

483 Supplementary Material

484 Table of Contents

486	A Implementation Details	14
487	A.1 Task Construction	14
488	A.2 Retrieval Policy Training	14
489	A.3 Spatial-aware Representation Training	15
490	A.4 Sim-to-Real transfer	16
491	A.5 Baseline Implementation	17
492	A.6 Evaluation metric	19
493	A.7 Computing Resources	20
494	B Additional Experiments	20
495	B.1 Impact of Occlusion Rate	20

499 **A Implementation Details**

500 **A.1 Task Construction**

501 **Domain Randomization.** To enhance the robustness and generalization capability of our system, we
502 implement comprehensive domain randomization strategies during the environment reset phase. The
503 randomization encompasses multiple aspects of the environment:

- 504 • Object Mass Randomization: At the beginning of each episode, object masses are randomized
505 by scaling each object’s default mass with a random factor sampled from a uniform distribution
506 $U(1, 1.5)$ (units: kg):

$$m_{\text{curr}} = m_{\text{default}} \cdot \alpha, \quad \alpha \sim U(1, 1.5)$$

- 507 • Object Position Randomization: Small perturbations are applied to the initial positions of objects
508 to introduce variability (units: meters):

$$\begin{aligned} \Delta x &\sim U(-0.02, 0.02) \\ \Delta y &\sim U(-0.02, 0.02) \end{aligned}$$

- 509 • Target Position Randomization: The initial position of the target object within the box is randomized.
510 The random displacements are sampled from (units: meters):

$$\begin{aligned} \Delta x &\sim U(-0.15, 0.15) \\ \Delta y &\sim U(-0.2, 0.2) \end{aligned}$$

511 This ensures that the target object can be placed within 70% of the box’s area.

- 512 • Camera Mount Randomization: During data collection, the camera’s mounting position is perturbed
513 with small random displacements (units: meters):

$$\mathbf{p}_{\text{camera}} = \mathbf{p}_{\text{default}} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim U(-0.01, 0.01)^3$$

514 **A.2 Retrieval Policy Training**

515 **Reward Design.** We carefully design the reward for our system, which contain serval item including
516 distance reward, stir reward, proximity clearance reward, pixel emergency reward and penalty. We
517 show the specific parameters setup for each reward item in Table 4.

Table 4: Reward Components and Definitions.

Reward Item	Definition	Parameters
Distance Reward	$r_{\text{dist}} = \exp(-5 \cdot \min(d - e_0, 0))$	$e_0 = 0.15$
Stir Reward	$r_{\text{stir}} = \alpha \ p_t^{\text{all}} - p_{t-1}^{\text{all}}\ _2$	$\alpha = 500.0$
Proximity Clearance Reward	$r_{\text{clean}} = \beta \cdot r_{\text{dist}} \cdot \sum_{i=1}^k f_i$	$\beta = 100.0$
Pixel Emergency Reward	$r_{\text{pixel}} = C/15$	-
Penalty	$\lambda_1 \ a_t^{\text{hand}} - a_{t+1}^{\text{hand}}\ _2^2 + \lambda_2 \ a_t^{\text{arm}}\ _2^2$	$\lambda_1 = 1.0, \lambda_2 = 1.0$

518 **Policy Training.** We employ the Proximal Policy Optimization (PPO) algorithm [28] to train a
519 continuous control policy using an actor-critic architecture. Detailed hyperparameters are provided
520 in Table 5. The policy network is parameterized as a multi-layer perceptron (MLP) with three
521 layers of sizes [256, 256, 256], utilizing the ELU activation function for improved gradient flow and
522 non-linearity. The standard deviation of the policy distribution is learned via a log-std representation,
523 enabling dynamic adjustment of exploration during training.

Table 5: Hyperparameters for PPO Training

Category	Parameter	Value	Description
<i>Model Architecture</i>	MLP Layers	[256, 256, 256]	Number of neurons per layer
	Activation Function	ELU	Non-linearity used in the network
<i>Training Parameters</i>	Learning Rate	3×10^{-4}	Step size for policy update
	Discount Factor (γ)	0.99	Reward discounting factor
	GAE Parameter (τ)	0.95	Smoothing factor for GAE
	Entropy Coefficient	0	Weight of entropy regularization
	Gradient Clipping	Norm 1, Truncation Enabled	Prevents gradient explosion
	Clip Range (ϵ)	0.1	PPO clipping threshold
	KL Threshold	0.02	KL divergence threshold for stopping training
	Minibatch Size	512	Batch size for optimization
	Mini Epochs	5	Number of updates per batch
	Horizon Length	8	Number of steps before update
	Max Training Epochs	50,000	Maximum number of training iterations
	Value Learning Rate	1×10^{-3}	Learning rate for value function

524 To ensure stable and efficient learning, we adopt adaptive learning rate scheduling, starting at
 525 3×10^{-4} . The advantage function is normalized to reduce variance in policy gradient updates, while
 526 the Generalized Advantage Estimation (GAE) [27] parameter τ is set to 0.95, striking a balance
 527 between bias and variance. Gradient clipping with a norm threshold of 1 is applied to prevent
 528 exploding gradients. To constrain policy updates, we employ a PPO clipping range of 0.1, which
 529 limits large deviations from the current policy, and enforce a KL divergence threshold of 0.02 to
 530 promote conservative updates and prevent policy collapse.

531 The training runs for a maximum of 50,000 epochs, with model checkpoints saved every 1,000
 532 epochs. The best-performing model is selected based on validation returns and retained after 200
 533 epochs to prevent overfitting. A separate centralized value function is used for advantage estimation,
 534 parameterized as an MLP with the same architecture as the policy network. The critic network
 535 employs a higher learning rate of 1×10^{-3} to facilitate faster convergence in value estimation, a choice
 536 informed by preliminary experiments indicating more stable critic updates with this configuration.

537 A.3 Spatial-aware Representation Training

538 **Model Architecture.** We design a hierarchical GNN model for estimating visibility and extracting
 539 structured representations of target objects in cluttered scenes. The model integrates a PointNet++-
 540 based encoder for each object category (target, hand, obstacles) and performs message passing via
 541 EdgeConv [31] layers. Specifically, three main components constitute the model: (1) **Point Feature**
 542 **Encoding:** Each point cloud (target, hand, obstacles) is encoded via a PointNet++ module [26],
 543 consisting of three Set Abstraction (SA) layers and two Feature Propagation (FP) layers. Global max
 544 pooling over the final layer provides a fixed-length feature representation. (2) **Graph Construction:**
 545 For each sample, a complete graph is built among all encoded nodes, including one target node,
 546 one hand node, and $K = 9$ nearest neighbor obstacles. (3) **Graph Message Passing:** Two layers
 547 of EdgeConv are applied to propagate features among nodes, followed by a third EdgeConv that
 548 aggregates pairwise edge features for visibility prediction.

549 **Implementation Details.** Before training this module, we collect 2000 trajectories generated by a
 550 basic reinforcement learning policy. The policy is trained using reward signals. We optimize the
 551 model using mean squared error (MSE) loss, training the network to predict the visibility label,
 552 defined as the number of pixels occupied by the target object in the top-view camera image. The
 553 dataset is randomly split into 90% for training and 10% for validation. We employ the Adam optimizer
 554 and apply a learning rate scheduler that reduces the learning rate upon plateauing of the validation
 555 loss. To monitor early-stage performance closely, validation is performed every 29 training iterations.

556 To support efficient parallel data loading and training, we design a function to batch dynamic obstacle
 557 sets and filter invalid data samples during runtime. Each episode in the dataset contains 209 time
 558 steps, and the model is trained either on the last step or all time steps, depending on the configuration.
 559 Table 6 summarizes all hyperparameters used for GNN training. The output dimensionality of the
 560 extracted structured feature (used in downstream RL) is set to $d = 32$ unless otherwise specified.

Table 6: Hyperparameters for GNN Training.

Category	Parameter	Value	Description
Model Architecture	PointNet++ Encoder Output Dim	64	Dimension of point cloud encoding
	GNN Hidden Dim	64	Hidden units in EdgeConv layers
	Final Output Dim	32	Output dim of GNN feature extractor
	Num of Neighbors (K)	9	Number of obstacle nodes in graph
PointNet++ Settings	SA Layer Config	[512, 128, 1]	Points per SA layer
	SA Radii	[0.2, 0.4, all]	Search radii for each SA layer
	MLPs per SA Layer	[64,64,128], [128,128,256], [256,512,out]	MLPs applied in SA layers
	Graph Convs	2 EdgeConv + 1 predictor	Number and type of GNN layers
Training Parameters	Optimizer	Adam	Optimization algorithm
	Learning Rate	1×10^{-4}	Initial step size
	LR Scheduler	ReduceLROnPlateau	Decay LR on validation plateau
	Scheduler Settings	factor 0.5, patience 5	Scheduler configuration
	Loss Function	$0.5 \times \text{MSE} + 0.5 \times \text{Huber}$	Objective function
	Batch Size	32	Samples per batch
Validation Settings	Epochs	30	Max training epochs
	Validation Interval	Every 29 steps	Frequency of validation
	Data Split	90% train / 10% val	Dataset partition ratio
	Time Sequence Mode	last (default)	Sampling strategy from sequence

561 A.4 Sim-to-Real transfer

562 For real-world deployment, we collect a set of trajectories generated by our RL expert policy.
 563 To ensure data quality and consistency, we first select successful trajectories. We then filter out
 564 trajectories where the finger’s z -coordinate lower 2 cm above the box, a threshold empirically chosen
 565 to prevent unstable behavior and reduce the risk of collision with the box during manipulation. To
 566 promote generalization, we balance the dataset across various target object positions within the box,
 567 ensuring uniform coverage of spatial configurations. This prevents the model from overfitting to
 568 specific object placements and enhances its adaptability to unseen scenarios.

569 Figure 5 shows the model architecture, which consists of a state
 570 encoder followed by a multi-head self-attention mechanism with six transformer layers, each containing six attention heads.
 571 This design captures complex temporal dependencies across
 572 historical state sequences of length 15, enabling accurate pre-
 573 diction of future actions over a 5 step horizon chunking. The
 574 hidden dimension of 384, paired with a feed-forward expansion
 575 ratio of 5.33 (2048/384), strikes a balance between model
 576 expressiveness and computational efficiency.

578 To effectively manage the continuous action space inherent
 579 in robotic control, we introduce a custom Negative Log Prod-
 580 uct Loss function, which penalizes trajectory deviations more
 581 sensitively than traditional mean squared error. This loss func-
 582 tion emphasizes multi-step consistency, enhancing the model’s
 583 predictive stability. Training is performed using the Adam opti-
 584 mizer with a learning rate of 1×10^{-4} over 10,000 iterations
 585 and a batch size of 512. Mixed-precision training accelerates
 586 computation without compromising accuracy, while gradient
 587 clipping at 1.0 maintains stable learning dynamics. Hyperpa-
 588 rameter selection was guided by cross-validation on a held-out
 589 dataset to optimize both performance and robustness. Detailed
 590 architectural specifications and hyperparameters are provided in Table 7. Despite strong simulation
 591 performance, real-world deployment introduces challenges such as sensor noise, domain discrep-
 592 acies, and dynamic environmental conditions. Our transformer model mitigates these issues by
 593 leveraging temporal patterns to predict smooth and consistent actions.

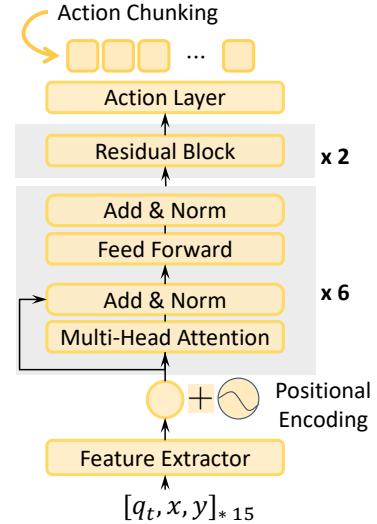


Figure 5: The architecture of student policy.

Table 7: Hyperparameters of Distilled Policy.

Category	Parameter	Value	Description
<i>Model Architecture</i>	Input State Dimension	9	Size of input state vector
	Action Dimension	13	Number of output actions
	History Frames	15	Past frames used as input
	Future Action Frames	5	Future actions predicted
	Transformer Hidden Size (d_{model})	384	Hidden layer size
	Number of Attention Heads	6	Transformer attention heads
	Number of Transformer Layers	6	Transformer depth
	Feed-forward Dimension	2048	FFN hidden size
	Dropout Rate	0.15	Dropout probability
<i>Training Parameters</i>	Batch Size	512	Training batch size
	Total Iterations	10,000	Training iterations
	Learning Rate	1e-4	Initial learning rate
	Optimizer	Adam	Optimization algorithm
	Loss Function	Negative Log Product Loss	Loss function used
	Gradient Clip Norm	1.0	Gradient clipping threshold

594 A.5 Baseline Implementation

595 In our simulation experiments, we compare our method against several baseline approaches. *Ours*
 596 (*SAC*) is a variant of our system where the PPO algorithm is replaced with Soft Actor-Critic (SAC)
 597 to evaluate the impact of different RL optimization strategies. *Ours (Gripper)* is a variant of our
 598 system where the dexterous hand is replaced with a UMI gripper to assess the importance of hand
 599 morphology in cluttered object retrieval. The other two baselines are *Visual-based Motion Planning*
 600 *Search (VMP)* and *Grasp-Pick*. VMP is a heuristic motion planning approach that uses target object
 601 segmentation masks to guide the robotic hand toward the target object and employs predefined rules
 602 for retrieval manipulation. Grasp-Pick involves sequentially grasping and placing objects based on
 603 the support relationships within the cluttered scene.

604 A.5.1 Ours (SAC)

605 SAC is an off-policy reinforcement learning algorithm grounded in the maximum entropy framework,
 606 which promotes exploration by encouraging stochastic action selection. The algorithm comprises a
 607 soft Q-function $Q_\theta(s, a)$ and a stochastic policy $\pi_\phi(a|s)$.

608 The soft Q-function is defined as the expected cumulative return:

$$Q_\theta(s_t, a_t) = \mathbb{E} \left[\sum_{t'=t}^T \gamma^{t'-t} r_{t'} \mid s_t, a_t \right], \quad (5)$$

609 where $\gamma \in [0, 1]$ denotes the discount factor. In maximum entropy RL, this return is typically
 610 augmented with an entropy bonus, though omitted here for clarity.

611 The Q-function parameters θ are optimized by minimizing the soft Bellman residual:

$$J_Q(\theta) = \mathbb{E}_{(s_t, a_t, s_{t+1}, r_t) \sim \mathcal{B}} \left[(Q_\theta(s_t, a_t) - r_t - \gamma \bar{V}(s_{t+1}))^2 \right], \quad (6)$$

612 where $\bar{V}(s_{t+1}) = \mathbb{E}_{a \sim \pi_\phi} [Q_{\bar{\theta}}(s_{t+1}, a) - \alpha \log \pi_\phi(a|s_{t+1})]$ uses a target Q-network with parameters
 613 $\bar{\theta}$. The temperature α controls the strength of the entropy regularization and can be fixed or learned
 614 during training.

615 The policy parameters ϕ are updated by minimizing:

$$J_\pi(\phi) = \mathbb{E}_{s_t \sim \mathcal{B}, a_t \sim \pi_\phi} [\alpha \log \pi_\phi(a_t|s_t) - Q_\theta(s_t, a_t)]. \quad (7)$$

616 SAC alternates between policy evaluation and improvement, enabling stable and sample-efficient
 617 learning. In this work, we adopt SAC as our base reinforcement learning algorithm due to its
 618 robustness in continuous control tasks.

Table 8: Hyperparameters for SAC Training

Category	Parameter	Value	Description
<i>Network Architecture</i>	Actor Hidden Dim	256	Hidden layer dimension of actor network
	Actor Hidden Depth	3	Number of hidden layers in actor network
	Critic Hidden Dim	256	Hidden layer dimension of critic network
	Critic Hidden Depth	3	Number of hidden layers in critic network
<i>Training</i>	Actor Learning Rate	5e-4	Learning rate for actor network
	Critic Learning Rate	5e-4	Learning rate for critic network
	Temperature Learning Rate	0.05	Learning rate for entropy coefficient
	Batch Size	8192	Number of samples per training batch
	Max Gradient Norm	0.5	Maximum gradient norm for clipping
	SAC Epochs	8	Number of training epochs per update
	Discount Factor	0.99	Future reward discount factor
	Initial Temperature	0.1	Initial entropy coefficient
	Critic Tau	0.05	Soft update coefficient for target critic
	Learnable Temperature	True	Whether to learn entropy coefficient
	Number of Environments	512	Parallel environments in IsaacGym
	Update Interval	1	Environment steps between updates
	Seed Steps	32	Initial random steps per environment
	TD n-step	3	Number of steps for TD learning
	Replay Buffer Size	5M	Maximum capacity of experience replay

619 We implement the baseline *Ours* (*SAC*), incorporating several key modifications to improve training
 620 efficiency and stability. The implementation is built on the IsaacGym physics engine, enabling
 621 large-scale parallelism with 512 simulated environments to accelerate data collection. The full set of
 622 hyperparameters is listed in Table 8. Both the actor and critic are modeled as three-layer multilayer
 623 perceptrons, each with 256 hidden units per layer. The actor outputs the mean and log standard
 624 deviation of a Gaussian action distribution, with the log standard deviation constrained between -5
 625 and 2 to ensure numerical stability. The critic adopts a double Q-learning architecture with target
 626 networks updated via a soft update rule using $\tau = 0.05$.

627 To stabilize training, we use a large batch size of 8192 and apply policy updates at every timestep.
 628 At the start of training, each environment performs 32 random actions to populate the replay buffer.
 629 The actor and critic are optimized using the AdamW optimizer with a learning rate of 5e-4, while the
 630 temperature parameter is updated with a learning rate of 0.05. Gradient clipping with a threshold of
 631 0.5 is applied to mitigate exploding gradients. The replay buffer stores up to 5 million transitions.
 632 For each policy update, we perform 8 epochs of optimization using a single mini-batch per epoch,
 633 promoting efficient data usage while reducing the risk of overfitting. The temperature parameter is
 634 initialized to 0.1 and is learned adaptively to maintain a balance between exploration and exploitation.

635 Our parallelized training framework significantly enhances sample throughput by synchronously
 636 simulating 512 environments, all of which share a single policy network. This design ensures
 637 consistent policy updates while leveraging high-throughput simulation. Additionally, we adopt 3-step
 638 temporal difference (TD) learning to reduce bias in value estimation and improve overall policy
 639 performance.

640 A.5.2 VMP

641 The VMP system implements a vision-guided manipulation framework for dexterous robotic retrieval
 642 tasks in cluttered environments. It integrates visual perception, motion planning, and control execution
 643 through a state machine architecture to ensure reliable object manipulation.

644 The vision module employs a top-down camera with a resolution of 1024×512 , capturing RGB,
 645 depth, and segmentation maps of the workspace. Target objects are identified using segmentation
 646 masks obtained from the segmentation map, with their IDs corresponding to known object labels. The
 647 center of the target mask's bounding box is extracted as the 2D image coordinate, which is projected
 648 into 3D space using depth data to obtain precise object localization. For motion planning, the robotic
 649 arm moves its end effector to the computed 3D coordinate and performs a scrape action to retrieve
 650 the object.

651 When the target object is completely occluded and its segmentation mask cannot be detected, the
 652 system employs an exploration strategy by randomly sampling four 3D coordinates within the
 653 cluttered bin area. The arm sequentially moves to these coordinates, performing scrape actions to
 654 uncover the target object. Specifically, the entire motion planning and scrape action process employs
 655 a four-stage approach to ensure reliable object retrieval.

656 **Pre-approach stage:** The end-effector moves to a predefined position ($h = 0.5$ m) above the target
 657 object. This configuration facilitates subsequent control of the hand to reach any position within the
 658 bounding box.

659 **Final approach stage:** Precise positioning is achieved using visual feedback combined with damped
 660 least squares inverse kinematics:

$$\tau = J^T(JJ^T + \lambda I)^{-1}\Delta x$$

661 where $\lambda = 0.05$ is the damping parameter, J is the Jacobian matrix, and Δx represents the positional
 662 error.

663 **Scraping stage:** The system executes a periodic motion pattern defined by:

$$x(t) = A \sin(2\pi ft + \phi) + O$$

664 where the amplitude $A = 2.0$, frequency $f = 20$ Hz, phase shift $\phi = \pi/4$, and offset $O = 0.5$.

665 **Reset stage:** The target object has been retrieved, so the robot arm will return the end effector to its
 666 initial position.

667 The control execution module utilizes position-based control for both arm and finger joints. Adaptive
 668 damping parameters are applied to ensure stable motion, while joint limits are strictly enforced
 669 throughout the execution process: $q_{\min} \leq q \leq q_{\max}$.

670 A.5.3 Grasp-Pick

671 This method relies on support relationships among cluttered objects to guide the grasping sequence.
 672 To ensure these assumptions hold, we designed tailored setups for both simulation and real-world
 673 experiments.

674 In simulation, object positions are directly accessible, enabling precise calculation of support rela-
 675 tionships. We employ a KD-Tree [29] to organize the coordinates of objects near the target. Based
 676 on Euclidean distance, we select the three to five nearest objects, depending on the scenario, and
 677 manipulate them sequentially to clear access to the target. While this approach offers computational
 678 simplicity, it assumes ideal sensing conditions and may not generalize to more complex spatial
 679 arrangements.

680 For robotic control, we implement damped least squares inverse kinematics:

$$\dot{\mathbf{q}} = \mathbf{J}^T(\mathbf{JJ}^T + \lambda \mathbf{I})^{-1}\dot{\mathbf{x}},$$

681 where λ is the damping coefficient, \mathbf{J} is the Jacobian matrix, and $\dot{\mathbf{x}}$ is the desired end-effector velocity.
 682 This formulation offers stable solutions near singularities but may limit the dexterity needed in
 683 cluttered environments.

684 In real-world experiments, sequentially grasping and removing multiple objects in stacked scenes
 685 with a dexterous hand remains challenging due to perception and control limitations. To address this,
 686 we employ predefined trajectories for each trial, simulating an idealized execution scenario. While
 687 this implementation provides an upper-bound estimate of this method’s efficiency, it does not reflect
 688 the challenges of autonomous execution in unstructured environments.

689 A.6 Evaluation metric

690 **Exposure Calculation.** The primary goal of object retrieval is to locate the target object and
 691 enhance its visibility within the camera’s field of view, facilitating subsequent manipulations. We
 692 define *exposure* as the proportion of unobstructed pixels of the target object in the imaging plane.
 693 Considering that changes in the object’s pose can affect the number of visible pixels, we proceed as
 694 follows:

695 At timestep t , we record the target object’s visible pixels p_t^{curr} and its 6D pose. Subsequently, all
 696 objects except the target are removed, the target object’s recorded 6D pose is reset, and its total visible

697 pixels are recorded as p_t^{all} . The exposure at time t is then computed as:

$$\text{exposure}_t = \frac{p_t^{\text{curr}}}{p_t^{\text{all}}}. \quad (8)$$



Figure 6: Examples of successful and failed object retrievals on the real robot.

698 **Success in Real-World Experiments.** To systematically evaluate the success rate of object retrieval
 699 on the real robot, we capture images before and after each task using a side-mounted RealSense D435
 700 camera. The success of the retrieval is determined by comparing the exposure of the target object in
 701 these images. As illustrated in Figure 6, we present examples of both successful and failed retrieval
 702 attempts.

703 A.7 Computing Resources

704 All experiments are conducted on a single
 705 NVIDIA RTX 4090 GPU. The reinforcement
 706 learning training takes approximately 6 hours,
 707 the GNN module requires around 20 minutes,
 708 and the student policy training completes within
 709 4 minutes.

710 B Additional Experiments

711 B.1 Impact of Occlusion Rate

712 We also explore the impact of clutter occlu-
 713 sion rate (i.e., $1 - \text{exposure}$) on target objects.
 714 Figure 7 presents the relationship between oc-
 715 clusion rates, retrieval success rate (RSR), and
 716 retrieval steps (RS) in cluttered environments,
 717 comparing smaller and larger target objects. As
 718 the occlusion rate rises from 30% to 100%, RSR decreases for both object sizes, with smaller objects
 719 consistently achieving higher RSR than larger ones. This indicates that larger objects, although

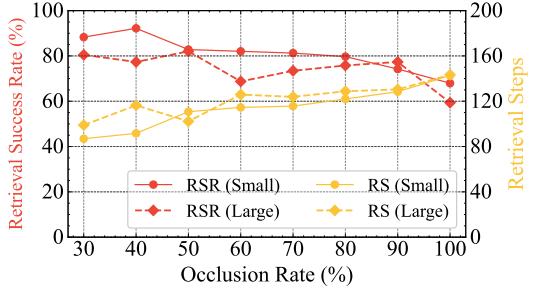


Figure 7: Impact of Occlusion Rate on Performance and Efficiency. We evaluate the retrieval success rate and retrieval steps of our policy for small and large target objects under varying occlusion levels.

720 generally easier to detect, are more prone to being significantly obstructed by substantial clutter.
721 In contrast, retrieval steps increase as occlusion rates rise, signifying reduced efficiency in highly
722 cluttered settings. Smaller objects generally require fewer retrieval steps than larger ones across most
723 occlusion levels, likely because retrieving larger objects involves removing more obstructing clutter.
724 These findings reveal that smaller objects are retrieved with higher success and efficiency, while
725 larger objects face greater challenges posed by occlusion.