

BEHAVIOR: Benchmark for Everyday Household Activities in Virtual, Interactive, and Ecological Environments

Sanjana Srivastava* Chengshu Li* Michael Lingelbach* Roberto Martín-Martín*
 Fei Xia Kent Vainio Zheng Lian Cem Gokmen Shyamal Buch C. Karen Liu
 Silvio Savarese Hyowon Gweon Jiajun Wu Li Fei-Fei

Stanford University

Abstract: We introduce BEHAVIOR, a benchmark for embodied AI with 100 activities in simulation, spanning a range of everyday household chores such as cleaning, maintenance, and food preparation. These activities are designed to be realistic, diverse and complex, aiming to reproduce the challenges that agents must face in the real world. Building such a benchmark poses three fundamental difficulties for each activity: definition (it can differ by time, place, or person), instantiation in a simulator, and evaluation. BEHAVIOR addresses these with three innovations. First, we propose an object-centric, predicate logic-based description language for expressing an activity’s initial and goal conditions, enabling generation of diverse instances for any activity. Second, we identify the simulator-agnostic features required by an underlying environment to support BEHAVIOR, and demonstrate its realization in one such simulator. Third, we introduce a set of metrics to measure task progress and efficiency, absolute and relative to human demonstrators. We include 500 human demonstrations in virtual reality (VR) to serve as the human ground truth. Our experiments demonstrate that even state-of-the-art embodied AI solutions struggle with the level of realism, diversity, and complexity imposed by the activities in our benchmark. We make BEHAVIOR publicly available at behavior.stanford.edu to facilitate and calibrate the development of new embodied AI solutions.

Keywords: Embodied AI, Benchmarking, Household Activities

1 Introduction

Embodied AI refers to the study and development of artificial agents that can perceive, reason, and interact with the environment with the capabilities and limitations of a physical body. Recently, significant progress has been made in developing solutions to embodied AI problems such as (visual) navigation [1–5], interactive Q&A [6–10], instruction following [11–15], and manipulation [16–22]. To calibrate the progress, several lines of pioneering efforts have been made towards benchmarking embodied AI in simulated environments, including Rearrangement [23, 24], TDW Transport Challenge [25], VirtualHome [26], ALFRED [11], Interactive Gibson Benchmark [27], MetaWorld [28], and RLBench [29], among others [30–32]). These efforts are inspiring, but their activities represent only a fraction of challenges that humans face in their daily lives. To develop artificial agents that can eventually perform and assist with everyday activities with human-level robustness and flexibility, we need a comprehensive benchmark with activities that are more **realistic**, **diverse**, and **complex**.

But this is easier said than done. There are three major challenges that have prevented existing benchmarks to accommodate more realistic, diverse, and complex activities:

- Definition: Identifying and defining meaningful activities for benchmarking;
- Realization: Developing simulated environments that realistically support such activities;
- Evaluation: Defining success and objective metrics for evaluating performance.

*indicates equal contribution
 correspondence to {sanjana2,chengshu}@stanford.edu

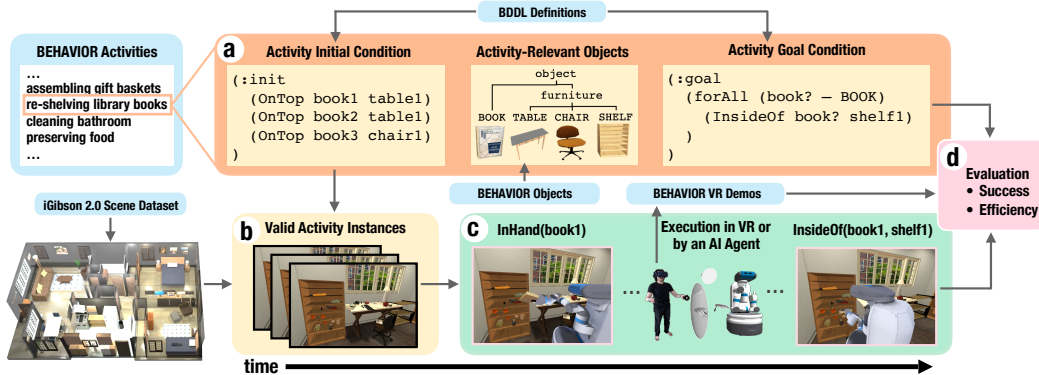


Figure 1: **Benchmarking Embodied AI with BEHAVIOR:** ① We define 100 realistic household activities from the American Time Use Survey [33] and define them with a set of relevant objects, organized with WordNet [34], and logic-symbolic initial and goal conditions in BDDL (Sec. 4). ② We provide an implementation of BEHAVIOR in iGibson 2.0 that generates potentially infinite diverse activity instances in realistic home scenes using the definition. ③ AI agents perform the activities in simulation through continuous physical interactions of an embodied avatar with the environment. Humans can perform the same activities in VR. BEHAVIOR includes a dataset of 500 successful VR demonstrations. ④ Changes in the scene are continuously mapped to their logic-symbolic equivalent representation in BDDL and checked against the goal condition; we provide intermediate success scores, metrics on agent’s efficiency, and a human-centric metric relative to the demonstrations.

We propose **BEHAVIOR** (Fig. 1)–**B**enchmark for **E**veryday **H**ousehold **A**ctivities in **V**irtual, **I**nteractive, and **ec**Ological **envi**Ronments, addressing the three key challenges aforementioned with three technical innovations. First, we introduce BEHAVIOR Domain Definition Language (BDDL), a representation adapted from predicate logic that maps simulated states to semantic symbols. It allows us to define 100 activities as initial and goal conditions, and further enables generation of potentially infinite initial states and solutions for achieving the goal states. Second, we facilitate its realization by listing environment-agnostic functional requirements for realistic simulation. With proper engineering, BEHAVIOR can be implemented in many existing environments; we provide a fully functional instantiation in iGibson 2.0 in this paper including the necessary object models (1217 models of 391 categories). Third, we provide a comprehensive set of metrics to evaluate agent performance in terms of success and efficiency. To make evaluation comparable across diverse activities, scenes, and instances, we propose a set of metrics relative to demonstrated human performance on each activity, and provide a large-scale dataset of 500 human demonstrations (758.5 min) in virtual reality, which serve as ground truth for evaluation and may also facilitate developing imitation learning solutions.

BEHAVIOR activities are realistic, diverse, and complex. They comprise of 100 activities often performed by humans in their homes (e.g., cleaning, packing or preparing food) and require long-horizon solutions for changing not only the position of multiple objects but also their internal states or texture (e.g., temperature, wetness or cleanliness levels). As we demonstrate by experimentally evaluating the performance of two state-of-the-art reinforcement learning algorithms (Section 7), these properties make BEHAVIOR activities extremely challenging for existing solutions. By presenting well-defined challenges beyond the capabilities of current solutions, BEHAVIOR can serve as a unifying benchmark that guides the development of embodied AI.

2 Related Work

Benchmarks and datasets have played a critical role in recent impressive advances in AI, particularly computer vision. Image [35–38] and video datasets [39–44] enable study and development of solutions for important research questions by providing both training data and fair comparison. These datasets, however, are passive observations, and therefore not well suited for development of embodied AI that must control and understand the consequences of their own actions.

Benchmarks for Embodied AI: Although real-world challenges [45–52] provide the ultimate testbed for embodied AI agents, benchmarks in simulated environments serve as useful alternatives with several advantages; simulation enables faster, safer learning, and supports more reproducible, accessible, and fair evaluation. However, in order to serve as a meaningful proxy for real-world

		BEHAVIOR	AI2THOR Vis. Room Rearr.																	ManipulatorTHOR ArmPointNav																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																															
			TDW Transport					Rearrangement T5 (Habitat)					Interactive Gibson Benchmark					VirtualHome					ALFRED					Rearrangement T2 (OCRTOC)					IKEA Furniture Assembly					RL Bench					Metaworld					Robosuite					SoftCym					DeepMind Control Suite					OpenAI Gym					Habitat 1.0					Gibson																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																										
			Mobile manipulation					Static manipulation					Navigation																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																						
Realism	Activity selections reflect human behavior	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗

¹ Estimate of a near-optimal, e.g. human, execution of the activity given the platform's action space

Table 1: Comparison of Embodied AI Benchmarks: BEHAVIOR activities are exceptionally realistic due to their grounding in human population time use [33] and realistic simulation (sensing, actuation, changes in environment) in iGibson 2.0. The activity set is diverse in topic, objects used, scenes done in, and state changes required. The diversity is reinforced by the ability to generate infinite new instances scene-agnostically. BEHAVIOR activities are complex enough to reflect real-world housework: many decision steps and objects in each activity. This makes BEHAVIOR uniquely well-suited to benchmark task-planning and control, and it is the only one to include human VR demonstrations (see Table A.1 for more detail).

performance, simulation benchmarks need to achieve high levels of 1) **realism** (in the activities, the models, the sensing and actuation of the agent), 2) **diversity** (of scenes, objects and activities benchmarked), and 3) **complexity** (length, number of objects, required skills and state changes). Below we review existing benchmarks based on these three criteria (see Table 1 for a summary).

Benchmarks for *visual navigation* [53, 54] provide high levels of visual realism and diversity of scenes, but they often lack interactivity or diversity of activities. The Interactive Gibson Benchmark [27] trades off some visual realism for physically realistic object manipulation in order to benchmark interactive visual navigation. While benchmarks for *stationary manipulation* [55, 29, 28, 30, 56, 31, 32] fare well on physical realism, they commonly fall short on diversity (of scenes, objects, tasks) and complexity (e.g., simple activities that take a few seconds). Benchmarks for *instruction following* [11, 26] provide diversity of scenes, objects and possible changes of the environment, but with low level of complexity; the horizon of the activities is shorter as the agents decide among a discrete set of predefined action primitives with full access to the state of the world.

Closer to BEHAVIOR, a recent group of benchmarks has focused on *rearrangement tasks* [23–25] in realistic simulation environments with diverse scenes. The initial Rearrangement position paper [23] poses critical questions such as how to define embodied AI tasks and measure solution quality. Importantly, however, most household activities go far beyond the scope of rearrangement (see comparison in Fig. A.2). While such focus can inspire new solutions for solving rearrangement tasks, these solutions may not generalize to activities that require more than physical manipulation of object coordinates. Indeed, the majority of household activities involve other state changes (cooking, washing, etc. (Fig. A.2, [33])). BEHAVIOR therefore incorporates 100 activities that humans actually spend time on at home [33] (Sec. 3). To express such diverse activities in a common language, we present a novel logic-symbolic representation that defines activities in terms of initial and goal states, inspired by but distinct from the Planning Domain Definition Language [57]. These yield in principle infinite instances per activity and accept any meaningful solution. We implement activity-independent metrics including a human-centric metric normalized to human performance; to facilitate comparison and development of new solutions, we also present a dataset of 500 successful VR demonstrations.

3 BEHAVIOR: Benchmarking Realistic, Diverse, Complex Activities

Building on the advances led by existing benchmarks, BEHAVIOR aims to reach new levels of realism, diversity, and complexity by using household activities as a domain for benchmarking AI. See Table 1 for comparisons between BEHAVIOR and existing benchmarks.

Realism in BEHAVIOR Activities: To effectively benchmark embodied AI agents in simulation, we need realistic activities that pose similar challenges to those in the real world. BEHAVIOR achieves this by using a data-driven approach to identify activities that approximate the true distribution of real household activities. To this end, we use the American Time Use Survey (ATUS, [33]): A survey from the U.S. Bureau of Labor Statistics on how Americans spend their time. BEHAVIOR activities come from, and are distributed similarly to, the full space of simulatable activities in ATUS (see Fig. A.2). The use of an independently curated source of real-world activities is a unique strength of BEHAVIOR as a benchmark that reflects natural behaviors of a large population.

BEHAVIOR also achieves realism by simulating these activities in reconstructions of real-world homes. We use iGibson 2.0, a simulation environment with realistic physics simulation from the Bullet [58] physics engine and high-quality virtual sensor signals (see Fig. A.7), which includes 15 ecological, fully interactive 3D models of real-world homes with furniture layouts that approximate their real counterparts. These scenes are further populated with object models created by professional artists from the new BEHAVIOR Object dataset, which includes 1217 models of 391 categories grounded in the WordNet [34] taxonomy. The dataset covers a data-driven selection of activity-related objects (see Fig. A.8). Figs. A.10 and A.9 illustrate examples of objects and taxonomic arrangement. The 100 BEHAVIOR activities, visualized in Fig. A.1, go beyond comparable benchmarks that evaluate a few hand-picked activities in less realistic setups (see Table 1 Realism).

Diversity in BEHAVIOR Activities: Benchmarks with diverse activities demand generalizable solutions. In real-world homes, agents encounter a range of activities that differ in 1) the capabilities required for achieving them, 2) the environments in which they occur (e.g., scenes, objects), and 3) the initial states of a particular scene. BEHAVIOR presents extensive diversity in all these dimensions. We include 100 activities that require a wide variety of state changes (e.g., moving objects, soaking materials, cleaning surfaces, heating/freezing food) demanding a broad set of agent capabilities (see Fig. A.2). To reflect the diversity in the ways humans encounter, understand, and accomplish these activities, we provide two example definitions per activity. BDDL, our novel representation for activity definition, allows new valid instances to be sampled from each definition, providing potentially infinite number of instances per activity. The resulting instances vary over scene, object models, and configuration, supported by implementation in iGibson 2.0 and BEHAVIOR Object dataset. Related benchmarks focus on fewer tasks, mostly limited to kinematic state changes and with scene- or position-constant instantiation (see Table 1 Diversity).

Complexity in BEHAVIOR Activities: Beyond diversity across activities, BEHAVIOR also raises the complexity of the activities themselves by benchmarking full household activities that parallel the length (number of steps an agent needs), the number of objects involved, and the number of required capabilities of real-world chores (see Fig. A.3, comparison in Table 1 Complexity). Compared to activities in existing benchmarks, these activities are very long-horizon with some requiring several thousand steps (even for humans in VR; see Fig. A.12), involve more objects (avg. 10.5), and require a heterogeneous set of capabilities (range: 2 - 8) to change various environment states.

4 Defining Realistic, Diverse, and Complex Household Activities with BDDL

BEHAVIOR challenges embodied AI agents to achieve a diverse set of complex long-horizon household activities through physical interactions in a realistically simulated home environment. Adopting the common formalism of partially-observable Markov decision processes (POMDP), each activity is represented by the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T}, \mathcal{R}, \gamma)$. Here, \mathcal{S} is the state space; \mathcal{A} is the action space; \mathcal{O} is the observation space; $\mathcal{T}(s'|s, a), s \in \mathcal{S}, a \in \mathcal{A}$, is the state transition model; $\mathcal{R}(s, a) \in \mathbb{R}$ is the reward function; γ is the discount factor. Based on a full representation of the physical state, \mathcal{S} , the simulation environment generates realistic transitions to embodied AI agents' actions, $a \in \mathcal{A}$, i.e., physical interactions, and close-to-real observations, $o \in \mathcal{O}$, e.g., virtual images.

We define an *activity* τ as two sets of states, $\tau = \{S_{\tau,0}, S_{\tau,g}\}$, where $S_{\tau,0}$ is a set of possible initial states and $S_{\tau,g}$ is a set of acceptable goal states. In an *activity instance*, the agent must change the world state from some concrete $s_0 \in S_{\tau,0}$ to any $s_g \in S_{\tau,g}$. However, describing activities in the physical state space generates scene- or pose-specific definitions (e.g., [23, 30, 29]) that are far more specific than how humans represent these activities, limiting the diversity and complexity of existing embodied AI benchmarks. To overcome this, we introduce *BEHAVIOR Domain Definition Language*

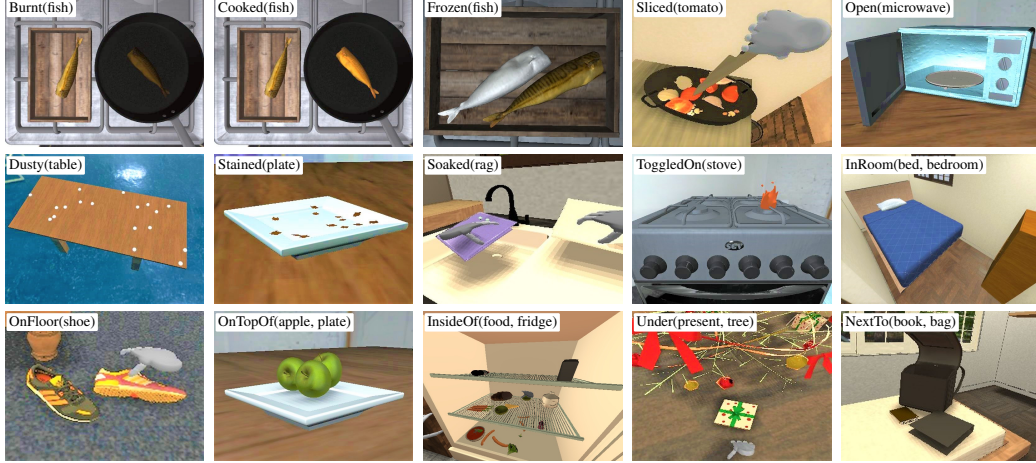


Figure 2: **Unary and Binary Predicates in BDDL:** We represent object states and relationships to other objects based on their kinematics, temperature, wetness level and other physical and functional properties, enabling a diverse and complex set of realistic activities

(BDDL), a predicate logic-based language that establishes a symbolic state representation built on predefined, meaningful predicates grounded in simulated physical states; its variables and constants represent object categories from the BEHAVIOR object dataset. Each activity is defined in BDDL as an initial and goal condition parametrizing sets of possible initial states and satisfactory goal states $\bar{S}_{\tau,0}$ and $\bar{S}_{\tau,g}$. BDDL predicates create symbolic counterparts of the physical state, \bar{S} (see Fig. 2).

BDDL overcomes limitations that hinder diversity through two mechanisms: first, an initial condition maps to infinite physical states in diverse scenes. Second, a goal condition detects all semantically satisfactory solutions, rather than limiting to a few or only those that obey semantically uninteresting geometric constraints (see Fig. A.6 for examples). This state-based definition is also entirely declarative, providing a true benchmark of planning ability. By comparison, other benchmarks are limited to scene- or pose-specific instantiation and solution acceptance, and/or have imperative plans. BEHAVIOR includes a systematic generation pipeline (see A.3.3) allowing unlimited definitions per activity and formalizing the inherent subjectivity and situationality of household activities. We include 200 definitions and 300 activity instances in simulation (see Sec. 5). BEHAVIOR is thus the only benchmark equipped to formalize unlimited human-defined versions of an activity and create practically infinite unique instantiations in any scene.

5 Instantiating BEHAVIOR in a Realistic Physics Simulator

While BEHAVIOR is not bounded to any specific simulation environment, there are a set of functional requirements that are necessary to simulate BEHAVIOR activities: 1) maintain an object-centric representation (object identities enriched with properties and states), 2) simulate physical forces and motion, and generate virtual sensor signals (images), 3) simulate additional, non-kinematic properties per object (e.g. temperature, wetness level, cleanliness level), 4) implement functionality to **generate** valid instances based on the literals defining an activity’s initial condition, e.g., instantiating an object `insideOf` another, and 5) implement functionality to **evaluate** the atomic formulae relevant to the goal condition, e.g. checking whether an object is `cooked` or `onTopOf` another.

Additionally, the simulator must provide an interface of the action space \mathcal{A} and the observation space \mathcal{O} of the underlying POMDP to embodied AI agents (Sec. 4). While BEHAVIOR activities are not tailored to a specific embodiment, we propose two concrete bodies to fulfill the activities (see Fig. 1): a *bimanual humanoid* avatar (24 degrees of freedom, DoF), and a *Fetch robot* (12/13 DoF), both capable of navigating, grasping and interacting with the hand(s). Humans in VR embody the bimanual humanoid. Agents trained with the Fetch embodiment could be directly tested with a real-world version of the hardware (see discussion on sim2real in Sec. A.8). Both embodiments receive sensor signals from the on-board virtual sensors, and perform actions at 30 Hz.

We provide a fully functional implementation of BEHAVIOR using iGibson 2.0, a new version of the open-source simulation environment iGibson that fulfills the requirements above. iGibson

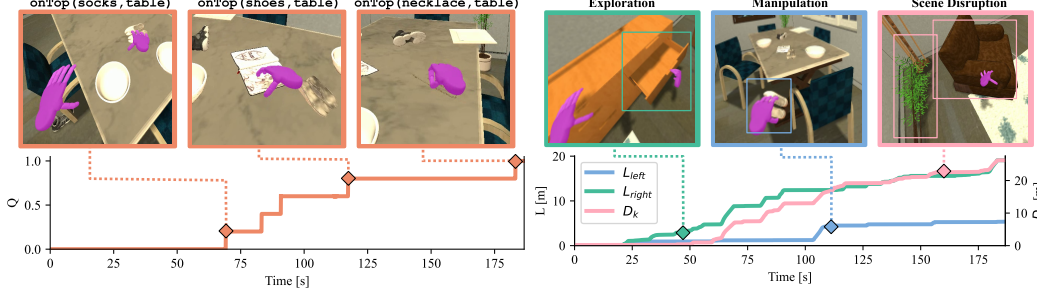


Figure 3: **Evaluation of human performance in `collect_misplaced_items`:** (Left) success score, Q ; (Right) efficiency metrics: kinematic disarrangement, (D_k , dotted), hand interaction displacement (L_{right} , green, and L_{left} , blue); frames at the top depict significant events detected by the metrics; the success score detects the completion of activity-relevant steps; exploration, manipulation and scene disruption events are captured by the efficiency metrics that provide complementary information about the performance of the agent

2.0 provides an object-centric representation with additional properties, support for sources of heat and water, dust and stain particles, and changes in object appearance based on extended states. We implement the two embodiments in iGibson 2.0: the agent receives proprioceptive information and has access to iGibson 2.0’s generated realistic signals: RGB, depth images, LiDAR, normals, flow (optical, spatial), and semantic and instance segmentation. While this control and sensing setup is standard in BEHAVIOR, we additionally implement a set of action primitives inspired by [25, 54, 59, 24] to facilitate solution prototyping and task-planning research. The primitives execute sequences of low-level actions resulting from a motion planning process (bilateral RRT* [60]) to `navigateTo`, `grasp`, `placeOnTop`, `placeInside`, `open`, and `close` the target object provided as arguments. Even though the agent only relies on sensory observations to decide on action primitive, the primitives themselves internally assume access to privileged information (e.g. object identities, poses, and geometric shapes for planning). Further details can be found in Sec. A.4 and in the cross-submission included in appendix. Our implementation of BEHAVIOR in iGibson 2.0 goes beyond the capabilities of existing benchmarks and amplifies realism, diversity, and complexity.

6 Evaluation Metrics: Success, Efficiency and Human-Centric Metric

BEHAVIOR provides evaluation metrics to quantify the performance of an embodied AI solution. Extending prior metrics suggested for Rearrangement [23], we propose a primary metric based on success and several secondary metrics for characterizing efficiency.

Primary Metric – Success Score Q : The main goal of an embodied AI agent in BEHAVIOR is to perform an activity successfully (i.e., all logical expressions in the goal condition are met). A binary definition of success, however, only signals the end of a successful execution and cannot assess interim progress. To provide more guidance to agents and enable comparisons of partial solutions, we propose **success score** as the primary metric, defined as the **maximum fraction of satisfied goal literals in a ground solution to the goal condition** at each step. More formally:

Given an activity τ with goal state set $\bar{S}_{\tau,g}$, its goal condition can be flattened to a set C of conjunctions C_i of ground literals l_{j_i} . For any $C_i \in C$, if all $l_{j_i} \in C_i$ are true then the goal condition is satisfied (see A.3.2 for definitions and technical details on flattening), i.e. for some current environment state s , we have $\bigvee_{C_i} \bigwedge_{l_{j_i}} l_{j_i} = \text{True} \implies s \in \bar{S}_{\tau,g}$. We compute the fraction of literals l_{j_i} that are True

for each C_i , and define the overall success score by taking the maximum: $Q = \max_C \frac{|\{l_{j_i} | l_{j_i} = \text{True}\}|}{|C_i|}$, where $|\cdot|$ is set cardinality.

An activity is complete when all literals in *at least one* C_i of its goal condition are satisfied, achieving $Q = 1$ (100%). Fig. 3 left depicts time evolution of Q during an activity execution. Q extends the fraction of objects in acceptable poses proposed as metric in [23], generalized to any type of activity.

Secondary Metrics – Efficiency: Beyond success, efficiency is critical to evaluation; a successful solution in real-world tasks may be ineffective if it takes too long or causes scene disruption. We propose six secondary metrics that complement the primary metric (see Fig. 3, right, for examples):

- *Simulated time, T_{sim}* : Accumulated time in simulation during execution as the number of simulated steps times the average simulated time per step. T_{sim} is independent of the computer used.
- *Kinematic disarrangement, D_K* : Displacement caused by the agent in the environment. This can be *accumulated* over time, or *differential*, i.e. computed between two time steps, e.g. initial, final.
- *Logical disarrangement, D_L* : Amount of changes caused by the agent in the logical state of the environment. This can be *accumulated* over time or *differential* between two time steps.
- *Distance navigated, L_{body}* : Accumulated distance traveled by the agent’s base body. This metric evaluates the efficiency of the agent in navigating the environment.
- *Displacement of hands, L_{left} and L_{right}* : Accumulated displacement of each of the agent’s hands while in contact with another object for manipulation (i.e., grasping, pushing, etc). This metric evaluates the efficiency of the agent in its interaction with the environment.

These efficiency metrics above can be quantified in absolute units (e.g., distance, time) for scene- and activity-specific comparisons (**general efficiency**). To enable fair comparisons cross diverse activities in BEHAVIOR, we also propose normalization relative to human performance (**human-centric efficiency**); given a human demonstration for an activity instance in VR, each secondary metric can be expressed as a *fraction of the maximum human performance* on that metric.

For this purpose, we present the BEHAVIOR Dataset of Human Demonstrations with 500 successful demonstrations of BEHAVIOR activities in VR (758.5 min). Humans are immersed in iGibson 2.0, controlling the same embodiment used by the AI agents (details in Sec. A.6). The dataset includes a complete record of human actions including manipulation, navigation, and gaze tracking data (Fig. A.12, Fig. A.14, and Fig. A.16), supporting analysis and subactivity segmentation (Fig. A.11). Sec. A.6.2 presents a comprehensive analysis of these data; we quantify human performance in BEHAVIOR efficiency metrics (see Fig. A.12), and Fig. A.13 provides a further decomposition of room occupancy and hand usage across each BEHAVIOR activity. To our knowledge, this is the largest available dataset of human behavior in VR; these data can facilitate development of new solutions for embodied AI (e.g., imitation learning) and also support studies of human cognition, planning, and motor control in ecological environments.

7 Evaluating Reinforcement Learning in BEHAVIOR

In this section, we aim to experimentally demonstrate the challenges imposed by BEHAVIOR’s realism, diversity, and complexity by evaluating the performance of some current state-of-the-art embodied AI solutions. While BEHAVIOR is a benchmark for all kinds of embodied AI methods, here we evaluate two reinforcement learning (RL) algorithms that have demonstrated excellent results in simpler embodied AI tasks with continuous or discrete action spaces [61, 62, 21, 63–67]: Soft-Actor Critic (SAC [16]) and Proximal-Policy Optimization (PPO [17]). We use SAC to train policies in the original low-level continuous action space of the agent, and PPO for experiments using our implemented action primitives (for details on the agents, see Sec. 5). Due to limited computational resources, we run our evaluation on the 12 most simple activities (based on involved types of state changes) until convergence. Reward is given by our staggered success score Q . We use as input to the policies a subset of the realistic agent’s observations, RGB, depth and proprioception (excluding LiDAR, segmentation, etc.). Sec. A.7 includes more experimental details.

Results in the original activities: The first row of Table 2 shows the results of SAC (mean Q at the end of training for 3 seeds) on the original 12 activities with the standard setup: realistic robot actions and onboard sensing. Even for these “simpler” activities, BEHAVIOR is too great a challenge: the training agents do not fulfill any predicate in the goal condition ($Q = 0$). In the following, we will analyze how each dimension of difficulty (realism, diversity, complexity) contributes to these results.

Effect of complexity (activity length): In the first experiment, we evaluate the impact of the activity complexity (time length) in robot learning performance. First, we evaluate the performance of an RL algorithm using our implemented action primitives based on motion planning. These are temporally extended actions that effectively shorten the horizon and length of the activity. The results of training with PPO are depicted in the second row of Table 2. Even in these simpler conditions, agents fail in all but one activity (`bringingInWood`, $Q = 0.13$). In a second oracle-driven experiment, we take a successful human demonstration for each activity from the BEHAVIOR Dataset and save the state of the environment a few seconds before its successful execution at T . We

use this as initial state and train agents with SAC: rows 3 to 6 of Table 2 show the mean success rate (SR , full accomplishment of the activity) in 100 evaluation episodes for the final policy resulting from training with three different random seeds (Q starts here close to 1 and is less informative). Even when starting 1 s away from a goal state, most learning agents fail to achieve the tasks. A few achieve better success but their performance decreases quickly as we start further away from the successful execution, being zero for all activities at 10 s. This indicates that the long-horizon of the activities in BEHAVIOR is in fact a paramount challenge for reinforcement learning. We hypothesize that Embodied AI solutions with a hierarchical structure such as hierarchical-RL or task-and-motion-planning (TAMP) may help to overcome the challenges of high complexity (length) of the BEHAVIOR activities [68–71].

Effect of realism (in sensing and actuation):

In a third experiment, we evaluate how much the realism in actuation and sensing affects the performance of embodied AI solutions. To evaluate the effect of realistic observability of the BEHAVIOR activities in the performance of robot learning approaches, we train agents with continuous motion control (SAC), and motion primitives (PPO) assuming full-observability of the state. Tables 2 (rows 7-8, subindex *FullObs*) depict the results. We observe that even with full observability the complexity dominates policies in the original action space and they do not accomplish any part of the activities. For policies selecting among action primitives, there is some partial success only in five of the activities indicating that the perceptual problem is part of the difficulty in BEHAVIOR. To evaluate the effect of realistic actuation, we train an agent using action primitives that execute without physics simulation, achieving their expected outcome (e.g. grasp an object, or place it somewhere). Tables 2 (row 9-10, subindex *noPhys*) shows the results, also in combination with unrealistic full-observability. We observe that without the difficulties of realistic physics and actuation, the learning agents achieve an important part of most activities, accomplishing consistently two of them ($Q = 1$) when full-observability of the state is also granted. This indicates that the generation of the correct actuation is a critical challenge for embodied AI solutions, even when they infer the right next step at the task-planning level, supporting the importance of benchmarks with realistically action execution over predefined action outcomes.

Effect of diversity (in activity instance and objects):

Another cause of the poor performance of robot learning solutions in the 12 BEHAVIOR activities may be the high diversity in multiple dimensions, such as scenes, objects, and initial states. This diversity forces embodied AI solutions to generalize to all possible conditions. In a second experiment, we evaluate the effect of BEHAVIOR’s diversity on performance. To present diversity across activities while alleviating their complexity, we train RL agents to complete five single-literal activities involving only one or two objects. Note that these activities are not part of BEHAVIOR. We evaluate training with RL (SAC) for each activity under diverse instantiations: initialization of the activity (object poses) and object instances. The results are shown in Table 3, where we report Q . First, we train without any diversity as baseline to understand the ground complexity of the single-literal activities. All agents achieve success. Then, we evaluate how well the RL policies train for a diverse set of instances of the activities, first changing objects’ initial pose, then changing the object. Performance in all activities decreases rapidly, especially in *sliced* and *stained*. These experiments

		bringingInWood collectingDisplacedItems movingBoxesToStorage organizingFridgeCabinet throwingAwayLeftovers puttingDishesAway puttingLeftoversAway ce-shelvingLibraryBooks layingTurfFloors settingUpCandles pickingUpTrash storingFood											
		Q^{ca}	Q^{ap}	$SR^{ca}@T=1s$	$SR^{ca}@T=2s$	$SR^{ca}@T=3s$	$SR^{ca}@T=10s$	$Q^{ca}_{FullObs}$	$Q^{ap}_{FullObs}$	Q^{ca}_{noPhys}	Q^{ap}_{noPhys}	$Q^{ca}_{FullObs,FullObs}$	$Q^{ap}_{FullObs,FullObs}$
complexity	Q^{ca}	0	0	0	0	0	0	0	0	0	0	0	0
	Q^{ap}	0.13	0	0	0	0	0	0	0	0	0	0	0
	$SR^{ca}@T=1s$	1	1	1	0	0	0	0	0	1	1	0.97	1
	$SR^{ca}@T=2s$	1	0.07	1	0	0	0	0	0	0	1	0.01	0
	$SR^{ca}@T=3s$	1	0.21	1	0	0	0	0	0	0	1	0.01	0
	$SR^{ca}@T=10s$	0	0	0	0	0	0	0	0	0	0	0	0
realism	$Q^{ca}_{FullObs}$	0	0	0	0	0	0	0	0	0	0	0	0
	$Q^{ap}_{FullObs}$	0.20	0.02	0.49	0	0	0	0.13	0	0.09	0	0	0
	Q^{ca}_{noPhys}	0.92	0.47	0.73	0	0.32	0.55	0.44	0.04	0	0.27	0	0.32
	Q^{ap}_{noPhys}	1.0	0.95	0.83	0	0.56	0.94	0.55	0.56	0	0.5	0.67	1.0
	$Q^{ca}_{FullObs,FullObs}$	0	0	0	0	0	0	0	0	0	0	0	0
	$Q^{ap}_{FullObs,FullObs}$	0.20	0.02	0.49	0	0	0	0.13	0	0.09	0	0	0

Table 2: **Evaluation of state-of-the-art RL algorithms on BEHAVIOR** Fully realistic, diverse and complex (row 1): SAC [16] for visuomotor continuous actions (superindex *ca*) performs poorly in all activities; *Complexity analysis* (rows 2-6): reducing complexity (horizon) with temporally extended action primitives (superindex *ap* and gray cells, trained with PPO [17]) or by starting few seconds away from a goal state, lead to some non-zero success rate (SR). *Realism analysis* (rows 7-10): Only by reducing realism in observations and physics, and complexity through action primitives, RL achieves significant success in a handful of the activities.

Diversity in...						
object poses	object instances	ontop	sliced	soaked	stained	cooked
✓	✗	1	0.15	1	1	1
✓	✗	0.825	0	0.935	0.28	0.66
✓	✓	0.46	0	0.925	0.11	0.265

Table 3: **Evaluation of the effect of BEHAVIOR’s diversity:** Results of training agents with SAC [16] in single-predicate activities of increasing diversity; Even in these simple activities, performance degrades quickly indicating that current state-of-the-art cannot cope with the dimensions of diversity spanned in BEHAVIOR

indicate that the diversity in BEHAVIOR goes beyond what current RL algorithms can handle even in simplified activities, and poses a challenge for generalization in embodied AI.

8 Conclusion and Future Work

We presented BEHAVIOR, a novel benchmark for embodied AI solutions of household activities. BEHAVIOR presents 100 realistic, diverse and complex activities with a new logic-symbolic representation, a fully functional simulation-based implementation, and a set of human-centric metrics based on the performance of humans on the same activities in VR. The activities push the state-of-the-art in benchmarking adding new types of state changes that the agent needs to be able to cause, such as cleaning surfaces or changing object temperatures. Our experiments with two state-of-the-art RL baselines shed light on the challenges presented by BEHAVIOR’s level of realism, diversity and complexity. BEHAVIOR will be open-source and free to use; we hope it facilitates participation and fair access to research tools, and paves the way towards a new generation of embodied AI.

Acknowledgments

We would like to thank Bokui Shen, Xi Jia Zhou, and Jim Fan for comments, ideas, and support in data collection. This work is in part supported by Toyota Research Institute (TRI), ARMY MURI grant W911NF-15-1-0479, Samsung, Amazon, and Stanford Institute for Human-Centered AI (SUHAI). S. S. and C. L. are supported by SUHAI Award #202521. S. S. is also supported by the National Science Foundation Graduate Research Fellowship Program (NSF GRFP). R. M-M. and S. B. are supported by SAIL TRI Center – Award # S-2018-28-Savarese-Robot-Learn. S. B. is also supported by a National Defense Science and Engineering Graduate (NDSEG) fellowship, SAIL TRI Center – Award # S-2018-27-Niebles, and SAIL TRI Center – Award # TRI Code 44. S. S. and S. B. are supported by a Department of Navy award (N00014-16-1-2127) issued by the Office of Naval Research (ONR). F. X. is supported by the Qualcomm Innovation Fellowship and Stanford Graduate Fellowship. This article solely reflects the opinions and conclusions of its authors and not any other entity.

References

- [1] Y. Zhu, R. Mottaghi, E. Kolve, J. J. Lim, A. Gupta, L. Fei-Fei, and A. Farhadi. Target-driven visual navigation in indoor scenes using deep reinforcement learning. In *2017 IEEE international conference on robotics and automation (ICRA)*, pages 3357–3364. IEEE, 2017.
- [2] N. Hirose, F. Xia, R. Martín-Martín, A. Sadeghian, and S. Savarese. Deep visual mpc-policy learning for navigation. *IEEE Robotics and Automation Letters*, 4(4):3184–3191, 2019.
- [3] E. Wijmans, A. Kadian, A. Morcos, S. Lee, I. Essa, D. Parikh, M. Savva, and D. Batra. Dd-ppo: Learning near-perfect pointgoal navigators from 2.5 billion frames. *arXiv preprint arXiv:1911.00357*, 2019.
- [4] S. Gupta, J. Davidson, S. Levine, R. Sukthankar, and J. Malik. Cognitive mapping and planning for visual navigation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2616–2625, 2017.
- [5] S. Bansal, V. Tolani, S. Gupta, J. Malik, and C. Tomlin. Combining optimal control and learning for visual navigation in novel environments. In *Conference on Robot Learning*, pages 420–429. PMLR, 2020.
- [6] A. Das, S. Datta, G. Gkioxari, S. Lee, D. Parikh, and D. Batra. Embodied question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–10, 2018.
- [7] L. Yu, X. Chen, G. Gkioxari, M. Bansal, T. L. Berg, and D. Batra. Multi-target embodied question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6309–6318, 2019.
- [8] A. Das, G. Gkioxari, S. Lee, D. Parikh, and D. Batra. Neural modular control for embodied question answering. In *Conference on Robot Learning*, pages 53–62. PMLR, 2018.

- [9] A. Mousavian, A. Toshev, M. Fišer, J. Košecká, A. Wahid, and J. Davidson. Visual representations for semantic target driven navigation. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 8846–8852. IEEE, 2019.
- [10] K. Marino, M. Rastegari, A. Farhadi, and R. Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3195–3204, 2019.
- [11] M. Shridhar, J. Thomason, D. Gordon, Y. Bisk, W. Han, R. Mottaghi, L. Zettlemoyer, and D. Fox. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10740–10749, 2020.
- [12] J. Fu, A. Korattikara, S. Levine, and S. Guadarrama. From language to goals: Inverse reinforcement learning for vision-based instruction following. In *International Conference on Learning Representations*, 2018.
- [13] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sünderhauf, I. Reid, S. Gould, and A. van den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [14] X. Wang, Q. Huang, A. Celikyilmaz, J. Gao, D. Shen, Y.-F. Wang, W. Y. Wang, and L. Zhang. Reinforced cross-modal matching and self-supervised imitation learning for vision-language navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6629–6638, 2019.
- [15] D. Fried, R. Hu, V. Cirik, A. Rohrbach, J. Andreas, L.-P. Morency, T. Berg-Kirkpatrick, K. Saenko, D. Klein, and T. Darrell. Speaker-follower models for vision-and-language navigation. *arXiv preprint arXiv:1806.02724*, 2018.
- [16] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, pages 1861–1870. PMLR, 2018.
- [17] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [18] N. Watters, L. Matthey, M. Bosnjak, C. P. Burgess, and A. Lerchner. Cobra: Data-efficient model-based rl through unsupervised object discovery and curiosity-driven exploration. *arXiv preprint arXiv:1905.09275*, 2019.
- [19] Y. Wu, W. Yan, T. Kurutach, L. Pinto, and P. Abbeel. Learning to manipulate deformable objects without demonstrations. *arXiv preprint arXiv:1910.13439*, 2019.
- [20] A. Billard and D. Kragic. Trends and challenges in robot manipulation. *Science*, 364(6446), 2019.
- [21] O. M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, et al. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20, 2020.
- [22] F. Xia, C. Li, R. Martín-Martín, O. Litany, A. Toshev, and S. Savarese. ReLMoGen: Leveraging motion generation in reinforcement learning for mobile manipulation. In *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020.
- [23] D. Batra, A. X. Chang, S. Chernova, A. J. Davison, J. Deng, V. Koltun, S. Levine, J. Malik, I. Mordatch, R. Mottaghi, M. Savva, and H. Su. Rearrangement: A challenge for embodied ai, 2020.
- [24] L. Weihs, M. Deitke, A. Kembhavi, and R. Mottaghi. Visual room rearrangement. *arXiv preprint arXiv:2103.16544*, 2021.

- [25] C. Gan, S. Zhou, J. Schwartz, S. Alter, A. Bhandwaldar, D. Gutfreund, D. L. Yamins, J. J. DiCarlo, J. McDermott, A. Torralba, et al. The threedworld transport challenge: A visually guided task-and-motion planning benchmark for physically realistic embodied ai. *arXiv preprint arXiv:2103.14025*, 2021.
- [26] X. Puig et al. Virtualhome: Simulating household activities via programs. In *IEEE CVPR*, 2018.
- [27] F. Xia, W. B. Shen, C. Li, P. Kasimbeg, M. E. Tchapmi, A. Toshev, R. Martín-Martín, and S. Savarese. Interactive gibbon benchmark: A benchmark for interactive navigation in cluttered environments. *IEEE Robotics and Automation Letters*, 5(2):713–720, 2020.
- [28] T. Yu, D. Quillen, Z. He, R. Julian, K. Hausman, C. Finn, and S. Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on Robot Learning*, pages 1094–1100. PMLR, 2020.
- [29] S. James, Z. Ma, D. Rovick Arrojo, and A. J. Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 2020.
- [30] Y. Zhu, J. Wong, A. Mandlekar, and R. Martín-Martín. robosuite: A modular simulation framework and benchmark for robot learning. *arXiv preprint arXiv:2009.12293*, 2020.
- [31] Y. Tassa, Y. Doron, A. Muldal, T. Erez, Y. Li, D. d. L. Casas, D. Budden, A. Abdolmaleki, J. Merel, A. Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.
- [32] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [33] U.S. Bureau of Labor Statistics. American Time Use Survey. <https://www.bls.gov/tus/>, 2019.
- [34] G. A. Miller. WordNet: a lexical database. *Communications of the ACM*, 38(11):39–41, 1995.
- [35] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [36] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [37] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.
- [38] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123(1):32–73, 2017.
- [39] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3354–3361. IEEE, 2012.
- [40] R. Goyal, S. Ebrahimi Kahou, V. Michalski, J. Materzynska, S. Westphal, H. Kim, V. Haenel, I. Fruend, P. Yianilos, M. Mueller-Freitag, et al. The "something something" video database for learning and evaluating visual common sense. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5842–5850, 2017.
- [41] G. A. Sigurdsson, A. Gupta, C. Schmid, A. Farhadi, and K. Alahari. Actor and observer: Joint modeling of first and third-person videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7396–7404, 2018.
- [42] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox. Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes. *arXiv preprint arXiv:1711.00199*, 2017.

- [43] R. Martín-Martín, M. Patel, H. Rezatofighi, A. Shenoi, J. Gwak, E. Frankel, A. Sadeghian, and S. Savarese. JrdB: A dataset and benchmark of egocentric robot visual perception of humans in built environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [44] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–970, 2015.
- [45] H. Kitano, M. Asada, Y. Kuniyoshi, I. Noda, E. Osawa, and H. Matsubara. Robocup: A challenge problem for ai. *AI magazine*, 18(1):73–73, 1997.
- [46] T. Wisspeintner, T. Van Der Zant, L. Iocchi, and S. Schiffer. Robocup@home: Scientific competition and benchmarking for domestic service robots. *Interaction Studies*, 10(3):392–426, 2009.
- [47] L. Iocchi, D. Holz, J. Ruiz-del Solar, K. Sugiura, and T. Van Der Zant. Robocup@ home: Analysis and results of evolving competitions for domestic and service robots. *Artificial Intelligence*, 229:258–281, 2015.
- [48] M. Buehler, K. Iagnemma, and S. Singh. *The DARPA urban challenge: autonomous vehicles in city traffic*, volume 56. springer, 2009.
- [49] E. Krotkov, D. Hackett, L. Jackel, M. Perschbacher, J. Pippine, J. Strauss, G. Pratt, and C. Orlowski. The darpa robotics challenge finals: Results and perspectives. *Journal of Field Robotics*, 34(2):229–240, 2017.
- [50] N. Correll, K. E. Bekris, D. Berenson, O. Brock, A. Causo, K. Hauser, K. Okada, A. Rodriguez, J. M. Romano, and P. R. Wurrman. Analysis and observations from the first amazon picking challenge. *IEEE Transactions on Automation Science and Engineering*, 15(1):172–188, 2016.
- [51] C. Eppner, S. Höfer, R. Jonschkowski, R. Martín-Martín, A. Sieverling, V. Wall, and O. Brock. Lessons from the amazon picking challenge: four aspects of building robotic systems. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 4831–4835, 2017.
- [52] M. A. Roa, M. Dogar, C. Vivas, A. Morales, N. Correll, M. Gerner, J. Rosell, S. Foix, R. Memmesheimer, F. Ferro, et al. Mobile manipulation hackathon: Moving into real world applications. *IEEE Robotics & Automation Magazine*, pages 2–14, 2021.
- [53] F. Xia, A. R. Zamir, Z. He, A. Sax, J. Malik, and S. Savarese. Gibson env: Real-world perception for embodied agents. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [54] Manolis Savva*, Abhishek Kadian*, Oleksandr Maksymets*, Y. Zhao, E. Wijmans, B. Jain, J. Straub, J. Liu, V. Koltun, J. Malik, D. Parikh, and D. Batra. Habitat: A Platform for Embodied AI Research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [55] Y. Lee, E. S. Hu, Z. Yang, A. Yin, and J. J. Lim. Ikea furniture assembly environment for long-horizon complex manipulation tasks. *arXiv preprint arXiv:1911.07246*, 2019.
- [56] X. Lin, Y. Wang, J. Olkin, and D. Held. Softgym: Benchmarking deep reinforcement learning for deformable object manipulation. In *Conference on Robot Learning*, 2020.
- [57] D. McDermott, M. Ghallab, A. Howe, C. Knoblock, A. Ram, M. Veloso, D. Weld, and D. Wilkins. Pddl - the planning domain definition language. Technical report, Technical Report 1165, Yale Computer Science, 1998.(CVC Report 98-003), 1998.
- [58] E. Coumans and Y. Bai. Pybullet, a python module for physics simulation for games, robotics and machine learning. *Technical Report*, 2016.
- [59] E. Kolve et al. AI2-THOR: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474*, 2017.

- [60] M. Jordan and A. Perez. Optimal bidirectional rapidly-exploring random trees. Technical Report MIT-CSAIL-TR-2013-021, Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA, August 2013. URL <http://dspace.mit.edu/bitstream/handle/1721.1/79884/MIT-CSAIL-TR-2013-021.pdf>.
- [61] R. Martín-Martín, M. A. Lee, R. Gardner, S. Savarese, J. Bohg, and A. Garg. Variable impedance control in end-effector space: An action space for reinforcement learning in contact-rich tasks. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1010–1017. IEEE, 2019.
- [62] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- [63] A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017.
- [64] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter. Learning agile and dynamic motor skills for legged robots. *Science Robotics*, 4(26), 2019.
- [65] I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino, M. Plappert, G. Powell, R. Ribas, et al. Solving rubik’s cube with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
- [66] T. Haarnoja, S. Ha, A. Zhou, J. Tan, G. Tucker, and S. Levine. Learning to walk via deep reinforcement learning. *arXiv preprint arXiv:1812.11103*, 2018.
- [67] B. R. Kiran, I. Sobh, V. Talpaert, P. Mannion, A. A. Al Sallab, S. Yogamani, and P. Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 2021.
- [68] A. G. Barto and S. Mahadevan. Recent advances in hierarchical reinforcement learning. *Discrete event dynamic systems*, 13(1):41–77, 2003.
- [69] A. S. Vezhnevets, S. Osindero, T. Schaul, N. Heess, M. Jaderberg, D. Silver, and K. Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. In *International Conference on Machine Learning*, pages 3540–3549. PMLR, 2017.
- [70] C. Li, F. Xia, R. Martín-Martín, and S. Savarese. Hrl4in: Hierarchical reinforcement learning for interactive navigation with mobile manipulators. In *Conference on Robot Learning*, pages 603–616. PMLR, 2020.
- [71] C. R. Garrett, T. Lozano-Pérez, and L. P. Kaelbling. Ffrob: An efficient heuristic for task and motion planning. In *Algorithmic Foundations of Robotics XI*, pages 179–195. Springer, 2015.
- [72] A. Aho and J. Ullman. *Foundations of Computer Science*. W. H. Freeman, 1992.
- [73] Upwork Global Inc. Upwork. <https://www.upwork.com/>, 2021. Accessed: 2021-06-16.
- [74] wikiHow, Inc. wikihow. <https://www.wikihow.com>, 2021. Accessed: 2021-06-16.
- [75] Google Alphabet. Blockly. <https://developers.google.com/blockly/>, 2021. Accessed: 2021-06-16.
- [76] B. Shen, F. Xia, C. Li, R. Martín-Martín, L. Fan, G. Wang, S. Buch, C. D’Arpino, S. Srivastava, L. P. Tchapmi, M. E. Tchapmi, K. Vainio, L. Fei-Fei, and S. Savarese. iGibson, a Simulation Environment for Interactive Tasks in Large Realistic Scenes, 2020.
- [77] F. Xiang, Y. Qin, K. Mo, Y. Xia, H. Zhu, F. Liu, M. Liu, H. Jiang, Y. Yuan, H. Wang, et al. SAPIEN: A simulated part-based interactive environment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11097–11107, 2020.

- [78] HTC Corporation. HTC Vive Pro Eye. <https://www.vive.com/us/product/vive-pro-eye/>, 2021. Accessed: 2021-06-16.
- [79] R. Kothari, Z. Yang, C. Kanan, R. Bailey, J. B. Pelz, and G. J. Diaz. Gaze-in-wild: A dataset for studying eye and head coordination in everyday activities. *Scientific reports*, 10(1):1–18, 2020.
- [80] J. Aloimonos, I. Weiss, and A. Bandyopadhyay. Active vision. *International journal of computer vision*, 1(4):333–356, 1988.
- [81] D. H. Ballard. Animate vision. *Artificial intelligence*, 48(1):57–86, 1991.
- [82] J. Bohg, K. Hausman, B. Sankaran, O. Brock, D. Kragic, S. Schaal, and G. S. Sukhatme. Interactive perception: Leveraging action in perception and perception in action. *IEEE Transactions on Robotics*, 33(6):1273–1291, 2017.
- [83] A. Sipatchin, S. Wahl, and K. Rifai. Accuracy and precision of the htc vive pro eye tracking in head-restrained and head-free conditions. *Investigative Ophthalmology & Visual Science*, 61(7):5071–5071, 2020.
- [84] S. Guadarrama, A. Korattikara, O. Ramirez, P. Castro, E. Holly, S. Fishman, K. Wang, E. Gonina, N. Wu, E. Kokiopoulou, L. Sbaiz, J. Smith, G. Bartók, J. Berent, C. Harris, V. Vanhoucke, and E. Brevdo. TF-Agents: A library for reinforcement learning in tensorflow. <https://github.com/tensorflow/agents>, 2018. URL <https://github.com/tensorflow/agents>. [Online; accessed 25-June-2019].
- [85] K. Kang, S. Belkhale, G. Kahn, P. Abbeel, and S. Levine. Generalization through simulation: Integrating simulated and real data into deep reinforcement learning for vision-based autonomous flight. In *2019 international conference on robotics and automation (ICRA)*, pages 6008–6014. IEEE, 2019.

Appendix for BEHAVIOR: Benchmark for Everyday Household Activities in Virtual, Interactive, and Ecological Environments

A.1 Visualizing 100 BEHAVIOR Activities

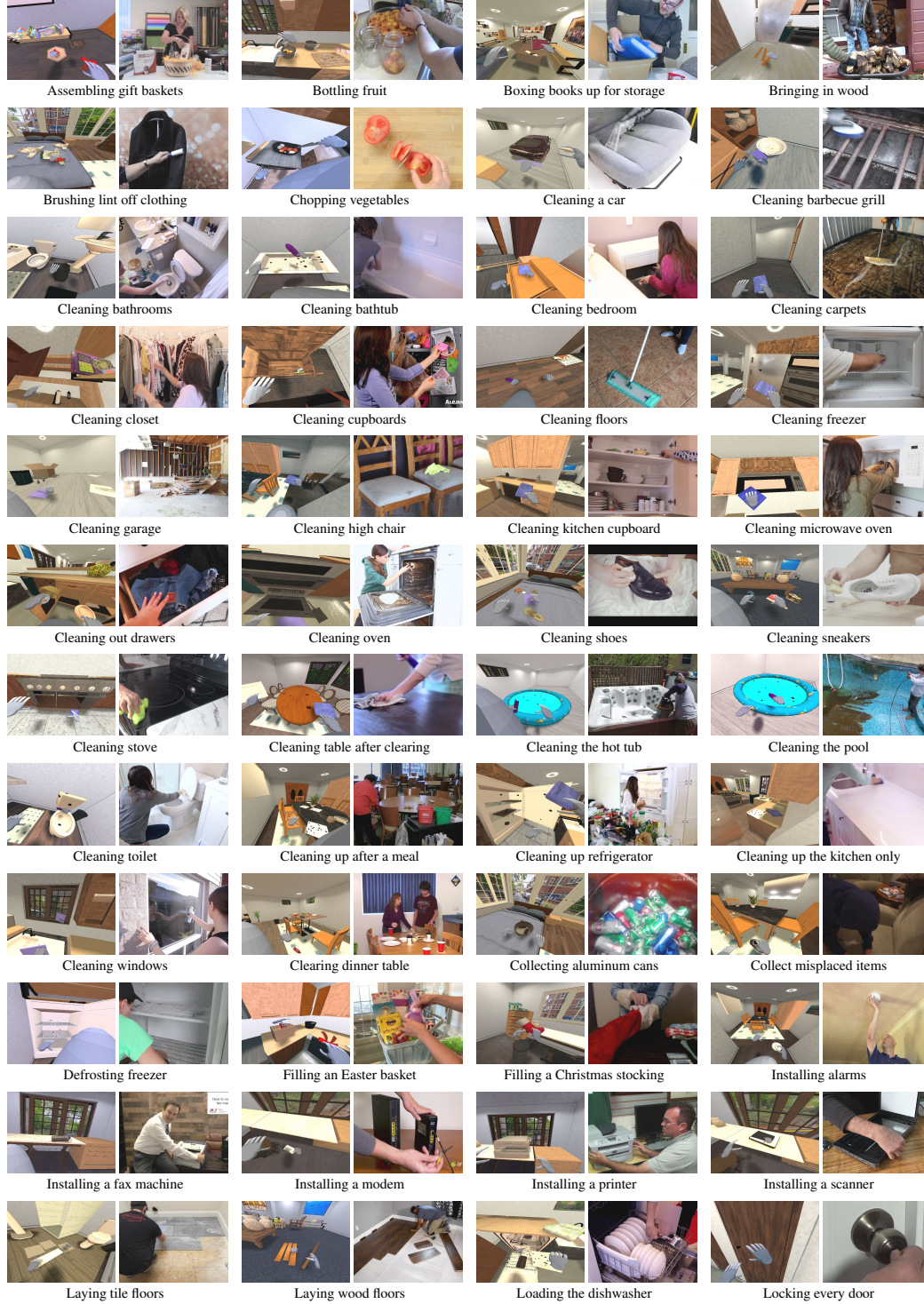


Figure A.1: **BEHAVIOR 100 activities:** Each pair of images depict a frame of the execution of the activity in BEHAVIOR from the agent's perspective in virtual reality (*left*) and the same activity in real-life from a YouTube video (*right*). All activities are selected from the American Time Use Survey [33], and correspond to simulatable household chores relevant in human's everyday life. The set of activities cover common areas like cleaning, maintenance, preparation for social activities, or household management.

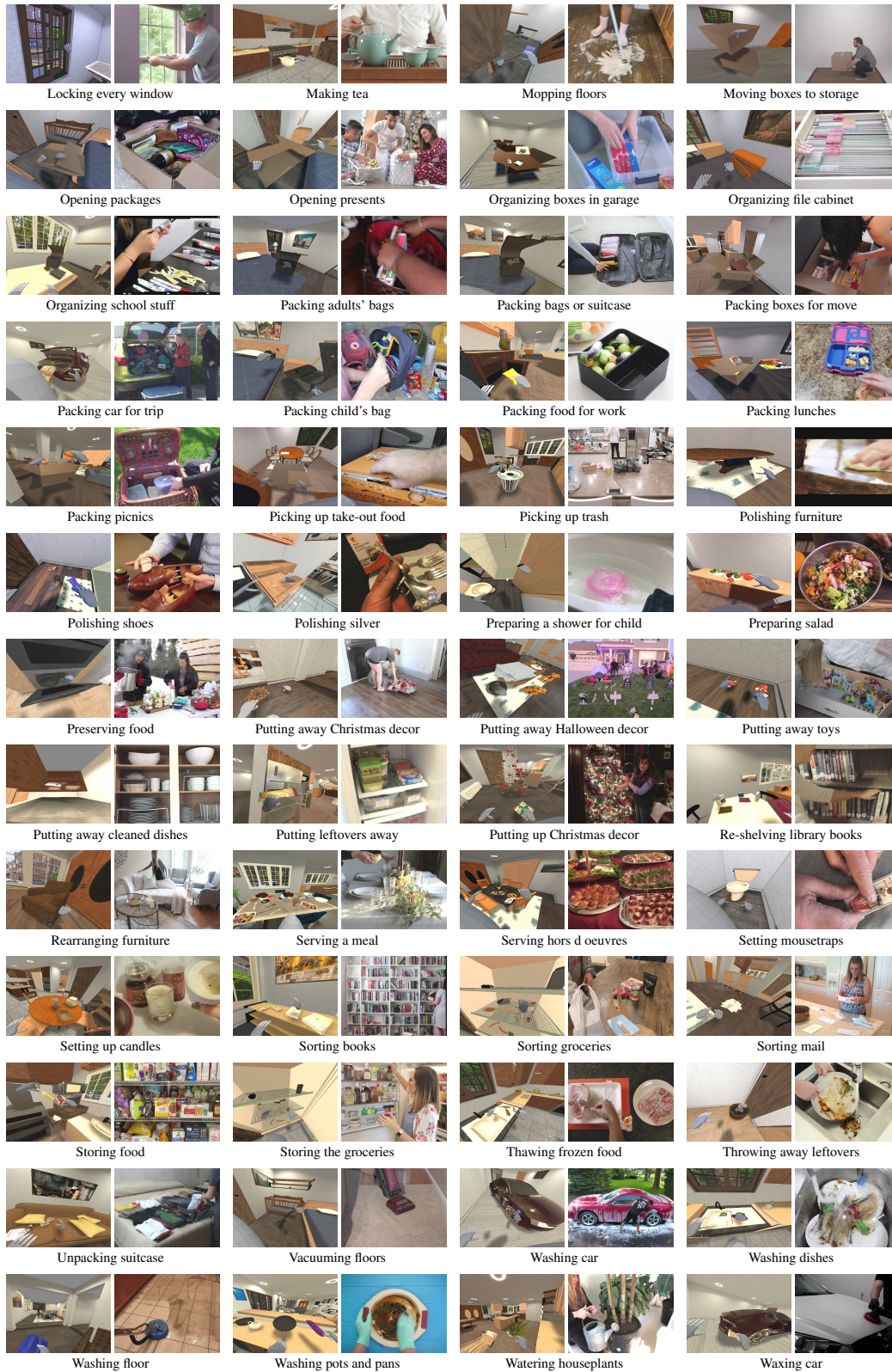


Figure A.1: **BEHAVIOR 100 activities** (cont.)

A.2 Additional Comparison between BEHAVIOR and other Embodied AI Benchmarks

		BEHAVIOR	A2THOR Visual Room Rearrangement Challenge																Gibson
			A2THOR	TDW	Rearrangement T5 (Habitat)	ManipulaTHOR	ArmPointNav	Interactive Gibson Benchmark	VirtualHome	ALFRED	OCRTOC	RLBench	Metaworld	IKEA Furniture Assembly	Robosuite	SoftGym	DespMoD Control Suite	OpenAI Gym Habitat 1.0	
Realism	Realistic activities	Activities match human time-use	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
	Realistic physics	Kinematics dynamics	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	Realistic embodied AI agents	Continuous temperature	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
		Flexible materials	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
	Realistic scenes	Realistic action execution	✓	✗	✓	✗	✓	✓	✗	✗	✗	✓	✓	✓	✓	✓	✓	✓	✓
		Realistic observations	✓	✓	✓	✓	✓	✓	✗	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓
	Realistic object models	Visually realistic	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✗	✗	✓	✗	✗	✓	✓
		Scenes reconstructed from real homes	✓	✗	✗	✓	✗	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✓	✓
	Diverse activities	Visually realistic	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✗	✗	✓	✗	✗	✓	N/A
		Weight, CoM, texture, cook temp	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	N/A
Diversity	Diverse activities	Activities	100	1	1	1	1	2	549	7	5	100	50	1	5	10	28	8	2
		Infinite scene-agnostic instantiation	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	N/A
	Diverse scenes and objects	# Object models	1217	118	112	YCB	150	152	84	101 + YCB	73+	28	80	10	4	4	4	Matterport + Gibson	N/A
		# Scenes / Rooms	15 / 100	- / 120	15 / 90-120	55 static / 30	- / 10	7 / 120	- / 1	- / 1	- / 1	1 / 1	1 / 1	1 / 1	1 / 1	1 / 1	1 / 1	Matterport + Gibson	572 static
	Diverse skills and activity reqs: Benchmark requires manipulating...	objects' pose	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
		agent's global pose	✓	✓	✓	✓	✓	✓	✓	✓	✓	✗	✗	✗	✗	✗	✗	✗	✓
		objects' joint config	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
		objects' geom.	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
		with two hands	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
		objects' functional state (ON/OFF)	✓	✗	✗	✗	✗	✗	✗	✗	✗	✓	✓	✓	✓	✓	✓	✓	✓
		with tools	✓	✗	✗	✗	✗	✗	✗	✗	✗	✓	✓	✓	✓	✓	✓	✓	✓
		object's surface	✓	✗	✗	✗	✗	✗	✗	✗	✗	✓	✓	✓	✓	✓	✓	✓	✓
Complexity	Activity length (steps)	objects' temp.	✓	✗	✗	✗	✗	✗	✗	✗	✗	✓	✓	✓	✓	✓	✓	✓	✓
		Activity length (steps)	300-2000	<100	100-1000	100-1000	<100	100-1000	<100	<100	100-1000	<1000	<100	<100	<100	<100	<100	<100	100-1000
		Obs. per activity	3-34	5	7-9	2-5	2-3	10	1-24	2	5-10	1-2	1-2	1	1-3	1-3	1-3	1	N/A
		# Obj. cats. in act.	2-17	1-5	7-10	2-5	1	1-18	2	1-10	1-2	1-2	1-2	1	1-2	1-3	1-3	1-2	N/A
		Diff. state changes required per activity (see A.2)	2-8	4	4	4	2	1-3	1-7	2-3	1	1-3	1-4	4	1	1-3	1-2	1-2	1
		Benchmark focus: Task-Planning and/or Control	TP+C	TP	TP+C	TP+C	TP+C	C	TP	TP	TP+C	C	TP+C	C	C	C	C	C	C
	# Human VR demos		400	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
			✗	✗	✗	✗	✗	✗	✗	✗	✗	✓	✗	✗	✗	✗	✗	✗	✗

Table A.1: Comparison between BEHAVIOR and other existing benchmarks for embodied AI. Expanded version of Table 1.

A.3 Defining BEHAVIOR Activities

This section includes additional information on how we define the 100 activities of BEHAVIOR, including details on 1) the process to select them from the American Time Use Survey [33] (ATUS), 2) BDDL, the predicate logic language to define them, 3) the crowdsourcing process to generate definitions (initial and goal conditions) for the activities, 4) and real BDDL examples of the generated definitions.

A.3.1 Selection of 100 Activities for BEHAVIOR

Our activities are extracted from the American Time Use Survey [33] that contains more than 2200 activities Americans spend their everyday time on. To select a subset for BEHAVIOR, we follow a set of criteria: **i) semantic diversity**: we select activities that span a wide range of semantic areas, from cleaning to food preparation, or repairing (see Fig. A.2a); **ii) diversity in the required state changes in the environment**: we select activities that requires manipulating different properties of the objects, their pose, temperature, cleanliness level, wetness level... (see Fig. A.2, b and c); and **iii) simulation feasibility**: given the current state of simulation environments, we select for BEHAVIOR activities that can be realistically simulated entirely in an indoor environment, involving only objects, most of them rigid or articulated, excluding activities outdoors, interactions with other humans or animals, or heavy simulation of flexible materials and fluids. The resulting full list of 100 BEHAVIOR activities selected can be visualized in Fig. A.1. They cover a large variety of activities such as cleaning (CleaningBathtub, CleaningTheKitchenOnly, WashingPotsAndPans), installing (InstallingAScanner, InstallingAlarms), waxing/polishing (PolishingSilver, WaxingCarsOrOtherVehicles), tidying (PuttingAwayToys, PuttingDishesAwayAfterCleaning), packing/assembling (PackingPicnics, AssemblingGiftBaskets), and preparing food (PreservingFood, ChoppingVegetables). Fig. A.2 depict statistics of the selected activities supporting that they approximate the semantic distribution of activities in the time use survey, and that they require a broad set changes in the environment. As comparison, Rearrangement tasks [23] and related benchmarks focus on activities that can be achieved by agent's pose (navigation), object's pose

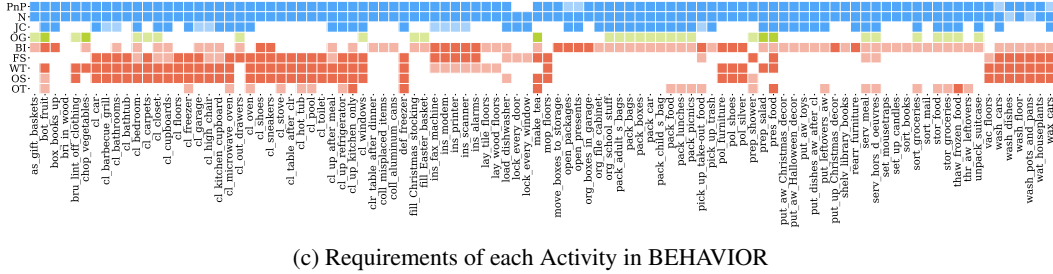
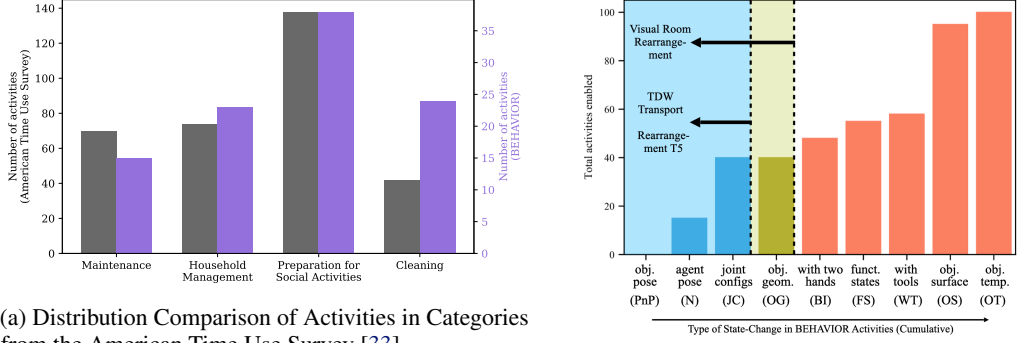


Figure A.2: **Statistics of the 100 activities in BEHAVIOR:** a) Distribution of simulatable activities in the American Time Use Survey (left axis) and BEHAVIOR (right axis) based on categories from the survey – BEHAVIOR covers a realistic distribution of activities; b) Cumulative visualization of activities enabled by different types of state changes in BEHAVIOR with comparison to recent prior work – based on requirements, some activities could be considered transport/rearrangement (blue) or visual-room rearrangement (blue and green), while others are out of their scope (red); c) We visualize the specific requirements for each of the BEHAVIOR activities, with the same coloring scheme as in b). Activities in BEHAVIOR present significantly more diverse requirements than prior work focused on transport/rearrangement tasks [23, 25, 24] enabling the evaluation of more general embodied AI solutions.

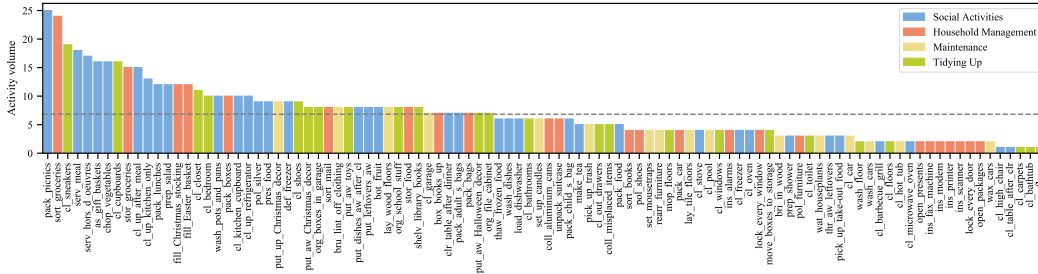


Figure A.3: **Activity volume in BEHAVIOR:** The number of literals in flattened goal conditions (volume, see Sec. A.3.2) provides a measure of the complexity of the activity and its length/horizon. The volume in BEHAVIOR activities span from one to 25 literals, very long horizon activities. Activities with one literal are often still long-horizon as they may require cleaning large surfaces (e.g. vacuumFloors or cleaningBathtub)

(pick-and-place), and joint configuration of articulated objects. VisualRoom Rearrangement [24] includes objects that can be broken (changing object geometry).

A.3.2 BDDL– BEHAVIOR Domain Definition Language

In BEHAVIOR, activities are defined using a new predicate logic language, BDDL, BEHAVIOR Domain Definition Language. BDDL creates a logic-symbolic counterpart to the physical state simulated by iGibson 2.0 through a set of logic functions (predicates). In this way, BDDL defines a set of symbols grounded into simulated objects and their states. The goal of BDDL is to enable defining activities in a unique, unifying language that connects to natural language to facilitate interpretability.

In this section, we provide additional information about the similarities and differences between BDDL and PDDL [57], a full description of the BDDL elements, syntax and grammar, and information about evaluation, grounding and “flattening” conditions, and the concept of “activity volume”.

BDDL vs. PDDL: While similar in name, BDDL is inspired by the Planning Domain Definition Language (PDDL) [57] but strongly divergent. Both are derived from predicate logic and share a common logic-symbolic structure. However, their goal and requirements are significantly different: while PDDL’s main objective is to define a complete space for symbolic planning without any necessary connection to a physical world, BDDL’s goal is to provide a diverse and fully-grounded symbolic representation of physical states to define activities as pairs of initial and goal logical conditions. Therefore, PDDL requires to define additional symbols for agent’s actions, while BDDL is only a representation of the state: agents act in the physical simulation to achieve the activities in BEHAVIOR. To facilitate the adoption of BDDL as standard language to define activities in embodied AI, we assume the well-known syntax of PDDL for states.

BDDL Syntax: In BDDL, we consider the following syntactic elements, a subset of the syntax of predicate logic defined in Aho and Ullman [72]:

- **Predicate:** logic function that takes as input one (unary) or two (binary) objects and returns a boolean value. Examples in BDDL: `ontop`, `stained`, `cooked`.
- **Variable:** element in a logical expression representing an object of the indicated category, always bound by a quantifier. Categories in BDDL are defined by WordNet [34] synsets (semantic meaning), indicated by the label structure `categoryName.n.synsetEntry`. A variable is then indicated by a character `?` followed by the category. Examples in BDDL: `?apple.n.01`, `?table.n.02`.
- **Constant:** ground term, i.e., variable linked to a specific instance of an object. In BDDL, constants are identified by a numerical id suffix (`_n`) appended to the variable name. Examples in BDDL: `apple.n.01_1`, `table.n.02_3`.
- **Category:** attribute of a constant or variable indicating the class of object it belongs to, and therefore which predicates it can be given as input to (e.g., `cooked`, `sliceable`). Examples in BDDL: `apple.n.01`, `table.n.02`.
- **Type:** synonymous with **category**, conventional for PDDL and therefore defined for BDDL.
- **Argument:** variable or constant used as input in a predicate.
- **Atomic formula:** single predicate with an appropriate number of arguments. Example in BDDL: `(ontop(apple.n.01_1, table.n.02_1))`
- **Logic operator:** Function mapping logical expressions to new logical expressions. In BDDL we include all four propositional logic operators: `and` (\wedge), `or` (\vee), `not` (\neg), `if` (\Rightarrow), and `iff` (\Leftrightarrow).
- **Quantifier:** Function of a variable to map existing logical expressions to new logical expressions. In BDDL we include the standard universal quantification (\forall), and existential quantification (\exists), and additional operators: `for_n`, `for_pairs`, `for_n_pairs` (definitions below).
- **Logical expression:** expression obtained by composing atomic formulas with logical operators. Example in BDDL: `(and (ontop(apple.n.01_1, table.n.02_1)) (forall (?apple.n.01 - apple.n.01) cooked(apple.n.01)))`
- **Initial condition:** set of atomic formulas that are guaranteed to be `True` at the beginning of all instances of the associated BEHAVIOR activity. See examples in Listings 1 and 2.
- **Goal condition:** logical expression that must be `True` for the associated BEHAVIOR activity to be considered successfully executed. See examples in Listings 1 and 2.
- **Literal:** atomic formula or negated atomic formula. Example in BDDL: `not(ontop(apple.n.01_1, table.n.02_1))`
- **Fact:** ground atomic formula evaluated on the current state of the simulated world and returning a Boolean. Example in BDDL: `ontop(apple.n.01_1, table.n.02_1) = True`.
- **State:** set of facts about the current state of the simulated world providing a logical representation that can be evaluated wrt. the goal condition.

Initial and final conditions for household activities could be expressed using the aforementioned first order logic syntax combined with BEHAVIOR’s predicates. However, our activities are defined by non-technical annotators through a crowdsourcing procedure. The annotators are not required to have background knowledge in formal logic or computer science. To facilitate their work, we include the following additional non-standard quantifiers:

- `for_n`: for some non-negative integer n and some object category C , the child condition must hold true for at least n instances of category C
- `for_pairs`: for two object categories $C1$ and $C2$, the child condition must hold true for some one-to-one mapping of object instances of $C1$ to object instances of $C2$ that covers all instances of at least one category
- `for_n_pairs`: for some non-negative integer n and two object categories $C1$ and $C2$, the child condition must hold true for at least n pairs of instances of $C1$ and instances of $C2$ that follow a one-to-one mapping.

Following the format of PDDL [57], in BDDL we consider two types of “files”: a domain file shared for all activities, and problem files for each activity. The domain file defines all possible predicates, including object categories (corresponding in BEHAVIOR to categories from WordNet) and semantic symbolic states. Each activity in BEHAVIOR is defined by a different problem file that includes the object instances involved in the activity (categorized), the conditions for initial and final states.

Evaluating Logical Expressions: For a logical expression to be evaluated, we first decompose recursively it into subcomponents at the operators and quantifiers until we obtain a hierarchical structure of atomic formulae. Each atomic formula is composed of a predicate and arguments, i.e., a mathematical relationship on the simulated object(s) properties passed as arguments. For example, the atomic formula `(cooked apple.n.01_1)` is evaluated by checking the relevant thermal information of the simulated object `apple.n.01_1`. For details on the implementation of each predicate, see the attached cross-submission on the simulator iGibson 2.0. Once the atomic formulae have been evaluated into facts with queries to the grounding simulated object states, we compose the facts through the logical operators to obtain the overall binary result of the whole expression. The BDDL symbolic definition of logical expressions creates flexibility: see Fig. A.6 bottom row for examples of multiple correct solutions accepted by the same BDDL specification.

Instantiating and Grounding Initial Conditions: The initial conditions of an activity in BEHAVIOR are defined at the beginning of each BDDL problem file. They include a list of object constants and a set of ground literals based on these constants. Instances of a BEHAVIOR activity are simulated physical states that fulfill all literals in the conditions. In our implementation of BDDL in iGibson 2.0, the initial conditions are instantiated in the simulated state by assigning all object constants to physical objects of the appropriate category, either matching to physical objects already in the simulated scene or instantiating new ones in the locations specified by the binary atomic formulae (e.g., `ontop`, `inside`, etc.). The ground unary literals are satisfied by setting the physical states of the simulated objects according to the value their associated constants as given in the initial condition (e.g., `(not (cooked (chicken.n.01_1)))` sets the temperature of the associated `chicken.n.01_1` instance to a value that corresponds to `uncooked`). Our instantiation of BDDL in iGibson 2.0 provides a sampling mechanism of unary and binary predicates that can generate potentially infinite variations of each set of initial conditions (more details in the iGibson 2.0 submission attached as supplementary) See Fig. A.6 top row for examples of multiple instantiations from the same BDDL specification.

“Flattening” a Goal Condition in an Activity Instance: BDDL provides a powerful mechanism to define the goal conditions in BEHAVIOR in their general form, e.g., `forall(?toy.n.01 - toy.n.01) inside(toy.n.01, box.n.01)`. As logical expression, BDDL goal conditions are independent of the concrete objects and the scene, and thus valid to all instances and capturing all variants of the solution. However, there are situations where grounding the goal conditions in the concrete instance of the activity at hand is helpful to understand the complexity (i.e., compute the activity volume), and the incremental progression towards the goal (i.e., compute the success score). Following on the previous example, for a possible goal condition of `PickingUpToys`, the activity’s complexity would be very different when the condition is applied on an activity instance (scene) with 100 toys or with only 1 toy. We call goal condition “flattening” in an activity instance to the process of generating possible ground states of a specific simulated world fulfilling a condition. Flattening involves decomposing the nested structure of operators and quantifiers in the logical expression into a flat structure of disjointed conjunctions C_i of ground literals l_{j_i} , $\bigvee_{C_i} \bigwedge l_{j_i}$, and grounding the literals

in all possible ways in the given instance. The final output of the flattening process is a list of *options*, each of which is a list of ground literals that would satisfy the goal condition. Because disjunctions, existential quantifiers, and `for_n`, `for_n_pairs` are satisfied as soon as one/ n of their children is/are satisfied, our implementation of the flattening process for BDDL in BEHAVIOR acts lazily,

generating only the minimal number of literals to fulfill each component of the goal condition. This prioritizes efficient solutions without losing any recall of possible solutions.

Activity Volume: The result of flattening a goal condition in an activity instance is a list of possible options to accomplish the activity, each option being a list of ground atomic literals. We define the *activity volume* as the length of the shortest flattened goal option for a given activity in a concrete instance. The activity volume provide a measure of the logical complexity of an activity, i.e., the number of atomic formulae that the agent needs to fulfill. For our previous example for `PickingUpToys`, the activity would have a volume of N for an activity instance with N toys, indicating the different complexity for an instance with 100 or with 1 toys.

A.3.3 Crowdsourcing the Annotation of Activities

Thanks to the connection in BDDL between the logical predicates and language semantics, BEHAVIOR activity definitions can be generated through crowdsourced annotation from non-experts workers, i.e., without background in computer science or logic. Through a visual interface, annotators can easily generate activity definitions in BDDL that reflect their idea of what the core of the activity is, and that are guaranteed to be simulatable in iGibson 2.0. We crowdsourced the generation of activity definitions to ensure that we do not introduce researcher biases in the design. The annotator pool was sourced from Upwork [73], limiting to Upwork freelancers based in the United States of America to maintain consistency and familiarity with ATUS activities. Each annotator was given a salary of \$15 per annotation, roughly \$20-30 per hour. Because the above process constitutes a complex annotation task, we developed a custom interface to guide and facilitate annotators’ work, and guarantee simulatable output.

Annotation Process and Interface: The annotation procedure is as follows. First, the annotator is presented with a BEHAVIOR activity label. When necessary, we modify the original labels to add a numerical context, e.g., “packing **four** lunches” for the original BEHAVIOR activity “packing lunches”. Then, the annotator reads the annotation instructions and enters the label into the interface. As response, the interface prompts the annotator to select one or more rooms that are relevant for the activity, and to choose objects already present in these rooms that are relevant to the activity (Fig. A.4 (a)). The annotators then select small objects from the BEHAVIOR Dataset of Objects organized in the WordNet hierarchy (Fig. A.4 (b)). To facilitate the annotation, instead of presenting the hierarchy for the entire BEHAVIOR Dataset, we preselect the most possible categories per activity based on a parsing procedure on how-to articles retrieved online, primarily from wikiHow [74] (see Sec. A.5). However, annotators can access the full hierarchy if the preselected items are not sufficient.

After this first phase to select activity-relevant objects, the annotator enters the second phase to annotate initial and goal conditions. First, they are introduced to a block-based, visual tool to generate BDDL (Fig. A.4 (c)) built on Blockly [75], which makes generating logical expressions intuitive and accessible to people without a programming background [26]. They use the tool to generate initial and goal conditions based on their concept of the activities (Fig. A.4 (d)). The resulting definitions have several guarantees: 1) they only use objects from the BEHAVIOR Dataset of Objects, 2) they only apply logical predicates to objects in a semantically meaningful manner (e.g., `cooked` can only be applied to `cookable` objects such as food). This is because blocks’ predicate fields are conditioned on entered categories that have been annotated with possible predicates in a separate manual WordNet annotation. 3) They will be in syntactically correct BDDL, through the implemented translation from Blockly. 4) They will not contain free variables or logically unsatisfiable conditions. 5) It will be possible to simulate them physically in at least three simulated scenes from iGibson 2.0. To guarantee feasibility, we assigned three possible home scenes from iGibson 2.0 to each activity and let the annotators evaluate the feasibility of their conditions at any point by clicking a button “Check feasibility” (Fig. A.4 (e)). The request will send the BDDL definition to up to three iGibson 2.0 simulators on a remote server that will attempt to sample the initial conditions and check if the goal conditions are feasible, returning real-time feedback to the annotators to correct any unfeasible condition. With the crowdsourcing procedure we obtain two alternative definitions per activity that are guarantee to be feasible in at least three simulated scenes.

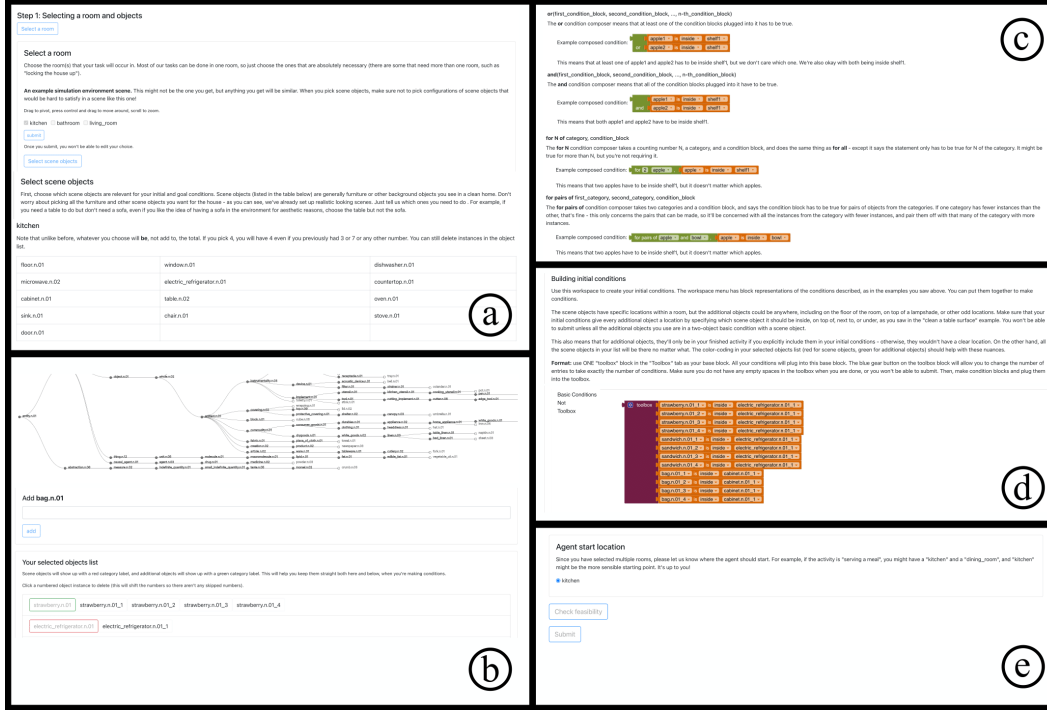


Figure A.4: Sections of the interface given to activity definition annotators. ④ shows selection of relevant rooms and scene objects. For the purpose of creating definitions compatible with multiple iGibson 2.0 scenes and likely to fit with new scenes, annotators were allowed to pick scene objects from the intersection of object sets in three pre-selected scenes. ⑤ shows selection of additional objects that would be added to the scene during activity instantiation, sourced from wikiHow [74] and taxonomized via WordNet [34]. ③ shows examples of the Blockly [75] version of BDDL, and ① shows the prompt for initial conditions and an example for a simple “packing lunches” definition. ② shows the decision of the agent’s start point and the interface for “checking feasibility”, i.e. confirming that the BDDL is syntactically correct, the initial and goal conditions are satisfiable, and the set-up can be physically simulated in iGibson 2.0 by attempting a sampling in an iGibson 2.0 instance on a remote server. Not shown: introductory instructions, goal condition prompt and example (similar to initial condition), some BDDL blocks, remote server communication. Full interface available online: <http://verified-states.herokuapp.com> (server currently disabled).

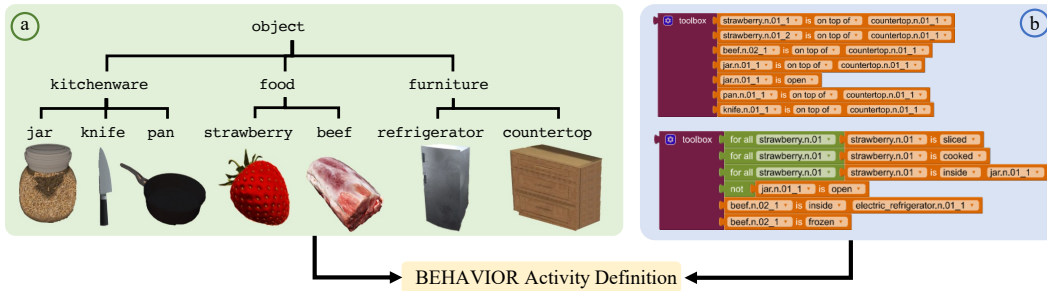


Figure A.5: **Activity annotation process for `preserving_food`**: a) annotators select objects from the WordNet organized BEHAVIOR Dataset of Objects; b) the selected objects are composed into logical predicates in BDDL for initial and final conditions using a visual interface derived from Blockly [75]; the result is a BDDL definition of the activity as logic predicates connected by logic operators and quantifiers, grounded in simulatable objects with physical properties

A.3.4 Example Definitions

```
(define
  (problem packing_lunches_1)
  (:domain igibson)

  (:objects
    shelf.n.01_1 - shelf.n.01
    water.n.06_1 - water.n.06
    countertop.n.01_1 - countertop.n.01
    apple.n.01_1 - apple.n.01
    electric_refrigerator.n.01_1 -
      electric_refrigerator.n.01
    hamburger.n.01_1 - hamburger.n.01
    basket.n.01_1 - basket.n.01
  )

  (:init
    (ontop water.n.06_1 countertop.n.01_1)
    (inside apple.n.01_1
      electric_refrigerator.n.01_1)
    (inside hamburger.n.01_1
      electric_refrigerator.n.01_1)
    (ontop basket.n.01_1 countertop.n.01_1)
    (inroom countertop.n.01_1 kitchen)
    (inroom electric_refrigerator.n.01_1
      kitchen)
    (inroom shelf.n.01_1 kitchen)
  )

  (:goal
    (and
      (for_n_pairs
        (1)
        (?hamburger.n.01 - hamburger.n.01)
        (?basket.n.01 - basket.n.01)
        (inside ?hamburger.n.01 ?basket.n.01)
      )
      (for_n_pairs
        (1)
        (?basket.n.01 - basket.n.01)
        (?water.n.06 - water.n.06)
        (inside ?water.n.06 ?basket.n.01)
      )
      (for_n_pairs
        (1)
        (?basket.n.01 - basket.n.01)
        (?apple.n.01 - apple.n.01)
        (inside ?apple.n.01 ?basket.n.01)
      )
      (forall
        (?basket.n.01 - basket.n.01)
        (ontop ?basket.n.01 ?countertop.n.01_1)
      )
    )
  )
)
```

Listing 1: packing_lunch

```
(define
  (problem serving_hors_doeuvres_1)
  (:domain igibson)

  (:objects
    tray.n.01_1 tray.n.01_2 - tray.n.01
    countertop.n.01_1 - countertop.n.01
    oven.n.01_1 - oven.n.01
    sausage.n.01_1 sausage.n.01_2 - sausage.n.01
    cherry.n.03_1 cherry.n.03_2 - cherry.n.03
    electric_refrigerator.n.01_1 -
      electric_refrigerator.n.01
  )

  (:init
    (ontop tray.n.01_1 countertop.n.01_1)
    (ontop tray.n.01_2 countertop.n.01_1)
    (inside sausage.n.01_1 oven.n.01_1)
    (inside sausage.n.01_2 oven.n.01_1)
    (inside cherry.n.03_1
      electric_refrigerator.n.01_1)
    (inside cherry.n.03_2
      electric_refrigerator.n.01_1)
    (inroom oven.n.01_1 kitchen)
    (inroom electric_refrigerator.n.01_1
      kitchen)
    (inroom countertop.n.01_1 kitchen)
  )

  (:goal
    (and
      (exists
        (?tray.n.01 - tray.n.01)
        (and
          (forall
            (?sausage.n.01 - sausage.n.01)
            (ontop ?sausage.n.01 ?tray.n.01)
          )
          (forall
            (?cherry.n.03 - cherry.n.03)
            (not
              (ontop ?cherry.n.03 ?tray.n.01)
            )
          )
        )
      )
      (exists
        (?tray.n.01 - tray.n.01)
        (and
          (forall
            (?cherry.n.03 - cherry.n.03)
            (ontop ?cherry.n.03 ?tray.n.01)
          )
          (forall
            (?sausage.n.01 - sausage.n.01)
            (not
              (ontop ?sausage.n.01 ?tray.n.01)
            )
          )
        )
      )
    )
  )
)
```

Listing 2: serving_hors_doeuvres

In Listings 1 and 2, we include two examples of activity definitions (initial and goal conditions) in BDDL. They are generating by mapping the input from crowdsourcing workers in our Blockly-like interface into BDDL language. The activities include several objects and predicates in the initial and goal specifications.

A.4 iGibson 2.0

While BEHAVIOR is agnostic to the underlying simulator, we provide a fully functional instantiation in iGibson 2.0. The details about iGibson 2.0 can be found in the cross-submission included in the supplementary material. Here we summarize its most important features in relation to BEHAVIOR.

Agents – Realistic Sensing and Actuation: We implement in iGibson 2.0 the two embodied agents mentioned in Sec. 5 to perform BEHAVIOR activities: a bimanual humanoid and a Fetch



Figure A.6: **BDDL Initial and Goal Conditions:** Our implementation in iGibson 2.0 can generate diverse valid activity instances from each BDDL definition (top row), and detect all successful variations of the solution (bottom row), promoting diversity and semantically-meaningful activities

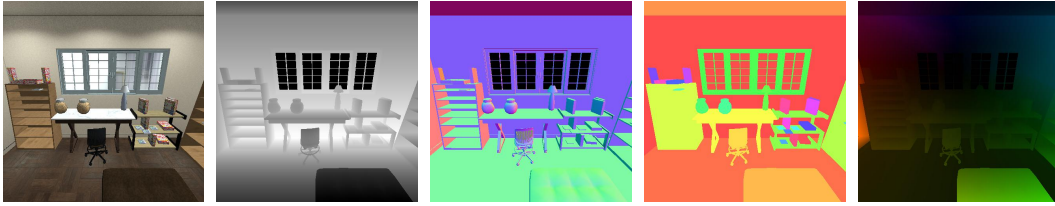


Figure A.7: **Virtual visual sensor signals generated by iGibson 2.0:** Color images are generated with a high-quality physics-based rendering procedure (PBR), exploiting the annotation of material (roughness, metallic) of all surfaces in our objects and scenes. iGibson 2.0 is able to generate RGB, depth, surface normals, semantic segmentation, instance segmentation, optical flow, scene flow and lidar (1-line and 16-line) sensors signals. Here we visualize a subset of those sensor signals, namely RGB, depth, surface normal, instance segmentation and optical flow.

robot. Agents embodying the bimanual humanoid must control 24 degrees of freedom (DoF) to navigate, move and grasp (1 continuous DoF) with the hands, and move the pose of the head that controls the camera point of view. This is the embodiment used by humans in VR. Agents embodying the Fetch robot control 12 or 13 DoF: the navigating motion of the base, the pose of the end-effector (6 DoF), or alternatively, the joint configuration of the arm (7 DoF), one prismatic joint to grasp and release, and pan/tilt motion of the head that moves the cameras.

The sensors used by humanoid agent and Fetch leverages the realistic sensor simulation from iGibson 2.0. iGibson 2.0 features a physically-based renderer that can generate highly photorealistic RGB camera images, as well as other modalities, including depth, surface normal, semantic segmentation, instance segmentation, lidars, scene flows and optical flows. Fig. A.7 highlights a subset of the generated sensor signals.

In terms of actuation, the actions are simulated accurately in pyBullet [58], the physics engine used by iGibson 2.0, with a very small physics simulation timestep of $\frac{1}{300}$ s. The small physics timestep can reduce physics simulation artifacts, such as objects clipping into each other, increasing realism.

Condition Checking and Sampling: The implementation of BEHAVIOR in iGibson 2.0 allows activities to be initialized, executed, and checked for completion. Given an activity definition in the BDDL, BEHAVIOR and iGibson 2.0 interface to generate a valid instance of the activity that satisfies the given object list and initial conditions. This mechanism can generate potentially infinite variation of scenes, objects and initial states to create different activity instances. In the generation of an activity instance, the goal conditions are checked for feasibility, avoiding the generation of activity

instances that cannot lead to successful executions (see Fig. A.6, top). iGibson 2.0 implements all necessary checking functionalities for the logical states. These checking functions execute in realtime together with the physical simulation and rendering, enabling live feedback to the agents for task completion and capturing all possible valid solutions (Fig. A.6, bottom). For more information about the condition checking and sampling, please refer to the concurrent submission iGibson 2.0 paper included as part of the supplementary material.

Implementation of the Action Primitives: To facilitate the development of solutions and to study the effect of the activity complexity on the performance of embodied AI algorithms, we provide action primitives implemented in iGibson 2.0 and that can be used in BEHAVIOR. The action primitives are temporally extended actions. We implemented six action primitives, namely `navigate_to(obj)`, `grasp(obj)`, `place_onTop(obj)`, `place_inside(obj)`, `open(obj)`, `close(obj)`. Each primitive can be applied relative to objects in the scene. For each action primitive, we implemented two variants. The first variant is “fully-simulated motion primitive”, where we first check the feasibility of the target configuration, and then plan a full valid path between the initial and the target configurations with a sampling based motion planner [60]. The second variant is “partially-simulated motion primitive”, where we only check for feasibility of the desired final configuration, and directly set the state of the world (agent and objects) to this desired configuration. This can be highly unrealistic as we do not verify if there is a valid path between the initial and the final configurations. The purpose of partially simulated motion primitive is to reduce the computation during RL training and to measure the relative complexity of generating full interactions vs. just finding the sequence of states to achieve an activity. Note that for both partially and fully simulated motion primitives, privileged information is given to the agent and the motion planner. For example, the agent knows how many activity-relevant objects are in the scene, and the motion planner knows the full geometry of the environment.

For the implementation of partially-simulated motion primitives, we only perform feasibility check when attempting to perform an action. For example, when trying to `navigate_to` an object, we will randomly sample points around the object and attempt to place the agent there: the goal is to find a collision-free location to place the agent. The second type of feasibility check is reachability: when attempting to `grasp`, `open` or `close` an object, we will check the distance from the hand to the closest point of the object is smaller than the arm length. When we `place` an object `inside` or `onTop` of an object, we use the sampling functionality available in iGibson 2.0.

For the implementation of fully-simulated motion primitives, in addition to the feasibility check, we attempt to plan and execute a collision free. We treat all objects as obstacles except for the objects given as argument for the primitive (e.g., objects that need to be picked up, receptacles that need to be opened), and plan a collision-free path from the start configuration to the target configuration. We use Bidirectional RRT [60] for motion planning and execute the motion with position control. In our experiments, we found that fully simulated motion primitives have much lower success rate than partially simulated motion primitives (Table 2) indicating that the difficulty in BEHAVIOR arises from solving the entire interaction rather than deciding on the strategy at a task-level. Our partially-simulated primitives and other benchmarks that do not simulate the full interaction, bypass this critical challenge.

Runtime performance of iGibson 2.0: iGibson 2.0 improved performance when compared to iGibson 1.0 [76], with optimizations on both physics and rendering. To evaluate the performance of iGibson 2.0 in BEHAVIOR activities, we benchmarked the different phases of each simulation step. We benchmark the activities in “idle” setting, which means we initialize the activity, and runs the simulation and condition checking loop. The agent applies zero actions and stays still. We benchmarked in two conditions using the same action time step of t_a but different physics time step of t_s , leading to slightly different reality in the physics simulation. The action step is the simulated-time between agent’s actions, while the physics time step is the simulated-time interval that the kinematics simulator (pyBullet) uses to integrate forces and compute the new kinematic states. We execute n_s queries to the simulator between agent actions, with $n_s = t_a/t_s$. The first condition we evaluate uses action time step $t_a = \frac{1}{30}$ s and physics time step of $t_s = \frac{1}{300}$ s, which creates high-fidelity physics simulation. The second condition uses action time step of $t_a = \frac{1}{30}$ s and physics time step of $t_s = \frac{1}{120}$ s, which has slightly lower physics fidelity, but has better performance and is sufficient for RL training. Both settings are benchmarked on a computer with Intel 5930k CPU and Nvidia GTX 1080 Ti GPU, in a single process setting, rendering 128×128 RGB-D images.

	bringing_in_wood	re-shelving_library_books	laying_tile_floors
Number of Objects	134	144	216
Simulation steps per second ($@t_s = \frac{1}{300}s / @t_s = \frac{1}{120}s$)	59 / 74	51 / 68	36 / 47
Kinematic State Update Time [ms] ($@t_s = \frac{1}{300}s / @t_s = \frac{1}{120}s$)	7.4 / 3.5	9.4 / 4.2	12.6 / 5.7
Non-kinematic State Update Time [ms]	3.4	3.4	5.2
Rendering Time [ms]	5.8	6.1	9.3
Logical Condition Checking Time [ms]	0.4	0.4	0.6

Table A.2: Benchmarking Simulation Time for BEHAVIOR Activities in iGibson 2.0

As shown in Table A.2, for the highest-fidelity physics setup, we can achieve 36-59 steps per second, 47-71 steps per second with larger simulated timestep, even in a very large scene with 100-200 movable objects, and with all the physical and logical states evaluated at each step. This frequencies provide pleasant experience in virtual reality. However, it only provide a $\times 2$ acceleration over clock-time to train RL agents. To increase the frequency in simulation and reduce the training time, we are exploring the parallelization of simulation and rendering and the more aggressive “sleep” of non-interacted objects.

A.5 BEHAVIOR Dataset of Objects

In order to instantiate BEHAVIOR activities in iGibson 2.0, we created a new dataset of everyday objects, the BEHAVIOR Dataset of Objects. To guide the selection of object categories, we analyze how-to articles, primarily WikiHow [74], explaining how to perform the activities included in BEHAVIOR. Specifically, we extract nouns of tangible objects from these articles that are activity-relevant, map them to WordNet synsets, and then purchase 3D models of these object categories from online marketplaces such as TurboSquid. This procedure allowed us to provide activity annotators and VR demonstrators with the most frequent objects necessary for the activities (see Fig. A.8).

The diversity of BEHAVIOR activities naturally leads to the diversity of the object dataset. In total, we curate 1217 object models across 391 object categories, to support 100 BEHAVIOR activities. The categories range from food items to tableware, from home decorations to office supplies, and from apparel to cleaning tools. In Fig. A.8, we observe that the BEHAVIOR Dataset of Objects cover a wide range of object categories.

To maintain high visual realism, all object models include material information (metallic, roughness) that can be rendered by iGibson 2.0 renderer. To maintain high physics realism, object models are annotated with size, mass, center of mass, moment of inertia, and also stable orientations. The collision mesh is a simplified version of the visual mesh, obtained with a convex decomposition using the VHACD algorithm. Object models with a shape close to a box are annotated with a primitive box collision mesh, much more efficient and robust for collision checking. For object categories that have the semantic property `openable` annotated, we make sure at least a subset of their object models have articulation, e.g. openable jars, backpacks, cars, etc. We either directly acquire them from the PartNet-Mobility Dataset [77] or acquire non-articulated models from TurboSquid, manually segment the models into parts, and then create the articulation in the URDF files. A subset of the object models are visualized in Fig. A.10.

We will publicly release the object dataset to be used for BEHAVIOR benchmarking. To preserve the rights of the model authors and the license agreement with TurboSquid, the 3D models are encrypted so that they can only be used within iGibson 2.0 and cannot be exported for other applications.

All models in the BEHAVIOR Dataset are organized following the WordNet [34], associating them to synsets. This structure allows us to define properties for all models of the same categories, but it also facilitates more general sampling of activity instances fulfilling initial conditions such as `onTop(fruit, table)` that can be achieved using any model within the branch `fruit` of WordNet. Fig. A.9 shows an example taxonomy of objects of the dataset organized in the WordNet taxonomy to perform a given household activity.

A.6 BEHAVIOR Dataset of Human Demonstrations in Virtual Reality

The main role of the VR demonstrations in BEHAVIOR is to provide a mechanism to normalize metrics, allowing to compare different embodied AI solutions between activity instances and scenes. However, we believe that the generated dataset of VR demos has the potential to be applied to other purposes, e.g., to generate AI solutions through imitation learning, or to study the mechanisms used

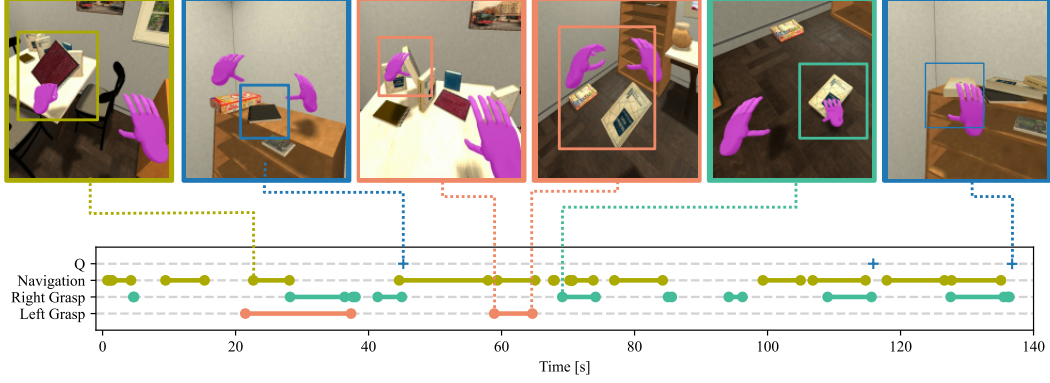


Figure A.11: **Sub-activity segmentation across activity execution for re-shelvingLibraryBooks:** We observe multiple cycles of long-range pick-and-place operations that eventually lead to activity success. In this figure, we show a sequence of snapshots of first-person view along with key frames (i.e. target objects placed on shelf, items dropped and picked up with alternating hands).

A.6.1 Collecting Human Demonstrations in Virtual Reality

To generate data, humans control a bimanual humanoid embodiment with a main body, two hands and a movable head based on stereo images displayed at 30 frames per second. The embodiment and the VR can be used with the most common VR hardware but for our dataset, we used a HTC Vive Pro Eye [78]. All recorded data can be deterministically replayed, achieving the same physical state transitions as reaction to the recorded physical interactions, which allows to generate any additional virtual sensor signal a posteriori. For more information about the VR interface, we provide the cross-submitted publication of iGibson 2.0 as part of the supplementary material.

We collect three different demonstrations of the same activity instance (same scene, same objects, same initialization) for each of the 100 activities in BEHAVIOR, 100 additional demonstrations, one for each activity for a different instance (different objects, different initialization) in the same scene, and 100 additional demonstrations, one for each activity in a different scene. This 500 demonstrations cover both the diversity in human execution, and the dimensions of variability in activity instances of BEHAVIOR. The data has been collected by voluntary participants and our own team.

A.6.2 Analysis and Statistics of Virtual Reality Demonstrations

The BEHAVIOR Dataset of Human Demonstrations in VR provides rich data of navigation, manipulation, and problem-solving from humans for long time-horizon and multi-step activities. Analyzing the statistical characteristics of the data (duration, hand use, room visitation, etc.) provides insights on how humans achieve their level of performance combining interaction and locomotion in the large BEHAVIOR scenes. Fig. A.11 depicts the segmentation of a VR demonstration into navigation and grasping phases while performing a pick-and-place rearrangement activity. This segmentation reveals multiple initial phases of navigation as the demonstrator observes the scene and locates activity relevant objects. For example, once the demonstrator reaches the table supporting the target objects (approx. at 18s), they pick up the target object with the non-dominant hand (approx. at 20s) and navigate to the goal location, before transferring the object to their dominant hand while positioning it (approx. at 30s). The demonstrator shows a preference for moving objects one-at-a-time, instead of stacking or carrying objects with each hand; this strategy will perform more poorly on the efficiency metrics T_{sim} , L_{body} , L_{right} , and L_{left} .

Fig. A.12 includes a) the duration of the VR demonstrations, b) the time spent in different room types, c) the hand used to interact and manipulate, and d) the complexity of the activities in logical representation vs. time. We observe that BEHAVIOR activities cover a wide range of time-horizons, from less than 2 minutes to more than 11 minutes. The activities show a bias towards living spaces (kitchen, living-room, bedroom), with the most prevalent room being the kitchen. A large portion of BEHAVIOR activities involve preparing food or cleaning appliances that are only supported in kitchens. Furthermore, as expected, the data reflect a bias towards dominant hand manipulation, followed by bimanual grasping, which is required for lifting and manipulating large objects. The

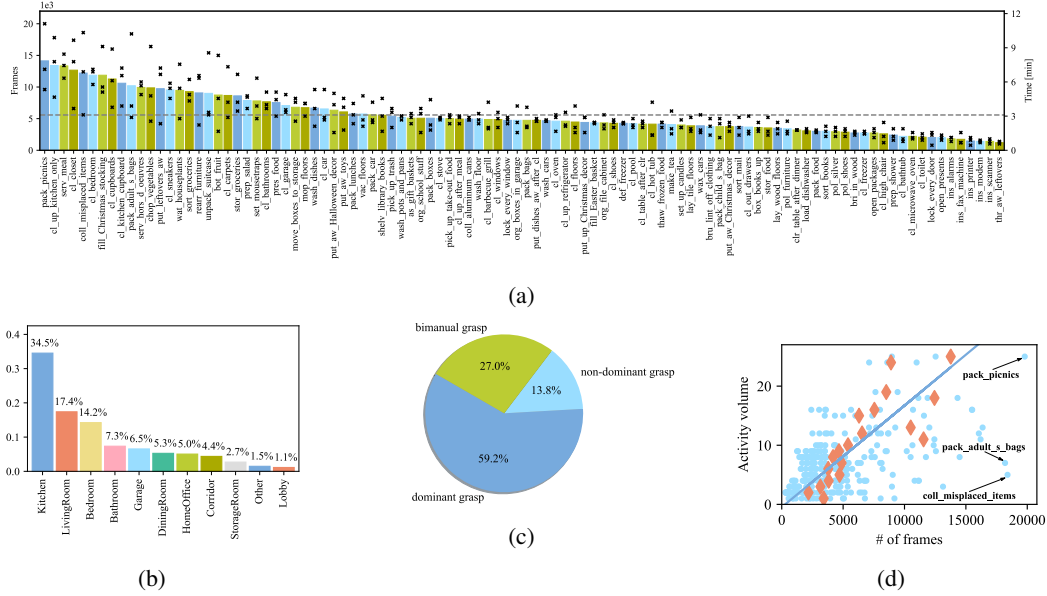


Figure A.12: **Analysis of human demonstrations of BEHAVIOR activities in virtual reality:** a) Duration of each successful demonstration (mean and individual trials, decreasing order); b) Fraction of total VR time spent in each type of room; c) Fraction of total VR time spent manipulating with the dominant, non-dominant, or both hands; d) Duration of each VR demonstration wrt. activity volume; blue dots denote individual demos and red diamonds denote the mean time for each number of ground literals (activity volume). Larger volume correlates with larger duration ($R^2 = 0.826$).

high use of two hands to manipulate correlate to the use in real-world; we hope that our dataset helps exploring this type of interaction that has been traditionally less studied in embodied AI. The total number of ground predicates (activity volume) is strongly correlated with the total activity time indicating that the volume is a good measure of the complexity of an activity. Outliers include activities with a high ratio of time to goal condition such as the ones that require cleaning a large area (cleaningCarpets, vacuumingFloors) or searching (collectMisplacedItems).

Analyzed individually, Fig. A.13 shows that room occupancy depends heavily on the type of activity. Room occupancy reflects common intuition about household activities; the ones associated with living-space decorations (puttingAwayChristmasDecorations, puttingAwayHalloweenDecorations) take place primarily in the living room, whereas cooking activities (preparingSalad, preservingFood) occur primarily in the kitchen. Similar activity preferences are observed in the grasping data; activities requiring installing unwieldy objects (layingWoodFloors, layingTileFloors) require the use of both hands, whereas simple cleaning activities (cleaningThePool) that require using a cleaning tool are performed with the dominant hand.

A.6.3 Gaze Tracking in Virtual Reality

Our preference on HTC Vive Pro to collect the BEHAVIOR Dataset of Human Demonstrations is motivated by its ability to track the gaze (pupil movement) of the demonstrator. We consider gaze information to be a valuable source to understand human performance in the activities. While other datasets of gaze are available [79], this is the largest dataset of active gaze attention during manipulation in simulation, providing synchronized ground-truth information of the object being observed, its state and full shape. Fig. A.15 depict examples of the tracked human gaze during activity execution, with the object attracting the gaze indicated in magenta. Fig. A.16 includes several statistics of the gaze attention over object categories in the entire dataset and for some example activities. Both figures indicate a clear correlation between the gaze data and the goal of the activities: we expect the dataset to be useful to study and predict human gaze attention, and to develop new embodied AI algorithms for active [80, 81] and interactive perception [82].



Figure A.15: Human gaze during activity execution; Red dot: human gaze point, Magenta: object gazed; The BEHAVIOR Dataset of Human Demonstrations in Virtual Reality includes 500 demonstrations (758.5 min) with gaze information while humans navigate and interact (accuracy: $\pm 4^\circ$ [83]); The gaze information correlates strongly with activity; We hope that this data can support new research in visual attention and active vision to control agent’s camera

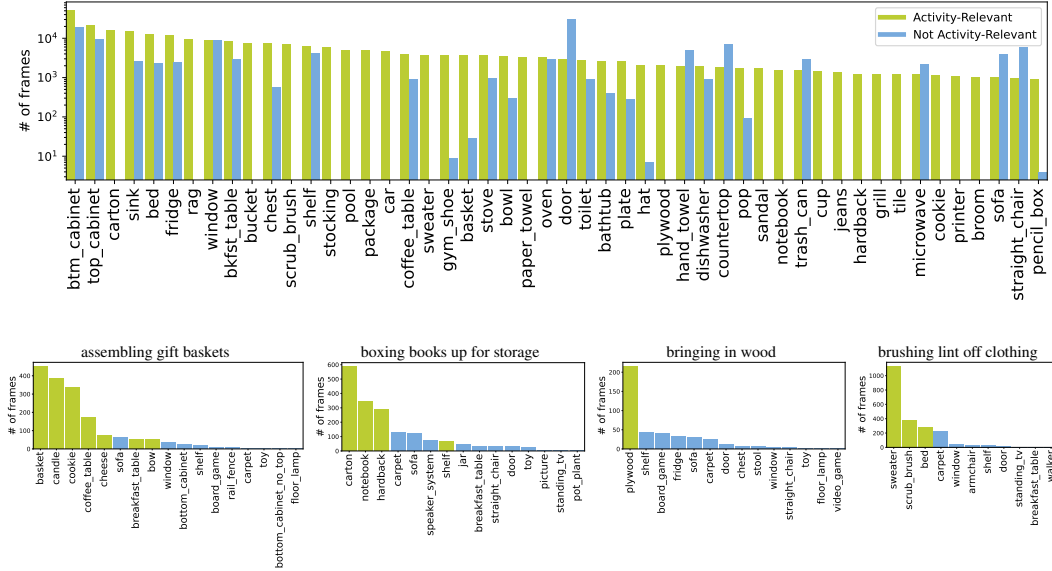


Figure A.16: Statistics of the attention over object instances of WordNet categories in the BEHAVIOR Dataset of Human Demonstrations in virtual reality aggregated for all demonstrations (top, logarithmic scale), and segregated for four activities (bottom row, linear scale), for activity-relevant (green) and not activity-relevant (blue) objects. In households, the aggregated of the visual attention goes to containers of objects (cabinets) and doors separating rooms; For individual activities, visual attention concentrate on specific activity-relevant objects

and RL with motion primitives. We will first elaborate the shared setup between the two, then go into their differences.

Shared Setup: In the normal “partial observability” setup, the observations include 128×128 RGB-D images from the onboard sensor on the agent’s head and proprioceptive information (Head pose in local frame, hand poses in local frame, and a fraction indicating how much each hand is closed). The proprioceptive information is 20 dimensional. For the experiments with “full observability”, the observations include the ground truth object poses for all the activity-relevant objects, the agent’s pose, and the proprioceptive information.

The agent receives a reward of 1 for every ground goal condition (literal) that it satisfies during the episode. The episode terminates if the agent achieves a success score Q of 1 (achieved all literals in the goal condition) or it times out.

The policy network architecture is largely shared in the following setups. With RGB-D images as input, we use a 3-layer convolutional neural network to encode the image into a 256 dimensional vector. Proprioceptive information and/or poses for all activity-relevant objects are also encoded into a 256 dimensional vector with an MLP, respectively. The features are concatenated and pass through another MLP to generate action representation, which could be a box action space or discrete action space depending on the setup (continuous actions or action primitives).

RL with Continuous Action Space: For this agent variant, we use Soft Actor-Critic (SAC) [16] implemented by TF-Agents [84]. The action space is continuous and has a dimensionality of 18. The first three dimensions represent the locomotion actions: desired x-y translation of the robot body and the desired rotation around the vertical axis. The next seven dimensions represent the linear and angular velocities of the left hand (in Cartesian space, 6 DoF) and 1 DoF closing/opening of the hand. The last seven dimensions is the same action but for the right hand. The maximum episode length depends on the experimental setup. For instance, if the initial state corresponds to 1 s away from a goal state, we will give the agent three times the amount of time (i.e. 3 s) to accomplish the activity. We train for 20K episodes, evaluate the final policy checkpoint and report the results in Table 2.

RL with Motion Primitives: For this agent variant, we use Proximal Policy Optimization (PPO) [17] implemented with TF-Agents [84]. The action space is discrete, with $n_r \times m$ choices, where n_r is the number of activity-relevant objects and m is the number of action primitives. Here we didn’t allow the agent to operate on all objects in the scene, but focus on activity-relevant objects to facilitate learning. Following our implementation of motion primitives, $m = 6$. Laying out the choices on a $n_r \times m$ grid, and i -th column j -th row means to apply j -th action primitive on i -th activity-relevant object. Not all combinations of action primitive and object are compatible and action that is not feasible is converted into no-ops. The maximum episode length is set to 100 for all activities. We experiment with partially simulated motion primitives and fully simulated motion primitives, as described in Sec. A.4). We train with partially simulated motion primitives until convergence, and evaluate and report the results on partially simulated motion primitives and fully simulated motion primitives, since training with motion planning in a complex scene is very time-consuming. In the experimental results shown in Table 2, generally fully simulated motion primitives results are much worse than partially simulated motion primitives, this is intuitive because motion planning performs more rigorous checks and complies with the physical model, highlighting the complexity of BEHAVIOR.

Experimental Setup for the Effect of Diversity: To evaluate diversity, we train for individual skills instead of full BEHAVIOR activities. Here, we adopt an easier experimental setup that allows us to study the effect of diversity; the results are reported in Table 3. Specifically, we use RL with continuous action space but with a more constrained action space: 6-dimensional representing the desired linear and angular velocities of the right hand (assuming the rest of the agent is stationary). For grasping, we adopt the “sticky mitten” simplification from other works [23]: we create a fixed constraint between the hand and the object as soon as they get in contact. We also use distance-based reward shaping to encourage the hand to approach activity-relevant objects. To evaluate the effect of diversity in object poses, we use the same object models and randomize their initial poses during training. To evaluate the effect of diversity in object instances, we randomize the object models during training. For example, for the `sliced` single-predicate activity, the agent will encounter different types of fruit (e.g. peach, strawberry, pineapple, etc) during training. We train for 10K episodes, evaluate the final policy checkpoint and report the results in Table 3.

A.8 Potential to Transfer to Real-World

BEHAVIOR is a benchmark in simulation. This facilitates a continuous evaluation of solutions, fair and equal conditions, and increased accessibility without expensive robot hardware. It is also instrumental for modern robot learning procedures that require generating large amount of experiences. However, the use of simulation introduces a gap between the activities in our benchmark and the equivalent activities in real world. We argue that, while not negligible, we have taken measures

to close this gap with the goal of providing a benchmark where the performance of embodied AI solutions is close to the performance they would have in a real world system.

Our instantiation of BEHAVIOR includes realistic scenes and object models, with high-quality visuals and close-to-real physical properties (mass, center of mass, friction) annotated in manual process assisted by the information obtained in the Internet. The underlying physics engine, pyBullet [58], is acknowledged as one of the standards in robotics and a high quality approximation of the underlying mechanical processes. The physics-based rendering from iGibson 2.0 generates high quality images to use as input in our evaluation. While our bimanual agent is not realistic, our second provided embodiment is a realistic robot model, a Fetch, with similar kinematics, actuation and sensing, facilitating the evaluation of solutions in BEHAVIOR that could act similarly on a real robot. Previous works have demonstrated good results developing solutions in iGibson 2.0 that could transfer to real world [2, 85], and evaluated the similarities between simulated and real-world sensor signals [22]. This indicates a high potential for the solutions evaluated in simulation in BEHAVIOR to perform similarly in the real world, a claim that we plan to evaluate experimentally after the pandemic.

A.9 Ethical Considerations

BEHAVIOR includes data (activity demonstrations) generated by humans. After evaluation, our institution’s Institutional Review Board (IRB) considered this project exempt from review: the data does not reveal any private information about the participants. Human demonstrations are collected in line with standard ethics practices, among lab members and volunteers. In terms of broader societal impacts, BEHAVIOR is aimed to facilitate research of autonomous robots performing activities of daily living. The potential impacts of this line of work, particularly on labor impacts of automation and the physical safety of humans interacting with autonomous robots, are far-reaching.