

506 Appendix

507 We divide the appendix into four different sections following the results section. Each section
 508 additionally provides hyper-parameters used for IPL in that section. The first section, setup, contains
 509 details information on the experimental setup and hyper-parameters used. The second section on
 510 benchmark results gives full learning curves for the experiments in Section 4.2. The third section
 511 provides full learning curves for the MetaWorld and Data-scaling experiments. The final Appendix
 512 section provides extended ablations.

513 A Setup

514 Here we provide the full algorithmic outline of IPL using Implicit Q-Learning [27] that
 515 mimics our implementation. While in practice the policy π could be extracted at the
 516 end of training, we do it simultaneously as in [27] in order to construct learning curves.

Algorithm 2: IPL Algorithm (IQL Variant)

Input : $\mathcal{D}_p, \mathcal{D}_o, \lambda, \alpha$
for $i = 1, 2, 3, \dots$ **do**
 517 Sample batches $B_p \sim \mathcal{D}_p, B_o \sim \mathcal{D}_o$
 Update Q : $\min_Q \mathbb{E}_{B_p} [\mathcal{L}_p(Q)] + \lambda \mathbb{E}_{B_p \cup B_o} [\mathcal{L}_r(Q)]$
 Update V : $\min_V \mathbb{E}_{B_p \cup B_o} [|\tau - 1(Q(s, a) - V(s))| (Q(s, a) - V(s))^2]$
 Update π : $\max_\pi \mathbb{E}_{\mathcal{D}_p \cup \mathcal{D}_o} [e^{\beta(Q(s, a) - V(s))} \log \pi(a|s)]$

518 Note that above we write the temperature parameter β as done in IQL, instead of how it is usually
 519 done, using α in the denominator [18, 42].

520 When sampling batches of preference data $B_p \sim \mathcal{D}_p$, we take sub-samples of each segment σ of
 521 length s . For a sampled data point $(\sigma^{(1)}, \sigma^{(2)}, y)$, we sample $\text{start} \sim \text{Unif}[0, 1, 2, \dots, k - s]$ and then
 522 let take $\sigma = s_{\text{start}}, a_{\text{start}}, \dots, s_{\text{start}+s}$. We use the same start value across the entire batch.

523 Given that we run experiments using MLPs, all of our experiments were run on CPU compute
 524 resources. Each seed for each method requires one CPU core (two hyper-threads) and 8 Gb of
 525 memory.

526 B Benchmark Results

527 Here we provide details for our experiments on the preference-based RL benchmark from Kim et al.
 528 [25]. We use the same hyperparameters as Kim et al. [25] and Kostrikov et al. [27] where applicable
 529 as shown in Table 4.

530 **Gym-Mujoco Locomotion.** Hopper and Walker2D agents are tasked with learning locomotion
 531 policies from datasets of varying qualities taken from the D4RL [16] benchmark. Preference datasets
 532 were constructed by Kim et al. [25] by uniformly sampling queries and labeling a subset of them. For
 533 all locomotion tasks the segment length. Preference datasets for “medium” quality offline datasets
 534 contain 500 queries, while preference datasets for “expert” quality offline datasets contain 100 queries.
 535 Segment lengths $k = 100$ for all datasets, and were subsampled to length $s = 64$ by IPL and our MR
 536 (reimplementation). Evaluation was performed over 10 episodes every 5000 steps. Full learning
 537 curves are shown in Fig. 3.

538 **RoboMimic.** The RoboMimic datasets contain interaction data of two types: ph — proficient human
 539 and mh — multihuman. The multi-human data was collected from human demonstrators of mixed
 540 quality. The robot is tasked with learning how to lift a cube (lift) or pick and place a can (can).
 541 Preference datasets were again taken directly from Kim et al. [25]. Preference datasets of size 100
 542 with segment lengths $k = 50$, randomly sub-sampled to length $s = 32$ were used for the ph datasets.
 543 Preference datasets of size 500 with segment lengths $k = 100$, randomly sub-sampled to length $s = 64$
 544 were used for the mh datasets. Evaluation was performed over 25 episodes every 50000 steps. Full
 545 learning curves are shown in Fig. 4.

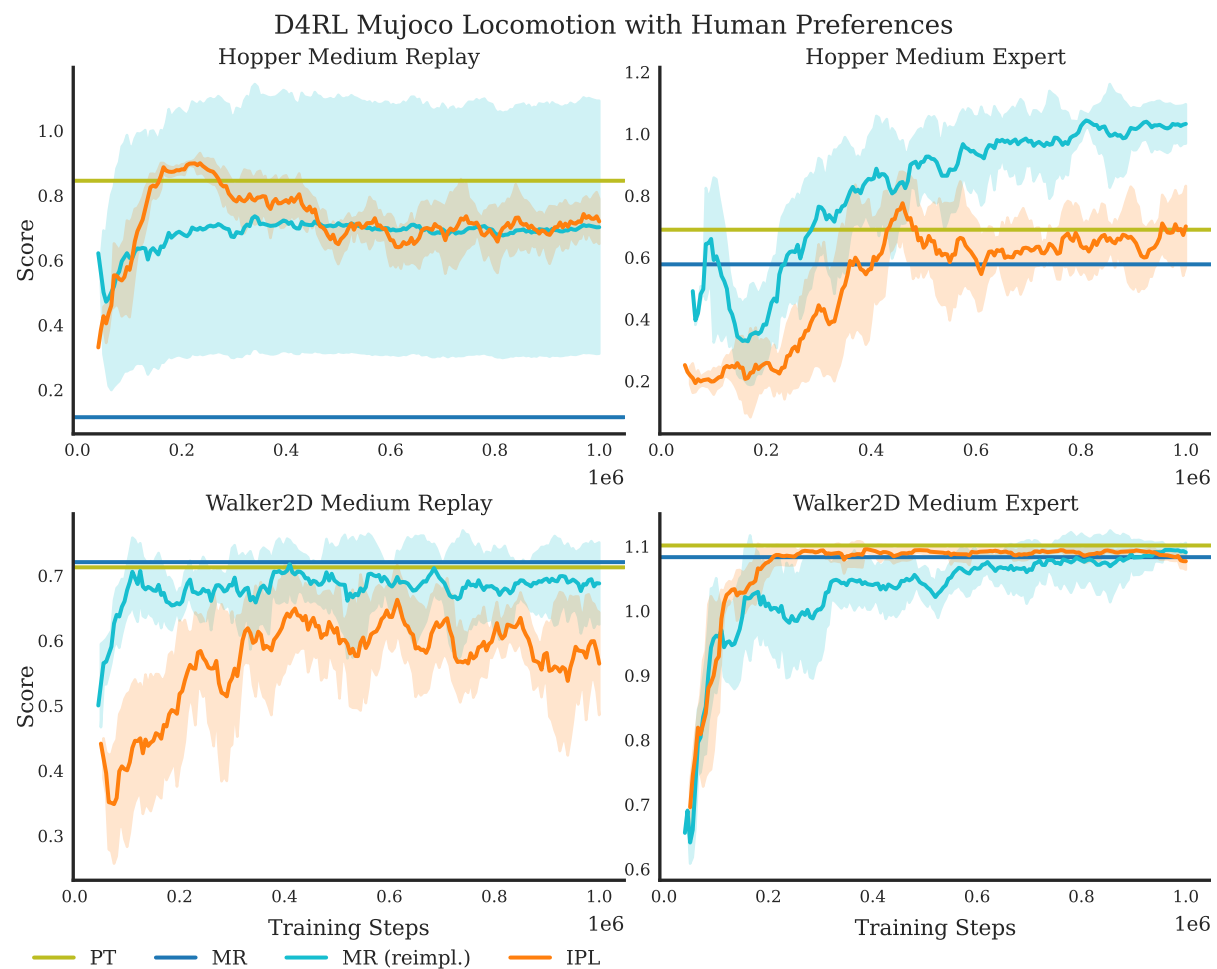


Figure 3: Full learning curves on the D4RL locomotion benchmark with human preferences.

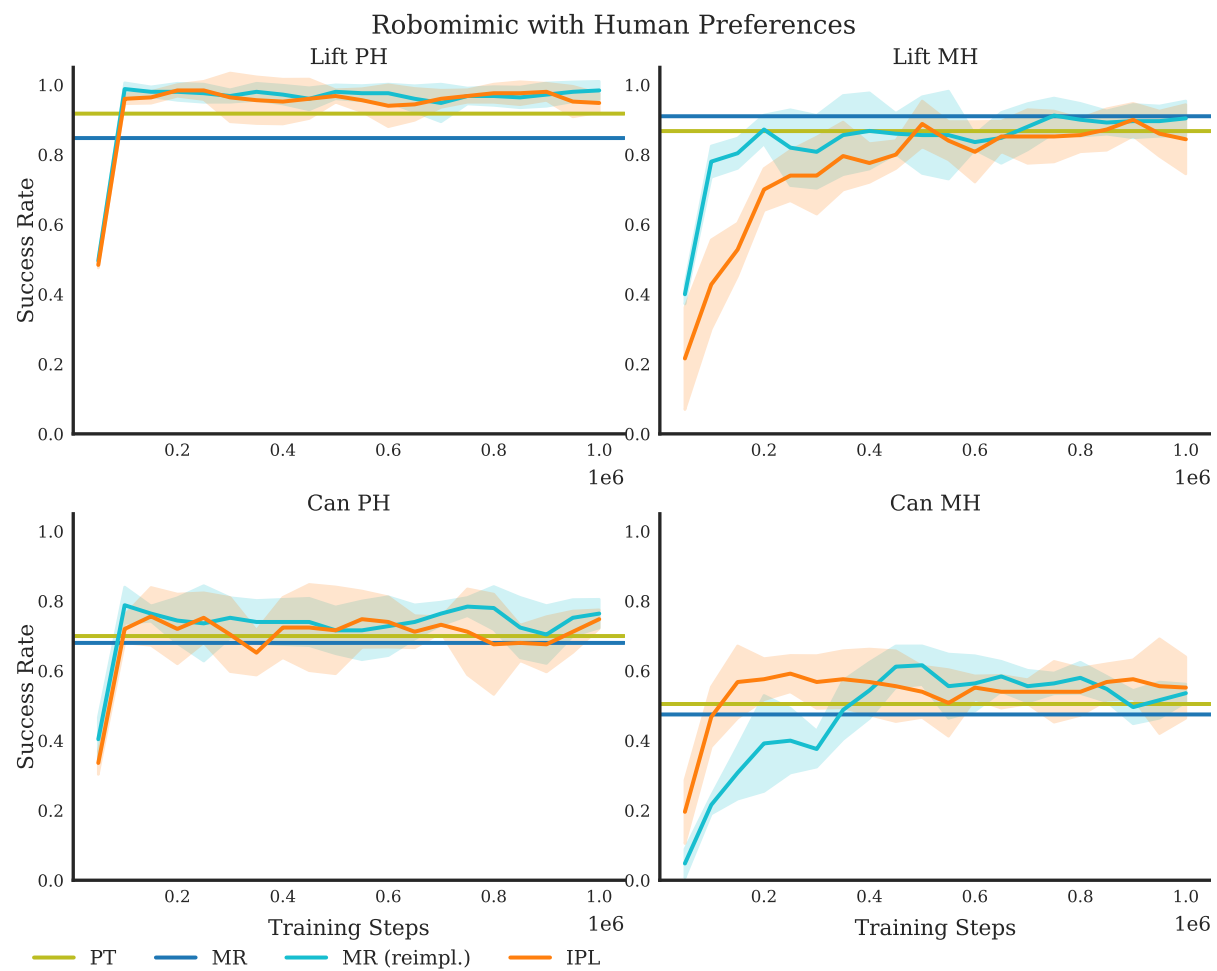


Figure 4: Full learning curves on the RoboMimic benchmark with human preferences.

Common Hyperparameters			MR Hyperparameters		
Parameter	Locomotion	Robomimic	Parameter	Locomotion	Robomimic
Q, V, π Arch	2x 256d	2x 256d	r_θ Arch	2x 256d	2x 256d
Learning Rate	0.0003	0.0003	r_θ LR	0.0003	0.0003
Optimizer	Adam	Adam	r_θ Optimizer	Adam	Adam
β	3.0	0.5	r_θ Steps	20k	20k
τ	0.7	0.7			
\mathcal{D}_o Batch Size	256	256			
\mathcal{D}_p Batch Size	8	8			
Training Steps	1 Mil	1 Mil			
k	100	100, 50			
Subsample s	64	64, 32			

IPL Hyperparameters		
Parameter	Locomotion	Robomimic
λ	0.5	4

Table 4: Hyperparameters used for the benchmark experiments. We can see that IPL has fewer hyperparameters.

Common Hyperparameters		MR Hyperparameters	
Parameter	Value	Parameter	Value
Q, V, π Arch	3x 256d	r_θ Arch	3x 256d
Learning Rate	0.0003	r_θ LR	0.0003
Optimizer	Adam	r_θ Optimizer	Adam
β	4.0	r_θ Steps	20k
τ	0.7		
\mathcal{D}_p Batch Size	16		
Training Steps	200k		
k	25		
Subsample s	16		

IPL Hyperparameters	
Parameter	Locomotion
λ	0.5

Table 5: Hyper-parameters used in the MetaWorld data scaling experiments.

546 C Data Scaling Results

547 Experiments for data scaling were conducted on the MetaWorld benchmark from Yu et al. [52]. Offline
548 datasets for five different MetaWorld tasks were constructed as follows: Collect 100 trajectories of
549 expert data on the target task using the built in ground truth policies with the addition of Gaussian
550 noise of standard deviation 1.0. Collect 100 trajectories of sub-optimal data by running the ground-
551 truth policy for a different randomization of the target task with Gaussian noise 1.0. Collect 100
552 trajectories of even more sub-optimal data by running the ground truth policy *of a different task*
553 with Gaussian noise standard deviation 1.0 in the target domain. Finally, collect 100 trajectories
554 with uniform random actions. As MetaWorld episodes are 500 steps long, this results in 200,000
555 time-steps of data. We then construct preference datasets by uniformly sampling segments from
556 the offline dataset and assigning labels y according to $\sum_t r(s_t^{(1)}, a_t^{(1)}) > \sum_t r(s_t^{(2)}, a_t^{(2)})$ where r is
557 the ground truth reward provided by metaworld. We then train using only the data from \mathcal{D} General
558 architecture hyper-parameters were taken from Lee et al. [30], Hejna and Sadigh [21] which also use
559 the MetaWorld benchmark, but for online preference-based RL. Full-hyper parameters are shown in
560 Table 5. We run 20 evaluation episodes every 2500 steps. Full learning curves are shown in
561 Fig. 5. When reporting values in Table 2, we choose the maximum point on the learning curves which
562 average across five seeds. This provides results as if early stopping was given by an oracle, which is
563 less optimistic than averaging the maximum of each seed as done in Mandlekar et al. [34].

564 D Ablations

565 In this section we provide additional ablations on both the benchmark datasets and MetaWorld
566 datasets. We keep the hyperparameters the same, except for the parameter-efficient experiments. We
567 run hyper-parameter sensitivity results for the human-preference benchmark datasets in Fig. 6. The
568 top row depicts the sensitivity for IPL to the value of λ . The bottom row depicts the sensitivity of MR
569 to the number of timesteps the reward function is trained for.

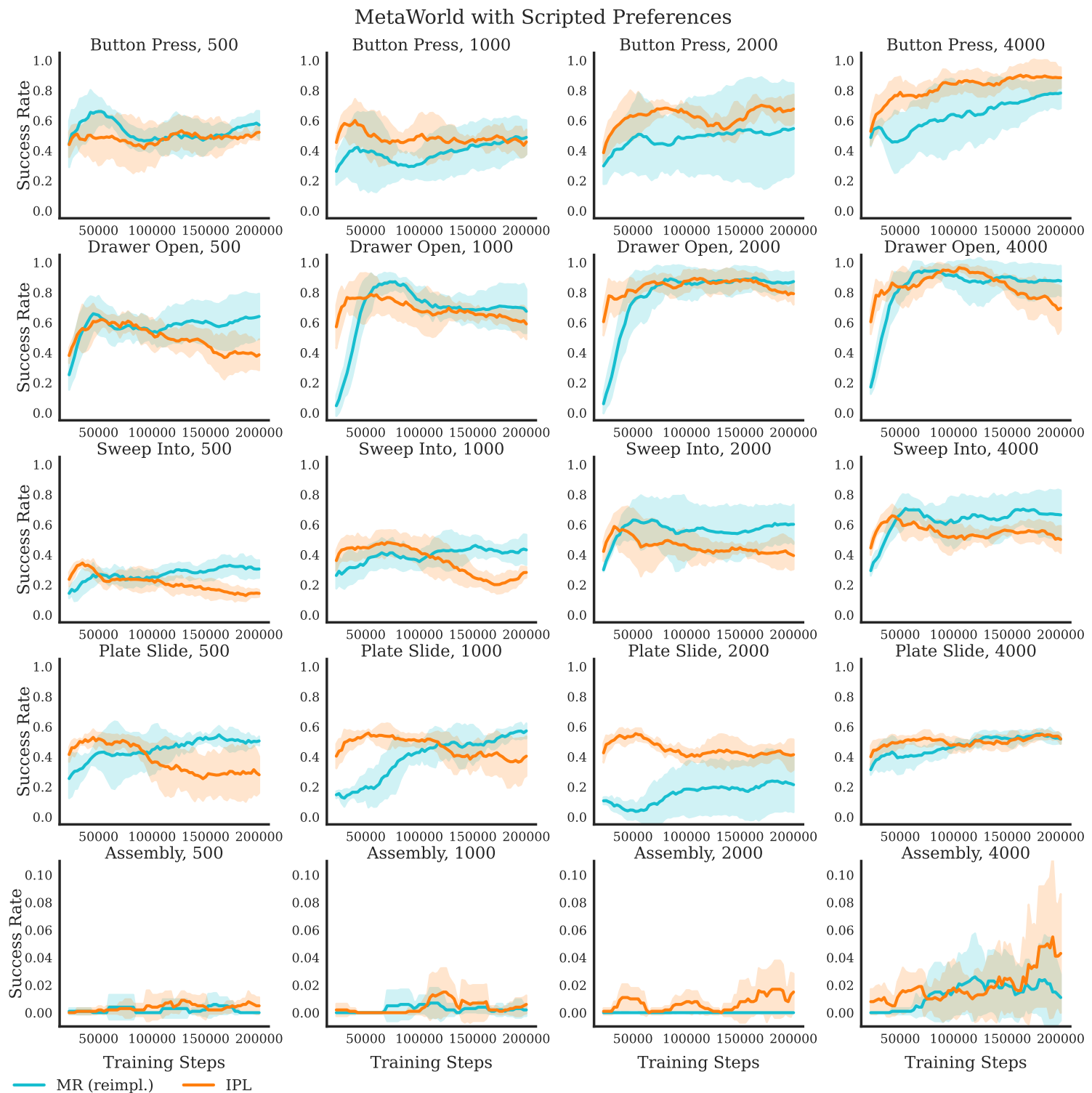


Figure 5: Full learning curves for the MetaWorld data scaling results with scripted preferences.

570 For the parameter-efficient experiments *only* we use an efficient version of IPL based on AWAC
571 [37] to additionally remove the value network. The outline of this variant is given below

Algorithm 3: IPL Algorithm (AWAC Variant)

Input: $\mathcal{D}_p, \mathcal{D}_o, \lambda, \alpha$
for $i = 1, 2, 3, \dots$ **do**
572 Sample batches $B_p \sim \mathcal{D}_p, B_o \sim \mathcal{D}_o$
 Estimate V as $Q(s, \pi(s))$
 Update Q : $\min_Q \mathbb{E}_{B_p} [\mathcal{L}_p(Q)] + \lambda \mathbb{E}_{B_p \cup B_o} [\mathcal{L}_r(Q)]$
 Update π : $\max_{\pi} \mathbb{E}_{\mathcal{D}_p \cup \mathcal{D}_o} [e^{\beta(Q(s,a) - Q(s, \pi(s)))} \log \pi(a|s)]$

573 For this version of IPL, we use $\lambda = 0.5$. All other hyper-parameters remain the same as in Table 6
574 except the architectures. For the parameter-efficiency experiments only we use MLPs consisting of
575 two dense layers with either dimension 64 or dimension 35. Running MR with a two-layer MLP of
576 dimension 35 has almost exactly the same number of parameters as IPL-AWAC with two-layer MLPs
577 of dimension 64. We include full results for the parameter-efficiency experiments in Table 6. We
578 find that on Drawer Open and Sweep Into, IPL outperforms both MR (64) and MR (35). In these
579 environments, performance increases from MR (35) to MR (64) indicating that the expressiveness
580 of the Q -function and policy are limiting performance. For the same budget, IPL is able to perform

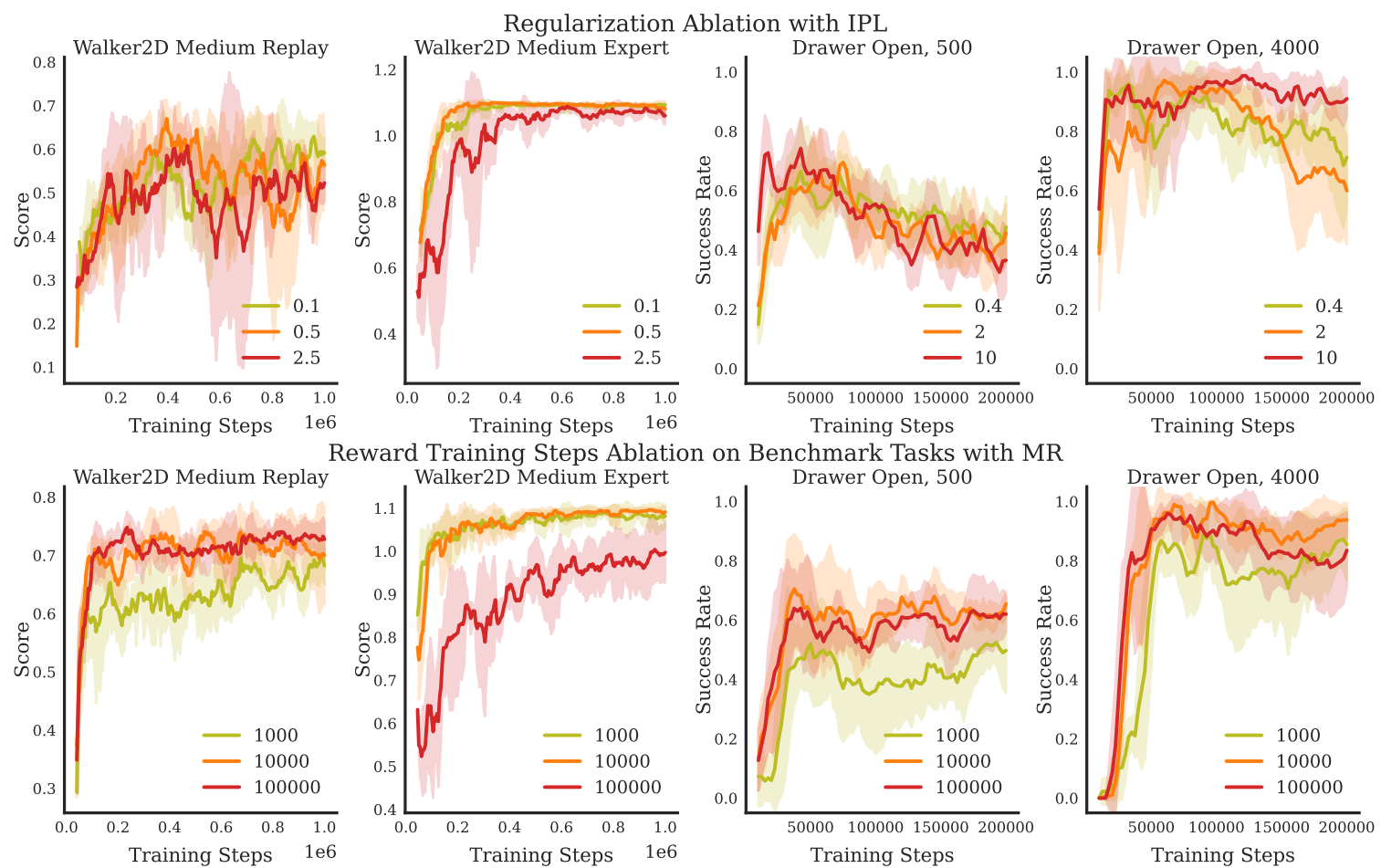


Figure 6: Ablations on regularization strength λ for IPL (top row) and the number of reward steps for MR (bottom row). We see that IPL is relatively consistent across different values of λ . MR on the other hand, can vary greatly if the reward function under or over fits. In Walker2D Medium Replay and Drawer Open, 500, we see that it can easily under-fit. In Walker2D Medium Expert it easily over-fits.

581 better. In Button Press, the simplest task, we find that MR (64) actually over-fits more than MR (35)
582 and MR (64) ends up performing worse. In Plate Slide, all methods perform similarly independent of
583 parameter count. We omit Assembly because of its low success rate at all data scales.

Preference Queries		500	1000	2000	4000
Button Press	MR (35)	73.9 \pm 8.9	86.8 \pm 8.2	89.9 \pm 14.4	99.0 \pm 1.0
	MR (64)	54.2 \pm 16.1	42.6 \pm 33.0	67.1 \pm 14.9	43.4 \pm 7.4
	IPL (64)	65.8 \pm 13.3	79.8 \pm 18.1	80.0 \pm 17.3	95.8 \pm 5.2
Drawer Open	MR (35)	13.4 \pm 13.9	12.6 \pm 21.9	15.5 \pm 20.1	18.4 \pm 25.6
	MR (64)	13.4 \pm 19.0	57.1 \pm 31.2	54.5 \pm 31.7	78.8 \pm 12.2
	IPL (64)	89.8 \pm 11.3	93.2 \pm 2.5	99.5 \pm 0.9	95.5 \pm 3.7
Sweep Into	MR (35)	35.1 \pm 8.9	42.4 \pm 9.9	45.9 \pm 9.6	35.9 \pm 4.1
	MR (64)	31.1 \pm 6.4	55.8 \pm 5.9	49.6 \pm 10.3	56.4 \pm 10.3
	IPL (64)	41.1 \pm 14.2	63.9 \pm 8.0	65.0 \pm 12.0	63.9 \pm 11.8
Plate Slide	MR (35)	55.2 \pm 6.1	51.1 \pm 4.4	53.0 \pm 2.0	48.9 \pm 3.3
	MR (64)	46.6 \pm 21.9	50.8 \pm 0.6	47.0 \pm 2.5	48.5 \pm 4.6
	IPL (64)	54.9 \pm 3.2	49.4 \pm 1.6	45.2 \pm 9.0	48.8 \pm 4.9

Table 6: Performance of different methods on the MetaWorld tasks under a limited parameter budget. MR (35) and IPL (64) have the same number of parameters. The Assembly task is omitted due to low success rate. On Button Press, fewer parameters appears to perform better as, due to the simplicity of the task, its easier for the bigger models to overfit. On Drawer Open and Sweep Into, we see consistent gains from increasing the number of parameters in the network, and IPL performs best overall. On the Plate Slide task, all methods at different parameter scales perform similarly.