

Hierarchical and Modular Network on Non-prehensile Manipulation in General Environments

Yoonyoung Cho*, Junhyek Han*, Jisu Han, Beomjoon Kim
 Korea Advanced Institute of Science and Technology

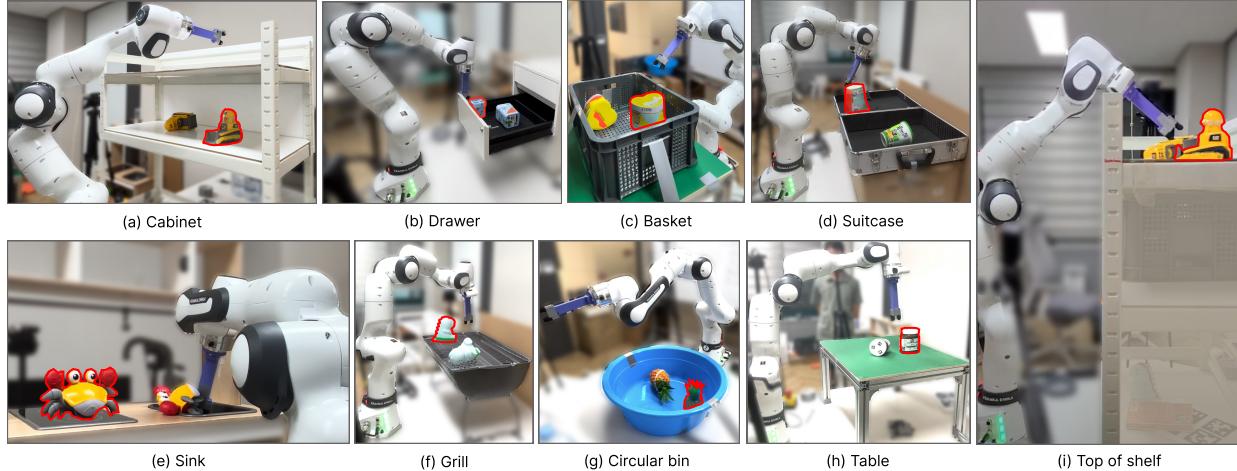


Fig. 1: Illustration of our real-world domains. Our approach generalizes to diverse and unseen objects and environments. The object outlined in red represents the goal pose. Note how the environment geometry limits object and robot motions: in (a), the robot must reach and maneuver the object while avoiding the ceiling; in (b) and (c), the robot must circumvent walls to access the object; in (d), the robot must move the object over a barrier; in (e), the robot must first traverse the sink to contact under the crab, lift it across the sink wall, then reorient the crab to face forward; and in (i), the robot must consider its kinematics while manipulating the object on a high shelf. Each object and its pose are chosen so that it cannot be simply pick-and-placed.

Abstract—For robots to operate in general environments like households, they must be able to perform non-prehensile manipulation actions such as toppling and rolling to manipulate ungraspable objects. However, prior works on non-prehensile manipulation cannot yet generalize across environments with diverse geometries. The main challenge lies in adapting to varying environmental constraints: within a cabinet, the robot must avoid walls and ceilings; to lift objects to the top of a step, the robot must account for the step’s pose and extent. While deep reinforcement learning (RL) has demonstrated impressive success in non-prehensile manipulation, accounting for such variability presents a challenge for the generalist policy, as it must learn diverse strategies for each new combination of constraints. To address this, we propose a modular architecture that uses different combinations of reusable modules based on task requirements. To capture the geometric variability in environments, we extend the contact-based object representation from CORN [17] to environment geometries, and propose a procedural algorithm for generating diverse environments to train our agent. Taken together, the resulting policy can zero-shot transfer to novel real-world environments despite training entirely within a simulator. We additionally release a simulation-based benchmark featuring nine digital twins of real-world scenes with 353 objects to facilitate non-prehensile manipulation research in realistic domains. Code, videos, and simulation benchmarks are available on the project website.

I. INTRODUCTION

Despite recent advances in robot manipulation, the practical deployment of robots in everyday environments like households remains challenging. One key reason is the robot’s inability to manipulate ungraspable objects. While much of prior work on manipulation centers around *prehensile* manipulation [57, 107, 27, 104], such approaches fall short in unstructured environments where objects are often ungraspable due to their geometry and the surrounding scene. To overcome this, robots must embrace non-prehensile manipulation, such as pushing, toppling, and rolling [44, 53, 17].

Recently, reinforcement learning (RL)-based approaches have achieved several successes in non-prehensile manipulation [45, 122, 17, 109]. However, these works have been limited to fixed objects in fixed scenes [45], general objects on flat tabletops [122, 17], or minor variations in objects and environments [109]. As such, no prior work addresses non-prehensile manipulation for novel objects *and* environments with arbitrary geometries as in Figure 1. Based on this

*Equal Contribution.

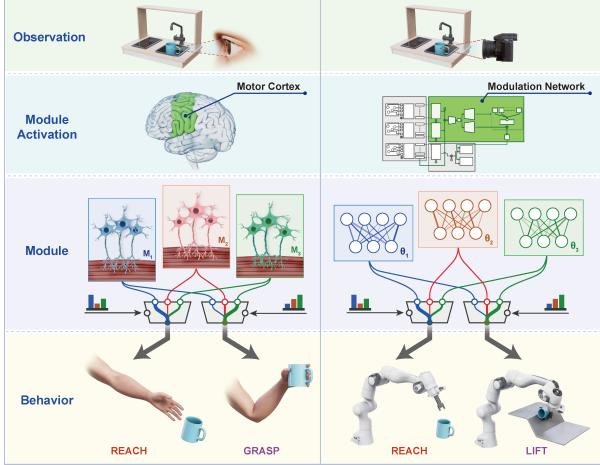


Fig. 2: Illustration of computational structure for biological motor control (left) and our architecture (right). Each row compares analogous components in (1) acquiring sensory observations, (2) determining module activations, (3) modules representing groups of co-activated neurons, and (4) composing the modules to construct specific behaviors. The bar graphs denote the activation weight of each module.

observation, our objective is to extend RL to enable non-prehensile manipulation in such diverse setups.

The key challenge here lies in training a policy that can adapt to the constraints imposed by the given scene. For instance, consider the scenarios in Figure 1: for each domain, the robot is presented with a unique set of constraints, which may also evolve during an episode (Figure 1e). This requires the policy to not only model multiple distinct behaviors, but also rapidly switch between behaviors in response to minor changes in state. For example, in the scenario depicted in Figure 1e, as soon as the toy crab is positioned above the sink, the robot must quickly switch from a *lifting* skill to a *translation* skill. However, standard networks struggle to learn such high-frequency functions, a phenomenon known as *spectral bias* [99, 4, 88].

Human brains, on the other hand, are extremely adaptive, and their computational structure for actions differs significantly from that of standard artificial neural networks. The computation of our brain is organized modularly [18], where a motor cortex orchestrates neural activities at the level of *motor modules*, a group of co-activating motor neurons¹ [103, 72], to produce an action based on the current task [73, 101, 43]. By invoking different sets of modules in response to the current context, the motor system produces disjoint behaviors such as reaching or grasping without interference [25]. This is illustrated in Figure 2, left.

Inspired by this, we propose a modular and hierarchical policy architecture (Figure 2, right). In our architecture, the *modulation network* assumes the role of the motor cortex which determines the activations of modules, each representing a group of co-activated network parameters. The weighted

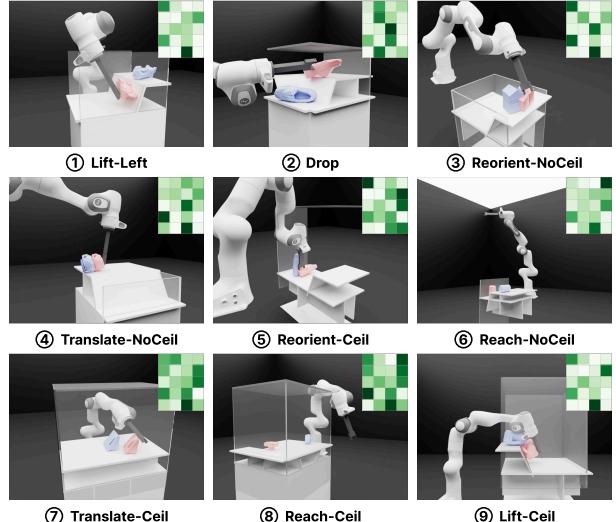


Fig. 3: Illustration of how distinct scenes map to distinct module activations that yield distinct behaviors, for a 5-layer network with 4 modules. The top-right colormap shows the activation of a particular module (column) for a particular layer (row): e.g., column 2, row 1 shows the activation of module 2 for layer 1. Opacity indicates the strength of module activation. The red object denotes the current object pose, and the blue object indicates the goal object pose. Ceil and NoCeil indicate the existence of a ceiling.

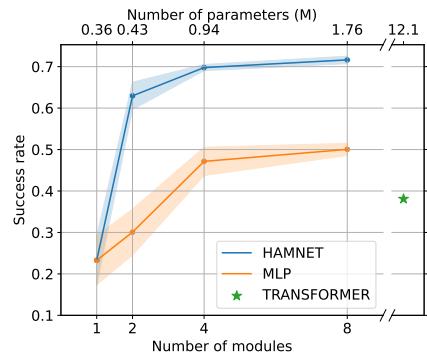


Fig. 4: Success rates per architecture by parameter counts for a monolithic architecture (MLP and transformer) (x-axis, top) and number of modules for HAMNET (x-axis, bottom)

combinations of these modules are then used to define the parameters of a *base network*, which has a fixed architecture but the parameters are determined online by the modulation network based on the current context, such as the environmental constraints. Unlike standard neural networks, this computational structure enables a single base network to model multiple functions, allowing the policy to produce distinct behaviors in response to even subtle changes in context.

We refer to our modular policy architecture as Hierarchical and Modular Network (HAMNET). We discovered that, when we train the policy with HAMNET, it autonomously discovers modules and their activations that correspond to distinct manipulation strategies, such as reaching, lifting, or reorienting objects (Figure 3). Furthermore, as shown in Figure 4, we

¹Neurons connected to muscle fibers that trigger muscular contractions.

found that HAMNET scales much more gracefully with the number of parameters compared to a standard neural network in our simulated domains.

Another challenge in generalizing across diverse environments and objects is acquiring geometric representations from high-dimensional point cloud observations. One option is to use end-to-end training, yet jointly learning to encode point clouds incurs significant memory footprint, limiting its use with large-scale GPU-based simulators [58], which have become an essential tool for training RL policies for robotics. Pre-training representation models can mitigate this issue, but their effectiveness depends on the pretext task capturing features relevant to manipulation. Consequently, large off-the-shelf models are often unsuitable, not only from their high computational cost but also due to spurious geometric features that are irrelevant for contact-rich manipulation [119, 17].

Instead, we build on CORN [17], which learns object representations tailored for manipulation based on a pretext task that predicts the presence and location of contact between a gripper and object. The key insight is for contact-rich manipulation, it is important to capture what forces and torques can be applied to an object in the given state, which in turn depends on the presence and location of the contact between the object and the gripper. While CORN restricts itself to gripper-object contacts, we also need object-environment contacts in our problems as we generalize over environments. So, we introduce Universal CORN (UNICORN) that generalizes to contact affordances between two arbitrary geometries, based on a Siamese pre-training pipeline where a single encoder learns the representations for both an object and environment, and a decoder predicts the contacts between each environment point cloud patch and the object.

Our last contribution concerns domain design. To obtain a policy that generalizes to diverse environments, training environments must encompass a range of geometric constraints, while affording fast simulation for practical training. This is a non-trivial problem: unlike objects, environment assets for simulated RL training are not widely available, and manually designing them is costly. Therefore, we propose a procedural generation algorithm for constructing environments based on cuboid primitives. Online rearrangement of cuboids at different poses and dimensions yields wide coverage of geometric features that exist in the real world, such as walls, ceilings, slopes, and bumps. Further, the convex geometry of cuboids affords efficient dynamics simulation. Figure 3 shows example environments from our environment generation algorithm.

We show that by leveraging our framework, we can train a non-prehensile manipulation policy that can operate in diverse and novel environments and objects in a data- and time-efficient manner. We train our policy entirely in simulation and zero-shot transfer to unseen environments and objects in the real world. Furthermore, we provide a simulation-based benchmark comprising 9 digital twins of real-world environments and 353 objects to serve as a benchmark for non-prehensile manipulation for general environments and objects.

TABLE I: Comparison of generalization capabilities for non-prehensile manipulation.

	Object generalization	General action space	Environment generalization
HACMAN [122]	O	X	X
CORN [17]	O	O	X
Wu et al. [109]	△	△	△
Ours	O	O	O

II. RELATED WORK

A. Nonprehensile Manipulation

1) *Planning-based approaches:* Prior works on planning-based non-prehensile manipulation use gradient-based optimization [67, 85, 68], graph-based search [56, 55, 66, 15, 49, 14], or a hybrid of both [10, 75]. To address the discontinuous dynamics arising from contact mode transitions, optimization-based works employ soft contact variables [67] or complementarity constraints [85, 68]. However, due to the imprecision of smoothed contact and the difficulty of precise constraint satisfaction, the resulting motions are difficult to realize in the real world.

On the other hand, graph-based methods [14, 40] can handle discrete dynamics transitions by representing the problem with a graph, where nodes represent robot states and contact modes, and edges encode the motion during the transition. This enables these methods to output more physically realistic motions for real-world deployment [15, 49]. However, to make the search tractable, these works restrict the diversity of motions, assuming quasi-static motions [14, 40] or predefined primitives [124, 49], limiting them to tasks with simple motions and sparse contact-mode transitions.

Other works combine optimization and sampling to accelerate planning [10, 75], yet remain too slow for online use due to the cost of searching large hybrid spaces with discontinuous dynamics. Further, most planning methods require knowledge of system parameters like mass and friction, which are difficult to estimate in real-world scenarios with varying objects and environments, harming practical real-world deployment.

2) *Learning-based approaches:* Recent works leverage reinforcement learning (RL) to bypass the limitations of traditional planners by learning a policy that maps actions directly from sensory inputs [52, 117, 81, 118, 31, 121, 45, 122, 17, 109]. While this circumvents the computational cost of planning or the requirement of full system parameters, most of these works suffer from limited generalization across object geometries, since the policy is only trained on a single object [52, 117, 81, 118, 31, 121, 45]. Recent works incorporate point-cloud inputs [122, 17] or employ contact re-targeting [109] to facilitate generalization across diverse object shapes, but none of these approaches adequately addresses the problem of generalizing across diverse environments using the full action space of the robot, as summarized in Table I.

To generalize across diverse objects, HACMan [122] predicts an object-centric affordance map on its point cloud. While sample-efficient, this restricts robot motion to a hand-designed poking primitive, limiting their applicability in diverse environments. CORN [17] learns a policy over the full joint space of the robot, and generalizes over objects with a

contact-based object representation tailored for manipulation. However, this work remains limited to a fixed tabletop due to the lack of environment representation. Wu et al. [109] retargets contacts from human demonstrations to determine robot actions for novel scenes, but their approach is limited to scenes similar to the original demonstration, as the actions are restricted to predefined skill sequences based on a primitive library. Like CORN, we train an RL policy over the full robot joint space to manipulate objects of general geometry. However, our approach generalizes across environments by leveraging our extended contact-based representation, UNICORN, and a modular network architecture, HAMNET.

B. Multi-task Neural Architectures

In multi-task learning, a single model learns to do a family of related tasks by leveraging task synergies for improved performance and training efficiency [9]. As a single model must distinguish multiple tasks, it additionally takes context variables (e.g., task IDs) as conditioning inputs. The simplest approach for multi-task learning uses a *monolithic* architecture, which incorporates context inputs simply by concatenating them with network inputs. However, this design suffers from interference among tasks, as the neurons must handle multiple functions, leading to performance degradation [115, 42].

Recently, *context-adaptive* architectures have been proposed, where a separate network g_ϕ determines the parameters of base network f_θ from context inputs. See, for example, Figure 5 (a) and (c) for a comparison of monolithic and context-adaptive architectures. This separation allows the neural network to define a different function for each context, which mitigates interference [42]. Our architecture, HAMNET, also falls into this category.

Representative context-adaptive architectures include conditional normalization [82, 80, 65, 7], hypernetworks [37, 42], and modular architectures [34, 112, 97]. In conditional normalization [82] (Figure 5b), g_ϕ takes the context variable z as input and applies feature-wise scale and bias $\{\gamma, \beta\} = g_\phi(z)$ to the intermediate features of the base network as $y = \gamma \odot f_\theta(x) + \beta$ where \odot denotes element-wise multiplication. While effective, the expressivity of these architectures is limited to affine transforms, restricting its capacity [89]. In hypernetworks [37] (Figure 5c), g_ϕ affords broader expressivity, as it generates the entire set of base network parameters, $\theta = g_\phi(z)$, but suffer from poor training stability due to its large decision space over densely interacting parameters [71].

HAMNET is an instance of modular architectures (Figure 5d), where g_ϕ only predicts sparse activation weights w of modules to determine θ . Specifically, in modular architectures, M denotes the number of modules, each of which is a network parameter, and the parameters of the base network are effectively formed as a weighted combination of these modules, such as $\theta = \sum_{i=1}^M w_i \theta_i$.

However, prior works on modular architectures for control assume that z is given by the user and is discrete, such as pre-defined task IDs [84, 97, 112, 39]. In contrast, to generalize to novel, real-world environments in our problem, we need to

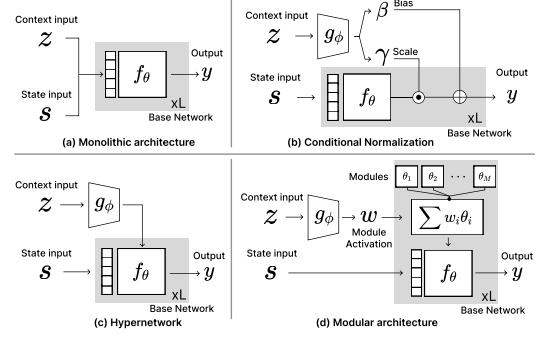


Fig. 5: Conceptual illustration of how different architectures incorporate context inputs. \oplus and \odot denote elementwise sum and multiplication, and L denotes number of layers. The base network that computes the output is indicated by f_θ , while g_ϕ is a separate network that determines θ from the context. In our work, we adopt a variant of the modular architecture (d).

infer z from sensory observations such as environment and object point clouds. To overcome this, we design geometric encoders that map high-dimensional sensory observations to z .

Our architecture is most similar to Soft Modularization (SM) [112], with three key differences. Unlike SM, which predicts *connectivities* between all pairs of modules between neighboring layers, we directly predict *module activations*, which simplifies computation and reduces the output dimensions from M^2 connections to M module activations. We additionally improve the computational efficiency by predicting all module activation weights in parallel, instead of predicting in series while conditioning on the preceding layer’s module activations. Lastly, we incorporate a gating mechanism [37] to enhance expressivity and boost policy performance.

C. Representation learning on point clouds

To accelerate RL training with high-dimensional sensory inputs, prior works use representational pre-training [114, 29] to bootstrap RL agents. Different pretext tasks have been proposed for this purpose, such as point completion [105], orientation and category estimation [12, 41], or contrastive learning [110]. Inspired by the advances in natural language processing [22, 87] and image analysis [38], recent works adopt self-supervised learning (SSL) on patch-based transformers [116, 11, 76, 120, 1] for point cloud representation learning. These works reconstruct unseen geometric patches via either autoregressive prediction [11] or masking [116, 76, 120, 1] to learn rich geometric representations, achieving state-of-the-art results in shape classification and segmentation [11].

Despite their success in general-purpose vision tasks, these representations are unsuitable for robot manipulation for two reasons. First, these models attempt to predict the missing patches in a point cloud, which forces the encoder to focus on encoding information about the object’s shape. However, knowing the exact shape of the object is often sufficient but unnecessary for manipulation. For example, manipulating a toy crab in Figure 1e does not require knowledge of the

exact shape between its legs, as that area is tightly confined and cannot be contacted by the robot or the environment. Second, capturing such spurious details requires a large model, which degrades training efficiency [28, 94] and policy performance [119]. In contrast, we extend CORN [17] to pre-train a representation to encode contact affordances among arbitrary geometry pairs, shown to be effective for robot manipulation.

D. Modularity in biological networks

Modularity in biological neural networks is a key principle underlying adaptation and learning [96]. In vertebrate motor systems, *muscle synergies* [103, 20, 72, 61, 23] serve as modules of movement that abstract muscle control, representing a coordinated contraction of a set of muscles to produce a desired behavior, such as the synchronized activation of the quadriceps and hamstrings for walking [23].

Modularizing motor control in this way provides several benefits. When adapting to a particular context, the central nervous system (CNS) can dictate behaviors using sparse, low-dimensional signals that activate specific muscle groups [5]. Compared to controlling individual motor neurons, this affords rapid switching between distinct motor skills like reaching and grasping depending on the context [73, 101]. Further, synergies can be reused across behaviors, producing diverse movements such as pinch- or power-grasps from a limited set of modules [86], which facilitates learning by recombining and adapting existing modules [18, 25]. In our work, we incorporate these principles to design our architecture.

E. Skill discovery in RL

In skill discovery, Unsupervised RL (URL) aims to find reusable skills using task-agnostic objectives like state coverage [79, 51, 64, 83, 113, 8] or skill diversity [35, 26, 93, 47, 77]. However, without task-specific priors, such intrinsically motivated methods often fail to make meaningful interactions in high-dimensional domains without engineered bias [3, 102]. While METRA [78] was shown to scale to pixel-based tasks, it remains prone to degenerate behaviors, such as lying still in varied poses in humanoid control domains [102]. On the other hand, our framework learns task-specific skills grounded in the training domain, yielding interpretable skills that transfer to real-world robot manipulation.

F. Mitigating Spectral Bias in RL

A central challenge in multi-task learning is enabling networks to adaptively *switch* functions under subtle shifts in its input. While our work focuses on the design of *network architecture*, one complementary line of work investigates improving *input representations* with high-frequency encodings, such as Learned Fourier Features (LFFs)[99], to enhance the network’s sensitivity to sharp transitions. While promising, the effectiveness of LFFs remains inconclusive in RL, with conflicting claims on the role of different frequency components [48, 111, 60]. Moreover, LFFs lead to an exponential increase in input dimensionality [6] and may degrade training stability from the increased variance, impairing generalization

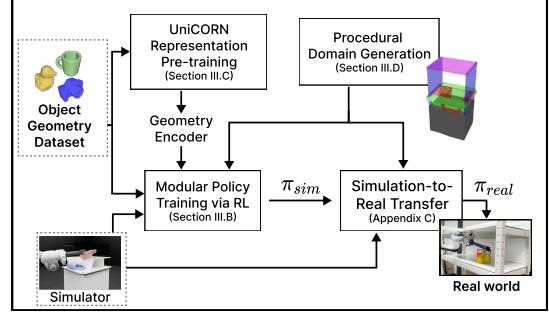


Fig. 6: Overall method overview. Our framework consists of four main components: a modular policy trained with RL, contact-based representation pre-training, a procedural domain generation scheme for environment geometries, and a simulation-to-real transfer method for real-world deployment. Dashed blocks indicate external inputs.

and robustness to noise [60]. These limitations suggest the need for an *architecture-level* solution, and not just the input representation. Still, one potentially promising future direction is to see whether high-frequency embeddings can complement the modular architecture by increasing network sensitivity to distinct contexts when determining module activations.

III. METHODOLOGY

We consider the non-prehensile manipulation problem, where a robot arm with a fixed base moves an object to a target pose in environments of general geometry, e.g., kitchen sinks, cabinets, and drawers (Figure 1). We model this problem as a Markov Decision Process (MDP), represented as a tuple (S, A, P, r, γ) denoting state space S , action space A , state transition model $P(s_{t+1}|a_t, s_t)$, reward model $r(s_t, a_t, s_{t+1})$, and discount factor γ . Our objective is to obtain policy π that maximizes the return $R_t = \mathbb{E}_{a_t \sim \pi(\cdot|s_t)}[\sum \gamma^t r(s_t, a_t, s_{t+1})]$ via a sequence of non-prehensile actions.

Figure 6 presents an overview of our framework. We leverage deep RL in a parallel GPU-based simulation [58] to train a modular policy (Section III-B) using the pre-trained point cloud representation (Section III-C) on procedurally generated domains (Section III-D). We distill the resulting policy for real-world deployment via teacher-student distillation (Appendix C). All pre-training, policy training, and distillation stages happen entirely in a simulation.

A. MDP Design

Our state space S consists of robot joint state x_t^q , end-effector pose x_t^{EE} , object geometry G_o , environment geometry G_e , and goal pose T_g . In the simulation, the agent additionally receives physics parameters ν and object state x_t^o . We represent all poses as 3D translation and 6D orientation to facilitate learning [123], and the goal is given as a relative pose from the current object pose. Object and scene geometries are given as surface-sampled point clouds.

Our action space A consists of joint residuals $\Delta q \in \mathbb{R}^7$ and controller gains, parameterized by proportional gains $k_p \in \mathbb{R}^7$ and the damping ratio $\rho \in \mathbb{R}^7$ that maps to the damping gain k_d as $\rho\sqrt{k_p}$, following [59, 45]. The resulting torque

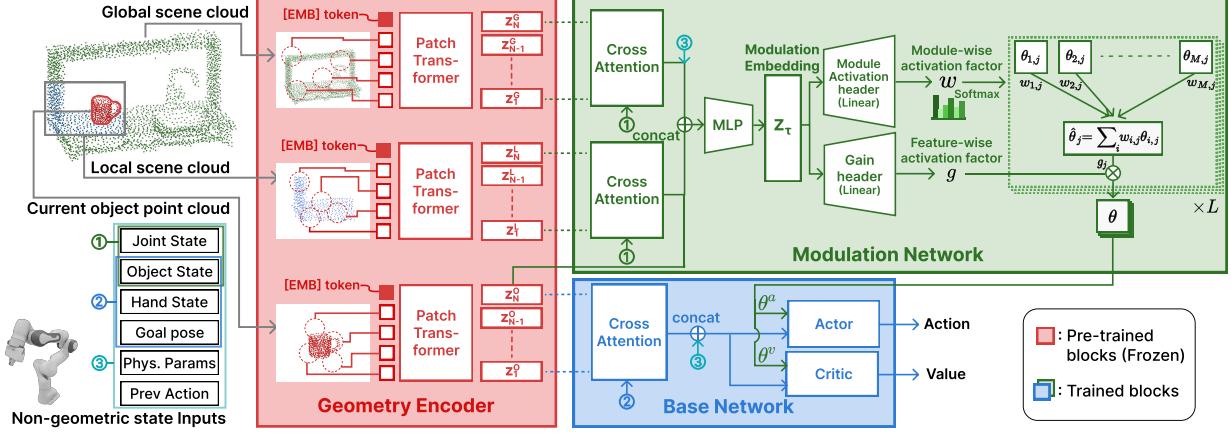


Fig. 7: Overall architecture. Our model comprises three components – the geometry encoder (red), the modulation network (green), and the base network (blue). The geometry encoder embeds the point clouds, and the modulation network maps the embeddings and non-geometric state inputs to the base network’s parameters θ . Conditioned on θ , the base network maps the state inputs and object geometry to actions and values. Input groups tagged with different numbers (①, ② and ③) indicate sets of non-geometric state inputs fed into different network parts. The inputs to cross-attention layers are concatenated and tokenized by a two-layer multi-layer perceptron (MLP).

for each joint is computed as $\tau = k_p \Delta q - k_d \dot{q}$. While prior works adopt Cartesian-space actions [122, 17], we adopt joint-space actions, which enables direct control of individual robot links to avoid collisions against the environment during manipulation.

The reward $r(s_t, a_t, s_{t+1})$ in our domain is defined as a sum of the task success reward r_s , goal-reaching reward r_r and the contact-inducing reward r_c : $r = r_s + \lambda_r r_r + \lambda_c r_c$, where λ_r and λ_c are scaling coefficients for the respective rewards. Since $r_s = \mathbb{1}_{suc}$ is sparsely given, we incorporate shaping rewards r_r and r_c as potential functions of the form $\gamma\phi(s') - \phi(s)$ with the discount factor $\gamma \in [0, 1]$, which preserves policy optimality [69]. Specifically, we have $\phi_r(s) = -\log(c_g \cdot d_{o,g}(s) + 1)$ for r_r , and $\phi_c(s) = -\log(c_r \cdot d_{h,o}(s) + 1)$ for r_c , where $c_g, c_r \in \mathbb{R}$ are scaling coefficients for the distance-based potential functions; $d_{o,g}(s)$ is the relative distance between the current object and the goal pose, based on the bounding-box distance [2]; $d_{h,o}(s)$ is the hand-object distance between the object and the tip of the end-effector. Task success is achieved when the object’s pose is within 0.1m and 0.1 radians of the target pose. The episode terminates if (1) the object reaches the goal, (2) the object is dropped from the workspace, or (3) the episode reaches the timeout of 300 simulation steps. Table VII summarizes our MDP design, and details on reward coefficients are in Table VIII.

B. HAMNET-based architecture

Our architecture, shown in Figure 7, consists of three main components: the geometry encoder (red), modulation network (green), and base network (blue). Our proposed modular architecture, HAMNET, consists of the modulation and base networks. Since we use PPO, an actor-critic algorithm, our base network outputs both value and action.

The geometry encoder processes three types of point cloud inputs: the *global* scene cloud, capturing the overall geometry

of the scene; the *local* scene cloud, detailing the nearby scene that surrounds the object; and the *object* point cloud, representing its surface geometry. Each cloud is patchified, tokenized, and embedded by the pre-trained geometry encoder (Section III-C), yielding latent geometric embeddings $z_{1:N}^{(G)}$, $z_{1:N}^{(L)}$, $z_{1:N}^{(O)}$, for global, local, and object embeddings, respectively. Details on point cloud acquisition are in Appendix E.

The role of the modulation network (Figure 7, green) is to output the parameters of the base network, θ . It takes the geometry embeddings $z_{1:N}^{(G)}$, $z_{1:N}^{(L)}$, $z_{1:N}^{(O)}$ and non-geometric states (① and ③) as input. To extract scene geometry information relevant to the policy’s current state, we apply cross-attention on the scene geometry embeddings $z_{1:N}^{(G)}$ and $z_{1:N}^{(L)}$, using the current robot and object states ① as queries. The resulting vector is concatenated with object geometry embedding $z_N^{(O)}$ and full non-geometric state inputs ③, and passed through an MLP to predict z_{τ} , the modulation embedding. Finally, the module activation and gain headers map z_{τ} to module activation weights w and gating values g respectively, for the L base network layers.

We then use w and g to build the base network parameters θ . For each layer $j \in [1 \dots L]$, $w = \{w_{i,j}\}_{j=1}^L \in \mathbb{R}^{L \times M}$ act as M module-wise weighting coefficients, passed through softmax to ensure $\sum_{i=1}^M w_{i,j} = 1$. The gating factor $g = \{g_j\}_{j=1}^L \in \mathbb{R}^{L \times D_j}$ is a feature-wise multiplier for each layer, with D_j denoting the number of output dimensions of layer j . Together, θ is constructed as a weighted composition of modules followed by gating, such that $\theta = \{(\sum_{i=1}^M w_{i,j} \theta_{i,j}) \odot g_j\}_{j=1}^L$.

The base network (Figure 7, blue) comprises actor and critic networks, where each network is an MLP. To produce the input for the base network, we first process the object embedding $z_{1:N}^{(O)}$ via cross-attention against input group ②, then concatenate the result with input group ③. The actor network outputs the action, and the critic network outputs the

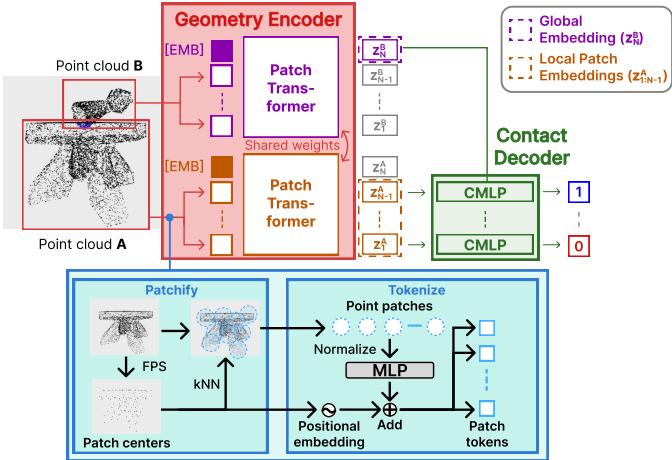


Fig. 8: Our pre-training architecture consists of a geometry encoder (red) and a contact decoder (green). The same geometry encoder operates on each point cloud A and B in a Siamese fashion to produce local patch embeddings $z_{1:N-1}^A$, $z_{1:N-1}^B$ and global embeddings z_N^A , z_N^B . The contact decoder (green) predicts contact between each patch $z_i^A \in z_{1:N-1}^A$ and z_N^B . The bottom block details the procedure to patchify and tokenize point clouds.

state values, but instead of a single scalar value, it uses three heads to predict the value for each reward component in our domain: r_s , r_r , and r_c . Since summing the rewards conflates the contributions from different reward terms, splitting the critic into multiple headers helps decrease the difficulty of value estimation [30, 54]. When training the actor network, we sum the advantages across reward terms to compute the policy gradients.

Note that since the base network has both the actor and critic, we keep separate sets of modules for each of them, denoted $\{\theta_{i,j}^{(a)}\}_{i=1}^M$, and $\{\theta_{i,j}^{(v)}\}_{i=1}^M$ for layer j . Our module activation weights and gating factor also consist of weights for a value and action, $w = \{w^{(v)}, w^{(a)}\}$ and $g = \{g^{(v)}, g^{(a)}\}$. To make a prediction, the base network gets instantiated twice for actor and critic; in the former case, the network uses the weight $\theta_j^{(a)} = \sum_{i=1}^M w_{i,j}^{(a)} \theta_{i,j} \odot g_j^{(a)}\}_{j=1}^L$, and in the latter, the network uses $\theta_j^{(v)} = \sum_{i=1}^M w_{i,j}^{(v)} \theta_{i,j} \odot g_j^{(v)}\}_{j=1}^L$. These details are omitted in the figure for brevity.

C. Training UNICORN

We design our representation pre-training task on estimating the presence and location of contact between two point clouds, A and B.

1) *Pre-training data generation* : To acquire data for pre-training, we generate a dataset containing pairs of objects represented as point clouds, and contact labels indicating the presence and location of contact. Using the objects from DexGraspNet dataset [106], we generate the data by (1) sampling *near-contact* object configurations, (2) creating the point clouds by sampling points from the surface of each object, and (3) labeling contact points based on whether they fall within the other object. To account for possible scale

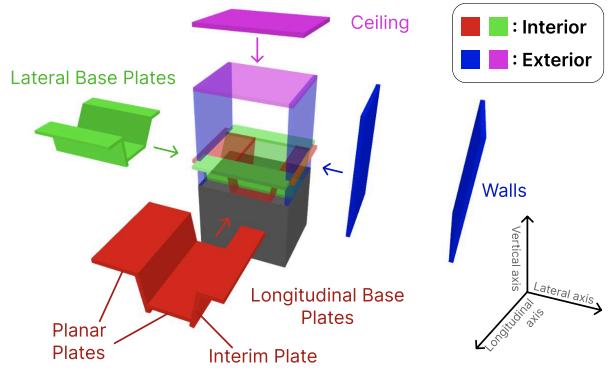


Fig. 9: Our pipeline for environment generation composes different environmental factors, such as walls, ceilings, and plates at different elevations for each axis, to construct geometrically diverse environments.

variations between geometries, we sample the point clouds at varying densities and scales during this process. For details on the data generation pipeline, see Appendix A.

2) *Network Architecture* : Figure 8 shows our pretraining network architecture, comprising the geometry encoder and the contact decoder. The encoder takes the point clouds of objects A and B, denoted x_A, x_B as inputs, mapping the patch-wise tokens from x_A and x_B and a learnable [EMB] token to local patch embeddings $z_{1:N-1}^A$, $z_{1:N-1}^B$ and global embeddings z_N^A, z_N^B . Afterward, the decoder takes $(z_i^A, z_N^B)_{i=1}^N$ and predicts the presence of contact at each of i -th local patch of object A with object B. The overall network is trained via binary cross-entropy against the patch-wise contact labels. During training, we alternate the roles of A and B (i.e., A-B and B-A) to ensure that we also use the global embedding of A and predict the contact at a patch of B.

Figure 8 (bottom) shows the procedure to tokenize the point clouds. In line with previous patch-based transformer architectures for point clouds [76, 11, 17], we first patchify the point cloud by gathering neighboring points from representative center points. These center points are selected via farthest-point sampling (FPS), and the points comprising the patches are determined as the k-nearest neighbors (kNN) of the patch center. These patches are normalized by subtracting their center coordinates, and a small MLP-based tokenizer [17] embeds the shape of each patch. Afterward, we add sinusoidal positional embeddings of the patch centers to the patch tokens to restore the global position information of each patch.

D. Procedural domain and curriculum generation

To create diverse environments and support curriculum learning for training our policy, we develop a procedural generation scheme for constructing environments as a composition of cuboidal primitives. Since we construct environments by dynamically rearranging existing geometric entities, it integrates well with most GPU-based simulators [58, 33] that prohibit spawning new assets after initialization.

Our procedural pipeline, shown in Figure 9, comprises two main components: domain *interior* and *exterior* generation.

TABLE II: Comparison between baselines regarding architecture and representation.

Model Name	Model Architecture	Representation
UNICORN-HAMNET (OURS)	HAMNET	UNICORN
UNICORN-HYPER	Hypernetwork	UNICORN
UNICORN-SM	Soft-Modularization [112]	UNICORN
UNICORN-TRANSFORMER	Transformer	UNICORN
UNICORN-MONO	MLP	UNICORN
POINTGPT-HAMNET	HAMNET	PointGPT [11]
E2E-HAMNET	HAMNET	End-to-end

The *interior* includes *planar* and *interim* plates arranged laterally and longitudinally, where *planar* plates form elevated surfaces and *interim* plates form sloped ramps. Their dimensions, elevations, and angles are randomly sampled to produce diverse topographies, yielding features like bumps, valleys, and steps (Figure 3). The *exterior* consists of walls and ceilings that impose accessibility constraints, where their presence, height, and ceiling type (*nominal* or *tight*) are randomly sampled. The proportion of ceiling types controls the difficulty of workspace accessibility, since the *nominal* ceiling is generated with sufficient clearance, whereas *tight* ceilings leave a narrow margin relative to the object’s height. Details on the procedural generation pipeline are in Appendix B.

As our procedural generation pipeline is fully parameterized, the sampling distributions of environmental parameters can be dynamically adjusted during training. This enables *curriculum learning*, where task complexity is incrementally increased throughout training. Specifically, we employ a curriculum for *robot initialization* and *ceiling types*. To facilitate this, we additionally introduce two types of robot initializations: *near* and *random*. In the *near* configuration, the end-effector begins within a 0.1m radius of the object based on collision-free inverse kinematics solutions from CuRobo [98]; in the *random* configuration, a collision-free joint configuration is uniformly sampled within the robot’s joint limits. Early in training, we preferentially sample *near* initializations and *nominal* ceilings to encourage interaction with the object. As training progresses, we linearly increase the proportion of *random* initializations and *tight* ceilings, encouraging the policy to develop obstacle-aware maneuvers for approaching objects from arbitrary configurations.

IV. EXPERIMENTAL RESULTS

A. Overview

Our goal is to evaluate the following claims: (1) our modular architecture, HAMNET, affords data-efficient training for a policy that generalizes over large domain diversity, compared to monolithic or hypernetwork architectures; (2) our contact-based representation, UNICORN, affords data-efficient training for a robot manipulation policy in geometrically rich domains compared to an off-the-shelf self-supervised representation; (3) our framework affords real-world transfer and generalization to novel environment geometries despite only training in a simulator with synthetic environments.

To evaluate our claims, we compare the performance of our proposed model (UNICORN-HAMNET) with the baselines

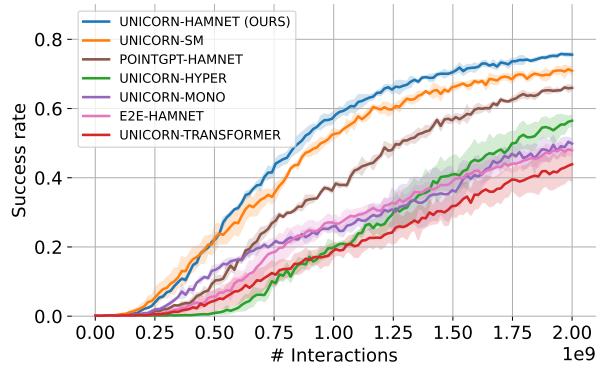


Fig. 10: Training progression. For each baseline, we show the mean (solid) and standard deviation (transparent) of success rates across three seeds. The interaction steps are aggregated across 1024 parallel environments.

summarized in Table II. These baselines explore alternative choices in network architecture or representation. UNICORN-HYPER uses a hypernetwork [50, 91] to predict base network parameters. UNICORN-SM is a variant of a modular architecture using Soft Modularization [112]. We include two variants of standard monolithic architectures, UNICORN-MONO and UNICORN-TRANSFORMER, respectively using an MLP and a transformer. POINTGPT-HAMNET and E2E-HAMNET consider alternative choices in the representation, where the former replaces the pretrained UNICORN with a PointGPT encoder [11], while the latter jointly trains the representation model end-to-end. All architectures are configured to have a similar number of trainable parameters up to the architectural constraints. Additional details on the baselines are in Appendix F.

B. Simulation experiment

To train our policy, we use a Franka Research 3 (FR3) arm manipulating a subset of 323 objects from DexGraspNet dataset [106] on the procedurally generated environments as in Section III-D. We train each baseline using PPO [92] with identical hyper-parameters, spanning 2 billion environment interactions across 1024 parallel environments in Isaac Gym [58]. Detailed hyperparameters for policy training are described in Appendix F. We consider two metrics: data efficiency and time efficiency.

To support our claim on training efficiency, we consider the training progression plot in Figure 10. Overall, modular architectures (UNICORN-HAMNET and UNICORN-SM) achieve the best data efficiency, with the mean success rates of 75.6% and 70.9% after training. In contrast, monolithic architectures show lower performance, regardless of whether conditioning is given by concatenation (UNICORN-MONO, 49.8%) or self-attention (UNICORN-TRANSFORMER, 43.9%). UNICORN-HYPER performs best among non-modular architectures at 56.5%, indicating the adaptivity of the network expedites policy training. UNICORN-SM (70.9%) further improves over hypernetworks, as its modularity affords reuse of network modules and reduces the learning complexity by predicting

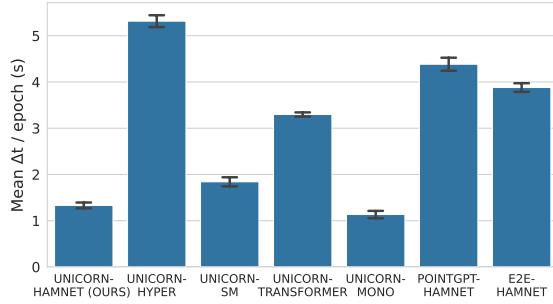


Fig. 11: Per-epoch training time comparison across all baselines, measured on identical hardware (NVIDIA A6000). Error bars represent two standard deviations.

sparse module activations rather than parameters of individual neurons. Lastly, UNICORN-HAMNET (75.6%) outperforms UNICORN-SM from the increased expressivity of the gating mechanism. To evaluate the representational efficacy of UNICORN, we also compare UNICORN-HAMNET to E2E-HAMNET and POINTGPT-HAMNET. End-to-end training (E2E-HAMNET) degrades performance (47.9%), emphasizing the utility of pre-training; while POINTGPT-HAMNET performs better (66.0%), it still underperforms UNICORN-HAMNET due to the overhead from spurious geometric details and increased embedding dimensions.

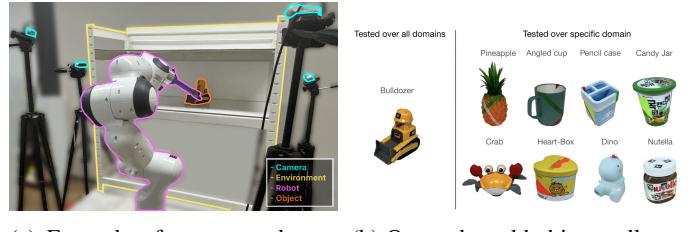
Since all our training happens in simulation, training time is also an important factor. To evaluate the time-efficiency of each baseline, we measure the per-epoch training time in Figure 11. While the monolithic MLP architecture (UNICORN-MONO) is the fastest (1.12s) due to its simplicity, modular architectures (UNICORN-HAMNET and UNICORN-SM) follow closely at just 1.33s and 1.84s, which shows that determining module-level activation adds minimal overhead; between the two, UNICORN-HAMNET achieves faster training than UNICORN-SM from the streamlined prediction of module activations. In contrast, UNICORN-HYPER takes the longest (5.31s) due to the cost of predicting the full set of base network parameters. To contextualize the overhead, a standard transformer (UNICORN-TRANSFORMER) takes around 3.30s. We also compare with representational baselines: POINTGPT-HAMNET and E2E-HAMNET. The large encoder of POINTGPT-HAMNET significantly lags training, averaging 4.40s per epoch, about 3.3 times slower than UNICORN-HAMNET. While E2E-HAMNET uses the same encoder architecture as UNICORN-HAMNET, it suffers from 2.9× slower training due to the overhead of co-training the representation model.

C. Real world experiment

To validate the real-world applicability and generalizability of our framework, we evaluate our policy in 9 real-world domains with novel everyday scenes and objects (Figure 1). We test two objects in each domain: one object (a toy bulldozer), shared across all domains, and one random object (Figure 12b), each with five trials at different initial and goal

TABLE III: Results on 9 unseen real-world domains.

Domain	Object	Success rate	Domain	Object	Success rate
Cabinet	Bulldozer	4/5	Top of cabinet	Bulldozer	3/5
	Heart-Box	3/5		Crab	4/5
Sink	Bulldozer	5/5	Basket	Bulldozer	3/5
	Angled Cup	4/5		Heart-Box	5/5
Drawer	Bulldozer	4/5	Grill	Bulldozer	5/5
	Pencil case	3/5		Dino	4/5
Circular bin	Bulldozer	4/5	Flat	Bulldozer	5/5
	Pineapple	3/5		Nutella	3/5
Suitcase	Bulldozer	4/5	Total		
	Candy Jar	5/5			78.9%



(a) Example of our scene layout in the *cabinet* domain.
(b) Our real-world objects, all unseen during policy training.

Fig. 12: Our real-world experimental setup.

poses.

For each scene, we mount four RealSense D435 cameras to observe the point clouds from multiple viewpoints, ensuring sufficient visibility of the object during execution (Figure 12a). To distinguish the object cloud from the environment cloud, we use SAM [46] to designate the initial object segmentation mask and utilize Cutie [13] to track the object during manipulation. We use FoundationPose [108] to estimate the object’s relative pose from the goal pose, using the view with the best visibility of the object (largest object segmentation mask) among the four cameras. We generate the environment point cloud by combining and filtering the point clouds from the depth cameras. We replaced the robot’s gripper to accommodate narrow environments, wrapped with a high-friction glove to reduce slipping. Further details on the real-world setup are in Appendix H.

Table III shows the results of our policy across 9 real-world domains. Overall, our policy demonstrates 78.9% success rate, indicating that our framework facilitates the policy to transfer to diverse, unseen real-world environments, despite only training in a simulation. The main failure modes of our policy are, in decreasing order of frequency: torque limit violation (5.56%); policy deadlock (4.44%); dropping objects (4.44%); getting blocked by the environment (3.33%); and perception error (3.33%). Detailed descriptions of these failure modes are in Appendix H2.

D. Emergence of skills in HAMNET

We show that HAMNET automatically discovers different manipulation skills and learns to sequence them. To do this, we inspect the modulation embedding z_τ (see Figure 7), which decides the activation weight of each module. We collect a dataset of z_τ by running a trained policy in 25,000 randomly sampled episodes in simulation. Since the high-dimensional z_τ is hard to interpret, we project z_τ into a three-dimensional manifold using UMAP [63] to visualize its structure.

To show that HAMNET discovers different skills, we apply

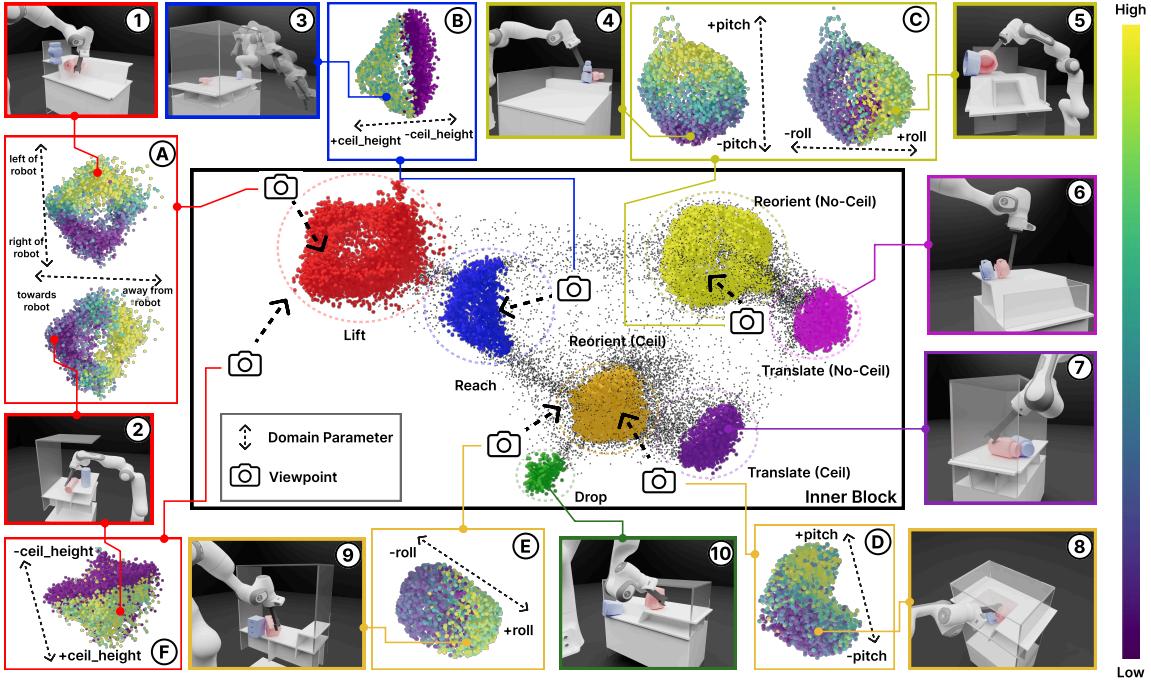
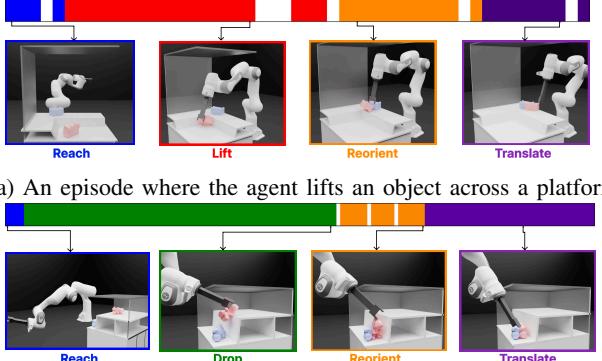


Fig. 13: (inner block) UMAP projection of the modulation embedding z_τ , colored by clusters from HDBScan. Unclustered points are in black. (outer block) Isolated view of each cluster, colored by a representative domain parameter, such as ceiling height or goal direction (to the left or right). The camera icon per each box denotes the viewpoint. The rendered scene shows a domain and state that generated an embedding in a cluster, with the red and blue objects indicating the current object pose and goal pose respectively.

HDBSCAN [62] to these projections of z_τ . Figure 13 shows the result. The *inner block* of Figure 13 shows the extracted clusters with different colors. We find that, without any manually designed bias or knowledge, these clusters naturally emerge and have semantically interpretable behaviors, such as *lifting* (red), *reaching* (blue), *reorienting* with (yellow) and without a ceiling (bright yellow), *translation* with (purple) and without a ceiling (bright purple), and *dropping* objects (green).

The *outer blocks* marked with numbers show the rendering of the situations in which these embeddings have been used. They show that our policy also learns *when* to use these skills based on the geometric constraint imposed by the environment, and the subgoal the robot is trying to achieve. For instance, to *lift* objects over platforms (Figure 13, ①), the policy must actively maintain contact between the object, wall, and the gripper. When *dropping* objects (Figure 13, ⑩), the robot must carefully prevent them from bouncing or rolling off the table. Similarly, the ceiling affects the policy’s *reaching* strategies: with the ceiling above the object, the robot must approach the object laterally (Figure 13, ⑧); in open environments, the robot can take overhand postures (Figure 13, ④) instead.

The *outer blocks* marked with alphabets in Figure 13 show that intra-cluster variation captures their behavioral variations within a skill. For example, as you move horizontally within the *lifting* cluster (Figure 13, ⑨) it models behavior that pulls the object towards or away from the robot ②, while the vertical direction maps to its left or right ①. While



(a) An episode where the agent lifts an object across a platform.
(b) An episode where the agent drops an object to a lower elevation.

Fig. 14: Illustration of how our architecture learns to use different skills. Color bar in each subfigure shows the cluster labels of z_τ at each step, and the bottom shows the domain rendering of representative keyframes. The red object is the current object pose, and the blue object is the goal pose.

subtler than the *categorical* differences across distinct skills like *lifting* and *reaching*, the emergence of such *intra-cluster* variations shows that HAMNET also learns to adjust the module activations to implement finer behavioral nuances.

To check if HAMNET can use these skills in sequence, we analyze how z_τ changes throughout a task. To highlight the transitions, we label z_τ at each step of an episode based on the precomputed HDBScan clustering shown in Figure 13. Figure 14 shows that the agent switches between behavioral

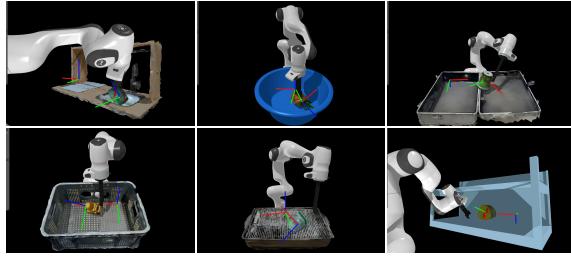


Fig. 15: Sample environments in our simulated benchmark. From the top left: sink, circular bin, suitcase, basket, grill, and cabinet. The axis on each object indicates its current pose, while the other axis represents the target pose.

clusters based on its internal subgoal: in Figure 14a, the robot initiates with a *reaching* skill to approach the object while avoiding obstacles. Afterward, the robot transitions to *lift* the object to the top platform. After a successful lift, the robot *reorients* the object to match the target orientation. Lastly, the robot *translates* the object to its target pose. The sequence of z_τ changes when the problem changes: when the robot has to drop an object to a lower platform instead (Figure 14b), the robot follows a different sequence (*reach-drop-reorient-translate*). This demonstrates that our architecture can (1) discover its own subgoals and (2) activate different modules to achieve different subgoals.

E. Simulated Benchmark in Realistic Domains

We release a simulated digital twin of our nine real-world setups as a benchmark for non-prehensile manipulation (Figure 15). The environment mesh is built using CAD, Nerfstudio [100], and Polycam. Our benchmark comprises 353 objects: 9 custom scans from the real world, 21 from GSO [24], and 323 from DGN [106]. For each domain and object pair, we sample 5 stable initial- and goal-poses and 128 random collision-free robot initializations to evaluate generality. Appendix G4 details domain configurations and provides baseline results.

V. CONCLUSION

In this work, we propose a novel framework for non-prehensile manipulation in general environments via deep reinforcement learning in a simulation. Our framework comprises a modular architecture (HAMNET), a contact-based object and environment representation (UNICORN), and a procedural domain generation algorithm for diverse environment geometries. Compared to conventional architectures and standard representations, our framework facilitates data- and time-efficient training of a policy that generalizes to diverse and unseen scenes. Despite solely training in synthetic environments in a simulation, our policy zero-shot generalizes to unseen real-world environments and objects. Overall, our combined framework achieves state-of-the-art performance in non-prehensile manipulation of general objects in general environments.

A. Limitations

Despite promising results, our approach has several limitations that can be addressed in future work:

Improved efficiency of HAMNET. In HAMNET, the parameters of the base network are updated at every step. However, our qualitative analysis (Section IV-D) shows that z_τ , and the resulting predictions for the module activation weights, remain stable for prolonged periods until a transition triggers a change, such as a successful lift. Thus, reusing the predicted network parameters over multiple steps can potentially reduce the computational overhead.

Dynamics-aware object representation. While UNICORN effectively represents object and environment geometries, it neglects dynamic properties like mass and inertia, which are critical for maneuvering objects with unusual dynamics, like roly-poly toys with non-uniform mass distributions. As such, one intriguing future research direction is to extend the representation to include dynamics information, potentially by incorporating memory [16, 36].

Generating fine-grained environment features. Our procedural generation pipeline relies on cuboidal primitives, limiting the diversity of fine-grained geometric features (e.g., textures, curvatures, small overhangs). While our experiments in the *Grill*, *Drawer*, and *Circular Bin* environments show that the policy can still adapt to uneven and curved surfaces, diversifying procedural generation through approaches like geometric generative models [95], may enhance the policy’s environmental generalization capability.

ACKNOWLEDGEMENTS

This work was supported by Institute of Information & communications Technology Planning & Evaluation (IITP) grant and National Research Foundation of Korea (NRF) funded by the Korea government(MSIT) (No.2019-0-00075, Artificial Intelligence Graduate School Program(KAIST)), (No.2022-0-00311, Development of Goal-Oriented Reinforcement Learning Techniques for Contact-Rich Robotic Manipulation of Everyday Objects), (No. 2022-0-00612, Geometric and Physical Commonsense Reasoning based Behavior Intelligence for Embodied AI), (No. RS-2024-00359085, Foundation model for learning-based humanoid robot that can understand and achieve language commands in unstructured human environments), (No. RS-2024-00509279, Global AI Frontier Lab).

REFERENCES

- [1] Karim Abou Zeid, Jonas Schult, Alexander Hermans, and Bastian Leibe. Point2Vec for self-supervised representation learning on point clouds. *German Conference on Pattern Recognition (GCPR)*, 2023.
- [2] Arthur Allshire, Mayank Mittal, Varun Lodaya, Viktor Makoviychuk, Denys Makoviichuk, Felix Widmaier, Manuel Wüthrich, Stefan Bauer, Ankur Handa, and Animesh Garg. Transferring dexterous manipulation from gpu simulation to a remote real-world trifinger. In *IEEE/RSJ International Conference on Intelligent*

- Robots and Systems (IROS)*, pages 11802–11809. IEEE, 2022.
- [3] Bowen Baker, Ingmar Kanitscheider, Todor Markov, Yi Wu, Glenn Powell, Bob McGrew, and Igor Mordatch. Emergent tool use from multi-agent autocurricula. In *International Conference on Learning Representations*, 2019.
- [4] Ronen Basri, Meirav Galun, Amnon Geifman, David Jacobs, Yoni Kasten, and Shira Kritchman. Frequency bias in neural networks for input of non-uniform density. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 685–694. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/basri20a.html>.
- [5] Emilio Bizzi and Vincent CK Cheung. The neural origin of muscle synergies. *Frontiers in Computational Neuroscience*, 7, 2013.
- [6] David Brellmann, David Filliat, and Goran Frehse. Fourier features in reinforcement learning with neural networks. *Transactions on Machine Learning Research*, 2023.
- [7] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alexander Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael S Ryoo, Grecia Salazar, Pannag R Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan H Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. RT-1: Robotics Transformer for Real-World Control at Scale. In *Proceedings of Robotics: Science and Systems*, Daegu, Republic of Korea, July 2023. doi: 10.15607/RSS.2023.XIX.025.
- [8] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=H1IJJnR5Ym>.
- [9] Rich Caruana. Multitask learning. *Machine learning*, 28:41–75, 1997.
- [10] Claire Chen, Preston Culbertson, Marion Lepert, Mac Schwager, and Jeannette Bohg. TrajectoTree: Trajectory optimization meets tree search for planning multi-contact dexterous manipulation. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 8262–8268. IEEE, 2021.
- [11] Guangyan Chen, Meiling Wang, Yi Yang, Kai Yu, Li Yuan, and Yufeng Yue. PointGPT: Auto-regressively generative pre-training from point clouds. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=rqE0fEQDqs>.
- [12] Tao Chen, Megha Tippur, Siyang Wu, Vikash Kumar, Edward Adelson, and Pulkit Agrawal. Visual dexterity: In-hand reorientation of novel and complex object shapes. *Science Robotics*, 8(84):eadc9244, 2023. doi: 10.1126/scirobotics.adc9244. URL <https://www.science.org/doi/abs/10.1126/scirobotics.adc9244>.
- [13] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Joon-Young Lee, and Alexander Schwing. Putting the object back into video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3151–3161, 2024.
- [14] Xianyi Cheng, Eric Huang, Yifan Hou, and Matthew T. Mason. Contact mode guided motion planning for quasidynamic dexterous manipulation in 3d. In *International Conference on Robotics and Automation (ICRA)*, pages 2730–2736, 2022. doi: 10.1109/ICRA46639.2022.9811872.
- [15] Xianyi Cheng, Sarvesh Patil, Zeynep Temel, Oliver Kroemer, and Matthew T Mason. Enhancing dexterity in robotic manipulation via hierarchical contact exploration. *IEEE Robotics and Automation Letters*, 9(1):390–397, 2023.
- [16] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, October 2014. URL <https://aclanthology.org/D14-1179>.
- [17] Yoonyoung Cho, Junhyek Han, Yoontae Cho, and Beomjoon Kim. CORN: Contact-based object representation for nonprehensile manipulation of general unseen objects. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=KTtEICH4TO>.
- [18] Jeff Clune, Jean-Baptiste Mouret, and Hod Lipson. The evolutionary origins of modularity. *Proceedings. Biological sciences / The Royal Society*, 280:20122863, 03 2013. doi: 10.1098/rspb.2012.2863.
- [19] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. In *Advances in Neural Information Processing Systems*, volume 35, pages 16344–16359, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/67d57c32e20fd0a7a302cb81d36e40d5-Paper-Conference.pdf.
- [20] Andrea d’Avella, Philippe Saltiel, and Emilio Bizzi. Combinations of muscle synergies in the construction of a natural motor behavior. *Nature neuroscience*, 6(3):300–308, 2003.

- [21] Dawson-Haggerty et al. trimesh. URL <https://trimesh.org/>.
- [22] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, June 2019. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- [23] Nadia Dominici, Yuri P. Ivanenko, Germana Cappellini, Andrea d’Avella, Vito Mondi, Marika Cicchese, Adele Fabiano, Tiziana Silei, Ambrogio Di Paolo, Carlo Giannini, Richard E. Poppele, and Francesco Lacquaniti. Locomotor primitives in newborn babies and their development. *Science*, 334(6058):997–999, 2011. doi: 10.1126/science.1210617. URL <https://www.science.org/doi/abs/10.1126/science.1210617>.
- [24] Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann, Thomas B. McHugh, and Vincent Vanhoucke. Google scanned objects: A high-quality dataset of 3d scanned household items. In *International Conference on Robotics and Automation (ICRA)*, pages 2553–2560, 2022. doi: 10.1109/ICRA46639.2022.9811809.
- [25] Kai Olav Ellefsen, Jean-Baptiste Mouret, and Jeff Clune. Neural modularity helps organisms evolve to learn new skills without forgetting old skills. *PLOS Computational Biology*, 11(4):1–24, 04 2015. doi: 10.1371/journal.pcbi.1004128. URL <https://doi.org/10.1371/journal.pcbi.1004128>.
- [26] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=SJx63jRqFm>.
- [27] Hao-Shu Fang, Chenxi Wang, Hongjie Fang, Minghao Gou, Jirong Liu, Hengxu Yan, Wenhui Liu, Yichen Xie, and Cewu Lu. AnyGrasp: Robust and efficient grasp perception in spatial and temporal domains. *IEEE Transactions on Robotics*, 39(5):3929–3945, 2023. doi: 10.1109/TRO.2023.3281153.
- [28] Zhiyuan Fang, Jianfeng Wang, Lijuan Wang, Lei Zhang, Yezhou Yang, and Zicheng Liu. {SEED}: Self-supervised distillation for visual representation. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=AHm3dbp7D1D>.
- [29] Jesse Farnbrother, Joshua Greaves, Rishabh Agarwal, Charline Le Lan, Ross Goroshin, Pablo Samuel Castro, and Marc G Bellemare. Proto-value networks: Scaling representation learning with auxiliary tasks. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=oGDKSt9JrZi>.
- [30] Mehdi Fatemi and Arash Tavakoli. Orchestrated value mapping for reinforcement learning. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=c87d0TS4yX>.
- [31] Juan Del Aguila Ferrandis, João Moura, and Sethu Vijayakumar. Nonprehensile planar manipulation through reinforcement learning with multimodal categorical exploration. *arXiv preprint arXiv:2308.02459*, 2023.
- [32] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-Aware Minimization for Efficiently Improving Generalization. *art. arXiv:2010.01412*, 2020.
- [33] C. Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier Bachem. Brax - a differentiable physics engine for large scale rigid body simulation, 2021.
- [34] Tomer Galanti and Lior Wolf. On the modularity of hypernetworks. *arXiv: Learning*, 2020.
- [35] Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. *Corr*, abs/1611.07507, 2016. URL <http://arxiv.org/abs/1611.07507>.
- [36] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces, 2024. URL <https://arxiv.org/abs/2312.00752>.
- [37] David Ha, Andrew M. Dai, and Quoc V. Le. Hypernetworks. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=rkpACe1lx>.
- [38] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16000–16009, June 2022.
- [39] Ahmed Hendawy, Jan Peters, and Carlo D’Eramo. Multi-task reinforcement learning with mixture of orthogonal experts. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=aZH1dM3GOX>.
- [40] Yifan Hou and Matthew T. Mason. Robust execution of contact-rich motion plans by hybrid force-velocity control. In *International Conference on Robotics and Automation (ICRA)*. IEEE, may 2019. doi: 10.1109/icra.2019.8794366. URL <https://doi.org/10.1109%2Ficra.2019.8794366>.
- [41] Wenlong Huang, Igor Mordatch, Pieter Abbeel, and Deepak Pathak. Generalization in dexterous manipulation via geometry-aware multi-task learning. *arXiv preprint arXiv:2111.03062*, 2021.
- [42] Siddhant M. Jayakumar, Wojciech M. Czarnecki, Jacob Menick, Jonathan Schwarz, Jack Rae, Simon Osindero, Yee Whye Teh, Tim Harley, and Razvan Pascanu. Multiplicative interactions and where to find them. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rylnK6VtDH>.
- [43] Shailesh Kantak, James Stinear, Ethan Buch, and

- Leonardo Cohen. Rewiring the brain: Potential role of the premotor cortex in motor control, learning, and recovery of function following brain injury. *Neurorehabilitation and neural repair*, 26:282–92, 09 2011. doi: 10.1177/1545968311420845.
- [44] Imin Kao, Kevin M. Lynch, and Joel W. Burdick. Contact Modeling and Manipulation. In *Springer Handbook of Robotics*, 2016.
- [45] Minchan Kim, Junhyek Han, Jaehyung Kim, and Beomjoon Kim. Pre-and post-contact policy decomposition for non-prehensile manipulation with zero-shot sim-to-real transfer. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10644–10651, 2023. doi: 10.1109/IROS55552.2023.10341657.
- [46] 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 Conference on Computer Vision*, pages 4015–4026, 2023.
- [47] Michael Laskin, Hao Liu, Xue Bin Peng, Denis Yarats, Aravind Rajeswaran, and Pieter Abbeel. Unsupervised reinforcement learning with contrastive intrinsic control. In *Advances in Neural Information Processing Systems*, volume 35, pages 34478–34491, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/debf482a7dbdc401f9052dbe15702837-Paper-Conference.pdf.
- [48] Alexander Li and Deepak Pathak. Functional regularization for reinforcement learning via learned fourier features. *Advances in Neural Information Processing Systems*, 34:19046–19055, 2021.
- [49] Jacky Liang, Xianyi Cheng, and Oliver Kroemer. Learning preconditions of hybrid force-velocity controllers for contact-rich manipulation. In *Proceedings of The 6th Conference on Robot Learning*, volume 205 of *Proceedings of Machine Learning Research*, pages 679–689. PMLR, 14–18 Dec 2023. URL <https://proceedings.mlr.press/v205/liang23a.html>.
- [50] Gidi Litwin and Lior Wolf. Deep meta functionals for shape representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1824–1833, 2019.
- [51] Hao Liu and Pieter Abbeel. Behavior from the void: Unsupervised active pre-training. In *Advances in Neural Information Processing Systems*, volume 34, pages 18459–18473, 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/99bf3d153d4bf67d640051a1af322505-Paper.pdf.
- [52] Kendall Lowrey, Svetoslav Kolev, Jeremy Dao, Aravind Rajeswaran, and Emanuel Todorov. Reinforcement Learning for Non-Prehensile Manipulation: Transfer from Simulation to Physical System. In *International Conference on Simulation, Modeling, and Programming for Autonomous Robots*, 2018.
- [53] Kevin M. Lynch and Matthew T. Mason. Dynamic nonprehensile manipulation: Controllability, planning, and experiments. *The International Journal of Robotics Research*, 18(1):64–92, 1999. doi: 10.1177/02783649901800105. URL <https://doi.org/10.1177/02783649901800105>.
- [54] James MacGlashan, Evan Archer, Alisa Devlic, Takuma Seno, Craig Sherstan, Peter Wurman, and Peter Stone. Value function decomposition for iterative design of reinforcement learning agents. In *Advances in Neural Information Processing Systems*, volume 35, pages 12001–12013, 2022.
- [55] Y. Maeda, H. Kijimoto, Y. Aiyama, and T. Arai. Planning of graspless manipulation by multiple robot fingers. In *International Conference on Robotics and Automation (ICRA)*, volume 3, pages 2474–2479, 2001. doi: 10.1109/ROBOT.2001.932994.
- [56] Yusuke Maeda and Tamio Arai. Planning of graspless manipulation by a multifingered robot hand. *Advanced Robotics*, 19(5):501–521, 2005.
- [57] Jeffrey Mahler, Jacky Liang, Sherdil Niyaz, Michael Laskey, Richard Doan, Xinyu Liu, Juan Aparicio, and Ken Goldberg. Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics. In *Proceedings of Robotics: Science and Systems*, Cambridge, Massachusetts, July 2017. doi: 10.15607/RSS.2017.XIII.058.
- [58] Viktor Makovychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, et al. Isaac gym: High performance gpu-based physics simulation for robot learning. *arXiv preprint arXiv:2108.10470*, 2021.
- [59] Roberto Martín-Martín, Michelle A. Lee, Rachel Gardner, Silvio Savarese, Jeannette Bohg, and Animesh Garg. Variable impedance control in end-effector space: An action space for reinforcement learning in contact-rich tasks. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1010–1017, 2019. doi: 10.1109/IROS40897.2019.8968201.
- [60] Augustine N Mavor-Parker, Matthew J Sargent, Caswell Barry, Lewis Griffin, and Clare Lyle. Frequency and generalisation of periodic activation functions in reinforcement learning. *arXiv preprint arXiv:2407.06756*, 2024.
- [61] David A. McCrea and Ilya A. Rybak. Organization of mammalian locomotor rhythm and pattern generation. *Brain Research Reviews*, 57(1):134–146, 2008. doi: <https://doi.org/10.1016/j.brainresrev.2007.08.006>. URL <https://www.sciencedirect.com/science/article/pii/S0165017307001798>.
- [62] Leland McInnes, John Healy, and Steve Astels. hdbSCAN: Hierarchical density based clustering. *The Journal of Open Source Software*, 2(11):205, 2017.
- [63] Leland McInnes, John Healy, Nathaniel Saul, and Lukas Grossberger. Umap: Uniform manifold approximation

- and projection. *The Journal of Open Source Software*, 3(29):861, 2018.
- [64] Russell Mendonca, Oleh Rybkin, Kostas Daniilidis, Danijar Hafner, and Deepak Pathak. Discovering and achieving goals via world models. In *Advances in Neural Information Processing Systems*, volume 34, pages 24379–24391, 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/cc4af25fa9d2d5c953496579b75f6f6c-Paper.pdf.
- [65] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [66] Kiyokazu Miyazawa, Yusuke Maeda, and Tamio Arai. Planning of graspless manipulation based on rapidly-exploring random trees. In *The International Symposium on Assembly and Task Planning: From Nano to Macro Assembly and Manufacturing*, pages 7–12. IEEE, 2005.
- [67] Igor Mordatch, Zoran Popović, and Emanuel Todorov. Contact-Invariant Optimization for Hand Manipulation. In *ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, 2012.
- [68] João Moura, Theodoros Stouraitis, and Sethu Vijayakumar. Non-prehensile planar manipulation via trajectory optimization with complementarity constraints. In *International Conference on Robotics and Automation (ICRA)*, pages 970–976, 2022. doi: 10.1109/ICRA46639.2022.9811942.
- [69] Andrew Y. Ng, Daishi Harada, and Stuart J. Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *Proceedings of the Sixteenth International Conference on Machine Learning*, page 278–287, 1999.
- [70] NVIDIA. TensorRT: A High-Performance Deep Learning Inference Library. <https://github.com/NVIDIA/TensorRT>, 2024.
- [71] Jose Javier Gonzalez Ortiz, John Guttag, and Adrian Dalca. Magnitude invariant parametrizations improve hypernetwork learning. *arXiv preprint arXiv:2304.07645*, 2023.
- [72] Simon A. Overduin, Andrea d’Avella, Jose M. Carmena, and Emilio Bizzi. Microstimulation activates a handful of muscle synergies. *Neuron*, 76(6):1071–1077, Dec 2012. doi: 10.1016/j.neuron.2012.10.018. URL <https://doi.org/10.1016/j.neuron.2012.10.018>.
- [73] Simon A. Overduin, Andrea d’Avella, Jinsook Roh, Jose M. Carmena, and Emilio Bizzi. Representation of muscle synergies in the primate brain. *Journal of Neuroscience*, 35(37):12615–12624, 2015. doi: 10.1523/JNEUROSCI.4302-14.2015. URL <https://www.jneurosci.org/content/35/37/12615>.
- [74] Jia Pan, Sachin Chitta, and Dinesh Manocha. FCL: A general purpose library for collision and proximity queries. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 3859–3866, 2012. doi: 10.1109/ICRA.2012.6225337.
- [75] Tao Pang, H. J. Terry Suh, Lujie Yang, and Russ Tedrake. Global planning for contact-rich manipulation via local smoothing of quasi-dynamic contact models. *IEEE Transactions on Robotics*, 39(6):4691–4711, 2023. doi: 10.1109/TRO.2023.3300230.
- [76] Yatian Pang, Wenxiao Wang, Francis EH Tay, Wei Liu, Yonghong Tian, and Li Yuan. Masked autoencoders for point cloud self-supervised learning. In *European conference on computer vision*, pages 604–621. Springer, 2022.
- [77] Seohong Park, Jongwook Choi, Jaekyeom Kim, Honglak Lee, and Gunhee Kim. Lipschitz-constrained unsupervised skill discovery. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=BGvt0ghNgA>.
- [78] Seohong Park, Oleh Rybkin, and Sergey Levine. METRA: Scalable unsupervised RL with metric-aware abstraction. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=c5pwL0Soay>.
- [79] Deepak Pathak, Dhiraj Gandhi, and Abhinav Gupta. Self-supervised exploration via disagreement. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pages 5062–5071. PMLR, 2019. URL <https://proceedings.mlr.press/v97/pathak19a.html>.
- [80] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4195–4205, October 2023.
- [81] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-Real Transfer of Robotic Control with Dynamics Randomization. In *International Conference on Robotics and Automation*, 2018.
- [82] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [83] Vitchyr Pong, Murtaza Dalal, Steven Lin, Ashvin Nair, Shikhar Bahl, and Sergey Levine. Skew-fit: State-covering self-supervised reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages 7783–7792. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/pong20a.html>.
- [84] Edoardo M Ponti, Alessandro Sordoni, Yoshua Bengio, and Siva Reddy. Combining modular skills in multitask learning. *arXiv preprint arXiv:2202.13914*, 2022.
- [85] Michael Posa, Cecilia Cantu, and Russ Tedrake. A Direct Method for Trajectory Optimization of Rigid Bodies Through Contact. *The International Journal of Robotics Research*, 2014.

- [86] Roberto Prevete, Francesco Donnarumma, Andrea d’Avella, and Giovanni Pezzulo. Evidence for sparse synergies in grasping actions. *Scientific Reports*, 8(1): 616, Jan 2018. doi: 10.1038/s41598-017-18776-y. URL <https://doi.org/10.1038/s41598-017-18776-y>.
- [87] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018.
- [88] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5301–5310. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/rahaman19a.html>.
- [89] Daniel Rebain, Mark J. Matthews, Kwang Moo Yi, Gopal Sharma, Dmitry Lagun, and Andrea Tagliasacchi. Attention beats concatenation for conditioning neural fields. *Transactions on Machine Learning Research*, 2023. URL <https://openreview.net/forum?id=GzqdMrFQsE>.
- [90] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Proceedings, 2011.
- [91] Elad Sarafian, Shai Keynan, and Sarit Kraus. Recomposing the reinforcement learning building blocks with hypernetworks. In *International Conference on Machine Learning*, pages 9301–9312. PMLR, 2021.
- [92] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [93] Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware unsupervised discovery of skills. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HJgLZR4KvH>.
- [94] Haizhou Shi, Youcai Zhang, Siliang Tang, Wenjie Zhu, Yaqian Li, Yandong Guo, and Yueting Zhuang. On the efficacy of small self-supervised contrastive models without distillation signals. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(2):2225–2234, Jun. 2022. doi: 10.1609/aaai.v36i2.20120. URL <https://ojs.aaai.org/index.php/AAAI/article/view/20120>.
- [95] Yawar Siddiqui, Antonio Alliegro, Alexey Artemov, Tatiana Tommasi, Daniele Sirigatti, Vladislav Rosov, Angela Dai, and Matthias Nießner. MeshGPT: Generating triangle meshes with decoder-only transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19615–19625, June 2024.
- [96] Olaf Sporns and Richard F. Betzel. Modular brain networks. *Annual Review of Psychology*, 67(Volume 67, 2016):613–640, 2016. doi: <https://doi.org/10.1146/annurev-psych-122414-033634>. URL <https://www.annualreviews.org/content/journals/10.1146/annurev-psych-122414-033634>.
- [97] Lingfeng Sun, Haichao Zhang, Wei Xu, and Masayoshi Tomizuka. PaCo: Parameter-compositional multi-task reinforcement learning. In *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=LYXTPNWJLr>.
- [98] Balakumar Sundaralingam, Siva Kumar Sastry Hari, Adam Fishman, Caelan Garrett, Karl Van Wyk, Valts Blukis, Alexander Millane, Helen Oleynikova, Ankur Handa, Fabio Ramos, Nathan Ratliff, and Dieter Fox. cuRobo: Parallelized collision-free minimum-jerk robot motion generation, 2023.
- [99] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. In *Advances in Neural Information Processing Systems*, volume 33, pages 7537–7547, 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/55053683268957697aa39fba6f231c68-Paper.pdf.
- [100] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Justin Kerr, Terrance Wang, Alexander Kristoffersen, Jake Austin, Kamyar Salahi, Abhik Ahuja, David McAllister, and Angjoo Kanazawa. Nerfstudio: A modular framework for neural radiance field development. In *ACM SIGGRAPH Conference Proceedings*, 2023.
- [101] Lena H Ting and J Lucas McKay. Neuromechanics of muscle synergies for posture and movement. *Current opinion in neurobiology*, 17(6):622–628, December 2007. doi: 10.1016/j.conb.2008.01.002. URL <https://europepmc.org/articles/PMC4350235>.
- [102] Andrea Tirinzoni, Ahmed Touati, Jesse Farenbrother, Mateusz Guzek, Anssi Kanervisto, Yingchen Xu, Alessandro Lazaric, and Matteo Pirotta. Zero-shot whole-body humanoid control via behavioral foundation models. In *International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=9sOR0nYLtz>.
- [103] Matthew C. Tresch, Philippe Saltiel, and Emilio Bizzi. The construction of movement by the spinal cord. *Nature Neuroscience*, 2(2):162–167, Feb 1999. doi: 10.1038/5721. URL <https://doi.org/10.1038/5721>.
- [104] Weikang Wan, Haoran Geng, Yun Liu, Zikang Shan, Yaodong Yang, Li Yi, and He Wang. UniDex-Grasp++: Improving dexterous grasping policy learning via geometry-aware curriculum and iterative generalist-specialist learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3891–3902, 2023.
- [105] Hanchen Wang, Qi Liu, Xiangyu Yue, Joan Lasenby, Modular

- and Matt J. Kusner. Unsupervised point cloud pre-training via occlusion completion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9782–9792, October 2021.
- [106] Ruicheng Wang, Jialiang Zhang, Jiayi Chen, Yinzhen Xu, Puhao Li, Tengyu Liu, and He Wang. Dexgraspnet: A large-scale robotic dexterous grasp dataset for general objects based on simulation. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 11359–11366. IEEE, 2023.
- [107] Shaochen Wang, Zhangli Zhou, and Zhen Kan. When transformer meets robotic grasping: Exploits context for efficient grasp detection. *IEEE Robotics and Automation Letters*, 7(3):8170–8177, 2022. doi: 10.1109/LRA.2022.3187261.
- [108] Bowen Wen, Wei Yang, Jan Kautz, and Stan Birchfield. Foundationpose: Unified 6d pose estimation and tracking of novel objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17868–17879, 2024.
- [109] Albert Wu, Ruocheng Wang, Sirui Chen, Clemens Eppner, and C Karen Liu. One-shot transfer of long-horizon extrinsic manipulation through contact retargeting. *arXiv preprint arXiv:2404.07468*, 2024.
- [110] Saining Xie, Jiatao Gu, Demi Guo, Charles R. Qi, Leonidas Guibas, and Or Litany. PointContrast: Unsupervised pre-training for 3d point cloud understanding. In *European Conference on Computer Vision (ECCV)*, 2020.
- [111] Ge Yang, Anurag Ajay, and Pulkit Agrawal. Overcoming the spectral bias of neural value approximation. In *International Conference on Learning Representations*, 2022.
- [112] Ruihan Yang, Huazhe Xu, Yi Wu, and Xiaolong Wang. Multi-task reinforcement learning with soft modularization. *Advances in Neural Information Processing Systems*, 33:4767–4777, 2020.
- [113] Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Reinforcement learning with prototypical representations. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pages 11920–11931. PMLR, 2021. URL <https://proceedings.mlr.press/v139/yarats21a.html>.
- [114] Lin Yen-Chen, Andy Zeng, Shuran Song, Phillip Isola, and Tsung-Yi Lin. Learning to see before learning to act: Visual pre-training for manipulation. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 7286–7293, 2020. doi: 10.1109/ICRA40945.2020.9197331.
- [115] Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Gradient surgery for multi-task learning. In *Advances in Neural Information Processing Systems*, volume 33, pages 5824–5836, 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/3fe78a8acf5fda99de95303940a2420c-Paper.pdf.
- [116] Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie Zhou, and Jiwen Lu. Point-BERT: Pre-training 3d point cloud transformers with masked point modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19313–19322, June 2022.
- [117] Weihao Yuan, Johannes A Stork, Danica Kragic, Michael Y Wang, and Kaiyu Hang. Rearrangement with Nonprehensile Manipulation Using Deep Reinforcement Learning. In *International Conference on Robotics and Automation*, 2018.
- [118] Weihao Yuan, Kaiyu Hang, Danica Kragic, Michael Y Wang, and Johannes A Stork. End-to-End Nonprehensile Rearrangement with Deep Reinforcement Learning and Simulation-to-Reality Transfer. *Robotics and Autonomous Systems*, 2019.
- [119] Amy Zhang, Rowan Thomas McAllister, Roberto Calandra, Yarin Gal, and Sergey Levine. Learning invariant representations for reinforcement learning without reconstruction. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=-2FCwDKRREu>.
- [120] Renrui Zhang, Ziyu Guo, Peng Gao, Rongyao Fang, Bin Zhao, Dong Wang, Yu Qiao, and Hongsheng Li. Point-M2AE: Multi-scale masked autoencoders for hierarchical point cloud pre-training. In *Advances in Neural Information Processing Systems*, volume 35, pages 27061–27074, 2022.
- [121] Wenxuan Zhou and David Held. Learning to grasp the ungraspable with emergent extrinsic dexterity. In *Conference on Robot Learning*, pages 150–160. PMLR, 2023.
- [122] Wenxuan Zhou, Bowen Jiang, Fan Yang, Chris Paxton, and David Held. HACMan: Learning hybrid actor-critic maps for 6d non-prehensile manipulation. In *Proceedings of The 7th Conference on Robot Learning*, volume 229 of *Proceedings of Machine Learning Research*, pages 241–265. PMLR, 06–09 Nov 2023. URL <https://proceedings.mlr.press/v229/zhou23a.html>.
- [123] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5745–5753, 2019.
- [124] Claudio Zito, Rustam Stolkin, Marek Kopicki, and Jeremy L Wyatt. Two-Level RRT Planning for Robotic Push Manipulation. In *International Conference on Intelligent Robots and Systems*, 2012.

APPENDIX

A. Details on training UNICORN

1) *Contact dataset generation:* Since predicting contacts between distant or fully overlapping objects is trivial, we prioritize sampling *near-contact* object configurations to encourage our model to learn informative representations. To achieve this, we follow CORN [17] and randomly sample two objects from the object dataset and position each object at a random SE(3) pose. Since this initial placement is unlikely to result in objects in colliding configurations, we first move the objects tangent to each other by measuring the shortest displacement between the two objects, and translating one of the objects by that amount so that they come into contact with each other. Afterward, we apply a small Gaussian noise to the poses so that the objects either slightly clip into each other or narrowly remain collision-free.

After positioning the objects, we generate the point clouds and contact labels. We sample the point clouds uniformly from the surface of each object, then label the points based on whether they fall within the other object. Since computing point-mesh intersection is often unreliable, we first apply convex decomposition to each object, then compute whether each point falls within any of the other object's convex parts, and vice versa. Overall, this procedure yields approximately half of the dataset comprising objects in colliding configurations, with the other half in near-contact configurations; the representation model must learn to distinguish the two scenarios. We iterate this procedure to generate a dataset comprising 500,000 point cloud pairs and their contact labels.

2) *Details on pretraining pipeline:* The contact decoder is a three-layer conditional MLP (CMLP), where each layer is a residual block with conditional batch normalization (CBN) [65]. CBN transforms each layer's output features by applying batch normalization, whose affine parameters are mapped from the conditioning input z_N^B with a single linear layer.

During training, we apply data augmentation by rotating, translating, and scaling both clouds, plus a small Gaussian noise. After patchifying the point clouds, we adjust the proportion of the inputs to the contact decoder so that approximately half of the input pairs are in contact (positive). This is done by resampling the positive patches with probability f/P and negative patches with probability $(1-f)/N$, where f denotes the target fraction of positive labels, P denotes the number of positive patches, and N denotes the number of negative patches. The hyperparameters for pretraining are in Table IV.

B. Details on Procedural Domain Generation

1) *Environment Geometry:* The flowchart in Figure 16 depicts our procedural generation pipeline, which samples environment geometries by composing a set of cuboidal primitives. Our pipeline is split into two main parts: the *interior* of the domain (Figure 16, purple), comprising the planar surface at various elevations and sloped hills formed by the *lateral* and *longitudinal* base plates (see Figure 9), and the *exterior*

TABLE IV: Pretraining pipeline hyperparameters.

Hyperparameter	Value
Batch size	1024
Optimizer	SAM [32]
Learning rate schedule	cosine
Base learning rate	0.0002
Min. learning rate	1e-6
Max. gradient norm	1,000
Weight decay	0.001
Rotational augmentation	$(-\pi, +\pi)$
Translational augmentation	$(-0.1, +0.1)$
Scale augmentation	(e^{-1}, e^{+1})
Noise augmentation	0.01
Positive patch fraction	0.5
Decoder size	(128, 128)

of the domain (Figure 16, cyan) comprising the ceiling and walls.

Our procedure begins by sampling the overall scene dimensions, i.e., whether a given table will be narrow or wide. This determines the overall scale of the cuboid primitives to fit within the bounds of the scene dimensions. Based on this, we compose the *interior* of the domain from a set of cuboidal plates. Each of the lateral and longitudinal axes comprises five plates: three *planar* plates form level surfaces at various heights, joined by two *interim* plates that form walls at various slopes.

To generate the *interior* of the domain (Figure 16, purple), we start by randomly sampling the plate dimensions, while ensuring the sum of the *planar* plate dimensions along each axis does not exceed the scene bounds. Afterward, the elevations of planar plates are randomly sampled by designating whether each plate is a *top* plate or a *bottom* plate, separated by a randomly sampled difference in elevation. Afterward, we sample the angles of the slopes corresponding to the interim plates that connect between two planar plates. Combining the planar and interim plates for both axes constructs the base surface of the domain, resulting in structural layouts such as sinks, bumps, or valleys (Figure 17).

Next, we generate the *exterior* of the domain, composed of walls and ceilings (Figure 16, cyan). Since these obstructions directly impact the accessibility of the workspace, we take a multi-step approach for positioning them. First, we randomly sample whether the ceiling should exist in the domain. Then, if the ceiling is present, the front-facing wall is deactivated to allow for the robot's entry. The remaining walls are randomly configured while ensuring at least one load-bearing wall is added on one of the four sides to support the ceiling. The disabled features are hidden beneath the tabletop to prevent interaction with other simulated entities.

Since the height of the ceiling heavily influences the accessibility of the object, we implement two different procedures to determine ceiling heights: *nominal* and *tight* (Figure 16,

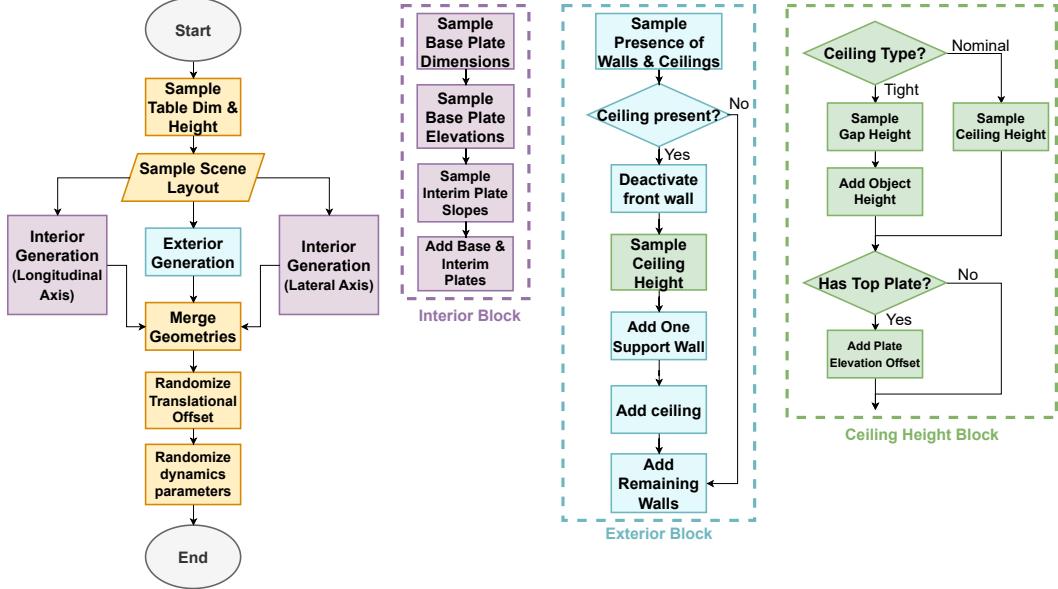


Fig. 16: Flowchart for our procedural generation pipeline. The leftmost flowchart describes the overall procedure; the expanded blocks (purple, cyan, and green) on the right show the subroutines in detail.

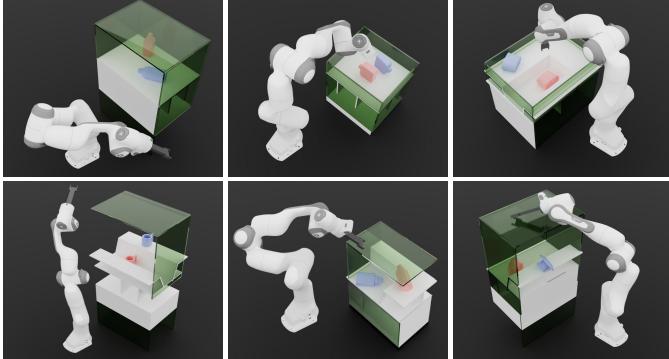


Fig. 17: Additional examples from our environment generation algorithm. While fully procedural and randomized, our pipeline yields geometries resembling those of real-world scenes, such as (from the top left): cabinet, basket, sink, valley, countertop, and step. Walls and ceilings are rendered with transparent green to distinguish them from the base plates. The red object indicates the object at initialization, and the blue object indicates the goal pose of the object.

green). The heights of *nominal*-type ceilings are simply sampled from a uniform distribution with a sufficient margin. On the other hand, heights of *tight*-type ceilings are determined by adding a small *gap* to the object height, where the height of the gap is drawn from a uniform distribution. Lastly, in both cases, if the interior of the domain has elevated platforms such as steps or bumps, we raise the ceiling heights by the height of the platform to prevent the object from intersecting with the ceiling.

After constructing the overall environment geometry, we further randomize the environment by shifting its position along each of the x, y, and z axes and adjusting its surface

friction coefficients. Table V lists the parameters for our procedural generation scheme.

TABLE V: Scene parameters and their ranges in our procedural generation pipeline. All angles are in degrees, and dimensions are in meters. U: uniform distribution; B: Bernoulli distribution.

Parameter	Value
table_dim.x	U(0.255, 0.51)
table_dim.y	U(0.325, 0.65)
table_dim.z	U(0.2, 0.4)
table_pos.x	U(0.0, 0.1)
table_pos.y	U(-0.15, 0.15)
table_pos.z	U(0.1, 0.8)
ramp_angle	U(0.0, 30) × 4
plate_height	U(0.0, 0.15) × 6
ceiling_height	U(0.3, 0.5)
gap_height	U(0.03, 0.05)
ceil_mask	B(0.5)
wall_mask	B(0.5) × 4
table_friction	U(0.2, 0.6)

2) *Object Placement*: Since our environment geometries change at the start of each episode, we must sample stable and collision-free object placements *online* to compute the initial and goal poses. Since this process is time-consuming, we precompute a set of stable object orientations as in [17] by dropping them in a simulation. Afterward, we also precompute the planar radius for each of the object’s stable orientations as the distance to the farthest point on the object from the object’s center.

For each episode, we sample object poses by combining one of the stable orientations with the position sampled from the horizontal plates in the environment. To compute collision-free and stable placements, we use the object’s precomputed radius to serve as the minimum distance away from the nearest wall or edge of the plates.

We first sample the goal poses for the objects, then sample their initial poses while ensuring sufficient separation from the goal in terms of both its position and orientation. This prevents the episode from terminating in success immediately. When sampling initial and goal poses in domains with height differences between plates, we bias the proportions to encourage the goal poses to be on elevated platforms compared to the initial poses, which favors sampling more challenging tasks where the robot must lift objects across a slope or a wall in the environment.

C. Simulation-to-Real Transfer Pipeline

The policy trained in simulation cannot be directly transferred to the real world. This is for two reasons. First, in contact-rich scenarios like non-prehensile object manipulation, the frequent contact between the robot and the object or the environment is prone to trigger hardware torque-limit violations. Second, the real-world policy does not have access to privileged information as in the simulation, such as the object’s mass and dynamics properties.

To overcome the torque-limit violations, we adopt two main strategies: *action magnitude curriculum* and *cartesian-space action clipping*, both of which reduce the scale of the policy actions to encourage conservative motions. To overcome the second issue of unavailable observations, we adopt *teacher-student distillation*, in which the student replicates the teacher’s actions solely from observations that would be available in the real world.

1) Action Magnitude Curriculum: In our domain, the robot frequently experiences contact with the object or the environment. However, when the policy operates at high velocity under mismatched robot dynamics due to the sim-to-real gap, such frequent contact may trigger hardware torque-limit violation when the impact is larger than expected.

As the robot must move the object through contact, the impact is inevitable. The policy cannot readily learn to prevent high-impact collisions either, since accurately reproducing the exact impact force in the simulation is challenging due to the sim-to-real gap. To circumvent this, we encourage the policy to perform generally conservative motions.

To this end, we adopt an action-magnitude curriculum inspired by the scheme from Kim et al. [45], where the maximum bounds of the subgoal residual that the policy can output is reduced gradually. During initial training, we start with the maximum joint-space residual magnitude $\xi = \xi_{\max}$ to facilitate the policy’s exploration. We then gradually reduce ξ to the target magnitude ξ^* , deemed safer for execution on the real robot. We apply different ξ_{\max} and ξ^* values for the large and small joints for the FR3 arm. Our reduction schedule follows a geometric sequence with ratio $\frac{\xi^*}{\xi_{\max}}^{\frac{N_s}{N_t}}$, where N_t

and N_s are hyperparameters that denote the total number of simulation steps for annealing and the interval between successive annealing steps, respectively. Detailed hyperparameters are in Table VI

TABLE VI: Sim2real hyperparameters.

Hyperparameter	Value
ξ^* (large joint)	0.16
ξ_{\max} (large joint)	0.26
ξ^* (small joint)	0.08
ξ_{\max} (small joint)	0.21
N_s	1024
N_t	2e6
ϵ_x	0.12
ϵ_x^{\max}	0.24
α	0.8

2) Cartesian-space Action Clipping: When using joint-space actions, multiple joints can simultaneously contribute to the same Cartesian direction. Despite reducing overall action magnitude, the sum of individual joint actions may still result in large end-effector movements, leading to high-force impacts that abort the robot. While further reducing the joint space could mitigate this, it would significantly degrade the robot’s dexterity and strength.

Algorithm 1

Cartesian-space action clipping algorithm.

Require: Policy π_θ , robot joint position q , cartesian action bound ϵ_x , damping parameter λ

Ensure: Clamped joint-space subgoal residual $\Delta q_{clamped}$

- 1: $J = \frac{\delta x}{\delta q}|_q$
 - 2: $\Delta q \sim \pi_\theta(s)$
 - 3: $\Delta x = J\Delta q$
 - 4: $\Delta x_{excess} = \Delta x \cdot \frac{\|\Delta x\| - \epsilon_x}{\|\Delta x\|}$
 - 5: $\Delta q_{excess} = J^T(JJ^T + \lambda^2 I)^{-1}\Delta x_{excess}$
 - 6: $\Delta q_{clamped} = \Delta q - \Delta q_{excess}$
-

Instead, we devise a scheme to clamp the joint residuals based on the projected end-effector space movement to reside within the bound ϵ_x with minimal change to the original action, shown in Algorithm 1. We first compute the Jacobian J of the robot in the current configuration (line 1). Then, we project the joint-space residuals Δq to the Cartesian-space end-effector movement Δx based on the Jacobian J (line 2-3). If the estimated end-effector movement exceeds the predefined Cartesian bound ϵ_x , we compute the excess movement Δx_{excess} compared to ϵ_x (line 4), then subtract the excess from the original joint residuals by re-projecting Δx_{excess} to Δq_{excess} based on damped least-squares method (line 5-6). Since clipping the action bounds potentially subjects the policy to a large behavioral change, we fine-tune the policy by starting with a large Cartesian bound ϵ_x^{\max} and gradually annealing it down to the target ϵ_x .

In addition to the magnitude scaling and Cartesian space

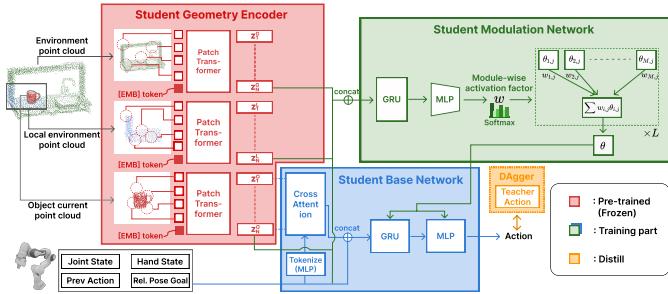


Fig. 18: Illustration of student policy architecture. As in the teacher network, the student architecture comprises a geometry encoder, a modulation network, and a base network.

clipping, we introduce a joint residual smoothing and energy-reducing loss for suppressing jerky motion during real-world deployment. Specifically, the smoothed residual $\Delta\bar{q}_t$ is computed from the original joint residual Δq_t with exponential moving average: $\Delta\bar{q}_t = \alpha q_t + (1-\alpha)\Delta\bar{q}_{t-1}$, where $\alpha \in [0, 1]$. The smoothed value $\Delta\bar{q}_t$ is then used to control the robot. The energy-reducing loss, computed by L2 norm of power $\|\sum_{i=1}^7 \tau_i \dot{q}_i\|_2$, is added as a regularizing loss during policy training.

3) *Teacher-Student Distillation*: During policy training, we utilize privileged information such as coefficients of friction, restitution, and the object’s inertial parameters and object velocity, which are not generally observable in the real world. As a result, the trained policy cannot operate without these quantities. To address this, we distill the trained *teacher* policy into a *student* policy that operates based on observations available in the real world. During distillation, we employ DAgger [90], where the student policy replicates the actions of the teacher solely based on the available observations through supervised learning during the simulation rollout.

As illustrated in Figure 18, the student policy shares a similar architecture with the teacher policy, consisting of a geometry encoder (red), a modulation network (green), and a base network (blue). As in the teacher, the student modulation network generates the weights for the student base network based on the outputs of the geometry encoder. The student base network then produces actions based on both the generated weights from the modulation network and the state inputs.

However, because the student policy must estimate the teacher policy’s actions with limited information, we incorporate a Gated Recurrent Unit (GRU) [16] into both the modulation network and the base network. This allows the student model to aggregate information from previous observations, helping it infer the teacher policy’s actions more effectively.

D. Details on MDP Design

The state, action, and reward components of MDP are summarized in Table VII. Hyperparameters of the reward terms are shown in Table VIII.

E. Point cloud sampling process

Our policy takes three types of point-cloud inputs: *global* scene cloud, *local* scene cloud, and *object* clouds. To obtain these inputs, we need to sample the points from the underlying meshes. For the *object* cloud, we can pre-sample the point clouds from the underlying mesh, then transform its point cloud to the current pose. However, obtaining environmental point clouds is non-trivial: since our scene is constructed dynamically, the corresponding point cloud must change across episodes. Thus, the point clouds cannot be pre-sampled, and an efficient online sampling procedure is necessary.

To sample the surface point clouds from a union of primitives, we must determine the subsection of the surfaces that form the exterior of the composed cuboids. While the simplest solution is to compute the boolean union of environment meshes, this operation is computationally costly as it cannot be parallelized. Instead, we subdivide the cuboid surfaces into a set of non-intersecting triangles, then cull the triangles that are contained by the cuboids. Afterward, we sample the point clouds proportional to the area of the remaining non-occluded triangles. By vectorizing this operation, we can efficiently sample the environmental point clouds by leveraging GPU-based acceleration.

Afterward, the *local* point cloud is sampled by selecting the points on the global cloud nearest to the object. For computational efficiency, we pre-sample a set of 64 keypoints on the object surface via FPS, then sub-sample 512 points from the global scene cloud with the lowest distance to the nearest object keypoint in the current pose.

F. Details on Baseline Architectures

UNICORN-SM uses UNICORN as the input representation (Section III-B), combined with the architecture from Soft Modularization [112]. The implementation is adapted from the author’s original code to operate with Isaac Gym [58]. For a fair comparison, we use the same number of modules and base network size as in HAMNET. The module activations are embedded with 128 dimensions, and passed through a two-layer MLP with a hidden dimension of 128 for each layer of the base network.

POINTGPT-HAMNET uses HAMNET with the PointGPT-S [11] representation model, utilizing the code and pre-trained weights released by the authors. However, since the original model is memory-intensive and computationally slow, we use TensorRT [70] to optimize the model.

UNICORN-MONO is a standard MLP using wider hidden layers (with dimensions [768, 384, 384, 192]) to match the number of trainable parameters in UNICORN-HAMNET. It uses Tanh activation and layer normalization in the hidden layers. When computing its input, we apply the same cross-attention as in HAMNET to the geometric embeddings, then concatenate the resulting embeddings with all non-geometric state inputs before feeding them into the network.

UNICORN-HYPER uses a hypernetwork [37] to output the parameters of the base network. As directly predicting all parameters requires an impossibly large hypernetwork,

TABLE VII: Summary of our MDP, in terms of state, action, and reward components. Each component is denoted by its name, shorthand symbol, dimensionality and a brief description. \dagger : only used in simulation.

State Component	Symbol	Dimension	Description
Object state \dagger	x_t^o	\mathbb{R}^{15}	Object pose and velocity
Robot state	x_t^q	\mathbb{R}^{14}	Joint positions and velocities
End-effector pose	x_t^{EE}	\mathbb{R}^9	Pose of the robot's end-effector
Physics parameters \dagger	ν	\mathbb{R}^6	Mass, friction, restitution of object and friction of robot and environment
Object geometry	G_o	$\mathbb{R}^{512 \times 3}$	Surface-sampled point cloud of the object
Environment geometry	G_e	$\mathbb{R}^{512 \times 3}$	Surface-sampled point cloud of the environment
Goal pose	T_g	\mathbb{R}^9	Target pose for the object, relative to current pose
Action Component	Symbol	Dimension	Description
Joint-space subgoal residuals	Δq	\mathbb{R}^7	Desired changes in joint positions
Proportional gains	k_p	\mathbb{R}^7	Joint-space proportional gains
Damping factors	ρ	\mathbb{R}^7	Factors for computing damping terms
Reward Component	Symbol	Dimension	Description
Task success reward	r_s	\mathbb{R}^1	Reward for task success
Goal-reaching reward	r_r	\mathbb{R}^1	Reward for moving object towards goal
Contact-inducing reward	r_c	\mathbb{R}^1	Reward for moving gripper towards object

TABLE VIII: Hyperparameters for the reward terms.

Parameter	Value	Description
λ_r	0.15	Goal-reaching reward coefficient
λ_c	0.03	Contact-inducing reward coefficient
c_g	3.0	Scale for goal-reaching distance potential
c_r	3.0	Scale for contact-inducing distance potential

we instead design UNICORN-HYPER to predict *low-rank* decompositions of base network parameters. For each layer of the base network of form $x_i = \sigma(Wx_{i-1} + b)$ where i denotes layer index, the hypernetwork ϕ outputs $\phi(z) = \{W_l, W_r, b\}$ and constructs $W = W_l W_r$ where $W \in \mathbb{R}^{N \times M}$, $W_l \in \mathbb{R}^{N \times k}$, and $W_r \in \mathbb{R}^{k \times M}$, where k denotes the rank, thus reducing the output dimensions from $N \times M + M$ to $(N + M) \times k + M$. In all our experiments, we configure the rank to be 16, and the hypernetwork is an MLP with hidden dimensions [256, 256] using Tanh activation and LayerNorm in the interim layers.

UNICORN-TRANSFORMER is a four-layer transformer, where each layer uses the embedding dimension of 512 and four attention heads. For computational efficiency, we leverage FlashAttention [19] in our implementation. The transformer receives learnable action and value tokens, geometric embeddings (generated by the same cross-attention used in HAMNET), and non-geometric state inputs tokenized with a linear layer. After passing these tokens through the transformer layers, a two-layer MLP with a 256-dimensional hidden layer maps the action and value token embeddings to the robot's action and state value, respectively. Detailed hyperparameters for the network architectures and PPO training are described in Table IX and Table X, respectively.

G. Ablation studies

1) *Effects of the gating mechanism:* To assess the benefits of the gating mechanism in HAMNET, we compare with UNICORN-HAMNET-WITHOUT-GATE, which omits the gating mechanism but is otherwise identical to UNICORN-HAMNET. Figure 19 illustrates the training progression for

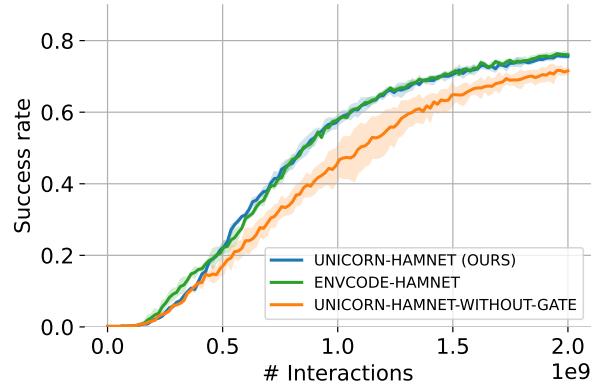


Fig. 19: Training progression for our ablation. Plots show the mean (solid) and standard deviations (transparent) for each baseline across three seeds. The interaction steps are reported as the total number of steps aggregated across 1024 parallel environments.

both. While initial trends are similar, UNICORN-HAMNET-WITHOUT-GATE trains slower and reaches lower final performance than UNICORN-HAMNET (71.6% vs. 75.6%). This result highlights the effectiveness of complementing the modules with the gating mechanism to afford additional expressivity.

2) *Representational quality of UNICORN :* To evaluate whether UNICORN sufficiently encodes geometric information about the environment and the object, we compare UNICORN-HAMNET with ENVCOde-HAMNET. This baseline uses the same architecture as UNICORN-HAMNET, but employs a hand-engineered oracle representation, ENVCOde. This is constructed by concatenating the procedural generation parameters as in Table XI, yielding a unique 25-dimensional real-valued vector for the scene representation.

As illustrated in Figure 19, UNICORN-HAMNET performs on par with ENVCOde-HAMNET, achieving a success rate of 75.6% compared to 76.0%. This indicates that UNICORN, despite operating from sensory observations, provides sufficient

TABLE IX: Network Hyperparameters.

Hyperparameter	Value	Hyperparameter	Value	Hyperparameter	Value
Num. points	512	Num. encoder layers	4	Modulation Network	MLP (256, 256)
Num. patches	16	Num. self-attn heads	4	Actor	MLP (256, 128, 128, 64)
Patch size	32	Cross-attn embedding dim.	64	Critic	MLP (256, 128, 128, 64)
Embedding dim.	128	Num. cross-attn heads (object / others)	8/4	Num. modules	8

TABLE X: PPO Hyperparameters.

Hyperparameter	Value	Hyperparameter	Value
Max Num. epoch	8	Base learning rate	0.0003
Early-stopping KL target	0.024	Adaptive-LR KL target	0.016
Entropy regularization	0	Learning rate schedule	KL-adaptive
Initial log std.	-0.4	log std. decay factor	-0.000367
Policy loss coeff.	2	Value loss coeff.	0.5
GAE parameter	0.95	Num. environment	1024
Discount factor	0.99	Episode length	300
PPO clip range	0.3	Update frequency	8
Bound loss coeff.	0.02	Energy loss coeff.	8e-5

TABLE XI: Content of ENVCODE.

Parameter	Dimensions	Description
ramp position	$\mathbb{R}^{2 \times 2}$	Position of each ramp
ramp slope	$\mathbb{R}^{2 \times 2}$	Angle of each ramp
plate elevations	$\mathbb{R}^{2 \times 3}$	Height of each base plate
wall heights	\mathbb{R}^4	Height of each wall
ceiling height	\mathbb{R}^1	Height of the ceiling
scene dimension	\mathbb{R}^3	Overall scene dimension
scene position	\mathbb{R}^3	Overall scene position

information about the environment, matching the performance of the hand-engineered *oracle* representation taken directly from the parameters of the underlying procedural generation pipeline.

3) *Details on the Parameter Scaling Experiment:* In our parameter scaling experiment (see Figure 4), we consider three baselines: HAMNET, MLP, and TRANSFORMER. In HAMNET, we omit the gating to isolate the effects of the modules. This ensures that MLP and HAMNET reduce to the same network architecture when the number of module is one, so that both baselines share the starting point. As we increase the number of modules, we keep the size of the base network the same in HAMNET, while increasing the width of the hidden layers in the MLP baseline to approximately match the parameter counts, as shown in Table XII. The TRANSFORMER model is considerably larger, with 8 transformer layers, each with the embedding dimension of 512 and 4 heads. All results are aggregated across three different seeds.

4) *Details on the Simulation Benchmark:* In our benchmark, we provide nine digital twins of real-world environments paired with a total of 353 objects in our benchmark, shown in Figure 21: nine custom-scanned objects from the real world (red box), 21 from GSO [24] (green box), and 323 from DGN [106] (blue box).

For the custom-scanned objects, stable poses are collected in the real world using FoundationPose [108]. For the GSO ob-

TABLE XII: MLP baseline configurations scaled to match the parameter count of the HAMNET network for each module count. Network size denotes the dimensions of hidden layers.

Corresponding Num. modules	Network size	Num. params
1	[256, 128, 128, 64]	0.36 M
2	[304, 144, 144, 64]	0.43 M
4	[512, 256, 256, 128]	0.94 M
8	[768, 384, 384, 192]	1.76 M

jects, we use Trimesh [21] to sample their stable orientations. For DGN objects, we use the precomputed stable orientations (see Appendix B2). To determine stable placements for GSO and DGN objects, each object is rotated to one of the pre-sampled stable orientations and randomly placed in a predefined region of the environment with a small vertical margin (0.005m). We resolve remaining collisions with FCL [74] to ensure that the object does not penetrate the environment.

From these stable poses, we randomly select pairs of initial and target object poses that maintain sufficient separation in both position and orientation to prevent trivial scenarios. For each domain-object pair, we sample 128 distinct episodes, each with different collision-free robot initializations within the robot’s joint limits, as shown in Figure 20.

To support the evaluation of new algorithms, we establish baseline results from representative monolithic and modular architectures: UNICORN-MONO and UNICORN-HAMNET, which are illustrated in Figure 22.

H. Details on the Real-World Setup

1) *Pose tracker implementation :* While the original FoundationPose [108] model only processes one image at a time, we operate with four cameras. To streamline computation, we batchify the model to allow multiple images to be processed at once. For additional robustness, we also generate multiple pose candidates by adding a small noise (0.02m, 0.15 radians). Each candidate undergoes the refinement procedure as in the original model, then we select the pose with the highest prediction score as input to the policy.

2) *Failure Modes in the Real World:* We describe our five main failure modes in the real world. In *torque-limit violation*, the robot aborts due to the robot exceeding the hardware’s safety limits. Despite the measures taken in Appendix C, the sim-to-real gap may still lead to spurious contact in domains with walls, such as the *drawer*, when the robot makes rapid movements to adjust contact sites. The policy may also get

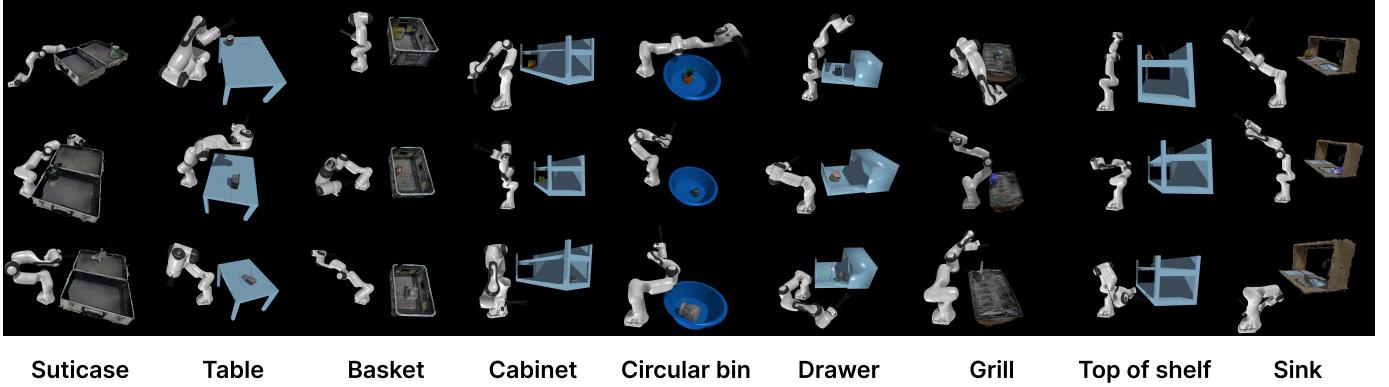


Fig. 20: Example configurations for each domain. Each row depicts a different object, and each column shows a different environment.

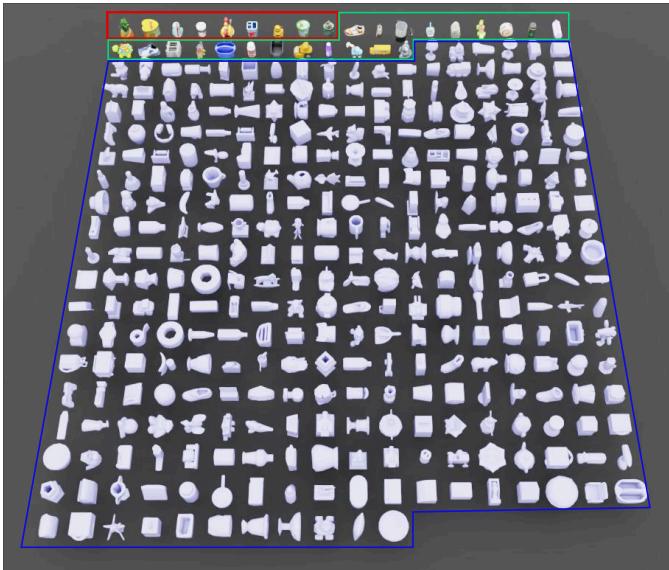


Fig. 21: Object meshes used for the benchmark (353 in total). The red box indicates nine custom-scanned objects from the real world; the green box contains 21 objects from GSO; and the blue boxes enclose 323 objects from DGN.

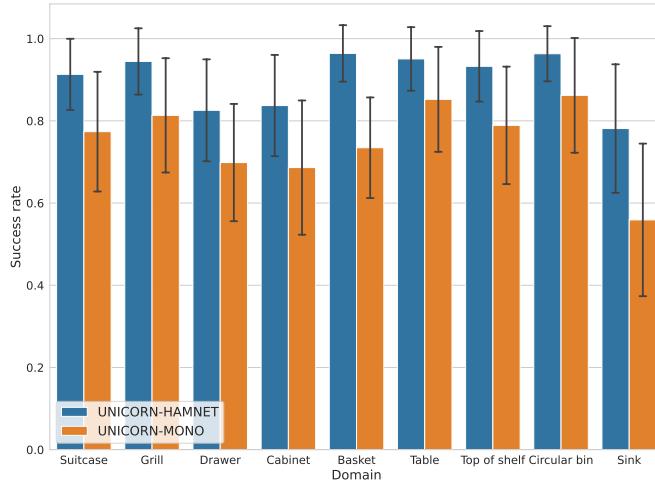


Fig. 22: Comparison of success rates of UNICORN-HAMNET (blue) and UNICORN-MONO (orange) on our simulated benchmark over 10000 simulation steps.

stuck in a *deadlock*, where it indefinitely repeats ineffective maneuvers. For instance, the policy may keep attempting a toppling maneuver for low-friction objects such as *heart-box*, which may not work due to slippage. In other cases, the robot may accidentally *drop* the object. For example, the object such as the *angled cup* may rapidly bounce or roll off the scene, and the limited dexterity of our hardware prevents catching such fast-moving objects. Another failure mode arises when the agent fails to *circumnavigate obstructions*, getting stuck against the environment. This also arises from the sim-to-real gap: while the simulation often allows the robot to move across a shallow barrier by pressing against it, real-world walls cannot be penetrated, which causes the robot to get stuck against the environment. Lastly, the remaining failures occur from the *perception stack*, where it loses track of the segmentation mask or the object pose. This most frequently occurs when key recognizable textures of the object are occluded by the robot or the environment.