

456 **A Proofs**

457 **A.1 Proof for Proposition 1**

458 For each  $k$ , we denote the stationary point of gradient descent as  $x_k^* = \lim_{n \rightarrow \infty} x_k^n$ . Using the  
 459 first-order optimality condition at the stationary point, we know that

$$\begin{aligned}
 0 &= \nabla_x d_{\sigma_k}(x; \mathcal{D})^2 \Big|_{x=x_k^*} \\
 &= -\nabla_x \sigma_k^2 \log \left[ \sum_i \exp \left[ -\frac{1}{2\sigma_k^2} \|x - x_i\|^2 \right] \right] \Big|_{x=x_k^*} \\
 &= \frac{\sigma_k^2 \sum_i \exp \left[ -\frac{1}{2\sigma_k^2} \|x - x_i\|^2 \right] \sigma_k^{-2} (x - x_i)}{\sum_i \exp \left[ -\frac{1}{2\sigma_k^2} \|x - x_i\|^2 \right]} \Big|_{x=x_k^*}.
 \end{aligned} \tag{13}$$

460 We then have that

$$\begin{aligned}
 \sum_i \exp \left[ -\frac{1}{2\sigma_k^2} \|x_k^* - x_i\|^2 \right] (x_k^* - x_i) &= 0 \\
 \left( \sum_i \exp \left[ -\frac{1}{2\sigma_k^2} \|x_k^* - x_i\|^2 \right] \right) x_k^* &= \sum_i \exp \left[ -\frac{1}{2\sigma_k^2} \|x_k^* - x_i\|^2 \right] x_i \\
 x_k^* &= \sum_i \frac{\exp \left[ -\frac{1}{2\sigma_k^2} \|x_k^* - x_i\|^2 \right]}{\sum_n \exp \left[ -\frac{1}{2\sigma_k^2} \|x_k^* - x_n\|^2 \right]} x_i.
 \end{aligned} \tag{14}$$

461 Let  $x_j \in \operatorname{argmin}_{x_i \in \mathcal{D}} \|x_k^* - x_i\|^2$  and  $\Delta_i(\sigma_k) = \sigma_k^{-2} (\|x_k^* - x_i\|^2 - \|x_k^* - x_j\|^2)$ .

462 **Lemma 4.** *When there is a unique minimizer  $x_j$ ,  $\lim_{k \rightarrow \infty} x_k^* = x_j$ .*

463 *Proof.* Since  $x_j$  is the unique closest data point to  $x_k^*$ ,  $\Delta_i > 0$  for  $i \neq j$  and  $\Delta_j = 0$ . With the  
 464 monotonically decreasing  $\sigma_k \rightarrow 0$ , we know that

$$\lim_{k \rightarrow \infty} \Delta_i(\sigma_k) = \begin{cases} 0 & \text{if } i = j \\ \infty & \text{otherwise.} \end{cases} \tag{15}$$

465 Therefore,

$$\lim_{k \rightarrow \infty} \exp \left[ -\frac{1}{2} \Delta_i(\sigma_k) \right] = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise.} \end{cases} \tag{16}$$

466 We then obtain

$$\begin{aligned}
 \lim_{k \rightarrow \infty} \frac{\exp \left[ -\frac{1}{2\sigma_k^2} (\|x_k^* - x_i\|^2 - \|x_k^* - x_j\|^2) \right]}{\sum_n \exp \left[ -\frac{1}{2\sigma_k^2} (\|x_k^* - x_n\|^2 - \|x_k^* - x_j\|^2) \right]} &= \lim_{k \rightarrow \infty} \frac{\exp \left[ -\frac{1}{2} \Delta_i(\sigma_k) \right]}{\sum_{n \neq j} \exp \left[ -\frac{1}{2} \Delta_n(\sigma_k) \right] + 1} \\
 &= \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise.} \end{cases}
 \end{aligned} \tag{17}$$

467 Therefore, by combining (14) and (17), we have that

$$\lim_{k \rightarrow \infty} x_k^* = x_j. \tag{18}$$

468 □

469 **Lemma 5.** *When there are multiple minimizers  $x_{j_1}, \dots, x_{j_m} \in \operatorname{argmin}_{x_i \in \mathcal{D}} \|x_k^* - x_i\|^2$ ,*  
 470  $\lim_{k \rightarrow \infty} x_k^* = \frac{1}{m} \sum_{l=1}^m x_{j_l}$ .

471 *Proof.* In contrast to (17) for the unique  $x_j$ , we have that

$$\begin{aligned}
& \lim_{k \rightarrow \infty} \frac{\exp \left[ -\frac{1}{2} (\|x_k^* - x_i\|^2 - \|x_k^* - x_{j_1}\|^2) \right]}{\sum_n \exp \left[ -\frac{1}{2} (\|x_k^* - x_n\|^2 - \|x_k^* - x_{j_1}\|^2) \right]} \\
&= \lim_{k \rightarrow \infty} \frac{\exp \left[ -\frac{1}{2} \Delta_i(\sigma_k) \right]}{\sum_{n \notin \{j_1, \dots, j_m\}} \exp \left[ -\frac{1}{2} \Delta_n(\sigma_k) \right] + m} \\
&= \begin{cases} \frac{1}{m} & \text{if } i \in \{j_1, \dots, j_m\} \\ 0 & \text{otherwise.} \end{cases}
\end{aligned} \tag{19}$$

472 Therefore,  $x_k^*$  converges to the geometric center of all the minimizers:

$$\lim_{k \rightarrow \infty} x_k^* = \frac{1}{m} \sum_{l=1}^m x_{j_l}. \tag{20}$$

473

□

474 **Lemma 6.** For a local minimizer  $x_k^*$ ,  $x_j$  is unique; for a local maximizer or saddle point  $x_k^*$ ,  $x_j$  is  
475 not unique.

476 *Proof.* We first prove that  $x_j$  is unique if  $x_k^*$  is a local minimizer by contradiction. If there are  
477 multiple minimizers  $x_{j_1}, \dots, x_{j_m} \in \operatorname{argmin}_{x_i \in \mathcal{D}} \|x_k^* - x_i\|^2$ , we have that

$$\begin{aligned}
\lim_{k \rightarrow \infty} d_{\sigma_k}(x_k^*; \mathcal{D})^2 &= \lim_{k \rightarrow \infty} -\sigma_k^2 \log \left[ \sum_i \exp \left[ -\frac{1}{2\sigma_k^2} \|x_k^* - x_i\|^2 \right] \right] \\
&= -\log \left[ \lim_{k \rightarrow \infty} \left( \sum_i \exp \left[ -\frac{1}{2\sigma_k^2} \|x_k^* - x_i\|^2 \right] \right)^{\sigma_k^2} \right] \\
&= \min_i \lim_{k \rightarrow \infty} \|x_k^* - x_i\|^2 \\
&= \lim_{k \rightarrow \infty} \|x_k^* - x_{j_1}\|^2 \\
&= \left\| \frac{1}{m} \sum_{l=1}^m x_{j_l} - x_{j_1} \right\|^2.
\end{aligned} \tag{21}$$

478 We observe that  $x_k^*$  is a local maximizer of  $\min_x d_{\sigma_k}(x; \mathcal{D})^2$ , which raises contradiction.

479 Similarly, we prove that  $x_j$  is not unique if  $x_k^*$  is a local maximizer or saddle point by contradiction.

480 Assume  $x_j$  is the unique minimizer of  $\min_{x_i \in \mathcal{D}} \|x_k^* - x_i\|^2$ , we have

$$\begin{aligned}
\lim_{k \rightarrow \infty} d_{\sigma_k}(x_k^*; \mathcal{D})^2 &= \min_i \lim_{k \rightarrow \infty} \|x_k^* - x_i\|^2 \\
&= \lim_{k \rightarrow \infty} \|x_k^* - x_j\|^2 \\
&= 0.
\end{aligned} \tag{22}$$

481 Hence,  $x_k^*$  can not be a local maximizer or saddle point of  $\min_x d_{\sigma_k}(x; \mathcal{D})^2$ , which raises contradic-  
482 tion. □

483 Since gradient descent converges to a local minimizer almost surely with random initialization [69]  
484 and appropriate step sizes, there is a unique  $x_j$  for  $x_k^*$  and (18) holds. We also note that a similar  
485 conclusion can be reached by using  $\Gamma$ -convergence of the softmin to min as  $\sigma_k \rightarrow 0$  [70].

486 **A.2 Proof for Proposition 2**

487 Let  $L_e$  be the local Lipschitz constant of the error  $e(x) := f(x) - f_\theta(x)$  over some domain  $\mathcal{Z} \subseteq \mathcal{X}$ ,  
 488 and define  $x_c := \arg \min_{x_i \in \mathcal{D}} \frac{1}{2} \|x - x_i\|_2^2$ , i.e., the closest data-point. Then, we have the following:

$$\begin{aligned}
 \|f(x) - f_\theta(x)\| &\leq \min_{x_i \in \mathcal{D}} [e(x_i) + L_e \|x - x_i\|_2] \\
 &\leq e(x_c) + L_e \|x - x_c\|_2 \\
 &= e(x_c) + \sqrt{2} L_e \sqrt{\frac{1}{2} \|x - x_c\|_2^2} \\
 &= e(x_c) + \sqrt{2} L_e \sqrt{\frac{1}{2} \min_{x_i \in \mathcal{D}} \|x - x_i\|_2^2} \\
 &= e(x_c) + \sqrt{2} L_e \sqrt{\min_{x_i \in \mathcal{D}} \frac{1}{2} \|x - x_i\|_2^2} \\
 &\leq e(x_c) + \sqrt{2} L_e \sqrt{-\sigma^2 \log \left( \sum_i \exp \left( -\frac{1}{2\sigma^2} \|x - x_i\|^2 \right) \right) + \sigma^2 \log N} \\
 &= e(x_c) + \sqrt{2} L_e \sqrt{d_\sigma(x; \mathcal{D})^2 + C_2}
 \end{aligned} \tag{23}$$

489 where  $C_2 := \sigma^2 \log N - C$ , and  $C$  is the constant defined in (4). In the second line, we use the  
 490 fact that as  $x_c$  is a feasible solution to the minimization in the first line, it is an upper bound on the  
 491 optimal value. In the sixth line, we have used the fact that for any vector  $v = [v_1, \dots, v_n]^\top \in \mathbb{R}^n$ ,  
 492  $\min\{v_1, \dots, v_n\} \leq -\frac{1}{t} \log \sum_{i=1}^n \exp(-tv_i) + \frac{\log n}{t}$  for some scaling  $t$ . In the final line, we have  
 493 applied the definition of  $d_\sigma(x; \mathcal{D})$  from (4).

494 **A.3 Proof for Proposition 3**

495 The perturbed data distribution can be written as a sum of Gaussians, since

$$\begin{aligned}
 p_\sigma(x') &:= \int \hat{p}(x) \mathcal{N}(x'; x, \sigma^2 \mathbf{I}) dx \\
 &= \int \left[ \frac{1}{N} \sum_i \delta(x_i) \right] \mathcal{N}(x'; x, \sigma^2 \mathbf{I}) dx \\
 &= \frac{1}{N} \sum_i \int \delta(x_i) \mathcal{N}(x'; x, \sigma^2 \mathbf{I}) \\
 &= \frac{1}{N} \sum_i \mathcal{N}(x_i; x, \sigma^2 \mathbf{I}) \\
 &= \frac{1}{N} \sum_i \mathcal{N}(x; x_i, \sigma^2 \mathbf{I})
 \end{aligned} \tag{24}$$

496 Then we consider the negative log of the perturbed data distribution multiplied by  $\sigma^2$ ,

$$\begin{aligned}
-\sigma^2 \log p_\sigma(x) &= -\sigma^2 \log \left[ \frac{1}{N} \sum_i \mathcal{N}(x; x_i, \sigma^2 \mathbf{I}) \right] \\
&= -\sigma^2 \log \left[ \sum_i \mathcal{N}(x; x_i, \sigma^2 \mathbf{I}) \right] + \log N \\
&= -\sigma^2 \log \left[ \frac{1}{\sqrt{(2\pi\sigma)^n}} \sum_i \exp \left[ -\frac{1}{2\sigma^2} \|x - x_i\|^2 \right] \right] + \sigma^2 \log N \\
&= -\sigma^2 \log \left[ \sum_i \exp \left[ -\frac{1}{2\sigma^2} \|x - x_i\|^2 \right] \right] + \sigma^2 \log N + \sigma^2 \frac{n}{2} \log(2\pi\sigma) \\
&= -\sigma^2 \text{LogSumExp}_i \left[ -\frac{1}{2\sigma^2} \|x - x_i\|^2 \right] + C(N, n, \Sigma) \\
&= \text{Softmin}_\sigma \left[ \frac{1}{2} \|x - x_i\|^2 \right] + C(N, n, \Sigma)
\end{aligned} \tag{25}$$

497 where we define  $C(N, n, \Sigma) := \sigma^2(\log N + n/2 \log(2\pi\sigma))$ .

## 498 B Details of the Planning algorithm

### 499 B.1 Computation of Gradients

500 Recall that the gradient of  $\log p_\sigma(x_i, u_i)$  cost w.r.t. the input variable  $u_j$  can be written as

$$501 \quad \nabla_{u_j} \log p(x_i, u_i) = \nabla_{x_i} \log p(x_i, u_i) \mathbf{D}_{u_j} x_i + \nabla_{u_i} \log p(x_i, u_i) \mathbf{D}_{u_j} u_i \tag{26}$$

502 where  $\mathbf{D}$  denotes the Jacobian. Writing the dependence on each variable more explicitly, we have

$$503 \quad \nabla_{u_j} \log p(x_i(u_j), u_i(u_j)) = \nabla_{x_i} \log p(x_i(u_j), u_i(u_j)) \mathbf{D}_{u_j} x_i(u_j) \tag{27}$$

$$504 \quad + \nabla_{u_i} \log p(x_i(u_j), u_i(u_j)) \mathbf{D}_{u_j} u_i(u_j) \tag{28}$$

505 where we note that  $\mathbf{D}_{u_j} u_i = 1$  if  $i = j$  and 0 otherwise. As long as  $i > j$ , we also note that  $x_i$   
506 has a dependence on  $u_j$ . Instead of computing this gradient explicitly, we first rollout the trajectory  
507 to compute  $x_i(u_j), u_i(u_j)$ , and compute the score function. Then we ask: which quantity do we  
508 need such that it gives us the above expression when differentiated w.r.t.  $u_j$ ? We use the following  
509 quantity,

$$510 \quad c_{ij} = s_x(x_i, u_i) x_i(u_j) + s_u(x_i, u_i) u_i(u_j) \tag{29}$$

511 where the score terms have been detached from the computation graph. Note that  $c_{ij}$  is a scalar and  
512 allows us to use reverse-mode automatic differentiation tools such as `pytorch` [71].

### 509 B.2 Noise-Annealing During Optimization

510 Additionally, we anneal the noise level during iterations of Adam. Given a sequence  $\sigma_k$  with  $K$  being  
511 the total number of annealing steps, we run Adam for  $\max \text{iter}_{\max}/K$  iterations, then run it with the  
512 next noise level.

### 513 B.3 Connection to Diffuser

514 We first lift the dynamics constraint into a quadratic penalty and write the penalty as  $\log p(x_{t+1}|x_t, u_t)$ .  
515 This equivalence is seen by considering a case where we fix  $x_t, u_t$  and perturb  $x_{t+1}$  with a Gaussian  
516 noise of scale  $\sigma$ . If  $(x_t, u_t, x_{t+1})$  is in the dataset, it obeys  $x_{t+1} = f(x_t, u_t)$  under real-world  
517 dynamics  $f$ . This allows us to write

$$\begin{aligned}
p_\sigma(x_{t+1}|x_t, u_t) &= \mathcal{N}(x_{t+1}|f(x_t, u_t), \sigma^2 \mathbf{I}) \\
\log p_\sigma(x_{t+1}|x_t, u_t) &= -\frac{1}{2} \|x_{t+1} - f(x_t, u_t)\|^2 + C
\end{aligned} \tag{30}$$

518 where  $C$  is some constant that does not effect the objective. Then, we rewrite our objective using the  
 519 factoring  $p(x_t, u_t) = p(u_t|x)p(x_t)$ . This allows us to rewrite the objective of Equation (11) as

$$\begin{aligned} & \sum_{t=1}^T r_t(x_t, u_t) + \beta \sum_{t=1}^T \log p(u_t|x_t) + \beta \sum_{t=1}^T \log p(x_t) + \beta \sum_{t=1}^T \log p(x_{t+1}|x_t, u_t) \\ & = V(x_{1:T}, u_{1:T}) + \beta \log p(x_{1:T}, u_{1:T}) + \sum_{t=1}^T \log p(x), \end{aligned} \quad (31)$$

520 where the first two terms are the objectives in Diffuser [25].

## 521 B.4 First-Order Policy Search

522 We note that our original method for gradient computation can easily be extended to the setting of  
 523 feedback first-order policy search, where we define the uncertainty-penalized value function as

$$\begin{aligned} & \max_{\alpha} \mathbb{E}_{x_1 \sim \rho} \left[ \sum_{t=1}^T r_t(x_t, u_t) + \beta \sigma^2 \sum_{t=1}^T \log p_{\sigma}(x_t, u_t; \mathcal{D}) \right] \\ & \text{s.t. } x_{t+1} = f_{\theta}(x_t, u_t), u_t = \pi_{\alpha}(x_t) \quad \forall t \in [1, T], \end{aligned} \quad (32)$$

524 where  $\rho$  is some distribution of initial conditions. We rewrite the objective with an explicit dependence  
 525 on  $\alpha$ , and use a Monte-Carlo estimator for the gradient of the stochastic objective,

$$\begin{aligned} & \nabla_{\alpha} \mathbb{E}_{x_1 \sim \rho} \left[ \sum_{t=1}^T r_t(x_t(\alpha), u_t(\alpha)) + \beta \sigma^2 \sum_{t=1}^T \log p_{\sigma}(x_t(\alpha), u_t(\alpha); \mathcal{D}) \right] \\ & = \mathbb{E}_{x_1 \sim \rho} \nabla_{\alpha} \left[ \sum_{t=1}^T r_t(x_t(\alpha), u_t(\alpha)) + \beta \sigma^2 \sum_{t=1}^T \log p_{\sigma}(x_t(\alpha), u_t(\alpha); \mathcal{D}) \right] \\ & \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\alpha} \left[ \sum_{t=1}^T r_t(x_t(\alpha), u_t(\alpha)) + \beta \sigma^2 \sum_{t=1}^T \log p_{\sigma}(x_t(\alpha), u_t(\alpha); \mathcal{D}) \quad \text{s.t. } x_1 = x_i \sim \rho \right], \end{aligned} \quad (33)$$

526 where the last equation denotes that fixing the initial condition to  $x_i$  sampled from  $\rho$ , and  $N$  is the  
 527 number of samples in the Monte-Carlo process. Since  $r_t$  and  $f_{\theta}$  are differentiable, we can obtain the  
 528 gradient

$$\nabla_{\alpha} \sum_{t=0}^T r_t(x_t(\alpha), u_t(\alpha)) \quad (34)$$

529 after rolling out the closed-loop system starting from  $x_i$  and using automatic differentiation w.r.t.  
 530 policy parameters  $\alpha$ . To compute the gradient w.r.t. the score function, we similarly use the chain  
 531 rule,

532 and differentiate it w.r.t  $\alpha$ , which lets us compute

$$\begin{aligned} \frac{\partial}{\partial \alpha} \left[ \sum_{t=1}^T \log p_{\sigma}(x_t, u_t) \right] &= \sum_{t=1}^T \frac{\partial}{\partial \alpha} \log p_{\sigma}(x_t(\alpha), u_t(\alpha)) \\ &= \sum_{t=1}^T \frac{\partial}{\partial x_t} \log p_{\sigma}(x_t, u_t) \frac{\partial x_t}{\partial \alpha} + \frac{\partial}{\partial u_t} \log p_{\sigma}(x_t, u_t) \frac{\partial u_t}{\partial \alpha} \\ &= \sum_{t=1}^T s_x(x_t, u_t; \sigma) \frac{\partial x_t}{\partial \alpha} + s_u(x_t, u_t; \sigma) \frac{\partial u_t}{\partial \alpha} \end{aligned} \quad (35)$$

533 where the last term is obtained by differentiating

$$\sum_{t=1}^T s_x(x_t, u_t; \sigma) x_t + s_u(x_t, u_t; \sigma) u_t \quad (36)$$

534 after detaching  $s_x$  and  $s_u$  from the computation graph.

## 535 B.5 Imitation Learning

536 We give more intuition for why maximizing the state-action likelihood leads to imitation learning.  
537 If the empirical data comes from an expert demonstrator, maximizing data likelihood leads to  
538 minimization of cross entropy between the state-action pairs encountered during planning and the  
539 state-action occupation measure of the demonstration policy, which is estimated with the perturbed  
540 empirical distribution  $p_\sigma(x_t, u_t)$ ,

$$\sum_t \log p_\sigma(x_t, u_t) = \sum_t \log p_\sigma(u_t|x_t) + \sum_t \log p_\sigma(x). \quad (37)$$

541 Note that the  $\log p_\sigma(u_t|x_t)$  is identical to the Behavior Cloning (BC) objective, while  $\log p_\sigma(x_t)$   
542 drives future states of the plan closer to states in the dataset. We note that Adversarial Inverse  
543 Reinforcement Learning (AIRL) [34] minimizes a similar objective as ours [12].

## 544 C Experiment Details

### 545 C.1 Cartpole with Learned Dynamics

546 **Environment.** We use the cart-pole dynamics model in [58, chapter 3.2], with the cost function  
547 being

$$c_t(x_t, u_t) = \begin{cases} \|x_t - x_g\|_{\mathbf{Q}}^2 & \text{if } t = T \\ 0 & \text{else.} \end{cases} \quad (38)$$

548  $\mathbf{Q} = \text{diag}(1, 1, 0.1, 0.1)$ . We choose the planning horizon  $T = 60$ .

549 **Training.** We randomly collected a dataset of size  $N = 1,000,000$  within the red box region in the  
550 state space. The dynamics model is an MLP with 3 hidden layers of width (64, 64, 32). For ensemble  
551 approach we use 6 different dynamics models, all with the same network structures.

552 We train a score function estimator, represented by an MLP with 4 hidden layers of width 1024. The  
553 network is trained form 400 epochs with a batch size of 2048.

554 **Parameters** During motion planning, for Adam optimizer we use a learning rate of 0.01. For CEM  
555 approach we use a population size of 10, with standard variance  $\sigma = 0.05$ , and we take the top 4  
556 seeds to update the mean in the next iteration.

557 **Data Distance Estimator.** To train a data distance estimator, we introduce a function approximator  
558  $d_\eta : \mathbb{R}^n \times \mathbb{R}_+ \rightarrow \mathbb{R}$  parametrized by  $\eta$  to predict the noise-dependent Softmin distance. The training  
559 objective is given by

$$\min_\eta \frac{1}{2} \mathbb{E}_{x \in \Omega, \sigma \in [0, \sigma_{max}]} \left[ \frac{1}{\sigma^2} \left| d(x, \sigma) - \text{Softmin}_{x_i \in \mathcal{D}} \frac{1}{2} \|x - x_i\|_{\sigma^{-2}\mathbf{I}}^2 \right| \right] \quad (39)$$

560 where  $\Omega$  is a large enough domain that covers the data distribution  $\mathcal{D}$ . For small datasets, it is possible  
561 to loop through all  $x_i$  in the dataset to compute this loss at every iteration. However, this training can  
562 get prohibitive as all of the training set needs to be considered to compute the loss, preventing batch  
563 training out of the training set.

### 564 C.2 D4RL Dataset

565 **Environment.** We directly use the D4RL dataset [43] Mujoco tasks [60] with 3 different envi-  
566 ronments of halfcheetah, walker2d, and hopper. We additionally use differnet sources of data with  
567 random, medium, and medium-expert.

568 **Training.** The dynamics and the score functions are both parametrized with MLP with 4 hidden  
569 layers of width 1024. The noise-conditioned score function is implemented by treating each level

570 of noise  $\sigma_k$  as an integer token, that gets embedded into a 1024 vector and gets multiplied with the  
 571 output of each layer. This acts similar to a masking of the weights depending on the level of noise.  
 572 The D4RL environment does not provide us with a differentiable reward function, so we additionally  
 573 train an estimator for the reward. We empirically saw that for score function estimation, wide shallow  
 574 networks performed better. We train both instances for 1000 iterations with Adam, with a learning  
 575 rate of  $1e-3$  and batch size of 2048.

576 We additionally set a noise schedule to be a cosine schedule that anneals from  $\sigma = 0.2$  to  $\sigma = 0.01$   
 577 for 10 steps in the normalized space of  $x, u$ .

578 **Parameters** We used a range of  $\beta$ s between  $1e^{-3}$  and  $1e^{-1}$  depending on the environment, where  
 579 in some cases it helped to be more reliant on reward, and in others it’s desirable to rely on imitation.  
 580 We use a MPC with  $T = 5$  and optimize it for 50 iterations with an aggressive learning rate of  $1e^{-1}$ .

### 581 C.3 Pixel Single Integrator

582 **Environment.** In this environment, we have a 2D single integrator  $f(x_t, u_t) = x_t + u_t$ ,  $x_t \in$   
 583  $\mathbb{R}^2$ ,  $u_t \in \mathbb{R}^2$  as the underlying ground-truth dynamics; however, instead of raw states  $x_t$ , we observe  
 584 a  $32 \times 32$  grayscale image  $y_t = h(x_t) \in \mathbb{R}^{32 \times 32}$ , which are top-down renderings of the robot. In  
 585 these observations, the position of the robot is represented with a dot. Moreover, we assume that we  
 586 do not directly assign the 2D control input  $u_t$ , but instead propose a  $32 \times 32$  grayscale image  $\hat{u}_t$ ,  
 587 where the value of the 2D control action  $u_t$  is extracted from the image via a spatial average:

$$u_t = \sum_{(p_x, p_y)} g_u(p_x, p_y) \hat{u}_t(p_x, p_y), \quad (40)$$

588 where the sum loops over each pixel  $(p_x, p_y) \in \{1, \dots, 32\}^2$ ,  $y_t(p_x, p_y)$  refers to the intensity of the  
 589 image at pixel  $(p_x, p_y)$ , and  $g_u : \mathbb{Z} \times \mathbb{Z} \rightarrow \mathbb{R}^2$  is a grid function mapping from pixel  $(p_x, p_y)$  to a  
 590 corresponding control action. The control image  $\hat{u}_t$  is normalized such that its overall intensity sums  
 591 to 1. Given some goal  $x_g$ , the reward is set to be the  $-c_t(x_t, u_t)$ , where the cost  $c_t(x_t, u_t)$  is

$$c_t(x_t, u_t) = \begin{cases} \|x_t - x_g\|_{\mathbf{Q}_t}^2 + \|u_t\|_{\mathbf{R}}^2 & \text{if } t = T \\ \|x_t - x_g\|_{\mathbf{Q}}^2 + \|u_t\|_{\mathbf{R}}^2 & \text{else,} \end{cases} \quad (41)$$

592 To evaluate this cost function for the planned sequence of image observations and control images, the  
 593 states  $x_t$  are also extracted from the image observations through a similar spatial averaging:

$$x_t = \sum_{(p_x, p_y)} g_{\text{im}}(p_x, p_y) y_t(p_x, p_y), \quad (42)$$

594 where  $g_{\text{im}} : \mathbb{Z} \times \mathbb{Z} \rightarrow \mathbb{R}^2$  is a grid function mapping from pixel  $(p_x, p_y)$  to a corresponding state.

595 In other words, we have running costs for the state and input, and a different terminal cost for the  
 596 state. We set  $\mathbf{R} = 6.5\mathbf{I}$ ,  $\mathbf{Q} = 500\mathbf{I}$ , and  $\mathbf{Q}_d = 1000\mathbf{I}$ , and plan with a horizon of  $T = 15$ .

597 **Training.** We collect a randomly collected dataset of size  $N = 200,000$ , with underlying 2D  
 598 data sampled from  $x_t \in [-1, 1]^2$  and  $u_t \in [-0.2, 0.2]^2$ . Both the dynamics and the score function  
 599 estimator are represented as U-Nets [72], with the architecture coming from [73].

600 **Parameters.** We found that  $\beta = 0.5$  is sufficient for the penalty parameter. Results are obtained  
 601 within 1450 iterations with a learning rate of 0.03. In representing the noise-conditioned score  
 602 function, we use 232 smoothing parameters  $\{\sigma_k\}_{k=1}^{232}$ , from  $\sigma_1 = 50$  to  $\sigma_{232} = 0.01$ .

603 **Baselines.** For gradient-based planning with ensembles, we set  $\beta = 0.5$  and use an ensemble of  
 604 size 10. For CEM with ensembles, we set  $\beta = 0.5$ , with an ensemble of size 5 (we only used the first  
 605 five networks in the original ensemble of size 10 due to RAM limitations).

606 **C.4 Box Pushing with Marker Dynamics**

607 **Environment.** We prepare a box-pushing environment where we assume that the box follows  
608 quasistatic dynamics, which allows us to treat the marker positions directly as state of the box that is  
609 bijective with its pose. We use 2D coordinates for each markers, and append the pusher position, also  
610 in 2D, resulting in  $x_t \in \mathbb{R}^{12}$ . The pusher is given a relative position command with a relatively large  
611 step size [68]. In addition, we give the robot knowledge of the pusher dynamics,  $x_{t+1}^{\text{pusher}} = x_t^{\text{pusher}} + u_t$ .  
612 The general goal of the task is to push the box and align the edge of the box with the blue tape line.

613 We formulate our cost as

$$c_t(x_t, u_t) = \begin{cases} \|x_t^{\text{marker}} - x_g^{\text{marker}}\|_{\mathbf{Q}_T}^2 & \text{if } t = T \\ \|u_t\|_{\mathbf{R}}^2 & \text{else,} \end{cases} \quad (43)$$

614 where  $\mathbf{Q}_T = \mathbf{I}$  and  $\mathbf{R} = 0.1\mathbf{I}$ . In order to get  $x_g^{\text{marker}}$ , we place the box where we want the goal to be  
615 and measure the position of the markers.

616 **Training.** We collect 100 demonstration trajectories resulting in 750 pairs of  $(x_t, u_t, x_{t+1})$ . The  
617 dynamics and the score functions are learned with a MLP of 4 hidden layers with size 1024, with the  
618 noise-conditioned score estimator being trained similar to the D4RL dataset with multiplicative token  
619 embeddings. We train for 500 iterations with a batch size of 32.

620 We additionally set a noise schedule to be a cosine schedule that anneals from  $\sigma = 0.2$  to  $\sigma = 0.01$   
621 for 10 steps.

622 **Parameters.** We observed that  $\beta = 1e^{-2}$  performs well for all the examples, with a learning rate  
623 of 0.1 and 50 iterations. We found that a horizon of  $T = 4$  was sufficient for our setup.