

---

# Gaussian Randomized Exploration for Semi-bandits with Sleeping Arms

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 This paper provides theoretical analyses of problem-independent upper and lower  
2 regret bounds for Gaussian randomized algorithms in semi-bandits with sleeping  
3 arms, where arms may be unavailable in certain rounds, and available arms satisfy  
4 combinatorial constraints can be played simultaneously. We first introduce  
5 the CTS-G algorithm, an adaptation of Thompson sampling with Gaussian priors,  
6 achieving an upper bound of  $\tilde{O}(m\sqrt{NT})$  over  $T$  rounds with  $N$  arms and up to  $m$   
7 arms played per round, where  $\tilde{O}$  hides the logarithmic factors. Next, we present  
8 CL-SG, which improves upon CTS-G by using a single Gaussian sample per round,  
9 achieving a near-optimal upper regret bound of  $\tilde{O}(\sqrt{mNT})$ . We also establish that  
10 both algorithms have lower regret bounds of  $\Omega(\sqrt{mNT \ln \frac{N}{m}})$  and  $\Omega(\sqrt{mNT})$ ,  
11 respectively.

## 12 1 Introduction

13 We consider a sleeping semi-bandit problem with a fixed set  $[N] = \{1, 2, \dots, N\}$  of  $N$  base arms  
14 and each base arm  $a \in [N]$  is associated with an unknown probability distribution  $p_a$  supported on  
15  $[0, 1]$  and mean  $r_a$ . Unlike standard combinatorial bandits (Kveton et al., 2015), where a learning  
16 agent, in each round  $t = 1, \dots, T$ , plays a super arm (combinations of base arms)  $A_t \in \Theta$ ,  
17 where  $\Theta \subseteq 2^{[N]}$  is a *feasible set* that satisfy certain constraints, sleeping semi-bandits involve a  
18 time-varying feasible set  $\Theta_t \subseteq \Theta$ , revealed at each round  $t$ . After observing the feasible set  $\Theta_t$  in  
19 round  $t$ , the learning agent selects a super arm  $A_t \in \Theta_t$ , observes rewards  $r_{a,t} \sim p_a$  for each base  
20 arm  $a \in A_t$ , and aims to minimize the  $T$ -round (pseudo)-regret defined as follows.

$$\mathcal{R}(T) := \sum_{t=1}^T \mathbf{E} \left[ \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} r_a \right], \quad (1)$$

21 where  $A_t^* := \arg \max_{A \in \Theta_t} \sum_{a \in A} r_a$  denotes the optimal super arm in round  $t$  and the expectation  
22 is taken over  $\Theta_t$ ,  $A_t$ , and  $A_t^*$ . Note that  $A_t^*$  is determined by  $\Theta_t$ . We further denote by  $m :=$   
23  $\max_{A \in \Theta} |A|$  the maximum number of base arms in any super arm.

24 The *upper confidence bound (UCB)* (Agrawal, 1995; Auer et al., 2002) and *Thompson sam-*  
25 *pling (TS)* (Thompson, 1933; Kaufmann et al., 2012; Agrawal & Goyal, 2012, 2017a) are two  
26 leading algorithmic families for addressing stochastic bandit problems. For semi-bandit settings, the  
27 minimax lower bound is established as  $\Omega(\sqrt{mNT})$  (Kveton et al., 2015; Merlis & Mannor, 2020), and  
28 UCB-based algorithms achieve an upper bound of  $O(\sqrt{mNT \ln T})$  (Kveton et al., 2015). Although  
29 TS-based algorithms have been analyzed for problem-dependent bounds in semi-bandits (Wang &  
30 Chen, 2018; Perrault et al., 2020), their results cannot be simply extended to reasonable problem-  
31 independent bounds because their bounds contain constant terms that grow exponentially with  $m$ .

32 While a substantial body of literature has explored the setting of sleeping semi-bandits (Hu et al.,  
 33 2019; Li et al., 2019; Wu & Li, 2024) using *upper confidence bound (UCB)*-based algorithms with an  
 34 upper bound of  $O(\sqrt{mNT \ln T})$ , the upper and lower bounds for *Thompson sampling (TS)*-based  
 35 algorithms for (sleeping) semi-bandits still remain an open problem. Since TS is highly competitive  
 36 with advanced UCB-based algorithms and widely used in large-scale applications (Chapelle & Li,  
 37 2011), investigating the theoretical performance of TS-based algorithms is crucial.

38 This work addresses long-standing gaps in the literature by introducing two algorithms with provable  
 39 theoretical guarantees. The first algorithm, CTS-G, is an adaptation of TS with Gaussian priors  
 40 specifically designed for sleeping semi-bandits, achieving an upper bound of  $\tilde{O}(m\sqrt{NT})$ , where  
 41  $\tilde{O}$  hides the logarithmic factors, and a lower bound of  $\Omega(\sqrt{mNT \ln \frac{N}{m}})$ . We further introduce  
 42 CL-SG, which improves upon CTS-G both theoretically and practically by employing only a single  
 43 Gaussian sample, resulting in tighter bounds: an upper bound of  $\tilde{O}(\sqrt{mNT})$  and a lower bound of  
 44  $\Omega(\sqrt{mNT})$ . CL-SG is minimax-optimal up to logarithmic factors compared to the known lower  
 45 bound for combinatorial bandits (Kveton et al., 2015; Merlis & Mannor, 2020).

## 46 2 Gaussian Randomized Algorithms

47 We first present some notations specific to this section. Let  $n_{a,t} := \sum_{\tau=1}^{t-1} \mathbf{1}[a \in A_\tau]$  denote the  
 48 total number of times that base arm  $a \in [N]$  has been pulled at the beginning of round  $t$ . Let  
 49  $\hat{r}_{a,n_{a,t}} := \frac{\sum_{\tau=1}^{t-1} \mathbf{1}[a \in A_\tau] \cdot r_{a,\tau}}{n_{a,t}}$  denote the empirical mean of base arm  $a$  at the beginning of round  
 50  $t$ , which is the average of  $n_{a,t}$  i.i.d. random variables according to reward distribution  $p_a$ . Let  $\mathcal{F}_t$   
 51 collect all the actions and observed rewards up to the end of round  $t$

52 In Sec. 2.1, we present CTS-G, an algorithm enjoying  $\tilde{O}(m\sqrt{NT})$  and  $\Omega(\sqrt{mNT \ln \frac{N}{m}})$  upper  
 53 and lower regret bounds. In Sec. 2.2, we present CL-SG, an algorithm enjoying  $\tilde{O}(\sqrt{mNT})$  and  
 54  $\Omega(\sqrt{mNT})$  upper and lower regret bounds. The practical performance of both algorithms is discussed  
 55 in Appendix A, and all the detailed proofs can be found in Appendix C to D.

### 56 2.1 Combinatorial Thompson Sampling with Gaussian Priors (CTS-G)

57 CTS-G presented in Alg. 1 is a direct adaptation of TS with Gaussian priors (Agrawal & Goyal, 2017b)  
 58 to the sleeping semi-bandit problems. The core idea is to use posterior distributions to model the mean  
 59 reward  $r_a$  of each base arm  $a \in [N]$ . In each round  $t$ , CTS-G draws a Gaussian posterior sample  
 60  $w_{a,t} \sim \mathcal{N}(\hat{r}_{a,n_{a,t}}, \frac{\gamma m \ln t}{n_{a,t} + 1})$  for each  $a \in [N]$ , where  $\gamma > 0$  is a constant parameter to control the  
 61 exploration level.<sup>1</sup> We can view the collection  $\mathbf{w}_t = \{w_{a,t}, \forall a \in [N]\}$  of all posterior samples as the  
 62 “sampled problem instance” based on which the learning agent conducts learning in round  $t$ . Then,  
 63 based on the revealed feasible set  $\Theta_t$ , CTS-G plays the super arm  $A_t \in \arg \max_{A \in \Theta_t} \sum_{a \in A} w_{a,t}$   
 64 with the highest aggregated value of posterior samples and observes each individual base arm’s  
 65 random reward.

66 **Theorem 1.** (1) The regret of CTS-G is  $O(m \ln(T) \sqrt{NT})$ . (2) There exists a problem instance  
 67 such that CTS-G suffers  $\Omega(\sqrt{mNT \ln(\frac{N}{m})})$  regret.

68 **Discussion.** Theorem 1 states that CTS-G is worst-case optimal up to an extra  $\ln(T) \sqrt{m}$  factor.  
 69 Compared with UCB-based algorithms for sleeping semi-bandits, our upper bound has an extra  
 70 factor of  $\sqrt{m \ln T}$  with the ones by Hu et al. (2019); Li et al. (2019), which are  $O(\sqrt{mNT \ln T})$ .  
 71 However, it is important to note a significant aspect of our model: unlike the assumptions in Hu et al.  
 72 (2019); Li et al. (2019), our bound is derived without relying on stochastic assumptions regarding  
 73 the availability of arms. Furthermore, the upper bound is minimax optimal up to an extra  $\ln(T) \sqrt{m}$   
 74 factor as compared to the  $\Omega(\sqrt{mNT})$  minimax lower bound for combinatorial bandits shown in  
 75 Merlis & Mannor (2020).

<sup>1</sup>In practice, we only need to draw posterior samples for available arms to improve efficiency.

76 **Upper bound proof sketch.** The theoretical analysis is non-trivial due to overlapping base arms  
77 among super arms, the dynamic nature of the optimal super arm  $A_t^*$ , and its unobservability, as only the  
78 played super arm  $A_t$  is observed in each round  $t$ . To decompose the regret, we define a high probability  
79 event for the empirical estimates. Let  $\mathcal{E}_t := \left\{ |r_a - \hat{r}_{a,n_{a,t}}| \leq \sqrt{\frac{3 \ln(Nt)}{n_{a,t}+1}}, \forall a \in [N] \right\}$  be the  
80 event that the empirical means are close to their true means by the beginning of round  $t$ . Let  
81  $t' = \max\{\sqrt{m}, 4\}$  and  $\mathbf{E}_{\Theta_t}[\cdot] := \mathbf{E}[\cdot \mid \Theta_t]$ . Then, we decompose the regret defined in (1) as

$$\mathcal{R}(T) \leq \underbrace{\sum_{t=t'}^T \mathbf{E} \left[ \sum_{a \in A_t^*} r_a - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \right]}_{=: I_1, \text{ optimism term}} + \underbrace{\sum_{t=t'}^T \mathbf{E} \left[ \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} (w_{a,t} - r_a) \mathbf{1}[\mathcal{E}_t] \right] \right]}_{=: I_2, \text{ deviation term}} + mt' + O(1).$$

82 The deviation term  $I_2$  is easy to analyze as we can observe  $A_t$ , and is upper bounded by  $\tilde{O}(m\sqrt{NT})$   
83 via using concentration bounds. The center question is how to upper bound the optimism term, which  
84 measures the gap between the *maximum amount of true reward*  $\sum_{a \in A_t^*} r_a$  the learning agent could  
85 achieve and the *expected maximum amount of reward*  $\sum_{a \in A_t} w_{a,t}$  the learning agent can observe  
86 in round  $t$ . Intuitively, if the learning agent is lucky, i.e., the history  $\mathcal{F}_{t-1}$  gives  $\sum_{a \in A_t^*} r_a \leq$   
87  $\mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right]$ , there is no regret in round  $t$  for this term. Let  $(\cdot)^+ := \max\{\cdot, 0\}$  be an  
88 activation function. Then, we have

$$\sum_{a \in A_t^*} r_a - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \leq \left( \sum_{a \in A_t^*} r_a - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \right)^+. \quad (2)$$

89 Let  $c(\gamma)$  be a constant only depending on  $\gamma$ . In our novel technical Lemma 1, inspired by Russo  
90 (2019), we show

$$\left( \sum_{a \in A_t^*} r_a - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \right)^+ \leq c(\gamma) \cdot \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \right)^+ \right], \quad (3)$$

91 which tackles the challenge brought by the unobservability of  $A_t^*$ .

92 Next, via introducing an independent ‘‘ghost’’ copy  $\tilde{w}_{a,t} \sim \mathcal{N}(\hat{r}_{a,n_{a,t}}, \frac{\gamma m \ln t}{n_{a,t}+1})$  of  $w_{a,t}$ , we show

$$\mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \right)^+ \right] \leq \mathbf{E}_{\Theta_t} \left[ \left| \sum_{a \in A_t} (w_{a,t} - \tilde{w}_{a,t}) \right| \right], \quad (4)$$

93 which gets rid of the introduced activation function.

94 Since  $w_{a,t} - \tilde{w}_{a,t} \sim \mathcal{N}\left(0, \frac{2\gamma m \ln t}{n_{a,t}+1}\right)$ , we only need to deal with Gaussian random variables and  
95 have

$$\sum_{t=t'}^T \mathbf{E} \left[ \left| \sum_{a \in A_t} (w_{a,t} - \tilde{w}_{a,t}) \right| \right] \leq O(m \ln T \sqrt{\gamma NT}). \quad (5)$$

96

97 **Lower bound proof sketch.** Inspired by Theorem 1.4 in Agrawal & Goyal (2017b), the lower bound  
98 is refined by constructing a path selection problem with  $N$  links (base arms) and  $K$  paths (super  
99 arms) of  $m$  links. This reduces the semi-bandits to  $K$  independent path selections, and the result  
100 follows using the anti-concentration inequality for Gaussian variables (Appendix C.7)

## 101 2.2 Combinatorial Learning with Single Gaussian Seed (CL-SG)

102 Since the upper bound of CTS-G still has an extra  $\ln(T)\sqrt{m}$  factor from the minimax lower  
103 bound Merlis & Mannor (2020) for combinatorial bandits, we are motivated to improve the up-  
104 per bound by controlling the amount of randomness injected within the learning algorithm.

105 Inspired by Xiong et al. (2022), we devise CL-SG which enjoys a  $\tilde{O}(\sqrt{mNT})$  regret bound. The  
106 key idea behind the removal of the extra  $\sqrt{m}$  factor as compared to the regret of CTS-G (Alg. 1)

107 is CL-SG uses a single random seed  $w_t \sim \mathcal{N}(0, 1)$  to perturb the empirical estimates of all the  
 108 base arms, as shown in Alg. 2. After drawing  $w_t$ , we construct  $\bar{r}_{a,t} = \hat{r}_{a,n_{a,t}} + w_t \cdot \sqrt{\frac{\gamma \ln t}{n_{a,t}+1}}$   
 109 for all the base arms  $a \in [N]$ , where constant  $\gamma > 0$  controls the exploration level. Then, we play  
 110  $A_t = \arg \max_{A \in \Theta_t} \sum_{a \in A} \bar{r}_{a,t}$  from the feasible set  $\Theta_t$  in round  $t$ .

111 **Theorem 2.** (1) The regret of CL-SG is  $O\left(\ln T \sqrt{mNT}\right)$ . (2) There exists a problem instance such  
 112 that CL-SG suffers  $\Omega\left(\sqrt{mNT}\right)$  regret.

113 **Discussion.** The extra  $\sqrt{m}$  in CTS-G comes from the  $m$  factor in the variance of the Gaussian  
 114 posterior sample  $w_{a,t}$ , necessary to keep  $c(\gamma)$  bounded by a constant. To bound  $c(\gamma)$ , we must  
 115 lower bound  $\Pr_{\Theta_t} \left( \sum_{a \in A_t^*} w_{a,t} - \hat{r}_{a,n_{a,t}} \geq \sum_{a \in A_t^*} \sqrt{\frac{4 \ln t}{n_{a,t}+1}} \right)$ , requiring the Cauchy-Schwarz  
 116 inequality to bring the summation inside the square root for the RHS term in the probability, which  
 117 scales with  $\sqrt{m}$ , i.e.,  $\sum_{a \in A_t^*} \sqrt{\frac{4 \ln t}{n_{a,t}+1}} \leq \sqrt{m} \sum_{a \in A_t^*} \frac{4 \ln t}{n_{a,t}+1}$ . This fact further results in an  
 118 extra  $m$  in the variance of CTS-G Gaussian samples for the probability to be lower bounded  
 119 by a constant. On the other hand, with CL-SG, using a single  $w_t$ , we lower bound a similar  
 120 probability,  $\Pr \left( \sum_{a \in A_t^*} w_t \sqrt{\frac{\gamma \ln t}{n_{a,t}+1}} \geq \sum_{a \in A_t^*} \sqrt{\frac{4 \ln t}{n_{a,t}+1}} \right)$ , allowing us to divide both sides by  
 121  $\sum_{a \in A_t^*} \sqrt{\frac{4 \ln t}{n_{a,t}+1}}$  and avoid the extra  $m$  in the variance.

122 The lower-bound proof considers the same problem instance to Theorem 1 but differs in addressing  
 123 the arms' dependency in CL-SG due to the common  $w_t$ . Let  $\Delta := \sqrt{\frac{N}{mT}}$  be the reward gap between  
 124 each suboptimal arm and the optimal super arm, and let  $Q_A(t)$  be the number of times that super  
 125 arm  $A$  has been played at the beginning of round  $t$ . Then, denote by  $B_t^* := \{Q_{A_1}(t) > t - cT\}$   
 126 the event that the optimal super arm  $A_1$  has been observed enough times at the beginning of  
 127 round  $t$ , where  $c \in (0, 1)$  is a constant. It is easy to prove that the regret is lower bounded  
 128 by  $cT \cdot m \cdot \Delta = \Omega(\sqrt{mNT})$  when  $B_t^*$  is false for some  $t \in [T]$ . The main challenge is  
 129 to demonstrate that, conditioned on the past histories  $F_{t-1}$  that lead to the happening of event  
 130  $B_t^*$ , the probability of playing a suboptimal super arm is at least a constant probability  $p_0$ , i.e.,  
 131  $\Pr(\exists A \in \Theta \setminus A_1 : A_t = A \mid \mathcal{F}_{t-1} = F_{t-1}) \geq p_0$ . This leads to lower bound a probability  
 132 that the empirical estimates of suboptimal arms are larger than the optimal super arm, i.e.,

$$\begin{aligned} & \Pr \left( \exists A \in \Theta \setminus A_1 : \sum_{a \in A} \hat{r}_{a,Q_A(t)} + w_t \sqrt{\frac{\gamma \ln t}{Q_A(t) + 1}} > \sum_{b \in A_1} \hat{r}_{b,Q_{A_1}(t)} + w_t \sqrt{\frac{\gamma \ln t}{Q_{A_1}(t) + 1}} \mid \mathcal{F}_{t-1} = F_{t-1} \right) \\ & \geq \Pr \left( \exists A \in \Theta \setminus A_1 : w_t \left( 1 - \sqrt{\frac{Q_A(t) + 1}{Q_{A_1}(t) + 1}} \right) > \Delta \sqrt{Q_A(t) + 1} \mid \mathcal{F}_{t-1} = F_{t-1} \right). \end{aligned}$$

133 This requires analyzing the play ratio between suboptimal and optimal super arms, while in the  
 134 lower-bound analysis of Theorem 1, we can avoid this situation by independently considering the  
 135 estimates of the optimal super arm is smaller than 0, and that of suboptimal arms is larger than 0.  
 136 The trick to address this ratio is to only consider the regret from  $\alpha T$  to  $T$ , with  $\alpha \in (0, 1)$  such that  
 137  $\frac{Q_A(t)+1}{Q_{A_1}(t)+1} \leq \frac{cT+1}{(\alpha-c)T+1}$  is a constant by tuning  $c$  and  $\alpha$ . Then, by applying the anti-concentration  
 138 bound for Gaussian variables and solving a non-trivial optimization problem, we can prove such a  $p_0$   
 139 exists, and regret is lower bounded by  $(1 - \alpha)T \cdot p_0 \cdot m\Delta = \Omega(\sqrt{mNT})$ .

### 140 3 Conclusion

141 In this paper, we have studied the problem of sleeping semi-bandits and presented CTS-G and CL-SG  
 142 with theoretical guarantees. Our results bridge the existing gap in the literature by providing upper and  
 143 lower bounds for TS-based algorithms in sleeping semi-bandits. Future work will focus on narrowing  
 144 the gap between these bounds, and studying the relationship between the number of random variables  
 145 and their variances.

146 **References**

- 147 Rajeev Agrawal. Sample Mean Based Index Policies by  $O(\log n)$  Regret for The Multi-armed Bandit  
148 Problem. *Advances in Applied Probability*, 27(4):1054–1078, 1995.
- 149 Shipra Agrawal and Navin Goyal. Analysis of Thompson Sampling for The Multi-armed Bandit  
150 Problem. In *Proc. Conference on Learning Theory (COLT)*, volume 23, pp. 39.1–39.26. PMLR,  
151 2012.
- 152 Shipra Agrawal and Navin Goyal. Near-optimal Regret Bounds for Thompson Sampling. *Journal of*  
153 *ACM*, 64(5):30:1–30:24, September 2017a. ISSN 0004-5411.
- 154 Shipra Agrawal and Navin Goyal. Near-optimal Regret Bounds for Thompson Sampling. <http://www.columbia.edu/~sa3305/papers/j3-corrected.pdf>, 2017b.
- 156 Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time Analysis of The Multiarmed Bandit  
157 Problem. *Machine Learning*, 47(2–3):235–256, 2002.
- 158 Olivier Chapelle and Lihong Li. An Empirical Evaluation of Thompson Sampling. In *Proc. Advances*  
159 *in Neural Information Processing Systems (NeurIPS)*, pp. 2249–2257, 2011.
- 160 Bingshan Hu, Yunjin Chen, Zhiming Huang, Nishant A. Mehta, and Jianping Pan. Intelligent Caching  
161 Algorithms in Heterogeneous Wireless Networks with Uncertainty. In *Proc. IEEE Conference on*  
162 *Distributed Computing Systems (ICDCS)*. IEEE, 2019.
- 163 Emilie Kaufmann, Nathaniel Korda, and Rémi Munos. Thompson Sampling: An Asymptotically  
164 Optimal Finite-time Analysis. In *Proc. International Conference on Algorithmic Learning The-*  
165 *ory (ALT)*, pp. 199–213. Springer, 2012.
- 166 Branislav Kveton, Zheng Wen, Azin Ashkan, and Csaba Szepesvari. Tight Regret Bounds for  
167 Stochastic Combinatorial Semi-Bandits. In *Proc. Artificial Intelligence and Statistics (AISTATS)*,  
168 pp. 535–543, 2015.
- 169 Fengjiao Li, Jia Liu, and Bo Ji. Combinatorial Sleeping Bandits with Fairness Constraints. *IEEE*  
170 *Transactions on Network Science and Engineering (TNSE)*, 7(3):1799–1813, 2019.
- 171 Nadav Merlis and Shie Mannor. Tight Lower Bounds for Combinatorial Multi-armed Bandits. In  
172 *Proc. Conference on Learning Theory (COLT)*, pp. 2830–2857. PMLR, 2020.
- 173 Pierre Perrault, Etienne Boursier, Michal Valko, and Vianney Perchet. Statistical Efficiency of  
174 Thompson Sampling for Combinatorial Semi-bandits. *Proc. Advances in Neural Information*  
175 *Processing Systems (NeurIPS)*, 33:5429–5440, 2020.
- 176 Daniel Russo. Worst-case Regret Bounds for Exploration via Randomized Value Functions. *Proc. Ad-*  
177 *vances in Neural Information Processing Systems*, 32, 2019.
- 178 William R Thompson. On the Likelihood that One Unknown Probability Exceeds Another in View  
179 of The Evidence of Two Samples. *Biometrika*, 25(3-4):285–294, 1933.
- 180 Siwei Wang and Wei Chen. Thompson Sampling for Combinatorial Semi-Bandits. In *Proc. Interna-*  
181 *tional Conference on Machine Learning (ICML)*, pp. 5101–5109, 2018.
- 182 Xiaoyi Wu and Bin Li. Achieving Regular and Fair Learning in Combinatorial Multi-Armed Bandit.  
183 In *Proc. IEEE Conference on Computer Communications (INFOCOM)*, pp. 361–370. IEEE, 2024.
- 184 Zhihan Xiong, Ruoqi Shen, Qiwen Cui, Maryam Fazel, and Simon Shaolei Du. Near-optimal  
185 Randomized Exploration for Tabular Markov Decision Processes. In *Proc. Advances in Neural*  
186 *Information Processing Systems*, 2022.

187 **A Numerical Experiments**

188 **A.1 Algorithm Description for CTS-G and CL-SG**

189 The algorithm descriptions for CTS-G and CL-SG are presented in Algs. 1 and 2, respectively.

---

**Algorithm 1** Combinatorial Thompson Sampling with Gaussian Priors (CTS-G)

---

**Require:** arm set  $[N]$ , exploration rate  $\gamma$   
Initialize  $n_{a,1} = 0$  and  $\hat{r}_{a,n_{a,1}} = 0$  for all base arms  $a \in [N]$   
**for**  $t = 1, 2, \dots$  **do**  
  Observe feasible set  $\Theta_t$   
  Draw  $w_{a,t} \sim \mathcal{N}(\hat{r}_{a,n_{a,t}}, \frac{\gamma m \ln t}{n_{a,t}+1})$  for each base arm  $a \in [N]$   
  Play super arm  $A_t = \arg \max_{A \in \Theta_t} \sum_{a \in A} w_{a,t}$   
  Observe  $r_{a,t} \sim p_a$  for all base arms  $a \in A_t$  and update  $n_{a,t}$  and  $\hat{r}_{a,n_{a,t}}$  for all  $a \in A_t$ .  
**end for**

---

---

**Algorithm 2** Combinatorial Learning with Single Gaussian Seed (CL-SG)

---

**Require:** arm set  $[N]$ , exploration rate  $\gamma$   
Initialize  $n_{a,1} = 0$  and  $\hat{r}_{a,0} = 0$  for all base arms  $a \in [N]$   
**for**  $t = 1, \dots$  **do**  
  Observe decision set  $\Theta_t$   
  Draw  $w_t \sim \mathcal{N}(0, 1)$   
  Construct  $\bar{r}_{a,t} = \hat{r}_{a,n_{a,t}} + w_t \cdot \sqrt{\frac{\gamma \ln t}{n_{a,t}+1}}$  for all base arms  $a \in [N]$   
  Play super arm  $A_t = \arg \max_{A \in \Theta_t} \sum_{a \in A} \bar{r}_{a,t}$   
  Observe  $r_{a,t} \sim p_a$  for all base arms  $a \in A_t$  and update  $n_{a,t}$  and  $\hat{r}_{a,n_{a,t}}$  for all  $a \in A_t$ .  
**end for**

---

190 **A.2 Combinatorial Learning with Least Gaussian Seed (CL-LG)**

191 We aim to explore whether further reducing the number of Gaussian samples in the algorithm can  
192 enhance the practical performance. To this end, we propose the *Combinatorial Learning with Least*  
193 *Gaussian Seed (CL-LG)* algorithm, as shown in Alg. 3. Different from CL-SG (see Alg. 2), which  
194 requires an independent Gaussian sample in each round, our approach only draws a single Gaussian  
195 sample  $w \sim \mathcal{N}(0, 1)$  at the beginning of the game.

---

**Algorithm 3** Combinatorial Learning with Least Gaussian Seed (CL-LG)

---

**Require:** arm set  $[N]$ , exploration rate  $\gamma$   
Initialize  $n_{a,1} = 0$  and  $\hat{r}_{a,0} = 0$  for all base arms  $a \in [N]$   
Draw  $w \sim \mathcal{N}(0, 1)$   
**for**  $t = 1, \dots$  **do**  
  Observe feasible set  $\Theta_t$   
  Construct  $\bar{r}_{a,t} = \hat{r}_{a,n_{a,t}} + w \cdot \sqrt{\frac{\gamma \ln t}{n_{a,t}+1}}$  for all base arms  $a \in [N]$   
  Play super arm  $A_t = \arg \max_{A \in \Theta_t} \sum_{a \in A} \bar{r}_{a,t}$   
  Observe  $r_{a,t} \sim p_a$  for all base arms  $a \in A_t$  and update  $n_{a,t}$  and  $\hat{r}_{a,n_{a,t}}$  for all  $a \in A_t$ .  
**end for**

---

196 **A.3 Experiment Settings**

197 We conduct experiments in two settings to show the performance of the proposed algorithms with  
198  $\gamma = 0.1$  to study the number of Gaussian seeds and the impact of different  $\gamma$ , which can be found in  
199 Appendix A.4 and A.5. All the experiment results are the average of 100 independent experiments  
200 conducted on a MacBook Pro with M1 Max and 32GB RAM using Numpy.

201 In Setting 1, we consider a simple environment with  $N = 10$  arms, and at most  $m = 3$  arms can  
202 be played in each round. The actual rewards for all the arms follow the Bernoulli distributions, while  
203 the first three arms have a mean reward of 0.9, and the rest of the arms have a mean reward of 0.8.  
204 In Setting 2, we consider a more complicated setting where  $N = 50$  and  $m = 15$ . In this setting,  
205 rewards are again based on Bernoulli distributions, where the first five arms have mean rewards

206 generated uniformly from  $[0.725, 0.75]$ , and the rest of the arms have mean rewards generated  
 207 uniformly from  $[0.7, 0.725]$ . For both settings, the availability of each arm is determined by a  
 208 Bernoulli distribution with a mean of 0.5. The reason we chose Bernoulli distributions for the  
 209 rewards is that we want to compare with the following CTS-B (which requires Bernoulli rewards)  
 210 and CombUCB algorithms. Both algorithms play arms  $A_t := \arg \max_{A \in \Theta_t} \sum_{a \in A} \theta_{a,t}$ , where  
 211  $\theta_{a,t}$  is defined differently as follows.

212 • CTS-B (Wang & Chen, 2018): In each round  $t$ , CTS-B draws random samples from Beta distribu-  
 213 tions for each available arm  $\theta_{a,t} \sim \mathbf{Beta}(\hat{r}_{a,n_{a,t}} n_{a,t} + 1, n_{a,t} - n_{a,t} \hat{r}_{a,n_{a,t}} + 1)$ , and plays  
 214 arms  $A_t := \arg \max_{A \in \Theta_t} \sum_{a \in A} \theta_{a,t}$ .

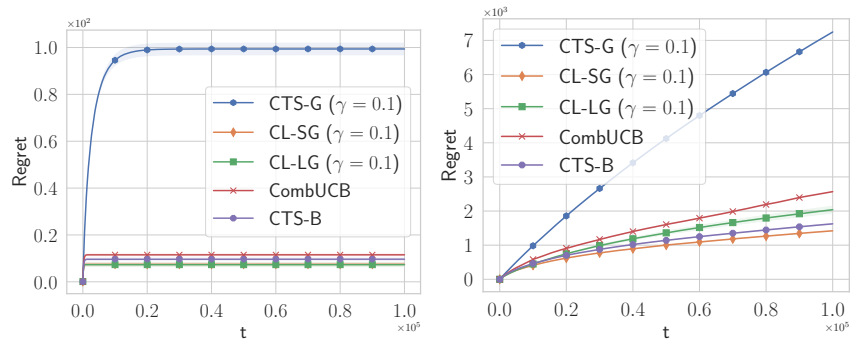
215 • CombUCB (Kveton et al., 2015): In each round  $t$ , CombUCB estimates the UCB values for each  
 216 arm  $\theta_{a,t} = \hat{r}_{a,n_{a,t}} + \sqrt{\frac{1.5 \ln t}{n_{a,t}}}$  in each round  $t$ .

#### 217 A.4 Impact of Number of Gaussian Seed

218 The regret results over  $T = 10^5$  rounds are shown in Fig. 1 with 97.5% confidence intervals.  
 219 In both settings, CTS-G performs worse than others, suffering the highest regret, because of the  
 220 algorithm’s reliance on Gaussian random samples, which are unbounded and result in an excessive  
 221 exploration rate. This overemphasis on exploration, at the expense of exploiting known rewarding  
 222 arms, fundamentally undermines the algorithm’s efficiency.

223 On the other hand, CL-SG demonstrates comparable performance to CL-LG in Setting 1, both  
 224 outperforming CTS-B. In Setting 2, CL-SG maintains its advantage, whereas CL-LG falls behind  
 225 CTS-G. This highlights the effectiveness of CL-SG’s design in optimizing the exploration-exploitation  
 226 trade-off more efficiently than its counterparts.

227 Notably, in both settings, CL-LG with  $\gamma = 0.1$  outperforms CombUCB, suggesting that the initial  
 228 randomness incorporated in CL-LG helps balance the trade-off between exploration and exploitation.  
 229 It remains an open question of how initial randomness helps.



(a) Setting 1 ( $N = 10$  and  $m = 3$ )      (b) Setting 2 ( $N = 50$  and  $m = 15$ )

Figure 1: The comparison of regret for both settings with  $\gamma = 0.1$ .

#### 230 A.5 Impact of Different $\gamma$

231 We performed experiments with the CTS-G, CL-SG, and CL-LG algorithms under Settings 1 and 2.  
 232 The experiments utilized  $\gamma$  values of 0.01, 0.1, 0.5, and 1. The results are illustrated in Figs. 2 and 3.

233 Regarding Setting 1, CTS-G performs worse as  $\gamma$  increases, as shown in Fig. 2a, because a higher  $\gamma$   
 234 corresponds to a higher exploration rate, which will over-explore the simple scenario. For CL-SG,  
 235 the performance with  $\gamma = 0.1$  is better than that with other  $\gamma$  values. CL-LG achieves the best  
 236 performance with 0.5, which indicates the performance of algorithms is not necessarily linear with  $\gamma$ .

237 When comparing the algorithms at their optimal  $\gamma$  values (see Fig. 2d), CTS-G shows the worst  
 238 performance, whereas CL-SG performs comparably to CL-LG.

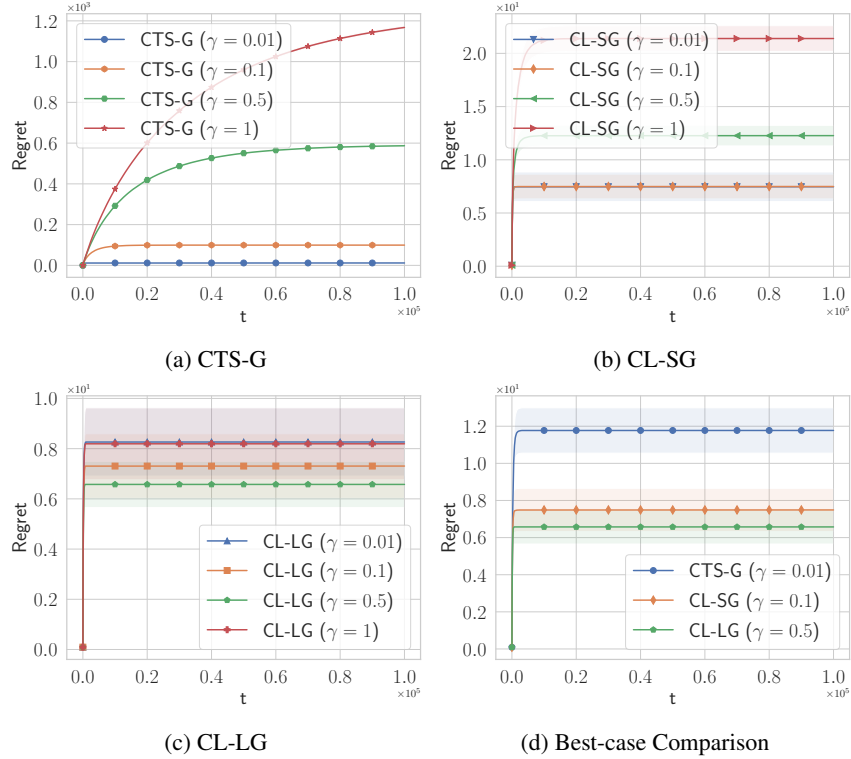


Figure 2: The comparison of different  $\gamma$  for CTS-G, CL-SG and CL-LG in Setting 1.

239 Additionally, CTS-G is very sensitive to the change of  $\gamma$ , and the performance of CTS-G with  $\gamma = 1$   
 240 is about 200 times worse than that of CTS-G with  $\gamma = 0.01$ . In contrast, CL-SG and CL-LG  
 241 demonstrate greater robustness to changes in  $\gamma$ , showing that fewer Gaussian samples may prevent  
 242 over-exploration.

243 Regarding the more complicated Setting 2, we can observe a change in CL-LG, where  $\gamma = 1$  leads  
 244 to the worst performance, while  $\gamma = 0.5$  achieves the best performance. This indicates that CL-LG  
 245 with  $\gamma = 1$  will over-explore. When comparing all the algorithms with their optimal  $\gamma$  values, we  
 246 can see that CL-SG with  $\gamma = 0.1$  achieves the best performance. More interestingly, we can observe  
 247 that algorithms with fewer Gaussian samples require higher  $\gamma$  to achieve better performance.

248 From this experiment, we can see that different Gaussian samples react differently to different  
 249 exploration rates. This observation raises an intriguing question for future research: what is the  
 250 relationship between the number of random variables and their variance, and what is the optimal  
 251 combination to achieve the best results?

## 252 A.6 Tightness of regret bound

253 We consider a setting of 100 arms, and at most 10 arms can be played in each round. The mean  
 254 rewards for the first 10 arms are 0.925, and the mean rewards for the rest suboptimal arms are 0.9.

255 We compare the regret of CTS-G with the lower regret bound  $0.1\sqrt{mNT} \ln(\frac{N}{m})$  in Fig. 4a. As we  
 256 can see, there are still gaps between the actual performance and the theoretical lower bound, and the  
 257 increasing rate of CTS-G is larger than the lower bound, which indicates that the lower bound may  
 258 still have room to be improved.

259 Similarly, we compared CL-SG with the lower bound of  $0.1\sqrt{mNT}$ , and we can see that the regret  
 260 of CL-SG increases faster than the lower bound, indicating that the lower bound can be improved.



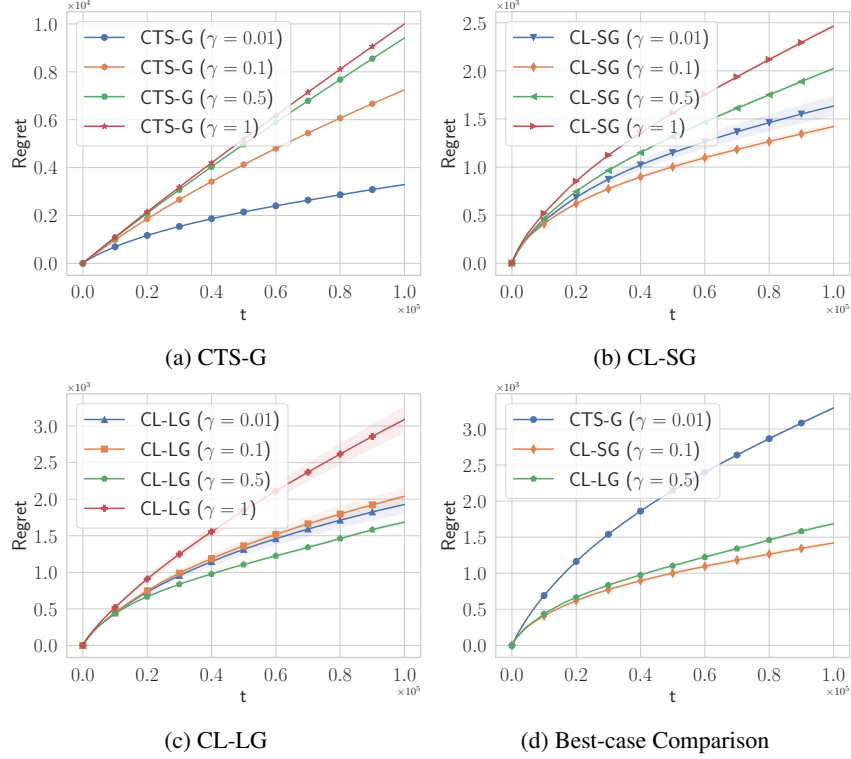
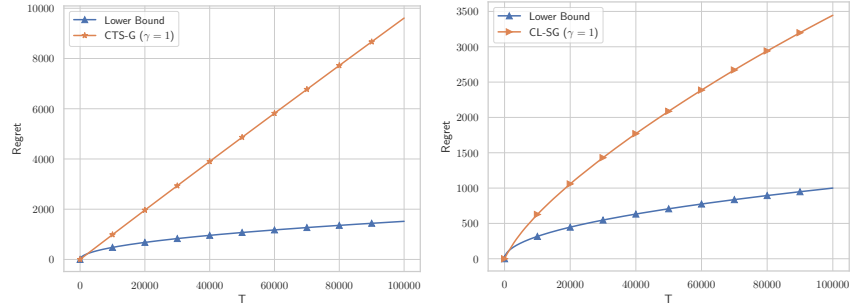


Figure 3: The comparison of different  $\gamma$  for CTS-G, CL-SG and CL-LG in Setting 2.



(a) Tightness of the lower bound for CTS-G (b) Tightness of the lower bound for CL-SG

Figure 4: Tightness of Regret Bound for both CTS-G and CL-SG.

## 261 B Notations and Facts

262 **Notations:** Let  $\mathcal{F}_{t-1}$  denote by the history of past actions and rewards until the end of  
 263 round  $t - 1$ . Recall that  $\mathbf{E}_{\Theta_t}[\cdot] := \mathbf{E}[\cdot \mid \Theta_t]$  and  $\Pr_{\Theta_t}(\cdot) := \Pr(\cdot \mid \Theta_t)$ . Denote by  $\mathcal{E}_t :=$   
 264  $\left\{ \forall a \in [N] : |r_a - \hat{r}_{a,n_{a,t}}| \leq \sqrt{\frac{3 \ln Nt}{n_{a,t}+1}} \right\}$  the high-probability event that the empirical mean  
 265 is close to the true mean reward for arm  $a$ , and by  $\bar{\mathcal{E}}_t$  the complementary event of  $\mathcal{E}_t$ . Recall that  
 266  $\tilde{w}_{a,t} \sim \mathcal{N}(\hat{r}_{a,n_{a,t}}, \frac{\gamma m \ln t}{n_{a,t}+1})$  is i.i.d. of  $w_{a,t}$  for CTS-G, and  $\tilde{r}_t := \hat{r}_{a,n_{a,t}} + \tilde{w}_t \sqrt{\frac{\gamma \ln t}{n_{a,t}+1}}$ , where  
 267  $\tilde{w}_t \sim \mathcal{N}(0, 1)$  is i.i.d. of  $w_t$  for CL-SG.

268 **Fact 1.** For a Gaussian distributed random variable  $Z$  with mean  $\mu$  and variance  $\delta^2$ , for any  $z$ , we  
 269 have that

$$\frac{1}{4\sqrt{\pi}} \cdot e^{-7z^2/2} \leq \Pr(|Z - \mu| > z\sigma) \leq \frac{1}{2} e^{-z^2/2}, \quad (6)$$

270 and for any  $z > 0$ ,

$$\Pr(Z - \mu > z\sigma) \geq \frac{1}{\sqrt{2\pi}} \frac{z}{z^2 + 1} e^{-\frac{z^2}{2}}. \quad (7)$$

271 **Fact 2.** Let  $X_1, \dots, X_N$  be  $N$  real random variables with  $X_i \sim \text{subG}(\sigma^2)$ ,  $i = 1, \dots, N$ ,  
 272 not necessarily independent. Then,

$$\mathbb{E} \left[ \max_{i=1, \dots, N} |X_i| \right] \leq \sigma \sqrt{2 \log(2N)}.$$

## 273 C Proofs for Theorem 1

### 274 C.1 Proof of Lemma 1

275 **Lemma 1.** In any round  $t \geq \max\{\sqrt{m}, 4\}$ , the optimism part in CTS-G satisfies that

$$\mathbf{E} \left[ \sum_{t=\max\{\sqrt{m}, 4\}}^T \left( \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} w_{a,t} \right) \right] \leq 8\sqrt{3\gamma} \Phi(-\sqrt{4/\gamma})^{-1} m \ln T \sqrt{NT}. \quad (8)$$

276 *Proof.* For each  $a \in [N]$ , we let  $\tilde{w}_{a,t} \sim \mathcal{N}(\hat{r}_{a,n_{a,t}}, \frac{m\gamma \ln t}{n_{a,t}+1})$  be an independent copy of  $w_{a,t}$ .  
 277 Let  $(\cdot)^+ := \max\{\cdot, 0\}$ . Let  $\mathbf{w}$  collect all the Gaussian random variables  $w_{a,t}$  for all  $a \in [N]$ .  
 278 Recall that  $\mathbf{E}_{\Theta_t}[\cdot] := \mathbf{E}[\cdot \mid \Theta_t]$ . There are three steps for the proofs.

279 **Step 1:** we show that in each round  $t \geq \max\{\sqrt{m}, 4\}$ , we have

$$\mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} w_{a,t} \right] \leq 2\Phi(-\sqrt{4/\gamma})^{-1} \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \right)^+ \right]. \quad (9)$$

280 **Step 2:** we further bound the expectation term in the RHS of (9) as follows.

$$\mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \right)^+ \right] \leq \mathbf{E}_{\Theta_t} \left[ \left| \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} \tilde{w}_{a,t} \right| \right]. \quad (10)$$

281 **Step 3:** summing over  $T$ , we show that (10) is upper bounded as follows.

$$\mathbf{E} \left[ \sum_{t=1}^T \left| \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} \tilde{w}_{a,t} \right| \right] \leq 4m \ln T \sqrt{3\gamma NT} \quad (11)$$

282 Combining these three steps, we have

$$\begin{aligned} & \mathbf{E} \left[ \sum_{t=\max\{\sqrt{m}, 4\}}^T \left( \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} w_{a,t} \right) \right] \\ & \leq 2\Phi(-\sqrt{4/\gamma})^{-1} \mathbf{E} \left[ \sum_{t=\max\{\sqrt{m}, 4\}}^T \left| \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} \tilde{w}_{a,t} \right| \right] \\ & \leq 2\Phi(-\sqrt{4/\gamma})^{-1} \mathbf{E} \left[ \sum_{t=1}^T \left| \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} \tilde{w}_{a,t} \right| \right] \\ & \leq 8\sqrt{3\gamma} \Phi(-\sqrt{4/\gamma})^{-1} m \ln T \sqrt{NT}. \end{aligned} \quad (12)$$

283 Now, we give the details for these three steps.

284 Let  $\alpha := \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} w_{a,t} \right]$ .

285 **Step 1 proof.** If  $\alpha = \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} w_{a,t} \right] \leq 0$ , the proof is trivial as the RHS of  
 286 (9) is non-negative. Note that  $2\Phi(-\sqrt{4/\gamma})^{-1} < +\infty$

287 For the case where  $\alpha > 0$ , we view  $(\sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} [\sum_{a \in A_t} w_{a,t}])^+ \geq 0$  as a non-  
 288 negative random variable and use Markov's inequality. We have

$$\mathbf{E}_{\Theta_t} \left[ (\sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} [\sum_{a \in A_t} w_{a,t}])^+ \right] \geq \alpha \Pr_{\Theta_t} (\sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} [\sum_{a \in A_t} w_{a,t}] \geq \alpha), \quad (13)$$

289 which gives

$$\begin{aligned} \alpha &\leq \frac{\mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \right)^+ \right]}{\Pr_{\Theta_t} \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \geq \alpha \right)} \\ &= \frac{\mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \right)^+ \right]}{\Pr_{\Theta_t} \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \geq \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} w_{a,t} \right] \right)} \\ &= \frac{\mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \right)^+ \right]}{\Pr_{\Theta_t} \left( \sum_{a \in A_t} w_{a,t} \geq \sum_{a \in A_t^*} r_a \right)} \stackrel{(a)}{\leq} \frac{\mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \right)^+ \right]}{\Pr_{\Theta_t} \left( \sum_{a \in A_t^*} w_{a,t} \geq \sum_{a \in A_t^*} r_a \right)} \\ &\stackrel{(b)}{\leq} 2\Phi(-\sqrt{4/\gamma})^{-1} \cdot \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \right)^+ \right], \end{aligned} \quad (14)$$

290 where step (a) is due to that  $A_t$  is the optimal super arm, and thus, we have  $\sum_{a \in A_t^*} w_{a,t} \leq$   
 291  $\sum_{a \in A_t} w_{a,t}$  and step (b) uses the result shown in Lemma 4.

292 **Step 2 proof.** Recall that  $w_{a,t}$  and  $\tilde{w}_{a,t}$  are i.i.d. according to  $\mathcal{N} \left( \hat{r}_{a,n_{a,t}}, \frac{m\gamma \ln t}{n_{a,t}+1} \right)$ ,  
 293 and  $A_t$  is the optimal super arm based on  $\Theta_t$  and  $\mathbf{w}$ . We have  $\mathbf{E}_{\Theta_t} [\sum_{a \in A_t} w_{a,t}] =$   
 294  $\mathbf{E}_{\Theta_t} [\max_{A \in \Theta_t} \sum_{a \in A} w_{a,t}] = \mathbf{E}_{\Theta_t} [\max_{A \in \Theta_t} \sum_{a \in A} \tilde{w}_{a,t}] \geq \mathbf{E}_{\Theta_t} [\sum_{a \in A_t} \tilde{w}_{a,t} \mid A_t] =$   
 295  $\mathbf{E}_{\Theta_t} [\sum_{a \in A_t} \tilde{w}_{a,t} \mid A_t, \mathbf{w}]$ . Then, we have

$$\begin{aligned}
\mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} w_{a,t} \right] \right)^+ \right] &\leq \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \tilde{w}_{a,t} \mid A_t \right] \right)^+ \right] \\
&= \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \tilde{w}_{a,t} \mid A_t, \mathbf{w} \right] \right)^+ \right] \\
&= \mathbf{E}_{\Theta_t} \left[ \left( \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} \tilde{w}_{a,t} \right) \mid A_t, \mathbf{w} \right] \right)^+ \right] \\
&\leq \mathbf{E}_{\Theta_t} \left[ \left| \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} \tilde{w}_{a,t} \right) \mid A_t, \mathbf{w} \right] \right| \right] \\
&\leq \mathbf{E}_{\Theta_t} \left[ \left| \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} \tilde{w}_{a,t} \mid A_t, \mathbf{w} \right| \right] \\
&= \mathbf{E}_{\Theta_t} \left[ \left| \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} \tilde{w}_{a,t} \right| \right], \tag{15}
\end{aligned}$$

296 where the last inequality is due to Jensen's inequality.

297 **Step 3 proof.** Since  $w_{a,t} - \tilde{w}_{a,t} \sim \mathcal{N} \left( 0, \frac{2\gamma m \ln t}{n_{a,t}+1} \right)$ , we can express  $w_{a,t} - \tilde{w}_{a,t}$  as  $\sqrt{2}\zeta_{a,t}\delta_{a,t}$ ,  
298 where  $\zeta_{a,t} \sim \mathcal{N}(0, 1)$  and  $\delta_{a,t} = \sqrt{\frac{\gamma m \ln t}{n_{a,t}+1}}$ . Thus, we have

$$\begin{aligned}
\mathbf{E} \left[ \sum_{t=1}^T \left| \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} \tilde{w}_{a,t} \right| \right] &\leq \sqrt{2} \mathbf{E} \left[ \sum_{t=1}^T \sum_{a \in A_t} |\zeta_{a,t} \delta_{a,t}| \right] \\
&\stackrel{(a)}{\leq} \sqrt{2} \mathbf{E} \left[ \max_{t \in [T], a \in [N]} |\zeta_{a,t}| \sum_{t=1}^T \sum_{a \in A_t} |\delta_{a,t}| \right] \\
&= \sqrt{2} \mathbf{E} \left[ \max_{t \in [T], a \in [N]} |\zeta_{a,t}| \sum_{t=1}^T \sum_{a \in A_t} \sqrt{\frac{\gamma m \ln t}{n_{a,t}+1}} \right] \tag{16} \\
&\stackrel{(b)}{\leq} 2m \sqrt{2\gamma NT \ln T} \mathbf{E} \left[ \max_{t \in [T], a \in [N]} \zeta_{a,t} \right] \\
&\stackrel{(c)}{\leq} 2m \sqrt{2\gamma NT \ln T} \cdot \sqrt{6 \ln T} \\
&\leq 4m \ln T \sqrt{3\gamma NT}.
\end{aligned}$$

299 where step (a) is due to Hölder's inequality. Step (b) is due to Lemma 5 such that  
300  $\sum_{t=1}^T \sum_{a \in A_t} \sqrt{\frac{1}{n_{a,t}+1}} \leq 2\sqrt{mNT}$ . Step (c) is due to the maximal inequality for Gaussian  
301 variables (Fact 2) such that  $\mathbf{E} \left[ \max_{t \in [T], a \in [N]} \zeta_{a,t} \right] \leq \sqrt{2 \ln 2NT} \leq \sqrt{6 \ln T}$  because  $2 \leq$   
302  $N \leq T$ .  $\square$

### 303 C.2 Proof of Lemma 2

**Lemma 2.** Let  $\mathcal{E}_t := \left\{ \forall a \in [N], \hat{r}_{a,n_{a,t}} - r_a \leq \sqrt{\frac{3 \ln Nt}{n_{a,t}+1}} \right\}$ . In CTS-G, the regret of the deviation part is

$$\mathbf{E} \left[ \sum_{t=1}^T \left( \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} r_a \right) \mathbf{1}[\mathcal{E}_t] \right] \leq 2m \ln T \sqrt{6\gamma NT} + 2\sqrt{6mNT \ln T}.$$

304 *Proof.* We can do decomposition as follows.

$$\begin{aligned}
& \mathbf{E} \left[ \sum_{t=1}^T \left( \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} r_a \right) \mathbf{1}[\mathcal{E}_t] \right] \\
&= \mathbf{E} \left[ \sum_{t=1}^T \left( \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} \hat{r}_{a,n_{a,t}} + \sum_{a \in A_t} \hat{r}_{a,n_{a,t}} - \sum_{a \in A_t} r_a \right) \mathbf{1}[\mathcal{E}_t] \right] \\
&= \mathbf{E} \left[ \sum_{t=1}^T \left( \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} \hat{r}_{a,n_{a,t}} \right) \mathbf{1}[\mathcal{E}_t] \right] + \mathbf{E} \left[ \sum_{t=1}^T \left( \sum_{a \in A_t} \hat{r}_{a,n_{a,t}} - \sum_{a \in A_t} r_a \right) \mathbf{1}[\mathcal{E}_t] \right] \\
&\stackrel{(a)}{\leq} \mathbf{E} \left[ \sum_{t=1}^T \sum_{a \in A_t} (w_{a,t} - \hat{r}_{a,n_{a,t}}) \right] + \mathbf{E} \left[ \sum_{t=1}^T \sum_{a \in A_t} \sqrt{\frac{6 \ln T}{n_{a,t} + 1}} \right] \\
&\stackrel{(b)}{\leq} \mathbf{E} \left[ \sum_{t=1}^T \sum_{a \in A_t} (w_{a,t} - \hat{r}_{a,n_{a,t}}) \right] + 2\sqrt{6mNT \ln T},
\end{aligned} \tag{17}$$

305 where step (a) is because event  $\mathcal{E}_t$  is true and  $\ln NT \leq 2 \ln T$  because of  $N \leq T$ , and step (b) is  
306 due to Lemma 5 such that  $\sum_{t=1}^T \sum_{a \in A_t} \sqrt{\frac{1}{n_{a,t}+1}} \leq 2\sqrt{mNT}$ .

307 We can represent each  $w_{a,t} - \hat{r}_{a,n_{a,t}}$  by  $\zeta_{a,t}\delta_{a,t}$ , where  $\zeta_{a,t} \sim \mathcal{N}(0, 1)$  and  $\delta_{a,t} = \sqrt{\frac{\gamma m \ln t}{n_{a,t}+1}}$ .  
308 Then, we can bound the first term on the RHS of the above equation as follows:

$$\begin{aligned}
\mathbf{E} \left[ \sum_{t=1}^T \sum_{a \in A_t} (w_{a,t} - \hat{r}_{a,n_{a,t}}) \right] &\leq \mathbf{E} \left[ \sum_{t=1}^T \sum_{a \in A_t} \zeta_{a,t} \delta_{a,t} \right] \\
&\stackrel{(a)}{\leq} \mathbf{E} \left[ \max_{t \in [T], a \in [N]} |\zeta_{a,t}| \cdot \sum_{t=1}^T \sum_{a \in A_t} |\delta_{a,t}| \right] \\
&= \mathbf{E} \left[ \max_{t \in [T], a \in [N]} |\zeta_{a,t}| \cdot \sum_{t=1}^T \sum_{a \in A_t} \sqrt{\frac{\gamma m \ln t}{n_{a,t} + 1}} \right],
\end{aligned} \tag{18}$$

309 where (a) is due to Hölder's inequality. By invoking Lemma 5 again, we have that

$$\sum_{t=1}^T \sum_{a \in A_t} \sqrt{\frac{\gamma m \ln t}{n_{a,t} + 1}} \leq \sqrt{\gamma m \ln T} \sum_{t=1}^T \sum_{a \in A_t} \sqrt{\frac{1}{n_{a,t} + 1}} \leq 2m\sqrt{\gamma NT \ln T}. \tag{19}$$

310 Then, using the maximal inequality (Fact 2), we have  $\mathbf{E} [\max_{t \in [T], a \in [N]} |\zeta_{a,t}|] \leq \sqrt{2 \ln 2NT} \leq$   
311  $\sqrt{6 \ln T}$ , where the last inequality is due to that  $2 \leq N \leq T$ . Thus, we have

$$\mathbf{E} \left[ \sum_{t=1}^T \sum_{a \in A_t} (w_{a,t} - \hat{r}_{a,n_{a,t}}) \right] \leq 2m \ln T \sqrt{6\gamma NT}. \tag{20}$$

312 Finally, by substituting (20) into (17), we complete the proof.  $\square$

### 313 C.3 Proof of Lemma 3

**Lemma 3.** *The probability that event  $\bar{\mathcal{E}}_t$  to happen satisfies that*

$$\sum_{t=1}^T \Pr(\bar{\mathcal{E}}_t) \leq \frac{\pi^2}{3}.$$

314 *Proof.* By a union bound and Hoeffding's inequality, we have that

$$\begin{aligned}
& \sum_{t=1}^T \Pr \left( \exists a \in [N] : |r_a - \hat{r}_{a,n_{a,t}}| > \sqrt{\frac{3 \ln Nt}{n_{a,t} + 1}} \right) \\
& \leq \sum_{t=1}^T \sum_{a \in [N]} \sum_{s=0}^{t-1} \Pr \left( |\hat{r}_{a,s} - r_a| > \sqrt{\frac{3 \ln Nt}{s + 1}} \right) \\
& = \sum_{a \in [N]} \sum_{t=1}^T \left( \Pr \left( r_a > \sqrt{3 \ln Nt} \right) + \sum_{s=1}^{t-1} \Pr \left( |\hat{r}_{a,s} - r_a| > \sqrt{\frac{3 \ln Nt}{s + 1}} \right) \right) \quad (21) \\
& \stackrel{(a)}{\leq} \sum_{a \in [N]} \left( 0 + \sum_{t=1}^T \sum_{s=1}^{t-1} \Pr \left( |\hat{r}_{a,s} - r_a| > \sqrt{\frac{3 \ln Nt}{2s}} \right) \right) \\
& \leq N \sum_{t=1}^{\infty} \sum_{s=1}^{t-1} \frac{2}{(Nt)^3} = \frac{\pi^2}{3N^2},
\end{aligned}$$

315 where step (a) is due to  $r_a \in [0, 1]$ ,  $\forall a \in [N]$  and  $3 \ln Nt > 1$  because  $N \geq 2$ , and that  
316  $s + 1 \leq 2s$  for any  $s \geq 1$ .

317

□

#### 318 C.4 Proof of Lemma 4

319 **Lemma 4.** In each round  $t \geq \max\{\sqrt{m}, 4\}$ , given any  $\Theta_t$ , we have

$$\frac{1}{\Pr_{\Theta_t} \left( \sum_{a \in A_t^*} w_{a,t} \geq \sum_{a \in A_t^*} r_a \right)} \leq 2\Phi \left( -\sqrt{4/\gamma} \right)^{-1},$$

320 where  $\Phi(\cdot)$  is the cdf of the standard Gaussian distribution.

321 *Proof.* Given  $\Theta_t$ ,  $A_t^*$  is determined. Define  $\mathcal{H}_t := \left\{ \forall a \in A_t^* : |r_a - \hat{r}_{a,n_{a,t}}| \leq \sqrt{\frac{4 \ln t}{n_{a,t} + 1}} \right\}$ .

322 Since  $t \geq \max\{\sqrt{m}, 4\}$ , we have that

$$\begin{aligned}
\Pr_{\Theta_t} (\mathcal{H}_t) & \geq 1 - \sum_{a \in A_t^*} \sum_{s_a=0}^{t-1} \Pr_{\Theta_t} \left( |r_a - \hat{r}_{a,s_a}| \geq \sqrt{\frac{4 \ln t}{s_a + 1}} \right) \\
& = 1 - \sum_{a \in A_t^*} \sum_{s_a=1}^{t-1} \Pr_{\Theta_t} \left( |r_a - \hat{r}_{a,s_a}| \geq \sqrt{\frac{4 \ln t}{s_a + 1}} \right) \\
& \geq 1 - \sum_{a \in A_t^*} \sum_{s_a=1}^{t-1} \Pr_{\Theta_t} \left( |r_a - \hat{r}_{a,s_a}| \geq \sqrt{\frac{4 \ln t}{2s_a}} \right) \quad (22) \\
& \geq 1 - mt \cdot 2 \cdot e^{-2 \cdot s_a \cdot 4 \ln t / (2s_a)} \\
& = 1 - \frac{2mt}{t^4} \\
& \geq 1 - \frac{2}{t} \\
& \geq 0.5 \quad .
\end{aligned}$$

323 We have

$$\begin{aligned}
\Pr_{\Theta_t} \left( \sum_{a \in A_t^*} w_{a,t} \geq \sum_{a \in A_t^*} r_a \right) &\geq \Pr_{\Theta_t} \left( \sum_{a \in A_t^*} w_{a,t} \geq \sum_{a \in A_t^*} r_a, \mathcal{H}_t \right) \\
&= \Pr_{\Theta_t} \left( \sum_{a \in A_t^*} w_{a,t} - \hat{r}_{a,n_{a,t}} \geq \sum_{a \in A_t^*} r_a - \hat{r}_{a,n_{a,t}}, \mathcal{H}_t \right) \\
&= \Pr_{\Theta_t}(\mathcal{H}_t) \cdot \Pr_{\Theta_t} \left( \sum_{a \in A_t^*} w_{a,t} - \hat{r}_{a,n_{a,t}} \geq \sum_{a \in A_t^*} r_a - \hat{r}_{a,n_{a,t}} \mid \mathcal{H}_t \right) \\
&\stackrel{(a)}{\geq} 0.5 \cdot \Pr_{\Theta_t} \left( \sum_{a \in A_t^*} w_{a,t} - \hat{r}_{a,n_{a,t}} \geq \sum_{a \in A_t^*} \sqrt{\frac{4 \ln t}{n_{a,t} + 1}} \right) \\
&\stackrel{(b)}{\geq} 0.5 \cdot \Pr_{\Theta_t} \left( \sum_{a \in A_t^*} w_{a,t} - \hat{r}_{a,n_{a,t}} \geq \sqrt{m \sum_{a \in A_t^*} \frac{4 \ln t}{n_{a,t} + 1}} \right) \\
&= 0.5 \cdot \Phi \left( -\sqrt{4/\gamma} \right) ,
\end{aligned} \tag{23}$$

324 where step (a) is due to (22) and the fact that event  $\mathcal{E}_t$  is true. Step (b) uses the Cauchy–Schwarz  
325 inequality, i.e., we have  $\sum_{a \in A_t^*} \sqrt{\frac{4 \ln t}{n_{a,t} + 1}} \leq \sqrt{m \cdot \sum_{a \in A_t^*} \frac{4 \ln t}{n_{a,t} + 1}}$ . The last equality is due to the  
326 standardization of Gaussian distribution.  $\square$

### 327 C.5 Proof of Lemma 5

328 **Lemma 5.** We have  $\sum_{t=1}^T \sum_{a \in A_t} \sqrt{\frac{1}{n_{a,t} + 1}} \leq 2\sqrt{mNT}$ .

329 *Proof.* Note that the LHS of the above inequality is a random variable. We provide an upper bound  
330 for this random variable.

331 Recall  $n_{a,t} := \sum_{\tau=1}^{t-1} \mathbf{1}[a \in A_\tau]$  is the number of times that arm  $a$  has been played at the be-  
332 ginning of round  $t$ . Let  $\tau_a(n)$  denote the round for arm  $a$  to be played for the  $n$ -th time, and thus  
333  $n_{a,\tau_a(n)} = n - 1$ .

$$\begin{aligned}
\sum_{t=1}^T \sum_{a \in A_t} \sqrt{\frac{1}{n_{a,t} + 1}} &= \sum_{t=1}^T \sum_{a \in [N]} \sqrt{\frac{1}{n_{a,t} + 1}} \mathbf{1}[a \in A_t] \\
&\stackrel{(a)}{=} \sum_{a \in [N]} \sum_{n=1}^{n_{a,T+1}} \sum_{t=\tau_a(n)}^{\tau_a(n+1)-1} \sqrt{\frac{1}{n_{a,t} + 1}} \mathbf{1}[a \in A_t] \\
&\stackrel{(b)}{=} \sum_{a \in [N]} \sum_{n=1}^{n_{a,T+1}} \sqrt{\frac{1}{n}} \leq \sum_{a \in [N]} \int_0^{n_{a,T+1}} \sqrt{\frac{1}{n}} dn \\
&= 2 \sum_{a \in [N]} \sqrt{n_{a,T+1}} \stackrel{(c)}{\leq} 2 \sqrt{N \sum_{a \in [N]} n_{a,T+1}} \\
&\stackrel{(d)}{=} 2\sqrt{mNT},
\end{aligned} \tag{24}$$

334 where step (a) partitions all  $T$  rounds into multiple intervals based on the arrivals of observations from  
335 arm  $a$ . Step (b) uses the fact that  $\sum_{t=\tau_a(n)}^{\tau_a(n+1)-1} \mathbf{1}[a \in A_t] \cdot \sqrt{\frac{1}{n_{a,t} + 1}} = \sqrt{\frac{1}{n-1+1}} = \sqrt{\frac{1}{n}}$ , be-  
336 cause  $n_{a,\tau_a(n)} = n - 1$  and  $\mathbf{1}[a \in A_t] = 0$  for all  $t \in \{\tau_a(n) + 1, \dots, \tau_a(n+1) - 1\}$ .  
337 Step (c) uses Cauchy-Schwarz inequality. Step (d) uses the fact that  $\sum_{a \in [N]} n_{a,T+1} \leq mT$ .

338  $\square$

339 **C.6 Proof of Upper bound**

340 *Upper Bound Proof of Theorem 1.* Denote by  $\mathcal{E}_t := \left\{ \forall a \in [N] : |r_a - \hat{r}_{a,n_{a,t}}| \leq \sqrt{\frac{3 \ln Nt}{n_{a,t}+1}} \right\}$   
 341 the high-probability event that the empirical mean reward is close to the true mean reward for arm  $a$ ,  
 342 and by  $\bar{\mathcal{E}}_t$  the complementary event of  $\mathcal{E}_t$ .

343 Let  $t' = \max\{\sqrt{m}, 4\}$ . We first decompose the regret as follows:

$$\begin{aligned}
 \mathcal{R}(T) &= \sum_{t=1}^{t'-1} \mathbf{E} \left[ \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} r_a \right] + \sum_{t=t'}^T \mathbf{E} \left[ \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} r_a \right] \\
 &\stackrel{(a)}{\leq} m \max\{\sqrt{m}, 4\} + \mathbf{E} \left[ \sum_{t=t'}^T \left( \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} r_a \right) \mathbf{1}[\mathcal{E}_t] \right] + \mathbf{E} \left[ \sum_{t=t'}^T \left( \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} r_a \right) \mathbf{1}[\bar{\mathcal{E}}_t] \right] \\
 &\stackrel{(b)}{\leq} m \max\{\sqrt{m}, 4\} + \mathbf{E} \left[ \sum_{t=t'}^T \left( \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} w_{a,t} + \sum_{a \in A_t} w_{a,t} - \sum_{a \in A_t} r_a \right) \mathbf{1}[\mathcal{E}_t] \right] + m \frac{\pi^2}{3N^2} \\
 &\leq \underbrace{\sum_{t=t'}^T \mathbf{E} \left[ \left( \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} w_{a,t} \right) \right]}_{=: I_1, \text{ optimism part}} + \underbrace{\sum_{t=t'}^T \mathbf{E} \left[ \sum_{a \in A_t} (w_{a,t} - r_a) \mathbf{1}[\mathcal{E}_t] \right]}_{=: I_2, \text{ deviation part}} + m \max\{\sqrt{m}, 4\} + \frac{\pi^2}{3},
 \end{aligned} \tag{25}$$

344 where step (a) is due to the fact that  $\sum_{a \in A_t^*} r_a - \sum_{a \in A_t} r_a \leq m$  by the definition of  $r_a$  and  $m$   
 345 and step (b) is due to Lemma 3.

346 Now, invoking Lemma 1 with proofs in Appendix C.1, we have term  $I_1$  bounded as follows:

$$I_1 \leq 8\sqrt{3\gamma}\Phi(-\sqrt{4/\gamma})^{-1}m \ln T\sqrt{NT}, \tag{26}$$

347 and  $I_2$  can be bounded by using Lemma 2 with proofs in Appendix C.2:

$$I_2 \leq 2m \ln T\sqrt{6\gamma NT} + 2\sqrt{6mNT \ln T}. \tag{27}$$

348 Thus, we have that

$$\begin{aligned}
 \mathcal{R}(T) &\leq \left( 2\sqrt{6\gamma} + 8\sqrt{3\gamma}\Phi(-\sqrt{4/\gamma})^{-1} \right) m \ln T\sqrt{NT} + 2\sqrt{6mNT \ln T} \\
 &\quad + m \left( \max\{\sqrt{m}, 4\} + \frac{\pi^2}{3} \right).
 \end{aligned} \tag{28}$$

349 Using numerical optimization methods searching from  $\gamma = 0.0001$  to  $\gamma = 100$ , we can find that  
 350 when  $\gamma = 6.4$ , the coefficient for the first item can achieve a minimum value of 175.74.

351 □

352 **C.7 Proof of Lower Bound**

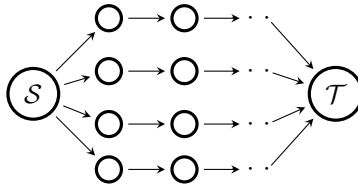


Figure 5: Problem instance for the lower-bound proof. Nodes  $S$  and  $T$  are the starting and ending points for each path.

353 *Lower bound Proof in Theorem 1.* Our proof uses similar ideas to the proofs of Theorem 1.4 in  
 354 Agrawal & Goyal (2017b).

355 We construct a path selection problem involving  $N$  links (each link corresponds to a base arm) and  $K$   
 356 paths (each path corresponds to a super arm), as illustrated in Fig. 5. Each path consists of  $m$  links,



357 and thus, the total number of base arms  $N = mK$ . We consider a fixed availability set throughout  
 358 all  $T$  rounds, i.e.,  $\Theta_t = \Theta := \{A_1, A_2, \dots, A_K\}$  for all rounds  $t \in [T]$  with each super arm  
 359  $A_k$  being a feasible path. We assume the first path  $A_1$  is the unique optimal one.

360 We construct the following Bernoulli reward distributions for each base arm. Let  $\Delta :=$   
 361  $\sqrt{K \ln K/T}$ . For any base arm in the optimal super arm  $A_1$ , we use a degenerate distribution  
 362 putting mass 1 on a single point  $\sqrt{\gamma}\Delta$ , i.e., if  $A_1$  is played, for any base arm in it, we always observe  
 363  $\sqrt{\gamma}\Delta$  as the random reward. Similarly, for the remaining base arms in the sub-optimal super arms, we  
 364 put mass 1 on a single point 0, i.e., the random reward is always 0 for any base arm in a sub-optimal  
 365 super arm.

366 Let  $Q_A(t)$  denote the number of times that super arm  $A \in \Theta$  has been played at the beginning  
 367 of round  $t$ . Since there are no overlapping base arms between two distinct super arms, we have  
 368  $Q_A(t) = n_{a,t}$  for all  $a \in A$ . Let  $c \in (0, 1)$  be some universal constant that will be tuned later.  
 369 Define  $B_t^* := \{Q_{A_1}(t) > t - cT\}$  as the event that the optimal super arm  $A_1$  has been observed  
 370 at least  $(t - cT)$  times by the beginning of round  $t$ .

371 We lower bound the total regret from round 1 to the end of round  $T$  by analyzing two cases that are  
 372 exhaustive and mutually exclusive based on events  $B_t^*$  for all rounds  $t \in [T]$ .

373 **If  $B_t^*$  is not true for some  $t \in [T]$ ,** we have the total number of times of playing sub-  
 374 optimal super arms by the beginning of round  $t$  is  $\sum_{A \in \Theta \setminus A_1} Q_A(t) = t - Q_{A_1}(t) \geq t - (t -$   
 375  $cT) \geq cT$ , which implies the total regret by the end of round  $T$  is at least  $cT \cdot m \cdot \sqrt{\gamma}\Delta =$   
 376  $\Omega(m\sqrt{KT \ln K}) = \Omega(\sqrt{mNT \ln(N/m)})$ . Note that the total regret from round 1 to round  $t$  is a  
 377 lower bound for the total regret over all  $T$  rounds.

378 **If  $B_t^*$  is true for all  $t \in [T]$ ,** we have the total number of times  $\sum_{A \in \Theta \setminus A_1} Q_A(t)$  of playing  
 379 sub-optimal super arms by the beginning of round  $t$  is upper bounded by

$$\sum_{A \in \Theta \setminus A_1} Q_A(t) = t - Q_{A_1}(t) \leq t - (t - cT) = cT \quad . \quad (29)$$

380 Due to the spread of a sub-optimal super arm's posterior distribution, the learning agent will make  
 381 mistakes when deciding which super arm to play. Formally, we show that with at least a constant  
 382 probability, the learning agent will play a sub-optimal super arm. Note that whether event  $B_t^*$  is true  
 383 or not is determined by the history information  $\mathcal{F}_{t-1}$ .

384 Recall  $w_{a,t} \sim \mathcal{N}\left(\hat{r}_{a,n_{a,t}}, \frac{\gamma m \ln t}{n_{a,t}+1}\right)$ . Now, we construct a lower bound for the probability of  
 385 selecting a sub-optimal arm in round  $t$  conditioned on instantiations  $F_{t-1}$  of  $\mathcal{F}_{t-1}$  such that  $B_t^*$  is  
 386 true. We have

$$\begin{aligned} & \Pr(\exists A \in \Theta \setminus A_1 : A_t = A \mid \mathcal{F}_{t-1} = F_{t-1}) \\ & \geq \Pr(\exists A \in \Theta \setminus A_1 : \sum_{a \in A} w_{a,t} > \sum_{a \in A_1} w_{a,t} \mid \mathcal{F}_{t-1} = F_{t-1}) \\ & \stackrel{(a)}{\geq} \Pr(\exists A \in \Theta \setminus A_1 : \sum_{a \in A} w_{a,t} \geq m\sqrt{\gamma}\Delta \mid \mathcal{F}_{t-1} = F_{t-1}) \\ & \quad \cdot \Pr(\sum_{a \in A_1} w_{a,t} < m\sqrt{\gamma}\Delta \mid \mathcal{F}_{t-1} = F_{t-1}) \\ & \stackrel{(b)}{=} \Pr(\exists A \in \Theta \setminus A_1 : \sum_{a \in A} w_{a,t} \geq m\sqrt{\gamma}\Delta \mid \mathcal{F}_{t-1} = F_{t-1}) \cdot \frac{1}{2} \quad , \end{aligned} \quad (30)$$

387 where step (a) uses the fact that all super arms are independent based on our construction of the path  
 388 selection problem. Step (b) uses the fact that the sum of multiple independent Gaussian random  
 389 variables is still Gaussian and  $(\sum_{a \in A_1} w_{a,t} - m\sqrt{\gamma}\Delta)$  is a zero-mean Gaussian distribution. Note  
 390 that the empirical mean of each base arm in the optimal super arm is exactly  $\sqrt{\gamma}\Delta$ .

391 Now, we construct a lower bound for (30) by using Gaussian anti-concentration bounds. We have

$$\begin{aligned}
& \Pr \left( \exists A \in \Theta \setminus A_1 : \sum_{a \in A} w_{a,t} \geq m\sqrt{\gamma}\Delta \mid \mathcal{F}_{t-1} = F_{t-1} \right) \\
& \geq \Pr \left( \exists A \in \Theta \setminus A_1 : \sum_{a \in A} w_{a,t} \geq m\sqrt{\gamma}\Delta\sqrt{\ln t} \mid \mathcal{F}_{t-1} = F_{t-1} \right) \\
& = 1 - \Pr \left( \sum_{a \in A} w_{a,t} < m\Delta\sqrt{\gamma \ln t}, \forall A \in \Theta \setminus A_1 \mid \mathcal{F}_{t-1} = F_{t-1} \right) \\
& = 1 - \prod_{A \in \Theta \setminus A_1} \left( 1 - \Pr \left( \sum_{a \in A} w_{a,t} \sqrt{(Q_A(t) + 1)} \geq m\Delta\sqrt{(Q_A(t) + 1)\gamma \ln t} \mid \mathcal{F}_{t-1} = F_{t-1} \right) \right) \\
& = 1 - \prod_{A \in \Theta \setminus A_1} \left( 1 - \Pr \left( \underbrace{\frac{\sum_{a \in A} w_{a,t} \sqrt{(Q_A(t) + 1)}}{m\sqrt{\gamma \ln t}}}_{\geq \frac{1}{8\sqrt{\pi}} e^{-\frac{7}{2}\Delta^2(Q_A(t)+1)}} \geq \Delta\sqrt{(Q_A(t) + 1)} \mid \mathcal{F}_{t-1} = F_{t-1} \right) \right) \\
& \geq 1 - \prod_{A \in \Theta \setminus A_1} \left( 1 - \frac{1}{8\sqrt{\pi}} e^{-\frac{7}{2}\Delta^2(Q_A(t)+1)} \right), \tag{31}
\end{aligned}$$

392 where the last inequality uses the fact that  $\frac{\sum_{a \in A} w_{a,t} \sqrt{Q_A(t)+1}}{m\sqrt{\gamma \ln t}} \sim \mathcal{N}(0, 1)$  and then the one-sided  
393 anti-concentration inequality shown in (6).

394 Tune constant  $c = 0.001$ . Then, we use the upper bound constructed in (29) to continue lower  
395 bounding (31). We have

$$\begin{aligned}
& 1 - \prod_{A \in \Theta \setminus A_1} \left( 1 - \frac{1}{8\sqrt{\pi}} e^{-\frac{7}{2}\Delta^2(Q_A(t)+1)} \right) \\
& \stackrel{(a)}{\geq} 1 - \prod_{A \in \Theta \setminus A_1} \left( 1 - \frac{1}{8\sqrt{\pi}} e^{-\frac{7c}{2}\Delta^2 \frac{\sqrt{KT \ln K}}{(K-1)\Delta} - \frac{7\Delta^2}{2}} \right) \\
& = 1 - \prod_{A \in \Theta \setminus A_1} \left( 1 - \frac{1}{8\sqrt{\pi}} e^{-\frac{7c}{2} \frac{K \ln K}{(K-1)} - \frac{7K \ln K}{2}} \right) \\
& = 1 - \prod_{A \in \Theta \setminus A_1} \left( 1 - \frac{1}{8\sqrt{\pi}} e^{-(\frac{7c}{2} \frac{K}{K-1} + \frac{7K}{2T}) \ln K} \right) \tag{32} \\
& \stackrel{(b)}{\geq} 1 - \prod_{A \in \Theta \setminus A_1} \left( 1 - \frac{1}{8\sqrt{\pi}} e^{-\ln K} \right) \\
& = 1 - \left( 1 - \frac{1}{8\sqrt{\pi}K} \right)^{K-1} \\
& \stackrel{(c)}{\geq} 1 - \left( e^{-\frac{1}{8\sqrt{\pi}K}} \right)^{K-1} \\
& \geq 1 - e^{-\frac{1}{16\sqrt{\pi}}},
\end{aligned}$$

396 where step (a) is due to the fact that, constrained on (29), i.e.,  $\sum_{A \in \Theta \setminus A_1} Q_A(t) \leq cT =$   
397  $c\frac{\sqrt{KT \ln K}}{\Delta}$ , the quantity  $\prod_{A \in \Theta \setminus A_1} \left( 1 - \frac{1}{4\sqrt{\pi}} e^{-\frac{7c}{2}\Delta^2 \frac{\sqrt{KT \ln K}}{(K-1)\Delta} - \frac{7\Delta^2}{2}} \right)$  is maximized when

398  $Q_A(t) = \frac{c\sqrt{KT \ln K}}{(K-1)\Delta}$  for all  $A \in \Theta \setminus A_1$ . Step (b) uses the fact that, when  $c = 0.001$ , we have  
399  $\frac{7c}{2} \frac{K}{K-1} + \frac{7K}{2T} \leq 1$  when  $T$  is sufficiently large, e.g.,  $T > 5K$ . Step (c) uses  $1 - x \leq e^{-x}$ .

400 Now, we are ready to complete the proof. Let  $p := \frac{1}{2} \left(1 - e^{-\frac{1}{16\sqrt{\pi}}}\right)$ . By plugging the lower bound  
 401 constructed in (32) into (30), we have  $\Pr(\exists A \in \Theta \setminus A_1 : A_t = A \mid \mathcal{F}_{t-1} = F_{t-1}) \geq p$ ,  
 402 which implies the total regret by the end of round  $T$  is at least  $Tpm\sqrt{\gamma}\Delta = \Omega(\sqrt{mNT} \ln(N/m))$ .  
 403  $\square$

## 404 D Proofs for Theorem 2

### 405 D.1 Proof of Lemma 6

406 **Lemma 6.** *The optimism part in CL-SG satisfies that*

$$\mathbf{E} \left[ \sum_{t=\max\{m,4\}}^T \left( \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} \bar{r}_{a,t} \right) \right] \leq 8\sqrt{2\gamma}\Phi(-\sqrt{4/\gamma})^{-1} \ln T\sqrt{mNT}. \quad (33)$$

407 *Proof.* Similar to the proof of Lemma 1. There are three steps for the proofs.

408 Step 1: Let  $t' = \max\{\sqrt{m}, 4\}$  we show that the following inequality holds for each round  $t$   
 409 conditioned on  $\Theta_t$ :

$$\mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} \bar{r}_{a,t} \right] \leq 2\Phi(-\sqrt{4/\gamma})^{-1} \cdot \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} \bar{r}_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \bar{r}_{a,t} \right] \right)^+ \right]. \quad (34)$$

410 Step 2: Let  $\tilde{r}_{a,t} = \hat{r}_{a,t} + \tilde{w}_t \sqrt{\frac{\gamma \ln t}{n_{a,t}+1}}$ , where  $\tilde{w}_t \sim \mathcal{N}(0, 1)$  is an independent copy of  $w_t$ . With  
 411  $\tilde{r}_{a,t}$ , we can further bound the last term in (34) as follows.

$$\mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} \bar{r}_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \bar{r}_{a,t} \right] \right)^+ \right] \leq \mathbf{E}_{\Theta_t} \left[ \left| \sum_{a \in A_t} \bar{r}_{a,t} - \sum_{a \in A_t} \tilde{r}_{a,t} \right| \right]. \quad (35)$$

412 Step 3: Summing over  $T$  rounds, we have that

$$\mathbf{E} \left[ \sum_{t=1}^T \left| \sum_{a \in A_t} \bar{r}_{a,t} - \sum_{a \in A_t} \tilde{r}_{a,t} \right| \right] \leq 4 \ln T \sqrt{2\gamma mNT} \quad (36)$$

413 Combining these three steps, we have

$$\begin{aligned} \mathbf{E} \left[ \sum_{t=\max\{\sqrt{m},4\}}^T \left( \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} \bar{r}_{a,t} \right) \right] &\leq 2\Phi(-\sqrt{4/\gamma})^{-1} \mathbf{E} \left[ \sum_{t=\max\{\sqrt{m},4\}}^T \left| \sum_{a \in A_t} \bar{r}_{a,t} - \sum_{a \in A_t} \tilde{r}_{a,t} \right| \right] \\ &\leq 2\Phi(-\sqrt{4/\gamma})^{-1} \mathbf{E} \left[ \sum_{t=1}^T \left| \sum_{a \in A_t} \bar{r}_{a,t} - \sum_{a \in A_t} \tilde{r}_{a,t} \right| \right] \\ &\leq 8\sqrt{2\gamma}\Phi(-\sqrt{4/\gamma})^{-1} \ln T\sqrt{mNT}. \end{aligned} \quad (37)$$

414 Now, we give the details for these three steps.

415 **Step 1 proof.** If  $\mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} \bar{r}_{a,t} \right] \leq 0$ , the proof is trivial as the RHS in (34) is  
 416 non-negative.

417 Recall  $(\cdot)^+ := \max\{\cdot, 0\}$ . For the case where  $\alpha := \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} \bar{r}_{a,t} \right] > 0$ , we use  
 418 Markov's inequality and have

$$\begin{aligned} \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} \bar{r}_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \bar{r}_{a,t} \right] \right)^+ \right] &\geq \alpha \cdot \Pr_{\Theta_t} \left( \left( \sum_{a \in A_t} \bar{r}_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \bar{r}_{a,t} \right] \right)^+ \geq \alpha \right) \\ &\geq \alpha \cdot \Pr_{\Theta_t} \left( \sum_{a \in A_t} \bar{r}_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \bar{r}_{a,t} \right] \geq \alpha \right), \end{aligned} \quad (38)$$

419 which gives

$$\begin{aligned} \alpha &= \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} \bar{r}_{a,t} \right] \\ &\leq \frac{\mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} \bar{r}_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \bar{r}_{a,t} \right] \right)^+ \right]}{\Pr_{\Theta_t} \left( \sum_{a \in A_t} \bar{r}_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \bar{r}_{a,t} \right] \geq \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} \bar{r}_{a,t} \right] \right)} \\ &= \frac{\mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} \bar{r}_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \bar{r}_{a,t} \right] \right)^+ \right]}{\Pr_{\Theta_t} \left( \sum_{a \in A_t} \bar{r}_{a,t} \geq \sum_{a \in A_t^*} r_a \right)} \leq \frac{\mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} \bar{r}_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \bar{r}_{a,t} \right] \right)^+ \right]}{\Pr_{\Theta_t} \left( \sum_{a \in A_t^*} \bar{r}_{a,t} \geq \sum_{a \in A_t^*} r_a \right)} \\ &= 2\Phi(-\sqrt{4/\gamma})^{-1} \cdot \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} \bar{r}_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \bar{r}_{a,t} \right] \right)^+ \right]. \end{aligned} \quad (39)$$

420 **Step 2 proof.** Since  $w_t$  and  $\tilde{w}_t$  are i.i.d., we have  $\mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \bar{r}_{a,t} \right] =$   
 421  $\mathbf{E}_{\Theta_t} \left[ \max_{A \in \Theta_t} \sum_{a \in A} \bar{r}_{a,t} \right] = \mathbf{E}_{\Theta_t} \left[ \max_{A \in \Theta_t} \sum_{a \in A} \tilde{r}_{a,t} \right] \geq \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \tilde{r}_{a,t} \mid A_t \right] =$   
 422  $\mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \tilde{r}_{a,t} \mid A_t, w_t \right]$ . Then, we have

$$\begin{aligned}
\mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} \bar{r}_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \bar{r}_{a,t} \right] \right)^+ \right] &\leq \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} \bar{r}_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \tilde{r}_{a,t} \mid A_t \right] \right)^+ \right] \\
&= \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} \bar{r}_{a,t} - \mathbf{E}_{\Theta_t} \left[ \sum_{a \in A_t} \tilde{r}_{a,t} \mid A_t, w_t \right] \right)^+ \right] \\
&= \mathbf{E}_{\Theta_t} \left[ \left( \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} \bar{r}_{a,t} - \sum_{a \in A_t} \tilde{r}_{a,t} \right) \mid A_t, w_t \right] \right)^+ \right] \\
&\leq \mathbf{E}_{\Theta_t} \left[ \left| \mathbf{E}_{\Theta_t} \left[ \left( \sum_{a \in A_t} \bar{r}_{a,t} - \sum_{a \in A_t} \tilde{r}_{a,t} \right) \mid A_t, w_t \right] \right| \right] \\
&\leq \mathbf{E}_{\Theta_t} \left[ \left| \sum_{a \in A_t} \bar{r}_{a,t} - \sum_{a \in A_t} \tilde{r}_{a,t} \right| \mid A_t, w_t \right] \\
&\leq \mathbf{E}_{\Theta_t} \left[ \left| \sum_{a \in A_t} \bar{r}_{a,t} - \sum_{a \in A_t} \tilde{r}_{a,t} \right| \right].
\end{aligned} \tag{40}$$

423 **Step 3 proof.** By Hölder's inequality, we have that

$$\begin{aligned}
\mathbf{E} \left[ \sum_{t=1}^T \left| \sum_{a \in A_t} \bar{r}_{a,t} - \sum_{a \in A_t} \tilde{r}_{a,t} \right| \right] &\leq \mathbf{E} \left[ \sum_{t=1}^T |w_t - \tilde{w}_t| \sum_{a \in A_t} \sqrt{\frac{\gamma \ln t}{n_{a,t} + 1}} \right] \\
&\leq \mathbf{E} \left[ \max_{t \in [T]} |w_t - \tilde{w}_t| \sum_{t=1}^T \sum_{a \in A_t} \sqrt{\frac{\gamma \ln t}{n_{a,t} + 1}} \right] \\
&\stackrel{(a)}{\leq} \mathbf{E} \left[ \max_{t \in [T]} |w_t - \tilde{w}_t| \right] \cdot 2\sqrt{\gamma m N T \ln T} \\
&\stackrel{(b)}{\leq} 2\sqrt{\ln 2T} \cdot 2\sqrt{\gamma m N T \ln T} \\
&\leq 4 \ln T \sqrt{2\gamma m N T},
\end{aligned} \tag{41}$$

424 where step (a) is due to Lemma 5, and step (b) is due to Fact 2 and  $w_t - \tilde{w}_t$  is a Gaussian variable  
425 with variance 2.  $\square$

## 426 D.2 Proof of Lemma 7

**Lemma 7.** *In CL-SG, the regret of the deviation part is*

$$\mathbf{E} \left[ \sum_{t=1}^T \left( \sum_{a \in A_t} \bar{r}_{a,t} - \sum_{a \in A_t} r_a \right) \mathbf{1}[\mathcal{E}_t] \right] \leq 4 \ln T \sqrt{\gamma m N T} + 2\sqrt{6m N T \ln T}.$$

427 *Proof.* Recall that  $\bar{r}_{a,t} = \hat{r}_{a,n_{a,t}} + w_t \sqrt{\frac{\gamma \ln t}{n_{a,t} + 1}}$ . When  $\mathcal{E}_t$  happens, we have that

$$\begin{aligned}
\mathbf{E} \left[ \sum_{t=1}^T \left( \sum_{a \in A_t} \bar{r}_{a,t} - \sum_{a \in A_t} r_a \right) \mathbf{1}[\mathcal{E}_t] \right] &= \mathbf{E} \left[ \sum_{t=1}^T \sum_{a \in A_t} \left( \hat{r}_{a,n_{a,t}} + w_t \sqrt{\frac{\gamma \ln t}{n_{a,t} + 1}} - \hat{r}_{a,n_{a,t}} + \sqrt{\frac{3 \ln N t}{n_{a,t} + 1}} \right) \mathbf{1}[\mathcal{E}_t] \right] \\
&\leq \sqrt{\gamma \ln T} \mathbf{E} \left[ \sum_{t=1}^T w_t \left( \sum_{a \in A_t} \sqrt{\frac{1}{n_{a,t} + 1}} \right) \right] + \sqrt{6 \ln T} \mathbf{E} \left[ \sum_{t=1}^T \sum_{a \in A_t} \sqrt{\frac{1}{n_{a,t} + 1}} \right],
\end{aligned} \tag{42}$$

428 where the last inequality is due to that  $N \leq T$ . Regarding the first item in RHS of (42), we can  
429 apply Hölder's inequality to have that

$$\begin{aligned}
\mathbf{E} \left[ \sum_{t=1}^T w_t \left( \sum_{a \in A_t} \sqrt{\frac{1}{n_{a,t} + 1}} \right) \right] &\leq \mathbf{E} \left[ \max_{1 \leq t \leq T} |w_t| \cdot \left| \sum_{t=1}^T \sum_{a \in A_t} \sqrt{\frac{1}{n_{a,t} + 1}} \right| \right] \\
&\leq \mathbf{E} \left[ \max_{1 \leq t \leq T} |w_t| \cdot 2\sqrt{mNT} \right] \\
&\leq 4\sqrt{mNT \ln T},
\end{aligned} \tag{43}$$

430 where the second inequality is due to Lemma 5, and the last inequality is due to the maximal  
431 inequality (Fact 2) for Gaussian variables such that  $\mathbf{E} [\max_{1 \leq t \leq T} |w_t|] \leq \sqrt{2 \ln 2T} \leq 2\sqrt{T}$ .

432 Regarding the second term in RHS of (42), we can invoke Lemma 5 again to give a bound of  
433  $2\sqrt{6mNT \ln T}$ .

434

□

### 435 D.3 Proof of Lemma 8

436 **Lemma 8.** *In each round  $t > \max\{\sqrt{m}, 4\}$ , given any  $\Theta_t$ , we have that for CL-SG:*

$$\frac{1}{\Pr_{\Theta_t} \left( \sum_{a \in A_t^*} \bar{r}_{a,t} \geq \sum_{a \in A_t^*} r_a \right)} \leq 2\Phi \left( -\sqrt{4/\gamma} \right)^{-1}. \tag{44}$$

437 *Proof of Lemma 8.* Given  $\Theta_t$ ,  $A_t^*$  is determined. Define  $\mathcal{H}_t :=$   
438  $\left\{ \forall a \in A_t^* : |r_a - \hat{r}_{a,n_{a,t}}| \leq \sqrt{\frac{4 \ln t}{n_{a,t} + 1}} \right\}$ . We have

$$\begin{aligned}
\Pr_{\Theta_t} (\mathcal{H}_t) &\geq 1 - \sum_{a \in A_t^*} \sum_{s_a=0}^{t-1} \Pr_{\Theta_t} \left( |r_a - \hat{r}_{a,s_a}| \geq \sqrt{\frac{4 \ln t}{s_a + 1}} \right) \\
&= 1 - \sum_{a \in A_t^*} \sum_{s_a=1}^{t-1} \Pr_{\Theta_t} \left( |r_a - \hat{r}_{a,s_a}| \geq \sqrt{\frac{4 \ln t}{s_a + 1}} \right) \\
&\geq 1 - \sum_{a \in A_t^*} \sum_{s_a=1}^{t-1} \Pr_{\Theta_t} \left( |r_a - \hat{r}_{a,s_a}| \geq \sqrt{\frac{4 \ln t}{2s_a}} \right) \\
&\geq 1 - mt \cdot 2 \cdot e^{-2 \cdot s_a \cdot 4 \ln t / (2s_a)} \\
&= 1 - \frac{2mt}{t^4} \\
&\geq 1 - \frac{2}{t} \\
&\geq 0.5,
\end{aligned} \tag{45}$$

439 where the last two inequalities are due to that  $t > \max\{\sqrt{m}, 4\}$ .

440 We have

$$\begin{aligned}
\Pr_{\Theta_t} \left( \sum_{a \in A_t^*} \bar{r}_{a,t} \geq \sum_{a \in A_t^*} r_a \right) &\geq \Pr_{\Theta_t} \left( \sum_{a \in A_t^*} \bar{r}_{a,t} \geq \sum_{a \in A_t^*} r_a, \mathcal{H}_t \right) \\
&= \Pr_{\Theta_t} \left( \sum_{a \in A_t^*} \bar{r}_{a,t} - \hat{r}_{a,n_{a,t}} \geq \sum_{a \in A_t^*} r_a - \hat{r}_{a,n_{a,t}}, \mathcal{H}_t \right) \\
&= \Pr_{\Theta_t}(\mathcal{H}_t) \cdot \Pr_{\Theta_t} \left( \sum_{a \in A_t^*} \bar{r}_{a,t} - \hat{r}_{a,n_{a,t}} \geq \sum_{a \in A_t^*} r_a - \hat{r}_{a,n_{a,t}} \mid \mathcal{H}_t \right) \\
&\stackrel{(a)}{\geq} 0.5 \cdot \Pr_{\Theta_t} \left( \sum_{a \in A_t^*} w_t \sqrt{\frac{\gamma \ln t}{n_{a,t} + 1}} \geq \sum_{a \in A_t^*} \sqrt{\frac{4 \ln t}{n_{a,t} + 1}} \right) \\
&= 0.5 \cdot \Pr_{\Theta_t} \left( w_t \geq \sqrt{4/\gamma} \right) \\
&= 0.5 \cdot \Phi \left( -\sqrt{4/\gamma} \right) \quad ,
\end{aligned} \tag{46}$$

441 where step (a) is due to (45) and the fact that event  $\mathcal{E}_t$  is true.  $\square$

#### 442 D.4 Proof of Upper Bound

443 *Proof.* Recall that  $\mathcal{E}_t := \left\{ \forall a \in [N] : |r_a - \hat{r}_{a,n_{a,t}}| \leq \sqrt{\frac{3 \ln Nt}{n_{a,t} + 1}} \right\}$  is the high-probability  
444 event that the empirical mean reward is close to the true mean reward for arm  $a$ , and  $\bar{\mathcal{E}}_t$  is the  
445 complementary event of  $\mathcal{E}_t$ .

446 Similar to the proof of the upper bound for CTS-G, we first let  $t' = \max\{\sqrt{m}, 4\}$ , and then  
447 decompose the regret as follows:

$$\begin{aligned}
\mathcal{R}(T) &\leq \underbrace{\sum_{t=t'}^T \mathbf{E} \left[ \left( \sum_{a \in A_t^*} r_a - \sum_{a \in A_t} \bar{r}_{a,t} \right) \right]}_{=: I_1, \text{ optimism part}} + \underbrace{\sum_{t=t'}^T \mathbf{E} \left[ \sum_{a \in A_t} (\bar{r}_{a,t} - r_a) \mathbf{1}[\mathcal{E}_t] \right]}_{=: I_2, \text{ deviation part}} \\
&\quad + m \left( \max\{\sqrt{m}, 4\} + \frac{\pi^2}{3} \right) ,
\end{aligned} \tag{47}$$

448 Now, invoking Lemma 6 with proofs in Appendix D.1, we have term  $I_1$  bounded as follows:

$$I_1 \leq 8\sqrt{2\gamma}\Phi(-\sqrt{4/\gamma})^{-1} \ln T\sqrt{mNT}, \tag{48}$$

449 and  $I_2$  can be bounded by using Lemma 7 with proofs in Appendix D.2:

$$I_2 \leq 4 \ln T\sqrt{\gamma mNT} + 2\sqrt{6mNT} \ln T. \tag{49}$$

450 Therefore, we have the regret bounded as follows:

$$\begin{aligned}
\mathcal{R}(T) &\leq \left( 4\sqrt{\gamma} + 8\sqrt{2\gamma}\Phi(-\sqrt{4/\gamma})^{-1} \right) \ln T\sqrt{mNT} + 2\sqrt{6mNT} \ln T \\
&\quad + m \left( \max\{\sqrt{m}, 4\} + \frac{\pi^2}{3} \right) ,
\end{aligned}$$

451 where the coefficient of the first term can be minimized to 144.43 at  $\gamma = 4.57$ .  $\square$

#### 452 D.5 Proof of Lower Bound

453 *Lower bound Proof in Theorem 2.* The main challenge in the proofs arises from the fact that all base  
454 arms share a single random Gaussian seed, creating dependencies between paths that are no longer  
455 independent. However, the lower-bound proof for Theorem 1 relies on the independence of each  
456 super arm. Therefore, this proof must manage these dependencies effectively.

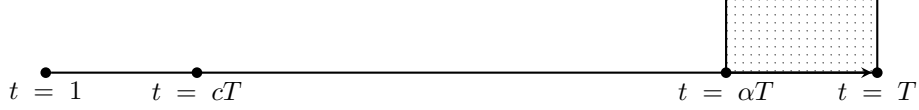


Figure 6: The regret of CL-SG is lower bounded by the regret from rounds  $\alpha T$  to  $T$ .

457 We construct a path selection problem involving  $N$  links (i.e., each link corresponds to a base arm)  
 458 and  $K$  paths (i.e., each path corresponds to a super arm). Each path consists of  $m$  links as illustrated  
 459 in Fig. 5 and the total number of base arms is  $N = mK$ . We use a fixed availability set throughout  
 460 all  $T$  rounds, i.e.,  $\Theta_t = \Theta := \{A_1, A_2, \dots, A_K\}$  for all rounds  $t \in [T]$  with each  $A_k$  being a  
 461 feasible path. We assume the first path is the unique optimal one.

462 We construct the following Bernoulli reward distributions for each base arm. Let  $\Delta := \sqrt{K/T}$ .  
 463 For any base arm in the optimal super arm  $A_1$ , we use a degenerate distribution putting mass 1 on a  
 464 single point  $\sqrt{\gamma}\Delta$ , i.e., if  $A_1$  is played, for each base arm in it, the observed random reward is always  
 465  $\sqrt{\gamma}\Delta$ . Similarly, for the remaining base arms in the sub-optimal super arms, we put mass 1 on a  
 466 single point 0, i.e., the random reward is always 0 for any base arm in a sub-optimal super arm.

467 Let  $Q_A(t)$  denote the total number of times that super arm  $A \in \Theta$  has been played at the beginning  
 468 of round  $t$ . Since there are no overlapping base arms between two distinct super arms, we have  
 469  $Q_A(t) = n_{a,t}$  for all  $a \in A$ , i.e., all base arms in a super arm have the same amount of observations.

470 Let  $c := \frac{1}{6}$ . Define  $B_t^* := \{Q_{A_1}(t) > t - cT\}$  as the event that the optimal super arm  $A_1$  has  
 471 been observed enough times at the beginning of round  $t$ .

472 We lower bound the total regret from round 1 to the end of round  $T$  by analyzing two cases that are  
 473 exhaustive and mutually exclusive based on events  $B_t^*$  for all rounds  $t \in [T]$ .

474 **If  $B_t^*$  is not true for some round  $t \in [T]$ ,** we have the total number of times of playing sub-optimal  
 475 super arms until the beginning of round  $t$  is  $\sum_{A \in \Theta \setminus A_1} Q_A(t) = t - Q_{A_1}(t) \geq cT$ . This lower  
 476 bound implies the total regret from round 1 to the end of round  $t - 1$  is at least  $cT \cdot m \cdot \sqrt{\gamma}\Delta =$   
 477  $\Omega(Tm\sqrt{K/T}) = \Omega(Tm\sqrt{N/(mT)}) = \Omega(\sqrt{mNT})$ . Note that this lower bound is also a regret  
 478 lower bound for the total regret from round 1 to the end of round  $T$ .

479 Let  $\alpha := \frac{5}{6}$ .

480 **If  $B_t^*$  is true for all rounds  $t \in [T]$ ,** the total regret from round 1 to the end of round  $T$  is lower  
 481 bounded by the total regret from round  $t = \alpha T$  to the end of round  $T$ , as shown in Fig. 6. In each  
 482 round  $t \geq \alpha T$ , we have the following inequalities:

$$\sum_{A \in \Theta \setminus A_1} Q_A(t) = t - Q_{A_1}(t) \leq t - (t - c \cdot T) = c \cdot T \quad , \quad (50)$$

483

$$Q_{A_1}(t) > t - c \cdot T \geq \alpha \cdot T - c \cdot T = (\alpha - c) \cdot T \quad , \quad (51)$$

484 and

$$Q_A(t) \leq c \cdot T, \quad \forall A \in \Theta \setminus A_1 \quad . \quad (52)$$

485 From (51) and (52), for each sub-optimal super arm  $A \in \Theta \setminus A_1$ , we have

$$\sqrt{\frac{Q_A(t)+1}{Q_{A_1}(t)+1}} \leq \sqrt{\frac{cT+1}{(\alpha-c)T+1}} = \sqrt{\frac{T/6+1}{4T/6+1}} \leq \frac{1}{2} \quad , \quad (53)$$

486 which gives

$$1 - \sqrt{\frac{Q_A(t)+1}{Q_{A_1}(t)+1}} \geq \frac{1}{2} \quad . \quad (54)$$

487 Let  $p_0 := \frac{1}{8\sqrt{\pi}}e^{-\frac{28}{3}}$ . In the following, we prove that, with at least a constant probability  $p_0$ , a  
 488 sub-optimal super arm is played in each round  $t \geq \alpha \cdot T$  conditioned on event  $B_t^*$  is true. Note  
 489 that whether event  $B_t^*$  is true or not is determined by the history information  $\mathcal{F}_{t-1}$ .



490 Conditioned on any instantiation  $F_{t-1}$  of  $\mathcal{F}_{t-1}$  such that event  $B_t^*$  is true, we have the probability of  
 491 playing a sub-optimal super arm is

$$\begin{aligned}
 & \Pr(\exists A \in \Theta \setminus A_1 : A_t = A \mid \mathcal{F}_{t-1} = F_{t-1}) \\
 & \geq \Pr\left(\exists A \in \Theta \setminus A_1 : \sum_{a \in A} \bar{r}_{a,t} > \sum_{b \in A_1} \bar{r}_{b,t} \mid \mathcal{F}_{t-1} = F_{t-1}\right) \\
 & = \Pr\left(\exists A \in \Theta \setminus A_1 : \sum_{a \in A} \hat{r}_{a,n_{a,t}} + w_t \sqrt{\frac{\gamma \ln t}{n_{a,t} + 1}} > \sum_{b \in A_1} \hat{r}_{b,n_{b,t}} + w_t \sqrt{\frac{\gamma \ln t}{n_{b,t} + 1}} \mid \mathcal{F}_{t-1} = F_{t-1}\right) \\
 & \stackrel{(a)}{=} \Pr\left(\exists A \in \Theta \setminus A_1 : \sum_{a \in A} w_t \sqrt{\frac{\gamma \ln t}{n_{a,t} + 1}} > \sum_{b \in A_1} \left(\sqrt{\gamma} \Delta + w_t \sqrt{\frac{\gamma \ln t}{n_{b,t} + 1}}\right) \mid \mathcal{F}_{t-1} = F_{t-1}\right) \\
 & \stackrel{(b)}{=} \Pr\left(\exists A \in \Theta \setminus A_1 : \sum_{a \in A} w_t \sqrt{\frac{\gamma \ln t}{Q_A(t) + 1}} > \sum_{b \in A_1} \left(\sqrt{\gamma} \Delta + w_t \sqrt{\frac{\gamma \ln t}{Q_{A_1}(t) + 1}}\right) \mid \mathcal{F}_{t-1} = F_{t-1}\right) \\
 & = \Pr\left(\exists A \in \Theta \setminus A_1 : w_t \sqrt{\frac{\ln t}{Q_A(t) + 1}} > \Delta + w_t \sqrt{\frac{\ln t}{Q_{A_1}(t) + 1}} \mid \mathcal{F}_{t-1} = F_{t-1}\right) \\
 & \geq \Pr\left(\exists A \in \Theta \setminus A_1 : w_t \left(1 - \sqrt{\frac{Q_A(t) + 1}{Q_{A_1}(t) + 1}}\right) > \Delta \sqrt{Q_A(t) + 1} \mid \mathcal{F}_{t-1} = F_{t-1}\right) \\
 & \stackrel{(c)}{\geq} \Pr\left(\exists A \in \Theta \setminus A_1 : w_t \cdot \frac{1}{2} > \Delta \sqrt{Q_A(t) + 1} \mid \mathcal{F}_{t-1} = F_{t-1}\right) \\
 & = 1 - \underbrace{\Pr\left(w_t \leq 2\Delta \sqrt{Q_A(t) + 1}, \forall A \in \Theta \setminus A_1 \mid \mathcal{F}_{t-1} = F_{t-1}\right)}_{\lambda} ,
 \end{aligned} \tag{55}$$

492 where step (a) uses the fact that for any base arm in a sub-optimal super arm, the empirical mean is 0,  
 493 whereas for any base arm in the optimal super arm, the empirical mean is  $\sqrt{\gamma} \Delta$  based on our reward  
 494 distribution construction. Step (b) uses the fact that all base arms in a super arm have the same number  
 495 of observations. Step (c) uses the lower bound constructed in (54), i.e.,  $1 - \sqrt{\frac{Q_A(t)+1}{Q_{A_1}(t)+1}} \geq \frac{1}{2}$ . Note  
 496 that for  $\lambda$ , the only randomness is  $w \sim \mathcal{N}(0, 1)$  as all  $Q_A(t)$  are determined by the history.

497 To construct an upper bound for  $\lambda$  above, we construct an optimization problem first using the  
 498 constraint shown in (50). Recall (50) is  $\sum_{A \in \Theta \setminus A_1} Q_A(t) \leq cT$ . We construct the optimization  
 499 problem with the objective function shown in the following (56)

$$\max_{x_1, x_2, \dots, x_{K-1}} \Pr_{w \sim \mathcal{N}(0,1)} \left( w \leq 2\Delta \sqrt{x_a + 1}, \forall a \in [K-1] \right) , \tag{56}$$

500 and constraints shown in (57)

$$x_a \geq 0, \forall a \in [K-1] \quad \text{and} \quad \sum_{a=1}^{K-1} x_a \leq c \cdot T . \tag{57}$$

501 Note that the optimal solution to (56) is the same as the optimal solution to the following objective  
 502 function (58):

$$\max_{x_1, x_2, \dots, x_{K-1}} \Pr_{w \sim \mathcal{N}(0,1)} \left( w \leq \min_{a \in [K-1]} x_a \right) . \tag{58}$$

503 It is not hard to verify that the objective function shown in (58) is maximized when  $x_a = \frac{cT}{K-1} =$   
 504  $\frac{c\sqrt{KT}}{(K-1)\Delta}$  for all  $a \in [K-1]$ . Therefore,  $x_a = \frac{cT}{K-1} = \frac{c\sqrt{KT}}{(K-1)\Delta}$  for all  $a \in [K-1]$  is also  
 505 the optimal solution to (56) and the maximum value of the objective function shown in (56) is  
 506  $\Pr_w \left( w \leq 2\Delta \sqrt{\frac{c\sqrt{KT}}{(K-1)\Delta} + 1} \right)$ .

507 Now, we are ready to construct an upper bound for  $\lambda$  and have

$$\begin{aligned}
\lambda &= \Pr_w \left( w \leq 2\Delta \sqrt{Q_A(t)} + 1, \forall A \in \Theta \setminus A_1 \mid \mathcal{F}_{t-1} = F_{t-1} \right) \\
&\leq \max_{x_1, x_2, \dots, x_{K-1}} \Pr_w \left( w \leq 2\Delta \sqrt{x_a + 1}, \forall a \in [K-1] \mid \mathcal{F}_{t-1} = F_{t-1} \right) \\
&= \max_{x_1, x_2, \dots, x_{K-1}} \Pr_w \left( w \leq 2\Delta \sqrt{x_a + 1}, \forall a \in [K-1] \right) \\
&\stackrel{(a)}{\leq} \Pr_w \left( w \leq 2\Delta \sqrt{\frac{c\sqrt{KT}}{(K-1)\Delta} + 1} \right) \\
&\stackrel{(a)}{\leq} 1 - \frac{1}{8\sqrt{\pi}} \cdot e^{-\frac{7}{2} \cdot 4\Delta^2 \cdot \left(\frac{c\sqrt{KT}}{(K-1)\Delta} + 1\right)} \tag{59} \\
&\stackrel{(c)}{\leq} 1 - \frac{1}{8\sqrt{\pi}} \cdot e^{-\frac{7}{2} \cdot 4\Delta^2 \cdot \frac{2c\sqrt{KT}}{0.5K\Delta}} \\
&= 1 - \frac{1}{8\sqrt{\pi}} \cdot e^{-\frac{7}{2} \cdot 8\sqrt{\frac{K}{T}} \cdot \frac{1}{6} \cdot \frac{\sqrt{KT}}{0.5K}} \\
&= 1 - \frac{1}{8\sqrt{\pi}} e^{-\frac{28}{3}} \\
&= 1 - p_0 \quad ,
\end{aligned}$$

508 where step (a) uses the fact that  $\Pr_w \left( w \leq 2\Delta \sqrt{\frac{c\sqrt{KT}}{(K-1)\Delta} + 1} \right)$  is the maximum value of the  
509 objective function shown in (56). Step (b) uses the one-sided anti-concentration inequality shown  
510 in (6). Step (c) uses the fact that  $K - 1 > 0.5K$ ,  $\Delta = \sqrt{K/T}$ , and when  $T$  is large enough, we  
511 have  $\frac{c\sqrt{KT}}{(K-1)\Delta} \geq 1$ .

512 By plugging the upper bound for  $\lambda$  into (55), we have  
513  $\Pr(\exists A \in \Theta \setminus A_1 : A_t = A \mid \mathcal{F}_{t-1} = F_{t-1}) \geq p_0$ , which concludes the proof for  
514 the statement that with at least a constant probability  $p_0$ , a sub-optimal super arm is played in round  $t$ .

515 To complete the proof, we use the fact that the total regret from round  $t = \alpha T$  to round  $T$  is  
516 at least  $(1 - \alpha)T \cdot p_0 \cdot m \cdot \sqrt{\gamma}\Delta = \Omega(Tm\Delta) = \Omega(Tm\sqrt{K/T}) = \Omega(Tm\sqrt{N/(mT)}) =$   
517  $\Omega(\sqrt{mNT})$ , which is also a regret lower bound for the total regret from round 1 to round  $T$ .  $\square$