

AllenAct: A Framework for Embodied AI Research

Luca Weihs^{*1}, Jordi Salvador^{*1}, Klemen Kotar^{*1}, Unnat Jain², Kuo-Hao Zeng³
 Roozbeh Mottaghi^{1,3}, Aniruddha Kembhavi^{1,3}

¹ PRIOR @ Allen Institute for AI ² UIUC ³ University of Washington

<https://allenact.org>

Abstract: The domain of Embodied AI, in which agents learn to complete tasks through interaction with their environment from egocentric observations, has experienced substantial growth with the advent of deep reinforcement learning and increased interest from the computer vision, NLP, and robotics communities. This growth has been facilitated by the creation of a large number of simulated environments (such as AI2-THOR, Habitat and CARLA), tasks (like point navigation, instruction following, and embodied question answering), and associated leader-boards. While this diversity has been beneficial and organic, it has also fragmented the community: a huge amount of effort is required to do something as simple as taking a model trained in one environment and testing it in another. This discourages good science. We introduce AllenAct, a modular and flexible learning framework designed with a focus on the unique requirements of Embodied AI research. AllenAct provides first-class support for a growing collection of embodied environments, tasks and algorithms, provides reproductions of state-of-the-art models and includes extensive documentation, tutorials, start-up code, and pre-trained models. We hope that our framework makes Embodied AI more accessible and encourages new researchers to join this exciting area.

1 Introduction

In recent years we have witnessed a surge of interest within the computer vision, natural language, and robotics communities towards the domain of *Embodied AI* (E-AI) - learning, while situated within some animate body (*e.g.* a robot), to perform tasks in environments through interaction. This has led to the development of a multitude of simulated environments employing photorealistic images (such as Gibson [1] and AI Habitat [2]), involving robot-object interaction (such as AI2-THOR [3] and Virtual Home [4]), focused on manipulation (such as RL-Bench [5], Sapien [6], and Meta-world [7]), using advanced physics simulations (such as MuJoCo [8] and ThreeDWorld [9]), and also physical counterparts to simulation environments (RoboTHOR [10] and iGibson [11]) to enable research in simulation-to-real transfer. Within these environments, research has progressed towards learning to interact: including visual navigation [12, 13, 14], question answering [15, 16], task completion [17], instruction following [18, 19], language grounding [20, 21], grasping [22, 23], object manipulation [24, 25], future prediction [26], and multi-agent collaboration [27, 28]; as well as using interaction as a tool to learn: environment representations [29], intuitive physics [30], and objects and attributes [31]. The rapidly growing list of publications in E-AI, see Fig. 1, as well as popularity of E-AI workshops and challenges in top computer vision and machine learning conferences over the past couple of years exemplify the growing interest in this domain.

As the domain of E-AI continues to grow, it faces several challenges: (a) *Replication across tasks and datasets* - While our community proposes a host of novel methods each publication cycle, these techniques are frequently evaluated on a single task and within a single simulation environment (see Fig. 1). Just as we now expect neural architectures to be evaluated across multiple tasks (*e.g.* computer vision tasks include classification, detection, and segmentation) and multiple datasets (such as ImageNet [33] and Places [34]), we must also start evaluating E-AI methods across tasks and datasets. Unfortunately, this currently requires large-scale changes to a code-base and thereby discourages comprehensive evaluation. (b) *Unravelling what matters* - As the field progresses via

* indicates equal contribution.

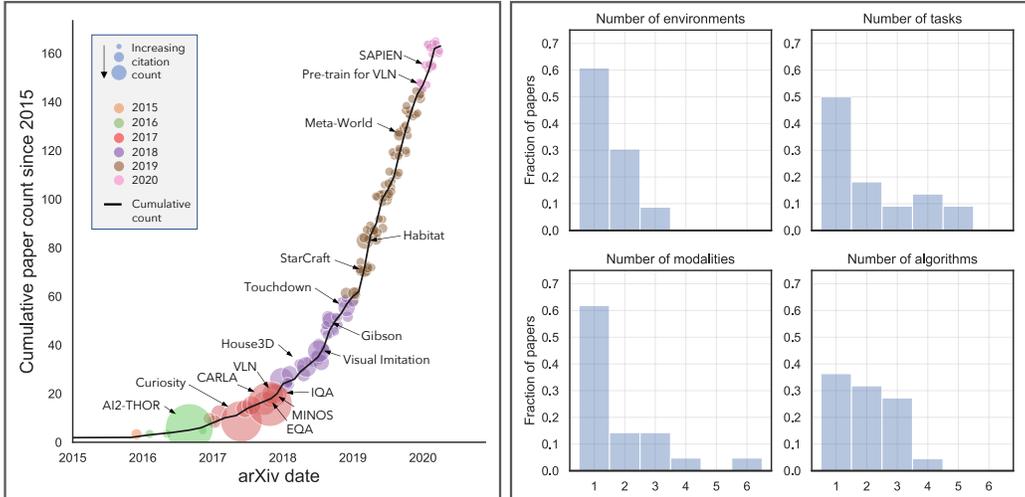


Figure 1: **Growth and fragmentation of E-AI.** *Left* - the cumulative number of papers published on arXiv since 2015 which were identified as being in the E-AI domain (see Appendix A for details). The number of publications has dramatically increased in recent years. *Right* - after manually annotating the 20+ most-cited E-AI papers, we plot histograms of the frequencies with which these papers ran experiments on multiple environments, with multiple tasks, etc. The large frequency with which only a single task, environment, and modality is evaluated suggest large barriers to comprehensive evaluation. While several papers evaluate multiple algorithms, we noticed little standardization - some compare imitation with reinforcement learning, some compare A3C and PPO, others try Q-learning. Moreover, these represent the most cited papers of the past years, likely making this analysis not representative of a randomly selected paper. This analysis used the S2ORC [32].

improvements on standard tasks and benchmarks, it is crucial to understand what components of systems matter most and which do not matter at all. This unravelling requires careful ablation studies and analyses. In addition to evaluations across tasks and datasets, this involves swapping out learning algorithms (e.g., on-policy and off-policy), losses (e.g., primary and auxiliary), model components (representation stacks, maps, etc.), and hyperparameters. These analyses are often critical for good science and also fast progress. Today’s frameworks and libraries can certainly be improved in this regard. Why should swapping a learning algorithm (e.g., PPO with A2C) be any more tedious than changing, e.g., the learning rate? (c) *Ramp up time* - Getting up to speed with E-AI algorithms takes significantly longer than ramping up to classical tasks in vision and NLP like image classification or sentiment analysis. Just as the early deep learning libraries like Caffe and Theano, and numerous online tutorials, lowered entry barriers and ushered in a new wave of researchers towards deep learning, E-AI can benefit from modularized coding frameworks, comprehensive tutorials, and ample startup code. (d) *Large training cost* - E-AI is expensive. Today’s state of the art reinforcement learning (RL) methods are sample inefficient and training competitive models for embodied tasks can cost tens of thousands of dollars - within the reach of industrial AI labs, but unaffordable for independent researchers and smaller organizations. The availability of large networks pre-trained on ImageNet (with accompanying code and models on standard libraries like PyTorch) significantly reduced the cost of training on downstream tasks. A similar centralized repository with a diverse set of E-AI code and models can greatly benefit our community.

As detailed in Section 2, there is no shortage of open-source reinforcement learning libraries and frameworks available today. While these frameworks excel in their particular domains, for research in E-AI we found that each individually lacked features we consider critical. In particular, no single framework simultaneously provides: support for a large number of E-AI environments and tasks, a variety of training algorithms, a capacity to construct custom training pipelines, the right balance between adding new and exploiting existing functionality, and a high likelihood of continued support and development. For this reason, we set out to develop a new framework focused on E-AI research.

We present the AllenAct framework, written in Python and using PyTorch [35], designed for research in E-AI with a focus on modularity, flexibility, and well encapsulated abstractions. It inherits

the best design principles and builds upon other AI libraries including `pytorch-a2c-ppo-acktr`¹ and Habitat-API [2]. While AllenAct will continue to improve, we highlight the following existing features: (1) *Environments* - we provide first-class support for the iTHOR [3], RoboTHOR [10], and Habitat [2] embodied environments and numerous tasks within, as well as for grid-worlds including MiniGrid [36]. Grid-worlds serve as excellent sand-boxes to evaluate new algorithms owing to their rendering speed and variable complexity. Swapping out environments, as well as adding new ones, is made simple. (2) *Task Abstraction* - tasks and environments are decoupled in AllenAct. This allows researchers to easily implement a large variety of tasks in the same environment. (3) *Algorithms* - we provide support for a variety of on-policy algorithms including PPO [37], DD-PPO [13], A2C [38], Imitation Learning (IL), and DAgger [39] as well as offline training such as offline IL. (4) *Sequential Algorithms* - AllenAct makes it trivial to experiment with different sequences of training routines, which are often the key to successful policies (example: IL followed by PPO). (5) *Simultaneous Losses* - AllenAct allows researchers to easily combine various losses while training models (for instance, use an external self-supervised loss while optimizing a PPO loss). While seemingly trivial, we found that present day RL libraries make this unnecessarily harder than it need be. (6) *Multi-agent support* - AllenAct provides support for multi-agent algorithms and tasks. (7) *Visualizations* - effective visualizations of embodied environments are critical for debugging and ideation. AllenAct provides out-of-the-box support to easily visualize first person and third person cameras for agents as well as intermediate model tensors and integrates these into Tensorboard. (8) *Pre-trained models* - AllenAct provides a number of models and accompanying code to train these models for standard E-AI tasks. (9) *Tutorials* - we provide start-up code to help ramp up new researchers to the field of embodied-AI as well as tutorials for performing common actions like adding new environments, tasks, and models.

The AllenAct framework will be made open source and freely available under the MIT License. We welcome and encourage contributions to AllenAct’s core functionalities as well as the addition of new environments, tasks, models, and pre-trained model weights. Our goal in releasing AllenAct is to make E-AI more accessible and encourage thorough, reproducible, research.

2 Related Work

Embodied AI platforms. AI research has benefited from platforms that enable agents to interact with, and obtain observations from, an environment. These platforms have been used as benchmarks to evaluate AI models on different types of tasks ranging from games [40] to performing tasks in indoor environments [3] to autonomous driving [41]. ALE [40], ViZDoom [42], and Malmo [43] are example game environments. Arena [44] provides a multi-agent platform for games. Several efforts have produced environments for navigation with virtual robotic agents [10, 2, 11]. AI2-THOR [3], CHAI [45], and Virtual Home [4] are examples of platforms that go beyond navigation and enable evaluation of agents on tasks that require interaction such as applying forces and/or changing object states. Platforms such as RL Bench [5], Sapien [6], and Meta-World [7] focus on manipulation tasks, while [46, 47] enable studying the task of grasping. The DeepMind Control Suite [48] provides a platform for continuous control tasks. CARLA [41] is designed to evaluate autonomous driving capabilities. Our goal is to provide a framework with general abstractions so researchers can easily plug-in their environment of interest and begin experimentation. We provide code for integrating multiple environments and tasks (see Sec. 3). We will continue to add more ourselves and encourage other researchers to do the same. OpenAI Gym [49] also provides a standard wrapper for a set of environments including the Atari games and MuJoCo control tasks. AllenAct differs from Gym in its abstractions, capabilities, and E-AI being its primary focus.

Embodied AI and reinforcement learning libraries. There have been several libraries developed over the years for E-AI, a few recent libraries that are most relevant to ours are discussed here. ML-Agents [50] enables defining new environments in the Unity game engine and provides a collection of example environments and RL algorithms. PyRoboLearn [51] provides a robot learning framework, where the idea is to disentangle the learning algorithms, models, robots, and their interface and to provide an abstraction for each of these components. Habitat-API [2] is a modular framework for defining embodied tasks and agent configurations and training and evaluating these agents. There are also RL frameworks without an embodied focus, for example, Garage² (also

¹<https://github.com/ikostrikov/pytorch-a2c-ppo-acktr-gail>

²Garage www.github.com/rlworkgroup/garage. Keras-RL www.github.com/keras-rl/keras-rl

known as `rllib`), OpenAI Gym [49], Dopamine [52], and Keras-RL². Each of these libraries offers a unique feature set well-suited to a particular research, or production, workflow. In contrast to these libraries, AllenAct is designed to provide first-class support (*e.g.*, including tutorials, starter-code, visualization, and pretrained models) for a wide range of E-AI tasks while also allowing for substantial flexibility in defining new training pipelines and integrating new environments and tasks.

3 The AllenAct framework

Designing software for AI tasks requires a delicate balance between the ease with which (a) new functionality can be added and (b) the existing functionality can be exploited. For instance, a framework designed only to train GRU based agents with the PPO algorithm to complete a navigation task within the AI2-THOR environment can narrow its API so that a user needs only to specify a small set of relevant hyperparameters before running a new experiment. This makes research within this domain extremely streamlined at the expense of flexibility: if a user now wants to try something beyond the scope of the design (*e.g.* train with the A2C loss) they will need to dive into the internals of the framework to understand what, often substantial, changes must be made. In our experience with E-AI research, the frequency with which we have had to modify our software to adapt to new experimental requirements has followed the following approximate pattern:

Daily-Weekly: Modify hyperparameters associated with the training loss (*e.g.* reward discount factor γ), model (*e.g.* RNN hidden state size), optimization (*e.g.* learning rate, batch size), and hardware (*e.g.* number of GPUs and training processes).

Weekly-Monthly: Modify model architectures, training strategies (*e.g.* warm-start a model with IL before training with PPO), sensor modalities (*e.g.* adding depth maps as input to the model).

Quarterly-Yearly: Adding new environments (*e.g.* SAPIEN), new tasks (*e.g.* a language and vision task such as ALFRED [19]), changes to the definition of an existing task (*e.g.* success in object-navigation requires an explicit stop signal in addition to proximity to the object), new losses to be used during training (*e.g.* auxiliary self-supervision), and incorporating new training paradigms (*e.g.* moving from asynchronous methods, *e.g.* A3C [38], to synchronous methods, *e.g.* PPO [37]).

In designing AllenAct, we have stressed modularity and flexibility while keeping the above in mind. Thus changing hyperparameters or model architectures is trivial and making more substantial changes, such as adding a new training paradigm (*e.g.* deep-Q learning), requires more knowledge of the framework’s internals but is still relatively straightforward. Following community standards, AllenAct is written in the Python programming language and heavily leverages the PyTorch library for designing deep-neural models and enabling their optimization. Next, we describe AllenAct’s API, features, documentation, and associated pre-trained E-AI models.

3.1 Abstractions and API

The API of AllenAct is defined by a collection of abstractions (each corresponding to a Python class) which, themselves, are best understood in context of their relationships. At a high level, an agent is defined by an `ACTORCRITICMODEL`, observes the world using `SENSORS`, and interacts with its `ENVIRONMENT` to complete a `TASK` which defines rewards and success criteria. New instances of a `TASK` (*e.g.* navigation starting from a different point) are created, sequentially, for the agent by a `TASKSAMPLER` and, during training, the agent’s parameters are updated to minimize some collection of `LOSSES`. Which `LOSSES` are used at a particular point in training is determined by the `TRAININGPIPELINE`. Rather than describing all of these abstractions in detail³ we instead highlight how these abstractions differ from those used in most RL libraries. Our code adapts and generalizes several abstractions from Habitat-API. For instance, their `DATASET` is generalized into our `TASKSAMPLER`, and while we both share a `TASK` abstraction, theirs is used within an Open AI gym `ENV` class while ours acts as an intermediary between the agent and the environment.

Experiments defined in code. In AllenAct, experiments are defined by creating an implementation of the abstract `EXPERIMENTCONFIG` class. Changing hyperparameters in such files is just as simple as doing so within text-based configuration files (a necessity, as noted above, as these types of changes occur daily to weekly) but with the added benefit that, at the cost of some additional boilerplate, it is trivial to add new hyperparameters, update model architectures, etc. Moreover, writing

³See AllenAct’s documentation for these comprehensive details.

configuration in code allows easy access to a wide range of productivity features provided by modern integrated development environments such as auto-completion and type hints. This hugely simplifies daily-weekly modifications and enables researchers to easily run several experiments. An example of how one might create an EXPERIMENTCONFIG implementation to train a navigation model in AI2-THOR can be seen in the documentation.

Flexible training pipelines. Training high-quality agents often requires a pipelined approach where, for example, an agent’s policy is given a warm-start by first training with IL after which reinforcement learning is used to further improve performance and generalization. While such training pipelines can be accomplished manually, AllenAct introduces a TRAININGPIPELINE class which makes the concept of a training pipeline a core concept within the framework. A TRAININGPIPELINE is defined by a collection of sequential PIPELINESTAGES which define: (a) the losses to be used, (b) the length of training and any early stopping criteria, and (c) whether or not to apply teacher forcing (see Sec. 3.2). During training, AllenAct moves through these stages and updates the agent accordingly. With this design, adding an IL warm-start to an experiment requires adding a single additional line of code. Thus the weekly-monthly change in training pipeline takes, at most, a few minutes and requires little additional bookkeeping.

Decoupling the environment from the task. A standard abstraction used within multiple RL frameworks is OpenAI Gym’s ENV. This ENV class defines (i) how an agent interacts with the environment, (ii) whether or not success criteria are met, (iii) the rewards returned after every action, (iv) observations available to the agent, and (v) how to reset itself. This abstraction is an excellent fit for many settings, especially those in which the environment (*e.g.* an Atari game) is intimately tied to the agent’s intended goal (*e.g.* beating the game). This abstraction is less natural in the setting of E-AI where the environment (*e.g.* AI2-THOR, Habitat, ThreeDWorld, etc.) has no innate goal and, in fact, a huge variety of distinct goals can be defined. Within the AI2-THOR environment alone, we are aware of nine unique tasks defined by various authors ranging from navigation [10] to multi-agent furniture moving [28]. Instead, in AllenAct we disentangle the TASK from the ENVIRONMENT. The ENVIRONMENT provides a means by which to modify environment state while the TASK encapsulates the agent’s goal and acts as an intermediary between the agent and the environment. The TASK defines the actions available to the agent, any success criteria, and the rewards returned to the agent. Once a TASK has been completed by the agent it is simply thrown away and, as described below, a new task is generated by the TASKSAMPLER. Beyond being conceptually appealing, this decoupling can make the quarterly-yearly updates to (and additions of) tasks far easier as, generally, changes are confined to the TASK class and large portions of code require no changes. This decoupling also simplifies the process of introducing new environments within the AllenAct framework.

Flexible task initialization. As previously noted, an OpenAI Gym’s ENV instance must be able to reset itself and so any implementation of ENV implicitly defines the stream of goals that the agent observes during training. This is well-suited for most RL research but is not a good fit for E-AI where one often needs more control over which goals are presented to the agent and their order. We instead use a TASKSAMPLER abstraction to allow complete control of how new instances of a task are, sequentially, generated for the agent. With this abstraction, enabling curriculum learning is as simple as defining a TASKSAMPLER that progressively samples more difficult instances of a task. TASKSAMPLERS enable quick experimentation with new training strategies which we require at the weekly-monthly frequency.

Together these changes allow for considerable flexibility and provide a useful mindset by which to approach E-AI problems. These abstractions are sufficiently general to be of use even beyond E-AI research, indeed AllenAct has been used in a grid-world-based study of reinforcement learning methodology [53].

3.2 Features

Environments and Tasks. A key goal of AllenAct is to provide first-class support for a diverse range of embodied environments and tasks. In this early release, we provide support for Habitat, iTHOR, and RoboTHOR and tasks within them (See Table 1). We also provide support for MiniGrid that serves as a fast sand-box for algorithm development. In future releases, we will extend support to the recently released SAPIEN [6] and ThreeDWorld [9] environments and associated tasks (*e.g.* robotic manipulation). A crucial advantage of AllenAct is the ease at which one may test the same model (or training pipeline, loss, etc.) across multiple environments and tasks. The

AllenAct documentation shows an example of the few changes required to an EXPERIMENTCONFIG to switch from one task to another in iTHOR and then move to the Habitat environment.

Algorithms. Our framework currently supports decentralized, distributed, synchronous, on-policy training with A2C, PPO, IL, or any user-defined loss. This builds upon the decentralized distributed proximal policy optimization (DD-PPO) [13] algorithm available in the Habitat-API with the additional flexibility that training can be done with arbitrary loss functions so that we can, for example, run DD-A2C or DD-curiosity-driven-exploration [56]. While

on-policy and synchronous RL algorithms have become very popular, with PPO the de-facto standard, and have been used with great success, on-policy methods are notoriously sample-inefficient and synchronous methods impose run-time limitations that are not appropriate for every problem. To this end, AllenAct currently supports two means by which to relax the on-policy assumption. *Teacher forcing* - in order to implement algorithms such as DAgger [39], an agent must be able to replace its action with the action of an expert with some probability (potentially decaying over training). In AllenAct, implementing teacher forcing (and thus DAgger) is as simple as defining a linear-decay function. *Training with a fixed, external, dataset* - it is frequently beneficial to be able to supervise an agent using a fixed dataset (e.g. IL from a dataset of human examples). AllenAct enables this type of supervision and also interleaving off-policy updates with on-policy ones. While we have found the above relaxations of the synchronous and on-policy assumptions to be sufficient for most prior work, we recognize that this will not be the case for all users of AllenAct. In light of this, our future roadmap includes incorporating deep Q-learning methods as well as capabilities for asynchronous execution and training.

Multiple Agents. A key facet of E-AI having received relatively little attention is collaboration and communication among multiple agents. To support research in this direction we have natively enabled training multi-agent systems in AllenAct. As seen in Table 1, we will soon provide high-quality support for the multi-agent tasks recently developed for AI2-THOR [27, 28].

Visualization. In our experience, visualizations and qualitative evaluations of the policies learned by E-AI agents are critical to debugging and understanding the limitations of current systems. Unfortunately, producing such visualizations can be time consuming, especially so if these visualizations are meant to be of sufficient quality to be used in presentations or publications. To lower the burden of visualization, the AllenAct framework contains a number of utilities (including ego-centric views, top-down views, third party views and tensor visualizations) for environments with first-class support (recall Table 1). Some of these visualizations, which can be automatically logged during inference, are presented in Figure 2. The range and scope of these visualization utilities will grow as further embodied environments and tasks are incorporated into our framework.

Tutorials, Documentation, and Typing. Beginning to work with a new framework can be a daunting and frustrating experience as one must internalize a large number of new abstractions, frequently with little documentation, written in an unfamiliar style. For this reason, we have made tutorials and documentation a high priority within AllenAct and, in our first release, we have several tutorials such as training a PointNav model in RoboTHOR or how to switch environments for a particular task. Moreover, we have added type hints throughout our codebase to enable IDE auto-completion and warnings.

Pre-trained models. To encourage reproducibility, we include several pre-trained model checkpoints (reproducing, within error, published results) for tasks with first-class support. This includes

Environment	Tasks
iTHOR [3]	PointNav [54], ObjectNav [12], ALFRED [19] FurnLift [27], FurnMove [28]
RoboTHOR [10]	PointNav [54], ObjectNav [12]
Habitat [2]	PointNav [54], ObjectNav [12], VLN [18]
MiniGrid [36, 55]	All ⁴

Table 1: Environments and tasks supported in AllenAct. This support includes starter-code, visualization tools, and pre-trained models. Tasks and environments in purple have support planned for upcoming releases.

⁴As an easily composable grid-world, MiniGrid can be used to quickly generate a huge number of different tasks. Taken together, MiniGrid and BabyAI have ≥ 30 unique tasks by default.

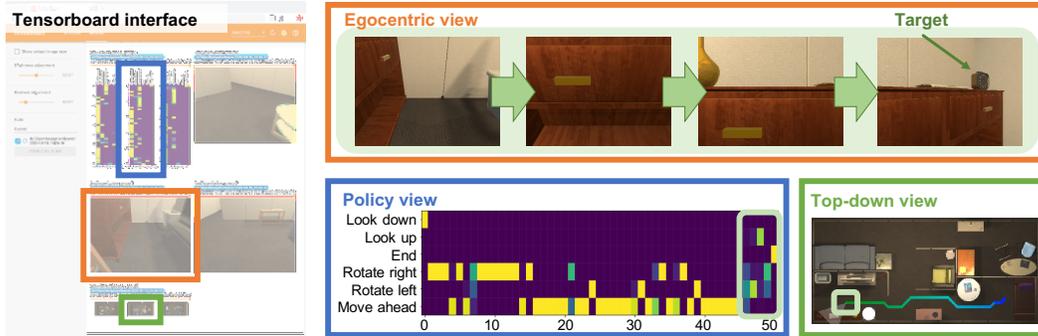


Figure 2: **Visualizations.** A simple plug-in-based interface allows generating different types of visualizations in Tensorboard for a particular task as shown in the screen capture on the left. A navigation task is shown in this example. We display ego-centric and top-down views as well as the policy over time. The top-down view enables observing the overall behavior of the agent. We highlight one segment of the trajectory with a green overlay. The ego-centric view enables interpreting the policy and visually assessing task success. The policy enables analyzing and debugging the probability of each action at each point in time.

all models trained towards the experiments in Sec. 4. As these models were trained within AllenAct, we also provide training and inference code for these models.

Future development and support. An important consideration when deciding whether or not to adopt a new framework is that of future support. Our team is committed to research in the domain of E-AI and expect to continue AllenAct’s development for, at least, several more years. Indeed we currently have a number of ongoing projects using AllenAct, in executing these projects we expect to obtain robust feedback that will be used to improve AllenAct. We will also encourage and make it easy for E-AI researchers to contribute code and models to AllenAct.

4 Experiments

We now highlight the capabilities of AllenAct by reproducing results from the (embodied) RL literature along with evocative ablations. Code and model checkpoints for all experiments can be easily accessed within AllenAct and serve as strong comparative baselines for future work.

Support for embodied environments and tasks. Using AllenAct, we have reproduced a number of results for navigation in the iTHOR, RoboTHOR, and Habitat environments, see Fig. 3a. Two variants of navigation are commonly considered in the literature: PointNav (navigating with a displacement vector to the goal) and ObjectNav (navigating to goal specified by a category label). We train DD-PPO [13] for 75 million frames and obtain a validation accuracy of 92.5%. This accuracy is within error of, and indeed slightly outperforming, the model in [2]. [13] demonstrated that if trained for 2.5 billion frames using ≈ 4608 GPU hours DD-PPO can reach 99.9% validation set accuracy. For our aim of demonstrating reproducibility and functionality, we restrict ourselves to 75 million frames, the same as [2] with whom we compare. When training to complete PointNav in iTHOR and RoboTHOR we obtain similarly high performance.

We demonstrate the ObjectNav task on iTHOR and RoboTHOR. In these experiments, we use a ResNet and LSTM based architecture and train our models for 200 million steps. Within RoboTHOR, our model outperforms the best model submitted to the recent CVPR’20 RoboTHOR challenge⁵. While no such challenge was undertaken for iTHOR, our similarly high metrics suggest we have obtained a well-trained model that will serve as a strong baseline for future work. Implementation details are in Appendix B.1.

Support for online and offline algorithms. AllenAct supports a range of different built-in training algorithms beyond PPO including A2C and several varieties of IL (*e.g.* on-policy behavior cloning, DAgger, and purely offline training from a fixed dataset of demonstrations). Precise details of these IL-based methods are given in Appendix B.2. In Fig. 3b, we highlight the results of using these algorithms to train an agent to complete the GoToLocal task in the BabyAI [55] grid-world environ-

⁵<https://ai2thor.allenai.org/robothor/challenge/>

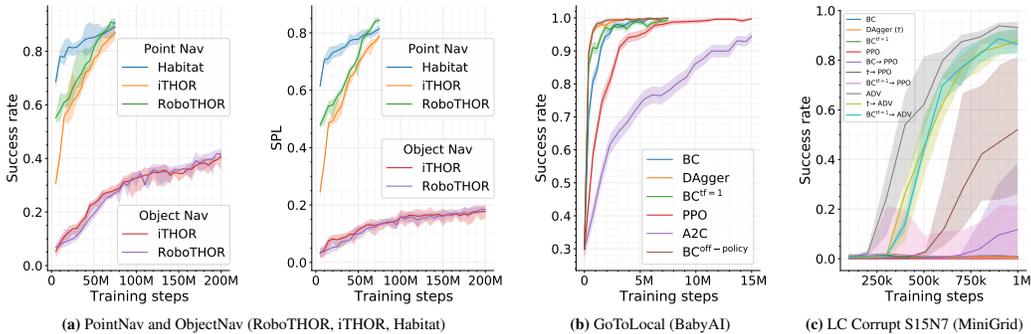


Figure 3: **Experimental results.** (a) Validation-set success rate and Success weighted by Path Length (SPL) [54] of models trained with DD-PPO to complete the PointNav & ObjectNav tasks in iTHOR, RoboTHOR & Habitat. Shaded regions indicate the range of metrics in five runs for PointNav in Habitat, three runs elsewhere. (b) Test-set success rate of models trained to complete BabyAI’s GoToLocal task (shaded regions indicate 95% confidence intervals around the mean); (c) Test-set success rate of models trained to complete the LC Corrupt S15N7 from [53] (shaded regions indicate inter-quartile ranges across 10 runs).

ment. Reproducing similar results as Chevalier-Boisvert et al. [55], we find that IL and RL-based algorithms can be used to train agents to nearly perfect test-set accuracy within a few million steps but that IL methods converge far faster in general. While Chevalier-Boisvert et al. trained their RL-based models only with PPO, our experiments suggest that A2C can be effective but appears to be substantially less sample efficient for this task.

Support for simultaneous and sequential algorithms. Beyond offering popular online and offline algorithms, AllenAct facilitates training with multiple losses in sequence (using a TRAINING-PIPELINE) or in parallel. Sequential training, particularly IL \rightarrow RL, is widely adopted in E-AI to *warm-start* agents [17, 16, 27] and has recently been formally studied [53]. Reproducing similar results as [53], we find that a IL \rightarrow RL combination of BC, DAgger, or BC^{cf=1} (as variants of IL) followed by PPO can significantly boost their individual (near-zero) performance (see Fig. 3c). Moreover ADV, an adaptive, parallel, combination of IL and RL losses, performs the best on the challenging, MiniGrid-based, LC CORRUPT task [53] (see Appendix B.3 for details). In line with [53], we choose hyperparameters via random search with 50 samples for each baseline. For each baseline, the best hyperparameter combination (out of the 50) is selected and we train 10 models with different seeds. We summarize the test-set performance across these runs by plotting their medians and inter-quartile ranges.

Support for multi-agent systems. We reproduce the Markov Stag Hunt [57]. In our re-implementation, two agents navigate in a 5×5 map with 1 Stag and 2 Plants, for 45 steps. Agents can move in any of the four cardinal directions to goal of collecting a Stag or Plant. An agent gains +1 reward when co-located with a Plant but can obtain a much larger reward (+5) if it coordinates with the other agent so that they occupy the same location as the Stag simultaneously. Alternatively, the agent receives a $-g$ penalty when it is co-located with the Stag but the other agent is not. A Stag or Plant disappears when collected and then re-spawns randomly. As in [57], the Stag moves towards to the closet agent at each step. Our reproduced model achieves 124.7 ± 2.47 , and 57.0 ± 4.21 reward as $g = \{0.5, 3.0\}$ at $3M$ training steps (roughly $66.7K$ training episodes). A similar baseline model converges to ≈ 120 reward at roughly $90K$ training episodes in [57]. As seen in prior work, qualitatively different behavior emerges depending on the choice of g (details in Appendix B.4).

Support for vision-language-embodiment. We re-implement ALFRED [19], which is an embodied language instruction following task. We are able to reproduce the results in the paper. More specifically, using their pre-trained model, we obtain a 4.0% task success rate and 9.4% on the goal condition success rate on the test-seen scenario.

5 Conclusion

We present AllenAct, a framework for reproducible and reusable research in the E-AI domain. Our framework provides a high degree of support (in the form of pre-trained models, starter code, and tutorials) for a growing collection of E-AI and general RL environments and tasks.

Appendix

A Generating Figure 1

Citation counts for papers in Fig 1-left were obtained using The Semantic Scholar Open Research Corpus (S2ORC) [32]. We restricted this analysis to papers submitted on arXiv within the past 6 years (2015 to 2020). Papers were determined to being in the E-AI domain if either the title or abstract contained at least 1 word from set A and at least 1 word from set B, shown below:

Set A: habitat, gibbon, igibson, ai2-thor, ai2thor, matterport, matter-port, matter-port3d, matterport3d, r2r, room-to-room, house3d, pybullet, vizdoom, chai, malmo, rlbench, procgen, touchdown, retouchdown, carla, minos, chalet, meta-world, sapien, mujoco, replica, house3d, embodiedqa

Set B: robot, agent, embodied, simulator, autonomous

The annotations for the histograms in Fig 1-right were produced manually by the authors of this submission. The top 23 papers (sorted by S2ORC citation counts) were used for this analysis.

B Experiment Details

In this section we provide additional details about the experiments listed in the main paper.

B.1 Support for embodied environments and tasks.

Methods. We trained all the models using DD-PPO. We trained the ObjectNav models for 200M steps and the PointNav models for 75M steps. For all the models we used a starting learning rate of $3e-4$ and have annealed it linearly to 0, over the course of the training run. We used rollout lengths of 30 and a γ of 0.99.

Task. We defined success on the PointNav task as taking the stop action within 0.2m from the target. We define success on the ObjectNav task as taking the stop action while looking at the target at a distance of no more than 1.0m. We used a turning angle of 30 degrees and a forward motion distance of 0.25m for all the experiments.

B.2 Support for online and offline algorithms.

In Section 4 we trained a number of different RL and IL baseline methods to complete the Go-ToLocal task in the BabyAI environment. We describe the details of these baseline methods, along with relevant hyperparameters below. For all of these experiments we use the same model used by Chevalier-Boisvert et al. [55].

Methods.

- *PPO* [37] – proximal policy optimization is an onpolicy, synchronous, RL algorithm which uses an easy-to-implement clipping methodology allowing for multiple gradient-updates with a single collection of rollouts from the agent’s policy to obtain better sample efficiency than several other popular approaches. For each update we use 1,536 rollouts of length 32 which are broken into batches of size 384×32 and iterated across 4 times (a total of 16 gradient updates per collection of rollouts). We use a fixed learning rate of 10^{-4} , a reward discount factor of $\gamma = 0.99$ and a clipping parameter of 0.1 (linearly decaying to 0 over training). Further training details can be found in our code base.
- *A2C* [58] – advantage actor critic is a synchronous variant of A3C ([38]) which, empirically, often results in better performance. While A2C has fallen out of favor, with PPO largely taking its place as the de facto onpolicy, synchronous, RL algorithm, it remains a strong comparative baseline. For each update we use 768 rollouts of length 16. With A2C, each such sample is used for only a single gradient update and the rollouts are together a single

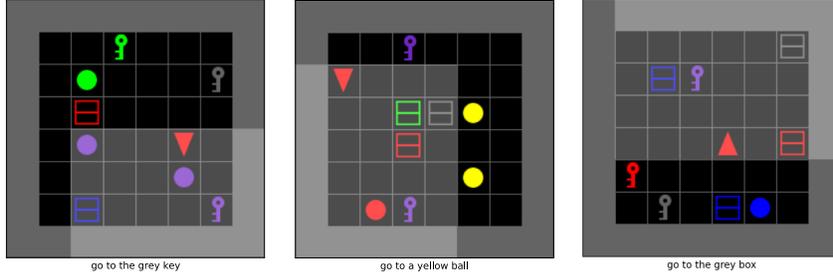


Figure 4: **BabyAI’s GoToLocal**. Three visualizations of the GoToLocal task in the BabyAI environment. At every step the agent (red triangle) obtains a 7×7 egocentric observation corresponding to the highlighted region and must follow a given “go” instruction. E.g. in the leftmost image, the agent must “go to the grey key”.

batch (this is not a limitation of AllenAct, we wished instead to remain faithful to the standard implementation of A2C). As A2C’s hyperparameters are a subset of those of PPO, we used the same hyperparameters as from PPO when applicable.

- *BC* - behavioral cloning is a straightforward variant of imitation learning in which: (i) rollouts are collected using the agent’s current policy, (ii) for each such action, we compute and store the, possibly different, expert’s action, and (iii) the agent’s policy is trained by minimizing a negative cross entropy loss between the agent’s policy and the expert action. For each such update we used 128 rollouts each of length 128. Similarly as for PPO, we use these rollouts to create 4 batches of size 32 and iterate over these batches 4 times for a total of 16 gradient updates per collection of rollouts. We use a learning rate of 10^{-3} which decays linearly to 0 throughout training.
- *Dagger* - training of DAgger is essentially identical to BC except where instead of always sampling actions from the agent’s policy we use the expert’s action with probability starting at 1 and decaying linearly to 0 during training. This means that the agent will initially see many successful rollouts which can improve training performance.
- $BC^{tf=1}$ - in this variant of behavioral cloning we *always* take the expert’s action. This is equivalent to training using a fixed dataset of expert trajectories but where this dataset is sufficiently large that no trajectory is seen more than once during training. All other hyperparameters are identical to as in BC.
- $BC^{off-policy}$ - this is imitation learning using a fixed dataset of 1 million expert demonstrations. Unlike in $BC^{tf=1}$, the agent will see the expert trajectory multiple times throughout training. While the training batches used in $BC^{off-policy}$ are of the same size as in the above IL methods, a single trajectory is not iterated over more than once until an entire epoch over the dataset is complete. In our experience this difference generally means that $BC^{off-policy}$ will obtain better results early in training (when compared to $BC^{tf=1}$) but, as $BC^{tf=1}$ is not limited to to a fixed number of experiences, with $BC^{tf=1}$ there is no fixed training set to which it can overfit.

Task. Some examples of the GoToLocal task are given in Fig. 4. At each step the agent is given a 7×7 egocentric observation and a five-word instruction defining the object to which it must navigate.

B.3 Support for simultaneous and sequential algorithms.

Methods. We test sequential IL→RL baselines comprising of the methods studied in Sec. B.2 and find that these IL→RL methods significantly improve over using only IL or RL alone. Moreover, we train ADVISOR [53] which adaptively combines imitation and rewards-based losses to bridge the ‘imitation gap’ between the expert and the agent. This is achieved via an auxiliary actor trained only via imitation loss, details of which can be found in [53]. In line with [53], we find that ADVISOR’s adaptive and parallel combination of IL and RL losses performs the best. Note that when referring to imitation in this study, the baselines learn from an expert policy. [53] lists three additional baselines where the agent learns from offline demonstrations.

Task. We deploy the above methods on the LAVA CROSSING (LC) CORRUPT (S15, N7) task

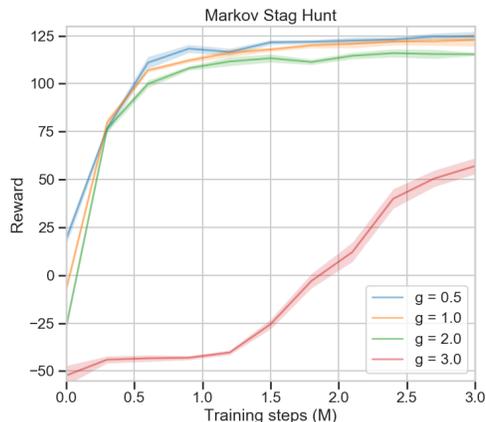


Figure 5: **Testing results for the Markov Stag Hunt over training.** Here g denotes the penalty, $g \in \{0.5, 1.0, 2.0, 3.0\}$, when only one is agent co-located with the Stag. The shaded regions indicate 95% confidence intervals around the mean over 5 different random seeds.

from [53]. This task is based on LAVACROSSING in MiniGrid environment, where an agent, only based on its egocentric view, needs to navigate to a goal while avoiding rivers of lava. Lava in a cell indicates that the episode will end if the agent steps on it. CORRUPT denotes that the (shortest-path) expert might provide corrupted supervision to the agent. Particularly, the expert policy becomes a random policy when the expert comes ≤ 10 steps from the goal. This tackles a realistic challenge of training agents which learn despite corruption or noise. S15 indicates that the grid is 15×15 in size. N7 marks that there are a total of 7 horizontal and vertical lava rivers in the environment.

B.4 Support for multi-agent systems.

We utilize MultiGrid⁶ to reproduce the Markov Stag Hunt [57]. MultiGrid, which is built within MiniGrid [36], is a grid-world environment developed for studying multi-agent reinforcement learning methods. On a 5×5 map, the agent and the Plant occupy 1×1 tiles, while the Stag occupies a 2×2 area. Following [57], we use the full map as the observation to the agents and this observation explicitly encodes entities’ attributes with indices predefined by the MiniGrid. Our model includes an embedding layer with hidden size 8 to encode the observation, a GRU with hidden size 512 to process long-term memory, and a linear layer for actor-critic output (i.e., distribution over 4 possible actions and value estimation). Both agents share the same embedding layer and GRU, while the linear layer’s parameters are not shared.

We train the agents for 3M steps (roughly $66.7K$ training episodes). We evaluate the learned model over 1000 testing trajectories with 5 different random seeds. As a result, we compute the average reward with its standard deviation over 5 different random seeds. Test-time performance over training is shown in Fig. 5. We observe that the agents converge to the payoff-dominant equilibrium when $g = \{0.5, 1.0, 2.0\}$ and risk-dominant equilibrium when $g = 3.0$. In other words, the trained agents learn to cooperate and collect the Stag when $g = \{0.5, 1.0, 2.0\}$ but, when $g = 3.0$, learn instead to individually focus on collecting the Plants. Thus, the payoff-dominant agents receive higher rewards than the risk-dominant agents.

⁶<https://github.com/ArnaudFickinger/gym-multigrid>

References

- [1] F. Xia, A. R. Zamir, Z.-Y. He, A. Sax, J. Malik, and S. Savarese. Gibson env: real-world perception for embodied agents. In *CVPR*, 2018.
- [2] M. Savva, A. Kadian, O. 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 *ICCV*, 2019. first three authors contributed equally.
- [3] E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Herrasti, D. Gordon, Y. Zhu, A. Gupta, and A. Farhadi. AI2-THOR: An Interactive 3D Environment for Visual AI. *arXiv*, 2017.
- [4] X. Puig, K. Ra, M. Boben, J. Li, T. Wang, S. Fidler, and A. Torralba. Virtualhome: Simulating household activities via programs. In *CVPR*, 2018.
- [5] S. James, Z. Ma, D. Rovick Arrojo, and A. J. Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 2020.
- [6] F. Xiang, Y. Qin, K. Mo, Y. Xia, H. Zhu, F. Liu, M. Liu, H. Jiang, Y. Yuan, H. Wang, L. Yi, A. X. Chang, L. J. Guibas, and H. Su. SAPIEN: A simulated part-based interactive environment. In *CVPR*, 2020.
- [7] 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 *CoRL*, 2019.
- [8] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. *IROS*, 2012.
- [9] C. Gan, J. I. Schwartz, S. Alter, M. Schrimpf, J. Traer, J. L. de Freitas, J. Kubilius, A. Bhandwaldar, N. Haber, M. Sano, K. Kim, E. Wang, D. Mrowca, M. Lingelbach, A. Curtis, K. T. Feigelis, D. M. Bear, D. Gutfreund, D. Cox, J. J. DiCarlo, J. H. McDermott, J. B. Tenenbaum, and D. L. K. Yamins. Threedworld: A platform for interactive multi-modal physical simulation. *arXiv*, 2020.
- [10] M. Deitke, W. Han, A. Herrasti, A. Kembhavi, E. Kolve, R. Mottaghi, J. Salvador, D. Schwenk, E. VanderBilt, M. Wallingford, L. Weihs, M. Yatskar, and A. Farhadi. RoboTHOR: An Open Simulation-to-Real Embodied AI Platform. In *CVPR*, 2020.
- [11] F. Xia, W. B. Shen, C. Li, P. Kasimbeg, M. Tchapmi, A. Toshev, R. Martn-Martn, and S. Savarese. Interactive gibbon benchmark: A benchmark for interactive navigation in cluttered environments. *IEEE Robotics and Automation Letters*, 2020.
- [12] W. Yang, X. Wang, A. Farhadi, A. Gupta, and R. Mottaghi. Visual semantic navigation using scene priors. In *ICLR*, 2019.
- [13] 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. In *ICLR*, 2020.
- [14] D. S. Chaplot, D. Gandhi, S. Gupta, A. Gupta, and R. Salakhutdinov. Learning to explore using active neural slam. In *ICLR*, 2020.
- [15] D. Gordon, A. Kembhavi, M. Rastegari, J. Redmon, D. Fox, and A. Farhadi. IQA: Visual question answering in interactive environments. In *CVPR*, 2018.
- [16] A. Das, S. Datta, G. Gkioxari, S. Lee, D. Parikh, and D. Batra. Embodied Question Answering. In *CVPR*, 2018.
- [17] Y. Zhu, D. Gordon, E. Kolve, D. Fox, L. Fei-Fei, A. Gupta, R. Mottaghi, and A. Farhadi. Visual semantic planning using deep successor representations. In *ICCV*, 2017.
- [18] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sunderhauf, I. Reid, S. Gould, and A. van den Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *CVPR*, 2018.

- [19] 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 *CVPR*, 2020.
- [20] H. Mehta, Y. Artzi, J. Baldridge, E. Ie, and P. Mirowski. Retouchdown: Adding touchdown to streetlearn as a shareable resource for language grounding tasks in street view. *arXiv*, 2020.
- [21] Y. Qi, Q. Wu, P. Anderson, M. Liu, C. Shen, and A. van den Hengel. Reverie: Remote embodied visual referring expression in real indoor environments. In *CVPR*, 2020.
- [22] S. Levine, P. Pastor, A. Krizhevsky, and D. Quillen. Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *IJRR*, 2018.
- [23] A. Gupta, A. Murali, D. Gandhi, and L. Pinto. Robot learning in homes: Improving generalization and reducing dataset bias. In *NeurIPS*, 2018.
- [24] L. Fan, Y. Zhu, J. Zhu, Z. Liu, O. Zeng, A. Gupta, J. Creus-Costa, S. Savarese, and L. Fei-Fei. Surreal: Open-source reinforcement learning framework and robot manipulation benchmark. In *CoRL*, 2018.
- [25] Y.-W. Lee, E. S. Hu, Z. Yang, A. C. Yin, and J. J. Lim. Ikea furniture assembly environment for long-horizon complex manipulation tasks. *arXiv*, 2019.
- [26] K.-H. Zeng, R. Mottaghi, L. Weihs, and A. Farhadi. Visual reaction: Learning to play catch with your drone. In *CVPR*, 2020.
- [27] U. Jain, L. Weihs, E. Kolve, M. Rastegari, S. Lazebnik, A. Farhadi, A. G. Schwing, and A. Kembhavi. Two body problem: Collaborative visual task completion. *CVPR*, 2019. first two authors contributed equally.
- [28] U. Jain, L. Weihs, E. Kolve, A. Farhadi, S. Lazebnik, A. Kembhavi, and A. G. Schwing. A cordial sync: Going beyond marginal policies for multi-agent embodied tasks. *ECCV*, 2020. first two authors contributed equally.
- [29] L. Weihs, A. Kembhavi, W. Han, A. Herrasti, E. Kolve, D. Schwenk, R. Mottaghi, and A. Farhadi. Artificial agents learn flexible visual representations by playing a hiding game. *arXiv*, 2019.
- [30] P. Battaglia, R. Pascanu, M. Lai, D. Jimenez Rezende, and K. Kavukcuoglu. Interaction networks for learning about objects, relations and physics. In *NeurIPS*, 2016.
- [31] M. Lohmann, J. Salvador, A. Kembhavi, and R. Mottaghi. Learning about objects by learning to interact with them. *arXiv*, 2020.
- [32] K. Lo, L. L. Wang, M. Neumann, R. Kinney, and D. S. Weld. S2ORC: The Semantic Scholar Open Research Corpus. In *ACL*, 2020.
- [33] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [34] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba. Places: A 10 million image database for scene recognition. *PAMI*, 2018.
- [35] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019.
- [36] M. Chevalier-Boisvert, L. Willems, and S. Pal. Minimalistic gridworld environment for openai gym. <https://github.com/maximecb/gym-minigrid>, 2018.
- [37] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv*, 2017.

- [38] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *ICML*, 2016.
- [39] S. Ross, G. J. Gordon, and J. A. Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *AISTATS*, 2011.
- [40] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 2013.
- [41] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun. CARLA: An open urban driving simulator. In *CORL*, 2017.
- [42] M. Kempka, M. Wydmuch, G. Runc, J. Toczek, and W. Jaskowski. Vizdoom: A doom-based ai research platform for visual reinforcement learning. *IEEE Conference on Computational Intelligence and Games*, 2016.
- [43] M. Johnson, K. Hofmann, T. Hutton, and D. Bignell. The malmo platform for artificial intelligence experimentation. In *IJCAI*, 2016.
- [44] Y. Song, J. Wang, T. Lukasiewicz, Z. Xu, M. Xu, Z. Ding, and L. Wu. Arena: A general evaluation platform and building toolkit for multi-agent intelligence. *arXiv*, 2019.
- [45] D. K. Misra, A. Bennett, V. Blukis, E. Niklasson, M. Shatkhin, and Y. Artzi. Mapping instructions to actions in 3d environments with visual goal prediction. In *EMNLP*, 2018.
- [46] B. León, S. Ulbrich, R. Diankov, G. Pucho, M. Przybylski, A. Morales, T. Asfour, S. Moio, J. Bohg, J. Kuffner, and R. Dillmann. Opengrasp: A toolkit for robot grasping simulation. In *SIMPAR*, 2010.
- [47] G. Kootstra, M. Popović, J. A. Jørgensen, D. Kragic, H. G. Petersen, and N. Krüger. Visgrab: A benchmark for vision-based grasping. *Paladyn*, 2012.
- [48] Y. Tassa, Y. Doron, A. Muldal, T. Erez, Y. Li, D. de Las Casas, D. Budden, A. Abdolmaleki, J. Merel, A. Lefrancq, T. Lillicrap, and M. Riedmiller. Deepmind control suite. *arXiv*, 2018.
- [49] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. OpenAI Gym. *arXiv*, 2016.
- [50] A. Juliani, V.-P. Berges, E. Teng, A. Cohen, J. Harper, C. Elion, C. Goy, Y. Gao, H. Henry, M. Mattar, and D. Lange. Unity: A general platform for intelligent agents. *arXiv*, 2018.
- [51] B. Delhaisse, L. D. Roza, and D. G. Caldwell. Pyrobolearn: A python framework for robot learning practitioners. In *CoRL*, 2019.
- [52] P. S. Castro, S. Moitra, C. Gelada, S. Kumar, and M. G. Bellemare. Dopamine: A Research Framework for Deep Reinforcement Learning. *arXiv*, 2018.
- [53] L. Weihs, U. Jain, J. Salvador, S. Lazebnik, A. Kembhavi, and A. Schwing. Bridging the imitation gap by adaptive insubordination. *arXiv*, 2020. first two authors contributed equally.
- [54] P. Anderson, A. Chang, D. S. Chaplot, A. Dosovitskiy, S. Gupta, V. Koltun, J. Kosecka, J. Malik, R. Mottaghi, M. Savva, and A. R. Zamir. On evaluation of embodied navigation agents. *arXiv*, 2018.
- [55] M. Chevalier-Boisvert, D. Bahdanau, S. Lahlou, L. Willems, C. Saharia, T. H. Nguyen, and Y. Bengio. BabyAI: First steps towards grounded language learning with a human in the loop. In *ICLR*, 2019.
- [56] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell. Curiosity-driven exploration by self-supervised prediction. In *ICML*, 2017.
- [57] A. Peysakhovich and A. Lerer. Prosocial learning agents solve generalized stag hunts better than selfish ones. In *AAMAS*, 2018.
- [58] Y. Wu, E. Mansimov, R. B. Grosse, S. Liao, and J. Ba. Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation. In *NeurIPS*, 2017.