
Efficient Algorithms for Contextual Apple Tasting with Log-Loss

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This paper introduces two novel algorithmic approaches designed for the Con-
2 textual Apple Tasting problem, where the learner faces an asymmetric feedback
3 structure by observing binary labels only upon an ‘Accept’ action. To address
4 the inherent decision bias and exploration challenges, we propose two distinct
5 but complementary strategies. First, we introduce LogCBPSide-AT, an algorithm
6 leveraging Confidence Bounds for Partial monitoring (CBP) to explicitly quantify
7 predictive uncertainty and effectively balance the exploration-exploitation trade-off.
8 Second, we present LogCB-AT, an approach that reduces the apple tasting problem
9 to an online regression oracle. This reduction-based strategy offers a computation-
10 ally efficient and scalable alternative that fundamentally bypasses the complex,
11 often intractable confidence bound constructions required by traditional methods.
12 Theoretically, we prove that both algorithms achieve sublinear regret bounds for
13 losses associated with the binary labels, guaranteeing robust performance even
14 under fundamentally restricted feedback. The practical utility of our methods is
15 empirically validated through adaptive Large Language Model (LLM) cascading,
16 where they effectively optimize the trade-off between inference cost and response
17 accuracy.

18 1 Introduction

19 In the field of sequential decision-making, the standard contextual bandit framework serves as a
20 foundational model for learning under uncertainty. Within this setup, a learner observes a context
21 and selects an action, receiving feedback only for the chosen arm. Achieving high performance
22 in such tasks depends on effectively managing the balance between exploration and exploitation.
23 The agent must decide whether to exploit its greedy action to maximize immediate rewards or
24 explore alternative actions to gather information that enables more accurate future decision-making.
25 While this standard framework has proven effective across a wide range of practical problems [Li
26 et al., 2010, Tewari and Murphy, 2017, Durand et al., 2018], more challenging settings arise when
27 the information structure is fundamentally restricted. This challenge is especially prominent in
28 environments characterized by partial monitoring, where feedback may be entirely absent for certain
29 choices. A particularly notable subclass of this setting is the Apple Tasting problem, originally
30 introduced by Helmbold et al. [2000]. In this framework, the learner faces an asymmetric information
31 structure where the true label of an instance (e.g., whether the apple is rotten or not) is revealed only
32 if the learner decides to ‘Accept’ (or ‘Taste’) it. If the learner chooses to remain ‘Reject’ (or ‘Pass’),
33 no feedback is provided, leaving the true status of the instance unknown. The learner’s goal is to
34 minimize the accumulated loss induced by the interaction between its actions (Taste vs. Pass) and the
35 binary labels (rotten vs non-rotten), e.g., reducing the number of untasted rotten apples and tasted
36 non-rotten apples. This objective faces a fundamental challenge where errors in rejection (false
37 negatives) are never explicitly corrected, potentially leading to a persistent bias where the model

38 remains stuck in a suboptimal state due to the continuous absence of corrective feedback.

39

40 In this paper, we study the contextual apple tasting problem, where action selection is informed by
41 contextual information. The decision to Accept reveals the binary label at a specific cost, while
42 the Reject arm does not provide any feedback, and potentially results in large loss according to the
43 value of y_t . This setup is motivated by modern AI application scenarios, for which we present two
44 representative examples below:

45 **LLM cascading** The system initially processes a query through a Small Language Model (SLM)
46 and must determine whether to cascade to a Large Language Model (LLM) for a more refined result.
47 While the SLM provides an efficient baseline response, complex queries may require the LLM’s
48 superior reasoning at the cost of significantly higher resources and latency. Since the true performance
49 gain is only observable upon invoking the LLM, the system operates under partial feedback: staying
50 with the SLM prevents learning whether the LLM would have produced a better answer. Optimizing
51 this trade-off is critical to avoid either unnecessary computational cost or degraded user satisfaction
52 from suboptimal SLM outputs. In this example, invoking LLM and staying with SLM correspond to
53 ‘Accept’ and ‘Reject’ actions respectively, and whether a performance gain was obtained corresponds
54 to binary label. The sequence of queries serves as contextual information.

55 **Digital Healthcare** Wearable devices monitor bio-signals (contextual information) to determine
56 if a patient needs a clinical visit. This is a quintessential contextual apple tasting problem because
57 the actual presence of a disease (binary label) can usually only be confirmed if the device has
58 notified the patient to see a doctor (‘Accept’). If the device remains silent (‘Reject’), the true medical
59 status remains unobserved. In this context, failing to notify a patient who is actually ill results
60 in an exceptionally large loss. Conversely, notifying a healthy patient leads to significant patient
61 inconvenience and unnecessary medical expenses.

62 The asymmetric feedback structure necessitates a sophisticated exploration strategy to prevent the
63 model from becoming permanently biased due to missing data. In this paper, we propose two distinct
64 but complementary algorithmic approaches that efficiently tackle this challenge: LogCBPSide-AT, an
65 confidence bound algorithm based on a logistic model assumption on the label, and LogCB-AT, a
66 reduction-to-regression-oracle approach that accommodates any model class that can be learned under
67 the log-loss objective. The two algorithms cater to different modeling assumptions and theoretical
68 guarantees: the former employs a refined exploration-exploitations strategy within a more restrictive
69 model class with possibly logarithmic regret, while the latter accommodates a broader class of models
70 with $O(\sqrt{T})$ regret, where T is the time horizon.

71 First, we develop LogCBPSide-AT, an algorithm leveraging Confidence Bounds for Partial monitoring
72 (CBP) to explicitly quantify predictive uncertainty under the logistic model. While this approach
73 provides a principled and theoretically rigorous mechanism for balancing exploration and exploitation
74 under partial feedback [Bartók and Szepesvári, 2012, Heuillet et al., 2024], constructing explicit
75 confidence bounds can become computationally demanding when scaling to complex function
76 classes. To provide a scalable alternative that accommodates a broader function class, including
77 the logistic model, we subsequently introduce LogCB-AT, a reduction-based approach that adapts
78 the SquareCB algorithm of Foster and Rakhlin [2020] to the binary, partial feedback setting. This
79 approach fundamentally bypasses the challenge of explicit bound construction by reducing the bandit
80 learning task to repeated calls to an online regression oracle with a log-loss objective. Consequently, it
81 eliminates the need for complex exploration terms and allows for straightforward integration of black-
82 box or complex predictive models. Importantly, both algorithms, LogCBPSide-AT and LogCB-AT,
83 adapt their respective base algorithms—CBPSide and SquareCB— to the one-sided feedback structure,
84 by exploiting the structural property of the Apple Tasting problem: the latent true label y_t is a shared
85 ground truth independent of the learner’s action.

86 2 Problem Setting and Preliminaries

87 We consider a contextual binary decision problem over T rounds. At each round $t \in [T] :=$
88 $\{1, \dots, T\}$, the learner observes a context vector $x_t \in \mathcal{X} \subset \mathbb{R}^d$ and must choose an action $a_t \in$
89 $\mathcal{A} := \{0, 1\}$, where 0 and 1 denote ‘Reject’ and ‘Accept’, respectively. After the action is chosen,
90 a latent binary label $y_t \in \{0, 1\}$ is realized according to an unknown conditional distribution. The

91 defining feature of apple tasting is its one-sided feedback structure. Specifically, the true label y_t
 92 is revealed only when the learner chooses the revealing action ($a_t = 1$); if the learner chooses the
 93 Reject action ($a_t = 0$), no feedback is observed.

94 Our basic assumption is realizability: the learner has access to a function class $\mathcal{F} \subseteq (\mathcal{X} \rightarrow [0, 1])$
 95 that contains the true conditional expectation of the binary label.

Assumption 1 (Realizability). *There exists a function $f^* \in \mathcal{F}$ such that, for all $x \in \mathcal{X}$,*

$$f^*(x) = \mathbb{E}[y \mid x].$$

96

97 As for LogCBPSide-AT, we further assume (in Section 4) that \mathcal{F} is the class of logistic functions,
 98 whereas for LogCB-AT, we impose no additional restriction beyond the existence of a regression
 99 oracle that can efficiently learn f^* under a log-loss objective (Assumption 2 in Section 5). We denote
 100 $f^*(x_t) = \mathbb{E}[y_t \mid x_t]$ by y_t^* , the conditional mean at round t . The learner maintains an estimate \hat{f}_{t-1}
 101 of f^* , and its prediction at round t is $\hat{y}_t := \hat{f}_{t-1}(x_t)$. The learner is evaluated through the asymmetric
 102 loss matrix

	$y = 0$	$y = 1$
$a = 0$	0	ℓ_{01}
$a = 1$	ℓ_{10}	ℓ_{11}

103 where $\ell_{01}, \ell_{10}, \ell_{11} \in [0, 1]$. Here, ℓ_{01} is the loss of rejecting an item with label 1, ℓ_{10} is the loss of
 104 accepting an item with label 0, and ℓ_{11} is the cost incurred when accepting an item with label 1.
 105 Throughout the paper, we assume the following cost conditions

$$\ell_{01} > \ell_{11} \geq 0 \quad \text{and} \quad \ell_{10} > 0.$$

106 The strict inequality $\ell_{01} > \ell_{11}$ is necessary to prevent the Reject action ($a = 0$) from trivially always
 107 being the optimal choice. This aligns naturally with practical scenarios like LLM cascading: the
 108 performance penalty of missing a better answer (ℓ_{01}) inherently outweighs the cost of invoking the
 109 LLM (ℓ_{11}). Additionally, $\ell_{11} \geq 0$ simply reflects that any invocation incurs a basic cost, while
 110 $\ell_{10} > 0$ ensures that an unnecessary call strictly incurs a computational penalty.

111 For any scalar $z \in [0, 1]$, we define the loss function as $\psi_a(z) := \ell_{01}z(1-a) + \{\ell_{10} + (\ell_{11} - \ell_{10})z\}a$.
 112 For label y , $\psi_0(y)$ and $\psi_1(y)$ represent the losses of rejection and acceptance, respectively. Also,
 113 since $\psi_a(\cdot)$ is linear in z for each $a \in \mathcal{A}$, we have,

$$\mathbb{E}[\psi_a(y_t) \mid x_t] = \psi_a(f^*(x_t)) = \psi_a(y_t^*).$$

114 The optimal action at round t is defined as

$$a_t^* \in \arg \min_{a \in \mathcal{A}} \psi_a(y_t^*),$$

115 i.e., the action with the smaller conditional expected loss under the true model. The instantaneous
 116 regret is then defined by

$$r_t := \psi_{a_t}(y_t) - \psi_{a_t^*}(y_t),$$

117 and the cumulative regret is $R(T) := \sum_{t=1}^T r_t$. In Sections 4 and 5, we present two novel algorithms,
 118 LogCBPSide-AT and LogCB-AT, respectively, both of which achieve $O(\sqrt{T})$ upper bounds on $R(T)$
 119 with high probability under their respective assumptions. As for LogCBPSide-AT, we also present
 120 instance-dependent regret bound with polylogarithmic order in T .

121 3 Related work

122 Partial monitoring (PM) studies sequential decision problems where the learner does not observe the
 123 loss directly, but only a feedback signal determined by the chosen action and the hidden label. In
 124 stochastic finite PM, games are classically categorized into trivial, easy, hard, and hopeless classes
 125 according to their minimax regret rates [Bartók and Szepesvári, 2012, Bartók et al., 2014]. Apple
 126 tasting is one of the canonical PM examples, and is known to be a two-action, two-outcome easy
 127 game, for which the minimax regret scales on the order of \sqrt{T} [Heuillet et al., 2024].

128 Helmbold et al. [2000] was the first to formalize the apple tasting model, demonstrating that an
 129 $O(\sqrt{T})$ mistake bound is achievable through randomized algorithms in the realizable setting, where
 130 the target hypothesis is assumed to reside within the hypothesis class. This foundational work
 131 established the essential framework for transforming standard online learning models into the apple
 132 tasting paradigm. Subsequently, Raman et al. [2024] extended this research from a combinatorial
 133 perspective, providing a comprehensive analysis across both realizable and agnostic settings. By
 134 introducing a new combinatorial parameter known as ‘Effective Width,’ they established a trichotomy
 135 for mistake bounds in the realizable setting, showing that the expected number of mistakes must fall
 136 into one of three categories: $\Theta(1)$, $\Theta(\sqrt{T})$, or $\Theta(T)$. Notably, they resolved a long-standing open
 137 question by characterizing agnostic learnability via the Littlestone dimension. In contrast to previous
 138 literature that primarily focused on randomized algorithms, Chase and Mehalé [2024] explored the
 139 feasibility of deterministic algorithms within the realizable setting. They proved that any hypothesis
 140 class learnable under Apple tasting feedback is necessarily learnable by a deterministic algorithm.
 141 By providing an optimal mistake bound of $O(\sqrt{L(\mathcal{H})T \log T})$, where $L(\mathcal{H})$ denotes the Littlestone
 142 dimension of the hypothesis class \mathcal{H} , this study completed the theoretical foundation for deterministic
 143 Apple tasting.

144 While the aforementioned literature established the theoretical foundation of apple tasting in non-
 145 contextual settings, subsequent research expanded this framework to contextual bandits, where side
 146 information is available for decision-making. Bartók and Szepesvári [2012] extended the confi-
 147 dence bound based method of Bartók et al. [2012] for partial monitoring to the contextual setting,
 148 introducing CBPSide. More recently, Heuillet et al. [2024] proposed RandCBP and RandCBPSide
 149 by incorporating randomization into these frameworks to improve exploration efficiency. Bayesian
 150 strategies have also been adapted to better balance the exploration-exploitation trade-off in partial
 151 monitoring. Tsuchiya et al. [2020] established theoretical regret bounds for Thompson Sampling
 152 (TS), while Kirschner et al. [2020] proposed Information-Directed Sampling (IDS) to minimize the
 153 information ratio. Specifically, Grant and Leslie [2021] revisited the apple tasting problem, demon-
 154 strating the empirical superiority of these Bayesian mechanisms in managing restricted feedback. In
 155 addition, Harris et al. [2023] proposed an EXP-based algorithm to address apple tasting in adversarial
 156 scenarios. While LogCBPSide-AT is adapted from CBPSide, it is simplified and specialized for
 157 the logistic contextual apple tasting setting. Beyond this, we establish an instance-dependent regret
 158 bound for LogCBPSide-AT, a guarantee that, to our knowledge, has not been previously derived for
 159 CBPSide. We note that the partial monitoring literature predominantly focuses on minimax bounds;
 160 among the few works that establish instance-dependent guarantees, the MED-based approach of
 161 Komiyama et al. does not prove a minimax bound, whereas LogCBPSide-AT achieves both.

162 Recent contextual bandit research has increasingly explored reductions to online regression oracles,
 163 beginning with SquareCB[Foster and Rakhlin, 2020]. Within this line of work, Zhang et al. [2023]
 164 proposed SquareCB.G, adapting SquareCB [Foster and Rakhlin, 2020] to the partial monitoring
 165 setting, including Apple Tasting. However, the SquareCB.G framework is restricted to the use of a
 166 square-loss oracle, which may not be well-suited for binary labels, as in the original Apple Tasting
 167 problem and the aforementioned modern applications in LLMs and Digital Healthcare. In this
 168 work, we adapt the SquareCB algorithm to the binary feedback setting using a log-loss oracle. The
 169 derivation of the action-selection policy fundamentally differs from that of SquareCB.G., resulting in
 170 a fully randomized algorithm, as opposed to the partially deterministic policy used in SquareCB.G.

171 4 Confidence bound-based algorithm: LogCBPSide-AT

172 In this section, we present LogCBPSide-AT (Algorithm 1), a deterministic algorithm for the logistic
 173 contextual apple tasting problem, adapted from CBPSide Bartók and Szepesvári [2012]. We analyze
 174 logistic function class, $f^*(x) = \mathbb{E}[y | x] = \sigma(x^T \theta^*)$. The algorithm maintains an estimate of θ^*
 175 (maximum quasi-likelihood estimator) alongside a confidence interval of width $\beta_{t-1}^{x_t}(\delta_t)$ (defined
 176 in Algorithm 1) for the unknown latent label y_t^* given context x_t and selects actions based on the
 177 relationship between this confidence interval and a decision threshold.

178 We define the flip point $\alpha \in (0, 1)$ as the value at which the expected losses of actions Accept and
 179 Reject are equal:

$$\psi_0(\alpha) = \psi_1(\alpha) \implies \alpha = \ell_{10} / (\ell_{10} + \ell_{01} - \ell_{11})$$

Algorithm 1 LogCBPSide-AT

Input: $\lambda, \alpha, \{\delta_t\}$
Observe the context x_1 , Play action $a_1 = 1$ (Accept) and Observe y_1
 $N_1 = 1$ (N_t is the number of times the action Accept ($a_t = 1$) is chosen up to time t)
 $V_1 = \lambda I + x_1 x_1^T$
Compute $\hat{\theta}_1$ such that $(y_1 - \sigma(x_1^T \hat{\theta}_1)) x_1 = 0$
for $t > 1$ **do**
 Observe the context x_t , Calculate $\hat{y}_t = \sigma(x_t^T \hat{\theta}_{t-1})$
 Let $\beta_{t-1}^{x_t}(\delta_t) = \frac{1}{2c_\sigma} \|x_t\|_{V_{t-1}^{-1}} \sqrt{\left(3 + 2 \log\left(1 + \frac{2}{\lambda}\right)\right) \cdot 2d \log(N_{t-1}) \cdot \log\left(\frac{d}{\delta_t}\right)}$
 if $|\hat{y}_t - \alpha| < \beta_{t-1}^{x_t}(\delta_t)$ **then**
 Exploration phase
 Play action $a_t = 1$ (Accept) and observe y_t
 else
 Exploitation phase
 if $\hat{y}_t \leq \alpha$ **then**
 Play action $a_t = 0$ (Reject)
 else
 Play action $a_t = 1$ (Accept) and observe y_t
 end if
 end if
 Update $V_t = V_{t-1} + x_t x_t^T \cdot \mathbb{1}\{a_t = 1\}$
 Compute $\hat{\theta}_t$ such that $\sum_{s=1}^t \mathbb{1}\{a_s = 1\} (y_s - \sigma(x_s^T \hat{\theta}_t)) x_s = 0$
 $N_t = N_{t-1} + \mathbb{1}\{a_t = 1\}$
end for

180 For any context x_t , action $a = 0$ (Reject) is the optimal action whenever $y_t^* < \alpha$, and action $a = 1$
181 (Accept) is optimal otherwise. Accordingly, when the confidence interval lies entirely to one side of
182 α , LogCBPSide-AT acts greedily by exploiting its current estimate \hat{y}_t to select the optimal action.
183 When the confidence interval overlaps with α , the algorithm resolves this uncertainty by selecting
184 Accept, thereby obtaining feedback to refine future estimates and confidence bounds. We present the
185 full procedure for LogCBPSide-AT in Algorithm 1.

186 The theoretical analysis of LogCBPSide-AT is deferred to Appendix C. We state the main regret guar-
187 antee below. LogCBPSide-AT achieves the minimax-optimal regret bound of order \sqrt{T} , consistent
188 with the known lower bounds for easy partial monitoring games with side information.

Theorem 1 (LogCBPSide-AT minimax regret bound). *Assume $\sup_{x \in \mathcal{X}} \|x\|_2 \leq 1$ and $\sup_{x \in \mathcal{X}} |x^T \theta^*| \leq C_{\max}$. For $\delta_t = \frac{1}{t}$, the expected cumulative regret of the LogCBPSide-AT algorithm is bounded as:*

$$\mathbb{E} [\text{Reg}_T] \leq O\left(\left(1 + e^{C_{\max}}\right)^2 e^{-C_{\max}} d \sqrt{T} \log^{\frac{3}{2}}(T)\right).$$

(Ignoring logarithmic factors in d)

189

190 Furthermore, an important advantage of confidence-based algorithms in online learning with limited
191 feedback is that they often enjoy instance dependent regret bound of order $\text{polylog}(T)$. We prove
192 that LogCBPSide-AT also enjoys a $\log^2 T$ regret bound (see Theorem 2), a guarantee that, to our
193 knowledge, has not been previously derived for CBPSide. Proof of Theorem 2 is deferred to Appendix
194 C.2.

Theorem 2 (LogCBPSide-AT instance dependent regret bound). *Assume $\sup_{x \in \mathcal{X}} \|x\|_2 \leq 1$, $\sup_{x \in \mathcal{X}} |x^T \theta^*| \leq C_{\max}$, and the existence of $\Delta > 0$ such that $\Delta \leq \min_{x \in \mathcal{X}} |\alpha - \sigma(x^T \theta^*)|$.*

195

For $\delta_t = \frac{1}{t}$, the expected cumulative regret of the LogCBPSide-AT algorithm is bounded as:

$$\mathbb{E}[\text{Reg}_T] \leq O\left(\frac{1}{\Delta} d^2 \log^2 T\right)$$

(ignoring doubly logarithmic factors and $\text{polylog}(\frac{1}{\Delta})$)

196

197 5 SquareCB-based algorithm: LogCB-AT

198 We now introduce **LogCB-AT**, our adaptation of SquareCB to the contextual apple tasting setting.

199 Our algorithm relies on an online regression oracle, denoted by Alg_{KL} , which produces a prediction

200 $\hat{y}_t = \hat{f}_{t-1}(x_t)$ at round t . Since the true label is revealed only when the learner takes the Accept

201 action ($a_t = 1$), the oracle is updated only on those rounds.

Assumption 2 (Logarithmic-loss oracle guarantee). *For any (possibly adaptive) sequence of contexts, actions, and labels, the online log loss regression oracle Alg_{KL} satisfies*

$$\sum_{t=1}^T \mathbf{1}\{a_t = 1\} \ell_{\log}(\hat{y}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \mathbf{1}\{a_t = 1\} \ell_{\log}(f(x_t), y_t) \leq \text{Reg}_{\text{KL}}(T_1) \leq \text{Reg}_{\text{KL}}(T),$$

where $\hat{y}_t = \hat{f}_{t-1}(x_t)$, ℓ_{\log} denotes the standard log-loss (binary cross-entropy),

$$\ell_{\log}(\hat{y}, y) = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})],$$

and $T_1 := \sum_{t=1}^T \mathbf{1}\{a_t = 1\}$ denotes the number of rounds on which the action $a_t = 1$ is taken.

202

203 Our proposed algorithm, LogCB-AT, is presented in Algorithm 2. It closely follows the structure of
 204 SquareCB. At each round t , the learner first computes the greedy action b_t and defines the action-
 205 selection policy which a learning rate $\gamma > 0$ and an exploration parameter μ . The learner then
 206 samples an action a_t , observes the label y_t only if $a_t = 1$, and updates Alg_{KL} only on those rounds.
 207 A notable difference from SquareCB is that (i) the oracle uses a log-loss objective rather than a
 208 square-loss objective, and (ii) the oracle is updated only on rounds with $a_t = 1$. These differences
 209 call for non-trivial regret analysis and a new formula for the optimal value of γ which enables an
 210 $O(\sqrt{T})$ regret bound under mild guarantees of the oracle. We elaborate on details in the following
 211 subsections.

212 We analyze the regret of LogCB-AT algorithm, and propose a new formula for the learning rate γ
 213 which effectively minimizes the upper bound of the regret under the challenge of one-sided, binary
 214 feedback. Our analysis begins by partitioning the time steps $[T]$ into two disjoint sets, \mathcal{I}_1 and \mathcal{I}_0 ,
 215 according to the greedy action b_t : $\mathcal{I}_1 := \{t \in [T] : b_t = 1\}$ and $\mathcal{I}_0 := \{t \in [T] : b_t = 0\}$. The regret
 216 is then analyzed separately over the rounds in \mathcal{I}_1 and \mathcal{I}_0 .

217 **Regret for $t \in \mathcal{I}_1$** For $t \in \mathcal{I}_1$, the regret is decomposed according to the following corollary, which
 218 originates from Foster and Rakhlin [2020].

Corollary 1. *Suppose Assumptions 1 hold. With probability at least $1 - \delta$, we have*

$$\sum_{t \in \mathcal{I}_1} r_t \leq \sum_{t \in \mathcal{I}_1} \sum_{a \in \mathcal{A}} p_{t,a} \left[(\psi_a(y_t^*) - \psi_{a_t^*}(y_t^*)) - \frac{\gamma}{4} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 \right]$$

$$+ \frac{\gamma}{4} \sum_{t \in \mathcal{I}_1} \sum_{a \in \mathcal{A}} p_{t,a} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 + \max(\ell_{10}, \ell_{01} - \ell_{11}) \sqrt{2T \log(2\delta^{-1})}.$$

219

220 The first term of the right-hand side is upper bounded by $4T/\gamma$ under the action selection policy
 221 specified in Algorithm 2, with specific choice of $\mu = 2$. The challenge arises in the second term,

Algorithm 2 LogCB-AT

Require: $\gamma > 0, \mu > 0$, online regression oracle for log loss Alg_{KL}

- 1: **for** $t = 1, \dots, T$ **do**
- 2: Observe context $x_t \in \mathcal{X}$
- 3: Obtain prediction $\hat{y}_t = \hat{f}_{t-1}(x_t)$ from Alg_{KL}
- 4: Set $b_t \in \arg \min_{a \in \{0,1\}} \psi_a(\hat{y}_t)$
- 5: Set

$$p_{t,a} = \frac{1}{\mu + \gamma(\psi_a(\hat{y}_t) - \psi_{b_t}(\hat{y}_t))} \text{ for } a \neq b_t, \quad p_{t,b_t} = 1 - \sum_{a \neq b_t} p_{t,a}$$

- 6: Sample $a_t \sim p_t$
 - 7: **if** $a_t = 1$ **then**
 - 8: Observe y_t and update Alg_{KL} with (x_t, y_t)
 - 9: **end if**
 - 10: **end for**
-

222 which represents the squared loss of the prediction \hat{y}_t over all rounds in \mathcal{I}_1 , whereas Assumption 2
223 guarantees an upper bound only on the cumulative log loss over rounds with $a_t = 1$. To circumvent
224 this issue, we exploit the structural property of Apple Tasting that, at each time t , both actions share a
225 common latent label y_t , and leverage a probability-shift argument between $p_{t,0}$ and $p_{t,1}$ for rounds
226 in \mathcal{I}_1 . Finally, we employ the conversion between squared loss and log-loss used in Foster and
227 Krishnamurthy [2021] to bound the second term as follows,

$$\frac{\gamma}{4} \sum_{t \in \mathcal{I}_1} \sum_{a \in \mathcal{A}} p_{t,a} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 \leq \frac{\gamma}{4} \Lambda_1 \left(8\text{Reg}_{\text{KL}}(T) + 4 \log(2\delta^{-1}) \right),$$

228 where $\Lambda_1 = (\ell_{10} - \ell_{11})^2 + \ell_{01}^2 / (\mu - 1)$. The complete proofs are deferred to Appendix A. We
229 note that choosing $\gamma = O(\sqrt{T/\text{Reg}_{\text{KL}}(T)})$ yields a bound of $O(\sqrt{T\text{Reg}_{\text{KL}}(T)})$ for the cumulative
230 regret over rounds in \mathcal{I}_1 .

231 **Regret for $t \in \mathcal{I}_0$** A similar technique to that used for analyzing the regret over rounds in \mathcal{I}_1
232 could also be employed to bound the regret for $t \in \mathcal{I}_0$. However, when $b_t = 0$, the probability
233 ratio $p_{t,0}/p_{t,1}$ introduces an additional dependence on γ in the second term. As a result, setting
234 $\gamma = O(\sqrt{T/\text{Reg}_{\text{KL}}(T)})$ leads to a regret bound of $O(T\text{Reg}_{\text{KL}}(T))$. The optimal choice of γ
235 for rounds in \mathcal{I}_0 instead turns out to be $O(\{T/\text{Reg}_{\text{KL}}(T)\}^{1/3})$, which yields a suboptimal regret
236 bound of $O(\{T\text{Reg}_{\text{KL}}(T)\}^{2/3})$. To achieve $O(\sqrt{T\text{Reg}_{\text{KL}}(T)})$ regret for rounds in \mathcal{I}_0 as well, we
237 devise a novel relationship between the regret and prediction loss that fundamentally differs from
238 Corollary 2. Specifically, we further decompose the rounds in \mathcal{I}_0 according to whether $a_t^* = 1$ or
239 $a_t^* = 0$, and derive the following inequalities for each case:

240 For $t \in \mathcal{I}_0 \cap \{t \in [T] : a_t^* = 1\}$,

$$\mathbb{E}[r_t] \leq C_\Delta \mathbb{E}[I(a_t = 1)|y_t^* - \hat{y}_t|] + \gamma C_\Delta^2 \mathbb{E}[I(a_t = 1)(y_t^* - \hat{y}_t)^2],$$

241 and for $t \in \mathcal{I}_0 \cap \{t \in [T] : a_t^* = 0\}$,

$$\mathbb{E}[r_t] \leq \frac{1}{\gamma} + C_\Delta \mathbb{E}[I(a_t = 1)|\hat{y}_t - y_t^*|],$$

242 where $C_\Delta = |\ell_{01} - \ell_{11} + \ell_{10}|$. We note that neither inequality involves a quadratic dependence on γ ,
243 and moreover, the prediction loss in the right-hand side is accumulated only over rounds with $a_t = 1$.
244 Consequently, the cumulated regret over rounds in \mathcal{I}_0 enjoys the same $O(\sqrt{T\text{Reg}_{\text{KL}}(T)})$ bound
245 with the same choice of γ . Detailed proof is provided in Appendix A.

246 Building upon the analysis above, we present the regret bound for the LogCB-AT algorithm.

Theorem 3 (LogCB-AT Regret Bound). Suppose Assumptions 1 and 2 hold. For the choice of $\gamma = O(\sqrt{T/\text{Reg}_{\text{KL}}(T)})$, the cumulative regret of the LogCB-AT algorithm is bounded by:

$$R(T) = O\left(\sqrt{T \cdot \text{Reg}_{\text{KL}}(T)}\right)$$

with probability at least $1 - \delta$.

247

248 **Remark.** LogCB-AT guarantees strictly positive, closed-form probabilities for all actions, satisfying
 249 $p_{t,a} \geq \frac{1}{2+\gamma \max(l_{01}, l_{11}, l_{10})}$ for all $a \in \{0, 1\}$. This fully randomized design contrasts with partially
 250 deterministic methods such as SquareCB.G and CBPSide, which may assign zero probability to
 251 non-greedy actions. This distinction is particularly important for Off-Policy Evaluation (OPE), where
 252 estimators such as IPW require non-zero action probabilities for valid counterfactual inference [Bian
 253 and Jun, 2022, Balagopalan and Jun, 2025]. By explicitly lower bounding the action probabilities,
 254 LogCB-AT avoids unstable IPW weights and enables robust offline evaluation.

255 **Comparison of LogCBPSide-AT and LogCB-AT** In this work, we proposed two complementary
 256 algorithms for the contextual Apple Tasting problem with provable guarantees under respective
 257 assumptions. In particular, LogCBPSide-AT exhibits strong theoretical guarantees, with the possibility
 258 of achieving logarithmic regret under suitable assumptions in the logistic-linear setting. On the
 259 other hand, LogCB-AT does not attain logarithmic regret, but instead offers greater flexibility by
 260 accommodating richer hypothesis classes beyond the logistic-linear setting. Therefore, when a good
 261 and rich representation is available and the logistic-linear assumption is appropriate, LogCBPSide-AT
 262 may be preferred due to its stronger guarantees. In contrast, when the logistic-linear assumption is
 263 hard to validate, LogCB-AT provides a more flexible alternative. In this sense, the two approaches
 264 serve as complementary alternatives, each being advantageous under different modeling assumptions
 265 and application scenarios.

266 6 Experiments

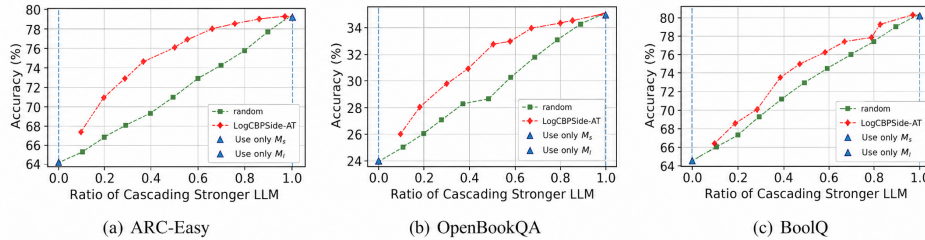


Figure 1: Accuracy vs. LLM Invocation Ratio: LogCBPSide-AT

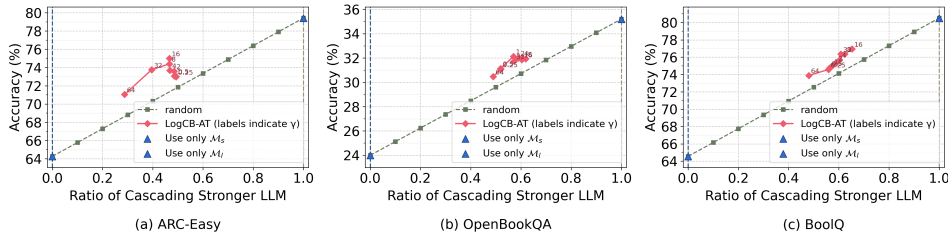


Figure 2: Accuracy vs. LLM Invocation Ratio: LogCB-AT

267 **Task Formulation** In this section, we empirically evaluate the performance of the proposed LogCB-
 268 AT and LogCBPSide-AT algorithms on the practical real-world application of test-time adaptive
 269 Large Language Model (LLM) cascading, which we frame as a contextual Apple Tasting problem,
 270 namely AppleTeA: Test-Time Adaptive LLM Cascading via Logistic Apple Tasting. The goal is to
 271 dynamically route streaming queries to either a computationally efficient Small Language Model
 272 (SLM) or a highly capable Large Language Model (LLM), thereby optimizing the trade-off between
 273 inference cost and response accuracy. At each time step t , the system receives a query q_t and first
 274 generates a response using the SLM (\mathcal{M}_s). The algorithm then makes a binary decision $a_t \in \{0, 1\}$:
 275 choosing $a_t = 0$ (Reject) finalizes the SLM response without additional cost, while choosing $a_t = 1$
 276 (Accept) invokes the stronger LLM (\mathcal{M}_l) at a higher computational cost. The ground-truth label is
 277 defined as $y_t = 1$ if the response generated by \mathcal{M}_s differs from that of \mathcal{M}_l , and $y_t = 0$ otherwise.
 278 Importantly, this label is observed only when $a_t = 1$, naturally inducing a partial-feedback apple
 279 tasting structure. We carefully construct the loss matrix to encode two competing cost components
 280 1) the invocation cost of the strong LLM (i.e., token-level inference cost) and 2) the price of error
 281 incurred when the weak SLM is invoked in contexts where prediction failure is likely. Further details
 282 on the experiments, including the construction of contexts, datasets, and the choice of SLM/LLM
 283 models, are provided in Appendix D.

284 **Baselines and Evaluation Metrics.** We compare our algorithms against three baselines: (i) \mathcal{M}_s -
 285 only, (ii) \mathcal{M}_l -only, and (iii) random cascading with a fixed probability $p_{random} \in [0, 1]$. In the figure,
 286 the x-axis represents the proportion of total queries routed to the stronger LLM. Moving to the right
 287 indicates an increase in invocations, which directly translates to higher operational costs. The y-axis
 288 reflects the final accuracy resulting from these routing decisions.

289 **Experimental Results** Figure 1 shows that LogCBPSide-AT forms a smooth, concave curve,
 290 demonstrating a highly efficient trade-off between accuracy and LLM invocations. Figure 2 illustrates
 291 that LogCB-AT, as a probabilistic approach, does not deterministically greedily invoke the LLM
 292 solely to maximize accuracy. Instead, it inherently preserves exploration, effectively managing the
 293 exploration-exploitation trade-off to maintain robust cost-efficiency.

294 7 Conclusion

295 In this paper, we introduced two novel algorithms, LogCBPSide-AT and LogCB-AT, designed to
 296 tackle the contextual Apple Tasting problem. Specifically, LogCBPSide-AT specializes the partial
 297 monitoring-based CBPSide framework for the Apple Tasting scenario by introducing a logistic setting.
 298 Concurrently, LogCB-AT adapts the contextual bandit-based SquareCB algorithm to operate within
 299 a log-loss setting. By leveraging these frameworks, we provided rigorous theoretical guarantees,
 300 demonstrating that our algorithms achieve sublinear regret bounds. Beyond the theoretical formu-
 301 lations, we empirically validated the practical utility of our approaches in a real-world application
 302 of test-time adaptive Large Language Model (LLM) cascading. Our experiments showed that dy-
 303 namically routing queries between a computationally efficient Small Language Model (SLM) and a
 304 highly capable LLM via these algorithms effectively optimizes the trade-off between inference cost
 305 and response accuracy. Building upon these foundational results, our future work aims to extend
 306 this framework to achieve first-order regret bounds. By adapting the FastCB algorithm [Foster and
 307 Krishnamurthy, 2021], we plan to develop a more advanced variant that scales with the loss of the
 308 optimal policy, yielding tighter regret guarantees and further accelerating convergence in practical
 309 environments.

310 **References**

- 311 Marc Abeille and Alessandro Lazaric. Linear Thompson Sampling Revisited. In Proceedings of the
312 International Conference on Artificial Intelligence and Statistics (AISTATS), volume 54, pages
313 176–184, 2017.
- 314 Kapilan Balagopalan and Kwang-Sung Jun. Minimum empirical divergence for sub-gaussian linear
315 bandits. In Yingzhen Li, Stephan Mandt, Shipra Agrawal, and Emtiyaz Khan, editors, Proceedings
316 of The 28th International Conference on Artificial Intelligence and Statistics, volume 258 of
317 Proceedings of Machine Learning Research, pages 1585–1593. PMLR, 03–05 May 2025.
- 318 Gábor Bartók and Csaba Szepesvári. Partial monitoring with side information. In International
319 Conference on Algorithmic Learning Theory, pages 305–319. Springer, 2012.
- 320 Gábor Bartók, Navid Zolghadr, and Csaba Szepesvári. An adaptive algorithm for finite stochastic
321 partial monitoring. ICML’12, page 1779–1786, Madison, WI, USA, 2012. Omnipress. ISBN
322 9781450312851.
- 323 Gábor Bartók, Dean P Foster, Dávid Pál, Alexander Rakhlin, and Csaba Szepesvári. Partial monitor-
324 ing—classification, regret bounds, and algorithms. Mathematics of Operations Research, 39(4):
325 967–997, 2014.
- 326 Jie Bian and Kwang-Sung Jun. Maillard sampling: Boltzmann exploration done optimally.
327 In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera, editors, Proceedings of
328 The 25th International Conference on Artificial Intelligence and Statistics, volume 151 of
329 Proceedings of Machine Learning Research, pages 54–72. PMLR, 28–30 Mar 2022. URL
330 <https://proceedings.mlr.press/v151/bian22a.html>.
- 331 Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric
332 Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al.
333 Pythia: A suite for analyzing large language models across training and scaling. In International
334 conference on machine learning, pages 2397–2430. PMLR, 2023.
- 335 Zachary Chase and Idan Mehal. Deterministic apple tasting. arXiv preprint arXiv:2410.10404,
336 2024.
- 337 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
338 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.
339 arXiv preprint arXiv:1803.05457, 2018.
- 340 Audrey Durand, Charis Achilleos, Demetris Iacovides, Katerina Strati, Georgios D Mitsis, and Joelle
341 Pineau. Contextual bandits for adapting treatment in a mouse model of de novo carcinogenesis. In
342 Machine learning for healthcare conference, pages 67–82. PMLR, 2018.
- 343 Sarah Filippi, Olivier Cappé, Aurélien Garivier, and Csaba Szepesvári. Parametric Bandits: The
344 Generalized Linear Case. In Advances in Neural Information Processing Systems (NeurIPS),
345 pages 586–594, 2010.
- 346 Dylan Foster and Alexander Rakhlin. Beyond ucb: Optimal and efficient contextual bandits with
347 regression oracles. In International conference on machine learning, pages 3199–3210. PMLR,
348 2020.
- 349 Dylan J Foster and Akshay Krishnamurthy. Efficient first-order contextual bandits: Prediction,
350 allocation, and triangular discrimination. Advances in Neural Information Processing Systems,
351 34:18907–18919, 2021.
- 352 James A Grant and David S Leslie. Apple tasting revisited: Bayesian approaches to partially
353 monitored online binary classification. arXiv preprint arXiv:2109.14412, 2021.
- 354 Keegan Harris, Chara Podimata, and Steven Z Wu. Strategic apple tasting. Advances in Neural
355 Information Processing Systems, 36:79918–79945, 2023.
- 356 David P Helmbold, Nicholas Littlestone, and Philip M Long. Apple tasting. Information and
357 Computation, 161(2):85–139, 2000.

- 358 Maxime Heuillet, Ola Ahmad, and Audrey Durand. Randomized confidence bounds for stochastic
359 partial monitoring. [arXiv preprint arXiv:2402.05002](#), 2024.
- 360 Kwang-Sung Jun and Jungtaek Kim. Noise-adaptive confidence sets for linear bandits and application
361 to bayesian optimization. [Proceedings of the International Conference on Machine Learning](#)
362 [\(ICML\)](#), 2024.
- 363 Johannes Kirschner, Tor Lattimore, and Andreas Krause. Information directed sampling for linear
364 partial monitoring. In [Conference on Learning Theory](#), pages 2328–2369. PMLR, 2020.
- 365 Junpei Komiyama, Junya Honda, and Hiroshi Nakagawa. Regret lower bound and optimal algorithm
366 in finite stochastic partial monitoring.
- 367 Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to
368 personalized news article recommendation. In [Proceedings of the 19th international conference](#)
369 [on World wide web](#), pages 661–670, 2010.
- 370 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct
371 electricity? a new dataset for open book question answering. In [Proceedings of the 2018 conference](#)
372 [on empirical methods in natural language processing](#), pages 2381–2391, 2018.
- 373 Vinod Raman, Unique Subedi, Ananth Raman, and Ambuj Tewari. Apple tasting: Combinatorial
374 dimensions and minimax rates. In [The Thirty Seventh Annual Conference on Learning Theory](#),
375 pages 4358–4380. PMLR, 2024.
- 376 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An
377 adversarial winograd schema challenge at scale. [Communications of the ACM](#), 64(9):99–106,
378 2021.
- 379 Ambuj Tewari and Susan A Murphy. From ads to interventions: Contextual bandits in mobile health.
380 In [Mobile health: sensors, analytic methods, and applications](#), pages 495–517. Springer, 2017.
- 381 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
382 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
383 and fine-tuned chat models. [arXiv preprint arXiv:2307.09288](#), 2023.
- 384 Taira Tsuchiya, Junya Honda, and Masashi Sugiyama. Analysis and design of thompson sampling for
385 stochastic partial monitoring. [Advances in Neural Information Processing Systems](#), 33:8861–8871,
386 2020.
- 387 Mengxiao Zhang, Yuheng Zhang, Olga Vrousseau, Haipeng Luo, and Paul Mineiro. Practical
388 contextual bandits with feedback graphs. [Advances in Neural Information Processing Systems](#),
389 36:30592–30617, 2023.

Appendix

390

391 A LogCB-AT Regret Analysis

392 In this section, we derive a high-probability regret bound for the LogCB-AT algorithm. We aim to
 393 show that with probability at least $1 - \delta$, the cumulative regret is bounded by $O(\sqrt{T \cdot \text{Reg}_{\text{KL}}(T)})$.

394 A.1 Definition of Regret

395 The instantaneous regret at round t is defined as the difference between the loss of the action taken a_t
 396 and the loss of the optimal action a_t^* .

$$r_t = \psi_{a_t}(y_t) - \psi_{a_t^*}(y_t). \quad (1)$$

397 The cumulative regret $R(T)$ over T rounds is the sum of these instantaneous regrets.

$$R(T) = \sum_{t=1}^T r_t. \quad (2)$$

Lemma 1. For any $\hat{y}, y^* \in [0, 1]$, the KL-divergence between two Bernoulli distributions with parameters \hat{y} and y^* satisfies:

$$d_{\text{KL}}(\hat{y}, y^*) \geq \frac{1}{2} \frac{(\hat{y} - y^*)^2}{\hat{y} + y^*}. \quad (3)$$

Consequently, the following pointwise inequality holds for all $\hat{y}, y^* \in [0, 1]$:

$$(\hat{y} - y^*)^2 \leq 4d_{\text{KL}}(\hat{y}, y^*). \quad (4)$$

398

399 *Proof.* The first inequality is a known lower bound for the KL-divergence of Bernoulli distributions
 400 (see, e.g., Proposition 5 of Foster and Krishnamurthy [2021]).

401 Since $\hat{y}, y^* \in [0, 1]$, their sum is at most 2 ($\hat{y} + y^* \leq 2$). Rearranging the first inequality, we obtain

$$(\hat{y} - y^*)^2 \leq 2(\hat{y} + y^*)d_{\text{KL}}(\hat{y}, y^*) \leq 4d_{\text{KL}}(\hat{y}, y^*).$$

402

□

Lemma 2. Let $C_1 = |\ell_{11} - \ell_{10}|$ be the constant for the Accept action ($a = 1$). From the definition of ψ_1 , it follows that:

$$(\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2 = C_1^2(\hat{y}_t - y_t^*)^2. \quad (5)$$

Then, under the Assumption 2 and Freedman's inequality, with probability at least $1 - \delta$, the following bound holds:

$$\sum_{t=1}^T p_{t,1}(\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2 \leq 8C_1^2 \text{Reg}_{\text{KL}}(T) + 4C_1^2 \log(2\delta^{-1}). \quad (6)$$

403

404 *Proof.* Let $\mathcal{F}_{t-1} = \sigma((x_1, a_1, y_1), \dots, (x_{t-1}, a_{t-1}, y_{t-1}), x_t)$ be the filtration. Define

$$M_t := \mathbb{1}\{a_t = 1\}(\hat{y}_t - y_t^*)^2.$$

405 Set $Z_t := M_t - \mathbb{E}[M_t \mid \mathcal{F}_{t-1}]$. Since $(\hat{y}_t - y_t^*)^2 \leq 1$ and $0 \leq \mathbb{1}\{a_t = 1\} \leq 1$, the range is
 406 $0 \leq M_t \leq 1$. Then,

$$\mathbb{E}[M_t \mid \mathcal{F}_{t-1}] = p_{t,1}(\hat{y}_t - y_t^*)^2, \quad \mathbb{E}[Z_t^2 \mid \mathcal{F}_{t-1}] \leq \mathbb{E}[M_t \mid \mathcal{F}_{t-1}].$$

407 Applying Freedman's inequality with range bound 1 and $\eta = \frac{1}{2}$, with probability at least $1 - \delta$

$$\sum_{t=1}^T \mathbb{E}[M_t \mid \mathcal{F}_{t-1}] \leq 2 \sum_{t=1}^T M_t + 4 \log(2\delta^{-1}). \quad (7)$$

408 By Lemma 1, we know $d_{\text{KL}}(\hat{y}_t, y_t^*) \geq \frac{1}{2} \frac{(\hat{y}_t - y_t^*)^2}{\hat{y}_t + y_t^*}$. Since $\hat{y}_t, y_t^* \in [0, 1]$, we have $\hat{y}_t + y_t^* \leq 2$, which
 409 implies $(\hat{y}_t - y_t^*)^2 \leq 4d_{\text{KL}}(\hat{y}_t, y_t^*)$. Therefore, based on the Log-loss oracle assumption, we can
 410 bound the sum of M_t :

$$\sum_{t=1}^T M_t = \sum_{t=1}^T \mathbb{1}\{a_t = 1\} (\hat{y}_t - y_t^*)^2 \leq 4 \sum_{t=1}^T \mathbb{1}\{a_t = 1\} d_{\text{KL}}(\hat{y}_t, y_t^*) \leq 4 \text{Reg}_{\text{KL}}(T).$$

411 Plugging this into (7), we get

$$\sum_{t=1}^T p_{t,1}(\hat{y}_t - y_t^*)^2 \leq 8 \text{Reg}_{\text{KL}}(T) + 4 \log(2\delta^{-1}).$$

412 Finally, multiplying both sides by C_1^2 and using $(\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2 = C_1^2(\hat{y}_t - y_t^*)^2$, we obtain

$$\sum_{t=1}^T p_{t,1}(\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2 \leq 8C_1^2 \text{Reg}_{\text{KL}}(T) + 4C_1^2 \log(2\delta^{-1}).$$

413

□

Lemma 3. Suppose Assumptions 1 and 2 hold. Let $\mathcal{A} = \{1, 0\}$ be the action set, where $a = 1$ denotes the Accept action and $a = 0$ denotes the Reject action. Then, with probability at least $1 - \delta$, the cumulative regret $R(T)$ is bounded as

$$R(T) \leq \sum_{t=1}^T \sum_{a \in \mathcal{A}} p_{t,a} \left(\psi_a(y_t^*) - \psi_{a_t^*}(y_t^*) \right) + R' \sqrt{2T \log(2\delta^{-1})} \quad (8)$$

where $R' = \max(\ell_{10}, \ell_{01} - \ell_{11})$.

414

415 *Proof.* Let $\mathcal{F}_{t-1} = \sigma((x_1, a_1, y_1), \dots, (x_{t-1}, a_{t-1}, y_{t-1}), x_t)$ be the filtration representing the his-
 416 tory up to round t . The conditional expectation of the instantaneous regret $r_t = \psi_{a_t}(y_t) - \psi_{a_t^*}(y_t)$
 417 with respect to \mathcal{F}_{t-1} is given by

$$\begin{aligned} \mathbb{E}_{t-1}[r_t] &= \mathbb{E} \left[\psi_{a_t}(y_t) - \psi_{a_t^*}(y_t) \mid \mathcal{F}_{t-1} \right] \\ &= \sum_{a \in \mathcal{A}} p_{t,a} \left(\psi_a(y_t^*) - \psi_{a_t^*}(y_t^*) \right). \end{aligned}$$

418 We define a martingale difference sequence Z_t as the deviation of the realized regret from its
 419 conditional expectation

$$Z_t = r_t - \mathbb{E}_{t-1}[r_t].$$

420 By construction, $\mathbb{E}_{t-1}[Z_t] = 0$. To apply Azuma-Hoeffding's inequality, we determine the range of
 421 r_t by considering the possible values of the true label $y_t \in \{0, 1\}$:

- 422 • **Case $y_t = 0$:** The losses are $\psi_0(0) = 0$ (Reject) and $\psi_1(0) = \ell_{10}$ (Accept). Since $\ell_{10} \geq 0$
 423 by Assumption 1, the optimal action is $a_t^* = 0$. Thus, $r_t = \psi_{a_t}(0) - \psi_0(0) = \psi_{a_t}(0)$.
 424 Since $a_t \in \{0, 1\}$, we have $r_t \in \{0, \ell_{10}\}$.

425 • **Case** $y_t = 1$: The losses are $\psi_0(1) = \ell_{01}$ (Reject) and $\psi_1(1) = \ell_{11}$ (Accept). Since
 426 $\ell_{01} \geq \ell_{11}$ by Assumption 2, the optimal action is $a_t^* = 1$. Thus, $r_t = \psi_{a_t}(1) - \psi_1(1)$. If
 427 $a_t = 1$, $r_t = 0$; if $a_t = 0$, $r_t = \ell_{01} - \ell_{11}$. This implies $r_t \in \{0, \ell_{01} - \ell_{11}\}$.

428 Combining both cases, the instantaneous regret is bounded within the interval $r_t \in [0, \max(\ell_{10}, \ell_{01} -$
 429 $\ell_{11})]$. Let $R' = \max(\ell_{10}, \ell_{01} - \ell_{11})$ be this upper bound. Since $0 \leq \mathbb{E}_{t-1}[r_t] \leq R'$, the range of the
 430 martingale difference Z_t is also bounded by R' .

431 Applying Azuma-Hoeffding's inequality to the sum $\sum_{t=1}^T Z_t$, we have that with probability at least
 432 $1 - \delta$

$$\sum_{t=1}^T (r_t - \mathbb{E}_{t-1}[r_t]) \leq R' \sqrt{2T \log(\delta^{-1})}.$$

433 we obtain

$$\begin{aligned} R(T) &= \sum_{t=1}^T r_t \leq \sum_{t=1}^T \mathbb{E}_{t-1}[r_t] + R' \sqrt{2T \log(2\delta^{-1})} \\ &= \sum_{t=1}^T \sum_{a \in \mathcal{A}} p_{t,a} \left(\psi_a(y_t^*) - \psi_{a_t^*}(y_t^*) \right) + R' \sqrt{2T \log(2\delta^{-1})}. \end{aligned}$$

434 This completes the proof. \square

435 A.2 Regret Decomposition

436 We partition the total rounds $[T]$ into two disjoint sets based on the learner's greedy decision
 437 $b_t = \arg \min_{a \in \{1,0\}} \psi_a(\hat{y}_t)$:

- 438 • B_1 **Rounds** (\mathcal{I}_1): $\{t \in [T] : b_t = 1\}$.
- 439 • B_0 **Rounds** (\mathcal{I}_0): $\{t \in [T] : b_t = 0\}$.

440 The regret in \mathcal{I}_0 requires an additional refined decomposition to address the lack of direct feedback.
 441 Specifically, we further categorize the rounds in \mathcal{I}_0 by comparing the greedy choice with the true
 442 optimal action a_t^* :

- 443 • $a_t^* = 1$ (Optimal is Accept): The optimal action is to Accept, but the greedy choice is to
 444 Reject. This represents a misclassification where the learner misses an opportunity to
 445 observe.
- 446 • $a_t^* = 0$ (Optimal is Reject): The optimal action is to Reject, and the greedy choice is correct.
 447 In this case, regret is only incurred due to the learner's stochastic exploration, specifically
 448 when the Accept action is sampled with probability $p_{t,1}$.

449 The resulting comprehensive decomposition of the cumulative regret $R(T)$ is as follows

$$\begin{aligned} R(T) &= \sum_{t \in \mathcal{I}_1} r_t + \sum_{t \in \mathcal{I}_0} r_t \\ &= \sum_{t \in \mathcal{I}_1} r_t + \sum_{t \in \mathcal{I}_0, a_t^* = 1} r_t + \sum_{t \in \mathcal{I}_0, a_t^* = 0} r_t \\ &= R_1(T) + R_{0,1}(T) + R_{0,0}(T). \end{aligned} \tag{9}$$

450 By applying Lemma 3, the realized cumulative regret $R(T)$ can be bounded by the sum of conditional
 451 expected regrets (denoted as \bar{R}) and the martingale concentration term. Specifically, with probability
 452 at least $1 - \delta$:

$$R(T) \leq \bar{R}_1(T) + \bar{R}_{0,1}(T) + \bar{R}_{0,0}(T) + R' \sqrt{2T \log(2\delta^{-1})}, \tag{10}$$

453 where $\bar{R}_1(T) = \sum_{t \in \mathcal{I}_1} \mathbb{E}_{t-1}[r_t]$, $\bar{R}_{0,1}(T) = \sum_{t \in \mathcal{I}_0, a_t^* = 1} \mathbb{E}_{t-1}[r_t]$, and $\bar{R}_{0,0}(T) =$
 454 $\sum_{t \in \mathcal{I}_0, a_t^* = 0} \mathbb{E}_{t-1}[r_t]$.

Corollary 2. Suppose Assumptions 1 and 2 hold. Following Lemma 3, with probability at least $1 - \delta$, the cumulative regret $R(T)$ can be decomposed and bounded as

$$R(T) \leq \sum_{t \in \mathcal{I}_1} \sum_{a \in \mathcal{A}} p_{t,a} \left[(\psi_a(y_t^*) - \psi_{a_t^*}(y_t^*)) - \frac{\gamma}{4} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 \right] \quad (11)$$

$$+ \frac{\gamma}{4} \sum_{t \in \mathcal{I}_1} \sum_{a \in \mathcal{A}} p_{t,a} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 \quad (12)$$

$$+ \bar{R}_0(T) + R' \sqrt{2T \log(2\delta^{-1})}. \quad (13)$$

455

456 *Proof.* From Lemma 3, with probability at least $1 - \delta$, the realized cumulative regret $R(T)$ is bounded
457 by the sum of conditional expected regrets and the martingale concentration term.

$$R(T) \leq \bar{R}_1(T) + \bar{R}_0(T) + R' \sqrt{2T \log(2\delta^{-1})},$$

458 where the conditional expected regrets are defined as $\bar{R}_1(T) = \sum_{t \in \mathcal{I}_1} \sum_{a \in \mathcal{A}} p_{t,a} (\psi_a(y_t^*) - \psi_{a_t^*}(y_t^*))$
459 and $\bar{R}_0(T) = \sum_{t \in \mathcal{I}_0} \sum_{a \in \mathcal{A}} p_{t,a} (\psi_a(y_t^*) - \psi_{a_t^*}(y_t^*))$.

460 To connect the regret in B_1 rounds ($\bar{R}_1(T)$) to the regression oracle, we employ an additive-subtractive
461 trick. We add and subtract the term $\frac{\gamma}{4} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2$ for each action in the $\bar{R}_1(T)$ summation.

$$\begin{aligned} \bar{R}_1(T) &= \sum_{t \in \mathcal{I}_1} \sum_{a \in \mathcal{A}} p_{t,a} \left[(\psi_a(y_t^*) - \psi_{a_t^*}(y_t^*)) - \frac{\gamma}{4} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 + \frac{\gamma}{4} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 \right] \\ &= \sum_{t \in \mathcal{I}_1} \sum_{a \in \mathcal{A}} p_{t,a} \left[(\psi_a(y_t^*) - \psi_{a_t^*}(y_t^*)) - \frac{\gamma}{4} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 \right] \\ &\quad + \frac{\gamma}{4} \sum_{t \in \mathcal{I}_1} \sum_{a \in \mathcal{A}} p_{t,a} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2. \end{aligned}$$

462

□

Lemma 4. Suppose Assumptions 1 and 2 hold. Let $C_1 = |\ell_{11} - \ell_{10}|$. By the definition of the loss function ψ , we have $(\psi_0(\hat{y}_t) - \psi_0(y_t^*))^2 = \ell_{01}^2 (\hat{y}_t - y_t^*)^2$. Then, with probability at least $1 - \delta$, the expected squared error of the loss estimates for rounds in \mathcal{I}_1 is bounded as

$$\sum_{t \in \mathcal{I}_1} \sum_{a \in \{1,0\}} p_{t,a} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 \leq \Lambda_1 \left(8\text{Reg}_{\text{KL}}(T) + 4 \log(2\delta^{-1}) \right), \quad (14)$$

where the constant Λ_1 is defined as $\Lambda_1 := C_1^2 + \frac{\ell_{01}^2}{\mu - 1}$.

463

464 *Proof.* For rounds $t \in \mathcal{I}_1$, the greedy action is to Accept ($b_t = 1$). We decompose the total squared
465 error sum into components for each action $a \in \{1,0\}$

$$\sum_{t \in \mathcal{I}_1} \sum_{a \in \{1,0\}} p_{t,a} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 = \sum_{t \in \mathcal{I}_1} p_{t,1} (\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2 + \sum_{t \in \mathcal{I}_1} p_{t,0} (\psi_0(\hat{y}_t) - \psi_0(y_t^*))^2. \quad (15)$$

466 For the first term (Accept action $a = 1$), we can apply the bound derived from Freedman's inequality
467 in Lemma 2, summing only over $t \in \mathcal{I}_1 \subseteq [T]$

$$\sum_{t \in \mathcal{I}_1} p_{t,1} (\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2 \leq \sum_{t=1}^T p_{t,1} (\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2 \leq C_1^2 \left(8\text{Reg}_{\text{KL}}(T) + 4 \log(2\delta^{-1}) \right). \quad (16)$$

468 For the second term (Reject action $a = 0$), we relate its error to the Accept action by bounding the
469 probability ratio. Since $b_t = 1$ for $t \in \mathcal{I}_1$, we know that $\psi_0(\hat{y}_t) > \psi_1(\hat{y}_t)$. Based on the SquareCB

470 sampling rule, the probabilities are defined as

$$p_{t,0} = \frac{1}{\mu + \gamma(\psi_0(\hat{y}_t) - \psi_1(\hat{y}_t))}, \quad p_{t,1} = 1 - p_{t,0}.$$

471 The probability ratio is then:

$$\frac{p_{t,0}}{p_{t,1}} = \frac{p_{t,0}}{1 - p_{t,0}} = \frac{1}{\mu + \gamma(\psi_0(\hat{y}_t) - \psi_1(\hat{y}_t)) - 1}.$$

472 Since this ratio is decreasing with respect to the difference $\psi_0(\hat{y}_t) - \psi_1(\hat{y}_t)$ (which is strictly positive),
473 the maximum occurs as the difference approaches 0. Thus, we have the upper bound

$$\frac{p_{t,0}}{p_{t,1}} < \frac{1}{\mu - 1}.$$

474 Using the identity $(\psi_0(\hat{y}_t) - \psi_0(y_t^*))^2 = \frac{\ell_{01}^2}{C_1^2}(\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2$, we relate the squared error terms

$$\begin{aligned} p_{t,0}(\psi_0(\hat{y}_t) - \psi_0(y_t^*))^2 &= \frac{p_{t,0}}{p_{t,1}} \cdot \frac{\ell_{01}^2}{C_1^2} \cdot p_{t,1}(\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2 \\ &< \frac{1}{\mu - 1} \cdot \frac{\ell_{01}^2}{C_1^2} \cdot p_{t,1}(\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2. \end{aligned}$$

475 Summing this over all $t \in \mathcal{I}_1$ and applying Lemma 2 yields

$$\begin{aligned} \sum_{t \in \mathcal{I}_1} p_{t,0}(\psi_0(\hat{y}_t) - \psi_0(y_t^*))^2 &\leq \frac{1}{\mu - 1} \frac{\ell_{01}^2}{C_1^2} \sum_{t=1}^T p_{t,1}(\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2 \\ &\leq \frac{\ell_{01}^2}{\mu - 1} \left(8\text{Reg}_{\text{KL}}(T) + 4 \log(2\delta^{-1}) \right). \end{aligned} \quad (17)$$

476 Substituting (16) and (17) back into (15), we arrive at the final bound

$$\begin{aligned} \sum_{t \in \mathcal{I}_1} \sum_{a \in \{1,0\}} p_{t,a}(\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 &\leq \left(C_1^2 + \frac{\ell_{01}^2}{\mu - 1} \right) \left(8\text{Reg}_{\text{KL}}(T) + 4 \log(2\delta^{-1}) \right) \\ &= \Lambda_1 \left(8\text{Reg}_{\text{KL}}(T) + 4 \log(2\delta^{-1}) \right). \end{aligned}$$

477 This completes the proof. □

Lemma 5. *Suppose Assumptions 1 and 2 hold. Let $\mathcal{I}_0 = \{t \in [T] : b_t = 0\}$ be the set of rounds where the greedy choice is the Reject action. Define the loss gap between the Reject and Accept actions as $\Delta(y) := \psi_0(y) - \psi_1(y)$, and let $C_\Delta := |\ell_{01} - \ell_{11} + \ell_{10}|$. By the definition of the linear loss functions, we have $|\Delta(y) - \Delta(y')| = C_\Delta |y - y'|$ for all $y, y' \in [0, 1]$. For the SquareCB exploration parameter $\mu = 2$, with probability at least $1 - \delta$, the sum of conditional expected regrets in B_0 rounds, $\bar{R}_0(T) = \bar{R}_{0,1}(T) + \bar{R}_{0,0}(T)$, is bounded by*

$$\begin{aligned} \bar{R}_0(T) &\leq \frac{T}{\gamma} + 2C_\Delta \sqrt{T \left(8\text{Reg}_{\text{KL}}(T) + 4 \log(2\delta^{-1}) \right)} \\ &\quad + \gamma C_\Delta^2 \left(8\text{Reg}_{\text{KL}}(T) + 4 \log(2\delta^{-1}) \right). \end{aligned} \quad (18)$$

478

479 *Proof.* For any round $t \in \mathcal{I}_0$, the greedy choice is $b_t = 0$, which implies that the estimated gap
480 is non-positive $\hat{\Delta}_t := \Delta(\hat{y}_t) \leq 0$. The true gap is denoted as $\Delta_t^* := \Delta(y_t^*)$. According to
481 the LogCB sampling rule, the exploration probability for the Accept action ($a = 1$) is given by
482 $p_{t,1} = \frac{1}{2 + \gamma(\psi_1(\hat{y}_t) - \psi_0(\hat{y}_t))} = \frac{1}{2 + \gamma(-\hat{\Delta}_t)}$.

483 We analyze the conditional expected regret $\mathbb{E}_{t-1}[r_t]$ for $t \in \mathcal{I}_0$. By definition, the conditional
 484 expected regret is the sum of instantaneous regrets over all possible sampled actions $a \in \{1, 0\}$,
 485 weighted by their exploration probabilities $p_{t,a}$:

$$\mathbb{E}_{t-1}[r_t] = \sum_{a \in \{1, 0\}} p_{t,a} (\psi_a(y_t^*) - \psi_{a_t^*}(y_t^*)). \quad (19)$$

486 Notice that when the sampled action matches the optimal action ($a = a_t^*$), the term becomes zero
 487 (e.g., $\psi_1(y_t^*) - \psi_1(y_t^*) = 0$). Thus, regret is strictly incurred only when the algorithm explores a
 488 suboptimal action ($a \neq a_t^*$). We consider two exhaustive cases based on the true optimal action
 489 $a_t^* \in \{1, 0\}$.

490 **Case 1:** $a_t^* = 1$ ($\bar{R}_{0,1}$) The optimal action is Accept, meaning the true gap is positive ($\Delta_t^* > 0$).
 491 Expanding the expected regret definition from (19), the regret for sampling the correct action ($a = 1$)
 492 vanishes:

$$\begin{aligned} \mathbb{E}_{t-1}[r_t] &= p_{t,1}(\psi_1(y_t^*) - \psi_1(y_t^*)) + p_{t,0}(\psi_0(y_t^*) - \psi_1(y_t^*)) \\ &= 0 + p_{t,0}\Delta_t^* = p_{t,0}\Delta_t^*. \end{aligned}$$

493 Using the fact that $p_{t,0} = 1 - p_{t,1}$, the probability ratio is

$$\frac{p_{t,0}}{p_{t,1}} = \frac{1 - p_{t,1}}{p_{t,1}} = \left(2 + \gamma(-\hat{\Delta}_t)\right) - 1 = 1 + \gamma(-\hat{\Delta}_t).$$

494 We can rewrite the expected regret as

$$\mathbb{E}_{t-1}[r_t] = p_{t,1} \frac{p_{t,0}}{p_{t,1}} \Delta_t^* = p_{t,1} (1 + \gamma(-\hat{\Delta}_t)) \Delta_t^*.$$

495 Since $\Delta_t^* > 0$ and $\hat{\Delta}_t \leq 0$, we bound the terms using the definition of Δ :

- 496 • $-\hat{\Delta}_t \leq \Delta_t^* - \hat{\Delta}_t \leq |\Delta_t^* - \hat{\Delta}_t| = C_\Delta |y_t^* - \hat{y}_t|$.
- 497 • $\Delta_t^* = \hat{\Delta}_t + (\Delta_t^* - \hat{\Delta}_t) \leq \Delta_t^* - \hat{\Delta}_t \leq |\Delta_t^* - \hat{\Delta}_t| = C_\Delta |y_t^* - \hat{y}_t|$ (because $\hat{\Delta}_t \leq 0$).

498 Substituting these upper bounds yields:

$$\mathbb{E}_{t-1}[r_t] \leq p_{t,1} (1 + \gamma C_\Delta |y_t^* - \hat{y}_t|) C_\Delta |y_t^* - \hat{y}_t| = C_\Delta p_{t,1} |y_t^* - \hat{y}_t| + \gamma C_\Delta^2 p_{t,1} (y_t^* - \hat{y}_t)^2.$$

499 **Case 2:** $a_t^* = 0$ ($\bar{R}_{0,0}$) The optimal action is Reject, meaning $\Delta_t^* \leq 0$. Even though the greedy
 500 action is correct ($b_t = a_t^* = 0$), regret is incurred due to stochastic exploration when $a = 1$ is
 501 sampled. Expanding the expected regret, the term for sampling the correct action ($a = 0$) vanishes:

$$\begin{aligned} \mathbb{E}_{t-1}[r_t] &= p_{t,1}(\psi_1(y_t^*) - \psi_0(y_t^*)) + p_{t,0}(\psi_0(y_t^*) - \psi_0(y_t^*)) \\ &= p_{t,1}(-\Delta_t^*) + 0 = p_{t,1}(-\Delta_t^*). \end{aligned}$$

502 We decompose $-\Delta_t^*$ by adding and subtracting $\hat{\Delta}_t$:

$$-\Delta_t^* = -\hat{\Delta}_t + (\hat{\Delta}_t - \Delta_t^*) \leq -\hat{\Delta}_t + |\hat{\Delta}_t - \Delta_t^*| = -\hat{\Delta}_t + C_\Delta |\hat{y}_t - y_t^*|.$$

503 Multiplying by $p_{t,1}$, we obtain:

$$\mathbb{E}_{t-1}[r_t] \leq p_{t,1}(-\hat{\Delta}_t) + C_\Delta p_{t,1} |\hat{y}_t - y_t^*|.$$

504 Notice that $p_{t,1}(-\hat{\Delta}_t) = \frac{-\hat{\Delta}_t}{2 + \gamma(-\hat{\Delta}_t)}$. Since the function $f(x) = \frac{x}{2 + \gamma x}$ is monotonically increasing
 505 for $x \geq 0$ and bounded by $\frac{1}{\gamma}$, and given $-\hat{\Delta}_t \geq 0$, we have $p_{t,1}(-\hat{\Delta}_t) \leq \frac{1}{\gamma}$. Therefore,

$$\mathbb{E}_{t-1}[r_t] \leq \frac{1}{\gamma} + C_\Delta p_{t,1} |\hat{y}_t - y_t^*|.$$

506 **Combining the Bounds** Summing the expected regrets from both cases over all rounds $t \in \mathcal{I}_0$, and
 507 relaxing the summation to all T rounds since $\mathcal{I}_0 \subseteq [T]$ and all terms are non-negative, we get:

$$\begin{aligned}\bar{R}_0(T) &= \sum_{t \in \mathcal{I}_0} \mathbb{E}_{t-1}[r_t] \leq \sum_{t \in \mathcal{I}_0} \left[\frac{1}{\gamma} + 2C_\Delta p_{t,1} |y_t^* - \hat{y}_t| + \gamma C_\Delta^2 p_{t,1} (y_t^* - \hat{y}_t)^2 \right] \\ &\leq \frac{T}{\gamma} + 2C_\Delta \sum_{t=1}^T p_{t,1} |y_t^* - \hat{y}_t| + \gamma C_\Delta^2 \sum_{t=1}^T p_{t,1} (y_t^* - \hat{y}_t)^2.\end{aligned}$$

508 Applying the Cauchy-Schwarz inequality to the middle term:

$$\sum_{t=1}^T p_{t,1} |y_t^* - \hat{y}_t| = \sum_{t=1}^T \sqrt{p_{t,1}} \cdot \sqrt{p_{t,1}} |y_t^* - \hat{y}_t| \leq \sqrt{\sum_{t=1}^T p_{t,1}} \sqrt{\sum_{t=1}^T p_{t,1} (y_t^* - \hat{y}_t)^2} \leq \sqrt{T \sum_{t=1}^T p_{t,1} (y_t^* - \hat{y}_t)^2}.$$

509 Substituting this back into the expected regret bound yields:

$$\bar{R}_0(T) \leq \frac{T}{\gamma} + 2C_\Delta \sqrt{T \sum_{t=1}^T p_{t,1} (y_t^* - \hat{y}_t)^2} + \gamma C_\Delta^2 \sum_{t=1}^T p_{t,1} (y_t^* - \hat{y}_t)^2.$$

510 From Lemma 2, we know that with probability at least $1 - \delta$, the sum of squared errors is bounded
 511 by $\sum_{t=1}^T p_{t,1} (\hat{y}_t - y_t^*)^2 \leq 8\text{Reg}_{\text{KL}}(T) + 4\log(2\delta^{-1})$. Applying this upper bound yields the final
 512 result. \square

Lemma 6 (Lemma 3 of SquareCB[Foster and Rakhlin, 2020]). *For any round t , the probability distribution p_t over the action set $\mathcal{A} = \{1, 0\}$ (where $K = 2$), chosen by the SquareCB algorithm, ensures that for any true loss vector defined by y_t^* :*

$$\sum_{a \in \mathcal{A}} p_{t,a} \left[(\psi_a(y_t^*) - \psi_{a_t^*}(y_t^*)) - \frac{\gamma}{4} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 \right] \leq \frac{4}{\gamma}. \quad (20)$$

Summing this inequality over all rounds $t \in \mathcal{I}_1$, we obtain:

$$\sum_{t \in \mathcal{I}_1} \sum_{a \in \mathcal{A}} p_{t,a} \left[(\psi_a(y_t^*) - \psi_{a_t^*}(y_t^*)) - \frac{\gamma}{4} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 \right] \leq \frac{4|\mathcal{I}_1|}{\gamma} \leq \frac{4T}{\gamma}. \quad (21)$$

513

514 *Proof.* The first inequality is a direct application of Lemma 3 in SquareCB[Foster and Rakhlin, 2020],
 515 adapted to our binary action setting with $K = 2$. By substituting the arbitrary loss f_a^* with our
 516 specific loss function $\psi_a(y_t^*)$ and the estimator \hat{y}_a with $\psi_a(\hat{y}_t)$, the bound holds for each individual
 517 round t . Summing this bound over all rounds $t \in \mathcal{I}_1$ and using the trivial upper bound $|\mathcal{I}_1| \leq T$
 518 yields the final cumulative result. \square

Theorem 3 (LogCB-AT Regret Bound). *Suppose Assumptions 1 and 2 hold. For the choice of $\gamma = O(\sqrt{T/\text{Reg}_{\text{KL}}(T)})$, the cumulative regret of the LogCB-AT algorithm is bounded by:*

$$R(T) = O\left(\sqrt{T \cdot \text{Reg}_{\text{KL}}(T)}\right)$$

with probability at least $1 - \delta$.

519

520 *Proof.* From Corollary 2,

$$\begin{aligned}R(T) &\leq \sum_{t \in \mathcal{I}_1} \sum_{a \in \{1,0\}} p_{t,a} \left[(\psi_a(y_t^*) - \psi_{a_t^*}(y_t^*)) - \frac{\gamma}{4} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 \right] \\ &\quad + \frac{\gamma}{4} \sum_{t \in \mathcal{I}_1} \sum_{a \in \{1,0\}} p_{t,a} (\psi_a(\hat{y}_t) - \psi_a(y_t^*))^2 \\ &\quad + \bar{R}_0(T) + R' \sqrt{2T \log(2\delta^{-1})}.\end{aligned}$$

521 By substituting the upper bounds derived in Lemma 4, Lemma 5, and Lemma 6 into the inequality,
 522 we obtain

$$\begin{aligned} R(T) &\leq \frac{4T}{\gamma} + \frac{\gamma}{4}\Lambda_1 \left(8\text{Reg}_{\text{KL}}(T) + 4\log(2\delta^{-1})\right) \\ &\quad + \frac{T}{\gamma} + 2C_\Delta \sqrt{T \cdot (8\text{Reg}_{\text{KL}}(T) + 4\log(2\delta^{-1}))} + \gamma C_\Delta^2 \left(8\text{Reg}_{\text{KL}}(T) + 4\log(2\delta^{-1})\right) \\ &\quad + R' \sqrt{2T \log(2\delta^{-1})}. \end{aligned}$$

523 Grouping the terms by $\frac{1}{\gamma}$ and γ , we have:

$$\begin{aligned} R(T) &\leq \frac{5T}{\gamma} + \gamma \left[\left(2\Lambda_1 + 8C_\Delta^2\right) \text{Reg}_{\text{KL}}(T) + \left(\Lambda_1 + 4C_\Delta^2\right) \log(2\delta^{-1}) \right] \\ &\quad + 2C_\Delta \sqrt{T \cdot (8\text{Reg}_{\text{KL}}(T) + 4\log(2\delta^{-1}))} + R' \sqrt{2T \log(2\delta^{-1})}. \end{aligned}$$

524 By setting the learning rate γ to balance the first two terms:

$$\gamma = \sqrt{\frac{5T}{(2\Lambda_1 + 8C_\Delta^2) \text{Reg}_{\text{KL}}(T) + (\Lambda_1 + 4C_\Delta^2) \log(2\delta^{-1})}}, \quad (22)$$

525 the sum of the $\frac{1}{\gamma}$ and γ terms simplifies to $2\sqrt{5T \left[(2\Lambda_1 + 8C_\Delta^2) \text{Reg}_{\text{KL}}(T) + (\Lambda_1 + 4C_\Delta^2) \log(2\delta^{-1}) \right]}$.

526 This yields the final bounded expression:

$$\begin{aligned} R(T) &\leq 2\sqrt{5T \left[(2\Lambda_1 + 8C_\Delta^2) \text{Reg}_{\text{KL}}(T) + (\Lambda_1 + 4C_\Delta^2) \log(2\delta^{-1}) \right]} \\ &\quad + 2C_\Delta \sqrt{T \cdot (8\text{Reg}_{\text{KL}}(T) + 4\log(2\delta^{-1}))} + R' \sqrt{2T \log(2\delta^{-1})}. \end{aligned}$$

527 Extracting the dominant terms with respect to T , we conclude

$$R(T) = O\left(\sqrt{T \cdot \text{Reg}_{\text{KL}}(T)}\right). \quad (23)$$

528

□

529 B SquareCB-AT Regret Analysis

530 In this section, we present the regret analysis for the Contextual Apple Tasting problem under
 531 the square loss setting, serving as an alternative to the log-loss framework used in LogCB-AT. To
 532 ensure compatibility with the standard SquareCB architecture, which inherently relies on square
 533 loss regression oracles, we adopt a **Linear Probability Model (LPM)** as our generative foundation.
 534 Under this setting, the conditional expectation is modeled linearly with a fixed offset:

$$y_t \mid x_t \sim \text{Bernoulli}\left(x_t^\top \theta^* + \frac{1}{2}\right). \quad (24)$$

535 Accordingly, we introduce our new algorithm (Algorithm 3) and assumption (Assumption 3).

Assumption 3 (Squared-loss oracle guarantee). *For any (possibly adaptive) sequence of contexts, actions, and labels, the regression oracle satisfies*

$$\sum_{t=1}^T \mathbf{1}\{a_t = 1\} (\hat{y}_t - f^*(x_t))^2 \leq \text{Reg}_{\text{Sq}}(T_1) \leq \text{Reg}_{\text{Sq}}(T),$$

536

Algorithm 3 SquareCB-AT

Require: Learning rate $\gamma > 0$, exploration parameter $\mu > 1$, online regression oracle Alg_{Sq}

- 1: **for** $t = 1, \dots, T$ **do**
 - 2: Observe context $x_t \in \mathcal{X}$
 - 3: Obtain prediction $\hat{y}_t = \hat{f}_{t-1}(x_t)$ from SqAlg
 - 4: Compute $\psi_0(\hat{y}_t)$ and $\psi_1(\hat{y}_t)$
 - 5: Set

$$b_t \in \arg \min_{a \in \{0,1\}} \psi_a(\hat{y}_t)$$
 - 6: **for each** $a \neq b_t$ **do**
 - 7: Set

$$p_{t,a} = \frac{1}{\mu + \gamma(\psi_a(\hat{y}_t) - \psi_{b_t}(\hat{y}_t))}$$
 - 8: **end for**
 - 9: Set

$$p_{t,b_t} = 1 - \sum_{a \neq b_t} p_{t,a}$$
 - 10: Sample $a_t \sim p_t$
 - 11: **if** $a_t = 1$ **then**
 - 12: Observe y_t and update Alg_{Sq} with (x_t, y_t)
 - 13: **end if**
 - 14: **end for**
-

where

$$T_1 := \sum_{t=1}^T \mathbf{1}\{a_t = 1\}$$

denotes the number of rounds on which the action A is taken.

537

538 Next, we introduce Lemma 7. This lemma serves a similar role to Lemma 2 in LogCB-AT.

Lemma 7. Let $C_1 = |\ell_{11} - \ell_{10}|$ be the constant for the Accept action ($a = 1$). From the definition of ψ_1 , it follows that:

$$(\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2 = C_1^2 (\hat{y}_t - y_t^*)^2. \quad (25)$$

Then, under Assumption 2 and Freedman's inequality, with probability at least $1 - \delta$, the following bound holds:

$$\sum_{t=1}^T p_{t,1} (\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2 \leq 2C_1^2 \text{Reg}_{\text{Sq}}(T) + 4C_1^2 \log(2\delta^{-1}). \quad (26)$$

539

540 *Proof.* Let $\mathcal{F}_{t-1} = \sigma((x_1, a_1, y_1), \dots, (x_{t-1}, a_{t-1}, y_{t-1}), x_t)$ be the filtration. Define

$$M_t := \mathbf{1}\{a_t = 1\} (\hat{y}_t - y_t^*)^2.$$

541 Set $Z_t := M_t - \mathbb{E}[M_t \mid \mathcal{F}_{t-1}]$. Since $(\hat{y}_t - y_t^*)^2 \leq 1$ and $0 \leq \mathbf{1}\{a_t = 1\} \leq 1$, the range is
 542 $0 \leq M_t \leq 1$. Then,

$$\mathbb{E}[M_t \mid \mathcal{F}_{t-1}] = p_{t,1} (\hat{y}_t - y_t^*)^2, \quad \mathbb{E}[Z_t^2 \mid \mathcal{F}_{t-1}] \leq \mathbb{E}[M_t \mid \mathcal{F}_{t-1}].$$

543 Applying Freedman's inequality with range bound 1 and $\eta = \frac{1}{2}$, with probability at least $1 - \delta$:

$$\sum_{t=1}^T \mathbb{E}[M_t \mid \mathcal{F}_{t-1}] \leq 2 \sum_{t=1}^T M_t + 4 \log(2\delta^{-1}). \quad (27)$$

544 By Assumption 2, $\sum_{t=1}^T M_t = \sum_{t=1}^T \mathbb{1}\{a_t = 1\}(\hat{y}_t - y_t^*)^2 \leq \text{Reg}_{\text{Sq}}(T)$. Plugging this into (27)

$$\sum_{t=1}^T p_{t,1}(\hat{y}_t - y_t^*)^2 \leq 2\text{Reg}_{\text{Sq}}(T) + 4\log(2\delta^{-1}).$$

545 Finally, multiplying both sides by C_1^2 and using $(\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2 = C_1^2(\hat{y}_t - y_t^*)^2$, we obtain

$$\sum_{t=1}^T p_{t,1}(\psi_1(\hat{y}_t) - \psi_1(y_t^*))^2 \leq 2C_1^2\text{Reg}_{\text{Sq}}(T) + 4C_1^2\log(2\delta^{-1}).$$

546

□

Theorem 4 (SquareCB-AT Regret Bound). *Suppose Assumptions 1 and 3 hold. Then, with probability at least $1 - \delta$, the cumulative regret of the SquareCB-AT algorithm is bounded by:*

$$R(T) = O\left(\sqrt{T \cdot \text{Reg}_{\text{Sq}}(T)}\right). \quad (28)$$

547

548 *Proof.* The proof structure is essentially identical to the regret analysis of LogCB-AT presented in
 549 Theorem 3. The core steps and the regret decomposition remain exactly the same, with the only
 550 modification being the substitution of Lemma 2 (used for the log-loss setting) with Lemma 7 (derived
 551 for the square-loss setting).

552 we configure the optimal learning rate γ as

$$\gamma = \sqrt{\frac{5T}{\left(\frac{\Lambda_1}{2} + 2C_\Delta^2\right)\text{Reg}_{\text{Sq}}(T) + (\Lambda_1 + 4C_\Delta^2)\log(2\delta^{-1})}}. \quad (29)$$

553 Substituting this γ simplifies the sum of the $\frac{1}{\gamma}$ and γ terms, yielding the explicit bounded expression:

$$\begin{aligned} R(T) &\leq 2\sqrt{5T \left[\left(\frac{\Lambda_1}{2} + 2C_\Delta^2\right)\text{Reg}_{\text{Sq}}(T) + (\Lambda_1 + 4C_\Delta^2)\log(2\delta^{-1}) \right]} \\ &\quad + 2C_\Delta\sqrt{T \cdot \left(2\text{Reg}_{\text{Sq}}(T) + 4\log(2\delta^{-1})\right)} + R'\sqrt{2T\log(2\delta^{-1})}. \end{aligned} \quad (30)$$

554 Finally, we conclude

$$R(T) = O\left(\sqrt{T \cdot \text{Reg}_{\text{Sq}}(T)}\right). \quad (31)$$

555

□

556 **C LogCBPSide-AT Regret Analysis**

557 In this section we provide the detailed proof for Theorem 1.

558 Our reward model is,

$$\mathbb{E} [y_t | x_t] = \sigma \left(x_t^T \cdot \theta^* \right)$$

559 Following the work of Filippi et al. [2010], we use *maximum quasi-likelihood estimator* for $\hat{\theta}_t$. That
560 is,

$$\sum_{s=1}^t \mathbb{1}\{a_s = 1\} \left(y_s - \sigma(x_s^T \hat{\theta}_t) \right) x_s = 0.$$

561 Since, obtained rewards y_s are conditionally independent of the context (x_s), we can modify the
562 Proposition 1 of Filippi et al. [2010] from Appendix A.2 by replacing t with N_{t-1} (number of times
563 action Reject was chosen).

Proposition 1. *Take any δ_t, N_{t-1} such that $0 < \delta_t < \min\{1, \frac{d}{e}\}, 1 \leq N_{t-1} \leq T$. Let x_t be any random variable. Let*

$$\beta_{t-1}^{x_t}(\delta_t) = \frac{2k_\sigma R_{\max}}{c_\sigma} \|x_t\|_{V_{t-1}^{-1}} \sqrt{\left(3 + 2 \log\left(1 + \frac{2}{\lambda}\right)\right) \cdot 2d \log(N_{t-1}) \cdot \log\left(\frac{d}{\delta_t}\right)}.$$

Then, with probability at least $1 - \delta_t$, it holds that

$$|\sigma(x_t^T \theta^*) - \sigma(x_t^T \hat{\theta}_{t-1})| \leq \beta_{t-1}^{x_t}(\delta_t).$$

where,

- $R_{\max} := \max_t y_t = 1$.
- $k_\sigma :=$ Lipschitz constant of σ .

$$\begin{aligned} \sigma'(z) &= \sigma(z)(1 - \sigma(z)) \\ &\leq \frac{1}{4} \\ k_\sigma &= \frac{1}{4}. \end{aligned}$$

- $C_{\max} = \sup_{x_t \in \mathcal{X}, \theta^* \in \Theta} x_t^T \cdot \theta^*$
- $c_\sigma := \inf \sigma'(z)$.

$$\begin{aligned} c_\sigma &= \inf \sigma(z)(1 - \sigma(z)) \\ &= \sup(|\sigma(z)|)(1 - \sup(|\sigma(z)|)) \\ &= \frac{e^{C_{\max}}}{(1 + e^{C_{\max}})^2}. \end{aligned}$$

564

565 Furthermore, we introduce 2 more Lemmas necessary for the regret analysis.

566 The Elliptical Potential Lemma (EPL),

Lemma 8 (Elliptical potential lemma adapted from Proposition 2 of Abeille and Lazaric [2017]). *Let $x_1, x_2, \dots, x_t \in \mathbb{R}^d$ be a sequence of vectors with $\|x_s\|_2 \leq 1, \forall s \in [t]$. Let*

567

$V_{t-1} = \lambda I + \sum_{s=1}^{t-1} x_s x_s^T$ for some $\lambda > 0$. Then,

$$\sum_{s=1}^t \|x_s\|_{V_{s-1}}^2 \leq 2d \log\left(1 + \frac{t}{d\lambda}\right).$$

Adapting this to our case:

Let $x_1 \cdot \mathbb{1}\{a_1 = 1\}, x_2 \cdot \mathbb{1}\{a_2 = 1\}, \dots, x_t \cdot \mathbb{1}\{a_t = 1\} \in \mathbb{R}^d$ for some $\{a_s\}_{s=1}^t \in \{0, 1\}$. Let $V_{t-1} = \lambda I + \sum_{s=1}^{t-1} x_s x_s^T \mathbb{1}\{a_s = 1\}$ for some $\lambda > 0$. Then,

$$\sum_{s=1}^t \|x_s\|_{V_{s-1}}^2 \mathbb{1}\{a_s = 1\} \leq 2d \log\left(1 + \frac{t}{d\lambda}\right).$$

568

569 and Elliptical Potential Count Lemma (EPC)

Lemma 9 (Elliptical potential count lemma adapted from Lemma C.2 of Jun and Kim [2024]). Let $x_1, x_2, \dots, x_t \in \mathbb{R}^d$ be a sequence of vectors with $\|x_s\|_2 \leq 1, \forall s \in [t]$. Let $V_{t-1} = \lambda I + \sum_{s=1}^{t-1} x_s x_s^T$ for some $\lambda > 0$. Let $J = \{s \in [t] : \|x_s\|_{V_{s-1}}^2 \geq L^2\}$ for some $L^2 \leq 1$. Then,

$$|J| \leq 3 \frac{d}{L^2} \ln\left(1 + \frac{2}{L^2 \lambda}\right).$$

Adapting this to our case:

Let $x_1 \cdot \mathbb{1}\{a_1 = 1\}, x_2 \cdot \mathbb{1}\{a_2 = 1\}, \dots, x_t \cdot \mathbb{1}\{a_t = 1\} \in \mathbb{R}^d$ for some $\{a_s\}_{s=1}^t \in \{0, 1\}$. Let $V_{t-1} = \lambda I + \sum_{s=1}^{t-1} x_s x_s^T \mathbb{1}\{a_s = 1\}$ for some $\lambda > 0$. Let $J = \{s \in [t] : \|x_s \mathbb{1}\{a_s = 1\}\|_{V_{s-1}}^2 \geq L^2\}$ for some $L^2 \leq 1$. Then,

$$|J| \leq 3 \frac{d}{L^2} \ln\left(1 + \frac{2}{L^2 \lambda}\right).$$

570

571 C.1 Minimax regret analysis

572 **Theorem 1** (LogCBPSide-AT minimax regret bound). For $\delta_t = \frac{1}{t}$, with $\forall t \quad \|x_t\|_2 \leq 1$ and
573 $\sup_{x_t \in \mathcal{X}} x_t^T \theta^* \leq C_{\max}$, the expected cumulative regret of the LogCBPSide-AT algorithm is bounded
574 as:

$$\mathbb{E}[\text{Reg}_T] \leq O\left(\frac{(1 + e^{C_{\max}})^2}{e^{C_{\max}}} d \sqrt{T} \log^{\frac{3}{2}}(T)\right).$$

(Ignoring logarithmic terms in d)

575 Throughout the analysis, we use $\hat{y}_t = \sigma(x_t^T \hat{\theta}_{t-1})$ and $y_t^* = \sigma(x_t^T \theta^*)$ interchangeably to denote the
576 estimated and true latent labels, respectively.

577 Let us define the good event \mathcal{G}_t as,

$$\mathcal{G}_t = \left\{ |\sigma(x_t^T \theta^*) - \sigma(x_t^T \hat{\theta}_{t-1})| \leq \beta_{t-1}^{x_t}(\delta_t) \right\}$$

578 Instantaneous regret,

$$r_t = \psi_{a_t}(y_t) - \psi_{a_t^*}(y_t)$$

579 Let $\mathcal{F}_{t-1} = \sigma((x_1, a_1, y_1), \dots, (x_{t-1}, a_{t-1}, y_{t-1}), x_t)$ be the filtration representing the history up
 580 to round t . Furthermore, use the shorthand $\mathbb{E}_{t-1}[\cdot]$ to denote $\mathbb{E}[\cdot \mid \mathcal{F}_{t-1}]$.

$$\mathbb{E}_{t-1}[r_t] = \psi_{a_t}(y_t^*) - \psi_{a_t^*}(y_t^*)$$

581 Cumulative regret,

$$\begin{aligned} \text{Reg}_T &= \sum_{t=1}^T r_t \\ \mathbb{E}[\text{Reg}_T] &= \mathbb{E}\left[\sum_{t=1}^T r_t\right] \\ \mathbb{E}[\text{Reg}_T] &= \underbrace{\mathbb{E}\left[\sum_{t=1}^T r_t \mathbb{1}\{\mathcal{G}_t\}\right]}_{F_1} + \underbrace{\mathbb{E}\left[\sum_{t=1}^T r_t \mathbb{1}\{\bar{\mathcal{G}}_t\}\right]}_{F_2} \end{aligned}$$

582 First, let us bound F_2

$$\begin{aligned} F_2 &= \mathbb{E}\left[\sum_{t=1}^T r_t \mathbb{1}\{\bar{\mathcal{G}}_t\}\right] \\ &= \max(\ell_{10}, \ell_{01} - \ell_{11}) \sum_{t=1}^T \mathbb{P}(\bar{\mathcal{G}}_t) \\ &= \max(\ell_{10}, \ell_{01} - \ell_{11}) \sum_{t=1}^T \delta_t \\ &= \max(\ell_{10}, \ell_{01} - \ell_{11}) \sum_{t=1}^T \frac{1}{t} && \text{(Proposition 1, set } \delta_t = \frac{1}{t}\text{)} \\ &\lesssim \max(\ell_{10}, \ell_{01} - \ell_{11}) \log T \\ &= O(\log T) \end{aligned}$$

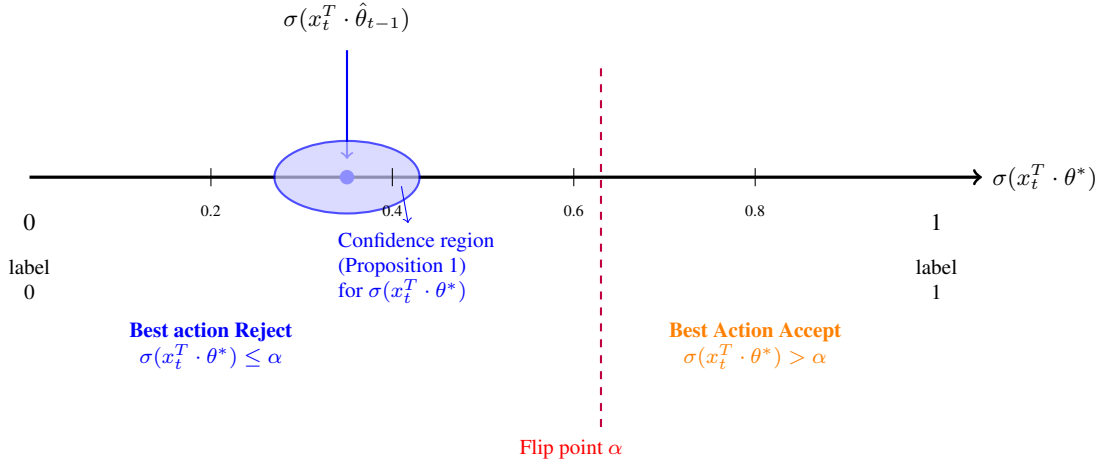
583 Now let us bound F_1 ,

$$\begin{aligned} F_1 &= \sum_{t=1}^T \mathbb{E}[r_t \mathbb{1}\{\mathcal{G}_t\}] \\ &= \underbrace{\sum_{t=1}^T \mathbb{E}\left[r_t \mathbb{1}\{\mathcal{G}_t\} \mathbb{1}\left\{\sigma(x_t^T \hat{\theta}_{t-1}) \leq \alpha - \beta_{t-1}^{x_t}(\delta_t)\right\}\right]}_{F_{11}} + \underbrace{\sum_{t=1}^T \mathbb{E}\left[r_t \mathbb{1}\{\mathcal{G}_t\} \mathbb{1}\left\{\sigma(x_t^T \hat{\theta}_{t-1}) \geq \alpha + \beta_{t-1}^{x_t}(\delta_t)\right\}\right]}_{F_{12}} \\ &\quad + \underbrace{\sum_{t=1}^T \mathbb{E}\left[r_t \mathbb{1}\{\mathcal{G}_t\} \mathbb{1}\left\{\alpha - \beta_{t-1}^{x_t}(\delta_t) < \sigma(x_t^T \hat{\theta}_{t-1}) < \alpha + \beta_{t-1}^{x_t}(\delta_t)\right\}\right]}_{F_{13}}. \end{aligned}$$

584 **Case F_{11} :** $\hat{y}_t = \sigma(x_t^T \hat{\theta}_{t-1}) \leq \alpha - \beta_{t-1}^{x_t}(\delta_t)$

585 First of all $|\hat{y}_t - \alpha| \geq \beta_{t-1}^{x_t}(\delta_t)$ and $\hat{y}_t \leq \alpha$, hence, $a_t = 0$ (Reject) by the behavior of Algorithm 1

586 This is the case where, estimated distribution \hat{y}_t is comfortably to the left of flip point, such that,
 587 based on Proposition 1, true mean of the distribution $\sigma(x_t^T \theta^*)$ itself will be on the left side of the flip
 588 point.



589

590 Since \mathcal{G}_t happens,

$$\begin{aligned}
 |\sigma(x_t^T \theta^*) - \sigma(x_t^T \hat{\theta}_{t-1})| &\leq \beta_{t-1}^{x_t}(\delta_t) \\
 \sigma(x_t^T \hat{\theta}_{t-1}) - \beta_{t-1}^{x_t}(\delta_t) &\leq \sigma(x_t^T \theta^*) \leq \sigma(x_t^T \hat{\theta}_{t-1}) + \beta_{t-1}^{x_t}(\delta_t) \\
 \sigma(x_t^T \theta^*) &\leq \alpha
 \end{aligned}$$

591 Hence, $a_t^* = 0$.

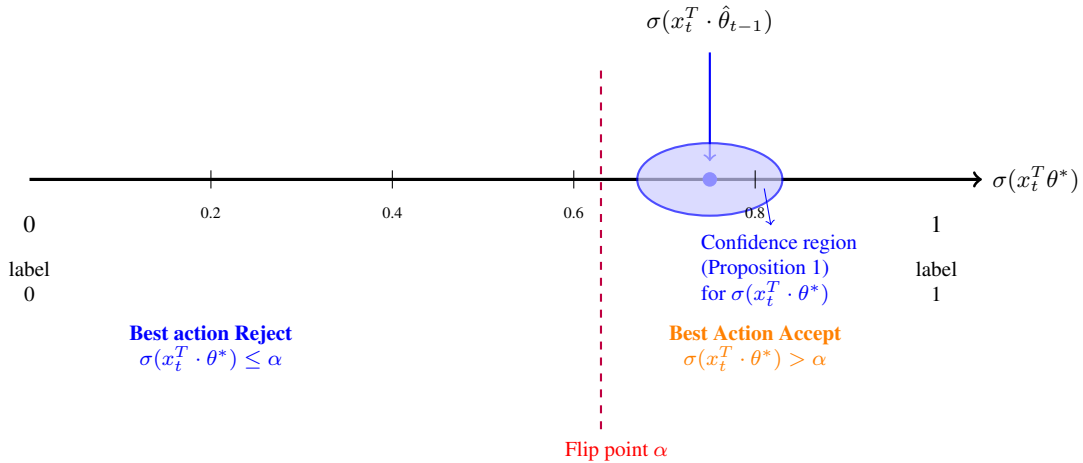
$$a_t = a_t^* = 0.$$

$$F_{11} = 0$$

592 **Case F_{12} :** $\hat{y}_t = \sigma(x_t^T \hat{\theta}_{t-1}) \geq \alpha + \beta_{t-1}^{x_t}(\delta_t)$

593 First of all $|\hat{y}_t - \alpha| \geq \beta_{t-1}^{x_t}(\delta_t)$ and $\hat{y}_t \geq \alpha$, hence, $a_t = 1$ (Accept)

594 This is the case where, estimated distribution \hat{y}_t is comfortably to the right of flip point, such that,
 595 based on Proposition 1, true mean of the distribution $\sigma(x_t^T \theta^*)$ itself will be on the right side of the
 596 flip point.



597

598 Since \mathcal{G}_t happens,

$$\begin{aligned}
|\sigma(x_t^T \theta^*) - \sigma(x_t^T \hat{\theta}_{t-1})| &\leq \beta_{t-1}^{x_t}(\delta_t) \\
\sigma(x_t^T \hat{\theta}_{t-1}) - \beta_{t-1}^{x_t}(\delta_t) &\leq \sigma(x_t^T \theta^*) \leq \sigma(x_t^T \hat{\theta}_{t-1}) + \beta_{t-1}^{x_t}(\delta_t) \\
\sigma(x_t^T \theta^*) &\geq \alpha
\end{aligned}$$

599 Hence $a_t^* = 1$ (Accept).

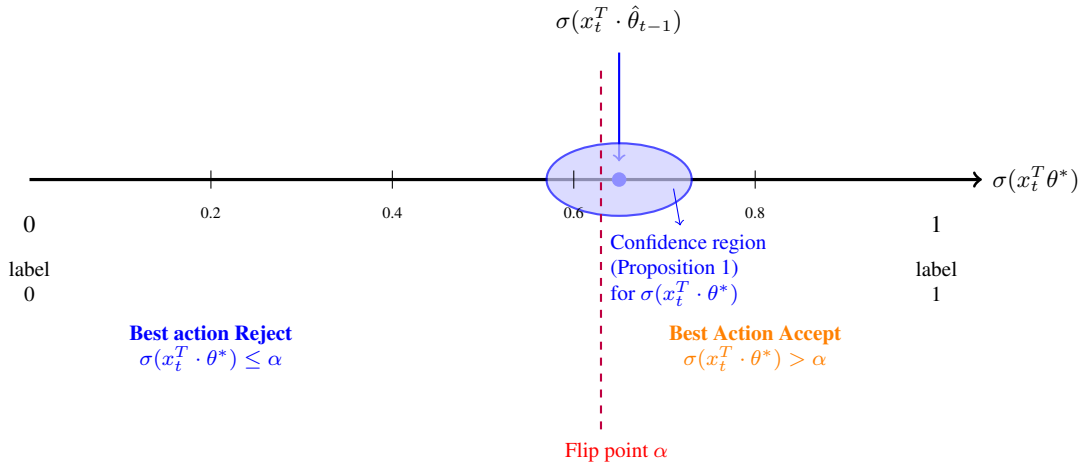
$$a_t = a_t^* = 1.$$

$$F_{12} = 0$$

600 **Case F_{13} :** $\alpha - \beta_{t-1}^{x_t}(\delta_t) < \sigma(x_t^T \hat{\theta}_{t-1}) < \alpha + \beta_{t-1}^{x_t}(\delta_t)$

601 First of all $|\hat{y}_t - \alpha| < \beta_{t-1}^{x_t}(\delta_t)$, hence, $a_t = 1$ (Accept)

602 This is an uncertain region which requires more exploration. However the nature of the confidence
603 region $\beta_{t-1}^{x_t}(\delta_t)$ is such that it shrinks with more exploration (because its dependence on $\|x_t\|_{V_{t-1}^{-1}}^2$
604 which is a monotonically decreasing function). Hence the regret could be controlled using Elliptical
605 potential lemma.



606

607 Since \mathcal{G}_t happens,

$$\begin{aligned}
|\sigma(x_t^T \theta^*) - \sigma(x_t^T \hat{\theta}_{t-1})| &\leq \beta_{t-1}^{x_t}(\delta_t) \\
\sigma(x_t^T \hat{\theta}_{t-1}) - \beta_{t-1}^{x_t}(\delta_t) &\leq \sigma(x_t^T \theta^*) \leq \sigma(x_t^T \hat{\theta}_{t-1}) + \beta_{t-1}^{x_t}(\delta_t) \\
\alpha - 2\beta_{t-1}^{x_t}(\delta_t) &\leq \sigma(x_t^T \theta^*) \leq \alpha + 2\beta_{t-1}^{x_t}(\delta_t) \\
-2\beta_{t-1}^{x_t}(\delta_t) &\leq \sigma(x_t^T \theta^*) - \alpha \leq 2\beta_{t-1}^{x_t}(\delta_t)
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_{t-1} [r_t] &= \mathbb{E}_{t-1} [\psi_{a_t}(y_t) - \psi_{a_t^*}(y_t)] \\
&= \mathbb{E}_{t-1} [\psi_1(y_t) - \psi_{a_t^*}(y_t)] \\
&\leq \left| \mathbb{E}_{t-1} [\psi_1(y_t) - \psi_0(y_t)] \right| \\
&= \left| \mathbb{E}_{t-1} [\ell_{10} + y_t(\ell_{11} - \ell_{10}) - \ell_{01}y] \right| \\
&= \left| \ell_{10} + \mathbb{E}_{t-1} [y_t] (\ell_{11} - \ell_{10}) - \ell_{01} \mathbb{E}_{t-1} [y_t] \right| \\
&= \left| \ell_{10} + \sigma(x_t^T \theta^*) (\ell_{11} - \ell_{10}) - \ell_{01} \sigma(x_t^T \theta^*) \right| \\
&= \left| \left(\sigma(x_t^T \theta^*) - \alpha \right) (\ell_{11} - \ell_{10}) - \ell_{01} \left(\sigma(x_t^T \theta^*) - \alpha \right) \right| \quad (\text{definition of } \alpha) \\
&\leq 2\beta_{t-1}^{x_t}(\delta_t) \max(\ell_{10}, \ell_{01} - \ell_{11})
\end{aligned}$$

$$\begin{aligned}
F_{13} &= \sum_{t=1}^T \mathbb{E} \left[r_t \mathbf{1} \{ \mathcal{G}_t \} \mathbf{1} \left\{ \alpha - \beta_{t-1}^{x_t}(\delta_t) < \sigma(x_t^T \hat{\theta}_{t-1}) < \alpha + \beta_{t-1}^{x_t}(\delta_t) \right\} \right] \\
&\leq \sum_{t=1}^T \mathbb{E} [r_t \mathbf{1} \{ a_t = 1 \}] \\
&= \sum_{t=1}^T \mathbb{E} [\mathbb{E}_{t-1} [r_t] \mathbf{1} \{ a_t = 1 \}] \\
&\leq 2 \max(\ell_{10}, \ell_{01} - \ell_{11}) \mathbb{E} \left[\sum_{t=1}^T \beta_{t-1}^{x_t}(\delta_t) \mathbf{1} \{ a_t = 1 \} \right] \\
&= 2 \max(\ell_{10}, \ell_{01} - \ell_{11}) \mathbb{E} \left[\sum_{t=1}^T \frac{2k_\sigma R_{\max}}{c_\sigma} \|x_t\|_{V_{t-1}^{-1}} \sqrt{\left(3 + 2 \log \left(1 + \frac{2}{\lambda}\right)\right) \cdot 2d \log(N_{t-1}) \cdot \log \left(\frac{d}{\delta_t}\right)} (\delta_t) \mathbf{1} \{ a_t = 1 \} \right] \\
&= 2 \max(\ell_{10}, \ell_{01} - \ell_{11}) \mathbb{E} \left[\sum_{t=1}^T \frac{2k_\sigma R_{\max}}{c_\sigma} \|x_t\|_{V_{t-1}^{-1}} \sqrt{\left(3 + 2 \log \left(1 + \frac{2}{\lambda}\right)\right) \cdot 2d \log(T) \cdot \log(dT)} (\delta_t) \mathbf{1} \{ a_t = 1 \} \right] \\
&= \frac{4 \max(\ell_{10}, \ell_{01} - \ell_{11}) k_\sigma R_{\max}}{c_\sigma} \sqrt{\left(3 + 2 \log \left(1 + \frac{2}{\lambda}\right)\right) \cdot 2d \log(T) \cdot \log(Td)} \mathbb{E} \left[\sum_{t=1}^T \|x_t\|_{V_{t-1}^{-1}} \mathbf{1} \{ a_t = 1 \} \right] \\
&\leq \frac{4\sqrt{T} \max(\ell_{10}, \ell_{01} - \ell_{11}) k_\sigma R_{\max}}{c_\sigma} \sqrt{\left(3 + 2 \log \left(1 + \frac{2}{\lambda}\right)\right) \cdot 2d \log(T) \cdot \log(Td)} \sqrt{\mathbb{E} \left[\sum_{t=1}^T \|x_t\|_{V_{t-1}^{-1}}^2 \mathbf{1} \{ a_t = 1 \} \right]} \\
&\hspace{15em} (\text{Cauchy-Schwartz})
\end{aligned}$$

608 By applying Lemma 8,

$$\begin{aligned}
F_{13} &\leq \left(\sqrt{2d \log\left(1 + \frac{T}{d\lambda}\right)} \right) \frac{4\sqrt{T} \max(\ell_{10}, \ell_{01} - \ell_{11}) k_\sigma R_{\max}}{c_\sigma} \sqrt{\left(3 + 2 \log\left(1 + \frac{2}{\lambda}\right)\right) \cdot 2d \log(T) \cdot \log(Td)} \\
&= O\left(\frac{(1 + e^{C_{\max}})^2}{e^{C_{\max}}} d\sqrt{T} \log^{\frac{3}{2}}(T)\right) \quad \text{(Ignoring logarithmic terms in } d)
\end{aligned}$$

609 Hence,

$$\begin{aligned}
\mathbb{E}[\text{Reg}_T] &\leq \left(\sqrt{2d \log\left(1 + \frac{T}{d\lambda}\right)} \right) \frac{4\sqrt{T} \max(\ell_{10}, \ell_{01} - \ell_{11}) k_\sigma R_{\max}}{c_\sigma} \sqrt{\left(3 + 2 \log\left(1 + \frac{2}{\lambda}\right)\right) \cdot 2d \log(T) \cdot \log(Td)} \\
&\quad + \log T \\
&= O\left(\frac{(1 + e^{C_{\max}})^2}{e^{C_{\max}}} d\sqrt{T} \log^{\frac{3}{2}}(T)\right) \\
&\quad \text{(Ignoring logarithmic terms in } d)
\end{aligned}$$

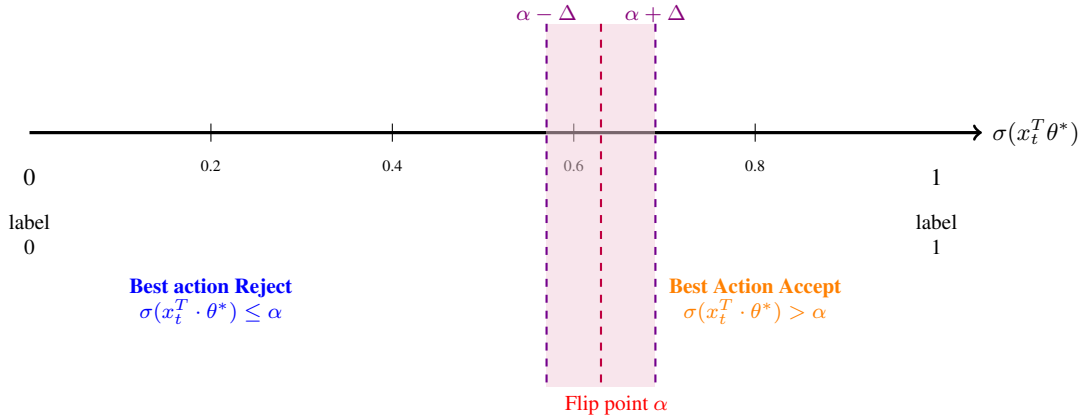
610 C.2 Instance dependent regret analysis

611 **Theorem 2** (LogCBPSide-AT instance dependent regret bound). For $\delta_t = \frac{1}{t}$, with $\forall x \in \mathcal{X} \|x\|_2 \leq 1$
612 and $\sup_{x \in \mathcal{X}} x^T \theta^* \leq C_{\max}$, and $0 < \Delta \leq \min_{x \in \mathcal{X}} |\alpha - \sigma(x^T \cdot \theta^*)|$ the expected cumulative regret
613 of the LogCBPSide-AT algorithm is bounded as:

$$\begin{aligned}
\mathbb{E}[\text{Reg}_T] &\leq O\left(\frac{1}{\Delta} d^2 \log^2 T\right) \\
&\quad \text{(ignoring doubly logarithmic factors and } \text{polylog}\left(\frac{1}{\Delta}\right))
\end{aligned}$$

614 Let define the problem dependent parameter $\Delta > 0$ as,

$$\Delta \leq \min_{x \in \mathcal{X}} |\alpha - \sigma(x^T \cdot \theta^*)|$$



615

616 This defines how close $y_t^* = \sigma(x_t^T \cdot \theta^*)$ is to the decision boundary α . The closer, the problem
617 instance become more difficult owing to the fact that we need more exploration to distinguish it. In
618 the figure above, there are no instances of x_t falls into the shaded region.

619 With this additional assumption, we can refine our previous analysis.

620 F_2, F_{11}, F_{12} all go though for this case as well. We skip them and move to the analysis of F_{13} .

621 **Case F_{13} :** $\alpha - \beta_{t-1}^{x_t}(\delta_t) < \sigma(x_t^T \hat{\theta}_{t-1}) < \alpha + \beta_{t-1}^{x_t}(\delta_t)$

622 First of all $|\hat{y}_t - \alpha| < \beta_{t-1}^{x_t}(\delta_t)$, hence, $a_t = 1$ (Accept)

623 Since \mathcal{G}_t happens,

$$\begin{aligned}
|\sigma(x_t^T \theta^*) - \sigma(x_t^T \hat{\theta}_{t-1})| &\leq \beta_{t-1}^{x_t}(\delta_t) \\
\sigma(x_t^T \hat{\theta}_{t-1}) - \beta_{t-1}^{x_t}(\delta_t) &\leq \sigma(x_t^T \theta^*) \leq \sigma(x_t^T \hat{\theta}_{t-1}) + \beta_{t-1}^{x_t}(\delta_t) \\
\alpha - 2\beta_{t-1}^{x_t}(\delta_t) &\leq \sigma(x_t^T \theta^*) \leq \alpha + 2\beta_{t-1}^{x_t}(\delta_t) \\
-2\beta_{t-1}^{x_t}(\delta_t) &\leq \sigma(x_t^T \theta^*) - \alpha \leq 2\beta_{t-1}^{x_t}(\delta_t).
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_{t-1}[r_t] &= \mathbb{E}_{t-1}[\psi_{a_t}(y_t) - \psi_{a_t^*}(y_t)] \\
&= \mathbb{E}_{t-1}[\psi_1(y_t) - \psi_{a_t^*}(y_t)] \\
&\leq \left| \mathbb{E}_{t-1}[\psi_1(y_t) - \psi_0(y_t)] \right| \\
&= \left| \mathbb{E}_{t-1}[\ell_{10} + y_t(\ell_{11} - \ell_{10}) - \ell_{01}y] \right| \\
&= \left| \ell_{10} + \mathbb{E}_{t-1}[y_t](\ell_{11} - \ell_{10}) - \ell_{01} \mathbb{E}_{t-1}[y_t] \right| \\
&= \left| \ell_{10} + \sigma(x_t^T \theta^*)(\ell_{11} - \ell_{10}) - \ell_{01} \sigma(x_t^T \theta^*) \right| \\
&= \left| \left[(\sigma(x_t^T \theta^*) - \alpha)(\ell_{11} - \ell_{10}) - \ell_{01}(\sigma(x_t^T \theta^*) - \alpha) \right] \right| \quad (\text{definition of } \alpha) \\
&\leq 2\beta_{t-1}^{x_t}(\delta_t) \max(\ell_{10}, \ell_{01} - \ell_{11})
\end{aligned}$$

$$\begin{aligned}
F_{13} &= \sum_{t=1}^T \mathbb{E} \left[r_t \mathbf{1} \{ \mathcal{G}_t \} \mathbf{1} \left\{ \alpha - \beta_{t-1}^{x_t}(\delta_t) < \sigma(x_t^T \hat{\theta}_{t-1}) < \alpha + \beta_{t-1}^{x_t}(\delta_t) \right\} \right] \\
&\leq \sum_{t=1}^T \mathbb{E} \left[r_t \mathbf{1} \{ a_t = 1 \} \mathbf{1} \left\{ -2\beta_{t-1}^{x_t}(\delta_t) \leq \sigma(x_t^T \theta^*) - \alpha \leq 2\beta_{t-1}^{x_t}(\delta_t) \right\} \right] \\
&= \sum_{t=1}^T \mathbb{E} \left[r_t \mathbf{1} \{ a_t = 1 \} \mathbf{1} \left\{ |\sigma(x_t^T \theta^*) - \alpha| \leq 2\beta_{t-1}^{x_t}(\delta_t) \right\} \right] \\
&= \sum_{t=1}^T \mathbb{E} \left[r_t \mathbf{1} \{ a_t = 1 \} \mathbf{1} \left\{ |\sigma(x_t^T \theta^*) - \alpha| \leq 2\beta_{t-1}^{x_t}(\delta_t) \right\} \mathbf{1} \left\{ \beta_{t-1}^{x_t}(\delta_t) < \frac{\Delta}{2} \right\} \right] \\
&\quad + \sum_{t=1}^T \mathbb{E} \left[r_t \mathbf{1} \{ a_t = 1 \} \mathbf{1} \left\{ |\sigma(x_t^T \theta^*) - \alpha| \leq 2\beta_{t-1}^{x_t}(\delta_t) \right\} \mathbf{1} \left\{ \beta_{t-1}^{x_t}(\delta_t) \geq \frac{\Delta}{2} \right\} \right] \\
&= \sum_{t=1}^T \mathbb{E} \left[r_t \mathbf{1} \{ a_t = 1 \} \mathbf{1} \left\{ |\sigma(x_t^T \theta^*) - \alpha| < \Delta \right\} \right] \\
&\quad + \sum_{t=1}^T \mathbb{E} \left[r_t \mathbf{1} \{ a_t = 1 \} \mathbf{1} \left\{ |\sigma(x_t^T \theta^*) - \alpha| \leq 2\beta_{t-1}^{x_t}(\delta_t) \right\} \mathbf{1} \left\{ \beta_{t-1}^{x_t}(\delta_t) \geq \frac{\Delta}{2} \right\} \right] \\
&= 0 \tag{By definition of \Delta} \\
&\quad + \sum_{t=1}^T \mathbb{E} \left[r_t \mathbf{1} \{ a_t = 1 \} \mathbf{1} \left\{ |\sigma(x_t^T \theta^*) - \alpha| \leq 2\beta_{t-1}^{x_t}(\delta_t) \right\} \mathbf{1} \left\{ \beta_{t-1}^{x_t}(\delta_t) \geq \frac{\Delta}{2} \right\} \right] \\
&= \sum_{t=1}^T \mathbb{E} \left[r_t \mathbf{1} \{ a_t = 1 \} \mathbf{1} \left\{ |\sigma(x_t^T \theta^*) - \alpha| \leq 2\beta_{t-1}^{x_t}(\delta_t) \right\} \mathbf{1} \left\{ \beta_{t-1}^{x_t}(\delta_t) \geq \frac{\Delta}{2} \right\} \right] \\
&\leq 2 \max(\ell_{10}, \ell_{01} - \ell_{11}) \sum_{t=1}^T \mathbb{E} \left[\mathbf{1} \{ a_t = 1 \} \beta_{t-1}^{x_t}(\delta_t) \mathbf{1} \left\{ \beta_{t-1}^{x_t}(\delta_t) \geq \frac{\Delta}{2} \right\} \right]
\end{aligned}$$

624 From Proposition 1,

$$\begin{aligned}
\beta_{t-1}^{x_t}(\delta_t) &= \frac{2k_\sigma R_{\max}}{c_\sigma} \|x_t\|_{V_{t-1}^{-1}} \sqrt{\left(3 + 2 \log\left(1 + \frac{2}{\lambda}\right)\right) \cdot 2d \log(N_{t-1}) \cdot \log\left(\frac{d}{\delta_t}\right)} \\
&\leq \frac{2k_\sigma R_{\max}}{c_\sigma} \|x_t\|_{V_{t-1}^{-1}} \sqrt{\left(3 + 2 \log\left(1 + \frac{2}{\lambda}\right)\right) \cdot 2d \log T \cdot \log(dT)} \\
&:= h(T) \cdot \|x_t\|_{V_{t-1}^{-1}} \\
\beta_{t-1}^{x_t}(\delta_t) \geq \frac{\Delta}{2} &\implies \|x_t\|_{V_{t-1}^{-1}} \geq \frac{\Delta}{2h(T)} \\
&\quad \text{(Let } h(T) := \frac{2k_\sigma R_{\max}}{c_\sigma} \sqrt{\left(3 + 2 \log\left(1 + \frac{2}{\lambda}\right)\right) \cdot 2d \log T \cdot \log(dT)} \text{)}
\end{aligned}$$

$$\begin{aligned}
F_{13} &\leq 2 \max(\ell_{10}, \ell_{01} - \ell_{11}) \sum_{t=1}^T \mathbb{E} \left[\mathbf{1} \{a_t = 1\} \beta_{t-1}^{x_t}(\delta_t) \mathbf{1} \left\{ \|x_t\|_{V_{t-1}^{-1}} \geq \frac{\Delta}{2h(T)} \right\} \right] \\
&\leq 2 \max(\ell_{10}, \ell_{01} - \ell_{11}) h(T) \sum_{t=1}^T \mathbb{E} \left[\mathbf{1} \{a_t = 1\} \|x_t\|_{V_{t-1}^{-1}} \mathbf{1} \left\{ \|x_t\|_{V_{t-1}^{-1}} \geq \frac{\Delta}{2h(T)} \right\} \right]
\end{aligned}$$

625 By applying peeling method,

$$\begin{aligned}
F_{13} &\leq 2 \max(\ell_{10}, \ell_{01} - \ell_{11}) h(T) \sum_{t=1}^T \mathbb{E} \left[\mathbf{1} \{a_t = 1\} \|x_t\|_{V_{t-1}^{-1}} \mathbf{1} \left\{ \|x_t\|_{V_{t-1}^{-1}} \geq \frac{\Delta}{2h(T)} \right\} \right] \\
&= 2 \max(\ell_{10}, \ell_{01} - \ell_{11}) h(T) \sum_{t=1}^T \sum_{\ell=1}^{\infty} \mathbb{E} \left[\mathbf{1} \{a_t = 1\} \|x_t\|_{V_{t-1}^{-1}} \mathbf{1} \left\{ 2^{\ell+1} \frac{\Delta}{2h(T)} \geq \|x_t\|_{V_{t-1}^{-1}} \geq \frac{\Delta}{2h(T)} \cdot 2^\ell \right\} \right] \\
&\leq 2 \max(\ell_{10}, \ell_{01} - \ell_{11}) h(T) \sum_{t=1}^T \sum_{\ell=1}^{\infty} \mathbb{E} \left[\mathbf{1} \{a_t = 1\} 2^{\ell+1} \frac{\Delta}{2h(T)} \mathbf{1} \left\{ 2^{\ell+1} \frac{\Delta}{2h(T)} \geq \|x_t\|_{V_{t-1}^{-1}} \geq \frac{\Delta}{2h(T)} \cdot 2^\ell \right\} \right] \\
&\leq \max(\ell_{10}, \ell_{01} - \ell_{11}) \sum_{t=1}^T \sum_{\ell=1}^{\infty} \mathbb{E} \left[2^{\ell+1} \Delta \mathbf{1} \left\{ 2^{\ell+1} \frac{\Delta}{2h(T)} \geq \|x_t\|_{V_{t-1}^{-1}} \geq \frac{\Delta}{2h(T)} \cdot 2^\ell \right\} \right] \\
&= \max(\ell_{10}, \ell_{01} - \ell_{11}) \sum_{\ell=1}^{\infty} 2^{\ell+1} \Delta \sum_{t=1}^T \mathbb{E} \left[\mathbf{1} \left\{ 2^{\ell+1} \frac{\Delta}{2h(T)} \geq \|x_t\|_{V_{t-1}^{-1}} \geq \frac{\Delta}{2h(T)} \cdot 2^\ell \right\} \right] \\
&\leq \max(\ell_{10}, \ell_{01} - \ell_{11}) \sum_{\ell=1}^{\infty} 2^{\ell+1} \Delta \sum_{t=1}^T \mathbb{E} \left[\mathbf{1} \left\{ \|x_t\|_{V_{t-1}^{-1}} \geq \frac{\Delta}{2h(T)} \cdot 2^\ell \right\} \right] \\
&= \max(\ell_{10}, \ell_{01} - \ell_{11}) \sum_{\ell=1}^{\infty} 2^{\ell+1} \Delta \mathbb{E} \left[\sum_{t=1}^T \mathbf{1} \left\{ \|x_t\|_{V_{t-1}^{-1}}^2 \geq \frac{\Delta^2}{4h^2(T)} \cdot 2^{2\ell} \right\} \right] \\
&\leq \max(\ell_{10}, \ell_{01} - \ell_{11}) \sum_{\ell=1}^{\infty} 2^{\ell+1} \Delta \mathbb{E} \left[\frac{12dh^2(T)}{2^{2\ell}\Delta^2} \log \left(1 + \frac{8h^2(T)}{\lambda 2^{2\ell}\Delta^2} \right) \right] \quad (\text{Lemma 9 (EPC)}) \\
&\leq \frac{1}{\Delta} 24dh^2(T) \max(\ell_{10}, \ell_{01} - \ell_{11}) d \log \left(1 + \frac{8h^2(T)}{\lambda 4\Delta^2} \right) \sum_{\ell=1}^{\infty} 2^{-\ell} \\
&\leq \frac{1}{\Delta} 24dh^2(T) \max(\ell_{10}, \ell_{01} - \ell_{11}) d \log \left(1 + \frac{8h^2(T)}{\lambda 4\Delta^2} \right) \\
&\quad \text{(here } h(T) = \frac{2k_\sigma R_{\max}}{c_\sigma} \sqrt{\left(3 + 2 \log \left(1 + \frac{2}{\lambda}\right)\right) \cdot 2d \log T \cdot \log(dT)})
\end{aligned}$$

626 Hence,

$$\begin{aligned}
\mathbb{E} [Reg_T] &\leq \frac{1}{\Delta} 24dh^2(T) \max(\ell_{10}, \ell_{01} - \ell_{11}) d \log \left(1 + \frac{8h^2(T)}{\lambda 4\Delta^2} \right) + O(\log(T)) \\
&\quad \text{(here } h(T) = \frac{2k_\sigma R_{\max}}{c_\sigma} \sqrt{\left(3 + 2 \log \left(1 + \frac{2}{\lambda}\right)\right) \cdot 2d \log T \cdot \log(dT)}) \\
&= O\left(\frac{1}{\Delta} d^2 \log^2 T\right) \quad \text{(ignoring doubly logarithmic factors and polylog}\left(\frac{1}{\Delta}\right)\text{)}
\end{aligned}$$

627 **D Experiment Details**

628 We call our framework AppleTeA: Test-Time Adaptive LLM Cascading via Logistic Apple Tasting,
 629 which is presented in Algorithm D

Algorithm 4 AppleTeA

Input: $\mathcal{M}_s, \mathcal{M}_\ell, n, \phi, \alpha, \gamma, \lambda, \mu$, algorithm Alg
if Alg = “LogCB-AT” **then**
 Instantiate $\mathcal{A} = \text{LogCB-AT}(\gamma, \mu)$
else
 Instantiate $\mathcal{A} = \text{LogCBPSide-AT}(\lambda, \alpha)$
end if
for $t = 1, \dots, n$ **do**
 Observe the query q_t
 Run SLM $r_t = \mathcal{M}_s(q_t)$
 Extract the answer: $\hat{y}_{w,t} = \text{answer}(r_t)$
 Compute the context: $x_t = \phi(q_t, r_t)$
 Run single step of algorithm \mathcal{A} and receive the action recommendation a_t
 if $a_t = 0$ (Reject) **then**
 output $\hat{y}_{w,t}$
 else
 Run the strong LLM: $r_t = \mathcal{M}_\ell(q_t)$
 Extract the answer: $\hat{y}_{s,t} = \text{answer}(r_t)$
 Compute the context: $x_t = \phi(q_t, r_t)$
 Update \mathcal{A} with context x_t and reward = $\mathbb{1}_{\{\hat{y}_{s,t} \neq \hat{y}_{w,t}\}}$
 output $\hat{y}_{s,t}$
 end if
end for

630 We have established that both LogCB-AT and LogCBPSide-AT attain sub-linear regret with respect
 631 to the optimal policy. This regret bound directly translates into a cumulative cost guarantee for
 632 AppleTeA. By carefully constructing the loss matrix to encode two competing cost components
 633 1) the invocation cost of the strong LLM (i.e., token-level inference cost) and 2) the price of error
 634 incurred when the weak SLM is invoked in contexts where prediction failure is likely; AppleTeA
 635 is equipped to navigate this cost-quality tradeoff in a principled manner. Consequently, AppleTeA
 636 achieves optimal performance in expectation, a guarantee that follows directly from the formulation
 637 of the model-selection problem as an Apple Tasting problem and the sub-linear regret guarantees of
 638 the underlying core algorithms.

639 **Context Construction.** To enable effective decision-making, we construct the context vector
 640 $x_t \in \mathbb{R}^d$ by capturing the internal state and the uncertainty score of the SLM model. Specifically, x_t
 641 is formed by concatenating the dimensionally-reduced hidden state vector of the final token from
 642 \mathcal{M}_s and its generative uncertainty score. The uncertainty score can be quantified using either the
 643 maximum probability of the generated tokens or the logit margin between the top two candidate
 644 answers. This context provides a rich signal for the the proposed methods to estimate the probability
 645 of the LLM model altering the final answer.

646 **Datasets and Models.** We adopt Pythia-2.8B [Biderman et al., 2023] and LLaMA2-13B [Touvron
 647 et al., 2023] as \mathcal{M}_s and \mathcal{M}_ℓ , respectively. We evaluate the effectiveness of our algorithms on the
 648 commonsense reasoning benchmark: ARC-Easy [Clark et al., 2018], OpenBookQA [Mihaylov et al.,
 649 2018] and BoolQ [Sakaguchi et al., 2021]

650 **NeurIPS Paper Checklist**

651 **1. Claims**

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

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

655 Justification: The abstract and introduction clearly state the contributions, scope, and
656 assumptions of the paper.

657 Guidelines:

- 658 • The answer [\[N/A\]](#) means that the abstract and introduction do not include the claims
659 made in the paper.
- 660 • The abstract and/or introduction should clearly state the claims made, including the
661 contributions made in the paper and important assumptions and limitations. A [\[No\]](#) or
662 [\[N/A\]](#) answer to this question will not be perceived well by the reviewers.
- 663 • The claims made should match theoretical and experimental results, and reflect how
664 much the results can be expected to generalize to other settings.
- 665 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
666 are not attained by the paper.

667 **2. Limitations**

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

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

670 Justification: We address the limitations of our approach in the Conclusion.

671 Guidelines:

- 672 • The answer [\[N/A\]](#) means that the paper has no limitation while the answer [\[No\]](#) means
673 that the paper has limitations, but those are not discussed in the paper.
- 674 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 675 • The paper should point out any strong assumptions and how robust the results are to
676 violations of these assumptions (e.g., independence assumptions, noiseless settings,
677 model well-specification, asymptotic approximations only holding locally). The authors
678 should reflect on how these assumptions might be violated in practice and what the
679 implications would be.
- 680 • The authors should reflect on the scope of the claims made, e.g., if the approach was
681 only tested on a few datasets or with a few runs. In general, empirical results often
682 depend on implicit assumptions, which should be articulated.
- 683 • The authors should reflect on the factors that influence the performance of the approach.
684 For example, a facial recognition algorithm may perform poorly when image resolution
685 is low or images are taken in low lighting. Or a speech-to-text system might not be
686 used reliably to provide closed captions for online lectures because it fails to handle
687 technical jargon.
- 688 • The authors should discuss the computational efficiency of the proposed algorithms
689 and how they scale with dataset size.
- 690 • If applicable, the authors should discuss possible limitations of their approach to
691 address problems of privacy and fairness.
- 692 • While the authors might fear that complete honesty about limitations might be used by
693 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
694 limitations that aren’t acknowledged in the paper. The authors should use their best
695 judgment and recognize that individual actions in favor of transparency play an impor-
696 tant role in developing norms that preserve the integrity of the community. Reviewers
697 will be specifically instructed to not penalize honesty concerning limitations.

698 **3. Theory assumptions and proofs**

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

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

702 Justification: All assumptions are explicitly stated in the main text and complete proofs are
703 provided in the Appendix.

704 Guidelines:

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

715 4. Experimental result reproducibility

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

719 Answer: [Yes]

720 Justification: All experimental details are in Experiments section.

721 Guidelines:

- 722 • The answer [N/A] means that the paper does not include experiments.
- 723 • If the paper includes experiments, a [No] answer to this question will not be perceived
724 well by the reviewers: Making the paper reproducible is important, regardless of
725 whether the code and data are provided or not.
- 726 • If the contribution is a dataset and/or model, the authors should describe the steps taken
727 to make their results reproducible or verifiable.
- 728 • Depending on the contribution, reproducibility can be accomplished in various ways.
729 For example, if the contribution is a novel architecture, describing the architecture fully
730 might suffice, or if the contribution is a specific model and empirical evaluation, it may
731 be necessary to either make it possible for others to replicate the model with the same
732 dataset, or provide access to the model. In general, releasing code and data is often
733 one good way to accomplish this, but reproducibility can also be provided via detailed
734 instructions for how to replicate the results, access to a hosted model (e.g., in the case
735 of a large language model), releasing of a model checkpoint, or other means that are
736 appropriate to the research performed.
- 737 • While NeurIPS does not require releasing code, the conference does require all submis-
738 sions to provide some reasonable avenue for reproducibility, which may depend on the
739 nature of the contribution. For example
 - 740 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
741 to reproduce that algorithm.
 - 742 (b) If the contribution is primarily a new model architecture, the paper should describe
743 the architecture clearly and fully.
 - 744 (c) If the contribution is a new model (e.g., a large language model), then there should
745 either be a way to access this model for reproducing the results or a way to reproduce
746 the model (e.g., with an open-source dataset or instructions for how to construct
747 the dataset).
 - 748 (d) We recognize that reproducibility may be tricky in some cases, in which case
749 authors are welcome to describe the particular way they provide for reproducibility.
750 In the case of closed-source models, it may be that access to the model is limited in
751 some way (e.g., to registered users), but it should be possible for other researchers
752 to have some path to reproducing or verifying the results.

753 5. Open access to data and code

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

757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807

Answer: [Yes]

Justification: The source code and datasets are provided in the supplementary material.

Guidelines:

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

6. Experimental setting/details

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

Answer: [Yes]

Justification: The core experimental settings are described in Experiments section.

Guidelines:

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

7. Experiment statistical significance

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

Answer: [Yes]

Justification: They are included in the experimental results.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- 808
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - 810
 - 811
 - 812
 - 813
 - 814
 - 815
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).
 - If error bars are reported in tables or plots, the authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

816 8. Experiments compute resources

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

820 Answer: [Yes]

821 Justification: These details are provided in the Experiments section.

822 Guidelines:

- 823 • The answer [N/A] means that the paper does not include experiments.
- 824 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- 825
- 826 • The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- 827
- 828 • The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 829
- 830

831 9. Code of ethics

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

834 Answer: [Yes]

835 Justification: The research conducted in this paper strictly conforms to the NeurIPS Code of Ethics.

837 Guidelines:

- 838 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- 839
- 840 • If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
- 841
- 842 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 843

844 10. Broader impacts

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

847 Answer: [N/A]

848 Justification: This work focuses on foundational theoretical algorithms and has no direct societal impact.

850 Guidelines:

- 851 • The answer [N/A] means that there is no societal impact of the work performed.
- 852 • If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
- 853
- 854 • Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- 855
- 856
- 857

- 858
- 859
- 860
- 861
- 862
- 863
- 864
- 865
- 866
- 867
- 868
- 869
- 870
- 871
- 872
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

873 **11. Safeguards**

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

877 Answer: [N/A]

878 Justification: This paper focuses on theoretical algorithms and does not release models or
879 datasets with a high risk of misuse.

880 Guidelines:

- 881
- 882
- 883
- 884
- 885
- 886
- 887
- 888
- 889
- 890
- The answer [N/A] means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

891 **12. Licenses for existing assets**

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

895 Answer: [Yes]

896 Justification: We appropriately cite all existing datasets and code used in our experiments
897 and adhere to their respective licenses.

898 Guidelines:

- 899
- 900
- 901
- 902
- 903
- 904
- 905
- 906
- 907
- 908
- 909
- The answer [N/A] means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- 910
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- 911
- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.
- 912
- 913

914 **13. New assets**

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

916 Answer: [Yes]

917 Justification: The code for the our proposed algorithms are provided in the supplementary material.

918 Guidelines:

- 919
- The answer [N/A] means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 920
- 921
- 922
- 923
- 924
- 925
- 926
- 927
- 928

929 **14. Crowdsourcing and research with human subjects**

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

931 Answer: [N/A]

932 Justification: This research does not involve crowdsourcing or human subjects.

933 Guidelines:

- 934
- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 935
- 936
- 937
- 938
- 939
- 940
- 941
- 942
- 943

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

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

946 Answer: [N/A]

947 Justification: This study does not involve human subjects or IRB-related research.

948 Guidelines:

- 949
- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
 - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- 950
- 951
- 952
- 953
- 954
- 955
- 956
- 957
- 958
- 959
- 960

- 961 • For initial submissions, do not include any information that would break anonymity (if
962 applicable), such as the institution conducting the review.

963 **16. Declaration of LLM usage**

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

968 Answer: [N/A]

969 Justification: LLMs were not used as a core component of the methodology in this research.

970 Guidelines:

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