

448 **A Proofs from Section 2**

---

**Algorithm 4:** Output  $\hat{\alpha} \in [G^{-1}(1 - \eta_1) - \frac{\varepsilon}{3}, G^{-1}(1 - \eta_1 + \eta_2) + \frac{\varepsilon}{3}]$  with probability  $1 - \frac{\delta}{2}$

---

- 1 **input:** arm set  $\mathcal{S} = (a_1, a_2, \dots)$  and parameters  $(\eta_1, \eta_2, \varepsilon, \delta) \in (0, 1)$  with  $\eta_2 < \eta_1$ .
  - 2 initialize:  $K = \frac{C\eta_1 \log(1/\delta)}{\eta_2^2}$ .
  - 3 **for**  $i = 1, 2, \dots, K$  **do**
  - 4     Collect  $n = \frac{C \log(1/\eta_2)}{\varepsilon^2}$  samples of arm  $i$ . Set  $\hat{p}_i = \hat{p}_i(n)$  to be the average observed reward.
  - 5 **end**
  - 6 Let  $\hat{\alpha}$  be the  $k$ -th largest value in  $\{\hat{p}_1, \dots, \hat{p}_K\}$  for  $k = \lceil K(\eta_1 - \frac{\eta_2}{2}) \rceil$ .
  - 7 Return  $\hat{\alpha}$
- 

449 We show the following generalization of Proposition 2.1.

450 **Proposition A.1.** Fix  $0 \leq \eta_1, \eta_2, \varepsilon, \delta \leq 1$  with  $\eta_2 \leq \eta_1$ . With probability at least  $1 - \frac{\delta}{2}$ , the output  
 451  $\hat{\alpha}$  of Alg. 4 satisfies

$$\hat{\alpha} \in \left[ G^{-1}(1 - \eta_1) - \frac{\varepsilon}{3}, G^{-1}(1 - \eta_1 + \eta_2) + \frac{\varepsilon}{3} \right].$$

452 Moreover, Alg. 4 has sample complexity

$$O\left(\frac{\eta_1 \log(1/\eta_2) \log(1/\delta)}{\eta_2^2 \varepsilon^2}\right).$$

453 *Proof.* The sample complexity is clear so we focus on the first statement. First observe that by a  
 454 Chernoff estimate, for each  $i \in [K]$ ,

$$\mathbb{P}\left[|p_i - \hat{p}_i| \geq \frac{\varepsilon}{3}\right] \leq \frac{\eta_2}{8}. \quad (\text{A.1})$$

455 Let  $N(\varepsilon)$  be the number of  $i \in [K]$  such that  $|p_i - \hat{p}_i| \geq \frac{\varepsilon}{3}$ . Applying a second Chernoff estimate  
 456 (of multiplicative form, see e.g. Theorem 4.5 in [MU17]) on these events as  $i$  varies and noting that  
 457  $K\eta_2 \geq C \log(1/\delta)$ , (A.1) implies

$$\mathbb{P}\left[N(\varepsilon) \leq \frac{K\eta_2}{6}\right] \geq 1 - \frac{\delta}{8}. \quad (\text{A.2})$$

458 We next show that with probability at least  $1 - \frac{\delta}{4}$ ,

$$\hat{\alpha} \leq \bar{\alpha} + \frac{\varepsilon}{3} \equiv G^{-1}(1 - \eta_1 + \eta_2) + \frac{\varepsilon}{3}. \quad (\text{A.3})$$

459 With  $p_i$  the (true) mean reward from arm  $a_i$ , let

$$N_{\bar{\alpha}} \equiv |\{i \in [K] : p_i > \bar{\alpha}\}|$$

460 denote the number of the  $K$  tested arms which satisfy  $p_i > \bar{\alpha}$ . By definition,  $N_{\bar{\alpha}}$  is stochastically  
 461 dominated by a Bin  $(K, \eta_1 - \frac{9\eta_2}{10})$  random variable, and  $\eta_1 - \frac{3\eta_2}{4} = \Theta(\eta_1)$  since  $\eta_2 \leq \eta_1$ . Note  
 462 that

$$\begin{aligned} \eta_1 - \frac{9\eta_2}{10} &\asymp \eta_1 - \frac{3\eta_2}{4} \asymp \eta_1, \\ \frac{\eta_1 - \frac{9\eta_2}{10}}{\eta_1 - \frac{3\eta_2}{4}} &\geq 1 + \frac{\eta_2}{20\eta_1}. \end{aligned}$$

463 Therefore another multiplicative Chernoff estimate implies

$$\mathbb{P}\left[N_{\bar{\alpha}} \leq K\left(\eta_1 - \frac{3\eta_2}{4}\right)\right] \geq e^{-\Omega(K\eta_2^2/\eta_1)} \geq 1 - \frac{\delta}{8}.$$

464 When both  $N(\varepsilon) \leq \frac{K\eta_2}{6}$  and  $N_{\bar{\alpha}} \leq K\left(\eta_1 - \frac{3\eta_2}{4}\right)$  hold, it follows by definition that  $\hat{\alpha} \leq \bar{\alpha} + \frac{\varepsilon}{3}$ .  
 465 Hence recalling (A.2) above, we conclude that

$$\mathbb{P}\left[\hat{\alpha} \leq \bar{\alpha} + \frac{\varepsilon}{3}\right] \geq 1 - \frac{\delta}{4},$$

466 establishing (A.3). The other direction is similar. With  $\alpha = G^{-1}(1 - \eta_1)$  as usual, we set

$$N_\alpha \equiv |\{i \in [K] : p_i \geq \alpha\}|. \quad (\text{A.4})$$

467 This time,  $N_\alpha$  stochastically dominates a  $\text{Bin}(K, \eta_1)$  random variable. Yet another Chernoff esti-  
468 mate yields

$$\mathbb{P} \left[ N_\alpha \geq K \left( \eta_1 - \frac{\eta_2}{4} \right) \right] \geq 1 - \frac{\delta}{8}.$$

469 Using (A.2) in the same way as above, we find

$$\mathbb{P} \left[ \hat{\alpha} \geq \alpha - \frac{\varepsilon}{3} \right] \geq 1 - \frac{\delta}{4}.$$

470 This concludes the proof.  $\square$

471 *Proof of Theorem 2.1.* First we analyze the expected sample complexity. On the event that

$$\hat{\alpha} \in \left[ G^{-1}(1 - \eta) - \frac{\varepsilon}{3}, G^{-1} \left( 1 - \frac{\eta}{2} \right) + \frac{\varepsilon}{3} \right] \quad (\text{A.5})$$

472 we claim that Alg. 2 terminates with probability  $\eta/4$  for each  $a_i$ . Indeed, if

$$\hat{p}_i \geq G^{-1} \left( 1 - \frac{\eta}{2} \right)$$

473 then termination always happens by definition. This has probability at least  $1/4$  if  $p_i \geq G^{-1} \left( 1 - \frac{\eta}{2} \right)$   
474 by Theorem 1 in [GM14], and the latter condition has probability at least  $\eta/2$  by definition. It  
475 follows that when (A.5) holds, the expected sample complexity of Alg. 2 is  $O \left( \frac{\log(1/\eta\delta)}{\eta\varepsilon^2} \right)$ . On  
476 the other hand, (A.5) fails to hold with probability less than  $\delta$ . Because of the explicit termina-  
477 tion condition in Alg. 2, this yields a additional sample complexity contribution of smaller order  
478  $O \left( \delta \log(1/\delta) \frac{\log(1/\eta\delta)}{\eta\varepsilon^2} \right)$ . Finally Alg. 4 has sample complexity

$$O \left( \frac{\log(1/\eta) \log(1/\delta)}{\eta\varepsilon^2} \right)$$

479 which clearly forms the dominant contribution. This completes the proof of the sample complexity  
480 bound and we now turn to proving correctness with probability  $1 - \delta$ . First, it is easy to see that  
481 Alg. 4 outputs some arm  $a_i$  with probability at least  $1 - \frac{\delta}{2}$ . It therefore suffices to show that for  
482 any fixed  $\hat{\alpha}$  satisfying (A.5), conditioned on the event  $\hat{p}_i \geq \hat{\alpha} - \frac{\varepsilon}{3}$ , the conditional probability that  
483  $p_i \geq \alpha - \varepsilon$  is at least  $1 - \frac{\delta}{2}$ .

484 We do this using Bayes' rule. If  $p_i \geq G^{-1}(1 - \frac{\eta}{2})$ , then as above Theorem 1 in [GM14] implies

$$\mathbb{P} \left[ \hat{p}_i \geq \hat{\alpha} - \frac{\varepsilon}{3} \right] \geq \mathbb{P}[\hat{p}_i \geq p_i] \geq 1/4.$$

485 This event hence contributes probability at least  $\eta/4$  to the event  $p_i \geq G^{-1}(1 - \eta)$ . On the other  
486 hand, if  $p_i \leq G^{-1}(1 - \eta) - \varepsilon \leq \hat{\alpha} - \frac{2\varepsilon}{3}$ , then

$$\mathbb{P} \left[ \hat{p}_i \geq \hat{\alpha} - \frac{\varepsilon}{3} \right] \leq \mathbb{P} \left[ \hat{p}_i \geq p_i + \frac{\varepsilon}{3} \right] \leq \eta\delta/8$$

487 for an absolute constant  $C$ . Combining these via Bayes' rule implies the desired result.  $\square$

## 488 B Lower Bound for Fixed Budget

### 489 Fixed Budget with Unknown $\alpha$

490 Before giving the proof, we give some qualitative discussion of the role of unknown  $\alpha$ . We consider  
491 Theorem 3.2 to be a definitive lower bound, since e.g. being given the value of  $\alpha$  only makes the  
492 result stronger. When  $\alpha$  is unknown, it is possible to give an essentially matching algorithm, but  
493 more care is required when stating the result. This is inherent and stems from the fact that the value  
494  $\alpha = G_\mu^{-1}(1 - \eta)$  can be difficult or even impossible to estimate, yet determines the constant  $c_{\alpha,\beta}$  in  
495 the desired rate.

496 Let us illustrate the issue by a counterexample. Consider  $\mu_N$  defined by:

$$\begin{aligned}\mathbb{P}^{p \sim \mu_N}[p = 0.4] &= \frac{1}{2} + e^{-10N}, \\ \mathbb{P}^{p \sim \mu_N}[p = 0.2] &= \frac{1}{2} - e^{-10N}.\end{aligned}\tag{B.1}$$

497 Similarly define  $\tilde{\mu}_N$  by:

$$\begin{aligned}\mathbb{P}^{p \sim \tilde{\mu}_N}[p = 0.4] &= \frac{1}{2} - e^{-10N}, \\ \mathbb{P}^{p \sim \tilde{\mu}_N}[p = 0.3] &= 2e^{-10N}, \\ \mathbb{P}^{p \sim \tilde{\mu}_N}[p = 0.2] &= \frac{1}{2} - e^{-10N}.\end{aligned}\tag{B.2}$$

498 Then  $\mu_N$  and  $\tilde{\mu}_N$  are not distinguishable using  $N$  samples, yet  $G_{\mu}^{-1}(1/2) = 0.4$  while  $G_{\tilde{\mu}}^{-1}(1/2) =$   
499  $0.3$ . Using non-distinguishability it follows that the lower bound of Theorem 3.2 applies to  $\tilde{\mu}_N$  with  
500 threshold  $\alpha = G_{\mu_N}^{-1}(1/2) = 0.4$ , as opposed to the direct application using  $G_{\tilde{\mu}_N}^{-1}(1/2) = 0.3$ . It is  
501 not hard to show using monotonicity of  $\frac{1}{\sqrt{x(1-x)}}$  that

$$c_{0.4,0.4-\varepsilon} < c_{0.3,0.3-\varepsilon}$$

502 for all  $\varepsilon \leq 0.3$ . As a result, it is information-theoretically **impossible** to achieve the rate (3.1) for  
503  $\tilde{\mu}_N$  if the target quantile value  $\alpha$  is not given. The core reason is that the value  $G_{\tilde{\mu}}^{-1}(1/2) = 0.3$  is  
504 too sensitive to the choice  $\eta = 1/2$  of quantile.

505 Fortunately, this issue is more of an annoyance than a real difficulty. It can be fixed in several ways.  
506 In Theorems B.1, B.2, and B.3 below we give three concrete formulations under which the guarantee  
507 (3.1) can be achieved, as mentioned in the main body.

508 **Theorem B.1.** *For fixed  $\eta_1, \eta_2, \varepsilon$ , there is a sequence  $(\mathcal{A}_N)_{N \geq 1}$  of  $N$ -sample algorithms outputting*  
509  *$a_{i^*}$  such that the following holds for any sequence  $(\mu_N)_{N \geq 1}$  of reservoir distributions. Letting*

$$\alpha_N = \frac{1}{\eta_1 - \eta_2} \cdot \int_{1-\eta_1}^{1-\eta_2} G_{\mu_N}^{-1}(x) dx$$

510 *be a quantile average of  $\mu_N$ , we have*

$$\limsup_{N \rightarrow \infty} \frac{(-\log \mathbb{P}[p_{i^*} < \alpha_N - \varepsilon]) \cdot \log^2 N}{c_{\alpha_N, \alpha_N - \varepsilon} N} \geq 1.\tag{B.3}$$

511 **Theorem B.2.** *For fixed  $\eta, \varepsilon$ , there is a sequence  $(\mathcal{A}_N)_{N \geq 1}$  of  $N$ -sample algorithms outputting  $a_{i^*}$*   
512 *such that for any sequence of reservoir distributions  $\mu_N$  satisfying*

$$\alpha_N \equiv G_{\mu_N}^{-1}(1 - \eta) \geq \frac{1 + \varepsilon}{2},$$

513 *we have*

$$\limsup_{N \rightarrow \infty} \frac{(-\log \mathbb{P}[p_{i^*} < G_{\mu_N}^{-1}(1 - \eta) - \varepsilon]) \cdot \log^2 N}{c_{\alpha_N, \alpha_N - \varepsilon} N} \geq 1.\tag{B.4}$$

514 **Theorem B.3.** *For any fixed  $\varepsilon_1 > \varepsilon$ , there is a sequence  $(\mathcal{A}_N)_{N \geq 1}$  of  $N$ -sample algorithms out-*  
515 *putting  $a_{i^*}$  such that for any fixed reservoir distribution  $\mu$  with  $\mu^* > \varepsilon$ ,*

$$\limsup_{N \rightarrow \infty} \frac{(-\log \mathbb{P}[p_{i^*} < \mu^* - \varepsilon_1]) \cdot \log^2 N}{N} \geq c_{\mu^*, \mu^* - \varepsilon}.\tag{B.5}$$

516 We emphasize that the rate (3.1) is optimal in all cases since the lower bound of Theorem 3.2 is  
517 for an easier problem. The first formulation above may be the most principled choice. The idea  
518 is that an averaged quantile depends continuously on  $\mu$ , and can in fact be estimated by applying  
519 Proposition A.1 for several pairs  $(\eta_1, \eta_2)$  and computing a Riemann sum. The second formulation  
520 requires only the mild condition that  $\alpha \geq \frac{1+\varepsilon}{2}$  and uses monotonicity of  $c_{\alpha, \alpha - \varepsilon}$  on this set. (In other  
521 words, if the average reward values  $p$  appearing in (B.1), (B.2) were larger than 0.5, there would be  
522 no counterexample.) The third formulation allows us to almost send  $\eta$  all the way down to 0. It uses  
523 the fact that

$$\mu^* - (\varepsilon_1 - \varepsilon) \leq G_{\mu}^{-1}(1 - \eta')$$

524 for some  $\eta' = \eta'(\mu, \varepsilon_1, \varepsilon) > 0$ . These results show that (3.1) is achievable even without knowledge  
525 of  $\alpha$ , up to a choice of technical modification to sidestep the counterexample discussed above.

526 **Remark B.1.** *In fact uniformity in  $(\alpha, \beta)$  holds in the following sense. For any sequence*  
 527  *$(\alpha_N, \beta_N)_{N \geq 1}$  of pairs with  $\min(\beta_N, \alpha_N - \beta_N, 1 - \alpha_N)$  uniformly bounded below, there is a*  
 528 *sequence  $(\mathcal{A}_N)_{N \geq 1}$  of  $N$ -sample algorithms such that for any  $\eta \in (0, 1)$  and any sequence of*  
 529 *reservoir distributions  $\mu_N$  with  $G_{\mu_N}^{-1}(1 - \eta) \geq \alpha_N$ ,*

$$\limsup_{N \rightarrow \infty} \frac{(-\log \mathbb{P}[p_{i^*} < \beta_N]) \cdot \log^2 N}{c_{\alpha_N, \beta_N} N} \leq 1. \quad (\text{B.6})$$

530 *This can be shown identically to Theorem 3.1, though we don't give the proof in this generality. It is*  
 531 *useful for the reduction arguments in Theorems B.1, B.2, and B.3.*

### 532 B.1 Preparation for the Proof

533 Here we prove Theorem 3.2. For any  $\alpha, \beta, \eta, \varrho > 0$  we construct a reservoir  $\mu = \mu_{\alpha, \beta, \eta, \varrho}$  such that

$$\liminf_{N \rightarrow \infty} \frac{(-\log \mathbb{P}^\mu[p_{i^*} < \beta]) \cdot \log^2 N}{N} \leq c_{\alpha, \beta} + \lambda(\varrho) \quad (\text{B.7})$$

534 holds for any sequence of  $N$ -sample algorithms  $\mathcal{A}_N$ , and where  $\lim_{\varrho \rightarrow 0} \lambda(\varrho) = 0$  for fixed  $\alpha, \beta, \eta$ .

### 535 B.2 Admissible Reservoirs and Bayesian Perspective

536 In proving Theorem 3.2, we will use reservoir distributions  $\mu$  of a specific form. Namely, we require  
 537 each  $\mu$  to be supported on an interval  $[\underline{\gamma}, \bar{\gamma}]$ , where

$$0 < \beta - \varrho < \underline{\gamma} < \beta < \alpha < \bar{\gamma} < \alpha + \varrho < 1.$$

538 In fact we define  $\underline{\gamma}, \bar{\gamma}$  explicitly (recall that  $\varrho > 0$  is a small constant which we eventually send to 0)  
 539 by

$$\begin{aligned} \theta(\underline{\gamma}) &= \theta(\beta) - \varrho^2; \\ \theta(\bar{\gamma}) &= \theta(\alpha) + \varrho^2. \end{aligned} \quad (\text{B.8})$$

540 We say  $\mu$  is  $(\underline{\gamma}, \bar{\gamma}, \underline{f}, \bar{f})$  *admissible* if  $\mu$  has density  $\mu(dx) = f(x)dx$  for a Borel measurable function  
 541  $f$  and satisfies for constants  $0 < \underline{f} < \bar{f} < \infty$ ,

$$f(x) \in [\underline{f}, \bar{f}], \quad \forall x \in [\underline{\gamma}, \bar{\gamma}].$$

542 Towards proving Theorem 3.2, we fix throughout this section some  $(\underline{\gamma}, \bar{\gamma}, \underline{f}, \bar{f})$  admissible  $\hat{\mu}$  such  
 543 that  $G_{\hat{\mu}}^{-1}(\alpha) = \eta$  holds, for appropriate constants  $(\underline{f}, \bar{f})$  depending only on  $(\eta, \varepsilon, \alpha, \beta, \underline{\gamma}, \bar{\gamma})$ . It is  
 544 easy to see that this is always possible.

545 An admissible  $\mu$  is roughly comparable to the uniform distribution on an interval. Using admissible  
 546 reservoirs gives each  $a_i$  the potential to slowly degrade in observed quality over time. We remark  
 547 that while it is more convenient to work with reservoirs supported away from the boundaries, i.e. in  
 548  $[\underline{\gamma}, \bar{\gamma}] \subseteq (0, 1)$ , we do not expect this to be essential.

549 It will be helpful throughout this section to take a Bayesian point of view. We treat  $\mu_N$  as known  
 550 to  $\mathcal{A}_N$ , since  $\mathcal{A}_N$  is in fact allowed to depend on  $\mu_N$ . Thus at each time  $t$ , each  $p_i$  has a posterior  
 551 probability distribution which we denote by  $\mu_{i,t}$ . Note that each  $\mu_{i,t}$  depends only on  $(n_{i,t}, \hat{p}_{i,t})$  and  
 552 is initialized at  $\mu_{i,0} = \mu$ . We denote by

$$\boldsymbol{\mu}^t = (\mu_{1,t}, \mu_{2,t}, \dots) \quad (\text{B.9})$$

553 the sequence of posterior distributions  $\mu_{i,t}$ . Since arms are independent,  $\boldsymbol{\mu}^t$  is the full time- $t$  poste-  
 554 rior of the algorithm.

### 555 B.3 Batched Algorithms and Adversaries

556 In pure exploration problems, it is possible to significantly simplify the structure of any algorithm  
 557 at the cost of a small multiplicative increase in the sample complexity. We carry this out using the  
 558 notion of a batch-compressed algorithm.

559 **Definition B.1.** Given an increasing sequence  $B = (b_1, b_2, \dots)$  of positive integers, an algorithm  
 560  $\mathcal{A}$  is  **$B$ -batch-compressed** if  $\mathcal{A}$  can only act by increasing the number of times  $n_i$  that  $a_i$  has been  
 561 sampled from  $b_k$  to  $b_{k+1}$ , so that  $n_i \in B$  holds at all times.  $B$  is  **$\varrho$ -slowly increasing** if

$$\frac{b_{k+1}}{b_k + 1} \leq 1 + \varrho, \quad \forall k \geq 1.$$

562 Finally if  $\mathcal{A}$  is  $B$ -batch-compressed and  $B$  is  $\varrho$ -slowly increasing, we say that  $\mathcal{A}$  is  $\varrho$ -batch-  
 563 compressed.

564 Unlike the batched algorithms studied in [PRCS16, GHRZ19], batch-compression is only important  
 565 for us as an analysis technique. Indeed the following proposition shows that it does not fundamen-  
 566 tally affect pure exploration algorithms.

567 **Proposition B.2.** If  $B$  is  $\varrho$ -slowly increasing, then for any  $N$ -sample algorithm  $\mathcal{A}$ , there exists an  
 568  $B$ -batch-compressed  $\lfloor N(1 + \varrho) \rfloor$  algorithm  $\mathcal{A}'$  with the same output.

569 *Proof.* We show how to simulate  $\mathcal{A}$  using the  $B$ -batch-compressed  $\mathcal{A}'$ , assuming that the sequence  
 570 of rewards for each  $a_i$  is fixed. Each time  $\mathcal{A}$  samples arm  $i$  for the  $n_i = (a_k + 1)$ -st time for  
 571  $a_k \in A$ ,  $\mathcal{A}'$  samples arm  $i$  until  $n_i = a_{k+1}$ . Then  $\mathcal{A}'$  has all the information of  $\mathcal{A}$  at all times, hence  
 572 can simulate the behavior and output of  $\mathcal{A}$ . Moreover by the definition of  $\varrho$ -slowly increasing, the  
 573 sample complexity of  $\mathcal{A}'$  is larger than that of  $\mathcal{A}$  by at most a factor  $(1 + \varrho)$ .  $\square$

574 We will use the above with  $\varrho \rightarrow 0$  slowly as  $N \rightarrow \infty$ . Then the sample complexity increase  $1 + \varrho$   
 575 is absorbed into the  $1 + o(1)$  factor in Theorem 3.2. As a result it suffices to establish (B.7) under  
 576 the additional assumption that  $\mathcal{A}_N$  is  $\varrho$ -batch-compressed.

#### 577 B.4 Fisher Information Distance

578 Determining the tight constant  $c_{\alpha, \beta}$  requires significant care. In particular the adversary must de-  
 579 crease the empirical average rewards  $\hat{p}_{i,t}$  at a precise rate depending on  $n_{i,t}$ . This rate turns out to  
 580 involve the *Fisher information distance*. For  $a, b \in [0, 1]$  we define the Fisher information distance  
 581  $d_F(a, b)$  between  $a$  and  $b$  to be

$$d_F(a, b) = \left| \int_a^b \frac{dx}{\sqrt{x(1-x)}} \right|.$$

582 This agrees with the more general Fisher information metric when each  $a \in [0, 1]$  is identified with  
 583 the corresponding Bernoulli distribution. We refer the reader to [Nie20] for a survey on informa-  
 584 tion geometry. In short, the Fisher information yields a natural Riemannian metric on families of  
 585 probability distributions which are parametrized by smooth manifolds. However we will use only  
 586 elementary properties of  $d_F$ .

587 We parametrize  $[0, 1]$  using the function  $\theta : [0, 1] \rightarrow [0, \pi]$  defined by

$$\theta(a) = d_F(0, a) = \int_0^a \frac{dx}{\sqrt{x(1-x)}} = \arccos(1 - 2a). \quad (\text{B.10})$$

588 In particular,

$$d_F(a, b) = |\arccos(1 - 2a) - \arccos(1 - 2b)| \geq 2|a - b|$$

589 and so  $d_F(0, 1) = \pi$ . The main property of  $\theta$  that we will use is the resulting differential equation

$$\theta'(a) = \frac{1}{\sqrt{\theta(a)(1 - \theta(a))}}. \quad (\text{B.11})$$

590 In our case,  $\theta^{-1}$  parametrizes a ‘‘constant speed’’ path through the space of Bernoulli variables, view-  
 591 ing the Fisher information. Correspondingly, our adversary will ensure that  $\theta(\hat{p}_i(n_{i,t}))$  decreases  
 592 linearly in  $\log(n_{i,t})$ .

593 **B.5 Preliminary Lemmas from Moderate Deviations**

594 Recall that for positive integers  $a$  and  $b$ , the  $\text{Beta}(a, b)$  distribution has probability density function

$$\frac{(a+b-1)!}{(a-1)!(b-1)!} x^{a-1} (1-x)^{b-1}$$

595 for  $x \in [0, 1]$ . We now recall a moderate deviations principle for the binomial distribution and a  
596 central limit theorem for the beta distribution.

597 **Lemma 4** (Theorem 2.2 in [DA92]). *For any  $0 < \underline{q} < \bar{q} < 1$  and constant  $\varrho > 0$  there exists*  
598  $\Delta_0(\underline{q}, \bar{q}, \varrho)$  *and  $M_0(\underline{q}, \bar{q}, \varrho)$  such that the following holds for all  $p \in [\underline{q}, \bar{q}]$ . For  $n \geq n_0(\cdot)$  sufficiently*  
599 *large and any  $\frac{1}{\Delta_0 \sqrt{n}} \leq \Delta \leq \Delta_0$  we have*

$$e^{\left(-\frac{\Delta^2}{2p(1-p)} - \varrho\right)n} \leq \mathbb{P} \left[ \frac{\text{Bin}(n, p)}{n} \leq p - \delta \right] \leq e^{\left(-\frac{\Delta^2}{2p(1-p)} + \varrho\right)n}.$$

600 **Lemma 5** (Lemma A.1 in [MNS16]). *Let  $\{a_n\}_{n \geq n_0}$  be a sequence satisfying*

$$\underline{\gamma} \leq \frac{a_n}{n} \leq \bar{\gamma}.$$

601 *Then the  $\text{Beta}(n - a_n + 1, a_n + 1)$  distribution on  $[0, 1]$  obeys a central limit theorem with mean  $\frac{a_n}{n}$*   
602 *and standard deviation  $\sqrt{\frac{(a_n/n)(1-(a_n/n))}{n}}$  in the sense that for any bounded sequence  $(w_n)_{n \geq n_0}$*   
603 *of real numbers and with  $\Phi$  the normal CDF,*

$$\lim_{n \rightarrow \infty} \left| \Phi(w_n) - \mathbb{P}^{x \sim \text{Beta}(n - a_n + 1, a_n + 1)} \left[ \left( x - (a_n/n) \right) \cdot \sqrt{\frac{n}{(a_n/n)(1 - (a_n/n))}} \leq w_n \right] \right| = 0.$$

604 In the next two lemmas, we lower bound the probability that  $\hat{p}_{i,t}$  changes significantly when the  
605 number  $n_{i,t}$  of samples for  $a_i$  increases by a factor  $(1 + \varrho)$ .

606 **Lemma 6.** *Assume  $\mu$  is  $(\underline{\gamma}, \bar{\gamma}, \underline{f}, \bar{f})$ -admissible. Suppose that arm  $i$ 's average reward  $\hat{p}_{i,t}$  after*  
607  *$n = n_{i,t}$  samples satisfies*

$$\hat{p}_{i,t} \in [\beta, \bar{\gamma}]. \tag{B.12}$$

608 *Then for  $n \geq C(\underline{\gamma}, \bar{\gamma}, \underline{f}, \bar{f}, \beta)$  sufficiently large,*

$$\mathbb{P}^{x \sim \mu_{i,n}} [x \leq \hat{p}_{i,t}] \geq \frac{\underline{f}}{3\bar{f}}. \tag{B.13}$$

609 *Proof.* Let  $R_{i,t} = n\hat{p}_{i,t}$  be the total reward from arm  $i$  so far. The posterior distribution  $\mu_{i,t}$  for  $p_i$   
610 takes the form

$$\mu_{i,t}(dx) = \frac{x^{R_{i,t}} (1-x)^{n-R_{i,t}} f(x) dx}{\int_{\underline{\gamma}}^{\bar{\gamma}} x^{R_{i,t}} (1-x)^{n-R_{i,t}} f(x) dx}.$$

611 For  $x \in [\underline{\gamma}, \bar{\gamma}]$  we estimate

$$\frac{x^{R_{i,t}} (1-x)^{n-R_{i,t}} f(x)}{\int_{\underline{\gamma}}^{\bar{\gamma}} x^{R_{i,t}} (1-x)^{n-R_{i,t}} f(x) dx} \geq (\underline{f}/\bar{f}) \cdot \frac{x^{R_{i,t}} (1-x)^{n-R_{i,t}}}{\int_0^1 x^{R_{i,t}} (1-x)^{n-R_{i,t}} dx}.$$

612 The right-hand side is the density of a beta variable with parameters  $(R_{i,t} + 1, n - R_{i,t} + 1)$ . We  
613 conclude that

$$\mathbb{P}^{x \sim \mu_{i,t}} [x \in [\underline{\gamma}, \hat{p}_{i,t}]] \geq (\underline{f}/\bar{f}) \cdot \mathbb{P}^{z \sim \text{Beta}(n - R_{i,t} + 1, R_{i,t} + 1)} [z \in [\underline{\gamma}, \hat{p}_{i,t}]]$$

614 For  $n$  sufficiently large, it follows from Lemma 5 and (B.12) that

$$\mathbb{P}^{z \sim \text{Beta}(n - R_{i,t} + 1, R_{i,t} + 1)} [z \in [\underline{\gamma}, \hat{p}_{i,t}]] \geq \frac{1}{3}.$$

615 Therefore  $\mathbb{P}^{\mu_{i,t}} [p_i \leq \hat{p}_{i,t}] \geq \frac{1}{3}$ , proving (B.13).  $\square$

616 **Lemma 7.** Assume  $\mu$  is  $(\underline{\gamma}, \bar{\gamma}, \underline{f}, \bar{f})$ -admissible and that (B.12) holds. For  $n = n_{i,t}$ , let  $\tilde{n} \geq 1$   
 617 satisfy  $|\tilde{n} - \varrho n| \leq 2$ . Let

$$\tilde{p}_i = \frac{R_{i,n+\tilde{n}} - R_{i,n}}{\tilde{n}}$$

618 be the average reward from the  $(n+1)$ -th through  $(n+\tilde{n})$ -th samples of arm  $i$ . Then as  $n \rightarrow \infty$ ,  
 619 for any sequence  $\Delta_n = \Theta(1/\log n)$ ,

$$\mathbb{P}^t[\tilde{p}_i \leq \theta^{-1}(\theta(\hat{p}_{i,t}) - \delta)] \geq \exp\left(-\frac{n\varrho\Delta_n^2(1+o_n(1))}{2}\right). \quad (\text{B.14})$$

620 *Proof.* Stochastic monotonicity implies that

$$\mathbb{P}\left[\frac{\text{Bin}(\tilde{n}, p)}{\tilde{n}} \leq \theta^{-1}(\theta(\hat{p}_{i,t}) - \Delta_n)\right]$$

621 is a decreasing function of  $p \in [0, 1]$ . Combining with Lemma 6, it follows that

$$\begin{aligned} \mathbb{P}^t[E] &= \int \mathbb{P}\left[\frac{\text{Bin}(\tilde{n}, x)}{\tilde{n}} \leq \theta^{-1}(\theta(\hat{p}_{i,t}) - \Delta_n)\right] d\mu_{i,t}(x) \\ &\geq \mathbb{P}^{\mu_{i,t}}[p_i \leq \hat{p}_{i,t}] \cdot \mathbb{P}\left[\frac{\text{Bin}(\tilde{n}, \hat{p}_{i,t})}{\tilde{n}} \leq \theta^{-1}(\theta(\hat{p}_{i,t}) - \Delta_n)\right] \\ &\geq \frac{f}{3\bar{f}} \cdot \mathbb{P}\left[\frac{\text{Bin}(\tilde{n}, \hat{p}_{i,t})}{\tilde{n}} \leq \theta^{-1}(\theta(\hat{p}_{i,t}) - \Delta_n)\right]. \end{aligned}$$

622 Since  $\theta$  is smooth with smooth inverse on  $[\underline{\gamma}, \bar{\gamma}]$  and  $\Delta_n \leq o_n(1)$ , we have

$$\begin{aligned} \hat{p}_{i,t} - \theta^{-1}(\theta(\hat{p}_{i,t}) - \Delta_n) &= (1 \pm o_n(1))\Delta_n \cdot (\theta^{-1})'(\theta(\hat{p}_{i,t})) \\ &= \frac{(1 \pm o_n(1)) \cdot \Delta_n}{\theta'(\theta^{-1}(\hat{p}_{i,t}))} \\ &= (1 \pm o_n(1)) \cdot \Delta_n \sqrt{\hat{p}_{i,t}(1 - \hat{p}_{i,t})}. \end{aligned}$$

623 The result now follows from Lemma 4, where we absorb the factor  $f/(3\bar{f})$  into the  $o_n(1)$ .  $\square$

## 624 B.6 Proof of Theorem 3.2

625 Recall the definition (B.8) of  $\underline{\gamma}$  and  $\bar{\gamma}$ . We require  $\mathcal{A}$  to be  $B$ -batch-compressed for  $B = B(N, \varrho)$   
 626 containing:

- 627 1. All positive integers at most  $N^{2\varrho}$ .
- 628 2. All positive multiples of  $\lfloor N^\varrho \rfloor$  at most  $N^{6\varrho}$ .
- 629 3. Integers of the form  $\lfloor N^{6\varrho}(1 + \varrho)^j \rfloor$  for  $j \geq 0$ .

630 It is easy to see that  $B$  thus defined is  $\varrho$ -slowly increasing for any  $\varrho > 0$  and  $N$  sufficiently large.  
 631 We denote  $b_k = \lfloor N^{6\varrho}(1 + \varrho)^k \rfloor$  so that  $|b_{k+1} - (1 + \varrho)b_k| \leq 2$ . (This choice of indexing differs  
 632 from that of Definition B.1, which will not be used in the sequel.)

633 We next construct our randomness distorting adversary  $\mathbb{A} = \mathbb{A}(N, \varrho)$ . For each arm  $i$ , the adversary  
 634  $\mathbb{A}$  acts as follows depending on the current number of samples  $n_{i,t}$ .

- 635 1. If  $n_{i,t} \leq N^{2\varrho}$ , then  $\mathbb{A}$  does nothing.
- 636 2. When  $N^{2\varrho} \leq n_{i,t} < N^{6\varrho}$  increases by  $N^\varrho$ ,  $\mathbb{A}$  declares that the average reward of this batch  
 637 of  $N^\varrho$  samples is at most  $\bar{\gamma} - N^{-\varrho}$ .
- 638 3. When  $n_{i,t}$  increases from  $b_k \geq N^{6\varrho}$  to  $b_{k+1}$ :  
 639 (a) If  $\hat{p}_i(b_k) > \beta$  holds, then  $\mathbb{A}$  declares that

$$\theta(\hat{p}_i(b_{k+1})) \leq \theta(\hat{p}_i(b_k)) - \frac{\varrho(1 + 10\varrho)d_F(\alpha, \beta)}{\log N}. \quad (\text{B.15})$$

640 (b) If  $\hat{p}_i(b_k) \leq \beta$  holds, then  $\mathbb{A}$  declares that

$$\hat{p}_i(b_{k+1}) \leq \beta.$$

641 4. When the  $\mathcal{A}$  chooses the arm  $a_{i^*}$  to output,  $\mathbb{A}$  declares that  $p_{i^*} < \beta$ .

642 Due to step 4, the declarations made by  $\mathbb{A}$  ensure that  $p_{i^*} < \beta$ . Recalling Lemma 4 and Proposi-  
643 tion B.2, it remains to show the upper bound

$$\text{strength}(\mathbb{A}) \leq \frac{(c_{\alpha, \beta} + C_* \varrho)N}{\log^2(N)}$$

644 for a constant  $C_* = C_*(\underline{\gamma}, \bar{\gamma}, \underline{f}, \bar{f}, \beta, \alpha)$  independent of  $\varrho$  (and  $N$ ). We show this bound in several  
645 parts. Recalling (3.3), we refer to the *cost* of a step above as the contribution to *Cost* from the  
646 corresponding declarations by  $\mathbb{A}$ . The most important parts are Lemmas 10 and 11, which bound  
647 the cost of the main step 3a and form the dominant contribution to *Cost*. Note that throughout the  
648 analysis below, all cost upper bounds hold almost surely and we **assume that all of  $\mathbb{A}$ 's declarations**  
649 **hold true**.

650 **Lemma 8.** *The total cost from step 2 is at most  $C_* N^{1-e}$ , for  $N \geq C(\underline{\gamma}, \bar{\gamma}, \underline{f}, \bar{f}, \beta, \alpha, \varrho)$  sufficiently*  
651 *large.*

652 *Proof.* The probability for each such declaration by  $\mathbb{A}$  is at least

$$\mathbb{P}[\text{Bin}(N^{2e}, \bar{\gamma}) \leq \bar{\gamma}N^{2e} - N^e] \tag{B.16}$$

653 since  $p_i \leq \bar{\gamma}$  almost surely. Recall that a  $\text{Bin}(N^{2e}, \bar{\gamma})$  random variable obeys a central limit theorem  
654 centered at  $\bar{\gamma}N^{2e}$  with standard deviation at least  $C(\bar{\gamma})N^e$ . Therefore the probability in (B.16) is  
655 at least  $\frac{1}{3}$  for  $N$  is sufficiently large depending on  $\varrho$ . Hence each such declaration costs at most  
656  $C_*$  for  $N$  sufficiently large. Moreover such declarations can occur only  $N^{1-e}$  times because each  
657 one involves  $N^e$  samples, and the base algorithm  $\mathcal{A}$  is an  $N$ -sample algorithm. This completes the  
658 proof.  $\square$

659 **Lemma 9.** *The total cost from step 3b is at most  $C_* N^{1-6e}$  as long as  $N \geq C(\underline{\gamma}, \bar{\gamma}, \underline{f}, \bar{f}, \varrho)$ .*

660 *Proof.* It suffices to show that the cost per step 3b declaration is at most  $C_*$ . This follows from  
661 (B.13) and stochastic monotonicity.  $\square$

662 **Lemma 10.** *The total cost from step 3a is at most*

$$\frac{N}{\log^2(N)} \cdot (c_{\alpha, \beta} + C_* \varrho + o_N(1)).$$

663 *Proof.* We claim that the cost from a single instance of step 3a when increasing from  $b_k$  to  $b_{k+1}$   
664 samples is at most

$$\left( \frac{(b_{k+1} - b_k)}{\log^2(N)} \right) (c_{\alpha, \beta} + C_* \varrho + o_N(1)).$$

665 This implies the desired result since  $\mathcal{A}_N$  is an  $N$ -sample algorithm. Taking  $\Delta = (1 +$   
666  $10\varrho)d_F(\alpha, \beta)/\log(N)$  in Lemma 7, we find that the declared event has probability at least

$$\exp\left(-\frac{(b_{k+1} - b_k)(1 + 10\varrho)^2 d_F(\alpha, \beta)^2 (1 + o_N(1))}{2\log^2(N)}\right) \geq \exp\left(-\frac{(b_{k+1} - b_k)}{\log^2(N)} (c_{\alpha, \beta} + C_* \varrho + o_N(1))\right).$$

667 This implies the desired claim and completes the proof.  $\square$

668 **Lemma 11.** *For any  $a_i$  sampled  $b_0 = \lfloor N^{6e} \rfloor$  times,  $\hat{p}_i(b_0) \leq \bar{\gamma}$ .*

669 *Proof.* By definition of  $\mathbb{A}$ ,

$$\begin{aligned} \hat{p}_i(b_0) &\leq \frac{N^{2e} + (N^{6e} - N^{2e})(\bar{\gamma} - N^{-e})}{N^{6e}} \\ &= \bar{\gamma} - \frac{1}{N^e} + \frac{(1 - \bar{\gamma})}{N^{4e}} + \frac{1}{N^{5e}} \\ &\leq \bar{\gamma}. \end{aligned}$$



670 In the last step we used the fact that

$$\frac{1}{N^\varrho} \geq \frac{(1 - \bar{\gamma})}{N^{4\varrho}} + \frac{1}{N^{5\varrho}}$$

671 for any  $\varrho > 0$  if  $N$  is sufficiently large.  $\square$

672 **Lemma 12.** For  $\varrho \in (0, 1/100)$ , if  $n_{i,t} \geq N^{1-\varrho}$  and the declarations of  $\mathbb{A}$  hold, then  $\hat{p}_{i,t} \leq \beta$ .

673 *Proof.* We analyze the rate at which the adversary forces  $\theta(\hat{p}_i(b_k))$  to decrease. From (B.15) and  
674 (11) it follows that for  $k$  with  $b_k \geq N^{1-\varrho}$ , we have

$$\begin{aligned} \theta(\hat{p}_i(b_k)) &\leq \theta(\bar{\gamma}) - \frac{\varrho(1 + 10\varrho)d_F(\alpha, \beta) \log_{1+\varrho}(N^{1-8\varrho})}{\log N} \\ &= \theta(\bar{\gamma}) - \frac{\varrho(1 + 10\varrho)(1 - 8\varrho)d_F(\alpha, \beta)}{\log(1 + \varrho)} \\ &\leq \theta(\bar{\gamma}) - (1 + \varrho)d_F(\alpha, \beta) \\ &\stackrel{\text{(B.8)}}{<} \theta(\beta). \end{aligned}$$

675 Here we used the fact that  $\log(1 + \varrho) \leq \varrho$  and  $(1 + 10\varrho)(1 - 8\varrho) \geq 1$  for  $\varrho \in (0, 1/100)$ . Since  $\theta$   
676 is increasing, this shows that  $\hat{p}_{i,t} = \hat{p}_i(b_k) < \beta$  for  $b_k \geq N^{1-\varrho}$ , completing the proof.  $\square$

677 **Lemma 13.** The cost from step 4 is at most  $C_*(N^{1-\varrho} + 1)$ .

678 *Proof.* First, if  $\hat{p}_{i^*,N} \leq \beta$  then the cost from step 4 is at most  $C_*$ . On the other hand if  $\hat{p}_{i^*,N} > \beta$ ,  
679 then Lemma 11 implies  $n_{i^*,N} \leq N^{1-\varrho}$ . Since the prior  $\mu$  is supported in  $[\underline{\gamma}, \bar{\gamma}]$ , the likelihood ratio  
680 of updates from  $N^{1-\varrho}$  samples is almost surely bounded by  $e^{C_*N^{1-\varrho}}$ . Therefore

$$\begin{aligned} \mathbb{P}^{x \sim \mu_{i^*,N}}[x < \beta] &\geq e^{-C_*N^{1-\varrho}} \mathbb{P}^{x \sim \mu}[x < \beta] \\ &\geq e^{-C_*N^{1-\varrho}} \frac{(\beta - \underline{\gamma})f}{\bar{f}}. \end{aligned}$$

681 This completes the proof.  $\square$

682 We now combine the lemmas above to conclude Theorem 3.1 via (B.7).

683 *Proof of Theorem 3.1.* Let  $C'_*$  be a larger constant depending on the same parameters. Then by  
684 Lemmas 8, 9, and 13, the total cost from Steps 2, 3b, 4 combines to  $C'_*N^{1-\varrho} \leq o_N(N/\log^2 N)$ .  
685 The main cost contribution of

$$\frac{N}{\log^2 N} (c_{\alpha,\beta} + C_*\varrho + o_N(1)).$$

686 comes from Lemma 10, and all other terms are of strictly smaller order. We have thus constructed a  
687 reservoir sequence  $(\mu_N(\varrho))_{N \geq 1}$  satisfying (B.7) for arbitrary  $\varrho > 0$ , completing the proof.  $\square$

## 688 C An Optimal Algorithm with Fixed Budget

689 Here we provide an asymptotically optimal algorithm which establishes Theorems B.1, B.2, and B.3.  
690 In the next subsection in which we show how to reduce the other results mentioned to Theorem 3.1  
691 (in which  $\alpha$  is given) using Proposition A.1. Our main focus will then be to prove Theorem 3.1.

692 We will fix  $\varrho > 0$  small and construct a sequence of  $N$ -sample algorithms  $(\mathcal{A}(N, \varrho))$  satisfying the  
693 slightly relaxed guarantee

$$\liminf_{N \rightarrow \infty} \frac{(-\log(\mathbb{P}^{\mu_N(\varrho)}[p_{i^*} < \beta])) \cdot \log^2 N}{N} \geq c_{\alpha,\beta} - \lambda(\varrho) \quad (\text{C.1})$$

694 for a (possibly different) function  $\lambda$  satisfying  $\lim_{\varrho \rightarrow 0} \lambda(\varrho) = 0$  (for fixed  $\alpha, \beta, \eta$ ). Here  $(\mu_N)_{N \geq 1}$  is  
695 any sequence of reservoir distributions satisfying  $G_{\mu_N}^{-1}(1 - \eta) = \alpha$ . An elementary diagonalization  
696 argument then implies Theorem 3.1. Thus it suffices to construct algorithms satisfying (C.1) for any  
697 desired  $\varrho > 0$ .

698 **C.1 Reduction to Known  $\alpha$**

699 We explain why Theorems B.1, B.2, and B.3 all follow from Theorem 3.1 (more precisely, the  
700 uniform statement given in Remark B.1). We begin with Theorem B.1, where

$$\alpha_N = \frac{1}{\eta_1 - \eta_2} \cdot \int_{1-\eta_1}^{1-\eta_2} G_{\mu_N}^{-1}(x) dx.$$

701 Let  $J = \lceil \frac{6}{\varepsilon(\eta_1 - \eta_2)} \rceil$  and define

$$\eta^{(j)} = \frac{(J-j)\eta_1 + j\eta_2}{J}, \quad j \in [J].$$

702 It is easy to see that  $\eta^{(j+1)} - \eta^{(j)} \leq \eta^{(j)}$  for all  $j$ . We next apply Alg. 4 on  $(\eta^{(j)}, \eta^{(j+1)} - \eta^{(j)}, \varepsilon', \delta')$   
703 for  $0 \leq j \leq J-1$ , with:

$$\begin{aligned} \varepsilon' &= \log^{-1/3}(N), \\ \delta' &= e^{-\frac{10N}{\log^2(N)}} / J. \end{aligned}$$

704 This requires sample complexity

$$N_A \leq \frac{C(\eta_1, \eta_2)N \log \log(N)}{\log(N)} \leq o_N(N). \quad (\text{C.2})$$

705 Let  $\hat{\alpha}_j$  be the resulting output. With probability  $1 - J\delta$ , we have for each  $0 \leq j \leq J-1$ ,

$$\hat{\alpha}_j \in \left[ G^{-1}(1 - \eta^{(j)}) - \frac{\varepsilon}{3}, G^{-1}\left(1 - \eta^{(j+1)}\right) + \frac{\varepsilon}{3} \right]. \quad (\text{C.3})$$

706 Note that the function  $G_{\mu}^{-1}$  is increasing and  $[0, 1]$ -valued. Therefore if (C.3) holds for each  $j$ , then

$$\left| \frac{1}{J} \cdot \sum_{j=0}^{J-1} \hat{\alpha}_j - \frac{1}{\eta_1 - \eta_2} \cdot \int_{1-\eta_1}^{1-\eta_2} G_{\mu_N}^{-1}(x) dx \right| \leq \frac{\varepsilon}{3} + \frac{1}{J} \leq \frac{\varepsilon}{2}.$$

707 Therefore the estimator

$$\hat{\alpha}_A = \frac{1}{J} \cdot \sum_{j=0}^{J-1} \hat{\alpha}_j$$

708 satisfies

$$\mathbb{P} \left[ \left| \hat{\alpha}_A - \frac{1}{\eta_1 - \eta_2} \cdot \int_{1-\eta_1}^{1-\eta_2} G_{\mu_N}^{-1}(x) dx \right| \leq \varepsilon/2 \right] \geq 1 - J\delta' = 1 - e^{-\frac{10N}{\log^2(N)}}.$$

709 Finally,  $c_{\alpha, \alpha - \varepsilon} \leq \pi < 10$  for any  $\alpha, \varepsilon \in [0, 1]$  (see (B.10)). Therefore the  $\delta' = e^{-\frac{10N}{\log^2(N)}}$  failure  
710 probability above has a negligible contribution in Theorem B.1. It follows that applying Theorem 3.1  
711 with  $\alpha = \hat{\alpha}_A$  as above and  $N' = N - N_A$  implies Theorem B.1.

712 We now turn to Theorem B.2, where  $\mu_N$  is required to satisfy  $G_{\mu_N}^{-1}(1 - \eta) \geq \frac{1+\varepsilon}{2}$ . We run Alg. 4  
713 with parameters

$$\begin{aligned} \eta_1 &= \eta, \\ \eta_2 &= \log^{-1/3}(N), \\ \varepsilon' &= \log^{-1/3}(N), \\ \delta' &= e^{-\frac{10N}{\log^2(N)}}. \end{aligned}$$

714 The sample complexity  $N_B$  again satisfies  $N_B \leq o(N)$  exactly as in (C.2). Let  $\hat{\alpha}_B + \varepsilon'$  be the  
715 resulting output. Then with probability at least  $1 - e^{-\frac{10N}{\log^2(N)}}$ ,

$$\hat{\alpha}_B \geq G_{\mu_N}^{-1}(1 - \eta) - 2\varepsilon'$$

716 and so with  $\varepsilon'' = \varepsilon - 2\varepsilon'$ , we have

$$\hat{\alpha}_B - \varepsilon'' \geq G_{\mu_N}^{-1}(1 - \eta) - \varepsilon.$$

717 Moreover, also with probability at least  $1 - e^{-\frac{10N}{\log^2(N)}}$ ,

$$\hat{\alpha}_B \leq G_{\mu_N}^{-1}(1 - \eta + \eta_2).$$

718 It follows that applying the algorithm of Theorem 3.1 with

$$(N, \alpha, \eta, \varepsilon) = (N - N_B, \hat{\alpha}_B, \eta - \eta_2, \varepsilon - 2\varepsilon')$$

719 suffices to recover Theorem B.2, since  $\eta_2$  and  $\varepsilon'$  tend to 0 as  $N \rightarrow \infty$ . As in our discussion of  
720 Theorem B.1 above, the failure probability  $e^{-\frac{10N}{\log^2(N)}}$  is negligible compared to the relevant rate in  
721 Theorem B.2.

722 Finally, Theorem B.3 relies on the simple fact

$$\lim_{\eta \rightarrow 0} G_{\mu}^{-1}(1 - \eta) = \mu \tag{C.4}$$

723 Recall that  $\mu^* \in [0, 1]$  denotes the maximum value in the support of  $\mu$ . We run Alg. 4 on  
724  $(\eta_1, \eta_2, \varepsilon', \delta')$  where:

$$\begin{aligned} \eta_1 &= \log^{-1/3}(N), \\ \eta_2 &= \eta_1/2, \\ \varepsilon' &= \varepsilon_1 - \varepsilon, \\ \delta' &= e^{-\frac{10N}{\log^2(N)}}. \end{aligned}$$

725 It follows from Proposition A.1 that the resulting output  $\hat{\alpha}_C + \frac{\varepsilon_1 - \varepsilon}{2}$  is computed using  
726  $O\left(\frac{N \log \log(N)}{\log(N)}\right) \leq o(N)$  samples as in the previous cases. Moreover for  $N$  sufficiently large:

$$\begin{aligned} \mathbb{P}\left[\hat{\alpha}_C + \frac{\varepsilon_1 - \varepsilon}{2} \geq \mu^* - \frac{\varepsilon'}{3} - o_N(1)\right] &\stackrel{\text{(C.4)}}{\geq} \mathbb{P}\left[\hat{\alpha}_C + \frac{\varepsilon_1 - \varepsilon}{2} \geq G_{\mu}^{-1}(1 - \eta_1) - \frac{\varepsilon'}{3}\right] \\ &\geq 1 - \delta' \\ &= 1 - e^{-\frac{10N}{\log^2(N)}}. \end{aligned}$$

727 Since  $\varepsilon_1 > \varepsilon$ , this means for  $N \geq N_0(\mu, \varepsilon', \dots)$  large enough,

$$\mathbb{P}[\hat{\alpha}_C \geq \mu^* - (\varepsilon_1 - \varepsilon)] \geq 1 - e^{-\frac{10N}{\log^2(N)}}.$$

728 Note that Alg. 4 also ensures that with probability  $1 - e^{-\frac{10N}{\log^2(N)}}$ ,

$$\begin{aligned} \hat{\alpha}_C &\leq \mu^* + \frac{\varepsilon'}{3} - \frac{\varepsilon_1 - \varepsilon}{2} = \mu^* - \frac{\varepsilon_1 - \varepsilon}{6} \\ &\leq G_{\mu}^{-1}(1 - \eta') \end{aligned}$$

729 for some  $\eta'(\mu, \varepsilon_1, \varepsilon) > 0$ . It follows that applying Theorem 3.1 with

$$(N, \alpha, \eta, \varepsilon) = (N - N', \hat{\alpha}_C, \eta', \varepsilon)$$

730 implies Theorem B.3.

## 731 C.2 The Fixed Budget Algorithm

732 We now present Algorithm 3 for the fixed budget problem (recall the informal discussion in Sec-  
733 tion 3). Algorithm 3 studies one arm  $a_i$  at a time, moving to  $a_{i+1}$  if  $a_i$  is rejected. Similarly to the  
734 previous section, some details are needed while  $n_{t,i}$  is small, since large deviation asymptotics may  
735 not have kicked in yet. As explained at the start of the section, we choose a small constant  $\varrho > 0$ .  
736 In fact, we will eventually choose small constants

$$0 < \varrho \ll \varrho_1 \ll \varrho_2 \ll \varrho_3 \ll \varrho_4 \ll \varrho_5 \ll 1$$

737 which all tend to 0 as  $\varrho \rightarrow 0$ . These constants will be defined throughout the proof. More formally,  
738 these values can be obtained by choosing  $\varrho_5 > 0$  arbitrarily small, then  $\varrho_4 > 0$  sufficiently small  
739 depending on  $\varrho_5$ , and so on.

740 Algorithm 3 operates in a batch-compressed way, for a sequence  $(b_1, b_2, \dots)$  defined as follows:

$$\begin{aligned}
b_0 &= \lceil \varrho_1 \log^2(N) \rceil, \\
k_0 &= \lceil \log_{1+\varrho}(\log^4(N)/b_0) \rceil \\
b_k &= b_0(1+\varrho)^k, \quad k \leq k_0 \\
b_{k_0+j} &= \lceil (1+\varrho)^j b_{k_0} \rceil, \quad j \geq 1 \\
\tau_k &= \alpha - \varrho - \frac{k}{\sqrt{\log N}}, \quad k \leq k_0 \\
\tau_{k_0+j} &= \theta(\alpha - 2\varrho) - j \cdot \frac{d_F(\alpha, \beta)\varrho(1-\varrho_2)}{\log N}, \quad j \geq 1.
\end{aligned}$$

741 Note in particular that  $b_{k_0} \geq \log^4(N)$ . We denote by  $\hat{p}_{i,t}$  the empirical average reward collected by  
742  $a_i$  from its first  $t$  samples.

---

**Algorithm 5:** Output arm with  $p_i \geq \beta$  using  $N$  samples with high probability

---

```

1 input: an infinite sequence of arms  $i = 1, 2, \dots$ 
2 initialize:  $i = 0$ 
3 while fewer than  $N$  samples have been collected do
4    $i \leftarrow i + 1$ 
5   Collect  $b_0$  samples of arm  $i$ .
6   if  $\hat{p}_{i,b_0} \leq \alpha - \varrho$  then
7     Reject arm  $i$ 
8   end
9   for  $k = 1, 2, \dots, k_0$  do
10    Collect  $b_k - b_{k-1}$  samples of arm  $i$  for a total of  $b_k$  samples.
11    if  $\hat{p}_{i,b_k} \leq \alpha - \varrho - \frac{k}{\sqrt{\log N}}$  then
12      Reject arm  $i$ ;
13    end
14  end
15  for  $j = 1, 2, \dots$  do
16    Collect  $b_{k_0+j} - b_{k_0+j-1}$  samples of arm  $i$  for a total of  $b_{k_0+j}$ .
17    if  $\theta(\hat{p}_{i,b_{k_0+j}}) \leq \theta(\alpha - 2\varrho) - j \cdot \frac{d_F(\alpha, \beta)\varrho(1-\varrho_2)}{\log N}$  then
18      Reject arm  $i$ 
19    end
20  end
21 end
22 Return arm  $i$ .

```

---

743 The role of the values  $b_j$  is as follows. When an arm  $a_i$  reaches  $b_k$  samples for some  $k \geq 0$ ,  
744 it is checked for possible rejection by comparing its empirical average reward to the threshold  $\tau_k$ .  
745 Algorithm 3 rejects arm  $i$  and moves to arm  $a_{i+1}$  if the empirical average  $\hat{p}_{i,b_k}$  of arm  $a_i$  drops below  
746 a moving threshold  $\tau_k$ . The threshold  $\tau_k$  begins close to  $\alpha$  and gradually decreases until reaching  
747  $\beta + \varrho$  by the time  $\tau_k \geq \Omega(N)$ .

748 So for, our informal description of Alg. 3 also applies to the algorithm proposed in [GM20]. We now  
749 highlight two important differences. The first is that our algorithm is defined more carefully during  
750 the “early” phases when an arm has been sampled at most  $N^{O(\varrho)}$  times. This is crucial for carrying  
751 out a rigorous analysis. The second difference is that in the main phase, we increase the sample size  
752 for a given arm in powers of  $1 + \varrho$  rather than powers of 2, and also move the rejection thresholds  $\tau_k$   
753 based on the Fisher information distance via the function  $\theta$ . The latter ingredients allow us to obtain  
754 the optimal constant factor.

755 We begin the analysis of Alg 3 by proving Lemma 3.

756 *Proof of Lemma 3.* Let  $M_j = \prod_{1 \leq i \leq j} Y_i$  and observe that  $M_j^c$  is a positive supermartingale with  
757  $M_0 = 0$ . The result follows by Doob’s maximal inequality.  $\square$

758 We will apply Lemma 3 in the following way. Let  $X_i$  be the number of samples used by arm  
759  $a_i$  before rejection, and  $I_i \in \{0, 1\}$  be the indicator of the event that  $a_i$  is ever rejected, even if  
760 Algorithm 3 were to continue past time  $N$  and sample arm  $i$  an infinite number of times. We set

$$Y_i = e^{X_i} \cdot I_i,$$

761 With  $M$  defined from  $(Y_i)_{i \geq 1}$  as in Lemma 3, it follows that  $\log(M)$  is at most the amount of  
762 time spent on eventual rejections before the first eventually accepted arm. Therefore if  $\log(M) \leq$   
763  $N(1 - \varrho)$ , we conclude that the last arm to be studied was sampled at least  $N\varrho$  times. Since it was  
764 not rejected during that time, we can conclude this arm has  $p_i \geq \beta$  with probability  $1 - e^{-\Omega_\varrho(N)}$ .  
765 The main contribution to the failure probability of Algorithm 3 comes from the event  $\{M \geq A\}$   
766 above, for suitable  $A$ . Correspondingly, the main work will be to verify  $\mathbb{E}[Y_i^c] \leq 1$  for suitable  $c$ .

767 Note that  $Y_i \in \{0\} \cup [1, \infty)$  almost surely for each  $i$ . Therefore a necessary first step in showing  
768  $\mathbb{E}[Y_i^c] \leq 1$  is to lower bound  $\mathbb{P}[Y_i = 0]$ , the probability that Algorithm 3 never rejects  $a_i$ . We now  
769 give a sufficient lower bound from the event  $p_i \geq \alpha$ .

770 **Proposition C.1.** *Let  $x_1, x_2, \dots$  be an i.i.d. Bernoulli( $p$ ) sequence for  $p \geq \alpha$ , and let  $S_k = \sum_{i=1}^k x_i$   
771 and set*

$$\underline{S} = \inf_{k \geq 1} S_k/k.$$

772 *Then  $\underline{S} \geq \alpha - \varrho$  holds with probability at least  $c(\alpha, \varrho) > 0$ . Thus  $\mathbb{E}[I_i] \leq 1 - c(\alpha, \varrho)$ .*

773 *Proof.* Since the probability that  $\underline{S} \geq \alpha - \varrho$  is increasing in  $p$  it suffices to take  $p = \alpha$  and show  
774 the probability is positive for any  $\varrho > 0$ . Assume not. Then by restarting the indexing every time  
775  $S_k \leq k(\alpha - \varrho)$  holds, we find that

$$\liminf_{n \rightarrow \infty} S_n/n \leq \alpha - \varrho.$$

776 This contradicts the strong law of large numbers, thus completing the proof of the first assertion.  
777 The second assertion follows since if  $S_k/k \geq \alpha - \varrho$  for all  $k$  where  $x_1, \dots$  are the rewards of arm  
778  $i$ , then arm  $i$  will never be rejected by Algorithm 3.  $\square$

779 Based on Proposition C.1 above, to show

$$\mathbb{E} \left[ e^{X_i \cdot \frac{c_{\alpha, \beta} - \varrho_3}{\log^2 N}} \cdot I_i \right] \leq 1$$

780 (which is essentially what we want in light of Lemma 3), it suffices to show that

$$\mathbb{E} \left[ \left( e^{X_i \cdot \frac{c_{\alpha, \beta} - \varrho_3}{\log^2 N}} - 1 \right) \cdot I_i \right] \leq c(\alpha, \varrho). \quad (\text{C.5})$$

781 We let  $I_i^t = I_i \cdot 1_{X_i=t}$  be the event that arm  $i$  was rejected after exactly  $t$  steps. Since Alg 3 can  
782 only reject after  $b_j$  samples, we have

$$I_i = \sum_{j=0}^{\infty} I_i^{b_j}$$

783 We use this to break the left-hand side of (C.5) into three separate parts and estimate the parts  
784 separately. The parts correspond to  $b_0, b_1$  through  $b_{k_0}$ , and  $b_{k_0+1}$  onward. The first two parts are  
785 easier and handled in Subsection C.3 below. The final term is the main contribution and is handled  
786 in Subsection C.4.

### 787 C.3 Analysis of Algorithm 3 in the Small and Medium Sample Phases

788 Proposition C.2 bounds the contribution to (C.5) from the *small sample phase*, i.e. the first rejection  
789 condition in line 7 of Alg 3.

790 **Proposition C.2.** *For any  $\alpha, \varrho$  there is  $\varrho_1 > 0$  sufficiently small that with  $b_0$  as defined above, and  
791 with  $N$  sufficiently large,*

$$\mathbb{E} \left[ \left( e^{X_i \cdot \frac{c_{\alpha, \beta} - \varrho_3}{\log^2 N}} - 1 \right) \cdot I_i^{b_0} \right] \leq c(\alpha, \varrho)/4$$

792 *Proof.* It suffices to observe that for fixed  $\alpha, \varrho$  and  $\varrho_1$  small and  $N$  sufficiently large, we have

$$e^{b_0 \cdot \frac{c_{\alpha, \beta} - \varrho_3}{\log^2 N}} - 1 \leq e^{\varrho_1} - 1 \leq 2\varrho_1.$$

793

□

794 Proposition C.3 bounds the contribution to (C.5) from the *medium sample phase*, i.e. the second  
795 rejection condition in line 12 of Alg 3.

796 **Proposition C.3.** For any  $\alpha, \varrho, \varrho_1$  and for  $N$  sufficiently large,

$$\sum_{k=1}^{k_0} \mathbb{E} \left[ \left( e^{X_i \cdot \frac{c_{\alpha, \beta} - \varrho_3}{\log^2 N}} - 1 \right) \cdot I_i^{b_k} \right] \leq c(\alpha, \varrho)/4$$

797 *Proof.* The event  $I_i^{b_k}$  requires  $|\hat{p}_{i, b_k} - \hat{p}_{i, b_{k-1}}| \geq \frac{1}{\sqrt{\log N}}$ . Hence by a standard Chernoff estimate,  
798 regardless of the true reward probability  $p_i$ ,

$$\mathbb{E}[I_i^{b_k}] \leq e^{-\Omega_{\alpha, \varrho, \varrho_1}(b_k / \log N)}.$$

799 Since by construction  $b_0 \geq \varrho_1 \log^2 N$ , we have

$$\begin{aligned} \mathbb{E} \left[ \left( e^{X_i \cdot \frac{c_{\alpha, \beta} - \varrho_3}{\log^2 N}} - 1 \right) \cdot I_i^{b_k} \right] &\leq e^{b_k \frac{c_{\alpha, \beta} - \varrho_3}{\log^2 N} - \Omega_{\alpha, \varrho, \varrho_1}(b_k / \log N)} \\ &\leq e^{-\Omega_{\alpha, \varrho, \varrho_1}(\log N)} \\ &= N^{-\Omega_{\alpha, \varrho, \varrho_1}(1)}. \end{aligned}$$

800 Since  $k_0 \leq O(\log N)$ , summing gives the desired conclusion. □

801 Propositions C.2 and C.3 imply that the total contribution from rejections in the small and medium  
802 sample phases is at most  $c(\alpha, \varrho)/2$ . It remains to analyze the large sample phase in the following  
803 subsection.

#### 804 C.4 Analysis of Algorithm 3 in the Large Sample Phase

805 Similarly to the previous section, the main part of the analysis concerns the large sample phases  
806  $b_{k_0+j}$  for  $j \geq 1$ . Our goal is to precisely estimate the rejection probability at each time  $b_{k_0+j}$ . Note  
807 that these estimates should not depend on the true average rewards  $p_i$ .

808 Our approach is based on exchangeability and avoids any consideration of  $p_i$ . For a given value  $j$   
809 and a large constant  $L = L(\varrho)$ , consider the sequence of times

$$b_{k_0+j-L}, b_{k_0+j-L+1}, \dots, b_{k_0+j}$$

810 and the associated sequence of empirical average rewards

$$\hat{p}_{i, b_{k_0+j-L}}, \hat{p}_{i, b_{k_0+j-L+1}}, \dots, \hat{p}_{i, b_{k_0+j}}. \quad (\text{C.6})$$

811 It follows from the algorithm description that for  $I_i^{b_{k_0+j}}$  to occur, we must have

$$\hat{p}_{i, b_{k_0+j}} - \hat{p}_{i, b_{k_0+j-\ell}} \geq \ell \cdot \frac{d_F(\alpha, \beta)\varrho(1-\varrho_2)}{\log N}, \quad \forall 1 \leq \ell \leq L. \quad (\text{C.7})$$

812 This is clear for  $j > L$ , but it holds also for  $0 \leq j \leq L$  as for  $N$  sufficiently large,

$$\alpha - \varrho - \frac{k_0}{\sqrt{\log N}} - L \cdot \frac{d_F(\alpha, \beta)\varrho(1-\varrho_2)}{\log N} \geq \alpha - 2\varrho.$$

813 By exchangeability, conditioned on the future values  $\hat{p}_{i, b_{k_0+j}}, \dots, \hat{p}_{i, b_{k_0+j-\ell}}$  the law of  $\hat{p}_{i, b_{k_0+j-\ell-1}}$   
814 depends only on  $\hat{p}_{i, b_{k_0+j-\ell}}$  and is given explicitly by a hypergeometric variable. Recalling that

815  $R_{i,t} = n_{i,t} \hat{\rho}_{i,t}$  is the total reward from the first  $n_{i,t}$  samples of arm  $i$ ,  $R_{i,b_{k_0+j-\ell-1}}$  has hypergeo-  
 816 metric conditional law given by:

$$\begin{aligned} \mathbb{P}\left[R_{i,b_{k_0+j-\ell-1}} = k \mid (\hat{\rho}_{i,b_{k_0+j}}, \dots, \hat{\rho}_{i,b_{k_0+j-\ell}})\right] &= \mathbb{P}\left[R_{i,b_{k_0+j-\ell-1}} = k \mid \hat{\rho}_{i,b_{k_0+j-\ell}}\right] \\ &= \frac{\binom{b_{k_0+j-\ell-1}}{k} \binom{b_{k_0+j-\ell}-b_{k_0+j-\ell-1}}{R_{k_0+j-\ell}-k}}{\binom{b_{k_0+j-\ell}}{R_{k_0+j-\ell}}}. \end{aligned} \quad (\text{C.8})$$

817 We will refer to this as the HyperGeom( $b_{k_0+j-\ell}, b_{k_0+j-\ell-1}, R_{k_0+j-\ell}$ ) distribution. Importantly,  
 818 this distribution is independent of  $\mu$ . We exploit this below to control the probability of a given  
 819 sequence  $(\hat{\rho}_{i,b_{k_0+j-L}}, \hat{\rho}_{i,b_{k_0+j-L+1}}, \dots, \hat{\rho}_{i,b_{k_0+j}})$  of empirical average rewards. The following  
 820 useful result states that hypergeometric variables automatically inherit tail bounds from the corre-  
 821 sponding binomial random variables.

822 **Lemma 1** ([LP14, Hoe94]). *Fix non-negative integers  $A \geq B, C$  and let  $X \sim$   
 823 HyperGeom( $A, B, C$ ) and  $Y \sim \text{Bin}(B, C/A)$ . Then for any convex function  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,*

$$\mathbb{E}[f(X)] \leq \mathbb{E}[f(Y)].$$

824 **Lemma 2.** *For any  $0 < \underline{q} < \bar{q} < 1$  and constants  $\varrho > 0$  there exists  $\Delta_0(\underline{q}, \bar{q}, \varrho)$  and  $N_0(\underline{q}, \bar{q}, \varrho)$   
 825 such that the following holds for all  $p \in [\underline{q}, \bar{q}]$ . For  $n \geq n_0$  sufficiently large and  $\frac{1}{\Delta_0 \sqrt{n}} \leq \Delta \leq \Delta_0$ ,*

$$\mathbb{P}\left[\frac{\text{HyperGeom}(n(1+\varrho), n, np(1+\varrho))}{n} \leq p - \Delta\right] \leq e^{(-\frac{\Delta^2}{2p(1-p)} + \varrho)n}.$$

826 *Proof.* The corresponding binomial result Lemma 4 is proved in Theorem 2.2 in [DA92] by upper  
 827 bounding an exponential moment. The same proof applies here by Lemma 1.  $\square$

828 It will be convenient to define a restricted set of *good* sequences  $(q_L, q_{L-1}, \dots, q_0)$ . These satisfy  
 829 the key properties of empirical average reward sequences (C.6) for which  $I_i^{b_{k_0+j}}$  holds. We say such  
 830 a length  $L+1$  sequence is good if the following conditions are satisfied:

831 1.  $q_0 \in [\underline{q}, \bar{q}] \subseteq (0, 1)$  for constants  $0 < \underline{q} < \bar{q} < 1$  depending only on  $\varrho, L$ .

832 2.

$$\max_{\ell_1, \ell_2} |q_{\ell_1} - q_{\ell_2}| \leq O(1/\sqrt{\log N}). \quad (\text{C.9})$$

833 3. For each  $1 \leq \ell \leq L$ :

$$\begin{aligned} \theta(q_0) &\leq \theta(\alpha - 2\varrho) - j \cdot \frac{d_F(\alpha, \beta)\varrho(1 - \varrho_2)}{\log N} \\ &\leq \theta(\alpha - 2\varrho) - (j - \ell) \cdot \frac{d_F(\alpha, \beta)\varrho(1 - \varrho_2)}{\log N} \\ &\leq \theta(q_\ell). \end{aligned}$$

833 The third condition above is necessary for  $I_i^{b_{k_0+j}, i} = 1$ , and these together imply the first condition.  
 834 Indeed for fixed  $\underline{q}, \bar{q}$  and small  $\varrho \in (0, 1/10)$  one always has

$$\frac{\hat{\rho}_{i,b_{k_0+j-1}}}{\hat{\rho}_{i,b_{k_0+j}}}, \frac{1 - \hat{\rho}_{i,b_{k_0+j-1}}}{1 - \hat{\rho}_{i,b_{k_0+j}}} \in [1 - 2\varrho, (1 - 2\varrho)^{-1}]$$

835 for large enough  $N$  and any  $j$ . Hence it suffices to take  $\underline{q} = \beta(1 - 2\varrho)^L$  and  $\bar{q} = 1 - (1 - \alpha)(1 - 2\varrho)^L$ .  
 836 With this choice, if

$$\hat{\rho}_{i,b_{k_0+j-L}}, \hat{\rho}_{i,b_{k_0+j-L+1}}, \dots, \hat{\rho}_{i,b_{k_0+j}}$$

837 is **not** good and  $I_i^{b_{k_0+j}} = 1$ , then the second condition must be the only violated one. The fol-  
 838 lowing easy lemma controls the failure probability of the second condition. Recall from (C.8) that  
 839 conditioning on  $\hat{\rho}_{i,b_{k_0+j}}$  determines the joint conditional law of the previous conditional rewards,  
 840 regardless of  $\mu$ .

841 **Lemma 3.** All sequences violating only the second condition (C.9) above have probability at most

$$e^{-\Omega_{L,\varrho}(b_{k_0+j}/\log N)},$$

842 even after conditioning on an arbitrary value for  $\hat{p}_{i,b_{k_0+j}}$ .

843 *Proof.* The claim follows by an elementary Chernoff estimate for hypergeometric variables, which  
 844 hold just as for binomial variables by Lemma 1. Indeed the assumption implies that some adjacent  
 845 difference  $|\hat{p}_{i,b_{k_0+j-\ell}} - \hat{p}_{i,b_{k_0+j-\ell+1}}|$  has size  $\Omega(1/\sqrt{\log N})$ . (Note for applying the Chernoff bound  
 846 that  $L$  is a constant independent of  $N$ , and so  $b_{k_0+j-L} \geq \Omega_{L,\varrho}(b_{k_0+j})$ .)  $\square$

847 We now focus on upper-bounding the probability of any good sequence  $(q_L, \dots, q_0)$  appearing,  
 848 conditionally on  $q_0$ .

849 **Lemma 4.** For any good sequence  $(q_L, q_{L-1}, \dots, q_0)$  and  $j \geq 0$ ,

$$\begin{aligned} & \mathbb{P}\left[\left(\hat{p}_{i,b_{k_0+j-L}}, \hat{p}_{i,b_{k_0+j-L+1}}, \dots, \hat{p}_{i,b_{k_0+j}}\right) = (q_L, q_{L-1}, \dots, q_0) \mid p_{i,b_{k_0+j}} = q_0\right] \\ & \leq \exp\left(-\frac{(1-O(\varrho))}{2q_0(1-q_0)\varrho} \sum_{\ell=0}^{L-1} b_{k_0+j-\ell}(q_\ell - q_{\ell+1})^2\right). \end{aligned}$$

850 *Proof.* It suffices to show that

$$\mathbb{P}[\hat{p}_{i,b_{k_0+j-\ell-1}} = q_{\ell+1} \mid q_\ell] \leq \exp\left(-\frac{(1-O(\varrho))}{2q_0(1-q_0)\varrho} b_{k_0+j-\ell}(q_\ell - q_{\ell+1})^2\right)$$

851 This follows by applying Lemma 2 to the hypergeometric random variable

$$\hat{p}_{i,b_{k_0+j-\ell}} \cdot b_{k_0+j-\ell} - \hat{p}_{i,b_{k_0+j-\ell-1}} \cdot b_{k_0+j-\ell-1} = R_{b_{k_0+j-\ell}} - R_{b_{k_0+j-\ell-1}}.$$

852 The fact that

$$b_{k_0+j-\ell+1} - b_{k_0+j-\ell} = \varrho \cdot b_{k_0+j-\ell} \pm O(1)$$

853 leads to the factor of  $\varrho$  in the denominator of the desired result.  $\square$

854 **Lemma 5.** For fixed problem parameters and  $N$  large, any good sequence  $(q_L, \dots, q_0)$  satisfies

$$q_\ell \geq q_0 + \frac{\ell \cdot d_F(\alpha, \beta)\varrho(1-2\varrho_2) \cdot \sqrt{q_0(1-q_0)}}{(\log N)}$$

855 *Proof.* Recall that  $\theta'(q) = \frac{1}{\sqrt{q(1-q)}}$  and that  $\theta$  is smooth on  $[q, \bar{q}] \subseteq (0, 1)$ . By Item 2 above, all  $q_\ell$   
 856 are within  $o_N(1)$  of each other, so the result follows from the inverse function theorem. (Notice that  
 857 the factor  $(1-\varrho_2)$  changed to  $(1-2\varrho_2)$  above.)  $\square$

858 **Lemma 6.** For  $1 \leq m \leq L$  and any good sequence  $(q_L, \dots, q_0)$ , we have

$$\sum_{\ell=0}^{m-1} (q_\ell - q_{\ell+1})^2 \geq \frac{m \cdot d_F(\alpha, \beta)^2 \varrho^2 (1-4\varrho_2) \cdot q_0(1-q_0)}{\log^2 N}.$$

859 *Proof.* The result follows from Lemma 5 and Cauchy-Schwarz in the form

$$\sum_{\ell=0}^{m-1} (q_\ell - q_{\ell+1})^2 \geq m^{-1} \left( \sum_{\ell=0}^{m-1} |q_\ell - q_{\ell+1}| \right)^2.$$

860  $\square$

861 **Lemma 7.** For any good sequence  $(q_L, \dots, q_0)$  and  $j \geq 0$ , we have

$$\sum_{\ell=0}^{L-1} b_{k_0+j-\ell}(q_\ell - q_{\ell+1})^2 \geq (1-O(\varrho_2)) \cdot \frac{b_{k_0+j}\varrho d_F(\alpha, \beta)^2 \cdot q_0(1-q_0)}{\log^2 N}.$$



862 *Proof.* We break the sum into parts and apply Lemma 6 to each one. We have:

$$\begin{aligned}
\sum_{\ell=0}^{L-1} b_{k_0+j-\ell} (q_\ell - q_{\ell+1})^2 &= b_{k_0+j-L+1} \sum_{\ell=0}^{L-1} (q_\ell - q_{\ell+1})^2 + \sum_{m=1}^{L-1} (b_{k_0+j-m+1} - b_{k_0+j-m}) \sum_{\ell=0}^{m-1} (q_\ell - q_{\ell+1})^2 \\
&\geq \sum_{m=1}^{L-1} b_{k_0+j} \cdot \frac{\varrho}{(1+\varrho)^{m+10}} \cdot (1-4\varrho_2) \frac{m\varrho^2 d_F(\alpha, \beta)^2 \cdot q_0(1-q_0)}{\log^2 N} \\
&\geq (1 - O(\varrho + \varrho_2)) \cdot b_{k_0+j} \cdot \frac{\varrho^3 d_F(\alpha, \beta)^2 \cdot q_0(1-q_0)}{\log^2 N} \cdot \sum_{m=1}^{L-1} \frac{m}{(1+\varrho)^m}.
\end{aligned}$$

863 For  $L = L(\varrho) = O(\varrho^{-1} \log(\varrho^{-1}))$  sufficiently large,

$$\begin{aligned}
\sum_{m=1}^{L-1} \frac{m\varrho}{(1+\varrho)^m} &\geq (1-\varrho) \sum_{m=1}^{\infty} \frac{m}{(1+\varrho)^m} \\
&= (1-\varrho) \left( \sum_{m=1}^{\infty} \frac{1}{(1+\varrho)^m} \right)^2 \\
&= \frac{1-\varrho}{\varrho^2}.
\end{aligned}$$

864 Substituting and recalling that  $\varrho \ll \varrho_2$  completes the proof.  $\square$

865 Combining with Lemma 4 yields the second inequality below (the first is trivial).

866 **Corollary C.4.** For any  $\mu$  and  $q_0$ , we have

$$\begin{aligned}
&\mathbb{P}^{p_i \sim \mu} \left[ (\hat{p}_{i, b_{k_0+j-L}}, \hat{p}_{i, b_{k_0+j-L+1}}, \dots, \hat{p}_{i, b_{k_0+j}}) = (q_L, q_{L-1}, \dots, q_0) \right] \\
&\leq \mathbb{P} \left[ (\hat{p}_{i, b_{k_0+j-L}}, \hat{p}_{i, b_{k_0+j-L+1}}, \dots, \hat{p}_{i, b_{k_0+j}}) = (q_L, q_{L-1}, \dots, q_0) \mid p_{i, b_{k_0+j}} = q_0 \right] \\
&\leq \exp \left( - (1 - O(\varrho_2)) \frac{b_{k_0+j} d_F(\alpha, \beta)^2}{2 \log^2 N} \right).
\end{aligned}$$

867 **Lemma 8.** Let  $j_0$  be the largest  $j$  such that  $b_{k_0+j} \leq N$ . Then for  $N$  sufficiently large,

$$\sum_{j=1}^{j_0} \mathbb{E} \left[ e^{X_i \cdot \frac{c_{\alpha, \beta} - \varrho_3}{\log^2 N}} \cdot I_i^{b_{k_0+j}} \right] \leq c(\alpha, \varrho)/4.$$

868 *Proof.* Recall that  $c_{\alpha, \beta} = \frac{d_F(\alpha, \beta)^2}{2}$ , and observe that the number of total sequences  $(q_L, \dots, q_0) \in$   
869  $[0, 1]^{L+1}$  with  $b_{k_0+j+\ell} q_\ell \in \mathbb{Z}$  is at most  $N^{L+1}$  for each  $j \leq j_0$ . Combining Lemma 3 and Corol-  
870 lary C.4 and noting that the latter always gives the main contribution, we find for each  $j \leq j_0$ ,

$$\begin{aligned}
\mathbb{E} \left[ e^{X_i \cdot \frac{c_{\alpha, \beta} - \varrho_3}{\log^2 N}} \cdot I_i^{b_{k_0+j}} \right] &\leq N^{L+1} \exp \left( \frac{b_{k_0+j}}{\log^2 N} \cdot ((c_{\alpha, \beta} - \varrho_3) - (1 - O(\varrho_2)) c_{\alpha, \beta}) \right) \\
&\leq \exp \left( -\Omega \left( \frac{\varrho_3 b_{k_0+j}}{\log^2 N} \right) \right)
\end{aligned}$$

871 so long as  $\varrho_3$  is chosen so that  $\varrho_3 \gg \max(\varrho, \varrho_2)$ . In the last line we used the fact that  $b_{k_0+j} \geq$   
872  $b_{k_0} \geq \log^4 N$  to absorb the factor  $N^{L+1} \leq e^{\varrho \log^{3/2} N}$  for large  $N$ . Summing over  $j$  gives the

873 desired result, since for  $\varrho_4 = \Omega(\varrho_3)$  and  $N$  sufficiently large,

$$\begin{aligned}
\sum_{j=1}^{\infty} e^{-\Omega\left(\frac{\varrho_3 b_{k_0+j}}{\log^2 N}\right)} &\leq \sum_{m=1}^{\infty} e^{-\frac{\varrho_4(m+b_{k_0})}{\log^2 N}} \\
&= e^{-\varrho_4 \log^2 N} \sum_{m=1}^{\infty} e^{-\frac{\varrho_4 m}{\log^2 N}} \\
&\leq e^{-\varrho_4 \log^2 N} \cdot O\left(\frac{\log^2 N}{\varrho_4}\right) \\
&\leq e^{-\frac{\varrho_4 \log^2 N}{2}} \\
&\leq c(\alpha, \varrho)/4.
\end{aligned}$$

874

□

875 We now use Lemma 3 to conclude.

876 *Proof that Algorithm 3 achieves the guarantee of Theorem 3.1.* By combining Lemma 8 with the  
877 previous Propositions C.2 and C.3, it follows that

$$\mathbb{E}\left[e^{X_i \cdot \frac{c_{\alpha, \beta} - \varrho_3}{\log^2 N}} \cdot I_i\right] \leq 1.$$

878 Lemma 3 now implies that the total amount of time spent on eventually rejected arms is at most  
879  $N(1 - \varrho)$  with probability

$$e^{-\frac{(c_{\alpha, \beta} - \varrho_3)(1 - \varrho)N}{\log^2 N}}.$$

880 On this event, the output arm  $i^*$  satisfies  $n_{i^*, N} \geq \varrho N$  by definition. Since  $i^*$  was not rejected, for  
881  $j_1$  be the largest value such  $b_{k_0+j_1} \leq \varrho N$  we have

$$\hat{p}_{i^*, b_{k_0+j_1}} \geq \beta + \varrho.$$

882 The probability for this to hold if  $p_i \leq \beta$  is at most  $e^{-\Omega_e(N)}$ . Altogether we find that

$$\mathbb{P}[p_{i^*} \geq \beta] \geq 1 - \exp\left(-\frac{(c_{\alpha, \beta} - \varrho_5)N}{\log^2 N}\right) - e^{-\Omega_e(N)} \quad (\text{C.10})$$

883 for  $\varrho_5$  arbitrarily small. This concludes the analysis of Algorithm 3 (since the last error term is  
884 negligible). □

### 885 C.5 Finding Many Good Arms with a Fixed Budget

886 In this final subsection we observe that Algorithm 3 can be modified to output as many as  $\log N$   
887 distinct arms each of which satisfies the same  $(\eta, \varepsilon, \delta)$ -PAC guarantee<sup>2</sup>, with no degradation in the  
888 asymptotic failure probability. With other parameters fixed, we denote the  $N$ -sample version of  
889 Algorithm 3 by  $\mathcal{A}_N$  to emphasize the dependence on  $N$ . In particular,  $N$  both equals the number of  
890 steps in  $\mathcal{A}_N$  and appears (via its logarithm) in the description of  $\mathcal{A}_N$ 's individual steps.

891 Let  $\tilde{N} = N + \lceil \frac{2N}{\log^{1/2}(N)} \rceil$ . We consider a modified algorithm  $\tilde{\mathcal{A}}_{\tilde{N}}$  which mimicks the behavior of  
892  $\mathcal{A}_N$  with two changes:

- 893 1.  $\tilde{\mathcal{A}}_{\tilde{N}}$  is a  $\tilde{N}$ -sample algorithm.
- 894 2. If an arm  $a_i$  has not yet been rejected after  $M = \lceil N/\log^{3/2}(N) \rceil$  samples, then  $\tilde{\mathcal{A}}_{\tilde{N}}$   
895 accepts  $a_i$  and continues to  $a_{i+1}$ . In particular,  $\tilde{\mathcal{A}}_{\tilde{N}}$  may accept several arms instead of just  
896 one.

897 **Theorem C.9.** *With probability  $1 - \exp\left(-\frac{(c_{\alpha, \beta} - \varrho_5 - o_N(1))N}{\log^2 N}\right)$ ,  $\tilde{\mathcal{A}}_{\tilde{N}}$  accepts at least  $\log(N)$  distinct  
898 arms  $a_i$ , all of which satisfy  $p_i \geq \beta$ .*

<sup>2</sup>In fact  $\log N$  can be replaced by anything  $o_N(\log^2 N)$  by more precisely defining  $M$  and  $\tilde{N}$ .

899 The change from  $N$  to  $\tilde{N}$  is almost irrelevant in the actual statement of Theorem C.9 since  $\log(N) \geq$   
900  $\log(\tilde{N}) - o_N(1)$ . In particular,  $\tilde{\mathcal{A}}_{\tilde{N}}$  is a  $\tilde{N}$ -sample algorithm which outputs at least  $\log(\tilde{N}) - 1$  arms  
901 with probability  $1 - \exp\left(-\frac{(c_{\alpha,\beta} - \varrho_5 - o_N(1))\tilde{N}}{\log^2 \tilde{N}}\right)$ . It is certainly not really necessary to use the value  
902  $\log(N)$  rather than  $\log(\tilde{N})$  to describe the individual steps taken by  $\tilde{\mathcal{A}}_{\tilde{N}}$ . However introducing  $\tilde{N}$   
903 streamlines the proof below by letting us treat  $\mathcal{A}_N$  as a blackbox.

904 *Proof.* To show that all accepted arms  $a_i$  satisfy  $p_i \geq \beta$  with sufficiently high probability, it suffices  
905 to consider (C.10) with the final term replaced by  $e^{-\Omega_e(N/\log^{3/2}(N))}$ . In particular, observe that the  
906 main term does not change, even after multiplying the failure probability by  $O(\log^{3/2}(N))$  (the  
907 maximum possible number of arms accepted by  $\tilde{\mathcal{A}}_{\tilde{N}}$ . Thus we focus on showing that  $\tilde{\mathcal{A}}_{\tilde{N}}$  outputs at  
908 least  $\log(N)$  arms with high probability.

909 Consider yet another  $N$ -sample algorithm  $\hat{\mathcal{A}}_N$  which deletes each arm independently with probabil-  
910 ity  $1/N$  and follows  $\mathcal{A}_N$  on the set of non-deleted arms in order of increasing index. (Like  $\mathcal{A}_N$ ,  $\hat{\mathcal{A}}_N$   
911 never accepts arms before time  $N$ .) We simulate  $\tilde{\mathcal{A}}_{\tilde{N}}$  and  $\hat{\mathcal{A}}_N$  on the same reward sequences, i.e. we  
912 couple them so that the  $t$ -th sample of arm  $a_i$  always gives the same result for each  $(t, i)$ . We **claim**  
913 that in this coupling, conditioned on  $\tilde{\mathcal{A}}_{\tilde{N}}$  failing to accept  $\log(N)$  arms within the first  $\tilde{N}$  samples,  
914  $\hat{\mathcal{A}}_N$  has probability  $\Omega(N^{-\log(N)})$  to fail (i.e. output  $a_i$  with  $p_i < \beta$ ) when run for  $N$  samples.

915 First let us assume the claim and deduce Theorem C.9. Denote by  $p(N)$  the probability for  $\mathcal{A}_N$  to  
916 fail. Note that  $\hat{\mathcal{A}}_N$  has the same failure probability  $p(N)$ , having in fact the same behavior as  $\mathcal{A}_N$   
917 in distribution (as the set of deleted arms is independent of everything else). Moreover let  $\tilde{p}(\tilde{N}, k)$   
918 denote the probability that  $\tilde{\mathcal{A}}_{\tilde{N}}$  fails to accept at least  $k$  arms. The claim above implies that

$$\begin{aligned} \tilde{p}(\tilde{N}, \log N) &\leq O(N^{\log N}) \cdot p(N, 1) \\ &\leq e^{o_N(N/\log^2 N)} \cdot p(N, 1) \\ &\leq \exp\left(-\frac{(c_{\alpha,\beta} - \varrho_5 - o_N(1))N}{\log^2 N}\right). \end{aligned}$$

919 It remains to prove the above claim. Let us say the infinite i.i.d. reward sequence  $(r_{i,n})_{n \geq 1}$  of arm  
920  $a_i$  is **acceptable** if  $\mathcal{A}_N$  would not reject  $a_i$  within  $M$  samples, i.e.  $\tilde{\mathcal{A}}_{\tilde{N}}$  will either accept  $a_i$  or run  
921 out of samples before doing so. We take the point of view that each  $a_i$  is either acceptable or not (by  
922 randomly fixing the reward sequences at the start). Then with probability  $\Omega(N^{-\log(N)})$ , the first  
923  $\log(N)$  acceptable arms are skipped by  $\hat{\mathcal{A}}_N$ , and the first  $\tilde{N}$  unacceptable arms are not skipped. On  
924 this event, the first  $\tilde{N} - M \geq N$  samples obtained by  $\hat{\mathcal{A}}_N$ , i.e. all  $N$  of its samples, are drawn from  
925 unacceptable arms. On this event,  $\tilde{\mathcal{A}}_{\tilde{N}}$  fails with constant probability, which establishes the claim  
926 and completes the proof.  $\square$