Offline Action-Free Learning of Ex-BMDPs by Comparing Diverse Datasets

Anonymous authors Paper under double-blind review

Keywords: representation learning, action-free RL, Ex-BMDP, controllable state representations

Summary

While sequential decision-making environments often involve high-dimensional observations, not all features of these observations are relevant for control. In particular, the observation space may capture factors of the environment which are not controllable by the agent, but which add complexity to the observation space. The need to ignore these "noise" features in order to operate in a tractably-small state space poses a challenge for efficient policy learning. Due to the abundance of video data available in many such environments, task-independent representation learning from action-free offline data offers an attractive solution. However, recent work has highlighted theoretical limitations in action-free learning under the Exogenous Block MDP (Ex-BMDP) model, where temporally-correlated noise features are present in the observations. To address these limitations, we identify a realistic setting where representation learning in Ex-BMDPs becomes tractable: when action-free video data from multiple agents with differing policies are available. Concretely, this paper introduces CRAFT (Comparisonbased Representations from Action-Free Trajectories), a sample-efficient algorithm leveraging differences in controllable feature dynamics across agents to learn representations. We provide theoretical guarantees for CRAFT's performance and demonstrate its feasibility on a toy example, offering a foundation for practical methods in similar settings.

Contribution(s)

1. We present a provably sample-efficient algorithm, CRAFT, that can learn high-accuracy latent state encoders under the Ex-BMDP model, when provided with two sets of offline observation trajectories, *without action labels*, that are collected by two agents with sufficiently-distinct policies.

Context: Misra et al. (2024) has shown that efficient representation learning in Ex-BMDPs using a single offline dataset without action labels is in general *not* possible. This work therefore represents to our knowledge the first *positive* theoretical result for this problem. The Ex-BMDP model was introduced by Efroni et al. (2022), who propose a provably sample-efficient algorithm for *online* representation learning in this model. Efroni et al. (2022) assume *near*-deterministic latent-state dynamics, while we make a strict determinism assumption on latent-state dynamics. However, the negative result given by Misra et al. (2024) applies even to the full-determinism variant.

- 2. We prove the correctness of CRAFT and prove sample-complexity bounds. **Context:** None.
- 3. We demonstrate the feasibility of CRAFT on a toy problem, and present the results. **Context:** None.

Offline Action-Free Learning of Ex-BMDPs by Comparing Diverse Datasets

Anonymous authors

Paper under double-blind review

Abstract

1	While sequential decision-making environments often involve high-dimensional obser-
2	vations, not all features of these observations are relevant for control. In particular, the
3	observation space may capture factors of the environment which are not controllable by
4	the agent, but which add complexity to the observation space. The need to ignore these
5	"noise" features in order to operate in a tractably-small state space poses a challenge
6	for efficient policy learning. Due to the abundance of video data available in many such
7	environments, task-independent representation learning from action-free offline data of-
8	fers an attractive solution. However, recent work has highlighted theoretical limitations
9	in action-free learning under the Exogenous Block MDP (Ex-BMDP) model, where
10	temporally-correlated noise features are present in the observations. To address these
11	limitations, we identify a realistic setting where representation learning in Ex-BMDPs
12	becomes tractable: when action-free video data from multiple agents with differing
13	policies are available. Concretely, this paper introduces CRAFT (Comparison-based
14	Representations from Action-Free Trajectories), a sample-efficient algorithm leverag-
15	ing differences in controllable feature dynamics across agents to learn representations.
16	We provide theoretical guarantees for CRAFT's performance and demonstrate its feasi-
17	bility on a toy example, offering a foundation for practical methods in similar settings.

18 1 Introduction

Many sequential decision-making settings, such as robotic navigation environments, involve highdimensional observations with many uncontrollable noise features. In order to efficiently learn policies for many downstream tasks, techniques for task-independent representation learning have been proposed. These techniques learn encoders that map the large observation space into a much smaller set of learned latent states, which can then be used to learn policies for downstream objectives more efficiently than learning from observations directly.

In some such settings, such as social navigation, large amounts of offline *video* data are available,
collected either with similar robots or with human agents. In video data, observations are available,
but *action* labels are not. Past work has shown that this offline data can be used to learn encodings
that can be leveraged for downstream tasks (Ma et al., 2023; Nair et al., 2023; Seo et al., 2022).

29 However, in recent work, Misra et al. (2024) has shown an important theoretical limitation to this 30 approach: for some important classes of environments, efficient action-free representation learning 31 is not possible. In particular, the Exogenous Block MDP (Ex-BMDP) model (Efroni et al., 2022) de-32 scribes a class of environments where observations depend both on a deterministic, action-controlled 33 latent state, and a potentially high-dimensional temporally-correlated noise factor, which is action-34 independent. The goal of representation learning in the Ex-BMDP setting is to learn a mapping from 35 the observation space to the much-smaller space of action-controlled latent states, while ignoring the 36 noise factor. Misra et al. (2024) show that, even with high coverage over latent states, representa-37 tion learning from Ex-BMDPs is not possible in general. This property of Ex-BMDPs is in contrast to *Block* MDPs, where the observation noise is not time-correlated, and which Misra et al. (2024) demonstrates *are* amenable to efficient action-free representation learning. At a high level, Misra et al. (2024)'s hardness result stems from the fact that, without action labels, action-controllable features are indistinguishable from uncontrollable features. (For example, if the observations capture both the controllable ego-agent's state and other uncontrollable "background" agents' states, it is ambiguous what state should be encoded in the learned representation.)

44 In this work, we describe a realistic setting where representation learning for Ex-BMDPs is in fact 45 tractable, and propose a provably sample-efficient algorithm for this setting. Specifically, we con-46 sider cases where videos are available of *multiple distinct agents* operating in the same environment. 47 Intuitively, the idea is that controllable latent features will differ in their dynamics between datasets 48 collected by different agents, while *uncontrollable* features will have *the same* dynamics in the two 49 datasets. Our main result is that, if two agents' policies sufficiently differ at every latent state, then, 50 under assumptions similar to those in Misra et al. (2024) and Efroni et al. (2022), sample-efficient 51 representation learning from offline action-free data collected by the two agents is possible. To show this fact, we propose a provably sample-efficient algorithm for Ex-BMDP representation 52 53 learning without action labels, which we call Comparison-based **R**epresentations from Action-Free 54 Trajectories, or **CRAFT**. CRAFT enjoys a sample complexity that depends only on the size of, 55 and coverage assumptions on, the controllable latent states of the environment, and, logarithmically, 56

on the size of the encoder hypothesis class. The sample complexity has no explicit dependence of the size of the space of exogenous noise. At a high level, CRAFT works by clustering sequential observation-pairs together based on how likely the pairs are to have been observed by each agent.

In this work, we introduce CRAFT, prove its correctness and sample-complexity, and validate its use on a toy example problem. To our knowledge, this is the first work to propose a provably sampleefficient algorithm for action-free offline representation learning in the Ex-BMDP setting. While this work is theoretical in nature due to some restrictive assumptions on the setting (which are inherited from the prior work upon which we build; see discussion in Section 6), we expect that the CRAFT algorithm can inspire practical methods that rely on the same principle of comparing action-free video datasets from diverse agents, in order to extract controllable feature representations.

66 2 Background

In this section, we define notation, formally introduce the action-free offline Ex-BMDP setting, and
 state our technical assumptions.

69 2.1 General Notation

We use [N] to denote the set $\{1, ..., N\}$. For a sequence $x_1, ..., x_N$, we use $x_{i:j}$ to denote the subsequence $x_i, x_{i+1}..., x_j$. For multisets A, B, we use $A \uplus B$ to denote the union of the two multisets, with multiplicities added.

73 2.2 Ex-BMDP Framework

The Ex-BMDP model describes a class of sequential decision-making environments where an agent's actions only operate on a hidden latent state, while the observations that the agent receives are also functions of a temporally-correlated exogenous "noise" factor, in addition to this controllable latent state. Following Efroni et al. (2022), we consider the *finite horizon* variant of this model, and also assume that the controllable latent dynamics are *deterministic*.¹ Formally, a (reward-free) Ex-BMDP can be described as a tuple, $\mathcal{M} = \langle H, \mathcal{A}, \chi_{1:H}, \mathcal{S}_{1:H}^*, \mathcal{E}_{1:H}, \mathcal{Q}_{1:H}, T_{2:H}, \mathcal{T}_{2:H}^e, S_1^*, P_1^e \rangle$, where *H* is the horizon (the number of steps per episode). At each timestep $h \in [H]$, the observation

¹Efroni et al. (2022) presents an algorithm for efficient online representation learning of Ex-BMDPs with *near*deterministic latent dynamics: the controllable latent state deviates from deterministic behavior with frequency \ll one time per episode. See Section 6 for further discussion.

81 $x_h \in \mathcal{X}_h$ is determined by two latent factors, $s_h^* \in \mathcal{S}_h^*$ and $e_h \in \mathcal{E}_h$. We assume that \mathcal{S}_h^* is finite, 82 while the \mathcal{E}_h and the observation space \mathcal{X}_h may be continuous.

The controllable latent state s_h^* evolves deterministically, depending on the action a_h taken by the agent: $s_{h+1}^* = T_{h+1}(s_h^*, a_h)$, where $T_{h+1} \in S_h^* \times \mathcal{A} \to S_{h+1}^*$ is a deterministic function, and \mathcal{A} is the set of possible actions, which we assume to be finite. Note that s_1^* is a constant. (Each episode starts at the same controllable latent state, so $S_1^* = \{s_1^*\}$.)

By contrast, the exogenous (noise) state evolves stochastically as a Markov chain, independent of actions. The initial exogenous state e_1 is sampled from the distribution $P_1^e \in \mathcal{P}(\mathcal{E}_1)$, and subsequent observations are sampled as $e_{h+1} \sim \mathcal{T}_{h+1}^e(e_h)$, where $\mathcal{T}_{h+1}^e \in \mathcal{E}_h \rightarrow \mathcal{P}(\mathcal{E}_{h+1})$. We can refer to the distribution of exogenous states e_h at time h as $P_h^e = \mathcal{T}_h^e(\mathcal{T}_{h-1}^e(...\mathcal{T}_2^e(P_1^e)...))$, and the *joint* distribution of exogenous states e_h and e_{h+1} as $P_{h:h+1}^e$.²

92 The observation x_h is then sampled as $x_h \sim Q_h(s_h^*, e_h)$, where $Q_h \in S_h^* \times \mathcal{E}_h \to \mathcal{P}(\mathcal{X}_h)$ is the 93 *emission function*. Under the Ex-BMDP model, we assume that the latent variables s_h^* and e_h can 94 always be inferred from x_h : that is, Q_h has a deterministic inverses ϕ_h^* and ϕ_h^e , such that if x_h is 95 sampled from $Q_h(s_h^*, e_h)$, then $\phi_h^*(x_h) = s_h^*$ and $\phi_h^e(x_h) = e_h$. (This assumption is the *block* 96 assumption referred to in the name "Exogenous Block MDP.")

97 The agent does not have access to s_h^* , e_h , or the "ground-truth" encoders ϕ_h^* , ϕ_h^e ; instead, it only 98 has access to the observations x_h . The goal of representation learning is to learn an encoder ϕ_h :

99 $\mathcal{X}_h \to \mathbb{N}$ for each timestep h, that approximates ϕ_h^* , up to label permutation. (We are *not* interested

100 in learning the exogenous encoder ϕ_h^e , because it is assumed that this noise factor is irrelevant for

101 control, and may be very large.)

102 2.3 Action-Free, Offline Setting

103 In this work, we consider a setting where the learner has access to multiple sets of offline trajectories 104 collected by different agents, but where only the observations x_h , and *not* the actions a_t , are avail-105 able. For simplicity, in this work we assume that there are only datasets from two distinct agents, 106 but the proposed method could be straightforwardly generalized to support more agents.

107 We refer to the two trajectory datasets as τ_A and τ_B . Each trajectory in τ_A (or τ_B) is a se-108 quence of observations $x_{1:H}$. We use $(\tau_A)_{h:h+i}$ to refer to the multiset of tuples of observations 109 $(x_h, x_{h+1}, ..., x_{h+i})$ for each trajectory in τ_A , and use $\tau_A[\{i, i', i''\}]$ to refer to the subset con-110 sisting of three *trajectories* (indexed i, i' and i'') in τ_A . We use $\mathcal{D}^*_A(s^*_h, s^*_{h+1})$ to refer to the 111 multiset consisting of all observation pairs $(x_h, x_{h+1}) \in (\tau_A)_{h:h+1}$ such that $\phi^*_h(x_h) = s^*_h$ and 112 $\phi^*_{h+1}(x_{h+1}) = s^*_{h+1}$; and $\mathcal{D}^*(s^*_h, s^*_{h+1}) := \mathcal{D}^*_A(s^*_h, s^*_{h+1}) \uplus \mathcal{D}^*_B(s^*_h, s^*_{h+1})$. We also define $\mathcal{D}^*(s^*_h)$ 113 as the multiset of observations in $x_h \in (\tau_A)_h \uplus (\tau_B)_h$ such that $\phi^*_h(x_h) = s^*_h$.

114 2.4 Technical Assumptions: Data Collection Method

As in previous works in offline representation learning in the Ex-BMDP setting (Misra et al., 2024; Islam et al., 2023; Levine et al., 2024; Lamb et al., 2023), we assume that the agents' actions a_h are chosen *independently* of the observations $x_{1:h}$, given the controllable latent states $s_{1:h}^*$. In other words, roughly speaking, we assume that the agents used to collect the offline data choose actions based *only* on the controllable latent state, not on the full observation. Misra et al. (2024) justifies this assumption by positing that the offline data are likely collected by expert agents which "would not make decisions based on noise." We discuss this assumption further in Section 6.

Beyond sharing this noise-independence assumption, our technical assumptions on the datacollection policies are otherwise significantly weaker that those in Misra et al. (2024). While Misra et al. (2024) assumes that each trajectory is generated by a *Markovian* policy (i.e, that $a_h \sim \pi(s_h^*)$, for some $\pi \in S_h^* \to \mathcal{P}(\mathcal{A})$), and furthermore that the policies used to generate each trajectory are

²Note that in general, $P_{h:h+1}^e \neq P_h^e \times P_{h+1}^e$.

- 126 chosen i.i.d., we make neither such assumption. In other words, we allow for both non-Markovian 127 behavioral policies – for example, we could have a policy in the form $a_h \sim \pi(s_{1:h}^*)$ – and non-i.i.d. 128 sampling of behavioral policies between episodes – for example: the data collector could evolve 129 over time between episodes, in order to, for instance, maximize the diversity of visited latent states.
- 130 Explicitly, our *only* assumption on the data-collection mechanism is that the process which generates 131 action sequences $a_{1:H}$ – and therefore, equivalently, controllable latent-state sequences $s_{1:H}^*$ – is 132 independent from observation noise over *both entire datasets*. Formally, for a dataset τ that consists 133 of trajectories $x_{1:H}$, let $\phi^*(\tau)$ denote the corresponding set of controllable latent state trajectories 134 $s_{1:H}^*$, and $\phi^e(\tau)$ denote the corresponding set of exogenous state trajectories $e_{1:H}$. Then, our only 135 requirement on the data collection mechanism is that it ensures:

$$\Pr(\tau_A, \tau_B) = \Pr(\phi^*(\tau_A), \phi^*(\tau_B)) \\ \cdot \Pr_{P_e^e, \mathcal{T}^e}(\phi^e(\tau_A)) \cdot \Pr_{Q}(\tau_A | \phi^*(\tau_A), \phi^e(\tau_A)) \cdot \Pr_{Q}(\tau_B | \phi^*(\tau_B), \phi^e(\tau_B)).$$
(1)

136 To see why Equation 1 is a sufficiently strong assumption to allow for useful analysis despite 137 its apparent generality, fix any two latent states $s_h^*, s_{h+1}^* \in \mathcal{S}_{h:h+1}^*$, and consider the multiset 138 $\mathcal{D}_A^*(s_h^*, s_{h+1}^*)$ as defined in Section 2.3; also let $n := |\mathcal{D}_A^*(s_h^*, s_{h+1}^*)|$. Then the marginal distri-139 bution of $\mathcal{D}_A^*(s_h^*, s_{h+1}^*)$ can be described as:

$$\mathcal{D}_{A}^{*}(s_{h}^{*}, s_{h+1}^{*}) \sim [(\mathcal{Q}(s_{h}^{*}, e_{h}), \mathcal{Q}(s_{h+1}^{*}, e_{h+1}))|e_{h}, e_{h+1} \sim P_{h:h+1}^{e}]^{n}.$$
(2)

140 We see that $\mathcal{D}_A^*(s_h^*, s_{h+1}^*)$ consists of i.i.d. samples from a fixed, policy-independent distribution. 141 Consequently, this property will frequently allow us to use standard concentration bounds in our 142 analysis, while still allowing for a wide class of non-Markovian, non-i.i.d. behavioral policies.

143 While we do not require the behavioral policies to be Markovian, it will be useful to refer to the 144 "empirical policies" $\pi_A^{emp.}$ and $\pi_B^{emp.}$, defined as:

$$\pi_A^{emp.}(s_{h+1}^*|s_h^*) := \frac{|\mathcal{D}_A^*(s_h^*, s_{h+1}^*)|}{\sum_{s' \in \mathcal{S}_{h+1}^*} |\mathcal{D}_A^*(s_h^*, s')|},\tag{3}$$

145 and likewise for π_B^{emp} . This is the empirical likelihood in the provided data that agent A (respec-146 tively, B) chooses an action that results in a transition from s_h^* to s_{h+1}^* .

147 2.5 Technical Assumptions: Coverage, Policy Diversity, and Realizability

148 In order to learn accurate latent state encoders $\phi_{1:H}$, we need to ensure adequate coverage over all 149 latent states s_h . For all timesteps h, and all pairs of latent states $(s_h^*, s_{h+1}^*) \in \mathcal{S}_h^* \times \mathcal{S}_{h+1}^*$ such that 150 $s_{h+1}^* = T_h(s_h^*, a)$ for some action a, we require that

$$\frac{|\mathcal{D}^*(s_h^*, s_{h+1}^*)|}{|\tau_A| + |\tau_B|} \ge \nu,\tag{4}$$

- 151 for some known lower-bound ν . This coverage assumption is presented in terms of the *actually* 152 *realized offline datasets* τ_A and τ_B . By contrast, Misra et al. (2024) assumes that trajectories in the 153 offline dataset are sampled i.i.d., and makes coverage assumptions on the *policies* used sample them.
- 154 We also require that the two agents, which produced datasets τ_A and τ_B , behaved sufficiently differ-
- 155 ently so that we can infer the latent dynamics from their differences. In particular, for some known
- 156 lower bound $\alpha > 0$, we require that, $\forall h \in [H-1], \forall s_h^* \in \mathcal{S}_h^*$, and for any two successor states 157 $s_{h+1}^*, s_{h+1}^{\prime*} \in \mathcal{S}_{h+1}^*$, such that s_h^* can transition to either s_{h+1}^* or $s_{h+1}^{\prime*}$ under T_h , we have, either:

$$e^{\alpha} \cdot \frac{\pi_{B}^{emp.}(s_{h+1}^{\prime*}|s_{h}^{*})}{\pi_{B}^{emp.}(s_{h+1}^{*}|s_{h}^{*})} \leq \frac{\pi_{A}^{emp.}(s_{h+1}^{\prime*}|s_{h}^{*})}{\pi_{A}^{emp.}(s_{h+1}^{*}|s_{h}^{*})}, \text{ or, } e^{\alpha} \cdot \frac{\pi_{A}^{emp.}(s_{h+1}^{\prime*}|s_{h}^{*})}{\pi_{A}^{emp.}(s_{h+1}^{*}|s_{h}^{*})} \leq \frac{\pi_{B}^{emp.}(s_{h+1}^{\prime*}|s_{h}^{*})}{\pi_{B}^{emp.}(s_{h+1}^{*}|s_{h}^{*})}.$$
(5)

158 In other words, the relative likelihood of transitioning to $s_{h+1}^{\prime*}$, versus transitioning to s_{h+1}^{*} , is dif-159 ferent in τ_A and τ_B by a multiplicative factor of at least e^{α} .

Finally, we also require that the difference in *total* state coverage between τ_A and τ_B for any pair of sequential latent states (s_h^*, s_{h+1}^*) is not *too* extreme. We require that, for a known lower-bound η :

$$\frac{|\mathcal{D}_{A}^{*}(s_{h}^{*}, s_{h+1}^{*})|}{|\mathcal{D}^{*}(s_{h}^{*}, s_{h+1}^{*})|} \ge \eta, \tag{6}$$

- and likewise for $\mathcal{D}_B^*(s_h^*, s_{h+1}^*)$. Our sample-complexity bound also depends on an additional param-
- 163 eter ν' , which does *not* need to be known a priori. This is the minimum single-state coverage ratio:

$$\nu' := \min_{h \in [H], s_h^* \in \mathcal{S}_h^*} \frac{|\mathcal{D}^*(s_h^*)|}{|\tau_A| + |\tau_B|}.$$
(7)

Function Approximation Assumptions. We assume access to hypothesis classes of encoder functions $\Phi_{1:H}$, as well as binary classification functions $\mathcal{G}_h \subseteq \mathcal{X}_h \to \{0, 1\}$. We make standard realizability assumptions (in brief, $\phi_h^* \in \Phi_h$, and $\forall s_h^*, s_h'^*, \exists g \in \mathcal{G}_h$ such that g can perfectly distinguish between observations of s_h^* and those of $s_h'^*$) and assume access to training oracles. See Appendix A for further information about these assumptions. We use $|\Phi|$ to denote $\max_h |\Phi_h|$. We also assume a known upper-bound N_s on $\max_h |\mathcal{S}_h|$; that is, the maximum output range of any encoder in Φ_h .

170 **3 Method**

171 In this section, we describe the CRAFT algorithm. First, however, we motivate its design by exam-172 ining a simpler version of the problem setting and of the algorithm.

173 3.1 Intuition: Single-step, Binary Action Case

174 In this section, we present a toy algorithm for an extremely simplified version of the Ex-BMDP

175 representation learning problem: an explanation of the toy algorithm captures some of the intu-

176 itions of CRAFT, while the limitations of the toy algorithm will motivate some of the less-intuitive

177 algorithmic details. We can call this naive, "first draft" form of the CRAFT algorithm "DRAFT."

178 We first consider the Ex-BMDP model with H = 2 and |A| = |S₂*| = 2. In this setting, the
179 representation learning problem reduces to the task of learning to distinguish the two latent states s^{*}₂ and s₂^{*} ∈ S₂* that can occur at h = 2. (See Figure 1a.)



(a) Latent-state transition graph.

(b) Composition of the datasets $(\tau_A)_2$ and $(\tau_B)_2$.

Figure 1: Dynamics and composition of the two-step Ex-BMDP example in Section 3.1.

180

181 We will also assume that $|\tau_A| = |\tau_B| = m$. Then, $(\tau_A)_1$ and $(\tau_B)_1$ both consist entirely of m i.i.d. 182 samples of the same distribution $\mathcal{Q}(s_1^*, P_1^e)$. The structures of $(\tau_A)_2$ and $(\tau_B)_2$ are (slightly) more 183 complex. If we let $\gamma_A := \pi_A^{emp.}(s_2^*|s_1^*)$ and $\gamma_B := \pi_B^{emp.}(s_2^*|s_1^*)$, then we see that the dataset $(\tau_A)_2$ 184 consists of $m \cdot \gamma_A$ i.i.d. samples of the distribution $\mathcal{Q}(s_2^*, P_2^e)$ and $m \cdot (1 - \gamma_A)$ i.i.d. samples of the 185 distribution $\mathcal{Q}(s_2'^*, P_2^e)$, while $(\tau_B)_2$ consists of $m \cdot \gamma_B$ i.i.d. samples of the distribution $\mathcal{Q}(s_2^*, P_2^e)$ 186 and $m \cdot (1 - \gamma_B)$ i.i.d. samples of the distribution $\mathcal{Q}(s_2'^*, P_2^e)$. (See Figure 1b.) 187 Now, by assumption, the agents generating datasets τ_A and τ_B are not behaviorally identical: this 188 means that $\gamma_A \neq \gamma_B$. Without loss of generality, assume that $\gamma_A > \gamma_B$. Our key insight is that, in 189 the limit as $m \to \infty$, the Bayes optimal classifier to distinguish a sample x_2 selected from $(\tau_A)_2$ 190 from a sample x_2 selected from $(\tau_B)_2$ is in fact (up to label permutation) the latent state encoder 191 $\phi_2^*(x_2)$. Concretely, consider a classifier ϕ_2' trained to minimize the 0-1 classification loss between 192 $(\tau_A)_2$ and $(\tau_B)_2$. In the limit of infinite data, we can define this classification loss function as:

$$\mathcal{L}_{pop}(\phi_2) := \lim_{m \to \infty} \mathop{\mathbb{E}}_{x \in (\tau_A)_2} \phi_2(x) + \mathop{\mathbb{E}}_{x \in (\tau_B)_2} 1 - \phi_2(x), \tag{8}$$

193 From Equation 8 and the composition of the datasets:

$$\mathcal{L}_{pop}(\phi_{2}) = \gamma_{A} \underset{x \sim \mathcal{Q}(s_{2}^{*}, P_{2}^{e})}{\mathbb{E}} \phi_{2}(x) + (1 - \gamma_{A}) \underset{x \sim \mathcal{Q}(s_{2}^{*}, P_{2}^{e})}{\mathbb{E}} \phi_{2}(x) + \gamma_{B} \underset{x \sim \mathcal{Q}(s_{2}^{*}, P_{2}^{e})}{\mathbb{E}} (1 - \phi_{2}(x)) + (1 - \gamma_{B}) \underset{x \sim \mathcal{Q}(s_{2}^{*}, P_{2}^{e})}{\mathbb{E}} (1 - \phi_{2}(x)) = (\gamma_{A} - \gamma_{B}) [\underset{x \sim \mathcal{Q}(s_{2}^{*}, P_{2}^{e})}{\mathbb{E}} \phi_{2}(x) + \underset{x \sim \mathcal{Q}(s_{2}^{*}, P_{2}^{e})}{\mathbb{E}} 1 - \phi_{2}(x)] + C = -(\gamma_{A} - \gamma_{B}) (\Pr(\phi_{2}(x) = 0 | \phi^{*}(x) = s_{2}^{*}) + \Pr(\phi_{2}(x) = 1 | \phi^{*}(x) = s_{2}^{*})) + C$$

$$(9)$$

194 where C is independent of ϕ_2 . Under the mapping $(0 \rightarrow s_2^*, 1 \rightarrow s_2'^*)$, we see that $\mathcal{L}_{pop}(\phi_2)$ scales

195 linearly with the rate that ϕ_2 produces incorrect encodings, with a δ -increase in \mathcal{L}_{pop} corresponding

196 to an $\mathcal{O}((\gamma_A - \gamma_b) \cdot \delta)$ increase in encoder failure. In particular, \mathcal{L}_{pop} is uniquely minimized by the

197 ground-truth encoder ϕ_2^* . We now examine how this simple algorithm functions with finite datasets:

Algorithm 1 "DRAFT" algorithm for
$$H = 2$$
 Ex-BMDPs.
Require: Datasets τ_A , τ_B with $H = 2$, hypothesis class $\Phi_2 \in \mathcal{X}_2 \to \{0, 1\}$.
Let $\phi'_1 := \mathcal{X}_1 \to 0$, and $\phi'_2 := \arg \min_{\phi_2 \in \Phi_2} \mathcal{L}(\phi_2)$, where:
 $\mathcal{L}(\phi_2) := \frac{1}{|\tau_A|} \sum_{x \in (\tau_A)_2} \phi_2(x) + \frac{1}{|\tau_B|} \sum_{x \in (\tau_B)_2} 1 - \phi_2(x).$ (10)

Return: ϕ'_1, ϕ'_2

198

199 With finite *m*, our main concern is overfitting: if Φ_2 is large enough such that some $\phi'_2 \in \Phi_2$ can 200 perfectly distinguish the samples that happen to fall into $(\tau_A)_2$ from those in $(\tau_B)_2$, then this ϕ'_2 will 201 attain a lower empirical loss than ϕ^*_2 , while being bad at distinguishing $\mathcal{Q}(s^*_2, P^e_2)$ from $\mathcal{Q}(s'^*_2, P^e_2)$ 202 in general. However, as long as $|\Phi_2|$ is controlled, we can use standard concentration inequalities to 203 limit this overfitting. In particular,

$$|\mathcal{L}(\phi_2) - \mathcal{L}_{pop}(\phi_2)| \approx \mathcal{O}(1/\sqrt{m}).$$
(11)

By combining Equations 11 and 9, we can determine how quickly $\phi'_2 = \arg \min \mathcal{L}(\cdot)$ will approach $\phi^*_2 = \arg \min \mathcal{L}_{pop}(\cdot)$ as *m* increases, in terms of *accuracy as a latent state encoder*. To ensure ϕ'_2 approximates ϕ^*_2 with a failure rate of at most ϵ , we need $m \approx \mathcal{O}\left(\frac{1}{(\gamma_A - \gamma_b)^2 \epsilon^2}\right)$ samples. Intuitively, the smaller the difference between behavior policies of the two agents $(\gamma_A - \gamma_B)$, the more samples are required to attain a given accuracy of encoder.

209 3.1.1 "DRAFT" doesn't generalize easily to the long-horizon setting

210 A naive first attempt to extend "DRAFT" to the H > 2 case might be to apply it *recursively*. That 211 is, once the two distinct latent states at h = 2 can be decoded, we can extract from $(\tau)_A$ and $(\tau)_B$ 212 the trajectories which contain (say) s_2^* , and then run DRAFT again on these samples, to obtain 213 an encoder that can separate the two states into which s_2^* can transition.³ We can then repeat this

³Here, we are continuing to assume $|\mathcal{A}| = 2$, and that the two actions have different effects from each other in each state.

- procedure for $s_2^{\prime*}$. If $s_2^{\prime*}$ and s_2^* both transition to the same latent state (say $s_3^{\prime*}$), we can easily detect
- this situation by attempting to learn binary classifiers between the observations of each state that
- succeeds s_2^* and each state that succeeds $s_2'^*$: if it is impossible distinguish these observations better than by random chance, then the two successor states are the same:



Figure 2: Illustration of recursive use of the "DRAFT" algorithm.

217

However, it turns out that it is difficult (and may be impossible) to prove an efficient samplecomplexity bound for this recursive algorithm. This is for two reasons:

220 1. After the first timestep, the input datasets to subsequent applications of DRAFT are corrupted 221 by mis-classified samples, such that the datasets are no longer mixtures of i.i.d. samples from 222 distributions $Q(s^*, P_h^e)$ for multiple values of $s^* \in S_h^*$.

223 2. DRAFT is highly sensitive to small changes in its input dataset.

To see (1), note that the encoder ϕ'_2 returned by the first application of DRAFT will misclassify some 224 225 $\mathcal{O}(\sqrt{m}/(\gamma_A - \gamma_B))$ samples. Moreover, these misclassified samples will *not* be chosen uniformly: the encoder ϕ'_2 may rely on some spurious features of the observations x_2 , which depend on e_2 , 226 227 to classify these observations. Consequently, because e_3 also depends on e_2 , the exogenous noise distributions of the observations x_3 of state $s_3^{\prime*}$ (for instance, in the dynamics example in Figure 2) 228 229 present in the datasets for the recursive DRAFT instances associated with s_2^* and $s_2'^*$ will differ 230 from each other, and each will differ from $\mathcal{Q}(s_3^{\prime*}, P_h^e)$, in a way that depends on the choice of ϕ_2^{\prime} . 231 Moreover, because ϕ'_2 is trained to distinguish τ_A from τ_B , this distribution shift may have different 232 effects on the distributions of observations from the two agents.

233 For (2), consider the (single step) DRAFT algorithm with some small number $\epsilon_{bad} \cdot m$ of samples 234 from \mathcal{X}_2 arbitrarily introduced to the datasets $(\tau_A)_2$ and $(\tau_B)_2$. For concreteness, we will replace some of the samples in $(\tau_A)_2$ for which $\phi_2^*(x) = s_2'^*$ with "bad" samples for which it *still* holds 235 that $\phi_2^*(x'_{bad}) = s'^*_2$, but which are not drawn i.i.d. from $\mathcal{Q}(s'^*_2, P^e_2)$. Instead, we assume that these 236 samples belong to some part of the support of $\mathcal{Q}(s_2'^*, P_2^e)$ which is typically sampled with negligible 237 probability; we can call their distribution \mathcal{D}'_{bad} . Similarly, we replace $\epsilon_{bad} \cdot m$ samples of $(\tau_B)_2$ for 238 which $\phi_2^*(x) = s_2^*$ with "bad" samples for which $\phi_2^*(x_{bad}) = s_2^*$, but which are drawn from \mathcal{D}_{bad} . 239 240 We consider the infinite-dataset limit. From Equation 8 and the composition of the datasets:

$$\begin{aligned} \mathcal{L}_{pop}(\phi_{2}) &= \gamma_{A} \mathop{\mathbb{E}}_{x \sim \mathcal{Q}(s_{2}^{*}, P_{2}^{e})} \phi_{2}(x) + (1 - \gamma_{A} - \epsilon_{bad}) \mathop{\mathbb{E}}_{x \sim \mathcal{Q}(s_{2}^{*}, P_{2}^{e})} \phi_{2}(x) + \epsilon_{bad} \mathop{\mathbb{E}}_{x \sim \mathcal{D}_{bad}} (\phi_{2}(x)) \\ &+ (1 - \gamma_{B}) \mathop{\mathbb{E}}_{x \sim \mathcal{Q}(s_{2}^{*}, P_{2}^{e})} (1 - \phi_{2}(x)) + (\gamma_{B} - \epsilon_{bad}) \mathop{\mathbb{E}}_{x \sim \mathcal{Q}(s_{2}^{*}, P_{2}^{e})} (1 - \phi_{2}(x)) + \epsilon_{bad} \mathop{\mathbb{E}}_{x \sim \mathcal{D}_{bad}} (1 - \phi_{2}(x)) \\ &= -(\gamma_{A} - \gamma_{B} + \epsilon_{bad}) \Big(\Pr(\phi_{2}(x) = 0 | x \sim \mathcal{Q}(s_{2}^{*}, P_{2}^{e})) + \Pr(\phi_{2}(x) = 1 | x \sim \mathcal{Q}(s_{2}^{*}, P_{2}^{e})) \Big) \\ &- \epsilon_{bad} \Big(\Pr(\phi_{2}(x) = 1 | x \sim \mathcal{D}_{bad}) + \Pr(\phi_{2}(x) = 0 | x \sim \mathcal{D}_{bad}') \Big) + C \end{aligned}$$

Note that the ground-truth encoder ϕ_2^* has loss $\mathcal{L}_{pop}(\phi_2^*) = -(\gamma_A - \gamma_B + \epsilon_{bad}) + C$. However, we can construct an encoder ϕ_2' that incorrectly encodes all samples in \mathcal{D}_{bad} as belonging to s_2^{**} (i.e., returns 1 on these samples), and incorrectly encodes all samples in \mathcal{D}'_{bad} as belonging to s_2^* . Furthermore, we can construct this ϕ_2' to also have an accuracy of only $1 - \epsilon_{bad}/(\gamma_A - \gamma_B + \epsilon_{bad})$ on the samples in the "natural" distributions $\mathcal{Q}(s_2^*, P_2^e)$ and $\mathcal{Q}(s_2'^*, P_2^e)$. Surprisingly, by the above expression for \mathcal{L}_{pop} , we see that this less-accurate encoder *also* has loss $\mathcal{L}_{pop}(\phi'_2) = -(\gamma_A - \gamma_B + \epsilon_{bad}) + C$. Therefore, if this encoder ϕ'_2 is included in the hypothesis class Φ_2 , then the ERM step of Algorithm 1 may just as easily return ϕ'_2 rather than ϕ^*_2 . (Furthermore, the realizability assumption does not guarantee that any lower-loss "third option" encoders exist.) Consequently, the misclassification rate can "blow up" from ϵ_{bad} to $\epsilon_{bad}/(\gamma_A - \gamma_B + \epsilon_{bad})$, that is, by a *multiplicative* factor of $1/(\gamma_A - \gamma_B + \epsilon_{bad})$ – even before accounting for finite datasets.

252 We can now see that DRAFT both (1) can make non-uniformly-distributed encoding errors, and (2) 253 given as input a dataset with (non-uniform) errors, can produce output "next-state" datasets with 254 a multiplicatively-increased error rate. Therefore, it seems difficult to derive a sample-complexity analysis of "recursive DRAFT" that does not require a number of samples exponential in the Ex-255 BMDP horizon H^4 . In the next section, we present our CRAFT algorithm, which is intention-256 257 ally designed to solve the multiple-agent action-free Ex-BMDP representation learning problem 258 while avoiding recursively training state classifiers on datasets derived from the output of previous-259 timestep state classifiers. We then can sidestep the issues with "Recursive DRAFT" shown here.

Note that, for Ex-BMDPs with deterministic latent dynamics, the issues with recursion seen here are unique to the offline, action-free setting. In the online setting, as in Efroni et al. (2022), once the dynamics up to timestep h have been learned, "fresh" samples of any given latent state s_h^* can then be constructed via closed-loop planning: there are no issues with compounding error.⁵ The action-free offline setting thus presents a new set of issues requiring a novel algorithmic solution.

265 3.2 CRAFT: High-Level Description of Method

Here, we give a high-level overview of the CRAFT algorithm. The complete algorithm is presented as Algorithm 2 in Appendix C. See also Figure 3 for a pictorial overview of the approach.

In CRAFT, we initially treat each trajectory as a sequence of observation *pairs*: (x_1, x_2) , (x_2, x_3) , ... (x_{H-1}, x_H) . (See Figure 3-a.) For each timestep-pair (h, h + 1), we train a model f_h to predict, given a sample (x_h, x_{h+1}) , whether the pair was collected by agent A or agent B: that is, whether (x_h, x_{h+1}) was selected from $(\tau_A)_{h,h+1}$ or $(\tau_B)_{h,h+1}$. However, unlike in "DRAFT", we do not treat this problem as hard binary classification. Instead, we train $f_h(x_h, x_{h+1})$ to predict the *log*odds ratio between the two possibilities, for a given (x_h, x_{h+1}) . That is, we train f_h to predict

$$\ln\left(\frac{\Pr[(x_h, x_{h+1}) \in (\tau_A)_{h,h+1} | (x_h, x_{h+1}) \in (\tau_A \uplus \tau_B)_{h,h+1}]}{\Pr[(x_h, x_{h+1}) \in (\tau_B)_{h,h+1} | (x_h, x_{h+1}) \in (\tau_A \uplus \tau_B)_{h,h+1}]}\right).$$
(12)

274 To accomplish this task, we can train f_h to minimize the following loss function:

$$\mathcal{L}(f_h) := \sum_{(x_h, x_{h+1}) \in (\tau_A)_{h,h+1}} \ln(1 + e^{-f(x_h, x_{h+1})}) + \sum_{(x_h, x_{h+1}) \in (\tau_B)_{h,h+1}} \ln(1 + e^{f(x_h, x_{h+1})}).$$
(13)

275 Note that in the limit of infinite data, the f_h that minimizes this loss will return

$$f_h^*(x_h, x_{h+1}) \to \ln\left(\frac{|\mathcal{D}_A^*(\phi_h^*(x_h), \phi_{h+1}^*(x_{h+1}))|}{|\mathcal{D}_B^*(\phi_h^*(x_h), \phi_{h+1}^*(x_{h+1}))|}\right).$$
(14)

276 Consequently (for sufficiently-large datasets) we expect the values of $f_h(x_h, x_{h+1})$ of all

observation-pairs (x_h, x_{h+1}) corresponding to the same latent-state pair $(\phi_h^*(x_h), \phi_{h+1}^*(x_{h+1})) =$

278 (s_h^*, s_{h+1}^*) to "cluster together" around the same value (See Figure 3-b): this effect can be quantified

⁴We are not claiming that "Recursive DRAFT" *actually does* require exponential samples in H, simply that there are clear obstacles to proving that it *does not*.

⁵Even with *offline* data, if action labels are available and the latent dynamics up to timestep h are known perfectly (which is achievable if the latent dynamics are deterministic), then "error-free" datasets can still be constructed for timestep h + 1.

a. <u>Trajectories are treated as observation pairs.</u>



Sets of observation pairs (xh, xh+1) which

succeed each inferred sh are labeled pairs(sh).

f2(X2,X3)

C.



Models fn predict log-odds that agent A or



Figure 3: Schematic of the CRAFT algorithm. See text of Section 3.2.

using standard concentration arguments. Note that the training of models f_h and resulting "clustering" of observation-pairs can be carried out *simultaneously and independently* for all time-steps h: there is no "recursion" here, and so each model f_h is trained on an "untainted" dataset.

Side note on realizability and discretization: To ensure that an f_h can be found that minimizes Equation 13, we need f_h to be chosen from a sufficiently-expressive hypothesis class \mathcal{F}_h . We can construct such an \mathcal{F}_h as $\Phi_h \times \Phi_{h+1} \times (N_s^2 \to \mathbb{R})$: by the realizability assumptions on Φ_h and Φ_{h+1} , we are ensured that this class contains the optimal predictor in Equation 14. However, the $(N_s^2 \to \mathbb{R})$ component of this hypothesis class makes it non-finite. In order to allow for a simple finite-hypothesis analysis, we instead construct \mathcal{F}_h as $\Phi_h \times \Phi_{h+1} \times (N_s^2 \to \Xi)$, where Ξ is a

- discrete space (roughly, every $(\alpha/4)$ -th interval on a range determined by η). It turns out that this discretization still ensures that a function "close enough" to f_{t}^{*} will always exist, and additionally
- discretization still ensures that a function "close enough" to f_h^* will always exist, and additionally greatly simplifies the identification of "clusters" in the output distribution of f_h on finite data.

While we expect the values of $f_h(x_h, x_{h+1})$ to cluster for sets of observations-pairs with the *same* latent-state-pair, it does not immediately follow that the values of $f_h(x_h, x_{h+1})$ and $f_h(x'_h, x'_{h+1})$ will *differ* if $(s_h, s_{h+1}) \neq (s'_h, s'_{h+1})$. In fact, this is not true in general: two distinct "clusters" may overlap entirely. This is where the assumption given in Equation 5 becomes useful: from Equation 5 and algebra, we can see that for any fixed s^*_h and distinct s^*_{h+1}, s'^*_{h+1} which can both follow s^*_h that:

$$\left| \ln \left(\frac{|\mathcal{D}_A^*(s_h^*, s_{h+1}^*)|}{|\mathcal{D}_B^*(s_h^*, s_{h+1}^*)|} \right) - \ln \left(\frac{|\mathcal{D}_A^*(s_h^*, s_{h+1}^{\prime*})|}{|\mathcal{D}_B^*(s_h^*, s_{h+1}^{\prime*})|} \right) \right| \ge \alpha.$$
(15)

In other words, the "clusters" associated with two pairs (s_h^*, s_{h+1}^*) and (s_h^*, s_{h+1}^*) are guaranteed to be distinct if $s_h^* = s_h^{\prime *}$. One consequence is that the observation-pairs (x_1, x_2) associated with each possible latent-state pair (s_1^*, s_2^*) must form distinct, well-separated clusters (because these pairs all share the same initial latent state s_1^*). Therefore, datasets of observations associated with each latent state s_2 can be immediately identified (as shown in green in Figure 3-b; note that we omit the asterisk, to indicate that these are *inferred*, rather than *ground-truth*, latent states.).

CRAFT then continues "recursively": once the trajectories which contain a particular $s'_2 \in S$ are 302 known, we can then examine the spectrum of values of $f_2(x_2, x_3)$ for only observation pairs (x_2, x_3) 303 304 from this *subset* of trajectories (referred to in Algorithm 2 as pairs (s'_2)). Because these observationpairs all (up to an error factor) share the same initial state $s_2^{\prime*}$, we expect to see well-separated clusters 305 for each latent state which can succeed s'_2 . Note that we do not retrain f_2 on only these samples in 306 307 pairs(s'_2). Therefore, any errors (missing or extra trajectories) in the construction of pairs(s'_2) can 308 only substantially affect the outcome of this step by compromising CRAFT's ability to recognize 309 distinct clusters in the *precomputed values* of $f_2(x_2, x_3)$. Due to the discretization of the range of 310 \mathcal{F}_h , this "cluster identification" is robust to even adversarial errors affecting a bounded number of 311 trajectories. The total number of misclassified trajectories then grows only linearly with H. (In 312 Figure 3-c, we show the spectrum of values of $f_2(x_2, x_3)$ for each subset pairs(s_2), pairs(s'_2), and pairs(s''_2); Figure 3-d shows the result of the cluster identification: the observation-pairs (x_2, x_3) 313 314 corresponding to each state in S_3 which can succeed each of s_2, s'_2 , and s''_2 have been identified).

Once we identify each latent state that can succeed each $s_2' \in \mathcal{S}_2$ individually, we now determine 315 whether or not any of these successor states to distinct states s_2 , s'_2 are in fact the same latent state 316 $s_3 \in S_3$. This can be accomplished easily, by attempting to learn binary classifiers between the ob-317 318 servations x_3 which are part of the observation-pair sets. If these observations are indistinguishable, 319 then the sets of observation-pairs represent the same latent state; if they are perfectly distinguish-320 able, then they represent different latent states. Figure 3-e illustrates this process. Note that while 321 there may be errors in these observation sets, each binary-classifier training ultimately produces a 322 boolean result (either the sets are distinguishable, or they are not) with a substantial allowance for error in the input sets: there is (with high probability) no accumulation of errors due to this process. 323

Finally, the observations corresponding to each unique latent state S_3 have been identified. (See Figure 3-f). We can then continue to timestep h = 4, and so on. As mentioned above, both the cluster-identification and state-merging processes are robust to bounded errors in their input data, so the total number of misclassified states grows only linearly in H. As a final step, the encoders ϕ'_h are trained on the assembled datasets for each timestep h.

329 3.3 Guarantees

330 We prove the following polynomial sample-complexity guarantee for CRAFT in Appendix C:

Theorem 3.1. Assume that CRAFT (Algorithm 2 in the Appendix) is given datasets τ_A and τ_B such that the assumptions given in Equations 1, 4,5, and 6 all hold. Then there exists an

$$f\left(H, |\Phi|, N_s, \frac{1}{\delta}, \frac{1}{\epsilon_0}, \frac{1}{\nu}, \frac{1}{\nu'}, \frac{1}{\eta}, \frac{1}{\alpha}\right) \in \mathcal{O}^*\left(\frac{H^2(\ln(|\Phi|/\delta) + N_s^2)}{\nu\eta^2\alpha^4} \cdot \max\left(\frac{1}{\nu^2}, \frac{1}{\epsilon_0^2\nu'^2}\right)\right), \quad (16)$$

333 where $\mathcal{O}^*(f(x)) := \mathcal{O}(f(x)\log^k(f(x)))$, such that for any given $\delta, \epsilon_0 \ge 0$, if $\forall s_h^*, s_{h+1}^*$ such that s_h^*

 $\text{ san transition to } s^*_{h+1}, |\mathcal{D}^*(s^*_h, s^*_{h+1})| \geq f\left(H, |\Phi|, N_s, \frac{1}{\delta}, \frac{1}{\epsilon_0}, \frac{1}{\nu}, \frac{1}{\nu'}, \frac{1}{\eta}, \frac{1}{\alpha}\right), \text{ then, with probability}$

at least $1 - \delta$, the encoders ϕ'_h returned by the algorithm will each have accuracy on at least $1 - \epsilon_0$, in the sense that, under some bijective mapping $\sigma_h : S_h \to S_h^*$,

$$\forall s^* \in \mathcal{S}_h^*, \quad \Pr_{x \sim \mathcal{Q}(s^*, P_h^e)}(\phi_h'(x) = \sigma_h^{-1}(\phi_h^*(x))) \ge 1 - \epsilon_0. \tag{17}$$

337

338 4 Simulation Results

339 We test CRAFT on a toy environment which captures CRAFT's ability to distinguish controllable 340 features in the observation space from time-correlated uncontrollable features. In the environment, $s_1^* = 0, \mathcal{A} = \mathcal{S}_{h>1}^* = \{0, 1\}$ and $s_{h+1}^* = a_h$; in other words, the agent can simply set the next latent 341 state using the action. The exogenous state consists of M - 1 factors: $e = (e^1, e^2, ..., e^{M-1})$. Each exogenous factor is a two-state Markov Chain: for $e^2, ... e^{M-1}$, the initial state distribution and state 342 343 transition probabilities are arbitrary parameters chosen uniformly at random for each chain, while e^1 344 has $Pr(e_1^1 = 0) = 0.5$ and transition probabilities of zero. The observation $x_h \in \{0, 1\}^M$ consists 345 of s_h^* concatenated with $(s_h^* \operatorname{XOR} e_h^i)$, for each $i \in [H-1]$. Additionally, at each timestep, the 346 347 order of s_h^* and the other factors is *permuted* by some arbitrary permutation which depends on h. The hypothesis classes are $\Phi_h := \{(x_h) \to (x_h)_i | i \in [M]\}$. The representation learning problem is 348 349 then to determine, for each h, which of the M components of the observation x_h is the controllable 350 factor s_h^* (or, failing at that, to find a component corresponding to a $(s_h^* \text{ XOR } e_h^i)$ where e_h^i is lowentropy, so the encoder imperfect but still useful). Agent A selects actions uniformly at random, 351 352 while for agent B, $Pr(a_h = s_h^*) = 3/4$.

353 Results are shown in Table 1. The setting is designed to prevent various "shortcuts" to learning an 354 encoder from working. Simply choosing the component of x_h that best predicts the policy ("Singleobservation classification" in Table 1) will not work, because at any sufficiently large timestep h, 355 356 the latent state distributions of the two policies are essentially identical (with a total-variation gap of 2^{-h}). Furthermore, given observations of a *pair* of sequential timesteps (x_h, x_{h+1}) , choosing 357 the components of x_h and x_{h+1} , respectively, that *together* best predict the agent also will not 358 359 work ("Paired-observation classification"). In particular, the "distractor" features $(s_h^* \text{ XOR } e_h^1)$ and $(s_{h+1} \text{ XOR } e_{h+1}^1)$ are, taken together, about as informative about the *agent's identity* as s_h^* and 360 361 s_{h+1}^* , but provide *no* information about the latent state s_h^* or s_{h+1}^* . In Table 1, we see that, given sufficient data (≥ 1000 trajectories for each agent), CRAFT is capable of learning highly-accurate 362 encoders in this setting, while these two "shortcut" techniques are not. In particular, while the 363 364 "Paired-observation classification" shortcut is about as effective as CRAFT in the very-low data regime, its performance plateaus (and even seems to drop) as more data becomes available. (The 365 366 *drop* in performance is likely because the adversarially-designed "distractor" features $(s_h^* \text{ XOR } e_h^1)$ 367 and $(s_{h+1}^* \text{ XOR } e_{h+1}^1)$ are more likely to be chosen by this method as more data becomes available.)

Table 1: Results of toy environment simulation, with H = 30, M = 128, averaged over 20 random seeds. See text of Section 5, and Appendix D for further details.

Technique	Avg. Encoder Acc. $(\tau_A = \tau_B = 500)$	" " 1000	" " 5000
CRAFT	86.4%	97.7%	>99.9%
Single-obs. classification	67.8%	68.7%	69.7%
Paired-obs. classification	87.4%	86.1%	82.1%

369 **5 Related Works**

370 Action-free representation learning Many prior works have tackled action-free representation 371 learning in practical scenarios, demonstrating empirically-validated methods. Common approaches 372 utilize observation reconstruction losses (Seo et al., 2022) or temporal contrastive losses (Nair et al., 373 2023). Some of these works infer "latent actions" by finding a compact representation that is highly 374 informative for predicting forward dynamics (Edwards et al., 2019; Menapace et al., 2021; Ye et al., 375 2023; Schmidt & Jiang, 2024) Another line of work augments large action-free offline datasets with 376 significantly smaller action-labeled datasets. For example, an offline dataset action-free dataset can 377 be used to train a goal-conditioned value function (Xu et al., 2022; Ma et al., 2023; Ghosh et al., 378 2023; Park et al., 2023). Alternatively, an inverse-dynamics model can be learned from the action-379 labelled data to "fill in" missing actions (Schmeckpeper et al., 2021; Zheng et al., 2023; Baker et al., 380 2022). By contrast, in this work we are interested in provable sample-efficiency of representation 381 learning, and assume no access to action-labeled data during pretraining.

Learning in Ex-BMDPs. As discussed throughout this work, numerous prior works consider the Ex-BMDP model (Efroni et al., 2022; Mhammedi et al., 2024), including in the offline setting (Islam et al., 2023; Lamb et al., 2023; Levine et al., 2024). Misra et al. (2024) in particular demonstrates a hardness result: that Ex-BMDP latent representations cannot be learned in general from offline action-free data. In this work, we demonstrate a special case where this representation learning problem is in fact tractable: the case where offline data from multiple diverse agents are available.

388 6 Discussion and Limitations

389 One major assumption of this work (as well as Misra et al. (2024); Islam et al. (2023) and other 390 prior works) is that offline data are collected by a policy which acts independently of observation 391 noise. This assumption stems from the fact that, if noise features influence the behavioral policy, 392 they (indirectly) influence the latent-state dynamics of the agent: these noise features may then be 393 erroneously captured in the learned representation. However, in real-world settings, it may actually 394 be beneficial to capture such features in the learned representation: if "expert" agents are relying 395 on some uncontrollable feature, this feature may be relevant to the expert agents' reward functions, 396 and may therefore also be relevant to downstream tasks for which our learned representations will 397 be used. Therefore, the noise-independent policy assumption might not be necessary in practice.

398 An additional restrictive assumption of this work is that the latent dynamics are deterministic, and 399 that each episode starts at the same latent state s_1^* . However, this assumption is also essentially 400 present even in the best-known result for provably sample-efficient Ex-BMDP representation learn-401 ing in the online setting (Efroni et al., 2022) – that work does allow for rare departures from deter-402 ministic dynamics, however, and it may be possible to adapt the analysis of CRAFT to that setting 403 as well, although we have focused on the strictly-deterministic case here for ease of presentation. 404 Mhammedi et al. (2024) proposes an online algorithm for learning Ex-BMDPs with nondeterminis-405 tic latent dynamics, but that work assumes "simulator access": the ability to reset the environment 406 to any previously-visited observation. Several works (Lamb et al., 2023; Levine et al., 2024; Islam 407 et al., 2023) consider learning Ex-BMDPs from offline data (with action labels) without assuming 408 restarts to s_1^* : these works are "practical" algorithms that do not provide sample-complexity guar-409 antees. A similar "practical" algorithm for the action-free, multiple-agent setting based on the ideas 410 presented in this work may also be possible.

411 The assumption that two policies differ substantially at *every* latent state may also be impractical.

412 One direction for future work may be to leverage data from *several* agents, such that it is more likely 413 that *some* agent has a distinct behavior at each latent state.

414 While access to training oracles is a common assumption in representation learning (Agarwal et al.,

415 2020; Efroni et al., 2022; Uehara et al., 2022), the optimization of Equation 13, on a discretized

416 domain, may be troublesome in practice. Additionally, the sample complexity bounds in Equation

417 16, while polynomial, may not be optimal: these issues are potential directions for future work.

418 **References**

Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. Flambe: Structural complexity and representation learning of low rank mdps. *Advances in neural information processing systems*, 33:20095–20107, 2020.

Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon
Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching
unlabeled online videos. *Advances in Neural Information Processing Systems*, 35:24639–24654,
2022.

Simon Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John Langford.
Provably efficient rl with rich observations via latent state decoding. In *International Conference on Machine Learning*, pp. 1665–1674. PMLR, 2019.

Ashley Edwards, Himanshu Sahni, Yannick Schroecker, and Charles Isbell. Imitating latent policies
 from observation. In *International conference on machine learning*, pp. 1755–1763. PMLR, 2019.

Yonathan Efroni, Dipendra Misra, Akshay Krishnamurthy, Alekh Agarwal, and John Langford.
Provably filtering exogenous distractors using multistep inverse dynamics. In *International Con- ference on Learning Representations*, 2022. URL https://openreview.net/forum?
id=RQLLzMCefQu.

Dibya Ghosh, Chethan Anand Bhateja, and Sergey Levine. Reinforcement learning from passive
data via latent intentions. In *International Conference on Machine Learning*, pp. 11321–11339.
PMLR, 2023.

Riashat Islam, Manan Tomar, Alex Lamb, Yonathan Efroni, Hongyu Zang, Aniket Rajiv Didolkar,
Dipendra Misra, Xin Li, Harm Van Seijen, Remi Tachet Des Combes, et al. Principled offline rl in
the presence of rich exogenous information. In *International Conference on Machine Learning*,
pp. 14390–14421. PMLR, 2023.

Alex Lamb, Riashat Islam, Yonathan Efroni, Aniket Rajiv Didolkar, Dipendra Misra, Dylan J Foster,
Lekan P Molu, Rajan Chari, Akshay Krishnamurthy, and John Langford. Guaranteed discovery
of control-endogenous latent states with multi-step inverse models. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?
id=TNocbXm5MZ.

- Alexander Levine, Peter Stone, and Amy Zhang. Multistep inverse is not all you need. *Reinforce- ment Learning Journal*, 2:884–925, 2024.
- Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy
 Zhang. VIP: Towards universal visual reward and representation via value-implicit pre-training.
 In *The Eleventh International Conference on Learning Representations*, 2023. URL https:
 //openreview.net/forum?id=YJ7o2wetJ2.
- Willi Menapace, Stephane Lathuiliere, Sergey Tulyakov, Aliaksandr Siarohin, and Elisa Ricci.
 Playable video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10061–10070, 2021.
- Zakaria Mhammedi, Dylan J Foster, and Alexander Rakhlin. The power of resets in online rein forcement learning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=7sACcaOmGi.
- Dipendra Misra, Mikael Henaff, Akshay Krishnamurthy, and John Langford. Kinematic state ab straction and provably efficient rich-observation reinforcement learning. In *International confer- ence on machine learning*, pp. 6961–6971. PMLR, 2020.

462 Dipendra Misra, Akanksha Saran, Tengyang Xie, Alex Lamb, and John Langford. Towards princi 463 pled representation learning from videos for reinforcement learning. In *The Twelfth International* 464 *Conference on Learning Representations*, 2024. URL https://openreview.net/forum?
 465 id=3mnWvUZIXt.

- Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation. In *Conference on Robot Learning*, pp. 892–909.
 PMLR, 2023.
- Seohong Park, Dibya Ghosh, Benjamin Eysenbach, and Sergey Levine. HIQL: Offline goal conditioned RL with latent states as actions. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=cLQCCtVDuW.
- Karl Schmeckpeper, Oleh Rybkin, Kostas Daniilidis, Sergey Levine, and Chelsea Finn. Reinforcement learning with videos: Combining offline observations with interaction. In Jens Kober, Fabio
 Ramos, and Claire Tomlin (eds.), *Proceedings of the 2020 Conference on Robot Learning*, volume 155 of *Proceedings of Machine Learning Research*, pp. 339–354. PMLR, 16–18 Nov 2021.
 URL https://proceedings.mlr.press/v155/schmeckpeper21a.html.
- 477 Dominik Schmidt and Minqi Jiang. Learning to act without actions. In *The Twelfth International* 478 *Conference on Learning Representations*, 2024. URL https://openreview.net/forum?
 479 id=rvUq3cxpDF.
- Younggyo Seo, Kimin Lee, Stephen L James, and Pieter Abbeel. Reinforcement learning with
 action-free pre-training from videos. In *International Conference on Machine Learning*, pp.
 19561–19579. PMLR, 2022.
- Masatoshi Uehara, Xuezhou Zhang, and Wen Sun. Representation learning for online and offline
 RL in low-rank MDPs. In *International Conference on Learning Representations*, 2022. URL
 https://openreview.net/forum?id=J4iSIR9fhY0.
- Haoran Xu, Li Jiang, Jianxiong Li, and Xianyuan Zhan. A policy-guided imitation approach
 for offline reinforcement learning. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and
 Kyunghyun Cho (eds.), Advances in Neural Information Processing Systems, 2022. URL
 https://openreview.net/forum?id=CKbqDtZnSc.
- Weirui Ye, Yunsheng Zhang, Pieter Abbeel, and Yang Gao. Become a proficient player with lim ited data through watching pure videos. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=Sy-o2N0hF4f.
- 493 Qinqing Zheng, Mikael Henaff, Brandon Amos, and Aditya Grover. Semi-supervised offline rein 494 forcement learning with action-free trajectories. In *International conference on machine learning*,
 495 pp. 42339–42362. PMLR, 2023.

496	Supplementary Materials
497 498	The following content was not necessarily subject to peer review.

499 A Hypothesis Classes and Realizability Assumptions

As mentioned in Section 2.5, we assume access to hypothesis classes of encoder functions $\Phi_{1:H}$. We make a realizability assumption: that is, the true encoder $\phi_h^* \in \Phi_h$. Moreover, we assume that for any arbitrary permutation σ , $\sigma(\phi_h^*(x)) \in \Phi_h$ – this allows us to train an encoder in Φ_h on datasets of observations representing each latent state without knowing the "correct" ordering of the latent states. Similar realizability assumptions are common in representation learning literature for structured MDPs (Du et al., 2019; Efroni et al., 2022; Misra et al., 2020; 2024; Uehara et al., 2022; Agarwal et al., 2020).

We use a second set of hypothesis classes, $\mathcal{G}_h \subseteq \mathcal{X}_h \to \{0, 1\}$, for which, for any *pair* of latent states $s_h^*, s_h'^*$, there exists some $g \in \mathcal{G}_h$ that can perfectly distinguish observations of s_h^* from observations of $s_h'^*$. In our sample-complexity results, we assume that $|\mathcal{G}_h| \le |\Phi_h|$. In this work, we are chiefly concerned with sample-complexity: we make use of training oracles for a variety of loss functions which may not be tractable to optimize in practice. See Section 6 for further discussion.

512 **B** Algorithm

513 The full CRAFT algorithm is presented as Algorithm 2.

514 C Proofs

- 515 In this section, we prove the correctness and sample complexity bounds of CRAFT presented in
- 516 Theorem 3.1. First, though, we prove various lemmas the will be helpful in proving the final result.
- 517 C.1 Preliminary Note
- 518 Recall Equation 1 in the main text:

$$\Pr(\tau_A, \tau_B) = \Pr(\phi^*(\tau_A), \phi^*(\tau_B)) \cdot \Pr_{\substack{P_1^e, \mathcal{T}^e}}(\phi^e(\tau_A)) \cdot \Pr_{\substack{P_1^e, \mathcal{T}^e}}(\phi^e(\tau_B))$$
$$\cdot \Pr_Q(\tau_A | \phi^*(\tau_A), \phi^e(\tau_A)) \cdot \Pr_Q(\tau_B | \phi^*(\tau_B), \phi^e(\tau_B))$$

Throughout our proofs, we will make use of this assumption in the following way: we will treat the controllable latent state trajectories $\phi^*(\tau_A)$, $\phi^*(\tau_B)$ as *fixed but arbitrary*, not as random variables, and treat the exogenous noise Markov chains $\phi^e(\tau_A)$, $\phi^e(\tau_B)$ and the emission function Q as the only random variables. Then, if the algorithm succeeds with high probability for any such *fixed*, *arbitrary* $\phi^*(\tau_A)$, $\phi^*(\tau_B)$, we can conclude by the independence assumption that it also succeeds with high probability under any data-generating process for which Equation 1 holds.

525 C.2 Concentration Lemmas

In this section, we present concentration bounds on the loss functions used in Algorithm 2. We start with the log-odds loss given in Equation 13:

Lemma C.1. Given m distributions $\mathcal{D}_1, ..., \mathcal{D}_m \in \mathcal{P}(\mathcal{X})$, each with two corresponding positive integers a_i, b_i , for $i \in [m]$, let $A_i \sim \mathcal{D}^{a_i}$ and $B_i \sim \mathcal{D}^{b_i}$ be two multi-sets consisting of a_i and b_i i.i.d. samples from \mathcal{D}_i , respectively. Then, for any $\xi > 0$ and $n_{\Xi} \in \mathbb{N}_+$ such that $\forall i, |\ln(a_i/b_i)| \leq \frac{n_{\Xi}\xi}{2}$, let $\Xi = \{-\frac{n_{\Xi}\xi}{2}, -\frac{n_{\Xi}\xi}{2} + \xi, -\frac{n_{\Xi}\xi}{2} + 2\xi, ..., \frac{n_{\Xi}\xi}{2}\}$. Further, let $\bar{c}_i \in \Xi$ be the smallest value in Ξ greater than or equal to $\ln(a_i/b_i)$, and $\underline{c}_i \in \Xi$ be the largest value in Ξ less than or equal to $\ln(a_i/b_i)$. Also, assume that $\forall i, \frac{m_i+b_i}{\sum_{i'=1}^m a_{i'}+b_{i'}} \geq \nu$.

534 *Given any function* $f \in \mathcal{X} \to \Xi$ *, define:*

$$\mathcal{L}(f) := \sum_{i=1}^{m} \left[\sum_{x \in A_i} \ln(e^{-f(x)} + 1) + \sum_{x \in B_i} \ln(e^{f(x)} + 1) \right].$$
 (18)

Algorithm 2 CRAFT

Require: Trajectory datasets τ_A , τ_B , known lower-bounds α , η , and ν , encoder function classes $\Phi_h \subseteq \mathcal{X}_h \rightarrow \mathcal{X}_h$ N_s and classification function class $\mathcal{G}_h \subseteq \mathcal{X}_h \to \{0, 1\}$. 1: $\alpha \leftarrow \min(1, \alpha)$. 2: Let $\xi := \alpha/4$; $n_{\Xi} := \lceil 8 \ln(\eta^{-1} - 1)/\alpha \rceil$. 3: $\eta \leftarrow 1/(1 + e^{n_{\Xi}\alpha/8})$. 4: Let $\Xi := \{-\frac{n_{\Xi}\xi}{2}, -\frac{n_{\Xi}\xi}{2} + \xi, -\frac{n_{\Xi}\xi}{2} + 2\xi, ..., \frac{n_{\Xi}\xi}{2}\}.$ 5: Initialize $S_1 := \{s_1\}, D_{A,1}(s_1) := [|\tau_A|], D_{B,1}(s_1) := [|\tau_B|].$ \\ First timestep should have a single state rep., associated with every trajectory index. 6: Let $\phi'_1 := \mathcal{X}_1 \to 0$. 7: for $h \in \{1, 2, ..., H - 1\}$ do Initialize $S_{h+1} := \{\}$. 8: Let the inverse-actor-prediction function class $\mathcal{F}_h \subseteq \mathcal{X} \times \mathcal{X} \to \Xi$ be composed as $\mathcal{F}_h = \Phi_h \times \Phi_{h+1} \times$ 9: $(N_s^2 \to \Xi).$ 10: Find the $f_h \in \mathcal{F}_h$ which minimizes Equation 13. 11: Let $q_{thresh.} := \frac{h\nu}{8H}$. 12: for $s_{pred} \in S_h$ do Initialize merged_already(s) := False for all $s \in S_{h+1}$. 13: 14: Initialize $S_{new} := \{\}$ Let pairs $(s_{pred}) := (\tau_A)_{h:h+1}[D_{A,h}(s_{pred})] \cup (\tau_B)_{h:h+1}[D_{B,h}(s_{pred})].$ 15: \\ All transitions which start at s_{prev} . $\forall j \in \{0, ..., n_{\Xi}\}$, Let pred_succ $[j] := \{(x_h, x_{h+1}) \in \text{pairs}(s_{pred}) | f(x_h, x_{h+1}) = j\xi - \frac{n_{\Xi}\xi}{2} \}$. 16: Initialize j := 017: while $j \leq n_{\Xi}$ do 18: if $|\text{pred_succ}[j]| \ge q_{thresh.}(|\tau_A| + |\tau_B|)$ then 19: Let j' be the minimum integer > j such that $|\text{pred}_\text{succ}[j']| < q_{thresh.}(|\tau_A| + |\tau_B|)$, or n_{Ξ} if 20: no such integer exists. Let $\mathcal{D}_{new_pairs} := \{x' | (x, x') \in \bigcup_{k=\max(0, i-1)}^{j'} \text{pred_succ}[k]\}, \mathcal{D}_{new} := \{x' | (x, x') \in U_{k} \}$ 21: \mathcal{D}_{new_pairs} } 22: Initialize new_state? \leftarrow True. for $s \in S_{h+1}$, such that merged_already(s) == False do 23: Let $\mathcal{D}_s := (\tau_A)_{h+1} [D_{A,h+1}(s)] \uplus (\tau_B)_{h+1} [D_{B,h+1}(s)]$ 24: Train a classifier $g \in \mathcal{G}$ to distinguish \mathcal{D}_{new} and \mathcal{D}_s , with loss $\mathcal{L}(g)$ given in Equation 53. 25: 26: if the loss $\mathcal{L}(g)$ on \mathcal{D}_{new} and \mathcal{D}_s is > 0.5 then Append to $D_{A,h+1}(s)$ indices of trajectories in τ_A that observations in \mathcal{D}_{new} are from. 27: 28: Append to $D_{B,h+1}(s)$ indices of trajectories in τ_B that observations in \mathcal{D}_{new} are from. merged_already?(s) \leftarrow True 29: $new_state? \gets False$ 30: break. 31: 32: end if 33: end for 34: if new_state? then Add new state s_{new} to S_{new} 35: Initialize $D_{A,h+1}(s_{new})$ as indices of trajectories in τ_A that observations in \mathcal{D}_{new} are from. 36: 37: Initialize $D_{B,h+1}(s_{new})$ as indices of trajectories in τ_B that observations in \mathcal{D}_{new} are from. 38: end if 39: $j \leftarrow j' + 2$ else 40: 41: $j \gets j + 1$ end if 42. end while 43: $\mathcal{S}_{h+1} := \mathcal{S}_{h+1} \cup \mathcal{S}_{new}$ 44: 45: end for $\arg\min_{\phi\in\Phi_{h+1}}\sum_{s\in\mathcal{S}_{h+1}}\Big[\frac{1}{|\mathcal{D}_s|}\sum_{x\in\mathcal{D}_s}(1-\mathbb{1}_{(\phi(x)=s)})\Big],\quad\text{where}\quad\mathcal{D}_s$ 46: ϕ_{h+1}' :=:= $(\tau_A)_{h+1}[D_{A,h+1}(s)] \uplus (\tau_B)_{h+1}[D_{B,h+1}(s)]$ 47: end for 48: **Return:** $\phi'_1, ..., \phi'_H$

$$\mathcal{L}_{ref} := \sum_{i=1}^{m} \left[\min_{c_i \in \{\bar{c}_i, \underline{c}_i\}} a_i \ln \left(e^{-c_i} + 1 \right) + b_i \ln \left(e^{c_i} + 1 \right) \right],\tag{19}$$

536 and let $\eta := (e^{n_{\Xi}\xi/2} + 1)^{-1}$. For any ϵ and δ , if:

$$\forall i \in [m], \ a_i + b_i \ge \frac{50 \ln(2/\delta) \ln^2(1/\eta)}{\nu \epsilon^2 \eta^2 \xi^4}$$
 (20)

537 then the probability that both $\mathcal{L} \leq \mathcal{L}_{ref}$, and:

$$\exists i \in [m]: \left| \{ x \in A_i \uplus B_i | f(x) \notin \{ \underline{c}_i, \overline{c}_i \} \right| > \epsilon(a_i + b_i)$$
(21)

538 is at most δ .

539 Proof. We can define

$$\mathcal{L}_{\text{ref}}^{i} := \min_{c_{i} \in \{\bar{c}_{i}, \underline{c}_{i}\}} a_{i} \ln\left(e^{-c_{i}} + 1\right) + b_{i} \ln\left(e^{c_{i}} + 1\right),$$
(22)

540 So that $\mathcal{L}_{ref} = \sum_{i=1}^{m} \mathcal{L}_{ref}^{i}$, and also define

$$\mathcal{L}^{i}_{pop}(f) := a_i \left(\underset{x \sim \mathcal{D}_i}{\mathbb{E}} \ln(e^{-f(x)} + 1) \right) + b_i \left(\underset{x \sim \mathcal{D}_i}{\mathbb{E}} \ln(e^{f(x)} + 1) \right)$$
(23)

- 541 and $\mathcal{L}_{pop} := \sum_{i=1}^{m} \mathcal{L}_{pop}^{i}$.
- 542 First, we consider the "population loss" for each distribution:

$$\mathcal{L}_{pop}^{i}(f) = a_{i} \left(\underset{x \sim \mathcal{D}_{i}}{\mathbb{E}} \ln(e^{-f(x)} + 1) \right) + b_{i} \left(\underset{x \sim \mathcal{D}_{i}}{\mathbb{E}} \ln(e^{f(x)} + 1) \right)$$

$$= a_{i} \left(\sum_{\zeta \in \Xi} \Pr_{x \sim \mathcal{D}_{i}} (f(x) = \zeta) \cdot \ln(e^{-\zeta} + 1) \right)$$

$$+ b_{i} \left(\sum_{\zeta \in \Xi} \Pr_{x \sim \mathcal{D}_{i}} (f(x) = \zeta) \cdot \ln(e^{\zeta} + 1) \right)$$

$$= \sum_{\zeta \in \Xi} \Pr_{x \sim \mathcal{D}_{i}} (f(x) = \zeta) \cdot (a_{i} \ln(e^{-\zeta} + 1) + b_{i} \ln(e^{\zeta} + 1))$$

$$= \sum_{\zeta \in \Xi} \Pr_{x \sim \mathcal{D}_{i}} (f(x) = \zeta) \cdot a_{i} \left(\left(1 + \frac{b_{i}}{a_{i}} \right) \ln(e^{\zeta} + 1) - \zeta \right)$$
(24)

543 We can define $h_{\gamma}(\zeta) := (1 + \gamma^{-1}) \ln(e^{\zeta} + 1) - \zeta$, so that

$$\mathcal{L}^{i}_{pop}(f) = a_{i} \sum_{\zeta \in \Xi} \Pr_{x \sim \mathcal{D}_{i}}(f(x) = \zeta) \cdot h_{a_{i/b_{i}}}(\zeta)$$
(25)

544 Now, note that:

$$h'_{\gamma}(\zeta) = \left(1 + \gamma^{-1}\right) \frac{e^{\zeta}}{e^{\zeta} + 1} - 1$$
(26)

545 and

$$h_{\gamma}''(\zeta) = \left(1 + \gamma^{-1}\right) \frac{e^{\zeta}}{(e^{\zeta} + 1)^2}.$$
(27)

546 Then, we see that $h_{\gamma}(\zeta)$ is a convex function, with a global minimum at $\zeta = \ln(\gamma)$, and second-547 derivative at least

$$(1+\gamma^{-1}) \frac{e^{n_{\Xi}\xi/2}}{(e^{n_{\Xi}\xi/2}+1)^2} \left(=(1+\gamma^{-1}) \frac{e^{-n_{\Xi}\xi/2}}{(e^{-n_{\Xi}\xi/2}+1)^2}\right)$$
(28)

- 548 everywhere on the interval $[-n_{\Xi}\xi/2, n_{\Xi}\xi/2]$.
- 549 Due to the convexity of $h_{a_{i/b_{i}}}(\zeta)$, we have, $\forall j > 0$,

$$a_i \cdot h_{a_i/b_i}(\ln(a_i/b_i)) \le \min(a_i \cdot h_{a_i/b_i}(\bar{c}_i), a_i \cdot h_{a_i/b_i}(\underline{c}_i)) (= \mathcal{L}_{ref}^i) \le a_i \cdot h_{a_i/b_i}(\bar{c}_i) < a_i \cdot h_{a_i/b_i}(\bar{c}_i+j\xi)$$

550 and, similarly, $\forall j > 0$:

$$a_i \cdot h_{a_i/b_i}(\ln(a_i/b_i)) \le \min(a_i \cdot h_{a_i/b_i}(\bar{c}_i), a_i \cdot h_{a_i/b_i}(\underline{c}_i)) (= \mathcal{L}_{ref}^i) \le a_i \cdot h_{a_i/b_i}(\underline{c}_i) < a_i \cdot h_{a_i/b_i}(\underline{c}_i - j\xi).$$

551 In particular, by a second-order Taylor bound, we have that, for j > 0:

$$\mathcal{L}_{ref}^{i} \leq a_{i} \cdot h_{a_{i/b_{i}}}(\bar{c}_{i}+j\xi) - a_{i} \cdot \left(1 + (a_{i}/b_{i})^{-1}\right) \frac{e^{n_{\Xi}\xi/2}}{(e^{n_{\Xi}\xi/2}+1)^{2}} \cdot \frac{(j\xi)^{2}}{2}$$

$$\leq a_{i} \cdot h_{a_{i/b_{i}}}(\bar{c}_{i}+j\xi) - (a_{i}+b_{i}) \frac{e^{n_{\Xi}\xi/2}}{(e^{n_{\Xi}\xi/2}+1)^{2}} \cdot \frac{\xi^{2}}{2}$$
(29)

552 and similarly for $a_i \cdot h_{a_i/b_i}(\underline{c}_i - j\xi)$:

$$\mathcal{L}_{ref}^{i} \le a_{i} \cdot h_{a_{i/b_{i}}}(\underline{c}_{i} - j\xi) - (a_{i} + b_{i}) \frac{e^{n_{\Xi}\xi/2}}{(e^{n_{\Xi}\xi/2} + 1)^{2}} \cdot \frac{\xi^{2}}{2}.$$
(30)

553 In particular, by Equation 25,

$$\mathcal{L}^{i}_{pop}(f) \ge \mathcal{L}^{i}_{ref} + \Pr_{x \sim \mathcal{D}_{i}}(f(x) \notin \{\underline{c}_{i}, \overline{c}_{i}\}) \cdot (a_{i} + b_{i}) \frac{e^{n_{\Xi}\xi/2}}{(e^{n_{\Xi}\xi/2} + 1)^{2}} \cdot \frac{\xi^{2}}{2}.$$
 (31)

554 In terms of η , this is:

$$\mathcal{L}_{pop}^{i}(f) \geq \mathcal{L}_{ref}^{i} + \Pr_{x \sim \mathcal{D}_{i}}(f(x) \notin \{\underline{c}_{i}, \overline{c}_{i}\}) \cdot (a_{i} + b_{i}) (\eta - \eta^{2}) \cdot \frac{\xi^{2}}{2}$$

$$\geq \mathcal{L}_{ref}^{i} + \Pr_{x \sim \mathcal{D}_{i}}(f(x) \notin \{\underline{c}_{i}, \overline{c}_{i}\}) \cdot (a_{i} + b_{i}) \cdot \frac{\eta\xi^{2}}{4}$$
(32)

555 where we use the fact that $\eta \leq 1/2$ in the last inequality. This gives us:

$$\Pr_{x \sim \mathcal{D}_i}(f(x) \notin \{\underline{c}_i, \overline{c}_i\}) \le \frac{4(\mathcal{L}_{pop}^i(f) - \mathcal{L}_{ref}^i)}{(a_i + b_i)\eta\xi^2}$$
(33)

556 Because $\mathcal{L}_{pop}^{i} - \mathcal{L}_{ref}^{i} \geq 0$, this implies:

$$\begin{aligned} \forall i \in [m], \\ (a_i + b_i) \cdot \Pr_{x \sim \mathcal{D}_i}(f(x) \notin \{\underline{c}_i, \overline{c}_i\}) &\leq \frac{4(\mathcal{L}_{pop}^i(f) - \mathcal{L}_{ref}^i)}{\eta \xi^2} \leq \frac{4\sum_{i=1}^m (\mathcal{L}_{pop}^i(f) - \mathcal{L}_{ref}^i)}{\eta \xi^2} \\ &= \frac{4(\mathcal{L}_{pop}(f) - \mathcal{L}_{ref})}{\eta \xi^2} \end{aligned}$$
(34)

557 Meanwhile, from Equations 18 and 23 applying (one-sided) Hoeffding's lemma gives us, with prob-558 ability at least $1 - \delta/2$:

$$\mathcal{L}(f) - \mathcal{L}_{pop}(f) + \sqrt{\sum_{i=1}^{m} (a_i + b_i) \ln(1/\eta) \sqrt{2\ln(2/\delta)}} \ge 0.$$
(35)

559 which implies, by assumption:

$$\forall i, \ \mathcal{L}(f) - \mathcal{L}_{pop}(f) + \sqrt{a_i + b_i} \sqrt{1/\nu} \ln(1/\eta) \sqrt{2\ln(2/\delta)} \ge 0.$$
(36)

560 Combining Equations 34 and 36 gives, with probability at least $1 - \delta/2$, we have that $\forall i \in [m]$,

$$(a_i + b_i) \Pr_{x \sim \mathcal{D}_i}(f(x) \notin \{\underline{c}_i, \overline{c}_i\}) \le \frac{4(\mathcal{L}(f) - \mathcal{L}_{ref})}{\eta \xi^2} + \frac{4(\sqrt{a_i + b_i})\ln(1/\eta)\sqrt{2\ln(2/\delta)}}{\sqrt{\nu}\eta \xi^2}.$$
 (37)

561 Then, with probability at least $1 - \delta/2$, the condition $\mathcal{L}(f) \leq \mathcal{L}_{ref}$ implies:

$$\forall i \in [m], \ (a_i + b_i) \Pr_{x \sim \mathcal{D}_i}(f(x) \notin \{\underline{c}_i, \overline{c}_i\}) \le \frac{4\sqrt{a_i + b_i}\ln(1/\eta)\sqrt{2\ln(2/\delta)}}{\sqrt{\nu}\eta\xi^2}.$$
 (38)

- 562 We can apply Hoeffding's lemma once for each $i \in [m]$, to the binary variable of whether on not
- 563 $f(x) \in \{\underline{c}_i, \overline{c}_i\}$, where x is sampled $(a_i + b_i)$ times to produce the dataset $A_i \cup B_i$. By union bound,
- 564 we have, with probability at least 1δ , $\mathcal{L}(f) \leq \mathcal{L}_{ref}$ implies:

$$\forall i \in [m], \quad \left| \left\{ x \in A_i \cup B_i | f(x) \notin \{\underline{c}_i, \overline{c}_i\} \right| \leq \\ \sqrt{(a_i + b_i) \ln(2m/\delta)/2} + \frac{4\sqrt{a_i + b_i} \ln(1/\eta) \sqrt{2 \ln(2/\delta)}}{\sqrt{\nu} \eta \xi^2} \leq \\ \sqrt{(a_i + b_i)} \sqrt{2} \left(\frac{\sqrt{\ln(2m/\delta)}}{2} + \frac{4\sqrt{\ln(2/\delta)} \ln(1/\eta)}{\sqrt{\nu} \eta \xi^2} \right) \leq \\ \sqrt{(a_i + b_i)} \sqrt{2} \frac{\ln(1/\eta)}{\eta \xi^2} \left(\frac{\sqrt{\ln(2/\delta)}}{2} + \frac{\sqrt{\ln(m)}}{2} + \frac{4\sqrt{\ln(2/\delta)}}{\sqrt{\nu}} \right),$$

$$(39)$$

where in the last line, we used triangle inequality and the fact that $\eta = 1/(e^{n_{\Xi}\xi/2}+1) \le 1/(e^{\xi/2}+1)$, which in turn implies:

$$\frac{\ln(1/\eta)}{\eta\xi^2} \ge \frac{(e^{\xi/2} + 1)\ln(e^{\xi/2} + 1)}{\xi^2} > 1 \quad (\forall \xi > 0).$$
(40)

Note that, because each $a_i + b_i$ contains at least a ν -fraction of the total $\sum_i^m a_i + b_i$, we must have $m \le 1/\nu$. Then:

$$\frac{\sqrt{\ln(m)}}{2} \le \frac{\sqrt{\ln(1/\nu)}}{2} \le \frac{1}{2\sqrt{e}\sqrt{\nu}} \le \frac{1}{2\sqrt{e}\sqrt{\ln(2)}} \frac{\sqrt{\ln(2/\delta)}}{\sqrt{\nu}} \le \frac{\sqrt{\ln(2/\delta)}}{2\sqrt{\nu}}$$
(41)

569 Therefore (and noting $\nu < 1$), we can combine terms in Equation 39 to conclude:

$$\forall i \in [m], \ \left| \left\{ x \in A_i \cup B_i | f(x) \notin \left\{ \underline{c}_i, \overline{c}_i \right\} \right| \le \frac{5\sqrt{2(a_i + b_i)\ln(2/\delta)}\ln(1/\eta)}{\sqrt{\nu}\eta\xi^2}.$$
(42)

570 Now, to ensure $\forall i \in [m]$, $\left| \{ x \in A_i \cup B_i | f(x) \notin \{\underline{c}_i, \overline{c}_i\} \right| \leq \epsilon(a_i + b_i)$, we need,

$$\forall i \in [m], \quad \frac{5\sqrt{2(a_i + b_i)\ln(2/\delta)}\ln(1/\eta)}{\sqrt{\nu}\eta\xi^2} \le \epsilon(a_i + b_i) \tag{43}$$

571 or:

$$\forall i \in [m], \ \frac{50\ln(2/\delta)\ln^2(1/\eta)}{\nu\epsilon^2\eta^2\xi^4} \le a_i + b_i$$
(44)

as provided by Equation 20. Note that because the implication

$$\mathcal{L}(f) \le \mathcal{L}_{ref} \to \forall i \in [m], \ \left| \{ x \in A_i \cup B_i | f(x) \notin \{ \underline{c}_i, \overline{c}_i \} \right| \le \epsilon(a_i + b_i)$$
(45)

573 holds with probability at least $(1 - \delta)$, this implication can only be *broken*, by the case that

$$\mathcal{L}(f) \le \mathcal{L}_{ref} \land \exists i \in [m] : \left| \{ x \in A_i \cup B_i | f(x) \notin \{ \underline{c}_i, \overline{c}_i \} \right| > \epsilon(a_i + b_i)$$
(46)

574 with probability at most δ .

Corollary C.2. Let $\mathcal{D}^*(s_h^*, s_{h+1}^*)$ be the multiset of observation pairs (x_h, x_{h+1}) from both τ_A and τ_B in Algorithm 1, such that $\phi_h^*(x_h) = s_h^*$ and $\phi_{h+1}^*(x_{h+1}) = s_{h+1}^*$, and let $\mathcal{D}_A^*(s_h^*, s_{h+1}^*)$ and $\mathcal{D}_B^*(s_h^*, s_{h+1}^*)$ be the elements of $\mathcal{D}^*(s_h^*, s_{h+1}^*)$ originating from τ_A and τ_B 578 respectively. Further, let $\bar{c}_{s_h^*, s_{h+1}^*} \in \Xi$ be the smallest value in Ξ greater than or equal to $\ln(|\mathcal{D}_A^*(s_h^*, s_{h+1}^*)|/|\mathcal{D}_B^*(s_h^*, s_{h+1}^*)|)$, and $\underline{c}_{s_h^*, s_{h+1}^*} \in \Xi$ be the largest value in Ξ less than or equal 580 to $\ln(|\mathcal{D}_A^*(s_h^*, s_{h+1}^*)|/|\mathcal{D}_B^*(s_h^*, s_{h+1}^*)|)$.

581 Further, assume the realizability condition that $\phi_h^* \in \Phi_h$ and $\phi_{h+1}^* \in \Phi_{h+1}$.

582 If $\forall s_h^*, s_{h+1}^*$, such that s_h^* can transition to s_{h+1}^* ,

$$|\mathcal{D}^*(s_h^*, s_{h+1}^*)| \ge \frac{50(\ln(2|\Phi|^2/\delta) + N_s^2 \ln(n_{\Xi} + 1))\ln^2(1/\eta)}{\nu\epsilon^2 \eta^2 \xi^4}$$
(47)

then with probability at least $1 - \delta$, the function $f(\cdot)$ found in Line 10 of Algorithm 2 will be such that, $\forall s_h^*, s_{h+1}^*$, such that s_h^* can transition to s_{h+1}^* ,

$$\left| \left\{ x \in \mathcal{D}^*(s_h^*, s_{h+1}^*) | f(x) \notin \left\{ \underline{c}_{s_h^*, s_{h+1}^*}, \overline{c}_{s_h^*, s_{h+1}^*} \right\} \right| \le \epsilon \left| \mathcal{D}^*(s_h^*, s_{h+1}^*) \right|.$$
(48)

585

Proof. By application of Lemma C.1 with $\delta' := \delta/(|\Phi|^2 \cdot (n_{\Xi} + 1)^{(N_s^2)}) \leq \delta/|\mathcal{F}_h|$, we have that, 586 for any fixed hypothesis f', the probability that $\mathcal{L}(f') \leq \mathcal{L}_{ref}$ and Equation 48 is violated is at most 587 $\delta/|\mathcal{F}_h|$. Then by union bound, the probability that any such f' exists in \mathcal{F}_h is at most $1-\delta$. However, 588 by the realizability assumption, we know that an f^* exists in \mathcal{F} which achieves loss $\mathcal{L}(f^*) = \mathcal{L}_{ref}$ 589 and also that respects Equation 48. (In particular, this f^* is simply (ϕ_h^*, ϕ_{h+1}^*) composed with a 590 591 mapping from the representations corresponding to each (s_h^*, s_{h+1}^*) to the corresponding $\underline{c}_{s_h^*, s_{h+1}^*}$ or $\bar{c}_{s_h^*,s_{h+1}^*}$ which minimizes Equation 22.) Therefore with probability at least $1-\delta$, the $f \in \mathcal{F}_h^{n+1}$ which 592 593 minimizes $\mathcal{L}(f)$ must respect Equation 48.

We now give two simple results for classification under corrupted data. First though, we prove a minor claim, which is simply some "deferred algebra" for the lemmas which follow:

Proposition C.3. Consider a multiset $Z = \{z_1, ..., z_m\}$ of items $z_i \in [0, 1]$, and a modified multiset 597 Z', also consisting of items in [0, 1], such that the symmetric difference between Z and Z' has size 598 at most k (that is, Z' can be constructed from Z by inserting and/or removing a total of at most k599 items). Then

$$\left|\sum_{\mathcal{Z}} \frac{z}{|\mathcal{Z}|} - \sum_{\mathcal{Z}'} \frac{z'}{|\mathcal{Z}'|}\right| \le \frac{k}{m} \tag{49}$$

600 Proof. Define $Z_{removed}, Z_{added}$, and Z_{kept} such that $Z = Z_{kept} + Z_{removed}$, and $Z' = Z_{kept} + Z_{added}$. Note that $k = |Z_{removed}| + |Z_{added}|$ and $m = |Z_{removed}| + |Z_{kept}|$. We first as-602 sume that $|Z| \ge |Z'|$. In other words, we assume $|Z_{removed}| \ge |Z_{added}|$. Then we can con-603 struct Z' from Z by (1) removing some arbitrary subset $Z'_{removed} \subseteq Z_{removed}$ from Z, such that 604 $|Z'_{removed}| = |Z_{added}|$; then (2) inserting the samples Z_{added} ; and then finally (3) removing the 605 multiset $Z''_{removed} = Z_{removed} \setminus Z'_{removed}$. Let the intermediate set constructed after step (2) be 606 $Z'' := (Z \setminus Z'_{removed}) \uplus Z_{added} = Z' \uplus Z''_{removed}$. Note that |Z| = |Z''|, and

$$\left|\sum_{\mathcal{Z}} \frac{z}{|\mathcal{Z}|} - \sum_{\mathcal{Z}''} \frac{z''}{|\mathcal{Z}''|}\right| = \left|\sum_{\mathcal{Z}} \frac{z}{|\mathcal{Z}|} - \sum_{\mathcal{Z}''} \frac{z''}{|\mathcal{Z}|}\right| = \frac{1}{|\mathcal{Z}|} \left|\sum_{\mathcal{Z}' removed} z - \sum_{\mathcal{Z}_{added}} z\right| \le \frac{|\mathcal{Z}_{added}|}{|\mathcal{Z}|}$$
(50)

607 Additionally, note that:

$$\left| \sum_{Z''} \frac{z''}{|\mathcal{Z}''|} - \sum_{Z'} \frac{z'}{|\mathcal{Z}'|} \right| = \frac{1}{|\mathcal{Z}''|} \left| \sum_{Z''} z'' - \frac{|\mathcal{Z}''|\sum_{Z'} z'}{|\mathcal{Z}'|} \right| = \frac{1}{|\mathcal{Z}'|} \left| \frac{|\mathcal{Z}'|\sum_{Z'} z'}{|\mathcal{Z}'|} + \frac{|\mathcal{Z}''_{removed}|\sum_{T''_{removed}} z_r}{|\mathcal{Z}''_{removed}|} - \frac{|\mathcal{Z}''|\sum_{Z'} z'}{|\mathcal{Z}'|} \right| = \frac{1}{|\mathcal{Z}|} \left| \frac{|\mathcal{Z}''_{removed}|\sum_{T''_{removed}} z_r}{|\mathcal{Z}''_{removed}|} - \frac{|\mathcal{Z}''_{removed}|\sum_{Z'} z'}{|\mathcal{Z}'|} \right| = \frac{|\mathcal{Z}''_{removed}|}{|\mathcal{Z}|} \left| \frac{\sum_{T''_{removed}} z_r}{|\mathcal{Z}''_{removed}|} - \frac{\sum_{Z'} z'}{|\mathcal{Z}'|} \right| \le \frac{|\mathcal{Z}''_{removed}|}{|\mathcal{Z}|}$$
(51)

608 Finally, by triangle inequality, we have that

$$\left|\sum_{\mathcal{Z}} \frac{z}{|\mathcal{Z}|} - \sum_{\mathcal{Z}'} \frac{z'}{|\mathcal{Z}'|}\right| \leq \left|\sum_{\mathcal{Z}} \frac{z}{|\mathcal{Z}|} - \sum_{\mathcal{Z}''} \frac{z''}{|\mathcal{Z}''|}\right| + \left|\sum_{\mathcal{Z}''} \frac{z''}{|\mathcal{Z}''|} - \sum_{\mathcal{Z}'} \frac{z'}{|\mathcal{Z}'|}\right|$$

$$\leq \frac{|\mathcal{Z}_{added}|}{|\mathcal{Z}|} + \frac{|\mathcal{Z}''_{removed}|}{|\mathcal{Z}|} \leq \frac{|\mathcal{Z}_{added}|}{|\mathcal{Z}|} + \frac{|\mathcal{Z}_{removed}|}{|\mathcal{Z}|} \leq \frac{k}{m}$$
(52)

as desired. A similar argument can be made for the case of $|\mathcal{Z}| < |\mathcal{Z}'|$.

610 We now give the classification lemmas:

611 **Lemma C.4.** *Given a finite hypothesis class* $\mathcal{G} \subseteq \mathcal{X} \to \{0,1\}$ *, and two datasets (multisets)* $A, B \subset \mathcal{X}$ *, let:*

$$g' := \arg\min_{g \in \mathcal{G}} \mathcal{L}(g) \mathcal{L}(g) := \frac{1}{|A|} \sum_{a \in A} (1 - g(a)) + \frac{1}{|B|} \sum_{b \in B} g(b).$$
(53)

- 613 Let $n = \min(|A|, |B|)$, and assume that A and B are constructed as follows:
- 614 $A' \sim \mathcal{D}_A^{|A'|}$ 615 • $B' \sim \mathcal{D}_B^{|B'|}$
- At most a total of m arbitrary (non-i.i.d.) samples are either added to or removed from A' or B',
- 617 or moved from A' to B' or vice-versa, to create A and B.

618 Then, if

$$n \ge 8m \text{ and } n \ge \frac{128}{7} \ln(2|\mathcal{G}|/\delta)$$
 (54)

619 *then with probability at least* $1 - \delta$ *,*

620 • If $D_A = D_B$, then $\mathcal{L}(g') > 1/2$

- 621 Conversely, if \mathcal{D}_A and \mathcal{D}_B have disjoint support, such that some $g^* \in \mathcal{G}$ maps all elements in the 622 support of \mathcal{D}_A to 1 and all elements in the support of \mathcal{D}_B to 0, then $\mathcal{L}(g') \leq 1/2$.
- 623 Proof. Define

$$\mathcal{L}_{clean}(g) := \frac{1}{|A'|} \sum_{a \in A'} (1 - g(a)) + \frac{1}{|B'|} \sum_{b \in B'} g(b).$$
(55)

624 Then, fix any $g \in \mathcal{G}$. From some algebra (see Proposition C.3), it can be shown that

$$\mathcal{L}_{clean}(g) - \frac{2m}{n} \le \mathcal{L}(g) \le \mathcal{L}_{clean}(g) + \frac{2m}{n}$$
(56)

Now, note that \mathcal{L}_{clean} is the sum of |A'| random variables bounded on [0, 1/|A'|], and |B'| random

variables bounded on [0, 1/|B'|], all of which are i.i.d. Then, by Hoeffding's Lemma and Equation 56, with probability $1 - \delta/|\mathcal{G}|$:

$$\mathbb{E}[\mathcal{L}_{clean}(g)] - \sqrt{\left(\frac{1}{|A'|} + \frac{1}{|B'|}\right) \ln(2|\mathcal{G}|/\delta)/2} - \frac{2m}{n} < \mathcal{L}_{clean}(g) - \frac{2m}{n} \le \mathcal{L}(g)$$

$$\le \mathcal{L}_{clean}(g) + \frac{2m}{n} < \mathbb{E}[\mathcal{L}_{clean}(g)] + \sqrt{\left(\frac{1}{|A'|} + \frac{1}{|B'|}\right) \ln(2|\mathcal{G}|/\delta)/2} + \frac{2m}{n}$$
(57)

Because $|A'|, |B'| \ge n - m$, and applying union bound over all $g \in \mathcal{G}$, we have, with probability $1 - \delta$:

$$\forall g \in \mathcal{G}, \ \mathbb{E}[\mathcal{L}_{clean}(g)] - \sqrt{\frac{\ln(2|\mathcal{G}|/\delta)}{n-m}} - \frac{2m}{n} < \mathcal{L}(g) < \mathbb{E}[\mathcal{L}_{clean}(g)] + \sqrt{\frac{\ln(2|\mathcal{G}|/\delta)}{n-m}} + \frac{2m}{n}.$$
(58)

630 Note that:

$$\forall g \in \mathcal{G}, \, \mathbb{E}[\mathcal{L}_{clean}(g)] = 1 - \mathbb{E}_{x \in \mathcal{D}_A}[g(x)] + \mathbb{E}_{x \in \mathcal{D}_B}[g(x)].$$
(59)

631 If $\mathcal{D}_A = \mathcal{D}_B$, then $\forall g \in \mathcal{G}$, $\mathbb{E}[\mathcal{L}_{clean}(g)] = 1$, so by Equation 58, we have, with probability $1 - \delta$:

$$\forall g \in \mathcal{G}, \ 1 - \sqrt{\frac{\ln(2|\mathcal{G}|/\delta)}{n-m}} - \frac{2m}{n} \le \mathcal{L}(g), \tag{60}$$

632 and in particular:

$$1 - \sqrt{\frac{\ln(2|\mathcal{G}|/\delta)}{n-m}} - \frac{2m}{n} < \mathcal{L}(g').$$
(61)

633 Conversely, if \mathcal{D}_A and \mathcal{D}_B have disjoint support, such that some $g^* \in \mathcal{G}$ maps all elements in the 634 support of \mathcal{D}_A to 1 and all elements in the support of \mathcal{D}_B to 0, then we have:

$$\mathbb{E}[\mathcal{L}_{clean}(g^*)] = 1 - \mathbb{E}_{x \in \mathcal{D}_A}[g^*(x)] + \mathbb{E}_{x \in \mathcal{D}_B}[g^*(x)] = 1 - 1 - 0 = 0.$$
(62)

635 Then, with probability at least $1 - \delta$:

$$\mathcal{L}(g') \le \mathcal{L}(g^*) < \sqrt{\frac{\ln(2|\mathcal{G}|/\delta)}{n-m}} + \frac{2m}{n}.$$
(63)

636 To complete the proof, we only need to show that

$$\sqrt{\frac{\ln(2|\mathcal{G}|/\delta)}{n-m}} + \frac{2m}{n} \le \frac{1}{2}.$$
(64)

637 With $m \le n/8$, this condition becomes:

$$\sqrt{\frac{8\ln(2|\mathcal{G}|/\delta)}{7n}} \le \frac{1}{4}.$$
(65)

638 or

$$n \ge \frac{128}{7} \ln(2|\mathcal{G}|/\delta) \tag{66}$$

639

640 **Lemma C.5.** Given a finite hypothesis class $\Phi \subset \mathcal{X} \to \mathbb{N}$, and N datasets (multisets) 641 $D_1, D_2, ..., D_N \subset \mathcal{X}$, let:

$$\phi' := \arg\min_{\phi \in \Phi} \mathcal{L}(\phi)$$

$$\mathcal{L}(\phi) := \sum_{i=1}^{N} \left[\frac{1}{|D_i|} \sum_{x \in D_i} (1 - \mathbb{1}_{(\phi(x)=i)}) \right]$$
(67)

642 Let $n = \min_i(|D_i|)$, and assume that each D_i is constructed as follows:

- 643 $\forall i, D'_i \sim \mathcal{D}_i^{|D'_i|}$
- At most a total of m arbitrary (non-i.i.d.) samples are arbitrarily moved between the datasets D'_i , to create the datasets D_i .
- 646 Additionally, assume that $\exists \phi^* \in \Phi : \forall i, x \sim \mathcal{D}_i \implies \phi^*(x) = i$. Then, if

$$n \ge \frac{8m}{\epsilon} \text{ and } n \ge \frac{64N\ln(2|\Phi|/\delta)}{7\epsilon^2}$$
 (68)

647 *then with probability at least* $1 - \delta$ *,*

$$\forall i \in [N], \quad \Pr_{x \sim \mathcal{D}_i}(\phi'(x) = i) \ge 1 - \epsilon.$$
(69)

648

649 Proof. Define

$$\mathcal{L}_{clean}(\phi) := \sum_{i=1}^{N} \left[\frac{1}{|D'_i|} \sum_{x \in D'_i} (1 - \mathbb{1}_{(\phi(x)=i)}) \right].$$
(70)

Then, fix any $\phi \in \Phi$. From Proposition C.3 (regarding each transfer of a sample as removing a sample into one multiset D'_i , and inserting a new sample into another) we see that:

$$\mathcal{L}_{clean}(\phi) - \frac{2m}{n} \le \mathcal{L}(\phi) \le \mathcal{L}_{clean}(\phi) + \frac{2m}{n}$$
(71)

Now, note that \mathcal{L}_{clean} is the sum of $|D'_1|$ random variables bounded on $[0, 1/|D'_1|]$, and $|D'_2|$ random variables bounded on $[0, 1/|D'_2|]$, et cetera, all of which are i.i.d. Then, by Hoeffding's Lemma and Equation 71, with probability $1 - \delta/|\Phi|$:

$$\mathbb{E}[\mathcal{L}_{clean}(\phi)] - \sqrt{\left(\sum_{i \in [N]} \frac{1}{|D'_i|}\right) \ln(2|\Phi|/\delta)/2} - \frac{2m}{n} < \mathcal{L}_{clean}(\phi) - \frac{2m}{n} \le \mathcal{L}(\phi)$$

$$\leq \mathcal{L}_{clean}(\phi) + \frac{2m}{n} < \mathbb{E}[\mathcal{L}_{clean}(\phi)] + \sqrt{\left(\sum_{i \in [N]} \frac{1}{|D'_i|}\right) \ln(2|\Phi|/\delta)/2} + \frac{2m}{n}$$
(72)

Because $\forall i, |D'_i| \ge n - m$, and applying union bound over all $\phi \in \Phi$, we have, with probability 1 - $\delta, \forall \phi \in \Phi$, :

$$\mathbb{E}[\mathcal{L}_{clean}(\phi)] - \sqrt{\frac{N\ln(2|\Phi|/\delta)}{2(n-m)}} - \frac{2m}{n} < \mathcal{L}(\phi) < \mathbb{E}[\mathcal{L}_{clean}(\phi)] + \sqrt{\frac{N\ln(2|\Phi|/\delta)}{2(n-m)}} + \frac{2m}{n}.$$
 (73)

657 Then we have:

$$\mathbb{E}[\mathcal{L}_{clean}(\phi')] < \mathcal{L}(\phi') + \sqrt{\frac{N\ln(2|\Phi|/\delta)}{2(n-m)}} + \frac{2m}{n} \leq \mathcal{L}(\phi^*) + \sqrt{\frac{N\ln(2|\Phi|/\delta)}{2(n-m)}} + \frac{2m}{n} \leq \mathbb{E}[\mathcal{L}_{clean}(\phi^*)] + \sqrt{\frac{2N\ln(2|\Phi|/\delta)}{(n-m)}} + \frac{4m}{n} = (74)$$

$$\sqrt{\frac{2N\ln(2|\Phi|/\delta)}{(n-m)}} + \frac{4m}{n}.$$

- 658 Where we use the fact that, by the definition of ϕ' as a minimizer, $\mathcal{L}(\phi') \leq \mathcal{L}(\phi^*)$, as well as the 659 fact that, by definition, $\mathcal{L}_{clean}(\phi^*) = 0$.
- Also, note that by the definition of \mathcal{L}_{clean} , we have that, for any ϕ ,

$$\mathbb{E}[\mathcal{L}_{clean}(\phi)] = \sum_{i \in [N]} 1 - \Pr_{x \sim \mathcal{D}_i}(\phi(x) = i)$$
(75)

661 Then, for any particular $i \in [N]$, we have that $1 - \Pr_{x \sim D_i}(\phi(x) = i) \leq \mathbb{E}[\mathcal{L}_{clean}(\phi)]$. Then, by 662 Equation 74, we have, $\forall i \in [N]$:

$$1 - \Pr_{x \sim \mathcal{D}_i}(\phi'(x) = i) < \sqrt{\frac{2N\ln(2|\Phi|/\delta)}{(n-m)}} + \frac{4m}{n}.$$
(76)

663 By algebra, our desired result (Equation 69) holds as long as:

$$\sqrt{\frac{2N\ln(2|\Phi|/\delta)}{(n-m)}} + \frac{4m}{n} \le \epsilon \tag{77}$$

664 Which follows from the given conditions on n.

665 C.3 Main Proof of Theorem 3.1

Here, we present the proof of Theorem 3.1. We first split out correctness proof of the main recursive step of the algorithm as a lemma:

668 **Lemma C.6.** In Algorithm 2, suppose that the ground-truth data coverage assumptions given in 669 Equations 4,5, and 6 all hold. Additionally, assume that the relative coverage lower-bound η can be 670 written in the form

$$\eta = \frac{e^{-n_{\Xi}\alpha/8}}{1 + e^{-n_{\Xi}\alpha/8}}$$
(78)

for some non-negative integer n_{Ξ} . Further, assume that for each $s^* \in S_h^*$, there exists some $s \in S_h^*$

- 672 S_h that represents approximately the same set of observations. In particular, each index in $[|\tau_A|]$ 673 appears in at most one set $D_{A,h}(s)$ for some s (and likewise for $[|\tau_B|]$ and $D_{B,h}(s)$), and there
- exists some bijective mapping $\sigma_h : S_h \to S_h^*$, such that for most indices j in $[|\tau_A|]$

$$j \in D_{A,h}(\sigma_h^{-1}(\phi_h^*((\tau_A)_h[j])))$$
(79)

675 and for most indices j in $[|\tau_B|]$

$$j \in D_{B,h}(\sigma_h^{-1}(\phi_h^*((\tau_B)_h[j]))),$$
(80)

- 676 with at most a combined $\beta(|\tau_A| + |\tau_B|)$ indices in either dataset for which this does not hold.
- 677 For any ϵ such that:

$$\epsilon < \frac{\nu}{8} - \beta,\tag{81}$$

678 and

$$\epsilon + \beta \le q_{thresh.} < \frac{\nu(1-\epsilon)}{2} - \beta, \tag{82}$$

679 where $q_{thresh.}$ is the threshold defined on Line 11 of the algorithm, assume that 680 $\forall s_{h}^{*}, s_{h+1}^{*}$, such that s_{h}^{*} can transition to s_{h+1}^{*} ,

$$|\mathcal{D}^*(s_h^*, s_{h+1}^*)| \ge \frac{12800(\ln(4|\Phi|^2/\delta) + N_s^2 \ln(n_{\Xi} + 1))\ln^2(1/\eta)}{\nu\epsilon^2 \eta^2 \alpha^4}.$$
(83)

Then, with high probability, for each $s^* \in S_{h+1}^*$, there exists some $s \in S_{h+1}$ that represents approximately the same set of observations. In particular, each index in $[|\tau_A|]$ appears in at most one set $D_A(s)$ for some s (and likewise for $[|\tau_B|]$ and $D_B(s)$), and there exists some bijective mapping $\sigma_{h+1}: S_{h+1} \to S_{h+1}^*$, such that for most indices j in $[|\tau_A|]$

$$j \in D_{A,h+1}(\sigma_{h+1}^{-1}(\phi_{h+1}^*(\tau_A)_{h+1}[j])))$$
(84)

685 and for most indices j in $[|\tau_A|]$

$$j \in D_{B,h+1}(\sigma_{h+1}^{-1}(\phi_{h+1}^*((\tau_B)_{h+1}[j]))),$$
(85)

686 with at most a combined $(\beta + \epsilon)(|\tau_A| + |\tau_B|)$ indices in either dataset for which this does not hold.

687 *Proof.* We first show that the datasets of observation pairs \mathcal{D}_{new_pairs} defined in Line 21 of the algo-688 rithm each correspond uniquely to a pair of ground truth latent states in $\mathcal{S}_h^* \times \mathcal{S}_{h+1}^*$, such that no pair 689 of observations is included in more than one such \mathcal{D}_{new_pairs} sets, and, with high probability, each 690 pair of observations x, x' is included in the correct \mathcal{D}_{new_pairs} corresponding to $(\phi_h^*(x), \phi_{h+1}^*(x'))$, 691 with up to at most $(\beta + \epsilon)(|\tau_A| + |\tau_B|)$ exceptions.

692 For any $s_{pred}^* \in \mathcal{S}_h^*$, consider any two distinct $s^*, s'^* \in \mathcal{S}_{h+1}^*$, such that 693 $|\mathcal{D}^*(s_{pred}^*, s^*)|, |\mathcal{D}^*(s_{pred}^*, s'^*)| > 0.$

694 Recall the assumption that, without loss of generality,

$$e^{\alpha} \cdot \frac{\pi_B^{emp.}(s'^*|s_{pred}^*)}{\pi_B^{emp.}(s^*|s_{pred}^*)} \le \frac{\pi_A^{emp.}(s'^*|s_{pred}^*)}{\pi_A^{emp.}(s^*|s_{pred}^*)},$$
(86)

695 Multiplying both sides by $\pi_A^{emp.}(s^*|s^*_{pred})/\pi_B^{emp.}(s'^*|s^*_{pred})$ yields

$$e^{\alpha} \cdot \frac{\pi_A^{emp.}(s^*|s_{pred}^*)}{\pi_B^{emp.}(s^*|s_{pred}^*)} \le \frac{\pi_A^{emp.}(s'^*|s_{pred}^*)}{\pi_B^{emp.}(s'^*|s_{pred}^*)},$$
(87)

696 From the definition of π^{emp} , this is:

$$e^{\alpha} \cdot \frac{|\mathcal{D}_{A}^{*}(s_{pred}^{*}, s^{*})| / |\mathcal{D}_{A}^{*}(s_{pred}^{*})|}{|\mathcal{D}_{B}^{*}(s_{pred}^{*}, s^{*})| / |\mathcal{D}_{B}^{*}(s_{pred}^{*})|} \leq \frac{|\mathcal{D}_{A}^{*}(s_{pred}^{*}, s^{\prime *})| / |\mathcal{D}_{A}^{*}(s_{pred}^{*})|}{|\mathcal{D}_{B}^{*}(s_{pred}^{*}, s^{\prime *})| / |\mathcal{D}_{B}^{*}(s_{pred}^{*})|}.$$
(88)

697 Multiplying both sides by $|\mathcal{D}_A^*(s_{pred}^*)|/|\mathcal{D}_B^*(s_{pred}^*)|$ and taking the logarithms yields:

$$\alpha + \ln\left(\frac{|\mathcal{D}_{A}^{*}(s_{pred}^{*}, s^{*})|}{|\mathcal{D}_{B}^{*}(s_{pred}^{*}, s^{*})|}\right) \le \ln\left(\frac{|\mathcal{D}_{A}^{*}(s_{pred}^{*}, s'^{*})|}{|\mathcal{D}_{B}^{*}(s_{pred}^{*}, s'^{*})|}\right).$$
(89)

698 By the definitions of $\underline{c}_{s_h^*, s_{h+1}^*}$ and $\overline{c}_{s_h^*, s_{h+1}^*}$ in Corollary C.2, and the fact that $\xi = \alpha/4$, we see that 699 there must be at least two values in Ξ between $\overline{c}_{s_{pred}^*, s^*}$ and $\underline{c}_{s_{pred}^*, s'^*}$; that is to say:

$$\underline{c}_{s_{pred}^*,s^*} \le \bar{c}_{s_{pred}^*,s^*} < \bar{c}_{s_{pred}^*,s^*} + \xi < \underline{c}_{s_{pred}^*,s'^*} - \xi < \underline{c}_{s_{pred}^*,s'^*} \le \bar{c}_{s_{pred}^*,s'^*}$$
(90)

- Therefore, by Corollary C.2 we have, with probability at least $1 \delta/2$, for any s^* such that $|\mathcal{D}^*(s^*_{pred}, s^*)| > 0$:
- At least $(1 \epsilon) |\mathcal{D}^*(s_{pred}^*, s^*)|$ of the samples in $\mathcal{D}^*(s_{pred}^*, s^*)$ will be mapped by f to $\underline{c}_{s_{pred}^*, s^*}$ or $\overline{c}_{s_{pred}^*, s^*}$
- for some choice of $\hat{c}_{s_{pred}^*,s^*} \in \{\underline{c}_{s_{pred}^*,s^*}, \overline{c}_{s_{pred}^*,s^*}\}$, at least $(1-\epsilon)/2 \cdot |\mathcal{D}^*(s_{pred}^*,s^*)|$ of the samples in $\mathcal{D}^*(s_{pred}^*,s^*)$ will be mapped to $\hat{c}_{s_{pred}^*,s^*}$.
- $^{706} By definition, {<u>c</u>_{s[*]_{pred},s[*]}, c̄_{s[*]_{pred},s[*]}} ⊂ {ĉ_{s[*]_{pred},s[*]} − ξ, ĉ_{s[*]_{pred},s[*]}, ĉ_{s[*]_{pred},s[*]} + ξ}.$
- Furthermore, for no two states s^* , s'^* , with $\hat{c}_{s_{pred}^*,s^*} \in \{\underline{c}_{s_{pred}^*,s^*}, \overline{c}_{s_{pred}^*,s^*}\}$ and $\hat{c}_{s_{pred}^*,s'^*} \in \{\underline{c}_{s_{pred}^*,s'^*}, \overline{c}_{s_{pred}^*,s^*}, \overline{c}_{s_{pred}^*,s^*}\}$ chosen arbitrarily, will the sets $\{\hat{c}_{s_{pred}^*,s^*} \xi, \hat{c}_{s_{pred}^*,s^*}, \hat{c}_{s_{pred}^*,s^*} + \xi\}$ and $\{\hat{c}_{s_{pred}^*,s'^*} \xi, \hat{c}_{s_{pred}^*,s^*}, \hat{c}_{s_{pred}^*,s^*} + \xi\}$ and $\{\hat{c}_{s_{pred}^*,s'^*} \xi, \hat{c}_{s_{pred}^*,s^*}, \hat{c}_{s_{pred}^*,s^*} + \xi\}$ overlap (By Equation 90).
- 710 Recall that by assumption, $|\mathcal{D}^*(s_{pred}^*, s^*)| \ge \nu(|\tau_A| + |\tau_B|)$. Therefore, at least $(\nu(1 \epsilon)/2)(|\tau_A| + |\tau_B|)$ of the samples in $\mathcal{D}^*(s_{pred}^*, s^*)$ will be mapped to $\hat{c}_{s_{pred}^*, s^*}$.
- The total number of samples in $\mathcal{D}^*(s_{pred}^*, s'^*)$, over all choices of s'^* , which are not mapped by f to a value in the respective set $\{\hat{c}_{s_{pred}^*, s'^*} - \xi, \hat{c}_{s_{pred}^*, s'^*}, \hat{c}_{s_{pred}^*, s'^*} + \xi\}$, is at most $\epsilon |\mathcal{D}^*(s_{pred}^*)|$.

714 •
$$\epsilon |\mathcal{D}^*(s_{pred}^*)| \le \epsilon (|\tau_A| + |\tau_B|).$$

Therefore, as long as $(\nu(1-\epsilon)/2) > \epsilon$, then among the pairs in $\mathcal{D}^*(s_{pred}^*, \cdot) := \bigcup_{s'^*} \mathcal{D}^*(s_{pred}^*, s'^*)$, 715 if there is any $z \in \Xi$ such that $> \epsilon(|\tau_A| + |\tau_B|)$ of the pairs are mapped by f to z, then we know that 716 the set of elements in $\mathcal{D}^*(s_{nred}^*, \cdot)$ which are mapped to $\{z - 1, z, z + 1\}$ contains at least $(1 - \epsilon)$ 717 of the elements of the set $\mathcal{D}^*(s_{pred}^*, s^*)$ for some s^* ; furthermore, such a z exists for each possible value of s^* where $|\mathcal{D}^*(s_{pred}^*, s^*)| > 0$, and, for distinct s^* and s'^* , these values ($\{z - \xi, z, z + \xi\}$ 718 719 and $\{z' - \xi, z', z' + \xi\}$) are non-overlapping. Consequently, by identifying subsets of $\mathcal{D}^*(s^*_{pred}, \cdot)$ 720 721 of size greater than $\epsilon(|\tau_A| + |\tau_B|)$ that f maps to the same value, and expanding these subsets to the elements in $\mathcal{D}^*(s_{pred}^*, \cdot)$ mapped to adjacent values in Ξ , we can partition $\mathcal{D}^*(s_{pred}^*, \cdot)$ into subsets 722 corresponding to each $\mathcal{D}^*(s_{pred}^*, s^*)$, with at most $\epsilon |\mathcal{D}^*(s_{pred}^*, \cdot)|$ errors. 723

Note however that we do not have access to $\mathcal{D}^*(s_{pred}^*, \cdot)$, only to pairs (s_{pred}) (where $\sigma_h(s_{pred}) = s_{pred}^*$). However, by assumption, $\mathcal{D}^*(s_{pred}^*, \cdot)$ and pairs (s_{pred}) differ (in terms of symmetric difference) by at most $\beta(|\tau_A| + |\tau_B|)$. Therefore, we claim that, if

$$\epsilon + \beta \le q_{thresh.} < \frac{\nu(1-\epsilon)}{2} - \beta$$
(91)

then, we can identify values of $\hat{c}_{s_{pred}^*,s'^*}$ (for some s'^*) as those values $j\xi - \frac{n \pm \xi}{2}$ for which (as shown in Line 20 of Algorithm 2):

$$\operatorname{pred_succ}[j] > q_{thresh.}(|\tau_A| + |\tau_B|), \tag{92}$$

729 and, conversely, if

$$\operatorname{pred_succ}[j] \le q_{thresh.}(|\tau_A| + |\tau_B|), \tag{93}$$

- 730 then $j\xi \frac{n_{\Xi}\xi}{2}$ does not correspond to some $\bar{c}_{s_{pred}^*,s'^*}$ or $\underline{c}_{s_{pred}^*,s'^*}$.
- 731 To validate this claim, note that if Equation 91 holds, then:

- # of samples (x, x') in pairs (s_{pred}) such that f(x, x') = z, if $\exists s^* : z \in \{\bar{c}_{s_{pred}^*, s^*}, \underline{c}_{s_{pred}^*, s^*}\} \leq$ # of samples (x, x') in $\mathcal{D}^*(s_{pred}^*)$ such that f(x, x') = z, if $\exists s^* : z \in \{\bar{c}_{s_{pred}^*, s^*}, \underline{c}_{s_{pred}^*, s^*}\} + |\operatorname{pairs}(s_{pred}) \setminus \mathcal{D}^*(s_{pred}^*)| \leq$ $\epsilon(|\tau_A| + |\tau_B|) + |\operatorname{pairs}(s_{pred}) \setminus \mathcal{D}^*(s_{pred}^*)| \leq$ $\epsilon(|\tau_A| + |\tau_B|) + |\operatorname{pairs}(s_{pred}) \setminus \mathcal{D}^*(s_{pred}^*)| \leq$ $(\operatorname{Note this line:}) \quad \epsilon(|\tau_A| + |\tau_B|) + \beta(|\tau_A| + |\tau_B|) \leq$ $q_{thresh}.(|\tau_A| + |\tau_B|) <$ $(\nu(1 - \epsilon)/2)(|\tau_A| + |\tau_B|) - \beta(|\tau_A| + |\tau_B|) \leq$ $(\nu(1 - \epsilon)/2)(|\tau_A| + |\tau_B|) - |\mathcal{D}^*(s_{pred}^*) \setminus \operatorname{pairs}(s_{pred})| \leq$ # of samples (x, x') in $\mathcal{D}^*(s_{pred}^*)$ such that f(x, x') = z, if $\exists s^* : z \in \{\bar{c}_{s_{pred}^*, s^*}, \underline{c}_{s_{pred}^*, s^*}\} - |\mathcal{D}^*(s_{pred}^*) \setminus \operatorname{pairs}(s_{pred})| \leq$ # of samples (x, x') in pairs (s_{pred}) such that f(x, x') = z, if $\exists s^* : z \in \{\bar{c}_{s_{pred}^*, s^*}, \underline{c}_{s_{pred}^*, s^*}\}$
- Therefore, for any s_{pred}^* , we can define:

$$\mathcal{D}_{new_pairs}^*(s_{pred}^*, j, j') := \{(x, x') | (x, x') \in \bigcup_{k=j-1}^{j'} \operatorname{pred_succ}^*[k]\}$$
(94)

733 where

pred_succ*[k] := {
$$(x_h, x_{h+1}) \in \mathcal{D}^*(s_{pred}^*, \cdot) | f(x_h, x_{h+1}) = k\xi - \frac{n_{\Xi}\xi}{2}$$
}. (95)

If j and j' are chosen as in Line 19 and 20 of Algorithm 2, then for any pair (s_{pred}^*, s^*) there 734 is a unique set $\mathcal{D}^*_{new_pairs}(s^*_{pred}, j, j')$ containing at least a $(1 - \epsilon)$ fraction of the samples in 735 736 $\mathcal{D}(s_{pred}^*, s^*)$. Furthermore, note that by assumption, summing over all pairs (s_{pred}^*, s^*) , the sets 737 $\mathcal{D}_{new\ pairs}^*(s_{pred}^*, j, j')$ and $\mathcal{D}_{new\ pairs}$ can differ by at most $\beta(|\tau_A| + |\tau_B|)$ members in total (because all datasets $\mathcal{D}^*(s_{pred}^*, \cdot)$ and pairs (s_{pred}) differ by at most this many members in total). There-738 fore, we have shown that, with high probability, each pair of observations x, x' is included in the 739 740 correct \mathcal{D}_{new_pairs} corresponding to $(\phi_h^*(x), \phi_{h+1}^*(x'))$, with up to at most $(\beta + \epsilon)(|\tau_A| + |\tau_B|)$ 741 exceptions. (Furthermore, different sets \mathcal{D}_{new_pairs} are non-overlapping by construction.) Then we only need to show that, with high probability, Line 26 of Algorithm 2 will only merge two sets \mathcal{D}_{new} 742 743 if these sets correspond to the same latent state. By Lemma C.4, taking a union bound over all pairs 744 of latent-state sequential latent-state pairs, we have, with probability at least $1 - \delta/2$ that, if

$$\nu \ge 8(\beta + \epsilon) \text{ and } |\mathcal{D}^*(s_h^*, s_{h+1}^*)| \ge \frac{128}{7} \ln(4 \cdot (\text{\# of latent state pairs}) \cdot |\mathcal{G}_{h+1}|/\delta)$$
(96)

then the classifiers trained in Line 25 will have loss greater than 1/2 if and only if the two datasets being compared correspond to the same latent state. Also note that (# of latent state pairs) $\leq 1/\nu$; then, under the assumption that $|\mathcal{G}_{h+1}| \leq |\Phi| (\leq |\Phi|^2)$, we have

$$\frac{128}{7}\ln(4 \cdot (\text{\# of latent state pairs}) \cdot |\mathcal{G}_{h+1}|/\delta) \le \frac{128}{7}(\ln(4|\Phi^2|/\delta) + \ln(1/\nu)).$$
(97)

Note that by Equation 40 (and $\xi = \alpha/4$), we have

$$\frac{16\ln(1/\eta)}{\eta\alpha^2} \ge 1,\tag{98}$$

749 so we can write:

$$\frac{128}{7}\ln(4\cdot(\#\text{ of latent state pairs})\cdot|\mathcal{G}_{h+1}|/\delta) \le \frac{128\cdot 256}{7}\frac{(\ln(4|\Phi^2|/\delta) + \ln(1/\nu))\ln^2(1/\eta)}{\eta^2\alpha^4}.$$
 (99)

Also noting that $\epsilon \leq 1$, and $\ln(4|\Phi|^2/\delta) \geq \ln(4) \geq 1$, and $1/\nu \geq \{\ln(1/\nu), 1\}$, we have:

$$\frac{128}{7}\ln(4\cdot(\text{\# of latent state pairs})\cdot|\mathcal{G}_{h+1}|/\delta) \le \frac{128\cdot256}{7}\frac{(\ln(4|\Phi^2|/\delta) + \ln(4|\Phi^2|/\delta))\ln^2(1/\eta)}{\nu\epsilon^2\eta^2\alpha^4}.$$
(100)

751 Then, because $N_s^2 \ln(n_{\Xi} + 1) > 0$, and $128 \cdot 256 \cdot 2/7 \le 12800$, we have

$$\frac{128}{7}\ln(4 \cdot (\# \text{ of latent state pairs}) \cdot |\mathcal{G}_{h+1}|/\delta) \le \frac{12800(\ln(4|\Phi^2|/\delta) + N_s^2\ln(n_{\Xi}+1))\ln^2(1/\eta)}{\nu\epsilon^2\eta^2\alpha^4}.$$
(101)

Therefore, the number of observations of each latent state pair $|\mathcal{D}^*(s_h^*, s_{h+1}^*)|$ assumed in Equation 83 is sufficient to ensure that all datasets \mathcal{D}_{new} will be merged correctly. Then, by union bound, with probability at least $1 - \delta$, the samples corresponding to the indices in $D_{A,h+1}(s)$ and $D_{B,h+1}(s)$ will correspond to the observations of a unique latent state s^* , with up to at most $(\beta + \epsilon)(|\tau_A| + |\tau_B|)$ exceptions.

757 The following lemma is essentially Theorem 3.1, with a minor additional assumption:

- **Lemma C.7.** Assume that Algorithm 2 is given datasets τ_A and τ_B such that the assumptions given
- in Equations 1, 4,5, and 6 all hold. Additionally, assume that the relative coverage lower-bound η

760 can be written in the form

$$\eta = \frac{e^{-n_{\Xi}\alpha/8}}{1 + e^{-n_{\Xi}\alpha/8}}$$
(102)

761 for some non-negative integer n_{Ξ} . Then, for any given $\delta, \epsilon_0 \geq 0$, if 762 $\forall s_h^*, s_{h+1}^*$, such that s_h^* can transition to s_{h+1}^* ,

$$|\mathcal{D}^*(s_h^*, s_{h+1}^*)| \ge \frac{819200H^2(\ln(8H|\Phi|^2/\delta) + N_s^2\ln(n_{\Xi}+1))\ln^2(1/\eta)}{\nu\eta^2\alpha^4} \cdot \max\left(\frac{1}{\nu^2}, \frac{1}{\epsilon_0^2\nu'^2}\right),\tag{103}$$

763 then, with probability at least $1 - \delta$, the encoders ϕ'_h returned by the algorithm will each have 764 accuracy on at least $1 - \epsilon_0$, in the sense that, under some bijective mapping $\sigma_h : S_h \to S_h^*$,

$$\forall s^* \in \mathcal{S}_h^*, \quad \Pr_{x \sim \mathcal{Q}(s^*, P_h^e)}(\phi_h'(x) = \sigma_h^{-1}(\phi_h^*(x))) \ge 1 - \epsilon_0. \tag{104}$$

765

Proof. Note that the conclusion applies at timestep h = 1 vacuously: there is only one latent state, and ϕ'_1 returns a constant value. Further, $D_{A,1}(s_1)$ and $D_{B,1}(s_1)$ contain exactly the sets of trajec-

768 tories in τ_A and τ_B which visit s_1^* at step 1.

769 For subsequent steps, we apply Lemma C.6 recursively, with:

770 •
$$\epsilon = \min(\frac{\nu}{8H}, \frac{\nu'\epsilon_0}{8H})$$

771 •
$$\beta = (h-1)\epsilon$$

- 772 $\delta_{\text{Lemma C.6}} := \delta/(2H).$
- 773 Note that the assumptions in Equations 81 and 82 are met, because:

$$\epsilon = \epsilon + \beta - \beta = h\epsilon - \beta < H\epsilon - \beta = H\min(\frac{\nu}{8H}, \frac{\nu'\epsilon_0}{8H}) - \beta < \frac{H\nu}{8H} - \beta < \frac{\nu}{8} - \beta$$
(105)

774 and

$$\epsilon + \beta = h \min\left(\frac{\nu}{8H}, \frac{\nu'\epsilon_0}{8H}\right) \le \frac{h\nu}{8H} = q_{thresh.}$$
(106)

775 and

$$q_{thresh.} - q_{thresh.} - q_{thresh.} - q_{thresh.} + \beta + \frac{\epsilon\nu}{2} - \beta - \frac{\epsilon\nu}{2} = \frac{h\nu}{8H} + (h - 1 + \nu/2) \min\left(\frac{\nu}{8H}, \frac{\nu'\epsilon_0}{8H}\right) - \beta - \frac{\epsilon\nu}{2} \leq \frac{h\nu}{8H} + \frac{(h - 1 + \nu/2)\nu}{8H} - \beta - \frac{\epsilon\nu}{2} < \frac{2H\nu}{8H} - \frac{\epsilon\nu}{2} - \beta < \frac{\nu(1 - \epsilon)}{2} - \beta.$$

$$(107)$$

Also, note that the assumption of Equation 83 is met (by comparison to Equation 103, with $\epsilon = \min(\frac{\nu}{8H}, \frac{\nu'\epsilon_0}{8H})$ and $\delta_{\text{Lemma C.6}} := \delta/(2H)$.) Finally, the inductive hypothesis, that $D_{A,h}(s)$ and $D_{B,h}(s)$ correspond to observations of some state s^* , with at most $\beta(|\tau_A| + |\tau_B|)$ exceptions, can be shown to hold. In particular, at iteration $h \ge 2$, we have that the input dataset has at most $\beta_h(|\tau_A| + |\tau_B|) = (\beta_{h-1} + \epsilon)(|\tau_A| + |\tau_B|)$ errors: we can confirm that $\beta_h = (h-1)\epsilon = (h-2)\epsilon + \epsilon = \beta_{h-1} + \epsilon$ for all $h \ge 2$, with $\beta_1 = 0$ (because there are no errors in $D_{A,1}(s_1)$ and $D_{B,1}(s_1)$).

Therefore, by induction and union bound, we can conclude that, with probability at least $1 - \delta/2$, for each $h \in [H]$ and each $s^* \in S_h^*$, there exists some $s \in S_h$ that represents approximately the same set of observations. In particular, each index in $[|\tau_A|]$ appears in at most one set $D_A(s)$ for some s(and likewise for $[|\tau_B|]$ and $D_B(s)$), and there exists some bijective mapping $\sigma_h : S_h \to S_h^*$, such that for most indices j in $[|\tau_A|]$

$$j \in D_{A,h}(\sigma_h^{-1}((\phi^*(\tau_A)_h[j])))$$
(108)

787 and for most indices j in $[|\tau_A|]$

$$j \in D_{B,h}(\sigma_h^{-1}(\phi^*((\tau_B)_h[j]))), \tag{109}$$

with at most a combined $(H-1)\min(\frac{\nu}{8H}, \frac{\nu'\epsilon_0}{8H})(|\tau_A| + |\tau_B|)$ indices in either dataset for which this does not hold. Note in particular that fewer than $\frac{\nu'\epsilon_0}{8}(|\tau_A| + |\tau_B|)$ indices will be mis-categorized at any timestep. Then, by application of Lemma C.5 with $n \ge \nu'(|\tau_A| + |\tau_B|)$, $m = \frac{\nu'\epsilon_0}{8}(|\tau_A| + |\tau_B|)$,

791 $\delta_{\text{Lemma C.5}} = \delta/(2H)$ and $\epsilon = \epsilon_0$, we have that as long as:

$$\nu'(|\tau_A| + |\tau_B|) \ge \frac{64 \cdot \max_i |\mathcal{S}_i| \cdot \ln(4H|\Phi|/\delta)}{7\epsilon_0^2},$$
(110)

then, by union bound, with probability at least $1 - \delta$, the encoders learned on line 46 of Algorithm 2 will each have accuracy at least $1 - \epsilon_0$ as in Equation 104, as desired. All that remains to be shown is that Equation 110 holds. Note that this equation can be re-written as:

$$|\tau_A| + |\tau_B| \ge \frac{64 \cdot \max_h |\mathcal{S}_h| \cdot \ln(4H|\Phi|/\delta)}{7\nu'\epsilon_0^2}.$$
(111)

795 Then, we have:

$$\begin{aligned} |\tau_{A}| + |\tau_{B}| &\geq \\ |\mathcal{D}^{*}(s_{h}^{*}, s_{h+1}^{*})| &\geq \\ \frac{819200H^{2}(\ln(8H|\Phi|^{2}/\delta) + N_{s}^{2}\ln(n_{\Xi} + 1))\ln^{2}(1/\eta)}{\nu\eta^{2}\alpha^{4}} \cdot \max\left(\frac{1}{\nu^{2}}, \frac{1}{\epsilon_{0}^{2}\nu'^{2}}\right) &\geq \\ \frac{819200H^{2}(\ln(8H|\Phi|^{2}/\delta) + N_{s}^{2}\ln(n_{\Xi} + 1))\ln^{2}(1/\eta)}{\epsilon_{0}^{2}\nu'^{2}\nu\eta^{2}\alpha^{4}} &\geq \\ (\text{by Equation 98)} \quad \frac{3200H^{2}(\ln(8H|\Phi|^{2}/\delta) + N_{s}^{2}\ln(n_{\Xi} + 1))}{\epsilon_{0}^{2}\nu'^{2}\nu} &\geq \\ (\text{log. of integer } \geq 0) \quad \frac{3200H^{2}\ln(8H|\Phi|^{2}/\delta)}{\epsilon_{0}^{2}\nu'^{2}\nu} &\geq \\ (\text{By definition, } (1/\nu' \geq \max_{h} |\mathcal{S}_{h}|) \quad \frac{3200H^{2}\max_{h} |\mathcal{S}_{h}|\ln(8H|\Phi|^{2}/\delta)}{\epsilon_{0}^{2}\nu'\nu} &\geq \\ \frac{64 \cdot \max_{h} |\mathcal{S}_{h}| \cdot \ln(4H|\Phi|/\delta)}{7\nu'\epsilon_{0}^{2}} &= \\ \end{aligned}$$

completing the proof. 796

797 Finally, we prove Theorem 3.1:

798 **Theorem 3.1.** Assume that CRAFT (Algorithm 2 in the Appendix) is given datasets τ_A and τ_B such 799 that the assumptions given in Equations 1, 4,5, and 6 all hold. Then there exists an

$$f\left(H, |\Phi|, N_s, \frac{1}{\delta}, \frac{1}{\epsilon_0}, \frac{1}{\nu}, \frac{1}{\nu'}, \frac{1}{\eta}, \frac{1}{\alpha}\right) \in \mathcal{O}^*\left(\frac{H^2(\ln(|\Phi|/\delta) + N_s^2)}{\nu\eta^2\alpha^4} \cdot \max\left(\frac{1}{\nu^2}, \frac{1}{\epsilon_0^2\nu'^2}\right)\right), \quad (16)$$

where $\mathcal{O}^*(f(x)) := \mathcal{O}(f(x)\log^k(f(x)))$, such that for any given $\delta, \epsilon_0 \ge 0$, if $\forall s_h^*, s_{h+1}^*$ such that s_h^* 800

801

can transition to s_{h+1}^* , $|\mathcal{D}^*(s_h^*, s_{h+1}^*)| \ge f\left(H, |\Phi|, N_s, \frac{1}{\delta}, \frac{1}{\epsilon_0}, \frac{1}{\nu}, \frac{1}{\nu'}, \frac{1}{\eta}, \frac{1}{\alpha}\right)$, then, with probability at least $1 - \delta$, the encoders ϕ'_h returned by the algorithm will each have accuracy on at least $1 - \epsilon_0$, 802 803 in the sense that, under some bijective mapping $\sigma_h : S_h \to S_h^*$,

$$\forall s^* \in \mathcal{S}_h^*, \quad \Pr_{x \sim \mathcal{Q}(s^*, P_h^e)}(\phi_h'(x) = \sigma_h^{-1}(\phi_h^*(x))) \ge 1 - \epsilon_0.$$
(17)

Proof. This final theorem follows close-to-directly from Lemma C.7, with the caveat that we no 804 longer assume that $\eta = (e^{-n \equiv \alpha/8})/(1 + e^{-n \equiv \alpha/8})$ for some non-negative integer n_{Ξ} . To do this, it 805 806 is important to note that the provided η is a *lower bound*: if we replace η in the algorithm with any arbitrary $\eta' \leq \eta$, then Lemma C.7 will still apply, with a sample-complexity in terms of η' rather 807 than η . Similarly, α is a lower-bound, and so Lemma C.7 will apply for any smaller α' . Our task is 808 then to replace η and α with some η' and α' , such that the *asymptotic* sample complexity as given 809 810 by Equation 16 still applies. For simplicity, we can write the sample-complexity given in Equation 811 103 as:

$$|\mathcal{D}^*(s_h^*, s_{h+1}^*)| \ge \frac{(C_1 + C_2 \ln(n_{\Xi} + 1)) \ln^2(1/\eta')}{\eta'^2 \alpha'^4}, \tag{113}$$

where C_1 and C_2 are independent of α , η , and n_{Ξ} . Now, we must choose η' such that: 812

$$\eta \ge \eta' = \frac{e^{-n_{\Xi}\alpha/8}}{1 + e^{-n_{\Xi}\alpha/8}} \left(= \frac{1}{1 + e^{n_{\Xi}\alpha/8}} \right).$$
(114)

An obvious choice is to take (recalling that by definition, $\eta \le 1/2$, so $\ln(\eta^{-1} - 1) > 0$): 813

$$n_{\Xi} := \left\lceil 8\ln(\eta^{-1} - 1)/\alpha \right\rceil \tag{115}$$

814 so that:

$$\eta' = \frac{1}{1 + e^{\lceil 8 \ln(\eta^{-1} - 1)/\alpha \rceil \alpha/8}}$$
(116)

815 Then we have:

$$\eta \ge \eta' \ge \frac{1}{1 + e^{(8\ln(\eta^{-1} - 1)/\alpha + 1)\alpha/8}} \ge \eta \cdot e^{-\alpha/8}$$
(117)

so that we have sufficient samples for Lemma C.7 if:

$$|\mathcal{D}^*(s_h^*, s_{h+1}^*)| \ge \frac{(C_1 + C_2 \ln(8\ln(\eta^{-1} - 1)/\alpha' + 2))(\ln(1/\eta) + \alpha'/8)^2 e^{\alpha'/4}}{\eta^2 \alpha'^4}, \tag{118}$$

Strictly speaking, Equation 118 with $\alpha' = \alpha$ satisfies the "big-O" asymptotic complexity given in Equation 16 in terms of $1/\alpha$ and $1/\eta$ as these quantities approach infinity. However, if we just take $\alpha' = \alpha$, notice that Equation 118 seems to require an exponentially large number of samples for large α . Recall though that α is a *lower bound*, so we can simply select an arbitrarily lower α' in the case of large α . In particular, if we take $\alpha' = \min(1, \alpha)$, then $\alpha' \leq \alpha$ as needed, and (ignoring lower-order polynomial terms and all logarithmic factors), the dependence of our sample complexity on α becomes:

$$\min(e^{\alpha/4}/\alpha^4, e^{1/4}/1^4) = \min(e^{\alpha/4}/\alpha^4, e^{1/4}) \le C \cdot 1/\alpha^4$$
(119)

so that the sample complexity is bounded even for large α .

These modifications to η and α are performed on Lines 1-3 of Algorithm 2, so the overall asymptotic sample complexity given in Equation 16 holds for the algorithm overall, with the input α and η .

827 **D** Experiment Details

For the hyperparameters η , ν and α of CRAFT, we use the "population" values based on the groundtruth dynamics and policy. In other words, we set:

$$e^{\alpha} = \min_{s_{h}^{*}, s_{h+1}^{*}, s_{h+1}^{*}} \max\left[\left(\frac{\Pr_{\pi_{A}}(s_{h+1}^{'*}|s_{h}^{*}) / \Pr_{\pi_{A}}(s_{h+1}^{*}|s_{h}^{*})}{\Pr_{\pi_{B}}(s_{h+1}^{'*}|s_{h}^{*}) / \Pr_{\pi_{B}}(s_{h+1}^{*}|s_{h}^{*})} \right), \\ \left(\frac{\Pr_{\pi_{B}}(s_{h+1}^{'*}|s_{h}^{*}) / \Pr_{\pi_{B}}(s_{h+1}^{*}|s_{h}^{*})}{\Pr_{\pi_{A}}(s_{h+1}^{*}|s_{h}^{*}) / \Pr_{\pi_{A}}(s_{h+1}^{*}|s_{h}^{*})} \right) \right] = \frac{0.75/0.25}{0.5/0.5} = 3$$
(120)

830 so $\alpha = \ln(3)$, and

$$\nu = \min_{s_h^*, s_{h+1}^*} \frac{\Pr_{\pi_A}(s_h^*, s_{h+1}^*) + \Pr_{\pi_B}(s_h^*, s_{h+1}^*)}{2} = \frac{1/4 + 1/16}{2} = \frac{5}{32},$$
(121)

831 and

$$\eta = \min_{s_h^*, s_{h+1}^*} \frac{\Pr_{\pi_B}(s_h^*, s_{h+1}^*)}{\Pr_{\pi_A}(s_h^*, s_{h+1}^*) + \Pr_{\pi_B}(s_h^*, s_{h+1}^*)} = \frac{1/16}{1/4 + 1/16} = \frac{1}{5}.$$
 (122)

For the "Single observation classification" and "Paired observation classification" baselines, we select the feature ϕ_h (or feature-pair ϕ_h, ϕ_{h+1}) such that the mutual information between $\phi_h(x_h)$ (respectively, $(\phi_h(x_h), \phi_{h+1}(x_{h+1}))$) and the agent's identity is maximized on the collected trajectories.

The "Average Encoder Accuracy" was computed based on the "population" behavior of the environment: that is, the accuracy of encoder ϕ'_h which extracts the feature $(s_h^* \text{ XOR } e_h^i)$ from x_h is computed as $\max(\Pr(e_h^i = 1), \Pr(e_h^i = 0))$, which can be determined analytically from the parameters of Markov chain e^i . For a given algorithm, this quantity was then averaged over timesteps for the returned encoder.

For the "Paired observation classification" baseline, note that for timesteps h = 2 through h = H-1, the suggested baseline could refer two distinct encoders: the encoder ϕ_h such that $\phi_h(x_h)$ and some 843 $\phi_{h+1}(x_{h+1})$ are together most informative at predicting the agent observing (x_h, x_{h+1}) ; or the

encoder ϕ'_h such that $\phi'_h(x_h)$ and some $\phi_{h-1}(x_{h-1})$ are together most informative at predicting the agent observing (x_{h-1}, x_h) . In reporting the final encoder accuracies, took the average accuracy of

846 these two encoders.