

Scene Editing as Teleoperation: A Case Study in 6DoF Kit Assembly

Yulong Li^{*1}, Shubham Agrawal^{*1,2}, Jen-Shuo Liu¹, Steven K. Feiner¹, Shuran Song¹

¹Columbia University ²Samsung AI Center NY

<https://seat.cs.columbia.edu/>

Abstract—Studies in robot teleoperation have been centered around action specifications—from continuous joint control to discrete end-effector pose control. However, these “robot-centric” interfaces often require skilled operators with extensive robotics expertise. To make teleoperation accessible to non-expert users, we propose the framework “Scene Editing as Teleoperation” (SEaT), where the key idea is to transform the traditional “robot-centric” interface into a “scene-centric” interface—instead of controlling the robot, users focus on specifying the task’s goal by manipulating digital twins of the real-world objects. As a result, a user can perform teleoperation without any expert knowledge of the robot hardware. To achieve this goal, we utilize a category-agnostic scene-completion algorithm that translates the real-world workspace (with unknown objects) into a manipulable virtual scene representation and an action-snapping algorithm that refines the user input before generating the robot’s action plan. To train the algorithms, we procedurally generated a large-scale, diverse kit-assembly dataset that contains object-kit pairs that mimic real-world object-kitting tasks. Our experiments in simulation and on a real-world system demonstrate that our framework improves both the efficiency and success rate for 6DoF kit-assembly tasks. A user study demonstrates that SEaT framework participants achieve a higher task success rate and report a lower subjective workload compared to an alternative robot-centric interface.

I. INTRODUCTION

The vast majority of robot-teleoperation research has focused on how to better specify robot actions: from continuous joint control to discrete end-effector pose control. However, most of these “robot-centric” interfaces require skilled operators (with robotics expertise), complex input devices, or low-latency connections, which are hard to guarantee in practice.

To address these issues, we propose the framework of “Scene Editing as Teleoperation” (SEaT), where the key idea is to transform the traditional *robot-centric* interface into a *scene-centric* interface—instead of specifying robot actions, users focus on specifying task goals by manipulating digital twins of real-world objects. As a result, non-expert users, users who have a high-level understanding of the task but no experience of working with the robot, can perform teleoperation without knowledge of the robot hardware, control mechanisms, or current state—users do not even see the robot during teleoperation. In addition, by removing the need of continuous control, the system is able to gracefully handle variable network latency.

While SEaT is applicable for general “object rearrangement” tasks, we use 6DoF unknown object kit assembly as the case study in this paper. This task is selected because of its high requirements in precision and flexibility. Through

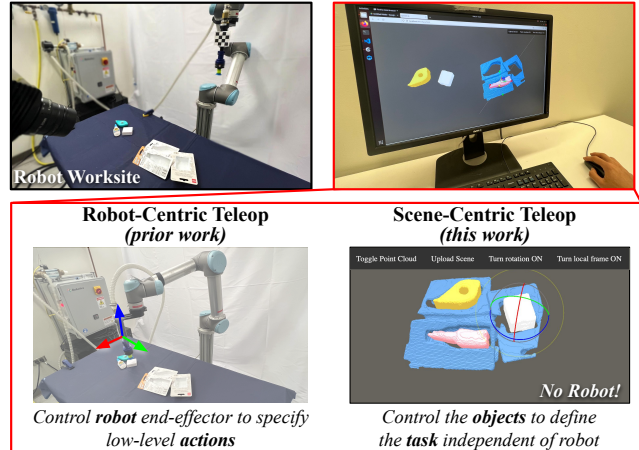


Fig. 1. **Scene Editing as Teleoperation.** With a scene-centric interface, our framework allows the user to efficiently specify the task goal without expert knowledge of the robot hardware or control, making this framework accessible to non-expert users. By removing the need for continuous control, the system is able to gracefully handle variable network latency.

this task, we hope to demonstrate the useful capabilities of SEaT that could not be achieved by either a traditional teleoperation system (struggles to produce precise actions in 6DoF space [1]) or an automated system (struggles to generalize to new objects and tasks [2]).

While there are many existing “scene editing” tools for manipulating virtual objects [3]–[5], the decisive challenge for our task is how to reliably translate between the real and virtual scene representations, specifically:

- How to translate the realworld workspace filled with *unknown* objects into an editable virtual scene.
- How to translate *imprecise* user edits (i.e., objects’ rearrangements) to the realworld with the robot’s actions.

To obtain the digital twins of unknown objects, we propose a category-agnostic scene-completion algorithm that segment and complete individual objects from depth images. To handle imprecise user inputs, we propose a 6DoF action-snapping algorithm that automatically refines user inputs and corrects object-kit alignment using a 3D shape matching network. Finally, virtual operations on object poses are translated by a sequence of robot actions generated by the robot planner. Learning from a large-scale kit-assembly dataset, our framework (both scene-completion and action-snapping algorithms) can generalize to unseen object-kit pairs, allowing quick adaptation to new assembly tasks.

In summary, our primary contribution is the framework of SEaT that allows non-expert end users to perform complex and precise 6DoF kit-assembly tasks over a high-latency

* indicates equal contributions

internet connection. This framework is enabled by the following technical contributions:

- A category-agnostic scene-completion algorithm that translates the real-world robot workspace (with unknown objects) into a virtual editable scene representation.
- An action-snapping algorithm that automatically refines user inputs and improves object-kit alignment using a 3D shape matching network.
- A large-scale kit-assembly dataset, KIT1000, that contains a diverse set of procedurally generated object-kit pairs that mimic real-world kitting tasks. This diverse training data allows the algorithm to generalize to new objects.

Extensive experiments suggest that SEaT improves both the efficiency and success rate of 6DoF kit-assembly tasks, while achieving a lower subjective workload compared to an alternative robot-centric interface. Please see our [project website](#) for more system videos. Code and data will be made publicly available.

II. RELATED WORK

Teleoperation. Early investigations in robot teleoperation focused on specifying a continuous motion trajectory [6]–[14], which often requires a low-latency connection between the teleoperator and robot or a complex input device for the operator. To reduce these requirements, other systems allow the operator to specify only the robot end-effector target poses [1], [15]–[17], and allow asynchronous execution to mitigate high communication latency. However, regardless of the levels of control, all these systems still focus on specifying the **robot’s action**, requiring expert teleoperators with knowledge and intuition of the robot embodiment. For example, the user needs to understand the robot kinematics to specify a reachable and collision-free arm trajectory or understand the robot gripper mechanism to specify a valid grasp pose. Training human operators with this expertise can be expensive and difficult to scale. In contrast, our system focus on specifying the **task goal** regardless of robot hardware. This idea of task-driven teleoperation has been studied in simple scenarios such as point-goal navigation [18] or manipulation with known objects [19]. However, how to enable precise and efficient task specification for a complex assembly task with unknown object parts is still an open research question, hence the focus of this paper.

Vision-based kit assembly. Traditional vision-based assembly approaches require strong prior knowledge of target objects (e.g., detailed CAD models) to perform object-pose estimation and motion planning [2], [20]. As a result, these approaches often cannot generalize to new objects without extensive data collection. Recent methods explore the idea of shape-informed assembly [2], [21], [22], where the task of assembly is formulated as a shape-matching problem between the object and its target location. This formulation allows the algorithms to generalize toward unseen objects by directly analyzing their 3D geometry. However, these algorithms are still limited to simpler tasks, such as 3DoF assembly [21], only predicting single object assembly [2], [22], only rotation prediction [2] or require precise demonstrations

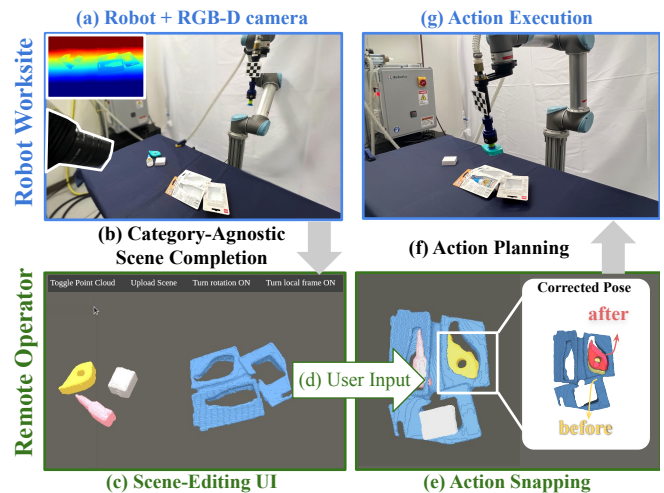


Fig. 2. **Overview.** Given a depth image, the scene-completion algorithm converts the workspace into a virtual scene (a–b §III-A). The user then specifies a target object pose by editing the virtual scene using the 3D UI (c–d, §III-B). Our action-snapping algorithm refines the object pose to improve object-kit alignment (e, §III-C). Finally, the system computes and executes the actions to assemble the objects (f–g, §III-D).

on the exact object-kit pair [22]. While top-down kits (3DoF assembly) may seem ubiquitous, most do not have a flat bottom and hence cannot stand vertically on their own on an assembly belt. Handling multiple objects simultaneously is required for kitting tasks involving packaging multiple related objects together (e.g., toothpaste and toothbrush or bundle of pens). Our approach is able to handle multi-unknown-object 6DoF kitting from imprecise user input, where user input helps reduce potential ambiguities and reduce search space, and the 3D shape-matching network further refines imprecise user input.

Creating digital twins of 3D scenes. Many 3D scene-understanding algorithms have been developed to produce high-quality digital models of real-world environments for teleoperation. These include algorithms for 3D object detection [23]–[31] and shape completion [32]–[37]. Unlike traditional 3D scene-understanding tasks that focus on *common* object categories (e.g., tables or chairs), in assembly tasks, a system often encounters a large number of new objects and parts that cannot be categorized into predefined categories. To address this issue, we propose a category-agnostic scene-completion algorithm that generalizes to unseen objects or parts without their 3D CAD model, allowing quick adaptation to new assembly tasks.

III. METHOD: SCENE EDITING AS TELEOPERATION

We study the task of **6DoF kit-assembly with multiple unknown objects**. To perform the task, the robot need to precisely place the object into their corresponding kit location with correct 6DoF poses. This task presents a set of unique challenges compared to general object rearrangement tasks: 1) High precision requirement – making it particularly challenging for human teleoperators with single view observation, hence, motivates our action snapping network with shape completed objects. 2) Ambiguities in object-kit correspondence. The ambiguities can be caused by similar or

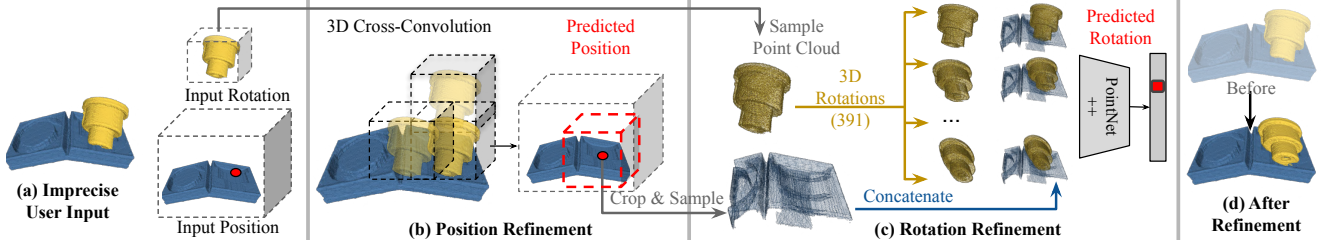


Fig. 3. **6DoF Action Snapping with SnapNet.** SnapNet uses 3D shape matching to refine the alignments between objects and their kits. Given the user’s imprecise input (a), the algorithm first refines the object position by using a 3D cross-convolution network between the geometric features computed from the object and kit volume (b). The cross-convolution is computed only in the local area around the user inputs (b). The algorithm then samples point-clouds from the object volume and the cropped kit volume centered at the predicted position and predicts the refined rotation from 391 rotations using a PointNet++ based classifier (c). Finally the algorithm outputs the refined position and rotation as the target pose.

symmetrical shapes, requiring human inputs to disambiguate. 3) Large search space — compared to top-down kit-assembly tasks [21], the possible object poses in 6DoF is significantly higher, making uniform search approach impractical. 4) Despite the ubiquity of the kit-assembly applications, a large-scale dataset is not yet available for this task, which is a key bottleneck for enabling learning-based approaches. In the following sections, we will discuss our approach to address above challenges.

A. Category-Agnostic Scene Completion

Given a single depth image I of the workspace with objects on one side and the kit on the other, the algorithm generates shape-completed geometries for individual objects using the following two steps:

Object-Instance Segmentation: The algorithm first detects and segments all object instances using SD-MaskRCNN [38]: a variant of MaskRCNN [39] using only depth for better sim2real generalization. Since the assembly task involves a large number of object parts that cannot be categorized into predefined categories, we train this algorithm in a category-agnostic manner with only a binary objectness label.

3D Shape Completion: Given an object’s instance mask M and the depth image I , the algorithm estimates the object’s full 3D geometry. This shape-completion step provides two benefits: 1) it aids the user during teleoperation by better visualization of the objects and provides more context for successful kitting, and 2) it helps in achieving better action-snapping results as shown in Tab. I.

To perform 3D shape completion, we first transform partial object geometry information from masked depth image $MD = I \times M$ into a 128^3 TSDF volume [40] representation $V_{partial}$ with voxel size 0.89 mm. This volume is then fed into our shape-completion network SC_{θ} to obtain the shape-completed 3D volume $V_{completed}$. SC_{θ} follows a 3D encoder-decoder style architecture with skip connections [34]. The network is trained to minimize voxel-wise MSE loss. We train a separate network for kits with same architecture as for object shape completion.

Both models are trained on the simulation data generated from objects and kits from our dataset (see §III-E) and then directly tested on unseen real world data.

B. Scene-Editing Interface

Given the 3D models for each object, the next step is to specify the task goal by changing their 3D poses in a

virtual scene. This interface (Fig. 2 c) is implemented as a 3D UI in a standard web browser using the three.js library [41]. The user can observe the 3D scene from an arbitrary viewpoint and select, translate, and rotate individual objects. The user sends the target poses to the robot by clicking the *Upload Scene* button. Our user study demonstrates that being able to directly manipulate objects at their target kits significantly reduces subjective workload as compared to traditional methods. Moreover, our interface does not require specialized hardware or a fast internet connection, making it accessible to common users (see video for interface demo).

C. SnapNet: 6DoF Action-Snapping Network

Specifying perfect 6DoF kitting poses is challenging. As supported by our study, allowing users to be imprecise greatly reduces their mental burden and task time as they can roughly align an object near its respective kit.

To make use of imprecise user inputs, we designed the SnapNet algorithm (Fig. 3) that refines the objects’ pose based on their 3D geometry. Concretely, the goal for SnapNet is to predict correct relative pose T_{gt} between object and kit given input volumes of object V_o , a kit $V_{k_{ws}}$, and user input $T_{user} \equiv (P_{user}, Q_{user}) \in SE(3)$. Here, we assume user input is within range: $\max_{i \in \{x, y, z\}} |P_{i, user} - P_{i, gt}| < \delta_{position}$ and $Q_{user} \cdot Q_{gt}^{-1} < \delta_{orientation}$ where $T_{gt} \equiv (P_{gt}, Q_{gt})$ is the ground-truth kitting pose. We train our system to handle poses up to $\delta_{position} = 2.8$ cm error along each translational axis and $\delta_{orientation} = 27.5^\circ$ quaternion difference.

To reduce the combinatorial search space, SnapNet predicts translation and rotation sequentially, which reduces the search space from $O(\theta_{xyz} \times \theta_{rpy})$ to $O(\theta_{xyz} + \theta_{rpy})$ where θ_{xyz} , θ_{rpy} represents discretization of translational and rotational search space.

Position prediction: Given V_o , $V_{k_{ws}}$ and P_{user} , the goal of position prediction is to infer P_{snap} . We first crop kit workspace volume $V_{k_{ws}}$ centered around P_{user} and of size $(2\delta_{position})^3$ to receive V_k . We then encode V_o and V_k via object and kit encoders (fully convolutional neural networks) to obtain deep feature embeddings $\phi(V_o)$ and $\psi(V_k)$ respectively. The algorithm then computes cross-convolution between $\phi(V_o)$ and $\psi(V_k)$ by treating $\phi(V_o)$ as convolution kernel. The output shares the same size as kit features $\psi(V_k)$. P_{snap} is chosen as position that corresponds to maximum feature correlation, i.e., *argmax* of cross convolution output. Both encoders are trained jointly to minimize voxel-wise BinaryCrossEntropy loss with label 1 at P_{gt} and 0 elsewhere.

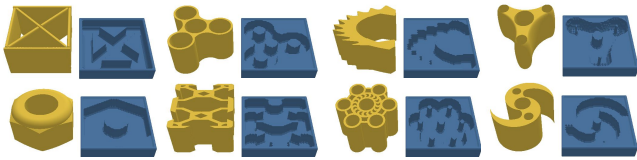


Fig. 4. **KIT1000 Dataset.** Examples of objects and generated kits.

Rotation prediction: Given V_k , V_o , user orientation Q_{user} , and position prediction P_{snap} , the goal of the Rotation module is to predict Q_{snap} . Directly regressing quaternions [2] fails to generalize (see Tab. I) and volume-based representations are susceptible to information loss under rotations. To address these issues, we use a **point-cloud-based** representation for rotation refinement. Using the refined position P_{snap} , V_k is further cropped down at center with size $(128)^3$. Both V_o and V_k volumes are converted to point-cloud representation ($N_o = 2048$ and $N_k = 4096$ points $\in \mathbb{R}^3$ respectively) to support rotation operations. We uniformly sample $N - 1$ rotations within $\delta_{orientation}$ from the user input Q_{user} . Q_{gt} is added to the set of rotations ($N = 391$) during training. For each rotation r in the set, we rotate the object point-cloud by r and concatenate it with the kit point-cloud. An additional fourth dimension is utilized to distinguish between object (1) and kit (-1) points. A PointNet++ based encoder [42] followed by fully connected layers is used to get binary classification score. We train the network using cross-entropy loss with 1 for Q_{gt} rotation and 0 otherwise.

All the modules are trained on the simulation data generated from objects and kits from our dataset (see §III-E) and then directly tested on unseen real world data.

D. Robot Planning and Execution

Picking and placing an object at specific goal pose is a challenging problem as the object may not initially be oriented such that the robot can grasp and then immediately place them in specific goal pose. Such manipulation systems are still an active research area [43], [44] and not the focus of this work. To test our system in real-world, we make a simplifying assumption that the object is top-down graspable, and the grasping surface is opposite to the kit insertion direction. No such assumptions are made for training and evaluation of scene completion and 6DoF pose prediction algorithms (Tab. I). To move the object from its current location to kitting location ${}^{robot}T_{snap}$, we pick the object via a suction-gripper-based top-down immobilizing grasp. The object is inserted into kit following a two-step primitive: (a) The robot first “hovers” at some fixed height and final orientation above the kitting location defined as ${}^{robot}T_{hover} = {}^{robot}T_{snap} \times {}^{snap}T_{hover}$, where ${}^{snap}T_{hover} \equiv ({}^{snap}P_{hover} = [0, 0, 0.1]$ m, ${}^{snap}Q_{hover} = [0, 0, 0, 1])$. (b) The robot follows a straight-line path from ${}^{robot}T_{hover}$ to final pose ${}^{robot}T_{snap}$ before releasing the suction. More details on the grasp pose estimation and trajectory computation can be found on the [webpage](#).

E. Dataset and Automatic Kit-Generation Procedure

Despite the ubiquity of kits in the packaging and transport industry, most kits are manually designed and no large-scale object-kit dataset exists. Given a 3D object geometry,

a typical kit (a) maximally confirms the object geometry and (b) allows the object to be inserted following a straight-line path at least along one direction. Our method neatly accounts for both of these: we capture an orthographic depth image of the object, which removes any artifacts that are not parallel to the insertion direction. The orthographic depth image is then converted to an occupancy grid. To allow some margin between kit and object geometry, the object 3D volume is then merged with replicas of itself after translating by margin distance along the horizontal direction. This creates a scaled version of the object geometry while preserving the centers of offset cavities. This scaled object geometry is then subtracted from the kit block to produce kit geometry.

We use objects from ABC Dataset [45], a large-scale CAD model dataset that contains a diverse set of mechanical parts. Each object is scaled to fit a $(5cm)^3$ box and a corresponding kit is generated as described above (see Fig. 4). To create 6DoF kits, we arbitrarily link 2–5 kits together using angle brackets with angles $\in [10^\circ, 45^\circ]$. We call this KIT1000 dataset and it will be made available.

IV. EXPERIMENTS

We first evaluate the action-snapping module (§IV-A) followed by a full system evaluation on a real-world platform (§IV-B) and a real-world user study (§IV-C).

A. Action-Snapping Evaluation

Metrics: We evaluate 6DoF pose prediction $T_{snap} \equiv (P_{snap}, Q_{snap})$ using two metrics: positional error $\delta_{pos} = \|P_{snap} - P_{gt}\|_2$. Rotational error δ_{rot} is computed as the geodesic distance $\arccos(2(Q_{snap} \cdot Q_{gt})^2 - 1)$.

Comparison with alternative approaches: We compare our algorithm with TransporterNet [22] and KitNet [2]. Since both algorithms are trained without user input, we modify our algorithm to also work without user input: For position prediction, instead of cropping V_{kws} around user input P_{user} , we directly use V_{kws} as V_k . For rotation prediction, we uniformly sample $roll, pitch \in [-15^\circ, 15^\circ]$, and $yaw \in [-180^\circ, 180^\circ]$. TransporterNet [22] consists of a pick and a place module. In our evaluation, we use the groundtruth pick position and retrain its place module with extensions to 6DoF actions. When user input is available, we filter out predictions that is far from provided pose, i.e., $T_{user} \pm (\delta_{position}, \delta_{orientation})$. KitNet [2] predicts only the rotation of the object via regression, so there is no straightforward way to incorporate user inputs. Thus, we only evaluate the rotation predictions of KitNet without user input.

Tab. I shows that both baselines fail to give accurate predictions. We hypothesize that without full geometry estimation, they do not have enough information to infer a 3D pose. By leveraging full 3D geometry and efficiently searching the SE(3) space, our model outperforms the baselines both with and without user input.

Effects of shape completion: To study the effect of shape completion on action snapping, we compare our approach without this step. *SnapNet-PartialVol* uses partial volume $V_{partial}$ to perform shape matching. Tab. I shows that our

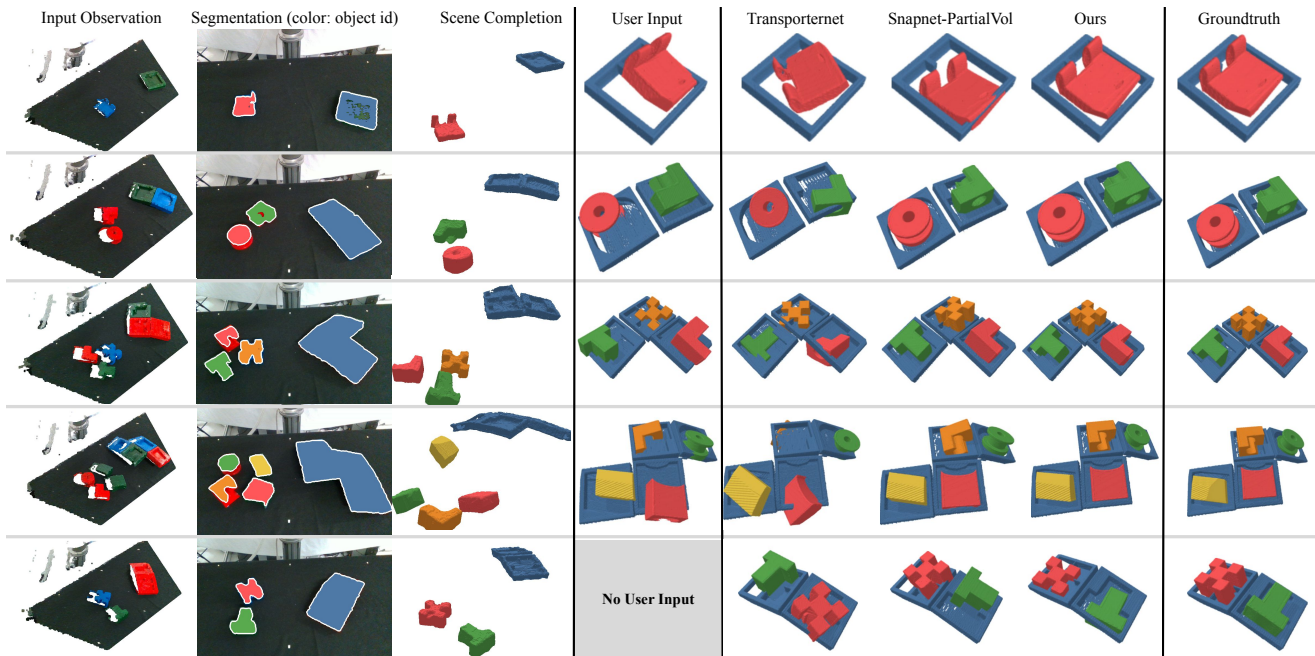


Fig. 5. **Comparisons to Alternative Approaches** We compare SEaT with 6DoF kitting baselines on novel object and kit geometries. TransporterNet fails to generalize to unseen object and kit geometries. SnapNet-PartialVol works for simple objects (row 2) but fails for objects with complex geometries (rows 3–4). When given no user input, both baselines frequently place objects at the wrong kits (row 5). In the last five columns, we use ground truth meshes to visualize poses. For more results, see the [project webpage](#).

TABLE I
ACTION-SNAPPING RESULTS AND COMPARISON

	With user input		Without user input	
	δ_{pos} (mm)	δ_{rot} (deg)	δ_{pos} (mm)	δ_{rot} (deg)
KitNet [2]	-	-	-	49.2
TransporterNet [22]	15.3	18.3	41.5	45.1
SnapNet-PartialVol	5.1	5.7	49.4	53.2
SnapNet (Ours)	3.9	4.9	10.8	29.6
SnapNet-GTVol	3.7	4.61	8.1	28.9

TABLE II
SYSTEM EVALUATION ON THE REAL-WORLD DATASET

Segmentation mIoU	Obj. Completion mIoU	Kit Completion Chamfer	Action Snapping pos	Action Snapping rot
69.1%	92.4%	6.3 mm	99.1 %	8.0 mm 7.2 mm 6.0°

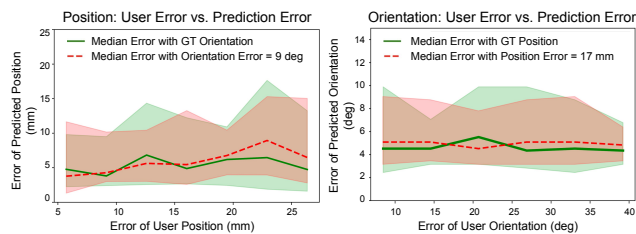


Fig. 6. **Robustness to User Input** with [20, 80] percentile region shaded. The left graph shows an analysis of error in position prediction, keeping error in user orientation fixed. As user position error increases, SnapNet maintains its low prediction error. Moreover, even with a large error in user orientation (dotted-red), SnapNet can predict position with low error. Similar results for predicted orientation, keeping the error in user position fixed, are shown on the right.

model *SnapNet* achieves better performance than *SnapNet-PartialVol*. We believe that this is because partial volumes lack of clear and precise object boundaries that shape matching crucially depends on. With ground-truth shape, *SnapNet-GTVol* can further improve action-snapping performance. This result indicates that the scene-completion module is essential for achieving accurate action snapping.

Robustness against user errors: We also test the algorithm’s robustness to different levels of user-input error. For a controlled experiment, we analyze error in position and

rotation prediction one-by-one by keeping the error in user orientation and user position fixed respectively. Fig. 6 breaks down the performance of our model by plotting prediction errors δ_{pos} , δ_{rot} against user-input errors. The plot shows that as user error increases, the model error remains roughly constant, demonstrating the robustness of the algorithm.

B. System Evaluation on Real-World Platform

Finally, we evaluate our algorithm on a real-world platform using a UR5 robot, an XYZ Robotics suction gripper [46], and a calibrated Intel RealSense D415 RGB-D camera. To account for RealSense camera precision (5 mm depth error [47], for pick-place task, the error would be 10 mm), we 3D-printed the kits from our test set with a larger object-kit margin of 1 cm as compared to 2.5 mm margin in simulation.

For systematic evaluation, we collect and label 23 scenes (7 of 1-kit, 7 of 2-kit, 4 of 3-kit, and 5 of 4-kit tasks), with ground-truth object target poses. We directly tested all our models (trained on simulation) with this real-world benchmark. To eliminate small holes in shape completed object volumes $V_{\text{completed}}$ due to sensor noise in input V_{partial} , we extend all the object voxels till the ground plane. To mimic user input, we randomly sample position and orientation in the vicinity (δ_{position} , $\delta_{\text{orientation}}$) of the ground-truth pose. Fig. 5 shows qualitative results on this real-world benchmark. Tab. II shows quantitative results for each individual component. The resulting average position and ro-

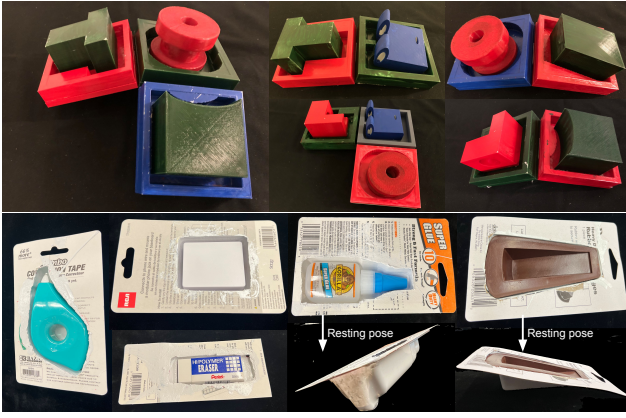


Fig. 7. **Kits for real-world experiments.** Top: 3D-printed kits from test dataset are connected at arbitrary angles to create 6DoF kits. Bottom: real-world kits. Arrows show the resting pose for a few kits which require non top-down object insertion.

tation error are comparable with the algorithm’s performance in simulation (Tab. I). Moreover, our model has similar level performance on training and test dataset with unseen shapes, which shows that our model is generalizable by leveraging a large simulated dataset.

In addition to 3D printed objects, we also evaluate the system on real-world object-kits (Fig. 7-bottom). Since these kits have a tighter object-kit margin, we use Photoneo Scanner with higher depth precision of 0.5 mm [48]. Fig. 8 shows the qualitative evaluation. We refer readers to supplementary video for real-world demonstration of our system.

C. User Study on Real-World Platform

Our user study aims to test the hypothesis that the SEaT interface would be easier to use than traditional teleoperation interfaces. We conducted a user study, approved by our institution’s IRB with 10 non-expert users.

Task and Procedure: Participants completed four kit-assembly tasks per interface (two 2-kit and two 3-kit tasks). For each n -kit task, we randomly attached n kits from a set of six unseen 3D-printed kits using randomly chosen angle brackets $\{10^\circ, 20^\circ, 30^\circ\}$ (see Fig 7). The study used a within-subjects design, where all participants performed both tasks using both interfaces in random order. Participants performed the 2-kit tasks first and then the 3-kit tasks for each interface.

Comparisons: We compared with EE-Control, a representative teleoperation interface where a user can specify 6DoF pick-and-place pose of the end-effector on the point-cloud representation of the scene. In the EE-Control interface, the user specifies a single pick-and-place pose followed by robot execution. Once the robot executes, the user scene is updated with the new scene and the user repeats the process. In SEaT, the user specifies the goal poses of all objects at once.

Dependent Measures: Our objective dependent measures were *a. Success rate*: the number of kits successfully assembled over the total number of kits, *b. specification time*: the time the user spent interacting with the interface for specifying goals, and *c. execution time*: the total system time minus the specification time. We also had a subjective dependent measure *d. unweighted NASA Task Load Index*

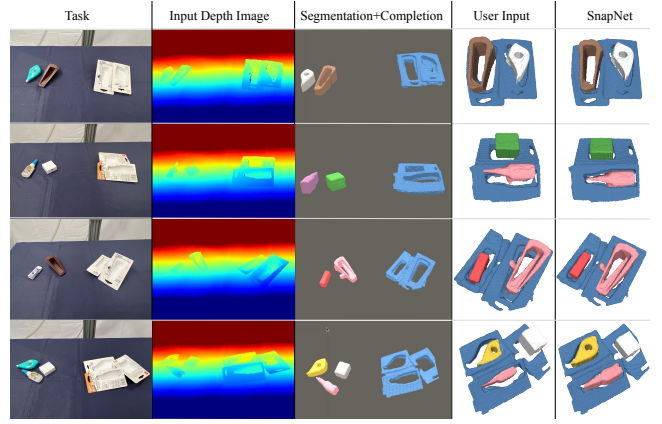


Fig. 8. **Qualitative results on real-world kits.** See video for more results.

(NASA-TLX) [49], which includes values for MentalDemand, PhysicalDemand, TemporalDemand, Performance, Effort, and Frustration. Since a user is allowed to operate on the updated scene in the EE-Control interface, in theory they can always assemble all the objects given infinite time. Therefore, for both interfaces, a user can only start an update for an n -kit task if the time already spent is less than n minutes. Users are informed about this time limit beforehand. We hypothesized that each of these dependent measures would differ between the SEaT and EE-Control interfaces.

Results: We evaluated the hypotheses for significance with $\alpha = .05$. While the *execution time* of SEaT is longer (+12s) due to model inference, the users spent significantly shorter *specification time* ($-27s, p < .001$), and achieved significantly higher *task success rate* (+33.1%, $p < .001$). For subjective measures (NASA-TLX), the participants reported significantly lower MentalDemand ($-39.2\%, p = .003$), lower TemporalDemand ($-43.1\%, p < .001$), lower Effort ($-32.0\%, p = .002$), and lower Frustration ($-40.7\%, p = .024$). The reported differences in PhysicalDemand and Performance are not significant between these two methods.

The shorter specification time and lower mental load of SEaT indicates a potential possibility of using SEaT to simultaneously operate multiple robots. In this case, a user can continue specify tasks (for another robot) during model inference and robot execution time, which will further improve the system’s overall efficiency.

V. CONCLUSION

We introduced “Scene Editing as Teleoperation”, which allows non-expert end users to perform precise multi-unknown-object 6DoF kitting tasks. Experiments demonstrated that SEaT improves efficiency, success rate, and subjective workload for 6DoF kit-assembly tasks.

Since our teleoperation interface assumes rigid objects, it cannot be directly applied to tasks involving articulated objects (e.g., opening a drawer). It would be interesting to discover articulation via RGB-D images [50], [51] and integrate it with our system. Planning the grasp and a set of sequential 6DoF robot actions for general 6DoF kitting tasks would also be an interesting future direction, where the robot might need to plan a place-driven grasp [52] or reorient the object before kitting [53].

REFERENCES

- [1] D. Kent, C. Saldanha, and S. Chernova, "Leveraging depth data in remote robot teleoperation interfaces for general object manipulation," *The International Journal of Robotics Research*, vol. 39, no. 1, pp. 39–53, 2020.
- [2] S. Devgon, J. Ichnowski, M. Danielczuk, D. S. Brown, A. Balakrishna, S. Joshi, E. Rocha, E. Solowjow, and K. Goldberg, "Kit-Net: Self-supervised learning to kit novel 3D objects into novel 3D cavities," *arXiv preprint arXiv:2107.05789*, 2021.
- [3] C. Reinhard and P.-F. Breton, "Experimental validation of 3ds Max Design 2009 and Daysim 3.0," in *Proceedings 11th International IBPSA Conference*, 2009, pp. 1514–1521.
- [4] "Unity character animation," <http://video.unity3d.com/video/4655480/unity-character-animation-gdc>.
- [5] "<https://www.solidworks.com/>," 2021.
- [6] e. a. Billard, Aude, "Robot programming by demonstration." 2008.
- [7] T. Ren, Y. Dong, and K. C. Dan Wu, "Design of direct teaching behavior of collaborative robot based on force interaction."
- [8] S. Hayati and S. Venkataraman, "Design and implementation of a robot control system with traded and shared control capability," in *ICRA*, 1989. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=100161>
- [9] M. Oda, N. Inaba, Y. Takano, S. Nishida, M. Kayashi, and Y. Sugano, "Onboard local compensation on ETS-W space robot teleoperation," in *IEEE/ASME Intl. Conf. on Advanced Intelligent Mechatronics*, 1999. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=803253>
- [10] P. Michelman and P. Allen, "Shared autonomy in a robot hand teleoperation system," in *IROS*, 1994. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=407383>
- [11] A. E. Leeper, K. Hsiao, M. Ciocarlie, L. Takayama, and D. Gossow, "Strategies for human-in-the-loop robotic grasping," in *Proceedings of the seventh annual ACM/IEEE international conference on Human-Robot Interaction*, 2012, pp. 1–8.
- [12] D. Kent, C. Saldanha, and S. Chernova, "Leveraging depth data in remote robot teleoperation interfaces for general object manipulation," *The International Journal of Robotics Research*, vol. 39, no. 1, pp. 39–53, 2020.
- [13] M. Ciocarlie, K. Hsiao, A. Leeper, and D. Gossow, "Mobile manipulation through an assistive home robot," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 5313–5320.
- [14] M. Gualtieri, A. Ten Pas, K. Saenko, and R. Platt, "High precision grasp pose detection in dense clutter," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2016, pp. 598–605.
- [15] S. Y. Gadre, E. Rosen, G. Chien, E. Phillips, S. Tellex, and G. Konidaris, "End-user robot programming using mixed reality," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 2707–2713.
- [16] D. Gossow, A. Leeper, D. Hershberger, and M. Ciocarlie, "Interactive markers: 3-D user interfaces for ROS applications [ROS topics]," *IEEE Robotics & Automation Magazine*, vol. 18, no. 4, pp. 14–15, 2011.
- [17] H.-I. Lin and Y.-H. Lin, "A novel teaching system for industrial robots," in *Sensors*, 2014.
- [18] J. K. Lee and O. C. Jenkins, "Goal-based teleoperation for robot manipulation," in *Artificial Intelligence for Human-Robot Interaction: Papers from the 2014 AAAI Fall Symposium.*, 2014.
- [19] M. Ciocarlie, K. Hsiao, A. Leeper, and D. Gossow, "Mobile manipulation through an assistive home robot," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 5313–5320.
- [20] A. B. Y. Litvak and A. Bar-Hillel, "Learning pose estimation for high-precision robotic assembly using simulated depth images." 2018.
- [21] K. Zakka, A. Zeng, J. Lee, and S. Song, "Form2Fit: Learning shape priors for generalizable assembly from disassembly," in *International Conference on Robotics and Automation (ICRA)*, 2020.
- [22] A. Zeng, P. Florence, J. Tompson, S. Welker, J. Chien, M. Attarian, T. Armstrong, I. Krasin, D. Duong, V. Sindhwani, and J. Lee, "Transporter networks: Rearranging the visual world for robotic manipulation," *Conference on Robot Learning (CoRL)*, 2020.
- [23] S. Song and J. Xiao, "Deep sliding shapes for amodal 3D object detection in RGB-D images," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 808–816.
- [24] M. Schwarz, A. Milan, A. S. Periyasamy, and S. Behnke, "RGB-D object detection and semantic segmentation for autonomous manipulation in clutter," *The International Journal of Robotics Research*, vol. 37, no. 4-5, pp. 437–451, 2018.
- [25] J. Taylor, J. Shotton, T. Sharp, and A. Fitzgibbon, "The vitruvian manifold: Inferring dense correspondences for one-shot human pose estimation," in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE, 2012, pp. 103–110.
- [26] R. A. Güler, N. Neverova, and I. Kokkinos, "Densepose: Dense human pose estimation in the wild," *arXiv preprint arXiv:1802.00434*, 2018.
- [27] E. Brachmann, A. Krull, F. Michel, S. Gumhold, J. Shotton, and C. Rother, "Learning 6D object pose estimation using 3D object coordinates," in *European conference on computer vision*. Springer, 2014, pp. 536–551.
- [28] S. Gupta, P. Arbeláez, R. Girshick, and J. Malik, "Aligning 3D models to RGB-D images of cluttered scenes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 4731–4740.
- [29] R. Guo, *Scene understanding with complete scenes and structured representations*. University of Illinois at Urbana-Champaign, 2014.
- [30] J. Papon and M. Schoeler, "Semantic pose using deep networks trained on synthetic RGB-D," in *Computer Vision (ICCV), 2015 IEEE International Conference on*. IEEE, 2015, pp. 774–782.
- [31] M. Braun, Q. Rao, Y. Wang, and F. Flohr, "Pose-RCNN: Joint object detection and pose estimation using 3D object proposals," in *Intelligent Transportation Systems (ITSC), 2016 IEEE 19th International Conference on*. IEEE, 2016, pp. 1546–1551.
- [32] S. Song, F. Yu, A. Zeng, A. X. Chang, M. Savva, and T. Funkhouser, "Semantic scene completion from a single depth image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1746–1754.
- [33] A. Dai, D. Ritchie, M. Bokeloh, S. Reed, J. Sturm, and M. Nießner, "ScanComplete: Large-scale scene completion and semantic segmentation for 3D scans," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4578–4587.
- [34] Z. Xu, Z. He, J. Wu, and S. Song, "Learning 3d dynamic scene representations for robot manipulation," in *Conference on Robotic Learning (CoRL)*, 2020.
- [35] Y. Liang, B. Chen, and S. Song, "Sscnav: Confidence-aware semantic scene completion for visual semantic navigation," in *Proc. of The International Conference in Robotics and Automation (ICRA)*, 2021.
- [36] F. Liu, C. Shen, and G. Lin, "Deep convolutional neural fields for depth estimation from a single image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5162–5170.
- [37] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab, "Deeper depth prediction with fully convolutional residual networks," in *3D Vision (3DV), 2016 Fourth International Conference on*. IEEE, 2016, pp. 239–248.
- [38] M. Danielczuk, M. Matl, S. Gupta, A. Li, A. Lee, J. Mahler, and K. Goldberg, "Segmenting unknown 3D objects from real depth images using Mask R-CNN trained on synthetic data," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 7283–7290.
- [39] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
- [40] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon, "KinectFusion: Real-time dense surface mapping and tracking," in *2011 10th IEEE international symposium on mixed and augmented reality*. IEEE, 2011, pp. 127–136.
- [41] B. Danchilla, "Three.js framework," in *Beginning WebGL for HTML5*. Springer, 2012, pp. 173–203.
- [42] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "PointNet++: Deep hierarchical feature learning on point sets in a metric space," *arXiv preprint arXiv:1706.02413*, 2017.
- [43] K. Wada, S. James, and A. J. Davison, "Reorientbot: Learning object reorientation for specific-posed placement," *arXiv preprint arXiv:2202.11092*, 2022.
- [44] N. Chavan-Dafle, R. Holladay, and A. Rodriguez, "In-hand manipulation via motion cones," *arXiv preprint arXiv:1810.00219*, 2018.
- [45] S. Koch, A. Matveev, Z. Jiang, F. Williams, A. Artemov, E. Burnaev, M. Alexa, D. Zorin, and D. Panozzo, "ABC: A big CAD model dataset for geometric deep learning," in *Proceedings of the IEEE/CVF*

Conference on Computer Vision and Pattern Recognition, 2019, pp. 9601–9611.

- [46] “XYZ Robotics,” <https://en.xyzrobotics.ai/>, accessed: 2021-09-13.
- [47] “Depth Resolution of Intel® RealSense™ Depth Camera D435 and Intel® RealSense™ Camera SR300,” <https://www.intel.com/content/www/us/en/support/articles/000026260/emerging-technologies/intel-realsense-technology.html>, accessed: 2021-09-13.
- [48] “Datasheet for Photoneo PhoXi 3D Scanner M,” <https://www.photoneo.com/products/phoxi-scan-m/>, accessed: 2022-02-28.
- [49] S. G. Hart, “NASA-task load index (NASA-TLX); 20 years later,” in *Proceedings of the human factors and ergonomics society annual meeting*, vol. 50, no. 9. Sage publications Sage CA: Los Angeles, CA, 2006, pp. 904–908.
- [50] S. Y. Gadre, K. Ehsani, and S. Song, “Act the part: Learning interaction strategies for articulated object part discovery,” *ICCV*, 2021.
- [51] Z. Xu, Z. He, and S. Song, “UMPNet: Universal manipulation policy network for articulated objects,” 2021.
- [52] K. Fang, Y. Zhu, A. Garg, A. Kurenkov, V. Mehta, F-F. Li, and S. Savarese, “Learning task-oriented grasping for tool manipulation from simulated self-supervision,” *The International Journal of Robotics Research*, vol. 39, no. 2-3, pp. 202–216, 2020.
- [53] N. Chavan-Daffe, M. T. Mason, H. Staab, G. Rossano, and A. Rodriguez, “A two-phase gripper to reorient and grasp,” in *2015 IEEE International Conference on Automation Science and Engineering (CASE)*. IEEE, 2015, pp. 1249–1255.