Slot-Level Robotic Placement via Visual Imitation from Single Human Video

Anonymous submission

Paper ID

Abstract

001 The majority of modern robot learning methods focus on 002 learning a set of pre-defined tasks with limited or no gen-003 eralization to new tasks. Extending the robot skillset to novel tasks involves gathering an extensive amount of train-004 ing data for additional tasks. In this paper, we address 005 the problem of teaching new tasks to robots using human 006 007 demonstration videos for repetitive tasks (e.g., packing). This task requires understanding the human video to iden-008 tify which object is being manipulated (the pick object) and 009 where it is being placed (the placement slot). In addition, 010 it needs to re-identify the pick object and the placement 011 012 slots during inference along with the relative poses to en-013 able robot execution of the task. To tackle this, we propose SLeRP, a modular system that leverages several advanced 014 visual foundation models and a novel slot-level placement 015 detector Slot-Net, eliminating the need for expensive video 016 demonstrations for training. We evaluate our system using 017 018 a new benchmark of real-world videos. The evaluation results show that SLeRP outperforms several baselines and 019 020 can be deployed on a real robot.

1. Introduction

Humans demonstrate exceptional skill in performing fine-022 grained manipulation tasks with high precision in their daily 023 024 lives. From arranging eggs in an egg carton to sorting utensils in an organizer, we excel at tasks that require identify-025 ing and reasoning about which objects to pick up and how 026 to place them into confined slots. Cognitive and motor de-027 velopment theories suggest that we develop such skills at 028 029 a young age, based on early experiences like playing with shape sorter toys [49]. However, current robotic and auto-030 mated systems are not yet as adept as humans at perceiving 031 and performing these fine-grained manipulation tasks. 032

Slot-level manipulation is crucial in various industrial,
 logistics, and domestic contexts. For example, in industrial
 settings, machine tending [57] requires placing components
 precisely into machine slots for assembly or processing. In



Figure 1. We introduce the novel problem of imitating slot-level robotic placement from a single human video. Given a human demonstration video showing an object being placed in a slot, and a new robot-view image captured by the robot wrist camera (may feature varied camera and object poses, changed scenes), SLeRP is able to find the corresponding object and similar slots in the robot view, and provide the 6-DoF transformation matrix for each detected slot to guide the robot in placing the object accurately.

logistics, sorting and packaging tasks, such as organizing parcels in a warehouse or placing products into shipping containers, demand efficient and precise placement to optimize space and minimize damage. In domestic environments, future home assistant robots will need to perform slot-level manipulation tasks such as organizing items in cabinets, placing dishes in a dishwasher, and even preparing meals by accurately arranging ingredients in a pan.

The task of programming robots to perform slot-level 045 placement remains arduous. Traditional methods [20, 67] 046 often require manual programming with domain expertise 047 and assume that the object models and slot locations are 048 known beforehand. Learning-based approaches [9, 30] 049 show promise in alleviating the burden of programming; 050 however, collecting robot data through tele-operation re-051 mains tedious and inefficient and can be particularly brittle 052 for high-precision tasks due to embodiment gaps. Learning 053 from human demonstration videos has recently emerged as 054 a promising approach due to its ease, speed of collection, 055 and potential to capture slot-level details. However, previ-056 ous research [3, 6, 26, 76, 86] has generally been limited to 057 coarser object-level tasks and often requires large amounts 058 of training data to learn how to parse human demonstrations 059 and translate them into robot policies. 060

037

038

039

040

041

042

043

^{*} Work done during internship at Nvidia. [†] Primary mentor.

123

124

125

126

127

128



Figure 2. **Method Overview.** The system begins by analyzing the input human video, tracking the object (highlighted in yellow) throughout the sequence and identifying the placement slot (highlighted in red). Next, we re-identify the object and the slot in the robot's view by correlating the human-view and robot-view images. Using depth images, we reconstruct the observations in 3D and compute a single 6-DoF object transformation T in the robot's view, enabling the robot to transfer the object into the slot. If more than one slot is present, we detect all applicable slots and compute one 6-DoF object transformation for each slot. Finally, such 6-DoF object transformations are sent to the downstream robot planning and control pipeline for real robot pick-and-place execution.

061 In this paper, we study the novel problem of recognizing slot-level object placement from a single human video, and 062 063 estimating 6-DoF transformations for robot imitation. As 064 shown in Fig. 1, the task takes two visual inputs: (1) a sin-065 gle human RGB-D video in which a person demonstrates picking up an object (e.g., a muffin) and precisely placing it 066 into a slot within a placement object (e.g., a tray), and (2) a 067 single RGB-D image captured from the robot's wrist cam-068 era, representing the new setup for the robot to operate in 069 with possible varying camera and object poses compared to 070 071 the human video. The outputs aim to detect the object and all empty slots in the robot's view similar to the placement 072 slot in the human video, as well as compute the 6-DoF ob-073 ject transformations necessary for the robot to transfer the 074 object from its initial position to each of the slots. 075

076 We propose a novel modular approach called SLeRP 077 (*i.e.*, Slot-Level Robotic Placement), to tackle the problem. As shown in Fig. 2, SLeRP starts by analyzing the input hu-078 man demonstration video, tracking the manipulated object 079 080 across the video frames and identifying the placement slot. Next, within the robot's view, SLeRP re-identifies both the 081 object and the slot by correlating the human-view images 082 with the robot-view images. By lifting the observations in 083 3D using the depth sensing and camera parameters, SLeRP 084 calculates a 6-DoF transformation matrix for the robot to 085 transfer the object from its initial location to the desired slot 086 087 in the robot's view. If multiple slots are present, SLeRP de-088 tects all slots that are similar to the one in the human video 089 and computes the object transformations for all of them. Finally, the computed 6-DoF object transformations are sent 090 to the downstream robot planning and control pipeline for 091 092 robot pick-and-place execution.

A key component of SLeRP is the detection of placement
 slots. Currently, no existing method is specifically designed
 for this task, and simple image differencing or change de-

tection [18] does not effectively solve the problem. There-096 fore, we propose a new slot-level placement detector, Slot-097 Net, that takes two image frames from a human demonstra-098 tion video-one before and one after placement-and out-099 puts a 2D mask outlining the placement slot on the images. 100 Unlike common vision tasks, collecting a sizable training 101 dataset for slot-level placement detection is challenging. To 102 address this challenge, we introduce a generative AI-based 103 data creation pipeline that expands the training set by boot-104 strapping from a small set of images. 105

For evaluation, we introduce a new dataset compris-106 ing 288 real-world videos targeted at studying this novel 107 problem. We compare our method against several base-108 line approaches, including ORION [86], a state-of-the-109 art method for object-level pick-and-place from a sin-110 gle human video; CLIPort [61], an end-to-end imitation-111 learning-based language-conditioned policy for tabletop 112 tasks; adapted versions of both for the novel slot-level 113 task; and a custom baseline leveraging cutting-edge vision-114 language models like GPT-40 [25]. Our results demon-115 strate that SLeRP outperforms baselines in accurately pre-116 dicting placement slots and computing 6-DoF transforma-117 tions across diverse real-world tasks. Our ablation studies 118 further validate several key components and design choices 119 in our system. Finally, we conduct real-robot experiments 120 that successfully apply the system in real world scenarios. 121

In summary, the core contributions of this paper are:

- Studying the novel task of slot-level object placement by learning from a single human demonstration video;
- Designing the modular approach SLeRP and the slot-level placement detector Slot-Net to tackle this problem;
- Introducing a new benchmark and several baseline methods to systematically evaluate system performance;
- Demonstrating that SLeRP achieves strong performance 129 in real-world and real-robot evaluations. 130

131 2. Related Work

Object Placement in Robotics. Identifying where and how 132 133 to place an object after picking it up is a crucial step in robotic pick-and-place tasks [37]. Early works [5, 22, 23] 134 analytically search for flat features on the object and the 135 placement surface. Modern learning-based methods es-136 timate placement locations and poses using learned fea-137 138 tures, focusing mainly on flat surfaces [46], such as table-139 top [13, 43, 82] and furniture shelves [44]. Researchers have also explored tabletop object placement under spatial 140 and semantic constraints given other objects [35, 42, 52]. 141

In the more challenging case of placing an object on an-142 other non-flat object, prior work [16, 28, 50, 59, 62, 63, 66, 143 144 73, 86] has explored tasks like putting one spoon in a cup or hanging a mug on a rack. Our work extends these studies 145 by focusing on placing objects into all empty, fine-grained, 146 tight-fitting slots (e.g., all egg slots in a carton), a task that 147 demands greater precision in recognition and prediction, as 148 well as handling multiple placement locations. Addition-149 ally, unlike previous work [24, 61, 75, 83, 84] that addresses 150 a 2D planar setting and requires task-specific training from 151 a few robot demonstrations, our approach tackles this prob-152 153 lem in 3D, learning from a single human demonstration to enable one-shot generalization to novel tasks. 154

Imitation Learning from Human Videos. Human videos 155 serve as a natural, information-rich, and easily accessible 156 source of data for learning robotic manipulation. Previ-157 158 ous work has explored diverse methods to extract, represent, and apply knowledge from human videos to support 159 robot manipulation learning. These approaches include pre-160 training latent visual representations [15, 39, 45, 58], infer-161 ring action trajectories or plans [3, 32, 47, 71, 79], learning 162 163 value or reward functions [10, 38], reconstructing human 164 hand or hand-object interaction [41, 51, 54, 60, 64], parsing interaction goals and affordance [4, 29, 36, 78], learning 165 point tracks for human-to-robot transfer [7, 74, 76, 77, 81], 166 etc. While the primary goals of these works are typi-167 168 cally learning robot trajectories or manipulation policies, 169 our work explores a novel perspective by recognizing finegrained placement slots as visual imitation targets. 170

Additionally, we tackle robot imitation learning from a 171 172 single human video. Previous work has investigated one-173 shot [14, 26, 27, 40, 80] and even zero-shot learning from human videos [6]; however, these approaches often re-174 quire extensive human video datasets, sometimes paired 175 with robot videos, to span multiple tasks during training. 176 In contrast, our approach leverages existing visual founda-177 178 tion models, eliminating the need for large-scale training videos. A notably similar work ORION [86] relies on text 179 to recognize task-relevant objects and primarily focuses on 180 object-level pick-and-place. In contrast, our method exclu-181 sively extracts information from a single human video to 182 perform more fine-grained slot-level placement tasks. 183

3. Problem Formulation

We formulate the novel problem of recognizing slot-level185object placement from a single human video, and estimating1866-DoF transformations for downstream robot imitation.187

Inputs. The task takes the following inputs:

- a single RGB-D human demonstration video with n frames, denoted as $\mathcal{H} = {\mathbf{H}_1, \mathbf{H}_2, \cdots, \mathbf{H}_n}$, recording a person picking up an object O from the scene and placing it in a slot S within a placement object;
- a single RGB-D robot-view image **R** that captures the robot's observation, often taken from the robot wrist camera and possibly with different camera and object poses, or scene layouts.

Outputs. The task outputs, in the robot's view, are:

- an object mask $\mathbf{M}_{\mathbf{R}}^{O}$ over the robot image \mathbf{R} that segments the object O to pick;
- a list of slot masks $\{\mathbf{M}_{R}^{S_{0}}, \mathbf{M}_{R}^{S_{1}}, \cdots, \mathbf{M}_{R}^{S_{k}}\}$ over the robot image **R** that marks all empty slots on the placement object similar to the demonstrated placement slot in the human video \mathcal{H} ;
- a list of 3D 6-Degree-of-Freedom (DoF) object transformation matrices {T₀, T₁, ..., T_k | T_i ∈ SE(3)} in the robot's coordinate frame, for the robot to transfer the object O from its initial position to all the detected slots.

Passing the detected object and slot masks, as well as the calculated 6-DoF object transformatrion matrices, down-stream robot pick-and-place pipeline is able to execute slot-level object placement as shown in Fig. 1.

4. Method

In this section, we present the technical designs of SLeRP. We begin with an overview (Sec. 4.1) and then dive into more details in parsing the input human video (Sec. 4.2) and correlating to the robot's view image (Sec. 4.3).

4.1. System Overview

Taking as inputs a human demonstration video \mathcal{H} and a 218 robot-view image R, our method SLeRP (Fig. 2) starts with 219 parsing the input human video (Sec. 4.2) by tracking the 220 object O throughout the video frames and identifying the 221 placement slot S. After this process, we obtain an object 222 mask $\mathbf{M}_{\mathbf{H}_1}^O$ and a slot mask $\mathbf{M}_{\mathbf{H}_1}^S$ over the first frame of 223 the human video H_1 . Next, by leveraging this informa-224 tion, SLeRP correlates the human-view and robot-view im-225 ages (Sec. 4.3), and re-identify the object mask $\mathbf{M}_{\mathbf{R}}^{O}$ and the 226 slot mask $\mathbf{M}^{S_0}_{\mathbf{R}}$ in the robot-view image \mathbf{R} , as observed in 227 the first human frame. If multiple similar slots are present, the system detects other empty slots $\{\mathbf{M}_{R}^{S_{1}}, \cdots, \mathbf{M}_{R}^{S_{k}}\}$ as 228 229 well. Then, the system lifts the human and robot observa-230 tions in 3D using the depth sensing and camera intrinsics, 231 and computes a single 6-DoF object transformation matrix 232 $\mathbf{T}_i \in SE(3)$ for each detected slot $\mathbf{M}_R^{S_i}$. 233

184 185

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

213 214 215

212



Figure 3. **Parse Human Video.** Given the input human video (bottom), we run state-of-the-art hand-object detector (yellow) and tracker (blue) to obtain the pick object mask (yellow) and train a novel network Slot-Net (red) to identify the slot mask (red).

4.2. Parsing the Human Demonstration Video

The input human video precisely demonstrates what is the pick object O and where is the placement slot S. As shown in Fig. 3, our method utilizes powerful hand-object detection and tracking systems to identify the object mask $\mathbf{M}_{\mathbf{H}_1}^O$ and proposes a novel network Slot-Net for estimating the slot mask $\mathbf{M}_{\mathbf{H}_1}^S$ over the first frame of the human video.

Object detection and tracking. We use a hand-object de-241 242 tector [11] to detect frame-wise hands and in-contact ob-243 jects, enabling us to locate the pick object O in the human video. As the detector operates on a per-frame ba-244 245 sis, there may be temporally inconsistent predictions. To refine the detection results, we apply MASA's matching al-246 247 gorithm [34] to generate smooth trajectories for the hand 248 and pick object across the hand-object contact frames. We then identify a confident key frame, when the hand and ob-249 250 ject first interacts, and use SAM2 [55] to track through the video, producing per-frame object segmentation $\mathbf{M}_{\mathbf{H}_{i}}^{O}$. 251

Placement slot detection (Slot-Net). Since no prior work 252 has studied the problem of detecting the placement slot 253 254 given a human pick-and-place video, we propose our own novel network Slot-Net for this purpose. We leverage the 255 SAM architecture [31] as the backbone given its powerful 256 capability in segmentation. Slot-Net takes the starting frame 257 \mathbf{H}_1 of the pick-place video as the input, together with the 258 259 absolute image difference in gray-scale between the starting and end frame $|\mathbf{H}_1 - \mathbf{H}_n|$ as the visual prompt, and is 260 tasked to output a slot segment in the starting human frame 261 image $\mathbf{M}_{\mathbf{H}_1}^S$. We leverage SAM's large-scale pretraining 262 263 by preserving most of its designs (e.g., the image encoder, the mask decoder) and we use the same image encoder as 264 265 the prompt embedder to process image difference prompt. Since we find that the SAM pretraining does not directly 266 work on such customized new task, finetuning over such 267 slot-level placement data is necessary. 268

269 Slot-Net data generation. Training our SAM-based
270 SLeRP requires a lot of data, yet collecting fine-grained
271 slot-level placement data in the real world is expensive.

Figure 4. **Slot-Net Data Generation.** Given an object-centric image (top middle), we inpaint to remove an object and reveal its slot (top left) and manually annotate the slot mask (top right). We then outpaint these images with a scene background (bottom) to create a starting and end image pair with a ground-truth slot mask.

However, recent generative models have demonstrated great 272 capabilities in generating realistic images [53], excelling at 273 tasks such as object removal and image outpainting. We 274 therefore propose a semi-automatic synthetic data gener-275 ation pipeline (Fig. 4). Given a collected object-centric 276 image of a placement object with many slots, we utilize 277 a state-of-the-art object removal model (SDXL [70] and 278 Cleanup.pictures [12]) to remove one pick object from one 279 slot and manually annotate the slot mask for the removed 280 object using TORAS [65]. Then, we employ a powerful im-281 age outpainting generative model (Hugging Face Outpaint-282 ing script [1]) to expand the image canvas, generating 100 283 images in diverse backgrounds, prompted with Llama [69] 284 generated text prompts, for each object-centric image. 285

In this manner, we obtain a large number of annotated starting and end image pairs to train SLeRP. We crowdsourced and collected 2,138 object-centric images of items with slots, spanning 67 object categories, by capturing them in everyday environments. We applied 100 augmentations for each slot on the object-centric image, resulting in 156K images for training, with the rest left for testing and validation. See supplementary for more details.

4.3. Correlating to the Robot-view Image

After we obtain the pick object mask $\mathbf{M}_{\mathbf{H}_1}^O$ and the slot mask $\mathbf{M}_{\mathbf{H}_1}^S$ from the human video, the next step is to correlate this information to the robot's view \mathbf{R} . As shown in Fig. 5, SLeRP first re-identifies the object mask $\mathbf{M}_{\mathbf{R}}^O$ and a list of empty slot masks $\{\mathbf{M}_{\mathbf{R}}^{S_0}, \mathbf{M}_{\mathbf{R}}^{S_1}, \cdots, \mathbf{M}_{\mathbf{R}}^{S_k}\}$ similar to the human demonstrated placement slot. Then, 2D keypoint matching and 3D lifting enable the calculation of a single 6-DoF object transformation matrix $\mathbf{T}_i \in SE(3)$ for each detected slot $\mathbf{M}_{\mathbf{R}}^{S_i}$ for downstream robotic pick-and-place.

Object and slot re-identification. Taking as input a short304video with only two frames $\{H_1, R\}$, SAM2 [55] is employed to output the object mask M^O_R and one best-matched305slot mask $M^{S_0}_R$ over the robot image R, given the detected307

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

Figure 5. **Correlate with robot view.** Given the object and slot mask detected in the human video, we first re-identify the corresponding object and slot in robot view, and also find all similar empty slots. With corresponding object masks and slot masks, we first compute 2D keypoint matching among the detected object and mask local patches and then lift the observations to 3D to compute 6-DoF transforms.

3082D object mask $\mathbf{M}_{\mathbf{H}_1}^O$ and the slot mask $\mathbf{M}_{\mathbf{H}_1}^S$ on the human309first frame image \mathbf{H}_1 . If multiple similar slots are present in310the robot image, we leverage SAM [31] to propose segment311candidates and use DINOv2 [48] to collect additional slot312masks { $\mathbf{M}_{\mathbf{R}}^{S_1}, \cdots, \mathbf{M}_{\mathbf{R}}^{S_k}$ } that share similar DINOv2 fea-313tures with the detected slot mask $\mathbf{M}_{\mathbf{R}}^{S_0}$. Empirically, we find314that SAM2 and DINOv2 provide good enough performance315on our data.

2D keypoint matching. With two corresponding masks in 316 the human view and robot view, we use MASt3R [33] to de-317 tect 2D keypoint correspondences by expanding the masks 318 into local 2D bounding boxes. As illustrated in Fig. 5 (mid-319 320 dle left), we compute the 2D keypoint matching on two 321 pairs of object local patches (between the object mask $\mathbf{M}_{\mathbf{R}}^{O}$ in the robot frame and the object mask $\mathbf{M}^{O}_{\mathbf{H}_{1}}$ in the initial 322 human frame, and between the object mask $\mathbf{M}^{O}_{\mathbf{H}_{1}}$ in the 323 initial human frame and the object mask $\mathbf{M}_{\mathbf{H}_{n}}^{O}$ in the last 324 human frame) and one pair of slot local patches (between 325 the slot $\mathbf{M}_{\mathbf{H}_1}^S$ in the initial human frame and the slot $\mathbf{M}_{\mathbf{R}}^{S_i}$ 326 in the robot view for any slot S_i to place). 327

328 3D lifting and transformation calculation. Using the depth sensing and the camera intrinsic parameters, we can 329 lift all human and robot view images into 3D point cloud 330 331 observations. Then, we are able to lift the 2D keypoint correspondences into 3D correspondences. Equipped with the 332 3D correspondences, we use Procrustes analysis [21] with 333 334 RANSAC [19] to calculate three 6-DoF transformation matrices for the aforementioned three local patch pair match-335 ings. We denote the three computed 6-DoF transformations 336 as $T^O_{\mathbf{R}\to\mathbf{H}}$, $T^O_{\mathbf{H}}$, and $T^{S_i}_{\mathbf{H}\to\mathbf{R}}$ respectively. Fig. 5 (middle right) illustrates their geometric meanings: the object trans-337 338 339 formation from the robot scene to the human scene at the start of the human video, the transformation applied by the 340 person to the picked object in the human video, and the slot 341 transformation from the human scene to the robot scene. 342

Final object placement transformations. As clearly illustrated in Fig. 5 (middle right and rightmost), by chaining up
the three 6-DoF transformation matrix explained above, we

can compute the final desired 6-DoF transformation matrix for the robot to execute in order to transform the pick object O from its initial position to any target slot i in the robot coordinate frame as the following

$$\mathbf{\Gamma}_{i} = \mathbf{T}_{\mathbf{H}\to\mathbf{R}}^{S_{i}} \mathbf{T}_{\mathbf{H}}^{O} \mathbf{T}_{\mathbf{R}\to\mathbf{H}}^{O}.$$
 (1) 350

5. Experiments

We propose a new dataset and present an extensive evaluation of our system in Sec. 5.1, where our system, SLeRP, outperforms the baselines by a large margin. In Sec. 5.2, we present an in-depth ablation over Slot-Net and additional design choices in SLeRP. In Sec. 5.3, we show that SLeRP is effective with real-world robots.

5.1. System Evaluation

Given the novelty of the problem we address, existing evaluation benchmarks are unavailable, and there are no baseline methods with which to make direct comparisons. Consequently, we have curated a dataset comprising real-world videos and established a benchmark specific to this problem by developing suitable baselines and metrics.

Dataset. We collected 288 real-world RGB-D videos span-366 ning 9 different object-in-slot task scenarios. For each sce-367 nario, variations were introduced in the background and the 368 inclusion of distractor objects, camera positions, and slot 369 occupancy conditions. The suite of tasks includes challeng-370 ing, common daily activities such as putting bread into a 371 toaster, arranging eggs in an egg steamer, and setting mugs 372 on coasters. All the objects are unseen to Slot-Net during 373 training, and 3 out of the 9 tasks encompass previously un-374 seen task categories. Visualizations of the tasks and their 375 varying settings are provided in the supplementary material. 376

Benchmark setup. Given that our task necessitates paired377data comprising a human demonstration video and a novel378image for the robot's view, we construct test pairs by repairing the videos in our dataset. Each pair comprises380videos depicting the same object being placed into the same381slot, albeit with potential variations in background, camera angle, and initial slot occupancy. We designate the first383

346

347

348

349

352

353

354

355

356

357

358

359

360

361

362

363

364

Figure 6. **Qualitative Comparison.** We compare our method to baselines and present side-by-side results on three examples. For each example, the first column shows the input human video at the top and robot-view image in the bottom. The top row displays 2D reidentification results (object in yellow, slot in red), while the bottom row shows 6-DoF relative pose predictions by projecting the object point cloud onto the slots. Unlike the baselines that can only predict one exact slot, our approach can also identify multiple slots. These results clearly demonstrate that our system outperforms the baselines, achieving accurate slot and transformation predictions.

video in each pair as the human demonstration and employ
the initial frame of the second video as the robot's view, as
illustrated in Fig. 6. During video data collection, we ensure
that human hands are absent in the first frame. We generate
three distinct test splits, introducing variations in viewpoint
(288 video pairs), background (720 video pairs), and slot
occupancy (288 video pairs).

391 Metrics. We evaluate the accuracy of the predicted 2D 392 masks and the 6-DoF transformation matrix using five dis-393 tinct metrics. For the evaluation of 2D masks, we calculate (1) the intersection-over-union (IoU) for the object mask 394 and (2) the IoU of the exact slot mask in the robot view, 395 comparing them to their respective ground-truth masks. We 396 assess the accuracy of the 6-DoF transformation by trans-397 398 forming and projecting the object's point cloud onto the robot view, then measuring (3) the precision of the mask 399 against the ground-truth mask. For 3D evaluation, we com-400 pute (4) the Chamfer Distance and (5) Earth Mover's Dis-401 402 tance [17] between the transformed object point cloud and the ground-truth object point cloud at placement. If no mask 403 or transformation output is predicted for any method, we 404 use a default empty mask and identity matrix as the fall-405 back predictions. To establish ground truths, we annotate 406 the 2D masks of the object and exact slot in the start video 407 408 frame (i.e., robot's view), alongside the object's mask postplacement in the end frame.

Baselines. We design four baselines for comparison.

- *ORION.* Zhu *et al.* [86] perform vision-based human-torobot imitation learning focused on object-level placement. We adapt ORION to our slot-level setting by providing the required ground-truth object and slot names. In contrast, our method automatically detects in-contact objects and slots without the need of explicit name inputs.
- *ORION*++. We leverage the object and slot detection results from SLeRP to enhance ORION, thereby establishing a stronger baseline.
- *CLIPort*++. CLIPort [61] is an end-to-end imitationlearning-based language-conditioned policy for tabletop tasks. However, the original method requires videos with action labels for training, whereas ours does not. To construct a comparison, we randomly split the tasks into training and test sets, ensuring that tasks, objects, and scenes are unseen during testing, and use the training split to train CLIPort (more details in Supplementary).
- *VideoCap+FMs.* We test whether our proposed tasks can be effectively solved using state-of-the-art foundation models. We utilize Qwen2VL [72], a video captioning model, to summarize the video, and then employ GPT40 [2] to identify the object and slot names. Subsequently, we use grounding-SAM2 [56] to generate a

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

Method		Different view					Different background					Different slot occupancy				
	Obj↑	Slot↑	Prec. [↑]	CD↓	$\text{EMD}{\downarrow}$	Obj↑	Slot↑	Prec. [↑]	CD↓	$EMD{\downarrow}$	Obj↑	Slot↑	$\text{Prec.}\uparrow$	CD↓	EMD↓	
ORION [86]	0.00	0.11	0.41	0.0949	0.0552	0.45	0.13	0.00	0.0932	0.0540	0.40	0.12	0.82	0.0952	0.0559	
ORION++ [86]	10.21	8.89	2.10	0.1058	0.0584	7.53	5.97	0.75	0.1113	0.0597	8.83	7.44	2.95	0.1021	0.0576	
CLIPort++ [61]	1.54	0.47	18.75	0.3887	0.1663	1.51	0.35	1.71	0.1152	0.0615	3.06	0.34	12.50	0.1348	0.0681	
VideoCap+FMs	2.35	8.61	13.45	0.1918	0.0987	2.12	6.28	13.61	0.1743	0.0917	2.64	9.84	7.25	0.1508	0.0837	
Ours	73.85	54.37	36.40	0.0282	0.0182	70.27	44.70	25.39	0.0573	0.0323	68.12	47.04	30.30	0.0334	0.0223	

Table 1. Quantitative System Evaluation. We compare SLeRP with baselines and report the 2D detection and the 3D object transformation accuracy: IoU for the object mask prediction (Obj); IoU for the slot mask prediction (Slot); mask precision of the predicted object after placement onto the slot projected to the camera plane (*Prec.*); and Chamfer distance (*CD*) and Earth-Mover distance (*EMD*) between the predicted and ground-truth target object point clouds after placement. We evaluate in three different settings with the robot's views having different camera viewpoints, scene backgrounds, and initial states of the placement slot occupancy compared to the input human videos. We find that SLeRP substantially outperforms the baselines by large margins across all the metrics in all the three evaluation settings.

Mathad	Syn	thetic	Real (se	en tasks)	Real (un	Real (unseen tasks		
Method	F1 ↑	IoU↑	F1 ↑	IoU ↑	F1 ↑	IoU ↑		
Image difference	44.10	31.34	32.69	20.04	32.20	19.53		
Change detection [18]	2.27	1.58	0.06	0.03	6.30	4.19		
Object mask	60.16	48.73	65.90	52.10	58.67	42.98		
Object-box mask	-	-	43.39	37.74	39.37	34.65		
GPT40+SAM	0.91	0.49	4.82	2.84	4.45	2.68		
Slot-Net (end image)	83.44	74.57	48.83	38.68	38.19	28.37		
Slot-Net (ours)	82.89	74.62	73.27	61.59	66.50	54.26		

Table 2. **Slot Segmentation Results.** We compare Slot-Net against various alternative approaches on slot detection. We evaluate on test synthetic images and our collected real-world images (seen and unseen tasks). Dashes note that the method cannot be evaluated for synthetic data given no video inputs. We can observe cleary that our Slot-Net performs the best.

434 435 436 bounding box and employ SAM2 [55] to produce the object and slot mask, with the transformation matrix computed using modules in the same way as in our system.

437 **Results.** Table1 presents a quantitative evaluation comparing SLeRP to four baseline methods. The results indi-438 cate that SLeRP significantly outperforms all baselines by 439 considerable margins. Fig.6 offers qualitative comparisons, 440 showcasing 2D object and slot mask predictions alongside 441 3D object transformation estimations. We observe that 442 443 SLeRP generates more accurate 2D object and slot detection results, as baseline methods such as ORION and Video-444 445 Cap+FMs frequently struggle to describe slot names in natural language for subsequent visual recognition (e.g., incor-446 rectly detecting the entire placement object or table). Ad-447 448 ditionally, SLeRP achieves more precise 3D object transformations compared to baselines like CLIPort, which are 449 450 primarily designed for top-down 2D predictions.

Lastly, our method can fill multiple slots, whereas other 451 methods generate output for only a single slot. In a sub-452 453 set of the data with ground-truth annotations for multiple slots, SLeRP achieves IoU scores of 67.70 and 42.62 for 454 2D object and slot segmentation. Additionally, it attains 455 scores of 23.14, 0.0575, and 0.0350 for 3D transformation 456 predictions in terms of slot projection precision, Chamfer 457 Distance (CD), and Earth Mover's Distance (EMD), respec-458 459 tively. These results are comparable to the single-slot place-

Method (diff. view)	SlotNet	SAM2	Mast3r	Obj ↑	Slot \uparrow	Prec. \uparrow	$\mathrm{CD}\downarrow$	$\text{EMD}\downarrow$
Base design	×	X	×	28.53	28.23	24.72	0.2289	0.1183
Ours w/o SlotNet	×	1	1	73.85	29.63	27.84	0.0486	0.0289
Ours w/o SAM2	1	×	1	31.05	35.98	23.75	0.1219	0.0667
Ours w/o Mast3r	1	1	×	73.85	54.37	32.74	0.0321	0.0205
Ours	1	1	1	73.85	54.37	36.40	0.0282	0.0182

Table 3. **Ablation Study.** All metrics follow Table 1. Results show all the key modules help. See supplementary for the full table.

ment evaluations reported in Table 1. See the supplementary materials for further details. 461

5.2. Ablation Study

We evaluate the effectiveness of Slot-Net for placement slot detection in Sec. 5.2.1 and additional ablations on other components like object and slot re-identification and keypoint matching in Sec. 5.2.2.

5.2.1. Slot-Net Ablations

Baselines. We consider the following alternative and ablation approaches to replace Slot-Net.

- *Image difference*. We use the difference image between the gray-scale start and end frame, then apply threshold-ing to the difference image to obtain a mask.
- *Change detection.* We use an off-the-shelf change detection model [18] given two frames for the masks.
- *Object mask.* We directly use the ground-truth pick object mask as the slot mask prediction.
- *Object-box mask.* We take the pick object bounding box detected in the tracking procedure and query SAM for a proxy placement slot mask.
- *GPT4o+SAM*. We query GPT4o with start and end frames for slot bounding boxes and query SAM for masks.
- *Slot-Net (end image)*. We use the end frame to replace the difference image as the prompt for Slot-Net.

Benchmark and metrics. We use our newly proposed484real-world video dataset along with held-out synthetic images for evaluation. For the real images, we evaluate two486splits: seen tasks, which involve seen object categories during training but novel object instances, and unseen tasks,488featuring object categories not encountered during training.489

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

Figure 7. Slot Detection Comparisons. We compare Slot-Net to many alternative approaches, and results show that ours performs better.

Figure 8. **Real-Robot Experiments.** We show real-robot experiments for "block into a container" and "strawberry into an organizer". See supplementary for videos and more examples.

For each real-world video, we pair the starting and ending
frames as input and evaluate predictions against the groundtruth slot mask from the starting frame. The slot mask prediction performance is assessed using IoU and F1 scores.

Results. Table 2 shows that Slot-Net, when trained on syn-494 thetic data, generalizes effectively to real images, outper-495 forming alternative approaches. Fig. 7 provides side-by-496 side comparisons of different methods, revealing that Slot-497 Net excels in identifying slot boundaries of various shapes. 498 This underscores the necessity of training a custom model 499 for slot detection and demonstrates that our model is both 500 well-designed and effective. 501

502 5.2.2. Other Ablations

Beyond ablating Slot-Net, we further validate two addi-503 504 tional key design elements in our system: utilizing SAM2 505 for object and slot re-identification and employing MASt3R for keypoint matching. In the absence of SAM2, we rely 506 on DINOv2 feature similarity between the slot mask and 507 all SAM-generated masks in the robot image. To re-508 place MASt3R, we employ DINOv2 features for Hungarian 509 510 matching. Table 9 shows the necessity of these modules.

5.3. Real-Robot Experiments

As shown in Fig. 8, we perform real-robot experiments with 512 a Franka robot and show that SLeRP is effective for real 513 robots. The manipulation system employs a wrist-mounted 514 RGB-D camera (Realsense D415) and an external RGB-D 515 camera (Realsense L515). The camera intrinsics and ex-516 trinsics relative to the robot are known. The wrist-mounted 517 camera provides observations for SLeRP, while the exter-518 nal camera observes the entire scene and aids in planning 519 collision-free trajectories. The system utilizes Contact-520 Graspnet[68] to generate grasps and plans collision-free tra-521 jectories using methods described in [13, 44]. 522

#

511

523

6. Conclusion

In this paper, we address the novel problem of slot-level 524 object placement by learning from a single human demon-525 stration video. We propose a modular system to tackle this 526 problem, which operates without requiring additional train-527 ing video data and features a unique slot-level placement 528 detector. To evaluate the system's performance, we intro-529 duce a new benchmark consisting of real-world videos and 530 compare our system against key baseline methods. Our re-531 sults demonstrate that SLeRP outperforms these baselines 532 and functions effectively in real-robot experiments. 533

Limitations and future work. Given the novel problem 534 formulation, there is potential for further research in fine-535 grained slot-level object placement with minimal or no hu-536 man demonstrations. Future work could focus on relax-537 ing current system assumptions, such as the static camera, 538 single-handed interaction, and minimal motion of the place-539 ment object. Moreover, advancements in visual foundation 540 models could enhance the robustness of our system, as they 541 play a crucial role in this work. 542

References 543

- 544 [1] Advanced inference: Outpainting. https : 545 / / huggingface . co / docs / diffusers / en / 546 advanced inference / outpaint. Accessed: 547 2010-09-30. 4. 13. 14
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ah-548 549 mad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, 550 Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 551 Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 552 2023. 6
- 553 [3] Shikhar Bahl, Abhinav Gupta, and Deepak Pathak. 554 Human-to-robot imitation in the wild. arXiv preprint 555 arXiv:2207.09450, 2022. 1, 3
- [4] Shikhar Bahl, Russell Mendonca, Lili Chen, Unnat Jain, and Deepak Pathak. Affordances from human videos as a versatile representation for robotics. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 13778-13790, 2023. 3 560
- 561 [5] Johannes Baumgartl, Tim Werner, Per Kaminsky, and Do-562 minik Henrich. A fast, gpu-based geometrical placement 563 planner for unknown sensor-modelled objects and placement 564 areas. In 2014 IEEE International Conference on Robotics 565 and Automation (ICRA), pages 1552–1559. IEEE, 2014. 3
- 566 [6] Homanga Bharadhwaj, Abhinav Gupta, Vikash Kumar, and 567 Shubham Tulsiani. Towards generalizable zero-shot manip-568 ulation via translating human interaction plans. In 2024 569 IEEE International Conference on Robotics and Automation 570 (ICRA), pages 6904-6911. IEEE, 2024. 1, 3
- [7] Homanga Bharadhwaj, Roozbeh Mottaghi, Abhinav Gupta, and Shubham Tulsiani. Track2act: Predicting point tracks 573 from internet videos enables generalizable robot manipulation. In European Conference on Computer Vision (ECCV), 575 2024. 3
- 576 [8] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot trans-577 578 fer by combining relative and metric depth. arXiv preprint 579 arXiv:2302.12288, 2023. 13
- 580 [9] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakr-582 ishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. 583 Rt-1: Robotics transformer for real-world control at scale. 584 arXiv preprint arXiv:2212.06817, 2022. 1
- [10] Annie S Chen, Suraj Nair, and Chelsea Finn. Learning gener-585 alizable robotic reward functions from" in-the-wild" human 586 587 videos. arXiv preprint arXiv:2103.16817, 2021. 3
- 588 [11] Tianyi Cheng, Dandan Shan, Ayda Sultan Hassen, Richard 589 Ely Locke Higgins, and David Fouhey. Towards a richer 2d 590 understanding of hands at scale. In Thirty-seventh Conference on Neural Information Processing Systems, 2023. 4
- [12] Cleanup.pictures. https://cleanup.pictures. 4, 592 593 12.13
- 594 [13] Michael Danielczuk, Arsalan Mousavian, Clemens Eppner, 595 and Dieter Fox. Object rearrangement using learned implicit collision functions. In 2021 IEEE International Con-596 597 ference on Robotics and Automation (ICRA), pages 6010-598 6017. IEEE, 2021. 3, 8

- [14] Sudeep Dasari and Abhinav Gupta. Transformers for oneshot visual imitation. In Conference on Robot Learning, pages 2071-2084. PMLR, 2021. 3
- [15] Sudeep Dasari, Mohan Kumar Srirama, Unnat Jain, and Abhinav Gupta. An unbiased look at datasets for visuo-motor pre-training. In Conference on Robot Learning, pages 1183-1198. PMLR, 2023. 3
- [16] Ben Eisner, Yi Yang, Todor Davchev, Mel Vecerik, Jonathan Scholz, and David Held. Deep se (3)-equivariant geometric reasoning for precise placement tasks. In The Twelfth International Conference on Learning Representations. 3
- [17] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 605-613, 2017. 6
- [18] Sheng Fang, Kaiyu Li, and Zhe Li. Changer: Feature interaction is what you need for change detection. arXiv preprint arXiv:2209.08290, 2022. 2, 7
- [19] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. Communications of the ACM, 24(6):381-395, 1981. 5
- [20] Caelan Reed Garrett, Rohan Chitnis, Rachel Holladay, Beomioon Kim, Tom Silver, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. Integrated task and motion planning. Annual Review of Control, Robotics, and Autonomous Systems, 4(1):265-293, 2021. 1
- [21] Colin Goodall. Procrustes methods in the statistical analysis of shape. Journal of the Royal Statistical Society: Series B (Methodological), 53(2):285-321, 1991. 5
- [22] Kensuke Harada, Tokuo Tsuji, Kazuyuki Nagata, Natsuki Yamanobe, and Hiromu Onda. Validating an object placement planner for robotic pick-and-place tasks. Robotics and Autonomous Systems, 62(10):1463-1477, 2014. 3
- [23] Joshua A Haustein, Kaiyu Hang, Johannes Stork, and Danica Kragic. Object placement planning and optimization for robot manipulators. In 2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 7417-7424. IEEE, 2019. 3
- [24] Haojie Huang, Dian Wang, Robin Walters, and Robert Platt. Equivariant transporter network. In Robotics: Science and Systems, 2022. 3
- [25] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-40 system card. arXiv preprint arXiv:2410.21276, 2024. 2
- [26] Vidhi Jain, Maria Attarian, Nikhil J Joshi, Ayzaan Wahid, Danny Driess, Quan Vuong, Pannag R Sanketi, Pierre Sermanet, Stefan Welker, Christine Chan, et al. Vid2robot: End-to-end video-conditioned policy learning with crossattention transformers. arXiv preprint arXiv:2403.12943, 2024. 1, 3
- [27] Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Fred-651 erik Ebert, Corey Lynch, Sergey Levine, and Chelsea Finn. 652 Bc-z: Zero-shot task generalization with robotic imitation 653 learning. In Conference on Robot Learning, pages 991-654 1002. PMLR, 2022. 3 655

#

556

557

558

559

571

572

574

581

591

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

#

- [28] Yun Jiang, Marcus Lim, Changxi Zheng, and Ashutosh Saxena. Learning to place new objects in a scene. *The In- ternational Journal of Robotics Research*, 31(9):1021–1043,
 2012. 3
- 660 [29] Aditya Kannan, Kenneth Shaw, Shikhar Bahl, Pragna Mannam, and Deepak Pathak. Deft: Dexterous fine-tuning for
 662 hand policies. In *Conference on Robot Learning*, pages 928–
 663 942. PMLR, 2023. 3
- [30] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao,
 Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan
 Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An
 open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024. 1
- [31] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao,
 Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Con- ference on Computer Vision*, pages 4015–4026, 2023. 4, 5
- [32] Jangwon Lee and Michael S Ryoo. Learning robot activities
 from first-person human videos using convolutional future
 regression. In *Proceedings of the IEEE Conference on Com- puter Vision and Pattern Recognition Workshops*, pages 1–2,
 2017. 3
- [33] Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r. *arXiv preprint arXiv:2406.09756*, 2024. 5
- [34] Siyuan Li, Lei Ke, Martin Danelljan, Luigi Piccinelli, Mattia
 Segu, Luc Van Gool, and Fisher Yu. Matching anything by
 segmenting anything. In *Proceedings of the IEEE/CVF Con- ference on Computer Vision and Pattern Recognition*, pages
 18963–18973, 2024. 4
- [35] Weiyu Liu, Chris Paxton, Tucker Hermans, and Dieter Fox.
 Structformer: Learning spatial structure for language-guided
 semantic rearrangement of novel objects. In 2022 Inter-*national Conference on Robotics and Automation (ICRA)*,
 pages 6322–6329. IEEE, 2022. 3
- [36] YuXuan Liu, Abhishek Gupta, Pieter Abbeel, and Sergey
 Levine. Imitation from observation: Learning to imitate
 behaviors from raw video via context translation. In 2018 *IEEE international conference on robotics and automation*(*ICRA*), pages 1118–1125. IEEE, 2018. 3
- [37] Tomás Lozano-Pérez, Joseph L. Jones, Emmanuel Mazer,
 and Patrick A. O'Donnell. Task-level planning of pick-andplace robot motions. *Computer*, 22(3):21–29, 1989. 3
- [38] Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy Zhang. Vip: Towards universal visual reward and representation via value-implicit pre-training. In *The Eleventh International Conference on Learning Representations*. 3
- [39] Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Jason Ma, Claire Chen, Sneha Silwal, Aryan Jain, Vincent-Pierre Berges, Tingfan Wu, Jay Vakil, et al. Where are we in the search for an artificial visual cortex for embodied intelligence? *Advances in Neural Information Processing Systems*, 36:655–677, 2023. 3
- [40] Zhao Mandi, Fangchen Liu, Kimin Lee, and Pieter Abbeel.
 Towards more generalizable one-shot visual imitation learn-

ing. In 2022 International Conference on Robotics and Automation (ICRA), pages 2434–2444. IEEE, 2022. 3 714

- [41] Priyanka Mandikal and Kristen Grauman. Dexvip: Learning dexterous grasping with human hand pose priors from video.
 In *Conference on Robot Learning*, pages 651–661. PMLR, 2022. 3
- [42] Oier Mees, Alp Emek, Johan Vertens, and Wolfram Burgard. Learning object placements for relational instructions by hallucinating scene representations. In 2020 IEEE International Conference on Robotics and Automation (ICRA), pages 94– 100. IEEE, 2020. 3
- [43] Chaitanya Mitash, Rahul Shome, Bowen Wen, Abdeslam Boularias, and Kostas Bekris. Task-driven perception and manipulation for constrained placement of unknown objects. *IEEE Robotics and Automation Letters*, 5(4):5605–5612, 2020. 3
- [44] Adithyavairavan Murali, Arsalan Mousavian, Clemens Eppner, Adam Fishman, and Dieter Fox. Cabinet: Scaling neural collision detection for object rearrangement with procedural scene generation. In 2023 IEEE International Conference on Robotics and Automation (ICRA), pages 1866–1874. IEEE, 2023. 3, 8
- [45] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation. In *Conference on Robot Learning*, pages 892–909. PMLR, 2023. 3
- [46] Rhys Newbury, Kerry He, Akansel Cosgun, and Tom Drummond. Learning to place objects onto flat surfaces in upright orientations. *IEEE Robotics and Automation Letters*, 6(3): 4377–4384, 2021. 3
- [47] Anh Nguyen, Dimitrios Kanoulas, Luca Muratore, Darwin G Caldwell, and Nikos G Tsagarakis. Translating videos to commands for robotic manipulation with deep recurrent neural networks. In 2018 IEEE International Conference on Robotics and Automation (ICRA), pages 3782–3788. IEEE, 2018. 3
- [48] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*. 5
- [49] Helena Örnkloo and Claes von Hofsten. Fitting objects into holes: on the development of spatial cognition skills. *Devel*opmental psychology, 43(2):404, 2007. 1
- [50] Chuer Pan, Brian Okorn, Harry Zhang, Ben Eisner, and David Held. Tax-pose: Task-specific cross-pose estimation for robot manipulation. In *Conference on Robot Learning*, pages 1783–1792. PMLR, 2023. 3
- [51] Austin Patel, Andrew Wang, Ilija Radosavovic, and Jitendra Malik. Learning to imitate object interactions from internet videos. arXiv:2211.13225, 2022. 3
- [52] Chris Paxton, Chris Xie, Tucker Hermans, and Dieter Fox. Predicting stable configurations for semantic placement of novel objects. In 5th Annual Conference on Robot Learning. 3
- [53] Dustin Podell, Zion English, Kyle Lacey, AndreasBlattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and769

804

805

806

- [54] Yuzhe Qin, Yueh-Hua Wu, Shaowei Liu, Hanwen Jiang, Ruihan Yang, Yang Fu, and Xiaolong Wang. Dexmv: Imitation learning for dexterous manipulation from human videos. In *European Conference on Computer Vision*, pages 570–587.
 Springer, 2022. 3
- [55] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang
 Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman
 Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2:
 Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 4, 7
- [56] Tianhe Ren, Qing Jiang, Shilong Liu, Zhaoyang Zeng, Wenlong Liu, Han Gao, Hongjie Huang, Zhengyu Ma, Xiaoke Jiang, Yihao Chen, Yuda Xiong, Hao Zhang, Feng Li, Peijun Tang, Kent Yu, and Lei Zhang. Grounding dino 1.5: Advance the "edge" of open-set object detection, 2024. 6
- [57] Brian Rooks. Machine tending in the modern age. *Industrial Robot: An International Journal*, 30(4):313–318, 2003. 1
- [58] Jinghuan Shang, Karl Schmeckpeper, Brandon B May,
 Maria Vittoria Minniti, Tarik Kelestemur, David Watkins,
 and Laura Herlant. Theia: Distilling diverse vision foundation models for robot learning. In 8th Annual Conference
 on Robot Learning. 3
- [59] Aditya Sharma, Luke Yoffe, and Tobias Höllerer. Octo+:
 A suite for automatic open-vocabulary object placement
 in mixed reality. In 2024 IEEE International Conference
 on Artificial Intelligence and eXtended and Virtual Reality
 (AIxVR), pages 157–165. IEEE, 2024. 3
- [60] Kenneth Shaw, Shikhar Bahl, and Deepak Pathak. Videodex:
 Learning dexterity from internet videos. In *Conference on Robot Learning*, pages 654–665. PMLR, 2023. 3
 - [61] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic manipulation. In *Proceedings of the 5th Conference on Robot Learning (CoRL)*, 2021. 2, 3, 6, 7
- 807 [62] Anthony Simeonov, Yilun Du, Andrea Tagliasacchi,
 808 Joshua B Tenenbaum, Alberto Rodriguez, Pulkit Agrawal,
 809 and Vincent Sitzmann. Neural descriptor fields: Se (3)810 equivariant object representations for manipulation. In
 811 2022 International Conference on Robotics and Automation
 812 (ICRA), pages 6394–6400. IEEE, 2022. 3
- [63] Anthony Simeonov, Yilun Du, Yen-Chen Lin, Alberto Rodriguez Garcia, Leslie Pack Kaelbling, Tomás Lozano-Pérez, and Pulkit Agrawal. Se (3)-equivariant relational rearrangement with neural descriptor fields. In *Conference on Robot Learning*, pages 835–846. PMLR, 2023. 3
- [64] Himanshu Gaurav Singh, Antonio Loquercio, Carmelo Sferrazza, Jane Wu, Haozhi Qi, Pieter Abbeel, and Jitendra Malik. Hand-object interaction pretraining from videos. *arXiv preprint arXiv:2409.08273*, 2024. 3
- 822[65] TorontoAnnotationSuite.823https://aidemos.cs.toronto.edu/toras. 4, 12, 13
- [66] Yu Sun, Shaogang Ren, and Yun Lin. Object–object interaction affordance learning. *Robotics and Autonomous Systems*, 62(4):487–496, 2014. 3

- [67] Balakumar Sundaralingam, Siva Kumar Sastry Hari, Adam Fishman, Caelan Garrett, Karl Van Wyk, Valts Blukis, Alexander Millane, Helen Millane, Ankur Handa, Fabio Ramos, Nathan Ratliff, and Dieter Fox. CuRobo: Parallelized collision-free robot motion generation. In *ICRA*, 2023. 1
 827
 828
 829
 830
 831
 832
- [68] Martin Sundermeyer, Arsalan Mousavian, Rudolph Triebel, and Fox Dieter. Contact-graspnet: Efficient 6-dof grasp generation in clutteredscenes. *IEEE International Conference on Robotics and Automation (ICRA)*, 2021. 8
- [69] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023. 4, 13, 14, 18, 19
- [70] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu, and Thomas Wolf. Diffusers: State-of-the-art diffusion models. https://github.com/huggingface/ diffusers, 2022. 4, 12, 13
- [71] Chen Wang, Linxi Fan, Jiankai Sun, Ruohan Zhang, Li Fei-Fei, Danfei Xu, Yuke Zhu, and Anima Anandkumar. Mimicplay: Long-horizon imitation learning by watching human play. *arXiv preprint arXiv:2302.12422*, 2023. 3
- [72] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. arXiv preprint arXiv:2409.12191, 2024. 6
- [73] Yixuan Wang, Mingtong Zhang, Zhuoran Li, Tarik Kelestemur, Katherine Rose Driggs-Campbell, Jiajun Wu, Li Fei-Fei, and Yunzhu Li. D3 fields: Dynamic 3d descriptor fields for zero-shot generalizable rearrangement. In 8th Annual Conference on Robot Learning, 2024. 3
- [74] Chuan Wen, Xingyu Lin, John So, Kai Chen, Qi Dou, Yang Gao, and Pieter Abbeel. Any-point trajectory modeling for policy learning. In *Robotics: Science and Systems*, 2024. 3
- [75] Hongtao Wu, Jikai Ye, Xin Meng, Chris Paxton, and Gregory S Chirikjian. Transporters with visual foresight for solving unseen rearrangement tasks. In 2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 10756–10763. IEEE, 2022. 3
- [76] Haoyu Xiong, Quanzhou Li, Yun-Chun Chen, Homanga Bharadhwaj, Samarth Sinha, and Animesh Garg. Learning by watching: Physical imitation of manipulation skills from human videos. In 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 7827–7834. IEEE, 2021. 1, 3
- [77] Mengda Xu, Zhenjia Xu, Yinghao Xu, Cheng Chi, Gordon Wetzstein, Manuela Veloso, and Shuran Song. Flow as the cross-domain manipulation interface. In 8th Annual Conference on Robot Learning. 3
- [78] Yezhou Yang, Yi Li, Cornelia Fermuller, and Yiannis Aloimonos. Robot learning manipulation action plans by" watch 883

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

970

971

972

973

974

975

884 885 886

ceedings of the AAAI conference on artificial intelligence, 2015. 3

ing" unconstrained videos from the world wide web. In Pro-

- [79] Seonghyeon Ye, Joel Jang, Byeongguk Jeon, Sejune Joo, 887 Jianwei Yang, Baolin Peng, Ajay Mandlekar, Reuben Tan, 888 889 Yu-Wei Chao, Bill Yuchen Lin, et al. Latent action pretrain-890 ing from videos. arXiv preprint arXiv:2410.11758, 2024. 3
- 891 [80] Tianhe Yu, Chelsea Finn, Sudeep Dasari, Annie Xie, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot im-892 893 itation from observing humans via domain-adaptive meta-894 learning. Robotics: Science and Systems XIV, 2018. 3
- [81] Chengbo Yuan, Chuan Wen, Tong Zhang, and Yang Gao. 895 896 General flow as foundation affordance for scalable robot 897 learning. arXiv preprint arXiv:2401.11439, 2024. 3
- [82] Wentao Yuan, Adithyavairavan Murali, Arsalan Mousavian, 898 and Dieter Fox. M2t2: Multi-task masked transformer for 899 object-centric pick and place. In 7th Annual Conference on 900 901 Robot Learning. 3
- [83] Kevin Zakka, Andy Zeng, Johnny Lee, and Shuran Song. 902 903 Form2fit: Learning shape priors for generalizable assembly from disassembly. In 2020 IEEE International Confer-904 905 ence on Robotics and Automation (ICRA), pages 9404-9410. 906 IEEE, 2020. 3
- 907 [84] Andy Zeng, Pete Florence, Jonathan Tompson, Stefan 908 Welker, Jonathan Chien, Maria Attarian, Travis Armstrong, 909 Ivan Krasin, Dan Duong, Vikas Sindhwani, et al. Transporter 910 networks: Rearranging the visual world for robotic manipu-911 lation. In Conference on Robot Learning, pages 726-747. 912 PMLR, 2021. 3
- 913 [85] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding 914 conditional control to text-to-image diffusion models. In Proceedings of the IEEE/CVF International Conference on 915 Computer Vision, pages 3836-3847, 2023. 13 916
- 917 [86] Yifeng Zhu, Arisrei Lim, Peter Stone, and Yuke Zhu. Visionbased manipulation from single human video with open-918 world object graphs. arXiv preprint arXiv:2405.20321, 919 920 2024. 1, 2, 3, 6, 7

Contents 921

922	1. Introduction	1
923	2. Related Work	3
924	3. Problem Formulation	3
925	4. Method	3
926	4.1. System Overview	3
927	4.2. Parsing the Human Demonstration Video	4
928	4.3. Correlating to the Robot-view Image	4
929	5. Experiments	5
930	5.1. System Evaluation	5
931	5.2. Ablation Study	7
932	5.2.1. Slot-Net Ablations	7
933	5.2.2. Other Ablations	8
934	5.3 Real-Robot Experiments	8

6. Conclusion	8	935
A Slot-Net Data Generation	12	936
A.1. Details, Statistics and Examples	12	937
A.2 Synthetic Data Generation Process	13	938
A.3 Limitations	14	939
B Evaluation Videos	14	940
B.1. Visualization of 9 tasks	14	941
B.2. Visualization of 3 settings	14	942
C Additional Implementation Details	14	943
C.1. Slot-Net	14	944
C.2. CLIPort++	14	945
D Additional Results of SLeRP	14	946
D.1. Full Ablation Study Table	14	947
D.2 Multi-slots analysis.	14	948
D.3. More baseline comparison results	15	949
D.4. More ablation comparison results	15	950
D.5 Limitations	15	951
E Additional Results of Slot-Net	15	952
E.1. Additional results of Slot-Net	15	953
E.2. Limitations	15	954
A. Slot-Net Data Generation		955
A.1. Details, Statistics and Examples		956
Object-with-slots image collection We collect obj	iect-	957

Object-with-slots image collection. We collect object centric images of items with slots by photographing them 958 in everyday environments. These images serve as the end 959 object crops, and we generate corresponding start images 960 by removing objects, thus creating (start, end) image pairs 961 for a pick-and-place action. During data collection, partic-962 ipants were instructed to partially fill the slots. This ap-963 proach ensures that slots are not entirely empty, facilitating 964 object removal, while avoiding fully filled slots, which can 965 complicate inpainting during object removal. In total, we 966 have collected 2,138 object-centric images spanning 67 cat-967 egories, as detailed in Tab. 4. 968

Pick-object removal. The captured image serves as an object-centric end image. To create an object-centric start image, where objects are absent from the slot, we employ the SDXL [70] inpainting model to remove the pick-object. For cases where SDXL does not successfully remove the object, we use Cleanup.pictures [12] in a subsequent round to achieve more effective removal results.

Slot mask annotation. After obtaining the (start, end) 976 object-centric images, we annotate the slot mask by com-977 paring each pair using TORAS [65]. We perform mask an-978 notation prior to augmentation and automatically transform 979 these masks onto the canvas, enabling us to generate labels 980 for a large dataset. 981

Task Name	#(obj centric)	#synthetic
battery in batter compartment	20	50
bead in bead organizer	41	4,291
block in a toy vehicle	12	102
book in book holder	22	2,207
bottle in box	28	2927
bottle in organizer	87	7,810
bread in toaster	84	3,896
can in box	5	50
can on board	15	557
candle in organizer	5	450
capsule in tray	20	1 966
cards in organizer	5	51
choclate in organier	22	200
coin in organizer	6	200 451
cup in cardboard trav	101	4.51
cup in cardboard tray	5	245
cup on coaster	5	243
cup on coaster	5	199
cylinder in organizer	49	1,705
egg in egg carton	205	36,758
egg in egg steamer	13	649
egg in tray	22	2,116
flower in organier	15	250
food in organizor	123	9,117
food in tray	10	1,250
fruit in box	5	450
fruit in food organizer	35	3,146
fruit in organizer	155	11,303
glass in organizer	5	400
glass on coaster	7	200
glove in organizer	12	400
ice cube in ice tray	31	661
jewerly in organizer	18	960
key in organizer	6	900
lip stick in organizer	14	256
muffin in muffin tray	63	4,150
mug on board	10	310
mug on coaster	66	2,453
notepad in box	15	901
peach in box	12	903
peg in wood base	8	1,413
pen in basket	15	793
pen in cup	20	994
pen in organizer	103	8.690
pen in pen holder	76	3.305
pen on book	30	1.335
pepper in trav	10	492
pill in pill organizer	15	462
plant in vase	12	387
rectangle box in tray	35	366
tangerine in muffin tray	33 24	1 700
tool in organizer	2 4 61	400
tou ar in organizer	14	400
toy car in organizer	14	003
tube in tray	15	1,400
utensil in bowl	20	850
utensil in cup	19	2,311
utensil in utensil organizer	146	14,976
wood in organizer	9	1,100

Table 4. Statistics for object-centric images. This table shows the task names for the object-centric images we collected (col 1), 13 the amount of object-centric images (col 2), and the amount of corresponding generated images (col 3).

Figure 9. **Data augmentation.** We show 9 synthetic data augmentation examples for the "egg in egg carton" task. In training data generation, we use x100 augmentation for each slot in the objectcentric images.

Augmentation with outpainting. For each object-centric 982 slot object, we create a larger scene by first sampling ran-983 dom locations on a 1024x1024 canvas and then apply-984 ing outpainting with various generated prompts. Using a 985 short category name, such as "bread in toaster", we em-986 ploy Llama [69] to enrich the text prompts with descriptions 987 of the environments. We apply outpainting to the back-988 ground following the Hugging Face Outpainting script [1] 989 with enriched texts to create the end image. The outpaint-990 ing script incorporates ControlNet [85], SDXL [70], and 991 ZoeDepth [8]. Subsequently, the start object crop images 992 and slot masks are subjected to the same transformation to 993 create the outpainted (start, end) image pairs and the cor-994 responding ground-truth slot masks. In total, we apply 100 995 augmentations for each slot on the object-centric image, re-996 sulting in 156,000 images for training. Examples of "egg in 997 egg carton" augmentation are shown in Fig. 9. 998

A.2. Synthetic Data Generation Process

We present examples of our synthetic generation process in Tab. 5 and Tab. 6. Each row in the table, from left to right, illustrates the step-by-step process for one example.

Beginning with a task name (col 1) and an object-1003 centric image (col 3), we first enhance the description us-1004 ing Llama [69] (col 2). Next, we perform object removal 1005 from the slots through two rounds of image inpainting us-1006 ing SDXL [70] and Clean.pictures [12] (col 4) to ensure 1007 more effective results. We annotate the slot mask on the 1008 object-centric image utilizing the Toronto Annotation Suite 1009 (TORAS) [65] (col 5). 1010

For generating diverse daily backgrounds, we transform 1011

999

1000

1001

1066

1070

1071

#

1012 the object-centric image (col 3) onto a 1024x1024 can-1013 vas and outpaint the background using the Hugging Face Outpainting script [1] with the enriched text generated by 1014 Llama [69] (col 2) as a prompt to create the end image. Fi-1015 1016 nally, the start image and ground-truth mask undergo the same transformation onto a 1024x1024 canvas to produce 1017 the (start, end, mask) triplet samples for training. 1018 A.3. Limitations 1019 1020 Our data generation pipeline for creating (start, end) pairs with minimal data collection and annotation is highly effi-1021

cient, enabling us to train SAM with manageable effort. 1022 Despite its efficiency, the generated images have cer-1023 tain limitations. First, the diversity of the generated image 1024 styles is insufficient. Although we employ Llama to enrich 1025 text prompts for backgrounds and request diverse styles, the 1026 generated images lack sufficient diversity and realism. Sec-1027 ond, the generated images sometimes defy physical plausi-1028 bility. In the process of outpainting the object-centric im-1029 1030 ages, we sample random locations and rotations on the canvas. Occasionally, these locations or rotations are challeng-1031 ing to outpaint while adhering to physical principles, result-1032 ing in images that appear unrealistic. 1033

1034The purpose of outpainting data augmentation is to gen-1035erate more varied backgrounds, thereby aiding Slot-Net in1036achieving generalization. In this context, the appearance of1037the background does not adversely affect our use case. We1038anticipate that advances in generative AI models will miti-1039gate these limitations.

1040 B. Evaluation Videos

In this paper, we developed a pick-and-place video datasetcontaining 288 videos spanning 9 tasks for evaluatingboth Slot-Net and SLeRP.

1044For Slot-Net evaluation, we utilize the *start* and *end* im-1045ages from each video. Among the 9 tasks, 3 are categorized1046as unseen for Slot-Net evaluation: "bottle in organizer",1047"cup on saucer", and "egg in egg steamer".

For SLeRP evaluation, we select two videos of the same 1048 task, one serving as the human demonstration and the other 1049 as the robot-view video. The robot-view video has its final 1050 frame depicting what the robot's view will look like upon 1051 1052 completion of the action, providing the ground-truth end frame appearance. The three different settings we use to 1053 pair the videos are: (1) different views (captured from dif-1054 ferent camera angles), (2) different backgrounds (captured 1055 in different environments with setups on distinct tables), 1056 1057 and (3) different slot occlusions (captured with varying slot 1058 occlusions, where one video has all other slots empty while the other has some slots filled). 1059

In this section, we provide visualizations of the 9 tasksand 3 settings in Tab. 7 and Tab. 8.

B.1. Visualization of 9 tasks

We show one video for each of the 9 tasks in Tab. 7. For 1063 each video, we out the (*start, pick, place, end*) frames to 1064 illustrate the action. 1065

B.2. Visualization of 3 settings

We show two example videos for each of the 3 settings (different views, different backgrounds, and different slot occlusions) in Tab. 8.106710681068

C. Additional Implementation Details

C.1. Slot-Net

We follow SAM's training recipe for finetuning the ViT-
based encoder using a combination of Dice Loss and Cross-
Entropy Loss. We train the model on an A100 GPU for 72K
iterations. We use a learning rate of 1×10^{-5} with a batch
size of 8.1072
10731074
10751074
1075

C.2. CLIPort++ 1077

We start with the released code from the authors. To pre-1078 pare labels for training from our RGBD videos, we treat 1079 our videos as a one-step action with the corresponding task 1080 name as the language goal. And we treat the center of the 1081 object mask as the pick action location and the center of the 1082 slot mask as the place action location. We use 80 videos for 1083 training and 64 videos for evaluation, ensuring there is no 1084 overlap in tasks, objects, and backgrounds. 1085

D. Additional Results of SLeRP

D.1. Full Ablation Study Table

Tab. 9 presents the full ablation study results evaluated all1088all the three different settings.1089

D.2. Multi-slots analysis.

Tab. 10 further provides 1-to-another slot placement results,1091placing into a different slot as human videos, and shows1092ours generalizes well.1093

	Diff Slot Placement												
	∣Obj↑	Slot↑	Prec.↑	CD↓	EMD↓								
Ours	67.70	42.62	23.14	0.0575	0.0350								

Table 10. Object placement into a different slot as in the human video.

To evaluate DINOv2+SAM 1-to-N slot re-identification,1094we annotated all N slots, and Table 11 shows that our1095method does reliably well in finding all equivalent slots.1096

	1	Diff	View	Dif	f Bg	Diff Slot Occ.			
		mIoU	AP	mIoU	AP		mIoU	AP	
Ours		68.3	46.3	66.8	44.2		66.0	43.3	

Table 11. Multi-slot mask identification evaluation.

1086

1087

1097 D.3. More baseline comparison results

1098 We show more qualitative comparison results for baselines1099 in Fig. 10.

1100 D.4. More ablation comparison results

We show more qualitative comparison results for ablationsin Fig. 11.

1103 D.5. Limitations

Our system, SLeRP, demonstrates superior performance 1104 compared to baseline methods. However, there are areas 1105 that require improvement. Firstly, SLeRP is modular, with 1106 each module relying on the output of the preceding one, 1107 which may result in compounding errors. If a preceding 1108 module fails to provide accurate outputs, subsequent results 1109 may be adversely affected. Secondly, while the modules 1110 currently employed are not without flaws, they can be eas-1111 ily replaced with more advanced methods. 1112

1113 For instance, Slot-Net necessitates that the start and end images have minimal changes aside from the pick object; 1114 future methodologies could relax this requirement, enabling 1115 slot detection in more general settings. Additionally, SAM2 1116 is a tracking method repurposed for re-identification, and 1117 methods specifically designed for re-identification would 1118 1119 likely yield better results. Lastly, while Mast3r is effective for object-level matching, it struggles with slot-level match-1120 ing, indicating that improved slot-level matching algorithms 1121 could enhance accuracy. 1122

1123 E. Additional Results of Slot-Net

1124 E.1. Additional results of Slot-Net

We provide more Slot-Net comparison results in Tab. 12 and random results in Tab. 13.

1127 E.2. Limitations

1128 Slot-Net has limitations in slot segmentation across varying conditions due to its training data constraints: (1) It requires 1129 that the input (start, end) images have minimal changes in 1130 viewpoint. (2) It necessitates that the pick object is the sole 1131 changing element between (start, end) images, and neither 1132 humans nor human hands are present. Slot-Net may exhibit 1133 the following issues: (1) While it can generalize to unseen 1134 tasks, its performance may degrade. (2) The predicted slot 1135 may encompass multiple slots or extend beyond the ground 1136 truth. (3) It may incorrectly identify pick objects as slots. 1137

Figure 10. **Qualitative Comparison for baselines.** We compare our method to baselines and present side-by-side results on six examples. For each example, the top row displays 2D slot prediction results, while the bottom row shows 6-DoF relative pose predictions by projecting the object point cloud onto the slots.

Figure 11. **Qualitative Comparison for ablations.** We compare our method to ablations and present side-by-side results on six examples. For each example, the top row displays 2D slot prediction results, while the bottom row shows 6-DoF relative pose predictions by projecting the object point cloud onto the slots. 17

#

Table 5. Data generation process (Part 1/2). Start from a task name (col 1) and an object-centric image (col 3), we first get an enriched detailed description using [69] (col 2), then remove one object from the slots (col 4) and annotate the slot mask on the object-centric image (col 5). To put the object-centric image into various daily backgrounds, we transform the object-centric image (col 3) onto a 1024x1024 canvas and outpaint the background with the enriched text (col 2) to create *End*. Finally, *Start* and GT mask follow the same transformation onto the canvas.

Task	Enrich text	Obj-centic img	Obj removal	Annotation	Start	End	GT mask
mug on coaster	'up', 'umg,' 'botteet', 'constar', 'cyle'. Tigja- quaity and photorealistic, 2020s, modern, beight, moming time, 43-degree. A ceranic may with a textured, matte finish, having a simple vel degant to it, and a wooden constar- with a natural, nuistic texture, having a simple second. 'An examine model's and a textured, and a wooden constar- ter daily supple second, such as a gar of coffee and fee daily supple second, such as a gar of coffee and leve coffee corps, with a wooden table finish and a few coffee corps, with a wooden table finish and a constant a super of coffee and the constant and the second wall'.			Ö	3	3	
peach in box	'top': 'peach', 'bottom'; 'boa', 'style', 'mid-century modern, high-guadry and peace and the style of the style cosy, morning, 45-degree bod-down vere', 'dogtet' 'The peach has a smooth, bod-down vere', 'botta has a soft hown color', 'back- urg occurre, 'the botta has a worden texture, with a nut- ar soft hown color', 'back- man to the style and chairs, and a few daily applies in data a to atostet, with fresh flowers.'			•			•
pen in orga- nizer	'top': 'pen', 'bottom': 'organizer, 'tsyle': 'high- modern Encogenesity, bright and cory, morning, 5 d-signer look-down view', 'dejects': The pen is to solve the second second second instrument with a silver finish and a rubbeized pin. The cognitor is a diver finish and a rubbeized pin. The cognitor is a rubbin base and a me hocket.' 'background is a rubbin base and a me hocket.' 'background is a rubbin office, with a wooden desk and a controllable office chair. There are secret organizer, including a notebook, a stapler, and a cop of coffice:						
pepper in tray	'top': 'pepper', 'bottom'; 'tray', 'kyjk': 'kgjk-gaality contemporary, bright, lumi- noss, morting, 4S-degree lood-down view', 'dojecti' nure, placed in the center for a rectangular, stanless steel arry with a smooth' a reveal kut kom labe with a few daily supplies, a warm, 'a vessel kut kut hen wall, cabinets, and a window that the six in mannal light'						
					219 (100) (10) (10) (10) (10) (10) (10) (10) (10)		•

#

No.	Task	#Total	Start Frame	Pick Frame	Place Frame	End Frame
1	bread in toaster	32				
2	bottle in or-	32				
3	fruit in tray	32	Ø	0 3 2		
4	cup in saucer	32				
5	mug in coaster	32				
6	egg in egg steamer	32				Contraction of the second seco
7	fruit in or-	32				
8	utensil in utensil organizer	32				
9	muffin in muffip trav	32				633

Table 7. Evaluation videos for 9 tasks. We show one example video for each of the 9 tasks from our evaluation video dataset. The 4 frames showing for each video are the start frame, pick frame, place frame, and end frame in the video. We get the pick and place frame indexes from the in-contact object trajectory, where the pick frame index is when the in-contact object is first detected, and the place frame index is when the in-contact object is last detected in the video.

#

Settings	Video 1 - Start Frame	Video 1 - End Frame	Video 2 - Start Frame	Video 2 - End Frame
Different view		a.y	C. Q.	
Different Background		đ.Q		
Different Slot Occlustion		đ.Q		6.9
Settings	Video 1 - Start Frame	Video 1 - End Frame	Video 2 - Start Frame	Video 2 - End Frame
Different view			CONTRACT OF	Contraction of the second seco
Different Background				
	1 620		655	655

Table 8. Evaluation videos in 3 settings. We show two example videos ("bread in toaster", "muffin in muffin tray") in 3 different setting. For each task of our evaluation videos, we captured videos for 3 different settings to evaluate system performance when the human demonstration video and its paired robot-view video were in different conditions. The 3 different settings we captured are (1) different views (two videos are captured from different camera views), (2) different backgrounds (two videos are captured from different backgrounds where we set up the two scenes on two different tables), and (3) different slot occlusions (two videos are captured from different slot occlusions, where one video with all other slots empty while the other video with some slots full).

Mathad	ClatNat	SAM2	Mast3r		Diff. View					Diff. Background					Diff. Slot Occu.				
Wiethou	Slouver	SAM2		Obj ↑	Slot \uparrow	Prec. \uparrow	$\mathbf{CD}\downarrow$	$\text{EMD} \downarrow$	Obj ↑	$\operatorname{Slot}\uparrow$	Prec. \uparrow	$\mathbf{CD}\downarrow$	$\text{EMD} \downarrow$	Obj ↑	Slot \uparrow	Prec. \uparrow	$\mathbf{CD}\downarrow$	$EMD \downarrow$	
Base design	X	×	X	28.53	28.23	24.72	0.2289	0.1183	33.27	23.40	22.22	0.1217	0.0643	27.71	27.05	26.47	0.1499	0.0809	
Ours w/o SlotNet	×	1	1	73.85	29.63	27.84	0.0486	0.0289	70.27	23.91	18.80	0.0630	0.0367	68.12	26.33	25.06	0.0520	0.0315	
Ours w/o SAM2	1	×	1	31.05	35.98	23.75	0.1219	0.0667	33.98	30.87	16.99	0.0870	0.0505	28.62	37.48	22.97	0.1012	0.0575	
Ours w/o Mast3r	1	1	×	73.85	54.37	32.74	0.0321	0.0205	70.27	44.70	28.93	0.0540	0.0318	68.12	47.04	27.38	0.0675	0.0374	
Ours	1	1	1	73.85	54.37	36.40	0.0282	0.0182	70.27	44.70	25.39	0.0573	0.0323	68.12	47.04	30.30	0.0334	0.0223	

Table 9. System Ablation Studies. We conduct ablation studies on the key modules in SLeRP. All metrics and test splits follow Table 1. Results show all the key modules help, and our full system performs the best on average across different settings.

Table 12. More comparison results for Slot-Net on seen tasks. We show more examples for comparison results between Slot-Net and other baselines or heuristics.

Table 13. Random results for Slot-Net on seen tasks. We show randomly sampled results from Slot-Net on seen tasks.