

What Matters for Adversarial Imitation Learning?

Manu Orsini,^{*} Anton Raichuk,^{*} Léonard Huszenot,^{*†}
 Damien Vincent, Robert Dadashi, Sertan Girgin,
 Matthieu Geist, Olivier Bachem, Olivier Pietquin, Marcin Andrychowicz[‡]

Google Research, Brain Team

Abstract

Adversarial imitation learning has become a popular framework for imitation in continuous control. Over the years, several variations of its components were proposed to enhance the performance of the learned policies as well as the sample complexity of the algorithm. In practice, these choices are rarely tested all together in rigorous empirical studies. It is therefore difficult to discuss and understand what choices, among the high-level algorithmic options as well as low-level implementation details, matter. To tackle this issue, we implement more than 50 of these choices in a generic adversarial imitation learning framework and investigate their impacts in a large-scale study ($>500k$ trained agents) with both synthetic and human-generated demonstrations. While many of our findings confirm common practices, some of them are surprising or even contradict prior work. In particular, our results suggest that artificial demonstrations are not a good proxy for human data and that the very common practice of evaluating imitation algorithms only with synthetic demonstrations may lead to algorithms which perform poorly in the more realistic scenarios with human demonstrations.

1 Introduction

Reinforcement Learning (RL) has shown its ability to perform complex tasks in contexts where clear reward functions can be set-up (e.g. +1 for winning a chess game) [16, 38, 41, 44] but for many real-world applications, designing a correct reward function is either tedious or impossible [21], while demonstrating a correct behavior is often easy and cheap. Therefore, imitation learning (IL, [4, 7]) might be the key to unlock the resolution of more complex tasks, such as autonomous driving, for which reward functions are much harder to design.

The simplest approach to IL is Behavioral Cloning (BC, [2]) which uses supervised learning to predict the expert’s action for any given state. However, BC is often unreliable as prediction errors compound in the course of an episode. Adversarial Imitation Learning (AIL, [15]) aims to remedy this using inspiration from Generative Adversarial Networks (GANs, [9]) and Inverse RL [3, 5, 6]: the policy is trained to generate trajectories that are indistinguishable from the expert’s ones. As in GANs, this is formalized as a two-player game where a discriminator is co-trained to distinguish between the policy and expert trajectories (or states). See App. B for a brief introduction to AIL.

A myriad of improvements over the original AIL algorithm were proposed over the years [18, 32, 42, 45, 47], from changing the discriminator’s loss function [18] to switching from on-policy to off-policy agents [32]. However, their relative performance is rarely studied in a controlled setting, and never these changes were all compared together. The performance of these high-level choices may also depend on the low-level implementation details, which might be not even mentioned in the

^{*}Equal contribution.

[†]Univ. de Lille, CNRS, Inria Scool, UMR 9189 CRISyAL.

[‡]Corresponding author. E-mail: marcina@google.com.

publications [20, 30, 37, 43], as well as the hyperparameters (HPs) used. Thus, assessing whether the proposed changes are the reason for the presented improvements becomes extremely difficult. This lack of proper comparisons slows down the overall research in imitation learning and the industrial applicability of these methods.

We investigate such high- and low-level choices in depth and study their impact on the algorithm performance. Hence, as **our key contributions**, we (1) implement a highly-configurable generic AIL algorithm, with various axes of variation (>50 HPs), including 4 different RL algorithms and 7 regularization schemes for the discriminator, (2) conduct a large-scale ($>500k$ trained agents) experimental study on 10 continuous-control tasks⁴ and (3) analyze the experimental results to provide practical insights and recommendations for designing novel and using existing AIL algorithms.

Most surprising finding #1: regularizers. While many of our findings confirm common practices in AIL research, some of them are surprising or even contradict prior work. In particular, we find that standard regularizers from Supervised Learning — dropout [10] and weight decay [1] often perform similarly to the regularizers designed specifically for adversarial learning like gradient penalty [19]. Moreover, for easier environments (which were often the only ones used in prior work), we find that it is possible to achieve excellent results without using any explicit discriminator regularization.

Most surprising finding #2: human demonstrations. Not only does the performance of AIL heavily depend on whether the demonstrations were collected from a human operator or generated by an RL algorithm, but the relative performance of algorithmic choices also depends on the demonstration source. Our results suggest that artificial demonstrations are not a good proxy for human data and that the very common practice of evaluating IL algorithms only with synthetic demonstrations may lead to algorithms which perform poorly in the more realistic scenarios with human demonstrations.

Paper outline. In Sec. 2, we describe our experimental setting and the performance metrics used. We then present and analyze the results related to the agent (Sec. 3) and discriminator (Sec. 4) training. Afterwards, we compare RL-generated and human-collected demonstrations (Sec. 5) and analyze the choices influencing the computational cost of running the algorithm (Sec. 6). The appendices contain the details of the different algorithmic choices in AIL (App. C) as well as the raw results of the experiments (App. F–H).

2 Experimental design

Environments. We focus on continuous-control tasks as robotics appears as one of the main potential applications of IL and a vast majority of the IL literature thus focuses on it. In particular, we run experiments with five widely used environments from OpenAI Gym [14]: HalfCheetah-v2, Hopper-v2, Walker2d-v2, Ant-v2, and Humanoid-v2 and three manipulation environments from Adroit [22]: pen-v0, door-v0, and hammer-v0. All the environments are shown in Fig. 1. The Adroit tasks consist in aligning a pen with a target orientation, opening a door and hammering a nail with a 5-fingered hand. These two benchmarks bring orthogonal contributions. The former focuses on locomotion but has 5 environments with different state/action dimensionality. The latter, more varied in term of tasks, has an almost constant state-action space.

Demonstrations. For the Gym tasks, we generate demonstrations with a SAC [29] agent trained on the environment reward. For the Adroit environments, we use the “expert” and “human” datasets from D4RL [46],⁵ which are, respectively, generated by an RL agent and collected from a human operator. As far as we know, our work is the first to solve these tasks with human datasets in the imitation setup (most of the prior work concentrated on Offline RL). For all environments, we use 11 demonstration trajectories. Following prior work [15, 32, 47], we subsample expert demonstrations by only using every 20th state-action pair to make the tasks harder.

Adversarial Imitation Learning algorithms. We researched prior work on AIL algorithms and made a list of commonly used design decisions like policy objectives or discriminator regularization techniques. We also included a number of natural options which we have not encountered in literature

⁴A task is defined by an environment and the demonstrator type (either human or RL agent).

⁵For pen, we only use the “expert” dataset, the “human” one consisting of a single (yet very long) trajectory.

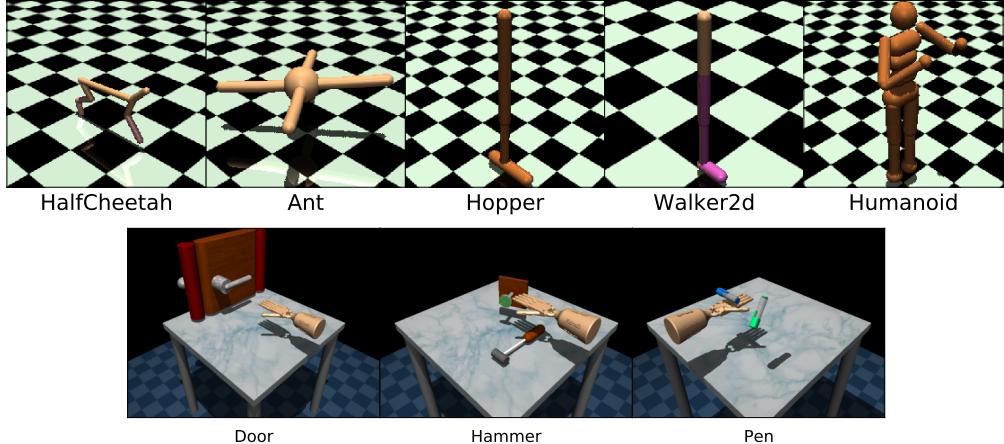


Figure 1: Environments: OpenAI Gym (top) and Adroit (bottom).

(e.g. dropout [10] in the discriminator or clipping rewards bigger than a threshold). All choices are listed and explained in App. C. Then, we implemented a single highly-configurable AIL agent which exposes all these choices as configuration options in the Acme framework [50] using JAX [26] for automatic differentiation and Flax [48] for neural networks computation. The configuration space is so wide that it covers the whole family of AIL algorithms, in particular, it mostly covers the setups from AIRL⁶ [18] and DAC [32]. We plan to open source the agent implementation.

Experimental design. We created a large HP sweep (57 HPs swept, >120k agents trained) in which each HP is sampled uniformly at random from a discrete set and independently from the other HPs. We manually ensured that the sampling ranges of all HPs are appropriate and cover the optimal values. Then, we analyzed the results of this initial experiment (called *wide*, detailed description and results in App. F), removed clearly suboptimal options and ran another experiment with the pruned sampling ranges (called *main*, 43 HPs swept, >250k agents trained, detailed description and results in App. G). The latter experiment serves as the basis for most of the conclusions drawn in this paper but we also run a few additional experiments to investigate some additional questions (App. H and App. I).

This pruning of the HP space guarantees that we draw conclusions based on training configurations which are highly competitive (training curves can be found in Fig. 27) while using a large HP sweep (including, for example, multiple different RL algorithms) ensures that our conclusions are robust and valid not only for a single RL algorithm and specific values of HPs, but are more generally applicable. Moreover, many choices may have strong interactions with other related choices, for example we find a surprisingly strong interaction between the discriminator regularization scheme and the discriminator learning rate (Sec. 4). This means that such choices need to be tuned together (as it is the case in our study) and experiments where only a single choice is varied but the interacting choices are kept fixed may lead to misleading conclusions.

Performance measure. For each HP configuration and each of the 10 environment-dataset pairs we train a policy and evaluate it 10 times through the training by running it for 50 episodes and computing the average undiscounted return using the environment reward. We then average these scores to obtain a single performance score which approximates the area under the learning curve. This ensures we assign higher scores to HP configurations that learn quickly.

Analysis. We consider two different analyses for each choice⁷:

Conditional 95th percentile: For each potential value of that choice (e.g., RL Algorithm = PPO), we look at the performance distribution of sampled configurations with that value. We report the 95th

⁶ Seminal AIRL uses TRPO [13] to train the policy, not supported in our implementation (PPO [23] used here).

⁷ This analysis is based on a similar type of study focused on on-policy RL algorithms [43].

percentile of the performance as well as error bars based on bootstrapping.⁸ This corresponds to an estimate of the performance one can expect if all other choices were tuned with random search and a limited budget of roughly 13 HP configurations⁹. All scores are normalized so that 0 corresponds to a random policy and 1 to the expert performance (expert scores can be found in App. E).

Distribution of choice within top 5% configurations. We further consider for each choice the distribution of values among the top 5% HP configurations. In particular, we measure the ratio of the frequency of the given value in the top 5% of HP configurations with the best performance to the frequency of this value among all HP configurations. If certain values are over-represented in the top models (ratio higher than 1), this indicates that the specific choice is important for good performance.

3 What matters for the agent training?

Summary of key findings. The AIRL reward function perform best for synthetic demonstrations while $-\ln(1 - D)$ is better for human demonstrations. Using explicit absorbing state is crucial in environments with variable length episodes. Observation normalization strongly affects the performance. Using an off-policy RL algorithm is necessary for good sample complexity while replaying expert data and pretraining with BC improves the performance only slightly.

Implicit reward function. In this section, we investigate choices related to agent training with AIL, the most salient of which is probably the choice of the implicit reward function. Let $D(s, a)$ be the probability of classifying the given state-action pair as *expert* by the discriminator¹⁰. In particular, we run experiments with the following reward functions: $r(s, a) = -\log(1 - D(s, a))$ (used in the original GAIL paper [15]), $r(s, a) = \log D(s, a) - \log(1 - D(s, a))$ (called the AIRL reward [18]), $r(s, a) = \log \bar{D}(s, a)$ (a natural choice we have not encountered in literature), and the FAIRL [47] reward function $r(s, a) = -h(s, a) \cdot e^{h(s, a)}$, where $h(s, a)$ is the discriminator logit¹¹. It can be shown that, under the assumption that all episodes have the same length, maximizing these reward functions corresponds to the minimization of different divergences between the marginal state-action distribution of the expert and the policy. See [47] for an in-depth discussion on this topic. We also consider clipping the rewards with absolute values bigger than a threshold which is a HP.

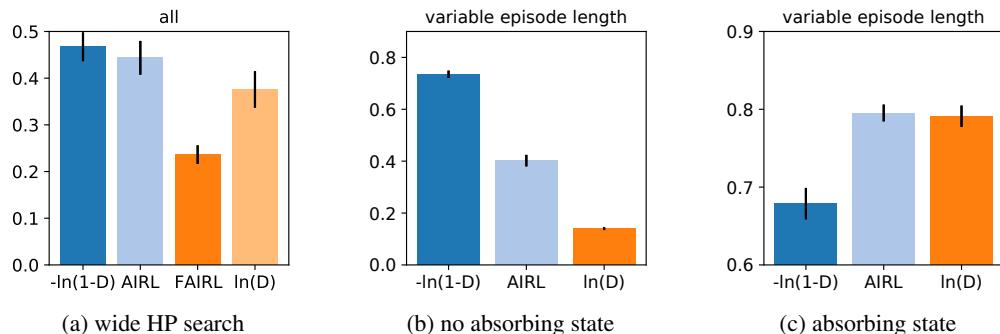


Figure 2: Comparison of different reward functions. The bars show the 95th percentile across HPs sampling of the *average* policy performance during training. Plot (a) shows the results averaged across all 10 tasks. Plots (b) and (c) show the performance on the subset of environments with variable length episodes when the absorbing state is disabled (b) or enabled (c). See Fig. 15 and Fig. 75 for the individual results in all environments.

The FAIRL reward performed much worse than all others in the initial wide experiment (Fig. 2a) and therefore was not included in our main experiment. This is mostly caused by its inferior performance with off-policy RL algorithms (Fig. 25). Moreover, reward clipping significantly helps the FAIRL

⁸We compute each metric 20 times based on a randomly selected half of all training runs, and then report the mean of these 20 measurements while the error bars show mean-std and mean+std.

⁹The probability that all 13 configurations score worse than the 95th percentile is equal $0.95^{13} \approx 50\%$.

¹⁰Some prior works, including GAIL[15], use the opposite notation, with $D(s, a)$ the *non-expert* probability.

¹¹It can also be expressed as $h(s, a) = \log D(s, a) - \log(1 - D(s, a))$.

reward (Fig. 26) while it does not help the other reward functions apart from some small gains for $-\ln(1 - D)$ (Fig. 85). Therefore, we suspect that the poor performance of the FAIRL reward function may be caused by its exponential term which may have very high magnitudes. Moreover, the FAIRL paper [47] mentions that the FAIRL reward is more sensitive to HPs than other reward functions which could also explain its poor performance in our experiments.

Fig. 28 shows that the $\ln(D)$ reward function performs a bit worse than the other two reward functions in the main experiment. Five out of the ten tasks used in our experiments have variable length episodes with longer episodes correlated with better behaviour¹² (Hopper, Walker2d, Ant, Humanoid, pen) — on these tasks we can notice that $r(s, a) = -\ln(1 - D(s, a))$ often performs best and $r(s, a) = \ln D(s, a)$ worst. This can be explained by the fact that $-\ln(1 - D(s, a)) > 0$ and $\ln D(s, a) < 0$ which means that the former reward encourages longer episodes and the latter one shorter ones [32]. Absorbing state (described in App. C.2) is a technique introduced in the DAC paper [32] to mitigate the mentioned bias and encourage the policy to generate episodes of similar length to demonstrations. In Fig. 2b-c we show how the performance of different reward functions compares in the environments with variable length episodes depending on whether the absorbing state is used. We can notice that without the absorbing state $r(s, a) = -\ln(1 - D(s, a)) > 0$ performs much better in the environments with variable episode length which suggests that the learning is driven to a large extent by the reward bias and not actual imitation of the expert behaviour [32]. This effect disappears when the absorbing state is enabled (Fig. 2c).

Fig. 75 shows the performance of different reward functions in all environments conditioned on whether the absorbing state is used. If the absorbing state is used, the AIRL reward function performs best in all the environments with RL-generated demonstrations, and $\ln(D)$ performs only marginally worse. The $-\ln(1 - D)$ reward function underperforms on the Humanoid and pen tasks while performing best with human datasets. We provide some hypothesis for this behaviour in Sec. 5, where we discuss human demonstrations in more details.

Observation normalization. We consider observation normalization which is applied to the inputs of all neural networks involved in AIL (policy, critic and discriminator). The normalization aims to transform the observations so that each observation coordinate has mean 0 and standard deviation 1. In particular, we consider computing the normalization statistics either using only the expert demonstrations so that the normalization is *fixed* throughout the training, or using data from the policy being trained (called *online*). See App. C.6 for more details. Fig. 3 shows that input normalization significantly influences the performance with the effects on performance being often much larger than those of algorithmic choices like the reward function or RL algorithm used. Surprisingly, normalizing observations can either significantly improve or diminish performance and whether the fixed or online normalization performs better is also environment dependent.

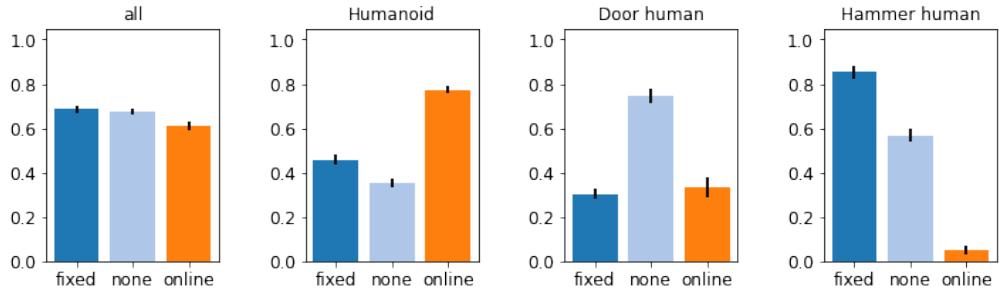


Figure 3: Comparison of observation normalization schemes. The bars show the 95th percentile of performance. The leftmost plot shows the results averaged across all 10 tasks. See Fig. 29 for the results on all environments.

Replaying expert data. When demonstrations as well as external rewards are available, it is common for RL algorithms to sample batches for off-policy updates from the demonstrations in addition to the replay buffer [31, 40]. We varied the ratio of the policy to expert data being replayed

¹²The episodes are terminated earlier if the simulated robot falls over or if the pen is dropped.

but found only very minor gains (Fig. 86). Moreover, in the cases when we see some benefits, it is usually best to replay 16–64 times more policy than expert data. On some tasks (Humanoid) replaying even a single expert transitions every 256 agent ones significantly hurts performance. We suspect that, in contrast to RL with demonstrations, we see little benefit from replaying expert data in the setup with learned rewards because (1) replaying expert data mostly helps when the reward signal is sparse (not the case for discriminator-based rewards), and (2) discriminator may overfit to the expert demonstrations which could result in incorrectly high rewards being assigned to expert transitions.

Pretraining with BC. We also experiment with pretraining a policy with Behavioral Cloning (BC, [2]) at the beginning of training. Despite starting from a much better policy than a random one, we usually observe that the policy quality deteriorates quickly at the beginning of training (see the pen task in Fig. 6) due to being updated using randomly initialized critic and discriminator networks, and the overall gain from pretraining is very small in most environments (Fig. 30).

RL algorithms. We run experiments with four different RL algorithms, three of which are off-policy algorithms (SAC [29], TD3 [27] and D4PG [25]), as well as PPO [23] which is nearly on-policy. Fig. 4 shows that the sample complexity of PPO is significantly worse than that of the off-policy algorithms while all off-policy algorithms perform overall similarly.

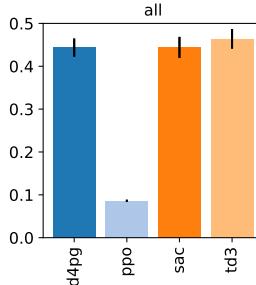


Figure 4: Comparison of RL algorithms (wide HP search). See Fig. 10 for the results on individual environments.

RL algorithms HPs. Fig. 11 shows that the discount factor is one of the most important HPs with the values of 0.97 – 0.99 performing well on all tasks. Fig. 32 shows that in most environments it is better not to erase any data from the RL replay buffer and always sample from all the experience encountered so far. It is common in RL to use a noise-free version of the policy during evaluation and we observe that it indeed improves the performance (Fig. 33). The policy MLP size does not matter much (Figs. 34–35) while bigger critic networks perform significantly better¹³ (Figs. 12–13). Regarding activation functions¹⁴, relu performs on par or better than tanh in all environments apart from door in which tanh is significantly better (Fig. 36). Our implementation of TD3 optionally applies gradient clipping¹⁵ but it does not affect the performance much (Fig. 37). D4PG can use n-step returns, this improves the performance on the Adroit tasks but hurts on the Gym suite (Fig. 38).

4 What matters for the discriminator training?

Summary of key findings. MLP discriminators perform on par or better than AIL-specific architectures. Explicit discriminator regularization is only important in more complicated environments (Humanoid and harder ones). Spectral norm is overall the best regularizer but standard regularizers from supervised learning often perform on par. Optimal learning rate for the discriminator may be 2–2.5 orders of magnitude lower than the one for the RL agent.

Discriminator input. In this section we look at the choices related to the discriminator training. Fig. 52 shows how the performance depends on the discriminator input. We can observe that while

¹³We thus only include critics with at least two hidden layers with the size at least 128 in the main experiment.

¹⁴We use the same activation function in the policy and critic networks.

¹⁵The reason for that is that the DAC paper [32] uses TD3 with gradient clipping.

it is beneficial to feed actions as well as states to the discriminator, the state-only demonstrations perform almost as well. Interestingly, on the door task with human data, it is better to ignore the expert actions. We explore the results with human demonstrations in more depth in Sec. 5.

Discriminator architecture. Regarding the discriminator network, our basic architecture is an MLP but we also consider two modifications introduced in AIRL [18]: a reward shaping term and a log $\pi(a|s)$ logit shift which introduces a dependence on the current policy (only applicable to RL algorithms with stochastic policies, which in our case are PPO and SAC). See App. C.3 for a detailed description of these techniques. Fig. 16 shows that the logit shift significantly hurts the performance. This is mainly due to the fact that it does not work well with SAC which is off-policy (Fig. 24). Fig. 53 shows that the shaping term does not affect the performance much. While the modifications from AIRL does not improve the sample complexity in our experiments, it is worth mentioning that they were introduced for another purpose, namely the recovery of transferable reward functions.

Regarding the size of the discriminator MLP(s), the best results on all tasks are obtained with a single hidden layer (Fig. 54), while the size of the hidden layer is of secondary importance (if it is not very small) with the exception of the tasks with human data where fewer hidden units perform significantly better (Fig. 55). All tested discriminator activation functions perform overall similarly while sigmoid performs best with human demonstrations (Fig. 56).

Discriminator training. Fig. 57 shows that it is best to use as large as possible replay buffers for sampling negative examples (i.e. agent transitions). As noticed in prior work, the initialization of the last *policy* layer can significantly influence the performance in RL [43], thus we tried initializing the last *discriminator* layer with smaller weights but it does not make much difference (Fig 58).

Discriminator regularization. An overfitting or too accurate discriminator can make agent’s training challenging, and therefore it is common to use additional regularization techniques when training the AIL discriminator (or GANs in general). We run experiments with a number of regularizers commonly used with AIL, namely Gradient Penalty [19] (GP, used e.g. in [32]), spectral norm [36] (e.g. in [45]), Mixup [24] (e.g. in [53]), as well as using the PUGAIL loss [42] instead of the standard cross entropy loss to train the discriminator. Apart from the above regularizers, we also run experiments with regularizers commonly used in Supervised Learning, namely dropout [10], the weight decay [1] variant from AdamW [34] as well as the entropy bonus of the discriminator output treated as a Bernoulli distribution. The detailed description of all these regularization techniques can be found in App. C.5.

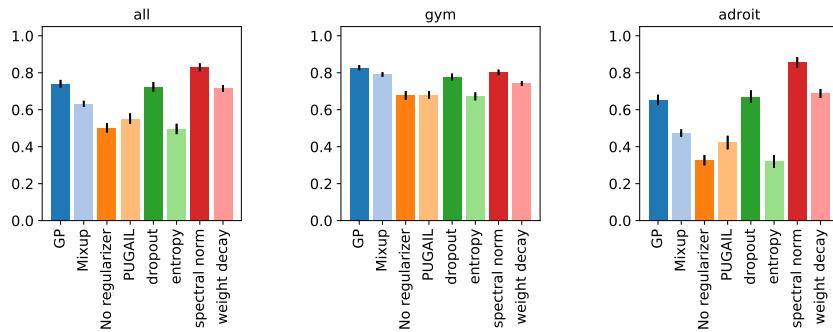


Figure 5: The 95th percentile of performance for different discriminator regularizers. The central plot shows the average performance across 5 tasks from OpenAI gym and the right one the average performance across 5 tasks from the Adroit suite. See Fig. 59 for the plots for individual environments.

Fig. 5 shows how the performance depends on the regularizer. Spectral normalization performs overall best, while GP, dropout and weight decay all perform on par with each other and only a bit worse than spectral normalization. We find this conclusion to be quite surprising given that we have not seen dropout or weight decay being used with AIL in literature. We also notice that the regularization is generally more important on harder tasks like Humanoid or the tasks in the Adroit suite (Fig. 59).

Most of the regularizers investigated in this section have their own HPs and therefore the comparison of different regularizers depends on how these HPs are sampled. As we randomly sample the regularizer-specific HPs in this analysis, our approach favours regularizers that are not too sensitive to

their HPs. At the same time, there might be regularizers that are sensitive to their HPs but for which good settings may be easily found. Fig. 70 shows that even if we condition on choosing the optimal HPs for each regularizer, the relative ranking of regularizers does not change.

Moreover, there might be correlations between the regularizer and other HPs, therefore their relative performance may depend on the distribution of all other HPs. In fact, we have found two such surprising correlations. Fig. 76 shows the performance conditioned on the regularizer used *as well* as the discriminator learning rate. We notice that for PUGAIL, entropy and no regularization, the performance significantly increases for lower discriminator learning rates and the best performing discriminator learning rate (10^{-6}) is in fact 2–2.5 orders of magnitude lower than the best learning rate for the RL algorithm (0.0001–0.0003, Figs. 17, 41, 42, 44, 50).¹⁶ On the other hand, the remaining regularizers are not too sensitive to the discriminator learning rate. This means that the performance gap between PUGAIL, entropy and no regularization and the other regularizers is to some degree caused by the fact that the former ones are more sensitive to the learning rate and may be smaller than suggested by Fig. 5 if we adjust for the appropriate choice of the discriminator learning rate. We can notice that PUGAIL and entropy are the only regularizers which only change the discriminator loss but do not affect the internals of the discriminator neural network. Given that they are the only two regularizers benefiting from very low discriminator learning rate, we suspect that it means that a very low learning rate can play a regularizing role in the absence of an explicit regularization inside the network.

Another surprising correlation is that in some environments, the regularizer interacts strongly with observation normalization (described App. C.6) employed on discriminator inputs (see Fig. 77 for an example on Ant). These two correlations highlight the difficulty of comparing regularizers, and algorithmic choices more broadly, as their performance significantly depends on the distribution of other HPs.

We also supplement our analysis by comparing the performance of different regularizers for the *best* found HPs. More precisely, we choose the best value for each HP in the main experiment (listed in App. D) and run them with different regularizers. To account for the mentioned correlations with the discriminator learning rate and observation normalization, we also include these two choices in the HP sweep and choose the best performing variant (as measure by the area under the learning curve) for each regularizer and each environment.

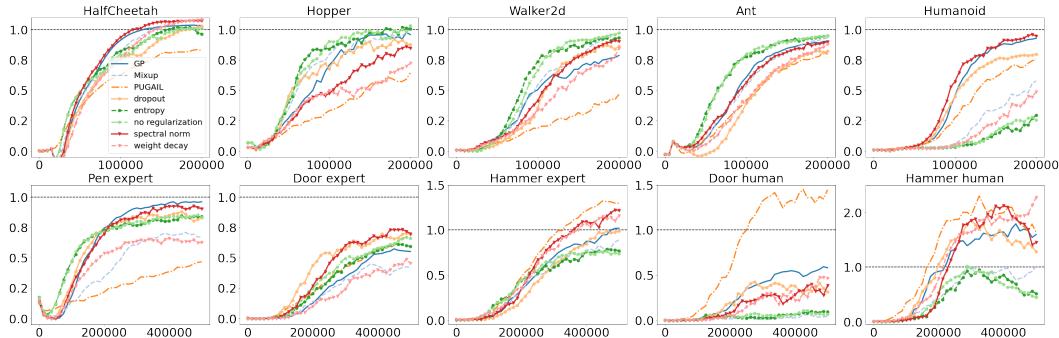


Figure 6: Learning curves for different discriminator regularizers when the other HPs are set to the best performing value across all tasks. The y-axis shows the average policy return normalized so that 0 corresponds to a random policy and 1 to the expert. See App. D for the HPs used. The plots shows the averages across 30 random seeds. Best seen in color.

While it is not guaranteed that the performance is going to be good at all because we greedily choose the best performing value for each HP and there might be some unaccounted HP correlations, we find that the performance is very competitive (Fig. 6). Notice that we use the same HPs in *all* environments¹⁷ and the performance can be probably improved by varying some HPs between the environments, or at least between the two environment suites.

We notice that on the four easiest tasks (HalfCheetah, Hopper, Walker2d, Ant), investigated discriminator regularizers provide no, or only minor performance improvements and excellent results

¹⁶The optimal learning rate for those regularizers was the smallest one included in the main experiment. We also run an additional sweep with smaller rates but found that even lower ones do not perform better (Fig. 84).

¹⁷Apart from the discriminator learning rate and observation normalization used.

can be achieved without them. On the tasks where regularization is beneficial, we usually see that there are multiple regularizers performing similarly well, with spectral normalization being one of the best regularizers in all tasks apart from the two tasks with human data where PUGAIL performs better.

Regularizers-specific HPs. For GP, the target gradient norm of 1 is slightly better in most environments but the value of 0 is significantly better in `hammer-human` (Fig. 60), while the penalty strength of 1 performs best overall (Fig. 61). For dropout, it is important to apply it not only to hidden layers but also to inputs (Fig. 62) and the best results are obtained for 50% input dropout and 75% hidden activations dropout (Figs. 62, 63 and 70). For weight decay, the optimal decay coefficient in the AIL setup is much larger than the values typically used for Supervised Learning, the value $\lambda = 10$ performs best in our experiments (Fig. 64). For Mixup, $\alpha = 1$ outperforms the other values on almost all tested environments (Fig. 65). For PUGAIL, the unbounded version performs much better on the Adroit suite, while the bounded version is better on the gym tasks (Fig. 66), and positive class prior of $\eta = 0.7$ performs well on most tasks (Fig. 67). For the discriminator entropy bonus, the values around 0.03 performed best overall (Fig. 68). All experiments with spectral normalization enforce the Lipschitz constant of 1 for each weight matrix.¹⁸

5 Are synthetic demonstrations a good proxy for human data?

Summary of key findings Human demonstrations significantly differ from synthetic ones. Learning from human demonstrations benefits more from discriminator regularization and may work better with different discriminator inputs and reward functions than RL-generated demonstrations.

Using a dataset of human demonstrations comes with a number of additional challenges. Compared to synthetic demonstrations, the human policy can be multi-modal in that for a given state different decisions might be chosen. A typical example occurs when the human demonstrator remains idle for some time (for example to think about the next action) before taking the actual relevant action: we have two modes in that state, the relevant action has a low probability while the idle action has a very high probability. The human policy might not be exactly markovian either. Those differences are significant enough that the conclusions on synthetic datasets might not hold anymore.

In this section, we focus on the Adroit door and `hammer` environments for which we run experiments with human as well as synthetic demonstrations.¹⁹ Note that on top of the aforementioned challenges, the setup with the Adroit environments using human demonstrations exhibits a few additional specifics. The demonstrations were collected letting the human decide when the task is completed: said in a different way, the demonstrator is offered an additional action to jump directly to a terminal state and this action is not available to the agent imitating the expert. The end result is a dataset of demonstrations of variable length while the agent can only generate episodes consisting of exactly 200 transitions. Note that there was no time limit imposed on the demonstrator and some of the demonstrations have a length greater than 200 transitions. Getting to the exact same state distribution as the human expert may be impossible, and imitation learning algorithms may have to make some trade-offs. The additional specificity of that setup is that the reward of the environment is not exactly what the human demonstrator optimized. In the door environment, the reward provided by the environment is the highest when the door is *fully* opened while the human might abort the task slightly before getting the highest reward. However, overall, we consider the reward provided by the environment as a reasonable metric to assess the quality of the trained policies. Moreover, in the `hammer` environment, some demonstrations have a low return and we suspect those are not successful demonstrations.²⁰

Discriminator regularization. When comparing the results for RL-generated (`adroit-expert`²¹) and human demonstrations (`adroit-human`) we can notice differences on a number of HPs related to the discriminator training. Human demonstrations benefit more from using discriminator regularizers (Fig. 59) and they also work better with smaller discriminator networks (Fig. 55) trained with lower

¹⁸There may be different Lipschitz constants for different networks depending on those of related activations.

¹⁹For `pen`, we only use the “expert” dataset, the “human” one consists of a single (yet very long) trajectory.

²⁰D4RL datasets [46] contain only the policy observations and not the simulator states and therefore it is not straightforward to visualize the demonstrations.

²¹We do not include `pen` in the `adroit-expert` plots so that both `adroit-expert` and `adroit-human` show the results averages across the `door` and `hammer` tasks and differ only in the demonstrations used.

learning rates (Fig. 69). The increased need for regularization suggest that it is easier to overfit to the idiosyncrasies of human demonstrations than to those of RL policies.

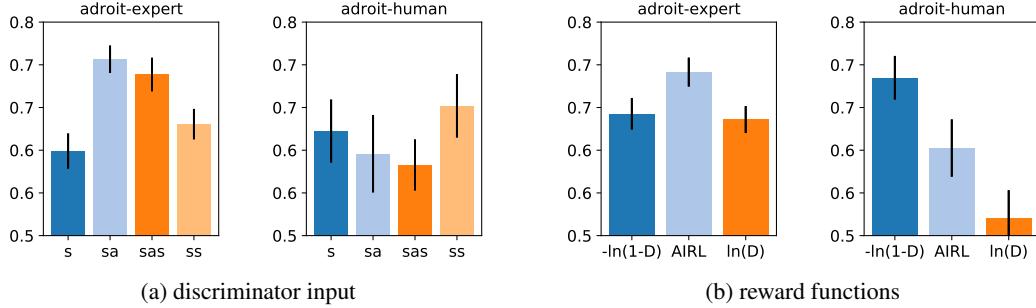


Figure 7: Comparison of discriminator inputs (a) and reward functions (b) for environments with human demonstrations. See Fig. 52 and Fig. 28 for the individual results in all environments.

Discriminator input. Fig. 7a shows the performance given the discriminator input depending on the demonstration source. For most tasks with RL-generated demonstrations, feeding actions as well as states improves the performance (Fig. 52). Yet, the opposite holds when human demonstrations are used. We suspect that it might be caused by the mentioned issue with demonstrations lengths which forces the policy to repeat a similar movement but with a different speed than the demonstrator.

Reward functions. Finally, we look at how the relative performance of different reward functions depends on the demonstration source. Fig. 7b shows that for RL-generated demonstrations the best reward function is AIRL while $-\ln(1 - D)$ performs better with human demonstrations. Under the assumption that the discriminator is optimal, these two reward functions correspond to the minimization of different divergences between the state (or state-action depending on the discriminator input) occupancy measures of the policy and the expert — See Table 1 for the details.

Table 1: Reward functions and corresponding divergences. π and E denote the state occupancy measures of, respectively, the policy and the expert. The proofs can be found in [15] and [47]. The *bounded* and *symmetric* column show whether the given *divergence* is bounded or symmetric. D denotes the probability of being classified as *expert*.

Paper	Reward	Divergence	Bounded	Symmetric
GAIL [15]	$-\ln(1 - D)$	$D_{JS}(\pi E)$	✓	✓
AIRL [18]	$\ln(D) - \ln(1 - D)$	$D_{KL}(\pi E)$	✗	✗

The reward function performing best with human demonstrations ($-\ln(1 - D)$) corresponds to the minimization of the Jensen-Shannon divergence (proof in [15]).²² Interestingly, this divergence is symmetric ($D_{JS}(\pi||E) = D_{JS}(E||\pi)$) and bounded ($0 \leq D_{JS}(\pi||E) \leq \ln(2)$). For AIRL, the symmetry means that it penalizes the policy for doing things the expert never does with exactly the same weight as for not doing some of the things the expert does while the boundedness means that the penalty for not visiting a single state is always finite. We suspect that this boundedness is beneficial for learning with human demonstrations because it may not be possible to exactly match the human distribution for the reasons explained earlier.

In contrast to Jensen-Shannon, the $D_{KL}(\pi||E)$ divergence which is optimized by the AIRL reward (proof in [47]) is neither symmetric, nor bounded — it penalizes the policy much more heavily for doing the things the expert never does than for not doing all the things the expert does and the penalty for visiting a single state the expert never visits is infinite (assuming a perfect discriminator).

While it is hard to draw any general conclusions only from the two investigated environments for which we had access to human demonstrations, our analysis shows that the differences between

²² $D_{JS}(P||Q) = D_{KL}(P||M) + KL(Q||M)$, where $M = \frac{P+Q}{2}$.

synthetic and human-generated demonstrations can influence the relative performance of different algorithmic choices. This suggests that RL-generated data are not a good proxy for human demonstrations and that the very common practice of evaluating IL only with synthetic demonstrations may lead to algorithms which perform poorly in the more realistic scenarios with human demonstrations.

6 How to train efficiently?

So far we have analysed how HPs affect the performance of AIL algorithms measured after fixed numbers of environment steps. Here we look at the HPs which influence sample complexity as well as the computational cost of running an algorithm. Raw experiment report can be found in App. H.

Batch size and replay ratio. One of the main factors influencing the throughput of a particular imitation algorithm is the number of times each transition is replayed on average and the batch size used.²³ See App. C.1 for the detailed description of the HPs involved. Fig. 79 shows that smaller batches perform overall better (given a fixed replay ratio) and increasing the replay ratio improves the performance, at least up to some threshold depending on the environment (Fig. 80). There is a very strong correlation between the two HPs — Fig. 83 shows that for most batch sizes, the optimal replay ratio is equal to the batch size, which corresponds to replaying exactly one batch of data per environment step. If we compare different batch sizes under the ratio of batches to environment steps fixed to one, the performance is mostly independent of the batch size (Fig. 83).

While in most of our experiment the discriminator and the RL agent are trained with exactly the same number of batches, we also tried doubling the number of discriminator batches. Fig. 81 shows that it improves the performance slightly on the Adroit suite.

Combining multiple batches. We also consider processing multiple batches at once for improved accelerator (GPU or TPU) utilization. In particular, we sample an N -times larger batch from a replay buffer, split it back into N smaller/proper batches on an accelerator, and process them sequentially. In order to keep the replay ratio unaffected, we decrease the frequency of updates accordingly, e.g. instead of performing one gradient update for every environment step, we perform N gradients updates every N environment steps. We apply this technique to the discriminator as well as the RL agent training. The effect on the sample complexity of the algorithm can be seen in Fig. 82. There is a small negative effect for values larger or equal to 16. The effect of this parameter on the throughput of our system could be observed in Fig. 8. The value of 8 provides a good compromise: almost no noticeable sample complexity regression while decreasing the training time by 2–3 times.

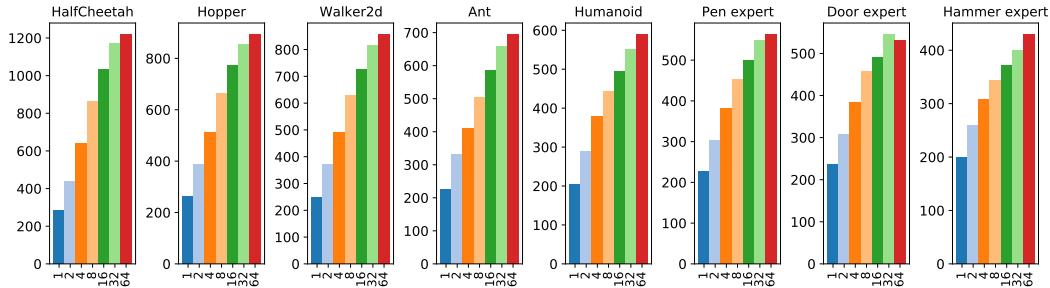


Figure 8: Training speed (in terms of environment steps per second) for combining multiple batches. The x-axis denotes the number of batches combined. Other HPs are set to the best performing value across all tasks (listed in App. D). The plots shows the averages across 10 random seeds.

²³We use the same batch size for the policy and actor networks while the discriminator batch size is effectively two times larger because its batches contain always `batch_size` (C7) demonstration transitions and `batch_size` (C7) policy transitions. The replay ratio is the same for all networks with the exception of the discriminator which can have its replay ratio doubled depending on the value of `discriminator_to_RL_updates_ratio` (C44). See App. C.4 for details.

7 Related work

The most similar work to ours is probably [45] which compares the performance of different discriminator regularizers and concludes that gradient penalty is necessary for achieving good performance with off-policy AIL algorithms. In contrast to [45], which uses a single HP configuration, we run large-scale experiments with very wide HP sweeps which allows us to reach more robust conclusions. In particular, we are able to achieve excellent sample complexity on all the environments used in [45]²⁴ without using any explicit discriminator regularizer (Fig. 6).

Another empirical study of IL algorithms is [54], which investigates the problem of HP selection in IL under the assumption that the reward function is not available for the HP selection.

The methodology of our study is mostly based on [43] which analyzed the importance of different choices for on-policy actor-critic methods. Our work is also similar to other large-scale studies done in other fields of Deep Learning, e.g. model-based RL [39], GANs [35], NLP [51], disentangled representations [33] and convolution network architectures [52].

8 Conclusions

In this empirical study, we investigate in depth many aspects of the AIL framework including discriminator architecture, training and regularization as well as many choices related to the agent training. Our key findings can be divided into three categories: (1) Corroborating prior work, e.g. for the underlying RL problem, off-policy algorithms are more sample efficient than on-policy ones; (2) Adding nuances to previous studies, e.g. while the regularization schemes encouraging Lipschitzness improve the performance, more classical regularizers like dropout or weight decay often perform on par; (3) Raising concerns: we observe a high discrepancy between the results for RL-generated and human data. We hope this study will be helpful to anyone using or designing AIL algorithms.

Acknowledgments

We thank Kamyar Ghasemipour for the discussions related to the FAIRL reward function and Lucas Beyer for the feedback on an earlier version of the manuscript.

²⁴[45] evaluated gradient penalty in the off-policy setup in the following environments: HalfCheetah, Hopper, Walker2d and Ant, as well as InvertedPendulum which we did not use due to its simplicity.

References

- [1] Stephen Hanson and Lorien Pratt. “Comparing biases for minimal network construction with back-propagation”. In: *Advances in neural information processing systems* 1 (1988), pp. 177–185.
- [2] Dean A Pomerleau. “Efficient training of artificial neural networks for autonomous navigation”. In: *Neural computation* 3.1 (1991), pp. 88–97.
- [3] Stuart Russell. “Learning agents for uncertain environments”. In: *Conference on Computational learning theory*. 1998.
- [4] Stefan Schaal. “Is imitation learning the route to humanoid robots?” In: *Trends in Cognitive Sciences* 3.6 (1999), pp. 233–242. ISSN: 1364-6613. DOI: [https://doi.org/10.1016/S1364-6613\(99\)01327-3](https://doi.org/10.1016/S1364-6613(99)01327-3). URL: <http://www.sciencedirect.com/science/article/pii/S1364661399013273>.
- [5] Andrew Y Ng, Stuart J Russell, et al. “Algorithms for inverse reinforcement learning.” In: *Icml*. Vol. 1. 2000, p. 2.
- [6] Brian D Ziebart et al. “Maximum entropy inverse reinforcement learning.” In: *Aaaai*. Vol. 8. Chicago, IL, USA. 2008, pp. 1433–1438.
- [7] Brenna D Argall et al. “A survey of robot learning from demonstration”. In: *Robotics and autonomous systems* 57.5 (2009), pp. 469–483.
- [8] Brian D Ziebart. “Modeling purposeful adaptive behavior with the principle of maximum causal entropy”. In: (2010).
- [9] Ian J Goodfellow et al. “Generative adversarial networks”. In: *arXiv preprint arXiv:1406.2661* (2014).
- [10] Nitish Srivastava et al. “Dropout: a simple way to prevent neural networks from overfitting”. In: *The journal of machine learning research* 15.1 (2014), pp. 1929–1958.
- [11] Diederik P. Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization”. In: *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*. 2015. URL: <http://arxiv.org/abs/1412.6980>.
- [12] John Schulman et al. “High-dimensional continuous control using generalized advantage estimation”. In: *arXiv preprint arXiv:1506.02438* (2015).
- [13] John Schulman et al. “Trust region policy optimization”. In: *International conference on machine learning*. 2015, pp. 1889–1897.
- [14] Greg Brockman et al. “Openai gym”. In: *arXiv preprint arXiv:1606.01540* (2016).
- [15] Jonathan Ho and Stefano Ermon. “Generative adversarial imitation learning”. In: *arXiv preprint arXiv:1606.03476* (2016).
- [16] David Silver et al. “Mastering the game of Go with deep neural networks and tree search”. In: *nature* 529.7587 (2016), p. 484.
- [17] Marc G Bellemare, Will Dabney, and Rémi Munos. “A distributional perspective on reinforcement learning”. In: *International Conference on Machine Learning*. PMLR. 2017, pp. 449–458.
- [18] Justin Fu, Katie Luo, and Sergey Levine. “Learning robust rewards with adversarial inverse reinforcement learning”. In: *arXiv preprint arXiv:1710.11248* (2017).
- [19] Ishaan Gulrajani et al. “Improved training of wasserstein gans”. In: *arXiv preprint arXiv:1704.00028* (2017).
- [20] Riashat Islam et al. “Reproducibility of benchmarked deep reinforcement learning tasks for continuous control”. In: *arXiv preprint arXiv:1708.04133* (2017).
- [21] Ivaylo Popov et al. “Data-efficient deep reinforcement learning for dexterous manipulation”. In: *arXiv preprint arXiv:1704.03073* (2017).
- [22] Aravind Rajeswaran et al. “Learning complex dexterous manipulation with deep reinforcement learning and demonstrations”. In: *arXiv preprint arXiv:1709.10087* (2017).
- [23] John Schulman et al. “Proximal policy optimization algorithms”. In: *arXiv preprint arXiv:1707.06347* (2017).
- [24] Hongyi Zhang et al. “mixup: Beyond empirical risk minimization”. In: *arXiv preprint arXiv:1710.09412* (2017).
- [25] Gabriel Barth-Maron et al. “Distributed distributional deterministic policy gradients”. In: *arXiv preprint arXiv:1804.08617* (2018).

- [26] James Bradbury et al. *JAX: composable transformations of Python+NumPy programs*. Version 0.2.5. 2018. URL: <http://github.com/google/jax>.
- [27] Scott Fujimoto, Herke Hoof, and David Meger. “Addressing function approximation error in actor-critic methods”. In: *International Conference on Machine Learning*. PMLR. 2018, pp. 1587–1596.
- [28] Tuomas Haarnoja et al. “Soft actor-critic algorithms and applications”. In: *arXiv preprint arXiv:1812.05905* (2018).
- [29] Tuomas Haarnoja et al. “Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor”. In: *International Conference on Machine Learning*. 2018, pp. 1861–1870.
- [30] Peter Henderson et al. “Deep reinforcement learning that matters”. In: *Thirty-Second AAAI Conference on Artificial Intelligence*. 2018.
- [31] Todd Hester et al. “Deep q-learning from demonstrations”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 32. 1. 2018.
- [32] Ilya Kostrikov et al. “Discriminator-actor-critic: Addressing sample inefficiency and reward bias in adversarial imitation learning”. In: *arXiv preprint arXiv:1809.02925* (2018).
- [33] Francesco Locatello et al. “Challenging common assumptions in the unsupervised learning of disentangled representations”. In: *arXiv preprint arXiv:1811.12359* (2018).
- [34] Ilya Loshchilov and Frank Hutter. “Fixing weight decay regularization in adam”. In: (2018).
- [35] Mario Lucic et al. “Are gans created equal? a large-scale study”. In: *Advances in neural information processing systems*. 2018, pp. 700–709.
- [36] Takeru Miyato et al. “Spectral normalization for generative adversarial networks”. In: *arXiv preprint arXiv:1802.05957* (2018).
- [37] George Tucker et al. “The mirage of action-dependent baselines in reinforcement learning”. In: *arXiv preprint arXiv:1802.10031* (2018).
- [38] Christopher Berner et al. “Dota 2 with large scale deep reinforcement learning”. In: *arXiv preprint arXiv:1912.06680* (2019).
- [39] Eric Langlois et al. “Benchmarking model-based reinforcement learning”. In: *arXiv preprint arXiv:1907.02057* (2019).
- [40] Tom Le Paine et al. “Making efficient use of demonstrations to solve hard exploration problems”. In: *arXiv preprint arXiv:1909.01387* (2019).
- [41] Oriol Vinyals et al. “Grandmaster level in StarCraft II using multi-agent reinforcement learning”. In: *Nature* 575.7782 (2019), pp. 350–354.
- [42] Danfei Xu and Misha Denil. “Positive-unlabeled reward learning”. In: *arXiv preprint arXiv:1911.00459* (2019).
- [43] Marcin Andrychowicz et al. “What matters in on-policy reinforcement learning? a large-scale empirical study”. In: *arXiv preprint arXiv:2006.05990* (2020).
- [44] OpenAI: Marcin Andrychowicz et al. “Learning dexterous in-hand manipulation”. In: *The International Journal of Robotics Research* 39.1 (2020), pp. 3–20.
- [45] Lionel Blondé, Pablo Strasser, and Alexandros Kalousis. “Lipschitzness Is All You Need To Tame Off-policy Generative Adversarial Imitation Learning”. In: *arXiv preprint arXiv:2006.16785* (2020).
- [46] Justin Fu et al. “D4rl: Datasets for deep data-driven reinforcement learning”. In: *arXiv preprint arXiv:2004.07219* (2020).
- [47] Seyed Kamyar Seyed Ghasemipour, Richard Zemel, and Shixiang Gu. “A divergence minimization perspective on imitation learning methods”. In: *Conference on Robot Learning*. PMLR. 2020, pp. 1259–1277.
- [48] Jonathan Heek et al. *Flax: A neural network library and ecosystem for JAX*. Version 0.3.3. 2020. URL: <http://github.com/google/flax>.
- [49] Matteo Hessel et al. *Optax: composable gradient transformation and optimisation, in JAX!* Version 0.0.1. 2020. URL: <http://github.com/deepmind/optax>.
- [50] Matt Hoffman et al. “Acme: A research framework for distributed reinforcement learning”. In: *arXiv preprint arXiv:2006.00979* (2020).
- [51] Jared Kaplan et al. “Scaling laws for neural language models”. In: *arXiv preprint arXiv:2001.08361* (2020).

- [52] Ilija Radosavovic et al. “Designing Network Design Spaces”. In: *arXiv preprint arXiv:2003.13678* (2020).
- [53] Annie Chen et al. “Batch exploration with examples for scalable robotic reinforcement learning”. In: *IEEE Robotics and Automation Letters* (2021).
- [54] Leonard Hussenot et al. *Hyperparameter Selection for Imitation Learning*. 2021. arXiv: 2105.12034 [cs.LG].

Contents

1	Introduction	1
2	Experimental design	2
3	What matters for the agent training?	4
4	What matters for the discriminator training?	6
5	Are synthetic demonstrations a good proxy for human data?	9
6	How to train efficiently?	11
7	Related work	12
8	Conclusions	12
A	Reinforcement Learning Background	17
B	Adversarial Imitation Learning Background	17
C	List of Investigated Choices	17
C.1	Reinforcement Learning algorithms	17
C.2	Imitation-specific changes to RL	18
C.3	Discriminator parameterization	19
C.4	Discriminator training	19
C.5	Discriminator regularization	19
C.6	Observation normalization	20
C.7	Combining multiple batches	20
D	Best hyperparameter values	22
E	Expert and random policy scores	23
F	Experiment wide	24
F.1	Design	24
F.2	Results	25
G	Experiment main	35
G.1	Design	35
G.2	Results	36
H	Experiment trade-offs	63
H.1	Design	63
H.2	Results	64
I	Additional experiments	68

A Reinforcement Learning Background

We consider the standard reinforcement learning formalism consisting of an agent interacting with an environment. To simplify the exposition we assume in this section that the environment is fully observable. An environment is described by a set of states \mathcal{S} , a set of actions \mathcal{A} , a distribution of initial states $p(s_0)$, a reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, transition probabilities $p(s_{t+1}|s_t, a_t)$ (t is a timestep index explained later), termination probabilities $T(s_t, a_t)$ and a discount factor $\gamma \in [0, 1]$.

A policy π is a mapping from state to a distribution over actions. Every episode starts by sampling an initial state s_0 . At every timestep t the agent produces an action based on the current state: $a_t \sim \pi(\cdot|s_t)$. In turn, the agent receives a reward $r_t = r(s_t, a_t)$ and the environment's state is updated. With probability $T(s_t, a_t)$ the episode is terminated, and otherwise the new environments state s_{t+1} is sampled from $p(\cdot|s_t, a_t)$. The discounted sum of future rewards, also referred to as the *return*, is defined as $R_t = \sum_{i=t}^{\infty} \gamma^{i-t} r_i$. The agent's goal is to find the policy π which maximizes the expected return $\mathbb{E}_{\pi}[R_0|s_0]$, where the expectation is taken over the initial state distribution, the policy, and environment transitions accordingly to the dynamics specified above.

B Adversarial Imitation Learning Background

See App. A for a very brief introduction to RL and the notation used in this section.

Drawing inspiration from Inverse Reinforcement Learning [5, 8] and Generative Adversarial Networks (GANs, [9]), adversarial imitation learning [15] aims at learning a behavior similar to that of the expert given a set of expert demonstrations \mathcal{D}_{expert} and the ability to interact with the environment.

To do so, the agent with policy π is initialized randomly and interacts with the environment. A discriminator network D is trained to distinguish between samples coming from the agent $(s_t, a_t, s_{t+1}) \sim \mathcal{D}_{\pi}$ and samples coming from the expert dataset $(s_t, a_t, s_{t+1}) \sim \mathcal{D}_{expert}$ with a cross-entropy loss. A reward function for the policy is then defined based on the discriminator prediction, e.g. $r(s, a) = -\ln(1 - D(s, a))$, where $D(s, a)$ denotes the probability of classifying the state-action pair as expert by the discriminator. The agent is then trained with an RL algorithm to maximize this reward and thus fool the discriminator. As in GANs, the training of the discriminator and that of the agent (here playing the role of the *generator*) are interleaved. Therefore, at the high level, the algorithm repeats the following steps in a loop: (1) interact with the environment using the current policy and store the experience in a replay buffer, (2) update the discriminator, (3) perform an RL update accordingly to the RL algorithm used.

C List of Investigated Choices

In this section we list all algorithmic choices which we consider in our experiments. See App. B for an introduction to adversarial imitation and the notation used in this section. For convenience, we mark each of the choices with a number (e.g., C8) and a fixed name (e.g. RL Algorithm (C8)) that can be easily used to find a description of the choice in this section.

C.1 Reinforcement Learning algorithms

In all experiments we use MLPs for the policy and critic/value networks and sample the following HPs controlling the networks architectures: policy MLP depth (C1) (the number of *hidden* layers), policy MLP width (C2), critic MLP depth (C3), critic MLP width (C4), RL activation (C5), as well as discount γ (C6) and batch size (C7). All networks are optimized with the Adam [11] optimizer.

We sample RL Algorithm (C8) from the following options:

Proximal Policy Optimization (PPO, [23]) For PPO, batch size (C7) denotes the number of experience fragments, each of consisting PPO unroll length (C9) transitions, collected in each policy update step. In each policy update step, we perform PPO number of epochs (C10) passes over the gathered data when in each pass the data is split into PPO number of minibatches (C11) minibatches. We use the PPO loss with the clipping threshold set by PPO clipping ϵ (C12)

and add an entropy loss with the coefficient specified by PPO entropy cost (C13). We also sample PPO learning rate (C14), and the GAE [12] returns mixing coefficient GAE λ (C15).

Soft Actor Critic (SAC, [29]) We use a version of SAC with a policy entropy constraint [28]. In particular, we choose SAC entropy per dimension (C16) and that set the entropy constraint so that the policy entropy is not lower than the number of action dimensions times this value. We also sweep SAC learning rate (C17) and the target network polyak averaging coefficient SAC polyak τ (C18) (the target network is updated after each minibatch).

Twin Delayed Deep Deterministic Policy Gradient (TD3, [27]) For TD3, we sweep TD3 policy learning rate (C19) and TD3 critic learning rate (C20) separately, as well as sample behavioral policy noise (C21). Following the original publication, we update the actor only using every other minibatch while the critic networks uses all minibatches. The target network is updated after every minibatch with the polyak coefficient fixed to 0.005. Following DAC [32], we clip actor gradients with magnitudes bigger than TD3 gradient clipping (C22).

Distributed Distributional Deterministic Policy Gradients (D4PG, [25]) This algorithm is similar to TD3 but uses a distributional C51-style critic [17] outputting distributions over number of atoms (C23) atoms spaced equally between -VMax (C24) and VMax (C24) as well as N-step returns (C25) returns. In contrast to the original D4PG [25], we use a single actor and do not use prioritized replay. The target network is fully updated every 100 training batches. As usual, we also sweep D4PG learning rate (C26).

Moreover, for off-policy algorithm (SAC, TD3 and D4PG) we sample replay ratio (C27) which denotes the average number of times each transition is replayed. This is achieved in the following way — if replay ratio (C27) \geq batch size (C7) than we replay replay ratio (C27) / batch size (C7) batches (each with batch size (C7) transitions) after every environment step. If batch size (C7) $>$ replay ratio (C27), we replay a single batch every batch size (C7) / replay ratio (C27) transitions. The transitions for replay are sampled uniformly from a FIFO replay buffer of size RL replay buffer size (C28) and we start training whenever we have at least 10k transition in the buffer.

For the RL algorithms which train stochastic policies (PPO and SAC) we use a Gaussian distribution followed by tanh to squash actions into the $[-1, 1]$ range.²⁵ More precisely, the policy network output is split into two parts — μ and ρ , and the action distribution used during training is $\tanh(\mathcal{N}(\mu, \text{softplus}(\rho) + 0.001))$. For policy evaluation, we choose evaluation behavior policy type (C29) from the following options:

- *stochastic*: sample from the distribution (same as behavioral policy used during training),
- *mode*: use the mode of the Gaussian instead of sampling,
- *average*: sample five action from the distribution and take the average of them.

C.2 Imitation-specific changes to RL

Reward function Let D denote the probability that a state-action pair (s, a) is classified as *expert* by the discriminator while h is the discriminator logit, i.e. $D = \sigma(h)$ where σ denotes the sigmoid function. Depending on the value of reward function (C30) we use one of the following reward functions (for completeness we write the formulas as a function of D as well as h):

- $r(s, a) = -\ln(1 - D) = \text{softplus}(h)$ (used in the original GAIL paper²⁶ [15]),
- $r(s, a) = \ln D - \ln(1 - D) = h$ (introduced in AIRL [18]).
- $r(s, a) = \ln D = -\text{softplus}(-h)$,
- $r(s, a) = -he^h$ (introduced in FAIRL [47]).

We also clip rewards with the absolute values higher than max reward magnitude (C31).

²⁵The action coordinates are scaled to $[-1, 1]$ regardless of the RL algorithm used.

²⁶The GAIL paper uses the inverse convention in which D denotes the probability as being classified as *non-expert*.

Absorbing state We optionally (if `absorbing state` (C32)=True) apply the absorbing state technique from DAC [32]. This technique encourages the agent to generate episodes of similar length to the ones of the expert. In particular, the demonstration and agent episodes are processed in the following way: for each terminal transition, we replace it with a non-terminal transition to a special absorbing state²⁷ and also add a transition from the absorbing state to itself with a zero action.

Replaying demonstrations For off-policy RL algorithms, we optionally (if `policy-to-expert replay ratio` (C33) $\neq \infty$) sample batches for RL training not only from the replay buffer, but also from the demonstrations. In particular, the ratio of policy to expert data in each minibatch is equal to `policy-to-expert replay ratio` (C33).

Initialization with behavior cloning We optionally (if `BC pretraining` (C34)=True) pre-train the policy network offline at the beginning of training using Behavior Cloning [2]. In particular, we perform 100k gradient steps with Adam on the MSE loss, using learning rate 10^{-4} and batch size 256.

C.3 Discriminator parameterization

Depending on the value of `discriminator input` (C35), the discriminator is fed single states, state-action pairs, state-state pairs or state-action-state tuples.

Our basic discriminator architecture is an MLP with `discriminator MLP depth` (C36) hidden layers, each of size `discriminator MLP width` (C37) with the activation function specified by `discriminator activation` (C38). Its output is interpreted as the logit of the probability of being classified as expert, i.e. for a state-action-state tuple (s, a, s') we have $D(s, a, s') = \sigma(f(s, a, s'))$, where D is the probability of classifying the tuple (s, a, s') as expert, σ denotes the sigmoid function, and f is a learnable function represented as an MLP.

We also consider two modifications introduced in the AIRL [18] paper. The first one (enabled if `reward shaping` (C39)=True) adds a reward shaping term where the f function is parameterized in the following way: $f(s, a, s') = g(s, a, s') + \gamma h(s') - h(s)$ where g and h are MLPs parameterized as described above, and γ is the RL discount factor.²⁸ The second modification (enabled if `subtract log-pi` (C40)=True) parameterizes the discriminator as $D(s, a, s') = \frac{\exp(f(s, a, s'))}{\exp(f(s, a, s')) + \pi(a|s)}$, where π is the current agent policy. It can be easily shown that it is equivalent to $D(s, a, s') = \sigma((f(s, a, s') - \log \pi(a|s))$ so this just shifts the logits by $\log \pi(a|s)$.

C.4 Discriminator training

All discriminator weight matrices use the `lecun_uniform` initializer from JAX [26]. The last discriminator layer initialization is additionally multiplied by `discriminator last layer init scale` (C41).

The discriminator is trained with the Adam [11] optimizer, the learning rate specified by `discriminator learning rate` (C42) and the cross-entropy loss. Each data batch contains exactly `batch size` (C7) expert transitions and `batch size` (C7) policy transitions. The policy transitions are sampled uniformly from a FIFO replay buffer of size `discriminator replay buffer size` (C43).

We perform `discriminator to RL updates ratio` (C44) discriminator gradient steps for each RL gradient step. More precisely, after each environment step, we compute the number of RL gradient steps as described in App. C.1, and perform `discriminator to RL updates ratio` (C44) that many discriminator gradient steps *before* performing the RL update.

C.5 Discriminator regularization

Depending on the value of `discriminator regularizer` (C45), we optionally apply one of the following regularizers to the discriminator:

²⁷In practice, this is done by adding a special bit to every observation which is set to zero for normal observations and one for the absorbing state. The remaining bits of the absorbing state are all zeros.

²⁸The inputs fed to g are specified by `discriminator input` (C35).

Gradient Penalty (GP, [19]) Gradient penalty is parameterized with gradient penalty k (C46) and gradient penalty λ (C47). This regularizer adds an extra term in the discriminator loss that encourages the discriminator gradient to be close to k on a convex combination of positive (expert) and negative (policy) data. In particular, for an expert data $x \sim \mathcal{D}_{expert}$ and policy data $\tilde{x} \sim \mathcal{D}_\pi$, the gradient penalty is defined as $\lambda(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - k)^2$, where \hat{x} is a convex combination of x and \tilde{x} , i.e. $\hat{x} := \epsilon x + (1 - \epsilon)\tilde{x}$ and ϵ follows a uniform distribution: $\epsilon \sim U[0, 1]$. In practice, k is usually chosen to be 0 (penalty for high gradients) or 1 (penalty for gradients with norms far from 1). Our gradient penalty implementation uses the gradient of the discriminator logit instead of the classification probability.

Spectral normalization [36] Spectral normalization guarantees that the discriminator is 1-Lipschitz: $|D(x_2) - D(x_1)| \leq \|x_2 - x_1\|$. It does so by dividing each dense layer matrix by its highest eigenvalue which can be efficiently computed with the power iteration method. See [36] for details.

Mixup [24] Mixup is parameterized with mixup α (C48) and relies on training the discriminator on a convex combination of positive (expert) and negative (policy) data. With expert data $x \sim \mathcal{D}_{expert}$ and policy data $\tilde{x} \sim \mathcal{D}_\pi$, let ϵ follow a Beta distribution: $\epsilon \sim Beta(\alpha, \alpha)$. Instead of training the discriminator on x and \tilde{x} separately, we only train it on the convex combination of them $\hat{x} := \epsilon x + (1 - \epsilon)\tilde{x}$ with the label being the convex combinations of the labels, i.e. expert with probability ϵ and non-expert with probability $1 - \epsilon$, so that the loss is $-\epsilon \ln D(\hat{x}) - (1 - \epsilon) \ln(1 - D(\hat{x}))$.

Positive Unlabeled GAIL (PUGAIL, [42]) Normally the discriminator is trained under the assumption that expert trajectories are positive examples and policy trajectories are negative examples. The PUGAIL loss assumes instead that policy trajectories are a mix of positive and negative examples.

With PUGAIL η (C49) denoting the assumed proportion of positive samples in the policy data and PUGAIL β (C50) being a clipping threshold, the discriminator is trained with the following loss:

$$\eta \hat{\mathbb{E}}_{x \sim \mathcal{D}_{expert}}[-\ln(D(x))] + \max\left(-\beta, \hat{\mathbb{E}}_{x \sim \mathcal{D}_\pi}[-\ln(1 - D(x))] - \eta \hat{\mathbb{E}}_{x \sim \mathcal{D}_{expert}}[-\ln(1 - D(x))]\right).$$

Dropout [10] We apply dropout to the hidden layers (dropout hidden rate (C51)) as well as inputs (dropout input rate (C52)). See [10] for the description of dropout.

Weight decay [1, 34] Weight decay is parameterized with a parameter controlling its strength weight decay λ (C53). Normally, weight decay is applied by adding a sum of the squares of the network parameters to the loss. However, this may interact negatively with an adaptive gradient optimizer like Adam [11] unless the optimizer is modified appropriately [34]. In our experiments, we use a version of Adam with weight decay called AdamW [34] from the Optax library [49]. See [34] for the details.

Entropy bonus Similarly to entropy bonus in RL, we also experiment with adding to the discriminator loss a term proportional to the entropy of the discriminator output treated as a Bernoulli distribution: $\lambda(D \ln D + (1 - D) \ln(1 - D))$ where entropy λ (C54) is a HP.

C.6 Observation normalization

We optionally apply input normalization (choice observation normalization (C55)) which transforms linearly the observations to all neural networks (in the RL algorithm as well as the discriminator) so that each coordinate has approximately mean zero and standard deviation equal one. This is done by subtracting from each observation μ and dividing by $\max(\rho, 0.001)$, where μ and ρ are the empirical mean and standard deviation of either all demonstrations (we call it *fixed* normalization because it does not change during training) or the empirical mean and standard deviation of all the observations encountered by the policy being trained so far (called *online* because it changes during training).

C.7 Combining multiple batches

We consider processing multiple batches at once for improved accelerator (GPU or TPU) utilization (choice number of combined batches (C56)). In particular, we sample an N -times larger batch

from a replay buffer, split it back into N smaller/proper batches on an accelerator, and process them sequentially. In order to keep the replay ratio unaffected, we decrease the frequency of updates accordingly, e.g. instead of performing one gradient update for every environment step, we perform N gradients updates every N environment steps. We apply this technique to the discriminator as well as the RL agent training.

D Best hyperparameter values

Table 2 shows the best value found for each HP in the main experiment. See App. G for the full experimental report. The sample complexity can be slightly improved by decreasing number of combined batches (C56) and increasing discriminator to RL updates ratio (C44). We used the suboptimal values from Table 2 because they give a good trade-off between sample complexity and runtime. discriminator learning rate (C42) equal 10^{-6} is better when PUGAIL, entropy or no discriminator regularizer is used, and $3 \cdot 10^{-5}$ is better otherwise. The performance of observation normalization schemes depends heavily on the environment and discriminator regularization used. For completeness, we present the best HPs for all discriminator regularizers.

Table 2: Best hyperparameter configuration.

Choice	Name	Best value
C1	policy MLP depth	2
C2	policy MLP width	256
C3	critic MLP depth	2
C4	critic MLP width	256
C5	RL activation	ReLU
C6	discount γ	0.97
C7	batch size	256
C8	RL Algorithm	SAC
C16	SAC entropy per dimension	-0.5
C17	SAC learning rate	$3 \cdot 10^{-4}$
C18	SAC polyak τ	0.01
C27	replay ratio	256
C28	RL replay buffer size	$3 \cdot 10^6$
C29	evaluation behavior policy type	mode
C30	reward function	AIRL
C31	max reward magnitude	∞
C32	absorbing state	True
C33	policy-to-expert replay ratio	∞
C34	BC pretraining	True
C35	discriminator input	(s, a)
C36	discriminator MLP depth	1
C37	discriminator MLP width	64
C38	discriminator activation	ReLU
C39	reward shaping	False
C40	subtract log-pi	False
C41	discriminator last layer init scale	1
C42	discriminator learning rate	10^{-6} or $3 \cdot 10^{-5}$
C43	discriminator replay buffer size	$3 \cdot 10^6$
C44	discriminator to RL updates ratio	1
C45	discriminator regularizer	spectral normalization
C46	gradient penalty k	0
C47	gradient penalty λ	1
C48	mixup α	1
C49	PUGAIL η	0.7
C50	PUGAIL β	∞
C51	dropout hidden rate	75%
C52	dropout input rate	50%
C53	weight decay λ	10
C54	entropy λ	0.03
C55	observation normalization	depends on the environment
C56	number of combined batches	8

E Expert and random policy scores

Table 3: Expert and random policy scores used to normalize the performance for all tasks.

Task	Random policy score	Expert score
HalfCheetah-v2	-282	8770
Hopper-v2	18	2798
Walker2d-v2	1.6	4118
Ant-v2	-59	5637
Humanoid-v2	123	9115
pen-expert-v0	94	3078
door-expert-v0	-56	2882
door-human-v0	-56	796
hammer-expert-v0	-274	12794
hammer-human-v0	-274	3071

F Experiment wide

F.1 Design

For each of the 10 tasks, we sampled 12083 choice configurations where we sampled the following choices independently and uniformly from the following ranges:

- RL Algorithm (C8): {d4pg, ppo, sac, td3}
 - For the case “RL Algorithm (C8) = sac”, we further sampled the sub-choices:
 - * SAC learning rate (C17): {0.0001, 0.0003, 0.001}
 - * SAC entropy per dimension (C16): {-2.0, -1.0, -0.5, 0.0}
 - * SAC polyak τ (C18): {0.001, 0.003, 0.01, 0.03}
 - * subtract log-pi (C40): {False, True}
 - * batch size (C7): {256.0}
 - For the case “RL Algorithm (C8) = d4pg”, we further sampled the sub-choices:
 - * D4PG learning rate (C26): {3e-05, 0.0001, 0.0003}
 - * behavioral policy noise (C21): {0.1, 0.2, 0.3, 0.5}
 - * VMax (C24): {150.0, 750.0, 1500.0}
 - * number of atoms (C23): {51.0, 101.0, 201.0, 401.0}
 - * N-step returns (C25): {1.0, 3.0, 5.0}
 - * batch size (C7): {256.0}
 - For the case “RL Algorithm (C8) = td3”, we further sampled the sub-choices:
 - * TD3 policy learning rate (C19): {0.0001, 0.0003, 0.001}
 - * TD3 critic learning rate (C20): {0.0001, 0.0003, 0.001}
 - * TD3 gradient clipping (C22): {40.0, ∞ }
 - * behavioral policy noise (C21): {0.1, 0.2, 0.3, 0.5}
 - * batch size (C7): {256.0}
 - For the case “RL Algorithm (C8) = ppo”, we further sampled the sub-choices:
 - * PPO learning rate (C14): {3e-05, 0.0001, 0.0003}
 - * PPO number of epochs (C10): {2.0, 5.0, 10.0, 20.0}
 - * PPO entropy cost (C13): {0.0, 0.001, 0.003, 0.01, 0.03, 0.1}
 - * PPO number of minibatches (C11): {8.0, 16.0, 32.0, 64.0}
 - * PPO unroll length (C9): {4.0, 8.0, 16.0, 32.0}
 - * PPO clipping ϵ (C12): {0.1, 0.2, 0.3}
 - * GAE λ (C15): {0.8, 0.9, 0.95, 0.99}
 - * subtract log-pi (C40): {False, True}
 - * batch size (C7): {64.0, 128.0, 256.0}
- RL replay buffer size (C28): {300000.0, 1000000.0, 3000000.0}
- policy MLP depth (C1): {1, 2, 3}
- policy MLP width (C2): {64, 128, 256, 512}
- critic MLP depth (C3): {1, 2, 3}
- critic MLP width (C4): {64, 128, 256, 512}
- RL activation (C5): {relu, tanh}
- discount γ (C6): {0.9, 0.97, 0.99, 0.997}
- BC pretraining (C34): {False, True}
- absorbing state (C32): {False, True}
- discriminator replay buffer size (C43): {300000, 1000000, 3000000}
- reward shaping (C39): {False, True}
- discriminator input (C35): {s, sa, sas, ss}
- discriminator MLP depth (C36): {1, 2, 3}
- discriminator MLP width (C37): {16, 32, 64, 128, 256, 512}

- **discriminator activation** (C38): {elu, leaky_relu, relu, sigmoid, swish, tanh}
- **discriminator last layer init scale** (C41): {0.001, 1.0}
- **discriminator regularizer** (C45): {GP, Mixup, No regularizer, PUGAIL, dropout, entropy, spectral norm, weight decay}
 - For the case “**discriminator regularizer** (C45) = GP”, we further sampled the sub-choices:
 - * **gradient penalty** λ (C47): {0.1, 1.0, 10.0}
 - * **gradient penalty** k (C46): {0.0, 1.0}
 - For the case “**discriminator regularizer** (C45) = Mixup”, we further sampled the sub-choices:
 - * **mixup** α (C48): {0.1, 0.4, 1.0}
 - For the case “**discriminator regularizer** (C45) = PUGAIL”, we further sampled the sub-choices:
 - * **PUGAIL** η (C49): {0.25, 0.5, 0.7}
 - * **PUGAIL** β (C50): {0.0, 0.7, ∞ }
 - For the case “**discriminator regularizer** (C45) = entropy”, we further sampled the sub-choices:
 - * **entropy** λ (C54): {0.0003, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3}
 - For the case “**discriminator regularizer** (C45) = weight decay”, we further sampled the sub-choices:
 - * **weight decay** λ (C53): {0.3, 1.0, 3.0, 10.0, 30.0}
 - For the case “**discriminator regularizer** (C45) = dropout”, we further sampled the sub-choices:
 - * **dropout input rate** (C52): {0.0, 0.25, 0.5, 0.75}
 - * **dropout hidden rate** (C51): {0.25, 0.5, 0.75}
- **observation normalization** (C55): {fixed, none}
- **evaluation behavior policy type** (C29): {average, mode, stochastic}
- **discriminator learning rate** (C42): {1e-06, 3e-06, 1e-05, 3e-05, 0.0001, 0.0003}
- **max reward magnitude** (C31): {0.5, 1.0, 2.0, 5.0, 10.0, 50.0, ∞ }
- **reward function** (C30): {-ln(1-D), AIRL, FAIRL, ln(D)}
- **replay ratio** (C27): {256}
- **discriminator to RL updates ratio** (C44): {1}
- **number of combined batches** (C56): {8}

F.2 Results

For each of the sampled choice configurations we compute the performance metric as described in Section 2. We report aggregate statistics of the experiment in Tables 4–7 as well as training curves in Figure 9. We further provide per-choice analyses in Figures 10–23.

Table 4: Quantiles of the *final* agent performance across HP configurations for OpenAI Gym tasks.

	Ant	HalfCheetah	Hopper	Humanoid	Walker2d
90%	0.18	0.80	0.99	0.06	0.56
95%	0.56	0.98	1.15	0.30	0.85
99%	0.92	1.10	1.20	0.79	0.99
Max	1.10	1.39	1.32	1.02	1.06

Table 5: Quantiles of the *final* agent performance across HP configurations for Adroit tasks.

	Door expert	Door human	Hammer expert	Hammer human	Pen expert
90%	0.12	0.07	0.16	0.12	0.28
95%	0.42	0.28	0.67	0.47	0.46
99%	0.90	1.20	1.26	2.03	0.77
Max	1.11	2.82	1.42	5.39	1.12

Table 6: Quantiles of the *average* agent performance during training across HP configurations for OpenAI Gym tasks.

	Ant	HalfCheetah	Hopper	Humanoid	Walker2d
90%	0.13	0.54	0.62	0.05	0.31
95%	0.31	0.66	0.80	0.20	0.49
99%	0.62	0.85	0.98	0.49	0.71
Max	0.94	0.99	1.08	0.84	0.92

Table 7: Quantiles of the *average* agent performance during training across HP configurations for Adroit tasks.

	Door expert	Door human	Hammer expert	Hammer human	Pen expert
90%	0.11	0.11	0.11	0.15	0.21
95%	0.27	0.24	0.34	0.33	0.35
99%	0.55	0.61	0.78	0.85	0.59
Max	0.87	1.65	1.01	1.97	0.84

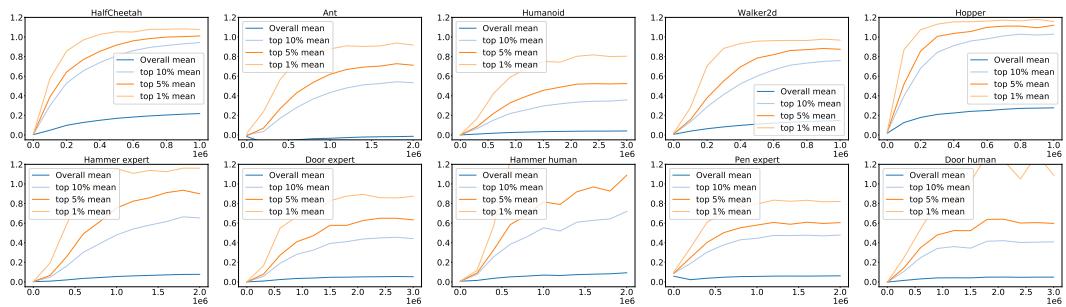


Figure 9: Training curves.

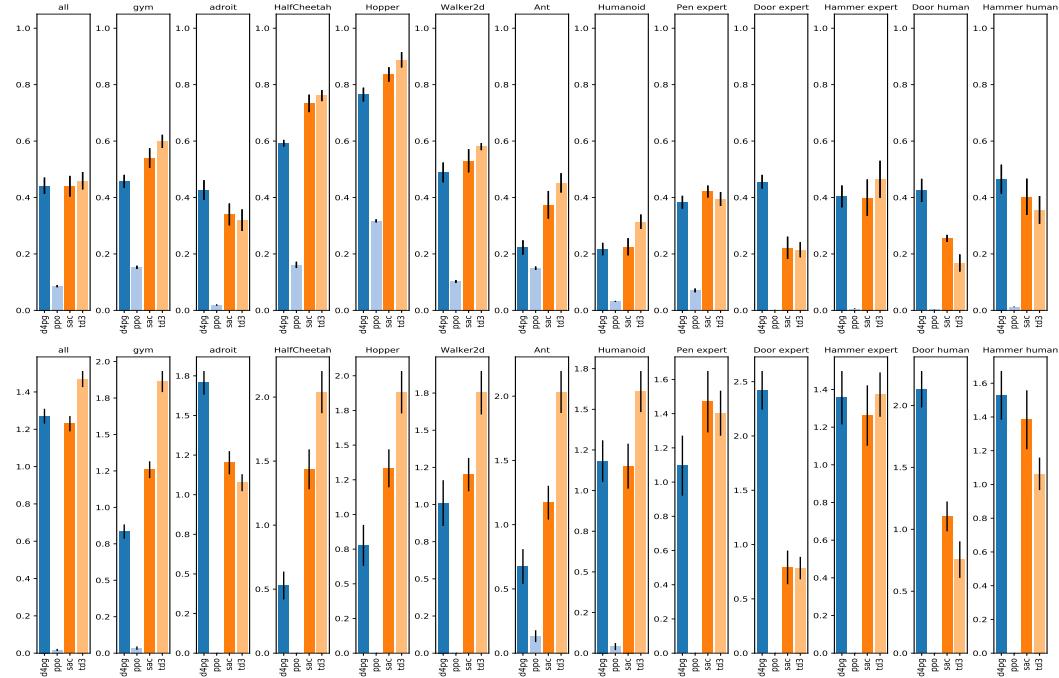


Figure 10: Analysis of choice RL Algorithm (C8): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

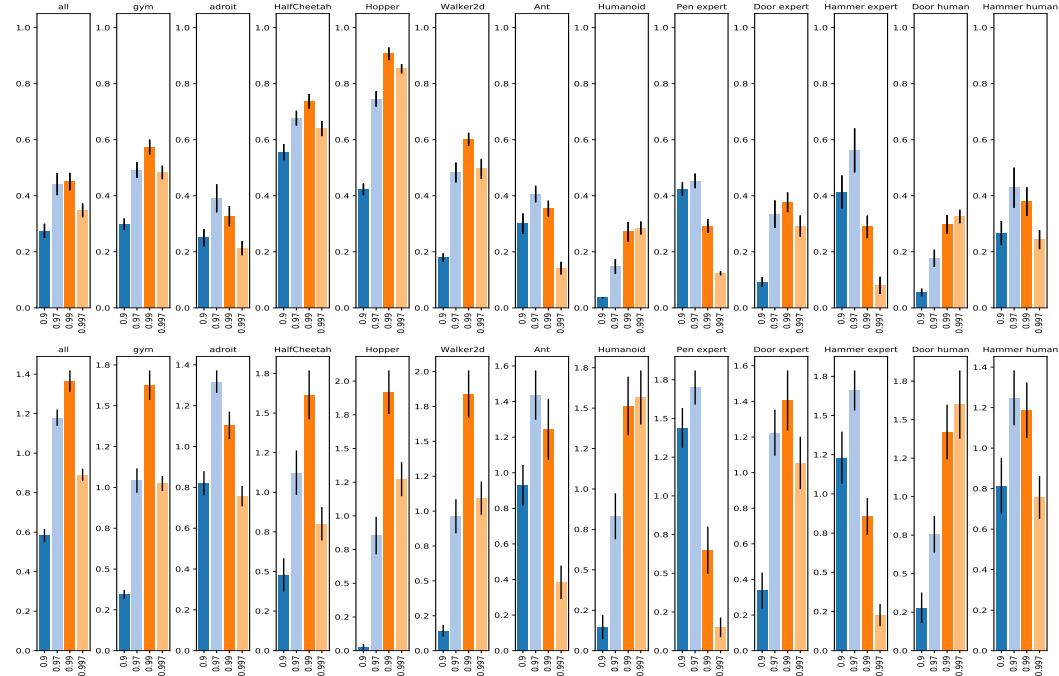


Figure 11: Analysis of choice discount γ (C6): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

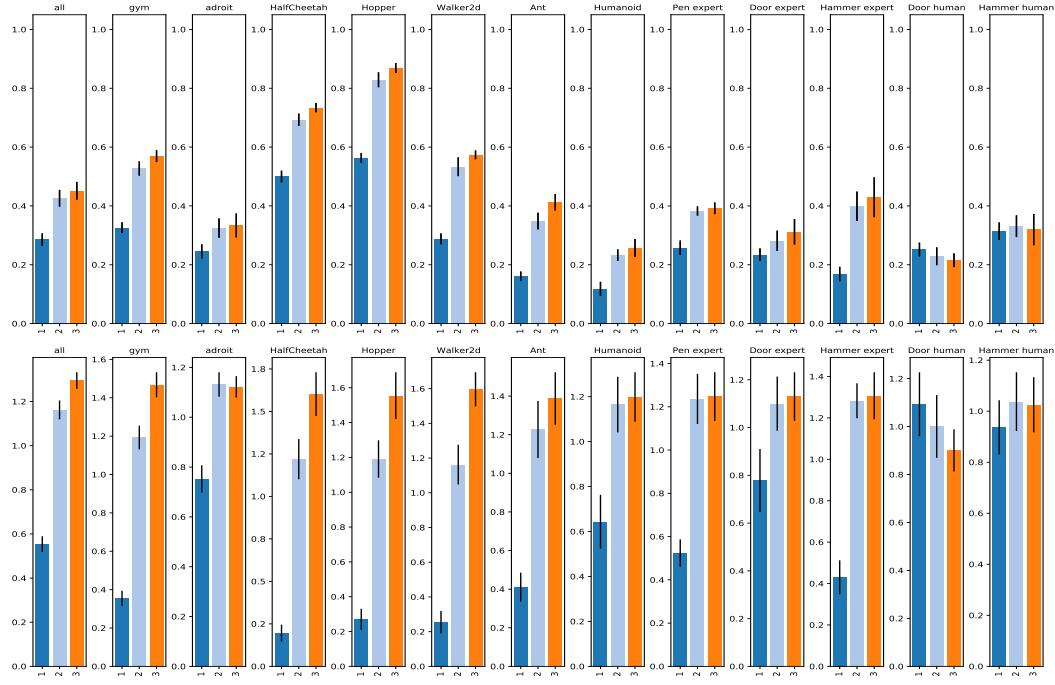


Figure 12: Analysis of choice critic MLP depth (C3): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

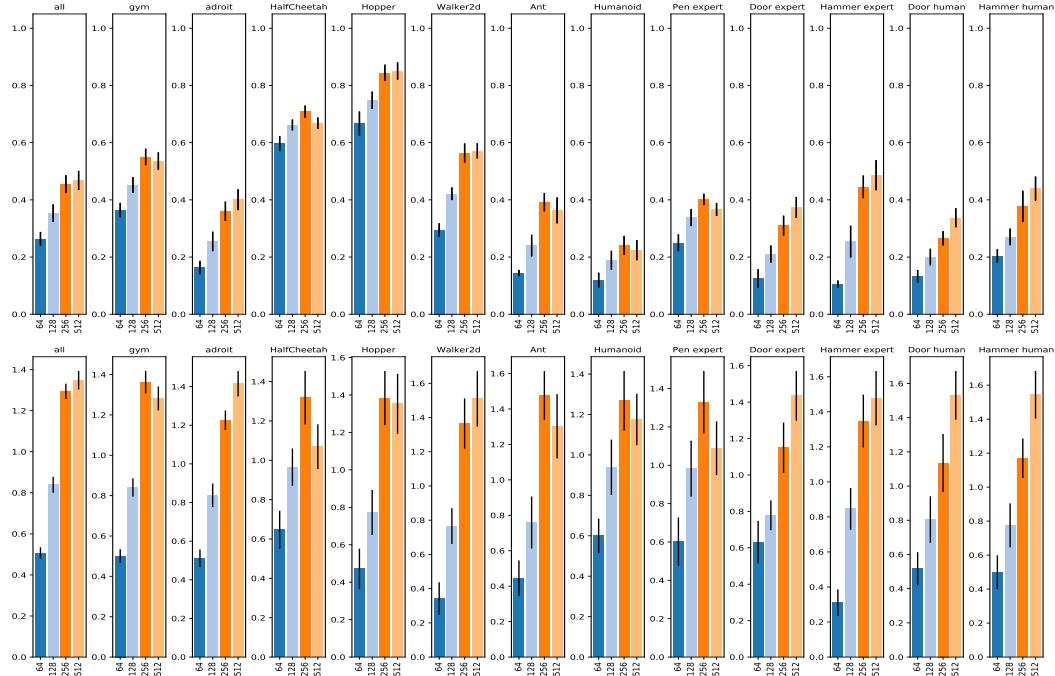


Figure 13: Analysis of choice critic MLP width (C4): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

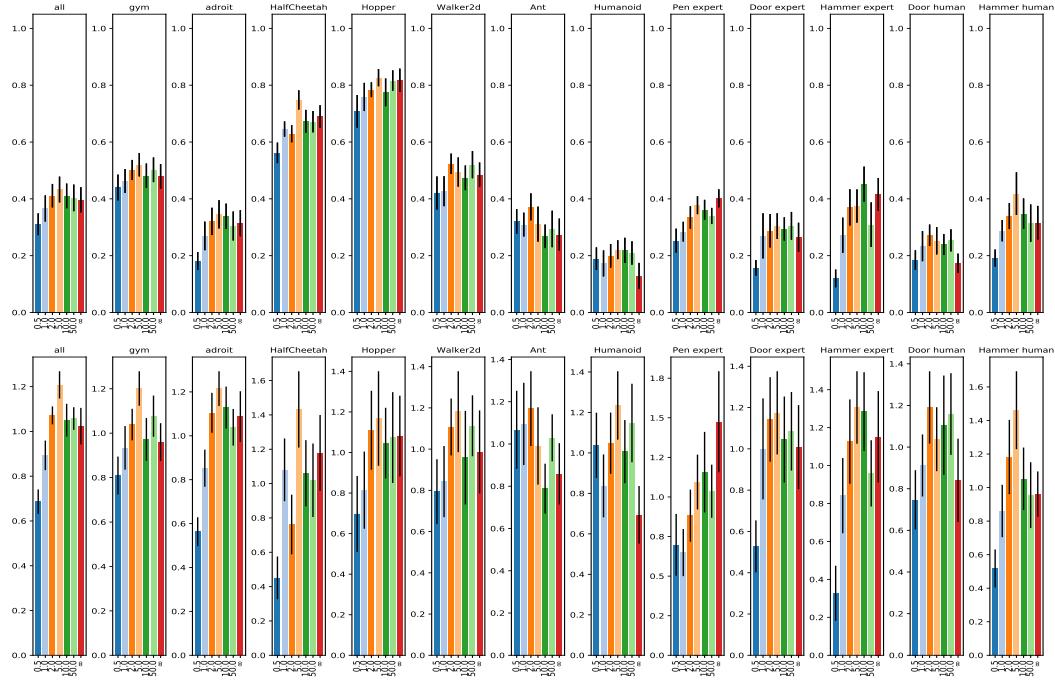


Figure 14: Analysis of choice max reward magnitude (C31): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

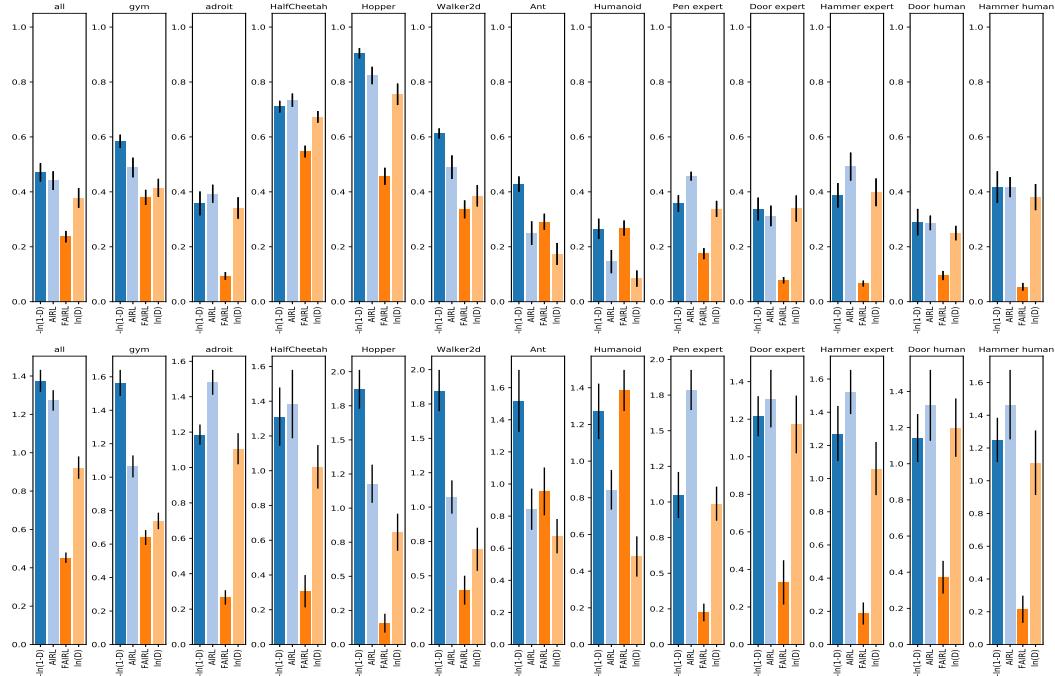


Figure 15: Analysis of choice reward function (C30): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

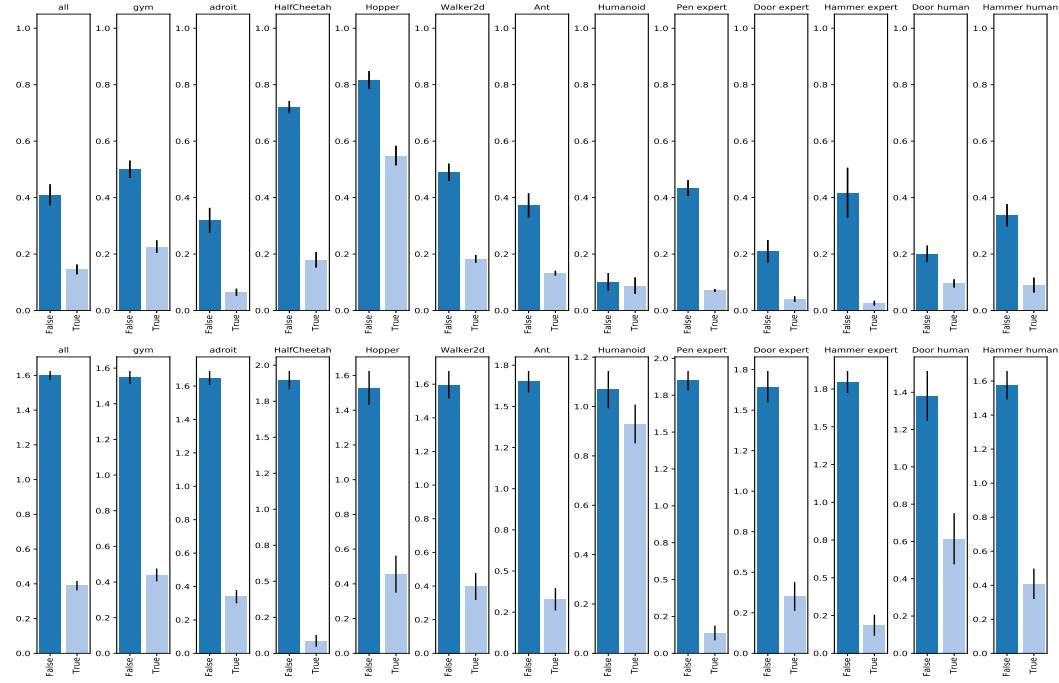


Figure 16: Analysis of choice subtract log-pi (C40): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

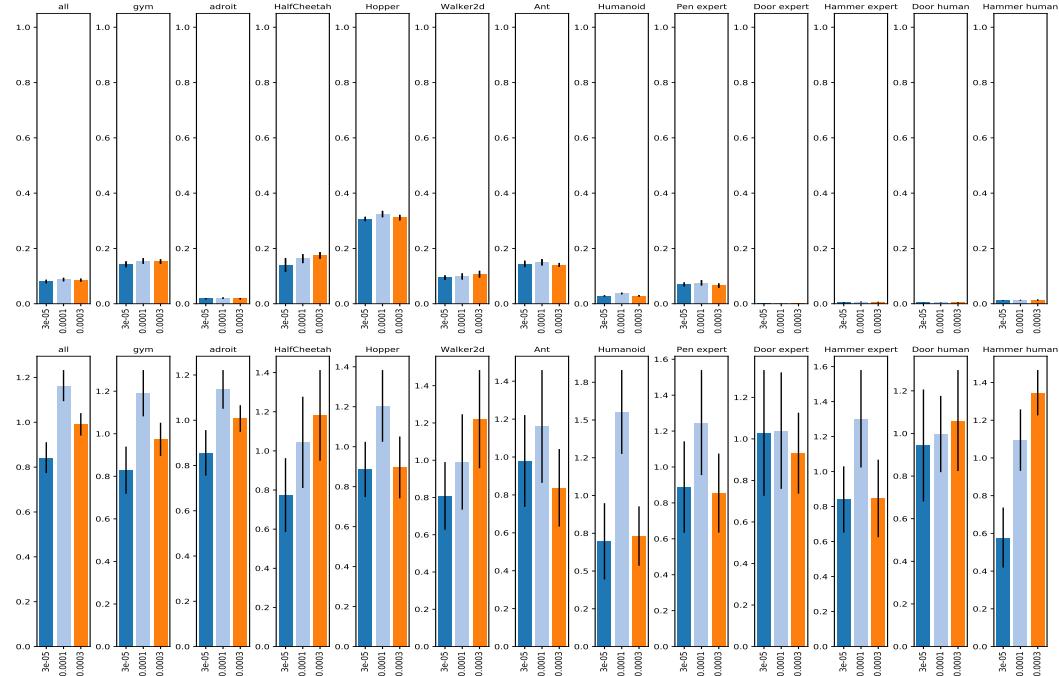


Figure 17: Analysis of choice PPO learning rate (C14): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

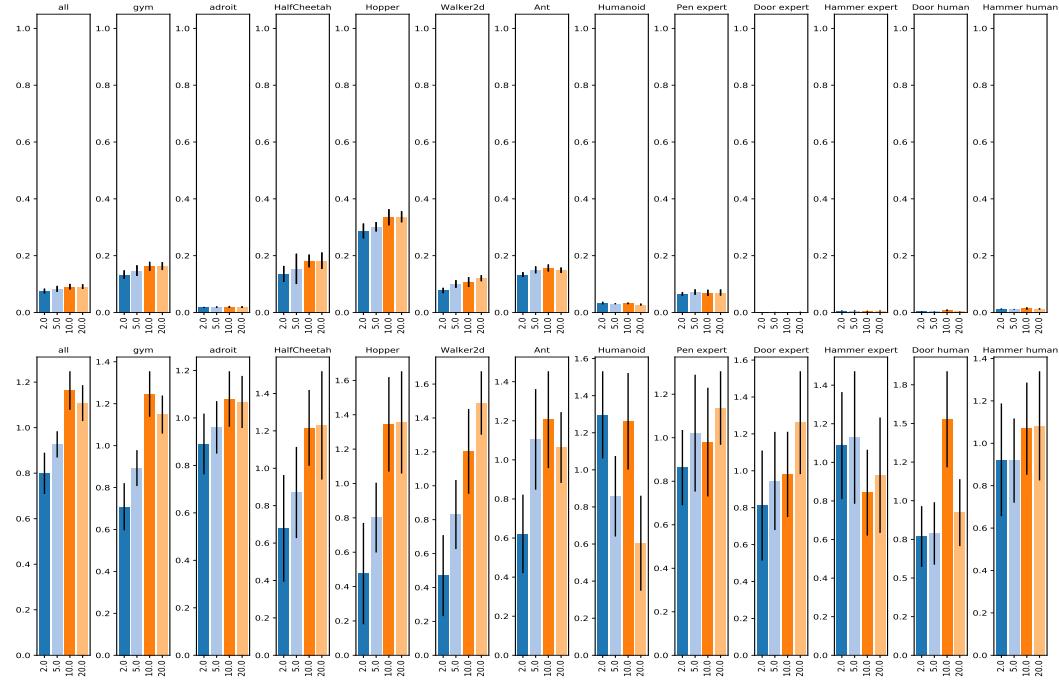


Figure 18: Analysis of choice PPO number of epochs (C10): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

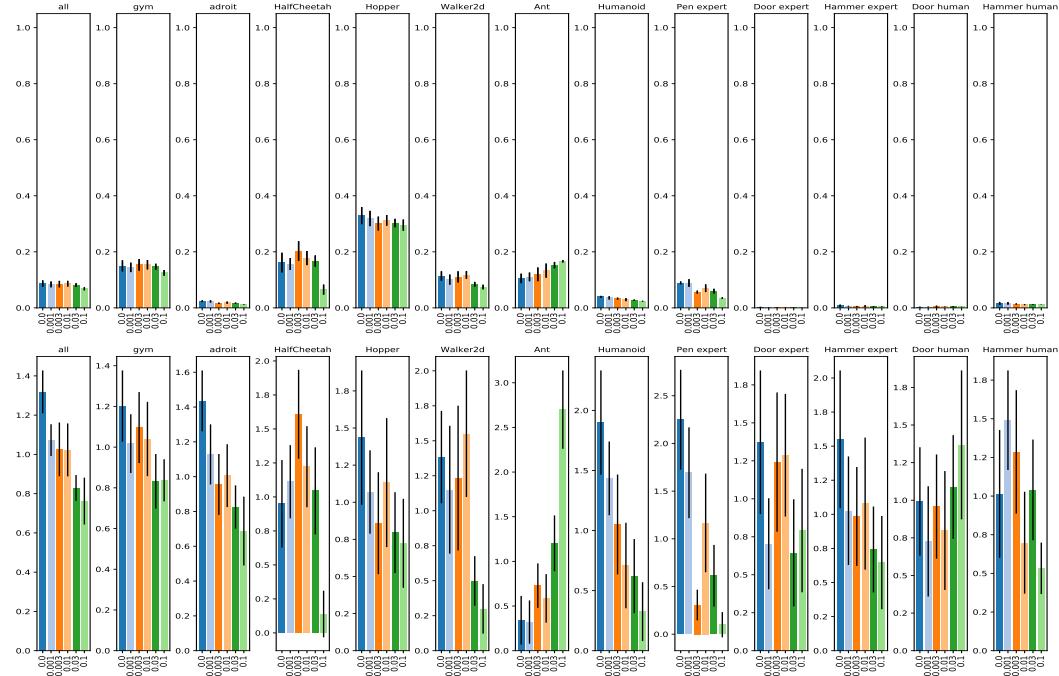


Figure 19: Analysis of choice PPO entropy cost (C13): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

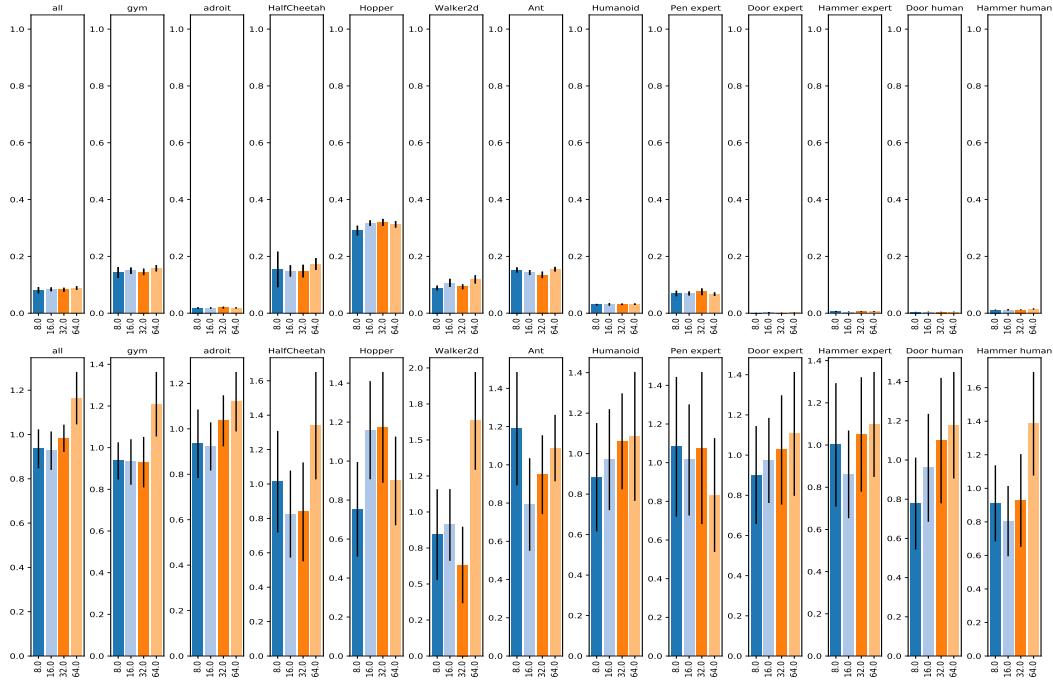


Figure 20: Analysis of choice PPO number of minibatches (C11): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

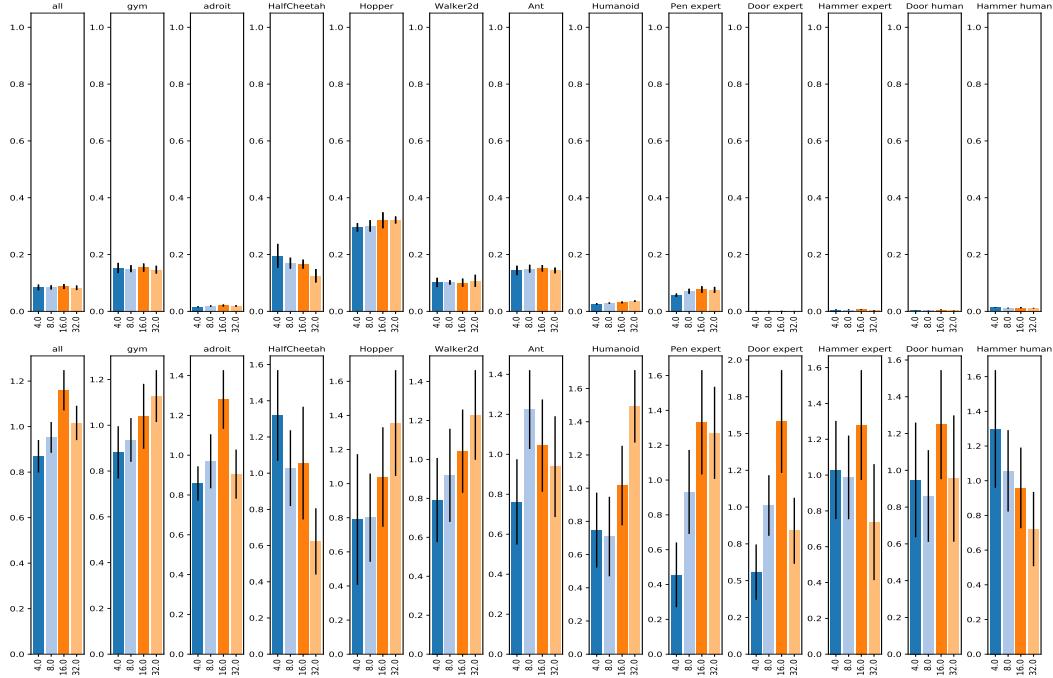


Figure 21: Analysis of choice PPO unroll length (C9): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

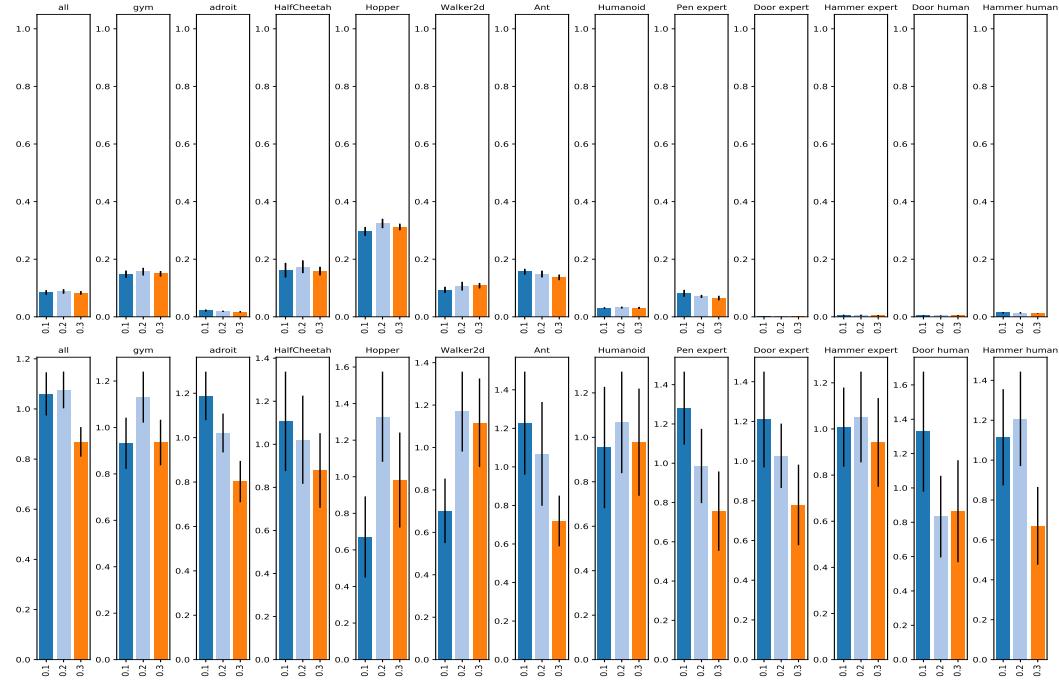


Figure 22: Analysis of choice PPO clipping ϵ (C12): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

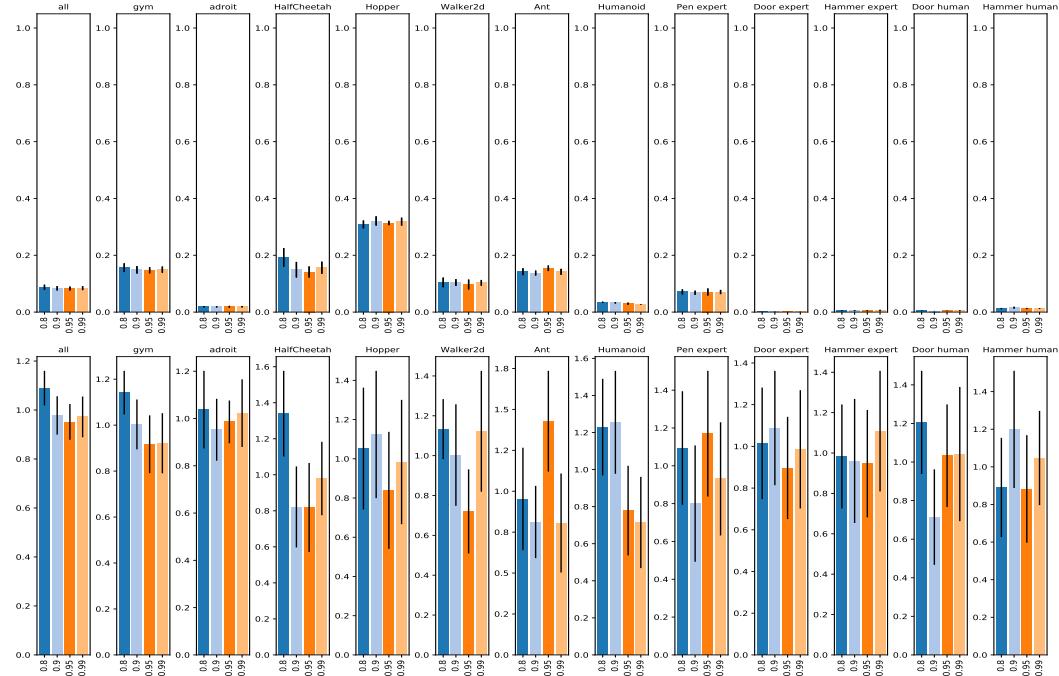


Figure 23: Analysis of choice GAE λ (C15): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

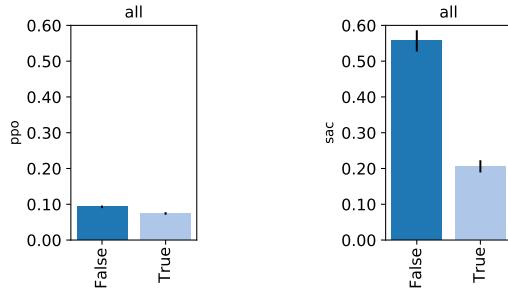


Figure 24: 95th percentile of performance scores conditioned on RL Algorithm (C8)(subplots) and subtract log-pi (C40)(bars).

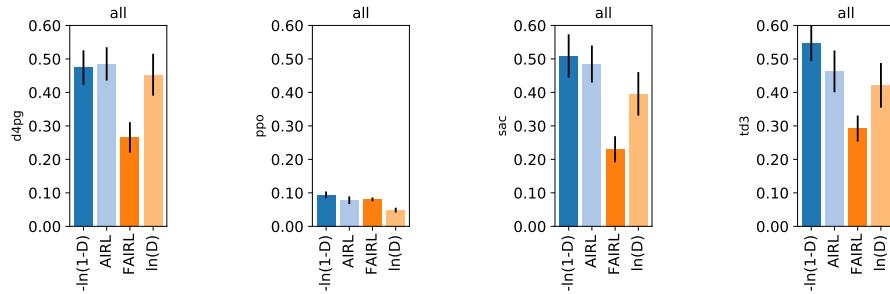


Figure 25: 95th percentile of performance scores conditioned on RL Algorithm (C8)(subplots) and reward function (C30)(bars).

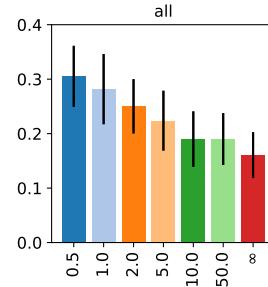


Figure 26: 95th percentile of performance scores conditioned on max reward magnitude (C31) and reward function (C30)=FAIRL.

G Experiment main

G.1 Design

For each of the 10 tasks, we sampled 25334 choice configurations where we sampled the following choices independently and uniformly from the following ranges:

- RL Algorithm (C8): {d4pg, sac, td3}
 - For the case “RL Algorithm (C8) = sac”, we further sampled the sub-choices:
 - * SAC learning rate (C17): {0.0001, 0.0003, 0.001}
 - * SAC entropy per dimension (C16): {-2.0, -1.0, -0.5, 0.0}
 - * SAC polyak τ (C18): {0.001, 0.003, 0.01, 0.03}
 - For the case “RL Algorithm (C8) = d4pg”, we further sampled the sub-choices:
 - * D4PG learning rate (C26): {3e-05, 0.0001, 0.0003}
 - * behavioral policy noise (C21): {0.1, 0.2, 0.3, 0.5}
 - * VMax (C24): {150.0, 750.0, 1500.0}
 - * number of atoms (C23): {51.0, 101.0, 201.0, 401.0}
 - * N-step returns (C25): {1.0, 3.0, 5.0}
 - For the case “RL Algorithm (C8) = td3”, we further sampled the sub-choices:
 - * TD3 policy learning rate (C19): {0.0001, 0.0003, 0.001}
 - * TD3 critic learning rate (C20): {0.0001, 0.0003, 0.001}
 - * TD3 gradient clipping (C22): {40.0, ∞ }
 - * behavioral policy noise (C21): {0.1, 0.2, 0.3, 0.5}
- RL replay buffer size (C28): {300000, 1000000, 3000000}
- policy MLP depth (C1): {1, 2, 3}
- policy MLP width (C2): {64, 128, 256, 512}
- critic MLP depth (C3): {2, 3}
- critic MLP width (C4): {256, 512}
- RL activation (C5): {relu, tanh}
- discount γ (C6): {0.97, 0.99}
- BC pretraining (C34): {False, True}
- absorbing state (C32): {False, True}
- discriminator replay buffer size (C43): {300000, 1000000, 3000000}
- reward shaping (C39): {False, True}
- discriminator input (C35): {s, sa, sas, ss}
- discriminator MLP depth (C36): {1, 2, 3}
- discriminator MLP width (C37): {16, 32, 64, 128, 256, 512}
- discriminator activation (C38): {elu, leaky_relu, relu, sigmoid, swish, tanh}
- discriminator last layer init scale (C41): {0.001, 1.0}
- discriminator regularizer (C45): {GP, Mixup, No regularizer, PUGAIL, dropout, entropy, spectral norm, weight decay}
 - For the case “discriminator regularizer (C45) = GP”, we further sampled the sub-choices:
 - * gradient penalty λ (C47): {0.1, 1.0, 10.0}
 - * gradient penalty k (C46): {0.0, 1.0}
 - For the case “discriminator regularizer (C45) = Mixup”, we further sampled the sub-choices:
 - * mixup α (C48): {0.1, 0.4, 1.0}
 - For the case “discriminator regularizer (C45) = PUGAIL”, we further sampled the sub-choices:

- * PUGAIL η (C49): {0.25, 0.5, 0.7}
- * PUGAIL β (C50): {0.0, 0.7, ∞ }
- For the case “discriminator regularizer (C45) = entropy”, we further sampled the sub-choices:
 - * entropy λ (C54): {0.0003, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3}
- For the case “discriminator regularizer (C45) = weight decay”, we further sampled the sub-choices:
 - * weight decay λ (C53): {0.3, 1.0, 3.0, 10.0, 30.0}
- For the case “discriminator regularizer (C45) = dropout”, we further sampled the sub-choices:
 - * dropout input rate (C52): {0.0, 0.25, 0.5, 0.75}
 - * dropout hidden rate (C51): {0.25, 0.5, 0.75}
- observation normalization (C55): {fixed, none, online}
- evaluation behavior policy type (C29): {average, mode, stochastic}
- discriminator learning rate (C42): {1e-06, 3e-06, 1e-05, 3e-05, 0.0001, 0.0003}
- reward function (C30): {-ln(1-D), AIRL, ln(D)}
- batch size (C7): {256}
- replay ratio (C27): {256}
- discriminator to RL updates ratio (C44): {1}
- number of combined batches (C56): {8}

G.2 Results

For each of the sampled choice configurations we compute the performance metric as described in Section 2. We report aggregate statistics of the experiment in Tables 8–11 as well as training curves in Figure 27. We further provide per-choice analyses in Figures 40–74.

Table 8: Quantiles of the *final* agent performance across HP configurations for OpenAI Gym tasks.

	Ant	HalfCheetah	Hopper	Humanoid	Walker2d
90%	0.90	1.07	1.18	0.51	0.99
95%	0.99	1.11	1.20	0.87	1.01
99%	1.07	1.17	1.23	1.01	1.04
Max	1.18	1.37	1.34	1.06	1.21

Table 9: Quantiles of the *final* agent performance across HP configurations for Adroit tasks.

	Door expert	Door human	Hammer expert	Hammer human	Pen expert
90%	0.72	0.25	1.08	0.46	0.74
95%	0.91	0.83	1.26	1.15	0.89
99%	1.04	2.29	1.37	3.04	1.11
Max	1.16	3.73	1.45	5.55	1.44

Table 10: Quantiles of the *average* agent performance during training across HP configurations for OpenAI Gym tasks.

	Ant	HalfCheetah	Hopper	Humanoid	Walker2d
90%	0.61	0.82	0.93	0.29	0.70
95%	0.72	0.87	0.98	0.53	0.76
99%	0.85	0.94	1.06	0.79	0.84
Max	0.96	1.05	1.10	0.92	0.92

Table 11: Quantiles of the *average* agent performance during training across HP configurations for Adroit tasks.

	Door expert	Door human	Hammer expert	Hammer human	Pen expert
90%	0.42	0.30	0.59	0.42	0.56
95%	0.57	0.56	0.77	0.70	0.66
99%	0.74	1.04	0.96	1.23	0.84
Max	0.92	2.08	1.18	3.42	1.09

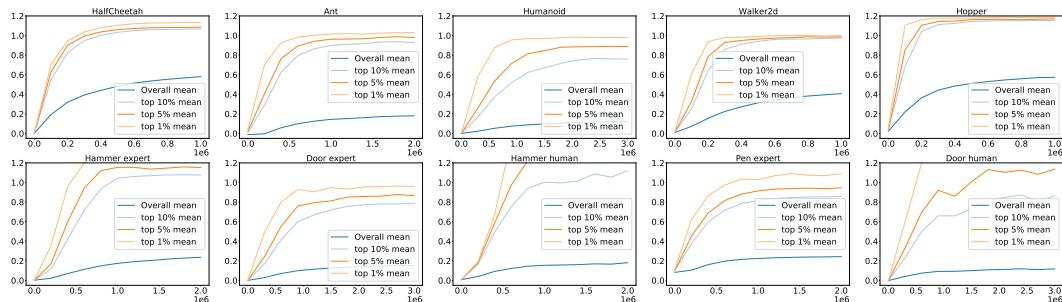


Figure 27: Training curves.

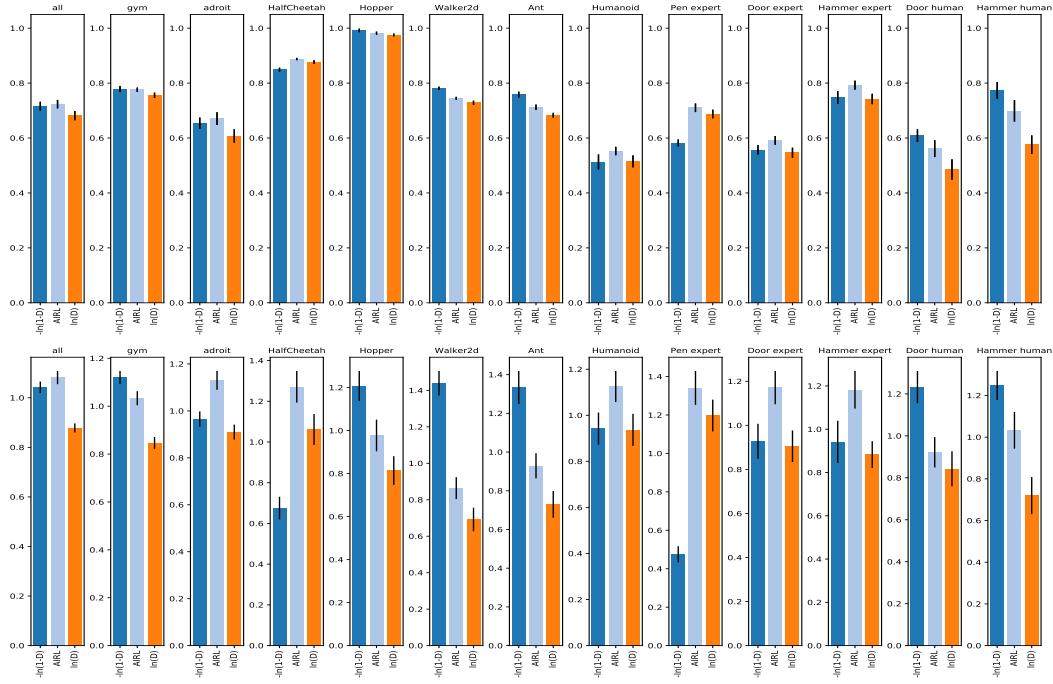


Figure 28: Analysis of choice reward function (C30): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

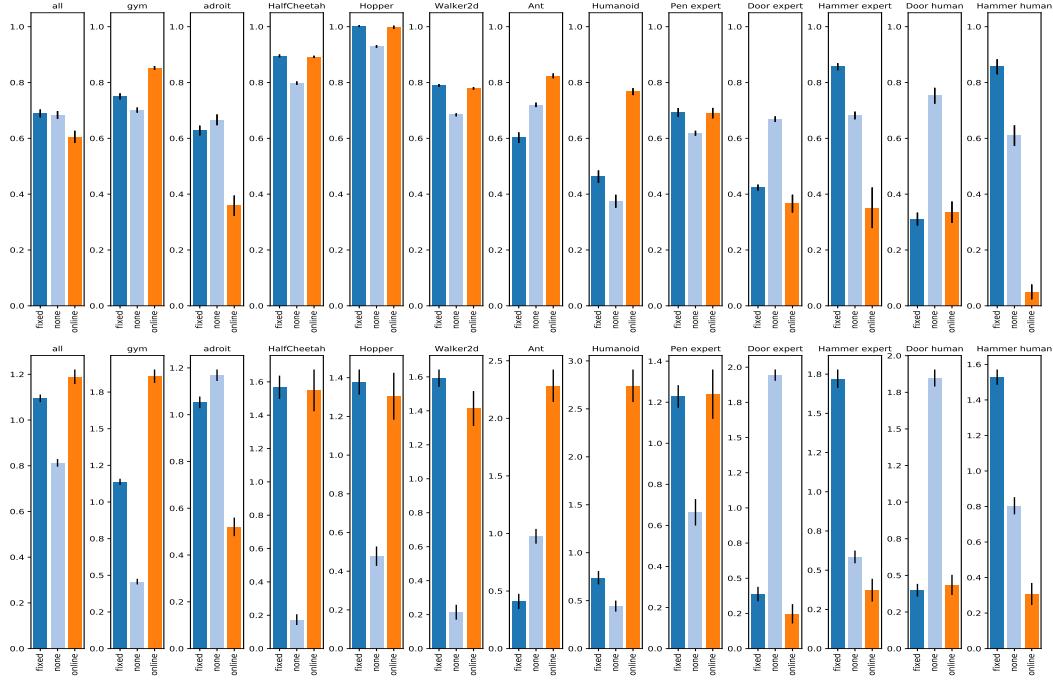


Figure 29: Analysis of choice observation normalization (C55): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

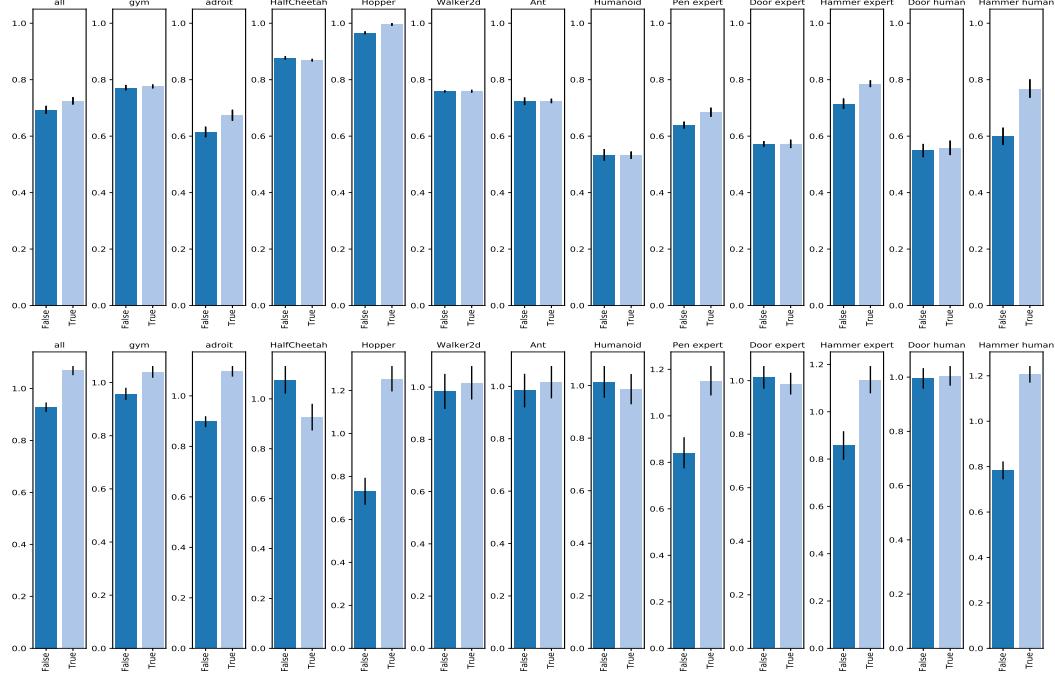


Figure 30: Analysis of choice BC pretraining (C34): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

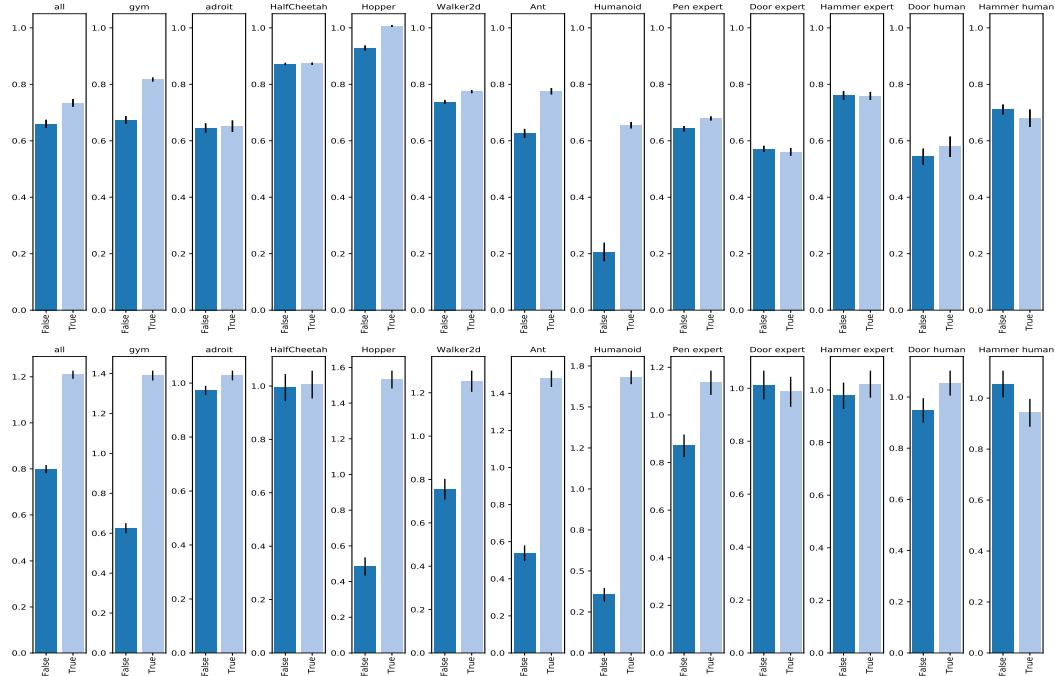


Figure 31: Analysis of choice absorbing state (C32): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

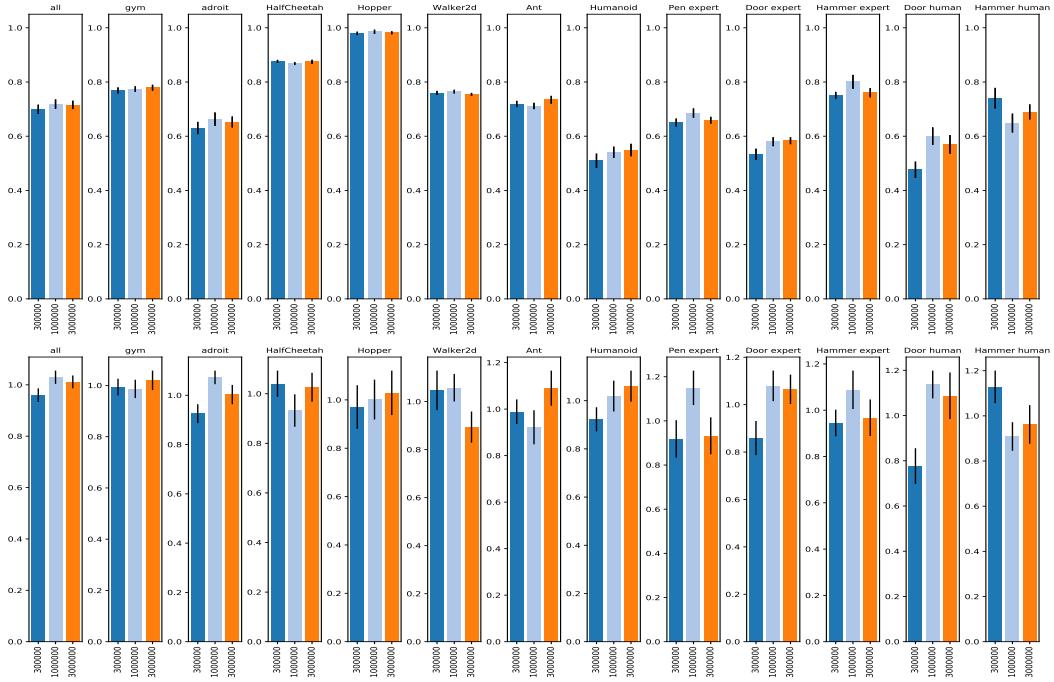


Figure 32: Analysis of choice RL replay buffer size (C28): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

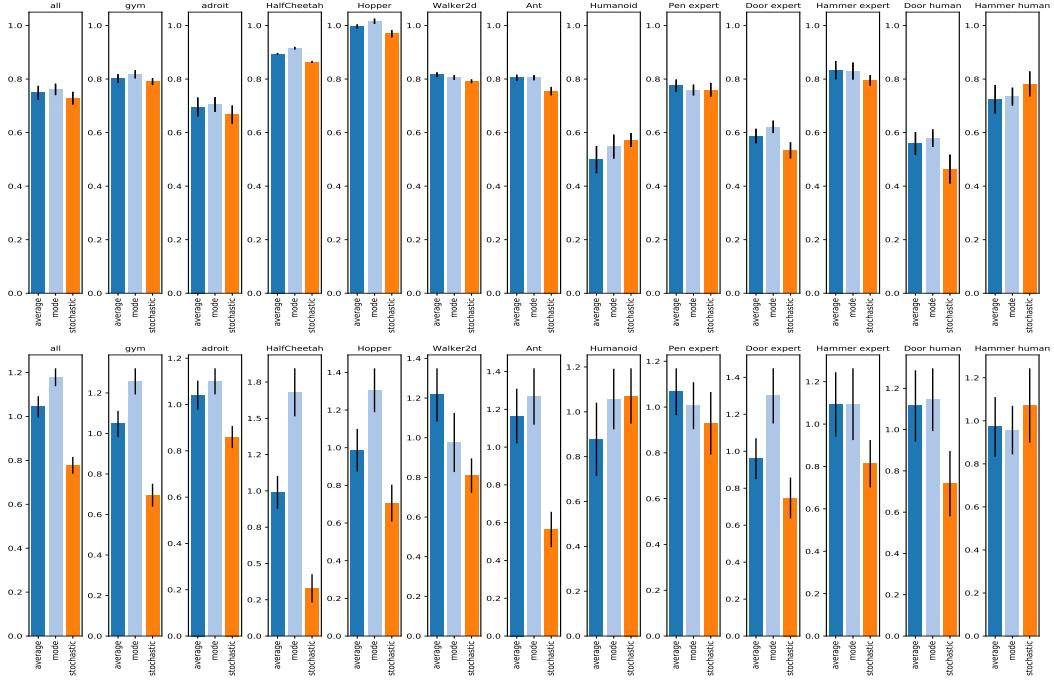


Figure 33: Analysis of choice evaluation behavior policy type (C29): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

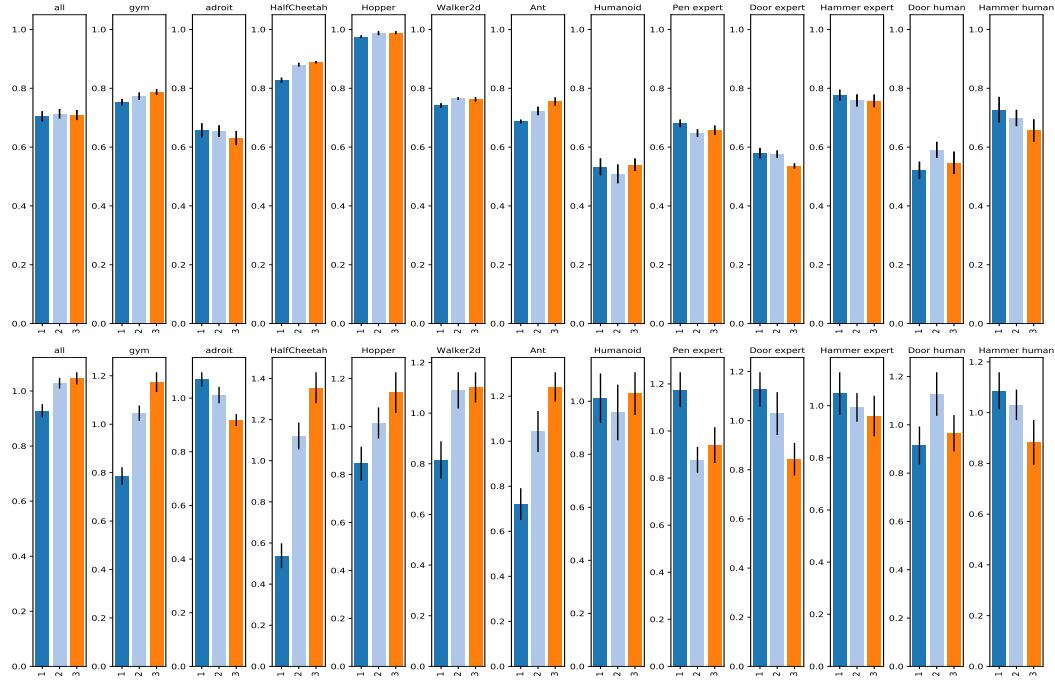


Figure 34: Analysis of choice policy MLP depth (C1): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

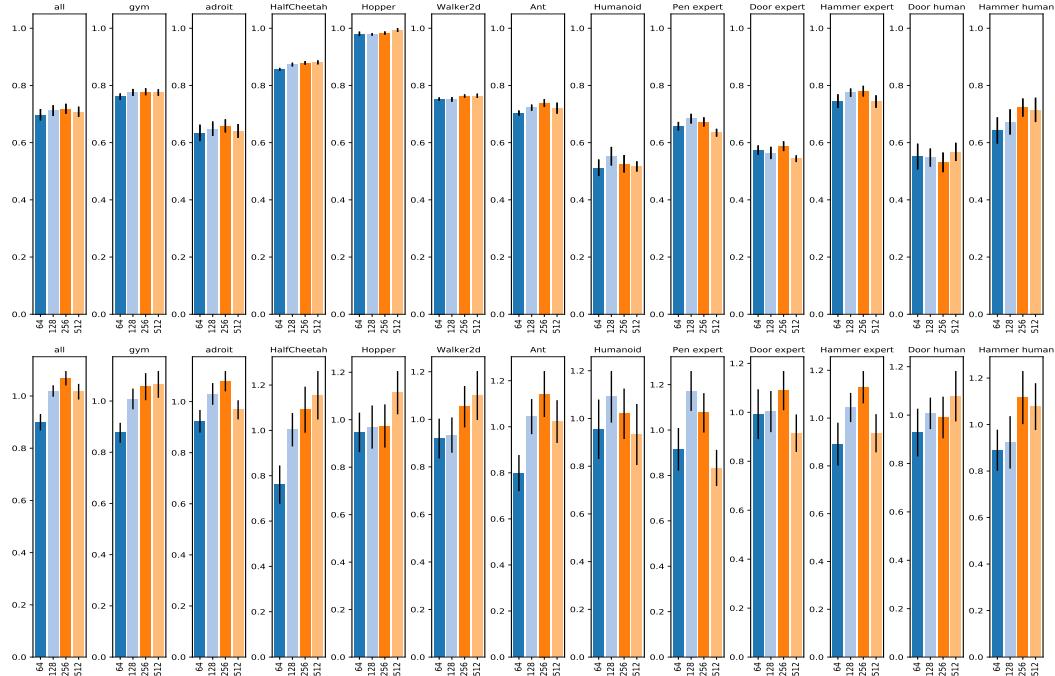


Figure 35: Analysis of choice policy MLP width (C2): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

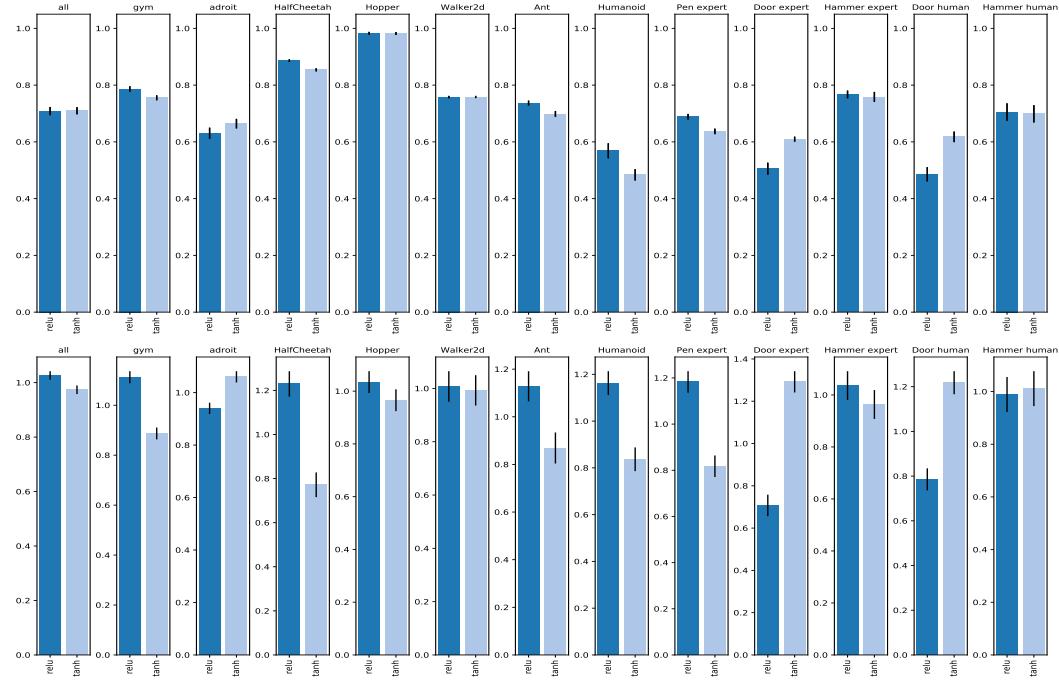


Figure 36: Analysis of choice RL activation (C5): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

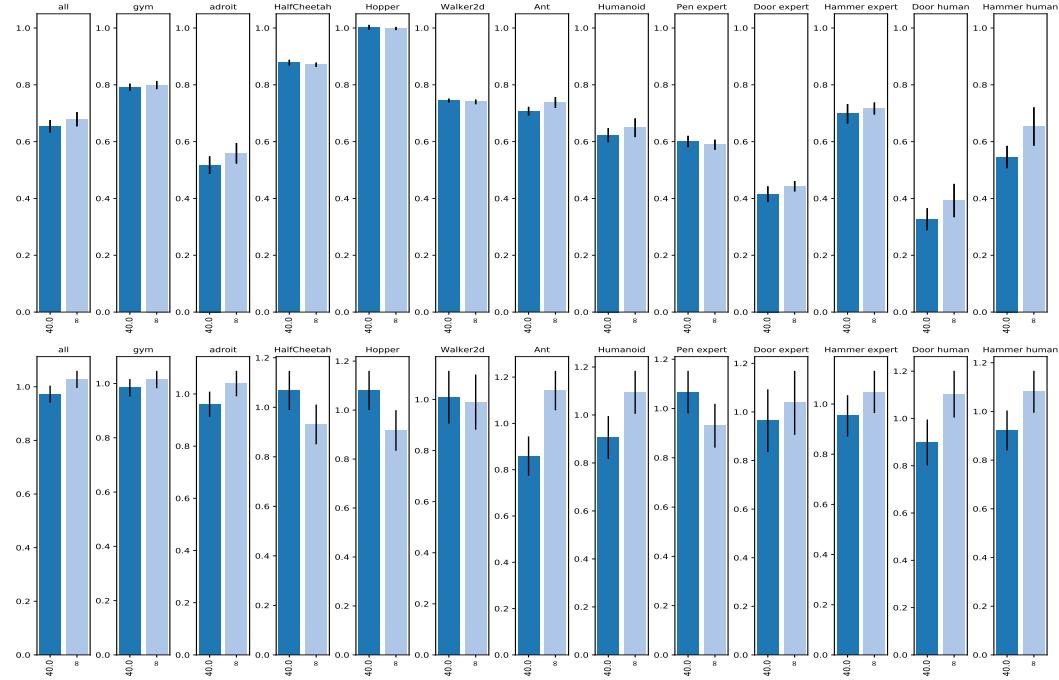


Figure 37: Analysis of choice TD3 gradient clipping (C22): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

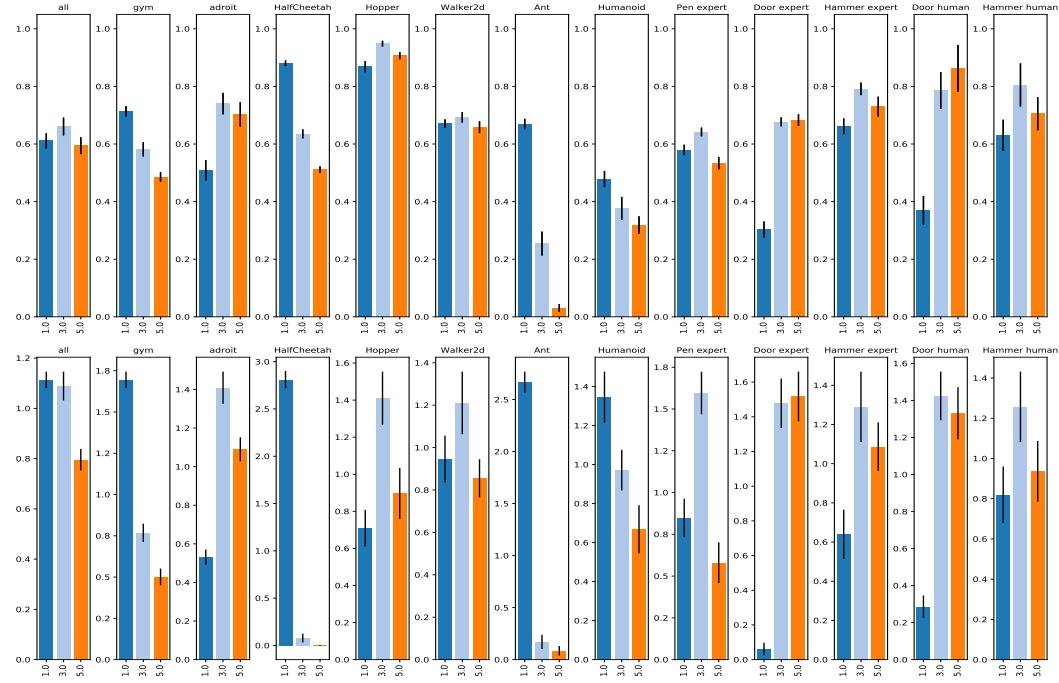


Figure 38: Analysis of choice N-step returns (C25): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

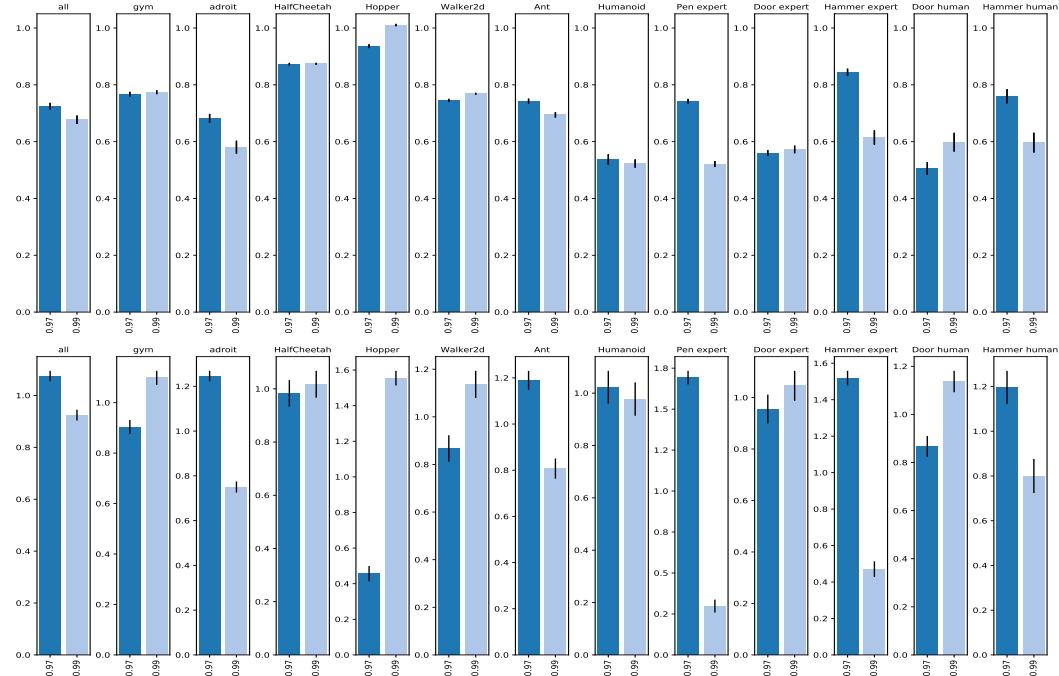


Figure 39: Analysis of choice discount γ (C6): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

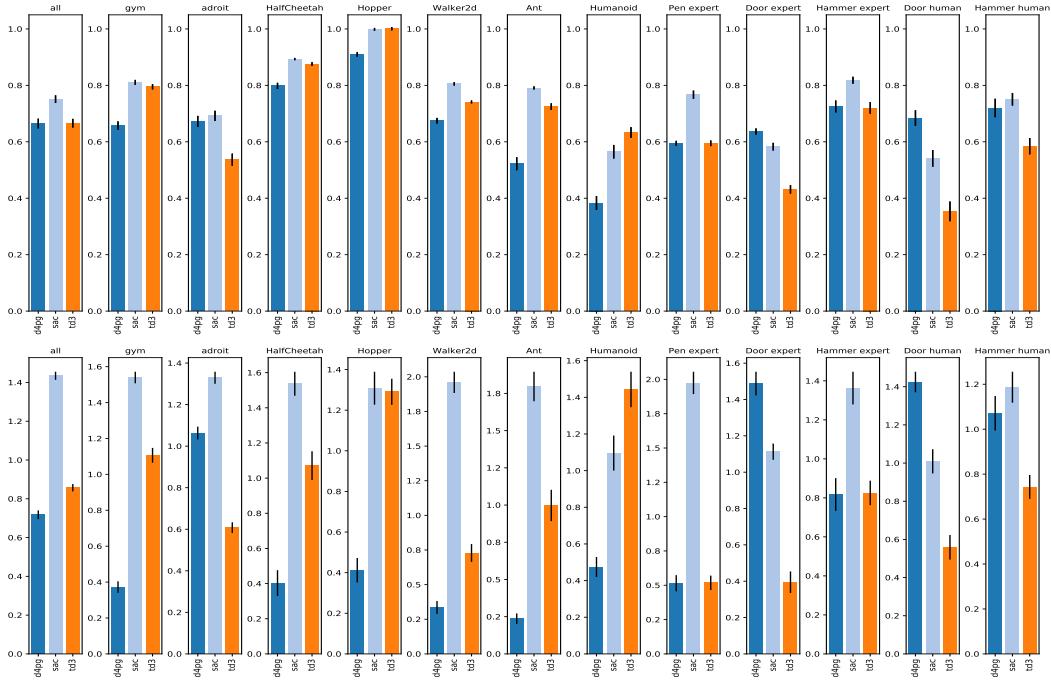


Figure 40: Analysis of choice RL Algorithm (C8): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

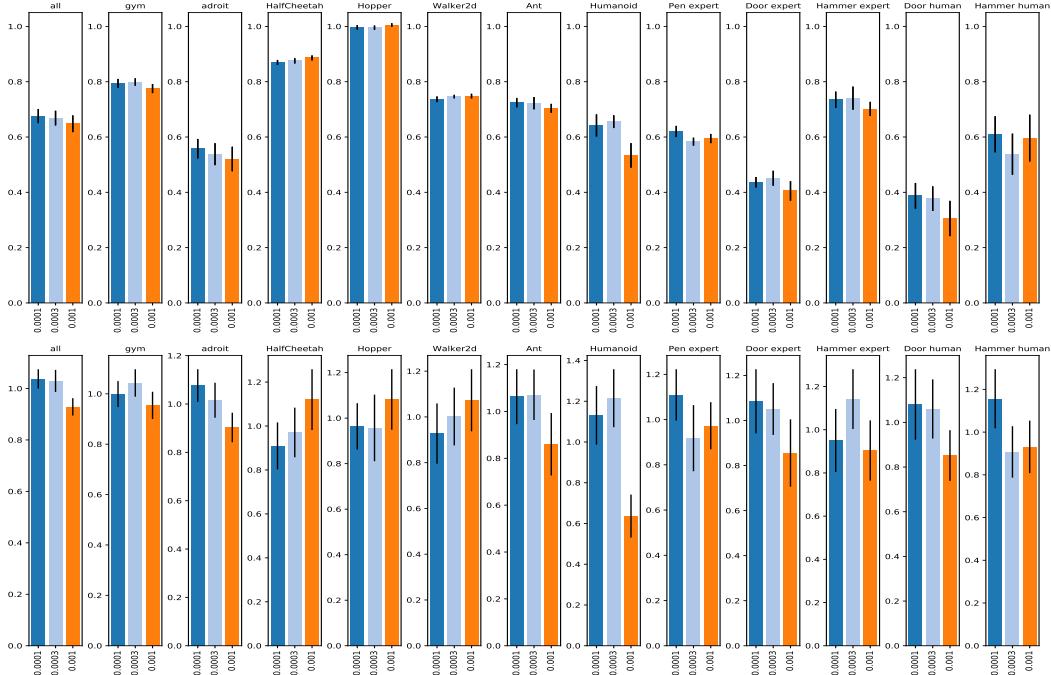


Figure 41: Analysis of choice TD3 policy learning rate (C19): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

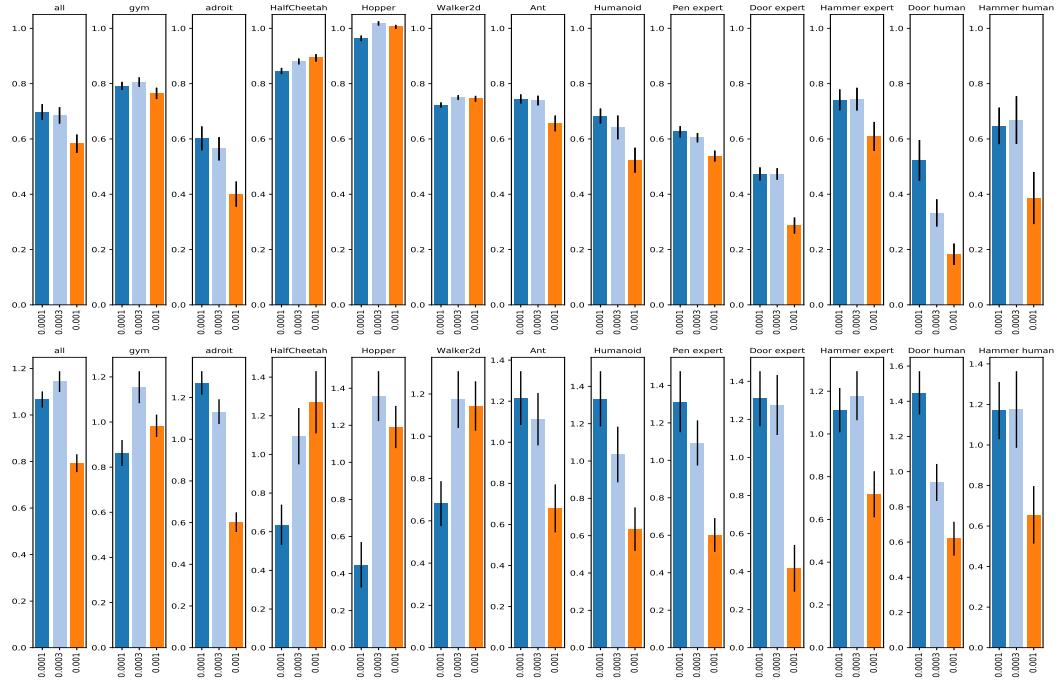


Figure 42: Analysis of choice TD3 critic learning rate (C20): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

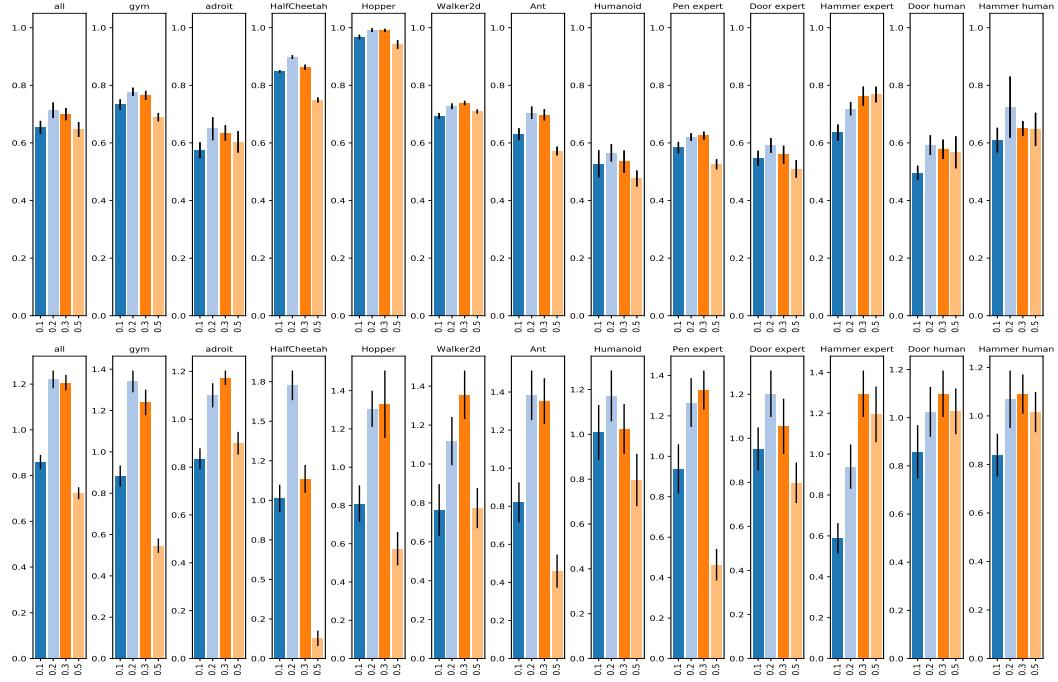


Figure 43: Analysis of choice behavioral policy noise (C21): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

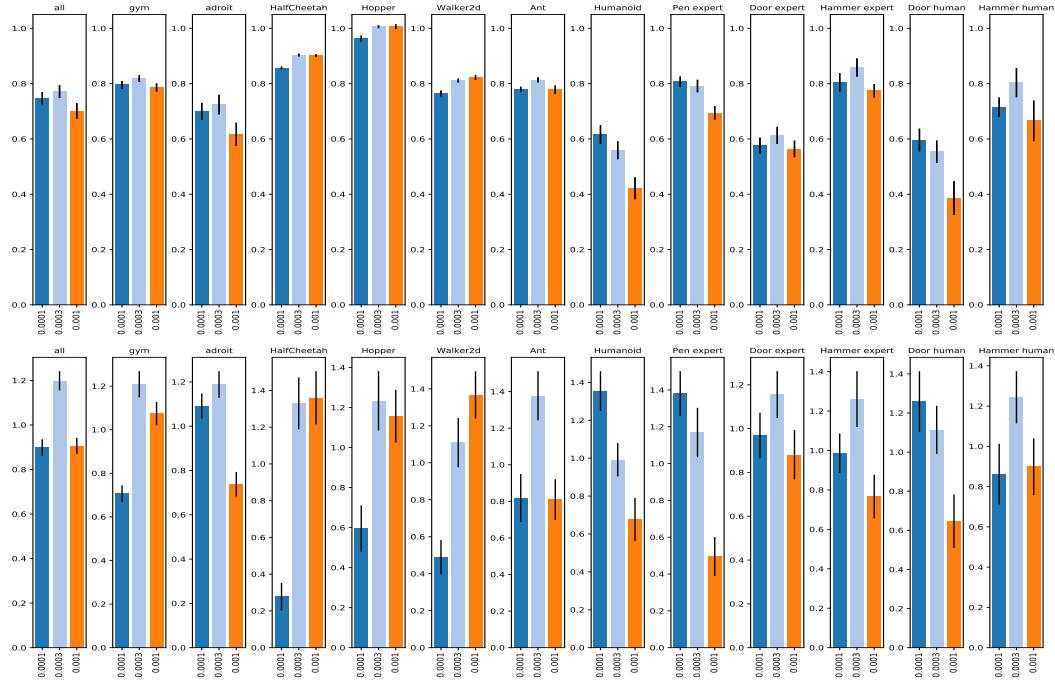


Figure 44: Analysis of choice SAC learning rate (C17): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

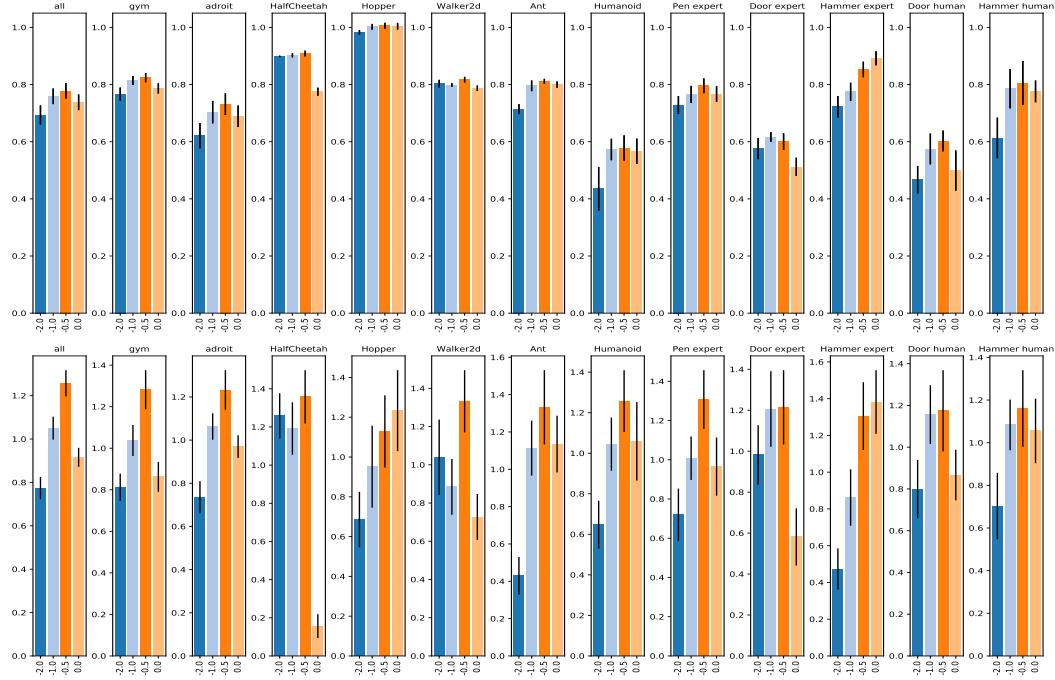


Figure 45: Analysis of choice SAC entropy per dimension (C16): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

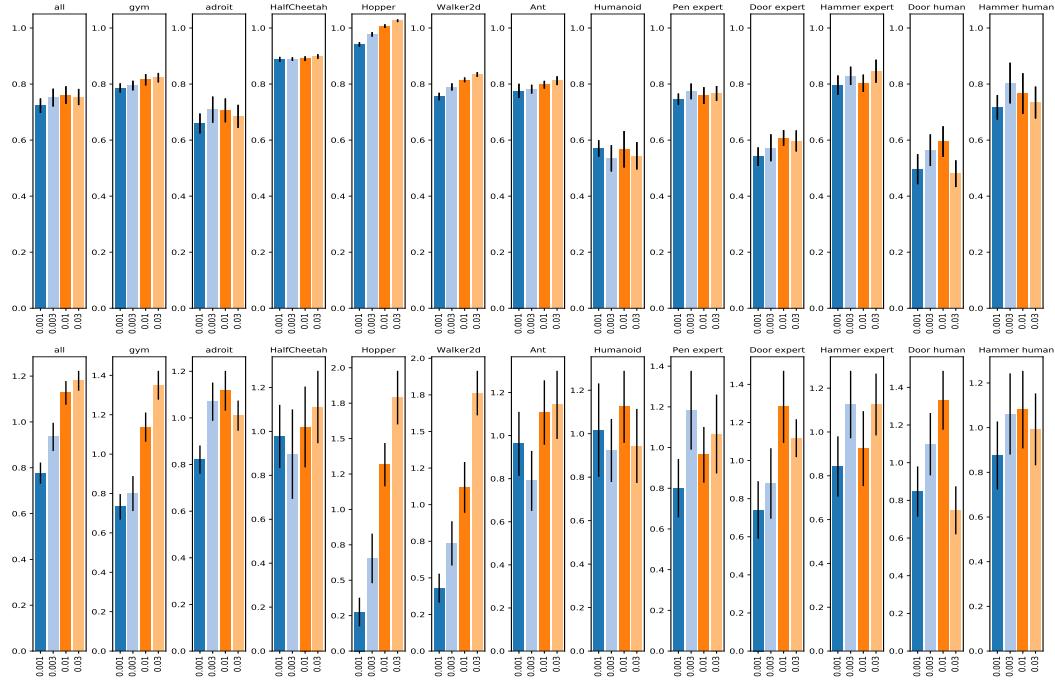


Figure 46: Analysis of choice SAC polyak τ (C18): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

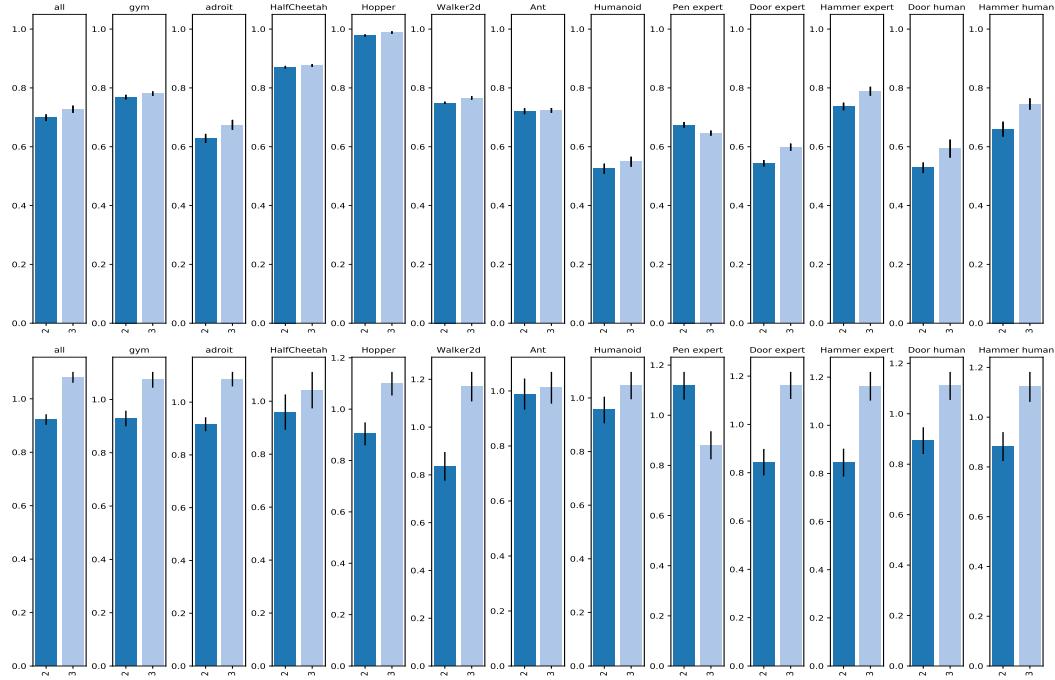


Figure 47: Analysis of choice critic MLP depth (C3): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

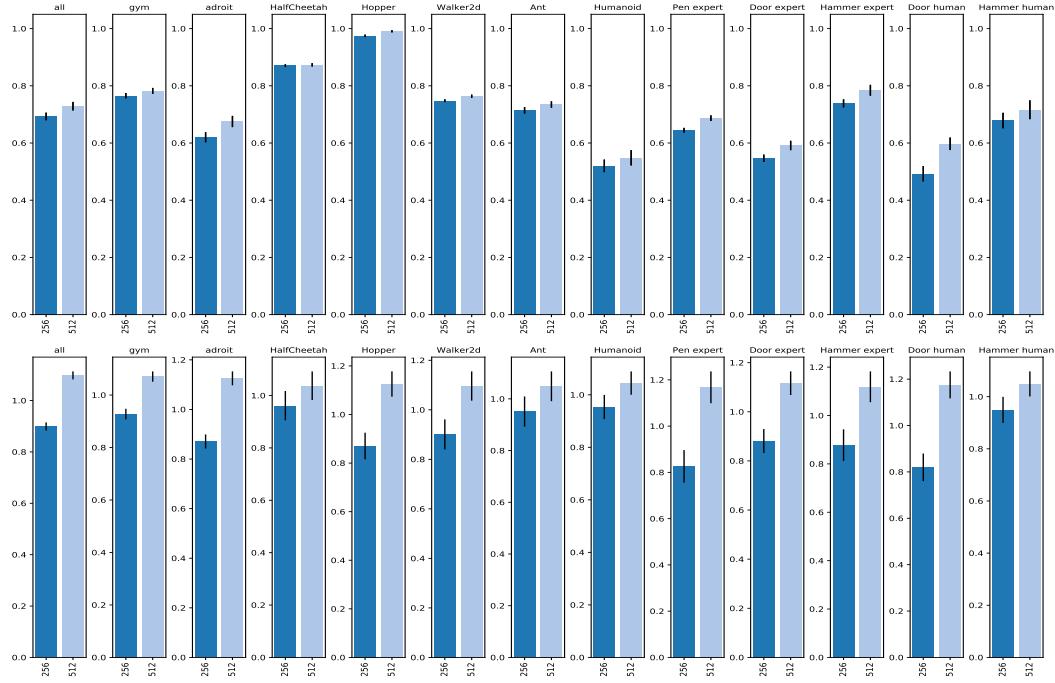


Figure 48: Analysis of choice critic MLP width (C4): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

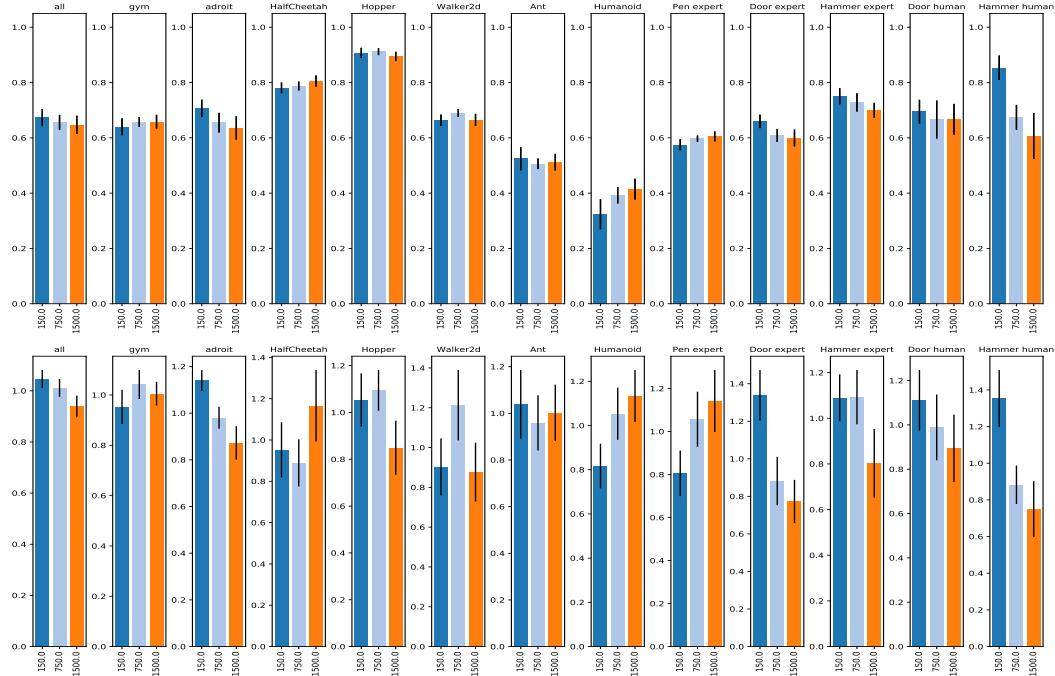


Figure 49: Analysis of choice VMax (C24): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

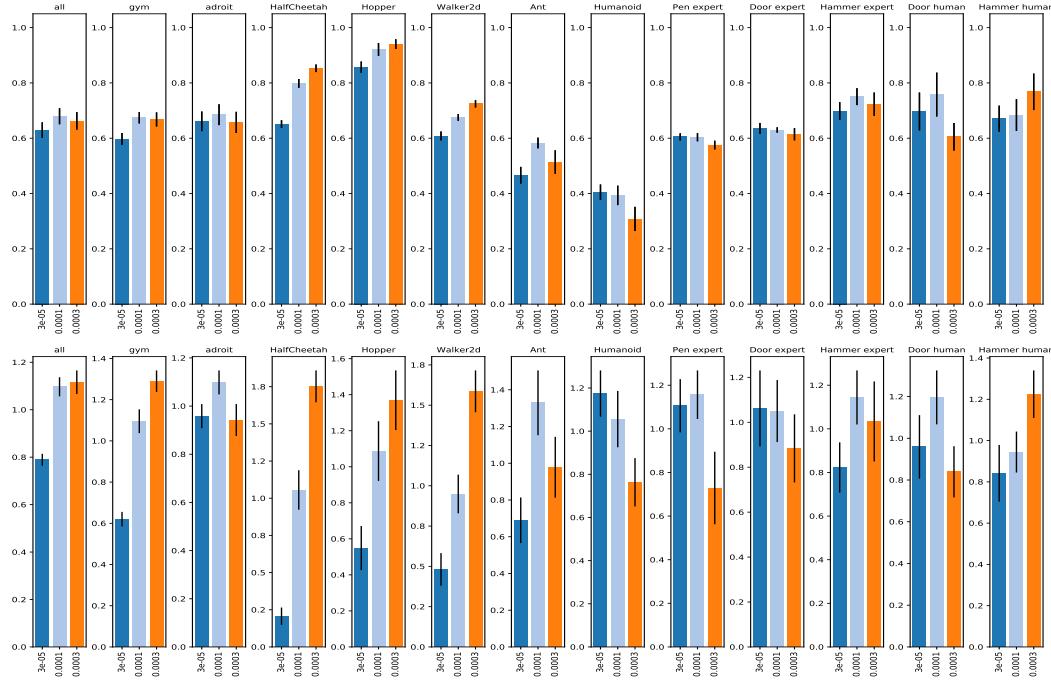


Figure 50: Analysis of choice D4PG learning rate (C26): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

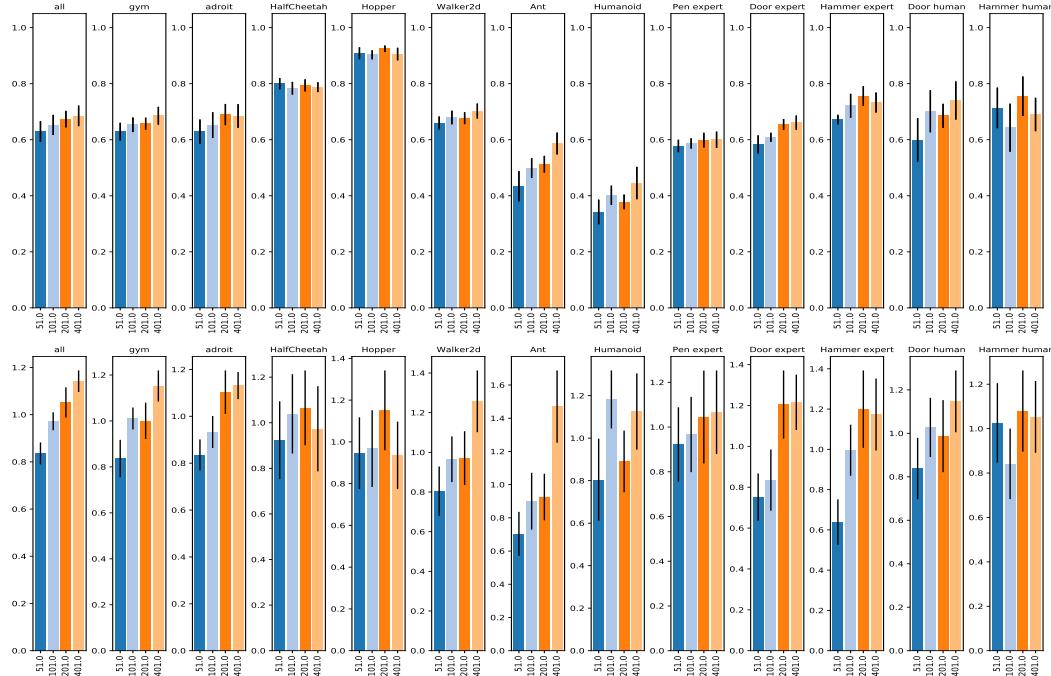


Figure 51: Analysis of choice number of atoms (C23): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

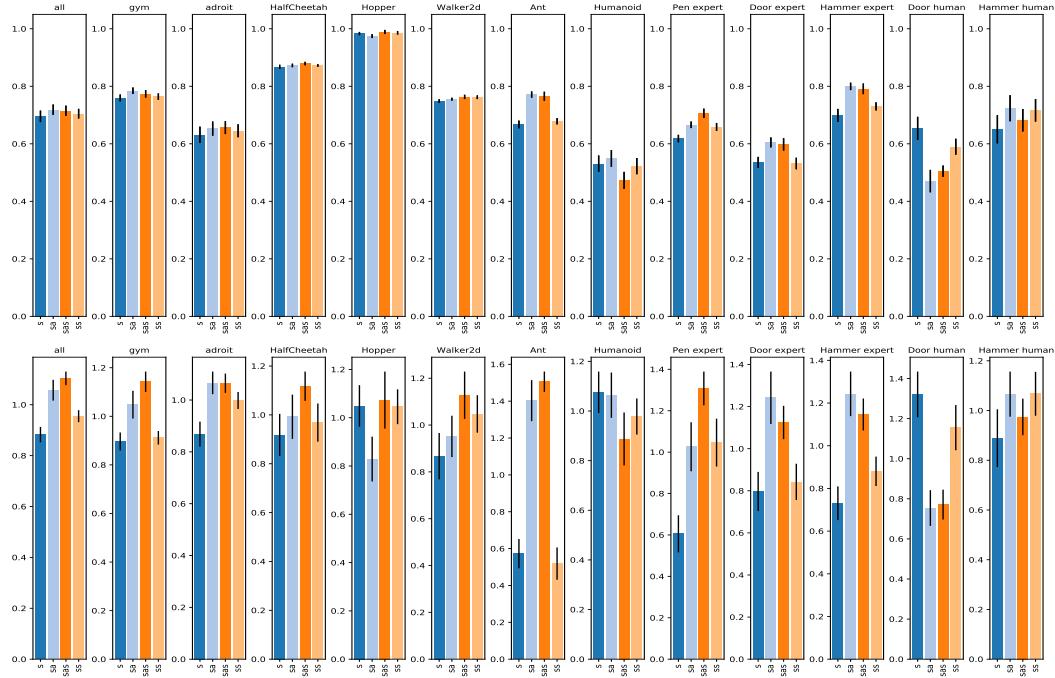


Figure 52: Analysis of choice discriminator input (C35): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

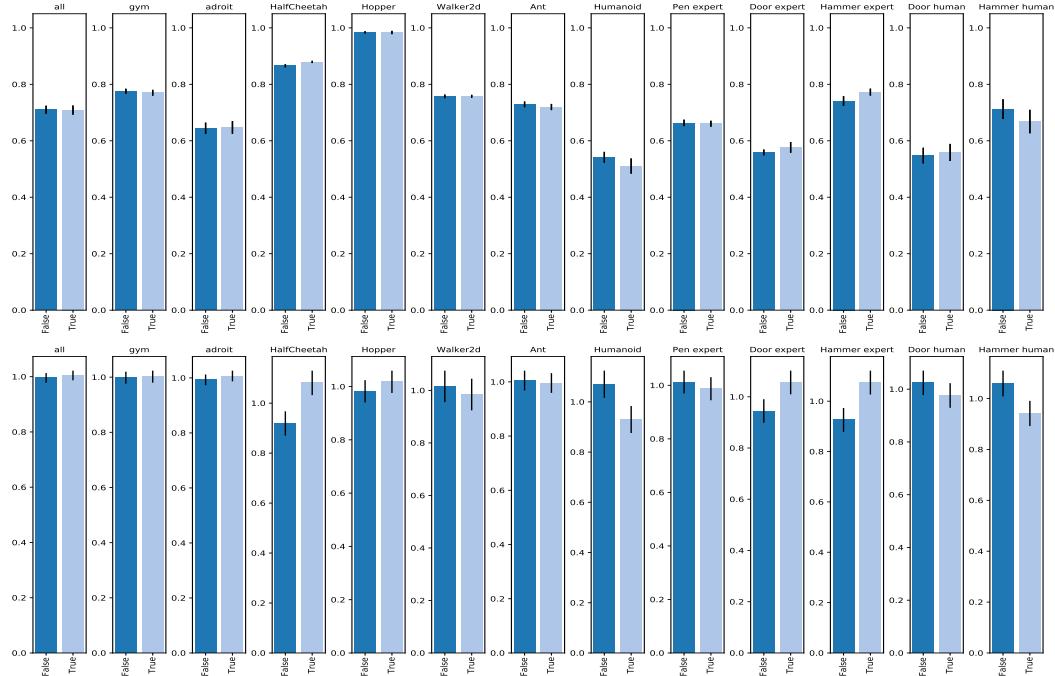


Figure 53: Analysis of choice reward shaping (C39): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

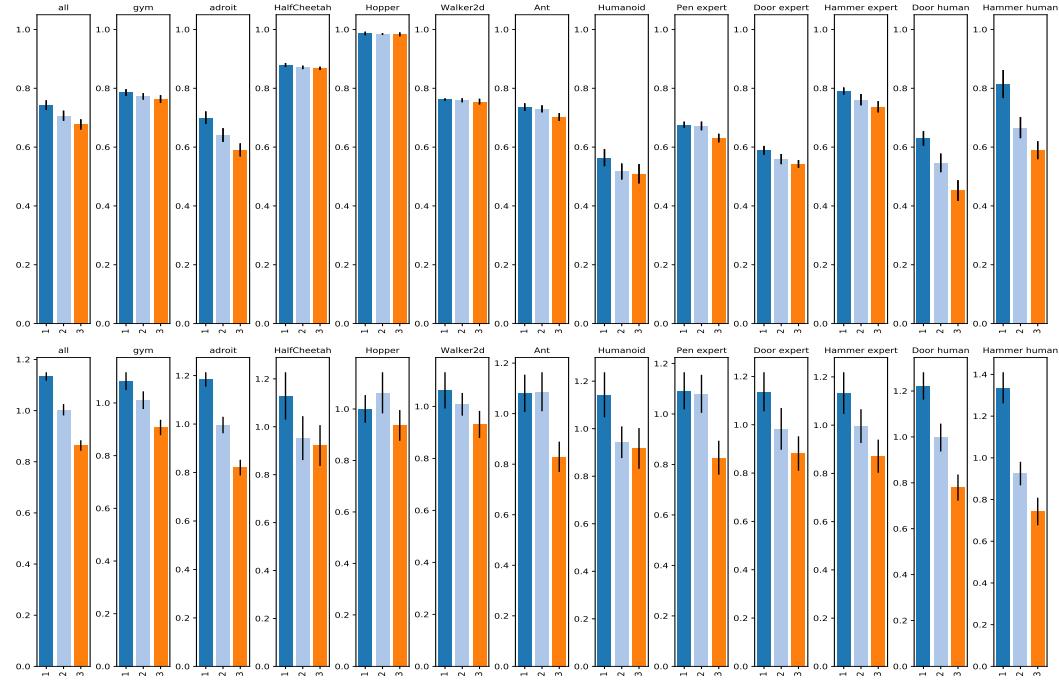


Figure 54: Analysis of choice discriminator MLP depth (C36): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

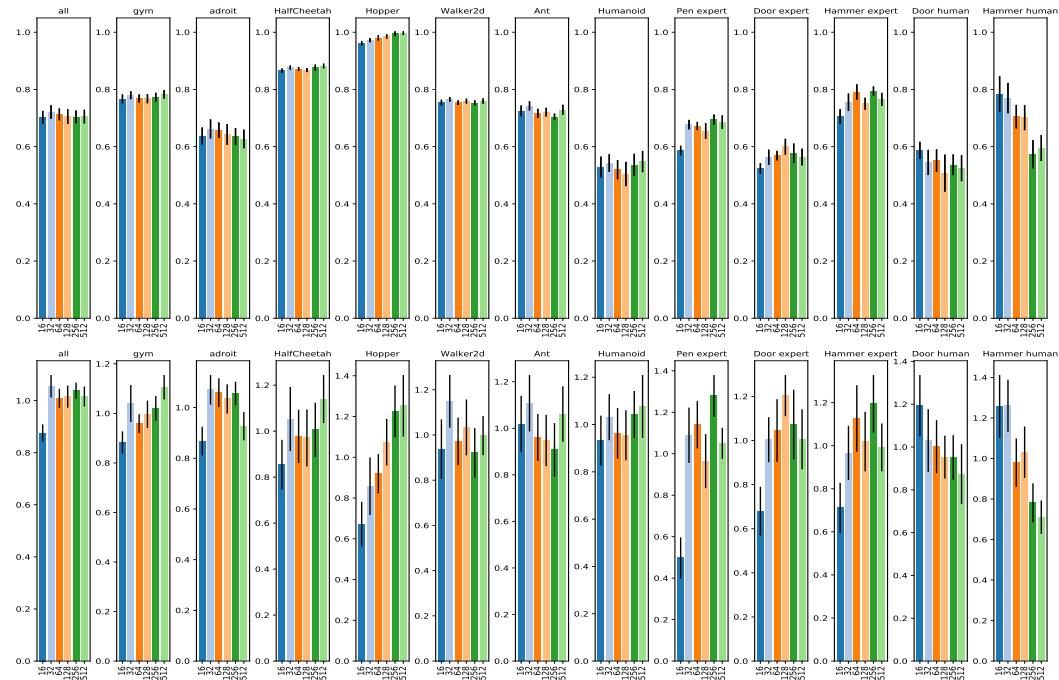


Figure 55: Analysis of choice discriminator MLP width (C37): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

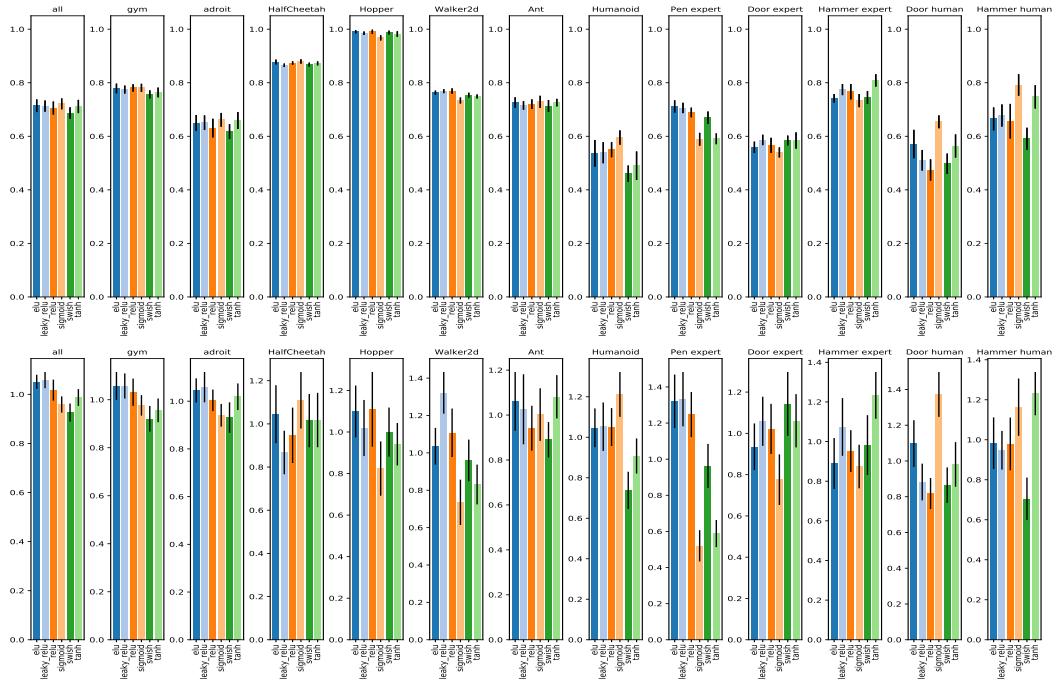


Figure 56: Analysis of choice discriminator activation (C38): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

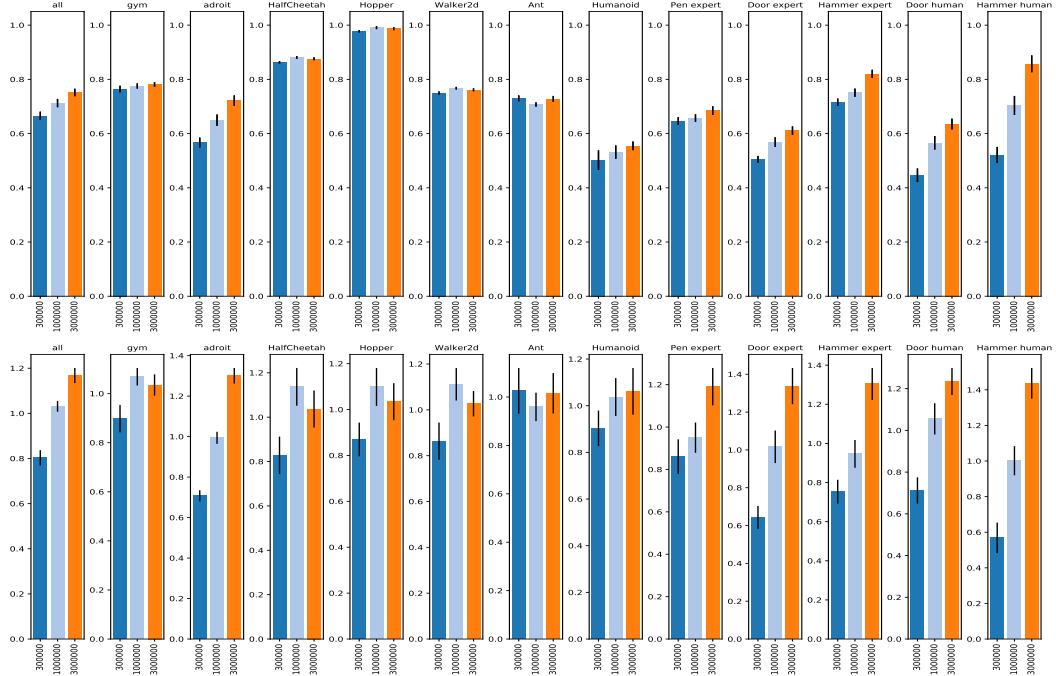


Figure 57: Analysis of choice discriminator replay buffer size (C43): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

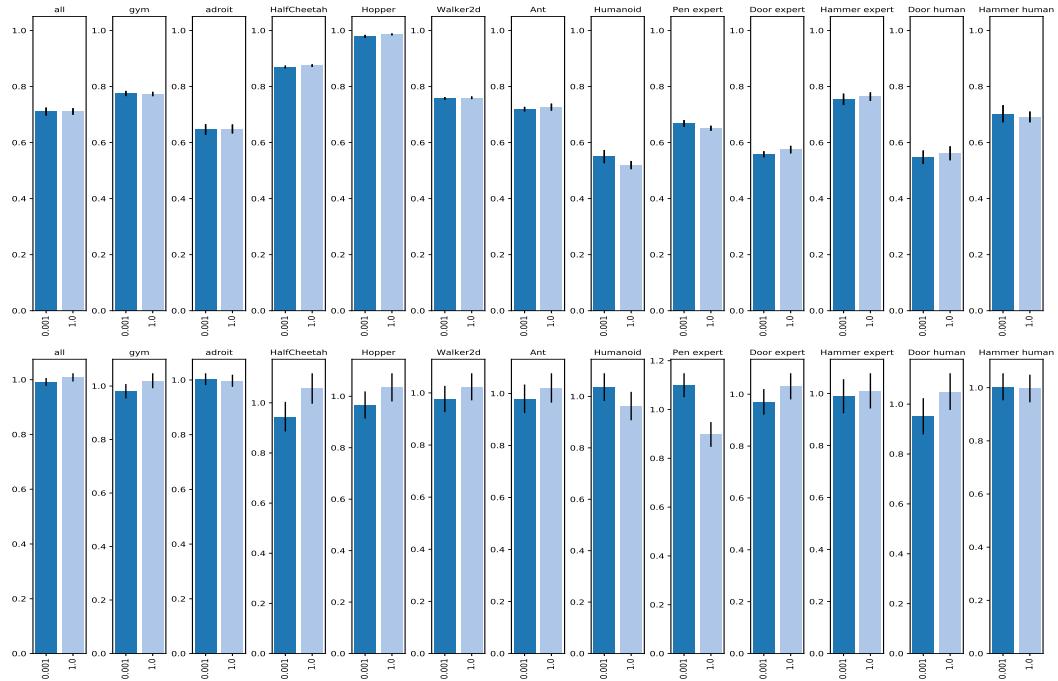


Figure 58: Analysis of choice discriminator last layer init scale (C41): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

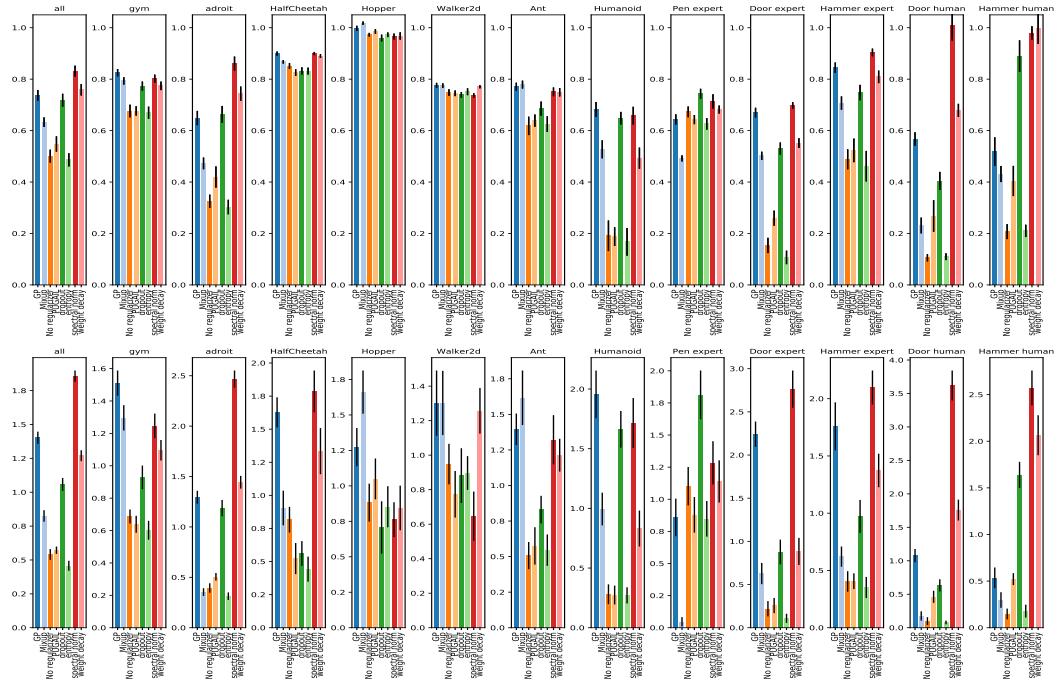


Figure 59: Analysis of choice discriminator regularizer (C45): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

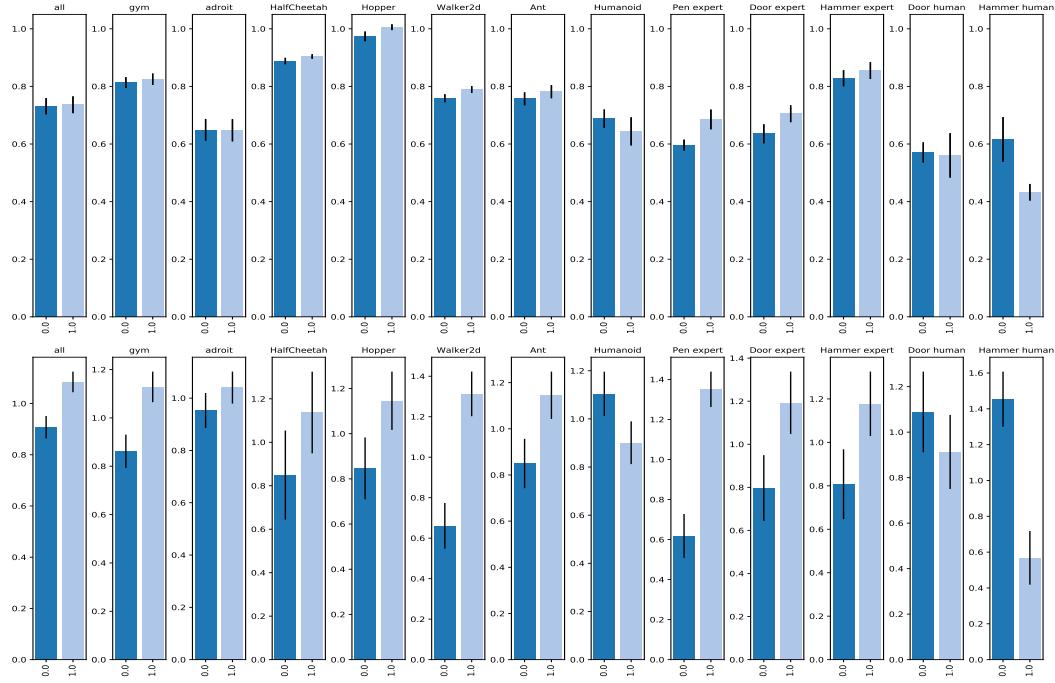


Figure 60: Analysis of choice gradient penalty k (C46): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

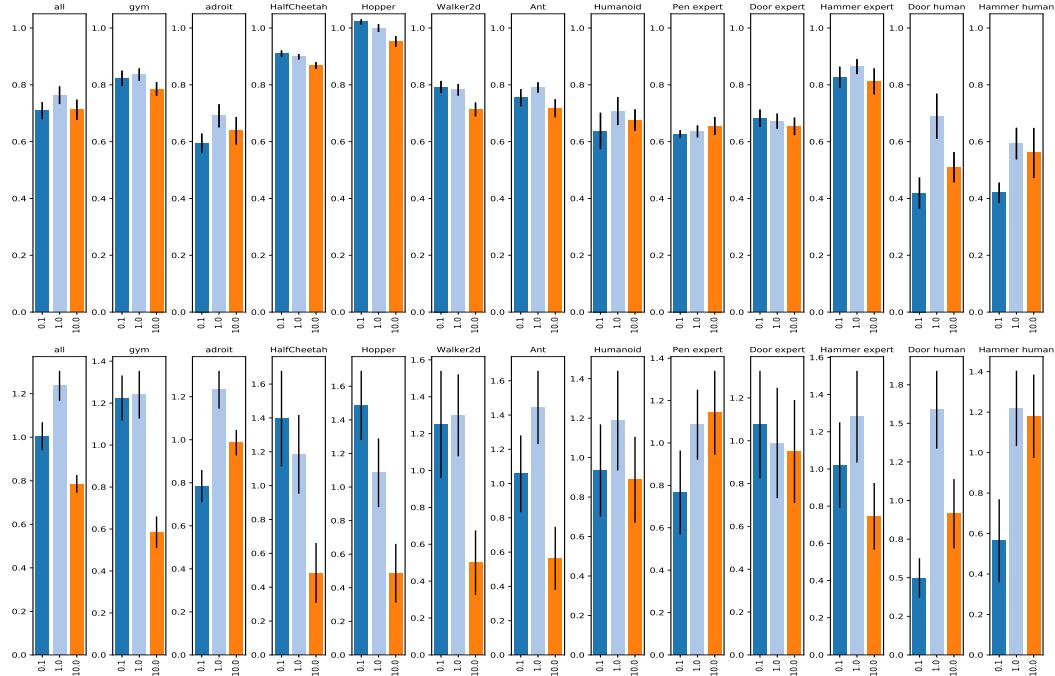


Figure 61: Analysis of choice gradient penalty λ (C47): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

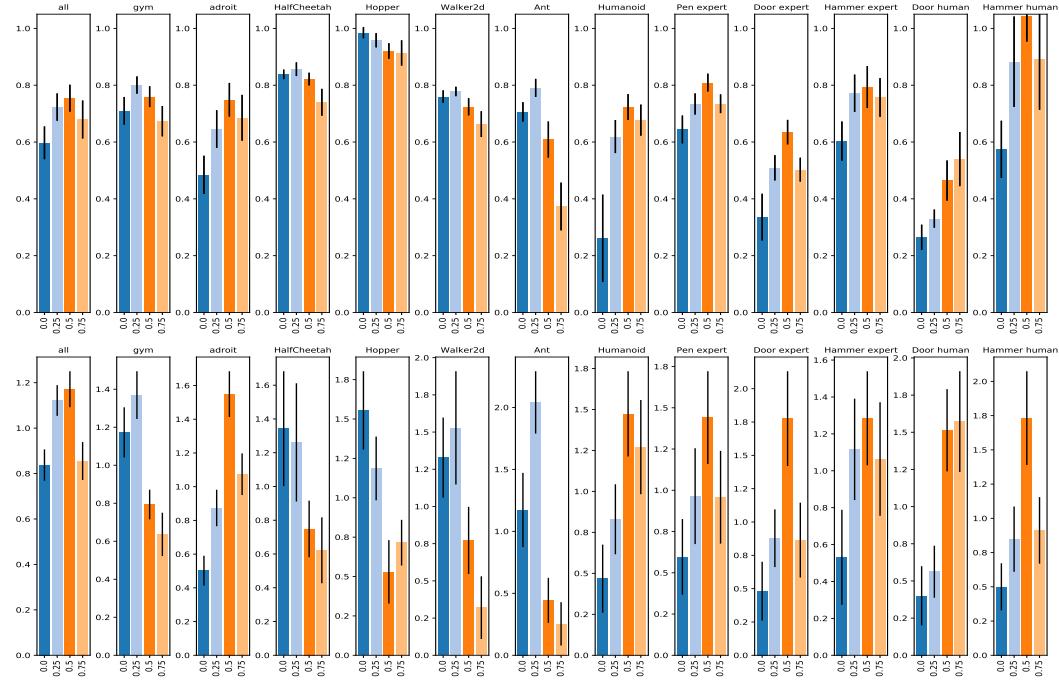


Figure 62: Analysis of choice dropout input rate (C52): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

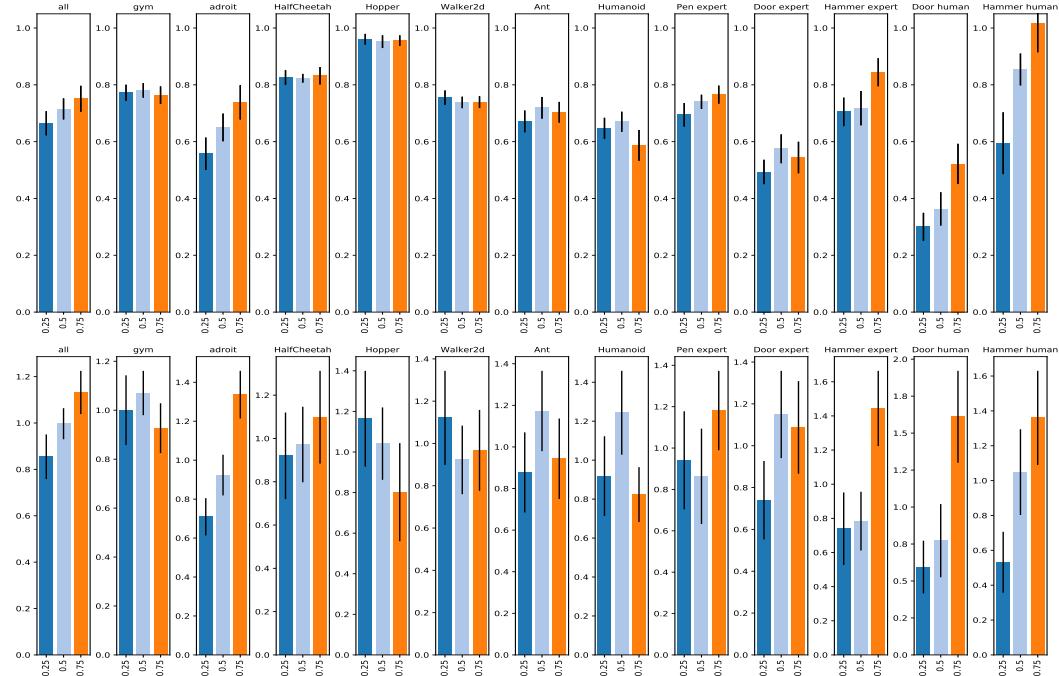


Figure 63: Analysis of choice dropout hidden rate (C51): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

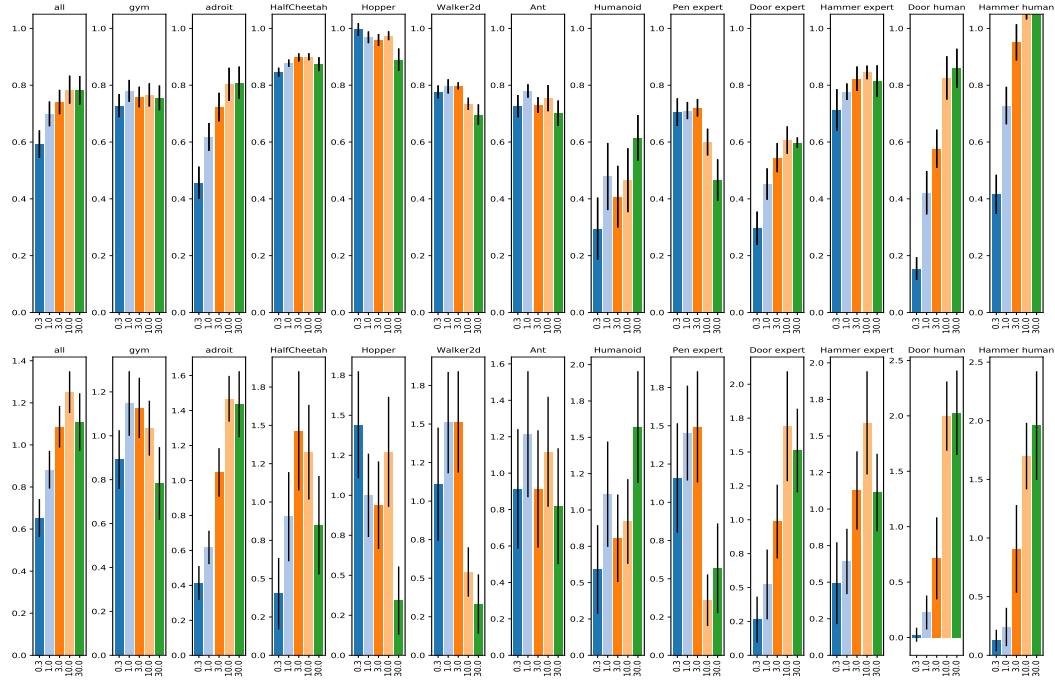


Figure 64: Analysis of choice weight decay λ (C53): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

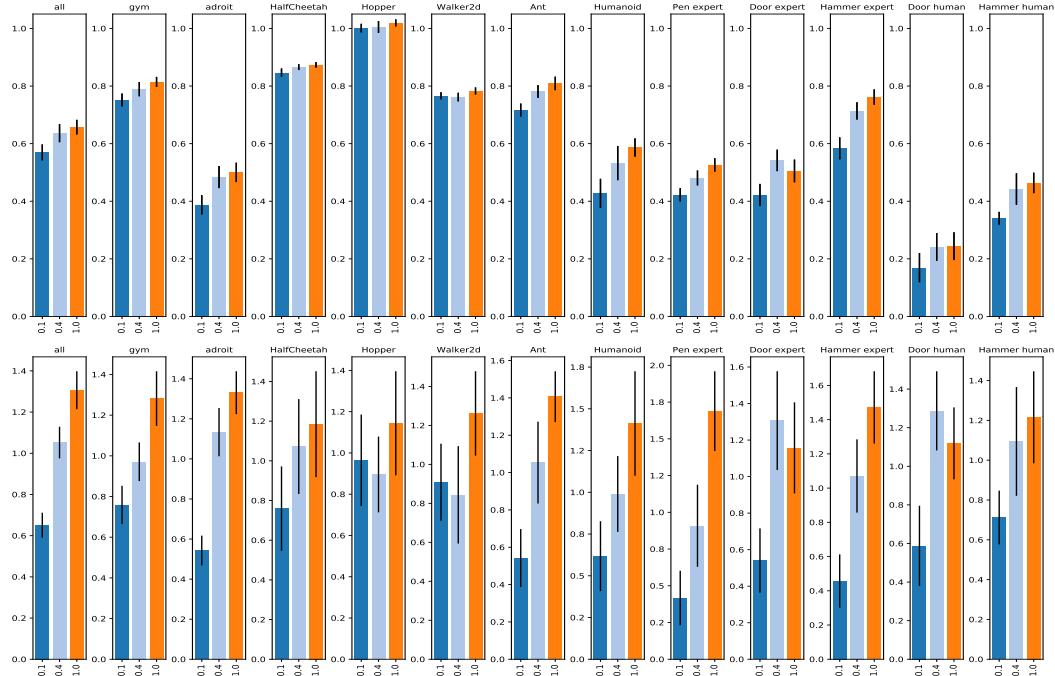


Figure 65: Analysis of choice mixup α (C48): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

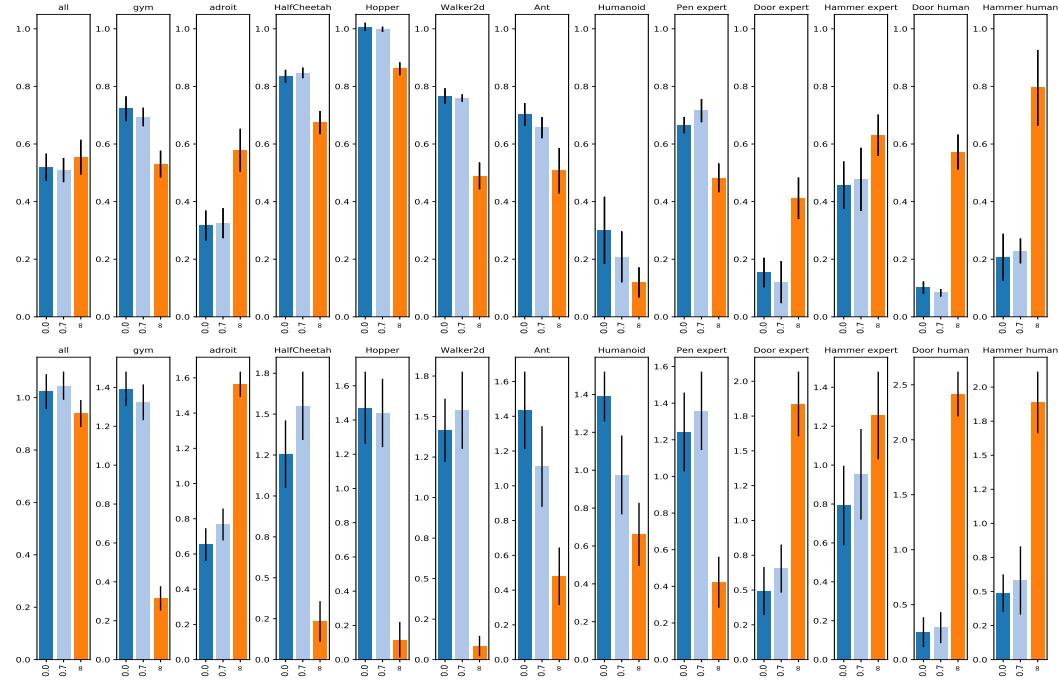


Figure 66: Analysis of choice PUGAIL β (C50): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

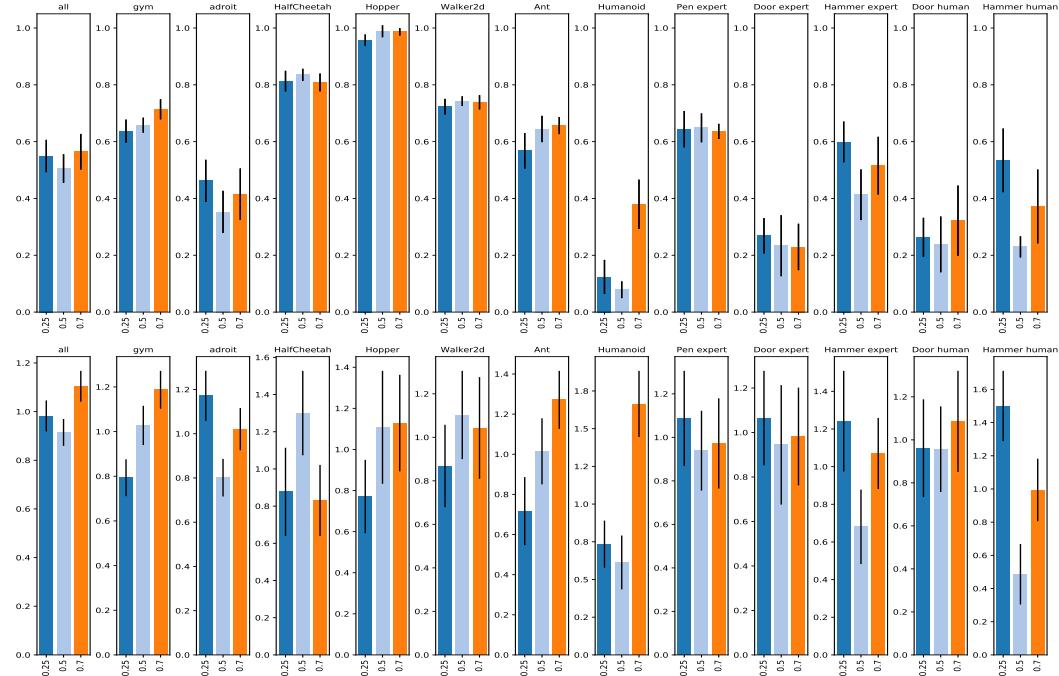


Figure 67: Analysis of choice PUGAIL η (C49): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

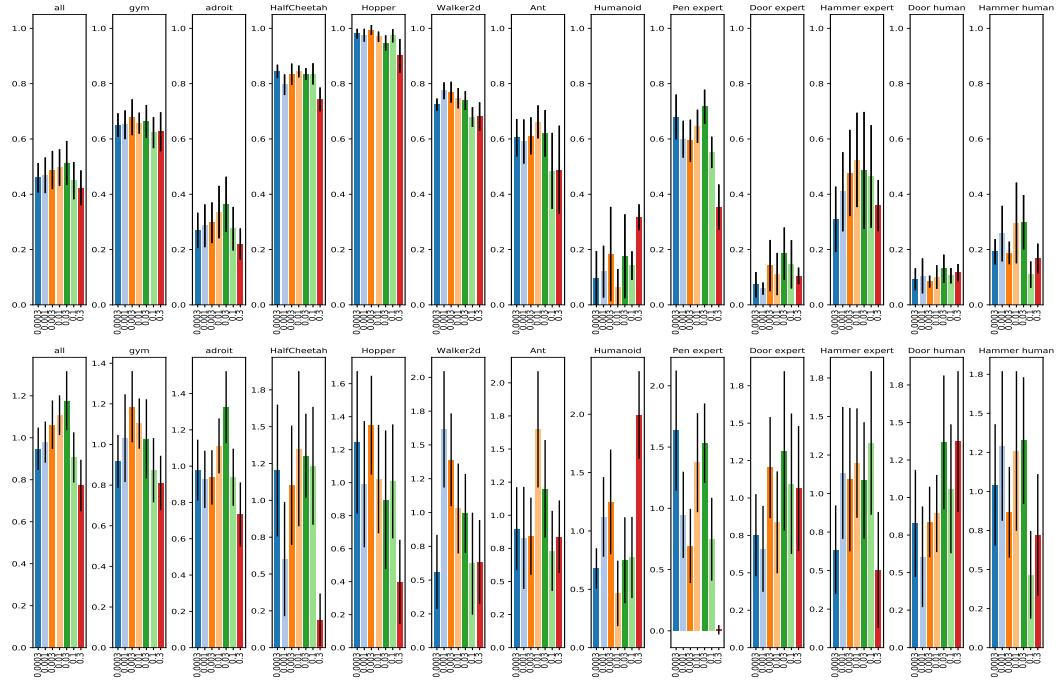


Figure 68: Analysis of choice entropy λ (C54): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

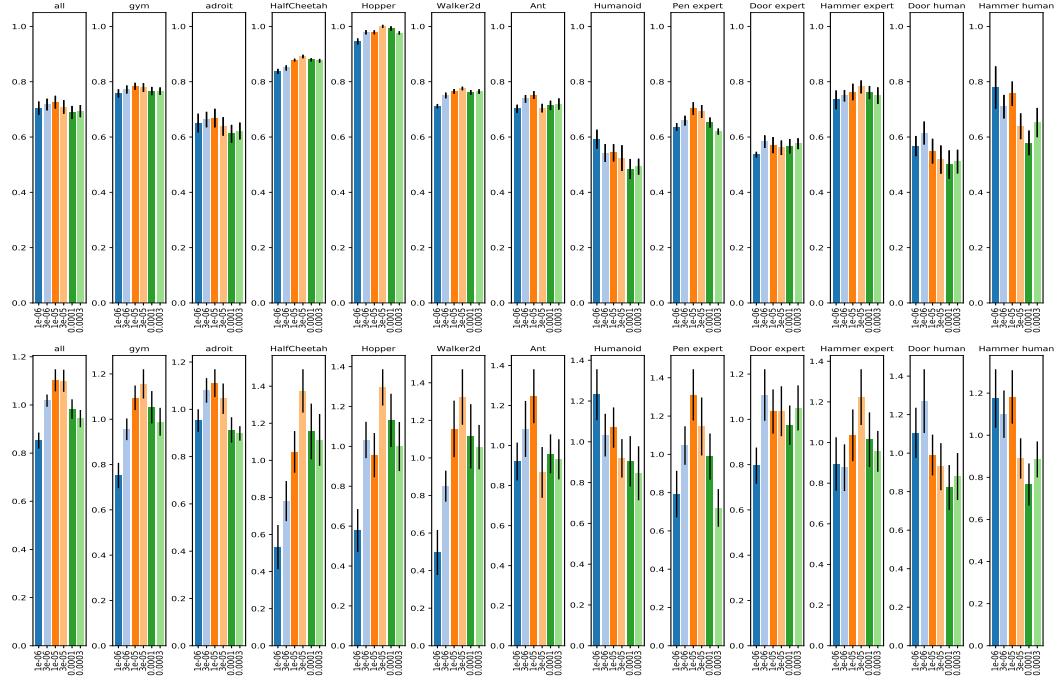


Figure 69: Analysis of choice discriminator learning rate (C42): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

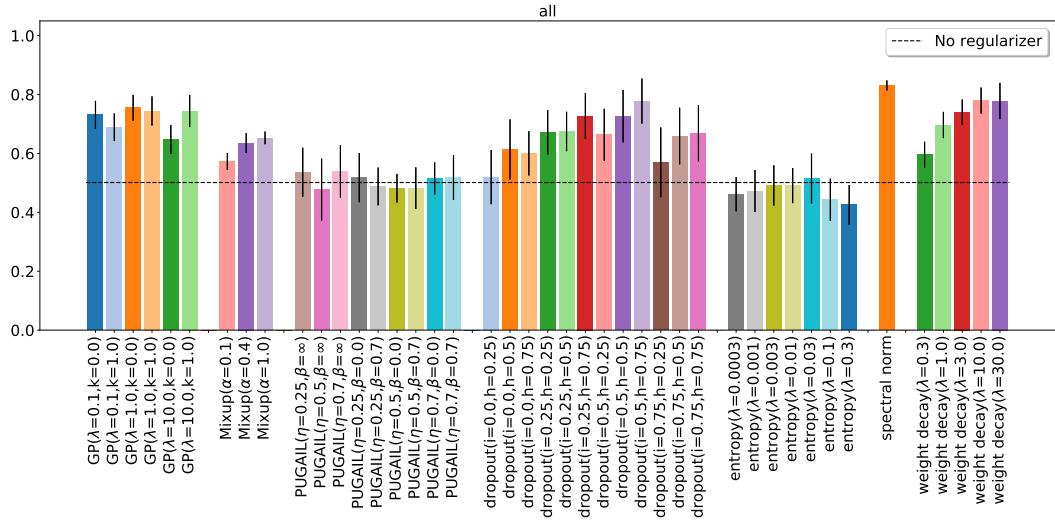


Figure 70: 95th percentile of performance scores conditioned on discriminator regularizer (C45) and regularizers' HPs averaged across all environments.

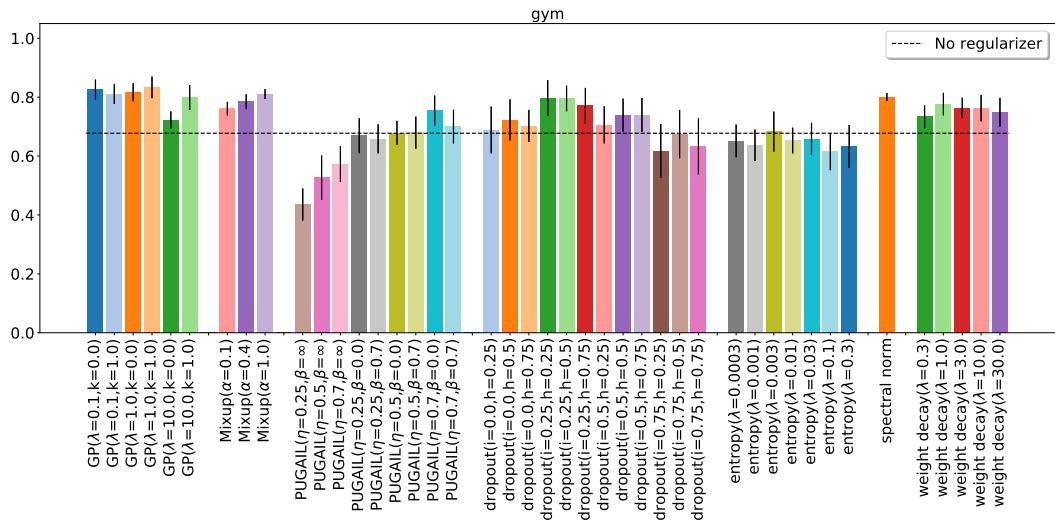


Figure 71: 95th percentile of performance scores conditioned on discriminator regularizer (C45) and regularizers' HPs averaged across OpenAI Gym environments.

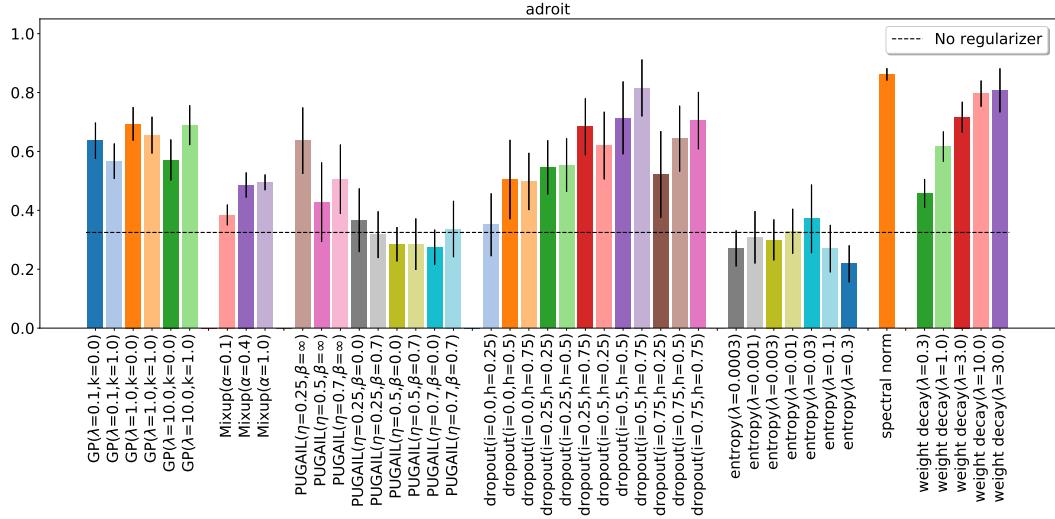


Figure 72: 95th percentile of performance scores conditioned on discriminator regularizer (C45) and regularizers' HPs averaged across Adroit environments.

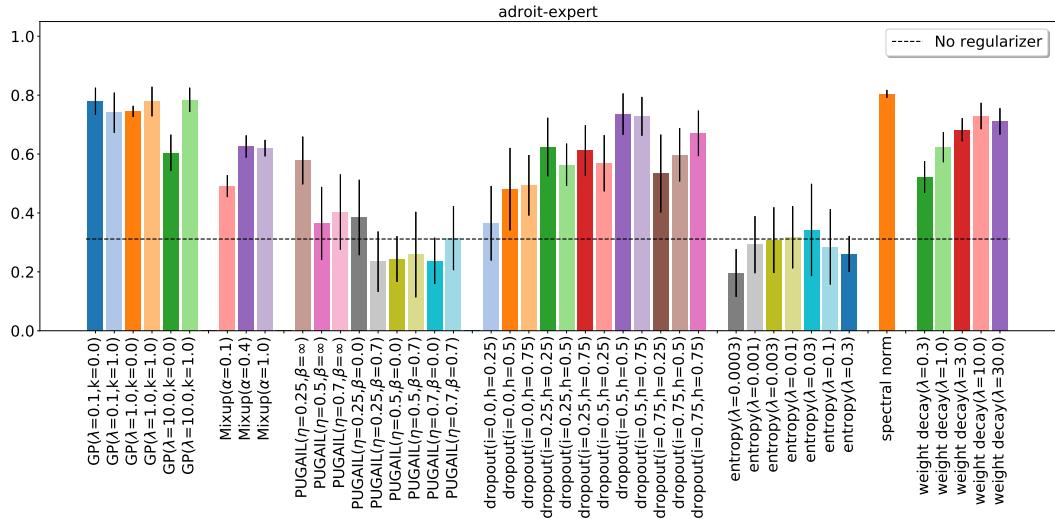


Figure 73: 95th percentile of performance scores conditioned on discriminator regularizer (C45) and regularizers' HPs averaged across door-expert and hammer-expert tasks.

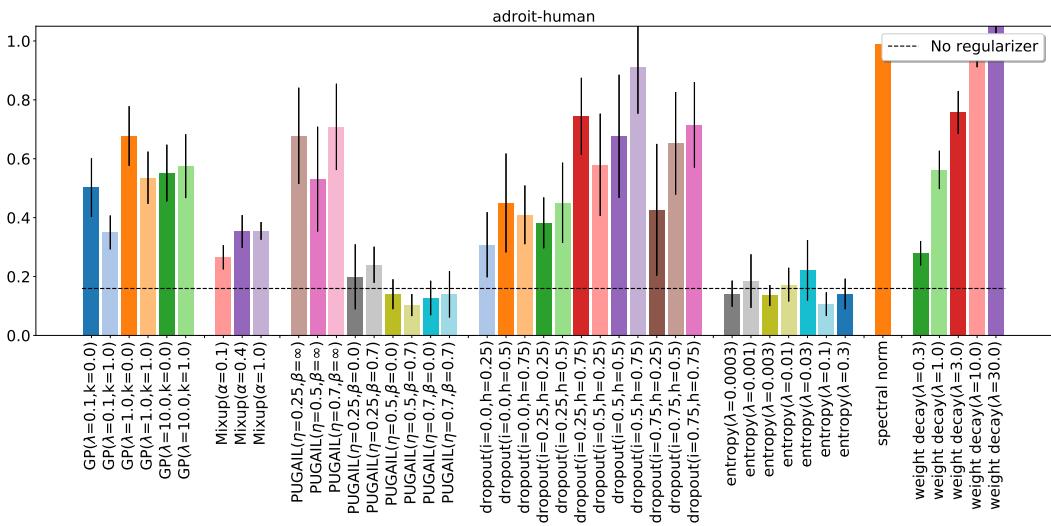


Figure 74: 95th percentile of performance scores conditioned on discriminator regularizer (C45) and regularizers' HPs averaged across door-human and hammer-human tasks.

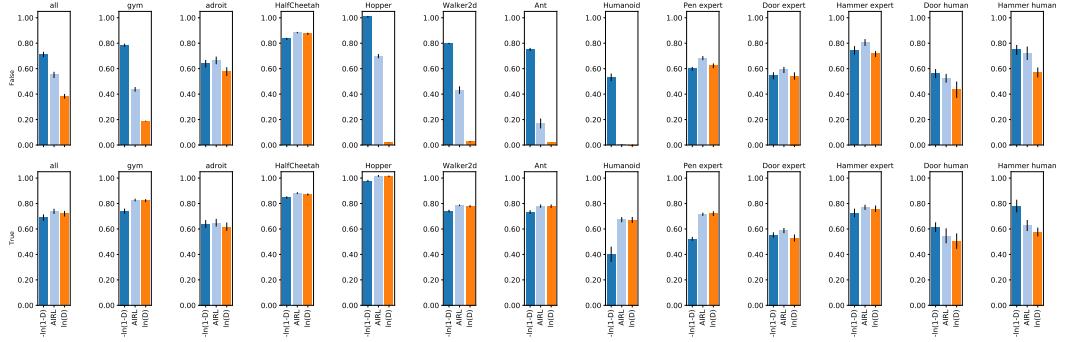


Figure 75: 95th percentile of performance scores conditioned on absorbing state (C32)(rows) and reward function (C30)(bars).

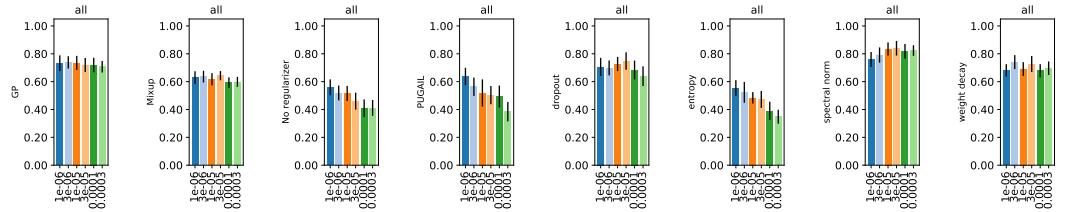


Figure 76: 95th percentile of performance scores conditioned on discriminator regularizer (C45)(subplots) and discriminator learning rate (C42)(bars).

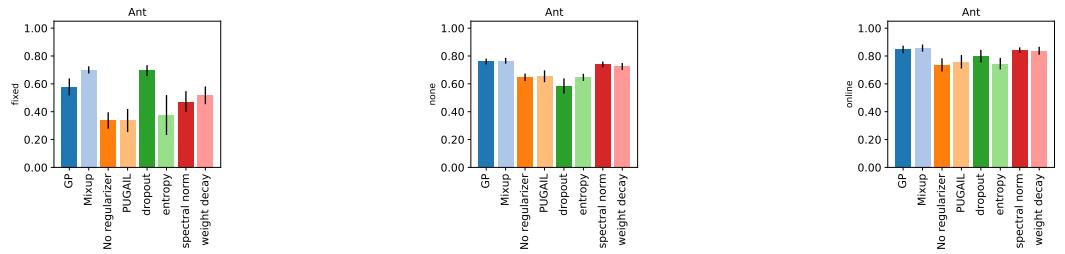


Figure 77: 95th percentile of performance scores conditioned on observation normalization (C55)(subplots) and discriminator regularizer (C45)(bars) in the Ant environment.

H Experiment trade-offs

H.1 Design

For each of the 10 tasks, we sampled 7991 choice configurations where we sampled the following choices independently and uniformly from the following ranges:

- RL Algorithm (C8): {d4pg, sac, td3}
 - For the case “RL Algorithm (C8) = sac”, we further sampled the sub-choices:
 - * SAC learning rate (C17): {0.0001, 0.0003, 0.001}
 - * SAC entropy per dimension (C16): {-2.0, -1.0, -0.5, 0.0}
 - * SAC polyak τ (C18): {0.001, 0.003, 0.01, 0.03}
 - For the case “RL Algorithm (C8) = d4pg”, we further sampled the sub-choices:
 - * D4PG learning rate (C26): {3e-05, 0.0001, 0.0003}
 - * behavioral policy noise (C21): {0.1, 0.2, 0.3, 0.5}
 - * VMax (C24): {150.0, 750.0, 1500.0}
 - * number of atoms (C23): {51.0, 101.0, 201.0, 401.0}
 - * N-step returns (C25): {1.0, 3.0, 5.0}
 - For the case “RL Algorithm (C8) = td3”, we further sampled the sub-choices:
 - * TD3 policy learning rate (C19): {0.0001, 0.0003, 0.001}
 - * TD3 critic learning rate (C20): {0.0001, 0.0003, 0.001}
 - * TD3 gradient clipping (C22): {40.0, ∞ }
 - * behavioral policy noise (C21): {0.1, 0.2, 0.3, 0.5}
- RL replay buffer size (C28): {300000, 1000000, 3000000}
- policy MLP depth (C1): {1, 2, 3}
- policy MLP width (C2): {64, 128, 256, 512}
- critic MLP depth (C3): {2, 3}
- critic MLP width (C4): {256, 512}
- RL activation (C5): {relu, tanh}
- discount γ (C6): {0.97, 0.99}
- BC pretraining (C34): {False, True}
- absorbing state (C32): {False, True}
- discriminator replay buffer size (C43): {300000, 1000000, 3000000}
- reward shaping (C39): {False, True}
- discriminator input (C35): {s, sa, sas, ss}
- discriminator MLP depth (C36): {1, 2, 3}
- discriminator MLP width (C37): {16, 32, 64, 128, 256, 512}
- discriminator activation (C38): {elu, leaky_relu, relu, sigmoid, swish, tanh}
- discriminator last layer init scale (C41): {0.001, 1.0}
- discriminator regularizer (C45): {GP, Mixup, No regularizer, PUGAIL, dropout, entropy, spectral norm, weight decay}
 - For the case “discriminator regularizer (C45) = GP”, we further sampled the sub-choices:
 - * gradient penalty λ (C47): {0.1, 1.0, 10.0}
 - * gradient penalty k (C46): {0.0, 1.0}
 - For the case “discriminator regularizer (C45) = Mixup”, we further sampled the sub-choices:
 - * mixup α (C48): {0.1, 0.4, 1.0}
 - For the case “discriminator regularizer (C45) = PUGAIL”, we further sampled the sub-choices:

- * PUGAIL η (C49): {0.25, 0.5, 0.7}
- * PUGAIL β (C50): {0.0, 0.7, ∞ }
- For the case “discriminator regularizer (C45) = entropy”, we further sampled the sub-choices:
 - * entropy λ (C54): {0.0003, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3}
- For the case “discriminator regularizer (C45) = weight decay”, we further sampled the sub-choices:
 - * weight decay λ (C53): {0.3, 1.0, 3.0, 10.0, 30.0}
- For the case “discriminator regularizer (C45) = dropout”, we further sampled the sub-choices:
 - * dropout input rate (C52): {0.0, 0.25, 0.5, 0.75}
 - * dropout hidden rate (C51): {0.25, 0.5, 0.75}
- observation normalization (C55): {fixed, none}
- evaluation behavior policy type (C29): {average, mode, stochastic}
- discriminator learning rate (C42): {1e-06, 3e-06, 1e-05, 3e-05, 0.0001, 0.0003}
- replay ratio (C27): {64, 128, 256, 512, 1024}
- batch size (C7): {64, 128, 256, 512, 1024}
- discriminator to RL updates ratio (C44): {1, 2}
- number of combined batches (C56): {1, 2, 4, 8, 16, 32, 64}
- reward function (C30): {-ln(1-D), AIRL, ln(D)}

H.2 Results

For each of the sampled choice configurations we compute the performance metric as described in Section 2. We report aggregate statistics of the experiment in Tables 12–15 as well as training curves in Figure 78. We further provide per-choice analyses in Figures 79–82.

Table 12: Quantiles of the *final* agent performance across HP configurations for OpenAI Gym tasks.

	Ant	HalfCheetah	Hopper	Humanoid	Walker2d
90%	0.81	1.04	1.18	0.14	0.97
95%	0.94	1.08	1.19	0.62	1.00
99%	1.04	1.15	1.22	0.98	1.03
Max	1.15	1.41	1.31	1.05	1.16

Table 13: Quantiles of the *final* agent performance across HP configurations for Adroit tasks.

	Door expert	Door human	Hammer expert	Hammer human	Pen expert
90%	0.71	0.25	1.03	0.45	0.70
95%	0.89	0.71	1.25	1.19	0.86
99%	1.04	2.12	1.36	2.95	1.07
Max	1.15	3.79	1.44	5.27	1.34

Table 14: Quantiles of the *average* agent performance during training across HP configurations for OpenAI Gym tasks.

	Ant	HalfCheetah	Hopper	Humanoid	Walker2d
90%	0.48	0.75	0.90	0.12	0.63
95%	0.63	0.83	0.98	0.32	0.72
99%	0.77	0.92	1.06	0.62	0.83
Max	0.89	1.00	1.10	0.85	0.92

Table 15: Quantiles of the *average* agent performance during training across HP configurations for Adroit tasks.

	Door expert	Door human	Hammer expert	Hammer human	Pen expert
90%	0.38	0.26	0.54	0.39	0.50
95%	0.53	0.49	0.71	0.65	0.63
99%	0.74	1.02	0.91	1.21	0.82
Max	0.94	2.05	1.17	2.13	1.01

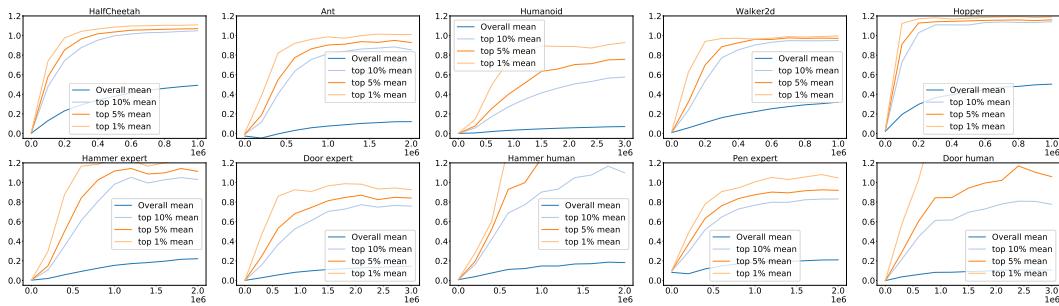


Figure 78: Training curves.

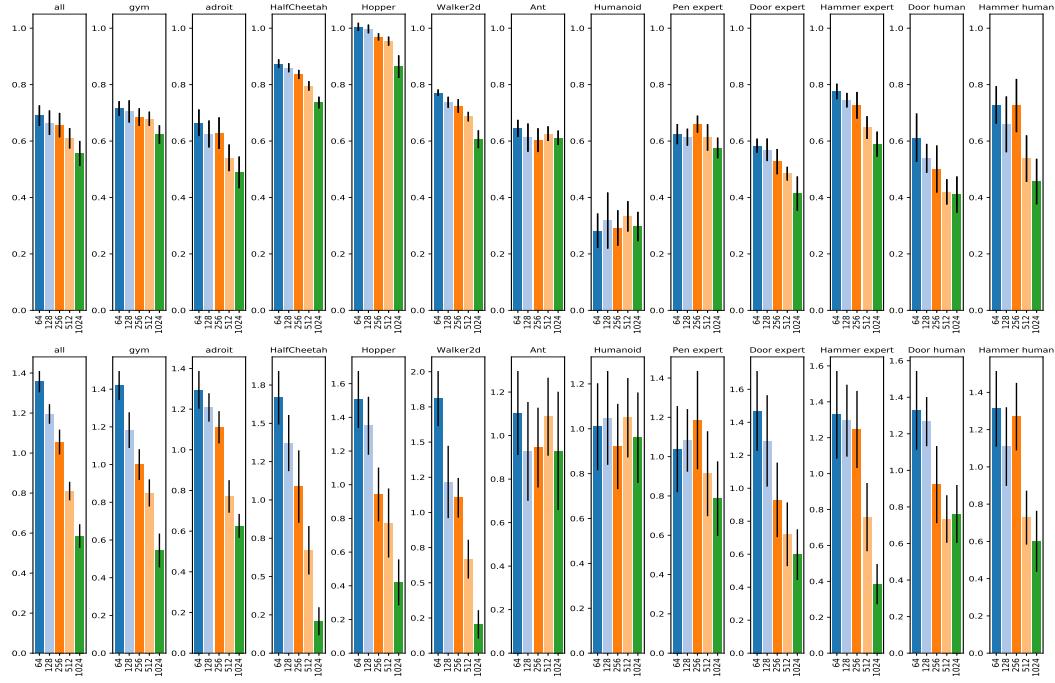


Figure 79: Analysis of choice batch size (C7): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

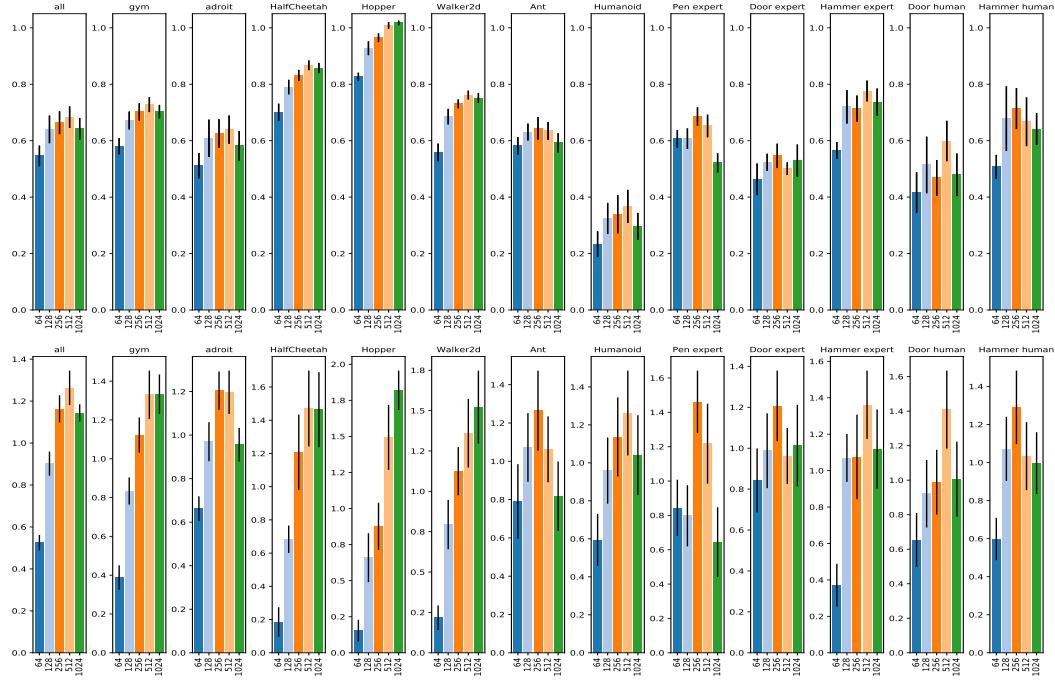


Figure 80: Analysis of choice replay ratio (C27): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

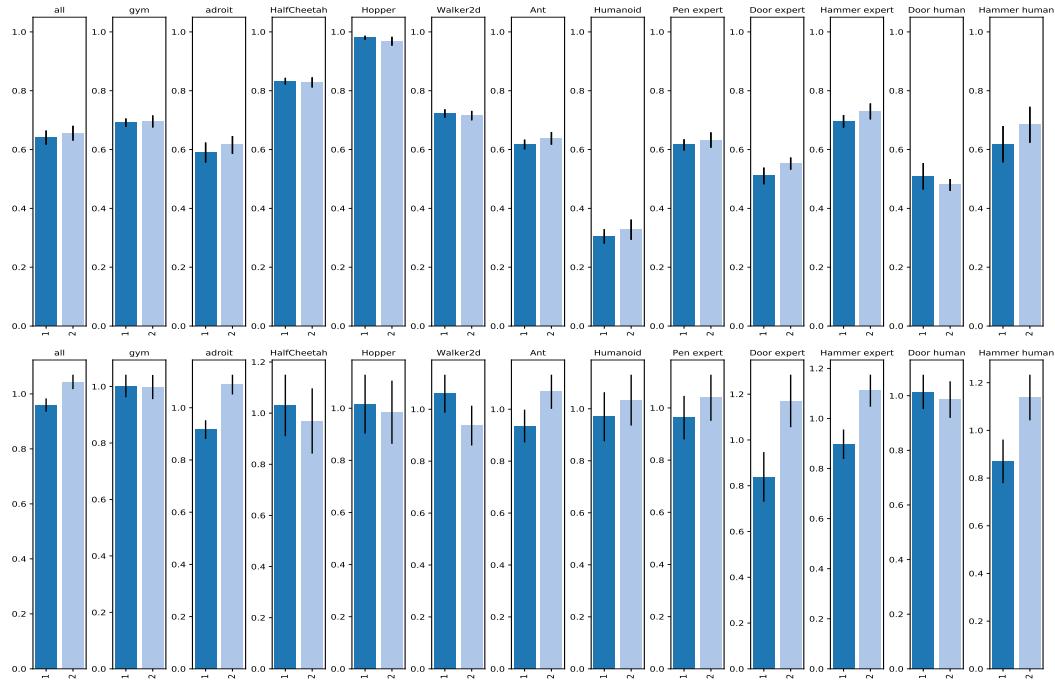


Figure 81: Analysis of choice discriminator to RL updates ratio (C44): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

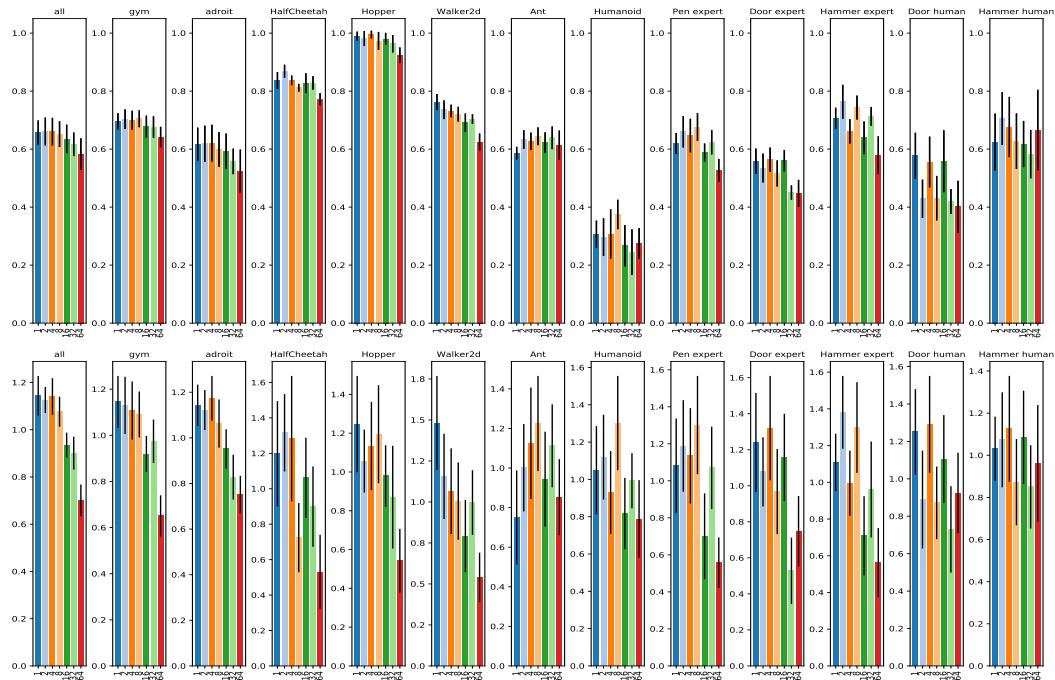


Figure 82: Analysis of choice number of combined batches (C56): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom).

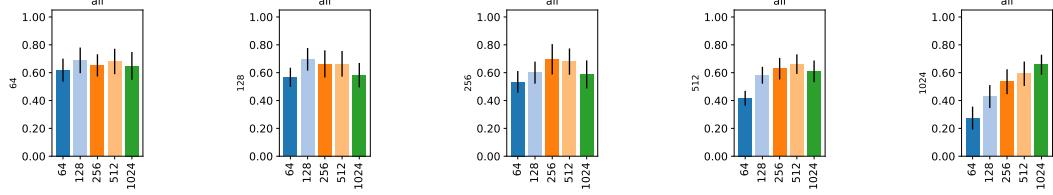


Figure 83: 95th percentile of performance scores conditioned on batch size (C7)(subplots) and replay ratio (C27)(bars).

I Additional experiments

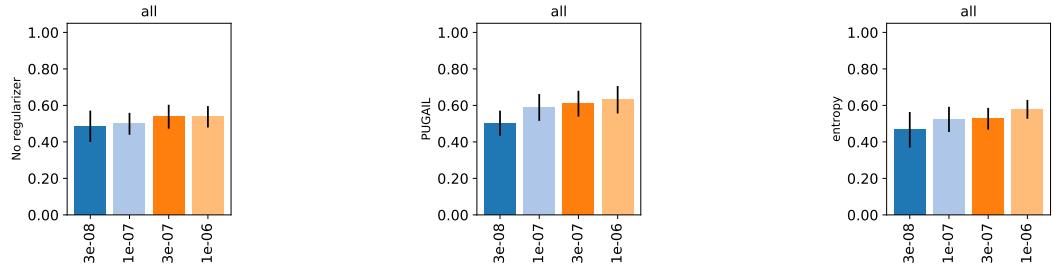


Figure 84: 95th percentile of performance scores conditioned on discriminator regularizer (C45)(rows) and discriminator learning rate (C42)(bars). The data comes from an experiment similar to the main one but with smaller values of discriminator learning rate (C42).

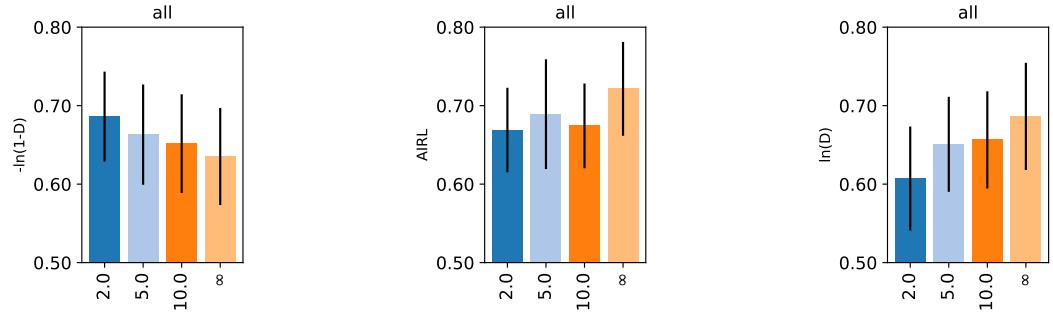


Figure 85: 95th percentile of performance scores conditioned on reward function (C30)(subplots) and max reward magnitude (C31)(bars). The data comes from an experiment similar to the main one but with max reward magnitude (C31) swept.

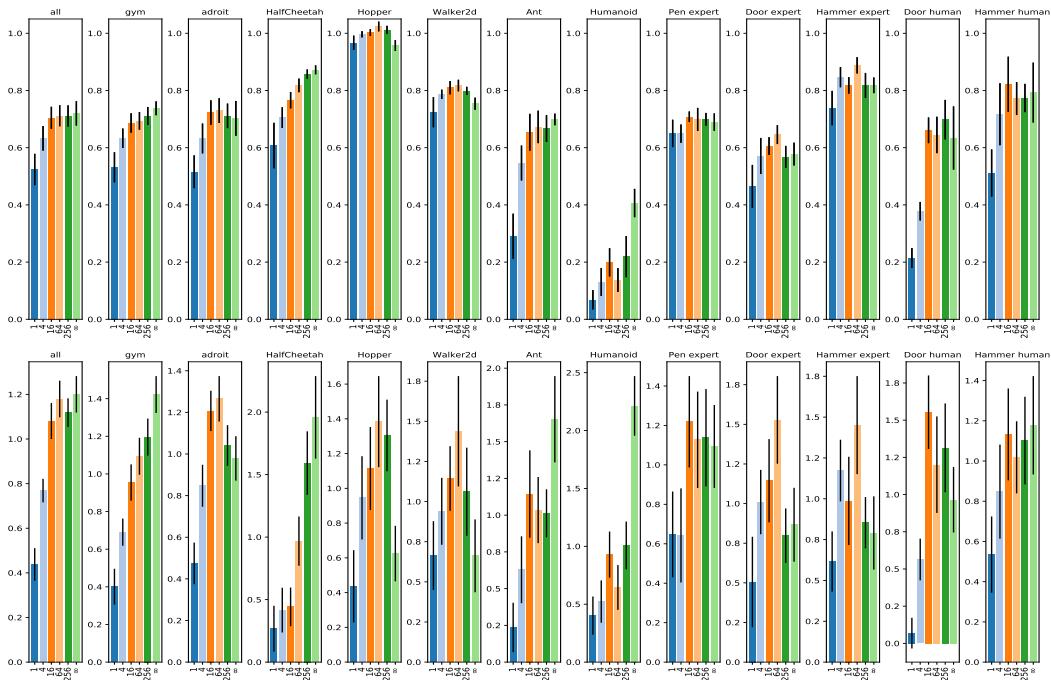


Figure 86: Analysis of choice policy-to-expert replay ratio (C33): 95th percentile of performance scores conditioned on choice (top) and distribution of choices in top 5% of configurations (bottom). policy-to-expert replay ratio (C33)=None means that expert transitions are not replayed. The data comes from an experiment similar to the main one but in which we also sweep policy-to-expert replay ratio (C33). All other experiments do not replay expert data.