
Behavior Cloning is Not All You Need: The Optimality of On-Policy Distillation for Noisy Expert Feedback

Anonymous Author(s)
Affiliation
Address
email

Abstract

1 Imitation Learning (IL) is a natural framework for learning in sequential decision-
2 making systems and has emerged as the dominant paradigm through which we
3 understand language model training. A central puzzle is that, while in theory offline
4 IL can be horizon-free and optimal, in practice online methods such as on-policy
5 distillation (OPD) often outperform offline methods such as supervised fine-tuning
6 (SFT). We propose a noisy expert model to explain this gap, in which the learner
7 only has access to a noisy version of the expert’s policy, but wishes to compete
8 against the reward achieved by a clean expert, motivated by the fact that in many
9 applications, e.g. training language models to perform long chains of thought, the
10 expert is often imperfect. In this setting, we show a sharp separation between offline
11 and online IL. Offline learning from noisy trajectories is fundamentally hard: to
12 compete with the clean expert, the sample complexity must grow exponentially, in
13 contradistinction to the clean expert setting where no explicit horizon dependence
14 exists. In contrast, we prove that online interaction with the noisy expert via a novel
15 variant of OPD enables horizon-free guarantees in some settings and polynomial
16 dependence on horizon in general. Our analysis leads to an alternative loss function
17 form that is commonly considered empirically for LM training. We further provide
18 algorithms and lower bounds, and extend our results to the more realistic setting
19 of unknown corruption when the clean expert is deterministic, thereby providing
20 a theoretical foundation for why on-policy distillation can outperform standard
21 supervised fine-tuning when training language models from imperfect teachers.
22 We complement our theoretical results with experiments on synthetic and natural-
23 language tasks, showing that the OPD variant suggested by our theory outperforms
24 both offline BC and existing OPD objectives under noisy expert feedback.

25 1 Introduction

26 Imitation learning (IL) is a powerful framework for training agents to perform complex tasks in
27 challenging environments by learning from expert demonstrations [7, 44, 47, 49]. In recent years, IL
28 has gained significant attention due to its success in a wide range of applications, including robotics
29 [5, 13], autonomous driving [11], and natural language processing [2, 7, 9, 41]. In particular, IL has
30 been proposed as a natural framework through which to understand Language Models (LMs), where
31 the training corpus is thought to consist of ‘expert demonstrations’ of human behavior. Due to the
32 increasing importance and cost of LMs, this raises the question of how to best leverage expert access
33 to train LMs, and more generally, how to best leverage expert access in IL.

34 Prior work has distinguished between two paradigms for IL: *offline* IL, where the learner only has
35 access to a fixed dataset of expert demonstrations, and *online* IL, where the learner can interact with
36 the expert and query for demonstrations in states visited by the learner [17, 45, 47, 49]. The canonical
37 algorithm in offline IL is Behavioral Cloning (BC), where the learner simply performs supervised

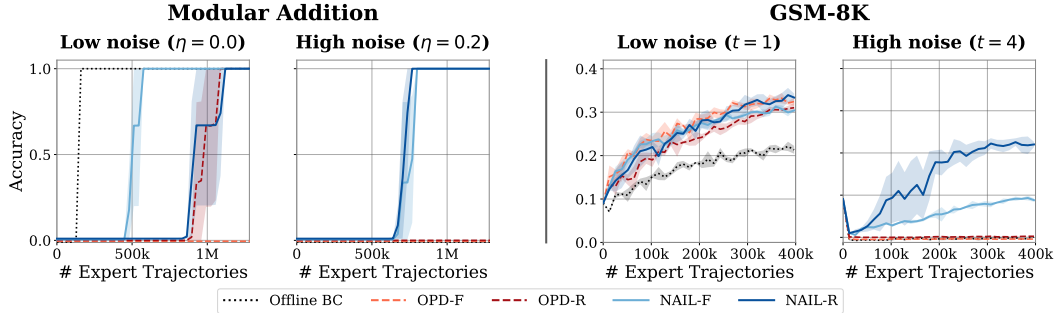


Figure 1: Comparison of offline BC, standard OPD from [2, 41], and our proposed NAIL for forward (F) and reverse (R) KL losses with clean and noisy experts. **Left:** Modular addition task: adding 31 numbers mod 7, where the expert is corrupted with probability $\eta \in \{0, 0.2\}$. **Right:** GSM-8K: math reasoning where the expert generates with temperature $t \in \{1, 4\}$. In both cases, NAIL is competitive with offline methods and standard OPD in the low noise regime and strongly outperforms these baselines in the high noise regime.

38 learning on the expert demonstrations in an effort to predict the action from the observation. While
 39 it was classically thought that BC suffers from compounding errors and thus has poor performance
 40 in long-horizon tasks, with the suboptimality of the learned policy growing quadratically with the
 41 horizon [47, 49], more recent work has shown that BC, when instantiated with the correct loss
 42 function, is both horizon-free and optimal in a minimax sense, even when comparing to online IL
 43 algorithms [17]. While this theoretical result may seem to suggest that BC is the best way to leverage
 44 expert access, in practice, many works have found that online IL algorithms, such as Dagger [49]
 45 and On-Policy Distillation (OPD) [2, 41], can significantly outperform offline IL in a wide range of
 46 settings. Recent work has suggested that this discrepancy between theory and practice may be due
 47 to the fact that the expert is often misspecified, meaning that the expert’s policy is not realizable by
 48 the learner’s policy class [45]; such misspecification can occur in robotics and autonomous driving
 49 problems when the expert and learner have different observation spaces (e.g. a human demonstrator
 50 has more sensory input than a robot learning) [5] and in LMs when we are trying to distill a large
 51 teacher model into a smaller student model [2, 24]. While misspecified experts are a natural model in
 52 many settings, they do not apply to one of the chief successes of IL for LMs, which is to use OPD to
 53 fine-tune a student LM from a teacher LM that has been trained with reinforcement learning (RL) but
 54 where both student and teacher share the same architecture [2, 41]. In this setting, the expert is not
 55 misspecified, but online IL (OPD) still outperforms offline IL (SFT) empirically.

56 In this work, we propose a *noisy expert* model to explain the discrepancy between theory and practice
 57 in IL. In particular, we suppose that there is a true expert policy π^* that achieves good expected
 58 reward when rolled out in an environment, but we only have access to a noisy version of this expert,
 59 $\pi_\eta^* = (1 - \eta) \cdot \pi^* + \eta \cdot \nu$, where $0 \leq \eta < 1$ is a noise level and ν is some noise distribution. This
 60 model captures the fact that in many real-world settings, such as math, code, and reasoning, the true
 61 expert against which we wish to compete could be approximately deterministic, but for computational
 62 reasons we are training a stochastic expert. Moreover, in the LM setting, the teacher model that we
 63 are distilling from is often a large model that is not fully converged and thus can be thought of as a
 64 noisy expert that occasionally makes mistakes [29, 61]. Such mistakes can also arise from human text
 65 that introduces minor errors, typos, and inconsistencies into datasets, which is known to be harmful
 66 for training LMs [35, 42, 60]. We thus ask: *Given access to a noisy expert, how should we leverage*
 67 *this expert to train a learner policy that performs well when rolled out in the environment?*

68 We provide an answer to this question in both the offline and online settings. For the sake of simplicity,
 69 in the introduction we discuss a special case of the noisy expert model (Definition 4), which is always
 70 satisfied when π^* is deterministic. We now summarize our main theoretical contributions. We first
 71 consider the setting where both η and ν are known to the learner; while unrealistic in practice, this
 72 setting allows us to cleanly characterize the fundamental limits of learning from a noisy expert
 73 without adding identifiability concerns (cf. Section 5). Our first result is that in offline settings, unlike

74 the clean expert setting where horizon-free guarantees are possible, we must pay *exponentially* in
 75 horizon H in order to guarantee learning whenever $\eta \gtrsim 1/H$.

76 **Theorem 1** (Informal version of [Theorem 6](#) and [Proposition 1](#)). *Suppose an expert π^* is contained in*
 77 *a policy class Π and a learner has access to n trajectories rolled out from a noisy expert π_η^* . In order*
 78 *to learn a policy $\hat{\pi}$ whose expected reward when rolled out is within ε of that of π^* , BC requires*
 79 *$n \gtrsim (1-\eta)^{-(H+2)} \cdot \log(|\Pi|)/\varepsilon$ and this is optimal.*

80 This result stands in marked contrast to the clean expert setting, where BC can achieve horizon-free
 81 guarantees, and suggests that offline IL is fundamentally intractable when the expert is noisy. In the
 82 *online* setting, however, we show that we can circumvent this exponential dependence on the horizon
 83 via *On-Policy Distillation*. Our second main result shows that a natural analogue of OPD, where
 84 the learner queries the noisy expert for demonstrations in states visited by the learner, can achieve
 85 horizon-free guarantees even when the expert is noisy.

86 **Theorem 2** (Informal version of [Theorem 7](#)). *Let $\mathbb{P}^{\pi, \pi'}$ denote the law of an augmented trajectory*
 87 *of states and actions generated by rolling out π in an environment, but at each step querying π' for*
 88 *an auxiliary action. Then, the suboptimality in expected reward of a policy $\hat{\pi}$ can be bounded by*
 89 *$(1-\eta)^{-2} \cdot (\text{D}_{\text{KL}}(\mathbb{P}^{\hat{\pi}, \pi_\eta^*} \parallel \mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta}) \wedge \text{D}_{\text{KL}}(\mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta} \parallel \mathbb{P}^{\hat{\pi}, \pi_\eta^*}))$.*

90 Note that $\mathbb{P}^{\hat{\pi}, \pi_\eta^*}$ is only available through online feedback: we roll out policy $\hat{\pi}$ and, at each step,
 91 ask for what the (noisy) expert π_η^* would do in that state. *This result suggests that rolling out the*
 92 *student as will be done at deployment (e.g. greedily) and then using the noisy expert to score the*
 93 *actions taken by the student is the right way to leverage online access to a noisy expert, providing*
 94 *a principled justification for OPD.* We emphasize that the suggested loss function is *different* from
 95 that considered in Agarwal et al. [2], Lu and Lab [41], which do not distinguish between rollout
 96 distribution (e.g. greedy) and student policy (e.g. temperature 1) during training. Our final main
 97 theoretical result concerns our ability to control this forward KL divergence.

98 **Theorem 3** (Informal version of [Theorems 8](#) and [9](#)). *Suppose we have access to a policy class Π*
 99 *that contains π^* . If η and ν are known, there is an OPD-like algorithm that can find a policy $\hat{\pi}$*
 100 *with suboptimality in expected reward at most ε with $n \gtrsim H^2 \cdot \log(|\Pi|)/\varepsilon^2 (1-\eta)^2$ rounds of interaction*
 101 *with the noisy expert π_η^* and this is optimal up to a factor of H . Moreover, if η and ν are unknown*
 102 *and π^* is deterministic, then as long as η is not too large and ν puts mass at most ρ on any action,*
 103 *$n \gtrsim \log(|\Pi|)/\varepsilon(1-\eta(1+\rho))^2$ rounds of interaction suffice and this is optimal.*

104 In essence, we show that horizon-free guarantees on suboptimality with respect to the *clean* expert
 105 are possible when given online access to a noisy expert, when the corruption is unknown but the
 106 expert is deterministic, while we pay only polynomially in the horizon in general when given online
 107 access to a noisy expert.

108 We empirically validate our theory in [Figure 1](#) on both a synthetic task of modular addition [38] and
 109 a natural language task of math reasoning (GSM-8K) [14] by minimizing a loss function motivated
 110 by [Theorem 7](#). Across both settings, online distillation substantially outperforms LogLossBC under
 111 noisy expert feedback, while remaining competitive under clean expert feedback, in line with our
 112 theoretical findings. Moreover, our NAIL variants consistently improve over standard OPD in the
 113 noisy expert setting, suggesting that the loss function suggested by our theory is indeed more effective
 114 for leveraging online access to a noisy expert. We defer further experimental details to [Section 6](#).

115 In [Section 2](#) we formally introduce the problem, as well as some prerequisite notions. In [Section 3](#),
 116 we analyze the *offline* setting before discussing the *online* setting in [Section 4](#). In [Section 5](#), we
 117 extend to the setting where η and ν are unknown and the expert is deterministic. Finally, in [Section 6](#)
 118 we present some preliminary empirical results validating our theoretical findings. We conclude with
 119 a discussion of related work and future directions in [Appendix A](#) and [Table 1](#) summarizes our results.

120 2 Formal Problem Setup and Preliminaries

121 We consider Imitation Learning in the standard episodic Markov Decision Process (MDP) setting
 122 [17, 47, 49]. In particular, we have an MDP M with horizon H consisting of a state space \mathcal{S} , an
 123 action space \mathcal{A} , transition kernels $P_h : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ for each step $h \in [H]$, and a reward function
 124 $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$. We assume that the initial state distribution is given by some $\rho \in \Delta(\mathcal{S})$. We
 125 will consider learning *policies* $\pi : \mathcal{S} \times [H] \rightarrow \Delta(\mathcal{A})$ that can be non-stationary across the horizon,

126 mapping states to distributions over actions. Given a policy π we denote by \mathbb{P}^π the distribution over
 127 trajectories induced by rolling out π in the MDP M , i.e., $\tau = (s_1, a_1, \dots, s_H, a_H)$ with $s_1 \sim \rho$,
 128 $a_h \sim \pi_h(\cdot | s_h)$ and $s_{h+1} \sim P_h(\cdot | s_h, a_h)$. We will also use \mathbb{E}^π to denote the expectation over
 129 trajectories induced by π . We are primarily concerned with finding policies π that have small *regret*
 130 with respect to the expert policy π^* : $\text{Reg}(\pi) = J(\pi^*) - J(\pi)$, where $J(\pi) = \mathbb{E}^\pi \left[\sum_{h=1}^H r(s_h, a_h) \right]$.
 131 As in Foster et al. [17], we will decouple the reward range from the horizon by assuming that
 132 $0 \leq \sum_{h=1}^H r(s_h, a_h) \leq R$ for some R , for any trajectory τ . While naively, R can scale with the
 133 horizon H , in many settings of interest we can have R be a constant independent of H , which allows
 134 for horizon-free regret bounds, e.g. for LMs learning to solve math problems where the trajectory
 135 reward is binary [23, 54].

136 Unlike in reinforcement learning, imitation learning assumes access only to trajectories τ without any
 137 reward signal. While in classical IL, one assumes access to expert actions, in this work we consider
 138 instead a *noisy expert*, where there exists a corruption policy $\nu : \mathcal{S} \times [H] \rightarrow \Delta(\mathcal{A})$ and a corruption
 139 level $\eta \in [0, 1)$ such that the learner has access to π_η^* defined as

$$\pi_{\eta,h}^*(a | s) = (1 - \eta) \cdot \pi_h^*(a | s) + \eta \cdot \nu_h(a | s). \quad (1)$$

140 As in more traditional IL, we distinguish between two different ways in which to interact with the
 141 (noisy) expert. In the *offline* setting [17, 44, 47], we assume access only to expert trajectories and
 142 attempt to learn without additional interaction, formalized as follows.

143 **Definition 1** (Offline Imitation Learning). In the *offline* setting, a learner is given n trajectories
 144 $\tau^1, \dots, \tau^n \sim \mathbb{P}^{\pi_\eta^*}$ with π_η^* as in (1) and must output a policy $\hat{\pi}$ without further interaction.

145 Meanwhile, in *online* imitation learning [49], the learner is instead allowed to query the expert on
 146 trajectories rolled out by the learner itself, formalized as follows.

147 **Definition 2.** In the *online IL* setting, a learner interacts with the environment in rounds. In each
 148 round t , the learner deploys a policy π_t and observes a trajectory $\tau^{(t)}$ drawn from \mathbb{P}^{π_t} . The learner
 149 can then query the (noisy) expert π_η^* at each of the states visited in $\tau^{(t)}$ to get the expert's action at
 150 those states, forming $\tau^{(t)'} = (s_1, a'_1, \dots, s_H, a'_H)$, where $a'_h \sim \pi_{\eta,h}^*(\cdot | s_h)$.

151 When $\eta = 0$, this recovers the online IL setting studied in Ross et al. [49]. A priori, the online setting
 152 is more powerful than the offline setting and classically, when $\eta = 0$ (i.e., we are in the clean expert
 153 setting) this additional power was reflected in an improved sample complexity of online algorithms
 154 like DAgger [49] compared to offline algorithms like BC. More recently, however, Foster et al.
 155 [17] showed that in the realizable setting, BC is actually optimal and achieves horizon-free regret
 156 bounds. In order to state these bounds, we recall the definition of the Hellinger distance between two
 157 distributions P and Q over the same space \mathcal{X} :

158 **Definition 3** (Hellinger Distance). The Hellinger distance between two distributions P and Q over
 159 the same space \mathcal{X} is defined as $D_{\text{H}^2}(P, Q) = 1/2 \cdot \int_{\mathcal{X}} (\sqrt{dP/d\mu(x)} - \sqrt{dQ/d\mu(x)})^2 d\mu(x)$, where μ
 160 is any distribution that dominates both P and Q .

161 While we defer to Polyanskiy and Wu [43] for a detailed discussion of the properties of the Hellinger
 162 distance, we provide a brief exposition in Appendix C.1. In particular, the Hellinger distance is
 163 intimately related to the more commonly used total variation distance up to a quadratic factor, and
 164 it is upper bounded by the KL divergence via Pinsker's inequality. The following result shows that
 165 regret is intimately related to the Hellinger distance between the trajectory distributions of the expert
 166 and learned policies.

167 **Theorem 4** (Theorem 2.1 & 3.1 from [17]). *Let π^* and $\hat{\pi}$ be any two policies in an MDP with reward*
 168 *range R . Then it holds that*

$$J(\pi^*) - J(\hat{\pi}) \lesssim R \cdot \sqrt{D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}})} + R \cdot D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}).$$

169 *Moreover, if π^* is deterministic, then it holds that $J(\pi^*) - J(\hat{\pi}) \lesssim R \cdot D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}})$.*

170 In addition, as we recall in Appendix C.3, Foster et al. [17] showed that the above reduction is tight,
 171 in the sense that for any two policies π^* and $\hat{\pi}$, there exists a reward function such that the Hellinger
 172 distance lower bounds the regret up to constant factors. Theorem 4 thus motivates us to focus on

173 controlling the Hellinger distance between the trajectory distributions of the expert and learned
 174 policies in order to get good regret guarantees; because of this, we will often abuse nomenclature and
 175 refer to $D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}})$ as *regret*. Thus, the question we answer in this work is the following:

176 *How much interaction with a noisy expert π_η^* is required in order to find a policy $\hat{\pi}$*
 177 *such that $D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) \leq \varepsilon$ and how does online access to the expert affect this*
 178 *sample complexity?*

179 We will at times focus on the case where the corruption distribution ν satisfies a certain *domination*
 180 condition with respect to the policy class Π , defined as follows.

181 **Definition 4** (κ -Domination). We say that a corruption distribution ν is κ -dominated by policies
 182 π, π' if for any time step $h \in [H]$, state $s \in \mathcal{S}$, and action $a \in \text{supp}(\pi_h(\cdot | s)) \cup \text{supp}(\pi'_h(\cdot | s))$, it
 183 holds that $\nu_h(a|s) \leq \kappa \cdot (\pi_h(a|s) + \pi'_h(a|s))$. We will say that ν is κ -dominated by a policy class Π
 184 if for fixed $\pi^* \in \Pi$, ν is κ -dominated by π^* and π for any $\pi \in \Pi$.

185 While this condition may appear somewhat unmotivated,¹ we will show that it is necessary in order
 186 to provide horizon-free guarantees in the noisy expert setting. The condition is weaker than it may
 187 first appear, as the likelihood domination only needs to hold on the support of the policies π and π' ,
 188 which can be much smaller than the entire action space. Indeed, note that it is always satisfied with
 189 $\kappa = 1$ for deterministic policies π , which we will study in detail in [Section 5](#).

190 3 Offline Imitation Learning with a Noisy Expert

191 We begin by considering the offline setting, where the learner only has access to a dataset of
 192 trajectories generated by the noisy expert π_η^* , i.e. $\tau^1, \dots, \tau^n \sim \mathbb{P}^{\pi_\eta^*}$. In this and the next sections,
 193 we will focus on the *known* corruption setting, where both η and ν are known to the learner and
 194 generalize to *unknown* corruptions in [Section 5](#). While the known corruption setting is often not
 195 realistic in practice, it allows the main ideas to be presented with significantly greater clarity, as well
 196 as providing a useful benchmark for the more difficult unknown corruption setting. In the classical
 197 regime, wherein the learner receives clean expert feedback, Foster et al. [[17](#)] showed:

198 **Theorem 5** (Theorem 2.1 Foster et al. [[17](#)]). *Let $\pi^* \in \Pi$ and suppose $\tau^1, \dots, \tau^n \sim \mathbb{P}^{\pi^*}$. If*

$$\hat{\pi} \in \underset{\pi \in \Pi}{\operatorname{argmin}} \sum_{i=1}^n \sum_{h=1}^H -\log \pi(a_h^{(i)} | s_h^{(i)}), \quad (2)$$

199 *then with probability at least $1 - \delta$, it holds that $D_{\text{H}^2}(\mathbb{P}^{\hat{\pi}}, \mathbb{P}^{\pi^*}) \lesssim \log(|\Pi|/\delta)/n$.*

200 This analysis is based on the key observation that (2) is the *Maximum Likelihood Estimator* (MLE)
 201 of the trajectory distribution of the expert π^* based on the observed trajectories, and thus we can
 202 apply classical results on the convergence of MLEs in Hellinger distance [[20](#), [64](#)]. While this is a
 203 powerful result when given clean expert feedback, in the noisy expert setting considered here, we
 204 only have access to trajectories from π_η^* rather than π^* , and thus, even when η and ν are known, we
 205 can guarantee closeness in Hellinger distance only to $\mathbb{P}^{\pi_\eta^*}$ rather than \mathbb{P}^{π^*} . Indeed, it would thus be
 206 natural to replace (2) with the following estimator:

$$\hat{\pi} \in \underset{\pi \in \Pi}{\operatorname{argmin}} \sum_{i=1}^n \sum_{h=1}^H -\log \pi_\eta(a_h^{(i)} | s_h^{(i)}), \quad (3)$$

207 where $\pi_\eta = (1 - \eta) \cdot \pi + \eta \cdot \nu$ is the noisy version of π . Critically, while the identical analysis as in
 208 [Theorem 5](#) shows that $D_{\text{H}^2}(\mathbb{P}^{\pi_\eta^*}, \mathbb{P}^{\hat{\pi}_\eta})$ is guaranteed to be small, our ultimate goal is to compete
 209 with the *clean* expert π^* rather than the noisy expert π_η^* , and thus we need to understand how this
 210 noisy Hellinger distance controls the clean Hellinger distance $D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}})$.

211 **Theorem 6.** *Let π, π' be policies that κ -dominate ν . Then, $D_{\text{H}^2}(\mathbb{P}^\pi, \mathbb{P}^{\pi'}) \lesssim$
 212 $(1 + \eta \cdot \kappa) (1 - \eta)^{-H-2} \cdot D_{\text{H}^2}(\mathbb{P}^{\pi_\eta}, \mathbb{P}^{\pi'_\eta})$. In particular, if $\hat{\pi}$ is as in (3) then with probability at*

¹This assumption arises naturally due to the singularity of the map $\eta \mapsto \sqrt{\eta}$ near $\eta = 0$ (cf. [Remark 1](#)).

213 least $1 - \delta$, it holds that $D_{\text{H}^2}(\mathbb{P}^{\hat{\pi}}, \mathbb{P}^{\pi^*}) \lesssim (1 + \eta \cdot \kappa)(1 - \eta)^{-H-2} \cdot \log(|\Pi|/\delta)/n$. In general, absent
 214 κ -domination it holds that $D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) \lesssim \sqrt{((1-\eta)^{-H}-1)/\eta(1-\eta)} \cdot D_{\text{H}^2}(\mathbb{P}^{\hat{\pi}_\eta}, \mathbb{P}^{\pi^*_\eta})$.

215 The proof is deferred to [Appendix D.1](#), where the key idea is to apply backward induction by peeling
 216 off the trajectory distributions one step at a time, and then to use the κ -domination condition to
 217 control the error at each step. Note that in the absence of κ -domination, as $\eta \downarrow 0$, we do not recover a
 218 tight dependence as a polynomial in H factor appears, which is tight.

219 While [Theorem 6](#) shows that as long as η and ν are known, a natural analogue of BC is consistent, the
 220 sample complexity is exponential in the horizon H whenever $\eta \gg 1/H$, even when the clean expert
 221 π^* is deterministic, in contradistinction to the clean setting where no explicit horizon dependence
 222 exists. Thus in the long horizon regime, e.g. LMs generating long chains of thought in order to solve
 223 a hard math problem, this upper bound becomes vacuous. Unfortunately, our next result shows that
 224 this exponential dependence is necessary.

225 **Proposition 1.** *For any $H \geq 2$, $\kappa \geq 1$, and $\eta < 1$ there exists a horizon H MDP with 3 actions,
 226 policies π^* , $\hat{\pi}$, and κ -dominated corruption distribution ν such that $D_{\text{H}^2}(\mathbb{P}^{\hat{\pi}}, \mathbb{P}^{\pi^*}) \gtrsim \eta \cdot \kappa \cdot (1 -$
 227 $\eta)^{-H-1} \cdot D_{\text{H}^2}(\mathbb{P}^{\hat{\pi}_\eta}, \mathbb{P}^{\pi^*_\eta})$. Moreover, if $\kappa = 1$, we can take π^* to be deterministic. In the absence
 228 of κ -domination, we may even force $D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) \gtrsim \sqrt{((1-\eta)^{-H}-1)/1-\eta} \cdot D_{\text{H}^2}(\mathbb{P}^{\hat{\pi}_\eta}, \mathbb{P}^{\pi^*_\eta})$.*

229 The constructions witnessing the above result can be found in [Appendix D.2](#). This result demonstrates
 230 that even in the easiest possible setting of noisy IL, where the contamination is known and the expert
 231 is deterministic, any offline IL algorithm must have sample complexity that depends exponentially
 232 on the horizon H in the worst case, presenting a fundamental barrier in the noisy expert setting.
 233 Moreover, the result demonstrates the fundamental necessity of κ -domination. Even worse, as we
 234 show in [Proposition 15](#), a corollary of this result demonstrates that *any* offline IL algorithm must
 235 suffer from this exponential dependence in sample complexity, making offline IL fundamentally
 236 intractable in the noisy expert setting. In the next section, we show that online IL can reduce this
 237 exponential dependence on the horizon to at most polynomial.

238 4 On-Policy Distillation Allows for Improved Guarantees

239 In the previous section, we showed that in the noisy expert setting, offline IL algorithms such as BC
 240 must suffer from an exponential dependence on the horizon H , which is unacceptable in the kinds
 241 of long-horizon tasks that are ubiquitous in modern AI applications [[23](#), [42](#), [54](#)]. In this section, we
 242 demonstrate that *online* access to the expert can circumvent this exponential dependence, in direct
 243 contradistinction to the classical setting where Foster et al. [[17](#)] showed that when given clean expert
 244 feedback, BC is optimal even in the presence of *online expert interaction*.

245 In order to motivate our main result of the section showing that online access can greatly improve the
 246 dependence on the horizon in the noisy expert setting, we first emphasize that the primary obstacle to
 247 offline IL in the noisy expert regime is precisely the tightness of [Theorem 6](#): in order to recover the
 248 clean expert’s trajectory distribution in Hellinger distance, we must have exponentially better control
 249 on the Hellinger distance between the trajectory distributions of the noisy expert and noisy learned
 250 policy. While this relationship is unfortunately tight, it is natural to wonder if a stronger guarantee
 251 on the closeness of the noisy policy trajectories can lead to one paying a smaller constant factor in
 252 the relationship with the clean policy trajectories. In order to show that this is indeed the case, we
 253 introduce the notation $\mathbb{P}^{\pi, \pi'}$ for augmented trajectory distributions. More precisely, for policies π, π'
 254 we say that $\tau' = (s_1, a'_1, \dots, s_H, a'_H) \sim \mathbb{P}^{\pi, \pi'}$ if $\tau = (s_1, a_1, \dots, s_H, a_H) \sim \mathbb{P}^\pi$ and $a'_h \sim \pi'(\cdot | s_h)$.
 255 In other words, $\mathbb{P}^{\pi, \pi'}$ is the distribution over trajectories obtained by rolling out π in the MDP but
 256 then taking actions according to π' instead of π ; note that $\mathbb{P}^{\pi, \pi}$ is in general not the same as \mathbb{P}^π if
 257 π is stochastic due to the resampling of actions, but they do have the same marginal distributions.
 258 The following result shows that if we change our notion of divergence from Hellinger to KL and we
 259 consider augmented trajectory distributions, then we can get a much tighter relationship.

260 **Theorem 7.** *Let π^* , $\hat{\pi}$ κ -dominate the corruption ν . Then for any $0 \leq \eta < 1$ it holds that*

$$D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) \lesssim (1 + \eta \cdot \kappa)/(1 - \eta)^2 \cdot (D_{\text{KL}}(\mathbb{P}^{\hat{\pi}, \pi^*_\eta} \| \mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta}) \wedge D_{\text{KL}}(\mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta} \| \mathbb{P}^{\hat{\pi}, \pi^*_\eta})). \quad (4)$$

Algorithm 1 NAIL: Noise-robust Aggregation for Imitation Learning

Require: Number of rounds n , policy class Π , noisy expert π_η^* , corruption level η , corruption distribution ν .

- 1: Initialize $w_1 = \text{Unif}(\Pi)$.
 - 2: **for** $t = 1$ to n **do**
 - 3: Define $\mu_t = \sum_{\pi \in \Pi} w_t(\pi) \cdot \pi$ and deploy μ_t to get trajectory $\tau^{(t)} \sim \mathbb{P}^{\mu_t}$.
 - 4: Query noisy expert π_η^* on $\tau^{(t)}$ to obtain augmented trajectory $\tau^{(t)'}$.
 - 5: Update $w_{t+1}(\pi) \propto w_t(\pi) \cdot \left(\prod_{h=1}^H \pi_\eta(a_h^{(t)' } | s_h^{(t)}) \right)^{1/H}$.
 - 6: **end for**
 - 7: **return** $\hat{\pi} = \mu_T$ where $T \sim \text{Unif}([n])$.
-

261 *Absent the κ -domination condition, it holds that*

$$D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) \lesssim (1 - \eta)^{-1} \cdot \sqrt{H \cdot (D_{\text{KL}}(\mathbb{P}^{\hat{\pi}, \pi_\eta^*} \| \mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta}) \wedge D_{\text{KL}}(\mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta} \| \mathbb{P}^{\hat{\pi}, \pi_\eta^*})}. \quad (5)$$

262 We defer a proof of this result to [Appendix E.1](#). The key idea is to use the *subadditivity* of the
263 Hellinger distance from Foster et al. [16, 18] to control the trajectory Hellinger distance by the sum of
264 the per-step Hellinger distances, and then demonstrate that these per-step Hellinger distances can be
265 controlled by the KL divergence of the noisy augmented trajectory distributions. Critically, we have
266 to use KL divergence instead of Hellinger distance in the upper bound in order to stitch the per-step
267 distances back into a divergence over the *trajectory distributions* because KL divergence satisfies
268 the *chain rule* (cf. [Proposition 4](#)) while Hellinger distance does not. Moreover, in [Proposition 17](#) we
269 show that (5) is tight up to a polynomial dependence on $(1 - \eta)$, and thus the improvement from
270 [Theorem 7](#) is indeed significant in the presence of κ -domination.

271 Note that, unlike the offline setting, [Theorem 7](#) is fully *independent of horizon* in the presence of
272 κ -domination and scales only polynomially with H in general. Moreover, the upper bounds in
273 [Theorem 7](#) are closely related to the popular *On-Policy Distillation* (OPD) algorithm [2, 41], which
274 is used to finetune a student policy to match the behavior of a teacher policy by minimizing a loss
275 closely related to the right hand side of (4). Indeed, a common approach is to roll out a policy
276 π_t and then use a target π_η^* to score the actions taken by π_t , in particular attempting to minimize
277 $D_{\text{KL}}(\mathbb{P}^{\pi_t} \| \mathbb{P}^{\pi_t, \pi_\eta^*})$. While our analysis suggests instead minimizing divergence on the *augmented*
278 *trajectory distributions*, this provides some explanation for why OPD can be so much more successful
279 than the offline BC (i.e., SFT) in practice.

280 While variants of OPD remain the recommended empirical approach ([Section 6](#)) and [Theorem 7](#)
281 provides theoretical backing for the success thereof, we continue in this section by introducing a
282 simple new algorithm, NAIL ([Algorithm 1](#)). At its core, much like the classical DAGger algorithm
283 [49], NAIL uses online learning to aggregate policies that are trained on the noisy expert’s behavior
284 on trajectories rolled out by the learner. Unlike some instantiations of DAGger, however, NAIL
285 is *always on-policy*, in the sense that it rolls out the current policy mixture μ_t at each round [49].
286 Moreover, NAIL learns the policy at a *trajectory level*, thereby circumventing the horizon dependence
287 that DAGger necessarily incurs by learning a different policy at each time step h [17, 49]. More
288 precisely, NAIL begins with a uniform distribution w_1 over policies in Π and at each round t rolls out
289 the mixture policy $\mu_t = \sum_{\pi \in \Pi} w_t(\pi) \cdot \pi$ to get a trajectory $\tau^{(t)}$, queries the noisy expert π_η^* on $\tau^{(t)}$
290 to get an augmented trajectory $\tau^{(t)'}$, and then updates w_{t+1} using the standard exponential weights
291 update [8] (cf. [Appendix C.4](#)) with the loss given by the negative log-likelihood of $\tau^{(t)'}$ under each
292 policy π . Finally, NAIL returns a random policy across the time steps. We show that NAIL achieves
293 the following regret guarantee in the noisy expert setting.

294 **Theorem 8.** *Let $\pi^* \in \Pi$ be an arbitrary policy, ν be an arbitrary corruption distribution, and*
295 *$0 \leq \eta < 1$ be a corruption level. If $\hat{\pi}$ is the policy returned by NAIL ([Algorithm 1](#)) after n rounds*
296 *with feedback from π_η^* , then $\mathbb{E} [D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}})] \lesssim H/1-\eta \cdot \sqrt{\log(|\Pi|)/n}$.*

297 We defer a proof of this result to [Appendix E.3](#), which proceeds by using the *mixability* (cf. [Ap-](#)
298 [pendix C.4](#)) of the trajectory-level log-loss to control the sum of on-policy trajectory-level KL
299 divergences across rounds $1 \leq t \leq n$, before applying [Theorem 7](#) and the convexity of the Hellinger

300 distance to get the desired bound. We reiterate that NAIL should be thought of as a theoretical
 301 algorithm that is designed to be optimal in the worst case and that directly minimizing the right hand
 302 side of (4) is a more practical approach to OPD. We further remark that, while NAIL uses exponential
 303 weights, really any online learning algorithm with small regret for the trajectory-level log-loss could
 304 be used in its place, as Theorem 7 allows for a generic reduction to online learning in the spirit of
 305 Ross et al. [49]. We now show that NAIL is optimal up to factors in H and η , although we do not
 306 rule out a sharper bound under κ -domination: the correct rate under this condition remains one of the
 307 chief theoretical questions left open by our work.

308 **Proposition 2.** *For any $H \geq 2$, $0 < \eta < 1$, and small ε there exists a horizon H MDP with 3 actions,
 309 a policy class Π of size $|\Pi| = 2$ containing π^* , and known corruption ν such that any algorithm with
 310 online access to the noisy expert π_η^* requires $n \gtrsim \eta \cdot H / (1-\eta)^2 \cdot \varepsilon^2$ in order to obtain regret ε .*

311 We defer a proof to Appendix E.4. To summarize: there is a marked distinction between the clean and
 312 noisy expert settings: in contradistinction to standard IL, when the expert is noisy, online interaction
 313 can exponentially improve over offline access to the expert.

314 5 Imitating Deterministic Experts with Unknown Corruptions

315 In the previous sections, we focused on the setting of noisy IL with a *known corruption*. While this
 316 setting is simple, it is somewhat unrealistic; we now relax this assumption. The first problem we run
 317 into is one of *identifiability*: without additional assumptions on the corruption, it may well be the case
 318 that it is impossible to identify the expert policy π^* from the noisy expert π_η^* , even absent stochasticity
 319 in π^* ; a concrete example of this phenomenon is provided in Appendix F.1. While this does not
 320 present an issue in standard IL, where closeness to the observed actions is the key desideratum, in
 321 the noisy expert setting, we care about learning a policy that is close to the *clean* expert and thus
 322 identifiability is critical. This issue does not arise when ν and η are known, as clean experts can be
 323 recovered from their noisy analogues, but recovery is not possible absent known corruption.

324 While many identifiability assumptions are possible, we focus on a natural one, satisfied in many
 325 of our motivating applications. We first restrict our focus to *deterministic experts*, motivated by
 326 applications to reasoning, math, or coding in LMs, where we think of the true expert as producing a
 327 fixed ‘correct’ output [23, 54]. We further assume that the corruption distribution ν is ‘smooth.’

328 **Definition 5.** We say that a corruption distribution ν is ρ -smooth if for all time steps $h \in [H]$, states
 329 $s \in \mathcal{S}$, and actions $a \in \mathcal{A}$, it holds that $\nu_h(a|s) \leq \rho$.

330 Smoothness is a natural assumption when the action space is large, for example when we think of
 331 corruption as arising from typos or other minor errors in the output of an LM or human demonstrator.
 332 In particular, if ν is a uniform distribution over a large action space, then ν is ρ -smooth with
 333 $\rho = |\mathcal{A}|^{-1}$. Finally, we assume a bound on the noise level such that $\eta \leq \alpha < 1$ to ensure that a
 334 *margin* exists between the expert and the noise, which is necessary for identifiability.²

335 We propose an algorithm, NAILGUN (Algorithm 3), similar to NAIL, but that is capable of achieving
 336 *horizon-free* regret guarantees in the noisy expert setting with unknown corruption under the above
 337 assumptions. While we defer a detailed description of the algorithm to Appendix F.2, we note that
 338 it is very similar to NAIL, except that we replace the mixture policy μ_t by its greedy analogue $\bar{\mu}_t$,
 339 which plays the action with the highest probability under μ_t at each state and time step; moreover, the
 340 precise form of the update is different in order to account for the fact that we do not know η and ν .

341 **Theorem 9.** *Suppose that $\pi^* \in \Pi$ is deterministic, ν is ρ -smooth, and $\eta \leq \alpha$ for some α satisfying
 342 $\alpha(1 + \rho) < 1$. There exists an algorithm (Algorithm 3) without knowledge of η or ν that returns a
 343 policy $\hat{\pi}$ after n rounds of online interaction such that $\mathbb{E} [D_{H^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}})] \lesssim \log(|\Pi|) / (1 - \alpha(1 + \rho))^2 \cdot n$.*

344 We defer a proof of Theorem 9 to Appendix F.3, which proceeds through an intricate analysis of the
 345 update rule and how the margin condition that $\alpha(1 + \rho) < 1$ allows us to control the error without
 346 knowledge of η or ν . Note that unlike in the previous section, the particular form of the update rule is
 347 critical to achieving this guarantee, and thus we do not prove a generic reduction to online learning.
 348 We now conclude this section with a matching lower bound.

²We expect that we could replace determinism and smoothness with a more general margin condition and recover the same results, but we leave this for future work.

349 **Proposition 3.** For any $0 < \alpha < 1$ and any ρ satisfying $\alpha(1 + \rho) < 1$, there exists a horizon $H = 2$
 350 MDP and a deterministic policy class Π of size $|\Pi| \lesssim 1/\rho$ such that in order to achieve expected
 351 regret at most ε , the learner requires $n \gtrsim \alpha\rho/\varepsilon(1-\alpha(1+\rho))^2$ rounds of interaction with an online noisy
 352 expert satisfying $\eta \leq \alpha$ and ν is ρ -smooth and unknown.

353 The proof of this lower bound can be found in [Appendix F.4](#). In particular, in the noisy expert regime
 354 where α and ρ are treated as constants, NAILGUN is optimal, even for horizon $H = 2$.

355 6 Experiments

356 We complement our theory with a small suite of experiments on synthetic and natural language tasks,
 357 showing that the OPD variant suggested by our theory (NAIL) can outperform both offline BC and
 358 existing OPD objectives under noisy expert feedback while remaining competitive under clean expert
 359 feedback. In particular, given an expert π_η^* , we compare five algorithms: (i) `LogLossBC` (SFT), which
 360 minimizes the negative log-likelihood of the noisy expert trajectories; (ii) `OPD-F`, the standard OPD
 361 algorithm minimizing the forward KL divergence between student and teacher trajectory distributions
 362 [2]; (iii) `OPD-R`, the OPD algorithm that minimizes the reverse KL divergence between student and
 363 teacher trajectory distributions [2, 41]; (iv) `NAIL-F`, the variant of NAIL minimizing forward KL
 364 divergence between the augmented trajectory distributions; and (v) `NAIL-R`, the variant of NAIL
 365 minimizing reverse KL divergence. In particular, the primary difference between OPD and NAIL
 366 is that NAIL always generates *greedily* from the student policy, whereas OPD samples trajectories
 367 from the student policy; we emphasize that we do not instantiate the exponential weights update in
 368 [Algorithm 1](#) but instead use gradient descent to minimize the upper bound in [Theorem 7](#). We consider
 369 two tasks: (i) modular addition and (ii) `GSM-8K` [14]. Additional details are in [Appendix B.1](#).

370 6.1 Modular Addition

371 We begin with a synthetic task from Li et al. [38], that of modular addition. Namely, for fixed p and
 372 m , given a sequence of m integers from $[p]$, the goal is to compute the sum of these integers modulo
 373 p . In the cited work, the authors demonstrated that low-depth transformers empirically benefit from a
 374 chain of thought (CoT) reasoning process in this task, making it a suitable testbed for studying the
 375 effect of horizon. All models are nanoGPT [30]. We train our expert π^* on a large corpus of data
 376 with $p = 7$ and $m = 31$ until its accuracy (with CoT) is approximately perfect. We then construct a
 377 noisy expert π_η^* by letting ν be uniform over tokens and corrupting the expert with probability η at
 378 each token if we have not yet seen an error, and corrupting all remaining tokens with probability 1 if
 379 we have already seen an error. We then train the student models using the five algorithms described
 380 above and evaluate the accuracy of the final learned policy $\hat{\pi}$ on a held-out test set of examples. We
 381 defer further details to [Appendix B.1.3](#). [Figure 1](#)(left) summarizes the results of this experiment. For
 382 the clean expert ($\eta = 0$) regime, we see that most methods work, with offline BC (i.e., SFT) learning
 383 the fastest, which is consistent with the fact that offline BC is optimal in the clean expert setting [17].
 384 For the noisy expert ($\eta = 0.2$) regime, we see a sharp separation between our proposed methods
 385 (`NAIL-F`, `NAIL-R`) and the baselines (`LogLossBC`, `OPD-F`, `OPD-R`), with the former achieving
 386 perfect accuracy and the latter failing to learn anything, even after over 1M expert interactions. We
 387 defer further comments and ablations to [Appendix B.2](#) and [Appendix B.3](#).

388 6.2 Mathematical Reasoning

389 We also consider a standard mathematical reasoning task, `GSM-8K` [14]. Here, we move beyond the
 390 realizable setting of our theory and use `GEMMA3-1B-IT` as the teacher and `GEMMA3-270M-IT` as
 391 the student [29]. To instantiate clean and noisy expert feedback with the same underlying model, we
 392 consider temperature 1 sampling as the ‘clean’ expert, which achieves 52.69% accuracy on `GSM-8K`
 393 test, and temperature 4 sampling as the ‘noisy’ expert, which achieves 0% accuracy. All training
 394 prompts are drawn from `TinyGSM` [40] and we evaluate zero-shot, exact-match accuracy on the
 395 `GSM-8K` test set with greedy decoding, reporting results over three random seeds. We defer further
 396 details to [Appendix B.1.4](#). We report results in [Figure 1](#)(right), where we see that all online methods
 397 substantially outperform `LogLossBC` even in the clean expert setting, consistent with results from
 398 Agarwal et al. [2], Lu and Lab [41]. In the noisy expert setting, only `NAIL-F` and `NAIL-R` remain
 399 effective, achieving nontrivial performance while `LogLossBC`, `OPD-F`, and `OPD-R` all fail to learn
 400 anything, maintaining 0 accuracy, like the noisy expert, throughout training.

References

- 401
- 402 [1] Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar,
403 Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. Phi-4 technical
404 report. *arXiv preprint arXiv:2412.08905*, 2024.
- 405 [2] Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos Garea,
406 Matthieu Geist, and Olivier Bachem. On-policy distillation of language models: Learning from
407 self-generated mistakes. In *The twelfth international conference on learning representations*,
408 2024.
- 409 [3] Awni Altabaa, Omar Montasser, and John Lafferty. Cot information: Improved sample com-
410 plexity under chain-of-thought supervision. In *The Thirty-ninth Annual Conference on Neural*
411 *Information Processing Systems*, 2026.
- 412 [4] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,
413 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- 414 [5] Jose Barreiros, Andrew Beaulieu, Aditya Bhat, Rick Cory, Eric Cousineau, Hongkai Dai, Ching-
415 Hsin Fang, Kunimatsu Hashimoto, Muhammad Zubair Irshad, Masha Itkina, et al. A careful
416 examination of large behavior models for multitask dexterous manipulation. *Science Robotics*,
417 11(113):eaea6201, 2026.
- 418 [6] Adam Block, Ali Jadbabaie, Daniel Pfrommer, Max Simchowitz, and Russ Tedrake. Provable
419 guarantees for generative behavior cloning: Bridging low-level stability and high-level behavior.
420 *Advances in Neural Information Processing Systems*, 36:48534–48547, 2023.
- 421 [7] Adam Block, Dylan J Foster, Akshay Krishnamurthy, Max Simchowitz, and Cyril Zhang.
422 Butterfly effects of sgd noise: Error amplification in behavior cloning and autoregression. In
423 *The Twelfth International Conference on Learning Representations*, 2024.
- 424 [8] Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university
425 press, 2006.
- 426 [9] Jonathan Chang, Kianté Brantley, Rajkumar Ramamurthy, Dipendra Misra, and Wen Sun.
427 Learning to generate better than your llm. In *NeurIPS 2023 Workshop on Instruction Tuning*
428 *and Instruction Following*, 2023.
- 429 [10] Fan Chen, Audrey Huang, Noah Golowich, Sadhika Malladi, Adam Block, Jordan T. Ash,
430 Akshay Krishnamurthy, and Dylan J Foster. The coverage principle: How pre-training enables
431 post-training. In *The Fourteenth International Conference on Learning Representations*, 2026.
- 432 [11] Jianyu Chen, Bodi Yuan, and Masayoshi Tomizuka. Deep imitation learning for autonomous
433 driving in generic urban scenarios with enhanced safety. In *2019 IEEE/RSJ international*
434 *conference on intelligent robots and systems (IROS)*, pages 2884–2890. IEEE, 2019.
- 435 [12] Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning
436 converts weak language models to strong language models. In *International Conference on*
437 *Machine Learning*, pages 6621–6642. PMLR, 2024.
- 438 [13] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ
439 Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion.
440 *The International Journal of Robotics Research*, 44(10-11):1684–1704, 2025.
- 441 [14] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
442 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to
443 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 444 [15] Dylan J Foster and Akshay Krishnamurthy. Efficient first-order contextual bandits: Prediction,
445 allocation, and triangular discrimination. *Advances in Neural Information Processing Systems*,
446 34:18907–18919, 2021.
- 447 [16] Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity
448 of interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021.

- 449 [17] Dylan J Foster, Adam Block, and Dipendra Misra. Is behavior cloning all you need? understanding horizon in imitation learning. *Advances in Neural Information Processing Systems*, 37: 450 120602–120666, 2024.
- 452 [18] Dylan J Foster, Yanjun Han, Jian Qian, and Alexander Rakhlin. Online estimation via offline 453 estimation: An information-theoretic framework. *Advances in Neural Information Processing 454 Systems*, 37:42840–42898, 2024.
- 455 [19] Yuqian Fu, Haohuan Huang, Kaiwen Jiang, Jiakai Liu, Zhuo Jiang, Yuanheng Zhu, and Dongbin 456 Zhao. Revisiting on-policy distillation: Empirical failure modes and simple fixes. *arXiv preprint 457 arXiv:2603.25562*, 2026.
- 458 [20] Sara A Geer. *Empirical Processes in M-estimation*, volume 6. Cambridge university press, 459 2000.
- 460 [21] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. Knowledge distillation: A 461 survey. *International journal of computer vision*, 129(6):1789–1819, 2021.
- 462 [22] Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. MiniLLM: Knowledge distillation of large 463 language models. In *The Twelfth International Conference on Learning Representations*, 2024.
- 464 [23] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, 465 Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. DeepSeek-R1: Incentivizing reasoning capability in 466 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 467 [24] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. 468 *arXiv preprint arXiv:1503.02531*, 2015.
- 469 [25] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural 470 information processing systems*, 29, 2016.
- 471 [26] Jonas Hübotter, Frederike Lübeck, Lejs Behric, Anton Baumann, Marco Bagatella, Daniel Marta, 472 Ido Hakimi, Idan Shenfeld, Thomas Kleine Buening, Carlos Guestrin, et al. Reinforcement 473 learning via self-distillation. *arXiv preprint arXiv:2601.20802*, 2026.
- 474 [27] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun 475 Liu. Tinybert: Distilling bert for natural language understanding. In *Findings of the association 476 for computational linguistics: EMNLP 2020*, pages 4163–4174, 2020.
- 477 [28] Nirmal Joshi, Gal Vardi, Adam Block, Surbhi Goel, Zhiyuan Li, Theodor Misiakiewicz, and 478 Nathan Srebro. A theory of learning with autoregressive chain of thought. In *The Thirty Eighth 479 Annual Conference on Learning Theory*, pages 3161–3212. PMLR, 2025.
- 480 [29] Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, 481 Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, et al. Gemma 3 482 technical report. *arXiv preprint arXiv:2503.19786*, 2025.
- 483 [30] Andrej Karpathy. NanoGPT. <https://github.com/karpathy/nanoGPT>, 2022.
- 484 [31] Yoon Kim and Alexander M Rush. Sequence-level knowledge distillation. In *Proceedings of 485 the 2016 conference on empirical methods in natural language processing*, pages 1317–1327, 486 2016.
- 487 [32] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd 488 International Conference on Learning Representations, ICLR 2015*, 2015.
- 489 [33] Jongwoo Ko, Sungnyun Kim, Tianyi Chen, and Se-Young Yun. DISTILLM: towards streamlined 490 distillation for large language models. In *Proceedings of the 41st International Conference on 491 Machine Learning*, pages 24872–24895, 2024.
- 492 [34] Michael Laskey, Jonathan Lee, Roy Fox, Anca Dragan, and Ken Goldberg. Dart: Noise injection 493 for robust imitation learning. In *Conference on robot learning*, pages 143–156. PMLR, 2017.

- 494 [35] Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Gadre, Hritik Bansal,
495 Etash Guha, Sedrick Keh, Kushal Arora, et al. Datacomp-lm: In search of the next generation
496 of training sets for language models. *Advances in Neural Information Processing Systems*, 37:
497 14200–14282, 2024.
- 498 [36] Yichen Li and Chicheng Zhang. Agnostic interactive imitation learning: New theory and
499 practical algorithms. *arXiv preprint arXiv:2312.16860*, 2023.
- 500 [37] Yichen Li and Chicheng Zhang. Interactive and hybrid imitation learning: Provably beating
501 behavior cloning. *arXiv preprint arXiv:2412.07057*, 2024.
- 502 [38] Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers
503 to solve inherently serial problems. In *The Twelfth International Conference on Learning*
504 *Representations*, 2024.
- 505 [39] Nick Littlestone. Learning quickly when irrelevant attributes abound: A new linear-threshold
506 algorithm. *Machine learning*, 2(4):285–318, 1988.
- 507 [40] Bingbin Liu, Sebastien Bubeck, Ronen Eldan, Janardhan Kulkarni, Yuanzhi Li, Anh Nguyen,
508 Rachel Ward, and Yi Zhang. TinyGSM: achieving 80% on GSM8k with one billion parameters.
509 In *The 3rd Workshop on Mathematical Reasoning and AI at NeurIPS'23*, 2023.
- 510 [41] Kevin Lu and Thinking Machines Lab. On-policy distillation. *Thinking Machines Lab: Con-*
511 *nectionism*, 2025. doi: 10.64434/tml.20251026. [https://thinkingmachines.ai/blog/on-policy-](https://thinkingmachines.ai/blog/on-policy-distillation)
512 [distillation](https://thinkingmachines.ai/blog/on-policy-distillation).
- 513 [42] Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita
514 Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, et al. 2 OLMo 2 furious. *arXiv preprint*
515 *arXiv:2501.00656*, 2024.
- 516 [43] Yury Polyanskiy and Yihong Wu. *Information theory: From coding to learning*. Cambridge
517 university press, 2025.
- 518 [44] Dean A Pomerleau. Alvin: An autonomous land vehicle in a neural network. *Advances in*
519 *neural information processing systems*, 1, 1988.
- 520 [45] Dhruv Rohatgi, Adam Block, Audrey Huang, Akshay Krishnamurthy, and Dylan J Foster.
521 Computational-statistical tradeoffs at the next-token prediction barrier: Autoregressive and
522 imitation learning under misspecification. In *The Thirty Eighth Annual Conference on Learning*
523 *Theory*, pages 4831–4837. PMLR, 2025.
- 524 [46] Elvezio M Ronchetti and Peter J Huber. *Robust statistics*. John Wiley & Sons Hoboken, NJ,
525 USA, 2009.
- 526 [47] Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. In *Proceedings of*
527 *the thirteenth international conference on artificial intelligence and statistics*, pages 661–668.
528 JMLR Workshop and Conference Proceedings, 2010.
- 529 [48] Stéphane Ross and J Andrew Bagnell. Reinforcement and imitation learning via interactive
530 no-regret learning. *arXiv preprint arXiv:1406.5979*, 2014.
- 531 [49] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and
532 structured prediction to no-regret online learning. In *Proceedings of the fourteenth interna-*
533 *tional conference on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and
534 Conference Proceedings, 2011.
- 535 [50] Andrei A Rusu, Sergio Gomez Colmenarejo, Caglar Gulcehre, Guillaume Desjardins, James
536 Kirkpatrick, Razvan Pascanu, Volodymyr Mnih, Koray Kavukcuoglu, and Raia Hadsell. Policy
537 distillation. *arXiv preprint arXiv:1511.06295*, 2015.
- 538 [51] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. DistilBERT, a distilled
539 version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.

- 540 [52] Ayush Sekhari, Karthik Sridharan, Wen Sun, and Runzhe Wu. Contextual bandits and imitation
541 learning with preference-based active queries. *Advances in Neural Information Processing*
542 *Systems*, 36:11261–11295, 2023.
- 543 [53] Ayush Sekhari, Karthik Sridharan, Wen Sun, and Runzhe Wu. Selective sampling and imitation
544 learning via online regression. *Advances in Neural Information Processing Systems*, 36:67213–
545 67268, 2023.
- 546 [54] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
547 Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical
548 reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 549 [55] Mingyang Song and Mao Zheng. A survey of on-policy distillation for large language models.
550 *arXiv preprint arXiv:2604.00626*, 2026.
- 551 [56] Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1.
552 MIT press Cambridge, 1998.
- 553 [57] Gokul Swamy, Sanjiban Choudhury, Drew Bagnell, and Steven Wu. Causal imitation learning
554 under temporally correlated noise. In *International Conference on Machine Learning*, pages
555 20877–20890. PMLR, 2022.
- 556 [58] Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xiong-Hui
557 Chen, Jianxin Yang, Zhenru Zhang, Yuqiong Liu, An Yang, Andrew Zhao, Yang Yue, Shiji
558 Song, Bowen Yu, Gao Huang, and Junyang Lin. Beyond the 80/20 rule: High-entropy minority
559 tokens drive effective reinforcement learning for LLM reasoning. In *The Thirty-ninth Annual*
560 *Conference on Neural Information Processing Systems*, 2026.
- 561 [59] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep
562 self-attention distillation for task-agnostic compression of pre-trained transformers. *Advances*
563 *in neural information processing systems*, 33:5776–5788, 2020.
- 564 [60] Maurice Weber, Daniel Y. Fu, Quentin Anthony, Yonatan Oren, Shane Adams, Anton Alexan-
565 drov, Xiaozhong Lyu, Huu Nguyen, Xiaozhe Yao, Virginia Adams, Ben Athiwaratkun, Rahul
566 Chalamala, Kezhen Chen, Max Ryabinin, Tri Dao, Percy Liang, Christopher Ré, Irina Rish, and
567 Ce Zhang. RedPajama: an open dataset for training large language models. *NeurIPS Datasets*
568 *and Benchmarks Track*, 2024.
- 569 [61] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
570 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*
571 *arXiv:2505.09388*, 2025.
- 572 [62] Chenxu Yang, Chuanyu Qin, Qingyi Si, Minghui Chen, Naibin Gu, Dingyu Yao, Zheng Lin,
573 Weiping Wang, Jiaqi Wang, and Nan Duan. Self-distilled rlvr. *arXiv preprint arXiv:2604.03128*,
574 2026.
- 575 [63] Ruixiang Zhang, Richard He Bai, Huangjie Zheng, Navdeep Jaitly, Ronan Collobert, and
576 Yizhe Zhang. Embarrassingly simple self-distillation improves code generation. *arXiv preprint*
577 *arXiv:2604.01193*, 2026.
- 578 [64] Tong Zhang. From ε -entropy to kl-entropy: Analysis of minimum information complexity
579 density estimation. *The Annals of Statistics*, pages 2180–2210, 2006.
- 580 [65] Siyan Zhao, Zhihui Xie, Mengchen Liu, Jing Huang, Guan Pang, Feiyu Chen, and Aditya
581 Grover. Self-distilled reasoner: On-policy self-distillation for large language models. *arXiv*
582 *preprint arXiv:2601.18734*, 2026.
- 583 [66] Binbin Zheng, Xing Ma, Yiheng Liang, Jingqing Ruan, Xiaoliang Fu, Kepeng Lin, Benchang
584 Zhu, Ke Zeng, and Xunliang Cai. SCOPE: Signal-calibrated on-policy distillation enhancement
585 with dual-path adaptive weighting. *arXiv preprint arXiv:2604.10688*, 2026.

586	Contents	
587	1 Introduction	1
588	2 Formal Problem Setup and Preliminaries	3
589	3 Offline Imitation Learning with a Noisy Expert	5
590	4 On-Policy Distillation Allows for Improved Guarantees	6
591	5 Imitating Deterministic Experts with Unknown Corruptions	8
592	6 Experiments	9
593	6.1 Modular Addition	9
594	6.2 Mathematical Reasoning	9
595	A Discussion	15
596	A.1 Limitations	15
597	A.2 Related Work	15
598	A.3 Future Directions	17
599	A.4 Broader Impacts	18
600	B Further Empirical Results	18
601	B.1 Further Experimental Details	18
602	B.1.1 Methods and Implementations	18
603	B.1.2 Gradient estimators.	20
604	B.1.3 Synthetic Task: Modular Addition	21
605	B.1.4 Mathematical Reasoning: GSM-8K	22
606	B.2 Further Discussion of Empirical Results in Section 6	23
607	B.3 Further Experiments	23
608	B.3.1 Ablating student rollout temperature	24
609	B.3.2 Interpolating between NAIL-F and NAIL-R	25
610	C Technical Tools	25
611	C.1 Information Theory	26
612	C.2 Maximum Likelihood Estimation	27
613	C.3 Imitation Learning	28
614	C.4 Online Learning	29
615	D Proofs from Section 3	30
616	D.1 Proof of Theorem 6	30
617	D.2 Proof of Proposition 1	34
618	D.3 Implications for Offline Imitation Learning	37

619	E Proofs from Section 4	37
620	E.1 Proof of Theorem 7	37
621	E.2 Necessity of Horizon Dependence in the Absence of κ -Domination	39
622	E.3 Proof of Theorem 8	39
623	E.4 Proof of Proposition 2	40
624	F Proofs from Section 5	42
625	F.1 The Problem of Identifiability	42
626	F.2 Description of the Algorithm	43
627	F.3 Proof of Theorem 9	43
628	F.4 Proof of Proposition 3	46

629 A Discussion

630 In this work we investigated the problem of imitation learning with noisy experts, where the learner
631 only has access to a corrupted version of the expert’s policy. We showed theoretically that while
632 offline imitation learning can be consistent in this setting, it suffers from an exponential dependence
633 on the horizon, making it fundamentally intractable in the long horizon regime. On the other hand,
634 we showed that online imitation learning can circumvent this exponential dependence and achieve
635 horizon-free guarantees in some cases and polynomial in horizon guarantees in general, in direct
636 contradistinction to the more classical IL setting with clean experts. Moreover, we provided formal
637 theoretical justification for the benefits of the commonly used *On-Policy Distillation* (OPD) algorithm
638 and validated our theory with a small suite of experiments on synthetic and natural language tasks. A
639 summary of our theoretical results is in Table 1. In this section, we discuss related work and future
640 directions for research.

641 A.1 Limitations

642 While our work presents a clean theoretical picture, there are several limitations in immediately
643 translating our results to practice. First, consistent with the literature in RL and IL, our results are
644 all in the finite Π setting [17, 45, 47, 49]; standard arguments can be used to extend our results to
645 infinite Π under covering conditions. Second, we adopt an idealized corruption model, which may
646 not capture all the nuances of real-world expert noise; we defer to future work the task of developing
647 more sophisticated corruption models that may better capture the types of noise observed in practice.
648 Third, many of our results are in the known corruption setting, and we only provide results in the
649 unknown corruption setting when the expert is deterministic; while this setting is relevant for many
650 of our motivating applications, e.g. LMs learning to perform long chains of thought, it would be
651 interesting to understand the unknown corruption setting more generally. Fourth, there remains a
652 linear in horizon gap between our upper bound for arbitrary experts and our lower bound; closing this
653 gap would be an interesting direction for future work; similarly, whether or not fully horizon-free
654 rates are possible for online IL in the presence of κ -domination is an interesting question. Finally,
655 due to resource constraints, our experiments are limited in scope and scale; it would be interesting to
656 conduct a more comprehensive empirical investigation of the phenomena we identify here, especially
657 in the context of training larger language models for longer.

658 A.2 Related Work

659 **Imitation Learning.** Behavior cloning was first introduced empirically in Pomerleau [44] and
660 has since engendered the field of Imitation Learning, with a large body of work considering both
661 theoretical and empirical aspects of the problem [5, 6, 13, 25, 47]. Of particular note is Ross et al.
662 [49], which introduced the online IL algorithm DAGger; while NAIL is conceptually similar (asking
663 for expert feedback to correct potential student mistakes), the analysis and guarantees are quite
664 different, with DAGger attempting to reduce horizon dependence in a clean expert setting.

Table 1: Summary of the main theoretical results in the general noisy-expert setting. Here $\text{KL}_{\text{aug}}(\pi)$ denotes the smaller of the two augmented-trajectory KL divergences appearing in Theorem 7. Bounds control $D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}})$ and are stated up to universal constants and logarithmic factors when indicated.

Regime	Setting	Bound	Reference
Offline			
	No κ -domination.	$\sqrt{\frac{(1-\eta)^{-H}-1}{\eta(1-\eta)} \cdot \frac{\log(\Pi /\delta)}{n}}$	Theorem 6 and Proposition 1 ; this is tight.
	κ -domination.	$(1+\eta\cdot\kappa)(1-\eta)^{-H-2} \cdot \frac{\log(\Pi /\delta)}{n}$	Proposition 13 ; this is tight.
Online			
	Augmented-trajectory KL reduction without κ -domination.	$\frac{1}{1-\eta} \sqrt{H \cdot \text{KL}_{\text{aug}}(\hat{\pi})}$	Theorem 7 ; this is tight up to factors in H and η .
	Augmented-trajectory KL reduction with κ -domination.	$\frac{1+\eta\cdot\kappa}{(1-\eta)^2} \cdot \text{KL}_{\text{aug}}(\hat{\pi})$	Theorem 7
	NAIL / trajectory-level on-line aggregation	$\frac{H}{1-\eta} \sqrt{\frac{\log \Pi }{n}}$	Theorem 8 . The lower bound in Proposition 2 is a factor of H and η away from this upper bound.
	Deterministic π^* , ρ -smooth ν , and $\eta \leq \alpha$ with $\alpha(1+\rho) < 1$	$\frac{\log \Pi }{(1-\alpha(1+\rho))^2 n}$	Theorem 9 and Proposition 3 ; this is tight up to a log factor.

665 Other IL algorithms study different feedback and data-collection models. SMILe [47] works in a
666 similar fashion as DAgger to collect learner-induced states under clean expert feedback, but instead
667 learns a sequence of policies and mixes them to control distribution shift. AggreVate [48] strengthens
668 the feedback model by using expert cost-to-go or advantage information, and GAIL [25] instead
669 casts IL as adversarial occupancy-measure matching from fixed expert demonstrations. DART [34]
670 collects off-policy demonstrations from a noise-injected expert so that the demonstrations include
671 recovery behavior.

672 More recently, Foster et al. [17] provided a comprehensive analysis of behavior cloning in the
673 realizable, clean expert setting and Rohatgi et al. [45] extended this to the agnostic setting. In the
674 online setting, Li and Zhang [37] considered a variant of DAgger where the learner pays per-state
675 instead of per-trajectory in expert feedback, while Li and Zhang [36] also considers the agnostic setting
676 but with strong assumptions on the MDP structure and a horizon dependence in their guarantees.
677 In contrast to these works, we consider a more challenging setting with noisy experts and show
678 that online IL can achieve horizon-free guarantees in some settings and polynomial dependence on
679 horizon in general even in this more difficult setting, while offline IL suffers from an exponential
680 dependence on the horizon. On the empirical side, OPD has seen recent success in training large
681 language models (LLMs) to perform long chains of thought [2, 4, 41], and our theoretical results
682 provide formal justification for the benefits of OPD in this setting.

683 The connections between IL and LMs have been explored in a number of recent works [7, 9, 17, 45],
684 with the deterministic expert setting being particularly relevant to LMs learning long chains of thought
685 [3, 28]. Our results generalize beyond this setting, allowing for stochastic experts and corrupted
686 feedback.

687 **Noisy Expert Feedback.** Most prior work in IL assumes access to a clean expert, with the notable
688 exceptions of Sekhari et al. [52, 53]. While a promising start, and the former work provides an
689 exponential in horizon lower bound for offline learning in a particular setting, these works make strong
690 assumptions on the noise and MDP structure and only apply to specific policy classes; moreover,
691 they only consider a very specific type of corruption arising from preference-based feedback, in
692 contradistinction to our bounds which apply to arbitrary policy classes and MDPs. Another related
693 work that concerns a noisy-expert model is that of Swamy et al. [57], which focuses (under very strong
694 assumptions) on causal IL under temporally correlated noise, where action noise changes future
695 states and creates spurious state-action correlations that BC can learn. Our corruption model has the
696 same rolled-in character, but our focus is different: rather than debiasing offline demonstrations via
697 causal assumptions, we show that online local queries can avoid learning from fully contaminated
698 trajectories.

699 **Distillation.** Knowledge Distillation is typically the process of training a student model to match
700 a stronger teacher, emphasizing compression and efficient deployment. Early work emphasized
701 transferring softened teacher predictions from cumbersome ensembles to smaller students [24],
702 and the same idea was later extended to reinforcement-learning policies [50] and to autoregressive
703 sequence models through sequence-level distillation [31]. Critically, this continued as an off-policy
704 paradigm: the student is trained on a fixed transfer set, teacher-generated sequences, or teacher
705 logits, as in [27, 51, 59]; see also the general survey of [21]. While highly effective for compression,
706 these methods do not train on prefixes induced by the student’s own generations, and are therefore
707 susceptible to the train–test mismatch inherent in autoregressive generation, a well known problem in
708 imitation learning [44, 47].

709 On-policy distillation [2] addresses this mismatch by training the student on self-generated sequences
710 while receiving teacher feedback on those sequences using varying discrepancy measures between
711 teacher and student. One objective in distillation is the forward KL from the teacher to the student,
712 which corresponds in our terminology to OPD-F. Prior work of Agarwal et al. [2] has noted that this
713 objective can be problematic when the student is not expressive enough to fit the teacher distribution:
714 because forward KL is mode-covering, the resulting student may place probability on generations that
715 are unlikely under the teacher. Subsequent work [22, 41] has therefore explored alternative on-policy
716 objectives, advocating for sequence/trajectory-level reverse KL to avoid the mode-covering behavior
717 of forward KL in generative settings, while [33] smooths the loss and reduces rollout cost with an
718 adaptive off-policy reuse scheme.

719 Our results offer a more unifying perspective. In the noisy-expert setting, the difficulty is not only
720 which KL direction is used, but also which prefix distribution the teacher is queried on and which
721 feedback is informative. Sampling from the student can itself induce noisy or off-distribution prefixes,
722 reducing the usefulness of local teacher feedback. This helps explain empirical findings that sampled-
723 token OPD can be brittle, and that stability can improve by restricting both the rollout distribution
724 and the teacher signal used for local updates [19]. It is also consistent with recent work that makes
725 OPD updates more selective, for example by down-weighting teacher KL on failed trajectories where
726 the teacher itself is uncertain [66]. More generally, work on pivotal or “forking” tokens suggests that
727 a small number of CoT tokens can determine which reasoning branch the model follows, and hence
728 whether downstream feedback remains informative [1, 58]. Our theory gives a principled prescription:
729 matching teacher feedback along the learner’s deployed prefix distribution; in our experiments with
730 deterministic experts, using greedy student rollouts together with augmented trajectory-level KL
731 objectives helps suppress sampling-induced noise and learn more effectively from imperfect teachers.
732 We leave a more detailed study of when forward versus reverse KL is preferable to future work. More
733 recently, self-play and self-distillation methods have shown that on-policy data generation can itself
734 be a source of improvement even without a separate external teacher [12, 26, 62, 63, 65]. For a more
735 detailed overview, we refer to the recent survey of Song and Zheng [55].

736 A.3 Future Directions

737 While we provided a near-tight theoretical analysis of the problems studied, there remain many
738 interesting open questions for future work. First, in the unknown corruption setting, we focused on
739 the case of deterministic experts and smooth corruptions with margin, but it would be interesting to
740 understand the extent to which these assumptions can be relaxed while still achieving horizon-free
741 guarantees, or even consistency. Second, while we considered the natural linear corruption setting,

742 motivated by classical statistical theory [46] and applications to LMs, alternative corruption models
743 are also interesting to study, for example trajectory-level, multiplicative, or token-adaptive corruptions
744 that target semantically pivotal “fork” tokens in long CoT reasoning [1, 10, 58, 66]. Third, our online
745 IL guarantees are *in expectation* due to the inherited guarantees from online learning with log loss,
746 whereas the offline guarantees are *with high probability* and resolving this discrepancy could provide
747 additional insight. Finally, scaling up the experimental results to more complex environments and
748 real-world tasks remains an important direction for future work.

749 **A.4 Broader Impacts**

750 This work is primarily theoretical and aims to improve our understanding of imitation learning from
751 imperfect expert feedback. We hope that showing fundamental separation between offline and online
752 methods (such as on-policy distillation) will lead to more reliable and sample-efficient post-training
753 methods, especially in long-horizon settings.

754 While our methods could be seen as enabling effective learning from imperfect teachers, they could
755 also make it easier to train models for harmful applications if deployed without appropriate safeguards.
756 Our experiments are limited to synthetic tasks and standard mathematical reasoning benchmarks, and
757 we do not release deployed systems or high-risk models. We therefore view the main impact of this
758 work as methodological, with downstream risks depending on how future systems built using these
759 ideas are trained, evaluated, and released.

760 **B Further Empirical Results**

761 This section provides additional details and discussion for the experiments in [Section 6](#). [Appendix B.1](#)
762 describes the implementation of the methods in the autoregressive LM setting, as well as the synthetic
763 and reasoning tasks setup. [Appendix B.2](#) gives a more detailed interpretation of the main empirical
764 results, including the differences between clean and noisy feedback and between greedy and sampled
765 rollouts. Finally, [Appendix B.3](#) reports additional experiments and ablations.

766 **B.1 Further Experimental Details**

767 In this section, we provide implementation details for the empirical results in [Section 6](#). [Ap-](#)
768 [pendix B.1.1](#) defines the offline and online learning objectives that we compare, including behavior
769 cloning, standard OPD, and our NAIL-inspired variants; in [Appendix B.1.2](#) we spell out their gradient
770 estimators. [Appendix B.1.3](#) describes the setup of the synthetic CoT tasks, including the data-
771 generation procedure, corruption model, architectures, and training details. Finally, [Appendix B.1.4](#)
772 gives the corresponding details for the GSM8K experiments, including the teacher and student models,
773 sampling procedure, training setup, and evaluation protocol.

774 We used NVIDIA RTX 6000 Pro GPUs for the modular addition experiments, with each run taking 1
775 hour. In total, the modular addition experiments used 50 GPU hours. For the GSM8K experiments, we
776 used NVIDIA A100 40GB GPUs. Each online-method run took approximately 60 hours, while each
777 offline-method run took approximately 3 hours; this disparity in runtime is an artifact of our highly
778 suboptimal hardware setup and a naïve implementation and we do not believe such a discrepancy
779 would exist in a more optimal setup. In total, the GSM8K experiments used roughly 800 GPU hours.

780 **B.1.1 Methods and Implementations**

781 We compare offline behavior cloning against several on-policy distillation objectives. The key
782 distinction among the online methods is that there are two separate choices: the policy used to
783 generate prefixes, and the divergence used to compare the student and teacher distributions on those
784 prefixes. This distinction is motivated by [Theorem 7](#), which controls the clean-expert trajectory error
785 through KL divergences between *augmented* trajectory laws. In particular, the rollout distribution
786 determines which prefixes are visited, while the per-prefix student distribution determines the next-
787 token distribution being matched to the noisy expert.

788 Let π_θ denote the trainable student policy and let $\bar{\pi}_\theta$ denote the greedy rollout policy induced by π_θ .
789 Thus $y_{1:T} \sim \bar{\pi}_\theta(\cdot | x)$ denotes the trajectory obtained by greedy decoding from the student, while
790 $\pi_\theta(\cdot | x, y_{<t})$ denotes the student’s full next-token distribution on the resulting prefix. We use π_η^* to

791 denote the noisy teacher. In the synthetic experiments, π_η^* is obtained by corrupting the clean expert
 792 token with the specified corruption law. In the language-model experiments, it is implemented by
 793 sampling from the teacher at the specified temperature.

794 **NAIL-F.** The forward-KL version of our method (which minimizes the first term in (4)) uses greedy
 795 student prefixes and trains the student to match noisy-teacher next-token feedback on those prefixes.
 796 Ideally, this corresponds to minimizing

$$\mathcal{L}_{\text{NAIL-F}}(\theta) = \mathbb{E}_{x \sim \mathcal{X}} \mathbb{E}_{y_{1:T} \sim \bar{\pi}_\theta(\cdot|x)} \left[\sum_{t=1}^T \text{D}_{\text{KL}}(\pi_\eta^*(\cdot | x, y_{<t}) \| \pi_\theta(\cdot | x, y_{<t})) \right],$$

797 Equivalently, up to the entropy of the noisy teacher, this is a soft-target cross-entropy objective:

$$-\mathbb{E}_{x \sim \mathcal{X}} \mathbb{E}_{y_{1:T} \sim \bar{\pi}_\theta(\cdot|x)} \left[\sum_{t=1}^T \sum_{a \in \Sigma} \pi_\eta^*(a | x, y_{<t}) \log \pi_\theta(a | x, y_{<t}) \right].$$

798 Since our implementation queries the teacher by sampling, we estimate this loss by drawing an
 799 independent teacher token $\tilde{y}_t \sim \pi_\eta^*(\cdot | x, y_{<t})$ at each greedy student prefix and minimizing

$$\hat{\mathcal{L}}_{\text{NAIL-F}}(\theta) = \mathbb{E}_{x \sim \mathcal{X}} \mathbb{E}_{y_{1:T} \sim \bar{\pi}_\theta(\cdot|x)} \mathbb{E}_{\tilde{y}_t \sim \pi_\eta^*(\cdot|x, y_{<t})} \left[\sum_{t=1}^T -\log \pi_\theta(\tilde{y}_t | x, y_{<t}) \right].$$

800 Thus, NAIL-F can be thought of as “local behavior cloning” on the learner’s own greedy prefixes:
 801 rather than imitating complete noisy teacher trajectories, it repeatedly asks what the noisy teacher
 802 would do at states actually reached by greedily rolling out the current student, and then trains the
 803 student to predict the noisy teacher’s sample via cross-entropy. This is the empirical analogue of the
 804 forward KL augmented-trajectory objective suggested by our theory. When η and ν are known, the
 805 literal theoretical objective replaces π_θ above by the noisy student policy $\pi_{\theta, \eta} = (1 - \eta)\pi_\theta + \eta\nu$; in
 806 our experiments we use the simpler student-matching surrogate above.

807 **NAIL-R.** We also consider the reverse KL analogue suggested by minimizing the second term in (4).
 808 As in NAIL-F, prefixes are generated by the greedy student rollout, but the KL direction is reversed:

$$\mathcal{L}_{\text{NAIL-R}}(\theta) = \mathbb{E}_{x \sim \mathcal{X}} \mathbb{E}_{y_{1:T} \sim \bar{\pi}_\theta(\cdot|x)} \left[\sum_{t=1}^T \text{D}_{\text{KL}}(\pi_\theta(\cdot | x, y_{<t}) \| \pi_\eta^*(\cdot | x, y_{<t})) \right].$$

809 We estimate the student expectation in the reverse KL by drawing an auxiliary token $\hat{y}_t \sim \pi_\theta(\cdot |$
 810 $x, y_{<t})$ at each greedy prefix:

$$\hat{\mathcal{L}}_{\text{NAIL-R}}(\theta) = \mathbb{E}_{x \sim \mathcal{X}} \mathbb{E}_{y_{1:T} \sim \bar{\pi}_\theta(\cdot|x)} \mathbb{E}_{\hat{y}_t \sim \pi_\theta(\cdot|x, y_{<t})} \left[\sum_{t=1}^T \log \pi_\theta(\hat{y}_t | x, y_{<t}) - \log \pi_\eta^*(\hat{y}_t | x, y_{<t}) \right].$$

811 We note that the auxiliary token \hat{y}_t is not necessarily the greedy rollout token y_t . The greedy rollout
 812 tokens determine the prefix distribution, while the auxiliary samples estimate the student expectation
 813 in the reverse KL. This separation is needed to match the augmented-trajectory viewpoint: the rollout
 814 policy and the next-token distribution being compared need not be the same object.

815 **OPD-F.** As an ablation, we also consider the forward-KL objective with sampled student rollouts
 816 rather than greedy rollouts:

$$\mathcal{L}_{\text{OPD-F}}(\theta) = \mathbb{E}_{x \sim \mathcal{X}} \mathbb{E}_{y_{1:T} \sim \pi_\theta(\cdot|x)} \left[\sum_{t=1}^T \text{D}_{\text{KL}}(\pi_\eta^*(\cdot | x, y_{<t}) \| \pi_\theta(\cdot | x, y_{<t})) \right],$$

817 This isolates the effect of the rollout distribution: comparing OPD-F to NAIL-F tests whether greedy
 818 learner prefixes are important, while comparing OPD-F to OPD-R tests the effect of KL direction
 819 under the same sampled-prefix distribution.

820 **OPD-R.** Finally, we include the standard on-policy distillation baseline using reverse KL [41]. This
 821 method samples trajectories from the student, queries the teacher on those same student-generated
 822 prefixes, and updates the student using the teacher log-probability of the sampled student tokens:

$$\mathcal{L}_{\text{OPD-R}}(\theta) := \mathbb{E}_{x \sim \mathcal{X}} \mathbb{E}_{y \sim \pi_\theta(\cdot|x)} \left[\sum_{t=1}^T \text{D}_{\text{KL}}(\pi_\theta(\cdot | x, y_{<t}) \| \pi_\eta^*(\cdot | x, y_{<t})) \right].$$

823 Using the sampled rollout token $y_t \sim \pi_\theta$, this becomes:

$$\widehat{\mathcal{L}}_{\text{OPD-R}}(\theta) = \mathbb{E}_{x \sim \mathcal{X}} \mathbb{E}_{y_{1:T} \sim \pi_\theta(\cdot | x)} \left[\sum_{t=1}^T \left(\log \pi_\theta(y_t | x, y_{<t}) - \log \pi_\eta^*(y_t | x, y_{<t}) \right) \right].$$

824 **LogLossBC.** The offline baseline is log-loss behavior cloning on a fixed dataset of noisy expert
825 rollouts [17]. Let \mathcal{D}_η denote trajectories generated by rolling out the noisy expert π_η^* . We train

$$\mathcal{L}_{\text{LogLossBC}}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}_\eta} \left[- \sum_{t=1}^{|y|} \log \pi_\theta(y_t | x, y_{<t}) \right].$$

826 Unlike the OPD variants, behavior cloning never queries the teacher on prefixes induced by the
827 current student. Consequently, if an early token in an offline teacher trajectory is corrupted, the
828 learner is also trained on all downstream prefixes induced by that corruption.

829 B.1.2 Gradient estimators.

830 We now spell out the stochastic gradient estimators used to optimize the empirical losses defined
831 above. Let $s_t = (x, y_{<t})$ denote a visited prefix, and write $p_{\theta,t}(\cdot) = \pi_\theta(\cdot | s_t)$ and $q_t(\cdot) = \pi_\eta^*(\cdot | s_t)$.
832 Throughout, $p_{\theta,t}$ denotes the temperature-one student distribution whose KL is optimized. Rollout-
833 temperature ablations only change the distribution used to collect prefixes. In particular, for rollout
834 temperature τ , let $\rho_{\theta,\tau}$ denote the stopped prefix-collection policy, with

$$\rho_{\theta,0} = \bar{\pi}_\theta, \quad \rho_{\theta,\tau}(\cdot | s_t) = \text{softmax}(z_\theta(s_t)/\tau) \quad \text{for } \tau > 0.$$

835 All gradients below treat prefixes sampled from $\rho_{\theta,\tau}$ as fixed; the gradient update is taken only
836 through the next-token distribution $p_{\theta,t}$ at those visited prefixes. The default settings are $\tau = 0$ for
837 NAIL-F and NAIL-R and $\tau = 1$ for OPD-F and OPD-R, while the rollout-temperature ablations
838 vary this prefix-collection temperature.

839 For $\widehat{\mathcal{L}}_{\text{NAIL-F}}(\theta)$, after drawing prefixes $y_{1:T} \sim \rho_{\theta,\tau}(\cdot | x)$ and independent teacher tokens $\tilde{y}_t \sim q_t$, we
840 use

$$\widehat{\mathcal{G}}_\tau^{\text{F}} = - \sum_{t=1}^T \nabla_\theta \log p_{\theta,t}(\tilde{y}_t).$$

841 This is an unbiased stochastic gradient estimator for the stopped-prefix forward-KL objective because
842 the entropy term of q_t is independent of θ . The methods differ only in how the prefixes are collected:
843 $\tau = 0$ gives the default NAIL-F estimator with greedy prefixes, while $\tau = 1$ gives the default
844 sampled-prefix OPD-F estimator. Other rollout temperatures use the same next-token estimator on
845 prefixes sampled from $\rho_{\theta,\tau}$. The offline LogLossBC baseline is ordinary cross-entropy on realized
846 noisy-expert trajectories; it has the same token-level gradient form, but its prefixes and labels come
847 from an offline dataset rather than from student rollouts.

848 For $\widehat{\mathcal{L}}_{\text{NAIL-R}}(\theta)$, prefix collection and reverse-KL token sampling are separate. After drawing prefixes
849 $y_{1:T} \sim \rho_{\theta,\tau}(\cdot | x)$, the implementation draws an auxiliary token from a stopped copy of the current
850 student distribution at the visited prefix: $\hat{y}_t \sim \text{sg}(p_{\theta,t})$, where $\text{sg}(\cdot)$ denotes stop-gradient. Thus, the
851 rollout temperature τ affects which prefixes are visited, but the auxiliary token for the reverse-KL
852 estimator is sampled from $p_{\theta,t}$, not from the rollout distribution. The stopped-prefix score-function
853 estimator is

$$\widehat{\mathcal{G}}_\tau^{\text{R}} = \sum_{t=1}^T (\text{sg}(\log p_{\theta,t}(\hat{y}_t)) - \log q_t(\hat{y}_t)) \nabla_\theta \log p_{\theta,t}(\hat{y}_t), \quad \hat{y}_t \sim \text{sg}(p_{\theta,t}).$$

854 Equivalently, this is implemented by backpropagating through the surrogate

$$\widehat{\ell}_t^{\text{R}}(\theta) = - \exp(\log p_{\theta,t}(\hat{y}_t) - \text{sg}(\log p_{\theta,t}(\hat{y}_t))) (\log q_t(\hat{y}_t) - \text{sg}(\log p_{\theta,t}(\hat{y}_t))).$$

855 The exponential factor is numerically equal to one at the sampled parameter value, but its numerator
856 still carries gradient, so differentiating gives exactly the estimator above.

857 For $\widehat{\mathcal{L}}_{\text{OPD-R}}(\theta)$, as in standard OPD, the sampled rollout token itself is used as the reverse-KL sample.
 858 Therefore, if the rollout token is sampled from the current temperature-one student distribution, this
 859 coincides with the on-policy reverse-KL estimator above. If one instead reuses rollout tokens sampled
 860 at a different temperature, then the token distribution no longer matches $p_{\theta,t}$, and the resulting update
 861 should be viewed as a temperature-mismatched surrogate. Note that this mirrors the recipe described
 862 by Lu and Lab [41].

863 For the interpolated objective in Appendix B.3, where

$$\widehat{\mathcal{L}}_{\beta}(\theta) = (1 - \beta)\widehat{\mathcal{L}}_{\text{NAIL-F}}(\theta) + \beta\widehat{\mathcal{L}}_{\text{NAIL-R}}(\theta),$$

864 we use the corresponding convex combination of the two stopped-prefix estimators,

$$\widehat{g}_{\beta} = (1 - \beta)\widehat{g}_0^{\text{F}} + \beta\widehat{g}_0^{\text{R}}.$$

865 B.1.3 Synthetic Task: Modular Addition

866 We use the modular-addition task of Li et al. [38], denoted C_p . Given $p \in \mathbb{N}$, the vocabulary
 867 is $\{0, \dots, p - 1, =\}$. An input has the form $x = (x_1, \dots, x_m, =)$, where each x_i is sampled
 868 independently from $\{0, \dots, p - 1\}$. The final answer is $f^*(x) = \sum_{i=1}^m x_i \bmod p$ and the chain of
 869 thought consists of the running partial sums $(\sum_{i=1}^t x_i \bmod p)_{t=1}^m$. Modular addition is not intended
 870 to be an inherently serial hardness instance; indeed, Li et al. [38] note that it is parallelizable and, in
 871 principle, can be solved by constant-depth transformers with sufficient precision. Nevertheless, their
 872 experiments show that low-depth transformers without CoT can perform near chance on C_7 , while
 873 CoT substantially improves performance, especially at longer sequence lengths. We include it as a
 874 simple controlled setting in which the target CoT has a transparent step-by-step structure, allowing us
 875 to study how noisy intermediate feedback affects offline and online imitation-learning methods.

876 **Training and Evaluation Details.** We use the nanoGPT codebase [30] for all modular-addition
 877 experiments. For the tokenizer, we use task-specific integer token IDs from the aforementioned
 878 symbolic vocabulary. Unless otherwise stated, all runs use a depth-one transformer with 8 attention
 879 heads, embedding dimension 512, dropout 0, no bias terms, and block size set to the CoT sequence
 880 length for the corresponding m . We train with batch size 64 using Adam [32] with learning rate
 881 10^{-5} , warmup for 2000 iterations, weight decay 0, $\beta_1 = 0.9$, $\beta_2 = 0.95$, gradient clipping at 1.0,
 882 and evaluation every 500 iterations.

883 We fix $p = 7$ and $m = 31$. All methods use the same prompt bank, containing 15 million training
 884 prompts and 5000 validation prompts. We train on a fixed 3 million-prompt subset in a fixed order.
 885 The clean expert is fixed across all runs and is synthetically modified using the corruption law
 886 described below. It was trained for 10,000 optimizer steps with batch size 64. The model reached
 887 perfect validation accuracy after roughly 2500 optimizer steps and stayed at this level, but we use the
 888 final checkpoint for all experiments.

889 We compare the five methods for two noise levels, $\eta \in \{0, 0.2\}$, and across three seeds. For each run
 890 seed, we use the same seed for student optimization and for the stochastic components of the noisy
 891 teacher law. For LogLossBC, the seed also determines the independently rendered noisy dataset.
 892 For the online methods, there is no pre-rendered trajectory; the same seed controls the online teacher
 893 sampling.

894 The noisy teacher is an absorbing instantiation of the state-dependent noise model in (1). The state is
 895 augmented with a binary flag indicating whether the prefix has already made a semantic error. In the
 896 unpoisoned state, at each target token the corruption distribution is

$$\nu(\cdot \mid x, y_{<t}) = \text{Unif}(\{0, \dots, p - 1\}),$$

897 so that

$$\pi_{\eta}^*(\cdot \mid x, y_{<t}) = (1 - \eta)\pi^*(\cdot \mid x, y_{<t}) + \eta\nu(\cdot \mid x, y_{<t}).$$

898 Here the uniform distribution includes the clean token. Since every modular-addition target token is a
 899 digit representing a running sum, every target token is treated as semantic. Once the sampled token
 900 differs from the clean running-sum token, the trajectory enters the poisoned state. In the poisoned
 901 state, both π^* and ν are defined to be uniform over $\{0, \dots, p - 1\}$, so all subsequent semantic
 902 feedback is independent of the clean computation while remaining syntactically valid.

903 This absorbing corruption law models a worst-case form of semantic error propagation in long
904 CoT reasoning: after an early pivotal mistake, the suffix may remain syntactically plausible while
905 providing minimal information about the clean continuation. This is motivated by recent work on
906 pivotal or forking tokens in reasoning models, where a small number of tokens can strongly affect the
907 success of the final completion [1, 58]. In our synthetic task, the running-sum tokens play this role, so
908 corrupting one makes downstream semantic feedback uninformative about the correct computation.
909 This pessimistic noise model provides a starting point for a controlled setting in which offline training
910 on corrupted rollouts can suffer from the trajectory-level contamination predicted by our theory, and
911 we leave the consideration of more nuanced noise models to future work.

912 For **LogLossBC**, training trajectories are rendered autoregressively from this noisy teacher law. Thus,
913 in the offline setting, an example is fully informative only if no visible semantic corruption occurs.
914 Under the idealized independent-trigger calculation, the probability of this event is $(1 - \eta + \eta/p)^m \approx$
915 2.9×10^{-3} in our setting.

916 For the online methods, the poisoned flag is inferred from the student-generated prefix: if a previous
917 running-sum token in the student prefix differs from the clean target, then later teacher feedback is
918 uniform over residues. Consequently, $\eta = 0$ has different semantics for offline and online runs under
919 this law: offline rendering is clean at $\eta = 0$, while online feedback can still become uninformative
920 after a student prefix error.

921 Online methods are trained for one pass over the fixed 3 million-prompt subset, giving
922 $3,000,000/64 = 46,875$ iterations. **LogLossBC** is also trained for one pass over the corresponding
923 rendered 3 million-example dataset. All methods are evaluated on the same clean validation prompt
924 bank.

925 **B.1.4 Mathematical Reasoning: GSM-8K**

926 We evaluate on the GSM-8K test set [14] using zero-shot greedy decoding with max generation
927 length 512 and seed 42. Each problem is placed in the Gemma instruction-tuned chat template with
928 the following prefix.

929 `Please reason step by step, and put your final answer within \boxed{}`.

930 For example, an evaluation prompt takes the form

```
931 <bos><start_of_turn>user
932 Please reason step by step, and put your final answer within \boxed{}.
933
934 Janet’s ducks lay 16 eggs per day. She eats three for breakfast every morning
935 and bakes muffins for her friends every day with four. She sells the remainder
936 at the farmers’ market daily for $2 per fresh duck egg. How much in dollars
937 does she make every day at the farmers’ market?<end_of_turn>
938 <start_of_turn>model
```

939 We extract the final answer in the `\boxed{}` and compute exact-match accuracy after standard
940 answer normalization.

941 **Training data.** We construct the training prompt set from TinyGSM [40]. Specifically, we first filter
942 out examples whose code answers contain more than 1024 characters, and then remove examples
943 whose code answers are not executable. This filtering yields approximately 11M examples. We then
944 sample a 400K-example subset that is used for all training runs.

945 **Training details.** For the **LogLossBC** baseline, we generate fixed teacher rollouts on the 400K
946 TinyGSM prompts using GEMMA3-1B-IT. We generate clean rollouts with temperature 1 in the
947 low-noise setting and noisy rollouts with temperature 4 in the high-noise setting, both using random
948 seed 42 and a maximum generation length of 512 new tokens. We measure the accuracy of these clean
949 and noisy teacher rollouts on the 400K TinyGSM subset by executing the reference code answers
950 and treating the execution outputs as ground truth. The clean teacher obtains 50.1% accuracy on
951 the training subset, while the noisy teacher obtains 0% accuracy. We then supervised fine-tune the
952 GEMMA3-270M-IT student separately on the clean and noisy teacher rollouts.

953 For the online reverse-KL variants, OPD-R and NAIL-R, we follow the recipe of Lu and Lab [41].
954 We generate students’ next token using temperature-1 sampling for OPD-R and greedy decoding
955 for NAIL-R. Along these student-generated trajectories, we compute the student and teacher log
956 probabilities conditioned on the same prefixes, using either the clean temperature-1 teacher or the
957 noisy temperature-4 teacher. The student is then updated using the corresponding importance-
958 sampling loss.

959 For the forward-KL variants, OPD-F and NAIL-F, we instead query the teacher locally on learner-
960 visited prefixes. At each such prefix, we sample the teacher’s next token, using the temperature-1
961 teacher in the low-noise setting and the temperature-4 teacher in the high-noise setting, and treat the
962 sampled token as a hard next-token label. The student is then updated using the standard next-token
963 log-loss. Subsequent learner prefixes are generated with temperature-1 sampling for OPD-F and
964 greedy decoding for NAIL-F.

965 All methods are trained on the same 400K TinyGSM prompts with a maximum generation length
966 of 512 new tokens. For the two online methods, we follow the recommended recipe from Lu and
967 Lab [41]. For the offline method, we tune the learning rate slightly in preliminary experiments, while
968 keeping all other hyperparameters fixed. Specifically, unless otherwise stated, all methods are trained
969 for one epoch using AdamW with default parameters, learning rate $1e-4$, linear warmup followed by
970 cosine decay, batch size 64, and bf16 precision. We use LoRA with rank 128 and $\alpha = 256$, applied
971 to all modules. For each method, we run three random seeds, 42, 43, and 44.

972 B.2 Further Discussion of Empirical Results in Section 6

973 All configurations are run with three random seeds. In the plots, each curve shows the mean across
974 seeds, and the shaded region corresponds to one standard deviation. We note that some shaded regions
975 are not visible as the runs are nearly identical across seeds. In particular, for Modular Addition, the
976 maximum standard deviation is less than 10^{-2} for LogLossBC and OPD-F in the low-noise setting,
977 and for OPD-F, OPD-R, and LogLossBC in the high-noise setting (as the latter all fail to learn).

978 For clarity of presentation, the main text shows accuracy only over the first 1M expert trajectories,
979 where the relevant separations between methods are most visible. In Figure 2, we provide the
980 corresponding validation-loss curves over the full 3M-trajectory training horizon. These longer-
981 horizon curves confirm the same qualitative picture: in the low noise setting, the NAIL variants
982 rapidly drive validation loss close to zero, while OPD-F plateaus early; in the high noise setting, the
983 gap becomes more pronounced, with NAIL-F and NAIL-R remaining substantially more stable than
984 the offline and OPD baselines.

985 To expand on the discussion in Section 6, Figure 2 shows a sharp contrast between the clean and
986 noisy regimes for the modular addition task. When $\eta = 0.2$, LogLossBC fails because an early
987 corrupted CoT token makes the remaining suffix random, so most offline trajectories contain little
988 usable signal. In contrast, our NAIL-F and NAIL-R reach perfect accuracy, consistent with the idea
989 that querying on student-induced prefixes avoids imitating fully corrupted trajectories. The failure
990 of sampled-rollout OPD further suggests that controlling the rollout distribution, here via greedy
991 prefixes, is important for keeping teacher feedback informative.

992 When $\eta = 0.0$, since offline traces are clean, LogLossBC learns fastest, as predicted by Foster
993 et al. [17]. Among online methods, NAIL-F learns much faster than OPD-F, suggesting that greedy
994 rollouts help suppress sampling-induced noise that hinders access to uncorrupted teacher feedback.
995 The gap between NAIL-F and the reverse-KL methods here points to the fact that teacher-sampled
996 cross-entropy gives a direct positive signal for the correct next token, whereas reverse KL only scores
997 student-sampled tokens and is therefore less directly corrective. Finally, when there is no expert
998 noise, greedy rollouts offer little benefit for NAIL-R over OPD-R because both methods update using
999 student-sampled tokens.

1000 B.3 Further Experiments

1001 We include two additional ablations on the modular addition task to better understand which parts of
1002 the online objective matter. In Appendix B.3.2, we fix greedy student rollouts and interpolate between
1003 the forward- and reverse-KL augmented-trajectory losses, testing whether performance depends on a
1004 single KL direction.

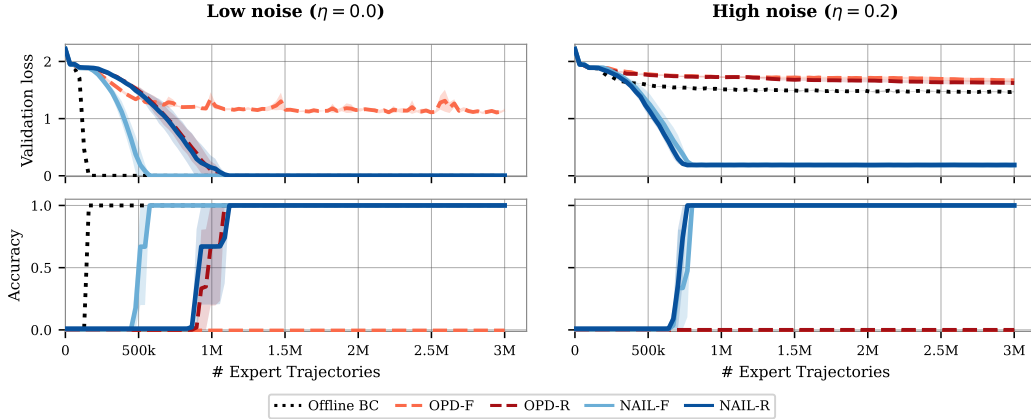


Figure 2: Full modular-addition results over 3M expert trajectories. **Top:** validation loss; **Bottom:** accuracy. Curves show the mean over three random seeds, with shaded regions indicating one standard deviation. In the low-noise setting ($\eta = 0$), the NAIL variants drive validation loss to zero and reach perfect accuracy, while OPD-F plateaus. In the high-noise setting ($\eta = 0.2$), the separation is more pronounced: NAIL-F and NAIL-R remain stable and reach perfect accuracy, whereas offline LogLossBC and the OPD baselines fail to solve the task.

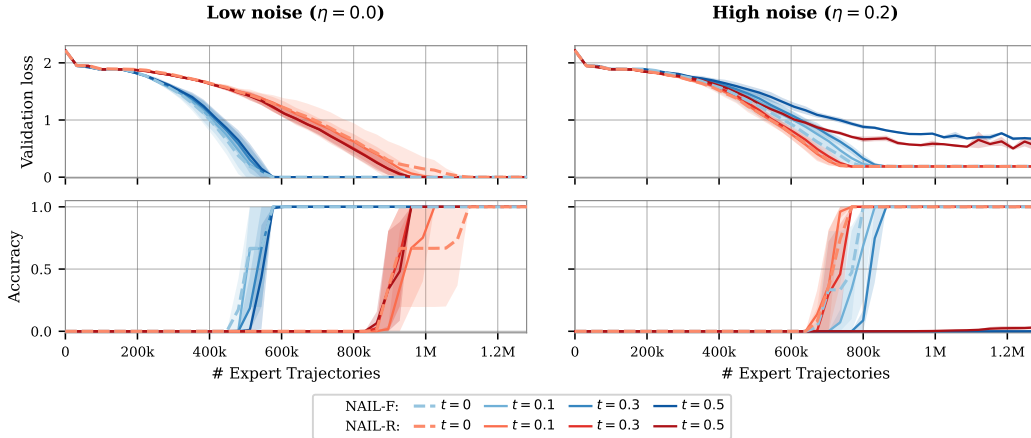


Figure 3: Ablation of student rollout temperature for NAIL-F and NAIL-R on Modular Addition. The parameter t controls the student sampling temperature used during training rollout. Curves show the mean over three random seeds, with shaded regions indicating one standard deviation. **Left:** low-noise setting, NAIL-F learns faster and is relatively insensitive to temperature, while NAIL-R is substantially slower. **Right:** high-noise setting, moderate student temperatures improve robustness for both objectives, with NAIL-F solving the task across temperatures and NAIL-R degrading at larger t .

1005 B.3.1 Ablating student rollout temperature

1006 We next ablate the temperature used when sampling student rollouts. In the main experiments, the
 1007 student prefixes are generated *greedily*, i.e. temperature $t = 0$. Here, we instead sample from the
 1008 student policy with various temperatures $t \in \{0.1, 0.3, 0.5\}$. We evaluate both NAIL-F and NAIL-R
 1009 on Modular Addition under the same low- and high-noise settings as in the main body. All variants
 1010 are trained using the same recipe as in [Appendix B.1](#); for readability, [Figure 3](#) shows only the first
 1011 1.25M expert-query trajectories.

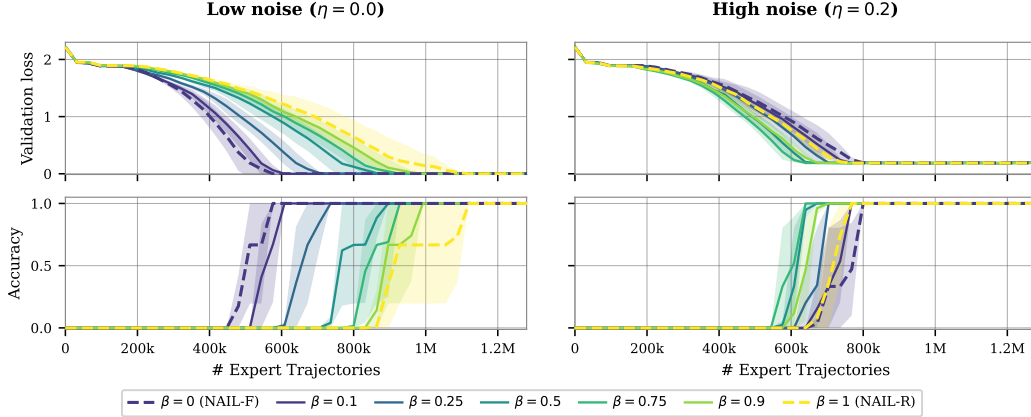


Figure 4: Interpolation between NAIL-F and NAIL-R on Modular Addition. The parameter β interpolates between the forward-KL ($\beta = 0$) and the reverse-KL ($\beta = 1$) losses. Curves show the mean over three random seeds, with shaded regions indicating one standard deviation. All interpolated variants eventually solve the task in both noise regimes, but the learning speed depends strongly on β . **Left:** in the low-noise setting, forward-KL-heavy objectives learn fastest. **Right:** in the high-noise setting, intermediate and reverse-KL-heavy objectives are competitive, suggesting that (i) the robustness of NAIL is primarily driven by querying the expert on learner-induced prefixes rather than by a single KL direction, and (ii) the best KL direction is task- and noise-dependent.

1012 In the low-noise setting ($\eta = 0$), NAIL-F solves the task at roughly the same sample complexity
 1013 across rollout temperatures, indicating that forward-KL training is fairly robust to this choice. In
 1014 contrast, NAIL-R remains slower across temperatures, consistent with the main results in Figure 1. In
 1015 the high-noise setting ($\eta = 0.2$), small to moderate rollout temperatures can still solve the task, but
 1016 large temperature substantially hurts both NAIL-F and NAIL-R. This suggests that some stochasticity
 1017 in student rollouts is tolerable, but excessive exploration can produce prefixes that are too noisy or
 1018 off-distribution for effective expert querying.

1019 B.3.2 Interpolating between NAIL-F and NAIL-R

1020 To further probe the role of the KL direction, we fix greedy student rollouts and interpolate between
 1021 the forward- and reverse-KL losses. For $\beta \in [0, 1]$, we minimize

$$J_{\beta}(\pi_{\theta}) = (1 - \beta) \cdot D_{\text{KL}} \left(\mathbb{P}^{\bar{\pi}_{\theta}, \pi_{\eta}^*} \parallel \mathbb{P}^{\bar{\pi}_{\theta}, \pi_{\theta}} \right) + \beta \cdot D_{\text{KL}} \left(\mathbb{P}^{\bar{\pi}_{\theta}, \pi_{\theta}} \parallel \mathbb{P}^{\bar{\pi}_{\theta}, \pi_{\eta}^*} \right).$$

1022 Thus $\beta = 0$ recovers NAIL-F, while $\beta = 1$ recovers NAIL-R. We train each variant for 3M expert-
 1023 query trajectories using the same recipe as in Appendix B.1; for readability, Figure 4 shows only
 1024 the first 1.25M trajectories. All interpolated variants eventually reach perfect accuracy in both the
 1025 low-noise and high-noise regimes, but their learning dynamics differ.

1026 When $\eta = 0$, performance changes smoothly from that of NAIL-F to NAIL-R as β increases,
 1027 matching the behavior in Figure 1: forward-KL-heavy objectives learn fastest. When $\eta = 0.2$,
 1028 however, intermediate and reverse-KL-heavy mixtures learn slightly faster than either endpoint, with
 1029 larger values of β reaching perfect accuracy roughly 100K trajectories earlier. These results suggest
 1030 that the optimal KL mixture is task- and noise-dependent. A more systematic study of when to prefer
 1031 forward KL, reverse KL, or mixtures of the two is an interesting direction for future work.

1032 C Technical Tools

1033 In this section, we recall some technical tools that are used throughout the proofs. We begin in
 1034 Appendix C.1, where we recall some basic definitions and properties of KL divergence and Hellinger
 1035 distance. We proceed in Appendix C.2 by stating some classical results on the performance of
 1036 maximum likelihood estimators. In Appendix C.3, we recall some key results from the theory of

1037 imitation learning that relate regret to Hellinger distance between trajectory distributions. Finally, in
 1038 [Appendix C.4](#) we recall some standard results from online learning that are used in the analysis of
 1039 our online algorithms.

1040 C.1 Information Theory

1041 In this section we recall some basic results from information theory that are used throughout the
 1042 paper. For a more complete introduction to the topic, see Polyanskiy and Wu [43].

1043 We first recall the definitions of KL divergence.

1044 **Definition 6.** Let P, Q be two distributions over the same space \mathcal{X} . The KL divergence between P
 1045 and Q is defined as

$$D_{\text{KL}}(P\|Q) = \mathbb{E}_{X \sim P} \left[\log \frac{dP}{dQ}(X) \right]$$

1046 with $D_{\text{KL}}(P\|Q) = \infty$ if P is not absolutely continuous with respect to Q .

1047 We use the following classical properties of KL divergence repeatedly throughout the paper; see
 1048 Polyanskiy and Wu [43] for details.

1049 **Proposition 4.** Let P, Q be distributions over the same space \mathcal{X} . Then it holds that $D_{\text{KL}}(P\|Q) \geq 0$
 1050 with equality if and only if $P = Q$. Moreover, $(P, Q) \mapsto D_{\text{KL}}(P\|Q)$ is jointly convex in its
 1051 arguments, i.e. for any $\lambda \in [0, 1]$ and any distributions P_1, P_2, Q_1, Q_2 ,

$$D_{\text{KL}}(\lambda P_1 + (1 - \lambda)P_2\|\lambda Q_1 + (1 - \lambda)Q_2) \leq \lambda \cdot D_{\text{KL}}(P_1\|Q_1) + (1 - \lambda) \cdot D_{\text{KL}}(P_2\|Q_2).$$

1052 Furthermore, if $P, Q \in \Delta(\mathcal{X}_1 \times \dots \times \mathcal{X}_n)$, then the KL divergence satisfies a chain rule: if P_{X_1, \dots, X_n}
 1053 and Q_{X_1, \dots, X_n} are the distributions of (X_1, \dots, X_n) under P and Q respectively, then

$$D_{\text{KL}}(P_{X_1, \dots, X_n}\|Q_{X_1, \dots, X_n}) = \sum_{i=1}^n \mathbb{E}_{X_1, \dots, X_{i-1} \sim P} [D_{\text{KL}}(P_{X_i|X_1, \dots, X_{i-1}}\|Q_{X_i|X_1, \dots, X_{i-1}})].$$

1054 While KL divergence is a fundamental notion of distance between distributions, it is infinite when the
 1055 two distributions are not absolutely continuous with respect to each other, which can be problematic
 1056 in many settings in IL. We thus also consider the *Hellinger distance*, defined as follows.

1057 **Definition 7.** Let P, Q be two distributions over the same space \mathcal{X} . The Hellinger distance between
 1058 P and Q is defined as

$$\begin{aligned} D_{\text{H}^2}(P, Q) &= 1 - \mathbb{E}_{X \sim P} \left[\sqrt{\frac{dQ}{dP}}(X) \right] \\ &= 1 - \mathbb{E}_{X \sim Q} \left[\sqrt{\frac{dP}{dQ}}(X) \right] \\ &= \frac{1}{2} \cdot \int_{\mathcal{X}} \left(\sqrt{\frac{dP}{d\mu}}(x) - \sqrt{\frac{dQ}{d\mu}}(x) \right)^2 d\mu(x), \end{aligned}$$

1059 where μ is any measure such that P and Q are absolutely continuous with respect to μ , e.g. $\mu =$
 1060 $(P+Q)/2$.

1061 While Hellinger distance is also nonnegative and equal to zero if and only if the two distributions are
 1062 equal, it is a weaker notion of distance than KL divergence, as it is always bounded by 1. In addition,
 1063 it satisfies a Pinsker-type inequality that relates it to KL divergence.

1064 **Proposition 5** (Pinsker's inequality). Let P, Q be two distributions over the same space \mathcal{X} . Then,

$$D_{\text{H}^2}(P, Q) \leq D_{\text{KL}}(P\|Q).$$

1065 The Hellinger distance is also intimately related to the total variation distance $\text{TV}(P, Q) =$
 1066 $\sup_{A \subseteq \mathcal{X}} |P(A) - Q(A)|$ up to a quadratic factor.

1067 **Proposition 6.** Let P, Q be two distributions over the same space \mathcal{X} . Then,

$$D_{\text{H}^2}(P, Q) \leq \text{TV}(P, Q) \leq \sqrt{2 \cdot D_{\text{H}^2}(P, Q)}.$$

1068 Much like KL divergence, Hellinger is jointly convex in its arguments. In contradistinction to KL
1069 divergence, however, Hellinger distance does not satisfy a chain rule, but it does satisfy two weaker
1070 properties that will be sufficient for our purposes.

1071 **Proposition 7.** Let P, Q be two distributions over $\mathcal{X}_1 \times \cdots \times \mathcal{X}_n$. Let $B_{n+1}(x_{1:n}) = 1$ and for all
1072 $i \leq n$, let for some common dominating measure μ ,

$$B_i(x_{1:i-1}) = \int_{\mathcal{X}_i} \sqrt{\frac{dP_{X_i|X_{1:i-1}=x_{1:i-1}}(x_i)}{d\mu} \cdot \frac{dQ_{X_i|X_{1:i-1}=x_{1:i-1}}(x_i)}{d\mu}} \cdot B_{i+1}(x_{1:i}) d\mu(x_i).$$

1073 Then,

$$D_{\text{H}^2}(P, Q) = 1 - B_1.$$

1074 In particular, if P, Q are product distributions, i.e. $P = P_1 \times \cdots \times P_n$ and $Q = Q_1 \times \cdots \times Q_n$, then

$$D_{\text{H}^2}(P, Q) = 1 - \prod_{i=1}^n (1 - D_{\text{H}^2}(P_i, Q_i)) \leq \sum_{i=1}^n D_{\text{H}^2}(P_i, Q_i).$$

1075 The second property is more similar in form to the chain rule for KL divergence, but it is an inequality
1076 rather than an equality.

1077 **Proposition 8** (Lemma D.2 Foster et al. [18]). Let P, Q be two distributions over $\mathcal{X}_1 \times \cdots \times \mathcal{X}_n$.
1078 Then,

$$D_{\text{H}^2}(P, Q) \leq \sum_{h=1}^H \mathbb{E}_P [D_{\text{H}^2}(P_{X_h|X_{1:h-1}}, Q_{X_h|X_{1:h-1}})].$$

1079 C.2 Maximum Likelihood Estimation

1080 Maximum Likelihood Estimation (MLE) is a fundamental statistical estimation technique that returns
1081 the density in a given class that maximizes the likelihood of the observed data. More precisely if Π'
1082 is a class of conditional distributions $\mathcal{S} \mapsto \Delta(\mathcal{A})$ and $(s^{(i)}, a^{(i)})_{i=1}^n$ are samples, then the MLE $\hat{\pi}$ is
1083 defined as

$$\hat{\pi} \in \operatorname{argmax}_{\pi \in \Pi'} \sum_{i=1}^n \log \pi(a^{(i)} | s^{(i)}) = \operatorname{argmin}_{\pi \in \Pi'} \sum_{i=1}^n -\log \pi(a^{(i)} | s^{(i)}). \quad (6)$$

1084 While the following result is due to Geer [20], Zhang [64], we state the version from Foster et al. [17]
1085 that is most relevant to our setting.

1086 **Theorem 10** (Proposition B.1 from [17]). Let Π' be a finite class of conditional distributions
1087 $\mathcal{S} \mapsto \Delta(\mathcal{A})$ and let $\pi^* \in \Pi'$. Let $(s^{(i)}, a^{(i)})_{i=1}^n$ be i.i.d. samples from \mathbb{P}^{π^*} , where \mathbb{P}^{π^*} is the
1088 distribution over $\mathcal{S} \times \mathcal{A}$ induced by sampling $s \sim \rho$ and $a \sim \pi^*(\cdot | s)$. If $\hat{\pi}$ is the MLE in (6), then
1089 with probability at least $1 - \delta$,

$$D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) \leq 12 \cdot \frac{\log(|\Pi'|/\delta)}{n}.$$

1090 An immediate consequence of [Theorem 10](#) is that Behavior Cloning with the logarithmic loss achieves
1091 Hellinger distance that scales in a horizon-free manner with the number of samples.

1092 **Corollary 1** (Proposition 2.1 from [17]). Let Π be a finite class of policies $\mathcal{S} \times [H] \rightarrow \Delta(\mathcal{A})$ and let
1093 $\pi^* \in \Pi$. Let $\tau^{(i)} = (s_1^{(i)}, a_1^{(i)}, \dots, s_H^{(i)}, a_H^{(i)})$ be i.i.d. trajectories from \mathbb{P}^{π^*} and let

$$\hat{\pi} \in \operatorname{argmin}_{\pi \in \Pi} \sum_{i=1}^n \sum_{h=1}^H -\log \pi(a_h^{(i)} | s_h^{(i)}).$$

1094 Then with probability at least $1 - \delta$,

$$D_{\text{H}^2} \left(\mathbb{P}^{\hat{\pi}}, \mathbb{P}^{\pi^*} \right) \lesssim \frac{\log(|\Pi|/\delta)}{n}.$$

1095 In particular,

$$\mathbb{E} \left[D_{\text{H}^2} \left(\mathbb{P}^{\hat{\pi}}, \mathbb{P}^{\pi^*} \right) \right] \lesssim \frac{\log(n \cdot |\Pi|)}{n}.$$

1096 Indeed, [Corollary 1](#) follows from observing that the policy does not affect the transition densities and
1097 thus

$$\begin{aligned} \hat{\pi} &= \operatorname{argmin}_{\pi \in \Pi} \sum_{i=1}^n \sum_{h=1}^H -\log \left(\pi \left(a_h^{(i)} \mid s_h^{(i)} \right) \right) = \operatorname{argmin}_{\pi \in \Pi} \sum_{i=1}^n \sum_{h=1}^H -\log \left(\pi \left(a_h^{(i)} \mid s_h^{(i)} \right) \cdot P_h \left(s_{h+1}^{(i)} \mid a_h^{(i)}, s_h^{(i)} \right) \right) \\ &= \operatorname{argmin}_{\pi \in \Pi} \sum_{i=1}^n -\log \mathbb{P}^{\pi}(\tau^{(i)}). \end{aligned}$$

1098 The second statement is immediate from the first and the fact that Hellinger distance is bounded by 1.

1099 C.3 Imitation Learning

1100 Recent work has revealed the fundamental importance of the Hellinger distance in the theory of
1101 interactive decision making [[15](#), [16](#), [45](#)]. In this section we recall the key fact that, at least in a
1102 minimax sense, Imitation Learning is essentially equivalent to learning the trajectory distribution
1103 of the expert in Hellinger distance. Following Foster et al. [[17](#)], we consider the deterministic and
1104 stochastic cases separately. In the deterministic case, we have the following result.

1105 **Theorem 11** (Theorem 2.1 from [[17](#)]). *Let π^* be a deterministic policy and let $\hat{\pi}$ be an arbitrary*
1106 *(possibly stochastic) policy. Then,*

$$J(\pi^*) - J(\hat{\pi}) \leq 4R \cdot D_{\text{H}^2} \left(\mathbb{P}^{\hat{\pi}}, \mathbb{P}^{\pi^*} \right),$$

1107 where R is such that

$$0 \leq \sum_{h=1}^H r(s_h, a_h) \leq R$$

1108 for all trajectories $\tau = (s_1, a_1, \dots, s_H, a_H)$.

1109 In the stochastic case, we recall the following result, which is a consequence of the more general
1110 results in Foster et al. [[17](#)].

1111 **Theorem 12** (Theorem 3.1 and Proposition 3.1 from Foster et al. [[17](#)]). *Let π^* be a (possibly*
1112 *stochastic) policy and let $\hat{\pi}$ be any policy. Then,*

$$J(\pi^*) - J(\hat{\pi}) \lesssim \sqrt{R^2 \cdot D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}})} + R \log \left(\frac{R}{D_{\text{H}^2}(\mathbb{P}^{\hat{\pi}}, \mathbb{P}^{\pi^*})} \right) \cdot D_{\text{H}^2}(\mathbb{P}^{\hat{\pi}}, \mathbb{P}^{\pi^*}).$$

1113 Note that Foster et al. [[17](#)] proves a tighter result, replacing the R^2 under the square root with

$$\sigma_{\pi^*}^2 = \sum_{h=1}^H \mathbb{E}^{\pi^*} \left[\left(\mathbb{E} \left[Q_h^{\pi^*}(s_h, a_h) \mid s_h \right] - Q_h^{\pi^*}(s_h, a_h) \right)^2 \right],$$

1114 where Q^{π^*} is the Q -function of π^* (cf. e.g. Sutton et al. [[56](#)]). The second cited result above, Foster
1115 et al. [[17](#), Proposition 3.1], controls $\sigma_{\pi^*}^2$ by R^2 leading to the version stated above. While $\sigma_{\pi^*}^2$ is a
1116 significantly more refined quantity that can be much smaller than R^2 in many settings of interest,
1117 such as when π^* is near-deterministic, we use the version stated above for simplicity. We leave
1118 to future work the interesting problem of understanding the precise role of $\sigma_{\pi^*}^2$ in the noisy expert
1119 setting and whether it can be used to obtain improved guarantees in certain regimes.

1120 Note that Foster et al. [[17](#), Theorem G.3] shows that the above reductions are tight up to constants in
1121 the sense that for any MDP and pair of policies, there exist reward functions that achieve the reverse
1122 inequalities up to constant factors. Thus, in a minimax sense, IL is essentially equivalent to learning
1123 the trajectory distribution of the expert in Hellinger distance. More formally, the result states the
1124 following.

Algorithm 2 Exponential Weights

Require: Number of rounds n , loss function $\ell : \Pi \times \mathcal{Y} \rightarrow \mathbb{R}$, learning rate $\lambda > 0$.

- 1: Set $w_1 = \text{Unif}(\Pi)$.
 - 2: **for** $t = 1$ to n **do**
 - 3: Observe y_t and suffer loss $\mathbb{E}_{\pi \sim w_t}[\ell(\pi, y_t)]$.
 - 4: Update $w_{t+1}(\pi) \propto w_t(\pi) \cdot e^{-\lambda \cdot \ell(\pi, y_t)}$.
 - 5: **end for**
-

1125 **Theorem 13** (Theorem G.3 from Foster et al. [17]). *Let M be an MDP. Then for any pair of policies*
1126 *$\pi^*, \hat{\pi}$ and any $\sigma > 0$ there exists a reward function r such that $\sigma_{\pi^*}^2 \leq \sigma^2$ and*

$$J(\pi^*) - J(\hat{\pi}) \gtrsim \sqrt{\sigma^2 \cdot D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}})}.$$

1127 *Moreover, there exists a reward function r such that*

$$J(\pi^*) - J(\hat{\pi}) \gtrsim R \cdot D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}})$$

1128 *and the same conclusion applies even if we assume π^* is deterministic.*

1129 By [Theorem 13](#), for any of our lower bounds on regret, it suffices to construct instances where
1130 learning the trajectory distribution of the expert in Hellinger distance is hard, which is what we do.

1131 C.4 Online Learning

1132 Several of our results rest on the use of online learning algorithms, and in particular the exponential
1133 weights algorithm, to learn policies in an online fashion. In this section we recall the key definitions
1134 and results from online learning that are used throughout the paper. For a more complete introduction
1135 to the topic, see Cesa-Bianchi and Lugosi [8].

1136 We are only concerned with online learning over finite classes of experts in this work. For a finite class
1137 Π , the online learning problem proceeds in rounds $t \in [n]$ as follows. The learner at the beginning
1138 is informed of a loss function $\ell : \Pi \times \mathcal{Y} \rightarrow \mathbb{R}$ and must choose a distribution $w_t \in \Delta(\Pi)$ over the
1139 experts. Then, an outcome $y_t \in \mathcal{Y}$ is revealed and the learner suffers loss $\mathbb{E}_{\pi \sim w_t}[\ell(\pi, y_t)]$. The goal
1140 of the learner is to minimize regret, defined as

$$\text{Reg}_n = \sum_{t=1}^n \mathbb{E}_{\pi \sim w_t}[\ell(\pi, y_t)] - \min_{\pi \in \Pi} \sum_{t=1}^n \ell(\pi, y_t).$$

1141 The exponential weights algorithm is a simple and classical online learning algorithm that achieves
1142 optimal regret guarantees in the adversarial setting, given in [Algorithm 2](#).

1143 We make use of the following definition.

1144 **Definition 8.** A loss $\ell : \Pi \times \mathcal{Y} \rightarrow \mathbb{R}$ is β -exp-concave if for all $y \in \mathcal{Y}$, the function $\pi \mapsto e^{-\beta \cdot \ell(\pi, y)}$
1145 is concave.

1146 We recall the following result about the regret of the exponential weights algorithm when the loss is
1147 exp-concave.

1148 **Proposition 9** (Proposition 3.1 from Cesa-Bianchi and Lugosi [8]). *If ℓ is β -exp-concave, then the*
1149 *exponential weights algorithm with learning rate $\lambda = \beta$ achieves regret*

$$\text{Reg}_n \leq \frac{\log(|\Pi|)}{\beta}.$$

1150 We emphasize that this regret is independent of the precise choice of y_1, \dots, y_n , and thus holds even
1151 if the outcomes are chosen by an adversary that observes the learner's distribution w_t at each round.
1152 Finally, we will recall a more general definition that is similar to exp-concavity but allows for a better
1153 regret bound.

1154 **Definition 9.** A loss $\ell : \Pi \times \mathcal{Y} \rightarrow \mathbb{R}$ is β -mixable if for all $y \in \mathcal{Y}$ and all distributions $w \in \Delta(\Pi)$,
1155 there exists $\pi_w \in \Pi$ such that

$$e^{-\beta \cdot \ell(\pi_w, y)} \geq \mathbb{E}_{\pi \sim w} \left[e^{-\beta \cdot \ell(\pi, y)} \right].$$

1156 Mixability is a classical notion in online learning and a more complete discussion of it can be found in
 1157 Cesa-Bianchi and Lugosi [8]. We recall the following guarantee for the exponential weights algorithm
 1158 when the loss is mixable.

1159 **Proposition 10** (Proposition 3.2 from Cesa-Bianchi and Lugosi [8]). *If ℓ is β -mixable, then the*
 1160 *exponential weights algorithm with learning rate $\lambda = \beta$ achieves regret*

$$\text{Reg}_n \leq \frac{\log(|\Pi|)}{\beta}.$$

1161 D Proofs from Section 3

1162 In this appendix, we prove the two results in Section 3 involving offline imitation learning with a noisy
 1163 expert. We begin by proving the upper bound that scales exponentially in horizon in Appendix D.1
 1164 before proving that this exponential dependence is necessary in the worst case in Appendix D.2. We
 1165 conclude in Appendix D.3 by showing that any offline IL algorithm must suffer from this exponential
 1166 dependence, making offline IL fundamentally intractable in the noisy expert setting.

1167 D.1 Proof of Theorem 6

1168 We first state a slightly tighter version of the main theorem, which recovers Theorem 6. We first
 1169 conclude the proof under κ -domination and then show how to get a weaker guarantee in the absence
 1170 of κ -domination.

1171 **Theorem 14.** *Let π, π' be policies that κ -dominate the corruption ν . For any $0 \leq \eta < 1$, it holds*
 1172 *that*

$$D_{\text{H}^2} \left(\mathbb{P}^\pi, \mathbb{P}^{\pi'} \right) \leq \frac{2(1 + \eta(2\kappa - 1))}{(1 - \eta)^{H+2}} \cdot D_{\text{H}^2} \left(\mathbb{P}^{\pi_\eta}, \mathbb{P}^{\pi'_\eta} \right).$$

1173 Noting that $\eta < 1$ immediately shows that the first statement of Theorem 6 follows from Theorem 14.
 1174 To prove the latter, we fix policies π, π' and introduce some notation in order to help with the proof.
 1175 First, for any state $s \in \mathcal{S}$, define

$$\mathbb{P}_h^\pi(s) = \mathbb{P}^\pi(a_h, \tau_{h+1:H} | s_h = s)$$

1176 to be the conditional distribution of the action a_h and the future trajectory of states and actions
 1177 conditioned on the event that $s_h = s$. We then define

$$D_h(s) = D_{\text{H}^2} \left(\mathbb{P}_h^\pi(s), \mathbb{P}_h^{\pi'}(s) \right) \quad \text{and} \quad D_h^\eta(s) = D_{\text{H}^2} \left(\mathbb{P}_h^{\pi_\eta}(s), \mathbb{P}_h^{\pi'_\eta}(s) \right).$$

1178 We have the following recursion.

1179 **Lemma 1.** *Let π, π' be any two policies. Then $D_{H+1}(s) = 0$ for all $s \in \mathcal{S}$ and for any $s \in \mathcal{S}$ and*
 1180 *$h \leq H$, it holds that*

$$D_h(s) = D_{\text{H}^2}(\pi_h(\cdot|s), \pi'_h(\cdot|s)) + \sum_{a \in \mathcal{A}} \sqrt{\pi_h(a|s) \cdot \pi'_h(a|s)} \cdot \int D_{h+1}(s') dP_h(s'|s, a).$$

1181 *Proof.* We use Proposition 7. Indeed, we compute

$$\begin{aligned} 1 - D_h(s) &= \mathbb{E}^\pi \left[\sqrt{\frac{\pi'_h(a_h|s)}{\pi_h(a_h|s)}} (1 - D_{h+1}(s_{h+1})) \right] \\ &= \sum_{a \in \mathcal{A}} \sqrt{\pi_h(a|s) \cdot \pi'_h(a|s)} \cdot \mathbb{E}_{s_{h+1} \sim P_h(\cdot|s, a)} [1 - D_{h+1}(s_{h+1})]. \end{aligned}$$

1182 The result follows by rearranging and using the definition of the Hellinger squared distance. \square

1183 Lemma 1 reduces the problem of bounding the contraction in Hellinger distance of trajectory
 1184 distributions to bounding the contraction on a per-step basis. We thus prove that in the case of
 1185 κ -dominated corruption distributions ν , we can control this contraction at each time step h .

1186 **Lemma 2.** Let π, π' be arbitrary policies and suppose that ν is κ -dominated, i.e., for all $1 \leq h \leq H$,
 1187 all $s \in \mathcal{S}$ and $a \in \text{supp}(\pi_h(\cdot|s)) \cup \text{supp}(\pi'_h(\cdot|s))$, it holds that

$$\nu_h(a|s) \leq \kappa (\pi_h(a|s) + \pi'_h(a|s)).$$

1188 Then for any $0 \leq \eta < 1$, any state $s \in \mathcal{S}$ and time $1 \leq h \leq H$, it holds that

$$D_{\text{H}^2}(\pi_h(\cdot|s), \pi'_h(\cdot|s)) \leq \frac{2(1 + \eta(2\kappa - 1))}{(1 - \eta)^2} \cdot D_{\text{H}^2}(\pi_{\eta,h}(\cdot|s), \pi'_{\eta,h}(\cdot|s)).$$

1189 *Proof.* For the sake of notational simplicity, we will suppress the h in this proof; in addition,
 1190 because we have fixed a state s , we will write $\pi(a)$ for $\pi_h(a|s)$. We observe that for any $a \in$
 1191 $\text{supp}(\pi) \cup \text{supp}(\pi')$, it holds that

$$\begin{aligned} \pi_\eta(a) + \pi'_\eta(a) &= (1 - \eta)(\pi(a) + \pi'(a)) + 2\eta \cdot \nu(a) \\ &\leq (1 - \eta)(\pi(a) + \pi'(a)) + 2\eta\kappa \cdot (\pi(a) + \pi'(a)) \\ &= (1 + \eta(2\kappa - 1))(\pi(a) + \pi'(a)). \end{aligned}$$

1192 Thus,

$$\left(\sqrt{\pi_\eta(a)} + \sqrt{\pi'_\eta(a)} \right)^2 \leq 2(\pi_\eta(a) + \pi'_\eta(a)) \leq 2(1 + \eta(2\kappa - 1)) \left(\sqrt{\pi(a)} + \sqrt{\pi'(a)} \right)^2,$$

1193 where we used the fact that $a + b \leq (\sqrt{a} + \sqrt{b})^2 \leq 2(a + b)$ for nonnegative a, b . Plugging into
 1194 the definition of Hellinger distance, we see that

$$\begin{aligned} D_{\text{H}^2}(\pi_\eta, \pi'_\eta) &= \frac{1}{2} \sum_{a \in \mathcal{A}} \left(\sqrt{\pi_\eta(a)} - \sqrt{\pi'_\eta(a)} \right)^2 \\ &= \frac{1}{2} \sum_{a \in \mathcal{A}} \frac{(\pi_\eta(a) - \pi'_\eta(a))^2}{\left(\sqrt{\pi_\eta(a)} + \sqrt{\pi'_\eta(a)} \right)^2} \\ &\geq \frac{(1 - \eta)^2}{2(1 + \eta(2\kappa - 1))} \cdot \frac{1}{2} \sum_{a \in \mathcal{A}} \frac{(\pi(a) - \pi'(a))^2}{\left(\sqrt{\pi(a)} + \sqrt{\pi'(a)} \right)^2} \\ &= \frac{(1 - \eta)^2}{2(1 + \eta(2\kappa - 1))} \cdot D_{\text{H}^2}(\pi, \pi'). \end{aligned}$$

1195 The result follows. \square

1196 **Remark 1.** Note that it is precisely in [Lemma 2](#) that κ -domination is used and in the proof one
 1197 can understand the necessity of such an assumption. Indeed, the problem is precisely that the map
 1198 $\eta \mapsto \sqrt{\eta}$ is not Lipschitz near $\eta = 0$. For a simple example of what can go wrong, let

$$\pi = \text{Bernoulli}(0) = \delta_0, \quad \pi' = \text{Bernoulli}(\varepsilon), \quad \text{and} \quad \nu = \text{Bernoulli}(1/2).$$

1199 An elementary computation then shows that

$$D_{\text{H}^2}(\pi, \pi') \asymp \varepsilon \quad \text{but} \quad D_{\text{H}^2}(\pi_\eta, \pi'_\eta) \asymp \varepsilon^2,$$

1200 for $\eta > 0$. Thus, in order to prevent such a quadratic blowup, we need to ensure that ν does not
 1201 put too much mass on actions that receive small, but positive, probability under π or π' , which is
 1202 precisely what κ -domination ensures.

1203 **Remark 2.** Note that even with κ -domination, the $(1 - \eta)^2$ dependence in the comparison is real.
 1204 Indeed, suppose that $\pi^* = \delta_{a_1}$ and $\hat{\pi} = \delta_{a_2}$ for some $a_1 \neq a_2$, and $\nu = 1/2(\pi^* + \hat{\pi})$. Then we have
 1205 κ -domination with $\kappa = 1$. On the other hand, we have $D_{\text{H}^2}(\pi^*, \hat{\pi}) = 1$, whereas

$$D_{\text{H}^2}(\pi_\eta^*, \hat{\pi}_\eta) = 1 - \sqrt{\eta(2 - \eta)} \asymp \frac{(1 - \eta)^2}{2}.$$

1206 Thus the $(1 - \eta)^2$ dependence is tight up to constant factors.

1207 We are now ready to conclude the proof of the result.

1208 *Proof of Theorem 14.* We prove the following claim by reverse induction from $H, \dots, 1$: it holds for
 1209 any $s \in \mathcal{S}$ that

$$D_h^\eta(s) \geq (1 - \eta)^{H+1-h} \cdot \frac{(1 - \eta)^2}{2(1 + \eta(2\kappa - 1))} \cdot D_h(s). \quad (7)$$

1210 The case $h = H$ follows immediately from Lemma 2. We thus suppose that (7) holds for $h + 1$. By
 1211 Lemma 1, it holds that

$$\begin{aligned} D_h^\eta(s) &= D_{\text{H}^2}(\pi_{\eta,h}(\cdot|s), \pi'_{\eta,h}(\cdot|s)) + \sum_{a \in \mathcal{A}} \sqrt{\pi_{\eta,h}(a|s) \cdot \pi'_{\eta,h}(a|s)} \cdot \int D_{h+1}^\eta(s') dP_h(s'|s, a) \\ &\geq \frac{(1 - \eta)^2}{2(1 + \eta(2\kappa - 1))} \cdot D_{\text{H}^2}(\pi_h(\cdot|s), \pi'_h(\cdot|s)) + \sum_{a \in \mathcal{A}} \sqrt{\pi_{\eta,h}(a|s) \cdot \pi'_{\eta,h}(a|s)} \cdot \int D_{h+1}^\eta(s') dP_h(s'|s, a) \\ &\geq \frac{(1 - \eta)^2}{2(1 + \eta(2\kappa - 1))} \cdot D_{\text{H}^2}(\pi_h(\cdot|s), \pi'_h(\cdot|s)) + (1 - \eta) \cdot \sum_{a \in \mathcal{A}} \sqrt{\pi_h(a|s) \cdot \pi'_h(a|s)} \cdot \int D_{h+1}^\eta(s') dP_h(s'|s, a) \\ &\geq \frac{(1 - \eta)^2}{2(1 + \eta(2\kappa - 1))} \cdot D_{\text{H}^2}(\pi_h(\cdot|s), \pi'_h(\cdot|s)) \\ &\quad + (1 - \eta) \cdot (1 - \eta)^{H-h} \cdot \frac{(1 - \eta)^2}{2(1 + \eta(2\kappa - 1))} \cdot \sum_{a \in \mathcal{A}} \sqrt{\pi_h(a|s) \cdot \pi'_h(a|s)} \cdot \int D_{h+1}^\eta(s') dP_h(s'|s, a) \\ &\geq \frac{(1 - \eta)^{H+3-h}}{2(1 + \eta(2\kappa - 1))} \cdot \left(D_{\text{H}^2}(\pi_h(\cdot|s), \pi'_h(\cdot|s)) + \sum_{a \in \mathcal{A}} \sqrt{\pi_h(a|s) \cdot \pi'_h(a|s)} \cdot \int D_{h+1}^\eta(s') dP_h(s'|s, a) \right) \\ &= \frac{(1 - \eta)^{H+3-h}}{2(1 + \eta(2\kappa - 1))} \cdot D_h(s), \end{aligned}$$

1212 where the first inequality used Lemma 2, the second inequality used the fact that

$$\begin{aligned} \sqrt{\pi_{\eta,h}(a|s) \cdot \pi'_{\eta,h}(a|s)} &= \sqrt{((1 - \eta)\pi_h(a|s) + \eta\nu_h(a|s))((1 - \eta)\pi'_h(a|s) + \eta\nu_h(a|s))} \\ &\geq (1 - \eta) \cdot \sqrt{\pi_h(a|s) \cdot \pi'_h(a|s)}, \end{aligned}$$

1213 the third inequality used the inductive hypothesis, and the final inequality again used Lemma 1. The
 1214 result follows immediately. \square

1215 Note that in the case that π, π' are both *deterministic* policies, then any ν is κ -dominated with $\kappa = 1$
 1216 and thus

$$D_{\text{H}^2}(\mathbb{P}^\pi, \mathbb{P}^{\pi'}) \lesssim (1 - \eta)^{-H-2} \cdot D_{\text{H}^2}(\mathbb{P}^{\pi_\eta}, \mathbb{P}^{\pi'_\eta}) \leq \frac{e^{\eta H/1-\eta}}{(1 - \eta)^2} \cdot D_{\text{H}^2}(\mathbb{P}^{\pi_\eta}, \mathbb{P}^{\pi'_\eta}).$$

1217 We now show how to get a weaker guarantee in the absence of κ -domination. We restate the result
 1218 with constants made explicit now.

1219 **Proposition 11.** *Let M be a horizon H MDP and let π, π' be two policies. For any choice of ν and
 1220 any $0 < \eta < 1$, it holds that*

$$D_{\text{H}^2}(\mathbb{P}^\pi, \mathbb{P}^{\pi'}) \leq \sqrt{2 \cdot \frac{(1 - \eta)^{-H} - 1}{\eta(1 - \eta)}} \cdot D_{\text{H}^2}(\mathbb{P}^{\pi_\eta}, \mathbb{P}^{\pi'_\eta}).$$

1221 We first prove the following general result on the contraction of Hellinger distance under arbitrary
 1222 corruptions.

1223 **Lemma 3.** *Let P, Q be two distributions over a common space \mathcal{X} and let ν be an arbitrary distribution
 1224 over \mathcal{X} . For any $0 < \eta < 1$, letting $P_\eta = (1 - \eta)P + \eta\nu$ and $Q_\eta = (1 - \eta)Q + \eta\nu$, it holds that*

$$D_{\text{H}^2}(P, Q) \leq \frac{\sqrt{2}}{1 - \eta} \cdot \sqrt{D_{\text{H}^2}(P_\eta, Q_\eta)}.$$

1225 *Proof.* By [Proposition 6](#), it holds that

$$\begin{aligned}
\sqrt{2D_{\mathbb{H}^2}(P_\eta, Q_\eta)} &\geq \text{TV}(P_\eta, Q_\eta) \\
&= \sup_{0 \leq f \leq 1} |\mathbb{E}_{P_\eta}[f] - \mathbb{E}_{Q_\eta}[f]| \\
&= (1 - \eta) \cdot \sup_{0 \leq f \leq 1} |\mathbb{E}_P[f] - \mathbb{E}_Q[f]| \\
&= (1 - \eta) \cdot \text{TV}(P, Q) \\
&\geq (1 - \eta) \cdot D_{\mathbb{H}^2}(P, Q),
\end{aligned}$$

1226 where the second equality comes from the linearity of expectation and the final inequality again
1227 comes from [Proposition 6](#). \square

1228 We can now prove the proposition.

1229 *Proof of [Proposition 11](#).* We will apply backward induction and [Lemma 1](#). Indeed, we use the
1230 identical notation as that used in the proof thereof. We suppose that there are constants C_h for
1231 $h = H, \dots, 1$ satisfying $C_H = 2^{-1} \cdot (1 - \eta)^2$ and

$$\frac{1}{C_h} = \frac{2}{(1 - \eta)^2} + \frac{1}{(1 - \eta) \cdot C_{h+1}}$$

1232 such that

$$D_h^\eta(s) \geq C_h \cdot D_h(s)^2 \quad \text{for all } s \in \mathcal{S}.$$

1233 By [Lemma 3](#), the statement holds for $h = H$. Now, by [Lemma 1](#), it holds that

$$\begin{aligned}
D_h^\eta(s) &= D_{\mathbb{H}^2}(\pi_{\eta,h}(\cdot|s), \pi'_{\eta,h}(\cdot|s)) + \sum_{a \in \mathcal{A}} \sqrt{\pi_{\eta,h}(a|s) \cdot \pi'_{\eta,h}(a|s)} \cdot \int D_{h+1}^\eta(s') dP_h(s'|s, a) \\
&\geq 2^{-1} \cdot (1 - \eta)^2 \cdot D_{\mathbb{H}^2}(\pi_h(\cdot|s), \pi'_h(\cdot|s))^2 \\
&\quad + (1 - \eta) \cdot \sum_{a \in \mathcal{A}} \sqrt{\pi_h(a|s) \cdot \pi'_h(a|s)} \cdot \int D_{h+1}^\eta(s') dP_h(s'|s, a) \\
&\geq 2^{-1} \cdot (1 - \eta)^2 \cdot D_{\mathbb{H}^2}(\pi_h(\cdot|s), \pi'_h(\cdot|s))^2 \\
&\quad + (1 - \eta) \cdot \sum_{a \in \mathcal{A}} \sqrt{\pi_h(a|s) \cdot \pi'_h(a|s)} \cdot \int C_{h+1} D_{h+1}(s')^2 dP_h(s'|s, a) \\
&\geq 2^{-1} \cdot (1 - \eta)^2 \cdot D_{\mathbb{H}^2}(\pi_h(\cdot|s), \pi'_h(\cdot|s))^2 \\
&\quad + C_{h+1}(1 - \eta) \cdot \left(\sum_{a \in \mathcal{A}} \sqrt{\pi_h(a|s) \cdot \pi'_h(a|s)} \cdot \int D_{h+1}(s') dP_h(s'|s, a) \right)^2,
\end{aligned}$$

1234 where the second inequality follows from the inductive hypothesis and the final inequality follows
1235 from Jensen's inequality and the fact that $\sum_a \sqrt{\pi_h(a|s) \cdot \pi'_h(a|s)} \leq 1$. By AM-GM, it holds that for
1236 $a, b, x, y \geq 0$,

$$a \cdot x^2 + b \cdot y^2 \geq \frac{1}{\frac{1}{a} + \frac{1}{b}} \cdot (x + y)^2.$$

1237 Thus, applying this to the final expression in the above display, we see that

$$\begin{aligned}
D_h^\eta(s) &\geq C_h \cdot \left(D_{\mathbb{H}^2}(\pi_h(\cdot|s), \pi'_h(\cdot|s)) + \sum_{a \in \mathcal{A}} \sqrt{\pi_h(a|s) \cdot \pi'_h(a|s)} \cdot \int D_{h+1}(s') dP_h(s'|s, a) \right)^2 \\
&= C_h \cdot D_h(s)^2,
\end{aligned}$$

1238 where the equality follows from [Lemma 1](#). Thus, the inductive step is complete.

1239 It remains to bound C_h itself. Letting $c_h = C_h^{-1}$, we see that

$$c_H = 2(1 - \eta)^{-2}, \quad c_h = 2(1 - \eta)^{-2} + (1 - \eta)^{-1} \cdot c_{h+1}$$

1240 for $h < H$. Thus by induction, it holds that

$$c_1 = 2 \cdot \frac{(1 - \eta)^{-H} - 1}{\eta(1 - \eta)}.$$

1241 The result follows. \square

1242 Finally, we note that [Theorem 6](#) follows immediately from [Theorem 14](#) and [Proposition 11](#), concluding
1243 the proof.

1244 D.2 Proof of Proposition 1

1245 We will prove three lower bounds: one for arbitrary κ , one for $\kappa = 1$ with a deterministic expert, and
1246 one that holds absent κ -domination. We begin by stating the common construction for the first two,
1247 before proving these first two results separately.

1248 We will take an MDP that has $H + 1$ states s_1, \dots, s_H and an absorbing state \perp . We will suppose
1249 the action space $\mathcal{A} = \{a_1, a_2, a_3\}$. Let the transition functions for $h < H$ be

$$P_h(s'|s, a) = \begin{cases} \delta_{s_{h+1}} & a = a_1 \text{ and } s = s_h \\ \delta_{\perp} & a \in \{a_2, a_3\} \text{ or } s = \perp \end{cases}.$$

1250 In other words, \perp is an absorbing state that is reached by taking a ‘wrong’ action at any state,
1251 otherwise we deterministically transition to the next state.

1252 We now state the formal proposition for the case of $\kappa = 1$ and deterministic experts.

1253 **Proposition 12.** *For any $H \geq 1$ and $0 \leq \eta < 1$, there exists an MDP, a policy class Π of size 2, and
1254 a corruption distribution ν such that both policies in Π are deterministic, but*

$$D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) = 1 \quad \text{but} \quad D_{\text{H}^2}(\mathbb{P}^{\pi_\eta^*}, \mathbb{P}^{\hat{\pi}_\eta}) \leq (1 - \eta)^{H+1}.$$

1255 *Proof.* We will consider the MDP described above with policy class $\Pi = \{\pi^*, \hat{\pi}\}$ of size 2 such
1256 that $\hat{\pi}(\cdot|s_h) = \pi^*(\cdot|s_h) = \delta_{a_1}$ (an atom on a_1) for $h < H$, $\pi^*(s_H) = \delta_{a_1}$, $\hat{\pi}(s_H) = \delta_{a_2}$, and both
1257 policies take the same action at \perp . Finally, suppose that

$$\nu(\cdot|s_H) = \frac{\delta_{a_1} + \delta_{a_2}}{2} \quad \text{and} \quad \nu(\cdot|s_h) = \delta_{a_3} \text{ for } h < H.$$

1258 We claim that for any $0 \leq \eta < 1$, the result of the proposition holds for this construction. Note that
1259 the $\eta = 0$ case is trivial, so suppose that $\eta > 0$. We now observe that $D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) = 1$ because
1260 the policies deterministically differ in the final time step, which is reached with probability 1. Thus
1261 we focus on upper bounding $D_{\text{H}^2}(\mathbb{P}^{\pi_\eta^*}, \mathbb{P}^{\hat{\pi}_\eta})$.

1262 Note first that for $\pi \in \{\pi^*, \hat{\pi}\}$, with probability at least $1 - (1 - \eta)^{H-1}$, it holds that the trajectory
1263 under π will transition to \perp before the final time step; moreover, conditional on this event, both
1264 corrupted policies induce the same distribution. On the complementary event of no contamination up
1265 to step H , it holds that

$$\pi_\eta^* = \left(1 - \frac{\eta}{2}\right) \cdot \delta_{a_1} + \frac{\eta}{2} \cdot \delta_{a_2} \quad \text{and} \quad \hat{\pi}_\eta = \frac{\eta}{2} \cdot \delta_{a_1} + \left(1 - \frac{\eta}{2}\right) \cdot \delta_{a_2}.$$

1266 Thus,

$$D_{\text{H}^2}(\pi_\eta^*(\cdot|s_H), \hat{\pi}_\eta(\cdot|s_H)) = 1 - \sqrt{\eta(2 - \eta)} = \frac{(1 - \eta)^2}{1 + \sqrt{\eta(2 - \eta)}} \leq (1 - \eta)^2.$$

1267 Combining these observations, we see that

$$D_{\text{H}^2}(\mathbb{P}^{\pi_\eta^*}, \mathbb{P}^{\hat{\pi}_\eta}) \leq (1 - \eta)^{H-1} \cdot (1 - \eta)^2 = (1 - \eta)^{H+1}.$$

1268 The result follows. \square

1269 We now state the formal proposition for the case of arbitrary κ .

1270 **Proposition 13.** *For any $H \geq 1$, $\kappa \geq 1$, and $0 \leq \eta < 1$, there exists an MDP, a policy class Π of*
 1271 *size 2, and a corruption distribution ν such that both policies in Π κ -dominate ν , but*

$$D_{\text{H}^2}(\mathbb{P}^\pi, \mathbb{P}^{\pi'}) \geq 4\eta \cdot \kappa \cdot (1 - \eta)^{-H-1} \cdot D_{\text{H}^2}(\mathbb{P}^{\pi_\eta}, \mathbb{P}^{\pi'_\eta}).$$

1272 *Proof.* We will consider the MDP described above but now let for $h < H$,

$$\pi_h(\cdot | s_h) = \pi'_h(\cdot | s_h) = (1 - \varepsilon) \cdot \delta_{a_1} + \varepsilon \cdot \delta_{a_3}$$

1273 and

$$\pi_H(\cdot | s_H) = (1 - \varepsilon) \cdot \delta_{a_1} + \varepsilon \cdot \delta_{a_2} \quad \text{and} \quad \pi'_H(\cdot | s_H) = (1 - \varepsilon) \cdot \delta_{a_1} + \varepsilon \cdot \delta_{a_3}.$$

1274 Let both policies take action a_2 on \perp . Finally, let

$$\nu_h(\cdot | s_h) = \delta_{a_3} \text{ for } h < H \quad \text{and} \quad \nu_H(\cdot | s_H) = \frac{\delta_{a_2} + \delta_{a_3}}{2},$$

1275 and let $\nu(\cdot | \perp) = \delta_{a_2}$. Suppose that $\varepsilon = 1/2\kappa$. We first note that ν is κ -dominated. Indeed, for $h < H$
 1276 this is immediate for action a_1 , and for a_3 ,

$$\nu_h(a_3 | s_h) = 1 = \kappa \cdot 2\varepsilon = \kappa \cdot (\pi_h(a_3 | s_h) + \pi'_h(a_3 | s_h)).$$

1277 At s_H , for a_2 we have

$$\nu_H(a_2 | s_H) = \frac{1}{2} = \kappa \cdot \varepsilon = \kappa \cdot (\pi_H(a_2 | s_H) + \pi'_H(a_2 | s_H)),$$

1278 and similarly for a_3 . Thus ν is κ -dominated.

1279 We now compute the Hellinger distance between the *clean* trajectory distributions. Let \mathcal{E} denote the
 1280 event that the trajectory reaches s_H . Note that the two laws are identical on \mathcal{E}^c . On the other hand,

$$\mathbb{P}^\pi(\mathcal{E}) = \mathbb{P}^{\pi'}(\mathcal{E}) = (1 - \varepsilon)^{H-1} \quad \text{and} \quad D_{\text{H}^2}(\pi_H(\cdot | s_H), \pi'_H(\cdot | s_H)) = \varepsilon$$

1281 Thus,

$$D_{\text{H}^2}(\mathbb{P}^\pi, \mathbb{P}^{\pi'}) = (1 - \varepsilon)^{H-1} \cdot \varepsilon.$$

1282 We now compute the Hellinger distance between the corrupted trajectory distributions. For $h < H$
 1283 the policies coincide and the probability of reaching s_H is $(1 - \eta)^{H-1} \cdot (1 - \varepsilon)^{H-1}$. At the final
 1284 time step,

$$D_{\text{H}^2}(\pi_{\eta, H}(\cdot | s_H), \pi'_{\eta, H}(\cdot | s_H)) = \left(\sqrt{\eta/2 + (1 - \eta)\varepsilon} - \sqrt{\eta/2} \right)^2 \leq \frac{(1 - \eta)^2 \varepsilon^2}{2\eta}.$$

1285 Thus,

$$\begin{aligned} D_{\text{H}^2}(\mathbb{P}^{\pi_\eta}, \mathbb{P}^{\pi'_\eta}) &\leq ((1 - \eta)(1 - \varepsilon))^{H-1} \cdot \frac{(1 - \eta)^2 \cdot \varepsilon^2}{2\eta} \\ &\leq \frac{(1 - \eta)^{H+1} \varepsilon}{2\eta} \cdot D_{\text{H}^2}(\mathbb{P}^\pi, \mathbb{P}^{\pi'}) \\ &= \frac{(1 - \eta)^{H+1}}{4\eta \cdot \kappa} \cdot D_{\text{H}^2}(\mathbb{P}^\pi, \mathbb{P}^{\pi'}). \end{aligned}$$

1286 The result follows. □

1287 We now provide a lower bound that holds absent κ -domination.

1288 **Proposition 14.** *For any $H \geq 1$ and $0 < \eta < 1$, there exists a horizon H MDP with 3 actions,*
 1289 *policies π^* , $\hat{\pi}$, and a corruption ν such that*

$$D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) \gtrsim \sqrt{\frac{(1 - \eta)^{-H} - 1}{1 - \eta}} \cdot D_{\text{H}^2}(\mathbb{P}^{\pi_\eta^*}, \mathbb{P}^{\hat{\pi}_\eta}).$$

1290 *Proof.* Let M have $H + 1$ states s_1, \dots, s_H and \perp , where

$$P_h(\cdot | s_h, a_2) = \delta_{s_{h+1}} \quad \text{and} \quad P_h(\cdot | s_h, a_1) = P_h(\cdot | s_h, a_3) = P_h(\cdot | \perp, a) = \delta_{\perp}.$$

1291 Let $t < 1/4$, let

$$S = \sum_{h=0}^{H-1} (1-\eta)^{-h} = \frac{(1-\eta)((1-\eta)^{-H} - 1)}{\eta} \quad \text{and} \quad \lambda = t/S.$$

1292 Let

$$u_h = \lambda(1-\eta)^{-h+1}.$$

1293 Note that $u_h \leq 1/4$ for every h . Now, let

$$\pi_h^*(\cdot | s_h) = u_h \cdot \delta_{a_1} + (1-u_h) \cdot \delta_{a_2} \quad \text{and} \quad \hat{\pi}_h(\cdot | s_h) = u_h \cdot \delta_{a_3} + (1-u_h) \cdot \delta_{a_2};$$

1294 let $\pi^*, \hat{\pi}$ agree on \perp and define

$$\nu_h(\cdot | s_h) = \frac{1}{2} \cdot \delta_{a_1} + \frac{1}{2} \cdot \delta_{a_3} \quad \text{and} \quad \nu_h(\cdot | \perp) = \delta_{a_1}.$$

1295 We will let

$$D_h = D_{\text{H}^2} \left(\mathbb{P}_h^{\pi^*}(s_h), \mathbb{P}_h^{\hat{\pi}}(s_h) \right) \quad \text{and} \quad D_h^\eta = D_{\text{H}^2} \left(\mathbb{P}_h^{\pi_\eta^*}(s_h), \mathbb{P}_h^{\hat{\pi}_\eta}(s_h) \right).$$

1296 Applying [Lemma 1](#), we see that

$$D_h = u_h + (1-u_h) \cdot D_{h+1}.$$

1297 Thus, by induction, we have that

$$1 - e^{-t} \leq D_1 = 1 - \prod_{h=1}^H (1-u_h) \leq t.$$

1298 Because $t \leq 1/4$, it thus holds that $t/2 \leq D_1 \leq t$. On the other hand, a direct computation shows that

$$D_{\text{H}^2} \left(\pi_{\eta,h}^*(\cdot | s_h), \hat{\pi}_{\eta,h}(\cdot | s_h) \right) = 1 - (1-\eta)(1-u_h) - \sqrt{\eta^2 + 2\eta(1-\eta)u_h} \leq \frac{(1-\eta)^2}{2\eta} \cdot u_h^2.$$

1299 By [Lemma 1](#) again we see that

$$\begin{aligned} D_h^\eta &= 1 - (1-\eta)(1-u_h) - \sqrt{\eta^2 + 2\eta(1-\eta)u_h} + (1-\eta)(1-u_h) \cdot D_{h+1}^\eta \\ &\leq \frac{(1-\eta)^2}{2\eta} \cdot u_h^2 + (1-\eta)(1-u_h) \cdot D_{h+1}^\eta. \end{aligned}$$

1300 Thus, by induction, we have that

$$D_1^\eta \leq \frac{(1-\eta)^2}{2\eta} \cdot \sum_{h=1}^H (1-\eta)^{h-1} \cdot u_h^2 = \frac{(1-\eta)^2}{2\eta} \cdot \lambda^2 S.$$

1301 Substituting in the definition of λ and S , we have that

$$D_1^\eta \leq \frac{1-\eta}{2((1-\eta)^{-H} - 1)} \cdot t^2 \leq \frac{2(1-\eta)}{(1-\eta)^{-H} - 1} \cdot D_1^2.$$

1302 The result follows immediately. □

1303 We now prove the main lower bound.

1304 *Proof of [Proposition 1](#).* This follows from combining [Propositions 12](#) to [14](#). □

1305 **D.3 Implications for Offline Imitation Learning**

1306 We conclude this appendix by demonstrating that the exponential dependence on horizon in [Theorem 6](#)
 1307 must appear in any offline IL algorithm through a similar construction to the one used in [Appendix D.2](#).

1308 **Proposition 15.** *Let $\kappa \geq 1$ and $0 \leq \eta < 1$. For any $H \geq 2$ and*

$$\varepsilon < \frac{(1 - 1/2\kappa)^{H-1}}{128\kappa},$$

1309 *there exists an MDP, a policy class Π of size 2, and a corruption distribution ν such that both policies*
 1310 *in Π κ -dominate ν , but any offline IL algorithm achieving $D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) \leq \varepsilon$ requires*

$$n \gtrsim \frac{\eta \cdot \kappa}{(1 - \eta)^{H+1} \cdot \varepsilon}$$

1311 *samples from the corrupted expert.*

1312 *Proof.* We use the identical construction as the proof in [Proposition 13](#), except we suppose that

$$\pi_H(\cdot | s_H) = (1 - u) \cdot \delta_{a_2} + u \cdot \delta_{a_1}, \quad \pi'_H(\cdot | s_H) = (1 - u) \cdot \delta_{a_2} + u \cdot \delta_{a_3},$$

1313 and

$$\nu_H(\cdot | s_H) = \kappa \cdot u \cdot \delta_{a_1} + (1 - 2\kappa \cdot u) \cdot \delta_{a_2} + \kappa \cdot u \cdot \delta_{a_3},$$

1314 for some $u \leq 1/2\kappa$. This construction is clearly κ -dominated. Moreover,

$$D_{\text{H}^2}(\mathbb{P}^\pi, \mathbb{P}^{\pi'}) = (1 - 1/2\kappa)^{H-1} \cdot u \quad \text{and} \quad D_{\text{H}^2}(\mathbb{P}^{\pi_\eta}, \mathbb{P}^{\pi'_\eta}) \leq (1 - \eta)^{H-1} \cdot (1 - 1/2\kappa)^{H-1} \cdot u \cdot \frac{(1 - \eta)^2}{4\eta\kappa}.$$

1315 Choosing

$$u = \frac{64 \cdot \varepsilon}{(1 - 1/2\kappa)^{H-1}},$$

1316 which allows $u \leq 1/2\kappa$ by the assumption on ε , we see that

$$D_{\text{H}^2}(\mathbb{P}^\pi, \mathbb{P}^{\pi'}) = 64 \cdot \varepsilon \quad \text{and} \quad D_{\text{H}^2}(\mathbb{P}^{\pi_\eta}, \mathbb{P}^{\pi'_\eta}) \leq (1 - \eta)^{H+1} \cdot \frac{16 \cdot \varepsilon}{\eta \cdot \kappa}.$$

1317 A standard two point argument concludes the proof. □

1318 **E Proofs from Section 4**

1319 In this appendix, we prove the results related to online Imitation Learning stated in the main body.
 1320 We begin in [Appendix E.1](#) with the proof of [Theorem 7](#), before continuing to prove the upper bound
 1321 on NAIL in [Appendix E.3](#), and concluding with a proof of the lower bound in [Appendix E.4](#).

1322 **E.1 Proof of Theorem 7**

1323 We first state a slightly tighter version of the result, albeit with a more complicated expression in the
 1324 upper bound. We break this into two separate results, one using forward KL and one using reverse
 1325 KL, which together imply [Theorem 7](#). First, under κ -domination, we have the following result.

1326 **Theorem 15.** *Let π, π' be any two policies that κ -dominate the corruption ν . Then, for any $0 \leq \eta < 1$*
 1327 *it holds that*

$$D_{\text{H}^2}(\mathbb{P}^\pi, \mathbb{P}^{\pi'}) \leq \frac{2(1 + \eta(2\kappa - 1))}{(1 - \eta)^2} \cdot \left(D_{\text{KL}}(\mathbb{P}^{\pi', \pi_\eta} \parallel \mathbb{P}^{\pi', \pi'_\eta}) \wedge D_{\text{KL}}(\mathbb{P}^{\pi', \pi'_\eta} \parallel \mathbb{P}^{\pi', \pi_\eta}) \right).$$

1328 Note that the first conclusion of [Theorem 7](#) can be recovered immediately.

1329 *Proof of Theorem 15.* We make use of the subadditivity of Hellinger squared divergence ([Proposi-](#)
1330 [tion 8](#)). Indeed, it holds that

$$D_{\text{H}^2} \left(\mathbb{P}^\pi, \mathbb{P}^{\pi'} \right) \leq \mathbb{E}^{\pi'} \left[\sum_{h=1}^H D_{\text{H}^2} \left(\pi(\cdot|s_h), \pi'(\cdot|s_h) \right) \right].$$

1331 Now, applying [Lemma 2](#), we have that for any s_h ,

$$D_{\text{H}^2} \left(\pi(\cdot|s_h), \pi'(\cdot|s_h) \right) \leq \frac{2(1 + \eta(2\kappa - 1))}{(1 - \eta)^2} \cdot D_{\text{H}^2} \left(\pi_\eta(\cdot|s_h), \pi'_\eta(\cdot|s_h) \right).$$

1332 Combining this fact with the preceding display and applying [Proposition 5](#), we have

$$\begin{aligned} D_{\text{H}^2} \left(\mathbb{P}^\pi, \mathbb{P}^{\pi'} \right) &\leq \frac{2(1 + \eta(2\kappa - 1))}{(1 - \eta)^2} \cdot \mathbb{E}^{\pi'} \left[\sum_{h=1}^H D_{\text{H}^2} \left(\pi_{\eta,h}(\cdot|s_h), \pi'_{\eta,h}(\cdot|s_h) \right) \right] \\ &\leq \frac{2(1 + \eta(2\kappa - 1))}{(1 - \eta)^2} \cdot \mathbb{E}^{\pi'} \left[\sum_{h=1}^H D_{\text{KL}} \left(\pi_{\eta,h}(\cdot|s_h) \parallel \pi'_{\eta,h}(\cdot|s_h) \right) \right]. \end{aligned} \quad (8)$$

1333 We may now apply the Chain rule for KL divergence ([Proposition 4](#)) to observe that

$$\mathbb{E}^{\pi'} \left[\sum_{h=1}^H D_{\text{KL}} \left(\pi_{\eta,h}(\cdot|s_h) \parallel \pi'_{\eta,h}(\cdot|s_h) \right) \right] = D_{\text{KL}} \left(\mathbb{P}^{\pi', \pi_\eta} \parallel \mathbb{P}^{\pi', \pi'_\eta} \right).$$

1334 The first result follows. The second follows by the same argument by observing that the Hellinger
1335 distance is symmetric in (8). \square

1336 We now provide a similar comparison between the clean and noisy trajectory distributions absent
1337 κ -domination.

1338 **Proposition 16.** *Let $\pi^*, \hat{\pi}$ be two policies and let ν be an arbitrary corruption distribution. Then for*
1339 *any $0 \leq \eta < 1$ it holds that*

$$D_{\text{H}^2} \left(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}} \right) \leq \frac{1}{1 - \eta} \cdot \sqrt{2H \cdot D_{\text{KL}} \left(\mathbb{P}^{\hat{\pi}, \pi_\eta^*} \parallel \mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta} \right) \wedge D_{\text{KL}} \left(\mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta} \parallel \mathbb{P}^{\hat{\pi}, \pi_\eta^*} \right)}.$$

1340 *Proof.* As in the proof of [Theorem 7](#), we use the subadditivity of the Hellinger distance ([Proposition 8](#))
1341 to get that

$$\begin{aligned} D_{\text{H}^2} \left(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}} \right) &\leq \mathbb{E}^{\hat{\pi}} \left[\sum_{h=1}^H D_{\text{H}^2} \left(\pi_h^*(\cdot|s_h), \hat{\pi}_h(\cdot|s_h) \right) \right] \\ &\leq \frac{\sqrt{2}}{1 - \eta} \cdot \mathbb{E}^{\hat{\pi}} \left[\sum_{h=1}^H \sqrt{D_{\text{H}^2} \left(\pi_{\eta,h}^*(\cdot|s_h), \hat{\pi}_{\eta,h}(\cdot|s_h) \right)} \right] \\ &\leq \frac{\sqrt{2H}}{1 - \eta} \cdot \sqrt{\mathbb{E}^{\hat{\pi}} \left[\sum_{h=1}^H D_{\text{H}^2} \left(\pi_{\eta,h}^*(\cdot|s_h), \hat{\pi}_{\eta,h}(\cdot|s_h) \right) \right]} \\ &\leq \frac{\sqrt{2H}}{1 - \eta} \cdot \sqrt{D_{\text{KL}} \left(\mathbb{P}^{\hat{\pi}, \pi_\eta^*} \parallel \mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta} \right)}, \end{aligned}$$

1342 where the second inequality comes from [Lemma 3](#), the third inequality comes from Cauchy-Schwarz,
1343 and the final inequality comes from applying Pinsker's inequality ([Proposition 5](#)) and the chain rule
1344 for KL divergence ([Proposition 4](#)). The first claimed bound follows. The second comes from the
1345 identical argument but using the symmetry of the Hellinger distance to apply Pinsker's inequality in
1346 the reverse direction. The result follows. \square

1347 Finally, we note that [Theorem 7](#) follows immediately from [Theorem 15](#) and [Proposition 16](#), concluding
1348 the proof of the main result.

1349 **E.2 Necessity of Horizon Dependence in the Absence of κ -Domination**

1350 We now show that the polynomial dependence in horizon that appears in [Proposition 16](#) is necessary
 1351 with the following lower bound. In particular, [Proposition 16](#) is tight up to a polynomial dependence
 1352 on $(1 - \eta)$.

1353 **Proposition 17.** *For any $H \geq 1$ and $0 < \eta < 1$, there exists a horizon H MDP with 2 actions,
 1354 policies π^* , $\hat{\pi}$, and corruption distribution ν such that*

$$D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) \geq 8^{-1} \cdot \sqrt{\frac{H \cdot \eta}{(1 - \eta)}} \cdot D_{\text{KL}}(\mathbb{P}^{\hat{\pi}, \pi_\eta^*} \parallel \mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta}).$$

1355 *Proof.* Let the MDP M have two actions, $\mathcal{A} = \{0, 1\}$ and H steps such that s_h transitions to s_{h+1}
 1356 deterministically, independent of the action. Let

$$\pi_h^*(0|s_h) = 1, \quad \hat{\pi}_h(0|s_h) = 1 - \varepsilon, \quad \text{and} \quad \nu_h(0|s_h) = 0$$

1357 for every h . Thus actions are independent across time. We thus compute

$$D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) = 1 - (1 - \varepsilon)^{H/2}$$

1358 and thus, if $\varepsilon \leq 1/4H$, it holds that

$$\frac{H\varepsilon}{4} \leq D_{\text{H}^2}(\mathbb{P}^{\pi^*}, \mathbb{P}^{\hat{\pi}}) \leq H\varepsilon.$$

1359 On the other hand, it holds that

$$\begin{aligned} D_{\text{KL}}(\mathbb{P}^{\hat{\pi}, \pi_\eta^*} \parallel \mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta}) &= H \cdot D_{\text{KL}}(\pi_\eta^*(\cdot|s) \parallel \hat{\pi}_\eta(\cdot|s)) \\ &= H \cdot \left((1 - \eta) \cdot \log\left(\frac{1}{1 - \varepsilon}\right) + \eta \cdot \log\left(\frac{\eta}{\eta + (1 - \eta)\varepsilon}\right) \right). \end{aligned}$$

1360 Using the fact that $\log(1 - x) \geq -x - x^2$ and $\log(1 + x) \geq x - x^2/2$ for $x \in (0, 1)$, we get that if
 1361 $\varepsilon \leq \eta/e(1 - \eta)$, then

$$D_{\text{KL}}(\mathbb{P}^{\hat{\pi}, \pi_\eta^*} \parallel \mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta}) \leq \frac{3}{2} \cdot H \cdot \frac{1 - \eta}{\eta} \cdot \varepsilon^2.$$

1362 Thus for sufficiently small ε , the claimed relation holds. □

1363 **E.3 Proof of Theorem 8**

1364 The main content of the proof is to show the following result, which is that the on-policy KL
 1365 divergence $\mathbb{E} \left[D_{\text{KL}}(\mathbb{P}^{\hat{\pi}, \pi_\eta^*} \parallel \mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta}) \right]$ for $\hat{\pi}$ returned by [Algorithm 1](#) is small. The result will then
 1366 follow from [Theorem 7](#). Indeed, we have the following result.

1367 **Lemma 4.** *Let Π be a finite policy class, ν be an arbitrary corruption distribution, and $0 \leq \eta < 1$
 1368 be a corruption level. Let $\hat{\pi}$ denote the policy returned by [Algorithm 1](#). Then it holds that*

$$\mathbb{E} \left[D_{\text{KL}}(\mathbb{P}^{\hat{\pi}, \pi_\eta^*} \parallel \mathbb{P}^{\hat{\pi}, \hat{\pi}_\eta}) \right] \leq \frac{H \cdot \log(|\Pi|)}{n}.$$

1369 *Proof.* Let $\mu_t = \sum_{\pi \in \Pi} w_t(\pi) \cdot \pi$ denote the mixture policy defined in [Algorithm 1](#). We first claim
 1370 that for any sequence of trajectories $\tau^{(t)}, \tau^{(t)'}$, it holds that

$$\sum_{t=1}^n \sum_{h=1}^H \log \left(\frac{\pi_{\eta, h}^*(a_h^{(t)' | s_h^{(t)})}}{\mu_{\eta, h}^{(t)}(a_h^{(t)' | s_h^{(t)})}} \right) \leq H \cdot \log(|\Pi|). \quad (9)$$

1371 To see this, let

$$\ell(\pi, \tau') = - \sum_{h=1}^H \log(\pi_{\eta, h}(a_h' | s_h))$$

1372 denote a loss function. We claim that ℓ is $1/H$ -mixable (Definition 9) and that Algorithm 1 amounts
 1373 to running the exponential weights algorithm (Algorithm 2) with respect to ℓ over the policy class

$$\Pi_\eta = \{\pi_\eta | \pi \in \Pi\}.$$

1374 Indeed, fix τ' and let $w \in \Delta(\Pi)$ and

$$\mu_w = \sum_{\pi \in \Pi} w(\pi) \cdot \pi.$$

1375 Then, we note that $(\mu_w)_\eta = \sum_{\pi \in \Pi} w(\pi) \cdot \pi_\eta$ by linearity and thus

$$\begin{aligned} e^{-H^{-1} \cdot \ell(\mu_w, \tau')} &= \left(\prod_{h=1}^H (\mu_w)_{\eta, h}(a'_h | s_h) \right)^{1/H} \\ &= \prod_{h=1}^H \left(\sum_{\pi \in \Pi} w(\pi) \cdot \pi_{\eta, h}(a'_h | s_h) \right)^{1/H} \\ &\geq \sum_{\pi \in \Pi} w(\pi) \cdot \left(\prod_{h=1}^H \pi_{\eta, h}(a'_h | s_h) \right)^{1/H} \\ &= \sum_{\pi \in \Pi} w(\pi) \cdot e^{-H^{-1} \cdot \ell(\pi, \tau')} \\ &= \mathbb{E}_{\pi \sim w} \left[e^{-H^{-1} \cdot \ell(\pi, \tau')} \right], \end{aligned}$$

1376 where the inequality follows from Hölder's inequality. Thus ℓ is $1/H$ -mixable. Moreover, it is
 1377 immediate that the update in Algorithm 1 corresponds to the exponential weights update with respect
 1378 to ℓ over the policy class Π_η with learning rate $\lambda = H^{-1}$. Thus it holds by Proposition 10 that (9)
 1379 holds for any sequence $\tau^{(t)}$.

1380 Now, for fixed $1 \leq t \leq n$ and temporarily suppressing the notational dependence on t , note that

$$\begin{aligned} \mathbb{E} \left[\sum_{h=1}^H \log \left(\frac{\pi_{\eta, h}^*(a'_h | s_h)}{\mu_{\eta, h}(a'_h | s_h)} \right) \right] &= \mathbb{E}^\mu \mathbb{E}_{a'_h \sim \pi_{\eta, h}^*(\cdot | s_h)} \left[\sum_{h=1}^H \log \left(\frac{\pi_{\eta, h}^*(a'_h | s_h) \cdot P_h(s_{h+1} | s_h, a_h)}{\mu_{\eta, h}(a'_h | s_h) \cdot P_h(s_{h+1} | s_h, a_h)} \right) \right] \\ &= \text{D}_{\text{KL}} \left(\mathbb{P}^{\mu, \pi_\eta^*} \parallel \mathbb{P}^{\mu, \mu_\eta} \right), \end{aligned} \quad (10)$$

1381 where the last equality follows from the chain rule for KL divergence (Proposition 4). Combining (9)
 1382 and (10) and renormalizing, we have that

$$\mathbb{E} \left[\frac{1}{n} \sum_{t=1}^n \text{D}_{\text{KL}} \left(\mathbb{P}^{\mu^{(t)}, \pi_\eta^*} \parallel \mathbb{P}^{\mu^{(t)}, \mu_\eta^{(t)}} \right) \right] \leq \frac{H \cdot \log(|\Pi|)}{n}.$$

1383 Since $\hat{\pi} = \mu_T$ for $T \sim \text{Unif}([n])$, independent of the training randomness, the left-hand side is
 1384 exactly the average over t . The desired bound follows. \square

1385 We can now prove the main result.

1386 *Proof of Theorem 8.* The result follows immediately by combining Lemma 4 and Theorem 7. \square

1387 E.4 Proof of Proposition 2

1388 In this section we prove the lower bound Proposition 2 demonstrating that even with online access,
 1389 noisy, stochastic experts necessitate the linear in horizon dependence in sample complexity that is
 1390 present in Theorem 8. We first state a more formal version of the lower bound, which is stated in the
 1391 main body as Proposition 2.

1392 **Proposition 18.** For any $H \geq 2$, corruption level $0 < \eta < 1$, and $\varepsilon < 1/32$, there exists a horizon H
1393 MDP with three actions, deterministic transitions, a known corruption distribution ν , and a policy
1394 class Π of size $|\Pi| = 2$ such that ν is κ -dominated by Π with $\kappa \asymp H/\varepsilon$ and any online IL algorithm
1395 must observe at least

$$n \geq \frac{\eta \cdot H}{512(1-\eta)^2 \cdot \varepsilon^2}$$

1396 trajectories of interaction with the noisy expert π_η^* in order to achieve regret $J(\pi^*) - J(\hat{\pi}) \leq \varepsilon$.

1397 *Proof.* By Foster et al. [17], it suffices to show a lower bound against learning in trajectory-wise
1398 Hellinger distance to order ε . Let an MDP have states s_1, \dots, s_H with transitions $s_h \rightarrow s_{h+1}$
1399 deterministically independent of the action. Let the action space be $\mathcal{A} = \{a_1, a_2, a_3\}$ and the policy
1400 space be $\Pi = \{\pi_+, \pi_-\}$. For some u to be determined such that $Hu \leq 1/2$ and $u \leq 1/4$, let

$$\pi_{+,h}(\cdot|s_h) = u \cdot \delta_{a_1} + (1-u) \cdot \delta_{a_2}, \quad \pi_{-,h}(\cdot|s_h) = u \cdot \delta_{a_3} + (1-u) \cdot \delta_{a_2}, \quad \text{and} \quad \nu_h = \frac{1}{4} \cdot (\delta_{a_1} + \delta_{a_3}) + \frac{1}{2} \cdot \delta_{a_2}.$$

1401 Note that ν_h is κ -dominated by Π with $\kappa = 1/4u$.

1402 Now observe that

$$D_{\text{H}^2}(\pi_{+,h}(\cdot|s_h), \pi_{-,h}(\cdot|s_h)) = 1 - \sqrt{(1-u)^2} = u.$$

1403 Moreover, because the action distributions are independent,

$$D_{\text{H}^2}(\mathbb{P}^{\pi_+}, \mathbb{P}^{\pi_-}) = 1 - (1-u)^H \geq \frac{Hu}{2},$$

1404 where the inequality used the assumption that $2Hu \leq 1$. On the other hand, an elementary computa-
1405 tion reveals that

$$\begin{aligned} \pi_{+, \eta, h} &= \left(\frac{\eta}{4} + (1-\eta)u\right) \cdot \delta_{a_1} + \left(\frac{\eta}{2} + (1-\eta)(1-u)\right) \cdot \delta_{a_2} + \frac{\eta}{4} \cdot \delta_{a_3} \\ \pi_{-, \eta, h} &= \frac{\eta}{4} \cdot \delta_{a_1} + \left(\frac{\eta}{2} + (1-\eta)(1-u)\right) \cdot \delta_{a_2} + \left(\frac{\eta}{4} + (1-\eta)u\right) \cdot \delta_{a_3}. \end{aligned}$$

1406 Thus,

$$\begin{aligned} D_{\text{KL}}(\pi_{+, \eta, h} \| \pi_{-, \eta, h}) &= \left(\frac{\eta}{4} + (1-\eta)u\right) \cdot \log\left(1 + \frac{4(1-\eta)u}{\eta}\right) + (1-\eta)u \cdot \log\left(\frac{\eta}{4(1-\eta)u + \eta}\right) \\ &= (1-\eta)u \cdot \log\left(1 + \frac{4(1-\eta)u}{\eta}\right) \\ &\leq \frac{4(1-\eta)^2 u^2}{\eta}. \end{aligned}$$

1407 Thus, after n rounds of interaction, the chain rule for KL divergence ([Proposition 4](#)) ensures that the
1408 KL divergence between the distributions over trajectories induced by π_+ and π_- is at most

$$nH \cdot \frac{4(1-\eta)^2 u^2}{\eta}$$

1409 and thus if

$$n \leq \frac{\eta}{8H(1-\eta)^2 u^2},$$

1410 then by Le Cam's inequality no algorithm can identify π^* with probability greater than $3/4$, which
1411 is required in order to achieve regret at most ε if $Hu/2 \geq 4\varepsilon$. Thus, setting $u = 8\varepsilon/H$, we see that as
1412 long as $\varepsilon \leq 1/32$, then $Hu \leq 1/2$ and $u \leq 1/4$. Plugging in concludes the proof. \square

1413 We now show the necessity of κ -domination in order to shave off the quadratic dependence on ε in
1414 the sample complexity. We have the following result.

1415 **Proposition 19.** For any $\kappa \geq 1$, any corruption level $0 < \eta < 1$, and any $\varepsilon < 2^{-4} \cdot \kappa^{-1}$, there exists
1416 a horizon $H = 1$ MDP with three actions, a policy class Π of size $|\Pi| = 2$, and a known corruption
1417 distribution ν such that:

Algorithm 3 NAILGUN: Noise-robust Aggregation for Imitation Learning with Greedy UNcertainty

Require: Number of rounds n , deterministic policy class Π , noisy expert π_η^* , noise ceiling $\alpha < 1$, contamination ceiling $\rho > 0$.

- 1: Set $r = \sqrt{(1-\alpha)/\alpha\rho}$ and $w_1 = \text{Unif}(\Pi)$.
 - 2: **for** $t = 1$ to n **do**
 - 3: Define $\mu_t = \sum_{\pi \in \Pi} w_t(\pi) \cdot \pi$ and $\bar{\mu}_t(\cdot|s) = \text{argmax}_a \mu_t(a|s)$.
 - 4: Deploy $\bar{\mu}_t$ to get trajectory $\tau^{(t)} \sim \mathbb{P}^{\bar{\mu}_t}$.
 - 5: Query noisy expert π_η^* on $\tau^{(t)}$ to obtain augmented trajectory $\tau'^{(t)}$.
 - 6: Update $w_{t+1}(\pi) \propto w_t(\pi) \cdot r^{\sum_{h=1}^H \mathbb{I}\{\pi(s_h^{(t)})=a_h^{(t)'}\}}$.
 - 7: **end for**
 - 8: **return** $\hat{\pi} = \bar{\mu}_T$ for $T \sim \text{Unif}([n])$.
-

1418 (a) ν is κ -dominated by Π and $\pi^* \in \Pi$, and

1419 (b) in order for an algorithm to achieve regret $J(\pi^*) - J(\hat{\pi}) \leq \varepsilon$, the learner must observe at
 1420 least $n \gtrsim \eta \cdot \kappa / \varepsilon (1-\eta)^2$ trajectories of interaction with the noisy expert π_η^* .

1421 *Proof.* Suppose there is a single state s and three actions a_1, a_2, a_3 . Fix some $0 < u \leq 1/2\kappa$ and
 1422 suppose that $\Pi = \{\pi_+, \pi_-\}$ where

$$\pi_+ = u \cdot \delta_{a_1} + (1-u) \cdot \delta_{a_2}, \quad \pi_- = u \cdot \delta_{a_3} + (1-u) \cdot \delta_{a_2}, \quad \text{and} \quad \nu = \kappa u \cdot \delta_{a_1} + \kappa u \cdot \delta_{a_3} + (1-2\kappa u) \cdot \delta_{a_2}.$$

1423 It is immediate that ν is κ -dominated by Π from the construction. Moreover, observe that

$$D_{\text{H}^2}(\mathbb{P}^{\pi_+}, \mathbb{P}^{\pi_-}) = 1 - \sqrt{(1-u)^2} = u.$$

1424 Thus, if $u > 4\varepsilon$, then any algorithm returning a policy $\hat{\pi}$ such that $D_{\text{H}^2}(\mathbb{P}^{\hat{\pi}}, \mathbb{P}^{\pi^*}) \leq \varepsilon$ must identify
 1425 π^* . We now observe that an elementary computation reveals that

$$D_{\text{KL}}(\mathbb{P}^{\pi_+, \eta} \| \mathbb{P}^{\pi_-, \eta}) = u(1-\eta) \cdot \log\left(1 + \frac{1-\eta}{\eta\kappa}\right) \leq u \cdot \frac{(1-\eta)^2}{\eta\kappa}.$$

1426 Thus, if we were to set $u = 8\varepsilon$, then we would have by Le Cam's inequality that if

$$n \cdot \frac{8\varepsilon \cdot (1-\eta)^2}{\eta\kappa} \leq \frac{1}{8}$$

1427 then with constant probability, $\hat{\pi} \neq \pi^*$. If $\varepsilon \leq 1/16\kappa$, then we can find such a u and thus the stated
 1428 sample complexity is required to achieve regret at most ε . The result follows. \square

1429 F Proofs from Section 5

1430 In this appendix, we provide the proofs of the results stated in Section 5. We begin by describing the
 1431 algorithm NAILGUN and providing intuition for its design in Appendix F.2. We then provide the
 1432 proof of Theorem 9 in Appendix F.3. Finally, we provide the proof of the lower bound Proposition 3
 1433 in Appendix F.4.

1434 F.1 The Problem of Identifiability

1435 In the unknown corruption setting, identifiability can be a major concern, i.e., it may be the case
 1436 that there are multiple policies π that are consistent with the observed noisy expert π_η^* and thus it is
 1437 impossible to identify the clean expert π^* . We now provide a simple example of this phenomenon.

1438 Let M be a horizon $H = 1$ MDP with a single state and action space \mathcal{A} . Let $\pi^* = \delta_{a^*}$ be a
 1439 deterministic expert and let ν be an arbitrary policy. For $\eta \leq \alpha$, let

$$\pi = \frac{1-\alpha}{1-\eta} \cdot \delta_{a^*} + \frac{\alpha-\eta}{1-\eta} \cdot \nu.$$

1440 Then it holds that $\pi_\alpha^* = \pi_\eta$, so the noisy expert π_α^* is consistent with both π^* and π as the underlying
 1441 expert policy, and thus it is impossible to distinguish between π^* and π . In particular, if $J(\pi^*) \gg$
 1442 $J(\pi)$, then it is impossible to learn a policy with good performance without additional assumptions
 1443 on the corruption or the feedback.

1444 **F.2 Description of the Algorithm**

1445 In [Theorem 8](#), we saw that we could get regret bounds that scale polynomially in the horizon H
 1446 by appealing to exponential weights and in particular the mixability of the trajectory-level log-loss.
 1447 Unfortunately, this strategy fundamentally relies on the learner having access to the loss at each round
 1448 of interaction, which depends on knowing both η and ν . In the unknown corruption setting, we do
 1449 not have access to this information and thus we cannot directly apply this strategy. Instead we use a
 1450 surrogate loss function specifically adapted to the assumption that the clean expert π^* is *deterministic*.
 1451 In particular, we use the following loss function:

$$\ell(\pi, \tau') = \log(r) \cdot \sum_{h=1}^H \mathbb{I}\{\pi(s_h) \neq a'_h\}. \quad (11)$$

1452 While this loss function itself can grow linearly in horizon, we demonstrate that its induced behavior
 1453 on the probability of error is substantially nicer. We eliminate this concern by rolling out our estimated
 1454 policies *greedily*. This leads naturally to [Algorithm 3](#), which we call **NAILGUN**. At each step, we
 1455 roll out our current policy greedily, collect noisy expert feedback, and then run [Algorithm 2](#) with
 1456 the loss function defined in (11). The key step in the analysis, provided below, is the observation
 1457 that the loss in (11) has the property that when we make a mistake at any point along the trajectory,
 1458 the weight we place on a policy that is not the expert goes down by a constant factor in expectation,
 1459 where critically this constant factor is independent of horizon. This analysis is thus substantially
 1460 different from the standard analysis of exponential weights, which relies on the mixability of the loss
 1461 function to control the sum of losses across rounds, and is more similar to the analysis of the *halving*
 1462 *algorithm* for online learning with expert advice [39].

1463 **F.3 Proof of Theorem 9**

1464 We begin by restating the theorem with a slightly tighter dependence on the parameters in question.

1465 **Theorem 16.** *Let Π be a class of deterministic policies and suppose that $\pi^* \in \Pi$. Let ν denote a*
 1466 *corruption distribution such that for all $h \in [H]$, $a \in \mathcal{A}$, and $s \in \mathcal{S}$ it holds that $\nu_h(a|s) \leq \rho$ for*
 1467 *some $0 < \rho < 1$. Let $\alpha > 0$ such that $(1 + \rho)\alpha < 1$ and suppose that the corruption noise in the*
 1468 *noisy expert η satisfies $\eta \leq \alpha$. If $\hat{\pi}$ is the policy returned by [Algorithm 3](#), then it holds that*

$$\mathbb{E} \left[\text{D}_{\text{H}^2} \left(\mathbb{P}^{\hat{\pi}}, \mathbb{P}^{\pi^*} \right) \right] \leq \frac{4 \cdot \log(|\Pi|)}{n (\sqrt{1 - \alpha} - \sqrt{\alpha\rho})^2}.$$

1469 We now derive [Theorem 9](#) directly from this result.

1470 *Proof of Theorem 9.* We compute

$$(\sqrt{1 - \alpha} - \sqrt{\alpha\rho})^2 = \frac{(1 - \alpha(1 + \rho))^2}{(\sqrt{1 - \alpha} + \sqrt{\alpha\rho})^2} \geq \frac{(1 - \alpha(1 + \rho))^2}{2}.$$

1471 Plugging this into [Theorem 16](#) concludes the proof. \square

1472 In order to prove [Theorem 16](#), we must introduce some notation. First, for any $1 \leq t \leq n$, let
 1473 \mathcal{F}_t denote the σ -algebra generated by all randomness before time t . Thus w_t, W_t, μ_t , and $\bar{\mu}_t$ are
 1474 \mathcal{F}_{t-1} -measurable. We recall that

$$r = \sqrt{\frac{1 - \alpha}{\alpha\rho}} \geq 1, \quad \mu_t = \sum_{\pi \in \Pi} w_t(\pi) \cdot \pi, \quad \text{and} \quad w_{t+1}(\pi) \propto w_t(\pi) r^{\sum_{h=1}^t \mathbb{I}\{\pi(s_h^{(t)}) \neq a_h^{(t)'}\}}.$$

1475 Moreover, $\bar{\mu}_t$ denotes the greedy policy associated with μ_t . We introduce the following notation:

$$R_t(\pi) = \frac{w_t(\pi)}{w_t(\pi^*)} \quad \text{and} \quad W_t = \frac{1}{w_t(\pi^*)} = 1 + \sum_{\pi \neq \pi^*} R_t(\pi). \quad (12)$$

1476 Note that if W_t is small, then μ_t places a lot of weight on the correct policy π^* and, moreover, if W_t
 1477 is very small, then we expect $\bar{\mu}_t$ to agree with π^* on a trajectory. We let

$$M_t = \left\{ \text{there exists } h \leq H \text{ such that } \bar{\mu}_{t,h}(s_h^{(t)}) \neq \pi_h^*(s_h^{(t)}) \right\} \quad (13)$$

1478 denote the event that the rolled out policy $\bar{\mu}_t$ disagrees with the *clean* expert π^* on at least one action.
 1479 The key lemma that we will show is that, in expectation, any time a mistake is made, W_t goes down
 1480 by a constant factor.

1481 **Lemma 5.** *Let W_t be as in (12) and let M_t be as in (13). Then it holds for any $1 \leq t \leq n$ that*

$$\mathbb{E}[W_{t+1}|\mathcal{F}_{t-1}] \leq W_t \left(1 - \frac{(\sqrt{1-\alpha} - \sqrt{\alpha\rho})^2}{2} \cdot \mathbb{E}[M_t|\mathcal{F}_{t-1}] \right).$$

1482 We will now prove [Theorem 16](#) assuming [Lemma 5](#) and return to the proof of this key result below.

1483 *Proof of [Theorem 16](#).* Let M_t be as in (13). Applying the fact that $\hat{\pi} = \mu_T$ for $T \sim \text{Unif}([n])$ as
 1484 well as Foster et al. [[17](#), Lemma D.3], we see that

$$\mathbb{E} \left[D_{\text{H}^2} \left(\mathbb{P}^{\hat{\pi}}, \mathbb{P}^{\pi^*} \right) \right] \leq \frac{1}{n} \sum_{t=1}^n \mathbb{E} \left[D_{\text{H}^2} \left(\mathbb{P}^{\bar{\mu}_t}, \mathbb{P}^{\pi^*} \right) \right] \leq \frac{2}{n} \sum_{t=1}^n \mathbb{P}(M_t).$$

1485 We now apply Jensen's inequality and [Lemma 5](#) to observe that

$$\begin{aligned} \mathbb{E}[\log(W_{t+1})|\mathcal{F}_{t-1}] &\leq \log(\mathbb{E}[W_{t+1}|\mathcal{F}_{t-1}]) \\ &\leq \log \left(W_t \left(1 - \frac{(\sqrt{1-\alpha} - \sqrt{\alpha\rho})^2}{2} \cdot \mathbb{E}[M_t|\mathcal{F}_{t-1}] \right) \right) \\ &= \log(W_t) + \log \left(1 - \frac{(\sqrt{1-\alpha} - \sqrt{\alpha\rho})^2}{2} \cdot \mathbb{P}(M_t|\mathcal{F}_{t-1}) \right) \\ &\leq \log(W_t) - \frac{(\sqrt{1-\alpha} - \sqrt{\alpha\rho})^2}{2} \cdot \mathbb{P}(M_t|\mathcal{F}_{t-1}). \end{aligned}$$

1486 The last inequality follows due to the fact that

$$0 \leq \frac{(\sqrt{1-\alpha} - \sqrt{\alpha\rho})^2}{2} \mathbb{P}(M_t|\mathcal{F}_{t-1}) \leq 1$$

1487 almost surely and the numerical inequality $\log(1-x) \leq -x$ for $0 < x < 1$. Rearranging, taking
 1488 expectations, and applying the tower property of conditional expectation, we see that

$$\begin{aligned} \frac{(\sqrt{1-\alpha} - \sqrt{\alpha\rho})^2}{2} \cdot \sum_{t=1}^n \mathbb{P}(M_t) &\leq \sum_{t=1}^n \mathbb{E}[\log(W_t) - \log(W_{t+1})] \\ &\leq \log(W_1) \\ &= \log(|\Pi|), \end{aligned}$$

1489 where the second inequality follows because $W_t \geq 1$ by construction in (12) and the equality follows
 1490 by assuming a uniform prior w_1 . Rearranging concludes the proof. \square

1491 Thus it remains to prove the key intermediate result, [Lemma 5](#).

1492 *Proof of [Lemma 5](#).* We fix a time t and omit the t from the notation of the trajectories τ and τ' in
 1493 order to simplify the presentation. Let \mathcal{G}_t denote the sigma-algebra generated by \mathcal{F}_{t-1} and the states
 1494 $s_1^{(t)}, \dots, s_H^{(t)}$, but not the noisy labels $a_1^{(t)'}, \dots, a_H^{(t)'}$. That is, $\mathcal{G}_t = \mathcal{F}_{t-1} \vee \sigma(s_{1:H}^{(t)})$. Observe that
 1495 for $\pi \neq \pi^*$, letting $a_h^* = \pi_h^*(s_h)$ and $\bar{a}_h = \pi_h(s_h)$, we have

$$R_{t+1}(\pi) = R_t(\pi) \cdot \exp \left(\log(r) \cdot \sum_{h=1}^H (\mathbb{I}\{\bar{a}_h = a_h'\} - \mathbb{I}\{a_h^* = a_h'\}) \right).$$

1496 In the event that $\bar{a}_h = a_h^*$, clearly $\mathbb{I}\{\bar{a}_h = a_h'\} - \mathbb{I}\{a_h^* = a_h'\} = 0$. On the other hand, if $\bar{a}_h \neq a_h^*$,
 1497 then by the margin assumption,

$$\mathbb{P}(a_h' = a_h^* | s_h) \geq 1 - \alpha \quad \text{and} \quad \mathbb{P}(a_h' = \bar{a}_h | s_h) \leq \alpha\rho.$$

1498 Thus,

$$\begin{aligned} \mathbb{E} \left[r^{\mathbb{I}\{\pi_h(s_h)=a'_h\}} - \mathbb{I}\{\pi_h^*(s_h)=a'_h\} \mid s_h \right] &= 1 - \mathbb{P}(a'_h = a_h^* \mid s_h) - \mathbb{P}(a'_h = \bar{a}_h \mid s_h) \\ &\quad + \frac{\mathbb{P}(a'_h = a_h^* \mid s_h)}{r} + r \cdot \mathbb{P}(a'_h = \bar{a}_h \mid s_h) \\ &= 1 - \left(1 - \frac{1}{r}\right) \cdot \mathbb{P}(a'_h = a_h^* \mid s_h) + (r-1) \cdot \mathbb{P}(a'_h = \bar{a}_h \mid s_h). \end{aligned}$$

1499 By the assumptions on α, ρ , it holds that $r > 1$ and thus combining the preceding two displays, we
1500 have

$$\begin{aligned} \mathbb{E} \left[r^{\mathbb{I}\{\pi_h(s_h)=a'_h\}} - \mathbb{I}\{\pi_h^*(s_h)=a'_h\} \mid s_h \right] &\leq 1 - \left(1 - \frac{1}{r}\right) \cdot (1 - \alpha) + (r-1) \cdot \alpha\rho \\ &= 1 - \left(\sqrt{1-\alpha} - \sqrt{\alpha\rho}\right)^2, \end{aligned}$$

1501 where we used the fact that $r = \sqrt{(1-\alpha)/\alpha\rho}$ in the last equality.

1502 Letting

$$N_t(\pi) = \sum_{h=1}^H \mathbb{I}\left\{\pi_h(s_h^{(t)}) \neq \pi_h^*(s_h^{(t)})\right\},$$

1503 we thus conclude that

$$\begin{aligned} \mathbb{E} [R_{t+1}(\pi) \mid \mathcal{G}_t] &\leq R_t(\pi) \cdot \left(1 - \left(\sqrt{1-\alpha} - \sqrt{\alpha\rho}\right)^2\right)^{N_t(\pi)} \\ &\leq R_t(\pi) \left(1 - \left(\sqrt{1-\alpha} - \sqrt{\alpha\rho}\right)^2 \cdot \mathbb{I}\{N_t(\pi) \geq 1\}\right). \end{aligned}$$

1504 Summing over $\pi \neq \pi^*$, we have that

$$\begin{aligned} \mathbb{E} [W_{t+1} - 1 \mid \mathcal{G}_t] &= \mathbb{E} \left[\sum_{\pi \neq \pi^*} R_{t+1}(\pi) \mid \mathcal{G}_t \right] \\ &\leq \sum_{\pi \neq \pi^*} R_t(\pi) \left(1 - \left(\sqrt{1-\alpha} - \sqrt{\alpha\rho}\right)^2 \cdot \mathbb{I}\{N_t(\pi) \geq 1\}\right) \\ &= W_t - 1 - \left(\sqrt{1-\alpha} - \sqrt{\alpha\rho}\right)^2 \cdot \sum_{N_t(\pi) \geq 1} R_t(\pi). \end{aligned} \tag{14}$$

1505 Now let $a_h = \bar{\mu}_{t,h}(s_h)$ be the rolled out action and suppose that M_t occurs; then there is some
1506 minimal $h \leq H$ such that $a_h \neq a_h^*$. Because $\bar{\mu}_t$ is the greedy policy, it must then hold that

$$\sum_{\substack{\pi \in \Pi \\ \pi(s_h)=a_h}} w_t(\pi) \geq \sum_{\substack{\pi \in \Pi \\ \pi(s_h)=a_h^*}} w_t(\pi)$$

1507 and thus $\mu_{t,h}(a_h^* \mid s_h) \leq 1/2$ and, in particular,

$$\sum_{N_t(\pi) \geq 1} w_t(\pi) \geq \frac{1}{2}.$$

1508 Dividing by $w_t(\pi^*)$, we see that

$$\sum_{N_t(\pi) \geq 1} R_t(\pi) \geq \frac{1}{2 \cdot w_t(\pi^*)} = \frac{W_t}{2}.$$

1509 By (14), it clearly holds that W_{t+1} is a supermartingale. Moreover, combining that equation with the
1510 preceding display, we see that the result holds. \square

1511 **F.4 Proof of Proposition 3**

1512 We now prove that NAILGUN is essentially optimal up to constants, the content of [Proposition 3](#). We
 1513 first state a more formal version of the result before giving the construction.

1514 **Proposition 20.** *For any $\varepsilon < 1/8$, any $0 < \alpha < 1$ and any $0 < \rho < 1$ such that $\alpha(1 + \rho) < 1$, there*
 1515 *exists a horizon $H = 2 \text{MDP}$ and a deterministic policy class Π of size $|\Pi| = M$ with $M \leq 1 + 1/\rho$*
 1516 *such that there exists a family of ρ -smooth measures $\{\nu_i\}_{i=1}^M$ such that for some $\eta \leq \alpha$, any online IL*
 1517 *algorithm with access to a noisy expert π_η^* must observe*

$$n \gtrsim \frac{1}{\varepsilon(1 - \alpha(1 + \rho)) \cdot \log\left(1 + \frac{1 - \alpha(1 + \rho)}{\alpha\rho}\right)}$$

1518 trajectories in order to guarantee expected regret $J(\pi^*) - J(\hat{\pi}) \leq \varepsilon$.

1519 *Proof.* Let $m = \lceil 1/\rho \rceil$ and let $\Pi = \{\pi^1, \dots, \pi^m\}$ be a class of m deterministic policies. Let
 1520 $\mathcal{A} = [m] \cup \{\perp\}$ and take the corruption level to be $\eta = \alpha$. We suppose there is a deterministic initial
 1521 state s_1 and all policies map to the same action, which transitions to state s_2 with probability q and
 1522 s'_2 with probability $1 - q$ for some q to be determined. For $1 \leq i \leq m$, let

$$\pi^i(s) = \begin{cases} \perp & s = s_1 \\ i & s = s_2, \\ \perp & s = s'_2 \end{cases}$$

1523 and let

$$\nu_i(i) = 0, \quad \nu_i(j) = \rho \text{ for } j \in [m] \setminus \{i\}, \quad \text{and} \quad \nu_i(\perp) = 1 - (m - 1)\rho.$$

1524 In other words, π^i is informative only on state s_2 and reveals the ‘correct’ action, whereas ν_i places
 1525 no mass on i and distributes mass evenly otherwise. Note that by construction ν_i is always ρ -smooth.

1526 We first note that π_i and π_j differ only upon transitioning to s_2 and thus

$$D_{\text{H}^2}(\mathbb{P}^{\pi_i}, \mathbb{P}^{\pi_j}) = q$$

1527 and so for any possibly randomized $\hat{\pi}$,

$$\mathbb{E}[D_{\text{H}^2}(\mathbb{P}^{\pi_i}, \mathbb{P}^{\hat{\pi}})] = q \cdot \mathbb{P}(\hat{\pi}(s_2) \neq i).$$

1528 Thus if $q = 4\varepsilon$, it suffices to show that for insufficiently many rounds of interaction, $\mathbb{P}(\hat{\pi}(s_2) \neq i) \geq$
 1529 $1/4$.

1530 Now, we will let Q_i denote the trajectory distribution obtained by rolling out $(1 - \alpha) \cdot \pi_i + \alpha \cdot \nu_i$
 1531 and we will denote by $p_i = (1 - \alpha)\delta_i + \alpha \cdot \nu_i$ and note that

$$p_i(j) = \begin{cases} 1 - \alpha & j = i \\ \alpha\rho & j \in [m] \setminus \{i\} \\ \alpha(1 - (m - 1)\rho) & j = \perp \\ 0 & \text{otherwise.} \end{cases}$$

1532 Thus,

$$D_{\text{KL}}(p_i \| p_j) = (1 - \alpha) \cdot \log\left(\frac{1 - \alpha}{\alpha\rho}\right) + \alpha\rho \cdot \log\left(\frac{\alpha\rho}{1 - \alpha}\right) = (1 - \alpha(1 + \rho)) \cdot \log\left(1 + \frac{1 - \alpha(1 + \rho)}{\alpha\rho}\right).$$

1533 Moreover, because p_i and p_j differ from each other only upon transitioning to s_2 , it holds by the
 1534 chain rule ([Proposition 4](#)) that

$$D_{\text{KL}}(Q_i \| Q_j) = q \cdot (1 - \alpha(1 + \rho)) \cdot \log\left(1 + \frac{1 - \alpha(1 + \rho)}{\alpha\rho}\right) = 4\varepsilon \cdot (1 - \alpha(1 + \rho)) \cdot \log\left(1 + \frac{1 - \alpha(1 + \rho)}{\alpha\rho}\right).$$

1535 Applying Le Cam’s two-point method to any pair $i \neq j$ and the chain rule again, it holds that in order
 1536 for $\hat{\pi}$ to have a better than $3/4$ probability of selecting the correct index, it must hold that

$$n \gtrsim \frac{1}{4\varepsilon \cdot (1 - \alpha(1 + \rho)) \cdot \log\left(1 + \frac{1 - \alpha(1 + \rho)}{\alpha\rho}\right)}.$$

1537 The result follows. □

1538

1539 Finally, we can prove the main result.

1540 *Proof of Proposition 3.* We observe that

$$\log \left(1 + \frac{1 - \alpha(1 + \rho)}{\alpha\rho} \right) \leq \frac{1 - \alpha(1 + \rho)}{\alpha\rho}$$

1541 and apply Proposition 20. □

1542 **NeurIPS Paper Checklist**

1543 **1. Claims**

1544 Question: Do the main claims made in the abstract and introduction accurately reflect the
1545 paper’s contributions and scope?

1546 Answer: **[Yes]**

1547 Justification: The abstract and introduction state the main theoretical claims, including the
1548 offline exponential separation, the online augmented-trajectory guarantee, the unknown-
1549 corruption extension, and the empirical scope. The assumptions and scope are stated in the
1550 problem setup and theorem statements.

1551 Guidelines:

- 1552 • The answer **[N/A]** means that the abstract and introduction do not include the claims
1553 made in the paper.
- 1554 • The abstract and/or introduction should clearly state the claims made, including the
1555 contributions made in the paper and important assumptions and limitations. A **[No]** or
1556 **[N/A]** answer to this question will not be perceived well by the reviewers.
- 1557 • The claims made should match theoretical and experimental results, and reflect how
1558 much the results can be expected to generalize to other settings.
- 1559 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1560 are not attained by the paper.

1561 **2. Limitations**

1562 Question: Does the paper discuss the limitations of the work performed by the authors?

1563 Answer: **[Yes]**

1564 Justification: We discuss limitations in the discussion, including the stylized corruption
1565 model, the focus on learning from deterministic experts in the unknown corruption setting,
1566 the gap between the theoretical algorithms and practical OPD implementations, and the
1567 limited scale of the empirical validation due to our compute budget.

1568 Guidelines:

- 1569 • The answer **[N/A]** means that the paper has no limitation while the answer **[No]** means
1570 that the paper has limitations, but those are not discussed in the paper.
- 1571 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 1572 • The paper should point out any strong assumptions and how robust the results are to
1573 violations of these assumptions (e.g., independence assumptions, noiseless settings,
1574 model well-specification, asymptotic approximations only holding locally). The authors
1575 should reflect on how these assumptions might be violated in practice and what the
1576 implications would be.
- 1577 • The authors should reflect on the scope of the claims made, e.g., if the approach was
1578 only tested on a few datasets or with a few runs. In general, empirical results often
1579 depend on implicit assumptions, which should be articulated.
- 1580 • The authors should reflect on the factors that influence the performance of the approach.
1581 For example, a facial recognition algorithm may perform poorly when image resolution
1582 is low or images are taken in low lighting. Or a speech-to-text system might not be
1583 used reliably to provide closed captions for online lectures because it fails to handle
1584 technical jargon.
- 1585 • The authors should discuss the computational efficiency of the proposed algorithms
1586 and how they scale with dataset size.
- 1587 • If applicable, the authors should discuss possible limitations of their approach to
1588 address problems of privacy and fairness.
- 1589 • While the authors might fear that complete honesty about limitations might be used by
1590 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
1591 limitations that aren’t acknowledged in the paper. The authors should use their best
1592 judgment and recognize that individual actions in favor of transparency play an impor-
1593 tant role in developing norms that preserve the integrity of the community. Reviewers
1594 will be specifically instructed to not penalize honesty concerning limitations.

1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: All theoretical results state their assumptions, including realizability, known or unknown corruption, deterministic versus stochastic experts, and domination/smoothness conditions. Full and correct proofs are provided in the appendix.

Guidelines:

- The answer [\[N/A\]](#) means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: The experimental setup, methods, task construction, noise model instantiation, architectures, training details, and evaluation metrics are described in [Section 6](#) and [Appendix B](#).

Guidelines:

- The answer [\[N/A\]](#) means that the paper does not include experiments.
- If the paper includes experiments, a [\[No\]](#) answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).

1648 (d) We recognize that reproducibility may be tricky in some cases, in which case
1649 authors are welcome to describe the particular way they provide for reproducibility.
1650 In the case of closed-source models, it may be that access to the model is limited in
1651 some way (e.g., to registered users), but it should be possible for other researchers
1652 to have some path to reproducing or verifying the results.

1653 5. Open access to data and code

1654 Question: Does the paper provide open access to the data and code, with sufficient instruc-
1655 tions to faithfully reproduce the main experimental results, as described in supplemental
1656 material?

1657 Answer: [No]

1658 Justification: We plan to release code and scripts for the camera ready version, but they are
1659 not included with the initial anonymous submission. The appendix provides implementation
1660 details sufficient to understand the main experimental results.

1661 Guidelines:

- 1662 • The answer [N/A] means that paper does not include experiments requiring code.
- 1663 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
1664 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1665 • While we encourage the release of code and data, we understand that this might not
1666 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
1667 including code, unless this is central to the contribution (e.g., for a new open-source
1668 benchmark).
- 1669 • The instructions should contain the exact command and environment needed to run to
1670 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
1671 //neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1672 • The authors should provide instructions on data access and preparation, including how
1673 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1674 • The authors should provide scripts to reproduce all experimental results for the new
1675 proposed method and baselines. If only a subset of experiments are reproducible, they
1676 should state which ones are omitted from the script and why.
- 1677 • At submission time, to preserve anonymity, the authors should release anonymized
1678 versions (if applicable).
- 1679 • Providing as much information as possible in supplemental material (appended to the
1680 paper) is recommended, but including URLs to data and code is permitted.

1681 6. Experimental setting/details

1682 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
1683 rameters, how they were chosen, type of optimizer) necessary to understand the results?

1684 Answer: [Yes]

1685 Justification: The paper specifies the training objectives, task setup, corruption process,
1686 model architecture, optimizer settings, dataset sizes, other hyperparameters, and evaluation
1687 protocol in [Section 6](#) and [Appendix B](#).

1688 Guidelines:

- 1689 • The answer [N/A] means that the paper does not include experiments.
- 1690 • The experimental setting should be presented in the core of the paper to a level of detail
1691 that is necessary to appreciate the results and make sense of them.
- 1692 • The full details can be provided either with the code, in appendix, or as supplemental
1693 material.

1694 7. Experiment statistical significance

1695 Question: Does the paper report error bars suitably and correctly defined or other appropriate
1696 information about the statistical significance of the experiments?

1697 Answer: [Yes]

1698 Justification: The plots report variability across independent runs using shaded error regions.
1699 The appendix specifies the number of seeds and how the error regions are computed.

1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).
- If error bars are reported in tables or plots, the authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The appendix reports the compute resources used for the synthetic and GSM8K experiments, including GPU type and approximate runtime ([Appendix B.1](#)).

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: The research conforms to the NeurIPS Code of Ethics. The work is primarily theoretical and uses synthetic tasks and standard public benchmarks, with no human-subject experiments or private data.

Guidelines:

- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

1752 Question: Does the paper discuss both potential positive societal impacts and negative
1753 societal impacts of the work performed?

1754 Answer: [Yes]

1755 Justification: The paper discusses potential positive impacts, such as more reliable learning
1756 from imperfect teachers, and potential risks, such as improving distillation methods for
1757 language models that could be misused if deployed without safeguards (Appendix A.4).

1758 Guidelines:

- 1759 • The answer [N/A] means that there is no societal impact of the work performed.
- 1760 • If the authors answer [N/A] or [No], they should explain why their work has no societal
1761 impact or why the paper does not address societal impact.
- 1762 • Examples of negative societal impacts include potential malicious or unintended uses
1763 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1764 (e.g., deployment of technologies that could make decisions that unfairly impact specific
1765 groups), privacy considerations, and security considerations.
- 1766 • The conference expects that many papers will be foundational research and not tied
1767 to particular applications, let alone deployments. However, if there is a direct path to
1768 any negative applications, the authors should point it out. For example, it is legitimate
1769 to point out that an improvement in the quality of generative models could be used to
1770 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
1771 that a generic algorithm for optimizing neural networks could enable people to train
1772 models that generate Deepfakes faster.
- 1773 • The authors should consider possible harms that could arise when the technology is
1774 being used as intended and functioning correctly, harms that could arise when the
1775 technology is being used as intended but gives incorrect results, and harms following
1776 from (intentional or unintentional) misuse of the technology.
- 1777 • If there are negative societal impacts, the authors could also discuss possible mitigation
1778 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1779 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1780 feedback over time, improving the efficiency and accessibility of ML).

1781 11. Safeguards

1782 Question: Does the paper describe safeguards that have been put in place for responsible
1783 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
1784 image generators, or scraped datasets)?

1785 Answer: [N/A]

1786 Justification: The paper does not release high-risk pretrained models, scraped datasets, or
1787 systems intended for deployment. The experiments use small synthetic tasks and standard
1788 benchmark data.

1789 Guidelines:

- 1790 • The answer [N/A] means that the paper poses no such risks.
- 1791 • Released models that have a high risk for misuse or dual-use should be released with
1792 necessary safeguards to allow for controlled use of the model, for example by requiring
1793 that users adhere to usage guidelines or restrictions to access the model or implementing
1794 safety filters.
- 1795 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1796 should describe how they avoided releasing unsafe images.
- 1797 • We recognize that providing effective safeguards is challenging, and many papers do
1798 not require this, but we encourage authors to take this into account and make a best
1799 faith effort.

1800 12. Licenses for existing assets

1801 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1802 the paper, properly credited and are the license and terms of use explicitly mentioned and
1803 properly respected?

1804 Answer: [Yes]

1805 Justification: Existing datasets, models, and codebases used in the experiments are cited
1806 in the paper and appendix. We exclusively used open-source code and datasets for all
1807 experiments.

1808 Guidelines:

- 1809 • The answer [N/A] means that the paper does not use existing assets.
- 1810 • The authors should cite the original paper that produced the code package or dataset.
- 1811 • The authors should state which version of the asset is used and, if possible, include a
1812 URL.
- 1813 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1814 • For scraped data from a particular source (e.g., website), the copyright and terms of
1815 service of that source should be provided.
- 1816 • If assets are released, the license, copyright information, and terms of use in the package
1817 should be provided. For popular datasets, paperswithcode.com/datasets has
1818 curated licenses for some datasets. Their licensing guide can help determine the license
1819 of a dataset.
- 1820 • For existing datasets that are re-packaged, both the original license and the license of
1821 the derived asset (if it has changed) should be provided.
- 1822 • If this information is not available online, the authors are encouraged to reach out to
1823 the asset’s creators.

1824 13. New assets

1825 Question: Are new assets introduced in the paper well documented and is the documentation
1826 provided alongside the assets?

1827 Answer: [N/A]

1828 Justification: The paper does not introduce a new benchmark dataset or pretrained model
1829 as a standalone released asset. Any released code will be documented with reproduction
1830 instructions.

1831 Guidelines:

- 1832 • The answer [N/A] means that the paper does not release new assets.
- 1833 • Researchers should communicate the details of the dataset/code/model as part of their
1834 submissions via structured templates. This includes details about training, license,
1835 limitations, etc.
- 1836 • The paper should discuss whether and how consent was obtained from people whose
1837 asset is used.
- 1838 • At submission time, remember to anonymize your assets (if applicable). You can either
1839 create an anonymized URL or include an anonymized zip file.

1840 14. Crowdsourcing and research with human subjects

1841 Question: For crowdsourcing experiments and research with human subjects, does the paper
1842 include the full text of instructions given to participants and screenshots, if applicable, as
1843 well as details about compensation (if any)?

1844 Answer: [N/A]

1845 Justification: The paper does not involve crowdsourcing, human-subject experiments, or
1846 collection of human participant data.

1847 Guidelines:

- 1848 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
1849 with human subjects.
- 1850 • Including this information in the supplemental material is fine, but if the main contribu-
1851 tion of the paper involves human subjects, then as much detail as possible should be
1852 included in the main paper.
- 1853 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1854 or other labor should be paid at least the minimum wage in the country of the data
1855 collector.

1856 **15. Institutional review board (IRB) approvals or equivalent for research with human**
1857 **subjects**

1858 Question: Does the paper describe potential risks incurred by study participants, whether
1859 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1860 approvals (or an equivalent approval/review based on the requirements of your country or
1861 institution) were obtained?

1862 Answer: [N/A]

1863 Justification: The paper does not involve crowdsourcing or human-subject research, so IRB
1864 approval is not applicable.

1865 Guidelines:

- 1866 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
1867 with human subjects.
- 1868 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1869 may be required for any human subjects research. If you obtained IRB approval, you
1870 should clearly state this in the paper.
- 1871 • We recognize that the procedures for this may vary significantly between institutions
1872 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1873 guidelines for their institution.
- 1874 • For initial submissions, do not include any information that would break anonymity (if
1875 applicable), such as the institution conducting the review.

1876 **16. Declaration of LLM usage**

1877 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1878 non-standard component of the core methods in this research? Note that if the LLM is used
1879 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
1880 scientific rigor, or originality of the research, declaration is not required.

1881 Answer: [N/A]

1882 Justification: LLMs are studied as part of the experimental domain, but no LLM is used as
1883 an important, original, or non-standard component of the research methodology beyond the
1884 models explicitly described in the experiments.

1885 Guidelines:

- 1886 • The answer [N/A] means that the core method development in this research does not
1887 involve LLMs as any important, original, or non-standard components.
- 1888 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
1889 be described.