou, A. Kar, M. Shugrina, K. Kreis, A. Geiger, and S. Fidler, "Atiss: Autoregressive transformers for indoor scene synthesis," in *NeurIPS*, 2021.

[81] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "Gans trained by a two time-scale update rule converge to a local nash equilibrium," in *NeurIPS*, 2017.

[82] T. Kynkäänniemi, T. Karras, M. Aittala, T. Aila, and J. Lehtinen, "The role of imagenet classes in fréchet inception distance," in *ICLR*, 2023.

[83] D. Coleman, I. Sucan, S. Chitta, and N. Correll, "Reducing the barrier to entry of complex robotic software: a moveit! case study," 2014. [Online]. Available: https://arxiv.org/abs/1404.3785

[84] E. Malis, "Complete closed-form and accurate solution to pose estimation from 3d correspondences," *RA-L*, vol. 8, no. 3, pp. 1786–1793, 2023.