
Learning When to Stop: Selective Imitation Learning Under Arbitrary Dynamics Shift

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Behavior cloning provides strong imitation learning guarantees when training and
2 test environments share the same dynamics. However, in many deployment set-
3 tings the test environment’s transitions differ from training, and classical offline
4 IL offers no recourse: the learner must commit to an action at every state, even
5 when its demonstrations are uninformative and could lead to arbitrary degradation
6 of performance. This motivates the study of *selective* imitation, where the learner
7 may choose to *stop* when it cannot act reliably. We introduce a model for selective
8 imitation under arbitrary dynamics shift: given labeled expert demonstrations from
9 a training environment and unlabeled state trajectories from the same expert in a
10 test environment, the learner outputs a *selective policy* that is *complete* (rarely stops
11 in training) and *sound* (incurs low regret before stopping in test). Our algorithm,
12 SeqRejectron, constructs a stopping rule using a small set of *validator policies*
13 whose size is independent of the horizon or policy class. For deterministic policies,
14 this yields horizon-free $\tilde{O}(\log |\Pi|/\epsilon^2)$ sample complexity, assuming sparse costs.
15 For stochastic policies, we obtain analogous horizon-free guarantees using a cumu-
16 lative Hellinger stopping time. We extend the framework to misspecified experts
17 and different expert policies across train and test and obtain results that gracefully
18 degrade with the amount of misspecification.

19 1 Introduction

20 A central challenge in deploying learned policies is *environment shift*: the dynamics governing state
21 transitions at test time may differ from those encountered during training. This arises ubiquitously: in
22 sim-to-real transfer for robotics [Tobin et al., 2017], in autonomous driving where conditions vary
23 between simulation and the road [Codevilla et al., 2019], and in any domain where a policy trained
24 on demonstrations must act in a changed environment.

25 Classical imitation learning offers strong guarantees when training and test environments coincide,
26 but these guarantees become vacuous under even modest dynamics shift. Even with infinite data, a
27 policy may perform arbitrarily poorly if the test environment’s state support differs from the training
28 environment’s. The fundamental issue is that it cannot recognize when its training data no longer
29 provides reliable guidance, and has no mechanism to stop before things go wrong. This motivates a
30 basic question: *can a policy learn not only what to do, but when to stop?*

31 We develop a framework for *learning to imitate with the option to abstain*, drawing on the *PQ learning*
32 paradigm from selective classification [Goldwasser et al., 2020]. We consider a practical setting where
33 the learner is given labeled (state, action) expert demonstrations from a training environment M , but
34 only unlabeled (state-only) expert trajectories from a test environment N . This asymmetry naturally
35 models settings where passive observation is cheap and abundant, but action telemetry is costly or
36 proprietary. For instance, in autonomous driving, a learner might train on rich kinematic data from a

37 sensor-equipped fleet, but must generalize to new cities using massive datasets of public dashcam
38 footage—which provide valid visual states of expert driving, but completely lack the underlying
39 steering commands [Xu et al., 2017, Torabi et al., 2018, Codevilla et al., 2019]. In healthcare, a
40 sepsis treatment policy trained on a fully instrumented ICU system may be deployed to a setting
41 where precise dosage decisions are unavailable [Komorowski et al., 2018, Finlayson et al., 2021]. In
42 algorithmic trading, a model trained on historical market conditions with full access to proprietary
43 executed trades must often navigate novel macroeconomic regimes where only public ticker data is
44 observable, leaving the exact buy and sell actions of top experts completely hidden.

45 The goal is to learn a *selective policy* that acts when its training data supports confident prediction
46 and stops execution otherwise. We require *completeness* (the policy rarely stops when deployed in
47 the training environment) and *soundness* (whenever the policy chooses to act in the test environment,
48 it matches the expert’s performance). A policy that always stops is trivially sound but useless; one
49 that never stops is trivially complete but unsafe. The tension makes the problem nontrivial.

50 1.1 Our approach and results

51 We certify trajectory prefixes using *validator policies*: competing explanations of the expert drawn
52 from the training-consistent version space. The learner executes a base policy π_0 , continuing only
53 while every validator agrees and abstaining at the first disagreement. Because a small validator set
54 always suffices, the resulting sample complexity depends on the horizon H only through the trajectory
55 cost scale C_{\max} and the complexity of the policy class. We call this approach SeqRejectron. It
56 generalizes the selective classification algorithm of Goldwasser et al. [2020] to sequential decision-
57 making, mirroring recent horizon-independent rates for imitation learning [Foster et al., 2024] while
58 accommodating dynamics shift. We analyze it in three increasingly general settings.

59 **Deterministic policies (Section 3).** We first develop our framework for finite classes of deterministic
60 policies. In this setting, each validator triggers at the first state where it disagrees with the base policy.
61 We prove that a surprisingly sparse set of these validators is always sufficient to safely stop execution.
62 Additionally, we demonstrate that this approach is oracle-efficient by providing an exact reduction to
63 multiple-instance learning [Maron and Lozano-Pérez, 1997].

64 **Stochastic policies (Section 4).** The binary disagreement test is too coarse for stochastic policies,
65 where two distributions may differ softly at every step without any single step being alarming. We
66 replace it with a cumulative Hellinger stopping time that tracks a trajectory-level divergence budget,
67 preserving the validator template. The main result gives finite-samples bounds on both the stopping
68 rate and the regret. This bound has both a *cost-driven* and *variance-driven* exponent. We prove
69 an $\Omega(1/\epsilon^3)$ labeled-sample lower bound already in the single-step stochastic setting, matching the
70 cost-driven part of the upper bound and leaving the variance-driven exponent open.

71 **Misspecified policy classes (Section 5).** Under misspecification, the exact-consistency version
72 space might not contain the expert. We replace it with a symmetrically regularized game that softly
73 penalizes source-data disagreement, ensuring the validator distribution remains well-defined and
74 sparse. We show that the guarantees degrade gracefully under misspecification.

75 1.2 Related work

76 Our work builds on three strands of literature.¹ First, it is closely related to imitation learning and
77 recent analyses of behavior cloning. Behavior cloning reduces imitation learning to supervised
78 prediction, but is classically limited by compounding error under distribution shift [Pomerleau, 1988,
79 Ross and Bagnell, 2010, Ross et al., 2011]. Recent theory has further characterized the statistical
80 limits of offline imitation learning and identified settings in which horizon-independent guarantees
81 are possible [Rajaraman et al., 2020, Swamy et al., 2022, Foster et al., 2024]. We study a modified
82 setting, where the learner is trained offline, but the test dynamics may shift arbitrarily.

83 Second, our framework is inspired by selective classification and learning under arbitrary distribution
84 shift. Closest to our setting is the PQ framework, which studies source-labeled, target-unlabeled
85 learning with abstention under arbitrary shift [Goldwasser et al., 2020, Kalai and Kanade, 2021, Goel
86 et al., 2024]. Our framework can be viewed as a sequential analogue of PQ learning, where the
87 learner abstains along a trajectory through a stopping rule rather than on a single prediction.

¹We defer a more comprehensive literature review to Section A.

88 Third, our work is related to robustness and uncertainty in reinforcement learning and robotics,
 89 including robust control, sim-to-real transfer, domain randomization, and OOD detection [Iyengar,
 90 2005, Nilim and El Ghaoui, 2005, Wiesemann et al., 2013, Tobin et al., 2017, Peng et al., 2018,
 91 Chebotar et al., 2019, Chae et al., 2022, Haider et al., 2023]. These approaches aim to make policies
 92 robust or adaptive to changing environments, but are typically heuristic, make strong assumptions
 93 about the nature of the shift, or require online access to the test distribution. By contrast, we study
 94 distribution-free guarantees under arbitrary dynamics shift using only labeled source demonstrations
 95 and unlabeled target trajectories, with abstention² as the safety mechanism.

96 2 Problem Formulation

97 **Markov decision processes.** A finite-horizon Markov decision process (MDP) is a tuple $M =$
 98 $(\mathcal{S}, \mathcal{A}, P, c, H)$, where \mathcal{S} is the state space, \mathcal{A} is the action space³, H is the horizon, $P = \{P_h\}_{h=0}^{H-1}$
 99 specifies the state evolution, and $c = \{c_h\}_{h=1}^H$ are per-step cost functions $c_h : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$. Here
 100 $P_0 \in \Delta(\mathcal{S})$ is the initial-state distribution, and for each $h = 1, \dots, H-1$, $P_h(\cdot | s, a) \in \Delta(\mathcal{S})$ is
 101 the transition kernel. A (possibly nonstationary, possibly randomized) policy is $\pi = \{\pi_h\}_{h=1}^H$ with
 102 $\pi_h : \mathcal{S} \rightarrow \Delta(\mathcal{A})$. Together, (M, π) induce a distribution over trajectories $(s_1, a_1), \dots, (s_H, a_H)$ via
 103 $s_1 \sim P_0, a_h \sim \pi_h(\cdot | s_h)$, and $s_{h+1} \sim P_h(\cdot | s_h, a_h)$ for $h = 1, \dots, H-1$. We denote the expected
 104 cumulative cost of π in M by $J_M(\pi; c) := \mathbb{E}^{M, \pi} \left[\sum_{h=1}^H c_h(s_h, a_h) \right]$.

105 **Stopping times.** A key feature of our framework is that the learner may stop execution mid-trajectory.
 106 We formalize this decision to abstain via a *stopping time* adapted to the observed history. Denote
 107 the canonical sample space of full state-action trajectories by $\Omega = (\mathcal{S} \times \mathcal{A})^H$. Let $\mathbb{G} = \{\mathcal{G}_h\}_{h=1}^H$
 108 be the pre-action filtration on Ω , where $\mathcal{G}_h = \sigma(s_1, a_1, \dots, s_{h-1}, a_{h-1}, s_h)$. Thus \mathcal{G}_h contains
 109 exactly the information available at step h , just before the learner chooses a_h . For convenience,
 110 we also let $\mathcal{G}_{H+1} = \sigma(s_1, a_1, \dots, s_H, a_H)$. We define the point of abstention as a stopping time
 111 $\tau : \Omega \rightarrow \{1, \dots, H, H+1\}$ adapted to the filtration \mathbb{G} . The event $\{\tau = h\}$ signifies that the
 112 learner stops execution at step h , having observed the history up to s_h . We adopt the convention that
 113 $\tau = H+1$ if the learner completes the full trajectory without abstaining.

114 **Problem statement.** Let $M := (\mathcal{S}, \mathcal{A}, P, c, H)$ and $N := (\mathcal{S}, \mathcal{A}, Q, c, H)$ be MDPs that share
 115 state space, action space, costs, and horizon, but may differ in their dynamics $P = \{P_h\}_{h=0}^{H-1}$ and
 116 $Q = \{Q_h\}_{h=0}^{H-1}$. We refer to M as the *training environment* and N as the *test environment*. For any
 117 environment $E \in \{M, N\}$ and policy π , we write \mathcal{P}_E^π for the law of the trajectory under (E, π) ,
 118 and $\mathcal{P}_{E, \text{state}}^\pi$ for the law of the state trajectory $s_{1:H}$. Let C_{\max} denote an upper bound on the total
 119 trajectory cost $\sum_{h=1}^H c_h(s_h, a_h)$.

120 Let Π be a finite policy class and let $\pi^* \in \Pi$ denote an expert policy whose behavior we seek to
 121 imitate. We do not assume that π^* is optimal for c . The learner observes demonstrations generated by
 122 π^* , but does not observe realized costs or expert actions in N :

- 123 • **Labeled training rollouts:** drawn i.i.d. from (M, π^*) , $\mathcal{D}_{\text{train}} = \{(s_{k,1:H}, a_{k,1:H})\}_{k=1}^m$.
- 124 • **Unlabeled test rollouts:** state-only trajectories drawn i.i.d. from (N, π^*) , $\mathcal{D}_{\text{test}} = \{t_{l,1:H}\}_{l=1}^n$.

125 Without the option to abstain, the problem is poorly posed under environment shift: if the expert-
 126 induced state supports of M and N are disjoint, even infinite labeled data from M cannot determine
 127 the expert’s actions on test-time states. To maintain low regret across all valid expert hypotheses, the
 128 learner is forced to *abstain* on unseen test states. We therefore allow the learner to stop execution
 129 mid-trajectory. We formalize this via the notion of a *selective policy*.

130 **Definition 1** (Selective policy, stopped regret, and stopping rate). *A selective policy is a pair*
 131 *(π, τ) , where $\pi \in \Pi$ is a policy and τ is a stopping time adapted to the pre-action filtration*
 132 *$\mathcal{G}_h = \sigma(s_1, a_1, \dots, s_{h-1}, a_{h-1}, s_h)$. When deployed, it incurs cost only on the time steps before τ .*
 133 *We define the stopped cost and stopped regret (in the test environment N) as*

$$J_N(\pi, \tau; c) := \mathbb{E}^{N, \pi} \left[\sum_{h=1}^{\tau-1} c_h(s_h, a_h) \right], \quad \text{Regret}_N(\pi, \tau; c) := J_N(\pi, \tau; c) - J_N(\pi^*, \tau; c).$$

²Throughout, we will interchangeably refer to *abstention* as *stopping*.

³To simplify presentation, we assume throughout that \mathcal{S} and \mathcal{A} are countable; the definitions and results extend to general spaces with the appropriate measure-theoretic treatment.

134 We define the stopping rate (in the training environment M) as $\alpha_M(\pi, \tau) := \Pr_{M, \pi}(\tau \leq H)$.

135 Intuitively, a selective policy can be thought of as an agent equipped with a fail-safe, executing its
 136 learned behavior only as long as it remains confident. For example, consider a self-driving car trained
 137 exclusively to navigate dry, sunny roads. If it is suddenly deployed in a blizzard where tire traction
 138 and visibility drastically change, a standard imitation policy might blindly attempt a normal turn and
 139 crash. A selective policy, by contrast, recognizes that these icy conditions fall outside its training
 140 support and safely pulls over before a catastrophic error occurs.

141 **Problem 1** (Selective imitation learning: stopped regret formulation). *Given labeled training rollouts*
 142 $\mathcal{D}_{\text{train}}$ *and unlabeled test rollouts* $\mathcal{D}_{\text{test}}$, *output a selective policy* (π, τ) *satisfying:*

- 143 • **Completeness:** $\alpha_M(\pi, \tau) \leq \epsilon$. (The policy rarely abstains in M).
- 144 • **Soundness:** $\text{Regret}_N(\pi, \tau; c) \leq \epsilon$. (The policy incurs low stopped regret in N).

145 3 Selective Execution for Deterministic Policies

146 We first present our results for the case when π^* and Π are deterministic policies. We first analyze
 147 a per-step baseline that suffers an avoidable horizon factor H by treating rollouts as independent
 148 decisions. We then resolve this via a trajectory-level construction based on *validator policies* that
 149 certifies entire execution prefixes, yielding an improved sample complexity.

150 **First attempt: Stepwise reduction to PQ-learning.** A naive baseline applies Rejection [Goldwasser
 151 et al., 2020] independently at each step. However, converting these per-step guarantees to a trajectory-
 152 level bound requires a union bound over H steps, yielding a suboptimal sample complexity of
 153 $\tilde{O}(H^2 C_{\text{max}}^2 \log |\Pi| / \epsilon^2)$. This introduces an avoidable H^2 penalty even when costs are sparse; we
 154 defer the formal analysis of this stepwise reduction to Section B.2.

155 **Trajectory-level validation via validator policies.** The per-step reduction treats a trajectory as H
 156 independent decisions, introducing an avoidable horizon factor. However, sequential rollouts have
 157 inherent structure: until the learner diverges from the expert, stopping incurs only a prefix of the
 158 expert’s cost. The danger is failing to stop before this first unrejected deviation. Therefore, our goal
 159 is to certify an entire execution prefix at once.

160 Fix a base policy $\pi_0 \in \Pi_{\text{version}}$. To determine when it is safe to continue executing π_0 , we compare it
 161 against other policies in the version space. For any $\pi \in \Pi_{\text{version}}$, the first state where π disagrees with
 162 π_0 marks a point where π_0 is no longer justified by that competing explanation of the expert. Since
 163 no single competitor reliably detects all problematic rollouts across all trajectories, we compare π_0
 164 against a *set* of validators, continuing only while every validator agrees with π_0 .

165 **Definition 2** (Validator-induced stopping time). *For any base policy* π_0 *and set of policies* $\Phi \subseteq$
 166 Π_{version} , *define the stopping time* $\tau_{\pi_0, \Phi}$ *as the first time along the realized trajectory at which at least*
 167 *one validator in* Φ *disagrees with* π_0 , *or* $H + 1$ *if no such disagreement occurs:*

$$\tau_{\pi_0, \Phi}(T) := \min \left(\{h \in [H] : \exists \pi \in \Phi, \pi_h(s_h) \neq \pi_0_h(s_h)\} \cup \{H + 1\} \right). \quad (1)$$

168 The key question is therefore how to choose the validator set Φ . Ideally, we would like Φ to be
 169 such that $\Pr_{\pi_0, N}[\tau_{\pi_0, \Phi} \leq \tau_{\pi_0, \{\pi^*\}}]$ is close to 1. That is, with high probability, Φ tells us to stop
 170 before π_0 deviates from π^* . This tells us that until the stopping time $\tau_{\pi_0, \Phi}$, the policies π_0 and π^*
 171 are approximately coupled, and hence incur similar cost. We would also like Φ to be sparse; this prevents
 172 us from rejecting too many trajectories, which is a form of overfitting.

173 The key difficulty, however, is that Φ must be computed *without knowledge of* π^* . We will thus
 174 ask: *Can we find a sparse validator set which stops before any policy in* Π_{version} *deviates from* π_0 ?
 175 There is no obvious reason that one small deterministic validator set should work uniformly against
 176 every policy in Π_{version} . Fortunately, however, the following result shows that it is always possible to
 177 compute a *distribution* over sparse validator sets which satisfies this guarantee with high probability.
 178 We thus sample from this distribution to obtain a sparse validator set for π^* .

179 **Definition 3** (ρ -valid validator distribution). *Fix* $\pi_0 \in \Pi_{\text{version}}$ *and unlabeled test trajectories*
 180 $\mathcal{D}_{\text{test}} = \{T_j\}_{j=1}^n$. *A distribution* q *over subsets* $\Phi \subseteq \Pi_{\text{version}}$ *is called* ρ -*valid for* π_0 *if*

$$181 \sup_{\pi \in \Pi_{\text{version}}} \mathbb{E}_{\Phi \sim q} \left[\frac{1}{n} \sum_{j=1}^n \mathbf{1}[\tau_{\pi_0, \Phi}(T_j) > \tau_{\pi_0, \{\pi\}}(T_j)] \right] \leq \rho.$$

182 In words, a ρ -valid distribution places mass on validator sets that, on almost all empirical test
 183 trajectories, stop no later than when *any* fixed competing policy in the version space diverges from π_0 .
 184 The next lemma shows that this averaged notion of coverage can always be achieved by distributions
 185 supported on *small* validator sets. Importantly, the support-size bound depends only on the target
 186 violation level ρ ; downstream, ρ will only depend on ϵ and C_{\max} .

187 **Lemma 1** (Existence of sparse valid validator distributions). *For every $\rho \in (0, 1)$, there exists a*
 188 *ρ -valid distribution q^* over subsets $\Phi \subseteq \Pi_{\text{version}}$ such that every $\Phi \in \text{supp}(q^*)$ satisfies $|\Phi| \leq \lceil \frac{1}{\rho} \rceil$.*

189 The proof uses only two ingredients: the minimax theorem and the exchangeability of i.i.d. draws.
 190 Importantly, there is no dependence on the state space, action space, horizon, or size of Π . This
 191 generality is important, and we will reuse it for stochastic policies (Section 4) and the misspecified
 192 setting (Section 5). We state the general version as Lemma 5 in the appendix.

193 **How do we compute q^* ?** To keep the statistical argument separate from the computation, we
 194 package Lemma 1 into a subroutine `SPARSEVALIDATORDIST`($\Pi_{\text{version}}, \pi_0, \mathcal{D}_{\text{test}}, \rho, \xi, \delta$), which with
 195 probability $1 - \delta$, returns a $(\rho + \xi)$ -valid distribution over subsets of Π_{version} with size $\leq \lceil 1/\rho \rceil$. Here
 196 ξ is a computational slack term that can be driven arbitrarily close to 0. A concrete, oracle-efficient
 197 implementation is given in Section B.6, where we realize this subroutine via no-regret dynamics.

198 **Algorithm and guarantees.** We are now ready to present our algorithm for selective imitation with
 199 deterministic policies. The algorithm (Algorithm 1) has a clean, simple structure. The labeled data
 200 define a version space and a base policy π_0 . The unlabeled test trajectories are used to construct a
 201 sparse validator distribution, from which we sample a validator set. The returned selector continues
 202 while every validator agrees with π_0 and abstains at the first disagreement.

Algorithm 1 SeqRejectron for deterministic policies

Require: deterministic policy class Π , labeled rollouts $\mathcal{D}_{\text{train}} = \{(s_{ih}, a_{ih})\}_{i \in [m], h \in [H]}$, unlabeled
 rollouts $\mathcal{D}_{\text{test}} = \{T_j\}_{j=1}^n$, tolerance η , computational slack ξ , confidence δ
 1: Compute the training-consistent version space

$$\Pi_{\text{version}} \leftarrow \{\pi \in \Pi : \pi_h(s_{ih}) = a_{ih} \text{ for all } i \in [m], h \in [H]\}.$$

- 2: Choose any base policy $\pi_0 \in \Pi_{\text{version}}$.
 - 3: Let $q^* \leftarrow \text{SPARSEVALIDATORDIST}(\Pi_{\text{version}}, \pi_0, \mathcal{D}_{\text{test}}, \eta/2, \xi, \delta/5)$.
 - 4: Sample $\Phi_1, \dots, \Phi_k \stackrel{\text{i.i.d.}}{\sim} q^*$ with $k \leftarrow \lceil \log_2(5/\delta) \rceil$ and set $\Phi \leftarrow \bigcup_{r=1}^k \Phi_r$.
 - 5: **return** the selective policy $(\pi_0, \tau_{\pi_0, \Phi})$.
-

203 **Theorem 1** (SeqRejectron guarantee for deterministic classes). *Run Algorithm 1 with hyperpa-*
 204 *rameters $\eta, \delta, \xi > 0$. Define $Z := (\lceil \log_2(5/\delta) \rceil \lceil 2/\eta \rceil + 1) \log |\Pi| + \log(5/\delta)$. Then, with*
 205 *probability at least $1 - \delta$, the selective policy $(\pi_0, \tau_{\pi_0, \Phi})$ satisfies $\alpha_M(\pi_0, \tau_{\pi_0, \Phi}) \leq \frac{2Z}{m}$ and*
 206 *$\text{Regret}_N(\pi_0, \tau_{\pi_0, \Phi}; c) \leq C_{\max} \left(\eta + 2\xi + \sqrt{\frac{2(\eta+2\xi)Z}{n}} + \frac{3Z}{n} \right)$.*

207 **Corollary 1.** *Suppose $m = n \geq 6 \lceil \log_2(5/\delta) \rceil \log |\Pi|$. Set $\eta := \sqrt{6 \lceil \log_2(5/\delta) \rceil \log |\Pi| / n}$. Then,*
 208 *with probability at least $1 - \delta$, both the stopping rate and $\text{Regret}_N / C_{\max}$ scale as $\tilde{O}(\sqrt{\log |\Pi| / n})$.*

209 *The role of weight sharing in $\log |\Pi|$.* The sample complexity of Theorem 1 scales with $\log |\Pi|$.
 210 When Π shares parameters across time steps (e.g., a single neural network mapping states to actions),
 211 $\log |\Pi|$ is independent of H and the bound is horizon-free in the strongest sense. In general, H enters
 212 only through C_{\max} and $\log |\Pi|$, consistent with Foster et al. [2024].

213 *Bounding the test-side abstention.* It is also natural to ask how often the selective policy abstains in
 214 N . When the environment shift is bounded in total variation, the test-side stopping rate $\alpha_N(\pi_0, \tau)$
 215 inherits the training-side guarantees up to a TV correction term, plus an *additional* sequential penalty
 216 that arises because the deployed policy's actions can steer the trajectory out of the safe prefix before
 217 the validator stops. See Section B.5 in the appendix for a precise statement.

218 **4 Selective Execution for Stochastic Policies**

219 Throughout this section, Π denotes a finite class of (possibly stochastic, possibly nonstationary)
 220 policies, and $\pi^* \in \Pi$ is the expert. Recall that the deterministic construction rests on a binary
 221 primitive: two policies either agree or disagree at each state. For stochastic policies this is too coarse,
 222 since two distributions may differ softly at every step without any single discrepancy being alarming.
 223 The right question is whether soft discrepancies *accumulate* enough to cause meaningful divergence
 224 from the expert. We therefore replace three ingredients from the deterministic construction.

225 First, we take π_0 to be the MLE and for $\gamma > 0$ define the version space as a log-loss ball,

$$\Pi_{\text{version}}^\gamma := \{ \pi \in \Pi : \text{LogLoss}(\pi, \mathcal{D}_{\text{train}}) \leq \text{LogLoss}(\pi_0, \mathcal{D}_{\text{train}}) + \gamma \},$$

226 where $\text{LogLoss}(\pi, \mathcal{D}_{\text{train}}) := -\frac{1}{m} \sum_{i=1}^m \sum_{h=1}^H \log \pi_h(a_{ih} \mid s_{ih})$. Second, we replace the first-
 227 disagreement stopping time with one based on cumulative squared Hellinger distance⁴: for $p, q \in$
 228 $\Delta(\mathcal{A})$, let $d_H^2(p, q) := 1 - \sum_a \sqrt{p(a)q(a)}$, and for $\theta > 0$ define

$$\tau_{\pi_0, \Phi}^\theta(T) := \min \left(\left\{ h \in [H] : \exists \pi \in \Phi, \sum_{k=1}^h d_H^2(\pi_k(\cdot \mid s_k), \pi_{0,k}(\cdot \mid s_k)) > \theta \right\} \cup \{H + 1\} \right).$$

229 Third, we define an analog to the SPARSEVALIDATORDIST subroutine: STOCHASTICSPARSEVAL-
 230 IDATORDIST, which when called with parameters ρ, θ , and δ returns a distribution q^* supported over
 231 subsets $\Phi \subseteq \Pi_{\text{version}}^\gamma$ of size at most $\lceil 1/\rho \rceil$, such that⁵, with probability $1 - \delta$, for every $\pi \in \Pi_{\text{version}}^\gamma$,
 232 it holds that $\mathbb{E}_{\Phi \sim q^*} \left[\frac{1}{n} \sum_{j=1}^n \mathbf{1} \left[\tau_{\pi_0, \Phi}^\theta(T_j) > \tau_{\pi_0, \{\pi\}}^\theta(T_j) \right] \right] \leq \rho$.

233 **Guarantees.** The stochastic SeqRejectron follows the same template as Algorithm 1, save for the
 234 changes mentioned above. To formalize our guarantees, let $\mathcal{P}_N^{\pi|\tau}$ denote the distribution of trajectories
 235 induced by π , stopped at time τ in environment N (i.e., $\mathcal{P}_N^{\pi|\tau}$ is supported over trajectory *prefixes*).⁶

236 **Theorem 2** (SeqRejectron guarantee for stochastic classes). *There exists an algorithm, the stochastic*
 237 *SeqRejectron, with hyperparameters $\eta, \theta, \gamma, \delta > 0$. Define $K_{\text{ens}} := \lceil \log_2(4/\delta) \rceil \lceil 2/\eta \rceil$ and $Z :=$
 238 $(K_{\text{ens}} + 1) \log |\Pi| + \log(4/\delta)$. If $n \geq 8Z/\eta$ and $m \geq (\log |\Pi| + \log(8/\delta))/\gamma$, then with probability
 239 at least $1 - \delta$, this algorithm returns a selective policy $(\pi_0, \tau_{\pi_0, \Phi}^\theta)$ satisfying $\alpha_M(\pi_0, \tau_{\pi_0, \Phi}^\theta) \leq$
 240 $K_{\text{ens}} \left(\frac{12}{\theta} + 48 \right) \gamma$ and $D_H^2(\mathcal{P}_N^{\pi_0|\tau}, \mathcal{P}_N^{\pi^*|\tau}) \leq 3(\theta + \eta)$.*

241 **Corollary 2.** *There exists a choice of parameters γ, η, θ for stochastic SeqRejectron (depending*
 242 *only on $m, n, \sigma_{\pi^*}^2, \log |\Pi|$, and δ) such that with probability at least $1 - \delta$, both $\alpha_M(\pi_0, \tau_{\pi_0, \Phi}^\theta)$ and*
 243 *Regret_N($\pi_0, \tau_{\pi_0, \Phi}^\theta; c$)/ C_{max} are $O \left(\max \left\{ \frac{(\log(1/\delta) + \log |\Pi|)^{2/5}}{m^{1/5}}, \frac{(\log(1/\delta) + \log |\Pi|)^{1/2}}{n^{1/4}} \right\} \right)$.*

244 To prove Corollary 2, we convert the guarantee of Theorem 2 into a total variation bound. A natural
 245 question is whether this bound can be sharpened, as Theorem 2 is stated as a guarantee on the *squared*
 246 *Hellinger* distance of the stopped trajectories. Indeed, Foster et al. [2024] use a variance-dependent
 247 Hellinger change-of-measure lemma to avoid this apparent looseness in the *full horizon* setting. Their
 248 bounds are stated in terms of an expert variance parameter $\sigma_{\pi^*}^2$; in our setting, however, this variance
 249 parameter would depend on the *learned* stopping time. Thus, it is not clear how to make full use
 250 of the squared Hellinger guarantee of Theorem 2 in our setting. It is an interesting question for
 251 future work whether or not a different algorithm or analysis can sharpen Corollary 2. However, in
 252 the following section, we present a natural variation of our objective which is more amenable to the
 253 sharper Hellinger change-of-measure, and prove sharper rates accordingly.

254 **Computational implementation and offline oracles.** Unlike the deterministic setting, the cumulative
 255 Hellinger thresholding required for stochastic policies does not yet admit a known oracle-efficient
 256 reduction to standard ERM. We leave the development of computational oracles for this cumulative
 257 objective—or rigorous guarantees for practical surrogates—to future work.

⁴For deterministic Π , d_H^2 reduces to the disagreement indicator, so $\tau_{\pi_0, \Phi}^\theta$ with $\theta < 1$ recovers (1).

⁵The existence of q^* follows the exact same combinatorial argument as Lemma 1, and is stated as Lemma 5 in the appendix.

⁶More formally, $\mathcal{P}_N^{\pi|\tau}$ is a probability measure over the space of finite histories $\mathcal{H} = \bigcup_{h=1}^{H+1} (\mathcal{S} \times \mathcal{A})^{h-1} \times \mathcal{S}$. A random draw from this measure takes the form $(s_1, a_1, \dots, s_{\tau-1}, a_{\tau-1}, s_\tau)$ on the event $\{\tau \leq H\}$, and is the full trajectory $(s_1, a_1, \dots, s_H, a_H)$ on the event $\{\tau = H + 1\}$.

258 4.1 Expert handoff interpretation and switched regret guarantees

259 So far, we have studied policies that can stop execution mid-trajectory. However, in collaborative
 260 autonomy—such as a self-driving car handing control back to a human driver upon encountering
 261 unfamiliar condition—the episode does not end at abstention. We would like to guarantee that the
 262 autonomous system does not drive the system into an unrecoverable state before handing off. Thus,
 263 one may consider the *switched policy* that follows the learner until τ and the expert thereafter.

264 **Definition 4** (Switched policy and switched regret). *For a base policy $\pi \in \Pi$, an expert policy*
 265 *$\pi^* \in \Pi$, and a stopping time τ , the switched policy $\pi^{\text{sw}}[\pi, \tau, \pi^*]$ executes π on steps $1, \dots, \tau - 1$*
 266 *and π^* from step τ onward. For costs $c = \{c_h\}_{h=1}^H$, we define the switched regret (in N) as*

$$\text{Regret}_N^{\text{sw}}(\pi, \tau; c) := J_N(\pi^{\text{sw}}[\pi, \tau, \pi^*]; c) - J_N(\pi^*; c). \quad (2)$$

267 **Problem 2** (Selective imitation learning: switched regret formulation). *Given labeled training*
 268 *rollouts $\mathcal{D}_{\text{train}}$ and unlabeled test rollouts $\mathcal{D}_{\text{test}}$, output a selective policy (π, τ) such that*

- 269 • **Completeness:** $\alpha_M(\pi, \tau) \leq \epsilon$. (The policy rarely abstains in M).
- 270 • **Soundness:** $\text{Regret}_N^{\text{sw}}(\pi, \tau; c) \leq \epsilon$. (The policy incurs low switched regret in N).

271 The switched regret is directly controlled by our Lemma 2, which generalizes the Hellinger regret
 272 decomposition of Foster et al. [2024]. Below, let $\sigma_{\pi^*}^2 := \sum_{h=1}^H \mathbb{E}_{\pi^*} [(V_h^{\pi^*}(x_h) - Q_h^{\pi^*}(x_h, a_h))^2] \leq$
 273 C_{max}^2 denote the expert-variance parameter of Foster et al. [2024] where (x_h, a_h) is the state-action
 274 pair at time h under π^* and $V_h^{\pi^*}, Q_h^{\pi^*}$ are the corresponding value and action-value functions.

275 **Lemma 2** (Hellinger-based regret decomposition for stopped trajectories; generalization of Theorem
 276 3.1 of Foster et al. [2024]). *Assume $C_{\text{max}} \geq 1$. Let π^* be the expert policy. For any learned policy*
 277 *$\hat{\pi}$ and stopping time τ , let π_{sw} denote the policy that executes $\hat{\pi}$ for $h < \tau$ and executes π^* for all*
 278 *remaining steps $h \geq \tau$. Then, for any $\epsilon \in (0, e^{-1})$, the expected regret of π_{sw} is bounded by:*

$$J(\pi_{\text{sw}}) - J(\pi^*) \leq \sqrt{6\sigma_{\pi^*}^2 \cdot D_H^2(\mathcal{P}_N^{\hat{\pi}|_{1:\tau}}, \mathcal{P}_N^{\pi^*|_{1:\tau}})} + O(C_{\text{max}} \log(C_{\text{max}} \epsilon^{-1})) \cdot D_H^2(\mathcal{P}_N^{\hat{\pi}|_{1:\tau}}, \mathcal{P}_N^{\pi^*|_{1:\tau}}) + \epsilon.$$

279 **Corollary 3.** *Set $\gamma := (\log |\Pi| + \log(8/\delta))/m$. Let $Z := \log |\Pi| + \log(8/\delta)$. There exists a*
 280 *choice of parameters $\eta = \theta$ (depending only on $m, n, \sigma_{\pi^*}/C_{\text{max}}, |\Pi|$, and δ) such that with*
 281 *probability at least $1 - \delta$: $\alpha_M(\pi_0, \tau_{\pi_0, \Phi}^\theta) = O\left(\max\left\{\frac{(\sigma_{\pi^*}/C_{\text{max}})^{4/5} Z}{m^{1/5}}, \frac{Z}{m^{1/3}}, \frac{\sigma_{\pi^*}}{C_{\text{max}}} \frac{Z}{n^{1/4}}, \frac{Z}{n^{1/2}}\right\}\right)$*
 282 *and $\text{Regret}_N^{\text{sw}}(\pi_0, \tau_{\pi_0, \Phi}^\theta; c) = O\left(\max\left\{\frac{(C_{\text{max}} \sigma_{\pi^*}^4)^{1/5} Z}{m^{1/5}}, \frac{C_{\text{max}} Z}{m^{1/3}}, \frac{\sigma_{\pi^*} Z}{n^{1/4}}, \frac{C_{\text{max}} Z}{n^{1/2}}\right\}\right)$.*

283 Corollary 3 has several parts. Both the stopping rate and regret bounds are a maximum over four
 284 terms. The first two terms are variance-driven and cost-driven contributions from the labeled sample,
 285 and the last two are their unlabeled-sample counterparts. In Section C.4, we show that the cost-driven
 286 labeled-data portion of this upper bound is tight.

287 *Switched regret vs. stopped regret.* As previously hinted, the switched regret is driven to zero at a
 288 faster rate than the stopped regret in the low $\sigma_{\pi^*}^2$ regime. This is perhaps unintuitive, as the switched
 289 regret bound provides a qualitatively strong guarantee: that the learner does not drive the system
 290 into an unrecoverable state before switching; on the other hand, the stopped regret guarantee only
 291 guarantees that the learner is competitive with the expert on the prefix for which it acts. Providing a
 292 resolution to this (for example, sharpening Corollary 2) is an interesting direction for future work.⁷

293 5 Extending to Misspecified Policy Classes

294 Previous sections assumed the expert π^* lies in the learner’s policy class Π . In practice, experts often
 295 employ richer representations, breaking this realizability assumption. Here, we show our validator
 296 framework degrades gracefully under misspecification. Throughout this section, let Π be a finite class
 297 of deterministic policies, and let π^* be a deterministic expert. We do not assume $\pi^* \in \Pi$.

⁷We suggest a partial resolution: define the *asymmetric stopped regret* as $J_N(\pi, \tau; c) - J_N(\pi^*; c)$ —physically, this compares the cost accumulated by the learner before handing over control to the cost the expert would have incurred driving the entire route. It is not hard to see that this is upper bounded by the switched regret, hence our algorithm controls the asymmetric stopped regret at the same rate as the switched regret (i.e., with the same rates as Corollary 3).

298 **Misspecification benchmark.** Let τ be as in (1). For any deterministic policy $\pi \in \Pi$, define
 299 $d_M(\pi) := \Pr_{M, \pi^*} [\tau_{\pi, \{\pi^*\}} \leq H]$ and $d_N(\pi) := \Pr_{N, \pi^*} [\tau_{\pi, \{\pi^*\}} \leq H]$. These are the probabilities
 300 that π deviates from the expert at any time during the trajectory (which has horizon H). We define the
 301 *policy-specific* and *best-in-class* misspecification, respectively, as $\Delta_\pi := \max(d_M(\pi), d_N(\pi))$ and
 302 $\Delta := \min_{\pi \in \Pi} \Delta_\pi$. Let $\tilde{\pi} = \arg \min_{\pi \in \Pi} \Delta_\pi$ be the policy which achieves the best-in-class misspecifi-
 303 cation. Finally, we define the empirical analogue of d_M as $\hat{d}_M(\pi) := \frac{1}{m} \sum_{i=1}^m \mathbf{1} [\tau_{\pi, \{\pi^*\}}(S_i) \leq H]$.

304 **Agnostic version of Algorithm 1.** When the expert lies outside the policy class ($\pi^* \notin \Pi$), the
 305 exact-consistency version space used in Algorithm 1 may fail to contain the expert. To ensure the
 306 validator set remains well-defined and sparse even under misspecification, we transition from a
 307 constrained optimization (forcing zero training disagreement) to a symmetrically regularized game.
 308 In this agnostic framework, we allow validators to disagree with the base policy on the source data,
 309 but we penalize this disagreement at a rate Λ . This allows the learner to "trade off" a small amount of
 310 source-side abstention for a significant reduction in target-side late-stop risk. This shift is captured by
 311 the following equilibrium result, which can be thought of as an "agnostic" analog of Lemma 1:

312 **Lemma 3.** *Let $\pi_0 \in \arg \min_{\pi \in \Pi} \hat{d}_M(\pi)$ be an empirical disagreement minimizer. For every*
 313 *penalty parameter $\Lambda > 0$ and integer $K \geq 1$, there exists a distribution q^* over validator sets*
 314 *$\Phi = (\phi_1, \dots, \phi_K) \in \Pi^K$ that simultaneously satisfies the following two properties:*

- 315 • $\mathbb{E}_{\Phi \sim q^*} \left[\sum_{k=1}^K \hat{d}_M(\phi_k) \right] \leq K \cdot \hat{d}_M(\pi_0) + \frac{1}{\Lambda},$
- 316 • $\mathbb{E}_{\Phi \sim q^*} \left[\frac{1}{n} \sum_{j=1}^n \mathbf{1} [\tau_{\pi_0, \Phi}(T_j) > \tau_{\pi_0, \{\pi^*\}}(T_j)] \right] \leq \Lambda \left(\hat{d}_M(\pi) - \hat{d}_M(\pi_0) \right) + \frac{1}{K}.$

317 Similar to Section 3, Lemma 3 is a non-constructive result, although it may be realized by no-regret
 318 play. We package the output Lemma 3 into an abstract subroutine, REGULARIZEDSPARSEVALIDA-
 319 TORDIST and leave the question of its efficient computational implementation to future work.

320 **Theorem 3.** *Let $\delta > 0$. Let $\pi_0 \in \arg \min_{\pi \in \Pi} \hat{d}_M(\pi)$. Fix a penalty $\Lambda > 0$ and committee size*
 321 *$K \geq 1$, and let $q^* \in \Delta(\Pi^K)$ be the equilibrium distribution from Lemma 3. Define the complexity*
 322 *measure $Z := (K + 1) \log |\Pi| + \log \frac{4}{\delta}$. Then, with probability at least $1 - \delta$, the randomized selective*
 323 *policy $(\pi_0, \tau_{\pi_0, \Phi})$ with validator set drawn $\Phi \sim q^*$ satisfies*

- 324 • $\mathbb{E}_{\Phi \sim q^*} [\alpha_M(\pi_0, \tau_{\pi_0, \Phi})] \leq (K + 1)\Delta + \frac{1}{\Lambda} + 2(K + 1)\sqrt{\frac{Z}{2m}},$
- 325 • $\mathbb{E}_{\Phi \sim q^*} [\text{Regret}_N(\pi_0, \tau_{\pi_0, \Phi}; c)] \leq C_{\max} \left(\Delta + \Lambda\Delta + \frac{1}{K} + 2\Lambda\sqrt{\frac{Z}{2m}} + \sqrt{\frac{Z}{2n}} \right).$

326 **Corollary 4.** *There exists a choice of parameters $K = \Lambda$ (depending on $m, n, \Delta, |\Pi|$,*
 327 *and δ) such that, with probability at least $1 - \delta$, both $\mathbb{E}_{\Phi \sim q^*} [\alpha_M(\pi_0, \tau_{\pi_0, \Phi})]$ and⁸
 328 $\mathbb{E}_{\Phi \sim q^*} [\text{Regret}_N(\pi_0, \tau_{\pi_0, \Phi}; c)]/C_{\max}$ are $\tilde{O} \left(\Delta^{1/2} + (\log |\Pi|/m)^{1/5} + (\log |\Pi|/n)^{1/3} \right).$*

329 **Extension to off-policy test trajectories.** While we have assumed a single expert across both
 330 environments, this regularized framework naturally extends to settings where the training demonstrator
 331 π^{tr} and test demonstrator π^{te} differ, and neither necessarily lies in Π . By generalizing the disagreement
 332 probabilities, we show that our algorithm degrades gracefully with the irreducible off-policy error
 333 Δ_{off} . We defer the formal setup and the proof of this guarantee to Section D.4.

334 6 Proof of Concept

335 We evaluate SeqRejectron in the LunarLander-v3 environment [Towers et al., 2024], with the goal of
 336 empirically characterizing the completeness-soundness tradeoff: the algorithm should stop minimally
 337 on train while reliably detecting out-of-distribution dynamics and stopping at test time.

338 **Experimental setup.** The step cost c_h is a quadratic penalty on distance to the landing pad, velocity
 339 and angular velocity with a cap at 1. If the lander ever crashes, it incurs the maximum step-wise
 340 cost for the remainder of the episode. Costs are then normalized to a maximum episode cost of 1.

⁸In contrast to Algorithm 1, the regret guarantee of the misspecified SeqRejectron is stated in expectation over the draw $\Phi \sim q^*$, and thus the algorithm must return the validator distribution q^* . This is because in the agnostic setting, taking the union of k validator sets (which we did to boost the success probability of Algorithm 1) would linearly amplify the irreducible misspecification error (scaling the penalty to $O(k\Delta)$).

341 For the dynamics shift, the source domain M is windless and uses a reduced initial random impulse,
 342 while the target domain N applies a rightward wind force. Thus the environment shift consists of
 343 the addition of wind and the stronger initial impulse. The expert π^* is a stochastic state-feedback
 344 controller, parameterized by a neural network on the raw LunarLander observation features, tuned to
 345 land in the target environment, but not necessarily optimal under the custom cost. The base learner
 346 π_0 is fit by MLE on labeled demonstrations $\mathcal{D}_{\text{train}}$ drawn only from M .

347 To ensure tractability, we replace the game-theoretic construction of Section B.6 with an adversarial
 348 posterior-sampling heuristic. First, using the labeled training demonstrations $\mathcal{D}_{\text{train}}$, we form a
 349 Bayesian posterior over the policy parameters and sample a pool of candidates approximately
 350 consistent with π^* on M . Next, we evaluate this pool on the unlabeled test trajectories $\mathcal{D}_{\text{test}}$ and
 351 greedily select the candidate that maximizes the cumulative squared Hellinger disagreement with
 352 the base learner π_0 (since our goal is to challenge our base learner on the test data). Repeating this
 353 selection $K = 3$ times yields a validator set Φ , which defines a stopping rule as in (2).

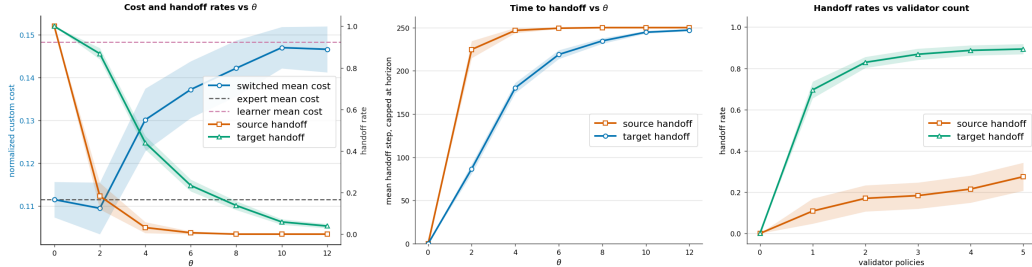


Figure 1: Left: normalized target switched cost and source/target handoff rates as functions of the safety threshold θ . Dashed horizontal lines show the expert and learner mean costs. Middle: mean handoff time as a function of θ , with trajectories that never hand off capped at the episode horizon. Right: source and target handoff rates as a function of the number of validator policies, with $\theta = 2$ fixed. In all plots, shaded bands indicate standard error across trials.

354 **Results.** We evaluate the framework over 20 independent trials. In each trial, we collect $m = 30$
 355 labeled training trajectories and $n = 30$ unlabeled test trajectories over a range of θ threshold values.
 356 Figure 1 summarizes the resulting tradeoffs. In the left panel, the source handoff rate collapses
 357 steeply to near 0%, confirming that the validators rarely flag in-distribution trajectories. The target
 358 handoff rate decays far more gradually with the trajectory cost mirroring its behavior. As handoff
 359 becomes less frequent, the switched policy’s cost approaches the learner’s cost. The dashed black
 360 line represents the expert’s average cost on the test environment. The red dashed line represents
 361 the expected cost of π_0 , and is the cost that would have been obtained by running vanilla behavior
 362 cloning, without any abstention mechanism. This shows that SeqRejectron (with a suitable choice of
 363 $\theta = 2$) is able to reliably detect and abstain when it encounters environment shift.

364 The middle panel shows that larger θ delays handoff, as expected, since the validator allows more
 365 cumulative disagreement before switching. Finally, the right panel shows that increasing the number
 366 of validators raises handoff rates, particularly in the target environment, by making the stopping
 367 rule more sensitive to disagreement. This effect quickly saturates, suggesting that SeqRejectron is
 368 somewhat robust to the choice of validator set size.

369 7 Conclusion

370 We introduced selective imitation learning, a framework where a learner abstains mid-trajectory
 371 when training data is uninformative about the test environment. Our algorithm, SeqRejectron,
 372 yields horizon-free sample complexity for deterministic classes under sparse costs, with analogous
 373 guarantees for stochastic and misspecified classes. We also established an $\Omega(1/\epsilon^3)$ lower bound that
 374 matches the labeled-sample cost-driven exponent. Several directions remain open. Computationally,
 375 developing practical surrogates for the cumulative Hellinger objective in stochastic settings is needed.
 376 Statistically, closing the gap between the variance-driven $\tilde{O}(\epsilon^{-5})$ rate and our lower bound remains
 377 open. Finally, incorporating partial structural knowledge about test dynamics (e.g., bounded TV
 378 distance, density ratio, or parametric shift) could tighten test-side stopping rates. Ultimately, this
 379 work establishes abstention as a principled response to unanticipated dynamics shift.

380 References

- 381 Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron,
382 Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, et al. Solving rubik’s cube with a
383 robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
- 384 Aaron D Ames, Samuel Coogan, Magnus Egerstedt, Gennaro Notomista, Koushil Sreenath, and
385 Paulo Tabuada. Control barrier functions: Theory and applications. In *2019 18th European control
386 conference (ECC)*, pages 3420–3431. Ieee, 2019.
- 387 OpenAI: Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew,
388 Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning
389 dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20,
390 2020.
- 391 Anastasios N Angelopoulos and Stephen Bates. Conformal prediction: A gentle introduction.
392 *Foundations and Trends in Machine Learning*, 16(4):494–591, 2023.
- 393 Somil Bansal, Mo Chen, Sylvia Herbert, and Claire J Tomlin. Hamilton-jacobi reachability: A brief
394 overview and recent advances. In *2017 IEEE 56th annual conference on decision and control
395 (CDC)*, pages 2242–2253. IEEE, 2017.
- 396 Shai Ben-David and Ruth Urner. On the hardness of domain adaptation and the utility of unlabeled
397 target samples. In *International Conference on Algorithmic Learning Theory*, pages 139–153.
398 Springer, 2012.
- 399 Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman
400 Vaughan. A theory of learning from different domains. *Machine learning*, 79(1):151–175, 2010.
- 401 Jongseong Chae, Seungyul Han, Whiyoung Jung, Myungsik Cho, Sungho Choi, and Youngchul Sung.
402 Robust imitation learning against variations in environment dynamics. In *International Conference
403 on Machine Learning*, pages 2828–2852. PMLR, 2022.
- 404 Gautam Chandrasekaran, Adam R Klivans, Vasilis Kontonis, Konstantinos Stavropoulos, and Arsen
405 Vasilyan. Efficient discrepancy testing for learning with distribution shift. *NeurIPS*, 2024.
- 406 Yevgen Chebotar, Ankur Handa, Viktor Makoviychuk, Miles Macklin, Jan Issac, Nathan Ratliff,
407 and Dieter Fox. Closing the sim-to-real loop: Adapting simulation randomization with real
408 world experience. In *2019 international conference on robotics and automation (ICRA)*, pages
409 8973–8979. IEEE, 2019.
- 410 Xiaoyu Chen, Jiachen Hu, Chi Jin, Lihong Li, and Liwei Wang. Understanding domain randomization
411 for sim-to-real transfer. *arXiv preprint arXiv:2110.03239*, 2021.
- 412 C Chow. On optimum recognition error and reject tradeoff. *IEEE Transactions on information theory*,
413 16(1):41–46, 2003.
- 414 Felipe Codevilla, Eder Santana, Antonio M López, and Adrien Gaidon. Exploring the limitations
415 of behavior cloning for autonomous driving. In *Proceedings of the IEEE/CVF international
416 conference on computer vision*, pages 9329–9338, 2019.
- 417 Siddharth Desai, Ishan Durugkar, Haresh Karnan, Garrett Warnell, Josiah Hanna, and Peter Stone.
418 An imitation from observation approach to transfer learning with dynamics mismatch. *Advances
419 in Neural Information Processing Systems*, 33:3917–3929, 2020.
- 420 Thomas G Dietterich, Richard H Lathrop, and Tomás Lozano-Pérez. Solving the multiple instance
421 problem with axis-parallel rectangles. *Artificial intelligence*, 89(1-2):31–71, 1997.
- 422 Anushri Dixit, Lars Lindemann, Skylar X Wei, Matthew Cleaveland, George J Pappas, and Joel W
423 Burdick. Adaptive conformal prediction for motion planning among dynamic agents. In *Learning
424 for Dynamics and Control Conference*, pages 300–314. PMLR, 2023.
- 425 Miroslav Dudík, Nika Haghtalab, Haipeng Luo, Robert E Schapire, Vasilis Syrgkanis, and Jen-
426 nifer Wortman Vaughan. Oracle-efficient online learning and auction design. *Journal of the ACM
427 (JACM)*, 67(5):1–57, 2020.

- 428 Ran El-Yaniv et al. On the foundations of noise-free selective classification. *Journal of Machine*
429 *Learning Research*, 11(5), 2010.
- 430 Nicolas Espinosa-Dice, Sanjiban Choudhury, Wen Sun, and Gokul Swamy. Efficient imitation under
431 misspecification. *arXiv preprint arXiv:2503.13162*, 2025.
- 432 Benjamin Eysenbach, Swapnil Asawa, Shreyas Chaudhari, Sergey Levine, and Ruslan Salakhutdinov.
433 Off-dynamics reinforcement learning: Training for transfer with domain classifiers. *arXiv preprint*
434 *arXiv:2006.13916*, 2020.
- 435 Arnaud Fickinger, Abderrahim Bendahi, and Stuart Russell. Statistical guarantees for offline domain
436 randomization. In *The Fourteenth International Conference on Learning Representations*, 2026.
- 437 Samuel G Finlayson, Adarsh Subbaswamy, Karandeep Singh, John Bowers, Annabel Kupke, Jonathan
438 Zittrain, Isaac S Kohane, and Suchi Saria. The clinician and dataset shift in artificial intelligence.
439 *New England Journal of Medicine*, 385(3):283–286, 2021.
- 440 Dylan J Foster, Adam Block, and Dipendra Misra. Is behavior cloning all you need? understanding
441 horizon in imitation learning. *Advances in Neural Information Processing Systems*, 37:120602–
442 120666, 2024.
- 443 Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse reinforce-
444 ment learning. *arXiv preprint arXiv:1710.11248*, 2017.
- 445 Tesshu Fujinami, Bruce D Lee, Nikolai Matni, and George J Pappas. Domain randomization is
446 sample efficient for linear quadratic control. *arXiv preprint arXiv:2502.12310*, 2025.
- 447 Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model
448 uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059.
449 PMLR, 2016.
- 450 Tanmay Gangwani and Jian Peng. State-only imitation with transition dynamics mismatch. *arXiv*
451 *preprint arXiv:2002.11879*, 2020.
- 452 Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. *Advances in*
453 *neural information processing systems*, 30, 2017.
- 454 Isaac Gibbs and Emmanuel Candes. Adaptive conformal inference under distribution shift. *Advances*
455 *in Neural Information Processing Systems*, 34:1660–1672, 2021.
- 456 Surbhi Goel, Abhishek Shetty, Konstantinos Stavropoulos, and Arsen Vasilyan. Tolerant algorithms
457 for learning with arbitrary covariate shift. *Advances in Neural Information Processing Systems*, 37:
458 124979–125018, 2024.
- 459 Shafi Goldwasser, Adam Tauman Kalai, Yael Kalai, and Omar Montasser. Beyond perturbations:
460 Learning guarantees with arbitrary adversarial test examples. *Advances in Neural Information*
461 *Processing Systems*, 33:15859–15870, 2020.
- 462 Tom Haider, Karsten Roscher, Felipe Schmoeller da Roza, and Stephan Günnemann. Out-of-
463 distribution detection for reinforcement learning agents with probabilistic dynamics models. In
464 *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*,
465 pages 851–859, 2023.
- 466 Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural*
467 *information processing systems*, 29, 2016.
- 468 Garud N Iyengar. Robust dynamic programming. *Mathematics of Operations Research*, 30(2):
469 257–280, 2005.
- 470 Stephen James, Paul Wohlhart, Mrinal Kalakrishnan, Dmitry Kalashnikov, Alex Irpan, Julian Ibarz,
471 Sergey Levine, Raia Hadsell, and Konstantinos Bousmalis. Sim-to-real via sim-to-sim: Data-
472 efficient robotic grasping via randomized-to-canonical adaptation networks. In *Proceedings of the*
473 *IEEE/CVF conference on computer vision and pattern recognition*, pages 12627–12637, 2019.

- 474 Gregory Kahn, Adam Villaflor, Vitchyr Pong, Pieter Abbeel, and Sergey Levine. Uncertainty-aware
475 reinforcement learning for collision avoidance. *arXiv preprint arXiv:1702.01182*, 2017.
- 476 Adam Tauman Kalai and Varun Kanade. Efficient learning with arbitrary covariate shift. In *Algorith-*
477 *mic Learning Theory*, pages 850–864. PMLR, 2021.
- 478 Adam R Klivans, Konstantinos Stavropoulos, and Arsen Vasilyan. Testable learning with distribution
479 shift. *37th Annual Conference on Learning Theory, COLT 2024 (to appear)*, 2024a.
- 480 Adam R Klivans, Konstantinos Stavropoulos, and Arsen Vasilyan. Learning intersections of halfspaces
481 with distribution shift: Improved algorithms and sq lower bounds. *37th Annual Conference on*
482 *Learning Theory, COLT 2024 (to appear)*, 2024b.
- 483 Matthieu Komorowski, Leo A Celi, Omar Badawi, Anthony C Gordon, and A Aldo Faisal. The
484 artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care.
485 *Nature medicine*, 24(11):1716–1720, 2018.
- 486 Ilya Kostrikov, Ofir Nachum, and Jonathan Tompson. Imitation learning via off-policy distribution
487 matching. *arXiv preprint arXiv:1912.05032*, 2019.
- 488 Ashish Kumar, Zipeng Fu, Deepak Pathak, and Jitendra Malik. Rma: Rapid motor adaptation for
489 legged robots. *arXiv preprint arXiv:2107.04034*, 2021.
- 490 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline
491 reinforcement learning. *Advances in neural information processing systems*, 33:1179–1191, 2020.
- 492 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive
493 uncertainty estimation using deep ensembles. *Advances in neural information processing systems*,
494 30, 2017.
- 495 Niklas Lauffer, Xiang Deng, Srivatsa Kundurthy, Brad Kenstler, and Jeff Da. Imitation learning for
496 multi-turn lm agents via on-policy expert corrections. *arXiv preprint arXiv:2512.14895*, 2025.
- 497 Brian Lee and Nikolai Matni. Single trajectory conformal prediction. In *2024 IEEE 63rd Conference*
498 *on Decision and Control (CDC)*, pages 3019–3024. IEEE, 2024.
- 499 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,
500 review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 501 Lars Lindemann, Matthew Cleaveland, Gihyun Shim, and George J Pappas. Safe planning in dynamic
502 environments using conformal prediction. *IEEE Robotics and Automation Letters*, 8(8):5116–5123,
503 2023a.
- 504 Lars Lindemann, Xin Qin, Jyotirmoy V Deshmukh, and George J Pappas. Conformal prediction
505 for stl runtime verification. In *Proceedings of the ACM/IEEE 14th International Conference on*
506 *Cyber-Physical Systems (with CPS-IoT Week 2023)*, pages 142–153, 2023b.
- 507 Björn Lütjens, Michael Everett, and Jonathan P How. Safe reinforcement learning with model
508 uncertainty estimates. In *2019 International Conference on Robotics and Automation (ICRA)*,
509 pages 8662–8668. IEEE, 2019.
- 510 Yishay Mansour, Mehryar Mohri, and Afshin Rostamizadeh. Domain adaptation: Learning bounds
511 and algorithms. *arXiv preprint arXiv:0902.3430*, 2009.
- 512 Oded Maron and Tomás Lozano-Pérez. A framework for multiple-instance learning. *Advances in*
513 *neural information processing systems*, 10, 1997.
- 514 Bhairav Mehta, Manfred Diaz, Florian Golemo, Christopher J Pal, and Liam Paull. Active domain
515 randomization. In *Conference on Robot Learning*, pages 1162–1176. PMLR, 2020.
- 516 Natinael Solomon Neggatu, Jeremie Houssineau, and Giovanni Montana. Evaluation-time policy
517 switching for offline reinforcement learning. *arXiv preprint arXiv:2503.12222*, 2025.
- 518 Arnab Nilim and Laurent El Ghaoui. Robust control of markov decision processes with uncertain
519 transition matrices. *Operations Research*, 53(5):780–798, 2005.

- 520 Kishan Panaganti and Dileep Kalathil. Sample complexity of robust reinforcement learning with
521 a generative model. In *International Conference on Artificial Intelligence and Statistics*, pages
522 9582–9602. PMLR, 2022.
- 523 Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of
524 robotic control with dynamics randomization. In *2018 IEEE international conference on robotics
525 and automation (ICRA)*, pages 3803–3810. IEEE, 2018.
- 526 Lerrel Pinto, James Davidson, Rahul Sukthankar, and Abhinav Gupta. Robust adversarial rein-
527 forcement learning. In *International conference on machine learning*, pages 2817–2826. PMLR,
528 2017.
- 529 Dean A Pomerleau. Alvin: An autonomous land vehicle in a neural network. *Advances in neural
530 information processing systems*, 1, 1988.
- 531 Stephen Prajna and Ali Jadbabaie. Safety verification of hybrid systems using barrier certificates. In
532 *International workshop on hybrid systems: Computation and control*, pages 477–492. Springer,
533 2004.
- 534 Mohit Prashant, Arvind Easwaran, Suman Das, and Michael Yuhas. Guaranteeing out-of-distribution
535 detection in deep rl via transition estimation. In *Proceedings of the AAAI Conference on Artificial
536 Intelligence*, 2025.
- 537 Nived Rajaraman, Lin Yang, Jiantao Jiao, and Kannan Ramchandran. Toward the fundamental limits
538 of imitation learning. *Advances in Neural Information Processing Systems*, 33:2914–2924, 2020.
- 539 Aravind Rajeswaran, Sarvjeet Ghotra, Balaraman Ravindran, and Sergey Levine. Epopt: Learning
540 robust neural network policies using model ensembles. *arXiv preprint arXiv:1610.01283*, 2016.
- 541 Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline rein-
542 forcement learning and imitation learning: A tale of pessimism. *Advances in Neural Information
543 Processing Systems*, 34:11702–11716, 2021.
- 544 Sahana Rayan and Ambuj Tewari. Learning to partially defer for sequences, 2025. URL <https://arxiv.org/abs/2502.01459>.
545
- 546 Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. In *Proceedings of the
547 thirteenth international conference on artificial intelligence and statistics*, pages 661–668. JMLR
548 Workshop and Conference Proceedings, 2010.
- 549 Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured
550 prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on
551 artificial intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Proceedings,
552 2011.
- 553 Andrei A Rusu, Matej Večerík, Thomas Rothörl, Nicolas Heess, Razvan Pascanu, and Raia Hadsell.
554 Sim-to-real robot learning from pixels with progressive nets. In *Conference on robot learning*,
555 pages 262–270. PMLR, 2017.
- 556 Fereshteh Sadeghi and Sergey Levine. Cad2rl: Real single-image flight without a single real image.
557 *arXiv preprint arXiv:1611.04201*, 2016.
- 558 Laixi Shi and Yuejie Chi. Distributionally robust model-based offline reinforcement learning with
559 near-optimal sample complexity. *Journal of Machine Learning Research*, 25(200):1–91, 2024.
- 560 Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-
561 likelihood function. *Journal of statistical planning and inference*, 90(2):227–244, 2000.
- 562 Masashi Sugiyama, Matthias Krauledat, and Klaus-Robert Müller. Covariate shift adaptation by
563 importance weighted cross validation. *Journal of Machine Learning Research*, 8(5), 2007.
- 564 Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1. MIT
565 press Cambridge, 1998.

- 566 Gokul Swamy, Sanjiban Choudhury, J Andrew Bagnell, and Steven Wu. Of moments and matching:
567 A game-theoretic framework for closing the imitation gap. In *International Conference on Machine*
568 *Learning*, pages 10022–10032. PMLR, 2021.
- 569 Gokul Swamy, Nived Rajaraman, Matt Peng, Sanjiban Choudhury, J Bagnell, Steven Z Wu, Jiantao
570 Jiao, and Kannan Ramchandran. Minimax optimal online imitation learning via replay estimation.
571 *Advances in Neural Information Processing Systems*, 35:7077–7088, 2022.
- 572 Jie Tan, Tingnan Zhang, Erwin Coumans, Atil Iscen, Yunfei Bai, Danijar Hafner, Steven Bohez, and
573 Vincent Vanhoucke. Sim-to-real: Learning agile locomotion for quadruped robots. *arXiv preprint*
574 *arXiv:1804.10332*, 2018.
- 575 Chen Tessler, Yonathan Efroni, and Shie Mannor. Action robust reinforcement learning and applica-
576 tions in continuous control. In *International Conference on Machine Learning*, pages 6215–6224.
577 PMLR, 2019.
- 578 Ryan J Tibshirani, Rina Foygel Barber, Emmanuel Candes, and Aaditya Ramdas. Conformal
579 prediction under covariate shift. *Advances in neural information processing systems*, 32, 2019.
- 580 Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain
581 randomization for transferring deep neural networks from simulation to the real world. In *2017*
582 *IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 23–30. IEEE,
583 2017.
- 584 Claire J Tomlin, John Lygeros, and S Shankar Sastry. A game theoretic approach to controller design
585 for hybrid systems. *Proceedings of the IEEE*, 88(7):949–970, 2000.
- 586 Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. *arXiv preprint*
587 *arXiv:1805.01954*, 2018.
- 588 Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U Balis, Gianluca De Cola, Tristan Deleu,
589 Manuel Goulão, Andreas Kallinteris, Markus Krimmel, Arjun KG, et al. Gymnasium: A standard
590 interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*, 2024.
- 591 Luca Viano, Yu-Ting Huang, Parameswaran Kamalaruban, Adrian Weller, and Volkan Cevher.
592 Robust inverse reinforcement learning under transition dynamics mismatch. *Advances in Neural*
593 *Information Processing Systems*, 34:25917–25931, 2021.
- 594 Wolfram Wiesemann, Daniel Kuhn, and Berç Rustem. Robust markov decision processes. *Mathe-*
595 *matics of Operations Research*, 38(1):153–183, 2013.
- 596 Huazhe Xu, Yang Gao, Fisher Yu, and Trevor Darrell. End-to-end learning of driving models from
597 large-scale video datasets. In *Proceedings of the IEEE conference on computer vision and pattern*
598 *recognition*, pages 2174–2182, 2017.
- 599 Wenhao Yu, Jie Tan, C Karen Liu, and Greg Turk. Preparing for the unknown: Learning a universal
600 policy with online system identification. *arXiv preprint arXiv:1702.02453*, 2017.
- 601 Huan Zhang, Hongge Chen, Chaowei Xiao, Bo Li, Mingyan Liu, Duane Boning, and Cho-Jui
602 Hsieh. Robust deep reinforcement learning against adversarial perturbations on state observations.
603 *Advances in neural information processing systems*, 33:21024–21037, 2020.
- 604 Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy inverse
605 reinforcement learning. In *Aaai*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.

606	Contents	
607	1 Introduction	1
608	1.1 Our approach and results	2
609	1.2 Related work	2
610	2 Problem Formulation	3
611	3 Selective Execution for Deterministic Policies	4
612	4 Selective Execution for Stochastic Policies	6
613	4.1 Expert handoff interpretation and switched regret guarantees	7
614	5 Extending to Misspecified Policy Classes	7
615	6 Proof of Concept	8
616	7 Conclusion	9
617	A Additional Related Work	17
618	A.1 Background on imitation learning and reinforcement learning	17
619	A.2 Background on selective classification and learning with distribution shift	17
620	A.3 Robustness, uncertainty quantification, and domain adaptation in RL	18
621	B Additional Results and Proofs from the Deterministic Section	20
622	B.1 Supporting technical lemmas	20
623	B.2 Analysis of stepwise reduction to PQ-learning	21
624	B.3 Proof of Lemma 1	22
625	B.4 Proof of Theorem 1	22
626	B.5 Transferring the completeness guarantee to the test environment	24
627	B.6 Constructing a validator distribution using no-regret dynamics	25
628	B.7 Proof of Proposition 1	26
629	B.8 Oracle-efficient constructions of a validator distribution	26
630	C Proofs from the Stochastic Section	29
631	C.1 Proof of Theorem 2	29
632	C.2 Supporting Hellinger lemmas	30
633	C.3 Proof of Lemma 2	33
634	C.4 Single-step lower bound for stochastic policies	34
635	C.5 Proof of Theorem 4	34
636	D Proofs from the Misspecification Section	35
637	D.1 Proof of Lemma 3	36
638	D.2 Proof of Theorem 3	36

639	D.3 Proof of Corollary 4	38
640	D.4 Off-policy test trajectories under misspecification	38
641	D.5 Proof of Proposition 4	39

642 A Additional Related Work

643 A.1 Background on imitation learning and reinforcement learning

644 **Reinforcement learning.** The standard framework for sequential decision-making is the Markov
645 decision process (MDP) [Sutton et al., 1998], where an agent maximizes cumulative reward by
646 interacting with an environment governed by unknown transition dynamics. In the offline setting,
647 where the agent must learn from a fixed dataset without further environment interaction, pessimism-
648 based approaches [Rashidinejad et al., 2021, Kumar et al., 2020, Levine et al., 2020] provide
649 guarantees when the dataset covers the target policy. Our work operates in a similarly offline
650 regime, but with a key difference: we observe datasets from *two* environments and provide PAC-style
651 guarantees without any online interaction with the test environment.

652 **Imitation learning.** The dominant paradigm for learning from demonstrations is *behavior cloning*
653 (BC), which reduces imitation learning to supervised prediction [Pomerleau, 1988, Ross and Bagnell,
654 2010]. A central limitation of BC is *compounding error*: small per-step mistakes cause the learner to
655 drift to states outside the expert’s distribution, where further errors compound. Interactive methods
656 address this by querying the expert on the learner’s own rollouts: DAGger [Ross et al., 2011] reduces
657 regret to $O(\epsilon H)$ by collecting on-policy corrections, and subsequent work gave game-theoretic
658 formulations [Swamy et al., 2021] and adversarial approaches via occupancy-measure matching [Ho
659 and Ermon, 2016, Ziebart et al., 2008]. Imitation from observation [Torabi et al., 2018] studies the
660 special case where only state sequences (not actions) are available, which is similar to our assumption
661 of unlabeled test trajectories, but does not provide theoretical guarantees.

662 **Sample complexity of imitation learning.** A precise theoretical understanding of BC’s sample
663 complexity has emerged over the past several years. Rajaraman et al. [2020] established fundamental
664 limits: offline BC requires $\Omega(H^2/\epsilon)$ trajectories in the worst case, while online methods can achieve
665 $O(H/\epsilon)$, and Swamy et al. [2022] showed that replay estimation achieves minimax-optimal online IL
666 by reducing empirical variance in the expert visitation distribution. This apparent gap between offline
667 and online IL was recently revisited by Foster et al. [2024], who showed that under weight-sharing,
668 log-loss BC achieves *horizon-independent* sample complexity via a trajectory-level Hellinger analysis,
669 whenever cumulative costs are controlled and the policy class shares parameters across time steps;
670 this matches the minimax rate even among interactive algorithms. Our work builds directly on this
671 Hellinger-based analysis, extending it to the environment-shift setting where test dynamics may
672 differ arbitrarily from training, and replacing the standard regret objective with a stopped regret that
673 accounts for the possibility of abstention.

674 A.2 Background on selective classification and learning with distribution shift

675 **Learning under distribution shift.** The classical approach to distribution shift is *covariate shift*
676 *adaptation*, where the marginal input distribution changes but the labeling rule is preserved. Shi-
677 modaira [2000] introduced importance-weighting as a principled correction, and Sugiyama et al.
678 [2007] gave an importance-weighted cross-validation procedure for model selection under this shift.
679 On the theoretical side, Ben-David et al. [2010] and Mansour et al. [2009] bound target error in terms
680 of source error and a divergence between domains. However, Ben-David and Urner [2012] establish
681 that without further assumptions on the relationship between source and target—such as bounded
682 density ratios or overlapping supports—successful domain adaptation is impossible in general. These
683 hardness results make clear that any distribution-free guarantee under arbitrary shift must either make
684 structural assumptions or give up on committing to an action everywhere.

685 **Selective classification and PQ learning.** *Selective classification* [Chow, 2003, El-Yaniv et al., 2010,
686 Geifman and El-Yaniv, 2017] offers a way out: by allowing a classifier to abstain on inputs it cannot
687 reliably label, one can trade coverage for accuracy. Recent work on partial deferral for sequence
688 prediction studies one-time handoff to an expert after a chosen point in the sequence, which is closely
689 related in spirit to our expert-handoff view of abstention [Rayan and Tewari, 2025]; however, they
690 do not consider environment shift, which is the central focus of this work. Closest to our formal
691 setup, Goldwasser et al. [2020] introduced the *PQ learning* framework, showing that abstention is
692 precisely the mechanism that circumvents the hardness of arbitrary distribution shift. Given labeled
693 samples from P and unlabeled samples from Q , PQ learning asks for a selective classifier with
694 low error on Q and low rejection rate on P —with no assumptions on the relationship between
695 P and Q . Their algorithm, Rejectron, iteratively builds validator classifiers that agree on P but

696 disagree on Q , and abstains wherever such disagreement exists. Subsequent work [Kalai and Kanade,
697 2021] showed that the computational complexity lies in between PAC learning and Agnostic PAC
698 learning. Another closely related model is the Testable Distribution Shift (TDS) Klivans et al. [2024a]
699 which allows for rejecting the entire test distribution instead of selectively predicting. There has
700 been a flurry of recent work designing computationally efficient algorithms for these models under
701 assumptions only on the train distribution [Klivans et al., 2024b, Chandrasekaran et al., 2024, Goel
702 et al., 2024]. Our formulation can be viewed as a sequential analogue of PQ learning (Section 2), and
703 the validator-based construction of SeqRejectron can be thought of as a generalization of Rejectron
704 to trajectory-level certification.

705 A.3 Robustness, uncertainty quantification, and domain adaptation in RL

706 **Sim-to-real transfer and domain randomization.** The *reality gap*—the mismatch between simu-
707 lated and real-world appearance, physics, and dynamics—is a central obstacle to deploying learned
708 policies on physical systems [Tobin et al., 2017].

709 The dominant mitigation is *domain randomization*: training over a broad distribution of simulated
710 environments so the real world falls within the training support. This has been applied to visual
711 parameters [Tobin et al., 2017, Sadeghi and Levine, 2016, James et al., 2019], physical parameters
712 such as friction, mass, and actuator delays [Peng et al., 2018, Tan et al., 2018], and scaled to dexterous
713 manipulation and Rubik’s-cube solving [Andrychowicz et al., 2020, Akkaya et al., 2019]. Comple-
714 menting these empirical milestones, recent literature has formalized the theoretical foundations of this
715 approach. Specifically, these works have bounded the sim-to-real gap for uniform domain randomiza-
716 tion [Chen et al., 2021], established the method’s optimal asymptotic sample efficiency [Fujinami
717 et al., 2025] for linear systems, and proved the statistical consistency of data-driven, offline domain
718 randomization [Fickinger et al., 2026]

719 A complementary line of work closes the gap *adaptively*: Chebotar et al. [2019] iteratively refine the
720 randomization distribution using real rollouts; Yu et al. [2017] and Kumar et al. [2021] condition
721 policies on online system-identification modules that infer dynamics at deployment; Rusu et al. [2017]
722 use progressive networks for visual transfer; and Mehta et al. [2020] steer randomization toward
723 maximally informative environments. Codevilla et al. [2019] document the resulting failure modes
724 of behavior cloning under sim-to-real shift, motivating principled responses to out-of-distribution
725 conditions.

726 All of these methods attempt to *bridge* the reality gap. Our work is complementary: rather than
727 correcting for an unknown shift, we learn to *detect* when dynamics have changed too much for
728 reliable execution and abstain, with distribution-free PAC guarantees that require no assumptions on
729 the magnitude or structure of the mismatch.

730 **Robust reinforcement learning.** Robust RL seeks policies that perform well under worst-case
731 dynamics drawn from a prespecified uncertainty set. The foundational framework is the robust MDP,
732 in which nature adversarially selects transitions within a rectangular uncertainty set; Iyengar [2005]
733 and Nilim and El Ghaoui [2005] independently showed that such problems admit tractable dynamic-
734 programming solutions, and Wiesemann et al. [2013] extended this to more general uncertainty
735 sets and ambiguity models. In the deep RL setting, Pinto et al. [2017] train a protagonist against
736 an adversary that applies destabilizing forces (RARL), and Rajeswaran et al. [2016] optimize over
737 an ensemble of source models (EPOpt) to guard against model error. Zhang et al. [2020] study
738 robustness to observation and action perturbations in deep RL, and Tessler et al. [2019] isolate
739 the action-perturbation setting. On the theoretical side, Panaganti and Kalathil [2022] and Shi and
740 Chi [2024] establish sample complexity bounds for learning robust policies from offline data under
741 (s, a) -rectangular and distributionally robust uncertainty sets, respectively. All of these methods
742 require the learner to commit ahead of time to an uncertainty set or dynamics prior—implicitly
743 assuming knowledge of the shift’s structure or magnitude.

744 **Imitation learning under dynamics mismatch.** Classical IL theory assumes that the learner and
745 expert share the same transition dynamics. A growing body of work relaxes this assumption.

746 Several approaches use adversarial methods to match occupancy measures across environments:
747 Ho and Ermon [2016] introduce GAIL, which frames IL as a distribution-matching problem and
748 provides a foundation for dynamics-robust extensions; Fu et al. [2017] show that adversarial IRL
749 (AIRL) can recover reward functions that transfer under moderate dynamics shift; and Kostrikov et al.

750 [2019] propose a distribution-correction method (DAC) that compensates for off-policy mismatch via
751 density-ratio reweighting.

752 Other works directly target the setting where expert and learner operate under different dynamics.
753 Gangwani and Peng [2020] and Desai et al. [2020] use adversarial state-distribution matching and
754 transition-dynamics alignment respectively. Viano et al. [2021] bound performance degradation as a
755 function of the ℓ_1 distance between transition kernels, and Chae et al. [2022] train policies that are
756 robust across a family of perturbed MDPs (RIME). Eysenbach et al. [2020] compensate for dynamics
757 mismatch by modifying the reward with a classifier that distinguishes source from target transitions,
758 though this requires online RL access to the target domain. Lauffer et al. [2025] address the covariate
759 shift that arises when a learned policy diverges into unfamiliar dynamics, using on-policy expert
760 corrections for multi-turn LM agents. Finally, Espinosa-Dice et al. [2025] study the misspecification
761 setting in which the expert falls outside the learner’s policy class.

762 All of these methods attempt to *correct for* or *tolerate* dynamics mismatch via reward modification,
763 distribution matching, or robust optimization, and most require either interactive access to the target
764 environment or assumptions on the magnitude of the shift. Our work is complementary: we make no
765 assumptions on the dynamics and instead learn, from offline data alone, to *detect* when the shift is
766 too large to act reliably and abstain, with distribution-free PAC guarantees on both the stopping rate
767 and the stopped regret.

768 **Control-theoretic safety and reachability.** Classical approaches to safe deployment ground safety
769 guarantees in known system dynamics. Hamilton-Jacobi reachability [Tomlin et al., 2000, Bansal
770 et al., 2017] computes the exact set of states from which a safety constraint can be guaranteed, and
771 control barrier functions (CBFs) [Ames et al., 2019] enforce forward invariance of a safe set via
772 Lyapunov-style certificates. Prajna and Jadbabaie [Prajna and Jadbabaie, 2004] give analogous barrier
773 certificates for hybrid systems. These methods offer strong formal guarantees, but require an accurate
774 model of the transition dynamics.

775 **Uncertainty-aware deployment.** A natural response to distribution shift is to monitor epistemic
776 uncertainty at deployment and abstain or intervene when it is high; MC dropout [Gal and Ghahramani,
777 2016] and deep ensembles [Lakshminarayanan et al., 2017] are the dominant practical tools for
778 obtaining such estimates. These estimates have been deployed directly in safety-critical settings:
779 Lütjens et al. [2019] use MC dropout and bootstrapped ensembles to modulate collision-avoidance
780 policies when epistemic uncertainty is high, and Kahn et al. [2017] propagate learned uncertainty
781 estimates through a model-based planner for robot navigation. However, the thresholds used to trigger
782 abstention in these methods are typically heuristic. Neggatu et al. [2025] use uncertainty estimates
783 to switch between an offline RL policy and a more conservative behavior-cloning policy, but do not
784 prove theoretical guarantees, unlike ours. Several works [Haider et al., 2023, Prashant et al., 2025]
785 frame out-of-distribution shift in RL as changes in the transition dynamics, but the proposed methods
786 do not rigorously bound the downstream regret.

787 Conformal prediction (see Angelopoulos and Bates [2023] for a foundational treatment) offers
788 a distribution-free alternative with rigorous finite-sample coverage guarantees. Tibshirani et al.
789 [2019] extend the framework to covariate shift via importance weighting; and Gibbs and Candes
790 [2021] develop online conformal methods that adapt to distribution shift as it occurs. Recent
791 work has extended these ideas to sequential and dynamical settings. Lindemann et al. [2023a] use
792 conformal prediction to construct probabilistic prediction regions for safe motion planning in dynamic
793 environments; Dixit et al. [2023] apply adaptive conformal prediction to quantify multi-step-ahead
794 uncertainty for MPC among dynamic agents; Lindemann et al. [2023b] give conformal runtime
795 verification guarantees for signal temporal logic specifications; and Lee and Matni [2024] establish
796 conformal coverage guarantees from a single trajectory of temporally correlated data generated by an
797 unknown stochastic dynamical system.

798 All of these methods share our high-level goal of detecting when a deployed agent has ventured
799 outside its training support and should stop. Ensemble and conformal methods differ from ours in
800 a fundamental way, however: ensemble approaches rely on heuristic variance thresholds to trigger
801 abstention, while conformal methods provide rigorous marginal coverage for *prediction sets* at
802 each step. Neither directly bounds the regret accumulated along a trajectory before abstention.
803 Our approach avoids both heuristics and prediction sets, instead providing frequentist PAC-style
804 guarantees directly on the stopping time that bound the stopped regret in the test environment.

805 B Additional Results and Proofs from the Deterministic Section

806 B.1 Supporting technical lemmas

807 **Lemma 4** (Prefix coupling for deterministic policies). *Let π and π^* be deterministic policies on*
 808 *a horizon- H MDP with pre-action filtration $\mathcal{G}_h = \sigma(s_1, a_1, \dots, s_{h-1}, a_{h-1}, s_h)$. Define the first-*
 809 *deviation time*

$$\tau_{\pi, \{\pi^*\}} := \min\{h \in [H] : \pi_h(s_h) \neq \pi_h^*(s_h)\},$$

810 *with the convention $\min \emptyset = H + 1$. Let τ be any stopping time adapted to $\{\mathcal{G}_h\}$, and define*
 811 *$h^* := \min(\tau, \tau_{\pi, \{\pi^*\}})$, and define the stopped sigma-algebra*

$$\mathcal{G}_{h^*} := \{A \subseteq \Omega : A \cap \{h^* = h\} \in \mathcal{G}_h \text{ for every } h \in [H + 1]\}.$$

812 *Intuitively, \mathcal{G}_{h^*} represents the information available up until the (random) time h^* . Then, for any*
 813 *environment M , every event $\mathcal{E} \in \mathcal{G}_{h^*}$ satisfies $\Pr_{M, \pi}[\mathcal{E}] = \Pr_{M, \pi^*}[\mathcal{E}]$. Consequently, for every*
 814 *bounded \mathcal{G}_{h^*} -measurable random variable Z , we also have $\mathbb{E}_{M, \pi}[Z] = \mathbb{E}_{M, \pi^*}[Z]$.*

815 *Proof.* We first check that h^* is a stopping time. For each $k \in [H]$, the event $\{\tau_{\pi, \{\pi^*\}} \leq k\} =$
 816 $\bigcup_{j=1}^k \{\pi_j(s_j) \neq \pi_j^*(s_j)\}$ belongs to \mathcal{G}_k because each s_j is $\mathcal{G}_j \subseteq \mathcal{G}_k$ -measurable and π_j, π_j^* are
 817 deterministic; for $k = H + 1$ the claim is trivial. Hence $\tau_{\pi, \{\pi^*\}}$ is a stopping time, and so is
 818 $h^* = \min(\tau, \tau_{\pi, \{\pi^*\}})$.

819 It remains to show that π and π^* induce the same distribution on \mathcal{G}_{h^*} -measurable events. Fix $k \in [H]$
 820 and any pre-action prefix $\mathbf{t}_k = (s_1, a_1, \dots, s_{k-1}, a_{k-1}, s_k)$ with $h^*(\mathbf{t}_k) = k$. On the event $\{h^* = k\}$
 821 we have $k \leq \tau_{\pi, \{\pi^*\}}$, so $\pi_j(s_j) = \pi_j^*(s_j)$ for all $j < k$, and therefore

$$\begin{aligned} \Pr_{M, \pi}[(s_1, a_1, \dots, s_{k-1}, a_{k-1}, s_k) = \mathbf{t}_k, h^* = k] &= P_0(s_1) \prod_{j=1}^{k-1} P_j(s_{j+1} \mid s_j, \pi_j(s_j)) \\ &= \Pr_{M, \pi^*}[(s_1, a_1, \dots, s_{k-1}, a_{k-1}, s_k) = \mathbf{t}_k, h^* = k]. \end{aligned}$$

822 When $k = H + 1$, the same argument applied to full trajectories shows that the probabilities also
 823 agree on the slice $\{h^* = H + 1\}$, since $\pi_j(s_j) = \pi_j^*(s_j)$ for every $j \in [H]$ there. For any $\mathcal{E} \in \mathcal{G}_{h^*}$,
 824 the slice $\mathcal{E} \cap \{h^* = k\}$ is \mathcal{G}_k -measurable, so its probability is a sum of prefix probabilities of the
 825 form above; summing over k gives $\Pr_{M, \pi}[\mathcal{E}] = \Pr_{M, \pi^*}[\mathcal{E}]$. Because the two laws are equal on \mathcal{G}_{h^*} ,
 826 their expectations agree for every bounded \mathcal{G}_{h^*} -measurable random variable. \square

827 **Lemma 5** (General statement of Lemma 1). *Let M be an positive integer and let $\mathcal{T} \subseteq [M]^n$ be any*
 828 *set of vectors. For any target $\epsilon \in (0, 1)$, set $K = \lceil 1/\epsilon \rceil$. There exists a probability distribution q^**
 829 *over subsets $S \subseteq \mathcal{T}$ of size at most K with the following two properties.*

830 1. (**Coverage**) For every $\tau \in \mathcal{T}$:

$$\mathbb{E}_{S \sim q^*} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{1} \left[\min_{v \in S} v_i > \tau_i \right] \right] < \epsilon.$$

831 2. (**Sparsity**) The support of q^* consists entirely of sets of size at most $\lceil 1/\epsilon \rceil$.

832 *Proof.* Define a finite zero-sum game in which the maximizer chooses $p \in \Delta(\mathcal{T})$ and the minimizer
 833 chooses $q \in \Delta(\mathcal{T}^K)$, with payoff

$$\phi(p, q) := \mathbb{E}_{\tau \sim p} \mathbb{E}_{S \sim q} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{1} \left[\min_{v \in S} v_i > \tau_i \right] \right].$$

834 Since ϕ is bilinear in (p, q) , von Neumann's minimax theorem gives

$$\min_q \max_p \phi(p, q) = \max_p \min_q \phi(p, q).$$

835 To bound the maximin, fix any $p \in \Delta(\mathcal{T})$ and let p^K denote the strategy of forming $S = \{v_1, \dots, v_K\}$
836 by drawing v_1, \dots, v_K i.i.d. from p . Then

$$\min_q \phi(p, q) \leq \phi(p, p^K) = \frac{1}{n} \sum_{i=1}^n \Pr_{\tau, v_1, \dots, v_K \sim p} \left(\min_{j \in [K]} (v_j)_i > \tau_i \right).$$

837 For a fixed coordinate i , the event $\{\min_j (v_j)_i > \tau_i\}$ holds if and only if τ_i is strictly the smallest
838 value among the $K + 1$ i.i.d. draws $\tau_i, (v_1)_i, \dots, (v_K)_i$. Since these draws are exchangeable, each is
839 equally likely to be the strict minimum, and since these $K + 1$ events are mutually exclusive, each
840 occurs with probability at most $1/(K + 1)$. Therefore

$$\phi(p, p^K) \leq \frac{1}{n} \sum_{i=1}^n \frac{1}{K + 1} = \frac{1}{K + 1}.$$

841 Since this holds for every $p \in \Delta(\mathcal{T})$, we conclude $\max_p \min_q \phi(p, q) \leq 1/(K + 1)$.

842 Combining with the minimax equality and using $K = \lceil 1/\epsilon \rceil$ so that $K + 1 > 1/\epsilon$:

$$\min_q \max_p \phi(p, q) = \max_p \min_q \phi(p, q) \leq \frac{1}{K + 1} < \epsilon.$$

843 Hence there exists $q^* \in \Delta(\mathcal{T}^K)$ with $\max_p \phi(p, q^*) < \epsilon$. Taking $p = \delta_\tau$ for any fixed $\tau \in \mathcal{T}$ yields
844 the coverage condition. The sparsity condition holds because every S in the support of q^* satisfies
845 $|S| \leq K = \lceil 1/\epsilon \rceil$. \square

846 B.2 Analysis of stepwise reduction to PQ-learning

847 Assume $\mathcal{A} = \{0, 1\}$ and $\pi^* \in \Pi$ is deterministic. For each $h \in [H]$, let μ_h^M and μ_h^N denote the
848 expert's step- h state marginals under M and N . We apply Rejection [Goldwasser et al., 2020] at each
849 step h with labeled examples $(s, \pi_h^*(s))$ for $s \sim \mu_h^M$ and unlabeled examples from $\bar{\mu}_h := \frac{1}{2}\mu_h^M + \frac{1}{2}\mu_h^N$,
850 yielding a predictor-selector pair $(\tilde{\pi}_h, g_h)$. Define $\tilde{\pi} = (\tilde{\pi}_1, \dots, \tilde{\pi}_H)$, $g = (g_1, \dots, g_H)$, and
851 $\tau_g = \min\{h : g_h(s_h) = 0\}$. We run each step- h learner with confidence δ/H and union-bound over
852 h , so all guarantees below hold simultaneously with probability at least $1 - \delta$.

853 For tolerance γ , the PQ guarantee gives for every $h \in [H]$:

$$\Pr_{s \sim \mu_h^M} [g_h(s) = 0] \leq \gamma, \tag{3}$$

$$\Pr_{s \sim \bar{\mu}_h} [g_h(s) = 1, \tilde{\pi}_h(s) \neq \pi_h^*(s)] \leq \gamma. \tag{4}$$

854 Since $\bar{\mu}_h$ averages μ_h^M and μ_h^N , both marginals inherit a 2γ accepted-error bound from (4).

855 **Abstention on M .** Let E_h be the event that step h is the *first* time the learner either abstains or makes
856 an unrejected mistake and let $A_h := \bigcap_{t=1}^{h-1} \{g_t(s_t) = 1, \tilde{\pi}_t(s_t) = \pi_t^*(s_t)\}$. Thus A_h is the event
857 that, before step h , the learner has neither abstained nor made an unrejected mistake. By definition,

$$E_h = A_h \cap \left(\{g_h(s_h) = 0\} \cup \{g_h(s_h) = 1, \tilde{\pi}_h(s_h) \neq \pi_h^*(s_h)\} \right).$$

858 Since the event inside the intersection is \mathcal{G}_h -measurable, by Lemma 4, we may change measure from
859 $\Pr_{M, \tilde{\pi}}$ to \Pr_{M, π^*} :

$$\begin{aligned} \Pr_{M, \tilde{\pi}}(E_h) &= \Pr_{M, \pi^*} \left(A_h \cap \left(\{g_h(s_h) = 0\} \cup \{g_h(s_h) = 1, \tilde{\pi}_h(s_h) \neq \pi_h^*(s_h)\} \right) \right) \\ &\leq \Pr_{M, \pi^*} \left(\{g_h(s_h) = 0\} \cup \{g_h(s_h) = 1, \tilde{\pi}_h(s_h) \neq \pi_h^*(s_h)\} \right) \\ &\leq \Pr_{M, \pi^*} (g_h(s_h) = 0) + \Pr_{M, \pi^*} (g_h(s_h) = 1, \tilde{\pi}_h(s_h) \neq \pi_h^*(s_h)) \\ &= \Pr_{s \sim \mu_h^M} [g_h(s) = 0] + \Pr_{s \sim \mu_h^M} [g_h(s) = 1, \tilde{\pi}_h(s) \neq \pi_h^*(s)] \leq 3\gamma. \end{aligned}$$

860 Since abstaining implies E_h occurs for some h , we have

$$\alpha_M(\tilde{\pi}, g) \leq \sum_{h=1}^H \Pr(E_h) \leq 3H\gamma.$$

861 **Regret on N .** Let F_h be the event that the first unrejected misstep occurs at step h prior to any
 862 abstention. Similar to the argument above, if $A_h := \bigcap_{t=1}^{h-1} \{g_t(s_t) = 1, \tilde{\pi}_t(s_t) = \pi_t^*(s_t)\}$, then

$$F_h = A_h \cap \{g_h(s_h) = 1, \tilde{\pi}_h(s_h) \neq \pi_h^*(s_h)\}.$$

863 We may first change measure from $\Pr_{N, \tilde{\pi}}$ to \Pr_{N, π^*} while retaining the factor A_h , and then drop A_h .
 864 Thus

$$\Pr_{N, \tilde{\pi}}(F_h) \leq \Pr_{N, \pi^*}(g_h(s_h) = 1, \tilde{\pi}_h(s_h) \neq \pi_h^*(s_h)) = \Pr_{s \sim \mu_h^N}[g_h(s) = 1, \tilde{\pi}_h(s) \neq \pi_h^*(s)] \leq 2\gamma.$$

865 Since the learner incurs zero stopped regret unless an unrejected misstep occurs—which costs at most
 866 C_{\max} —we obtain

$$\text{Regret}_N(\tilde{\pi}, g; c) \leq C_{\max} \sum_{h=1}^H \Pr(F_h) \leq 2C_{\max} H \gamma.$$

867 **Sample complexity.** Setting $\gamma = \Theta(\epsilon/(HC_{\max}))$ to achieve $\text{Regret}_N \leq \epsilon$, and using the PQ sample
 868 complexity of $\tilde{O}(\log |\Pi|/\gamma^2)$ per step, yields $m = n = \tilde{O}\left(\frac{H^2 C_{\max}^2 \log |\Pi|}{\epsilon^2}\right)$ trajectories (since one
 869 trajectory provides an example for each step). This H^2 penalty is unavoidable in this approach: the
 870 reduction certifies H isolated decisions rather than a single trajectory prefix, and the horizon factor
 871 enters through the per-step tolerance regardless of the cost structure.

872 B.3 Proof of Lemma 1

873 Let $M = H + 1$. For each policy $\pi \in \Pi_{\text{version}}$, let $v^\pi \in [M]^n$ be its empirical stopping-time
 874 vector on the n test trajectories, where the j -th coordinate is given by $(v^\pi)_j = \tau_{\pi_0, \{\pi\}}(T_j)$. Let
 875 $\mathcal{T} = \{v^\pi : \pi \in \Pi_{\text{version}}\}$ be the finite set of all such vectors.

876 Applying Lemma 5 to the set \mathcal{T} with target error $\epsilon = \rho$ implies the existence of a distribution q^* over
 877 subsets $\Phi \subseteq \Pi_{\text{version}}$ satisfying:

- 878 • For all $\pi \in \Pi_{\text{version}}$: $\mathbb{E}_{\Phi \sim q^*} \left[\frac{1}{n} \sum_{j=1}^n \mathbf{1}[\tau_{\pi_0, \Phi}(T_j) > \tau_{\pi_0, \{\pi\}}(T_j)] \right] < \rho$.
- 879 • For all $\Phi \in \text{supp}(q^*)$: $|\Phi| \leq \lceil 1/\rho \rceil$.

880 This is exactly the q^* required by Lemma 1.

881 B.4 Proof of Theorem 1

882 We prove the theorem through four steps: constructing the validator distribution, bounding the
 883 empirical risk, generalizing to the population, and bounding the regret on N .

884 Throughout, let $K = \lceil \log_2(5/\delta) \rceil \lceil 2/\eta \rceil$ denote the maximum ensemble size and

$$Z := (K + 1) \log |\Pi| + \log \frac{5}{\delta}$$

885 the complexity term appearing in generalization bounds.

886 **Step 1 (Existence of q^*).** By Proposition 1, Algorithm 2 yields a $(\eta/2 + \xi)$ -valid distribution q^* over
 887 subsets of Π_{version} such that every $\Phi' \in \text{supp}(q^*)$ satisfies $|\Phi'| \leq \left\lceil \frac{2}{\eta} \right\rceil$ with probability $1 - \delta/5$.

888 **Step 2 (Bound the empirical risk).** The output $\Phi = \bigcup_{i=1}^k \Phi_i$ is the union of $k = \lceil \log_2(5/\delta) \rceil$
 889 independent draws from q^* , so $|\Phi| \leq k \cdot \lceil 2/\eta \rceil = K$ deterministically. By Markov's inequality
 890 applied to the $\eta/2 + \xi$ coverage guarantee,

$$\begin{aligned} & \Pr_{\Phi_i \sim q^*} \left(\frac{1}{n} \sum_{j=1}^n \mathbf{1}[\tau_{\pi_0, \Phi_i}(T_j) > \tau_{\pi_0, \{\pi^*\}}(T_j)] > \eta + 2\xi \right) \\ & \leq \frac{\mathbb{E}_{\Phi_i \sim q^*} \left[\frac{1}{n} \sum_{j=1}^n \mathbf{1}[\tau_{\pi_0, \Phi_i}(T_j) > \tau_{\pi_0, \{\pi^*\}}(T_j)] \right]}{\eta + 2\xi} < \frac{\eta/2 + \xi}{\eta + 2\xi} = \frac{1}{2}. \end{aligned}$$

891 We refer to the term on the left hand side as the *risk*. Since $\tau_{\pi_0, \Phi}(T_j) = \min_i \tau_{\pi_0, \Phi_i}(T_j)$, then if any
 892 Φ_i achieves risk $\leq \eta$, then so does Φ . Equivalently, Φ only has high risk if all k independent draws
 893 have high risk. Since these draws are independent,

$$\Pr \left(\frac{1}{n} \sum_{j=1}^n \mathbf{1}[\tau_{\pi_0, \Phi}(T_j) > \tau_{\pi_0, \{\pi^*\}}(T_j)] > \eta + 2\xi \right) \leq \prod_{i=1}^k \frac{1}{2} \leq \left(\frac{1}{2} \right)^{\lceil \log_2(5/\delta) \rceil} \leq \delta/5.$$

894 For the training data, all policies in $\Phi \cup \{\pi^*\}$ belong to Π_{version} by construction, meaning they agree
 895 with π^* on every observed state-action pair in $\mathcal{D}_{\text{train}}$. Therefore, for all $i \in [m]$:

$$\tau_{\pi_0, \Phi}(S_i) = \tau_{\pi_0, \{\pi^*\}}(S_i) = H + 1,$$

896 where the first equality holds because no policy in Φ ever disagrees with π_0 on any training trajectory
 897 (since both agree with π^* on all observed states), and the second holds because π^* and π_0 both belong
 898 to Π_{version} and hence agree on all observed actions. Since $\tau_{\pi_0, \Phi}(S_i) = H + 1 > H$ for all i , the
 899 empirical stopping rate $\frac{1}{m} \sum_{i=1}^m \mathbf{1}[\tau_{\pi_0, \Phi}(S_i) \leq H] = 0$, and since $\tau_{\pi_0, \Phi}(S_i) = \tau_{\pi_0, \{\pi^*\}}(S_i)$ for all
 900 i , the empirical risk $\frac{1}{m} \sum_{i=1}^m \mathbf{1}[\tau_{\pi_0, \Phi}(S_i) > \tau_{\pi_0, \{\pi^*\}}(S_i)] = 0$ as well.

901 **Step 3 (Bound generalization error on the expert's distribution).** We use uniform convergence
 902 over $\mathcal{H}_K = \{(\pi_0, \Phi) : \pi_0 \in \Pi, \Phi \subseteq \Pi, |\Phi| \leq K\}$, which satisfies $|\mathcal{H}_K| \leq |\Pi|^{K+1}$ as π_0 and each
 903 of the at most K elements of Φ are chosen from Π . We identify three failure events each occurring
 904 with probability at most $\delta/5$.

905 *Test risk generalization.* Let $Y(\pi_0, \Phi) = \mathbf{1}[\tau_{\pi_0, \Phi}(T) > \tau_{\pi_0, \{\pi^*\}}(T)]$, so that $\hat{\mathbb{E}}[Y] = \frac{1}{n} \sum_{j=1}^n Y(T_j)$
 906 is the empirical test risk over $\mathcal{D}_{\text{test}}$ and $\mathbb{E}_N^{\pi^*}[Y] = \Pr_{N, \pi^*}[\tau_{\pi_0, \Phi}(T) > \tau_{\pi_0, \{\pi^*\}}(T)]$ is the true test
 907 risk. Since $Y \in \{0, 1\}$, $\widehat{\text{Var}}(Y) \leq \hat{\mathbb{E}}[Y^2] = \hat{\mathbb{E}}[Y]$. By the empirical Bernstein inequality applied to
 908 each fixed hypothesis in \mathcal{H}_K , then union-bounded over all $|\mathcal{H}_K|$ hypotheses, with probability at least
 909 $1 - \delta/4$, every $(\pi_0, \Phi) \in \mathcal{H}_K$ simultaneously satisfies:

$$\mathbb{E}_N^{\pi^*}[Y] \leq \hat{\mathbb{E}}[Y] + \sqrt{\frac{2\hat{\mathbb{E}}[Y] \log(5|\mathcal{H}_K|/\delta)}{n}} + \frac{3 \log(5|\mathcal{H}_K|/\delta)}{n}.$$

910 Since $|\mathcal{H}_K| \leq |\Pi|^{K+1}$, we have $\log(5|\mathcal{H}_K|/\delta) \leq (K+1) \log |\Pi| + \log(5/\delta) = Z$. Since the
 911 empirical step guarantees $\hat{\mathbb{E}}[Y] \leq \eta + 2\xi$, this implies:

$$\mathbb{E}_N^{\pi^*}[Y] \leq \eta + 2\xi + \sqrt{\frac{2(\eta + 2\xi)Z}{n}} + \frac{3Z}{n}.$$

912 *Train abstention generalization.* The true stopping rate of (π_0, Φ) on expert trajectories is
 913 $\Pr_{M, \pi^*}[\tau_{\pi_0, \Phi}(S) \leq H]$ and the empirical stopping rate over $\mathcal{D}_{\text{train}}$ is $\frac{1}{m} \sum_{i=1}^m \mathbf{1}[\tau_{\pi_0, \Phi}(S_i) \leq H]$.
 914 For any fixed $(\pi_0, \Phi) \in \mathcal{H}_K$ with $\Pr_{M, \pi^*}[\tau_{\pi_0, \Phi}(S) \leq H] > \frac{Z}{m}$, each of the m i.i.d. draws
 915 $S_i \sim (M, \pi^*)$ independently fails to trigger abstention with probability at most $1 - \frac{Z}{m}$, so the
 916 probability of observing zero empirical abstention is at most $(1 - \frac{Z}{m})^m \leq e^{-Z}$. A union bound over
 917 all $|\mathcal{H}_K| \leq |\Pi|^{K+1}$ hypotheses gives failure probability at most $|\mathcal{H}_K| e^{-Z} \leq \delta/5$.

918 *Train risk generalization.* The true train risk of $(\pi_0, \tau_{\pi_0, \Phi})$ is $\Pr_{M, \pi^*}[\tau_{\pi_0, \Phi}(S) > \tau_{\pi_0, \{\pi^*\}}(S)]$ and
 919 the empirical train risk is $\frac{1}{m} \sum_{i=1}^m \mathbf{1}[\tau_{\pi_0, \Phi}(S_i) > \tau_{\pi_0, \{\pi^*\}}(S_i)]$. Just like above, any hypothesis with
 920 $\Pr_{M, \pi^*}[\tau_{\pi_0, \Phi}(S) > \tau_{\pi_0, \{\pi^*\}}(S)] > \frac{Z}{m}$ achieves zero empirical train risk with probability at most
 921 $(1 - \frac{Z}{m})^m \leq e^{-Z}$, and a union bound over \mathcal{H}_K gives failure probability at most $\delta/5$.

922 Taking a union bound over all failure events, with probability at least $1 - \delta$, uniform convergence
 923 holds simultaneously for all three failure events across all hypotheses in \mathcal{H}_K .

924 **Step 4 (Bound stopping rate and regret on the learned policy's distribution).** We first address the
 925 stopping rate on M under $(\pi_0, \tau_{\pi_0, \Phi})$. The generalization steps for the train abstention and train risk,
 926 respectively, give $\Pr_{M, \pi^*}[\tau_{\pi_0, \Phi} \leq H] \leq \frac{Z}{m}$ and $\Pr_{M, \pi^*}[\tau_{\pi_0, \Phi} > \tau_{\pi_0, \{\pi^*\}}] \leq \frac{Z}{m}$. Define the safe
 927 event $B = \{\tau_{\pi_0, \Phi} \leq \tau_{\pi_0, \{\pi^*\}}\}$, so $\Pr_{M, \pi^*}[B^c] \leq \frac{Z}{m}$. Then,

$$\Pr_{M, \pi_0}[\tau_{\pi_0, \Phi} \leq H] \leq \Pr_{M, \pi_0}[\{\tau_{\pi_0, \Phi} \leq H\} \cap B] + \Pr_{M, \pi_0}[B^c].$$

928 To apply Lemma 4, let $h^* = \min(\tau_{\pi_0, \Phi}, \tau_{\pi_0, \{\pi^*\}})$. Both $\{\tau_{\pi_0, \Phi} \leq H\} \cap B$ and B^c lie in \mathcal{G}_{h^*} : for
 929 each $h \in [H]$, the slice $\{\tau_{\pi_0, \Phi} \leq H\} \cap B \cap \{h^* = h\}$ equals $\{\tau_{\pi_0, \Phi} = h\} \cap \{\tau_{\pi_0, \{\pi^*\}} \geq h\}$, which
 930 is \mathcal{G}_h -measurable since $\tau_{\pi_0, \Phi}$ is a stopping time and $\{\tau_{\pi_0, \{\pi^*\}} \geq h\}$ depends only on whether π_0
 931 and π^* agree on s_1, \dots, s_{h-1} (the case $h = H + 1$ is trivial as the slice is the empty set); the same
 932 argument applies to B^c . Lemma 4 therefore gives

$$\Pr_{M, \pi_0} [\{\tau_{\pi_0, \Phi} \leq H\} \cap B] = \Pr_{M, \pi^*} [\{\tau_{\pi_0, \Phi} \leq H\} \cap B] \leq \Pr_{M, \pi^*} [\tau_{\pi_0, \Phi} \leq H] \leq \frac{Z}{m},$$

933 and $\Pr_{M, \pi_0} [B^c] = \Pr_{M, \pi^*} [B^c] \leq \frac{Z}{m}$. Combining yields $\Pr_{M, \pi_0} [\tau_{\pi_0, \Phi} \leq H] \leq \frac{2Z}{m}$.

934 To bound the stopped regret on N , note that $\Pr_{N, \pi_0} [B^c] = \Pr_{N, \pi^*} [B^c] \leq \eta + 2\xi + \sqrt{\frac{2(\eta+2\xi)Z}{n}} + \frac{3Z}{n}$
 935 by the test generalization guarantee and Lemma 4. Decomposing the expected stopped cost of
 936 $(\pi_0, \tau_{\pi_0, \Phi})$ over B and B^c ,

$$\begin{aligned} \text{Regret}_N(\pi_0, \tau_{\pi_0, \Phi}; c) &= \mathbb{E}_N^{\pi_0} \left[\left(\sum_{h=1}^{\tau_{\pi_0, \Phi}-1} c_h \right) \mathbf{1}(B) \right] + \mathbb{E}_N^{\pi_0} \left[\left(\sum_{h=1}^{\tau_{\pi_0, \Phi}-1} c_h \right) \mathbf{1}(B^c) \right] - \mathbb{E}_N^{\pi^*} \left[\sum_{h=1}^{\tau_{\pi_0, \Phi}-1} c_h \right] \\ &\leq \mathbb{E}_N^{\pi^*} \left[\left(\sum_{h=1}^{\tau_{\pi_0, \Phi}-1} c_h \right) \mathbf{1}(B) \right] + C_{\max} \cdot \Pr_{N, \pi_0} [B^c] - \mathbb{E}_N^{\pi^*} \left[\sum_{h=1}^{\tau_{\pi_0, \Phi}-1} c_h \right] \\ &= -\mathbb{E}_N^{\pi^*} \left[\left(\sum_{h=1}^{\tau_{\pi_0, \Phi}-1} c_h \right) \mathbf{1}(B^c) \right] + C_{\max} \cdot \Pr_{N, \pi_0} [B^c] \\ &\leq C_{\max} \cdot \Pr_{N, \pi_0} [B^c] \quad (\text{costs are nonnegative}) \\ &\leq C_{\max} \left(\eta + 2\xi + \sqrt{\frac{2(\eta+2\xi)Z}{n}} + \frac{3Z}{n} \right), \end{aligned}$$

937 where the first inequality applies Lemma 4 on B (because $\tau_{\pi_0, \Phi} \leq \tau_{\pi_0, \{\pi^*\}}$ on B , the selective
 938 policy's expected prefix cost $\sum_{h=1}^{\tau_{\pi_0, \Phi}-1} c_h$ matches the expert's expected cost over that prefix) and
 939 bounds the total cost on B^c by C_{\max} .

940 B.5 Transferring the completeness guarantee to the test environment

941 While our objective does not require bounding the test-side stopping rate $\alpha_N(\pi_0, \tau) := \Pr_{N, \pi_0} [\tau \leq$
 942 $H]$ distribution-free, it is natural to ask how this rate scales if the environment shift is bounded. Let
 943 $d_{\pi^*}^{\text{state}}(M, N)$ denote the total variation distance between the state-trajectory distributions induced by
 944 the expert in M and N .

945 Define the safe event $B = \{\tau \leq \tau_{\pi_0, \{\pi^*\}}\}$ and let $h^* = \min(\tau, \tau_{\pi_0, \{\pi^*\}})$. Both $\{\tau \leq H\} \cap B$ and
 946 B^c lie in \mathcal{G}_{h^*} by the same measurability argument as in the proof of Theorem 1. By a union bound
 947 and two applications of Lemma 4,

$$\begin{aligned} \alpha_N(\pi_0, \tau_{\pi_0, \Phi}) &\leq \Pr_{N, \pi_0} [\{\tau_{\pi_0, \Phi} \leq H\} \cap B] + \Pr_{N, \pi_0} [B^c] \\ &= \Pr_{N, \pi^*} [\{\tau_{\pi_0, \Phi} \leq H\} \cap B] + \Pr_{N, \pi^*} [B^c] \\ &\leq \Pr_{N, \pi^*} [\tau_{\pi_0, \Phi} \leq H] + \Pr_{N, \pi^*} [\tau_{\pi_0, \Phi} > \tau_{\pi_0, \{\pi^*\}}]. \end{aligned}$$

948 For a fixed validator set, the abstention event $\{\tau \leq H\}$ is strictly a measurable function of the state
 949 trajectory. Thus, its probability under the expert can differ between the test environment (N) and
 950 training environment (M) by at most the TV distance:

$$\Pr_{N, \pi^*} [\tau_{\pi_0, \Phi} \leq H] \leq \Pr_{M, \pi^*} [\tau_{\pi_0, \Phi} \leq H] + d_{\pi^*}^{\text{state}}(M, N).$$

951 Finally, substituting the high-probability bounds from the proof of Theorem 1 (where the train-side
 952 abstention is bounded by Z/m and the late-stop risk by $\eta + \sqrt{2\eta Z/n} + 3Z/n$) yields

$$\alpha_N(\pi_0, \tau) \leq \underbrace{\frac{Z}{m} + d_{\pi^*}^{\text{state}}(M, N)}_{\text{PQ-style TV transfer}} + \underbrace{\eta + \sqrt{\frac{2\eta Z}{n}} + \frac{3Z}{n}}_{\text{Sequential penalty}}.$$

953 The first two terms are the sequential analogue of the rejection-rate transfer in batch selective
 954 classification [Goldwasser et al., 2020]: the abstention event is a fixed measurable set, so its probability
 955 shifts by at most the TV distance. The final term has no batch analogue; it arises because the deployed
 956 policy’s actions can steer the trajectory out of the safe prefix before the validators stop it.

957 B.6 Constructing a validator distribution using no-regret dynamics

958 The nonconstructive minimax proof of Lemma 1 can be made algorithmic via the well-known
 959 connection between approximate minimax equilibria and no-regret play. In Section B.6, we formulate
 960 the problem as a repeated game where a maximization player maintains a distribution over the
 961 version space using a no-regret algorithm, and the validation player plays an (approximate) best-
 962 response by sampling a sparse validator set. Averaging the validation player’s choices over T rounds
 963 yields a distribution that converges to the minimax optimum. This gives an explicit, oracle-efficient
 964 implementation of the SPARSEVALIDATORDIST subroutine.

965 Fix the version space Π_{version} , a base policy $\pi_0 \in \Pi_{\text{version}}$, and unlabeled test trajectories $\mathcal{D}_{\text{test}} =$
 966 $\{T_j\}_{j=1}^n$. Write each test trajectory as $T_j = (s_{j1}, \dots, s_{jH})$.

Algorithm 2 SPARSEVALIDATORDIST($\Pi_{\text{version}}, \pi_0, \mathcal{D}_{\text{test}}, \rho, \xi, \delta$)

Require: version space Π_{version} , base policy π_0 , unlabeled test trajectories $\mathcal{D}_{\text{test}} = \{T_j\}_{j=1}^n$, toler-
 ance $\rho \in (0, 1)$, computational slack ξ , confidence δ , no-regret algorithm \mathcal{A} with regret bound
 Reg_T .

- 1: Set $K \leftarrow \lceil 1/\rho \rceil$.
 - 2: Set T large enough that $\frac{\text{Reg}_T}{T} + \sqrt{\frac{\log(1/\delta)}{2T}} \leq \xi$.
 - 3: Initialize any distribution p^1 over Π_{version} .
 - 4: **for** $t = 1, \dots, T$ **do**
 - 5: Draw $\pi_1^t, \dots, \pi_K^t \stackrel{\text{i.i.d.}}{\sim} p^t$ and set $\Phi^t \leftarrow \{\pi_1^t, \dots, \pi_K^t\}$.
 - 6: For each $\pi \in \Pi_{\text{version}}$, define the payoff $u^t(\pi) := \frac{1}{n} \sum_{j=1}^n \mathbf{1} [\tau_{\pi_0, \Phi^t}(T_j) > \tau_{\pi_0, \{\pi\}}(T_j)]$.
 - 7: Update p^{t+1} using \mathcal{A} for maximizing the payoffs u^t .
 - 8: Return the empirical distribution $\bar{q} := \frac{1}{T} \sum_{t=1}^T \delta_{\Phi^t}$.
-

967 The interpretation of the payoff is straightforward: $u^t(\pi)$ is the fraction of empirical test trajectories
 968 on which the validator set Φ^t stops strictly later than the single competing policy π . Thus the
 969 maximization player tries to expose competitors against which the current validator set is weak, while
 970 the validation player responds by sampling a fresh set from the current maximizer distribution.

971 **Proposition 1** (Generic no-regret construction of q^*). *Let \bar{q} be the output of Algorithm 2. With*
 972 *probability at least $1 - \delta$, we have*

$$\sup_{\pi \in \Pi_{\text{version}}} \mathbb{E}_{\Phi \sim \bar{q}} \left[\frac{1}{n} \sum_{j=1}^n \mathbf{1} [\tau_{\pi_0, \Phi}(T_j) > \tau_{\pi_0, \{\pi\}}(T_j)] \right] \leq \rho + \xi.$$

973 *In particular, using Hedge yields $\text{Reg}_T \leq \sqrt{2T \log |\Pi|}$, meaning $T = O\left(\frac{\log |\Pi| + \log(1/\delta)}{\rho^2}\right)$ rounds*
 974 *suffice to return an $O(\rho)$ -valid distribution supported on sets of size at most $\lceil 1/\rho \rceil$.*

975 The no-regret construction above is explicit, but inefficient (when using Hedge) when the version
 976 space is large. A standard choice in online learning with large action spaces is to use Follow-the-
 977 Perturbed-Leader (FTPL) instead of Hedge. However, a naive FTPL implementation would perturb
 978 one coordinate per policy in Π_{version} , which is also intractable. Instead, we show that by perturbing
 979 the *dataset*, and then solving a multiple-instance learning problem [Dietterich et al., 1997, Maron
 980 and Lozano-Pérez, 1997] over this perturbed dataset, we can implement an *oracle-efficient* no-regret
 981 algorithm for the maximizer. This places our implementation in the setting of generalized FTPL
 982 [Dudík et al., 2020]. The full formulation and proofs of this reduction are deferred to Section B.8.

983 **B.7 Proof of Proposition 1**

984 By the definition of \bar{q} , the expected coverage for any fixed $\pi \in \Pi_{\text{version}}$ is exactly the average payoff
 985 over all rounds:

$$\mathbb{E}_{\Phi \sim \bar{q}} \left[\frac{1}{n} \sum_{j=1}^n \mathbf{1} [\tau_{\pi_0, \Phi}(T_j) > \tau_{\pi_0, \{\pi\}}(T_j)] \right] = \frac{1}{T} \sum_{t=1}^T u^t(\pi).$$

986 By the definition of external regret, this average payoff is bounded by the algorithm's average expected
 987 payoff plus the regret term:

$$\frac{1}{T} \sum_{t=1}^T u^t(\pi) \leq \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\pi' \sim p^t} [u^t(\pi')] + \frac{\text{Reg}_T}{T}.$$

988 It therefore remains to control the average expected payoff of the mixed strategy p^t against the
 989 sampled validator set Φ^t .

990 Conditioned on the history \mathcal{H}_{t-1} up to round $t-1$, the distribution p^t is fixed, and Φ^t is obtained by
 991 drawing $K = \lceil 1/\rho \rceil$ policies i.i.d. from p^t . The exchangeability argument from Lemma 1 bounds the
 992 conditional expected payoff:

$$\mathbb{E} \left[\mathbb{E}_{\pi' \sim p^t} [u^t(\pi')] \mid \mathcal{H}_{t-1} \right] \leq \frac{1}{K+1} < \rho.$$

993 Since the expected payoffs $\mathbb{E}_{\pi' \sim p^t} [u^t(\pi')]$ are bounded in $[0, 1]$, we can apply the Azuma–Hoeffding
 994 inequality. With probability at least $1 - \delta$, the empirical average concentrates around its conditional
 995 expectation upper bound:

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\pi' \sim p^t} [u^t(\pi')] \leq \rho + \sqrt{\frac{\log(1/\delta)}{2T}}.$$

996 Combining these bounds completes the proof.

997 **B.8 Oracle-efficient constructions of a validator distribution**

998 As mentioned in Section B.6, the generic no-regret construction of a sparse validator distribution is
 999 inefficient. In this section, we instantiate the generalized FTPL framework of Dudík et al. [2020] to
 1000 obtain an oracle-efficient implementation of the no-regret strategy required by Algorithm 2.

1001 The generalized FTPL framework guarantees efficient no-regret learning, provided an online game
 1002 can be reduced to a linearly perturbed offline optimization problem. To achieve this, we must map
 1003 our abstract policy class into a structured vector space via a translation matrix. Crucially, this matrix
 1004 must satisfy two structural conditions to ensure the random perturbations effectively smooth the
 1005 leader's choices: admissibility (the matrix must have bounded entries and separate distinct policies to
 1006 control the geometric complexity of the action space) and implementability (any payoff generated
 1007 by the adversary must be expressible as a linear combination of the matrix's columns). If these
 1008 conditions hold, FTPL can simulate the online game using synthetic datasets passed to a standard
 1009 offline optimization oracle.

1010 We now formulate the online learning problem by defining the strategy space X , the adversary action
 1011 space Y , and the translation matrix Γ . This is very similar to the online learning problem constructed
 1012 in Section B.6, except slightly more pedantic to ensure that our game satisfies the technical conditions
 1013 required by generalized FTPL.

- 1014 • **Strategy space X :** Quotient Π_{version} by empirical stopping-time equivalence: $\pi \sim \pi' \iff$
 1015 $\tau_{\pi_0, \{\pi\}}(T_j) = \tau_{\pi_0, \{\pi'\}}(T_j)$ for all $j \in [n]$. Let X denote the set of equivalence classes, and
 1016 for $\pi \in X$ write $x_\pi := (\tau_{\pi_0, \{\pi\}}(T_1), \dots, \tau_{\pi_0, \{\pi\}}(T_n)) \in \{1, \dots, H+1\}^n$.
- 1017 • **Adversary action space Y :** Let $Y := \{1, \dots, H+1\}^n$. For any validator set Φ , we define
 1018 its empirical cutoff vector $h^\Phi := (\tau_{\pi_0, \Phi}(T_1), \dots, \tau_{\pi_0, \Phi}(T_n))$, which is an element of Y .

1019 We define the payoff function $f: X \times Y \rightarrow [0, 1]$ as $f(\pi, h) := \frac{1}{n} \sum_{j=1}^n \mathbf{1}[x_{\pi,j} < h_j]$. Observe
 1020 that for any validator set Φ , the payoff evaluates to $f(\pi, h^\Phi) = \frac{1}{n} \sum_{j=1}^n \mathbf{1}[\tau_{\pi_0, \Phi}(T_j) > \tau_{\pi_0, \{\pi\}}(T_j)]$.
 1021 Thus, the validator-distribution problem is exactly equivalent to online learning over X against
 1022 adversary actions in Y with payoff function f .

1023 Finally, we define the generalized FTPL translation matrix Γ . Its rows are indexed by the strategy
 1024 space X , its columns by the environment time steps $(j, c) \in [n] \times \{1, \dots, H\}$, and its entries map
 1025 the policies into structured cumulative features:

$$\Gamma_{\pi, (j, c)} := \mathbf{1}[\tau_{\pi_0, \{\pi\}}(T_j) \leq c].$$

1026 We now prove that Γ satisfies the technical conditions necessary for generalized FTPL.

1027 **Lemma 6** (Admissibility and implementability of the cutoff matrix). *The matrix Γ is 1-admissible*
 1028 *and implementable with complexity 1, in the sense of Dudík et al. [2020].*

1029 *Proof.* We first prove admissibility. Let $\pi, \pi' \in X$ be distinct rows. Since they represent distinct
 1030 stopping-time vectors, there exists some trajectory j such that $\tau_{\pi_0, \{\pi\}}(T_j) \neq \tau_{\pi_0, \{\pi'\}}(T_j)$. Assuming
 1031 without loss of generality that $\tau_{\pi_0, \{\pi\}}(T_j) < \tau_{\pi_0, \{\pi'\}}(T_j)$, we set $c^* := \tau_{\pi_0, \{\pi'\}}(T_j) - 1$. Then
 1032 $\Gamma_{\pi, (j, c^*)} = 1$ and $\Gamma_{\pi', (j, c^*)} = 0$. Thus any two distinct rows are separated by some column. Because
 1033 all matrix entries lie in $\{0, 1\}$, distinct values in a column differ by exactly 1, making Γ 1-admissible.

1034 To prove implementability, we must show that each column can be simulated by a synthetic adversary
 1035 action. Fix a column (j, c) and define the adversary action $y^{(j, c)} \in Y$ by $y_j^{(j, c)} = c + 1$ and $y_{j'}^{(j, c)} = 1$
 1036 for all $j' \neq j$. Then for any $\pi \in X$:

$$\begin{aligned} f(\pi, y^{(j, c)}) &= \frac{1}{n} \left(\mathbf{1}[\tau_{\pi_0, \{\pi\}}(T_j) < c + 1] + \sum_{j' \neq j} \mathbf{1}[\tau_{\pi_0, \{\pi\}}(T_{j'}) < 1] \right) \\ &= \frac{1}{n} \mathbf{1}[\tau_{\pi_0, \{\pi\}}(T_j) \leq c] = \frac{1}{n} \Gamma_{\pi, (j, c)}, \end{aligned}$$

1037 since $\tau_{\pi_0, \{\pi\}}(T_{j'}) \geq 1$ always. Therefore, the weighted singleton dataset $S_{j, c} := \{(n, y^{(j, c)})\}$
 1038 implements column (j, c) , because for all $\pi, \pi' \in X$, $\sum_{(w, y) \in S_{j, c}} w(f(\pi, y) - f(\pi', y)) = \Gamma_{\pi, (j, c)} -$
 1039 $\Gamma_{\pi', (j, c)}$. Thus, each column is implementable by a single synthetic adversary action. \square

1040 By Corollary 2.11 of Dudík et al. [2020], the admissibility and implementability of Γ guarantee the
 1041 existence of an oracle-efficient, randomized no-regret maximization player. Specifically, because our
 1042 translation matrix Γ is 1-admissible (i.e., $\delta = 1$) and contains $N = nH$ columns, the generalized
 1043 FTPL framework yields an expected regret bound of $O(N\sqrt{T}/\delta) = O(nH\sqrt{T})$. It remains only to
 1044 identify the specific offline optimization problem required by the oracle on each round.

1045 It remains only to identify the specific offline optimization problem required by the oracle on each
 1046 round. According to the generalized FTPL algorithm Dudík et al. [2020], on each round t , the
 1047 maximization player draws independent, non-negative perturbations $\alpha_{j, c}^t$ for every column of the
 1048 translation matrix. The player must then compute the perturbed best response against the empirical
 1049 history of the adversary’s actions—which, in our game, are the cutoff vectors h^{Φ^s} of the validator
 1050 sets Φ^s sampled on previous rounds $s < t$. Formally, the oracle must solve:

$$\arg \max_{\pi \in \Pi_{\text{version}}} \left\{ \sum_{s=1}^{t-1} f(\pi, h^{\Phi^s}) + \sum_{j=1}^n \sum_{c=1}^H \alpha_{j, c}^t \Gamma_{\pi, (j, c)} \right\}, \quad (5)$$

1051 where the first term represents the cumulative historical payoff and the second term injects the shared
 1052 linear perturbations. We now connect this oracle optimization problem to the multiple instance
 1053 learning (MIL) framework of Maron and Lozano-Pérez [1997].

1054 **Problem (5) reduces to weighted MIL over the disagreement class.** To implement the oracle
 1055 efficiently, we show that the perturbed leader problem (5) can be cast as a (weighted) MIL problem.
 1056 The intuition behind this reduction is straightforward: a policy stops by step c on trajectory j if and
 1057 only if it disagrees with the base policy π_0 on *at least one* state in the trajectory prefix up to step c .
 1058 This “at least one” condition exactly mirrors the definition of a positive bag in MIL. By treating each

1059 trajectory prefix as a “bag” of states, and any disagreement with π_0 as a positive instance label, we
 1060 can formally map the FTPL objective to a realizable MIL objective.

1061 **Proposition 2.** *For the base policy π_0 and any policy $\pi \in \Pi_{\text{version}}$, define the disagreement hypothesis
 1062 $h_\pi(h, s) := \mathbf{1}[\pi_h(s) \neq \pi_{0,h}(s)]$, and let $\mathcal{H}_{\pi_0} := \{h_\pi : \pi \in \Pi_{\text{version}}\}$ be the corresponding
 1063 disagreement class. On round t , the perturbed leader problem in (5) is exactly equivalent to the
 1064 following weighted realizable multiple-instance disagreement problem over \mathcal{H}_{π_0} :*

$$\arg \max_{\pi \in \Pi_{\text{version}}} \sum_{j=1}^n \sum_{c=1}^H w_{j,c}^{t-1} \mathbf{1}[\exists(h, s_{jh}) \in B_{j,c} : h_\pi(h, s_{jh}) = 1], \quad (6)$$

1065 where the aggregated weights and prefix bags are defined as

$$w_{j,c}^{t-1} := |\{s < t : \tau_{\pi_0, \Phi^s}(T_j) - 1 = c\}| + n\alpha_{j,c}^t, \quad B_{j,c} := \{(1, s_{j1}), \dots, (c, s_{jc})\}.$$

1066 *Proof.* Multiplying the objective in (5) by n (which preserves the argmax) and expanding the
 1067 definitions of f and Γ yields:

$$\arg \max_{\pi \in \Pi_{\text{version}}} \left\{ \sum_{s=1}^{t-1} \sum_{j=1}^n \mathbf{1}[\tau_{\pi_0, \{\pi\}}(T_j) \leq \tau_{\pi_0, \Phi^s}(T_j) - 1] + n \sum_{j=1}^n \sum_{c=1}^H \alpha_{j,c}^t \mathbf{1}[\tau_{\pi_0, \{\pi\}}(T_j) \leq c] \right\}.$$

1068 Regrouping the first term by the cutoff value c and combining it with the perturbation term gives the
 1069 simplified objective $\sum_{j=1}^n \sum_{c=1}^H w_{j,c}^{t-1} \mathbf{1}[\tau_{\pi_0, \{\pi\}}(T_j) \leq c]$. We map this to the MIL framework using
 1070 the disagreement hypothesis $h_\pi(h, s) = \mathbf{1}[\pi_h(s) \neq \pi_{0,h}(s)]$. By definition, the event $\tau_{\pi_0, \{\pi\}}(T_j) \leq$
 1071 c occurs if and only if there is some step $h \leq c$ where $\pi_h(s_{jh}) \neq \pi_{0,h}(s_{jh})$. This is precisely
 1072 the condition that the prefix bag $B_{j,c}$ is labeled positive by h_π , establishing (6). Finally, since any
 1073 $\pi \in \Pi_{\text{version}}$ matches the expert on the labeled training data, h_π is strictly zero on those states,
 1074 satisfying the negative singleton constraints of realizable MIL. \square

1075 **Problem (5) reduces to weighted MIL over the original class.** Alternatively, for finite action
 1076 spaces \mathcal{A} , we can avoid the binary disagreement wrapper and reduce directly over Π_{version} via data
 1077 augmentation. By duplicating each state in a prefix for every alternative action $a \neq \pi_0(s)$, we create
 1078 an augmented bag that a standard multi-class MIL solver evaluates as positive if π outputs *any* of
 1079 those alternative actions. This inflates bag sizes by $|\mathcal{A}| - 1$.

1080 **Proposition 3** (Exact perturbed best response as multi-class MIL over Π_{version}). *Assume a finite
 1081 action space \mathcal{A} . On round t , the perturbed leader problem in (5) is exactly equivalent to the following
 1082 weighted realizable multi-class multiple-instance problem over Π_{version} :*

$$\arg \max_{\pi \in \Pi_{\text{version}}} \sum_{j=1}^n \sum_{c=1}^H w_{j,c}^{t-1} \mathbf{1}[\exists((h, s_{jh}), a) \in \tilde{B}_{j,c} : \pi_h(s_{jh}) = a], \quad (7)$$

1083 with weights $w_{j,c}^{t-1}$ as in Proposition 2 and augmented prefix bags:

$$\tilde{B}_{j,c} := \{((h, s_{jh}), a) \mid h \in \{1, \dots, c\}, a \in \mathcal{A} \setminus \{\pi_{0,h}(s_{jh})\}\}.$$

1084 *Proof.* The simplified FTPL objective matches Proposition 2; we only need to verify the positive bag
 1085 condition. A policy π disagrees with π_0 on the prefix up to step c if and only if $\pi_h(s_{jh}) = a$ for some
 1086 step $h \leq c$ and alternative action $a \in \mathcal{A} \setminus \{\pi_{0,h}(s_{jh})\}$. This exactly matches the condition that π
 1087 predicts at least one target label a in the augmented bag $\tilde{B}_{j,c}$. Finally, since all $\pi \in \Pi_{\text{version}}$ perfectly
 1088 match the expert on the labeled training data, they never predict alternative actions on those states,
 1089 naturally satisfying the negative singleton constraints of realizable MIL. \square

1090 **Agnostic oracle and computational-statistical tradeoff.** So far, we have established the existence of
 1091 an oracle-efficient subroutine to compute a validator distribution. However, the specific optimization
 1092 problem, (weighted) MIL, is somewhat of an exotic class of problem and it is not clear such oracles
 1093 may actually be implemented or approximated in practice. It turns out, however, that we can relax
 1094 the requirement of a MIL oracle for Π to a standard *agnostic learning* oracle for Π (in fact, we can

1095 relax this to a weaker oracle known as a *reliable learning* oracle [Kalai and Kanade, 2021]). This is
 1096 simply due to the stepwise reduction to PQ-learning of Section B.2.

1097 Of course, the cost is a degradation in our statistical guarantees. In particular, by moving from our
 1098 trajectory-wise problem to the stepwise reduction, we solve an easier computational problem but pick
 1099 up a dependence on H in our sample complexity. We leave the question of whether this tradeoff is
 1100 fundamental, or just a by-product of our algorithms, to future work.

1101 C Proofs from the Stochastic Section

Algorithm 3 SeqRejectron for stochastic policies

Require: policy class Π , labeled rollouts $\mathcal{D}_{\text{train}} = \{(s_{ih}, a_{ih})\}_{i \in [m], h \in [H]}$, unlabeled rollouts
 $\mathcal{D}_{\text{test}} = \{T_j\}_{j=1}^n$, tolerance η , confidence δ , threshold θ , version-space radius γ

- 1: Compute the MLE base policy $\pi_0 \leftarrow \arg \min_{\pi \in \Pi} \text{LogLoss}(\pi, \mathcal{D}_{\text{train}})$.
 - 2: Form the version space $\Pi_{\text{version}}^\gamma$.
 - 3: Let $q^* \leftarrow \text{STOCHASTICSPARSEVALIDATORDIST}(\Pi_{\text{version}}^\gamma, \pi_0, \mathcal{D}_{\text{test}}, \eta/2, \theta, \delta/4)$.
 - 4: Sample $\Phi_1, \dots, \Phi_k \stackrel{\text{i.i.d.}}{\sim} q^*$ with $k \leftarrow \lceil \log_2(4/\delta) \rceil$ and set $\Phi \leftarrow \bigcup_{r=1}^k \Phi_r$.
 - 5: **return** the selective policy $(\pi_0, \tau_{\pi_0, \Phi}^\theta)$.
-

1102 C.1 Proof of Theorem 2

1103 Let $(\pi_0, \tau_{\pi_0, \Phi}^\theta)$ be the output of Algorithm 3. We first control the abstention probability on M , and
 1104 then the stopped trajectory Hellinger distance on N .

1105 **Step 1 (Control abstention on M).** The quantity we want to bound is the probability that the
 1106 cumulative Hellinger distance exceeds θ . Define $\gamma^* := \frac{\log |\Pi| + \log(8/\delta)}{m}$, and let γ be the version space
 1107 radius which is used by Algorithm 3. By assumption on m , we have that $\gamma^* \leq \gamma$. Thus, by Lemma 7,
 1108 with probability at least $1 - \delta/4$, we have $\pi^* \in \Pi_{\text{version}}^\gamma$ and every $\pi \in \Pi_{\text{version}}^\gamma$ satisfies

$$D_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi^*}) \leq \frac{\gamma}{2} + \gamma^* \leq 1.5\gamma.$$

1109 Combining this with the Hellinger triangle inequality, we have that for every $\pi \in \Pi_{\text{version}}^\gamma$,

$$D_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi_0}) \leq 2D_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi^*}) + 2D_H^2(\mathcal{P}_M^{\pi_0}, \mathcal{P}_M^{\pi^*}) \leq 6\gamma.$$

1110 Now consider any validator $\pi \in \Phi \subseteq \Pi_{\text{version}}^\gamma$. Let E be the event that the cumulative Hellinger
 1111 distance between π and π_0 exceeds θ . To bound the probability of E under π_0 , we introduce the
 1112 geometric mixture policy $\bar{\pi} \propto \sqrt{\pi \pi_0}$, using the convention in Lemma 8. By Lemma 8, the trajectory
 1113 affinity satisfies

$$1 - D_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi_0}) \leq \mathbb{E}_{T \sim \mathcal{P}_M^{\bar{\pi}}} \left[e^{-\sum_{h=1}^H d_H^2(\pi_h, \pi_{0,h})} \right] \leq \frac{\Pr_{\mathcal{P}_M^{\bar{\pi}}}[E] e^{-\theta} + 1 - \Pr_{\mathcal{P}_M^{\bar{\pi}}}[E]}{2}.$$

1114 Rearranging this inequality yields $\Pr_{\mathcal{P}_M^{\bar{\pi}}}[E] \leq \frac{1}{1 - e^{-\theta}} D_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi_0}) \leq \frac{1 + \theta}{\theta} D_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi_0})$. To
 1115 shift this probability from the mixture policy to the deployed policy π_0 , we apply a Hellinger
 1116 change-of-measure inequality (e.g., Lemma 3.1 of Foster et al. [2024]) alongside Lemma 9:

$$\begin{aligned} \Pr_{\mathcal{P}_M^{\pi_0}}[E] &\leq 2 \Pr_{\mathcal{P}_M^{\bar{\pi}}}[E] + 4D_H^2(\mathcal{P}_M^{\bar{\pi}}, \mathcal{P}_M^{\pi_0}) \\ &\leq \frac{2(1 + \theta)}{\theta} D_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi_0}) + 6D_H^2(\mathcal{P}_M^{\bar{\pi}}, \mathcal{P}_M^{\pi_0}) \\ &= \left(\frac{2}{\theta} + 8 \right) D_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi_0}). \end{aligned}$$

1117 Substituting our version space bound $D_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi_0}) \leq 6\gamma$, the probability of E under π_0 is at
 1118 most $\left(\frac{12}{\theta} + 48 \right) \gamma$. Note that the algorithm abstains ($\tau_{\pi_0, \Phi}^\theta \leq H$) if and only if at least one of these
 1119 cumulative distances exceeds θ . By a union bound over the $|\Phi| \leq K_{\text{ens}}$ validators,

$$\alpha_M(\pi_0, g) \leq K_{\text{ens}} \left(\frac{12}{\theta} + 48 \right) \gamma.$$

1120 **Step 2 (Bound Hellinger Distance of stopped trajectories on N).** Let $M = H + 1$. For each policy
 1121 $\pi \in \Pi_{\text{version}}^\gamma$, let $v^\pi \in [M]^n$ be its empirical stopping-time vector on the n test trajectories, where the
 1122 j -th coordinate is given by $(v^\pi)_j = \tau_{\pi_0, \{\pi\}}^\theta(T_j)$. Let $\mathcal{T} = \{v^\pi : \pi \in \Pi_{\text{version}}^\gamma\}$ be the finite set of all
 1123 such vectors.

1124 By Lemma 5 applied to \mathcal{T} with target error $\eta/2$, there exists a distribution q^* over subsets $\Phi' \subseteq$
 1125 $\Pi_{\text{version}}^\gamma$ satisfying $|\Phi'| \leq \lceil 2/\eta \rceil$ such that, because $\pi^* \in \Pi_{\text{version}}^\gamma$:

$$\mathbb{E}_{\Phi' \sim q^*} \left[\frac{1}{n} \sum_{j=1}^n \mathbf{1} \left[\tau_{\pi_0, \Phi'}^\theta(T_j) > \tau_{\pi_0, \{\pi^*\}}^\theta(T_j) \right] \right] \leq \frac{\eta}{2}.$$

1126 For a validator set Φ , define the empirical risk $\hat{p}(\Phi) := \frac{1}{n} \sum_{j=1}^n \mathbf{1}[\tau_{\pi_0, \{\pi^*\}}^\theta(T_j) < \tau_{\pi_0, \Phi}^\theta(T_j)]$. By an
 1127 identical argument to Theorem 1 (Markov's inequality and multiple draws) we have $\hat{p}(\Phi) \leq \eta$ with
 1128 probability at least $1 - \delta/4$.

1129 To generalize to expert trajectories in N , note the final validator set satisfies $|\Phi| \leq K_{\text{ens}}$. Consider
 1130 the indicator class

$$\mathcal{Y}_{K_{\text{ens}}} := \{T \mapsto \mathbf{1}[\tau_{\pi_0, \{\pi^*\}}^\theta(T) < \tau_{\pi_0, \Phi}^\theta(T)] : \pi_0 \in \Pi, \Phi \subseteq \Pi, |\Phi| \leq K_{\text{ens}}\},$$

1131 which has cardinality at most $|\Pi|^{K_{\text{ens}}+1}$. Applying the empirical Bernstein inequality with a union
 1132 bound, with probability at least $1 - \delta/4$, simultaneously for all such (π_0, Φ) ,

$$p(\Phi) \leq \hat{p}(\Phi) + \sqrt{\frac{2\hat{p}(\Phi)Z}{n}} + \frac{3Z}{n} \leq \eta + \sqrt{\frac{2\eta Z}{n}} + \frac{3Z}{n},$$

1133 where $p(\Phi) := \Pr_{T \sim (\pi^*, N)} \left(\tau_{\pi_0, \{\pi^*\}}^\theta(T) < \tau_{\pi_0, \Phi}^\theta(T) \right)$. Under the assumption $n \geq 8Z/\eta$, each
 1134 correction term is at most $\eta/2$, yielding $p(\Phi) \leq 2\eta$.

1135 Finally, to transfer this to the deployed trajectories, let $\mathcal{E} := \{\tau_{\pi_0, \{\pi^*\}}^\theta(T) < \tau_{\pi_0, \Phi}^\theta(T)\}$ and define
 1136 $\tau'(T) := \tau_{\pi_0, \Phi \cup \{\pi^*\}}^\theta(T)$. Applying Lemma 10 with validator set $\Phi \cup \{\pi^*\}$ and comparison policy
 1137 π^* , we have deterministically $D_H^2(P_{|\tau'}^{\pi_0}, P_{|\tau'}^{\pi^*}) \leq \theta$, because the event $\tau_{\pi_0, \{\pi^*\}}^\theta(T) < \tau_{\pi_0, \Phi \cup \{\pi^*\}}^\theta(T)$
 1138 has probability zero under (π_0, N) . Since \mathcal{E} is measurable with respect to the trajectory stopped at τ' ,
 1139 the binary Hellinger change-of-measure inequality yields

$$\begin{aligned} \Pr_{(\pi_0, N)}(\mathcal{E}) &\leq \left(\sqrt{\Pr_{(\pi^*, N)}(\mathcal{E})} + \sqrt{2D_H^2 \left(\mathcal{P}_N^{\pi_0|\tau'}, \mathcal{P}_N^{\pi^*|\tau'} \right)} \right)^2 \\ &\leq 2 \Pr_{(\pi^*, N)}(\mathcal{E}) + 4D_H^2 \left(\mathcal{P}_N^{\pi_0|\tau'}, \mathcal{P}_N^{\pi^*|\tau'} \right) \\ &\leq 4\eta + 4\theta, \end{aligned}$$

1140 where the last step uses $p(\Phi) \leq 2\eta$. A final application of Lemma 10, now with validator set Φ ,
 1141 comparison policy π^* , and failure probabilities $4\eta + 4\theta$ and 2η , gives

$$\begin{aligned} D_H^2(P_N^{\pi_0|\tau}, P_N^{\pi^*|\tau}) &\leq \theta + \frac{(4\eta + 4\theta) + 2\eta}{2} \\ &= 3\theta + 3\eta. \end{aligned}$$

1142 The four failure events above each occur with probability at most $\delta/4$, so the claimed bounds hold
 1143 simultaneously with probability at least $1 - \delta$.

1144 C.2 Supporting Hellinger lemmas

1145 We collect the proofs of the four Hellinger lemmas used in the proof of Theorem 2.

1146 **Lemma 7** (Approximate minimizers are Hellinger-close). *Let π_0 be the MLE on a dataset of m*
 1147 *trajectories, and define $\gamma^* := (\log |\Pi| + \log(2/\delta))/m$. With probability at least $1 - \delta$:*

1148 *1. $\pi^* \in \Pi_{\text{version}}^{\gamma^*}$, i.e., $\text{LogLoss}(\pi^*) \leq \text{LogLoss}(\pi_0) + \gamma^*$.*

1149 2. Every $\pi \in \Pi$ with $\text{LogLoss}(\pi) \leq \text{LogLoss}(\pi_0) + \gamma$ satisfies $D_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi^*}) \leq \gamma/2 + \gamma^*$.

1150 *Proof.* Let $L(\pi) = \prod_{k=1}^m \prod_{h=1}^H \frac{\pi_h(a_h^{(k)} | s_h^{(k)})}{\pi_h^*(a_h^{(k)} | s_h^{(k)})}$ denote the likelihood ratio of the dataset under π versus
 1151 π^* . For a single trajectory, the likelihood ratio under $\mathcal{P}_M^{\pi^*}$ sums only the \mathcal{P}_M^π -mass on trajectories in
 1152 the support of $\mathcal{P}_M^{\pi^*}$, so its expectation is at most 1 (and exactly 1 whenever $\mathcal{P}_M^\pi \ll \mathcal{P}_M^{\pi^*}$). Since $L(\pi)$
 1153 is a product of m independent trajectory likelihood ratios, we have $\mathbb{E}[L(\pi)] \leq 1$ for all π . Thus, by
 1154 Markov's inequality and a union bound,

$$\Pr \left[\sup_{\pi \in \Pi} L(\pi) > \frac{2|\Pi|}{\delta} \right] \leq \frac{\delta}{2}.$$

1155 Taking the log and dividing by m , with probability at least $1 - \delta/2$ the log-loss gap satisfies
 1156 $\text{LogLoss}(\pi^*) \leq \text{LogLoss}(\pi_0) + \gamma^*$, proving Part 1.

1157 For Part 2, if π is a γ -approximate minimizer then $\text{LogLoss}(\pi) \leq \text{LogLoss}(\pi^*) + \gamma$ (since π_0 is the
 1158 MLE), which rearranges to

$$-\frac{1}{m} \sum_{k=1}^m \log \left(\prod_{h=1}^H \frac{\pi_h(a_h^{(k)} | s_h^{(k)})}{\pi_h^*(a_h^{(k)} | s_h^{(k)})} \right) \leq \gamma \implies \sum_{k=1}^m \log \sqrt{\prod_{h=1}^H \frac{\pi_h(a_h^{(k)} | s_h^{(k)})}{\pi_h^*(a_h^{(k)} | s_h^{(k)})}} \geq -\frac{m\gamma}{2} \quad (8)$$

$$\implies \prod_{k=1}^m \sqrt{\prod_{h=1}^H \frac{\pi_h(a_h^{(k)} | s_h^{(k)})}{\pi_h^*(a_h^{(k)} | s_h^{(k)})}} \geq e^{-m\gamma/2}. \quad (9)$$

1159 By definition, the expectation of each factor in the outer product equals $1 - D_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi^*})$, so the
 1160 LHS has expectation $(1 - D_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi^*}))^m \leq e^{-mD_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi^*})}$. For any π with $D_H^2(\mathcal{P}_M^\pi, \mathcal{P}_M^{\pi^*}) >$
 1161 $\gamma/2 + \gamma^*$, Markov's inequality applied to (9) gives failure probability at most $\delta/(2|\Pi|)$. A union
 1162 bound over Π completes the proof with overall probability at least $1 - \delta$. \square

1163 **Lemma 8** (Hellinger tensorization under geometric mixture). *Let \mathcal{P}^π and $\mathcal{P}^{\pi'}$ denote the trajectory*
 1164 *distributions under policies π and π' in an MDP M . Define the geometric mixture policy by*
 1165 $\bar{\pi}_h(a | s) := \sqrt{\pi_h(a | s) \pi_h'(a | s)} / Z_h(s)$, where $Z_h(s) := \sum_a \sqrt{\pi_h(a | s) \pi_h'(a | s)}$.⁹ Then

$$1 - D_H^2(\mathcal{P}^\pi, \mathcal{P}^{\pi'}) = \mathbb{E}_{T \sim \mathcal{P}^\pi} \left[\prod_{h=1}^H (1 - d_H^2(\pi_h(\cdot | s_h), \pi_h'(\cdot | s_h))) \right].$$

1166 *Proof.* By definition, the trajectory-level squared Hellinger affinity is

$$1 - D_H^2(\mathcal{P}^\pi, \mathcal{P}^{\pi'}) = \sum_T \sqrt{\mathcal{P}^\pi(T) \mathcal{P}^{\pi'}(T)}.$$

1167 Expanding each trajectory probability,

$$\sqrt{\mathcal{P}^\pi(T) \mathcal{P}^{\pi'}(T)} = P_0(s_1) \left(\prod_{h=1}^{H-1} P_h(s_{h+1} | s_h, a_h) \right) \left(\prod_{h=1}^H \sqrt{\pi_h(a_h | s_h) \pi_h'(a_h | s_h)} \right).$$

1168 Let $Z_h(s_h) := \sum_a \sqrt{\pi_h(a | s_h) \pi_h'(a | s_h)} = 1 - d_H^2(\pi_h(\cdot | s_h), \pi_h'(\cdot | s_h))$. With the convention
 1169 in the lemma statement, the identity $\sqrt{\pi_h(a_h | s_h) \pi_h'(a_h | s_h)} = \bar{\pi}_h(a_h | s_h) Z_h(s_h)$ holds for
 1170 every s_h, a_h : it is the definition when $Z_h(s_h) > 0$, and both sides are zero when $Z_h(s_h) = 0$.
 1171 Substituting,

$$\sqrt{\mathcal{P}^\pi(T) \mathcal{P}^{\pi'}(T)} = \mathcal{P}^{\bar{\pi}}(T) \prod_{h=1}^H Z_h(s_h).$$

1172 Summing over all trajectories T gives the result. \square

⁹If $Z_h(s) = 0$, define $\bar{\pi}_h(\cdot | s)$ arbitrarily. In this case $\pi_h(\cdot | s)$ and $\pi_h'(\cdot | s)$ have disjoint support, so
 $\sqrt{\pi_h(a | s) \pi_h'(a | s)} = 0$ for every a , and the identities below still hold.

1173 **Lemma 9** (Geometric mixture Hellinger bound). *Let $\pi, \pi', \bar{\pi}$ be as in Lemma 8 and let $\gamma =$*
 1174 $D_H^2(\mathcal{P}^\pi, \mathcal{P}^{\pi'})$. Then

$$D_H^2(\mathcal{P}^{\bar{\pi}}, \mathcal{P}^{\pi'}) \leq \frac{3}{2} D_H^2(\mathcal{P}^\pi, \mathcal{P}^{\pi'}).$$

1175 *Proof.* Define

$$P := \mathcal{P}^\pi, \quad Q := \mathcal{P}^{\pi'}, \quad \bar{P} := \mathcal{P}^{\bar{\pi}}.$$

1176 From the proof of Lemma 8, we have that pointwise $\bar{P}(T)Z(T) = \sqrt{P(T)Q(T)}$ where $Z(T) =$
 1177 $\prod_{h=1}^H Z_h(s_h) \leq 1$. Therefore $\bar{P}(T) \geq \sqrt{P(T)Q(T)}$, and

$$1 - D_H^2(\bar{P}, Q) = \sum_T \sqrt{\bar{P}(T)Q(T)} \geq \sum_T \sqrt{\sqrt{P(T)Q(T)} \cdot Q(T)} = \sum_T P(T)^{1/4} Q(T)^{3/4}.$$

1178 Applying Hölder's inequality with conjugate exponents $3/2$ and 3 , we obtain

$$\begin{aligned} \sum_T P(T)^{1/2} Q(T)^{1/2} &= \sum_T (P(T)^{1/4} Q(T)^{3/4})^{2/3} P(T)^{1/3} \\ &\leq \left(\sum_T [(P(T)^{1/4} Q(T)^{3/4})^{2/3}]^{3/2} \right)^{2/3} \left(\sum_T (P(T)^{1/3})^3 \right)^{1/3} \\ &= \left(\sum_T P(T)^{1/4} Q(T)^{3/4} \right)^{2/3} \left(\sum_T P(T) \right)^{1/3} \\ &= \left(\sum_T P(T)^{1/4} Q(T)^{3/4} \right)^{2/3}. \end{aligned}$$

1179 Hence $1 - D_H^2(\mathcal{P}^\pi, \mathcal{P}^{\pi'}) \leq (\sum_T P(T)^{1/4} Q(T)^{3/4})^{2/3}$. Raising to the $3/2$ power,
 1180 $\sum_T P(T)^{1/4} Q(T)^{3/4} \geq (1 - D_H^2(\mathcal{P}^\pi, \mathcal{P}^{\pi'}))^{3/2}$, so by Bernoulli's inequality:

$$D_H^2(\bar{P}, Q) \leq 1 - (1 - D_H^2(\mathcal{P}^\pi, \mathcal{P}^{\pi'}))^{3/2} \leq \frac{3}{2} D_H^2(\mathcal{P}^\pi, \mathcal{P}^{\pi'}).$$

1181 □

1182 **Lemma 10** (Stopped trajectory Hellinger bound). *Let $\pi_0, \pi' \in \Pi$ and $\Phi \subseteq \Pi$. Let $\tau := \tau_{\pi_0, \Phi}^\theta(T)$. If*
 1183 Φ satisfies

$$\Pr_{T \sim (\pi_0, N)}(\tau_{\pi_0, \{\pi'\}}^\theta(T) < \tau) < \eta \quad \text{and} \quad \Pr_{T \sim (\pi', N)}(\tau_{\pi_0, \{\pi'\}}^\theta(T) < \tau) < \eta'$$

1184 then, with probability at least $1 - \eta$ over $T \sim (\pi_0, N)$, the cumulative per-step Hellinger distance up
 1185 to τ is bounded by θ :

$$\sum_{h=1}^{\tau-1} d_H^2(\pi'_h(\cdot | s_h), \pi_{0,h}(\cdot | s_h)) \leq \theta.$$

1186 Furthermore, the stopped trajectory distributions satisfy

$$D_H^2(\mathcal{P}_N^{\pi_0 | \tau}, \mathcal{P}_N^{\pi' | \tau}) \leq \theta + \frac{\eta + \eta'}{2}.$$

1187 *Proof.* Let

$$\tau' := \tau_{\pi_0, \{\pi'\}}^\theta(T), \quad \mathcal{E} := \{\tau \leq \tau'\}.$$

1188 By hypothesis, $\Pr_{(\pi_0, N)}(\mathcal{E}^c) < \eta$, so $\Pr(\mathcal{E}) \geq 1 - \eta$. Conditioned on \mathcal{E} , the sum

$$\sum_{h=1}^{\tau-1} d_H^2(\pi'_h(\cdot | s_h), \pi_{0,h}(\cdot | s_h)) \leq \theta$$

1189 by the definition of τ' and monotonicity of the sum.

1190 For the second claim, define $\tau^*(T) := \min(\tau, \tau')$. Since $\tau^* \leq \tau'$ on every trajectory, we have
 1191 $\sum_{h < \tau^*} d_H^2(\pi'_h(\cdot | s_h), \pi_{0,h}(\cdot | s_h)) \leq \theta$ pathwise. The stopped affinity therefore satisfies

$$\begin{aligned} \sum_T \sqrt{\mathcal{P}_N^{\pi_{0|\tau^*}}(T) \mathcal{P}_N^{\pi'_{|\tau^*}}(T)} &= \mathbb{E}_{T \sim \mathcal{P}^\pi} \left[\prod_{h < \tau^*} (1 - d_H^2(\pi'_h(\cdot | s_h), \pi_{0,h}(\cdot | s_h))) \right] \\ &\geq \mathbb{E}_{T \sim \mathcal{P}^\pi} \left[1 - \sum_{h < \tau^*} d_H^2(\pi'_h(\cdot | s_h), \pi_{0,h}(\cdot | s_h)) \right] \geq 1 - \theta, \end{aligned}$$

1192 so $D_H^2(\mathcal{P}_N^{\pi_{0|\tau^*}}, \mathcal{P}_N^{\pi'_{|\tau^*}}) \leq \theta$.

1193 Write $A(P, Q) = \sum_T \sqrt{P(T)Q(T)}$ for the Hellinger affinity. On \mathcal{E} , the stopped distributions at τ
 1194 and τ^* agree, so

$$A(\mathcal{P}_N^{\pi_{0|\tau}}, \mathcal{P}_N^{\pi'_{|\tau}}) \geq \sum_{T \in \mathcal{E}} \sqrt{\mathcal{P}_N^{\pi_{0|\tau^*}}(T) \mathcal{P}_N^{\pi'_{|\tau^*}}(T)}.$$

1195 On the other hand,

$$A(\mathcal{P}_N^{\pi_{0|\tau^*}}, \mathcal{P}_N^{\pi'_{|\tau^*}}) = \sum_{T \in \mathcal{E}} \sqrt{\mathcal{P}_N^{\pi_{0|\tau^*}}(T) \mathcal{P}_N^{\pi'_{|\tau^*}}(T)} + \sum_{T \in \mathcal{E}^c} \sqrt{\mathcal{P}_N^{\pi_{0|\tau^*}}(T) \mathcal{P}_N^{\pi'_{|\tau^*}}(T)}.$$

1196 By AM-GM, the second sum satisfies

$$\sum_{T \in \mathcal{E}^c} \sqrt{\mathcal{P}_N^{\pi_{0|\tau^*}}(T) \mathcal{P}_N^{\pi'_{|\tau^*}}(T)} \leq \sum_{T \in \mathcal{E}^c} \frac{\mathcal{P}_N^{\pi_{0|\tau^*}}(T) + \mathcal{P}_N^{\pi'_{|\tau^*}}(T)}{2} \leq \frac{\eta + \eta'}{2}$$

1197 Therefore $A(\mathcal{P}_N^{\pi_{0|\tau}}, \mathcal{P}_N^{\pi'_{|\tau}}) \geq A(\mathcal{P}_N^{\pi_{0|\tau^*}}, \mathcal{P}_N^{\pi'_{|\tau^*}}) - \frac{\eta + \eta'}{2}$, which gives

$$D_H^2(\mathcal{P}_N^{\pi_{0|\tau}}, \mathcal{P}_N^{\pi'_{|\tau}}) \leq D_H^2(\mathcal{P}_N^{\pi_{0|\tau^*}}, \mathcal{P}_N^{\pi'_{|\tau^*}}) + (\eta + \eta')/2 \leq \theta + \frac{\eta + \eta'}{2}$$

1198 □

1199 C.3 Proof of Lemma 2

1200 Because π_{sw} and π^* are valid policies evaluated over the full-horizon MDP N , we can directly invoke
 1201 Theorem 3.1 of Foster et al. [2024] to bound their regret difference:

$$J(\pi_{\text{sw}}) - J(\pi^*) \leq \sqrt{6\sigma_{\pi^*}^2 \cdot D_H^2(\mathcal{P}^{\pi_{\text{sw}}}, \mathcal{P}^{\pi^*})} + O(C_{\max} \log(C_{\max} \epsilon^{-1})) \cdot D_H^2(\mathcal{P}^{\pi_{\text{sw}}}, \mathcal{P}^{\pi^*}) + \epsilon,$$

1202 where $\mathcal{P}^{\pi_{\text{sw}}}$ and \mathcal{P}^{π^*} are the full-horizon trajectory distributions of π_{sw} and π^* , respectively.

1203 To evaluate the full-horizon Hellinger distance, we expand the Hellinger affinity between $\mathcal{P}^{\pi_{\text{sw}}}$ and
 1204 \mathcal{P}^{π^*} . Recall that the squared Hellinger distance satisfies $1 - D_H^2(P, Q) = \sum_T \sqrt{P(T)Q(T)}$. Let $T_{|\tau}$
 1205 denote the prefix of the trajectory up to the stopping time τ , and let $T_{>\tau}$ denote the suffix. Therefore,
 1206 the full-horizon trajectory distribution for any policy π factors as $P^\pi(T) = P^\pi(T_{|\tau})P^\pi(T_{>\tau} | T_{|\tau})$.

1207 By construction, the mixed policy π_{sw} executes the learner $\hat{\pi}$ up to τ and the expert π^* for all steps
 1208 $h > \tau$. Thus, its prefix distribution is exactly the stopped learner distribution, $\mathcal{P}^{\pi_{\text{sw}}}(T_{|\tau}) = \mathcal{P}^{\hat{\pi}_{|\tau}}(T_{|\tau})$.
 1209 Crucially, because τ is a stopping time, the event $\{\tau = t\}$ depends exclusively on the trajectory prefix
 1210 $T_{1:t}$. Consequently, for any prefix where the switch occurs at step t , the conditional distribution of
 1211 the remaining trajectory is governed entirely by the expert policy, yielding $\mathcal{P}^{\pi_{\text{sw}}}(T_{t+1:H} | T_{1:t}) =$
 1212 $\mathcal{P}^{\pi^*}(T_{t+1:H} | T_{1:t})$. Thus,

$$\begin{aligned} 1 - D_H^2(\mathcal{P}^{\pi_{\text{sw}}}, \mathcal{P}^{\pi^*}) &= \sum_{T_{|\tau}} \sum_{T_{>\tau}} \sqrt{\mathcal{P}^{\pi_{\text{sw}}}(T_{|\tau}) \mathcal{P}^{\pi_{\text{sw}}}(T_{>\tau} | T_{|\tau}) \cdot \mathcal{P}^{\pi^*}(T_{|\tau}) \mathcal{P}^{\pi^*}(T_{>\tau} | T_{|\tau})} \\ &= \sum_{T_{|\tau}} \sqrt{\mathcal{P}^{\hat{\pi}_{|\tau}}(T_{|\tau}) \mathcal{P}^{\pi^*}(T_{|\tau})} \sum_{T_{>\tau}} \mathcal{P}^{\pi^*}(T_{>\tau} | T_{|\tau}). \end{aligned}$$

1213 Because $P^{\pi^*}(T_{>\tau} | T_{|\tau})$ is a valid probability distribution over the suffix trajectories, the inner sum
 1214 evaluates exactly to 1 for every prefix $T_{|\tau}$. The remaining sum is exactly the Hellinger affinity of the
 1215 stopped trajectories. Therefore, the full-horizon trajectory distance perfectly collapses to the distance
 1216 accumulated up to the stopping time:

$$D_H^2(\mathcal{P}^{\pi_{sw}}, \mathcal{P}^{\pi^*}) = D_H^2(\mathcal{P}^{\hat{\pi}_{|\tau}}, \mathcal{P}^{\pi^*}). \quad (10)$$

1217 Substituting (10) into the bound from Theorem 3.1 of Foster et al. [2024] concludes the proof.

1218 C.4 Single-step lower bound for stochastic policies

1219 How tight are these rates? The upper bound of Corollary 3 has two regimes: a variance-driven $\tilde{O}(\epsilon^{-5})$
 1220 term and a cost-driven $\tilde{O}(\epsilon^{-3})$ term. The next result shows that the cubic dependence is already
 1221 unavoidable in the single-step setting, matching the cost-driven part.

1222 **Theorem 4** (Stochastic PQ lower bound). *There exists a constant $c_0 > 0$ such that:*

1223 *For any dimension $d \geq 1$ and tolerance $\epsilon \in (0, 1/16]$, there exist a state space \mathcal{X} , a stochastic policy*
 1224 *class Π of log-cardinality d , and fixed known distributions P and Q over \mathcal{X} for which the following*
 1225 *holds: For any proper¹⁰ selective learner, if the labeled sample size is $m \leq c_0 \frac{d}{\epsilon^3}$, then there exists an*
 1226 *expert $\pi^* \in \Pi$ and a bounded cost function $c : \mathcal{X} \times \{0, 1\} \rightarrow [0, 1]$ such that*

$$\mathbb{E}[\alpha_P(\pi, \tau)] > \epsilon \quad \text{or} \quad \mathbb{E}[\text{Regret}_Q(\pi, \tau; c)] > \epsilon$$

1227 *where the outer expectation is over the labeled sample and the learner's internal randomization,*
 1228 *which in particular determine π and τ .*

1229 Theorem 4 shows that the ϵ^{-3} dependence in labeled sample complexity is already unavoidable in
 1230 the single-step stochastic setting. Crucially, because this lower bound is established in a single-step
 1231 environment ($H = 1$), there is no future horizon over which errors can accumulate after an abstention.
 1232 Consequently, the stopped regret and the switched regret are equal in this construction, meaning
 1233 this lower bound applies to both formulations. It therefore matches the C_{\max} -driven $\tilde{O}(\epsilon^{-3})$ portion
 1234 of Corollary 3, while leaving a gap to the current variance-driven $\tilde{O}(\epsilon^{-5})$ labeled-sample rate. For
 1235 deterministic classes, the lower bound of Goldwasser et al. [2020] for PQ learning implies that the
 1236 $\tilde{O}(1/\epsilon^2)$ labeled complexity is tight. Whether the sequential stochastic upper bound can be improved,
 1237 or whether stronger lower bounds are possible, remains open.

1238 C.5 Proof of Theorem 4

1239 Let the state space be $\mathcal{X} = \{1, 2, \dots, N\}$. Let P be uniform over \mathcal{X} , Q be uniform over a known
 1240 subset $S_Q = \{1, \dots, d\}$, and let $\Delta \in [0, 1/2]$ be a parameter to be chosen later. Let Π be the class
 1241 of stochastic policies parameterized by $\sigma \in \{-1, 1\}^d$: for $x \in S_Q$, the expert π^* takes action 1
 1242 with probability $1/2 + \sigma_x \Delta$; for $x \notin S_Q$, the probability is $1/2$. Similarly, define the cost function
 1243 $c_\sigma : \mathcal{X} \times \{0, 1\} \rightarrow [0, 1]$ by $c_\sigma(x, 1) = 1\{\sigma_x = -1\}$ and $c_\sigma(x, 0) = 1\{\sigma_x = 1\}$ for $x \in S_Q$ and 0
 1244 otherwise.

1245 We now apply Yao's minimax principle. It is enough to lower bound the expected performance
 1246 of a deterministic proper learner under the uniform prior $\sigma \sim \text{Unif}(\{-1, 1\}^d)$. Let D be the m
 1247 labeled samples from P the learner observes. Since the learner is deterministic, after observing D ,
 1248 it will, for each state x , choose an acceptance probability $\alpha_x(D) \in [0, 1]$ and a bernoulli parameter
 1249 $q_x(D) = \hat{\pi}(1|x)$. Crucially, because the learner is proper, it must output a policy in Π , meaning
 1250 $q_x(D) \in \{1/2 - \Delta, 1/2 + \Delta\}$.

1251 For $x \in S_Q$, define the one-step selective regret at state x as

$$\begin{aligned} R_x(D, \sigma) &:= \mathbb{E}_{a \sim \hat{\pi}_D(\cdot|x)}[c_\sigma(x, a)] - \mathbb{E}_{a \sim \pi_\sigma(\cdot|x)}[c_\sigma(x, a)] \\ &= \Delta + \sigma_x \left(\frac{1}{2} - q_x(D) \right). \end{aligned} \quad (\text{direct computation})$$

¹⁰By proper selective learner, we mean that the base policy π_0 always lies in Π . This holds for SeqRejectron. In the batch (i.e., single step) setting this holds, for example, for Rejectron [Goldwasser et al., 2020] and the Slice-and-Dice algorithm of Kalai and Kanade [2021].

1252 Now, define $\eta_x = \mathbb{E}[\sigma_x | D]$. By conditioning on D , we get

$$\mathbb{E}[R_x(D, \sigma) | D] = \Delta + \eta_x \left(\frac{1}{2} - q_x(D) \right) \geq \Delta - \Delta |\eta_x|$$

1253 where we used the proper learner constraint $|\frac{1}{2} - q_x(D)| = \Delta$.

1254 Multiplying by α_x and then taking the expectation over D gives

$$\mathbb{E}[\alpha_x(D) R_x(D, \sigma)] \geq \Delta \mathbb{E}[\alpha_x(D)] - \Delta \mathbb{E}[|\eta_x|]$$

1255 where the full expectation is now over σ and D .

1256 Let $n_x(D) \sim \text{Binomial}(m, 1/N)$ be the number of times state x is observed. The posterior bias
1257 satisfies

$$\begin{aligned} \mathbb{E}[|\eta_x(D)| | n_x(D) = k] &= \text{TV} \left(\text{Binomial} \left(k, \frac{1}{2} + \Delta \right), \text{Binomial} \left(k, \frac{1}{2} - \Delta \right) \right) \\ &\leq 2\sqrt{2}\Delta\sqrt{k} \end{aligned}$$

1258 by Pinsker's inequality. Taking expectations yields $\mathbb{E}[|\eta_x(D)|] \leq 2\sqrt{2}\Delta\mathbb{E}[\sqrt{n_x(D)}]$. Then,
1259 averaging over the possibilities of x gives

$$\mathbb{E}[\text{Regret}_Q(\hat{\pi}, \tau; c)] \geq \Delta \mathbb{E}[\bar{\alpha}(D)] - \frac{2\sqrt{2}\Delta^2}{d} \sum_{x=1}^d \mathbb{E}[\sqrt{n_x(D)}]$$

1260 where $\bar{\alpha}(D) = \frac{1}{d} \sum_{x=1}^d \alpha_x(D)$. Jensen's inequality on the second expectation then yields

$$\mathbb{E}[\text{Regret}_Q(\hat{\pi}, \tau; c)] \geq \Delta \mathbb{E}[\bar{\alpha}(D)] - 2\sqrt{2}\Delta^2 \sqrt{\frac{m}{N}}.$$

1261 We now bound the other expectation. Since P is uniform, $\mathbb{E}_{x \sim P}[1 - \alpha_x(D)] = 1 - \frac{1}{N} \sum_{x=1}^N \alpha_x(D)$.

1262 Thus, if the learner satisfies $\mathbb{E}[\mathbb{E}_{x \sim P}[1 - \alpha_x(D)]] \leq \epsilon$, then $\frac{1}{N} \sum_{x=1}^N \mathbb{E}[\alpha_x(D)] \geq 1 - \epsilon$. Even if
1263 the learner accepts with probability 1 on all states outside S_Q , this still forces $\mathbb{E}[\bar{\alpha}(D)] \geq 1 - \frac{N\epsilon}{d}$.

1264 Now, setting $N = \frac{d}{2\epsilon}$ and $\Delta = 4\epsilon$ gives us the bounds

$$E[\text{Regret}_Q(\hat{\pi}, \tau; c)] \geq 2\epsilon - 32\sqrt{2}\epsilon^2 \sqrt{\frac{2\epsilon m}{d}}.$$

1265 Requiring this expected regret to be at most ϵ yields:

$$1 \leq 32\sqrt{2}\epsilon \sqrt{\frac{2\epsilon m}{d}} \implies m \geq \frac{d}{4096\epsilon^3}.$$

1266 Setting $c_0 = 1/4096$ completes the proof.

1267 D Proofs from the Misspecification Section

1268 **Notation.** For a base policy $\pi_0 \in \Pi$, comparator $\pi \in \Pi$, and validator sequence Φ , define the
1269 expert-side late-stop risk

$$r_N(\pi_0, \pi, \Phi) := \Pr_{T \sim P_N^*} [\tau_{\pi_0, \Phi}(T) > \tau_{\pi_0, \{\pi\}}(T)],$$

1270 with empirical analogue

$$\hat{r}_N(\pi_0, \pi, \Phi) := \frac{1}{n} \sum_{j=1}^n \mathbf{1} [\tau_{\pi_0, \Phi}(T_j) > \tau_{\pi_0, \{\pi\}}(T_j)],$$

1271 where T_1, \dots, T_n are the unlabeled expert test trajectories. This is the expert-side late-stop quantity
1272 that appears in the equilibrium lemma and in the generalization argument.

1273 **D.1 Proof of Lemma 3**

1274 Define a zero-sum game in which the maximizer chooses $p \in \Delta(\Pi)$ and the minimizer chooses
 1275 $q \in \Delta(\Pi^K)$, with payoff

$$U(p, q) := \mathbb{E}_{\pi \sim p} \mathbb{E}_{\Phi \sim q} \left[\widehat{r}_N(\pi_0, \pi, \Phi) - \Lambda \cdot \widehat{d}_M(\pi) + \frac{\Lambda}{K} \sum_{k=1}^K \widehat{d}_M(\phi_k) \right].$$

1276 Since U is bilinear, von Neumann's minimax theorem implies $\min_q \max_p U(p, q) =$
 1277 $\max_p \min_q U(p, q)$. To bound the maximin, fix any $p \in \Delta(\Pi)$ and let p^K denote the strategy
 1278 of forming $\Phi = (\phi_1, \dots, \phi_K)$ by drawing ϕ_1, \dots, ϕ_K i.i.d. from p . Then,

$$\begin{aligned} \min_q U(p, q) &\leq U(p, p^K) = \mathbb{E}_{\pi \sim p, \Phi \sim p^K} [\widehat{r}_N(\pi_0, \pi, \Phi)] - \Lambda \cdot \mathbb{E}_{\pi \sim p} [\widehat{d}_M(\pi)] + \frac{\Lambda}{K} \sum_{k=1}^K \mathbb{E}_{\phi_k \sim p} [\widehat{d}_M(\phi_k)] \\ &= \mathbb{E}_{\pi \sim p, \Phi \sim p^K} [\widehat{r}_N(\pi_0, \pi, \Phi)] \quad (\text{the other two terms cancel}) \\ &\leq \frac{1}{K+1}. \quad (\text{same exchangeability argument as Lemma 5}) \end{aligned}$$

1279 Since this holds for every $p \in \Delta(\Pi)$, we conclude $\min_q \max_p U(p, q) = \max_p \min_q U(p, q) \leq$
 1280 $\frac{1}{K+1}$. Hence, there exists a minimizer strategy $q^* \in \Delta(\Pi^K)$ such that $\max_p U(p, q^*) \leq \frac{1}{K+1}$.
 1281 Because this bound holds against any p , it holds in particular for the pure strategies δ_π for any $\pi \in \Pi$.

1282 To prove the first part of Lemma 3, take $p = \delta_{\pi_0}$ to be the pure strategy playing π_0 . Since π_0 never
 1283 deviates from itself, the late-stop risk against any committee is exactly zero ($\widehat{r}_N = 0$). Thus,

$$\begin{aligned} U(\delta_{\pi_0}, q^*) &= 0 - \Lambda \cdot \widehat{d}_M(\pi_0) + \frac{\Lambda}{K} \mathbb{E}_{\Phi \sim q^*} \left[\sum_{k=1}^K \widehat{d}_M(\phi_k) \right] \leq \frac{1}{K+1} \\ \implies \mathbb{E}_{\Phi \sim q^*} \left[\sum_{k=1}^K \widehat{d}_M(\phi_k) \right] &\leq K \cdot \widehat{d}_M(\pi_0) + \frac{1}{\Lambda}. \end{aligned}$$

1284 To prove the second part of Lemma 3, take $p = \delta_\pi$ for an arbitrary target policy $\pi \in \Pi$. We have

$$\begin{aligned} U(\delta_\pi, q^*) &= \mathbb{E}_{\Phi \sim q^*} [\widehat{r}_N(\pi_0, \pi, \Phi)] - \Lambda \cdot \widehat{d}_M(\pi) + \frac{\Lambda}{K} \mathbb{E}_{\Phi \sim q^*} \left[\sum_{k=1}^K \widehat{d}_M(\phi_k) \right] \leq \frac{1}{K+1} \\ \implies \mathbb{E}_{\Phi \sim q^*} [\widehat{r}_N(\pi_0, \pi, \Phi)] &\leq \Lambda \left(\widehat{d}_M(\pi) - \frac{1}{K} \mathbb{E}_{\Phi \sim q^*} \left[\sum_{k=1}^K \widehat{d}_M(\phi_k) \right] \right) + \frac{1}{K+1} \\ \implies \mathbb{E}_{\Phi \sim q^*} [\widehat{r}_N(\pi_0, \pi, \Phi)] &\leq \Lambda \left(\widehat{d}_M(\pi) - \widehat{d}_M(\pi_0) \right) + \frac{1}{K}. \quad (\pi_0 \text{ minimizes } \widehat{d}_M) \end{aligned}$$

1285 **D.2 Proof of Theorem 3**

1286 The proof closely follows the proof of Theorem 1 (Section B.4); we highlight the key differences and
 1287 omit steps that are identical.

1288 **Existence of a sparse regularized validator distribution.** For parameters Λ, K , let q^* be the
 1289 distribution guaranteed by Lemma 3. Throughout, we let Λ, K be free and optimize them at the end.

1290 **Bounding the empirical risk.** Unlike the deterministic non-misspecified setting—which draws
 1291 k independent committees to suppress test risk and relies on the version space to prevent source
 1292 abstention—here we bound the expected empirical metrics over a single regularized sequence $\Phi \sim q^*$
 1293 (recall that this is done so as to not amplify the irreducible error).

1294 For the source data, the policy abstains if any validator in Φ disagrees with π_0 . Combining the union
 1295 bound with the the first part of Lemma 3 bounds the expected empirical abstention:

$$\mathbb{E}_{\Phi \sim q^*} [\widehat{\alpha}_M(\pi_0, \tau_{\pi_0, \Phi})] \leq \mathbb{E}_{\Phi \sim q^*} \left[\sum_{k=1}^K \widehat{d}_M(\phi_k) \right] \leq K \cdot \widehat{d}_M(\pi_0) + \frac{1}{\Lambda}. \quad (11)$$

1296 For the test data, evaluating against any comparator $\pi \in \Pi$, the second part of Lemma 3 directly
 1297 bounds the expected empirical late-stop risk by the empirical suboptimality of π . In particular, taking
 1298 our comparator to be $\tilde{\pi} = \arg \min_{\pi \in \Pi} \Delta_\pi$,

$$\mathbb{E}_{\Phi \sim q^*} [\hat{r}_N(\pi_0, \tilde{\pi}, \Phi)] \leq \Lambda \left(\hat{d}_M(\tilde{\pi}) - \hat{d}_M(\pi_0) \right) + \frac{1}{K}. \quad (12)$$

1299 **Bounding the generalization error on the expert's distribution.** We establish uniform convergence
 1300 over all $|\Pi|^{K+1}$ possible combinations of base policy π_0 and validator set Φ , exactly as in the
 1301 non-misspecified proof. Unlike that proof, the empirical source disagreements and target late-stop
 1302 risks are no longer small due to misspecification, so we use standard Hoeffding bounds rather than
 1303 fast-rate bounds. Let $Z := (K+1) \log |\Pi| + \log(4/\delta)$. By Hoeffding's inequality and a union bound,
 1304 with probability at least $1 - \delta$, simultaneously for all (π_0, Φ) :

$$|d_M(\pi_0) - \hat{d}_M(\pi_0)| \leq \sqrt{\frac{Z}{2m}}, \quad |r_N(\pi_0, \tilde{\pi}, \Phi) - \hat{r}_N(\pi_0, \tilde{\pi}, \Phi)| \leq \sqrt{\frac{Z}{2n}}. \quad (13)$$

1305 For the source data, the true stopping rate satisfies $\alpha_M(\pi_0, \tau_{\pi_0, \Phi}) \leq \sum_{k=1}^K d_M(\phi_k)$ by a union
 1306 bound. Taking the expectation over $\Phi \sim q^*$ and applying (13) for each $\phi_k \in \Phi$,

$$\begin{aligned} \mathbb{E}_{\Phi \sim q^*} [\alpha_M(\pi_0, \tau_{\pi_0, \Phi})] &\leq \mathbb{E}_{\Phi \sim q^*} \left[\sum_{k=1}^K \hat{d}_M(\phi_k) \right] + K \sqrt{\frac{Z}{2m}} && \text{(by (13))} \\ &\leq K \cdot \hat{d}_M(\pi_0) + \frac{1}{\Lambda} + K \sqrt{\frac{Z}{2m}} && \text{(by (11))} \\ &\leq K \cdot \hat{d}_M(\tilde{\pi}) + \frac{1}{\Lambda} + K \sqrt{\frac{Z}{2m}} && (\pi_0 \text{ minimizes } \hat{d}_M) \\ &\leq K\Delta + \frac{1}{\Lambda} + 2K \sqrt{\frac{Z}{2m}}. && \text{(by (13) and definition of } \Delta) \end{aligned}$$

1307 For the test data,

$$\begin{aligned} \mathbb{E}_{\Phi \sim q^*} [r_N(\pi_0, \tilde{\pi}, \Phi)] &\leq \mathbb{E}_{\Phi \sim q^*} [\hat{r}_N(\pi_0, \tilde{\pi}, \Phi)] + \sqrt{\frac{Z}{2n}} && \text{(by (13))} \\ &\leq \Lambda \left(\hat{d}_M(\tilde{\pi}) - \hat{d}_M(\pi_0) \right) + \frac{1}{K} + \sqrt{\frac{Z}{2n}} && \text{(by (12))} \\ &\leq \Lambda \left(d_M(\tilde{\pi}) - d_M(\pi_0) \right) + \frac{1}{K} + 2\Lambda \sqrt{\frac{Z}{2m}} + \sqrt{\frac{Z}{2n}} && \text{(by (13))} \\ &\leq \Lambda\Delta + \frac{1}{K} + 2\Lambda \sqrt{\frac{Z}{2m}} + \sqrt{\frac{Z}{2n}}, && \text{(definition of } \Delta, \text{ and } -d_M(\pi_0) \leq 0) \end{aligned}$$

1308 where the last step uses $d_M(\tilde{\pi}) - d_M(\pi_0) \leq \Delta_{\tilde{\pi}} \leq \Delta$.

1309 **Bounding the stopping rate and regret on the learned policy's distribution.** For the source data,
 1310 we apply the same union bound and prefix coupling argument as in the proof of Theorem 1. The
 1311 one difference is that π_0 may now deviate from π^* , contributing an additional $d_M(\pi_0)$ term. Since
 1312 $d_M(\pi_0) \leq d_M(\tilde{\pi}) + 2\sqrt{Z/2m} \leq \Delta + 2\sqrt{Z/2m}$ by (13) and the empirical minimality of π_0 , taking
 1313 the expectation over q^* yields:

$$\begin{aligned} \mathbb{E}_{\Phi \sim q^*} [\alpha_M(\pi_0, \tau_{\pi_0, \Phi})] &\leq \mathbb{E}_{\Phi \sim q^*} [\alpha_M(\pi^*, \tau_{\pi_0, \Phi})] + d_M(\pi_0) \\ &\leq \left(K\Delta + \frac{1}{\Lambda} + 2K \sqrt{\frac{Z}{2m}} \right) + \Delta + 2\sqrt{\frac{Z}{2m}} \\ &= (K+1)\Delta + \frac{1}{\Lambda} + 2(K+1) \sqrt{\frac{Z}{2m}}. \end{aligned}$$

1314 For the target data, let $B := \{\tau_{\pi_0, \{\pi^*\}} < \tau_{\pi_0, \Phi}\}$ be the event that the deployed policy π_0 deviates from
 1315 the true expert π^* before the validator set stops execution. As in the realizable case, the event B is
 1316 \mathcal{G}_{h^*} -measurable with $h^* = \min(\tau_{\pi_0, \Phi}, \tau_{\pi_0, \{\pi^*\}})$. Lemma 4 therefore gives $\Pr_{N, \pi_0}[B] = \Pr_{N, \pi^*}[B]$.

1317 To bound this probability on the expert’s distribution, we compare π^* to the best-in-class policy $\tilde{\pi}$.
 1318 On any expert trajectory, if π_0 deviates from π^* before Φ triggers (event B), then either (1) $\tilde{\pi}$ deviates
 1319 from π^* at some point during the trajectory, or (2) $\tilde{\pi}$ perfectly matches π^* everywhere, meaning the
 1320 deviation between π_0 and π^* is simultaneously a deviation between π_0 and $\tilde{\pi}$, and this deviation
 1321 occurs before Φ triggers. This implies the set containment $B \subseteq \{\tau_{\tilde{\pi}, \{\pi^*\}} \leq H\} \cup \{\tau_{\pi_0, \Phi} > \tau_{\pi_0, \{\tilde{\pi}\}}\}$
 1322 on expert trajectories. Applying a union bound, we have $\Pr_{N, \pi^*}[B] \leq d_N(\tilde{\pi}) + r_N(\pi_0, \tilde{\pi}, \Phi)$.

1323 As in the proof of Theorem 1, the stopped regret is bounded by C_{\max} times the probability of an
 1324 unrejected deviation from the expert, so taking the expectation over q^* yields:

$$\begin{aligned} \mathbb{E}_{\Phi \sim q^*} [\text{Regret}_N(\pi_0, \tau_{\pi_0, \Phi}; c)] &\leq C_{\max} \cdot \mathbb{E}_{\Phi \sim q^*} \left[\Pr_{N, \pi_0} [B] \right] \\ &= C_{\max} \cdot \mathbb{E}_{\Phi \sim q^*} \left[\Pr_{N, \pi^*} [B] \right] \\ &\leq C_{\max} \left(d_N(\tilde{\pi}) + \mathbb{E}_{\Phi \sim q^*} [r_N(\pi_0, \tilde{\pi}, \Phi)] \right) \\ &\leq C_{\max} \left(\Delta + \Lambda \Delta + \frac{1}{K} + 2\Lambda \sqrt{\frac{Z}{2m}} + \sqrt{\frac{Z}{2n}} \right). \end{aligned}$$

1325 Setting $\Lambda = K = \Theta(\Delta^{-1/2})$ balances the Δ -dependent and Δ -independent terms in both bounds,
 1326 yielding the rates in Theorem 3.

1327 D.3 Proof of Corollary 4

1328 Let $c_0 := \log |\Pi| + \log(4/\delta)$. We define the $\epsilon^* := \max \{(c_0/m)^{1/5}, (c_0/n)^{1/3}\}$ and set $K = \Lambda =$
 1329 $[(\Delta^{1/2} + \epsilon^*)^{-1}]$. Substituting this into Theorem 3,

$$K\Delta + \frac{1}{K} \leq \frac{\Delta}{\Delta^{1/2}} + (\Delta^{1/2} + \epsilon^*) = 2\Delta^{1/2} + \epsilon^*.$$

1330 For the sample estimation penalties, noting $C \leq (K+1)c_0$ and $K \leq (\epsilon^*)^{-1}$, the labeled penalty
 1331 is bounded by $(\epsilon^*)^{-3/2} \sqrt{c_0/m}$. By definition, $\epsilon^* \geq (c_0/m)^{1/5}$, which implies $\sqrt{c_0/m} \leq (\epsilon^*)^{5/2}$.
 1332 Thus, the labeled penalty is at most $(\epsilon^*)^{-3/2} (\epsilon^*)^{5/2} = \epsilon^*$. An identical balancing argument bounds
 1333 the unlabeled penalty: $K^{1/2} \sqrt{c_0/n} \leq (\epsilon^*)^{-1/2} (\epsilon^*)^{3/2} = \epsilon^*$.

1334 Because all terms in the abstention and regret bounds of Theorem 3 are upper-bounded by $O(\Delta^{1/2} +$
 1335 $\epsilon^*)$, expanding ϵ^* yields the claimed result.

1336 D.4 Off-policy test trajectories under misspecification

1337 So far, we have assumed that the expert π^* is the same across train and test environments. In practice,
 1338 however, train and test demonstrations often come from different sources—for example, using one
 1339 expert to pilot a simulated car, while using dashcam aggregated from multiple drivers. This creates a
 1340 challenging off-policy gap where the training behavior π^{tr} and test behavior π^{te} differ, and neither
 1341 necessarily lies in the policy class Π . Fortunately, the same idea we used for misspecified experts—
 1342 finding a validator distribution via a *regularized* game, which induces a stopping time—extends
 1343 readily to policy shift between environments.

1344 Concretely, suppose $\mathcal{D}_{\text{train}}$ is generated by a deterministic policy π^{tr} , while $\mathcal{D}_{\text{test}}$ is generated by a
 1345 different deterministic policy π^{te} , with neither policy necessarily in Π . Let d_M^{tr} and d_N^{te} denote the
 1346 corresponding disagreement probabilities in (M, π^{tr}) and (N, π^{te}) , respectively, and let $\text{Regret}_N^{\text{te}}$
 1347 denote regret against π^{te} in the test environment. Then, the agnostic SeqRejectron satisfies:

1348 **Proposition 4.** *There exists an algorithm, the agnostic SeqRejectron (which we described in*
 1349 *Section 5), with hyperparameters $K, \Lambda, \delta > 0$, with the following guarantee. Let $\Delta_{\text{off}} :=$
 1350 $\min_{\tilde{\pi} \in \Pi} \max \{d_M^{\text{tr}}(\tilde{\pi}), d_N^{\text{te}}(\tilde{\pi})\}$.¹¹ When run with a specific choice of parameters K, Λ (depending
 1351 on $m, n, \Delta_{\text{off}}, |\Pi|$, and δ), then with probability at least $1 - \delta$, the agnostic SeqRejectron will return
 1352 a pair (π_0, q^*) such that both $\mathbb{E}_{\Phi \sim q^*} [\alpha_M(\pi_0, \tau_{\pi_0, \Phi})]$ and $\mathbb{E}_{\Phi \sim q^*} [\text{Regret}_N^{\text{te}}(\pi_0, \tau_{\pi_0, \Phi}; c)]/C_{\max}$ are
 1353 $\tilde{O} \left(\Delta_{\text{off}}^{1/2} + (\log |\Pi|/m)^{1/5} + (\log |\Pi|/n)^{1/3} \right)$.*

¹¹Note that as long as there exists some policy $\pi \in \Pi$ which is able to behave similarly to π^{tr} on M and to π^{te} on N , then the irreducible error Δ_{off} will be small.

1354 **D.5 Proof of Proposition 4**

The proof is a direct adaptation of the proof of Theorem 3. We first redefine the suboptimality gaps and their empirical estimators with respect to the off-policy behavior policies π^{tr} and π^{te} instead of the expert π^* :

$$d_M^{\text{tr}}(\bar{\pi}) := \Pr_{M, \pi^{\text{tr}}} [\tau_{\bar{\pi}, \{\pi^{\text{tr}}\}} \leq H], \quad d_N^{\text{te}}(\bar{\pi}) := \Pr_{N, \pi^{\text{te}}} [\tau_{\bar{\pi}, \{\pi^{\text{te}}\}} \leq H]$$

$$\hat{d}_M^{\text{tr}}(\pi) := \frac{1}{m} \sum_{i=1}^m \mathbf{1}[\tau_{\pi, \{\pi^{\text{tr}}\}}(S_i) \leq H], \quad \hat{r}_N^{\text{te}}(\pi_0, \pi, \Phi) := \frac{1}{n} \sum_{j=1}^n \mathbf{1}[\tau_{\pi_0, \Phi}(T_j) > \tau_{\pi_0, \{\pi\}}(T_j)].$$

1355 Let $\pi_0 \in \arg \min_{\pi \in \Pi} \hat{d}_M^{\text{tr}}(\pi)$. Applying the symmetrically regularized equilibrium (Lemma 3) to
 1356 these modified losses guarantees the existence of a distribution q^* such that:

$$\mathbb{E}_{\Phi \sim q^*} \left[\sum_{k=1}^K \hat{d}_M^{\text{tr}}(\phi_k) \right] \leq K \hat{d}_M^{\text{tr}}(\pi_0) + \frac{1}{\Lambda},$$

$$\mathbb{E}_{\Phi \sim q^*} [\hat{r}_N^{\text{te}}(\pi_0, \pi, \Phi)] \leq \Lambda (\hat{d}_M^{\text{tr}}(\pi) - \hat{d}_M^{\text{tr}}(\pi_0)) + \frac{1}{K} \quad \text{for every } \pi \in \Pi.$$

1357 From here, the mechanics follow exactly as before. Applying Hoeffding bounds introduces the
 1358 $O(\sqrt{C/m})$ and $O(\sqrt{C/n})$ concentration penalties to these empirical terms. The prefix-coupling
 1359 argument bounds the true completeness α_M by the sum of training disagreements, and the test regret
 1360 by C_{\max} times the late-stop risk against π^{te} .

1361 Finally, to bound the regret against the training demonstrator π^{tr} in the test environment, coupling π^{tr}
 1362 to $\bar{\pi}$ and then $\bar{\pi}$ to π^{te} naturally introduces the discrepancy terms $d_N^{\text{tr}}(\bar{\pi})$ and $d_N^{\text{te}}(\bar{\pi})$ via the union
 1363 bound, yielding the secondary result.

1364 Finally, to obtain the optimal rates claimed in the proposition, we evaluate these bounds at the
 1365 comparator $\bar{\pi}_{\text{off}} \in \arg \min_{\bar{\pi} \in \Pi} \max\{d_M^{\text{tr}}(\bar{\pi}), d_N^{\text{te}}(\bar{\pi})\}$. Setting $K = \Lambda = \left[(\Delta_{\text{off}}^{1/2} + \epsilon^*)^{-1} \right]$ with
 1366 $\epsilon^* := \max\{(\log |\Pi|/m)^{1/5}, (\log |\Pi|/n)^{1/3}\}$ and applying the identical balancing argument from the
 1367 proof of Corollary 4 yields the final rate.