# Galactic: Scaling End-to-End Reinforcement Learning for Rearrangement at 100k Steps-Per-Second

Vincent-Pierre Berges*
Meta AI (FAIR)

Andrew Szot*
Georgia Tech

Devendra Singh Chaplot
Meta AI (FAIR)

Aaron Gokaslan
Cornell University

Roozbeh Mottaghi
Meta AI (FAIR)

Dhruv Batra
Meta AI (FAIR), Georgia Tech

Eric Undersander
Meta AI (FAIR)

## Abstract

*We present Galactic, a large-scale simulation and reinforcement-learning (RL) framework for robotic mobile manipulation in indoor environments. Specifically, a Fetch robot (equipped with a mobile base, 7DoF arm, RGBD camera, egomotion, and onboard sensing) is spawned in a home environment and asked to rearrange objects – by navigating to an object, picking it up, navigating to a target location, and then placing the object at the target location.*

*Galactic is fast. In terms of simulation speed (rendering + physics), Galactic achieves **over 421,000 steps-per-second (SPS)** on an 8-GPU node, which is 54x faster than Habitat 2.0 [55] (7699 SPS). More importantly, Galactic was designed to optimize the entire rendering+physics+RL interplay since any bottleneck in the interplay slows down training. In terms of simulation+RL speed (rendering + physics + inference + learning), Galactic achieves **over 108,000 SPS, which 88x faster than Habitat 2.0 (1243 SPS).***

*These massive speed-ups not only drastically cut the wall-clock training time of existing experiments, but also unlock an unprecedented scale of new experiments. First, Galactic can train a mobile pick skill to $>80\%$ accuracy in under 16 minutes, a 100x speedup compared to the over 24 hours it takes to train the same skill in Habitat 2.0. Second, we use Galactic to perform the largest-scale experiment to date for rearrangement using 5B steps of experience in 46 hours, which is equivalent to 20 years of robot experience. This scaling results in a single neural network composed of task-agnostic components achieving 85% success in GeometricGoal rearrangement, compared to 0% success reported in Habitat 2.0 for the same approach. The code is available at github.com/facebookresearch/galactic .*

## 1. Introduction

The scaling hypothesis posits that as general-purpose neural architectures are scaled to larger model sizes and training experience ever increasingly sophisticated intelligent behavior

*Equal contribution

emerges. These so-called 'scaling laws' appear to be driving many recent advances in AI, leading to massive improvements in computer vision [12, 42, 44] and natural language processing [3, 40]. But what about embodied AI? We contend that embodied AI experiments need to be scaled by *several orders of magnitude* to become comparable to the experiment scales of CV and NLP, and likely even further beyond given the multimodal interactive long-horizon nature of embodied AI problems.

Consider one of the largest-scale experiments in vision-and-language: the CLIP [39] model was trained on a dataset of 400 million images (and captions) for 32 epochs, giving a total of approximately 12 Billion frames seen during training. In contrast, most navigation experiments in embodied AI involve only 100-500M frames of experience [5, 30, 68]. The value of large scale training in embodied AI was demonstrated by Wijmans *et al.* [60] by achieving near-perfect performance on the PointNav task through scaling to 2.5 billion steps of experience. Since then, there have been several other examples of scaling experiments to 1B steps in navigation tasks [41, 65]. Curiously, as problems become more challenging, going from navigation to mobile manipulation object rearrangement, the scale of experiments have become smaller. Examples of training scales in rearrangement policies include [19, 55, 59] training skills in Habitat 2.0 for 100-200M steps and [58] in Visual Room Rearrangement for 75M steps. In this work, we demonstrate for the first time scaling training to 5 billion frames of experience for rearrangement in visually challenging environments.

Why is large-scale learning in embodied AI hard? Unlike in CV or NLP, data in embodied AI is collected through an agent *acting* in environments. This data collection process involves policy inference to compute actions, physics to update the world state, rendering to compute agent observations, and reinforcement learning (RL) to learn from the collected data. These separate systems for rendering, physics, and inference are largely absent in CV and NLP.

We present Galactic, a large-scale simulation+RL framework for robotic mobile manipulation in indoor environments. Specifically, we study and simulate the task of GeometricGoal

| Arcade RL Sims | Device | Res | Sensors | Train SPS | Sim SPS | Photoreal | Physics |
|---|---|---|---|---|---|---|---|
| VizDoom [26, 37] | 1x RTX 3090 | 28×72 | RGB | 18,900 | 38,100 | ✘ | ✓ |
| **Physics-only Sims** | **Device** | **Res** | **Sensors** | **Train SPS** | **Sim SPS** | **Photoreal** | **Physics** |
| Isaac Gym (Shadow Hand) [29] | 1x A100 | N/A | N/A | 150,000 | – | ✘ | ✔ |
| Brax (Grasp) [16] | 4x2 TPU v3 | N/A | N/A | 1,000,000 | 10,000,000 | ✘ | ✔ |
| ADPL Humanoid [64] | 1x TITAN X | N/A | N/A | 40,960 | 144,035 | ✘ | ✔ |
| **EAI Sims** | **Device** | **Res.** | **Sensors** | **Train SPS** | **Sim SPS** | **Photoreal** | **Physics** |
| iGibson [28, 51] | 1x GPU | 128×128 | RGB | – | 100 | ✔ | ✔ |
| AI2-THOR [11] | 8x RTX 2080TI | 224x224 | RGB | ≈300 | 2,860 | ✔ | ✔ |
| Megaverse [37] | 1x RTX 3090 | 128×72 | RGB | 42,700 | 327,000 | ✘ | ✓ |
| | 8x RTX 2080TI | 128×72 | RGB | 134,000 | 1,148,000 | | |
| LBS [48] | 1x RTX 3090 | 64×64 | RGB | 13,300 | 33,700 | ✔ | ✘ |
| | 1x Tesla V100 | 64×64 | RGB | 9,000 | – | | |
| | 8x Tesla V100 | 64×64 | RGB | 37,800 | – | | |
| Habitat 2.0 [55, 59] | 1x RTX 2080 Ti | 128×128 | RGBD | 367 | 1,660 | ✔ | ✔ |
| | 8x RTX 2080 Ti | 128×128 | RGBD | 1,243 | 7,699 | | |
| | 1x Tesla V100 | 128×128 | RGBD | 128 | 2,790 | | |
| | 8x Tesla V100 | 128×128 | RGBD | 945 | 17,465 | | |
| **Galactic (Ours)** | 1x Tesla V100 | 128×128 | RGBD | 14,807 | 54,966 | ✔ | ✔ |
| | 8x Tesla V100 | 128×128 | RGBD | 108,806 | 421,194 | | |

Table 1. High-level throughput comparison of different simulators. Steps-per-second (SPS) numbers are taken from source publications, and we don't control for all performance-critical variables including scene complexity and policy architecture. Comparisons should focus on orders of magnitude. We show Sim SPS (physics and/or rendering) and training SPS (physics and/or rendering, inference and learning) for various physics-only and Embodied AI simulators. We also describe VizDoom, an arcade simulator which has served as a classic benchmarks for RL algorithms due to its speed. The ✓ for Megaverse and VizDoom represent physics for abstract, non-realistic environments. Among EAI simulators that support realistic environments (photorealism and realistic physics), Galactic is $80\times$ faster than the existing fastest simulator, Habitat 2.0 (108,806 vs 1243 training SPS for 8 GPUs). Galactic's training speed is comparable to LBS, Megaverse, and VizDoom, even though LBS doesn't simulate physics and neither Megaverse nor VizDoom support realistic environments. We also compare to GPU-based physics simulators: while these are generally faster than Galactic, they entirely omit rendering, which significantly reduces their compute requirements. For Galactic, we observe near-linear scaling from 1 to 8 GPUs, with a 7.3x speedup.

Rearrangement [2], where a Fetch robot [43] equipped with a mobile base, 7DoF arm, RGBD camera, egomotion, and proprioceptive sensing must rearrange objects in the ReplicaCAD [55] environments by navigating to an object, picking up the object, navigating to a target location, and then placing the object at the target location.

Galactic is fast. In terms of simulation speed (rendering + physics), Galactic achieves *over 421,000 steps-per-second (SPS)* on an 8-GPU node, which is 54x faster than Habitat 2.0 [55] (7699 SPS). More importantly, Galactic was designed to optimize the entire rendering+physics+RL interplay since any bottleneck in the interplay slows down training. In terms of simulation+RL speed (rendering + physics + inference + learning), Galactic achieves *over 108,000 SPS*, which 88x faster than Habitat 2.0 (1243 SPS).

Our key technical innovations are: (1) integration of CPU-based batch physics with GPU-heavy batch rendering and inference, and (2) a new, approximate kinematic simulation targeted at EAI rearrangement tasks. Compared to a "one simulator, one environment, one process" paradigm, batching yields massive speedups due to memory savings, lower communication overhead, and greater parallelism. Meanwhile, we leverage our physics approximations and the reduced

complexity of kinematic simulation to reduce our CPU compute and yield further speedups.

These massive speed-ups not only drastically cut the wall-clock training time of existing experiments, but also unlock an unprecedented scale of new experiments. First, Galactic can train a mobile pick skill to $>80\%$ accuracy in under 16 minutes, a 100x speedup compared to the over 24 hours it takes to train the same skill in Habitat 2.0. Second, we use Galactic to perform the largest-scale experiment to date for rearrangement using 5B steps of experience in 46 hours, which is equivalent to 20 years of robot experience (assuming 8 actions per second). This scaling results in a single neural network composed of task-agnostic components (CNNs and LSTMs) achieving $85\%$ success in GeometricGoal rearrangement. This is impressive performance because (1) the task is extremely long horizon (involving navigation, picking, and placing), (2) the architecture is monolithic and has no mapping modules, task-planning, or motion-planning. For context, Habitat 2.0 reported $0\%$ success with a monolithic RL baseline in GeometricGoal rearrangement. We find the learned policies are able to efficiently navigate, avoid distractor objects, and synchronize base and arm movement for greater efficiency. Finally, we also show that models trained in Galactic are somewhat robust to zero-shot sim2sim generalization, *i.e.*

can achieve 26% success when deployed in Habitat 2.0 despite differences in rendering, physics, and underlying controller.

## 2. Related Work

**Scaling Approaches in Embodied AI.** There is a large body of work in speeding up training embodied agents in simulation. There are three general approaches for increasing efficiency of the overall system: distributing the policy training and inference, increasing the sample efficiency of the learning algorithms, and batch simulation.

*Distribution and parallelization:* The works in this domain [15, 24, 35, 54, 60, 67] achieve efficiency by distributing training across multiple GPUs or nodes and parallelization of the computation. Galactic's systems contributions are targeted at rollout computation (inference and environment-stepping) on a single CPU process and single GPU, so it is complementary to many of these approaches [59, 60, 67].

*Training algorithms:* Various techniques have been developed for the efficiency of training algorithms in interactive settings. Using auxiliary losses [25, 31, 36, 50, 66, 68], offline training [32, 33, 46, 49, 69, 70], and model-based training [6, 9, 20–22, 34] are some examples that lead to sample efficiency (and typically wall clock time) of the training algorithms. Similarly, since our experiments use a single DD-PPO [60] policy trained from scratch, Galactic can be combined with any of these techniques for further efficiency.

*Batch simulation:* Batch simulation refers to vectorized physics or rendering to compute updates across multiple environments with one operation in a batched fashion. This yields large speedups compared to the "one process, one simulator, one environment" paradigm. Our approach is closest to LBS [48] and Megaverse [37], but ours is the first work to combine batch physics and batch rendering to simulate realistic environments and vision sensors. [48] does not support physics, and instead only supports photorealistic, non-interactive scenes and cylinder agents. Megaverse [37] does not support physics simulation with articulated agents and realistic movable objects, and instead only supports "block worlds" with movable blocks and cylinder agents.

**Embodied AI Simulators.** There are various Embodied AI simulation platforms [10, 17, 27, 38, 45, 51, 55, 62, 63] for indoor environments that support tasks such as navigation [45, 61], object manipulation [14, 53], instruction following [1, 52], interactive question answering [8, 18], and object rearrangement [55, 58]. The efficiency of the simulators is important for these tasks since they typically involve long task horizons, and the state-of-the-art training algorithms require millions of iterations to converge. Galactic supports similar tasks in indoor environments.

**Kinematic Simulation.** Some recent work uses "kinematic simulation", in which the robot and objects are moved directly without simulating rigid body dynamics. [55] uses a "sticky mitten" abstraction for grasping rather than simulating contact physics. It also uses kinematic movement for the base rather than simulating base momentum and wheel forces. [57] explores training a navigation policy for a quadruped robot in simulation. They show that training with kinematic movement for the base transfers better to real compared to training with full quadruped dynamics. Both of these works use the Habitat 2.0 simulator with underlying CPU-based Bullet physics [7]. For Galactic, we introduce a new kinematic simulator optimized specifically for EAI rearrangement tasks.

## 3. Galactic System

From an ML systems perspective, RL training can be broken down into rollout computation (collecting experience by interacting with the EAI simulator) and learning (updating the policy). Our systems contributions focus on speeding up the experience collection. In particular, we optimize batch rollout computation for a single GPU and associated CPU process by (1) integrating CPU-based batch physics with GPU-based batch rendering and DNN inference (Section 3.1), and (2) introducing a new approximate kinematic simulation optimized for EAI rearrangement tasks (Section 3.2). For our RL experiments, we combine fast rollouts with DD-PPO [60] (Section 4.2), but other approaches to distributed RL are also compatible [59, 67]. We discuss throughput and scaling in Section 5.1. Galactic also supports adding assets from different sources (see Appendix H for details).

### 3.1. Batching

Consider non-batched EAI simulators like [27, 55]. They use a "one Python process, one simulator, one environment" paradigm. The simulator manages a single environment (including physics and rendering) and produces a single set of observations. Additionally, some task-specific logic is implemented directly in Python code, e.g. extracting a robot's pose from the sim and computing reward and other metrics.

Python's cooperative threading model and global interpreter lock make parallelism difficult. To scale this RL training to multiple environments, these approaches typically spawn multiple Python "env" processes, each hosting its own simulator instance. In-memory assets here include CPU data, such as physics collision geometry, and GPU data, such as meshes and textures, and unfortunately they must be duplicated across instances, not shared. Coordinating rollouts and gathering results requires significant interprocess communication, e.g. sending observations between the "env" and "main" processes. Observations are batched on the main process and fed to GPU-based batch DNN inference. The number of parallel environments is limited by total system memory, interprocess communication overhead, and other factors and is typically between 16 and 28 [27, 55].

Prior work has shown the benefit of batch simulation over this non-batched paradigm, both for CPU-based physics simulation [37], GPU-based physics simulation [16, 29, 64], and GPU-based rendering [37, 48]. For Galactic, we batch both CPU-based physics and GPU-based rendering. For physics,

the main Python process hosts a single C++ simulator instance. This instance steps physics for a batch of environments, sharing in-memory CPU assets across environments. Unlike [48] which moves nearly all task logic to C++, we retain the flexibility of Python for reward and other task-specific computation. This is implemented as efficient Numpy tensor operations, in contrast to the non-tensor Python code in the non-batched paradigm. We use the Python buffer protocol to ensure we have zero-copy conversion between C++ and Python.

For rendering, we use Bps3D [48]. A single Bps3D renderer instance renders all environments, sharing GPU memory assets across environments. Scene graph updates are communicated from our physics simulation to the Bps3D renderer efficiently in C++. The renderer outputs a single batch observation per camera sensor, essentially a stack of images from each environment in GPU memory. This is directly consumed by PyTorch GPU-based batch DNN inference, without ever copying pixels to the CPU. We visualize example observations from Galactic in Appendix D.

Figure 1 shows our integration of batch physics (orange), batch rendering (purple), and PyTorch DNN inference (blue) into the rollout computation loop. "Step post-processing" refers to reward calculation and other task-specific logic. CPU-based physics is computed in parallel with GPU-heavy rendering and inference. Because of this interleaving, we must accept a one-step-delay: as in [55], our policy's actions are computed not from the current step's observations, $o_t$, but rather those from the previous step, $o_{t-1}$. Our approximate kinematic sim is fast enough such that we are "GPU-bound": the bottleneck here is primarily the GPU, with large gaps on both the main and physics threads corresponding to idle CPU. This GPU-bound property is desirable and means Galactic will benefit greatly as faster GPUs become available.

Compared to the non-batched paradigm, Galactic rollouts have no interprocess communication overhead because all CPU compute happens in a single process. Because of memory savings from batching, our number of environments is larger than the non-batched paradigm (128 versus 16 to 32).
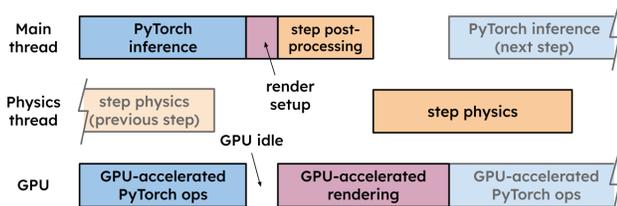


Figure 1. Rollout timeline for a single batch step, showing our integration of batch physics (orange), batch rendering (purple), and PyTorch DNN inference (blue). CPU-based physics is computed in parallel with GPU-heavy rendering and inference. Physics is fast enough for rollout computation to be primarily GPU-bound.

## 3.2. Approximate Kinematic Simulation

Batra et al. [2] reviews physics realism for EAI simulators, including agent embodiment and how it interacts with the scene. They describe a spectrum: at the most abstract end are simple cylinder embodiments, "magic pointer" grasping, and "virtual backpacks" [27, 45]. At the other end is full rigid-body dynamics simulation: [16, 29, 64].

Galactic's approximate kinematic simulation lies somewhere in the middle. It is kinematic: the robot and objects are moved directly without simulating rigid body dynamics. It uses approximations for detecting and resolving collisions and for simulating object-dropping. Regarding these approximation choices, an important design goal is feasible sim-to-real for policies trained in Galactic: our policies should produce detailed, physically-plausible trajectories for an articulated robot and movable objects, requiring only the addition of simple, low-level controllers on real hardware. As a proxy for sim-to-real in this work, we explore sim-to-sim transfer to the Habitat 2.0 dynamic simulator in Section 5.2.

Recent work has employed similar kinematic simulation [55, 57], but these are built on top of Habitat 2.0 and the existing CPU-based Bullet physics engine [7]. Meanwhile, we build a new simulator from scratch, leveraging both our physics approximations and the reduced complexity of kinematic simulation to reduce our compute.

Our action space for the articulated robot includes offsets for base forward/back, base rotation, and the other degrees of freedom, e.g. arm joint rotation. Inside the simulator, these offsets are applied directly to the robot state, then we use forward kinematics to compute the position of the robot's articulated links. This is kinematic movement: no velocities, forces, or momentum are simulated. This new candidate pose for the robot is tested for overlap (penetration) with the environment. A "collision" is resolved by either disallowing all robot movement for that step or by sliding the robot base ("allow sliding" is a training hyperparameter). Our sliding implementation is approximate and uses a heuristic search described in Appendix I.1.

Collision geometry is generally used in a physics simulator to perform contact and overlap tests. Common primitives include convex hulls and triangle meshes, which can accurately represent real-world shapes using sufficiently high vertex/triangle counts. For our rearrangement task, we assume sub-centimeter tolerances are not important, so we use alternate primitives which are less accurate but faster to query. In particular, we approximate the robot (and the current grasped object, if any) with a set of spheres (green and blue in Figure 3). We approximate movable objects with oriented bounding boxes (orange). Finally, we approximate the static (non-movable) parts of the environment with a voxel-like structure called a column grid (gray). It has limited precision in the lateral (XZ) direction (3 cm) but full floating-point precision in the vertical (Y) direction, so surface heights are accurately represented. Spheres are authored for the robot manually while bounding

Figure 2. Overview of the GeometricGoal rearrangement task. In this task, a Fetch robot must move an object from a start position, specified as a 3D coordinate, to a goal position, also specified as a 3D coordinate, all from egocentric sensing. Figure adapted from [56].

boxes and column grids are generated automatically from high-fidelity source meshes in a preprocessing step.

Sphere-versus-environment overlap tests are optimized in a couple ways: (1) we precompute a distinct column grid for each unique sphere radius, so sphere-versus-column-grid tests are simple lookups. (2) Resting movable objects are inserted into a regular-grid acceleration structure. Testing against the entire set of movable objects is fast because we need only iterate over the precomputed set of nearby (50 cm) objects.

Our simulated grasping approximates a real-world suction gripper kinematically without simulating contact or suction forces. Firstly, our action space includes a discrete grasp/release action. In the simulator, if the grasp action is active, we query a small (3 cm) sphere at the tip of the end effector and check for overlap (contact) with movable objects in the scene. If an overlapping object is found, the object is fixed to the end effector. The object moves with the end effector until the release action is performed.

On release, our "snap-to-surface" operation approximates the entire drop sequence for a released object: falling, landing, and settling. We use a sphere-cast query to find the nearest surface below the object. Instead of allowing movable objects to stack on top of each other, we use an additional heuristic search described in Appendix I.2 to find a collision-free resting position nearby.

We choose to implement our simulator in C++; it runs on the CPU, in contrast to recent GPU-based physics simulators [16, 29, 64]. While these offer greater throughput, CPU-based C++ code has advantages over GPU code: direct access to host memory/disk/network resources, a less restrictive memory model, and greater availability of third-party libraries. In addition, as shown in Figure 1, our kinematic sim ("step physics") is faster than GPU-based rendering and PyTorch inference. This makes it "free" within our GPU-bound rollout computation, as shown by the large idle gaps on the Physics thread. This is in contrast to GPU-based physics simulators, which despite their high throughput, would not be free in this GPU-bound scenario, which would negatively impact overall rollout compute time.
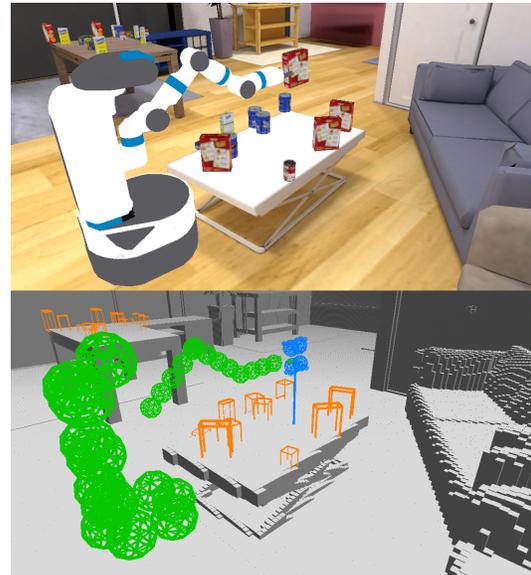


Figure 3. Visualization of an environment (top) and its collision geometry (bottom). In the bottom, the robot is represented with green spheres, and the held object with blue spheres. The other movable objects as oriented boxes (orange), and the static environment as a voxel-like "column grid" data structure (gray).

## 4. Experiment Setup: Object Rearrangement

### 4.1. Task Description

We study how Galactic can accelerate learning in *Geometric-Goal object rearrangement* [2] with end-to-end RL. We follow the GeometricGoal rearrangement task setup from [55], where a robot is tasked with moving an object from a start location to the desired goal entirely from onboard sensing consisting of an RGBD head camera and proprioceptive sensing. The task is specified via the 3D center of mass for the target object's start location and the 3D center of mass for the object's target location. To accomplish the task, the agent must navigate through an indoor environment, pick, place, and avoid distractor objects and clutter, all while operating from egocentric sensors. The agent is a simulated Fetch robot [43], with a wheeled base, a 7-DoF arm manipulator, a suction gripper, and a head-mounted RGBD camera ($128 \times 128$ pixels). The episode is successful if

the agent places the object within 15cm of the desired position within the episode horizon of 500 time steps and calls the stop action. We aim to transfer the policies learned in Galactic to the Rearrange Easy benchmark from the NeurIPS 2022 Habitat Rearrangement Challenge [55, 56]. We implement the same rearrangement task in Galactic to facilitate this transfer.

We use an 11-dimension action space in Galactic consisting of control for the base, arm, gripper, and episode termination. 2 actions control the linear and angular velocity of the robot base. 7 actions are for delta joint angles for each of the 7 arm joints. 1 binary action is for the grasping and 1 binary action indicates episode termination.

The observation space consists of visual sensors, proprioceptive sensing, and task specification. The visual sensor is a $128 \times 128$ RGBD, $90°$ FoV head-mounted camera. Proprioceptive sensing provides the angles in radians of all 7 joints and if the robot is holding an object. From a base egomotion sensor, we derive the relative position of the robot base and end-effector to the object start and goal position as observations. Additionally, we include an episode step counter in the observation.

We import the datasets and assets from the Rearrange Easy benchmark in Habitat 2.0 into Galactic. The Rearrange Easy benchmark uses the ReplicaCAD [55] scene dataset consisting of 105 interactive indoor home spaces. We load the train dataset of 10k episodes which specify rearrangement episodes in the train split of 63 room layouts in ReplicaCAD. We pre-compute the robot's starting position in each episode for greater efficiency. Each episode contains one target object to rearrange and 29 distractor objects.

### 4.2. Approach: End-to-End RL

Our approach relies on training a "sensors-to-actions" policy directly using RL with the task reward alone. Previous state-of-the-art approaches for mobile manipulation in Habitat 2.0 rely on decomposing the task into the separate skills of navigation, picking, and placing [19, 59]. However, such hierarchical approaches suffer from requiring a hand-specified task decomposition and "hand-off problems" where errors between skills compound. Our approach avoids these issues by training a single end-to-end policy via RL. Previously, such end-to-end approaches achieved little success in the Habitat 2.0 rearrangement tasks [55]. However, with the speed of the Galactic simulator, we can generate experience fast enough for the end-to-end approach to become a viable approach to rearrangement tasks since it allows gathering the billions of samples needed in only a few days with an 8-GPU compute node.

**Reward function.** We define the reward for the rearrangement task largely following the default reward in the Rearrange Easy benchmark, but with some modifications for smooth robot motion. The agent gets a sparse reward for completing the task and picking up the object. The robot is given a dense reward as the decrease in L2 distance between the end-effector and object as well as between the object and goal. The robot is

given a penalty for large differences in actions at subsequent steps. Furthermore, to speed up training convergence, we do not allow the robot to call the stop or drop action until the end-effector is within a cylinder of width 0.15m and height 0.3m around the goal position. We analyze this decision in detail in Section 5.3. Additional details of the rearrangement task and reward function in Galactic are in Appendix A.

**Architecture and Hyperparameters.** Our primary experiments use a ResNet18 [23] visual encoder, with a 2-layer 512 hidden unit LSTM, and then separate actor and critic network heads. We then train this with DD-PPO [60], a distributed version of PPO [47]. We run 128 environments per GPU across 8 GPUs giving 1024 environment instances in total. We use a policy rollout length of 64, 2 mini-batches per update, and 1 epoch over the rollout data per update. This setup runs at over 30,000 steps-per-second (SPS), where the policy contains 6 million trainable parameters, We provide all training hyperparameters in Appendix B.

## 5. Results

### 5.1. Throughput and Scaling

In Table 1, we report sim steps per second (physics and/or rendering) and training SPS (physics and/or rendering, plus inference and learning) for various physics-only and Embodied AI simulators. We also describe VizDoom, an arcade simulator that has served as a classic benchmark for RL algorithms due to its speed.

Among EAI simulators that support realistic environments (photorealism and realistic physics), Galactic is $80\times$ faster than the existing fastest simulator, Habitat 2.0 (108,806 vs 1243 training SPS for 8 GPUs). Galactic's training throughput is comparable to LBS, Megaverse, and VizDoom, even though LBS doesn't simulate physics, and neither Megaverse nor VizDoom supports realistic environments. We also compare to GPU-based physics-only simulators: while these are generally faster than Galactic, they entirely omit rendering, which significantly reduces the compute requirements of environment-stepping, inference, and training compared to visual simulators and vision-based policy training. As seen in the last two rows, for our distributed training, we observe near-linear scaling from 1 to 8 GPUs, achieving a 7.3x speedup.

In Figure 4, we show how the training SPS scales as a function of the number of environments (batch size) and visual encoder on a rearrangement task (described in detail in Section 4). SPS increases for simpler encoders and larger batch sizes.

### 5.2. Sim-to-Sim Results

In this section, we show that we can scale RL training in Galactic, achieve a high success rate in Galactic, and then zero-shot transfer the policy to Habitat 2.0.

**Training in Galactic.** First, we train a policy in Galactic using the end-to-end training setup and policy architecture described in Section 4.2. We leverage the fast simulation and easy scaling
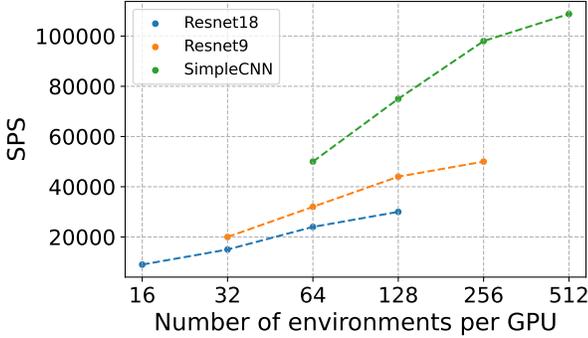
Figure 4. The steps-per-second throughput when varying visual encoder architecture and the number of environments per GPUs, using a single compute node and 8 GPUs. The rightmost point on each curve corresponds to the maximum number of environments we can use before running out of GPU memory.

of Galactic to train a neural policy with RL for 5 billion steps of experience. With 1024 environments across 8 GPUs, training runs at 30,000 SPS which takes 46 hours to train the policy to convergence at 5 billion steps.

**Transfer to Habitat 2.0.** Next, we transfer the policy trained in Galactic zero-shot to Habitat 2.0. The physics simulation in Galactic is kinematic, while the physics in Habitat 2.0 is dynamic. In the kinematic simulation of Galactic, the actions denote joint delta angles and base velocity, which are set by the simulator on the next step. However, in Habitat 2.0, the arm is controlled via joint motor torques at every step to achieve a target joint state. Hence, we bridge the dynamics simulation gap by using the policy trained in Galactic as a higher-level policy that outputs target joint offset angles in Habitat 2.0 at 10Hz. A controller then outputs the joint motor torques to achieve the target joint state set by the policy trained in Galactic operating at 120Hz. The policy takes as input the observations from Habitat 2.0 to output this target delta state. Observations between Galactic and Habitat 2.0 are similar since they use the same scenes and assets.

Other small differences in the simulators affect the sim-to-sim performance. Since Galactic does not simulate dynamics, the dropping mechanism "snaps" the object to the surface directly below the drop point of the object. This ignores the effects of the object rolling or bouncing off from the point of contact after the drop. However, as described in Section 4.2, the policy is trained to drop the object just above the surface. This low drop point minimizes the potential effects of bouncing or rolling due to high drop speeds, helping the sim-to-sim transfer. Furthermore, collisions are handled differently in the two simulators. In Galactic, we "allow sliding" for training like described in Section 3.2, meaning that the base of the robot will slide when the arm penetrates the environment. However, in Habitat 2.0, the arm can collide with the surrounding scene.

We compare to training the same policy architecture directly in Habitat 2.0 with RL. Specifically, we use an identical reward and policy architecture setup as described in Section 4.2 in

Habitat 2.0. We train this policy for 50 hours which leads to around 200M steps of training in Habitat 2.0 using the VER trainer [59]. Training in Habitat 2.0 for 5 billion steps like Galactic would take 46 days of training. Both the Galactic policy and the Habitat 2.0 baseline policy utilize comparable compute resources for training.

|  | Galactic | Galactic → H2.0 | H2.0 |
|---|---|---|---|
| **Train** | $95.30 \pm 0.67$ | $36.70 \pm 0.46$ | $0.00 \pm 0.00$ |
| **Eval** | $86.70 \pm 1.06$ | $26.40 \pm 0.43$ | $0.00 \pm 0.00$ |

Table 2. Success rates in the Rearrange task. The middle column shows the zero-shot success of the policy trained in Galactic on the fully dynamic Habitat 2.0 (H2.0). The rightmost column is the policy trained purely in H2.0. The leftmost column shows success for the Galactic trained policy evaluated in Galactic (grayed since it reports success in Galactic while other columns are in H2.0). Numbers represent mean and standard error across 1k episodes.

**Results.** We evaluate the success rate of the trained policies on unseen episodes in Table 2. The top row of Table 2 shows the success rate for 1k of the train episodes, and the bottom row shows the success rate over the evaluation dataset of 1k episodes (the "val" dataset from the Rearrange Easy challenge). The leftmost column (Galactic → Galactic) of Table 2 shows the performance of the policy trained in Galactic on the unseen episodes in Galactic, demonstrating that the policy can learn generalizable rearrangement behavior from 5B steps of training. In Habitat 2.0 → Habitat 2.0, we zero-shot transfer the policy trained in Galactic to Habitat 2.0. Despite differences in visuals and dynamics between the simulators, the policy can still achieve 26.40% success rate without any further training. This far outperforms training directly in Habitat 2.0 with a similar compute budget as demonstrated by the policy trained directly in Habitat 2.0 achieving no success. This result is consistent with [55, 56] where end-to-end policies also achieve no success. The policy in Habitat 2.0 learns too slow, and even after 200M steps, it only learns to pick the object up but struggles to place it. In Appendix C, we further analyze policy in Habitat 2.0 and the errors in the zero-shot transfer.

### 5.3. Analyzing Rearrangement Settings in Galactic

In this section, we analyze what properties of the rearrangement task structure are important for RL training. Galactic allows us to answer these questions at scale by training policies fully to convergence, even if it takes billions of steps.

We analyze three variations of the rearrange task with different conditions around the termination and drop action.
- *Rearrange* : The "default" version of rearrangement where the policy may drop the object at any point or call the stop action at any point in the episode.
- *Rearrange [Safe Drop/Stop]* : Same as *Rearrange* except the policy has a heuristic condition around when the stop and drop actions are executed. The drop action is

ignored except when the robot's end-effector is within a cylinder of radius 0.15m and height of 0.3m around the goal location. The stop action is ignored before the robot has dropped the object.

- *Rearrange [No Distractors]* : A version of Rearrange [Safe Drop/Stop] where there are no other objects in the scene other than the object the robot needs to move.

Policies are trained for 5 billion steps in the same training setup as Section 5.2. Figure 5 compares the learning curves for each setting where each point on the curve is the policy checkpoint evaluated for 100 validation episodes (unseen configurations of objects and scenes) from the Tidy House dataset. When evaluating the last checkpoint on 1000 evaluation episodes, we find that *Rearrange [Safe Drop/Stop]* achieves 84.73% success rate while *Rearrange* achieves 79.63% success rate, despite using the same sensor inputs. This demonstrates that the challenges of learning the stop and drop actions harm learning. Furthermore, the gap between *Rearrange [Safe Drop/Stop]* with a success rate of 84.73% and *Rearrange [No Distractors]* with a 99.5% success rate after only 500 million steps shows that the presence of distractor objects makes learning more difficult.
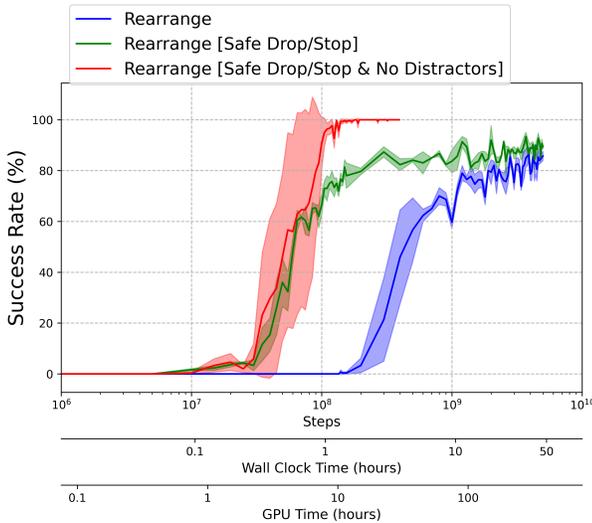


Figure 5. Evaluation success rate of different *Rearrange* tasks on 100 held-out episodes with mean and std across 3 seeds.

### 5.4. Speeding Up Mobile Pick

We also show that Galactic can rapidly train policies in an additional mobile pick task. In the mobile pick task, the robot is spawned within a 2-meter radius of the object to pick up and is provided the start coordinates of the object. The agent then controls the base and arm to navigate to and then pick up the object. The episode is successful if the agent picks up the correct object within the episode horizon of 300 steps. We use the policy architecture described in Section 4.2 with a ResNet18 visual encoder for policies in Galactic and Habitat 2.0.

The plot in Figure 6 shows that Galactic is capable of training

policies in Galactic $100\times$ faster than in Habitat 2.0. To reach an 80% success rate, training in Galactic requires 16 minutes of wall-clock time, while reaching the same success in Habitat 2.0 requires over 26 hours. This $100\times$ speed-up validates that Galactic's higher SPS translates directly to faster wall-clock time-to-convergence. It also demonstrates that Galactic can be used to accelerate even tasks less complex than full rearrangement.
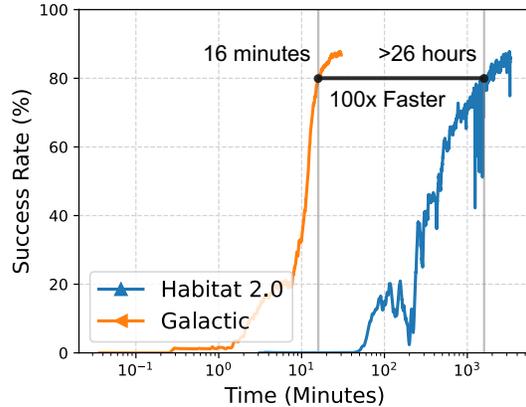


Figure 6. Wall-clock time versus success rate in the mobile pick task between training policies in Galactic and Habitat 2.0. Galactic takes 16 minutes to train a policy to 80% accuracy while Habitat 2.0 requires over 26 hours to train a policy to 80% accuracy.

## 6. Conclusion

We propose Galactic, a framework for rendering, physics, and RL training for embodied tasks at a massive scale. To achieve high throughput for experience collection, we develop an approximate kinematic simulator optimized for embodied rearrangement tasks and propose a batch processing approach to jointly render observations and simulate the physics of the world. Galactic processes more than 100,000 steps per second in RL training. Learning a mobile picking skill takes less than 16 minutes, while training for the same task in Habitat 2.0, one of the fastest embodied simulation frameworks, requires more than 26 hours. We also show that policies trained in Galactic generalize to zero-shot to Habitat 2.0. We hope that Galactic opens up new avenues in embodied AI research by enabling large-scale training.

# References

[1] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian D. Reid, Stephen Gould, and Anton van den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *CVPR*, 2018. 3

[2] Dhruv Batra, Angel X Chang, Sonia Chernova, Andrew J Davison, Jia Deng, Vladlen Koltun, Sergey Levine, Jitendra Malik, Igor Mordatch, Roozbeh Mottaghi, et al. Rearrangement: A challenge for embodied ai. *arXiv*, 2020. 2, 4, 5

[3] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *arXiv*, 2020. 1

[4] Berk Calli, Arjun Singh, Aaron Walsman, Siddhartha Srinivasa, Pieter Abbeel, and Aaron M Dollar. The ycb object and model set: Towards common benchmarks for manipulation research. In *ICAR*, 2015. 14

[5] Devendra Singh Chaplot, Dhiraj Prakashchand Gandhi, Abhinav Gupta, and Russ R Salakhutdinov. Object goal navigation using goal-oriented semantic exploration. In *NeurIPS*, 2020. 1

[6] Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. In *NeurIPS*, 2018. 3

[7] Erwin Coumans and Yunfei Bai. Pybullet, a python module for physics simulation for games, robotics and machine learning. http://pybullet.org. 3, 4

[8] Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied question answering. In *CVPR*, 2018. 3

[9] Marc Peter Deisenroth and Carl Edward Rasmussen. Pilco: A model-based and data-efficient approach to policy search. In *ICML*, 2011. 3

[10] Matt Deitke, Winson Han, Alvaro Herrasti, Aniruddha Kembhavi, Eric Kolve, Roozbeh Mottaghi, Jordi Salvador, Dustin Schwenk, Eli VanderBilt, Matthew Wallingford, Luca Weihs, Mark Yatskar, and Ali Farhadi. Robothor: An open simulation-to-real embodied ai platform. In *CVPR*, 2020. 3

[11] Matt Deitke, Eli VanderBilt, Alvaro Herrasti, Luca Weihs, Jordi Salvador, Kiana Ehsani, Winson Han, Eric Kolve, Ali Farhadi, Aniruddha Kembhavi, et al. Procthor: Large-scale embodied ai using procedural generation. In *NeurIPS*, 2022. 2

[12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 1

[13] 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 *ICRA*, 2022. 15

[14] Kiana Ehsani, Winson Han, Alvaro Herrasti, Eli VanderBilt, Luca Weihs, Eric Kolve, Aniruddha Kembhavi, and Roozbeh Mottaghi. ManipulaTHOR: A Framework for Visual Object Manipulation. In *CVPR*, 2021. 3

[15] Lasse Espeholt, Hubert Soyer, Rémi Munos, Karen Simonyan, Volodymyr Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, Shane Legg, and Koray Kavukcuoglu. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. In *ICML*, 2018. 3

[16] 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. In *NeurIPS Datasets and Benchmarks Track*, 2021. 2, 3, 4, 5

[17] Chuang Gan, Jeremy Schwartz, Seth Alter, Martin Schrimpf, James Traer, Julian De Freitas, Jonas Kubilius, Abhishek Bhandwaldar, Nick Haber, Megumi Sano, Kuno Kim, Elias Wang, Damian Mrowca, Michael Lingelbach, Aidan Curtis, Kevin T. Feigelis, Daniel Bear, Dan Gutfreund, David Cox, James J. DiCarlo, Josh H. McDermott, Joshua B. Tenenbaum, and Daniel L. K. Yamins. Threedworld: A platform for interactive multimodal physical simulation. In *NeurIPS (dataset track)*, 2021. 3

[18] Daniel Gordon, Aniruddha Kembhavi, Mohammad Rastegari, Joseph Redmon, Dieter Fox, and Ali Farhadi. Iqa: Visual question answering in interactive environments. In *CVPR*, 2018. 3

[19] Jiayuan Gu, Devendra Singh Chaplot, Hao Su, and Jitendra Malik. Multi-skill mobile manipulation for object rearrangement. In *ICLR*, 2023. 1, 6

[20] Shixiang Shane Gu, Timothy P. Lillicrap, Ilya Sutskever, and Sergey Levine. Continuous deep q-learning with model-based acceleration. In *ICML*, 2016. 3

[21] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. In *ICLR*, 2021. 3

[22] Danijar Hafner, Timothy P. Lillicrap, Ian S. Fischer, Ruben Villegas, David R Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In *ICML*, 2019. 3

[23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 6

[24] Dan Horgan, John Quan, David Budden, Gabriel Barth-Maron, Matteo Hessel, Hado van Hasselt, and David Silver. Distributed prioritized experience replay. In *ICLR*, 2018. 3

[25] Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z Leibo, David Silver, and Koray Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks. In *ICLR*, 2017. 3

[26] Michał Kempka, Marek Wydmuch, Grzegorz Runc, Jakub Toczek, and Wojciech Jaśkowski. ViZDoom: A Doom-based AI research platform for visual reinforcement learning. In *IEEE Conference on Computational Intelligence and Games*, 2016. 2

[27] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. Ai2-thor: An interactive 3d environment for visual ai. *arXiv*, 2017. 3, 4

[28] Chengshu Li, Fei Xia, Roberto Martín-Martín, Michael Lingelbach, Sanjana Srivastava, Bokui Shen, Kent Vainio, Cem Gokmen, Gokul Dharan, Tanish Jain, et al. igibson 2.0: Object-centric simulation for robot learning of everyday household tasks. *arXiv*, 2021. 2

[29] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac gym: High performance gpu-based physics simulation for robot learning. *arXiv*, 2021. 2, 3, 4, 5

[30] Oleksandr Maksymets, Vincent Cartillier, Aaron Gokaslan, Erik Wijmans, Wojciech Galuba, Stefan Lee, and Dhruv Batra. Thda: Treasure hunt data augmentation for semantic navigation. In *ICCV*, 2021. 1

[31] Piotr Wojciech Mirowski, Razvan Pascanu, Fabio Viola, Hubert Soyer, Andy Ballard, Andrea Banino, Misha Denil, Ross Goroshin, L. Sifre, Koray Kavukcuoglu, Dharshan Kumaran, and Raia Hadsell. Learning to navigate in complex environments. In *ICLR*, 2017. 3

[32] Rémi Munos, Tom Stepleton, Anna Harutyunyan, and Marc G. Bellemare. Safe and efficient off-policy reinforcement learning. In *NeurIPS*, 2016. 3

[33] Ofir Nachum, Shixiang Shane Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement learning. In *NeurIPS*, 2018. 3

[34] Anusha Nagabandi, Gregory Kahn, Ronald S. Fearing, and Sergey Levine. Neural network dynamics for model-based deep reinforcement learning with model-free fine-tuning. In *ICRA*, 2018. 3

[35] Arun Nair, Praveen Srinivasan, Sam Blackwell, Cagdas Alcicek, Rory Fearon, Alessandro De Maria, Vedavyas Panneershelvam, Mustafa Suleyman, Charlie Beattie, Stig Petersen, Shane Legg, Volodymyr Mnih, Koray Kavukcuoglu, and David Silver. Massively parallel methods for deep reinforcement learning. *arXiv*, 2015. 3

[36] Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *ICML*, 2017. 3

[37] Aleksei Petrenko, Erik Wijmans, Brennan Shacklett, and Vladlen Koltun. Megaverse: Simulating embodied agents at one million experiences per second. In *ICML*, 2021. 2, 3

[38] Xavier Puig, Kevin Kyunghwan Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Torralba. Virtualhome: Simulating household activities via programs. In *CVPR*, 2018. 3

[39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *arXiv*, 2021. 1

[40] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *arXiv*, 2019. 1

[41] Santhosh K Ramakrishnan, Aaron Gokaslan, Erik Wijmans, Oleksandr Maksymets, Alex Clegg, John Turner, Eric Undersander, Wojciech Galuba, Andrew Westbury, Angel X Chang, et al. Habitat-matterport 3d dataset (hm3d): 1000 large-scale 3d environments for embodied ai. *arXiv*, 2021. Matterport license available at http://kaldir.vc.in.tum.de/matterport/MP_TOS.pdf. 1, 15

[42] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, 2021. 1

[43] Fetch robotics. Fetch. http://fetchrobotics.com/, 2020. 2, 5, 14

[44] Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 1

[45] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. In *ICCV*, 2019. 3, 4

[46] Alexander Sax, Bradley Emi, Amir R. Zamir, Leonidas J. Guibas, Silvio Savarese, and Jitendra Malik. Mid-level visual representations improve generalization and sample efficiency for learning visuomotor policies. *arXiv*, 2018. 3

[47] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv*, 2017. 6

[48] Brennan Shacklett, Erik Wijmans, Aleksei Petrenko, Manolis Savva, Dhruv Batra, Vladlen Koltun, and Kayvon Fatahalian. Large batch simulation for deep reinforcement learning. In *ICLR*, 2021. 2, 3, 4

[49] Rutav M. Shah and Vikash Kumar. RRL: resnet as representation for reinforcement learning. In *ICML*, 2021. 3

[50] Evan Shelhamer, Parsa Mahmoudieh, Max Argus, and Trevor Darrell. Loss is its own reward: Self-supervision for reinforcement learning. *arXiv*, 2017. 3

[51] Bokui Shen, Fei Xia, Chengshu Li, Roberto Mart'in-Mart'in, Linxi (Jim) Fan, Guanzhi Wang, S. Buch, Claudia. Pérez D'Arpino, Sanjana Srivastava, Lyne P. Tchapmi, Micael Edmond Tchapmi, Kent Vainio, Li Fei-Fei, and Silvio Savarese. igibson, a simulation environment for interactive tasks in large realistic scenes. In *IROS*, 2021. 2, 3

[52] Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *CVPR*, 2020. 3

[53] Sanjana Srivastava, Chengshu Li, Michael Lingelbach, Roberto Mart'in-Mart'in, Fei Xia, Kent Vainio, Zheng Lian, Cem Gokmen, S. Buch, C. Karen Liu, Silvio Savarese, Hyowon Gweon, Jiajun Wu, and Li Fei-Fei. Behavior: Benchmark for everyday household activities in virtual, interactive, and ecological environments. In *CoRL*, 2021. 3

[54] Adam Stooke and Pieter Abbeel. Accelerated methods for deep reinforcement learning. *arXiv*, 2018. 3

[55] Andrew Szot, Alexander Clegg, Eric Undersander, Erik Wijmans, Yili Zhao, John Turner, Noah Maestre, Mustafa Mukadam, Devendra Singh Chaplot, Oleksandr Maksymets, Aaron Gokaslan, Vladimir Vondrus, Sameer Dharur, Franziska Meier, Wojciech Galuba, Angel Xuan Chang, Zsolt Kira, Vladlen Koltun, Jitendra Malik, Manolis Savva, and Dhruv Batra. Habitat 2.0: Training home assistants to rearrange their habitat. In *NeurIPS*, 2021. 1, 2, 3, 4, 5, 6, 7, 14

[56] Andrew Szot, Karmesh Yadav, Alex Clegg, Vincent-Pierre Berges, Aaron Gokaslan, Angel Chang, Manolis Savva, Zsolt Kira, and Dhruv Batra. Habitat rearrangement challenge 2022. https://aihabitat.org/challenge/rearrange_2022, 2022. 5, 6, 7

[57] Joanne Truong, Max Rudolph, Naoki Harrison Yokoyama, Sonia Chernova, Dhruv Batra, and Akshara Rai. Rethinking sim2real:

Lower fidelity simulation leads to higher sim2real transfer in navigation. In *CoRL*, 2022. 3, 4

[58] Luca Weihs, Matt Deitke, Aniruddha Kembhavi, and Roozbeh Mottaghi. Visual room rearrangement. In *CVPR*, 2021. 1, 3

[59] Erik Wijmans, Irfan Essa, and Dhruv Batra. Ver: Scaling on-policy rl leads to the emergence of navigation in embodied rearrangement. *arXiv*, 2022. 1, 2, 3, 6, 7

[60] Erik Wijmans, Abhishek Kadian, Ari S. Morcos, Stefan Lee, Irfan Essa, Devi Parikh, Manolis Savva, and Dhruv Batra. Dd-ppo: Learning near-perfect pointgoal navigators from 2.5 billion frames. In *ICLR*, 2020. 1, 3, 6

[61] Mitchell Wortsman, Kiana Ehsani, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Learning to learn how to learn: Self-adaptive visual navigation using meta-learning. In *CVPR*, 2019. 3

[62] Fei Xia, Amir R Zamir, Zhiyang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gibson env: Real-world perception for embodied agents. In *CVPR*, 2018. 3

[63] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, Li Yi, Angel X.Chang, Leonidas Guibas, and Hao Su. SAPIEN: A SimulAted Part-based Interactive ENvironment. In *CVPR*, 2020. 3

[64] Jie Xu, Viktor Makoviychuk, Yashraj Narang, Fabio Ramos, Wojciech Matusik, Animesh Garg, and Miles Macklin. Accelerated policy learning with parallel differentiable simulation. *ICLR*, 2022. 2, 3, 4, 5

[65] Karmesh Yadav, Ram Ramrakhya, Arjun Majumdar, Vincent-Pierre Berges, Sachit Kuhar, Dhruv Batra, Alexei Baevski, and Oleksandr Maksymets. Offline visual representation learning for embodied navigation. *arXiv*, 2022. 1

[66] Denis Yarats, Amy Zhang, Ilya Kostrikov, Brandon Amos, Joelle Pineau, and Rob Fergus. Improving sample efficiency in model-free reinforcement learning from images. In *AAAI*, 2021. 3

[67] Amir Yazdanbakhsh, Junchao Chen, and Yu Zheng. Menger: Massively large-scale distributed reinforcement learning. *NeurIPS, Beyond Backpropagation Workshop*, 2020. 3

[68] Joel Ye, Dhruv Batra, Erik Wijmans, and Abhishek Das. Auxiliary tasks speed up learning pointgoal navigation. In *CoRL*, 2020. 1, 3

[69] 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 *ICRA*, 2020. 3

[70] Albert Zhan, Philip Zhao, Lerrel Pinto, Pieter Abbeel, and Michael Laskin. Learning visual robotic control efficiently with contrastive pre-training and data augmentation. In *IROS*, 2022. 3

# A. Additional Task Details

The full reward function of the *Rearrange* task is described in Equation (1).

$$r_t = 10\mathbb{I}_{success} + 5\mathbb{I}_{pick} + \Delta^{obj}_{arm} + \Delta^{goal}_{obj} - 0.001C_t \quad (1)$$

Where:

- $\mathbb{I}_{success}$ is the indicator for task success.

- $\mathbb{I}_{pick}$ is the indicator if the agent just picked up the object.

- $\Delta^{obj}_{arm}$ is the change in Euclidean distance between the end-effector ($arm$) and the target object ($obj$). If $d_t$ is the distance between the two at timestep $t$, then $\Delta^{obj}_{arm} = d_{t-1} - d_t$).

- $\Delta^{goal}_{obj}$ is the change in Euclidean distance between the object (obj) and the goal position (goal).

- $C_t$ Is the squared difference in joint action values between the current and previous time step. If $a^k_t$ is the action for timestep t for moving joint $k$ then $C_t = \sum_k (a^k_t - a^k_{t-1})^2$

The reward signal for the *mobile pick* is identical, but the task ends in a success if the robot picks the correct object ($\mathbb{I}_{success} = \mathbb{I}_{pick}$).

The action space of the monolithic policy consist of 11 actions:

- 7 continuous actions controlling the change to the joints angles. These actions are normalized between $-1$ and $1$ with the minimum and maximum value corresponding to the maximum change to the joint angles allowed in each direction per step.

- 1 continuous action between $-1$ and $1$ corresponding to the robot moving forward. An action with value of $1$ corresponds to the robot moving forward by $10$cm and $-1$ to the robot moving backward by $10$cm in the simulation.

- 1 continuous action between $-1$ and $1$ corresponding to the robot rotating. An action with value of $1$ corresponds to the robot rotating in clockwise by $5°$ and $-1$ to the robot rotating counter-clockwise by $5°$.

- 1 discrete action with 2 options corresponding to the robot attempting to grasp or release an object. If the value is 0, and the robot is holding an object, the robot will attempt to release it. If the value is $1$ and the robot is not holding an object, the robot will attempt to grasp.

- 1 discrete action with 2 options corresponding to the robot attempting to terminate the episode. If the value is $0$ the robot will continue the task. If the value is $1$, the robot will signal that the task is completed.

# B. Method Details

More details about the method architecture here. Our Hyperparameters are described in Table 3.

| Hyperparameter | Value |
|---|---|
| start learning rate | $3.5 \times 10^{-4}$ |
| end learning rate | $0$ |
| learning rate schedule | $linear$ |
| entropy coefficient | $1 \times 10^{-3}$ |
| clip gradient norm | $2.0$ |
| time horizon | $64$ |
| number of epochs per updates | $1$ |
| number of mini batches per updates | $2$ |
| RGB and Depth image resolution | $128 \times 128$ |
| image encoder | $ResNet18$ |
| normalized advantage | $true$ |

Table 3. Hyperparameters used for DD-PPO training in Galactic

To calculate the entropy of this action space for entropy regularization in DD-PPO, we add the entropy of the discrete and continuous actions distributions together without any scaling.

The SimpleCNN model we use consists of 3 convolution layers followed by a fully connected layer. The kernel sizes for the three convolution layers are $8 \times 8$, $4 \times 4$ and $3 \times 3$, the strides are $4 \times 4$, $2 \times 2$ and $1 \times 1$ and there is no dilation nor padding. This is the same SimpleCNN visual encoder used in Habitat 2.0. The size of the models used are described in Table 4.
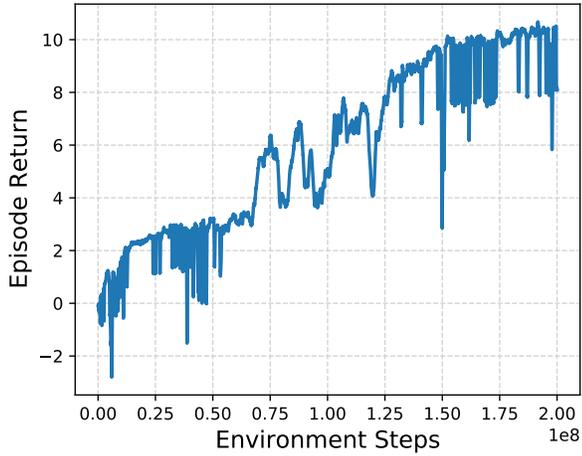
| Model | Total number of parameters |
|---|---|
| SimpleCNN | 4,046,999 |
| ResNet9 | 4,338,007 |
| ResNet18 | 5,906,647 |

Table 4. Model sizes for the different visual encoders used. This includes the visual encoder, the actor, and the critic.
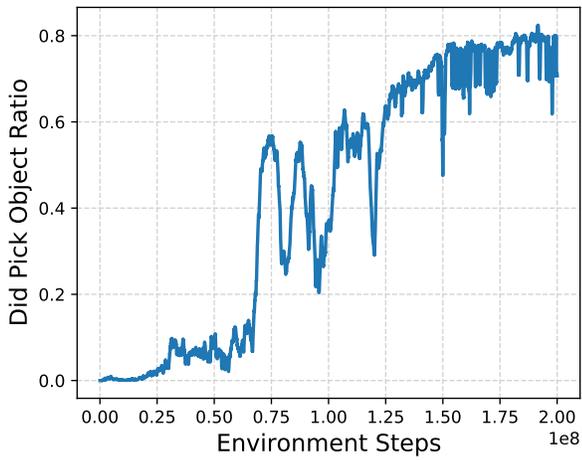
# C. Further Habitat 2.0 Results

First, we analyze the poor performance of the policy purely trained in Habitat 2.0, which achieves no success in the Table 2. Figure 7a shows the reward learning curve during training. This learning curve demonstrates that even after 200M steps of training, the reward is still increasing, which provides evidence for the necessity of Galactic to scale training. In this training time, the agent reliably learns to pick the object around $80\%$ of the time as shown by the training plot in Figure 7b comparing the fraction of the time the robot picked the object within an episode versus the number of training steps.

Next, we analyze the source of errors in the zero-shot transfer from Galactic to Habitat 2.0. We show the drop in performance is not due to the dynamic arm control by comparing to transferring to Habitat 2.0 with a kinematic arm controller instead of a dynamics-based torque controller. The agent with

(a) Reward



(b) Picked object ratio

Figure 7. Learning curves for the policy trained purely in Habitat 2.0. Figure 7a shows the episode reward does not saturate even after 200M training steps. Figure 7b shows that even though the agent is never successful, it still learns to pick the object.

the kinematic arm controller achieves a 29.7% success rate on the "Eval" dataset, barely any better than the 26.4% success rate the dynamics-based torque controller achieves.

## D. Additional Task visuals

In this section, we visualize observations rendered using the Galactic simulator. Figure 8 are examples of $128 \times 128$ RGB images used for training. Figure 9 are examples of $128 \times 128$ depth images used for training. We also visualize observations rendered using the Habitat 2.0 simulator also at $128 \times 128$ in Figure 10 and Figure 11.

## E. Training for 15 Billion steps

To show the usefulness of training for several billions of steps, we trained the *Rearrange* task defined in Section 4 for 15 billion steps. Training and validation success rates are still
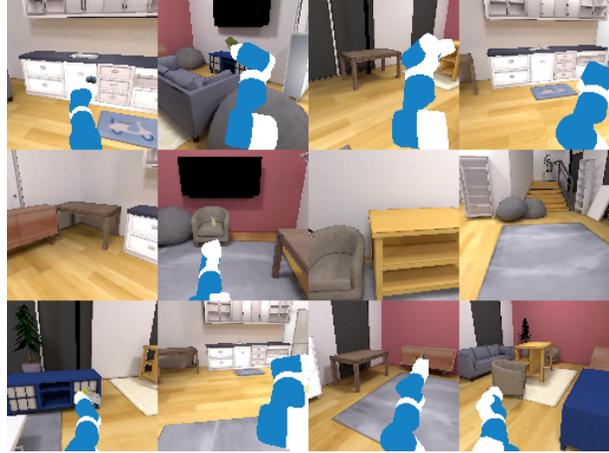


Figure 8. Samples of RGB observations collected in Galactic.



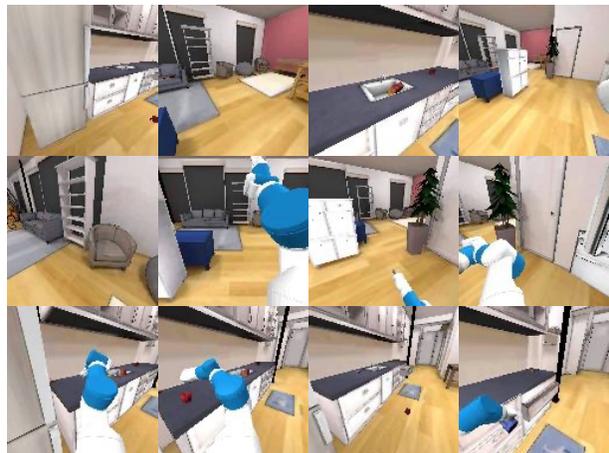Figure 9. Samples of Depth observations collected in Galactic.



Figure 10. Samples of RGB observations collected in Habitat 2.0.

improving, showing that training still hasn't converged, even after 15 billion steps.

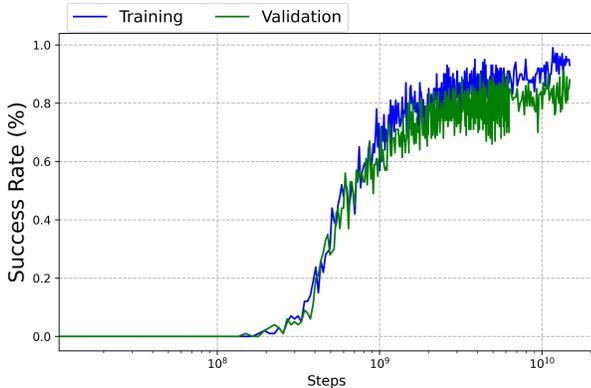Figure 11. Samples of Depth observations collected in Habitat 2.0.



Figure 12. Training and evaluation curves for a 15 billion steps training run. Each checkpoint is evaluated on 100 training or validation episodes.

| Visual Encoder | ResNet18 | SimpleCNN |
|---|---|---|
| Number of Envs | 128 | 512 |
| PyTorch Inference | 7.52 | 8.80 |
| Render Setup | 0.77 | 1.47 |
| Step Post-processing* | 1.84 | 2.05 |
| Step Physics* | 1.78 | 3.76 |
| GPU Rendering | 2.27 | 9.68 |
| Additional CPU | 0.74 | 1.00 |
| Total | 11.30 ms | 20.95 ms |

Table 5. Timing breakdown of a single batch rollout step, for two configurations in milliseconds. * Post-processing and Physics are interleaved with GPU rendering and PyTorch inference and don't contribute to total rollout step time. See also Figure 1. 1x Tesla V100, 10x Intel Xeon Gold 6230 CPU @ 2.10GHz, 128x128 RGBD sensors.

# F. Performance Timings

# G. Additional Collision-Detection Details

In this section, we'll expand on Section 3.2, in particular, we'll discuss our collision representations and collision-detection queries.

| Visual Encoder | ResNet18 | SimpleCNN |
|---|---|---|
| Number of Envs | 128 | 512 |
| Compute Rollouts | 726 | 1289 |
| Update Agent | 1381 | 856 |
| Total | 2204 ms | 2215 ms |
| Training SPS | 3716 SPS | 14791 SPS |

Table 6. Timing breakdown of a single train update for two configurations in milliseconds. 1x Tesla V100, 10x Intel Xeon Gold 6230 CPU @ 2.10GHz, 128x128 RGBD sensors, 64 batch rollout steps.

Galactic scenes include a Fetch robot [43], movable YCB objects [4], and 105 static (non-movable) ReplicaCAD scenes [55]. Note that scenes in the ReplicaCAD dataset include some interactive furniture (e.g. openable cabinet drawers and doors), but we don't simulate these in Galactic as they aren't required for the Rearrange Easy benchmark.

As discussed in Section 3.2, our approximate kinematic sim must perform collision queries between the robot (including grasped object, if any) and the environment (resting movable objects and the static scene). We represent each articulated link of the robot with a set of spheres (green in Figure 3). These are authored manually, with the goal to approximate the shape of the Fetch robot with a minimal number of spheres. We also represent each grasped object with a set of spheres (blue). These are generated offline using a space-filling heuristic. Rather than supporting arbitrary sphere radius, we limit ourselves to a sphere-radius "working set" of{1.5 cm, 5 cm, 12 cm}. This limitation is important as we'll see shortly.

We approximate a ReplicaCAD scene as a voxel-like structure called a column grid (gray in Figure 3). A column grid is generated offline for a particular scene and a particular sphere radius from our working set, so we generate three column grids per scene. A column grid is a dense 2D array of columns in the XZ (ground) plane, with 3-centimeter spacing. For each column, we represent vertical free space as a list of layers. For example, a column in an open area of the room would contain just one layer, storing two floating-point height values roughly corresponding to the height of the floor and the height of the ceiling. A column in the vicinity of a table, meanwhile, would contain two layers: one spanning from the floor to the underside of the table, and another spanning from the table surface to the ceiling. Finally, the stored height values don't actually represent the surface heights themselves, but rather the height of the query sphere (of known radius) in contact with the surface. For ReplicaCAD scenes, the maximum number of layers for any column is approximately 10 and corresponds to columns in the vicinity of a particular bookshelf with many shelves.

A column grid is generated offline using the ReplicaCAD scene's source triangle mesh and Habitat 2.0's sphere-query functionality. We load the scene in Habitat 2.0 and use the scene extents to derive the column grid's XZ (ground-plane) extents. We iterate over this region using our chosen 3-cm

spacing. For each column, we perform a brute-force search of the vertical region at the column's XZ position, using a series of sphere-overlap and vertical sphere-casts to find the free spans.

At runtime, to detect collisions between the robot (including grasped object, if any) and the static scene, we implement a fast sphere-versus-column-grid query. First, we select the appropriate column grid corresponding to the query sphere's radius. Second, we retrieve the nearest column corresponding to the sphere's XZ position. Finally, we linearly search the column's layers to determine whether the query sphere's Y position is in free versus obstructed space. This linear search is accelerated using caching: we start the search from the same layer index found in recent searches. This leverages spatial and temporal coherency, for example, consider the robot reaching under a table: if one sphere from the robot arm's link is found to be between the floor and the underside of a table, it's likely that other spheres from that same link or other queries from succeeding timesteps will also lie in that vertical layer.

Whereas a grasped movable object is represented with a set of spheres (blue in Figure 3), a resting movable object is approximated as an oriented box (orange). This is computed from the YCB object's triangle mesh in a preprocess. At runtime, to detect collisions between the robot (including grasped object, if any) and the resting movable objects, we perform sphere-versus-box queries. There are generally 30 resting movable objects in the environment (1 target object and 29 distractor objects) and we need to avoid performing all 30 sphere-versus-box queries. So, we use a "regular grid" acceleration structure to quickly retrieve a list of nearby resting objects.

Resting movable objects are inserted into a regular grid at episode initialization. This is a dense 2D array spanning the XZ (ground) plane, with each cell storing a list of objects that overlap it. Objects will generally overlap multiple cells and thus be present in the object lists of multiple cells. When an object is grasped by the robot, it is removed from all relevant cells in the regular grid, and if the object is later dropped, it is re-inserted into the regular grid at its new resting position.

Let's consider how to find the list of nearby resting objects for a given query sphere. The regular grid spacing is chosen such that cells are at least $4\times$ the largest radius in our sphere-radius working set (12 cm). A query sphere may overlap up to four adjacent cells in the regular grid, e.g. the sphere is centered near the shared edge of two cells or the shared corner of four cells. A naive approach here would be to merge and de-duplicate the object lists of the four cells. We avoid this expense and instead maintain *four separate regular grids*, all spanning the entire scene XZ extent, with carefully-chosen varying X and Z offsets for the cell boundaries. In this way, any query sphere is guaranteed to lie fully inside a single cell of one of these grids (not spanning a cell edge or corner). Thus, our list of nearby resting objects is simply the list stored in this cell; we don't have to merge or de-duplicate multiple lists. Note this approach of four somewhat-redundant regular grids comes at the expense of extra memory and added insertion/removal compute time.

## H. Simulator Flexibility to new Assets

Galactic can work with various assets (robots, scenes and objects) from different sources. We use a mostly-automated pipeline that includes optimizing assets for the batch renderer and generating collision geometry (see Appendix G). In Fig. 13 we added Stretch and Spot robots loaded in a scene from the MP3D dataset [41] with new objects.
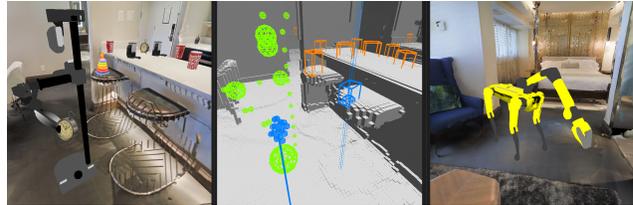


Figure 13. Galactic with an MP3D scene, Google Scanned Objects [13], Stretch robot (left), Stretch debug viz (center), and Spot robot (right).

## I. Description of Heuristics

### I.1. Sliding Heuristic

We implement the robot sliding heuristic as the following steps: (1) start from a candidate pose, (2) if the pose penetrates the scene, compute a jitter direction in the ground plane (3) jitter the robot base in the horizontal plane. Repeat from step 1 until a penetration-free pose is found, up to 3 times. If this fails, the robot does not move on that step.

### I.2. Object Placing Heuristic

The object placement heuristic is as follows: (1) start from a candidate pose, (2) if the pose penetrates the scene, compute a jitter direction in the ground plane, with some randomness, (3) jitter the dropped object and re-cast down to a support surface. Repeat from step 1 until a penetration-free pose is found, up to 6 times. If this fails, we restore the dropped object to its resting position prior to grasp. This approximates the dropped object bouncing or rolling away. Snap-to-surface is an instantaneous operation that resolves within one physics step; objects do not fall or settle over time.