# POIL: PREFERENCE OPTIMIZATION FOR IMITATION LEARNING

Anonymous authors

004

006

008 009 010

011 012 013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

Paper under double-blind review

#### Abstract

Imitation learning (IL) enables agents to learn policies by mimicking expert demonstrations. While online IL methods require interaction with the environment, which is costly, risky, or impractical, offline IL allows agents to learn solely from expert datasets without any interaction with the environment. In this paper, we propose Preference Optimization for Imitation Learning (POIL), a novel approach inspired by preference optimization techniques in large language model alignment. POIL eliminates the need for adversarial training and reference models by directly comparing the agent's actions to expert actions using a preferencebased loss function. We evaluate POIL on MuJoCo control tasks under two challenging settings: learning from a single expert demonstration and training with different dataset sizes (100%, 10%, 5%, and 2%) from the D4RL benchmark. Our experiments show that POIL consistently delivers superior or competitive performance against state-of-the-art methods in the past, including Behavioral Cloning (BC), IQ-Learn, DMIL, and O-DICE, especially in data-scarce scenarios, such as using one expert trajectory or as little as 2% of the full expert dataset. These results demonstrate that POIL enhances data efficiency and stability in offline imitation learning, making it a promising solution for applications where environment interaction is infeasible and expert data is limited.



Figure 1: Process overview of preference optimization for imitation learning (POIL). The agent, guided by an expert dataset, compares the agent's actions with expert actions and computes the preference loss for updating the policy parameters. This process does not require environmental interaction, as the red cross indicates, and will be discussed in more detail in Subsection 3.3

#### <sup>54</sup> 1 INTRODUCTION

055 056

*Reinforcement learning* (RL) (Sutton, 2018) has achieved remarkable success across various domains, including video games (Mnih et al., 2015; Schrittwieser et al., 2020), robotics (Kober et al., 2013), and even nuclear fusion control (Degrave et al., 2022). However, defining suitable reward functions remains a significant challenge (Eschmann, 2021), especially in tasks where desired behaviors are abstract or hard to specify, such as control problems (Kiumarsi et al., 2017). Poorly designed reward functions can lead to unintended or unsafe behaviors (Amodei et al., 2016), and deep RL algorithms are often sensitive to reward sparsity (Ladosz et al., 2022), complicating the development of effective reward signals.

064 Imitation learning (IL) (Zare et al., 2024; Hussein et al., 2017) offers an alternative by learning 065 policies directly from expert demonstrations without requiring explicit reward functions. In *online* 066 *IL*, the agent interacts with the environment to learn the expert's behavior. Prominent methods like 067 generative adversarial imitation learning (GAIL) (Ho & Ermon, 2016), adversarial inverse reinforce-068 ment learning (AIRL) (Fu et al., 2017), and discriminator actor-critic (DAC) (Kostrikov et al., 2018) 069 employ a generator (the policy) and a discriminator (distinguishing between expert and agent behaviors) in an adversarial setup to encourage the agent to mimic the expert closely (Garg et al., 2021). Despite their effectiveness, these methods face practical challenges: the adversarial optimization 071 process can be unstable and difficult to train, leading to biased, high-variance gradient estimators 072 and convergence issues (Garg et al., 2021). Moreover, the need for environment interaction makes 073 them impractical in real-world scenarios where such interaction is costly, risky, or infeasible (Lyu, 074 2024; Prudencio et al., 2023). 075

To address these limitations, *offline IL* methods have been developed to learn from pre-collected expert demonstrations without environment interaction. Behavior cloning (BC) (Pomerleau, 1991) is a straightforward approach that directly replicates expert actions through supervised learning. However, BC suffers from compounding errors due to distribution shifts and often requires large amounts of expert data (Ross et al., 2011). Recent advances aim to mitigate these issues by correcting for distribution discrepancies. The DICE family of algorithms, including ValueDICE (Kostrikov et al., 2019), DemoDICE (Kim et al., 2022), and O-DICE (Mao et al., 2024), improve upon BC by addressing distribution shift.

In this paper, we propose a novel offline imitation learning method called preference optimization for imitation learning (POIL), inspired by recent advances in preference-based alignment techniques 085 in large language models (LLMs) (Wang et al., 2024), and clearly different from previous offline IL 086 methods. An overview of the POIL process is illustrated in Figure 1, described in more detail 087 in Subsection 3.3. Specifically, we draw inspiration from direct preference optimization (DPO) 880 (Rafailov et al., 2024), contrastive preference optimization (CPO) (Xu et al., 2024), and self-play 089 fine-tuning (SPIN) (Chen et al., 2024). These methods have been successful in aligning LLMs with 090 human preferences but have specific requirements that limit their direct application to offline IL, 091 namely, DPO and CPO require preference datasets and, in the case of DPO, a reference model. 092 SPIN leverages an expert dataset but still relies on a reference model during training.

In contrast, POIL adapts these techniques to the offline IL setting by introducing a framework that directly compares the agent's actions to expert actions without the need for a discriminator or a reference model, as illustrated in Figure 1. By eliminating the need for preference datasets and reference models, POIL simplifies the learning process, avoids adversarial training instabilities, and enhances computational efficiency. This approach allows POIL to effectively overcome key constraints of existing methods in offline IL while improving overall performance.

We evaluate POIL on standard MuJoCo (Todorov et al., 2012) control tasks, including *HalfCheetah*,
 *Hopper*, and *Walker2d*. POIL consistently delivers superior or competitive results compared to state of-the-art methods across various dataset sizes, particularly excelling in data-scarce scenarios, e.g.,
 a single demonstration or small fractions of the dataset.

- Our contributions are summarized as follows:
- 105

• We introduce POIL, a novel offline imitation learning method that eliminates the need for

We introduce POIL, a novel offline imitation learning method that eliminates the need for adversarial training, preference datasets, and reference models by directly comparing agent and expert actions.

• We provide empirical evidences on MuJoCo tasks showing that POIL achieves superior or competitive performance against state-of-the-art methods in the past, especially in data-limited scenarios.

• We conduct ablation studies to analyze the impact of some key hyper-parameters on POIL's performance, providing insights into its robustness and applicability.

These results suggest that preference optimization techniques from LLM alignment are effectively
 adapted to offline imitation learning, opening new avenues for research and applications in control
 and robotics.

117 118

108

110

111

112

113

2 RELATED WORK

119 120

121

2.1 OFFLINE IMITATION LEARNING

Offline imitation learning methods (Zare et al., 2024; Hussein et al., 2017) learn from static datasets without needing interaction with the environment. An early approach, behavior cloning (Pomerleau, 1991), learns directly from expert demonstrations but has issues with compounding errors and distribution shift, especially with limited data (Ross et al., 2011).

126 To overcome these issues, the DICE (DIstribution Correction Estimation) family of algorithms pro-127 vides improvements. ValueDICE (Kostrikov et al., 2019) minimizes the KL divergence between stationary distributions, while SoftDICE (Sun et al., 2021) employs the Earth-Mover distance for 128 distribution matching. DemoDICE (Kim et al., 2022) can use demonstrations of varying quality, 129 and ODICE (Mao et al., 2024) adds orthogonal-gradient updates to handle conflicting gradients in 130 learning. Other methods include adaptations of *inverse reinforcement learning* for offline use (Zolna 131 et al., 2020; Yue et al., 2023), energy-based models (Jarrett et al., 2020), and OTR (Luo et al., 2023). 132 Despite these advances, challenges remain. Some methods, such as DemoDICE and SMODICE (Ma 133 et al., 2022), need extra data beyond expert demonstrations. Others struggle with limited datasets or 134 single demonstrations. 135

1362.2 PREFERENCE-BASED REINFORCEMENT LEARNING

138 Preference-based RL (PbRL) (Wirth et al., 2017) has emerged as a promising approach to address 139 the challenges of reward function design in traditional RL by incorporating human preferences into the learning process. Early work by Christiano et al. (2017) demonstrated the potential of PbRL 140 using deep learning techniques. This pioneering work opened new avenues for tackling challenging 141 domains but relied heavily on external preference feedback. Subsequent research focused on re-142 ducing this dependency. Lee et al. (2021) introduced PEBBLE, which improved feedback efficiency 143 through experience relabeling and unsupervised pre-training. Park et al. (2022) further improved this 144 with SURF, a semi-supervised approach leveraging data augmentation, yet neither fully eliminated 145 the need for human feedback. 146

A significant advancement in the field came with the inverse preference learning (IPL) proposed by Hejna & Sadigh (2024). IPL represents a novel approach specifically designed for learning from offline preference data. It leverages the key insight that for a fixed policy, the Q-function encodes all necessary information about the reward function. Using the Bellman operator, IPL eliminates the need for an explicit reward model, thereby simplifying the algorithm and improving parameter efficiency. This innovation marks a crucial step towards more efficient and scalable PbRL methods. However, IPL still requires an external preference dataset and remains a value-based method, leaving room for further improvements in data efficiency and algorithmic approach.

154 155

156

2.3 ALIGNMENT TECHNIQUES IN LARGE LANGUAGE MODELS

Reinforcement learning from human feedback (RLHF) (Kaufmann et al., 2023; Wang et al., 2024)
has emerged as a powerful approach for aligning large language models (LLMs) with human preferences. Pioneered by InstructGPT (Ouyang et al., 2022) and further developed by Bai et al. (2022),
RLHF involves training a reward model based on human preference data and then optimizing the
policy using reinforcement learning guided by this reward model. While effective, RLHF is computationally intensive and requires careful hyper-parameter tuning.

To address these challenges, a class of methods known as DPO-like methods has been proposed
 as simpler alternatives that bypass explicit reward modeling. DPO (Rafailov et al., 2024) directly
 optimized the policy to match preference data using a classification loss. The standard DPO loss
 function is defined as:

166 167 168

169

$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \left( \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right) \right], \tag{1}$$

where  $\pi_{\theta}$  is the agent's policy parameterized by  $\theta$ ,  $\pi_{ref}$  is the reference model,  $\sigma(z) = 1/(1 + e^{-z})$ is the sigmoid function,  $\beta$  is a scaling factor,  $\mathcal{D}$  is the dataset of preference pairs  $(x, y_w, y_l)$ , and  $y_w$ and  $y_l$  denote the preferred and less preferred responses given a prompt x, respectively.

Several variants of DPO, collectively referred to as DPO-like methods, have been developed to improve performance or address specific issues. For instance, identity preference optimization (IPO) (Azar et al., 2024) aimed to mitigate overfitting in preference learning, DPO-positive (DPOP) (Pal et al., 2024) introduced additional regularization to prevent reward degradation, and Kahneman-Tversky optimization (KTO) (Ethayarajh et al., 2024) incorporated insights from prospect theory to better model human decision-making.

Recent researches have focused on developing reference-free DPO-like methods to eliminate the need for a fixed reference model. Simple preference optimization (SimPO) (Meng et al., 2024)
introduced a loss function with length normalization and a reward margin, enabling reference-free optimization while addressing response length control. CPO (constrastive preference optimization) (Xu et al., 2024) showed that when the reference model perfectly aligns with the true data distribution of preferred data, the DPO loss is upper-bounded by a simpler loss function without a reference model by assuming a uniform distribution U. The preference part of the CPO loss function is given by:

187 188

189

$$\mathcal{L}_{\text{CPO}}^{Preference}(\pi_{\theta}) = \mathcal{L}_{\text{DPO}}(\pi_{\theta}; U) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \left(\log \pi_{\theta}(y_w | x) - \log \pi_{\theta}(y_l | x)\right)\right)\right].$$
(2)

These DPO-like methods not only enhance computational and memory efficiency but also retain
 comparable optimization performance to the standard DPO.

Most closely related to our work is SPIN (self-play fine-tuning) (Chen et al., 2024), which employed
a self-play mechanism to improve an LLM without additional human preference data iteratively.
SPIN generated its training data and refines itself by distinguishing between current and previous
outputs, continuously updating its reference model.

Our research bridges the gap between language model alignment and imitation learning, demonstrating how DPO-like alignment techniques, originally developed for text generation problems in LLMs, can be successfully adapted to control problems in reinforcement learning. This cross-domain application opens up new possibilities for improving imitation learning in complex, real-world tasks, and highlights the potential of adapting DPO-like methods.

202 203

#### **3** PREFERENCE OPTIMIZATION FOR IMITATION LEARNING

204 205 206

### 3.1 Adapting DPO-like Methods for Imitation Learning

Our approach leverages the strengths of DPO-like methods, and namely combines the self-play 207 mechanism from SPIN with the reference-free optimization from CPO, to enhance the imitation 208 learning process. In this method, we eliminate the need for a reference model in the self-play setup 209 by adopting CPO's reference-free loss function and allowing us to apply self-play in offline imitation 210 learning without the computational overhead of maintaining a reference model. Specifically, we 211 focus on iteratively refining the agent's policy by directly comparing its actions to those of experts, 212 thus enabling the model to align more closely with expert behavior while reducing computational 213 complexity. 214

In SPIN, a model generates synthetic data and refines itself by distinguishing between current and previous outputs using a continuously updated reference model. However, maintaining and updating

this reference model adds computational complexity. To address this, we adapt CPO's reference-free optimization, eliminating the need for a reference model while retaining the benefits of self-play.

In our approach, expert demonstrations serve as positive samples  $(y_w = a_E)$ , while the model's own actions during training are treated as negative samples  $(y_l = a)$ , where *a* denotes the agent's action and  $a_E$  denotes the expert's action. This direct comparison enables the model to prioritize actions that align more closely with expert behavior. By progressively learning from its own generated data in relation to expert demonstrations, the agent refines its policy, achieving a more stable and efficient imitation learning process without relying on predefined rewards or reference models.

3.2 POIL OBJECTIVE

225

226

231 232

233 234

235 236

237

238

239

240

248 249

250

251

252 253 254

265

Our goal is to adapt DPO-like methods to the offline IL setting by eliminating the need for a reference model and incorporating a BC regularization term to enhance performance. Inspired by the findings of CPO as described in Subsection 2.3, we consider the expert's policy as the true data distribution of preferred actions. First, we obtain a reference-free loss function for imitation learning as follows.

$$\mathcal{L}_{\text{POIL}}(\pi_{\theta}) = -\mathbb{E}_{(s, a_E, a) \sim \mathcal{D}_E} \left[ \log \sigma \left( \beta \left( \log \pi_{\theta}(a_E | s) - \log \pi_{\theta}(a | s) \right) \right) \right], \tag{3}$$

This loss function is designed to achieve two main objectives:

- 1. Align with expert behavior. By maximizing  $\log \pi_{\theta}(a_E|s)$ , we encourage the agent to assign higher probabilities to the expert's actions.
- 2. Discourage sub-optimal actions. By minimizing  $\log \pi_{\theta}(a|s)$ , we encourage the agent to move away from sub-optimal (agent's) behaviors.

The loss design of POIL (Equation 3) is to simultaneously minimize the divergence between the expert's behavior and the agent's behavior while maximizing the preference of expert actions over the agent's current actions.

To further encourage the policy to closely mimic expert actions, we incorporate a BC (behavior cloning) regularization term, similar to the approach in Xu et al. (2024). Specifically, we add the negative log-likelihood of expert actions under the agent's policy:

$$\mathcal{L}_{BC}(\pi_{\theta}) = -\mathbb{E}_{(s,a_E)\sim\mathcal{D}_E}\left[\log \pi_{\theta}(a_E|s)\right].$$
(4)

Our overall augmented POIL loss function then combines the preference optimization and the BC regularization:

$$\mathcal{L}_{\text{POIL}}^{\text{aug}}(\pi_{\theta}) = \mathcal{L}_{\text{POIL}}(\pi_{\theta}) + \lambda \cdot \mathcal{L}_{\text{BC}}(\pi_{\theta}), \tag{5}$$

where  $\lambda$  is a hyper-parameter that balances the trade-off between preference optimization and behavior cloning.

258 By incorporating  $\lambda$  as a tunable parameter, we allow for greater flexibility in balancing the influence 259 of the BC regularization term, which is crucial in scenarios with varying quality or quantity of expert 260 data.

Consider the provided by Xu et al. (2024), adapting it to the offline imitation learning setting and introducing  $\lambda$  to enhance the method's adaptability. This results in a loss function that effectively guides the policy towards aligning with the expert data distribution without the need for a reference model.

266 3.3 ALGORITHM

The POIL algorithm, detailed in Algorithm 1, iteratively refines the agent's policy to better align with expert behavior through preference-based optimization and behavior cloning regularization. As shown in Figure 1, the process begins by initializing the policy parameters randomly. During

each iteration, the algorithm samples state-action pairs from the expert demonstrations, which serve as the basis for learning.

To generate agent actions in a differentiable manner, we employ the *reparameterization trick* when sampling from the policy  $\pi_{\theta}(a|s)$ . Specifically, the policy outputs a mean  $\mu_{\theta}(s)$  and a standard deviation  $\zeta_{\theta}(s)$ . We sample actions using  $a = \tanh(\mu_{\theta}(s) + \zeta_{\theta}(s) \cdot \epsilon)$ , where  $\epsilon \sim \mathcal{N}(0, I)$ . This approach ensures that gradients flow through the sampling process during optimization, allowing effective updating of the policy parameters  $\theta$ .

size <i>m</i> , learning rate $\eta$ , number of iterations <i>T</i> . Randomly initialize policy parameters $\theta$ for iteration = 1 to <i>T</i> do Sample a batch of expert state action pairs $\{(c_1, q_{T_1})\}^m$ from $\mathcal{D}_T$
Randomly initialize policy parameters $\theta$ for iteration = 1 to T do Sample a batch of expert state action pairs $\{(c_1, a_2)\}^m$ from $\mathcal{D}_{\mathbb{T}}$
for iteration = 1 to T do Sample a batch of expert state action pairs $\{(a_1, a_2)\}^m$ from $\mathcal{D}_{\mathbb{T}}$
Sample a batch of expert state action pairs $\int (a_1, a_2) m$ from $\mathcal{D}_2$
Sample a batch of expert state-action pairs $\{(s_j, a_E, j)\}_{j=1}$ from $\mathcal{D}_E$
Generate agent actions $a_j$ via reparameterized sampling from $\pi_{\theta}(a s_j)$ for all j
for each $(s_i, a_{E,i}, a_i)$ in batch do
Compute the POIL loss: $\mathcal{L}_{POIL} = -\log \sigma \left(\beta \left(\log \pi_{\theta}(a_{E,i} s_i) - \log \pi_{\theta}(a_i s_i)\right)\right)$
Compute the BC regularization term: $\mathcal{L}_{BC} = -\log \pi_{\theta}(a_{E,j} s_j)$
Combine the losses: $\mathcal{L}_{POII}^{aug} = \mathcal{L}_{POIL} + \lambda \cdot \mathcal{L}_{BC}$
Update policy parameters: $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}_{poll}^{aug}$
end for
return $\theta$
end for

293 294 295

296

297

298

299

In Algorithm 1,  $\mathcal{D}_E$  represents the dataset of expert demonstrations, and N is the total number of expert state-action pairs. The regularization coefficient  $\lambda$  controls the weight of the BC regularization relative to the preference-based loss, while the scaling factor  $\beta$ , learning rate  $\eta$ , and number of iterations T are hyper-parameters that control the optimization process. The batch size m determines how many samples are used in each iteration to compute the gradient.

The algorithm proceeds by sampling batches of expert data and generating corresponding agent actions. The total loss  $\mathcal{L}_{\text{total}}$  is computed for each sample in the batch, balancing between preference optimization and behavior cloning. The policy parameters  $\theta$  are then updated using gradient descent to minimize the total loss, thereby improving the agent's policy to better match the expert's behavior.

304 305

306

#### 4 EXPERIMENTS

In this section, we evaluate the performance of POIL on control tasks in the MuJoCo environment (Todorov et al., 2012). Additionally, to demonstrate the scalability and robustness of POIL in more complex environments, we provide experimental results on the Adroit tasks from the D4RL benchmark (Fu et al., 2020) in Appendix A. We conduct several experiments<sup>1</sup> to demonstrate the effectiveness of POIL under different settings and compare it against some baseline methods, which include those state-of-the-art in the past.

## 313314 4.1 EXPERIMENTAL SETUP

 We conduct experiments on standard MuJoCo control tasks, specifically HalfCheetah-v2, Hopper-v2, and Walker2d-v2. These environments are widely used benchmarks in reinforcement learning, requiring agents to learn complex locomotion behaviors in high-dimensional state and action spaces.

For the single demonstration experiments, we utilize one expert trajectory per task, sourced from the same dataset as used in ValueDICE (Kostrikov et al., 2019). These expert trajectories are generated by well-trained policies and present a challenging setting for imitation learning due to the limited data available.

<sup>&</sup>lt;sup>1</sup>Refer to the supplementary to reproduce experiments in this section.

In the experiments involving different dataset sizes, we use datasets from the D4RL benchmark (Fu
 et al., 2020). We create subsets of the full dataset by randomly selecting different percentages of the
 data, namely 100%, 10%, 5%, and 2%. This approach allows us to assess the performance of POIL
 and baseline methods under different data availability conditions.

We compare the performance of POIL against several baseline methods, including BC (Pomerleau, 1991), IQ-Learn (Garg et al., 2021), DMIL (Zhang et al., 2023), and O-DICE (Mao et al., 2024), many of which were state-of-the-art in offline imitation learning in the past. These methods provide a comprehensive benchmark for evaluating POIL.

For fair comparison, we use the same neural network architecture for all methods. The policy network consists of two fully connected layers, each with 256 units, with ReLU activation functions applied after each layer. All models are trained for 100k timesteps. We use the Adam optimizer (Kingma, 2014) for optimization with default parameters. The experiments were conducted on a system equipped with 4 NVIDIA RTX A6000 GPUs, 128GB of RAM, and an AMD Threadripper PRO 5965WX processor featuring 24 cores and 48 threads.

338 339 340

#### 4.2 SINGLE DEMONSTRATION LEARNING

In this experiment, we evaluate the ability of POIL to learn effective policies from a single expert trajectory. The hyper-parameter  $\beta$  is set to 0.2 in this experiment.

This experiment tests the data efficiency of imitation learning methods when only minimal expert data is available. Detailed results, including standard deviations over multiple runs, are provided in Appendix B.

347 348

349

350

351

Table 1: The performance of various methods trained with a single expert demonstration on MuJoCo tasks. The results are averaged over 3 different runs, each using a unique random seed, and the scores represent the average over the last ten epochs. (The bold numbers represent the best, while the underscored numbers are the second best. Note that the scores are not normalized to expert data because we cannot directly get the expert scores from this dataset (Kostrikov et al., 2019)).

	Environment	Traj.	BC	IQ-Learn	DMIL	<b>O-DICE</b>	$\mathbf{POIL}_{\lambda=1}$	$POIL_{\lambda=0}$
-	HalfCheetah	traj 1	2775.39	3621.73	2764.81	1532.63	<u>3843.75</u>	4628.88
		traj2	3031.12	2957.77	<u>3546.40</u>	1801.12	2228.04	4833.13
		traj3	2927.38	<u>3805.18</u>	2802.43	1032.74	2645.44	4309.46
-	Hopper	traj l	1548.32	2431.73	2947.13	1672.45	2426.22	3501.67
		traj2	1687.93	2601.54	3195.28	1748.63	2534.35	3066.00
		traj3	1850.54	2520.18	3301.41	1797.43	<u>3372.09</u>	3499.08
		traj l	881.23	1011.89	1942.76	1543.67	<u>2645.29</u>	4722.79
	Walker2d	traj2	1023.85	1062.52	1686.98	1392.54	1826.08	5233.93
_		traj3	935.63	997.21	1864.76	1483.52	<u>2493.53</u>	3745.50

362 363 364

365

366

367

368

369

370

359 360 361

> As shown in Table 1, POIL with  $\lambda = 0$  achieves the highest performance across all trajectories in the HalfCheetah-v2 task, outperforming all baseline methods. In the Walker2d-v2 task, POIL with  $\lambda = 1$  wins on two trajectories, and POIL with  $\lambda = 0$  for one trajectory. For the Hopper-v2 task, POIL with  $\lambda = 1$  and  $\lambda = 0$  each achieves the best performance on one trajectory, with DMIL outperforming both POIL variants on the remaining one trajectory. Overall, our method only loses to DMIL on one trajectory and outperforms all baseline methods on the remaining tasks. Particularly, POIL with  $\lambda = 0$  performs the best on five trajectories, the second on three, and the third on one.

This demonstrates POIL's ability to effectively utilize limited expert data and suggests superior data efficiency compared to other methods. More discussion about  $\lambda$  is given inSubsection 4.4.2.

373

375

374 4.3 DIFFERENT DATASET SIZES

In this experiment, we assess the performance of POIL and baseline methods when trained on differ ent proportions of the D4RL datasets: 100%, 10%, 5%, and 2%. This evaluation tests the robustness of the methods under different amounts of training data, particularly in data-scarce scenarios.

As shown in Table 2, POIL demonstrates highly competitive performance across different dataset sizes and tasks. In the HalfCheetah-v2 and Hopper-v2 tasks, POIL with  $\lambda = 1$  performes the best, while POIL  $\lambda = 0$  performs the second best for all dataset sizes among all baseline methods. In the Walker2d-v2 task, POIL with  $\lambda = 0$  performs competitively. While IQ-Learn and O-DICE achieve the highest scores in this task, POIL remains very close to them in performance within 1% for all dataset sizes. Particularly for 10% dataset size, it achieves the best result.

Table 2: Performance of different methods trained on varying dataset sizes from D4RL on MuJoCo tasks. Results are normalized to the expert scores, and the scores represent the maximum scores achieved across the training process. (The bold numbers are the best, while the underscored numbers are the second best.)

Environment	Ratio	BC	IQ-Learn	DMIL	O-DICE	<b>POIL</b> $_{\lambda=1}$	$POIL_{\lambda=0}$
	100%	91.95	80.69	94.62	95.29	96.60	<u>95.75</u>
HalfChaatah	10%	90.64	61.48	95.37	94.67	96.14	<u>95.42</u>
папспестап	5%	82.90	58.42	93.70	91.79	96.47	<u>95.29</u>
	2%	23.58	44.17	89.89	92.73	95.51	<u>95.08</u>
	100%	95.06	110.41	112.50	114.93	116.23	118.86
Honnar	10%	83.52	104.22	112.66	114.33	117.66	<u>117.36</u>
порры	5%	73.35	100.13	112.93	<u>117.29</u>	115.45	119.01
	2%	53.54	98.35	112.81	113.67	<u>116.93</u>	119.16
	100%	107.35	113.67	109.24	111.85	110.01	<u>112.97</u>
Walker2d	10%	105.36	<u>111.64</u>	110.40	110.40	110.00	112.65
waikei 20	5%	103.21	111.84	109.53	110.79	110.15	<u>111.19</u>
	2%	58.34	<u>110.58</u>	109.32	110.64	110.46	110.34

These results indicate that POIL not only excels in challenging environments like HalfCheetah-v2 and Hopper-v2, but also performs competitively in Walker2d-v2. The consistent high performance across different dataset sizes underscores POIL's robustness and effectiveness in offline imitation learning, even in data-scarce scenarios. More discussion about  $\lambda$ s is given in Subsection 4.4.2.

410 411

406

407

408

409

384 385 386

387

388

389

396 397

412

413 414

415 416 4.4

4.4.1 Impact of Scaling Factor  $\beta$ 

ABLATION STUDY

417 We conduct an ablation study to analyze the impact of the scaling factor  $\beta$  on POIL's perfor-418 mance. Figure 2 shows the performance of POIL on the HalfCheetah-v2, Hopper-v2, and 419 Walker2d-v2 tasks with different values of  $\beta$ . We set  $\lambda$  to zero based on the observations from 420 Subsection 4.2 and Subsection 4.3, where POIL with  $\lambda = 0$  demonstrates better performance in 421 most cases.

422 Our observations indicate that a smaller value of  $\beta$ , specifically  $\beta = 0.2$ , yields the best performance 423 across most tasks in this setting. This suggests that choosing an appropriate scaling factor effectively 424 controls the sharpness of the preference function in the loss, influencing how strongly the model 425 distinguishes between expert and agent actions. A smaller  $\beta$  value tends to smooth the preference 426 function, which leads to more stable gradients and improved training dynamics, as we will further 427 explore when discussing different choices of  $\lambda$  in Subsection 4.4.2.

428 While  $\beta = 0.2$  performed well in most cases as above, it may not universally be the best choice 429 for all tasks or environments. Different environments may require different scaling factors based on 430 their complexity and the nature of the expert data. A consistent observation that smaller  $\beta$  values 431 yield better performance suggests that a smoother preference function be generally beneficial in 600 offline imitation learning with POIL.



Figure 2: Impact of the scaling factor  $\beta$  on POIL's performance in the HalfCheetah-v2, Hopper-v2, and Walker2d-v2 tasks.

#### 4.4.2 IMPACT OF BC COEFFICIENT $\lambda$

440

441

446

In addition to  $\beta$ , we also conduct further experiments to analyze the impact of the regularization coefficient  $\lambda$ , which controls the weight of the BC term. These experiments were conducted under the same conditions as described in Section 4.2, using only a single expert demonstration for each of the three environments, HalfCheetah-v2, Hopper-v2, and Walker2d-v2. The results are averaged over three different random seeds to ensure robust evaluation. Figure 3 illustrates the performance of POIL across the three tasks as  $\lambda$  varies from 0 to 3, specifically  $\lambda = [0, 0.2, 0.4, 0.6, 0.8, 1.0, 3.0]$ , for five different  $\beta = [0.1, 0.2, 0.5, 0.8, 1.0]$ .

454 Our findings consistently show that smaller values of  $\lambda$  lead to higher performance across all three 455 tasks for most  $\beta$ . In particular,  $\lambda = 0$  yields the highest returns for most  $\beta$  in all three environments. 456 As  $\lambda$  increases, performance decreases notably, especially for larger values of  $\lambda$ , such as 3.0, where 457 performance deteriorates significantly across the board.

This suggests that the contribution of the BC regularization term, controlled by  $\lambda$ , is not as beneficial in this offline imitation learning setup, particularly when smaller  $\beta$  values already balance the preference-based loss effectively. Higher values of  $\lambda$  seem to introduce a trade-off that limits the agent's ability to align its policy with the expert data, resulting in lower average returns.

462 However, as demonstrated in Subsection 4.3, when more demonstration data is available, both  $\lambda = 0$ 463 and  $\lambda = 1$  perform equally well across various dataset sizes, indicating that the regularization term's 464 effect varies depending on the amount of available expert data. In more data-rich environments, 465 moderate values of  $\lambda$  tend to help stabilize training without severely hindering performance, unlike 466 in the single demonstration scenario.

These results indicate that POIL performs best with minimal or no BC regularization ( $\lambda \approx 0$ ) in data-scarce settings. In contrast, with larger datasets, the impact of  $\lambda$  becomes less pronounced, and both  $\lambda = 0$  and  $\lambda = 1$  lead to competitive performance.





#### 487 4.4.3 COMPARISON BETWEEN DIFFERENT PREFERENCE OPTIMIZATION METHODS

In POIL, we utilize the CPO loss to achieve good performances in the tasks above. In addition to CPO, there are still many other reference-free preference optimization methods proposed in the context of preference learning in LLM that can be adapted to offline imitation learning, too. In this subsection, we include these methods for comparison, including SimPO (Meng et al., 2024), SLIC-HF (Zhao, 2023), RRHF (Yuan et al., 2023), and ORPO (Hong, 2024), to see whether any of them is better.

As shown in Figure 4, POIL (in olive green) significantly outperforms these methods across all tasks. This suggests that our approach of directly comparing agent actions with expert actions using the POIL loss function is more effective in guiding the policy learning process in offline imitation learning. The superior performance of POIL highlights its effectiveness over other reference-free methods.



Figure 4: Performance comparison between POIL and other reference-free preference optimization methods on single demonstration task.

#### 5 DISCUSSION

In this paper, we introduce POIL, a novel method inspired by preference optimization techniques
 from large language model alignment. POIL eliminates the need for adversarial training by directly
 comparing agent actions to expert actions. Through extensive experiments on MuJoCo tasks and the
 D4RL benchmark, we demonstrate that POIL performs best or competitively against state-of-the-art
 methods, particularly in data-scarce settings.

525 A key finding from our study is the role of the regularization coefficient  $\lambda$ . While CPO typically 526 sets  $\lambda = 1$ , our results suggest that  $\lambda = 0$  often leads to better performance, particularly in simpler 527 tasks or data-scarce settings where overfitting is less of a concern. However, we found that  $\lambda$  is more 528 useful in tasks where the expert data is harder to learn from, such as HalfCheetah-v2, where 529 the expert reaches a very high score. In such cases, the regularization from  $\lambda$  helps the agent stay 530 closer to expert actions, providing a stabilizing effect during training. This allows the agent to learn 531 more effectively in environments where the expert's performance sets a difficult target.

This insight opens potential opportunities for future research. For example, further investigation into tuning  $\lambda$  based on task complexity or the quality of the expert data provide more flexible and adaptive learning strategies. In tasks where the expert's performance is exceptionally high, increasing  $\lambda$  helps prevent the agent from diverging too far from expert behavior. In contrast, for simpler tasks or when expert data is more limited, reducing  $\lambda$  might enable the agent to explore a broader range of policies without over-regularization.

Overall, POIL offers a robust solution for offline imitation learning, especially when expert data is
 limited or challenging to learn from. Its flexibility in adapting to different dataset sizes and task
 difficulties makes it a promising direction for future research in imitation learning and related fields.

## 540 REFERENCES

558

570

571

572

578

579

- 542 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Con-543 crete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland,
  Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*,
  pp. 4447–4455. PMLR, 2024.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning
   converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*,
   2024.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- Jonas Degrave, Federico Felici, Jonas Buchli, Michael Neunert, Brendan Tracey, Francesco Carpanese, Timo Ewalds, Roland Hafner, Abbas Abdolmaleki, Diego de Las Casas, et al. Magnetic control of tokamak plasmas through deep reinforcement learning. *Nature*, 602(7897):414–419, 2022.
- Jonas Eschmann. Reward function design in reinforcement learning. *Reinforcement Learning Algorithms: Analysis and Applications*, pp. 25–33, 2021.
- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse reinforcement learning. *arXiv preprint arXiv:1710.11248*, 2017.
  - Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 573 Divyansh Garg, Shuvam Chakraborty, Chris Cundy, Jiaming Song, and Stefano Ermon. Iq-learn:
  574 Inverse soft-q learning for imitation. *Advances in Neural Information Processing Systems*, 34:
  575 4028–4039, 2021.
- Joey Hejna and Dorsa Sadigh. Inverse preference learning: Preference-based rl without a reward
   function. Advances in Neural Information Processing Systems, 36, 2024.
  - Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. Advances in neural information processing systems, 29, 2016.
- Lee N. Thorne J. Hong, J. Orpo: Monolithic preference optimization without reference model.
   *ArXiv Preprint ArXiv:2403.07691*, 2(5), 2024.
- Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Chrisina Jayne. Imitation learning: A survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35, 2017.
- Daniel Jarrett, Ioana Bica, and Mihaela van der Schaar. Strictly batch imitation learning by energy based distribution matching. *Advances in Neural Information Processing Systems*, 33:7354–7365, 2020.
- Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. A survey of reinforcement learning from human feedback. *arXiv preprint arXiv:2312.14925*, 2023.
- Geon-Hyeong Kim, Seokin Seo, Jongmin Lee, Wonseok Jeon, HyeongJoo Hwang, Hongseok Yang, and Kee-Eung Kim. Demodice: Offline imitation learning with supplementary imperfect demonstrations. In *International Conference on Learning Representations*, 2022.

- Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Bahare Kiumarsi, Kyriakos G Vamvoudakis, Hamidreza Modares, and Frank L Lewis. Optimal
   and autonomous control using reinforcement learning: A survey. *IEEE transactions on neural networks and learning systems*, 29(6):2042–2062, 2017.
- Jens Kober, J Andrew Bagnell, and Jan Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- Ilya Kostrikov, Kumar Krishna Agrawal, Debidatta Dwibedi, Sergey Levine, and Jonathan Tompson. Discriminator-actor-critic: Addressing sample inefficiency and reward bias in adversarial imitation learning. *arXiv preprint arXiv:1809.02925*, 2018.
- Ilya Kostrikov, Ofir Nachum, and Jonathan Tompson. Imitation learning via off-policy distribution
   matching. *arXiv preprint arXiv:1912.05032*, 2019.
- Pawel Ladosz, Lilian Weng, Minwoo Kim, and Hyondong Oh. Exploration in deep reinforcement learning: A survey. *Information Fusion*, 85:1–22, 2022.
- Kimin Lee, Laura Smith, and Pieter Abbeel. Pebble: Feedback-efficient interactive rein forcement learning via relabeling experience and unsupervised pre-training. *arXiv preprint arXiv:2106.05091*, 2021.
- 4 Yicheng Luo, Zhengyao Jiang, Samuel Cohen, Edward Grefenstette, and Marc Peter Deisenroth.
   4 Optimal transport for offline imitation learning. *arXiv preprint arXiv:2303.13971*, 2023.
- Ma X. Wan L. Liu R. Li X. Lu Z. Lyu, J. Seabo: A simple search-based method for offline imitation
   learning. *ArXiv Preprint ArXiv:2402.03807*, 2024.
- Yecheng Ma, Andrew Shen, Dinesh Jayaraman, and Osbert Bastani. Versatile offline imitation from observations and examples via regularized state-occupancy matching. In *International Conference on Machine Learning*, pp. 14639–14663. PMLR, 2022.
- Liyuan Mao, Haoran Xu, Weinan Zhang, and Xianyuan Zhan. Odice: Revealing the mystery of dis tribution correction estimation via orthogonal-gradient update. *arXiv preprint arXiv:2402.00348*, 2024.
- Yu Meng, Mengzhou Xia, and Danqi Chen. Simple preference optimization with a reference-free reward. *arXiv preprint arXiv:2405.14734*, 2024.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Belle mare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level
   control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-low instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.
- Arka Pal, Deep Karkhanis, Samuel Dooley, Manley Roberts, Siddartha Naidu, and Colin White.
  Smaug: Fixing failure modes of preference optimisation with dpo-positive. *arXiv preprint arXiv:2402.13228*, 2024.
- Jongjin Park, Younggyo Seo, Jinwoo Shin, Honglak Lee, Pieter Abbeel, and Kimin Lee. Surf:
   Semi-supervised reward learning with data augmentation for feedback-efficient preference-based
   reinforcement learning. *arXiv preprint arXiv:2203.10050*, 2022.
- Dean A Pomerleau. Efficient training of artificial neural networks for autonomous navigation. *Neural computation*, 3(1):88–97, 1991.
- Rafael Figueiredo Prudencio, Marcos ROA Maximo, and Esther Luna Colombini. A survey on offline reinforcement learning: Taxonomy, review, and open problems. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.

682

683

684

685

686

687

688 689

690

691

692

696

- 648 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea 649 Finn. Direct preference optimization: Your language model is secretly a reward model. Advances 650 in Neural Information Processing Systems, 36, 2024. 651 652 Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In Proceedings of the fourteenth international con-653 ference on artificial intelligence and statistics, pp. 627-635. JMLR Workshop and Conference 654 Proceedings, 2011. 655 656 Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon 657 Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari, 658 go, chess and shogi by planning with a learned model. Nature, 588(7839):604-609, 2020. 659 660 Kaustubh Sridhar, Souradeep Dutta, Dinesh Jayaraman, James Weimer, and Insup Lee. Memory-661 consistent neural networks for imitation learning. arXiv preprint arXiv:2310.06171, 2023. 662 Mingfei Sun, Anuj Mahajan, Katja Hofmann, and Shimon Whiteson. Softdice for imitation learning: 663 Rethinking off-policy distribution matching. arXiv preprint arXiv:2106.03155, 2021. 664 665 Richard S Sutton. Reinforcement learning: An introduction. A Bradford Book, 2018. 666 667 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. 668 In 2012 IEEE/RSJ international conference on intelligent robots and systems, pp. 5026–5033. 669 IEEE, 2012. 670 671 Zhichao Wang, Bin Bi, Shiva Kumar Pentyala, Kiran Ramnath, Sougata Chaudhuri, Shubham 672 Mehrotra, Xiang-Bo Mao, Sitaram Asur, et al. A comprehensive survey of Ilm alignment tech-673 niques: Rlhf, rlaif, ppo, dpo and more. arXiv preprint arXiv:2407.16216, 2024. 674 Christian Wirth, Riad Akrour, Gerhard Neumann, and Johannes Fürnkranz. A survey of preference-675 based reinforcement learning methods. Journal of Machine Learning Research, 18(136):1-46, 676 2017. 677 678 Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton 679 Murray, and Young Jin Kim. Contrastive preference optimization: Pushing the boundaries of llm 680 performance in machine translation. arXiv preprint arXiv:2401.08417, 2024. 681
  - Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. Rrhf: Rank responses to align language models with human feedback without tears. *arXiv preprint arXiv:2304.05302*, 2023.
  - Sheng Yue, Guanbo Wang, Wei Shao, Zhaofeng Zhang, Sen Lin, Ju Ren, and Junshan Zhang. Clare: Conservative model-based reward learning for offline inverse reinforcement learning. *arXiv* preprint arXiv:2302.04782, 2023.
  - 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.
- Wenjia Zhang, Haoran Xu, Haoyi Niu, Peng Cheng, Ming Li, Heming Zhang, Guyue Zhou, and
   Xianyuan Zhan. Discriminator-guided model-based offline imitation learning. In *Conference on Robot Learning*, pp. 1266–1276. PMLR, 2023.
- Joshi R. Liu T. Khalman M. Saleh M. Liu P. Zhao, Y. Slic-hf: Sequence likelihood calibration with
   human feedback. *ArXiv Preprint ArXiv:2305.10425*, 2023.
- Konrad Zolna, Alexander Novikov, Ksenia Konyushkova, Caglar Gulcehre, Ziyu Wang, Yusuf Aytar, Misha Denil, Nando de Freitas, and Scott Reed. Offline learning from demonstrations and unlabeled experience. *arXiv preprint arXiv:2011.13885*, 2020.

## A ADDITIONAL EXPERIMENTS ON ADROIT TASKS

To further demonstrate the scalability and robustness of POIL in more complex environments, we conducted additional experiments on the Adroit tasks from the D4RL benchmark (Fu et al., 2020).

707 A.1 EXPERIMENTAL SETUP

We focused on the expert datasets within the Adroit tasks, following the same settings as in MCNN (Sridhar et al., 2023). This ensures consistency and allows for a direct comparison between POIL and MCNN under similar conditions. All baseline results reported in our experiments are directly inherited from MCNN. Specifically, we evaluated POIL on the following tasks: pen-expert-v1, hammer-expert-v1, door-expert-v1, and relocate-expert-v1. For these experi-ments, we set the regularization coefficient  $\lambda = 0$  and the scaling factor  $\beta = 1$  based on our findings from the MuJoCo experiments. We trained our models using the same neural network architecture and optimization settings as in the main experiments to ensure fairness. 

A.2 RESULTS

The performance of POIL on the Adroit tasks is presented in Tables 3 and 4. In Table 3, we compare POIL against the methods reported in the MCNN paper when using the full expert dataset for each task. In Table 4, we evaluate the performance of POIL with varying numbers of demonstrations to assess its data efficiency.

Table 3: Performance comparison on Adroit tasks using the full expert dataset (5,000 demonstrations). The results are averaged over three runs with different random seeds. (The bold numbers indicate the best performance, while the underscored numbers are the second best.)

Task	BC	BeT-BC	Implicit BC	MCNN (Fixed)	MCNN (Tuned)	POIL
Pen	2633	$1853\pm117$	$2586\pm65$	$3285\pm209$	$3405\pm328$	$4077 \pm 66$
Hammer	16140	$2731\pm261$	$-132\pm25$	$16027\pm382$	$\textbf{16387} \pm \textbf{682}$	$16295\pm49$
Door	969	$356\pm35$	$361\pm67$	$3033\pm0.3$	$3035\pm7$	$\textbf{3040} \pm \textbf{12}$
Relocate	4289	$490\pm42$	-	$4566\pm47$	$\overline{4566\pm47}$	$\textbf{4606} \pm \textbf{45}$

As shown in Table 3, POIL outperforms all methods reported in the MCNN paper on pen-expert-v1, door-expert-v1, and relocate-expert-v1. On hammer-expert-v1, POIL achieves performance comparable to the best method (MCNN+MLP with tuned hyperparameters), with only a slight difference.

Table 4: Performance comparison on Adroit tasks with varying numbers of demonstrations. Results are averaged over three runs with different random seeds. (The bold numbers indicate the best performance for each dataset size, while the underscored numbers are the second best.)

Task	Demos	MCNN+MLP	POIL
	100	$2725\pm139$	$3015\pm3$
	500	$2931\pm32$	$\textbf{3025} \pm \textbf{1}$
Deer	1000	$\overline{2992\pm18}$	$\textbf{3027} \pm \textbf{7}$
Door	2000	$\overline{3017\pm10}$	$\textbf{3028} \pm \textbf{22}$
	4000	$\overline{3025\pm3}$	$\textbf{3033} \pm \textbf{3}$
	5000	$3035 \pm 7$	$\textbf{3041} \pm \textbf{12}$
	100		$4146 \pm 104$
	500	$3712\pm32$	$\textbf{4127} \pm \textbf{8}$
Pen	1000	$\overline{3808\pm6}$	$4141 \pm 38$
	2000	$\overline{3858\pm29}$	$\textbf{4197} \pm \textbf{136}$
	4000	$\overline{3934 \pm 42}$	$\textbf{4172} \pm \textbf{113}$
	5000	$\overline{4051\pm195}$	$\textbf{4078} \pm \textbf{66}$

Table 4 demonstrates that POIL consistently achieves superior performance across all tasks and dataset sizes. Notably, POIL shows strong performance even with a limited number of demonstra-tions (e.g., 100 demos), highlighting its data efficiency and robustness in data-scarce scenarios. 

#### Β ADDITIONAL RESULTS WITH ERROR BARS

In this appendix, we provide detailed results for the single demonstration learning experiments presented in Section 4.2, including standard deviations (error bars) over multiple runs. These results offer a more comprehensive view of the performance variability across different methods.

Table 5: Performance of BC, IQ-Learn, and DMIL trained with a single expert demonstration on MuJoCo tasks. The results are averaged over 3 runs with different random seeds, and the scores represent the average returns  $\pm$  standard deviation over the last ten epochs. The bold numbers represent the best performance, while the underscored numbers are the second best.

Environment	Traj.	BC	IQ-Learn	DMIL
	traj1	2775.39 ± 296.23	3621.73 ± 423.67	2764.81 ± 315.29
HalfCheetah	traj2	$3031.12 \pm 727.32$	$2957.77 \pm 803.21$	$3546.40 \pm 521.78$
	traj3	2927.38 ± 967.15	$3805.18 \pm 567.81$	$2802.43 \pm 648.34$
	traj1	$1548.32 \pm 108.74$	2431.73 ± 698.13	$2947.13 \pm 267.93$
Hopper	traj2	$1687.93 \pm 70.75$	$2601.54 \pm 532.47$	3195.28 ± 501.47
	traj3	$1850.54 \pm 870.87$	2520.18 ± 746.91	3301.41 ± 687.12
	traj1	881.23 ± 68.44	1011.89 ± 391.27	1942.76 ± 321.15
Walker2d	traj2	$1023.85 \pm 124.71$	$1062.52 \pm 287.91$	$1686.98 \pm 612.37$
	traj3	$935.63 \pm 76.50$	997.21 ± 159.43	1864.76 ± 491.72

Table 6: Performance of O-DICE and POIL trained with a single expert demonstration on MuJoCo tasks. The results are averaged over 3 runs with different random seeds, and the scores represent the average returns  $\pm$  standard deviation over the last ten epochs. The bold numbers represent the best performance, while the underscored numbers are the second best.

Environment	Traj.	O-DICE	<b>POIL</b> $_{\lambda=1}$	$\mathbf{POIL}_{\lambda=0}$
	traj1	$1532.63 \pm 716.74$	$3843.75 \pm 379.47$	4628.88 ± 183.47
HalfCheetah	traj2	$1801.12 \pm 697.67$	$2228.04 \pm 797.59$	$4833.13 \pm 30.95$
	traj3	$1032.74 \pm 1104.10$	$2645.44 \pm 876.50$	4309.46 ± 247.95
	traj1	$1672.45 \pm 441.29$	$2426.22 \pm 622.98$	3501.67 ± 469.42
Hopper	traj2	$1748.63 \pm 204.60$	$2534.35 \pm 398.42$	$3066.00 \pm 158.89$
	traj3	$1797.43 \pm 145.06$	$3372.09 \pm 52.41$	3499.08 ± 58.08
	traj1	$1543.67 \pm 445.67$	$2645.29 \pm 840.42$	4722.79 ± 521.90
Walker2d	traj2	1392.54 ± 1106.87	$1826.08 \pm 490.14$	5233.93 ± 234.21
	traj3	$1483.52 \pm 716.27$	$2493.53 \pm 729.79$	3745.50 ± 445.90