

## 481 A Pseudocode

---

### Algorithm 1: Dec\_UCB

---

**Input:**  $\mathbb{G}, T, C_{i,k}(t)$

1 **Initialization** Each agent samples each arm exactly once. Initialize  $z_{i,k}(0) = \bar{x}_{i,k}(0) = X_{i,k}(0)$ ,  
 $m_{i,k}(0) = n_{i,k}(0) = 1$ , and  $C_{i,k}(0) = 0$

2 **for**  $t = 0, \dots, T$  **do**

3      $\mathcal{A}_i = \emptyset$

4     **if**  $n_{i,k}(t) \leq m_{i,k}(t) - M$  **then**

5         Agent  $i$  puts  $k$  into a set  $\mathcal{A}_i$  # exploration consistency requirements

6     **end**

7     **if**  $\mathcal{A}_i = \emptyset$  **then**

8         **for**  $k = 1, \dots, M$  **do**

9              $Q_{i,k}(t+1) = z_{i,k}(t) + C_{i,k}(t)$  # belief update

10         **end**

11          $a_i(t+1) = \arg \max_k Q_{i,k}(t+1)$  # optimal arm in belief

12     **else**

13          $a_i(t+1)$  is randomly chosen from  $\mathcal{A}_i$

14     **end**

15     Agent  $i$  sends  $m_{i,k}(t)$  and  $z_{i,k}(t)$  to each agent  $j$  satisfying  $i \in \mathcal{N}_j$  # information transmission

16     Agent  $i$  receives  $m_{j,k}(t), z_{j,k}(t)$  from each neighbor  $j \in \mathcal{N}_i$

17      $n_{i,k}(t+1) = n_{i,k}(t), \forall k \in [M]$  # information updating

18      $n_{i,a_i(t+1)}(t+1) = n_{i,a_i(t+1)}(t) + 1$

19      $m_{i,k}(t+1) = \max\{n_{i,k}(t+1), m_{j,k}(t), j \in \mathcal{N}_i\}$

20      $z_{i,k}(t+1) = \sum_{j=1}^N w_{ij} z_{j,k}(t) + \bar{x}_{i,k}(t+1) - \bar{x}_{i,k}(t)$

21 **end**

---

## 482 B Analysis and Proofs

483 In this appendix, we provide the analysis of Dec\_UCB and proofs of Theorems 1 and 2.

484 We begin with some basic properties of sub-Gaussian random variables in B.1, and provide analysis  
485 on the exploration “consistency” of Dec\_UCB in B.2, which theoretically validates Remark 2. Based  
486 on the properties in B.1 and results in B.2, we prove Theorems 1 and 2 in B.3 and B.4, respectively.

### 487 B.1 Sub-Gaussian Random Variables

488 A random variable  $X$  with  $\mathbf{E}[X] = \mu$  is called  $\sigma^2$  sub-Gaussian if there is a positive  $\sigma$  such that

$$\mathbf{E}(e^{\lambda(X-\mu)}) \leq e^{\frac{\sigma^2 \lambda^2}{2}}, \quad \forall \lambda \in \mathbb{R}.$$

489 Such  $\sigma^2$  is called a variance proxy, and the smallest variance proxy is called the optimal variance  
490 proxy. Sub-Gaussian random variables have the following three properties.

491 **Lemma 1.** Let  $X$  be any  $\sigma^2$  sub-Gaussian random variable  $\mathbf{E}[X] = \mu$ . Then, for any  $a \geq 0$ ,

$$\mathbf{P}(X - \mu \geq a) \leq e^{-\frac{a^2}{2\sigma^2}}, \quad \mathbf{P}(\mu - X \geq a) \leq e^{-\frac{a^2}{2\sigma^2}}.$$

492 *Proof:* The proof can be found in Section 5.3 in [1]. ■

493 **Lemma 2.** Let  $X_1, \dots, X_n$  be  $n$  independent random variables such that  $X_i$  is  $\sigma_i^2$  sub-Gaussian  
494 random variable. Then,  $X_1 + \dots + X_n$  is  $(\sigma_1^2 + \dots + \sigma_n^2)$  sub-Gaussian.

495 *Proof:* The proof can be found in Section 5.3 in [1]. ■

496 **Lemma 3.** If a random variable  $X$  has a finite mean and  $a \leq X \leq b$  almost surely, then  $X$  is  
497  $\frac{1}{4}(b-a)^2$  sub-Gaussian.

498 *Proof:* The lemma is a direct consequence of Hoeffding’s Lemma in [2]. ■

## 499 B.2 Exploration Consistency of Dec\_UCB

500 In this subsection, we will show that after a finite number of pulls, for each agent  $i$  and each arm  
 501  $k$ , there holds  $n_{i,k}(t) \leq 2 \min_{j \in [N]} n_{j,k}(t)$ ; see Lemma 6. An immediate consequence of this  
 502 property is that  $\max_{j \in [N]} n_{j,k}(t) \leq 2 \min_{j \in [N]} n_{j,k}(t)$ ; in other words, the global maximal number  
 503 of pulls of each arm is no larger than twice the corresponding global minimum. This ensures that the  
 504 explorations of each arm among all the agents are consistent in the sense that no agent is behind too  
 505 much in arm exploration.

506 To proceed, we will make use of the concept of distance in a graph. For a strongly connected graph,  
 507 the distance from vertex  $i$  to another vertex  $j$  is the length of the shortest directed path from  $i$  to  $j$ .  
 508 The definition subsumes the special case of undirected graphs. For an undirected, connected graph,  
 509 the distance between two different vertices is the length of the shortest path connecting them. We use  
 510 the notation  $d_{i,j}$  to denote the distance from vertex  $i$  to vertex  $j$  in a graph  $\mathbb{G}$ , regardless of it being  
 511 directed or undirected. It is natural to define  $d_{i,i} = 0$  for any vertex  $i$ , and easy to see that  $d_{i,j} \leq d$ ,  
 512 where  $d$  is the diameter of  $\mathbb{G}$ .

513 It is worth emphasizing that the following three lemmas hold for both a directed, strongly connected  
 514 neighbor graph  $\mathbb{G}$  and an undirected, connected  $\mathbb{G}$ , using the same proofs.

515 For the purpose of analysis, we define  $n_{i,k}(t) = m_{i,k}(t) = 0$  for all  $i \in [N]$  and  $k \in [M]$  whenever  
 516  $t < 0$ .

517 **Lemma 4.** For any  $i \in [N]$  and  $k \in [M]$ ,

$$m_{i,k}(t) = \max_{j \in [N]} \{n_{j,k}(t - d_{j,i})\}. \quad (9)$$

518 *Proof:* We will prove the lemma by induction on  $t$ . For the basis step, suppose that  $t = 0$ . In this case,  
 519  $m_{i,k}(1) = \max\{n_{i,k}(0), m_{j,k}(0), j \in \mathcal{N}_i\} = 1$ . Note that  $\max_{j \in [N]} \{n_{j,k}(t - d_{j,i})\} = n_{i,k}(0) =$   
 520  $1$ . Thus, (9) holds when  $t = 0$ .

521 For the inductive step, assume (9) holds at time  $t$ , and now consider time  $t + 1$ . Note that

$$\begin{aligned} m_{i,k}(t+1) &= \max\{n_{i,k}(t+1), m_{j,k}(t), j \in \mathcal{N}_i\} \\ &= \max\{n_{i,k}(t+1), n_{h,k}(t - d_{j,h}), h \in [N], j \in \mathcal{N}_i\}. \end{aligned}$$

522 It is easy to see that  $d_{h,i} \leq d_{j,i} + d_{h,j} = 1 + d_{h,j}$ . Since  $n_{i,k}(t)$  is a non-decreasing function of  $t$  by  
 523 its definition,

$$m_{i,k}(t+1) \leq \max_{h \in [N]} \{n_{i,k}(t+1), n_{h,k}(t - d_{h,i} + 1)\} = \max_{j \in [N]} \{n_{j,k}(t - d_{j,i} + 1)\}. \quad (10)$$

524 Fix any vertex  $j \in [N]$  and let  $p = (j, v_{d_{j,i}}, \dots, v_2, i)$  be a shortest path from  $j$  to  $i$  in  $\mathbb{G}$ . From (2),

$$\begin{aligned} m_{i,k}(t+1) &\geq m_{v_2,k}(t) \geq \dots \geq m_{v_{d_{j,i}},k}(t - d_{j,i} + 2) \\ &\geq m_{j,k}(t - d_{j,i} + 1) \geq n_{j,k}(t - d_{j,i} + 1). \end{aligned} \quad (11)$$

525 Since  $j$  is arbitrarily chosen from  $[N]$ , we have  $m_{i,k}(t+1) \geq \max_{j \in [N]} \{n_{j,k}(t - d_{j,i} + 1)\}$ .  
 526 Combining with (10), we have

$$m_{i,k}(t+1) = \max_{j \in [N]} \{n_{j,k}(t - d_{j,i} + 1)\}.$$

527 Thus, (9) also holds at  $t + 1$ , which completes the induction. ■

528 **Lemma 5.** For any  $i \in [N]$  and  $k \in [M]$ ,

$$n_{i,k}(t) > m_{i,k}(t) - M(M + 2d).$$

*Proof:* We will prove the lemma by contradiction. Suppose that, to the contrary,  $\exists i, k_1$  such that  
 $n_{i,k_1}(t) \leq m_{i,k_1}(t) - M(M + 2d)$ . Let  $t'$  denote the first time at which the equality holds, i.e.,

$$n_{i,k_1}(t') = m_{i,k_1}(t') - M(M + 2d).$$

529 Here  $t'$  must exist, since when  $t = 0$ , we have  $n_{i,k}(0) > m_{i,k}(0) - M(M + 2d)$ , since both  
 530  $n_{i,k}(t)$  and  $m_{i,k}(t)$  increase by 0 and 1 at each time instance, if there exists some  $t$  such that

531  $n_{i,k_1}(t) < m_{i,k_1}(t) - M(M + 2d)$ , there must exist a  $t'$  between 0 and  $t$ , such that  $n_{i,k_1}(t') =$   
 532  $m_{i,k_1}(t') - M(M + 2d)$ . According to Lemma 4,  $\exists j \in [N]$  such that

$$m_{i,k_1}(t') = n_{j,k_1}(t' - d_{j,i}). \quad (12)$$

533 Then,

$$n_{j,k_1}(t' - d_{j,i}) - n_{i,k_1}(t') = M(M + 2d). \quad (13)$$

534 Since according to Lemma 4,  $m_{i,k_1}(t) \geq n_{j,k_1}(t - d_{j,i})$  always holds, so for  $t < t'$ , we have

$$n_{j,k_1}(t - d_{j,i}) - n_{i,k_1}(t) \leq m_{i,k_1}(t) - n_{i,k_1}(t) < M(M + 2d). \quad (14)$$

535 Since  $n_{i,k}(t)$  is non-decreasing for all  $i \in [N], k \in [M]$ , (13) and (14) imply that  $n_{j,k_1}(t' - d_{j,i}) >$   
 536  $n_{j,k_1}(t' - d_{j,i} - 1)$ . This further implies that at time  $t' - d_{j,i}$ , agent  $j$  pulls arm  $k_1$ .

537 Since each agent must pull an arm at each time, we have  $\sum_k n_{i,k}(t) = t, \forall i \in [N]$ . Then,

$$\sum_{k \in [M] \setminus k_1} n_{i,k}(t') - \sum_{k \in [M] \setminus k_1} n_{j,k}(t' - d_{j,i}) = M(M + 2d) + d_{j,i}.$$

538 Applying the Pigeonhole principle,  $\exists k_2 \in [M]$  such that

$$n_{i,k_2}(t') - n_{j,k_2}(t' - d_{j,i}) \geq \frac{M(M + 2d)}{M - 1} > M + 2d.$$

539 According to the definition of  $n_{i,k}(t)$ , it is non-decreasing and  $n_{i,k}(t + 1) \leq n_{i,k}(t) + 1$ , we obtain

$$n_{i,k_2}(t') = n_{i,k_2}(t' - d_{j,i} - d_{i,j} + d_{j,i} + d_{i,j}) \leq n_{i,k_2}(t' - d_{j,i} - d_{i,j}) + d_{j,i} + d_{i,j}.$$

540 Thus,

$$\begin{aligned} n_{i,k_2}(t' - d_{j,i} - d_{i,j}) - n_{j,k_2}(t' - d_{j,i}) &> n_{i,k_2}(t') - n_{j,k_2}(t' - d_{j,i}) - d_{j,i} - d_{i,j} \\ &> M + 2d - d_{j,i} - d_{i,j} > M. \end{aligned}$$

541 Using (11), we have  $m_{j,k_2}(t' - d_{j,i}) \geq n_{i,k_2}(t' - d_{j,i} - d_{i,j})$ . Thus,

$$m_{j,k_2}(t' - d_{j,i}) - n_{j,k_2}(t' - d_{j,i}) > M.$$

542 From the above analysis, agent  $j$  must pull arm  $k_1$  at time  $t' - d_{j,i}$ . According to the decision making  
 543 step of the algorithm, there holds

$$m_{j,k_1}(t' - d_{j,i}) - n_{j,k_1}(t' - d_{j,i}) \geq M > 0. \quad (15)$$

544 Note that from (11),

$$m_{i,k_1}(t') \geq m_{j,k_1}(t' - d_{j,i}). \quad (16)$$

545 Combining (12) – (16) together, we have

$$n_{j,k_1}(t' - d_{j,i}) = m_{i,k_1}(t') \geq m_{j,k_1}(t' - d_{j,i}) > n_{j,k_1}(t' - d_{j,i}),$$

546 which is a contradiction. Therefore, the statement of the lemma is true. ■

547 **Lemma 6.** For any  $i \in [N]$  and  $k \in [M]$ , if  $n_{i,k}(t) \geq 2(M^2 + 2Md + d)$ , then

$$n_{i,k}(t) \leq 2 \min_{j \in [N]} n_{j,k}(t).$$

548 *Proof:* From (11),  $m_{j,k}(t) \geq n_{i,k}(t - d_{i,j}), \forall i \in [N]$ . Combining with  $n_{j,k}(t + 1) \leq n_{j,k}(t) + 1$ ,  
 549 we have

$$m_{j,k}(t) \geq n_{i,k}(t) - d_{i,j} \geq n_{i,k}(t) - d.$$

550 From Lemma 5, we have

$$n_{j,k}(t) \geq n_{i,k}(t) - (M^2 + 2Md + d), \forall j, i \in [N].$$

551 Since  $j$  is arbitrarily chosen, when  $n_{i,k}(t) \geq 2(M^2 + 2Md + d)$ ,

$$\min_{j \in [N]} n_{j,k}(t) \geq n_{i,k}(t) - (M^2 + 2Md + d) \geq n_{i,k}(t) - \frac{1}{2}n_{i,k}(t) = \frac{1}{2}n_{i,k}(t),$$

552 which completes the proof. ■

### 553 B.3 Strongly Connected Graphs

554 To prove the performance of Dec\_UCB on strongly connected graphs as stated in Theorem 1, we need  
 555 to provide a tight bound of the variance proxy of  $z_{i,k}(t)$  in this case.

556 **Lemma 7.** For any  $i \in [N]$ ,  $k \in [M]$  and time  $t \geq 0$ , with  $W$  being defined in (5),  $z_{i,k}(t)$  is a  
 557 sub-Gaussian random variable, and when

$$n_{i,k}(t) \geq \max \{L, 2(M^2 + 2Md + d)\},$$

558 the optimal variance proxy of  $z_{i,k}(t)$  is no larger than  $\frac{1}{3n_{i,k}(t)}$ .

559 *Proof.* According to Lemma 3,  $z_{i,k}(t)$  is a sub-Gaussian random variable as it is bounded. From (1)  
 560 and (3), we know that  $z_{i,k}(t)$  is a linear combination of  $X_{j,k}(\tau)$ , for all  $j \in [N], \tau \in \{1, 2, \dots, t\}$ .  
 561 According to Lemma 3 and Lemma 2, in order to find the variance proxy of  $z_{i,k}(t)$ , we need to  
 562 estimate the coordinates of such  $X_{j,k}(\tau)$  first.

563 Note that from (4),

$$\begin{aligned} z_k(t) &= Wz_k(t-1) + \bar{x}_k(t) - \bar{x}_k(t-1) \\ &= W^t z_k(0) + \sum_{\tau=0}^{t-1} W^\tau (\bar{x}_k(t-\tau) - \bar{x}_k(t-\tau-1)) \\ &= \sum_{\tau=0}^{t-1} (W^{t-\tau} - W^{t-\tau-1}) \bar{x}_k(\tau) + \bar{x}_k(t). \end{aligned}$$

564 Thus,

$$z_{i,k}(t) = \sum_j \left\{ \sum_{\tau=0}^{t-1} [W^{t-\tau} - W^{t-\tau-1}]_{ij} \bar{x}_{j,k}(\tau) + [W^0]_{ij} \bar{x}_{j,k}(t) \right\}.$$

565 Denote  $\tau_{i,1}, \tau_{i,2}, \dots, \tau_{i,n_{i,k}(t)}$  as the ascending sequence of all time instances before time  $t$  at which  
 566 agent  $i$  pulls arm  $k$ . From the initialization step of the algorithm, it is clear that  $\tau_{i,1} = 0$ . According  
 567 to (1), if  $\tau_{i,m} \leq \tau < \tau_{i,m+1}$ , we have  $\bar{x}_{i,k}(\tau) = \bar{x}_{i,k}(\tau_{i,m})$ ,  $\forall i \in [N]$ . Then,

$$z_{i,k}(t) = \sum_j \left\{ \sum_{h=1}^{n_{j,k}(t)-1} [W^{t-\tau_{j,h}} - W^{t-\tau_{j,h+1}}]_{ij} \bar{x}_{j,k}(\tau_{j,h}) + [W^{t-\tau_{j,n_{j,k}(t)}}]_{ij} \bar{x}_{j,k}(\tau_{j,n_{j,k}(t)}) \right\}, \quad (17)$$

568 Let  $c_{i,k,j}^{(\tau)}(t)$  be the coefficient of  $X_{j,k}(\tau)$  in  $z_{i,k}(t)$ , it is not hard to see from the above equation that  
 569 when  $\tau \neq \tau_{j,1}, \tau_{j,2}, \dots, \tau_{j,n_{j,k}(t)}$ , we have  $c_{i,k,j}^{(\tau)}(t) = 0$ . Also, when  $\tau = \tau_{j,1}, \tau_{j,2}, \dots, \tau_{j,n_{j,k}(t)}$ ,  
 570 then  $X_{j,k}(\tau) \cdot \mathbb{1}(a_j(\tau) = k) = X_{j,k}(\tau)$  are i.i.d. random variables. Thus, from Lemma 2 and  
 571 Lemma 3,

$$\sigma_{i,k}^2 \triangleq \frac{1}{4} \sum_{j=1}^N \sum_{h=1}^{n_{j,k}(t)} \left| c_{i,k,j}^{(\tau_{j,h})}(t) \right|^2$$

572 is a variance proxy of  $z_{i,k}(t)$ . And we have

$$c_{i,k,j}^{(0)}(t) = \left[ \sum_{h=1}^{n_{j,k}(t)-1} \frac{W^{t-\tau_{j,h}} - W^{t-\tau_{j,h+1}}}{h} + \frac{W^{t-\tau_{j,n_{j,k}(t)}}}{n_{j,k}(t)} \right]_{ij} \quad (18)$$

573 which holds for all  $i \in [N], k \in [M]$ . The equation also can be written as

$$c_{i,k,j}^{(0)}(t) = \left[ W^t - \sum_{h=2}^{n_{j,k}(t)} \frac{W^{t-\tau_{j,h}}}{(h-1)h} \right]_{ij}. \quad (19)$$

574 From (6),  $c_{i,k,j}^{(0)}(t)$  satisfies

$$\begin{aligned} |c_{i,k,j}^{(0)}(t)| &\leq [W_\infty]_{ij} \left( 1 - \sum_{h=2}^{n_{j,k}(t)} \frac{1}{(h-1)h} \right) + c\rho_2^t + c \sum_{h=2}^{n_{j,k}(t)} \frac{\rho_2^{t-\tau_{j,h}}}{(h-1)h} \\ &= \frac{[W_\infty]_{ij}}{n_{j,k}(t)} + c\rho_2^t + \sum_{h=2}^{n_{j,k}(t)} \frac{c\rho_2^{t-\tau_{j,h}}}{(h-1)h}. \end{aligned}$$

575 Since  $0 < \rho_2 < 1$ , the smaller  $t - \tau_{j,h}$  is, the larger the right side of the inequality would be, so  
 576  $\rho_2^{t-\tau_{j,h}} \leq \rho_2^{n_{j,k}(t)-h}$ . Since

$$\begin{aligned} &\sum_{h=2}^{n_{j,k}(t)} \frac{\rho_2^{n_{j,k}(t)-h}}{(h-1)h} \\ &= \sum_{h=2}^{\frac{12N\lceil c \rceil}{12N\lceil c \rceil+1} n_{j,k}(t)} \frac{\rho_2^{n_{j,k}(t)-h}}{(h-1)h} + \sum_{\frac{12N\lceil c \rceil}{12N\lceil c \rceil+1} n_{j,k}(t)+1}^{n_{j,k}(t)} \frac{\rho_2^{n_{j,k}(t)-h}}{(h-1)h} \\ &= \rho_2^{n_{j,k}(t)-2} + (1-\rho_2) \sum_{h=2}^{\frac{12N\lceil c \rceil}{12N\lceil c \rceil+1} n_{j,k}(t)} \frac{\rho_2^{n_{j,k}(t)-h-1}}{h} + \sum_{\frac{12N\lceil c \rceil}{12N\lceil c \rceil+1} n_{j,k}(t)+1}^{n_{j,k}(t)} \frac{\rho_2^{n_{j,k}(t)-h}}{(h-1)h} \\ &\leq \rho_2^{n_{j,k}(t)-2} + (1-\rho_2) \sum_{h=2}^{\frac{12N\lceil c \rceil}{12N\lceil c \rceil+1} n_{j,k}(t)} \rho_2^{n_{j,k}(t)-h-1} + \sum_{\frac{12N\lceil c \rceil}{12N\lceil c \rceil+1} n_{j,k}(t)+1}^{n_{j,k}(t)} \frac{1}{(h-1)h} \\ &\leq \rho_2^{n_{j,k}(t)-2} + \rho_2^{\frac{n_{j,k}(t)}{12N\lceil c \rceil+1}-1} + \frac{1}{12N\lceil c \rceil n_{j,k}(t)}, \end{aligned}$$

577 we obtain that

$$\begin{aligned} |c_{i,k,j}^{(0)}(t)| &\leq \frac{[W_\infty]_{ij}}{n_{j,k}(t)} + c \left( \rho_2^t + \rho_2^{n_{j,k}(t)-2} + \rho_2^{\frac{n_{j,k}(t)}{12N\lceil c \rceil+1}-1} + \frac{1}{12N\lceil c \rceil n_{j,k}(t)} \right) \\ &\leq \frac{[W_\infty]_{ij}}{n_{j,k}(t)} + 3c\rho_2^{\frac{n_{j,k}(t)}{12N\lceil c \rceil+1}-1} + \frac{c}{12N\lceil c \rceil n_{j,k}(t)}. \end{aligned}$$

578 Recall that  $L$  is the smallest value such that when  $t \geq L$ , we have  $3\rho_2^{t/12N\lceil c \rceil} \leq \frac{\rho_2}{24N\lceil c \rceil t}$ . Since  
 579  $c \leq \lceil c \rceil$ , when  $n_{j,k}(t) \geq L$ , we have

$$|c_{i,k,j}^{(0)}(t)| \leq \frac{[W_\infty]_{ij}}{n_{j,k}(t)} + \frac{1}{8Nn_{j,k}(t)}.$$

580 Note that the expression in  $[\cdot]$  of (18) is a summation of  $n_{j,k}(t)$  terms. From the derivation of (18)  
 581 and the definition of  $c_{i,k,j}^{(\tau)}(t)$ , it is straightforward to verify that for each  $l \in \{2, \dots, n_{j,k}(t)\}$ ,

$$c_{i,k,j}^{(\tau_{j,l})} = \left[ \sum_{h=l}^{n_{j,k}(t)-1} \frac{W^{t-\tau_{j,h}} - W^{t-\tau_{j,h+1}}}{h} + \frac{W^{t-\tau_{j,n_{j,k}(t)}}}{n_{j,k}(t)} \right]_{ij},$$

582 in which  $[\cdot]$  is the summation of the last  $n_{j,k}(t) - l + 1$  terms in  $[\cdot]$  of (18). Then, following the same  
 583 steps as above, we can conclude that for all  $l \in \{2, \dots, n_{j,k}(t)\}$ , there holds

$$\left| c_{i,k,j}^{(\tau_{j,l})}(t) \right| \leq \frac{[W_\infty]_{ij}}{n_{j,k}(t)} + \frac{1}{8Nn_{j,k}(t)}$$

584 when  $n_{j,k}(t) \geq L$ . Using Lemma 6, when  $n_{j,k}(t) \geq \max\{L, 2(M^2 + 2Md + d)\}$ , we have

$$\begin{aligned}\sigma_{i,k}^2(t) &= \frac{1}{4} \sum_{j=1}^N \sum_{h=1}^{n_{j,k}(t)} \left| c_{i,k,j}^{(\tau_{j,h})} \right|^2 \\ &\leq \sum_{j=1}^N \frac{[W_\infty]_{ij}^2 + \frac{1}{4N} \cdot [W_\infty]_{ij} + \frac{1}{64N^2}}{4n_{j,k}(t)} \\ &\leq \frac{1}{2n_{i,k}(t)} \left( \sum_{j=1}^N [W_\infty]_{ij}^2 + \frac{1}{4N} \sum_{j=1}^N [W_\infty]_{ij} + \frac{1}{64N} \right).\end{aligned}\quad (20)$$

585 Since  $W$  is stochastic, so is  $W_\infty$ . Then,

$$\sum_{j=1}^N [W_\infty]_{ij} = 1, \quad (21)$$

586 and  $W_\infty \cdot \mathbf{1} = \mathbf{1}$ . Let  $a^\top$  be the “normalized” dominant left eigenvector of  $W$ , i.e.,  $a^\top \cdot \mathbf{1} = 1$  and  
587  $a^\top \cdot W = a^\top$ . Then,  $a^\top \cdot W_\infty = a^\top$  and moreover  $W_\infty = \mathbf{1} \cdot a^\top$  [3]. Since  $W$  is an irreducible  
588 nonnegative matrix, by the Perron-Frobenius Theorem, each entry of  $a$  is positive. With these facts,  
589 we have

$$\begin{aligned}\sum_{j=1}^N [W_\infty]_{ij}^2 &= a^\top a = a^\top W a = \text{tr}(a^\top W a) = \text{tr}(W \cdot a a^\top) \quad (\text{tr}(AB) = \text{tr}(BA)) \\ &\leq \sum_{i=1}^N \frac{\sum_{j=1}^N a_i a_j}{|\mathcal{N}_i|} \quad (a \text{ is positive}) \\ &= \sum_{i=1}^N \frac{a_i}{|\mathcal{N}_i|} \quad (\text{the sum of all entries of } a \text{ equals } 1) \\ &\leq \sum_{i=1}^N \frac{a_i}{2} \quad (|\mathcal{N}_i| \geq 2) \\ &= \frac{1}{2},\end{aligned}\quad (22)$$

590 where  $\text{tr}(\cdot)$  denotes the trace of a square matrix. Substituting (21) and (22) in (20), we have

$$\sigma_{i,k}^2(t) \leq \frac{1}{2n_{i,k}(t)} \cdot \left( \frac{1}{2} + \frac{1}{4N} + \frac{1}{64N} \right) < \frac{1}{3n_{i,k}(t)},$$

591 which completes the proof. ■

592 We are now in a position to prove Theorem 1.

593 **Proof of Theorem 1:** According to Lemma 1 and Lemma 7, when

$$n_{i,k}(t) \geq \max\{L, 2(M^2 + 2Md + d)\},$$

594 we have

$$\mathbf{P} \left( z_{i,k}(t) - \mu_k \geq \sqrt{\frac{4 \log t}{3n_{i,k}(t)}} \right) \leq \exp \left( -\frac{2 \log t}{3n_{i,k}(t) \sigma_{i,k}^2(t)} \right) \leq \frac{1}{t^2}.$$

595 Similarly,

$$\mathbf{P} \left( \mu_k - z_{i,k}(t) \geq \sqrt{\frac{4 \log t}{3n_{i,k}(t)}} \right) \leq \frac{1}{t^2}.$$

Now let us go back to the algorithm and set

$$C_{i,k}(t) = \sqrt{\frac{4 \log t}{3n_{i,k}(t)}}.$$

596 The Decision Making step of Dec\_UCB ensures that agent  $i$  chooses an arm  $k \neq 1$ , instead of the  
597 optimal arm 1, at time  $t$  only if one of the following four cases occurs:

598 Case 1:  $n_{i,k}(t) \leq m_{i,k}(t) - M$ ;

599 Case 2:  $z_{i,k}(t) - \mu_k \geq C_{i,k}(t)$ ;

600 Case 3:  $\mu_1 - z_{i,1}(t) \geq C_{i,1}(t)$ ;

601 Case 4:  $\mu_1 - \mu_k < 2C_{i,k}(t)$ .

602 We are then prepared to find a bound for  $\mathbf{E}(n_{i,k}(T))$ . First, it is easy to verify that when

$$n_{i,k}(t) \geq \frac{16}{3\Delta_k^2} \log T,$$

603 Case 4 does not hold. To proceed, we define  $t'$  as the first time instance, if any, that satisfies

$$n_{i,k}(t') = \max \left\{ \frac{16}{3\Delta_k^2} \log T, L, 2(M^2 + 2Md + d) \right\}.$$

604 In the case when there does not exist such a  $t' \leq T$ , it immediately follows that

$$n_{i,k}(T) < n_{i,k}(t') = \max \left\{ \frac{16}{3\Delta_k^2} \log T, L, 2(M^2 + 2Md + d) \right\}.$$

605 Next consider the case when  $t'$  exists and  $t' \leq T$ . Then,

$$\sum_{t > t'} \left[ \mathbf{P}(z_{i,k}(t) - \mu_k \geq C_{i,k}(t)) + \mathbf{P}(\mu_1 - z_{i,1}(t) \geq C_{i,1}(t)) \right] \leq \sum_{t > t'} \frac{2}{t^2} = \frac{\pi^2}{3},$$

606 which implies that after  $t'$ , the expected number of pulls of agent  $i$  on arm  $k$  due to Case 2 and  
607 Case 3 is no more than  $\frac{\pi^2}{3}$ . Since the difference between  $n_{i,k}(t)$  and  $m_{i,k}(t)$  is at most  $M^2 + 2Md$   
608 by Lemma 5 and  $n_{i,k}(t)$  keeps increasing until the difference is less than  $M$  according to the  
609 Decision Making step of Dec\_UCB, the expected number of pulls due to Case 1 must be no larger  
610 than  $\frac{\pi^2}{3} + M^2 + (2d - 1)M$ . Thus,

$$\begin{aligned} \mathbf{E}(n_{i,k}(T)) &\leq \mathbf{E}(n_{i,k}(T) \mid T \geq t') \\ &= n_{i,k}(t') + \frac{2\pi^2}{3} + M^2 + (2d - 1)M \\ &= \max \left\{ \frac{16}{3\Delta_k^2} \log T, L, 2(M^2 + 2Md + d) \right\} + \frac{2\pi^2}{3} + M^2 + (2d - 1)M. \end{aligned}$$

611 Now we can get an upper bound of agent  $i$ 's regret as follows:

$$\begin{aligned} R_i(T) &= T\mu_1 - \sum_{t=1}^T \mathbf{E}(X_{a_i(t)}) \\ &= \sum_{k: \Delta_k > 0} \mathbf{E}(n_{i,k}(T)) \cdot \Delta_k \\ &\leq \sum_{k: \Delta_k > 0} \left( \max \left\{ \frac{16}{3\Delta_k^2} \log T, L, 2(M^2 + 2Md + d) \right\} + \frac{2\pi^2}{3} + M^2 + (2d - 1)M \right) \Delta_k. \end{aligned}$$

612 This completes the proof. ■

## 613 B.4 Undirected and Connected Graphs

614 In this subsection, we prove the performance of Dec\_UCB on undirected and connected graphs as  
 615 stated in Theorem 2. The procedure is the same as that in B.3. We first provide a tight bound of the  
 616 variance proxy of  $z_{i,k}(t)$ .

617 **Lemma 8.** *For any  $i \in [N]$ ,  $k \in [M]$ , and time  $t$ , with  $W$  being defined in (8),  $z_{i,k}(t)$  is a  
 618 sub-Gaussian random variable, and when*

$$n_{i,k}(t) \geq \max \{L, 2(M^2 + 2Md + d)\},$$

619 *the optimal variance proxy of  $z_{i,k}(t)$  is no larger than  $\frac{3}{4|\mathcal{N}_i|n_{i,k}(t)}$ .*

620 *Proof.* The arguments in the proof of Lemma 7 until (20) still hold here, that is, when  $n_{j,k}(t) \geq$   
 621  $\max\{L, 2(M^2 + 2Md + d)\}$ , we still have

$$\sigma_{i,k}^2(t) \leq \frac{1}{2n_{i,k}(t)} \left( \sum_{j=1}^N [W_\infty]_{ij}^2 + \frac{1}{4N} \sum_{j=1}^N [W_\infty]_{ij} + \frac{1}{64N} \right),$$

622 where now  $W$  is a doubly stochastic matrix defined by (8). Since  $W$  is irreducible and doubly  
 623 stochastic, its dominant left and right eigenvectors of eigenvalue 1 are  $\mathbf{1}^\top$  and  $\mathbf{1}$ , respectively, and  
 624 thus  $W_\infty = \frac{1}{N} \mathbf{1} \mathbf{1}^\top$  [3]. Thus,

$$\sigma_{i,k}^2(t) \leq \frac{1}{2n_{i,k}(t)} \left( \frac{1}{N} + \frac{1}{4N} + \frac{1}{64N} \right) \leq \frac{3}{4Nn_{i,k}(t)} \leq \frac{3}{4|\mathcal{N}_i|n_{i,k}(t)},$$

625 which proves the lemma. ■

626 **Proof of Theorem 2:** According to Lemma 1 and Lemma 8, when

$$n_{i,k}(t) \geq \max \{L, 2(M^2 + 2Md + d)\},$$

627 we have

$$\mathbf{P} \left( z_{i,k}(t) - \mu_k \geq \sqrt{\frac{3 \log t}{|\mathcal{N}_i|n_{i,k}(t)}} \right) \leq \exp \left( -\frac{3 \log t}{2|\mathcal{N}_i|n_{i,k}(t)\sigma_{i,k}^2(t)} \right) \leq \frac{1}{t^2}.$$

628 Similarly,

$$\mathbf{P} \left( \mu_k - z_{i,k}(t) \geq \sqrt{\frac{3 \log t}{|\mathcal{N}_i|n_{i,k}(t)}} \right) \leq \frac{1}{t^2}.$$

629 Again, the Decision Making step of Dec\_UCB guarantees that agent  $i$  chooses an arm  $k \neq 1$  instead  
 630 of the optimal arm 1 at time  $t$  only if one of the four cases listed in the proof of Theorem 1 occurs.  
 631 First, it is easy to verify that when

$$n_{i,k}(t) \geq \frac{12}{|\mathcal{N}_i|\Delta_k^2} \log T,$$

632 Case 4 does not hold. Then, let  $t'$  be the first time instance, if any, that satisfies

$$n_{i,k}(t') = \max \left\{ \frac{12}{|\mathcal{N}_i|\Delta_k^2} \log T, L, 2(M^2 + 2Md + d) \right\}.$$

633 In the case when there does not exist such a  $t' \leq T$ , it immediately follows that

$$n_{i,k}(T) < n_{i,k}(t') = \max \left\{ \frac{12}{|\mathcal{N}_i|\Delta_k^2} \log T, L, 2(M^2 + 2Md + d) \right\}.$$

634 In the other case when  $t'$  exists and  $t' \leq T$ , using the same arguments as in the proof of Theorem 1,  
 635 the expected number of pulls due to Case 1 is no larger than  $\frac{\pi^2}{3} + M^2 + (2d - 1)M$ . Then,

$$\begin{aligned} \mathbf{E}(n_{i,k}(T)) &\leq \mathbf{E}(n_{i,k}(T) \mid T \geq t') \\ &= n_{i,k}(t') + \frac{2\pi^2}{3} + M^2 + (2d - 1)M \\ &= \max \left\{ \frac{12}{|\mathcal{N}_i|\Delta_k^2} \log T, L, 2(M^2 + 2Md + d) \right\} + \frac{2\pi^2}{3} + M^2 + (2d - 1)M. \end{aligned}$$

636 We thus get the following upper bound of agent  $i$ 's regret:

$$\begin{aligned}
R_i(T) &= T\mu_1 - \sum_{t=1}^T \mathbf{E}(X_{a_i(t)}) \\
&= \sum_{k:\Delta_k > 0} \mathbf{E}(n_{i,k}(T)) \cdot \Delta_k \\
&\leq \sum_{k:\Delta_k > 0} \left( \max \left\{ \frac{12}{|\mathcal{N}_i| \Delta_k^2} \log T, L, 2(M^2 + 2Md + d) \right\} + \frac{2\pi^2}{3} + M^2 + (2d-1)M \right) \Delta_k,
\end{aligned}$$

637 which completes the proof. ■

## 638 C Additional Simulations

639 In this appendix, we provide a set of additional simulations to complement those in Section 4. We  
640 compare the homogeneous and heterogeneous settings, and include simulations for special cases of  
641 interest.

### 642 C.1 Homogeneous vs. Heterogeneous Settings

643 Recall that a homogeneous setting requires that all agents observe the same reward distribution for  
644 any given arm, while in a heterogeneous setting, agents may observe different reward distributions  
645 for a given arm, as long as the mean reward for the arm is consistent across all agents. Below we  
646 present side by side homogeneous and heterogeneous simulations for otherwise identical settings,  
647 first for a case where only one class of reward distribution is present, and next in a case where three  
648 different classes of reward distributions are present.

#### 649 C.1.1 One Distribution

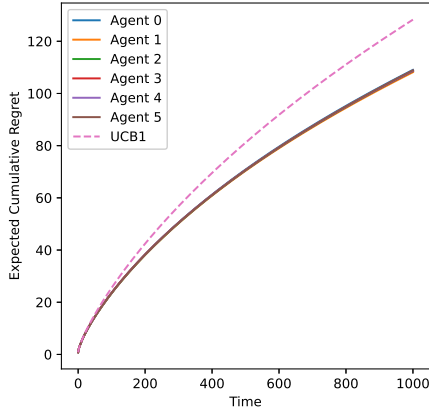
650 In this subsection we consider side by side homogeneous and heterogeneous simulations for strongly  
651 connected, undirected connected, and weakly connected graphs with 6 agents able to choose from  
652 a set of 6 arms for  $T = 1000$  time steps, and average the results over 100 trials. Arm means are  
653 chosen uniformly at random from  $[0.05, 0.95]$ . Reward distributions follow a Beta distribution with a  
654 standard deviation of either 0.01, 0.05, or 0.1. In the homogeneous case, each arm is assigned a Beta  
655 distribution with a certain standard deviation uniformly at random. In the heterogeneous case, each  
656 agent/arm pair is assigned a Beta distribution with a certain standard deviation uniformly at random.  
657 See Figures 5, 6, and 7.

#### 658 C.1.2 Three Distributions

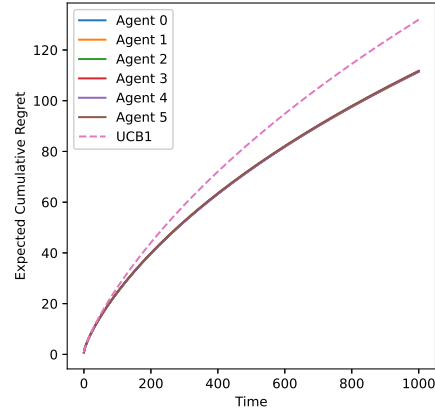
659 In this subsection we consider side by side homogeneous and heterogeneous simulations for strongly  
660 connected, undirected connected, and weakly connected graphs with 15 agents able to choose from  
661 a set of 10 arms for  $T = 1000$  time steps, and average the results over 100 trials. Arm means  
662 are chosen uniformly at random from  $[0.05, 0.95]$ . Reward distributions follow either a truncated  
663 normal distribution with standard deviation 0.2, a Bernoulli distribution, or a uniform distribution  
664 with the greatest width possible in  $[0, 1]$  given a particular mean. In the homogeneous case, each arm  
665 is assigned a distribution uniformly at random. In the heterogeneous case, each arm/agent pair is  
666 assigned a distribution uniformly at random. See Figures 8, 9, and 10.

#### 667 C.1.3 Discussion

668 Our simulations on homogeneous and heterogeneous settings reveal no significant differences between  
669 the two in terms of regret performance. To complement Section 4 with some variety we opted to  
670 use new graph sizes, new quantities of arms, and new distributions, namely, Beta distributions with  
671 standard deviations other than 0.05, a truncated normal distribution with a standard deviation four  
672 times greater than that of Section 4, and a uniform distribution. Momentarily digressing from the  
673 homogeneous vs. heterogeneous comparison, we notice that these different settings still agree with  
674 our results in Section 3.2. If any difference between the homogeneous and heterogeneous settings

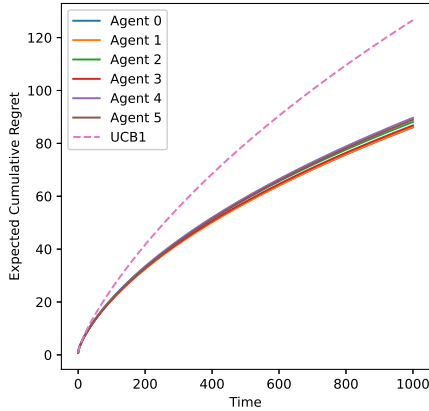


(a) Results for the homogeneous strongly connected setting. Each arm follows a Beta distribution with standard deviation chosen randomly from 0.01, 0.05, 0.1

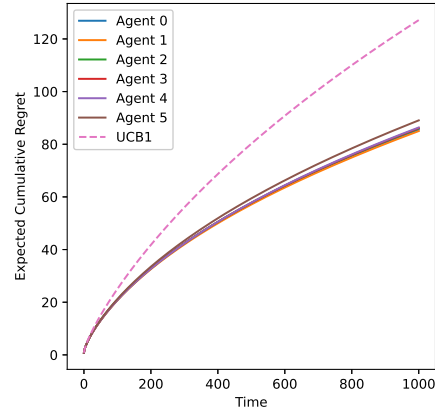


(b) Results for the heterogeneous strongly connected setting. Each agent/arm pair follows a Beta distribution with standard deviation chosen randomly from 0.01, 0.05, 0.1.

Figure 5: Homogeneous vs. heterogeneous regret plots for every agent running Dec\_UCB and the best agent running UCB1. Results for each plot are averaged over 100 different randomly generated (Erdős–Rényi) strongly connected graphs.



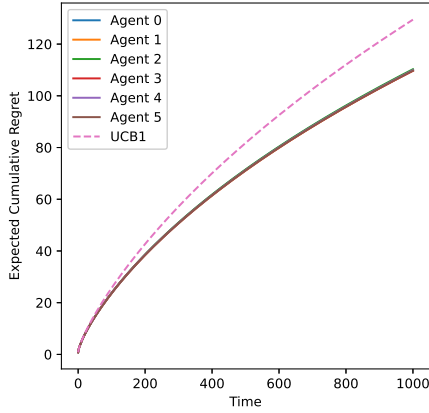
(a) Results for the homogeneous undirected connected setting. Each arm follows a Beta distribution with standard deviation chosen randomly from 0.01, 0.05, 0.1



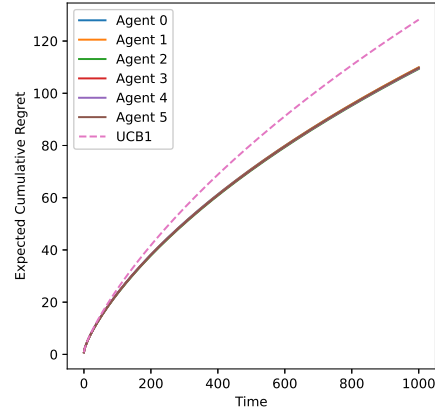
(b) Results for the heterogeneous undirected connected setting. Each agent/arm pair follows a Beta distribution with standard deviation chosen randomly from 0.01, 0.05, 0.1.

Figure 6: Homogeneous vs. heterogeneous regret plots for every agent running Dec\_UCB and the best agent running UCB1. Results for each plot are averaged over 100 different randomly generated (Erdős–Rényi) undirected connected graphs.

675 must be pointed out, we could observe from the figures in C.1.1 and C.1.2 that the homogeneous  
 676 regrets appear to be somewhat smaller than the heterogeneous regrets. However, the difference is  
 677 quite small, and without any theoretical backing we abstain from making any claims about differences  
 678 in the two settings.

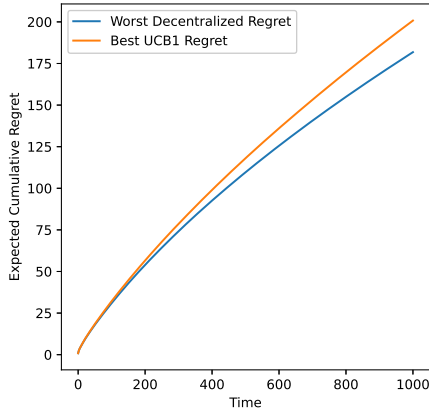


(a) Results for the homogeneous weakly connected setting. Each arm follows a Beta distribution with standard deviation chosen randomly from 0.01, 0.05, 0.1

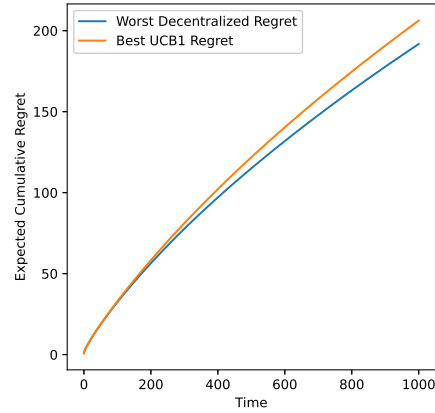


(b) Results for the heterogeneous weakly connected setting. Each agent/arm pair follows a Beta distribution with standard deviation chosen randomly from 0.01, 0.05, 0.1.

Figure 7: Homogeneous vs. heterogeneous regret plots for every agent running Dec\_UCB and the best agent running UCB1. Results for each plot are averaged over 100 different randomly generated (Erdős–Rényi) weakly connected graphs.



(a) Results for the homogeneous strongly connected setting. Each arm follows either a Bernoulli, uniform, or truncated normal ( $\sigma = 0.2$ ) distribution chosen at random.



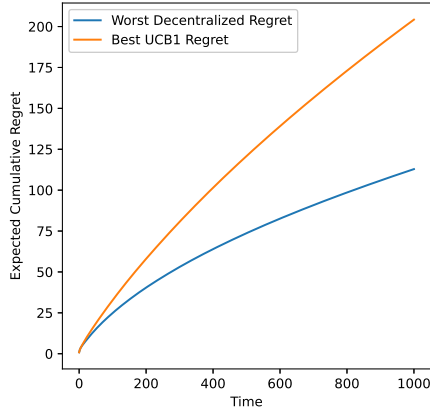
(b) Results for the heterogeneous strongly connected setting. Each agent/arm pair follows either a Bernoulli, uniform, or truncated normal ( $\sigma = 0.2$ ) distribution chosen at random.

Figure 8: Homogeneous vs. heterogeneous regret plots for both the worst performing agent of Dec\_UCB and best performing agent of UCB1 in otherwise identical settings. Results for each plot are averaged over 100 different randomly generated (Erdős–Rényi) strongly connected graphs.

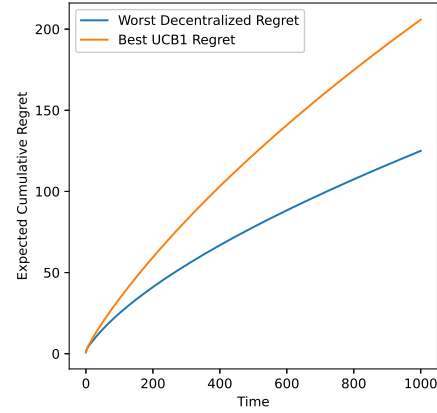
## 679 C.2 Selected Graphs

### 680 C.2.1 Corollary 2: Undirected vs. Strongly Connected

681 We first demonstrate the validity of Corollary 2. This is accomplished by first comparing the  
 682 performance of Dec\_UCB on an undirected, connected graph where all agents have at least two  
 683 neighbors with the performance of Dec\_UCB on a directed, strongly connected graph, keeping all  
 684 other parameters the same. Six agents and six arms were used, with set means of [0.10, 0.25, 0.45,

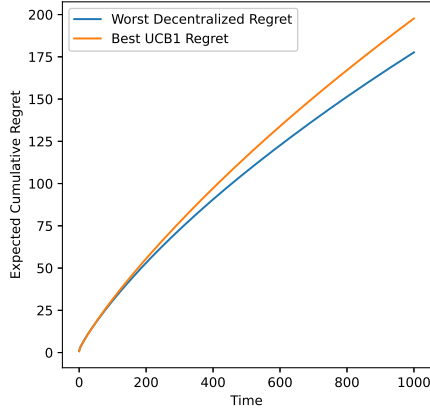


(a) Results for the homogeneous undirected connected setting. Each arm follows either a Bernoulli, uniform, or truncated normal ( $\sigma = 0.2$ ) distribution chosen at random.

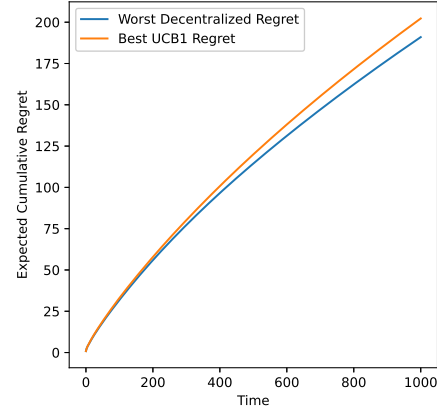


(b) Results for the heterogeneous undirected connected setting. Each agent/arm pair follows either a Bernoulli, uniform, or truncated normal ( $\sigma = 0.2$ ) distribution chosen at random.

Figure 9: Homogeneous vs. heterogeneous regret plots for both the worst performing agent of Dec\_UCB and best performing agent of UCB1 in otherwise identical settings. Results for each plot are averaged over 100 different randomly generated (Erdős–Rényi) undirected connected graphs.



(a) Results for the homogeneous weakly connected setting. Each arm follows either a Bernoulli, uniform, or truncated normal ( $\sigma = 0.2$ ) distribution chosen at random.



(b) Results for the heterogeneous weakly connected setting. Each agent/arm pair follows either a Bernoulli, uniform, or truncated normal ( $\sigma = 0.2$ ) distribution chosen at random.

Figure 10: Homogeneous vs. heterogeneous regret plots for both the worst performing agent of Dec\_UCB and best performing agent of UCB1 in otherwise identical settings. Results for each plot averaged over 100 different randomly generated (Erdős–Rényi) weakly connected graphs.

685 0.65, 0.75, 0.90] for the arms. The reward means  $\mu_k$  were the same for all agents on a given arm, with  
686 each  $\mu_k$  randomly chosen from a uniform distribution on  $[0.05, 0.95]$ . Possible reward distributions  
687 again included the Beta, Bernoulli, and truncated normal distributions, following the three distribution  
688 heterogeneous setting. Experiments were ran for  $T = 1000$  time steps for a total of 100 experiments.  
689 The used graphs and their performances are shown in Figures 11 and 12 respectively. As shown,  
690 every agent in the undirected graph performs better than the agents in the strongly connected graph,  
691 with the performance of an undirected agent being directly related to its number of neighbors.

Next we compare the performance of Dec\_UCB on an undirected, connected graph where all agents are not guaranteed to have at least two neighbors with the performance of Dec\_UCB on the same directed, strongly connected graph. This undirected graph is shown in Figure 13. As illustrated, the undirected agents with only one neighbor (Agents 0 and 5) perform approximately equivalent to the strongly connected agents. The undirected agents with two neighbors perform better than the strongly connected agents.

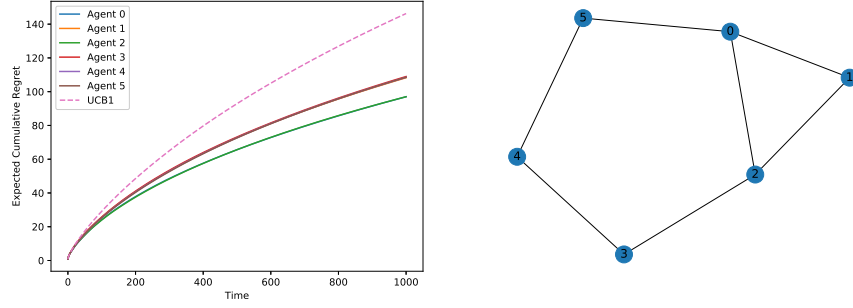


Figure 11: The used undirected graph, ensuring that all agents have at least two neighbors. Reward distributions used vary between agents for a given arm.

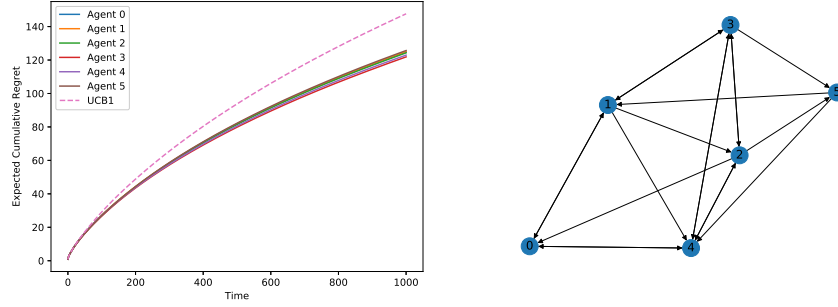


Figure 12: The used strongly connected graph for comparison. Reward distributions used vary between agents for a given arm.

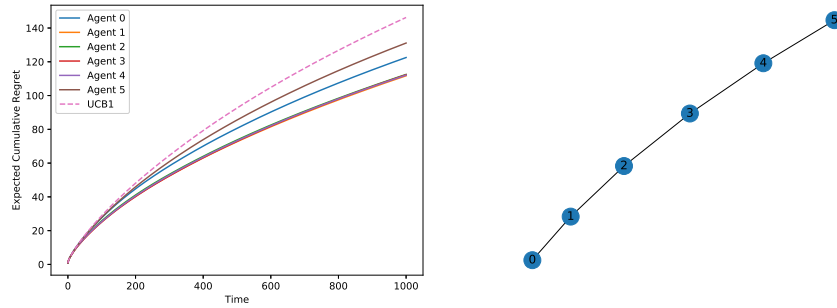


Figure 13: The used undirected graph, allowing for some agents to have only one neighbor. Reward distributions used vary between agents for a given arm.

### 698 C.2.2 (Weakly Connected) Directed Path Graph

699 We next illustrate results for a special weakly connected graph, directed path, for Dec\_UCB, compared  
700 with UCB1. Six agents and six arms were used, with results averaged over 100 experiments,

701 generating new random means  $\mu_k$  each time. The reward means  $\mu_k$  were the same for all agents on  
 702 a given arm, with each  $\mu_k$  randomly chosen from a uniform distribution on  $[0.05, 0.95]$ . Possible  
 703 reward distributions again included the Beta, Bernoulli, and truncated normal distributions, following  
 704 the three distribution heterogeneous setting. Each experiment was ran for  $T = 1000$  time steps. The  
 705 results are illustrated in Figure 14. These results further suggest that Dec\_UCB can perform better  
 706 than UCB1 for this special type of weakly connected graphs, though we are unable to definitively  
 707 confirm or deny this until theoretical backing is developed or a counterexample is discovered.

708 It is worth emphasizing that in such a directed path, the “root” agent (e.g., agent 0 in Figure 14) does  
 709 not have any incoming neighbor and thus receives no external information. The root agent therefore  
 710 solves the bandit problem as in the conventional single-agent case. Since the upper confidence bound  
 711 function of Dec\_UCB is smaller than that of UCB1, it is likely that there exists counterexamples,  
 712 with certain “worst-case” reward distributions, in which such a “root” agent cannot solve the bandit  
 problem using Dec\_UCB.

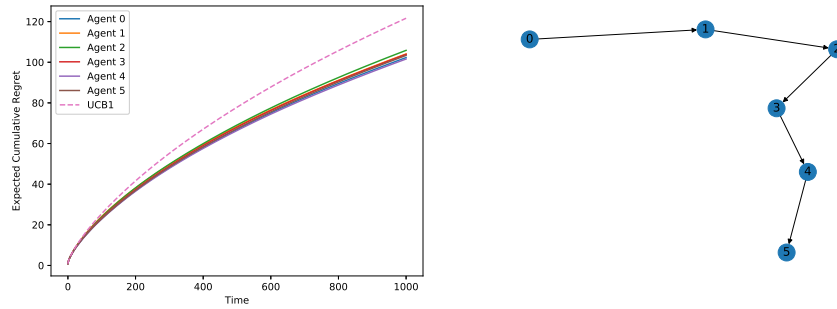


Figure 14: A plot of the regret for both Dec\_UCB and UCB1 of the weakly connected directed path graph, averaged over 100 experiments. Reward distributions used vary between agents for a given arm.

713

## 714 References

- 715 [1] T. Lattimore and C. Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2018.  
 716 [2] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the*  
 717 *American Statistical Association*, 58(301):13–30, 1963.  
 718 [3] L. Xiao and S. Boyd. Fast linear iterations for distributed averaging. *Systems & Control Letters*,  
 719 53(1):65–78, 2004.