
WHEN NOVICES TEACH BETTER: IMPROVING BEHAVIORAL CLONING WITH LOW-SKILL DATA

Anonymous authors

Paper under double-blind review

ABSTRACT

Behavioral cloning (BC), which trains models to replicate behavior from offline demonstrations, is a common approach in reinforcement learning. Several prior works argue that BC requires expert demonstrations and performs poorly when trained on low-skill or suboptimal data. We challenge this assumption by showing that, in certain regimes, training on low-skill demonstrations can yield models that outperform those trained on high-skill data. Since expert data is often costly and scarce, while low-skill data is cheaper and more abundant, this finding has important practical implications. To explain the result, we introduce a measure that quantifies the *resilience* of a policy—its ability to maintain reward under random perturbations—and show that resilience aligns with observed performance differences. Building on this insight, we introduce a new skill-based training curricula—structuring the training process according to policy skill levels—and show consistently improved BC performance compared to treating all data uniformly or filtering for experts. We show our findings in a custom synthetic environment and MuJoCo and validate using human data from Chess and Racing, showing consistency across domains.

1 INTRODUCTION

Behavioral cloning (BC) is a form of imitation learning that trains models to replicate behaviors from offline demonstrations (Bain and Sammut, 1995). It is widely used in reinforcement learning, particularly when explicit reward signals are unavailable or when online interaction is impractical due to safety or cost concerns (Zare et al., 2024). In some settings, BC can rival or even surpass offline RL methods (Kumar et al., 2022b; Qin et al., 2022). When trained on human data, it produces models that exhibit human-like behavior—a desirable property for AI systems that interact with people. It is commonly believed, however, that BC performs best with expert, high-skill demonstrations (Kurin et al., 2017) and performs poorly with low-skill demonstrations (Kumar et al., 2022c). As a result, practitioners often adopt Filtered BC (Mandlekar et al., 2021), discarding all but the highest-skill demonstrations. While effective when expert data is plentiful, this approach is limiting in practice, where datasets often span a wide range of skill levels and expert demonstrations are costly or scarce.

This paper revisits the common belief that BC works best with only expert data. We demonstrate that in certain types of low-data regimes, training on low-skill demonstrations can actually outperform training on expert data. Since expert data is often expensive and limited while low-skill data is cheaper and more abundant, this finding has significant practical implications. To explain this counterintuitive result, we introduce a central concept—*resilience*, which quantifies the ability of a policy (and the corresponding demonstrations generated by the policy) to withstand errors without severe degradation in performance. Our findings reveal that high-skill demonstrations often enter states where errors are disproportionately costly, whereas low-skill demonstrations more often avoid such states, resulting in more reliable models which suffer fewer errors.

Insights from psychology and education help motivate this perspective. Experts often rely on tacit knowledge and omit intermediate steps, making their behavior harder to model accurately (Nathan et al., 2001). High-skill strategies may operate closer to failure margins, amplifying the consequences of small modeling errors (Simon, 1997), and compounding errors in BC can cause learned policies to drift into unstable states (Ross et al., 2011). A racing game illustrates the point: expert drivers maximize speed and hug the track’s edges, achieving the fastest lap times but leaving little margin

054 for error, while novice drivers stay near the center and avoid fragile states. With unlimited training
055 data, BC can reproduce expert strategies, but with limited data, errors in expert demonstrations often
056 occur in fragile states, leading to catastrophic outcomes, whereas errors in novice data are more easily
057 absorbed. These perspectives suggest that the usefulness of a demonstration depends not only on
058 skill but also on the resilience of the demonstrating policy, motivating our study of whether and when
059 low-skill demonstrations can yield more reliable BC models than high-skill data and whether this
060 leads to useful insights in overall BC training.

061 More concretely, in this paper we investigate the role and impact of low-skill data in behavioral
062 cloning across three domains. We use expected reward as our measure of skill: policies with
063 higher rewards are considered high-skill, and those with lower rewards are low-skill. We begin
064 with controlled synthetic environments, where we manipulate rewards in a tree-based domain and
065 simulate policies across different skill levels. This setup allows us to systematically vary environment
066 characteristics and directly test how skill interacts with resilience. **We then replicate these results
067 across 12 tasks in MuJoCo.** We then turn to Chess, where public datasets provide billions of human
068 games across skill levels. Finally, we study Racing, using human demonstrations obtained through
069 collaboration with a simulator platform, allowing us to evaluate whether our findings generalize to a
070 complex real-world setting. Taken together, these three domains provide complementary perspectives
071 that reveal when and why low-skill data can be more useful than expert data.

072 Our experiments show that low-skill demonstrations are more resilient, BC models trained on the
073 low-skill data achieve higher performance at small budgets, and this effect disappears when high-skill
074 demonstrations are more resilient. These observations hold consistently across all **four** domains.

075 Motivated by these results, we also introduce a new method for training effectively with mixed-skill
076 datasets. Common strategies—such as filtering exclusively for expert demonstrations or training
077 uniformly over all data—do not fully leverage the diversity of available data. We instead evaluate
078 skill-based curricula that organize training by demonstrator skill level. Across both synthetic and
079 real-world domains, these curricula consistently improve behavioral cloning performance. While
080 secondary to our main contribution, these results suggest practical strategies for making better use of
081 demonstrations at varying skill levels.

082 **Contributions.** Overall, this work revisits the common belief that BC works best when trained only
083 on expert data and makes three contributions:

- 084 • We provide the counterintuitive empirical finding that in low-data regimes, BC trained on low-
085 skill demonstrations can outperform models trained on expert data, and show this effect in both
086 synthetic and real-world domains
- 087 • We introduce resilience as a metric for demonstration datasets, and show that resilience explains
088 when this phenomenon occurs.
- 089 • **To showcase the usefulness of this result, we apply the standard curriculum-learning framework**
090 **and demonstrate that leveraging low-skill data leads to consistent training improvements across**
091 **synthetic and real-world domains.**

092 093 094 2 RELATED WORK

095
096 **RL with Human Demonstrations.** A substantial body of work addresses reinforcement learning
097 using datasets of human demonstrations. Imitation learning trains models to replicate human policies
098 directly from demonstrations (Hussein et al., 2017). The simplest form, Behavioral Cloning (BC),
099 frames the task as supervised learning: given a state, the model predicts the demonstrated action (Bain
100 and Sammut, 1995). BC is widely applied in domains such as autonomous driving (Wang et al.,
101 2022), robotics (Florence et al., 2022), and gaming (Ruoss et al., 2024), since it avoids explicit
102 reward modeling or transition simulation. It can even outperform reward-based offline RL in some
103 settings (Kurenkov and Kolesnikov, 2022; Kumar et al., 2022b; Qin et al., 2022; Mandekar et al.,
104 2021). Alternatives to BC have also been developed. Offline RL incorporates reward signals during
105 training (Levine et al., 2020), which can improve task performance (Kumar et al., 2022c). Inverse
106 Reinforcement Learning aims to infer the underlying reward function of human behavior (Arora and
107 Doshi, 2021). Finally, online RL approaches can leverage demonstrations, especially when access to
a human oracle is available (Wagenmaker and Pacchiano, 2023).

Leveraging Low-Skill Human Demonstrations. We study whether low-skill human demonstrations can improve the training of machine learning models. In supervised learning, this idea has been widely explored in crowdsourcing (Howe et al., 2006; Vaughan, 2018). By recruiting laypeople to provide relatively low-skill annotations, large-scale datasets can be built without relying on costly expert input. This strategy has been highly successful, with prominent examples such as ImageNet (Deng et al., 2009) driving major progress in supervised learning.

In contrast, in the context of BC, prior work has shown that training on noisy or suboptimal demonstrations often degrades performance in BC (Fu et al., 2020; Mandelkar et al., 2021; Kanervisto et al., 2020), with follow-up theory reinforcing the view that BC typically requires expert data (Kumar et al., 2022c). To overcome this, extensions to BC have been proposed: filtered BC discards low-skill demonstrations using post-processing classifiers, improving performance in robotics (Chen et al., 2020; Wang et al., 2023), (Chen et al., 2025; Hejna et al., 2025), gaming (Kurin et al., 2017; Ruoss et al., 2024), and language modeling (Wang et al., 2024); weighted BC assigns importance weights to demonstrations (Ghosh et al., 2024); and conditional BC augments states with contextual variables such as modified rewards (Chen et al., 2021b) or goals (Ghosh et al., 2019). Different from this line of work, we focus on investigating whether and when low-skill data might improve BC training.

More broadly, reinforcement learning research has examined training with low-skill or mixed-skill human data, often using BC as a benchmark for evaluating alternative methods (Fu et al., 2020; Mandelkar et al., 2021; Kanervisto et al., 2020; Guss et al., 2019). Other studies have explored training with low-skill non-human data, typically generated by perturbing expert demonstrations with noise or using partially trained models (Brown et al., 2019; Chen et al., 2021a; Zhao et al., 2023a).

Critical States. The resilience metric in our work is related to the notion of critical states in BC and RL—the idea that certain regions of the state space are riskier or more important than others. For instance, Kumar et al. (2022a) define critical states as a subset of states from which taking a wrong action is especially costly to recover. Similarly, Zhang et al. (2020) propose an RL approach that involves first deploying a model in a less critical source environment (e.g., a traffic simulator) to learn representations of especially dangerous states, which can then be avoided once deployed in real-world environments. While related, our resilience concept differs in focus: it characterizes the *policy*, whereas the notion of critical states characterizes the *states*.

Curriculum Learning. Our approach builds on the idea of structuring training around demonstrations of varying skill levels. This is inspired by curriculum learning (Bengio et al., 2009), a technique modeled after human and animal pedagogy that organizes the training process to improve learning outcomes. Most research in curriculum learning structures training by gradually increasing difficulty, such as progressing from easy to hard examples. Curricula have also been designed around tasks, model architectures, or evaluation metrics (Soviany et al., 2022). Although first applied to supervised learning, curriculum learning has since gained traction in reinforcement learning as well (Narvekar et al., 2020). In our method, we structure the data by demonstrator skill level, beginning training with low-skill demonstrations and progressively incorporating higher-skill demonstrations.

3 SETTING AND METHODOLOGY

3.1 PROBLEM SETTING

Markov decision process (MDP). Behavioral cloning is formulated in the context of an MDP, defined by the tuple $\{\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma\}$, where \mathcal{S} is the state space, $\mathcal{A}(s)$ is the action space in state s , $\mathcal{P}(s, a)$ represents the state transition for taking action a in state s , $\mathcal{R}(s, a)$ specifies the reward for taking action a in state s , and γ is the discounting factor. In the remaining of the paper, our focus is on finite-time MDP and has set $\gamma = 1$, though our discussion is general to cases with $\gamma < 1$.

Policy and rewards. A policy is denoted by $\pi(a|s)$, indicating the probability for taking action a in state s . The policy reward, the cumulated rewards for deploying policy π in the MDP, is defined as $J(\pi) = \mathbb{E}[\sum_t R(s_t^\pi, a_t^\pi)]$, where s_t^π is the random variable denoting the state at t if the agent follows policy π and a_t^π is the random variable denoting the agent action in state s_t^π if the agent follows policy π . The expectation is over the randomness of the state transition and the policy. We can similarly define the value function $V^\pi(s) = \mathbb{E}[\sum_t R(s_t^\pi, a_t^\pi) | s_0 = s]$ as the reward of following policy π when starting with state s .

Skill levels. The skill level of a policy π is defined by the expected reward of the policy $J(\pi)$. We focus on situations when there exist policy functions at a range of skill levels. We denote these policies as $\{\pi_1, \pi_2, \dots, \pi_K\}$. Without loss of generality, we assume $J(\pi_1) \leq \dots \leq J(\pi_K)$. For low-skill (high-skill) data or demonstrations, we refer to the state-action pairs (s, a) generated by the policy with relatively low (high) rewards among the set of all policies. The total set of state-action pairs generated by the policy function π_k comprises a dataset \mathcal{D}_k .

Behavioral cloning. Behavioral cloning learns a policy function using supervised learning: Given demonstrations from \mathcal{D} , we train a model M to map each state to its corresponding action.

3.2 RESILIENCE

We define the *resilience* measure as the expected reward of a policy when a single action in the rollout (i.e., a trajectory generated by following the policy) is replaced with a random action, while all other actions follow the policy. Intuitively, resilience is a useful metric because BC models trained on finite datasets inevitably make errors, particularly under small data budgets. The resilience metric captures the extent to which a policy remains robust to these errors.

More formally, let $\zeta^\pi = \{(s_t, a_t, r_t)\}_{t=1, \dots, T}$ be the rollout of the policy π and (s_t, a_t, r_t) denote the (state, action, reward) at step t in the rollout. Also let $V^\pi(s)$ be the value function of state s if the agent takes policy π after state s , we can define resilience as follows.

$$\text{Res}(\pi) = \mathbb{E}_{\zeta^\pi} \left[\mathbb{E}_{k \sim U[1, T]} \left[\sum_{t=1}^{k-1} r_t + \mathbb{E}_{(a, s') \sim (\pi_{rand}, P(\cdot | s_k, a))} [R(s_k, a) + V^\pi(s')] \right] \right]$$

This measure requires knowledge of the underlying policy π , which is often unknown. Below, we define an empirical approximation that uses the demonstration dataset \mathcal{D}_π generated by π .¹

$$\overline{\text{Res}}(\mathcal{D}_\pi) = \frac{1}{|\mathcal{D}_\pi|} \sum_{\zeta \in \mathcal{D}_\pi} \frac{1}{T} \sum_{k=1}^T \left(\sum_{t=1}^{k-1} r_t + \frac{1}{|\mathcal{A}|} \sum_a \sum_{s'} P(s' | s_k, a) (R(s_k, a) + V(s')) \right) \quad (1)$$

3.2.1 CONNECTION BETWEEN RESILIENCE AND BC PERFORMANCE

Let $\hat{\pi}$ be the BC policy trained on demonstrations from π , and let ϵ be its per-timestep imitation error (e.g., probability of deviating from the demonstrator action). Using Theorem 2.1 from Ross et al. (2011) and adapting from a 0-1 loss formulation to a more general case, we obtain $J(\hat{\pi}) \approx J(\pi) - \epsilon T^2 C_\pi$, where $C_\pi = J(\pi) - \text{Res}(\pi)$ is the cost of making one error and following π afterwards². Thus, we get:

$$J(\hat{\pi}) \approx J(\pi) - \epsilon T^2 (J(\pi) - \text{Res}(\pi)) = (1 - \epsilon T^2) J(\pi) + \epsilon T^2 \text{Res}(\pi). \quad (2)$$

Implications on low-skill BC. Consider low- and high-skill demonstrators π_L, π_H with BC learners $\hat{\pi}_L, \hat{\pi}_H$. Applying Eq. (2) (under a shared- ϵ approximation) gives

$$J(\hat{\pi}_L) - J(\hat{\pi}_H) \approx (1 - \epsilon T)(J(\pi_L) - J(\pi_H)) + \epsilon T(\text{Res}(\pi_L) - \text{Res}(\pi_H)). \quad (3)$$

Since ϵ typically decreases as the demonstration budget increases, Eq. (3) predicts two regimes: (i) for high budgets (small ϵ), the skill term dominates and higher-skill demonstrations yield better BC; (ii) for low budgets (large ϵ), more resilient demonstrations can yield better BC even if they are lower-skill. In particular, if $\text{Res}(\pi_L) > \text{Res}(\pi_H)$, then there can exist a sufficiently large ϵT^2 regime where $J(\hat{\pi}_L) > J(\hat{\pi}_H)$.

¹This approximation requires a value function estimator, which could be obtained by training a value function on \mathcal{D}_π or leveraging known domain value functions (e.g., Stockfish in Chess)

²This approximation assumes that BC errors follow a uniform distribution – that $\hat{\pi}$ takes the correct action with $1 - \epsilon$ probability, and takes uniformly random actions with probability ϵ , aligning with our definition of resilience. In practice, BC errors might not be uniform, and accounting for more realistic BC error in our resilience metric would be an interesting future direction.

216 3.3 METHODOLOGY

217
218 In this work, we aim to better understand the effect of skill and resilience on behavioral cloning, and
219 how these insights can be leveraged to improve BC. Specifically, we first examine whether training
220 on low-skill data can outperform training on high-skill data, and whether resilience can explain this
221 phenomenon. We then explore how low-skill data can be leveraged to improve BC in practice.

222 3.3.1 USING RESILIENCE TO EXAMINE LOW-SKILL TRAINING IN BC

223
224 We begin by studying whether there are situations where training on low-skill data yields better models
225 in BC than training on high-skill data, and whether resilience can help explain this phenomenon. We
226 consider domains with policies π_1, \dots, π_K at varying skill levels, each generating a corresponding
227 demonstration dataset $\mathcal{D}_1, \dots, \mathcal{D}_K$. We assume the datasets are given and that the skill level of the
228 generating policy is known. Our goal is to examine whether, under different data sizes, training on a
229 low-skill dataset can yield better performance than training on a higher-skill dataset.

230 We use the empirical resilience measure in Equation 1 to compute the resilience of demonstrations
231 at different skill levels. In particular, we are interested in ordering resilience scores across skills to
232 identify positive or negative trends between resilience and skill. We denote environments where
233 high-skill agents are more resilient than low-skill agents as **Resilient High-Skill**, and those where
234 high-skill agents are less resilient than low-skill agents as **Fragile High-Skill**.

235 Given this classification, we then run BC on each dataset \mathcal{D}_k , consisting of B demonstrations
236 generated by π_k . Each dataset is split into training and validation sets. We instantiate a model M
237 and train it to minimize loss on the training data, repeating this process for multiple epochs until
238 validation loss converges. Finally, we evaluate the task performance of the trained models at each
239 budget and skill level.

240 3.3.2 TRAINING WITH SKILL-BASED CURRICULA

241
242 Next, we investigate how to leverage low-skill data in the common setting where a mixed-skill dataset
243 $\mathcal{D} = \{\mathcal{D}_1, \dots, \mathcal{D}_K\}$ has already been collected. We assume that, given a budget B , each dataset \mathcal{D}_k
244 contains $\frac{B}{K}$ demonstrations from π_k . Our task is to choose a training schedule that maximizes the
245 final model’s performance, especially under the common scenario where marginal compute costs
246 are much cheaper than data acquisition costs. We consider three possible schedules: two baselines
247 from prior literature and one proposed curriculum that orders training by skill level. Our goal is not
248 to prove that this proposed schedule is the optimal one, but simply show that it is possible to use
249 low-skill data to improve training scheduling in the first place.

- 251 • **Standard BC:** In this method, the model is trained using supervised learning to mimic the full
252 training dataset $\{\mathcal{D}_1, \dots, \mathcal{D}_K\}$ until convergence.
- 253 • **Filtered BC:** The dataset is first preprocessed to include only the data at the highest skill $\{\mathcal{D}_K\}$.
254 The model is then trained using the same BC methodology on the thresholded dataset.
- 255 • **Filtered BC-LSC³:** We first train a model on the full dataset $\{\mathcal{D}_1, \dots, \mathcal{D}_K\}$ until convergence
256 (Standard BC). We then iteratively remove the lowest-skill data and retrain on the remaining
257 datasets, ending with training only on \mathcal{D}_K until convergence (Filtered BC).

258
259 To assess the effect of different training curricula, we train a model under each schedule to obtain
260 three models and evaluate their task performance in the domain.

261 3.4 DOMAINS

262 3.4.1 TREE

263
264 We design a simplified class of tree environments to investigate the effect of data skill-level and
265 resilience on BC model performance.⁴ Each environment we construct is characterized by a depth-5

266 ³LSC stands for Low-Skill Curriculum

267 ⁴Our experiment code, including both tree and chess domains is available at the following anonymized link:
268 <https://anonymous.4open.science/r/iclr2026-C020>

270 binary tree structure, where each node in the tree represents a state. For each non-terminal node, the
271 agent’s action space consists of moving to either the left or right child node. The state transitions are
272 deterministic. Each node is associated with a reward, and the agent collects the reward from each
273 node it passes through.

274 To generate demonstrations at different skill levels, we model an agent’s skill based on how many
275 steps it can look ahead. Specifically, an agent with skill level 1 always chooses the direct child node
276 with the larger reward. An agent with skill level k evaluates all possible combinations of the next k
277 moves and selects the next move that leads to the optimal path over those k moves.

278 In order to better understand the relationship between resilience and BC performance, we specifically
279 design a set of environments where high-skill agents have lower resilience. Specifically, each
280 environment is randomly generated with different rewards assigned to the nodes. For each node, we
281 begin by drawing a reward value uniformly at random between 0 and 1. We then randomly designate
282 a subset of the nodes as *fragile*, with the proportion of fragile nodes varying based on the experiment.
283 For these nodes, we decrease their rewards by 8, and randomly select one of their child nodes to
284 increase its reward by 10. The intuition of this design is to provide a connection between the behavior
285 of low-skill and high-skill agents to our metric of *resilience*.

286 Using this definition of skill, we observe that low-skill and high-skill agents follow different strategies
287 in an environment. The lowest-skill agent will always avoid fragile nodes if possible, to avoid paying
288 the entry penalty. In contrast, the highest-skill agent will explicitly seek out these states, recognizing
289 that the additional reward obtained from them improves its overall score. However, this behavior
290 implies that high-skill agents are less resilient to fragile nodes than low-skill nodes.

291 Therefore, environments with a high proportion of fragile nodes can be classified as fragile high-skill,
292 while environments without any fragile nodes can be classified as resilient high-skill. We can tune
293 this parameter to understand the effect of resilience on BC performance. We provide more details
294 about the domain construction, as well as an example tree diagram in Appendix A.1.

296 3.4.2 MUJoCo

298 **MuJoCo is a physics engine used for simulating robotics and biomechanics and contains 11 different**
299 **environments for training agents. (Todorov et al., 2012). We use an existing set of pretrained models**
300 **which span several combinations of environment, architecture, and skill (Younis et al., 2024). In our**
301 **results, we evaluate on every environment/architecture combination with at least two models (except**
302 **for SAC-HumanoidStandup, which failed to run). We denote the highest skill model as "high-skill"**
303 **and the lowest-skill model as "low-skill". Overall, this gives us 12 additional tasks to evaluate.**

305 3.4.3 CHESS

306 The previous two examples allow us to examine training behavior in a simplified toy scenario, and a
307 more complex set synthetic environments. Next, we aim to show that these results hold in a real-world
308 setting, and first pick the domain of Chess. Chess is a useful domain to analyze because of the huge
309 human datasets available for analysis, a simple and well-defined MDP, and a highly tuned metric
310 of policy skill. We conduct our experiments using two datasets. First, we use the Lichess Player
311 Database (Lichess, 2025), which contains 6.2 billion human chess games at skill levels ranging from
312 beginners to grandmasters. Second, we construct a dataset of Stockfish games (the best chess engine
313 available), representing skill level as nodes searched. For all experiments in this domain, we use
314 the BC training framework provided by McIlroy-Young et al. (2020a). More details on the model
315 training and Stockfish game generation processes is available in Appendix C.

317 3.4.4 RACING

318 We also extend our exploration of low-skill BC to the domain of Racing. Through an academic
319 partnership with a commercial racing simulator platform, we obtained access to human driving
320 trajectories from a Formula 4 racing series. The dataset contains 10,000 complete laps on each of
321 four tracks: Summit Point Raceway, Circuito de Navarra, Autodromo Nazionale Monza, and Road
322 America (a diagram of each map is available in Appendix D). We approximate skill by lap time, with
323 the fastest 5,000 laps representing high-skill data and the slowest 5,000 laps representing low-skill
data. While our limited access to the platform prevents us from directly computing resilience scores,

our intuition strongly suggests that low-skill demonstrations may be more resilient, due to the nature of high-skill drivers to drive faster and cut corners close to the edge of the track.⁵

4 RESULTS

4.1 TREE

Following the methodology from Section 3.3.1, we train BC models on datasets generated by agents of each skill level from $k = 1$ to 4, using data budgets ranging from 2^8 to 2^{18} data points. More details on the BC training process is available in Appendix A.1.

After training, we evaluate each model by simulating 16,384 randomly generated episodes and report the average task performance. We repeat this process with proportion of fragile states ranging from 0 to 1. For each environment configuration, we train models with 200 random seeds and present the mean and standard error.

In particular, we highlight two state configurations - with 20% of tree nodes as fragile and with 0%, to better explore the effects of these states on low-skill BC. We plot the task performance of low-skill and high-skill BC at a variety of budgets for both experiment configurations in Figure 1.

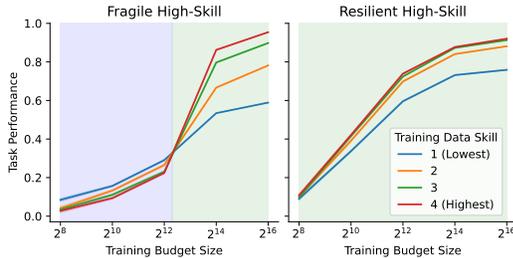


Figure 1: Effect of training BC models across budgets and skill levels in the tree domain without and with fragile nodes. Regions where low-skill training outperforms high-skill are highlighted in blue.

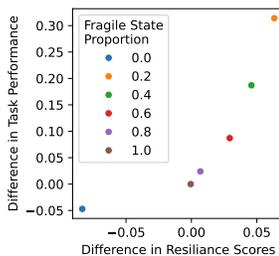


Figure 2: Comparing the difference in resilience and small-budget task performance by tree environment.

Overall, we see that in the fragile environment configurations, models trained on low-skill data outperform models trained on high-skill data at small budgets (under 2^{14}), while at larger budgets, training on high-skill data leads to improved task performance. However in the environment without fragile states, training on high-skill demonstrations improves performance at all data budgets. **These differences are significant with $p < .001$.**

We train each model until convergence, meaning that different models may use different amounts of compute. This is a reasonable assumption in domains where data collection is much more expensive than compute. However, one might ask if low-skill models are using more compute in order to get their higher performance. In Appendix A.1, we show that this is not true, and that low-skill

models actually take fewer epochs to converge.

To illustrate the potential connection between *resilience* and the benefit of low-skill training, we plot the differences in agent resilience and the difference in final model performance at low budgets (2^{10}) across environment configurations. We estimate agent resilience using Equation 1, and plot the difference in resilience and difference in task performance between the lowest and highest-skill agent for each environment configuration in Figure 2.

We see that a clear trend between the difference in resilience scores between the low and high-skill agents and BC model performance. When the low-skill agent has a higher resilience score, the corresponding low-skill BC model has a higher task performance during deployment.

4.2 MUJoCo

In each MuJoCo environment, we train BC models on the low-skill agent and the high-skill agent. For both models, we follow Equation 1 and compute the resilience score by rolling out the original policy for up to 1000 steps, and measure the accumulated reward over 5 timesteps after perturbations. We present the results in Table 1. Overall, we see that in 5 domains, the low-skill agent demonstrates higher resilience, and in the other 7 domains, the high-skill agent demonstrates more resilience.

⁵We are unable to release the dataset and code for this application due to our agreement with the provider.

Next, we compute the task performance of models trained on low-skill and high-skill data. Following the methodology from Section 3.3.1, we generate datasets consisting of 1000-step rollouts by each agent. We then train models at budgets ranging from 2^2 to 2^{10} until validation convergence. Finally, we evaluate the models by performing 10 rollouts, and reporting the cumulative reward attained. We repeat this procedure for 10 seeds, and present the results in the table above.

In Table 1, we report a summary of results in all 12 domains. The full set of training plots, measuring model performance at multiple training budgets is available in Appendix B.

In all five environments with low-skill resilience, we find that low-skill models indeed outperform high-skill models consistently at low budgets. Additionally, we find that in five of seven domains with high-skill resilience, low-skill models consistently outperform high-skill models at low budgets.

Overall, we find that the presence of low-skill resilience has 100% accuracy in predicting low-skill outperformance. Additionally, when low-skill models are less resilient, they generally fail to outperform high-skill models. The presence of the two exceptions (Pusher and Walker2d) does not invalidate our metric, but suggest that our approximate resilience calculation may not be accurate enough in these domains.

Domain	High-Skill Resilience	Low-Skill Resilience	Low-skill more resilient?	Low-Skill performs better?	Aligns
SAC-HalfCheetah	-0.32917	-0.211631	✓	✓	✓
SAC-Hopper	0.864271	0.906312	✓	✓	✓
SAC-Humanoid	4.807381	4.846941	✓	✓	✓
TQC-HalfCheetah	-0.242331	-0.231321	✓	✓	✓
TQC-Humanoid	4.827656	4.848168	✓	✓	✓
PPO-Swimmer	0.04955	0.044526	×	×	✓
SAC-Ant	-0.06562	-0.345851	×	×	✓
SAC-InvertedDbIPend	7.49	7.47	×	×	✓
SAC-InvPendulum	0.678	0.674	×	×	✓
SAC-Pusher	-0.41236	-0.420937	×	✓	×
SAC-Reacher	-0.14926	-0.149267	×	×	✓
SAC-Walker2d	0.904928	0.812948	×	✓	×

Table 1: Resilience scores and model performance summary across 12 architectures and environments in MuJoCo.

4.3 CHESS

We train two sets of models - human-like models, and engine-like models. For our human-like models, we select five skill levels from the Lichess dataset - moves by players rated at 800 (low), 1200, 1600, 2000, and 2400 (high). For our engine-like models, we select three skill levels, nodes=100 (low), nodes=1000, and nodes=10000 (high). Based on experiments conducted by Marco (2021), these three skill levels of Stockfish roughly correspond to human ratings of 1200, 1700, and 2500. In the main paper, we only focus on the highest and lowest skill levels, but show that the results are qualitatively similar for all skill levels in Appendix C.1.

First, we compute the resilience scores of the human and engine agents at each skill level. We sample 4096 actions from the training datasets at each skill level and use Equation 1 with the highest-depth Stockfish as a value function to compute resilience. We plot these scores in Table 2. In Appendix C.2, we show that these resilience scores are robust to weaker approximate value function and at multiple skill levels.

Skill	Human	Engine
Low	0.305	0.247
High	0.251	0.298

Table 2: Resilience scores for humans and engines for the lowest and highest skill levels.

We see that in the human dataset, resilience values decrease as skill level increases, with the highest skill level having the least resilience. On the other hand, in the engine dataset, resilience values increase as skill level increases, with the highest skill level having the most resilience.⁶ Thus, we denote the environment with human data as Fragile High-Skill, and the environment with engine data as Resilient High-Skill.

To validate our tree-domain results, we again follow the methodology from Section 3.3 and train BC models on datasets containing demonstrations at each skill level, with data budgets ranging from 2^{16} to 2^{24} data points. We evaluate

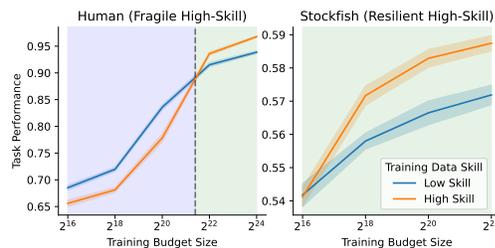


Figure 3: Effect of training BC models on human and engine demonstrations in chess by skill level. Budget regions where low-skill training outperforms high-skill is highlighted in blue.

⁶In Appendix C, we show an expanded plot showing this trend is monotonic across more skill levels.

the resulting models by playing 10240 chess games against a randomly-playing opponent.⁷ We repeat this process with 20 random seeds and present the mean and standard error for each experiment configuration in Figure 3.

In the human-trained models, we observe that with small data budgets, under 2^{22} , training on low-skill data achieves the highest performance, while at larger data budgets leads to the best performance. However, with the engine-trained models, we observe that training on the highest-skill data leads to the best task performance at all data budgets. **These differences are significant with $p < .001$.**

Both results match our intuition and our previous results. In the environment where low-skill agents are more resilient, models trained on low-skill data outperform high-skill models at small budgets, but when this isn't true, high-skill models outperform low-skill models at all budgets.

4.4 RACING

Again following methodology from Section 3.3, we train BC models on low and high-skill data on each of the four tracks, at data budgets ranging from 500 to 1000 laps. We evaluate each model using a virtual joystick (pyv, 2025) that forwards the controls to the simulator and measure track completion rate over 5 attempts. We present the results in Table 3.

We see that at the smallest budget, low-skill models complete 60% of races, while high-skill models complete 20% of races. This difference is significant with $p < 0.01$. We see that this trend holds in three out of four tracks. For example, in Summit Point, with a budget of 500 laps, a model trained on low-skill data can always complete the track, whereas a model trained on high-skill data can only do so reliably with 1000 or more laps. The only exception is Road America, where the model trained on high-skill data outperforms even at low budgets. To explain this, we find in Road America, high-skill models qualitatively perform more resiliently, which again aligns with our explanation. We provide more details of this phenomenon in Appendix D.

Track Name	Budget	500	750	1000
Summit Point	Low	1.0	1.0	1.0
	High	0.0	0.8	1.0
Circuito de Navarra	Low	0.4	0.8	1.0
	High	0.0	0.4	1.0
Autodromo Nazionale Monza	Low	1.0	1.0	1.0
	High	0.6	1.0	1.0
Road America	Low	0.0	0.0	1.0
	High	0.2	1.0	1.0

Table 3: Effect of training BC models in a racing simulator by track, budget and skill. In the first three tracks, low-skill training outperforms high-skill training, while in the fourth track, high-skill training always outperforms low-skill training. Budget regions where low-skill training outperforms high-skill is highlighted in blue.

5 SKILL-BASED CURRICULA

Next, we conduct experiments to evaluate if applying a standard curriculum schedule over demonstrator skill can improve performance when mixed-skill demonstrations are available.

5.1 TREE

We construct a tree dataset with a mixture of demonstrations from agents at four skill levels. We implement both the baselines, Standard BC and Filtered BC, as well as our proposed Filtered BC-LSC approach. We then train models with data budgets ranging from 2^8 to 2^{20} data points. After training, we evaluate each model by simulating 16,384 random episodes and report the task performance. We repeat this process with 200 random seeds and present the mean and standard error for each configuration in Figure 4.

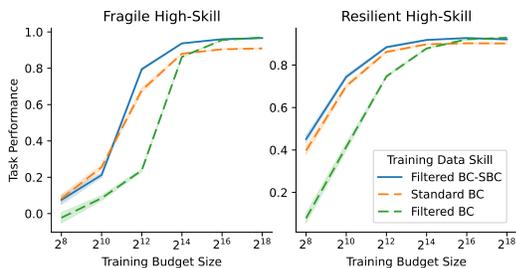


Figure 4: Effect of low-skill curricula on task performance in the tree domain, with environment configurations both containing and not containing fragile states.

⁷In Appendix C, we replicate our results with other evaluation approaches.

We observe a few key trends. In both environments, we see Standard BC outperform Filtered BC at small budgets (under 2^{14}), while Filtered BC performs better with larger budgets. This aligns with existing work which suggests Filtered BC is best when budgets are large enough to keep sufficient data after filtering. Second, we see that our proposed approach minimizes the downsides of both approaches at all budgets – Filtered BC-LSC matches or exceeds Standard BC at low budgets, while performing competitively at large budgets against Filtered BC. Overall, we see that simple curriculum strategy leveraging both low-skill and high-skill demonstrations allows us to perform well in both small-budget and large-budget regions.

5.2 CHESS

Following the methodology from Section 3.3.2, we train Standard BC and Filtered BC models, with and without a low-skill curriculum, with data budgets ranging from 2^{16} to 2^{24} , using the low- and high-skill agents in both the human and Stockfish datasets. We evaluate the resulting models by playing 10240 chess games against a randomly-playing opponent, and report the average win rate of the model. We repeat this process with 20 random seeds, and present the mean and standard error for each experiment configuration in Figure 5.

Overall, we observe similar trends to those seen with the tree domain results. Standard BC outperforms Filtered BC at lower budgets, while Filtered BC outperforms at higher budgets. Again, we see that Filtered BC-LSC acts as a combination of the two, performing as well as Standard BC at low budgets, while better taking advantage of high-skill data at high budgets. Overall, this gives us strong evidence that it is possible to leverage low-skill data to train BC models with high performance.

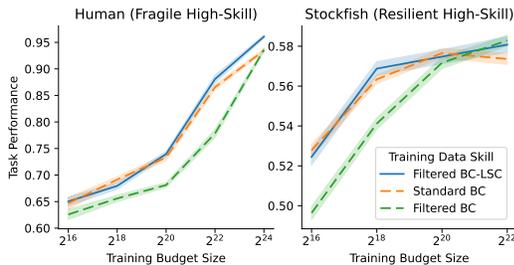


Figure 5: Effect of training BC models on human and engine demonstrations using skill-based curricula in chess by skill level.

6 CONCLUSION AND DISCUSSION

We demonstrate that low-skill demonstrations can be leveraged to improve the performance of BC models. Our findings show that BC trained on low-skill data can outperform BC trained on high-skill data in small-data regimes. This effect is strongly correlated with the *resilience* of the training datasets: it arises in domains where low-skill trajectories are more resilient to mistakes, but not in domains where this is not the case. We also show that a simple curriculum learning-like strategy that organizes data by skill level can improve BC performance compared to standard approaches.

These results lead to several useful takeaways. When practitioners have reason to believe that low-skill policies are more resilient, they should consider training on low-skill data. While this behavior is not guaranteed, our findings suggest that practitioners can benefit from testing whether it holds in their specific domain. In scenarios where only small, mixed-skill datasets are available, emphasizing low-skill data during training can be particularly valuable, especially when compute costs are much lower than data acquisition costs. By contrast, in settings with large data budgets, low-skill data is less helpful, and practitioners should instead focus on collecting the highest-quality demonstrations.

There are also several open directions for future work. **We showed that a simple curriculum schedule can lead to increased performance, but more complex schedules may show additional benefits, especially developed jointly with online algorithms that dynamically determine which data composition to acquire.** Second, methods for automatically estimating the complexity of human strategies in each domain would better guide practitioners in deciding which skill data to collect. Finally, it would be valuable to explore whether the benefits of training on low-skill data extend beyond BC to offline RL or other learning paradigms.

ETHICS STATEMENT

This work uses human-generated data from two primary sources: a public chess database and a private racing simulator dataset. The chess data is drawn from the Lichess Player Database, which consists of anonymized games that are publicly available for research. The racing data was obtained through an academic partnership with a commercial simulator platform. This dataset was handled in accordance with our data sharing agreement with the provider, and all data was anonymized to protect user privacy. The racing platform was not involved with the research direction or conclusions of this paper, and there are no conflicts of interest between the platform and this work.

Our research aims to create more robust and reliable models through behavioral cloning, and we do not foresee direct negative societal impacts. On the contrary, our finding that low-skill data can improve model resilience may contribute to the development of safer AI systems, as it encourages learning from a wider, more cautious range of behaviors rather than solely from "expert" demonstrations that may operate close to failure margins.

We have used large language models to assist with proofreading and polishing writing of this paper.

REPRODUCIBILITY STATEMENT

We have made every effort to ensure the reproducibility of our results. Our experiment code for both the synthetic tree and chess domains has been made available in an anonymized repository at <https://anonymous.4open.science/r/iclr2026-C020>. The core methodology, including the empirical definition of our resilience metric (Equation 1) and the design of our skill-based training curricula, is described in Section 3.3. For our experiments, detailed descriptions of the synthetic tree domain construction, agent generation, and model training procedures are provided in Section 4 and Appendix A. Similarly, implementation details for the chess experiments, including data sources, model architecture, and evaluation protocol, can be found in Section 5 and Appendix B. While the proprietary nature of the racing dataset used in Section 6 prevents its public release, we have detailed the experimental setup and provided qualitative analyses in Appendix C. We believe these resources provide the necessary details to reproduce our key findings. .

REFERENCES

- pyvjoy: Python bindings for vjoy, 2025. URL <https://github.com/tidzo/pyvjoy>.
- Saurabh Arora and Prashant Doshi. A survey of inverse reinforcement learning: Challenges, methods and progress. *Artificial Intelligence*, 297:103500, 2021.
- Michael Bain and Claude Sammut. A framework for behavioural cloning. In *Machine Intelligence 15*, pages 103–129, 1995.
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48, 2009.
- Daniel Brown, Wonjoon Goo, Prabhat Nagarajan, and Scott Niekum. Extrapolating beyond sub-optimal demonstrations via inverse reinforcement learning from observations. In *International conference on machine learning*, pages 783–792. PMLR, 2019.
- Annie S. Chen, Alec M. Lessing, Yuejiang Liu, and Chelsea Finn. Curating demonstrations using online experience, 2025. URL <https://arxiv.org/abs/2503.03707>.
- Letian Chen, Rohan Paleja, and Matthew Gombolay. Learning from suboptimal demonstration via self-supervised reward regression. In *Conference on robot learning*, pages 1262–1277. PMLR, 2021a.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021b.

594 Xinyue Chen, Zijian Zhou, Zheng Wang, Che Wang, Yanqiu Wu, and Keith Ross. Bail: Best-
595 action imitation learning for batch deep reinforcement learning. *Advances in Neural Information*
596 *Processing Systems*, 33:18353–18363, 2020.

597 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
598 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
599 pages 248–255. IEEE, 2009.

600
601 Pete Florence, Corey Lynch, Andy Zeng, Oscar A Ramirez, Ayzaan Wahid, Laura Downs, Adrian
602 Wong, Johnny Lee, Igor Mordatch, and Jonathan Tompson. Implicit behavioral cloning. In
603 *Conference on Robot Learning*, pages 158–168. PMLR, 2022.

604 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep
605 data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.

606
607 Dibya Ghosh, Abhishek Gupta, Ashwin Reddy, Justin Fu, Coline Devin, Benjamin Eysenbach,
608 and Sergey Levine. Learning to reach goals via iterated supervised learning. *arXiv preprint*
609 *arXiv:1912.06088*, 2019.

610 Udit Ghosh, Dripta S Raychaudhuri, Jiachen Li, Konstantinos Karydis, and Amit K Roy-Chowdhury.
611 Robust offline imitation learning from diverse auxiliary data. *arXiv preprint arXiv:2410.03626*,
612 2024.

613
614 William H Guss, Brandon Houghton, Nicholay Topin, Phillip Wang, Cayden Codel, Manuela Veloso,
615 and Ruslan Salakhutdinov. Minerl: A large-scale dataset of minecraft demonstrations. *arXiv*
616 *preprint arXiv:1907.13440*, 2019.

617
618 Joey Hejna, Suvir Mirchandani, Ashwin Balakrishna, Annie Xie, Ayzaan Wahid, Jonathan Tompson,
619 Pannag Sanketi, Dhruv Shah, Coline Devin, and Dorsa Sadigh. Robot data curation with mutual
620 information estimators, 2025. URL <https://arxiv.org/abs/2502.08623>.

621
622 Jeff Howe et al. The rise of crowdsourcing. *Wired magazine*, 14(6):176–183, 2006.

623
624 Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Chrisina Jayne. Imitation learning: A
625 survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35, 2017.

626
627 Anssi Kanervisto, Joonas Pussinen, and Ville Hautamäki. Benchmarking end-to-end behavioural
628 cloning on video games. In *2020 IEEE conference on games (CoG)*, pages 558–565. IEEE, 2020.

629
630 Aviral Kumar, Joey Hong, Anikait Singh, and Sergey Levine. When should we prefer offline
631 reinforcement learning over behavioral cloning? *ArXiv*, abs/2204.05618, 2022a. URL <https://api.semanticscholar.org/CorpusID:248118795>.

632
633 Aviral Kumar, Joey Hong, Anikait Singh, and Sergey Levine. When should we prefer offline
634 reinforcement learning over behavioral cloning? *arXiv preprint arXiv:2204.05618*, 2022b.

635
636 Aviral Kumar, Ilya Kostrikov, and Sergey Levine. Should i use offline rl or imitation learning? 2022c.
637 URL <https://bair.berkeley.edu/blog/2022/04/25/rl-or-bc/>.

638
639 Vladislav Kurenkov and Sergey Kolesnikov. Showing your offline reinforcement learning work:
640 Online evaluation budget matters. In *International Conference on Machine Learning*, pages
641 11729–11752. PMLR, 2022.

642
643 Vitaly Kurin, Sebastian Nowozin, Katja Hofmann, Lucas Beyer, and Bastian Leibe. The atari grand
644 challenge dataset. *arXiv preprint arXiv:1705.10998*, 2017.

645
646 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,
647 review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.

Lichess. Lichess database, 2025. URL [database.lichess.org](https://lichess.org).

Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni, Li Fei-
Fei, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in learning from offline
human demonstrations for robot manipulation. *arXiv preprint arXiv:2108.03298*, 2021.

648 Meloni Marco. Stockfish and lc0, test at different number of nodes,
649 2021. URL [https://www.melonimarco.it/en/2021/03/08/
650 stockfish-and-lc0-test-at-different-number-of-nodes/](https://www.melonimarco.it/en/2021/03/08/stockfish-and-lc0-test-at-different-number-of-nodes/).
651

652 Reid McIlroy-Young, Siddhartha Sen, Jon Kleinberg, and Ashton Anderson. Aligning superhuman
653 ai with human behavior: Chess as a model system. In *Proceedings of the 26th ACM SIGKDD
654 International Conference on Knowledge Discovery & Data Mining*, pages 1677–1687, 2020a.

655 Reid McIlroy-Young, Siddhartha Sen, Jon Kleinberg, and Ashton Anderson. Aligning superhuman
656 ai with human behavior: Chess as a model system. In *Proceedings of the 26th ACM SIGKDD
657 International Conference on Knowledge Discovery and Data Mining*, page 1677–1687, 2020b.
658

659 Sanmit Narvekar, Bei Peng, Matteo Leonetti, Jivko Sinapov, Matthew E Taylor, and Peter Stone.
660 Curriculum learning for reinforcement learning domains: A framework and survey. *Journal of
661 Machine Learning Research*, 21(181):1–50, 2020.

662 Mitchell J Nathan, Kenneth R Koedinger, Martha W Alibali, et al. Expert blind spot: When content
663 knowledge eclipses pedagogical content knowledge. In *Proceedings of the third international
664 conference on cognitive science*, volume 644648, pages 644–648, 2001.
665

666 Rong-Jun Qin, Xingyuan Zhang, Songyi Gao, Xiong-Hui Chen, Zewen Li, Weinan Zhang, and Yang
667 Yu. Neorl: A near real-world benchmark for offline reinforcement learning. *Advances in Neural
668 Information Processing Systems*, 35:24753–24765, 2022.

669 Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured
670 prediction to no-regret online learning. In *Proceedings of the fourteenth international conference
671 on artificial intelligence and statistics*, 2011.
672

673 Anian Ruoss, Grégoire Delétang, Sourabh Medapati, Jordi Grau-Moya, Li Kevin Wenliang, Elliot
674 Catt, John Reid, and Tim Genewein. Grandmaster-level chess without search. *arXiv preprint
675 arXiv:2402.04494*, 2024.

676 Herbert Alexander Simon. *Models of bounded rationality: Empirically grounded economic reason*,
677 volume 3. MIT press, 1997.

678 Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu Sebe. Curriculum learning: A survey.
679 *International Journal of Computer Vision*, 130(6):1526–1565, 2022.
680

681 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.
682 In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pages 5026–5033.
683 IEEE, 2012.

684 Jennifer Wortman Vaughan. Making better use of the crowd: How crowdsourcing can advance
685 machine learning research. *Journal of Machine Learning Research*, 18(193):1–46, 2018.
686

687 Andrew Wagenmaker and Aldo Pacchiano. Leveraging offline data in online reinforcement learning.
688 In *International Conference on Machine Learning*, pages 35300–35338. PMLR, 2023.
689

690 Lingguang Wang, Carlos Fernandez, and Christoph Stiller. High-level decision making for automated
691 highway driving via behavior cloning. *IEEE Transactions on Intelligent Vehicles*, 8(1):923–935,
692 2022.

693 Qiang Wang, Robert McCarthy, David Cordova Bulens, Francisco Roldan Sanchez, Kevin McGuin-
694 ness, Noel E O’Connor, and Stephen J Redmond. Identifying expert behavior in offline training
695 datasets improves behavioral cloning of robotic manipulation policies. *IEEE Robotics and Automa-
696 tion Letters*, 2023.

697 Ruiyi Wang, Haofei Yu, Wenxin Zhang, Zhengyang Qi, Maarten Sap, Graham Neubig, Yonatan Bisk,
698 and Hao Zhu. Sotopia- π : Interactive learning of socially intelligent language agents. *arXiv
699 preprint arXiv:2403.08715*, 2024.
700

701 Omar G. Younis, Rodrigo Perez-Vicente, John U. Balis, Will Dudley, Alex Davey, and Jordan K
Terry. Minari, September 2024. URL <https://doi.org/10.5281/zenodo.13767625>.

Maryam Zare, Parham M Kebria, Abbas Khosravi, and Saeid Nahavandi. A survey of imitation learning: Algorithms, recent developments, and challenges. *IEEE Transactions on Cybernetics*, 2024.

Jesse Zhang, Brian Cheung, Chelsea Finn, Sergey Levine, and Dinesh Jayaraman. Cautious adaptation for reinforcement learning in safety-critical settings. In *Proceedings of the 37th International Conference on Machine Learning, ICML'20*, 2020.

Tianxiang Zhao, Wenchao Yu, Suhang Wang, Lu Wang, Xiang Zhang, Yuncong Chen, Yanchi Liu, Wei Cheng, and Haifeng Chen. Skill disentanglement for imitation learning from suboptimal demonstrations. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 3513–3524, 2023a.

Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware, 2023b. URL <https://arxiv.org/abs/2304.13705>.

A TREE EXPERIMENTS

A.1 DOMAIN CONSTRUCTION

As described in Section 4, each tree task is composed of a binary tree structure. Each node in the tree is randomly assigned a reward between 0 and 1. We also designate a subset of the nodes as *fragile*. These nodes have an extra penalty for reaching them (-8), but an even larger reward (10) after making the correct action in this state. However, if the wrong action is taken in this state, no extra reward is given, and the agent simply suffers the initial penalty. An example initialization of a tree structure is visualized in Figure 6, with an example of a fragile state and extra-reward state highlighted in orange and green, respectively.

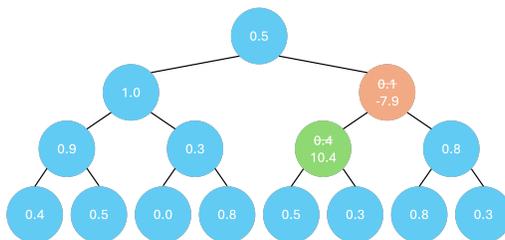


Figure 6: An example tree structure, with a fragile state highlighted in orange, and the extra reward state highlighted in green.

With this definition of state fragility, we can then characterize entire tree domains by the proportion of states which are randomly assigned as fragile, with this proportion ranging from 0 to 1.

A.2 MODEL TRAINING

In Section 4.1 we conduct all experiments with a tree depth of 5. To train our BC models, we construct datasets which contain data points on the action an agent would take at a randomly generated state. We then create multiple datasets at budgets ranging from 2^8 to 2^{16} , at skill levels ranging from $k = 1$ to 4, with fragility proportions in the set $\{0, 0.2, 0.4, 0.6, 0.8, 1.0\}$.

On each generated model, we train a BC model. This model consists of a simple fully-connected neural network with 1 layer, and 16 nodes. We allocated 75% of the dataset for training, and the rest for validation. We trained the neural network on the training dataset using the ADAM optimizer with the initial learning rate of 0.1, and trained until validation loss converged (measured using a patience level of 5 epochs). The full training code is available at <https://anonymous.4open.science/r/iclr2026-C020>.

While training until convergence is a reasonable assumption in domains with expensive data acquisition costs, one question may arise – are we spending too much in compute costs to get improved performance at the low-skill, low-budget region? In Table 4, we find that this is not true – in fact, low-skill models take fewer epochs to converge than high-skill models at the low-budget region.

Budget	Low-Skill Epochs	High-Skill Epochs
2^8	223.6	246.0
2^{10}	161.9	189.2
2^{12}	821.6	846.5
2^{14}	1375.6	635.6
2^{16}	1020.7	470.3

Table 4: Comparison of Low-Skill vs. High-Skill Epochs across budgets.

For each generated model, we evaluate its task performance by generating an additional set of 2^{14} tree environments using the same configuration that the model was originally trained on, and evaluate the expected performance of trained agent in these environments, by computing the average sum of rewards encountered by the agent.

The full set of results is available in Figure 7

In Section 4.2, we largely followed the same training paradigm as in Section 4.1, with two differences.

(1) First, we constructed training datasets which contained a mixture of actions across skill levels. For instance, with Standard BC with a budget of 2^{10} , the dataset actually contained 2^8 data points at $k = 1$, 2^8 data points at $k = 2$, 2^8 data points at $k = 3$, and 2^8 data points at $k = 4$, split into 75% training and 25% validation as before.

(2) We construct a training curricula which consists of a series of training datasets. In each step, we train the model until validation convergence. In the next step, we take the best model from the previous step, and re-run training on the next dataset, after entirely resetting the learning rate and optimizer state.

B MUJoCo EXPERIMENTS

In Section 4.2, we replicated our experiments from Tree in a more complex set of environments. In Table 1, we printed a summary of the results of low-skill and high-skill training. Below, we provide a more comprehensive set of results showing task performance for each environment, demonstrator skill, and training budget.

C CHESS EXPERIMENTS

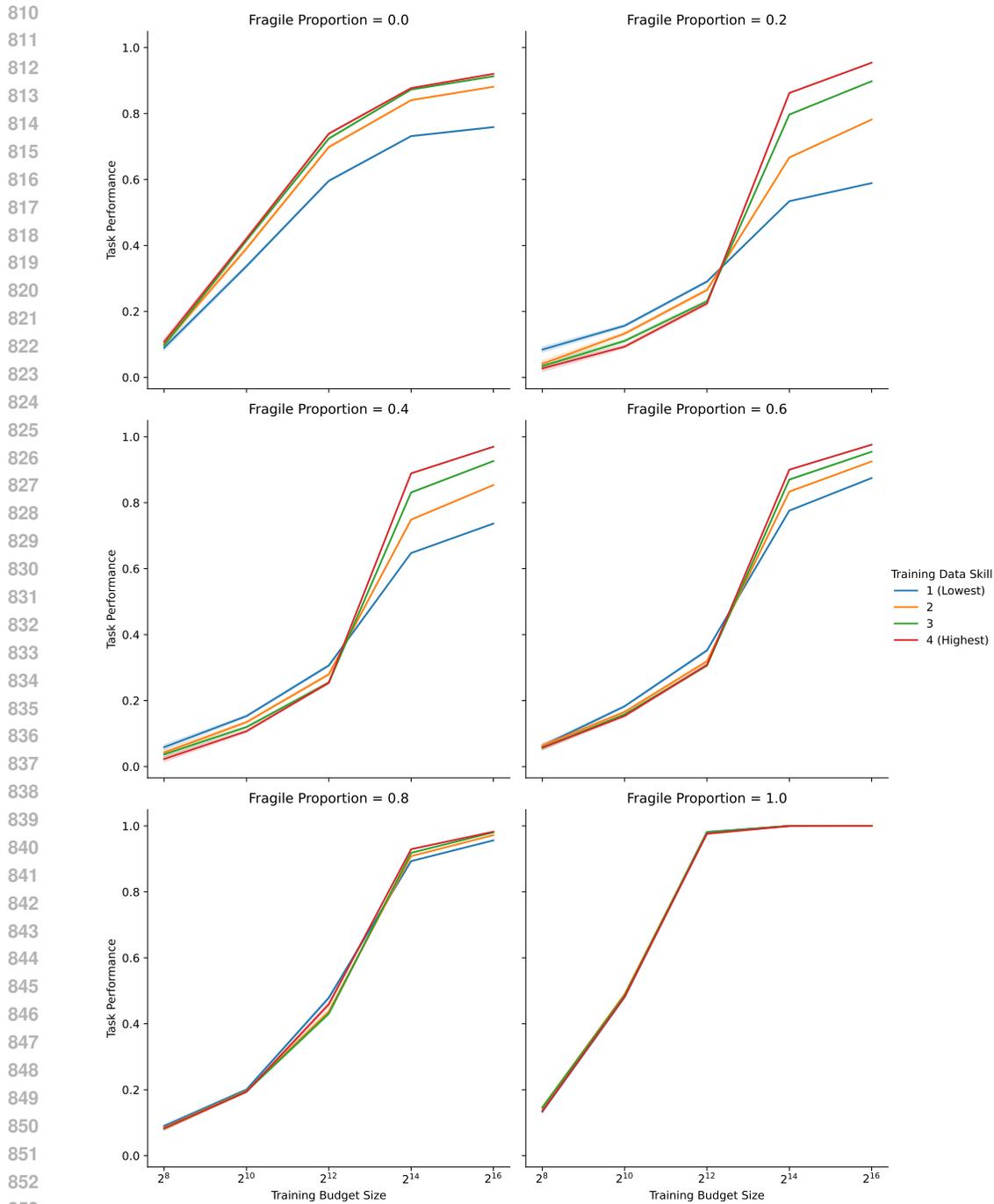
In Section 4.3, we repeat our experiments using Chess instead of a synthetic domain. Broadly, we follow the experimental decisions taken by McIlroy (McIlroy-Young et al., 2020b), which was itself inspired by the training process of the Leela chess engine. We collect human trajectories from the online platform Lichess.org. Each action is labeled according to the skill level of the human who made it, where 400 is the lowest possible skill, and 3400 is roughly the highest skill (though technically the range is uncapped).

In addition to the human data, we construct a synthetic demonstration dataset using the leading chess engine: Stockfish. We have Stockfish play matches against itself, varying the node depth of Stockfish search to simulate skill levels.

C.1 MODEL TRAINING

Each state in the training dataset is represented as a $112 \times 8 \times 8$ matrix, and each action is represented as a one-hot encoded vector of length 1858. Using the transformer-based architecture of the Leela T74 architecture, we randomly initialize the weights of the model before training the BC model on the human or Stockfish data until convergence using SGD with an initial learning rate of 0.1.

Once the model is trained, we evaluate its task performance by having the trained model play against a completely random (i.e untrained) chess model until the end of the game. Our task score is computed as the average result over a match of independent 10240 chess games.



854 Figure 7: Effect of training BC models across budgets and skill levels in the tree domain at varying
 855 proportions of fragile states
 856

857
 858
 859 The full set of results is available in Figure 9.
 860

861 In addition to evaluating the model against a completely random opponent, we also evaluate two
 862 other approaches, playing against a 800-rated opponent (a model trained to convergence on 256M
 863 chess positions at the 800 level), and evaluate the model on a set of puzzles. We present these results
 in Figure 10 and 11.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

Stockfish Depth	Low-Skill Resilience	High-Skill Resilience
1	0.31235	0.27070
10	0.31217	0.27074
100	0.30161	0.25938
1000	0.29605	0.24837
10000	0.30086	0.25364
100000	0.30203	0.25337

Table 5: Resilience scores at varying Stockfish depths.

C.2 RESILIENCE COMPUTATION

In Table 2, we showed resilience scores for low-skill and high-skill humans, evaluated using the highest depth stockfish. In Table 6, we show the resilience against other skill agents, and show that these trends still hold, where humans have higher resilience at lower skills, while engines have higher resilience at higher skills. In Table 5, we also show resilience scores for low-skill and high-skill humans, evaluated using multiple depths of stockfish, and we find that the trends hold at all depths.

Skill	Human	Engine
Low	0.305	0.247
	0.270	
↓	0.273	0.284
	0.263	
High	0.251	0.298

Table 6: Resilience scores for humans and engines across skill levels

D RACING EXPERIMENTS

D.1 MODEL TRAINING

Through an academic partnership with a commercial racing simulator platform, we obtained access to human driving trajectories from a Formula 4 racing series. Specifically, we use a dataset of 10000 complete laps on each of four different tracks: Summit Point Raceway, Road America, Circuito de Navarra, and Autodromo Nazionale Monza. A reference of these map diagrams is available in Figure 12.

We approximate skill by ordering the laps by time, with the fastest 5000 laps representing high-skill data and the slowest 5000 laps representing low-skill data. We evaluate our models using a virtual joystick (pyv, 2025) that forwards the controls to the iRacing simulator.

We train a deep MLP model on state variables (e.g., car speed, heading, etc.) to predict four driving controls (steering, throttle, brake, and gear) at 60Hz. We use action chunking (Zhao et al., 2023b) to ensure stability and accommodate inference latency, predicting 16 future actions at a time.

D.2 ADDITIONAL ANALYSIS

Figure 13 shows the failure experienced by the low-skill model on Road America. As the trajectory shows, the low-skill model makes a tighter turn and eases off the throttle too late, pushing it to the outer boundary of the track and resulting in a more fragile state. In contrast, the high-skill model takes a wider turn and eases the throttle early (dropping to 0.4), allowing it to complete the turn safely. Thus in this situation, the high-skill model is actually more resilient than the low-skill model, explaining why it performs better even at lower data budgets.

We also show a more typical example in Figure 14 where the low-skill model is more resilient than the high-skill model, in this case at the failure experienced by the high-skill models on Summit Point. Here, the high-skill model is too aggressive on the entry to the turn and does not brake enough (0.6) compared to the low-skill model (0.8).

918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971

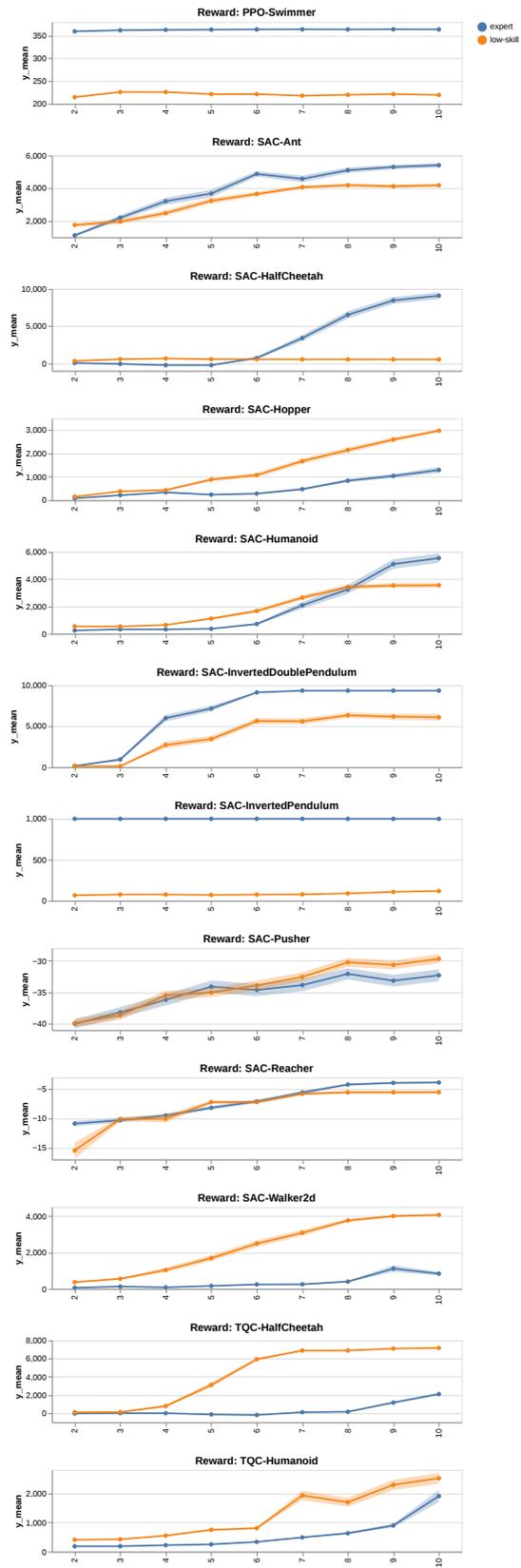


Figure 8: Effect of training BC models across budgets and skill levels in MuJoCo

972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025

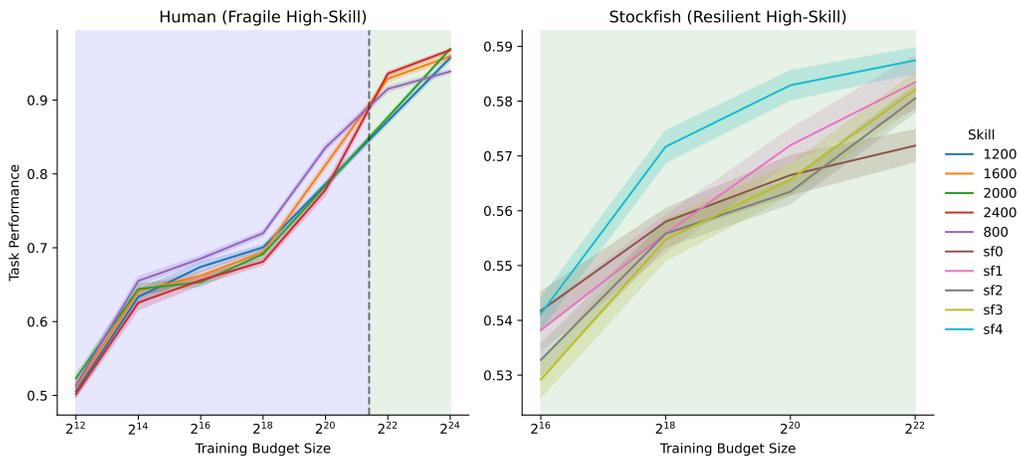


Figure 9: Effect of training BC models on human and engine demonstrations in chess by skill level. Budget regions where low-skill training outperforms high-skill is highlighted in blue.

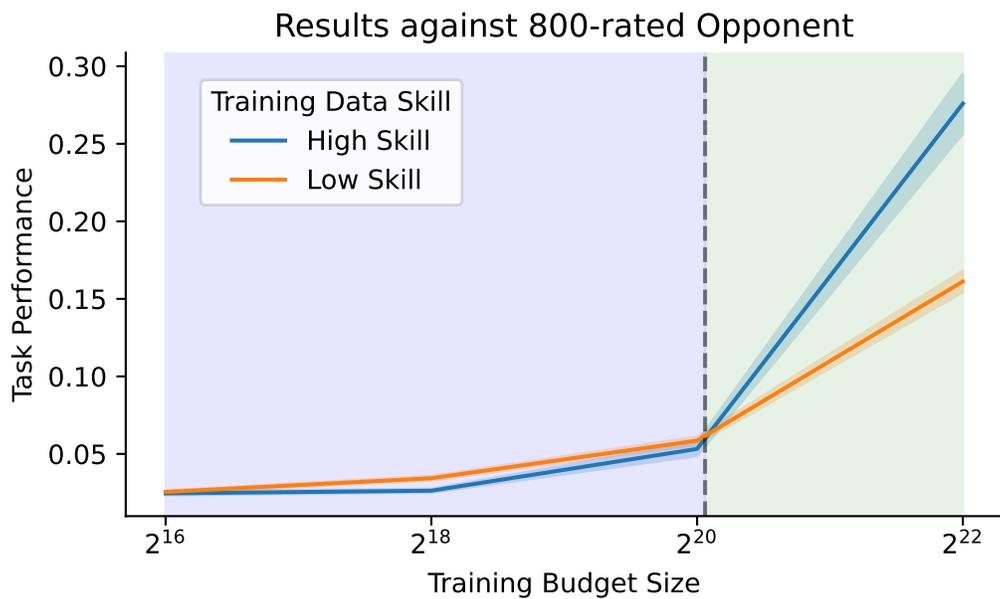


Figure 10: Effect of training BC models on human and engine demonstrations in chess by skill level, evaluating against an 800-rated opponent. Budget regions where low-skill training outperforms high-skill is highlighted in blue.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

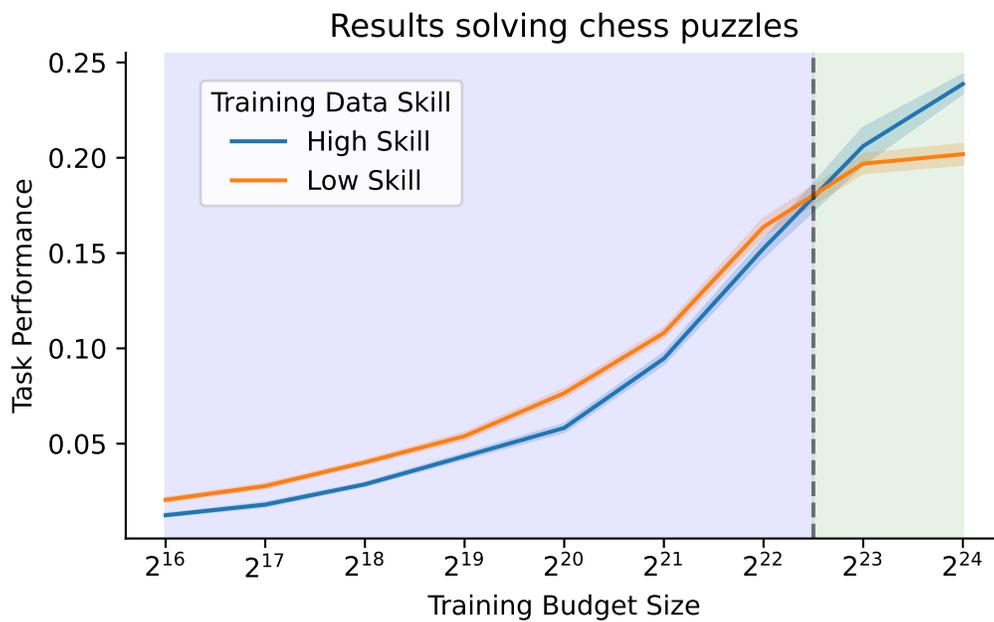


Figure 11: Effect of training BC models on human and engine demonstrations in chess by skill level, evaluated using puzzle solve rate. Budget regions where low-skill training outperforms high-skill is highlighted in blue.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

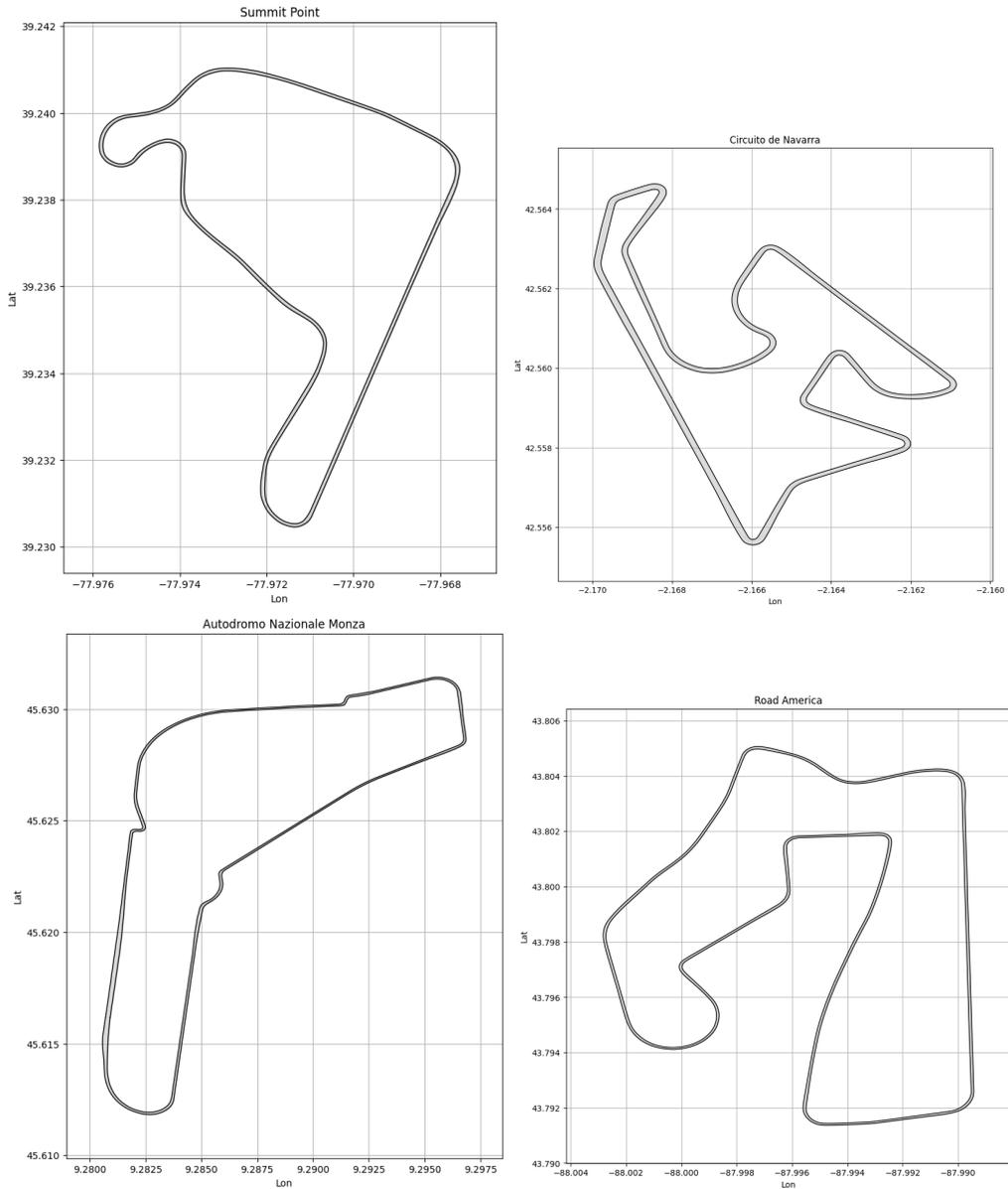


Figure 12: Diagrams of each of the four tracks we evaluate BC on.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

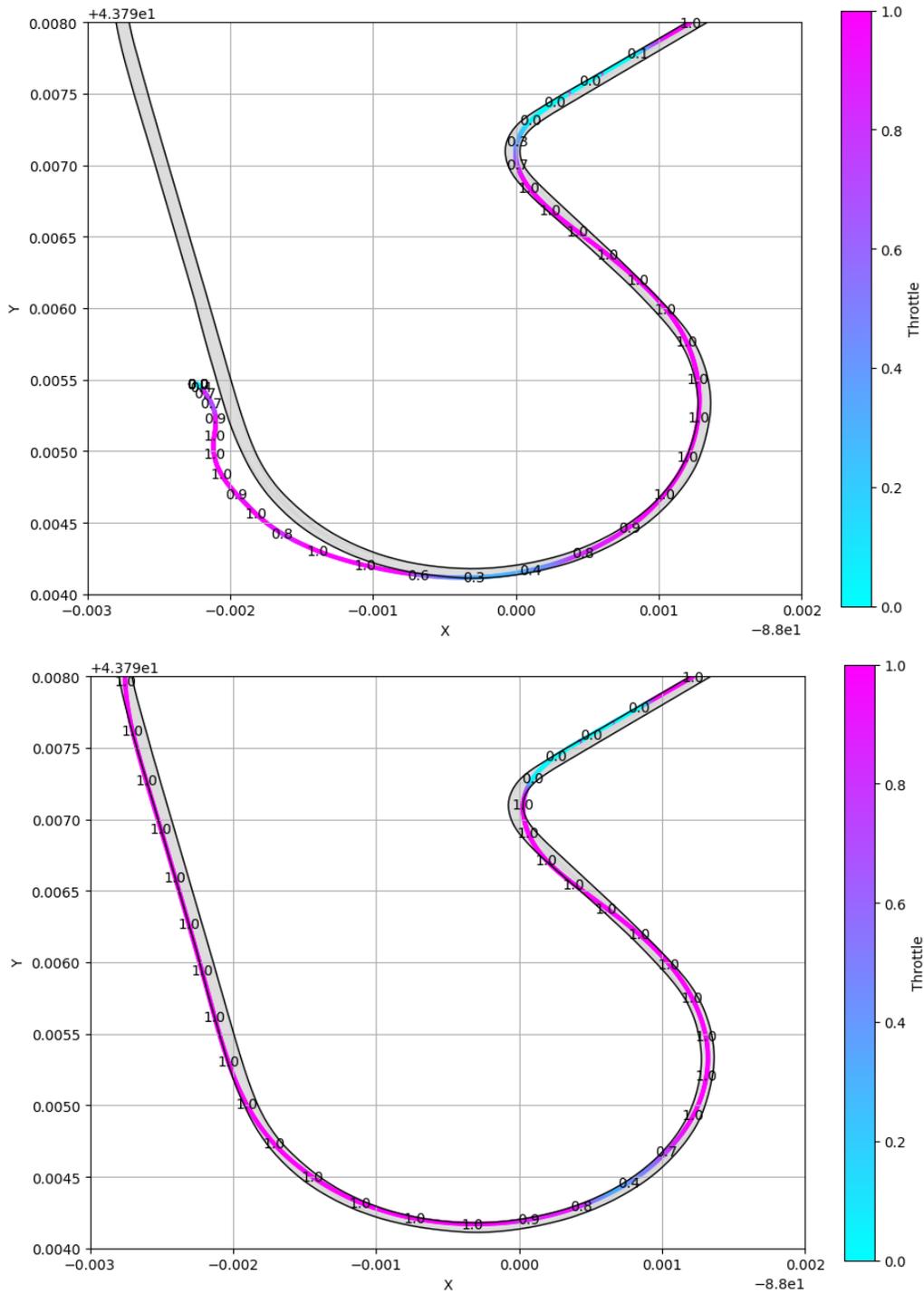


Figure 13: Failed low-skill trajectory and successful high-skill trajectory on Road America

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

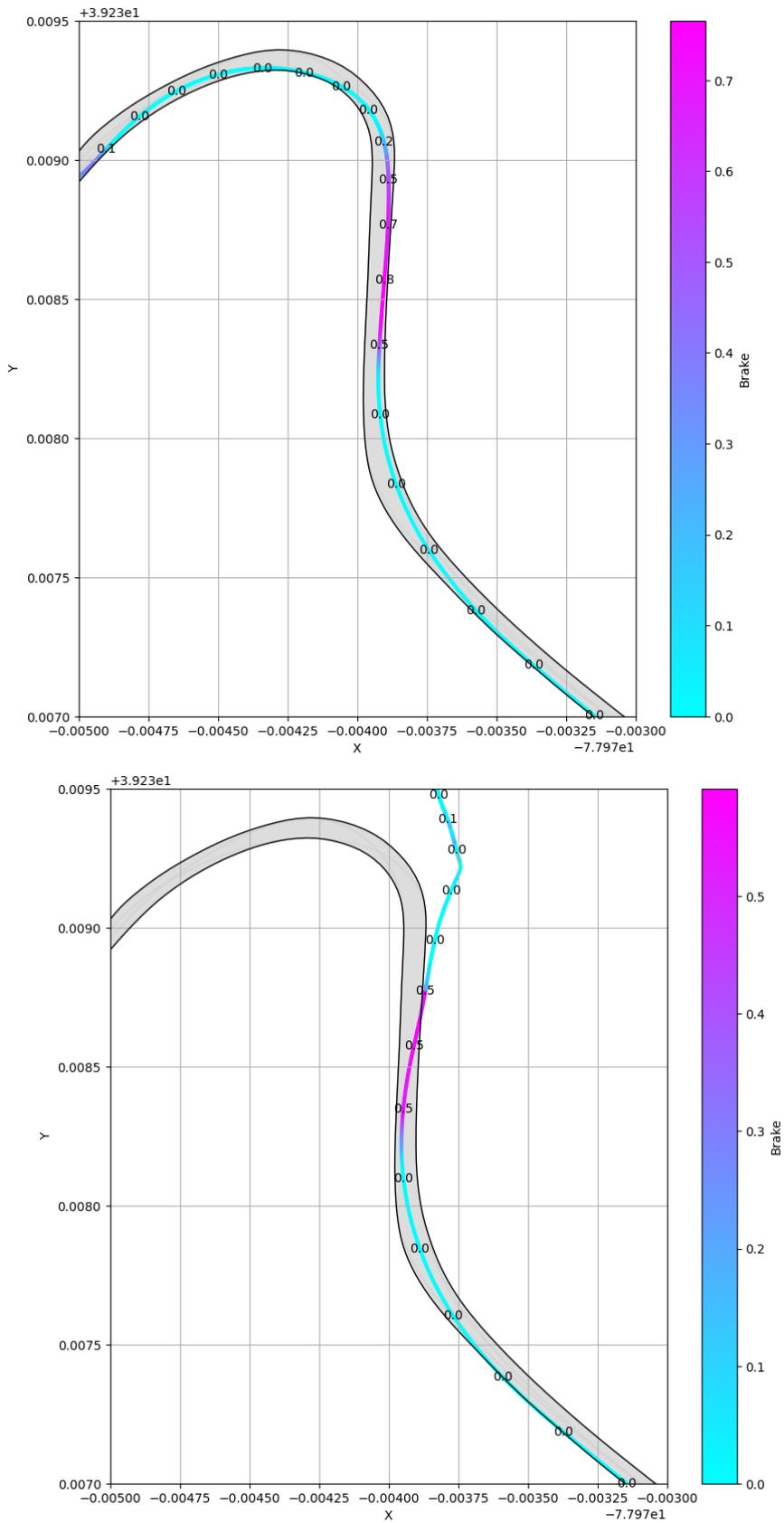


Figure 14: Successful low-skill trajectory and failed high-skill trajectory on Summit Point