



**HAL**  
open science

# Decentralized Learning in Online Queuing Systems

Flore Sentenac, Etienne Boursier, Vianney Perchet

► **To cite this version:**

Flore Sentenac, Etienne Boursier, Vianney Perchet. Decentralized Learning in Online Queuing Systems. NeurIPS, Dec 2021, Online, France. 10.5555/3540261.3541675 . hal-03324084

**HAL Id: hal-03324084**

**<https://hal.science/hal-03324084>**

Submitted on 12 Jul 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# Decentralized Learning in Online Queuing Systems

---

**Flore Sentenac\***  
 CREST, ENSAE Paris, Palaiseau, France  
 flore.sentenac@gmail.com

**Etienne Boursier\***  
 Centre Borelli, ENS Paris-Saclay, France  
 etienne.boursier1@gmail.com

**Vianney Perchet**  
 CREST, ENSAE Paris, Palaiseau, France  
 CRITEO AI Lab, Paris, France  
 vianney.perchet@normalesup.org

## Abstract

Motivated by packet routing in computer networks and resource allocation in radio networks, online queuing systems are composed of queues receiving packets at different rates. Repeatedly, they send packets to servers, each of them treating only at most one packet at a time. In the centralized case, the number of accumulated packets remains bounded (i.e., the system is *stable*) as long as the ratio between service rates and arrival rates is larger than 1. In the decentralized case, individual no-regret strategies ensures stability when this ratio is larger than 2. Yet, myopically minimizing regret disregards the long term effects due to the carryover of packets to further rounds. On the other hand, minimizing long term costs leads to stable Nash equilibria as soon as the ratio exceeds  $\frac{e}{e-1}$ . Stability with decentralized learning strategies with a ratio below 2 was a major remaining question. We first argue that for ratios up to 2, cooperation is required for stability of learning strategies, as selfish minimization of policy regret, a *patient* notion of regret, might indeed still be unstable in this case. We therefore consider cooperative queues and propose the first learning decentralized algorithm guaranteeing stability of the system as long as the ratio of rates is larger than 1, thus reaching performances comparable to centralized strategies.

## 1 Introduction

Inefficient decisions in repeated games can stem from both strategic and learning considerations. First, strategic agents selfishly maximize their own individual reward at others' expense. The *price of anarchy* [Koutsoupias and Papadimitriou, 1999] measures this inefficiency as the social welfare ratio between the best possible situation and the worst Nash equilibrium. Although reaching the best collective outcome might be illusory for selfish agents, considering the worst Nash equilibrium might be too pessimistic. In games with external factors, more complex interactions intervene and might lead the agents to the best equilibrium. Instead, the *price of stability* [Schulz and Moses, 2003] measures the inefficiency by the social welfare ratio between the best possible situation and the best Nash equilibrium.

Second, the agents also have to learn their environment, by repeatedly experimenting different outcomes. Learning equilibria in repeated games is at the core of many problems in computer science and economics [Fudenberg et al., 1998, Cesa-Bianchi and Lugosi, 2006]. The interaction between multiple agents can indeed interfere in the learning process, potentially converging to no or bad equilibria. It is yet known that in repeated games, if all agents follow no internal regret strategies, their actions converge in average to the set of correlated equilibria [Hart and Mas-Colell, 2000, Blum and Monsour, 2007, Perchet, 2014].

---

\*Equal contributions

Many related results are known in the classical repeated games [see e.g., Cesa-Bianchi and Lugosi, 2006, Roughgarden, 2010], where a single game is repeated over independent rounds (but the agents strategies might evolve and depend on the history). Motivated by packet routing in computer networks, Gaitonde and Tardos [2020a] introduced a repeated game with a *carryover* feature: the outcome of a round does not only depend on the actions of the agents, but also on the previous rounds. They consider heterogeneous queues sending packets to servers. If several queues simultaneously send packets to the same server, only the oldest packet is treated by the server.

Because of this carryover effect, little is known about this type of game. In a first paper, Gaitonde and Tardos [2020a] proved that if queues follow suitable no-regret strategies, a ratio of 2 between server and arrival rates leads to stability of the system, meaning that the number of packets accumulated by each queue remains bounded. However, the assumption of regret minimization sort of reflects a myopic behavior and is not adapted to games with carryover. Gaitonde and Tardos [2020b] subsequently consider a patient game, where queues instead minimize their asymptotic number of accumulated packets. A ratio only larger than  $\frac{e}{e-1}$  then guarantees the stability of the system, while a smaller ratio leads to inefficient Nash equilibria. As a consequence, going below the  $\frac{e}{e-1}$  factor requires some level of cooperation between the queues. This result actually holds with perfect knowledge of the problem parameters and it remained even unknown whether decentralized learning strategies can be stable with a ratio below 2.

We first argue that decentralized queues need some level of cooperation to ensure stability with a ratio of rates below 2. Policy regret can indeed be seen as a patient alternative to the regret notion. Yet even minimizing the policy regret might lead to instability when this ratio is below 2. An explicit decentralized cooperative algorithm called ADEQUA (A DEcentralized QUEuing Algorithm) is thus proposed. It is the first decentralized learning algorithm guaranteeing stability when this ratio is only larger than 1. ADEQUA does not require communication between the queues, but uses synchronisation between them to accurately estimate the problem parameters and avoid interference when sending packets. Our main result is given by Theorem 1 below, whose formal version, Theorem 5 in Section 4, also provides bounds on the number of accumulated packets.

**Theorem 1** (Theorem 5, informal). *If the ratio between server rates and arrival rates is larger than 1 and all queues follow ADEQUA, the system is strongly stable.*

The remaining of the paper is organised as follows. The model and existing results are recalled in Section 2. Section 3 argues that cooperation is required to guarantee stability of learning strategies when the ratio of rates is below 2. ADEQUA is then presented in Section 4, along with insights for the proof of Theorem 1. Section 5 finally compares the behavior of ADEQUA with no-regret strategies on toy examples and empirically confirms the different known theoretical results.

## 1.1 Additional related work

Queuing theory includes applications in diverse areas such as computer science, engineering, operation research [Shortle et al., 2018]. Borodin et al. [1996] for example use the stability theorem of Pemantle and Rosenthal [1999], which was also used by Gaitonde and Tardos [2020b], to study the problem of packet routing through a network. Our setting is the single-hop particular instance of throughput maximization in wireless networks. Motivated by resource allocation in multihop radio problem, packets can be sent through more general routing paths in the original problem. Tassioulas and Ephremides [1990] proposed a first stable centralized algorithm, when the service rates are known *a priori*. Stable decentralized algorithms were later introduced in specific cases [Neely et al., 2008, Jiang and Walrand, 2009, Shah and Shin, 2012], when the rewards  $X_k(t)$  are observed before deciding which server to send the packet. The main challenge then is coordination and queues aim at avoiding collisions with each other. The proposed algorithms are thus not adapted to our setting, where both coordination between queues and learning the service rates are required. We refer the reader to [Georgiadis et al., 2006] for an extended survey on resource allocation in wireless networks.

Krishnasamy et al. [2016] first considered online learning for such queuing systems model, in the simple case of a single queue. It is a particular instance of stochastic multi-armed bandits, a celebrated online learning model, where the agent repeatedly takes an action within a finite set and observes its associated reward. This model becomes intricate when considering multiple queues, as they interfere when choosing the same server. It is then related to the multiplayer bandits problem

which considers multiple players simultaneously pulling arms. When several of them pull the same arm, a *collision* occurs and they receive no reward [Anandkumar et al., 2010].

The collision model is here different as one of the players still gets a reward. It is thus even more closely related to competing bandits [Liu et al., 2020a,b], where arms have preferences over the players and only the most preferred player pulling the arm actually gets the reward. Arm preferences are here not fixed and instead depend on the packets’ ages. While collisions can be used as communication tools between players in multiplayer bandits [Bistritz and Leshem, 2018, Boursier and Perchet, 2019, Mehrabian et al., 2020, Wang et al., 2020], this becomes harder with an asymmetric collision model as in competing bandits. However, some level of communication remains possible [Sankararaman et al., 2020, Basu et al., 2021]. In queuing systems, collisions are not only asymmetric, but depend on the age of the sent packets, making such solutions unsuited.

While multiplayer bandits literature considers cooperative players, Boursier and Perchet [2020] showed that cooperative algorithms could be made robust to selfish players. On the other hand, competing bandits consider strategic players and arms as the goal is to reach a bipartite stable matching between them. Despite being cooperative, ADEQUA also has strategic considerations as the queues’ strategy converges to a correlated equilibrium of the patient game described in Section 2.

An additional difficulty here appears as queues are asynchronous: they are not active at each round, but only when having packets left. This is different from the classical notion of asynchronicity [Bonneto et al., 2017], where players are active at each round with some fixed probability. Most strategies in multiplayer bandits rely on synchronisation between the players [Boursier and Perchet, 2019] to avoid collisions. While such a level of synchronisation is not possible here, some lower level of synchronisation is still used to avoid collisions between queues.

## 2 Queuing Model

We consider a queuing system composed of  $N$  queues and  $K$  servers, associated with vectors of arrival and service rates  $\lambda, \mu$ , where at each time step  $t = 1, 2, \dots$ , the following happens:

- each queue  $i \in [N]$  receives a new packet with probability  $\lambda_i \in [0, 1]$ , that is marked with the timestamp of its arrival time. If the queue currently has packet(s) on hold, it sends one of them to a chosen server  $j$  based on its past observations.
- Each server  $j \in [K]$  attempts to clear the oldest packet it has received, breaking ties uniformly at random. It succeeds with probability  $\mu_j \in [0, 1]$  and otherwise sends it back to its original queue, as well as all other unprocessed packets.

At each time step, a queue only observes whether or not the packet sent (if any) is cleared by the server. We note  $Q_t^i$  the number of packets in queue  $i$  at time  $t$ . Given a packet-sending dynamics, the system is **stable** if, for each  $i$  in  $[N]$ ,  $Q_t^i/t$  converges to 0 almost surely. It is **strongly stable**, if for any  $r, t \geq 0$  and  $i \in [N]$ ,  $\mathbb{E}[(Q_t^i)^r] \leq C_r$ , where  $C_r$  is an arbitrarily large constant, depending on  $r$  but not  $t$ . Without ambiguity, we also say the policy or the queues are (strongly) stable. Naturally, a strongly stable system is also stable [Gaitonde and Tardos, 2020a].

In the following,  $x_{(i)}$  will denote the  $i$ -th order statistics of a vector  $\mathbf{x}$ , i.e.,  $\lambda_{(1)} \geq \lambda_{(2)} \geq \dots \geq \lambda_{(N)}$  and  $\mu_{(1)} \geq \dots \geq \mu_{(K)}$ . Without loss of generality, we assume  $K \geq N$  (otherwise, we simply add fictitious servers with 0 service rate). The key quantity of a system is its **slack**, defined as the largest real number  $\eta$  such that:

$$\sum_{i=1}^k \mu_{(i)} \geq \eta \sum_{i=1}^k \lambda_{(i)}, \quad \forall k \leq N.$$

We also denote by  $\mathcal{P}([K])$  the set of probability distributions on  $[K]$  and by  $\Delta$  the **margin** of the system defined by

$$\Delta := \min_{k \in [N]} \frac{1}{k} \sum_{i=1}^k (\mu_{(i)} - \lambda_{(i)}). \quad (1)$$

Notice that the alternative system where  $\tilde{\lambda}_i = \lambda_i + \Delta$  and  $\tilde{\mu}_k = \mu_k$  has a slack 1. In that sense,  $\Delta$  is the largest *margin* between service and arrival rates that all queues can individually get in the system. Note that if  $\eta > 1$ , then  $\Delta > 0$ . We now recall existing results for this problem, summarized in Figure 1 below.

**Theorem 2** (Marshall et al. 1979). *For any instance, there exists a strongly stable centralized policy if and only if  $\eta > 1$ .*

**Theorem 3** (Gaitonde and Tardos 2020a, informal). *If  $\eta > 2$ , queues following appropriate no regret strategies are strongly stable.*

*For each  $N > 0$ , there exists a system and a dynamic s.t.  $\eta > 2 - o(1/N)$ , all queues follow appropriate no-regret strategies, but they are not strongly stable.*

In the above theorem, an *appropriate no regret strategy* is a strategy such that there exists a partitioning of the time into successive windows, for which the incurred regret is  $o(w)$  with high probability on any window of length  $w$ . This for example includes the EXP3.P.1 algorithm [Auer et al., 2002] where the  $k$ -th window has length  $2^k$ .

The patient queuing game  $\mathcal{G} = ([N], (c_i)_{i=1}^n, \boldsymbol{\mu}, \boldsymbol{\lambda})$  is defined as follows. The strategy space for each queue is  $\mathcal{P}([K])$ . Let  $\mathbf{p}_{-i} \in (\mathcal{P}([K]))^{N-1}$  denote the vector of fixed distributions for all queues over servers, except for queue  $i$ . The cost function for queue  $i$  is defined as:

$$c_i(p_i, \mathbf{p}_{-i}) = \lim_{t \rightarrow +\infty} T_t^i / t,$$

where  $T_t^i$  is the age of the oldest packet in queue  $i$  at time  $t$ . Bounding  $T_t^i$  is equivalent to bounding  $Q_t^i$ .

**Theorem 4** (Gaitonde and Tardos 2020b, informal). *If  $\eta > \frac{e}{e-1}$ , any Nash equilibrium of the patient game  $\mathcal{G}$  is stable.*

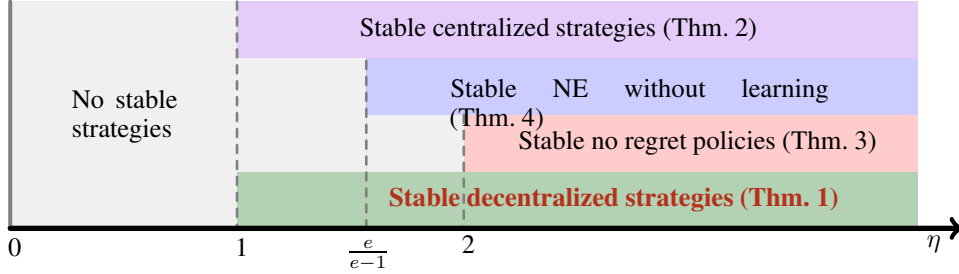


Figure 1: Existing results depending on the slack  $\eta$ . Our result is highlighted in red.

### 3 The case for a cooperative algorithm

According to Theorems 3 and 4, queues that are patient enough and select a fixed randomization over the servers are stable over a larger range of slack  $\eta$  than queues optimizing their individual regret. A key difference between the two settings is that when minimizing their regret, queues are myopic, which is formalized as follows. Let  $a_{1:t} = (a_1, \dots, a_t)$  be the vector of actions played by a queue up to time  $t$  and let  $\nu_t(a_{1:t})$  be the indicator that it cleared a packet at iteration  $t$ . Classical (external) regret of queue  $i$  over horizon  $T$  is then defined as:

$$R_i^{\text{ext}}(T) := \max_{p \in \mathcal{P}([K])} \sum_{t=1}^T \mathbb{E}_{\tilde{a}_t \sim p} [\nu_t(a_{1:t-1}, \tilde{a}_t)] - \sum_{t=1}^T \nu_t(a_{1:t}).$$

Thus minimizing the external regret is equivalent to maximizing the instant rewards at each iteration, ignoring the consequences of the played action on the state of the system. However, in the context of queuing systems, the actions played by the queues change the state of the system. Notably, letting other queues clear packets can be in the best interest of a queue, as it may give it priority in the subsequent iterations where it holds older packets. Since the objective is to maximize the total number of packets cleared, it seems adapted to minimize a *patient* version of the regret, namely the policy regret [Arora et al., 2012], rather than the external regret, which is defined by

$$R_i^{\text{pol}}(T) := \max_{p \in \mathcal{P}([K])} \sum_{t=1}^T \mathbb{E}_{\tilde{a}_{1:t} \sim \otimes_{i=1}^t p} [\nu_t(\tilde{a}_{1:t})] - \sum_{t=1}^T \nu_t(a_{1:t}).$$

That is,  $R_i^{\text{pol}}(T)$  is the expected difference between the number of packets queue  $i$  cleared and the number of packets it would have cleared over the whole period by playing a fixed (possibly random) action, taking into account how this change of policy would affect the state of the system.

However, as stated in Proposition 1, optimizing this patient version of the regret rather than the myopic one could not guarantee stability on a wider range of slack value. This suggests that adding only patience to the learning strategy of the queues is not enough to go beyond a slack of 2, and that any strategy beating that factor 2 must somewhat include synchronisation between the queues.

**Proposition 1.** *Consider the partition of the time  $t = 1, 2, \dots$  into successive windows, where  $w_k = k^2$  is the length of the  $k$ -th one. For any  $N \geq 2$ , there exists an instance with  $2N$  queues and servers, with slack  $\eta = 2 - \mathcal{O}\left(\frac{1}{N}\right)$ , s.t., almost surely, each queue’s policy regret is  $o(w_k)$  on all but finitely many of the windows, but the system is not strongly stable.*

To ease comparison, the formulation in the above proposition matches that of the counter-example used to prove Theorem 3 (Gaitonde and Tardos 2020a). In that counter-example, a set of system parameters, for which any no external regret policies were unstable, was exhibited. Whereas we exhibit in our case a *specific* strategy that satisfies the no policy regret condition, but is unstable.

*Sketch of proof.* Consider a system with  $2N$  queues and servers with  $\lambda_i = 1/2N$  and  $\mu_i = 1/N - 1/4N^2$  for all  $i \in [2N]$ . The considered strategy profile is the following. For each  $k \geq 0$ , the  $k^{\text{th}}$  time window is split into two stages. During the first stage, of length  $\lceil \alpha w_k \rceil$ , queues  $2i$  and  $2i + 1$  both play server  $2i + t \pmod{2N}$  at iteration  $t$ , for all  $i \in [N]$ . During the second stage of the time window, queue  $i$  plays server  $i + t \pmod{2N}$  at iteration  $t$ . This counter example, albeit very specific, illustrates well how when the queues are highly synchronised, it is better to remain synchronized rather than deviate, even if the synchronisation is suboptimal in terms of stability. The complete proof is provided in Appendix C.

Queues following this strategy accumulate packets during the first stage, and clear more packets than they receive during the second stage. The value of  $\alpha$  is tuned so that the queues still accumulate a linear portion of packets during each time window. For those appropriate  $\alpha$ , the system is unstable.

Now suppose that queue  $i$  deviates from the strategy and plays a fixed action  $p \in \mathcal{P}([K])$ . In the first stage of each time window, queue  $i$  can clear a bit more packets than it would by not deviating. However, during the second stage, it is no longer synchronised with the other queues and collides with them a large number of times. Because of those collisions, it will accumulate many packets. In the detailed analysis, we demonstrate that, in the end, for appropriate values of  $\alpha$ , queue  $i$  accumulates more packets than it would have without deviating.  $\square$

According to Theorem 4, the factor  $\frac{e}{e-1}$  can be seen as the price of anarchy of the problem, as for slacks below, the worst Nash equilibria might be unstable. On the other hand, it is known that for any slack above 1, there exists a centralized stable strategy. This centralized strategy actually consists in queues playing the same joint probability at each time step, independently from the number of accumulated packets. As a consequence, it is also a correlated equilibrium of the patient game and 1 can be seen as the correlated price of stability.

## 4 A decentralized algorithm

This section describes the decentralized algorithm ADEQUA, whose pseudocode is given in Algorithm 1. Due to space constraints, all the proofs are postponed to Appendix D. ADEQUA assumes all queues *a priori* know the number  $N$  of queues in the game and have a unique rank or *id* in  $[N]$ . Moreover, the existence of a shared randomness between all queues is assumed. The *id* assumption is required to break the symmetry between queues and is classical in multiplayer bandits without collision information. On the other side, the shared randomness assumption is equivalent to the knowledge of a common seed for all queues, which then use this common seed for their random generators. A similar assumption is used in multiplayer bandits [Bubeck et al., 2020].

ADEQUA is inspired by the celebrated  $\varepsilon$ -greedy strategy. With probability  $\varepsilon_t = (N + K)t^{-\frac{1}{3}}$ , at each time step, queues explore the different parameters  $\lambda_i$  and  $\mu_i$  as described below. Otherwise with probability  $1 - \varepsilon_t$ , they exploit the servers. Each queue  $i$  then sends a packet to a server following a policy solely computed from its local estimates  $\hat{\lambda}^i, \hat{\mu}^i$  of the problem parameters  $\lambda$  and  $\mu$ . The shared randomness is here used so that exploration simultaneously happens for all queues. If exploration/exploitation was not synchronized between the queues, an exploiting queue could collide with an exploring queue, biasing the estimates  $\hat{\lambda}^i, \hat{\mu}^i$  of the latter.

---

**Algorithm 1: ADEQUA**

---

**input :**  $N$  (number of queues),  $i \in [N]$  (queue  $id$ )

```
1 for  $t = 1, \dots, \infty$  do
2    $\hat{P} \leftarrow \phi(\hat{\lambda}, \hat{\mu})$  and  $\hat{A} \leftarrow \psi(\hat{P})$  where  $\phi$  and  $\psi$  are resp. defined by Equations (3) and (4)
3   Draw  $\omega_1 \sim \text{Bernoulli}((N + K)t^{-\frac{1}{5}})$  and  $\omega_2 \sim \mathcal{U}(0, 1)$  // shared randomness
4   if  $\omega_1 = 1$  then EXPLORE( $i$ ) // exploration
5   else Pull  $\hat{A}(\omega_2)(i)$  // exploitation
6 end
```

---

**Exploration.** When exploring, queues choose either to explore the servers' parameters  $\mu_k$  or the other queues' parameters  $\lambda_i$  as described in Algorithm 2 below. In the former case, all queues choose different servers at random (if they have packets to send). These rounds are used to estimate the servers means:  $\hat{\mu}_k^i$  is the empirical mean of server  $k$  observed by the queue  $i$  for such rounds. Thanks to the shared randomness, queues pull different servers here, making the estimates unbiased.

In the latter case, queues explore each other in a pairwise fashion. When queues  $i$  and  $j$  explore each other at round  $t$ , each of them sends their **most recent** packet to some server  $k$ , chosen uniformly at random, if and only if a packet appeared during round  $t$ . In that case, we say that *the queue  $i$  explores  $\lambda_j$*  (and vice versa). To make sure that  $i$  and  $j$  are the only queues choosing the server  $k$  during this step, we proceed as follows:

- queues sample a matching  $\pi$  between queues at random. To do so, the queues use the same method to plan an all-meet-all (or round robin) tournament, for instance Berger tables [Berger, 1899], and choose uniformly at random which round of the tournament to play. If the number of queues  $N$  is odd, in each round of the tournament, one queue remains alone and does nothing.
- the queues draw the same number  $l \sim \mathcal{U}([K])$  with their shared randomness. For each pair of queues  $(i, j)$  matched in  $\pi$ , associate  $k_{(i,j)} = l + \min(i, j) \pmod{K} + 1$  to this pair. The queues  $i$  and  $j$  then send to the server  $k_{(i,j)}$ .

As we assumed that the server breaks ties in the packets' age uniformly at random, the queue  $i$  clears with probability  $(1 - \frac{\lambda_j}{2})\bar{\mu}$ , where  $\bar{\mu} = \frac{1}{K} \sum_{k=1}^K \mu_k$ . Thanks to this,  $\lambda_j$  is estimated by queue  $i$  as:

$$\hat{\lambda}_j^i = 2 - 2\hat{S}_j^i / \tilde{\mu}^i, \quad (2)$$

where  $\tilde{\mu}^i = \frac{\sum_{k=1}^K N_k^i \hat{\mu}_k^i}{\sum_{k=1}^K N_k^i}$ ,  $N_k^i$  is the number of *exploration* pulls of server  $k$  by queue  $i$  and  $\hat{S}_j^i$  is the empirical probability of clearing a packet observed by queue  $i$  when exploring  $\lambda_j$ .

---

**Algorithm 2: EXPLORE**

---

**input :**  $i \in [N]$  // queue  $id$

```
1  $k \leftarrow 0$ 
2 Draw  $n \sim \mathcal{U}([N + K])$  // shared randomness
3 if  $n \leq K$  then // explore  $\mu$ 
4    $k \leftarrow n + i \pmod{K} + 1$ 
5   Pull  $k$ ; Update  $N_k$  and  $\hat{\mu}_k$ 
6 else // explore  $\lambda$ 
7   Draw  $r \sim \mathcal{U}([N])$  and  $l \sim \mathcal{U}([K])$  // shared randomness
8    $j \leftarrow r^{\text{th}}$  opponent in the all-meet-all tournament planned according to Berger tables
9    $k \leftarrow l + \min(i, j) \pmod{K} + 1$ 
10  if  $k \neq 0$  and packet appeared at current time step then // explore  $\lambda_j$  on server  $k$ 
11  | Pull  $k$  with most recent packet; Update  $\hat{S}_j$  and  $\hat{\lambda}_j$  according to Equation (2)
12  end
13 end
```

---

**Remark 1.** The packet manipulation when exploring  $\lambda_j$  strongly relies on the servers tie breaking rules (uniformly at random). If this rule was unknown or not explicit, the algorithm can be adapted: when queue  $i$  explores  $\lambda_j$ , queue  $j$  instead sends the packet generated at time  $t - 1$  (if it exists), while queue  $i$  still sends the packet generated at time  $t$ . In that case, the clearing probability for queue  $i$

is exactly  $(1 - \lambda_j)\bar{\mu}$ , allowing to estimate  $\lambda_j$ . Anticipating the nature of the round  $t$  (exploration vs. exploitation) can be done by drawing  $\omega_1 \sim \text{Bernoulli}(\varepsilon_t)$  at time  $t - 1$ . If  $\omega_1 = 1$ , the round  $t$  is exploratory and the packet generated at time  $t - 1$  is then kept apart by the queue  $j$ .

To describe the exploitation phase, we need a few more notations. We denote by  $\mathfrak{B}_K$  the set of doubly stochastic matrices (non-negative matrices such that each of its rows and columns sums to 1) and by  $\mathfrak{S}_K$  the set of permutation matrices in  $[K]$  (a permutation matrix will be identified with its associated permutation for the sake of cumbersomeness).

A **dominant mapping** is a function  $\phi : \mathbb{R}^N \times \mathbb{R}^K \rightarrow \mathfrak{B}_K$  which, from  $(\lambda, \mu)$ , returns a doubly stochastic matrix  $P$  such that  $\lambda_i < (P\mu)_i$  for any  $i \in [N]$  if it exists (and the identity matrix otherwise).

A **BvN** (Birkhoff von Neumann) **decomposition** is a function  $\psi : \mathfrak{B}_K \rightarrow \mathcal{P}(\mathfrak{S}_K)$  that associates to any doubly stochastic matrix  $P$  a random variable  $\psi(P)$  such that  $\mathbb{E}[\psi(P)] = P$ ; stated otherwise, it expresses  $P$  as a convex combination of permutation matrices. For convenience, we will represent this random variable as a function from  $[0, 1]$  (equipped with the uniform distribution) to  $\mathfrak{S}_K$ .

Informally speaking, those functions describe the strategies queues would follow in the centralized case: a dominant mapping gives adequate marginals ensuring stability (since the queue  $i$  clears in expectation  $(P\mu)_i$  packets at each step, which is larger than  $\lambda_i$  by definition), while a BvN decomposition describes the associated coupling to avoid collisions. Explicitly, the joint strategy is for each queue to draw a shared random variable  $\omega_2 \sim \mathcal{U}(0, 1)$  and to choose servers according to the permutation  $\psi(\phi(\lambda, \mu))(\omega_2)$

**Exploitation.** In a decentralized system, each queue  $i$  computes a mapping  $\hat{A}^i := \psi(\phi(\hat{\lambda}^i, \hat{\mu}^i))$  solely based on its own estimates  $\hat{\lambda}^i, \hat{\mu}^i$ . A shared variable  $\omega_2 \in [0, 1]$  is then generated uniformly at random and queue  $i$  sends a packet to the server  $\hat{A}^i(\omega_2)(i)$ . If all queues knew exactly the parameters  $\lambda, \mu$ , the computed strategies  $\hat{A}^i$  would be identical and they would follow the centralized policy described above.

However, the estimates  $(\hat{\lambda}^i, \hat{\mu}^i)$  are different between queues. The usual dominant mappings and BvN decompositions in the literature are non-continuous. Using those, even queues with close estimates could have totally different  $\hat{A}^i$ , and thus collide a large number of times, which would impede the stability of the system. Regular enough dominant mappings and BvN decompositions are required, to avoid this phenomenon. The design of  $\phi$  and  $\psi$  is thus crucial and appropriate choices are given in the following Sections 4.1 and 4.2. Nonetheless, they can be used in some black-box fashion, so we provide for the sake of completeness sufficient conditions for stability, as well as a general result depending on the properties of  $\phi$  and  $\psi$ , in Appendix A.

**Remark 2.** The exploration probability  $t^{-\frac{1}{5}}$  gives the smallest theoretical dependency in  $\Delta$  in our bound. A trade-off between the proportion of exploration rounds and the speed of learning indeed appears in the proof of Theorem 1. Exploration rounds have to represent a small proportion of the rounds, as the queues accumulate packets when exploring. On the other hand, if queues explore more often, the regime where their number of packets decreases starts earlier. A general stability result depending on the choice of this probability is given by Theorem 6 in Appendix A. Yet in Section 5, taking a probability  $t^{-\frac{1}{4}}$  empirically performs better as it speeds up the exploration.

#### 4.1 Choice of a dominant mapping

Recall that a dominant mapping takes as inputs  $(\lambda, \mu)$  and returns, if possible, a doubly stochastic matrix  $P$  such that

$$\lambda_i < \sum_{k=1}^K P_{i,k} \mu_k \text{ for all } i \in [N].$$

The usual dominant mappings sort the vector  $\lambda$  and  $\mu$  in descending orders [Marshall et al., 1979]. Because of this operation, they are non-continuous and we thus need to design a regular dominant mapping satisfying the above property. Inspired by the log-barrier method, it is done by taking the minimizer of a strongly convex program as follows

$$\phi(\lambda, \mu) = \arg \min_{P \in \mathfrak{B}_K} \max_{i \in [N]} -\ln \left( \sum_{j=1}^K P_{i,j} \mu_j - \lambda_i \right) + \frac{1}{2K} \|P\|_2^2. \quad (3)$$



Although the objective function is non-smooth because of the max operator, it enforces fairness between queues and leads to a better regularity of the arg min.

**Remark 3.** Computing  $\phi$  requires solving a non-smooth strongly convex minimization problem. This cannot be computed exactly, but a good approximation can be quickly obtained using the scheme described in Appendix B. If this approximation error is small enough, it has no impact on the stability bound of Theorem 5. It is thus ignored for simplicity, i.e., we assume in the following that  $\phi(\lambda, \mu)$  is exactly computed at each step.

As required,  $\phi$  always returns a matrix  $P$  satisfying that  $\lambda < P\mu$  if possible, since otherwise the objective is infinite (and in that case we assume that  $\phi$  returns the identity matrix). Moreover, the objective function is  $\frac{1}{K}$ -strongly convex, which guarantees some regularity of the arg min, namely local-Lipschitzness, leading to Lemma 1 below.

**Lemma 1.** For any  $(\lambda, \mu)$  with positive margin  $\Delta$  (defined in Equation (1)), if  $\|(\hat{\lambda} - \lambda, \hat{\mu} - \mu)\|_\infty \leq c_1 \Delta$ , for any  $c_1 < \frac{1}{2\sqrt{e+2}}$ , then

$$\|\phi(\hat{\lambda}, \hat{\mu}) - \phi(\lambda, \mu)\|_2 \leq \frac{c_2 K}{\Delta} \|(\hat{\lambda} - \lambda, \hat{\mu} - \mu)\|_\infty,$$

where  $c_2 = \frac{4}{(1-2c_1)/\sqrt{e}-2c_1}$ . Moreover, denoting  $\hat{P} = \phi(\hat{\lambda}, \hat{\mu})$ , it holds for any  $i \in [N]$ ,

$$\lambda_i \leq \sum_{k=1}^K \hat{P}_{i,k} \mu_k - \left( \frac{1-2c_1}{\sqrt{e}} - 2c_1 \right) \Delta.$$

The first property guarantees that if the queues have close estimates, they also have close doubly stochastic matrices  $\hat{P}$ . Moreover, the second property guarantees that any queue should clear its packets with a margin of order  $\Delta$ , in absence of collisions.

**Remark 4.** An alternative dominant mapping without the regularizing term in Equation (3) can also be proposed. Yet, its local Lipschitz bound would also depend on the smallest difference between the  $\lambda_i$  or the  $\mu_i$ , which can be arbitrarily small. If two parameters  $\lambda_i$  or  $\mu_i$  are equal, this choice of dominant mapping might lead to unstable policies. Using a regularization term in Equation (3) thus avoids this problem, although a smaller dependency might be possible without this regularization term when the parameters  $\lambda_i$  and  $\mu_i$  are very distinct.

## 4.2 Choice of a Birkhoff von Neumann decomposition

Given a doubly stochastic matrix  $\hat{P}$ , Birkhoff algorithm returns a convex combination of permutation matrices  $P[j]$  such that  $\hat{P} = \sum_j z[j] P[j]$ . The classical version of Birkhoff algorithm is non-continuous in its inputs and it even holds for its extensions as the one proposed by Dufossé et al. [2018]. Yet it can be smartly modified as in ORDERED BIRKHOFF, described in Algorithm 3, to get a regular BvN decomposition defined as follows for any  $\omega \in (0, 1)$ :

$$\psi(P)(\omega) = P[j_\omega] \tag{4}$$

where  $P = \sum_j z[j] P[j]$  is the decomposition returned by ORDERED BIRKHOFF algorithm

$$\text{and } j_\omega \text{ verifies } \sum_{j \leq j_\omega} z[j] \leq \omega < \sum_{j \leq j_\omega + 1} z[j].$$

For a matrix  $P$  in the following, its support is defined as  $\text{supp}(P) = \{(i, j) \mid P_{i,j} \neq 0\}$ .

---

### Algorithm 3: ORDERED BIRKHOFF

---

**input :**  $\hat{P} \in \mathfrak{B}_K$  (doubly stochastic matrix),  $C \in \mathbb{R}^{K \times K}$  (cost matrix)

- 1  $j \leftarrow 1$
- 2 **while**  $\hat{P} \neq \mathbf{0}$  **do**
- 3      $C_{i,k} \leftarrow +\infty$  for all  $(i, k) \notin \text{supp}(\hat{P})$                      // remove edge  $(i, k)$  in induced graph
- 4      $P[j] \leftarrow \text{HUNGARIAN}(C)$                                      // matching with minimal cost w.r.t.  $C$
- 5      $z[j] \leftarrow \min_{(i,k) \in \text{supp}(P[j])} \hat{P}_{i,k}$
- 6      $\hat{P} \leftarrow \hat{P} - z[j] P[j]$      and      $j \leftarrow j + 1$
- 7 **end**
- 8 **return**  $(z[j], P[j])_j$

---

Obviously  $\mathbb{E}_{\omega \sim \mathcal{U}(0,1)}[\psi(P)(\omega)] = P$  and permutations avoid collisions between queues. The difference with the usual Birkhoff algorithm happens at Line 4. Birkhoff algorithm usually computes any perfect matching in the graph induced by the support of  $\hat{P}$  at the current iteration. This is often done with the Hopcroft-Karp algorithm, while it is here done with the Hungarian algorithm with respect to some cost matrix  $C$ . Although using the Hungarian algorithm slightly increases the computational complexity of this step ( $K^3$  instead of  $K^{2.5}$ ), it ensures to output the permutation matrices  $P[j]$  according to a fixed order defined below.

**Definition 1.** A cost matrix  $C$  induces an order  $\prec_C$  on the permutation matrices defined, for any  $P, P' \in \mathfrak{S}_K$  by

$$P \prec_C P' \quad \text{iff} \quad \sum_{i,j} C_{i,j} P_{i,j} < \sum_{i,j} C_{i,j} P'_{i,j}.$$

This order might be non-total as different permutations can have the same cost. However, if  $C$  is drawn at random according to some continuous distribution, this order is total with probability 1. The order  $\prec_C$  has to be the same for all queues and is thus determined beforehand for all queues.

**Lemma 2.** Given matrices  $C \in \mathbb{R}^{K \times K}$  and  $P \in \mathfrak{B}_K$ , ORDERED BIRKHOFF outputs a sequence  $(z[j], P[j])_j$  of length at most  $K^2 - K + 1$ , such that

$$P = \sum_j z[j] P[j], \text{ where for all } j, z[j] > 0 \text{ and } P[j] \in \mathfrak{S}_K.$$

Moreover if the induced order  $\prec_C$  is total,  $z[j]$  is the  $j$ -th non-zero element of the sequence  $(z_l(P))_{1 \leq l \leq K!}$  defined by

$$z_j(P) = \min_{(i,k) \in \text{supp}(P_j)} \left( P - \sum_{l=1}^{j-1} z_l(P) P_l \right)_{i,k} \quad (5)$$

where  $(P_j)_{1 \leq j \leq K!}$  is a  $\prec_C$ -increasing sequence of permutation matrices, i.e.,  $P_j \prec_C P_{j+1}$  for all  $j$ .

Lemma 2 is crucial to guarantee the regularity of  $\psi$ , given by Lemma 3.

**Lemma 3.** Consider  $\psi$  defined as in Equation (4) with a cost matrix  $C$  inducing a total order  $\prec_C$ , then for any doubly stochastic matrices  $P, P'$

$$\int_0^1 \mathbb{1}(\psi(P)(\omega) \neq \psi(P')(\omega)) d\omega \leq 2^{2(K^2 - K + 1)} \|P - P'\|_\infty.$$

Lemma 3 indeed ensures that the probability of collision between two queues remains small when they have close estimates. Unfortunately, the regularity constant is exponential in  $K^2$ , which yields a similar dependency in the stability bound of Theorem 5. The existence of a BvN decomposition with polynomial regularity constants remains unknown, even without computational considerations. The design of a better BvN decomposition is left open for future work and would directly improve the stability bounds, using the general result given by Theorem 6 in Appendix A. The number of accumulated packets yet remain reasonably small in the experiments of Section 5, suggesting that the bound given by Lemma 3 is not tight and might be improved in future work.

### 4.3 Stability guarantees

This section finally provides theoretical guarantees on the stability of the system when all queues follow ADEQUA. The success of ADEQUA relies on the accurate estimation of all problem parameters by the queues, given by Lemma 9 in Appendix D.4. After some time  $\tau$ , the queues have tight estimations of the problem parameters. Afterwards, they clear their packets with a margin of order  $\Delta$ , thanks to Lemmas 1 and 3. This finally ensures the stability of the system, as given by Theorem 5.

**Theorem 5.** For any  $\eta > 1$ , the system where all queues follow ADEQUA, for any queue  $i$  and any  $r \in \mathbb{N}$ , there exists a constant  $C_r$  depending only on  $r$  such that

$$\mathbb{E}[(Q_t^i)^r] \leq C_r K N \left( \frac{N^{\frac{5}{2}} K^{\frac{5}{2}} 2^{5(K^2 - K + 1)}}{(\min(1, K\bar{\mu})\lambda)^{\frac{5}{4}} \Delta^5} \right)^r, \quad \text{for all } t \in \mathbb{N}.$$

As a consequence, for any  $\eta > 1$ , this decentralized system is strongly stable.

Despite yielding an exponential dependency in  $K^2$ , this anytime bound leads to a first decentralized stability result when  $\eta \in (1, \frac{e}{e-1})$ , which closes the stability gap left by previous works. Moreover it can be seen in the proof that the asymptotic number of packets is much smaller. It actually converges, in expectation, to the number of packets the queues would accumulate if they were following a stable centralized strategy from the beginning. As already noted by Krishnasamy et al. [2016] for a single queue, the number of packets first increases during the learning phase and then decreases once the queues have tight enough estimations, until reaching the same state as in the perfect knowledge centralized case. This is empirically confirmed in Section 5.

## 5 Simulations

Figures 2 and 3 compare on toy examples the stability of queues, when either each of them follows the no-regret strategy EXP3.P.1, or each queue follows ADeQUA. For practical considerations, we choose the exploration probability  $\varepsilon_t = (N + K)t^{-\frac{1}{4}}$  for ADeQUA, as the exploration is too slow with  $\varepsilon_t$  of order  $t^{-\frac{1}{5}}$ .

These figures illustrate the evolution of the average queue length on two different instances with  $N = K = 4$ . The code for the experiments is available at [gitlab.com/f\\_sen/queuing\\_systems](https://gitlab.com/f_sen/queuing_systems).

In the first instance shown in Figure 2, for all  $i \in [N]$ ,  $\lambda_i = (N + 1)/N^2$ . Moreover  $\mu_1 = 1$  and for all  $i \geq 2$ ,  $\mu_i = (N - 1)/N^2$ . Here  $\eta < 2$  and no-regret strategies are known to be unstable [Gaitonde and Tardos, 2020a]. It is empirically confirmed as the number of packets in each queue diverges when they follow EXP3.P.1. Conversely, when the queues follow ADeQUA, after a learning phase, the queues reach equilibrium and all succeed in clearing their packets.

In the second instance shown in Figure 3, for all  $i \in [N]$ ,  $\lambda_i = 0.55 - 0.1 \cdot i$  and  $\mu_i = 2.1\lambda_i$ . Here  $\eta > 2$  and both strategies are known to be stable, which is again empirically confirmed. However, ADeQUA requires more time to learn the different parameters, suggesting that individual no-regret strategies might be better on easy instances where  $\eta > 2$ .

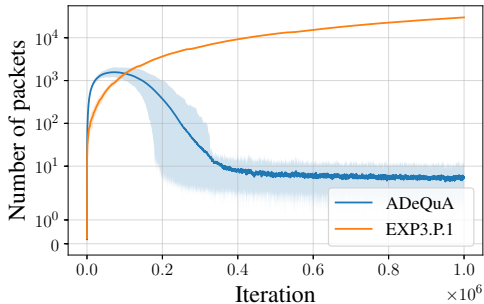


Figure 2: Hard instance,  $\eta < 2$ .

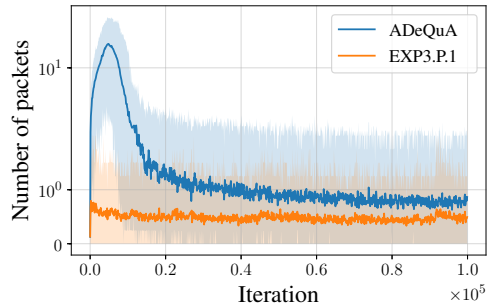


Figure 3: Easy instance,  $\eta > 2$ .

## 6 Conclusion

In this work, we showed that minimizing a more patient version of regret was not necessarily stable when the system’s slack is smaller than two and we argued that some level of cooperation was then required between learning queues to reach stability. We presented the first decentralized learning algorithm guaranteeing stability of any queuing system with a slack larger than 1. Our stability bound presents an exponential dependency in the number of queues and remains open for improvement, e.g., through a better dominant mapping/BvN decomposition or a tighter analysis of ours. The proposed algorithm relies heavily on synchronisation between the queues, which all start the game simultaneously and share a common time discretisation. In particular, the shared randomness assumption merely results from this synchronisation when the players use a common random seed. Stability of asynchronous queues thus remains open for future work, for which Glauber dynamics approaches used in scheduling problems might be of interest [see e.g., Shah and Shin, 2012].

## Acknowledgements

F. Sentenac was supported by IP PARIS' PhD Funding. E. Boursier was supported by an AMX scholarship. V. Perchet acknowledges support from the French National Research Agency (ANR) under grant number #ANR-19-CE23-0026 as well as the support grant, as well as from the grant "Investissements d'Avenir" (LabEx Ecodec/ANR-11-LABX-0047).

## References

- Animashree Anandkumar, Nithin Michael, and Ao Tang. Opportunistic spectrum access with multiple users: Learning under competition. In *2010 Proceedings IEEE INFOCOM*, pages 1–9. IEEE, 2010.
- Raman Arora, Ofer Dekel, and Ambuj Tewari. Online bandit learning against an adaptive adversary: from regret to policy regret. In *Proceedings of the 29th International Conference on International Conference on Machine Learning*, pages 1747–1754, 2012.
- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multi-armed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- Soumya Basu, Karthik Abinav Sankararaman, and Abishek Sankararaman. Beyond  $\log^2(t)$  regret for decentralized bandits in matching markets. *arXiv preprint arXiv:2103.07501*, 2021.
- Johann Berger. Schach-jahrbuch fur 1899/1900 : fortsetzung des schach-jahrbuches fur 1892/93. *Verlag von Veit*, 1899.
- Garrett Birkhoff. Tres observaciones sobre el algebra lineal. *Univ. Nac. Tucuman, Ser. A*, 5:147–154, 1946.
- Ilai Bistriz and Amir Leshem. Distributed multi-player bandits-a game of thrones approach. *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- Avrim Blum and Yishay Mansour. Learning, regret minimization, and equilibria. *Algorithmic Game Theory*, 2007.
- Rémi Bonnefoi, Lilian Besson, Christophe Moy, Emilie Kaufmann, and Jacques Palicot. Multi-armed bandit learning in iot networks: Learning helps even in non-stationary settings. In *International Conference on Cognitive Radio Oriented Wireless Networks*, pages 173–185. Springer, 2017.
- Allan Borodin, Jon Kleinberg, Prabhakar Raghavan, Madhu Sudan, and David P Williamson. Adversarial queueing theory. In *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*, pages 376–385, 1996.
- Etienne Boursier and Vianney Perchet. SIC-MMAB: Synchronisation involves communication in multiplayer multi-armed bandits. In *NIPS Proceedings*, 2019.
- Etienne Boursier and Vianney Perchet. Selfish robustness and equilibria in multi-player bandits. In *Conference on Learning Theory*, pages 530–581. PMLR, 2020.
- Sébastien Bubeck. Convex optimization: Algorithms and complexity. *arXiv preprint arXiv:1405.4980*, 2014.
- Sébastien Bubeck, Thomas Budzinski, and Mark Sellke. Cooperative and stochastic multi-player multi-armed bandit: Optimal regret with neither communication nor collisions. *arXiv preprint arXiv:2011.03896*, 2020.
- Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.
- Fanny Dufossé, Kamer Kaya, Ioannis Panagiotas, and Bora Uçar. Further notes on birkhoff–von neumann decomposition of doubly stochastic matrices. *Linear Algebra and its Applications*, 554: 68–78, 2018.

- Drew Fudenberg, Fudenberg Drew, David K Levine, and David K Levine. *The theory of learning in games*, volume 2. MIT press, 1998.
- Jason Gaitonde and Eva Tardos. Stability and learning in strategic queuing systems. In *Proceedings of the 21st ACM Conference on Economics and Computation*, pages 319–347, 2020a.
- Jason Gaitonde and Eva Tardos. Virtues of patience in strategic queuing systems. *arXiv preprint arXiv:2011.10205*, 2020b.
- Leonidas Georgiadis, Michael J Neely, and Leandros Tassiulas. *Resource allocation and cross-layer control in wireless networks*. Now Publishers Inc, 2006.
- Sergiu Hart and Andreu Mas-Colell. A simple adaptive procedure leading to correlated equilibrium. *Econometrica*, 68(5):1127–1150, 2000.
- Libin Jiang and Jean Walrand. A distributed csma algorithm for throughput and utility maximization in wireless networks. *IEEE/ACM Transactions on Networking*, 18(3):960–972, 2009.
- Elias Koutsoupias and Christos Papadimitriou. Worst-case equilibria. In *Annual Symposium on Theoretical Aspects of Computer Science*, pages 404–413. Springer, 1999.
- Subhashini Krishnasamy, Rajat Sen, Ramesh Johari, and Sanjay Shakkottai. Regret of queueing bandits. *Advances in Neural Information Processing Systems*, 29:1669–1677, 2016.
- Lydia T Liu, Horia Mania, and Michael Jordan. Competing bandits in matching markets. In *International Conference on Artificial Intelligence and Statistics*, pages 1618–1628. PMLR, 2020a.
- Lydia T Liu, Feng Ruan, Horia Mania, and Michael I Jordan. Bandit learning in decentralized matching markets. *arXiv preprint arXiv:2012.07348*, 2020b.
- Albert W Marshall, Ingram Olkin, and Barry C Arnold. *Inequalities: theory of majorization and its applications*, volume 143. Springer, 1979.
- Abbas Mehrabian, Etienne Boursier, Emilie Kaufmann, and Vianney Perchet. A practical algorithm for multiplayer bandits when arm means vary among players. In *International Conference on Artificial Intelligence and Statistics*, pages 1211–1221. PMLR, 2020.
- Michael J Neely, Eytan Modiano, and Chih-Ping Li. Fairness and optimal stochastic control for heterogeneous networks. *IEEE/ACM Transactions On Networking*, 16(2):396–409, 2008.
- Robin Pemantle and Jeffrey S Rosenthal. Moment conditions for a sequence with negative drift to be uniformly bounded in lr. *Stochastic Processes and their Applications*, 82(1):143–155, 1999.
- Vianney Perchet. Approachability, regret and calibration: Implications and equivalences. *Journal of Dynamics & Games*, 1(2):181, 2014.
- Tim Roughgarden. Algorithmic game theory. *Communications of the ACM*, 53(7):78–86, 2010.
- Abishek Sankararaman, Soumya Basu, and Karthik Abinav Sankararaman. Dominate or delete: Decentralized competing bandits with uniform valuation. *arXiv preprint arXiv:2006.15166*, 2020.
- Andreas S Schulz and Nicolás E Stier Moses. On the performance of user equilibria in traffic networks. In *SODA*, volume 3, pages 86–87, 2003.
- Devavrat Shah and Jinwoo Shin. Randomized scheduling algorithm for queueing networks. *The Annals of Applied Probability*, 22(1):128–171, 2012.
- John F Shortle, James M Thompson, Donald Gross, and Carl M Harris. *Fundamentals of queueing theory*, volume 399. John Wiley & Sons, 2018.
- Leandros Tassiulas and Anthony Ephremides. Stability properties of constrained queueing systems and scheduling policies for maximum throughput in multihop radio networks. In *29th IEEE Conference on Decision and Control*, pages 2130–2132. IEEE, 1990.
- Po-An Wang, Alexandre Proutiere, Kaito Ariu, Yassir Jedra, and Alessio Russo. Optimal algorithms for multiplayer multi-armed bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 4120–4129. PMLR, 2020.

## A General version of Theorem 5

ADEQUA is described for specific choices of the functions  $\phi$  and  $\psi$  given by Sections 4.1 and 4.2. It yet uses them in a black box fashion and different functions can be used, as long as they verify some key properties. This section provides a general version of Theorem 5, when the used dominant mapping and BvN decomposition respect the properties given by Assumptions 1 and 2.

**Assumption 1** (regular dominant mapping). *There are constants  $c_1, c_2 > 0$  and a norm  $\|\cdot\|$  on  $\mathbb{R}^{K \times K}$ , such that if  $\|(\hat{\lambda} - \lambda, \hat{\mu} - \mu)\|_\infty \leq c_1 \Delta$ , then*

$$\|\phi(\hat{\lambda}, \hat{\mu}) - \phi(\lambda, \mu)\| \leq L_\phi \cdot \|(\hat{\lambda} - \lambda, \hat{\mu} - \mu)\|_\infty.$$

Moreover,  $\hat{P} = \phi(\hat{\lambda}, \hat{\mu})$  is doubly stochastic and for any  $i \in [N]$ ,

$$\lambda_i \leq \sum_{k=1}^K \hat{P}_{i,k} \mu_k - c_2 \Delta.$$

**Assumption 2** (regular BvN decomposition). *Consider the same norm  $\|\cdot\|$  as Assumption 1 on  $\mathbb{R}^{K \times K}$ . For any doubly stochastic matrices  $P, P'$*

$$\int_0^1 \psi(P)(\omega) d\omega = P$$

$$\text{and } \int_0^1 \mathbb{1}(\psi(P)(\omega) \neq \psi(P')(\omega)) d\omega \leq L_\psi \cdot \|P - P'\|.$$

Lemmas 1 and 3 show that the functions described in Sections 4.1 and 4.2 verify Assumptions 1 and 2 with the constants  $L_\phi$  and  $L_\psi$  respectively of order  $\frac{K}{\Delta}$  and  $2^{2K^2}$  with the norm  $\|\cdot\|_\infty$ . Designing a dominant mapping and a BvN decomposition with smaller constants  $L_\phi$  and  $L_\psi$  is left open for future work. It would lead to a direct improvement of the stability bound, as shown by Theorem 6.

**Theorem 6.** *Assume all queues follow ADEQUA, using an exploration probability  $\varepsilon_t = xt^{-\alpha}$  with  $x > 0, \alpha \in (0, 1)$  and functions  $\phi$  and  $\psi$  verifying Assumptions 1 and 2 with the constants  $L_\phi, L_\psi$ . The system is then strongly stable and for any  $r \in \mathbb{N}$ , there exists a constant  $C_r$  such that:*

$$\mathbb{E}[(Q_t^i)^r] \leq C_r \left( \frac{x^{r/\alpha}}{\Delta^{r/\alpha}} + KN \left( \frac{N^2 K L_\phi^2 L_\psi^2}{\min(1, K \bar{\mu}) \lambda \Delta^2 x} \right)^{\frac{r}{1-\alpha}} \right), \quad \text{for all } t \in \mathbb{N}$$

The proof directly follows the lines of the proof of Theorem 5 in Appendix D.4 and is thus omitted here. From this version, it can be directly deduced that  $\alpha = \frac{1}{5}$  gives the best dependency in  $\Delta$  for ADEQUA. Moreover the best choice for  $x$  varies with  $r$ . When  $r \rightarrow \infty$ , it actually is  $x = N^{\frac{2}{5}} K^{\frac{3}{5}} 2^{\frac{4}{5}} K^2$  for ADEQUA. The choice  $x = N + K$  is preferred for simplicity and still yields quite similar problem dependent bounds.

## B Efficient computation of $\phi$

As mentioned in Section 4.1, computing exactly  $\phi(\hat{\lambda}, \hat{\mu})$  is not possible. Even efficiently approximating it is not obvious, as the function to minimize is neither smooth nor Lipschitz. We here describe how an approximation of  $\phi$  can be efficiently computed with guarantees on the approximation error.

First define the empirical estimate of the margin  $\Delta$ :

$$\hat{\Delta} := \min_{k \in [N]} \frac{1}{k} \left( \sum_{i=1}^k \hat{\mu}_{(i)} - \hat{\lambda}_{(i)} \right).$$

It can be computed in time  $\mathcal{O}(N \log(N))$  as it only requires to sort the vectors  $\hat{\lambda}$  and  $\hat{\mu}$ . If  $\hat{\Delta} \leq 0$ , then the value of the optimization problem is  $+\infty$  and any matrix can be returned. Assume in the following  $\hat{\Delta} > 0$ . Similarly to the proof of Lemma 1, it can be shown that the value of the optimization problem is smaller than  $-\ln(\hat{\Delta}/\sqrt{e})$ . Noting by  $\mathfrak{B}_K$  the set of  $K \times K$  doubly stochastic matrices, the optimization problem given by Equation (3) is then equivalent to

$$\arg \min_{P \in \mathfrak{B}_K} g(P), \tag{6}$$

where

$$\mathcal{X} = \left\{ P \in \mathfrak{B}_K \mid \forall i \in [N], \sum_{j=1}^K P_{i,j} \mu_j - \lambda_i \geq \frac{\hat{\Delta}}{\sqrt{e}} \right\},$$

and  $g(P) = \max_{i \in [N]} -\ln(\sum_{j=1}^K P_{i,j} \mu_j - \lambda_i) + \frac{1}{2K} \|P\|_2^2$ .

Thanks to this new constraint set, the objective function of Equation (6) is now  $(\frac{\sqrt{e}}{\Delta} + 1)$ -Lipschitz. We can now use classical results for Lipschitz strongly convex minimization to obtain convergence rates of order  $\frac{1}{t}$  for the projected gradient descent algorithm [see e.g., Bubeck, 2014, Theorem 3.9]. These results yet assume that the projection on the constraint set can be exactly computed in a short time. This is not the case here, but it yet can be efficiently approximated using interior point methods [see e.g., Bubeck, 2014, Section 5.3], which has a linear convergence rate. If this approximation is good enough, similar convergence guarantees than with exact projection can be shown similarly to the original proof.

Algorithm 4 then describes how to quickly estimate  $\phi(\hat{\lambda}, \hat{\mu})$ , where  $\hat{\Pi}_{\mathcal{X}}$  returns an approximation of the orthogonal projection on the set  $\mathcal{X}$  and  $\partial g$  is a sub-gradient of  $g$ . It uses an averaged value of the different iterates, as the last iterate does not have good convergence guarantees.

---

**Algorithm 4:** Compute  $\phi$

---

**input:** function  $g$ , constraint set  $\mathcal{X}$ ,  $P^0 \in \mathcal{X}$

```

1  $P, \hat{P} \leftarrow P^0$ 
2 for  $t = 1, \dots, n$  do
3    $P \leftarrow \hat{\Pi}_{\mathcal{X}} \left( P - \frac{2N}{(t+1)} \partial g(P) \right)$  // approximated projection
4    $\hat{P} \leftarrow \frac{t}{t+2} \hat{P} + \frac{2}{t+2} P$ 
5 end
6 return  $\hat{P}$ 

```

---

In practice, the approximation can even be computed faster by initializing  $P^0$  in Algorithm 4 with the solution of the previous round  $t - 1$ .

## C Unstable No-Policy regret system example

---

**Algorithm 5:** Unstable No-policy regret system example

---

**input:**  $w_k, N, \alpha, \lambda = (1/N, \dots, 1/N), \mu = (2(N-d)/N^2, \dots, 2(N-d)/N^2)$

```

1 for  $k = 1, \dots, \infty$  do
2   for  $t = 1, \dots, \lceil \alpha w_k \rceil$  do
3     Queues  $2i$  and  $2i + 1$  play server  $2i + t \pmod{N}$  // stage 1
4   end
5   for  $t = \lceil \alpha w_k \rceil + 1, \dots, w_k$  do
6     Queue  $i$  plays server  $i + t \pmod{N}$  // stage 2
7   end
8 end

```

---

**Lemma 4.** Consider the system where the queues play according to the policy described in Algorithm 5 over successive windows of length  $w_k = k^2$ . If  $\alpha > 1 - \frac{d}{N-d}$ , the system is not stable.

*Proof.* Note that the system is equivalent to a system where each queue or pair of queue would always pick the same server. For simplicity, the analysis deals with that equivalent system. Also, wlog, we analyse the subsystem with the two first queues and the two first servers. Let  $\{B_t^i\}_{i \in [n], t \geq 1}$  be the independent random variables indicating the arrival of a packet on queue  $i$  at time  $t$ ,  $\{S_t^i\}_{i \in [n], t \geq 1}$  be the indicators that server  $j$  would clear a packet at iteration  $\ell$  if one were sent to it. For each queue  $i \in [N]$  and  $t \geq 0$ , we have by Chernoff bound

$$\Pr \left( \left| \sum_{t=1}^{\ell} B_t^i - \lambda_i \ell \right| \geq \sqrt{\ell \ln(\ell)} \right) \leq \frac{2}{\ell^2}.$$

The same holds for each queue, thus the probability that this event happens for queue 1 or queue 2 is at most,  $\frac{4}{\ell^2}$ . As it is summable in  $\ell$ , The Borel-Cantelli Lemma implies that, for large enough  $\ell$ , almost surely, for any  $i \in [2]$ :

$$\sum_{\ell=1}^{\ell} B_t^i = \lambda_i \ell \pm \tilde{\mathcal{O}}(\sqrt{\ell}). \quad (7)$$

Let  $W_k = \sum_{i=1}^k w_i$ . Note that  $W_k = \Theta(k^3) = \Theta(w_k^{3/2})$ . Again by Chernoff bound and Borel-Cantelli, for large enough  $k$ , almost surely, for any  $i \in \{1, 2\}$ :

$$\sum_{t=W_{k-1}}^{W_{k-1} + \lceil \alpha w_k \rceil} S_t^i = \mu_i \alpha w_k \pm \tilde{\mathcal{O}}(\sqrt{w_k}), \quad \sum_{t=W_{k-1} + \lceil \alpha w_k \rceil}^{W_k} S_t^i = \mu_i (1 - \alpha) w_k \pm \tilde{\mathcal{O}}(\sqrt{w_k}). \quad (8)$$

Thus, for any large enough  $k$ , the total number of packet in both queues at time  $W_k$  is almost surely lower bounded as:

$$Q_{W_k}^1 + Q_{W_k}^2 \geq \sum_{t=1}^{W_k} (B_t^1 + B_t^2) - \sum_{t=1}^{W_k} S_t^1 - \sum_{t=1}^k \left( \sum_{t=W_{i-1} + \lceil \alpha w_i \rceil}^{W_i} S_t^2 \right) \quad (9)$$

$$\geq \left[ \frac{2}{N} - \frac{2(N-d)}{N^2} - (1-\alpha) \frac{2(N-d)}{N^2} \right] W_k - \tilde{\mathcal{O}}(W_k^{2/3}) \quad (10)$$

$$\geq \frac{2[\alpha(N-d) - (N-2d)]}{N^2} W_k - \tilde{\mathcal{O}}(W_k^{2/3}) \quad (11)$$

which is a diverging function of  $W_k$ . Note that this result also holds for any two pair of queues  $(2i-1, 2i)$ , with  $i \in [N/2]$ . □

**Lemma 5.** Consider the same setting as in Lemma 4. For any  $i \in [N]$ , for any large enough  $k$ , queue  $i$  clears

$$\left( \frac{N-d}{N^2} + (1-\alpha) \frac{N-d}{N^2} + o(1) \right) w_k$$

packets almost surely over window  $w_k$ .

*Proof.* The proof starts by showing that for any large enough  $t$ , all the queues hold roughly the same number of packets. Then, as they receive roughly the same number of packets over a time window and we can compute the approximate total number of packets cleared, the results follows.

Let  $T_i^t$  be the age of the oldest packet in queue  $i$  at time  $t$ . By Chernoff bound,

$$\mathbb{P}(|T_i^t - NQ_i^t| \geq N\sqrt{t \ln(t)}) \leq \frac{2}{t^2}.$$

Thus, using the Borel-Cantelli lemma, for any queue  $i$ , almost surely, for any large enough  $k$  and any  $t \in [W_{k-1} + 1, W_k]$ ,

$$|T_i^t - NQ_i^t| \leq N\sqrt{t \ln(t)} = \tilde{\mathcal{O}}(w_k^{3/4}). \quad (12)$$

For any  $(i, j) \in [N]^2$ , define

$$\phi_t^+(i, j) := \left( Q_t^i - Q_t^j - 2N\sqrt{t \ln(t)} \right)_+ \quad \text{and} \quad \phi_t^-(i, j) := \left( Q_t^i - Q_t^j + 2N\sqrt{t \ln(t)} \right)_-$$

Let  $C_t^i$  be the indicator function that queue  $i$  clears a packet at iteration  $t$ . Note that for any large enough  $t$ ,  $\phi_t^+(i, j)$  is a supermartingale. Indeed,

$$\begin{aligned} \mathbb{E}[\phi_{t+1}^+(i, j) | \phi_{1:t}^+(i, j)] &\leq \phi_t^+(i, j) + \mathbb{E}[B_t^i - B_t^j | \phi_{1:t}^+(i, j)] - \mathbb{E}[C_t^i - C_t^j | \phi_{1:t}^+(i, j)] \\ &\leq \phi_t^+(i, j). \end{aligned}$$



The second inequality comes from Equation (12), that implies that for any large enough  $t$ , if  $\phi_t^+(i, j)$  is strictly positive, queue  $i$  holds the oldest packet and thus clears one with higher probability than queue  $j$ . By the same arguments,  $\phi_t^-(i, j)$  a submartingale. Also,  $|\phi_{t+1}^+(i, j) - \phi_t^+(i, j)| \leq 2(N+1)$  for any  $t \geq 0$ , and the same holds for  $\phi_t^-(i, j)$ . Let  $\tau_{ij}$  be the stopping time of the smallest iteration after which Equation (12) always holds for queues  $i$  and  $j$ . By Azuma-Hoeffding's inequality,

$$\Pr \left( \phi_\ell^+(i, j) - \phi_{\tau_{ij}}^+(i, j) \geq 3(N+1)\sqrt{\ell \ln(\ell)} \right) \leq \frac{2}{\ell^2}$$

and

$$\Pr \left( \phi_\ell^-(i, j) - \phi_{\tau_{ij}}^-(i, j) \leq -3(N+1)\sqrt{\ell \ln(\ell)} \right) \leq \frac{2}{\ell^2}.$$

This, together with a union bound and Borel-Cantelli's Lemma implies that almost surely, for any large enough  $t$ , for any  $(i, j) \in [N]^2$

$$Q_t^i - Q_t^j = \tilde{O}(\sqrt{t}). \quad (13)$$

This with Equation (9) implies that for any large enough  $k$ , for any  $i \in [N]$ , almost surely,

$$Q_{W_k}^i \geq \frac{[\alpha(N-d) - (N-2d)]}{N^2} W_k - \tilde{O}(W_k^{2/3}).$$

This means that for any large enough  $k$ , every queue holds at least one packet over the whole window  $w_k$ . This and Equation (8) is already enough to show that for any time-window  $w_k$ , for any large enough  $k$ , the total number of packets cleared by any couple of queue  $(2i-1, 2i)$ ,  $i \in [N/2]$  is:

$$2 \left( \frac{N-d}{N^2} + (1-\alpha) \frac{N-d}{N^2} \right) w_k + \tilde{O}(\sqrt{w_k}).$$

During time window  $w_k$ , according to Equation (7), both every queue receives  $\alpha w_k / N + \tilde{O}(w_k^{3/4})$  packets almost surely for any large enough  $k$ . Equation (13) implies that for any  $i \in [N/2]$

$$Q_{W_k}^{2i-1} - Q_{W_k}^{2i} = \tilde{O}(w_k^{3/4}) \text{ and } Q_{W_{k-1}}^{2i-1} - Q_{W_{k-1}}^{2i} = \tilde{O}(w_k^{3/4}).$$

Therefore, over each time-window  $w_k$ , for any large enough  $k$ , each queue clears

$$\left( \frac{N-d}{N^2} + (1-\alpha) \frac{N-d}{N^2} + o(1) \right) w_k$$

packets almost surely.  $\square$

**Lemma 6.** *Consider again the system where the queues play according to the policy described in Algorithm 5 over successive windows of length  $w_k = k^2$ . If  $\alpha < 1 - \frac{1}{N-1}$ , the queues have  $o(w_k)$  policy regret in all but finitely many of the windows.*

Wlog, let us consider that queue 1 deviates, and plays at every iteration a server chosen from the probability distribution  $\mathbf{p} = (p_1, \dots, p_N)$ , with  $p_i$  the probability to play server  $i$ . To upper bound the number of packets queue 1 clears over each time window, we can assume it always has priority over queue 2 and ignore it in the analysis.

Before proving Lemma 6, we prove the following technical one.

**Lemma 7.** *Consider that a queue deviates from the strategy considered in Lemma 6 and plays at every iteration a server chosen from the probability distribution  $\mathbf{p} = (p_1, \dots, p_N)$ , with  $p_i$  the probability to play server  $i$ . For any large enough  $k$ , almost surely, the number of packets the deviating queue clears of the first stage of the  $k^{\text{th}}$  window is*

$$\left( \frac{1}{2} + \frac{1}{N} \right) \frac{2(N-d)}{N^2} \alpha w_k + \tilde{O}(w_k^{3/4}).$$

*Proof.* The proof starts by showing that for any large enough  $t$ , every non-deviating queue holds approximately the same number of packets.

First note that for any large enough  $t$ , Equation (12) still holds surely for any queue  $i$ . For any  $(i, j) \in \{3, \dots, N\}^2$ , define

$$\phi_\ell^+(i, j) := \left( Q_{\lceil \ell N \rceil}^i - Q_{\lceil \ell N \rceil}^j - 4N\sqrt{\lceil \ell N \rceil \ln(\lceil \ell N \rceil)} \right)_+$$

and

$$\phi_\ell^-(i, j) := \left( Q_{\lceil \ell N \rceil}^i - Q_{\lceil \ell N \rceil}^j + 4N\sqrt{\lceil \ell N \rceil \ln(\lceil \ell N \rceil)} \right)_-$$

For any interval  $[\lceil \ell N \rceil, \lceil (\ell + 1)N \rceil]$  where Equation (12) holds for queues 1,  $i$  and  $j$ , if  $\phi_\ell^+(i, j)$  is strictly positive, then

$$\mathbb{E} \left[ \sum_{t=\lceil \ell N \rceil}^{\lceil (\ell+1)N \rceil} C_t^j - C_t^i \middle| \phi_{1:t}^+(i, j) \right] \leq 0.$$

Indeed, if  $\phi_\ell^+(i, j)$  is strictly positive and Equation (12) holds, queue  $i$  holds the oldest packets throughout the interval. Also, queue  $i$  and queue  $j$  collide with queue 1 the same number of times over the interval in expectation, and if at one iteration of the interval, queue 1 holds an older packet than queue  $i$ , it holds an older packet than queue  $j$  over the whole interval. Thus  $\phi_\ell^+(i, j)$  is a submartingale. By the same arguments,  $\phi_\ell^-(i, j)$  is a supermartingale. Also,  $|\phi_{\ell+1}^+(i, j) - \phi_\ell^+(i, j)| \leq 4(N+1)^2$  and the same holds for  $\phi_\ell^-(i, j)$ . Finishing with the same arguments used to prove Equation (13), almost surely, for any  $(i, j) \in \{3, \dots, N\}^2$ ,

$$Q_t^i - Q_t^j = \tilde{O}(\sqrt{t}). \quad (14)$$

We now show that for any large enough  $t$ , queue 1 can not hold many more packets than the non-deviating queues. Define

$$\phi_t^+ := \left( Q_t^1 - \max_{i \geq 3} Q_t^i - 2N\sqrt{t \ln(t)} \right)_+.$$

Once again, at every iteration where  $\phi_t^+$  is strictly positive and Equation (12) holds, queue 1 holds the oldest packet and thus has priority on whichever server it chooses. This implies that for any large enough  $t$ ,  $\phi_t^+$  is a supermartingale. It also holds that for any  $t \geq 0$ ,  $|\phi_{t+1}^+ - \phi_t^+| \leq 2(N+1)$ . Thus, with the same arguments used to prove Equation (13), almost surely,

$$\left( Q_t^1 - \max_{i \geq 3} Q_t^i \right)_+ = \tilde{O}(\sqrt{t}). \quad (15)$$

With that at hand, we prove that for any large enough  $k$ , queue 1 does not get priority often over the other queues during the first stage of the  $k^{\text{th}}$  window. For any  $i \in \{2, \dots, N/2\}$ , pose:

$$\psi_\ell^i = \frac{1}{2} \left( Q_{\lceil \ell N \rceil}^{2i-1} + Q_{\lceil \ell N \rceil}^{2i} \right) - Q_{\lceil \ell N \rceil}^1 - \frac{2(N-d)}{N^3} (\lceil \ell N \rceil - W_{k-1})$$

For any  $\ell$  s.t.  $\{\lceil \ell N \rceil; \lceil (\ell + 1)N - 1 \rceil\}$  is included in the first phase of a window, we have

$$\begin{aligned} \sum_{t=\lceil \ell N \rceil}^{\lceil (\ell+1)N \rceil - 1} \mathbb{E} \left[ C_t^1 \middle| \psi_{1:t}^+(i, j) \right] &\geq \sum_{t=\lceil \ell N \rceil}^{\lceil (\ell+1)N \rceil - 1} \mathbb{E} \left[ S_t^i \mathbf{1}_{\{\text{queue 1 and only queue 1 picks server } i\}} \middle| \psi_{1:t}^+(i, j) \right] \\ &\geq \frac{N-d}{N} + \frac{2(N-d)}{N^2} \end{aligned}$$

as well as

$$\begin{aligned} \sum_{t=\lceil \ell N \rceil}^{\lceil (\ell+1)N \rceil - 1} \mathbb{E} \left[ \frac{1}{2} (C_t^{2i} + C_t^{2i-1}) \middle| \psi_{1:t}^+(i, j) \right] &\leq \sum_{t=\lceil \ell N \rceil}^{\lceil (\ell+1)N \rceil - 1} \mathbb{E} \left[ \frac{1}{2} S_{i+t \pmod{N}}^t \middle| \psi_{1:t}^+(i, j) \right] \\ &\leq \frac{N-d}{N}. \end{aligned}$$

Those two inequalities imply:

$$\begin{aligned}\mathbb{E}[\psi_{\ell+1}^i | \psi_{1:\ell}^+(i, j)] &= \psi_{\ell}^+(i, j) + \sum_{t=\lceil \ell N \rceil}^{\lceil (\ell+1)N \rceil - 1} \mathbb{E} \left[ \frac{1}{2} (B_t^{2i} - B_t^{2i-1}) - B_t^1 \middle| \psi_{1:\ell}^+(i, j) \right] \\ &\quad - \sum_{t=\lceil \ell N \rceil}^{\lceil (\ell+1)N \rceil - 1} \mathbb{E} \left[ \frac{1}{2} (C_t^{2i} - C_t^{2i-1}) - C_t^1 \middle| \psi_{1:\ell}^+(i, j) \right] - \frac{2(N-d)}{N^2} \\ &\geq \psi_{\ell}^+(i, j).\end{aligned}$$

Thus, for any  $\ell$  s.t.  $\{\lceil \ell N \rceil; \lceil (\ell+1)N - 1 \rceil\}$  is included in the first phase of a window,  $\psi_{\ell}^i$  is a submartingale. Moreover, for any  $\ell \geq 0$ ,  $|\psi_{\ell+1}^i - \psi_{\ell}^i| \leq 3N$ . Thus, by Azuma-Hoeffding's inequality, for any  $\ell$  s.t.  $\{\lceil \ell N \rceil; \lceil (\ell+1)N - 1 \rceil\} \subset [W_{k-1}, W_{k-1} + \alpha w_k]$ ,

$$\Pr \left( \psi_{\ell}^i - \psi_{W_k}^i \leq -6N \sqrt{\ell N \ln(\ell N)} \right) \leq \frac{1}{(\ell N)^2}.$$

Borel-Cantelli's lemma implies, that for any large enough  $\ell$  s.t.  $\{\lceil \ell N \rceil; \lceil (\ell+1)N - 1 \rceil\} \subset [W_{k-1}, W_{k-1} + \alpha w_k]$ , almost surely:

$$\psi_{\ell}^i \geq \psi_{W_k}^i - 6N \sqrt{\ell N \ln(\ell N)}.$$

This and Equation (15) applied at  $t = W_k$ , imply that for any large enough  $k$ , for any  $t \in [W_{k-1}, W_{k-1} + \alpha w_k]$ ,

$$\begin{aligned}\frac{1}{2} (Q_t^{2i-1} + Q_t^{2i}) &\geq Q_{\lceil \ell N \rceil}^1 + \frac{2(N-d)}{N^3} (t - W_{k-1}) + \psi_{W_k}^i - \tilde{O}(\sqrt{t}) \\ &\geq Q_{\lceil \ell N \rceil}^1 + \frac{2(N-d)}{N^3} (t - W_{k-1}) - \tilde{O}(w_k^{3/4}).\end{aligned}$$

This and Equation (12) imply that during the first stage of the time window, queue 1 holds younger packets than any other queues  $i \geq 3$  after at most  $\tilde{O}(w_k^{3/4})$  iterations.

By Chernoff bound and the Borel-Cantelli lemma again, for any large enough  $k$ , almost surely, the number of packets queue 1 clears during the first stage of the  $k^{\text{th}}$  window on servers where it does not collide with other queues is:

$$\sum_{t=W_{k-1}+\alpha w_k}^{W_{k-1}+\alpha w_k} \sum_{i=1}^N S_i^t \mathbb{1}_{\{\text{queue 1 and only queue 1 picks server } i\}} = \left(\frac{1}{2} + \frac{1}{N}\right) \frac{2(N-d)}{N^2} \alpha w_k + \tilde{O}(\sqrt{w_k}).$$

Since we have shown that for any large enough  $k$ , almost surely, queue 1 does not have priority over the other queues after at most  $\tilde{O}(w_k^{3/4})$  iterations, for any large enough  $k$ , almost surely, the number of packets queue 1 clears of the first stage of the  $k^{\text{th}}$  window is

$$\left(\frac{1}{2} + \frac{1}{N}\right) \frac{2(N-d)}{N^2} \alpha w_k + \tilde{O}(w_k^{3/4}).$$

□

We are now ready to prove Lemma 6.

*Proof.* By Chernoff bound and the Borel-Cantelli lemma, almost surely for any large enough  $k$ , the number of packets queue 1 clears during the second stage of the window on servers where it does not collide with other queues is:

$$\sum_{t=W_{k-1}+\alpha w_k}^{W_k-1} \sum_{i=1}^N S_i^t \mathbb{1}_{\{\text{queue 1 and only queue 1 picks server } i\}} = \frac{4(N-d)}{N^3} (1-\alpha) w_k + \tilde{O}(\sqrt{w_k}). \quad (16)$$

Suppose that during the second stage of the window, queue 1 never gets priority over another queue. In that case, according to Equation (16) and Lemma 7, for any large enough  $k$ , almost surely, the total number of packets cleared by queue 1 during the time window is

$$\left(\frac{\alpha}{2} + \frac{2-\alpha}{N}\right) \frac{2(N-d)}{N^2} w_k + \tilde{O}(w_k^{3/4}).$$

For any large enough  $k$ , if  $\alpha \leq 1 - \frac{1}{N-1}$  this is smaller than the number of packets queue 1 would have cleared had it not deviated, according to Lemma 5.

On the other hand, suppose that queue gets priority over some other queue  $i$  at some iteration  $\tau$  of the second stage of the window. In that case, at that iteration, queue 1 holds the oldest packets, which, according to Equation (12), implies

$$Q_1^\tau > Q_i^\tau - \tilde{O}(w_k^{3/4})$$

During the second stage of the window, for any  $i \geq 3$ ,  $\gamma_t^i := \left(Q_i^t - Q_1^t - 2N\sqrt{t \ln(t)}\right)_+$  is a supermartingale with bounded increments for any  $t$  where Equation (12) holds for queues 1 and  $i$ . Indeed, in that case, if  $\gamma_t^i$  is strictly positive, queue  $i$  holds an older packet than queue 1, and thus, whether they collide or not, it has a higher probability to clear a packet than queue 1. Thus, by Azuma-Hoeffding and the Borel-cantelli lemma again, for any large enough  $k$ , almost surely,

$$Q_i^{W_k} - Q_1^{W_k} \leq Q_i^\tau - Q_1^\tau + \tilde{O}(w_k^{3/4}).$$

Thus it holds that  $Q_1^{W_k} \geq Q_i^{W_k} - \tilde{O}(w_k^{3/4})$  for any  $i \geq 2$ . This and Equation (15) imply that all the queues clear approximately the same number of packets over those time windows for any large enough  $k$  almost surely. Thus queue 1 clears

$$\left[(2-\alpha)(N-2) + \left(\alpha + \frac{4-2\alpha}{N}\right)\right] \frac{(N-d)}{(N-1)N^2} w_k + \tilde{O}(w_k^{3/4})$$

packets almost surely, which again is smaller than the number of packets it would have cleared had it not deviated.

Thus, the deviating queue clears almost surely less packets by time window than it would have had it not deviated on all but finitely many of the time windows, which implies that it has no policy regret on all but finitely many of the time windows.  $\square$

## D Proofs of Section 4

### D.1 Proof of Lemma 1

We want to show that if  $\|(\hat{\lambda} - \lambda, \hat{\mu} - \mu)\|_\infty \leq c_1 \Delta$ , then

$$\|\phi(\hat{\lambda}, \hat{\mu}) - \phi(\lambda, \mu)\|_2 \leq \frac{c_2 K}{\Delta} \|(\hat{\lambda} - \lambda, \hat{\mu} - \mu)\|_\infty, \quad (17)$$

with the constants  $c_1, c_2$  given in Lemma 1.

Recall that  $\phi$  is defined as

$$\phi(\lambda, \mu) = \arg \min_{P \in \mathfrak{B}_K} f(P, \lambda, \mu),$$

where  $\mathfrak{B}_K$  is the set of  $K \times K$  doubly stochastic matrices and  $f$  is defined as:

$$f(P, \lambda, \mu) := \max_{i \in [N]} -\ln\left(\sum_{j=1}^K P_{i,j} \mu_j - \lambda_i\right) + \frac{1}{2K} \|P\|_2^2$$

Let  $P^*$  and  $\hat{P}^*$  be the minimizers of  $f$  with the respective parameters  $(\lambda, \mu)$  and  $(\hat{\lambda}, \hat{\mu})$ . They are uniquely defined as  $f$  is  $\frac{1}{K}$  strongly convex.

As the property of Lemma 1 is symmetric, we can assume without loss of generality that  $f(P^*, \lambda, \mu) \geq f(\hat{P}^*, \hat{\lambda}, \hat{\mu})$ .

Given the definition of  $\Delta$ , we have the bound

$$-\ln(\Delta) + \frac{1}{2} \geq f(P^*, \lambda, \mu) \geq -\ln(\Delta).$$

The lower bound holds because the term in the  $\ln$  is at most  $\Delta$  for at least one  $i$ . For the upper bound, some matrix  $P$  ensures that the term in the  $\ln$  is at least  $\Delta$  for all  $i$  and  $\|P\|_2^2 \leq K$ .

Similarly for  $\hat{P}^*$ , it follows:

$$-\ln((1 - 2c_1)\Delta) + \frac{1}{2} \geq f(\hat{P}^*, \hat{\lambda}, \hat{\mu}) \geq -\ln((1 + 2c_1)\Delta).$$

As a consequence, it holds for any  $i \in [N]$ :

$$\begin{aligned} -\ln\left(\sum_{j=1}^K \hat{P}_{i,j}^* \hat{\mu}_j - \hat{\lambda}_i\right) &\leq f(\hat{P}^*, \hat{\lambda}, \hat{\mu}) \\ &\leq -\ln((1 - 2c_1)\Delta/\sqrt{e}) \\ \sum_{j=1}^K \hat{P}_{i,j}^* \hat{\mu}_j - \hat{\lambda}_i &\geq (1 - 2c_1)\Delta/\sqrt{e}. \end{aligned}$$

Note that for any  $i \in [N]$ ,

$$\sum_{j=1}^K \hat{P}_{i,j}^* \hat{\mu}_j - \hat{\lambda}_i \leq \sum_{j=1}^K \hat{P}_{i,j}^* \mu_j - \lambda_i + 2\|(\hat{\lambda} - \lambda, \hat{\mu} - \mu)\|_\infty.$$

It then yields the second point of Lemma 1:

$$\sum_{j=1}^K \hat{P}_{i,j}^* \mu_j - \lambda_i \geq ((1 - 2c_1)/\sqrt{e} - 2c_1) \Delta$$

Moreover, it follows

$$\begin{aligned} -\ln\left(\sum_{j=1}^K \hat{P}_{i,j}^* \hat{\mu}_j - \hat{\lambda}_i\right) &\geq -\ln\left(\sum_{j=1}^K \hat{P}_{i,j}^* \mu_j - \lambda_i\right) - \ln\left(1 + \frac{2\|(\hat{\lambda} - \lambda, \hat{\mu} - \mu)\|_\infty}{\sum_{j=1}^K \hat{P}_{i,j}^* \mu_j - \lambda_i}\right) \\ &\geq -\ln\left(\sum_{j=1}^K \hat{P}_{i,j}^* \mu_j - \lambda_i\right) - \frac{2\|(\hat{\lambda} - \lambda, \hat{\mu} - \mu)\|_\infty}{((1 - 2c_1)/\sqrt{e} - 2c_1) \Delta} \end{aligned}$$

Recall that for a  $\frac{1}{K}$ -strongly convex function  $g$  of global minimum  $x^*$  and any  $x$ :

$$\|x - x^*\|_2 \leq 2K(g(x) - g(x^*))$$

As a consequence, it follows:

$$\begin{aligned} f(\hat{P}^*, \hat{\lambda}, \hat{\mu}) &\geq f(\hat{P}^*, \lambda, \mu) - \frac{2\|(\hat{\lambda} - \lambda, \hat{\mu} - \mu)\|_\infty}{((1 - 2c_1)/\sqrt{e} - 2c_1) \Delta} \\ &\geq f(P^*, \lambda, \mu) - \frac{2\|(\hat{\lambda} - \lambda, \hat{\mu} - \mu)\|_\infty}{((1 - 2c_1)/\sqrt{e} - 2c_1) \Delta} + \frac{1}{2K} \|P^* - \hat{P}^*\|_2. \end{aligned}$$

Equation (17) then follows.

## D.2 Proof of Lemma 2

The coefficient  $C_{i,j}$  is replaced by  $+\infty$  as soon as the whole weight  $P_{i,j}$  is exhausted. Thanks to this, the HUNGARIAN algorithm does return a perfect matching with respect to the bipartite graph with edges  $(i, j)$  where there remains some weight for  $P_{i,j}$ . Because of this, it can be shown following the usual proof of Birkhoff algorithm [Birkhoff, 1946] that the sequence  $(z[j], P[j])$  is indeed of length at most  $K^2 - K + 1$  and is a valid decomposition of  $P$ .

Now assume that  $\prec_C$  is a total order. At each iteration  $j$  of HUNGARIAN algorithm, denote  $\tilde{P}^j = P - \sum_{s=1}^{j-1} z[s]P[s]$  the remaining weights to attribute.

Let  $l_j$  be such that  $P[j] = P_{l_j}$  for any iteration  $j$  of HUNGARIAN algorithm.

It can now be shown by induction that

$$\tilde{P}^j = P - \sum_{l=1}^{l_j} z_l(P)P_l.$$

where  $z_l(P)$  are defined by Equation (5). Indeed, by definition

$$\begin{aligned} \tilde{P}^{j+1} &= \tilde{P}^j - z[j+1]P[j+1] \\ &= \tilde{P}^j - z[j+1]P_{l_{j+1}} \end{aligned}$$

The HUNGARIAN algorithm returns the minimal cost matching with respect to the modified cost matrix  $C$  where the coefficients  $i, k$  such that  $\tilde{P}_{i,k}^j = 0$  are replaced by  $+\infty$ . Thanks to this,  $P_{l_{j+1}}$  is the minimal cost permutation matrix  $P_l$  (for  $\prec_C$ ) such that  $\tilde{P}_{i,k}^j > 0$  for all  $(i, k) \in \text{supp}(P_l)$ .

This means that for any  $l < l_{j+1}$

$$\min_{(i,k) \in \text{supp}(P_l)} (\tilde{P}^j)_{i,k} = 0.$$

Using the induction hypothesis, this implies that  $z_l(P) = 0$  for any  $l_j < l < l_{j+1}$ . And finally, this also implies that  $z_{l_{j+1}}(P) = z[j+1]$ .

This finally concludes the proof as  $\tilde{P}^j = 0$  after the last iteration.

## D.3 Proof of Lemma 3

For  $z$  and  $z'$  the respective decompositions of  $P$  and  $P'$  defined in Lemma 2, then

$$\int_0^1 \mathbb{1}(\psi(P)(\omega) \neq \psi(P')(\omega)) d\omega = \mathbb{P}_{\omega \sim U(0,1)}(\psi(P)(\omega) \neq \psi(P')(\omega)).$$

In the following, note  $A = \psi(P)$  and  $A' = \psi(P')$ . It follows for  $P_n$  defined as in Lemma 2

$$\begin{aligned} \int_0^1 \mathbb{1}(\psi(P)(\omega) \neq \psi(P')(\omega)) d\omega &= \sum_{n=1}^{K!} \mathbb{P}(A = P_n \text{ and } A' \neq P_n) \\ &= \frac{1}{2} \sum_{n=1}^{K!} \mathbb{P}(A = P_n \text{ and } A' \neq P_n) + \frac{1}{2} \sum_{n=1}^{K!} \mathbb{P}(A' = P_n \text{ and } A \neq P_n) \\ &= \frac{1}{2} \sum_{n=1}^{K!} \text{vol} \left( \left[ \sum_{j=1}^{n-1} z_j(P), \sum_{j=1}^n z_j(P) \right] \ominus \left[ \sum_{j=1}^{n-1} z_j(P'), \sum_{j=1}^n z_j(P') \right] \right), \end{aligned}$$

where  $\text{vol}$  denotes the volume of a set and  $A \ominus B = (A \setminus B) \cup (B \setminus A)$  is the symmetric difference of  $A$  and  $B$ . The last equality comes from the expression of  $\psi$  with respect to the coefficients  $z_j(P)$ , thanks to Lemma 2.

It is easy to show that

$$\text{vol}([a, b] \ominus [c, d]) \leq (|c - a| + |d - b|) \mathbf{1}(b > a \text{ or } c > d).$$

The previous equality then leads to

$$\begin{aligned} \int_0^1 \mathbf{1}(\psi(P)(\omega) \neq \psi(P')(\omega)) \, d\omega &\leq \frac{1}{2} \sum_{n=1}^{K!} \left( \left| \sum_{j=1}^{n-1} z_j(P) - z_j(P') \right| + \left| \sum_{j=1}^n z_j(P) - z_j(P') \right| \right) \mathbf{1}(z_n(P) + z_n(P') > 0) \\ &\leq \sum_{n=1}^{K!} \left| \sum_{j=1}^n z_j(P) - z_j(P') \right| \mathbf{1}(z_n(P) + z_n(P') > 0). \end{aligned} \quad (18)$$

The last inequality holds because  $\sum_{j=1}^k z_j(P) - z_j(P')$  is counted twice when  $z_k(P) + z_k(P')$  is positive: when  $n = k$  and for the next  $n$  such that the elements are counted in the sum.

Thanks to Lemma 2, only  $2(K^2 - K + 1)$  elements  $z_j(P)$  and  $z_j(P')$  are non-zero. Let  $k_n$  be the index of the  $n$ -th non-zero element of  $(z_s(P) + z_s(P'))_{1 \leq s \leq K!}$ . Note that  $z_s(P')$  can be non-zero while  $z_s(P)$  is zero (or conversely). Let also

$$\begin{aligned} (i_{k_n}, j_{k_n}) &\in \arg \min_{(i,j) \in \text{supp}(P_{k_n})} P_{i,j} - \sum_{l < k_n} z_l(P) \mathbf{1}((i,j) \in \text{supp}(P_{k_n})), \\ (i'_{k_n}, j'_{k_n}) &\in \arg \min_{(i,j) \in \text{supp}(P_{k_n})} P'_{i,j} - \sum_{l < k_n} z_l(P') \mathbf{1}((i,j) \in \text{supp}(P_{k_n})). \end{aligned}$$

It then comes, thanks to Lemma 2

$$\begin{aligned} z_{k_n}(P) - z_{k_n}(P') &\leq P_{i'_{k_n}, j'_{k_n}} - P'_{i'_{k_n}, j'_{k_n}} - \sum_{l < k_n} (z_l(P) - z_l(P')) \mathbf{1}((i'_{k_n}, j'_{k_n}) \in \text{supp}(P_{k_n})) \\ &\leq P_{i'_{k_n}, j'_{k_n}} - P'_{i'_{k_n}, j'_{k_n}} - \sum_{l < n} (z_{k_l}(P) - z_{k_l}(P')) \mathbf{1}((i'_{k_n}, j'_{k_n}) \in \text{supp}(P_{k_n})) \end{aligned}$$

The second inequality holds, because for  $l' \notin \{k_l \mid l < 2K^2\}$ , the term in the sum is zero by definition of the sequence  $k_l$ .

A similar inequality holds for  $z_{k_n}(P') - z_{k_n}(P)$ , which leads to

$$|z_{k_n}(P) - z_{k_n}(P')| \leq \|P - P'\|_\infty + \sum_{l < n} |z_{k_l}(P) - z_{k_l}(P')|.$$

By induction, it thus holds

$$|z_{k_n}(P) - z_{k_n}(P')| \leq 2^{n-1} \|P - P'\|_\infty.$$

We finally conclude using Equation (18)

$$\begin{aligned} \int_0^1 \mathbf{1}(\psi(P)(\omega) \neq \psi(P')(\omega)) \, d\omega &\leq \sum_{n=1}^{K!} \left| \sum_{j=1}^n z_j(P) - z_j(P') \right| \mathbf{1}(z_n(P) + z_n(P') > 0) \\ &\leq \sum_{n=1}^{2(K^2 - K + 1) - 1} \left| \sum_{j=1}^{k_n} z_j(P) - z_j(P') \right| \\ &\leq \sum_{n=1}^{2(K^2 - K + 1) - 1} \left| \sum_{l=1}^n z_{k_l}(P) - z_{k_l}(P') \right| \\ &\leq \sum_{n=1}^{2(K^2 - K + 1) - 1} \sum_{j=1}^n 2^{j-1} \|P - P'\|_\infty \\ &\leq 2^{2(K^2 - K + 1)} \|P - P'\|_\infty. \end{aligned}$$

In the fourth inequality, the  $2(K^2 - K + 1)$ -th term of the sum is ignored. It is indeed 0 as  $z$  and  $z'$  both sum to 1.

#### D.4 Proof of Theorem 5

First recall below a useful version of Chernoff bound.

**Lemma 8.** For any independent variables  $X_1, \dots, X_n$  in  $[0, 1]$  and  $\delta \in (0, 1)$ ,

$$\mathbb{P}\left(\sum_{i=1}^n X_i \leq (1 - \delta) \sum_{i=1}^n \mathbb{E}[X_i]\right) \leq e^{-\frac{\delta^2 \sum_{i=1}^n \mathbb{E}[X_i]}{2}}.$$

We now prove the following concentration lemma.

**Lemma 9.** For any time  $t \geq (N + K)^5$  and  $\varepsilon \in (0, \frac{1}{4})$ ,

$$\begin{aligned} \mathbb{P}\left(|\hat{\mu}_k^i(t) - \mu_k| \geq \varepsilon\right) &\leq 3 \exp\left(-\lambda_i \left(t^{\frac{4}{5}} - 1\right) \varepsilon^2\right) \\ \mathbb{P}\left(|\hat{\lambda}_j^i(t) - \lambda_j| \geq \varepsilon\right) &\leq 6 \exp\left(-\lambda_i K \bar{\mu} \frac{t^{\frac{4}{5}} - 1}{145} \varepsilon^2\right). \end{aligned}$$

*Proof.*

**Concentration for  $\hat{\mu}$ .** Consider agent  $i$  in the following and denote by  $N_k(t)$  the number of exploratory pulls of this agent on server  $k$  at time  $t$ . By definition, the probability to proceed to an exploratory pull on the server  $k$  at round  $t$  is at least  $\lambda_i \min(t^{-\frac{1}{5}}, \frac{1}{N+K})$ . The term  $\lambda_i$  here appears as a pull is guaranteed if a packet appeared at the current time step. Yet the number of exploratory pulls might be much larger in practice as queues should accumulate a large number of uncleared packets at the beginning.

For  $t \geq (N + K)^5$ , it holds:

$$\begin{aligned} \sum_{n=1}^t \min(n^{-\frac{1}{5}}, \frac{1}{N+K}) &= \sum_{n=1}^{(N+K)^5} \frac{1}{N+K} + \sum_{n=(N+K)^5+1}^t n^{-\frac{1}{5}} \\ &\geq (N+K)^4 + \int_{(N+K)^5}^t x^{-\frac{1}{5}} dx - 1 \\ &\geq \frac{1}{4} \left(5t^{\frac{4}{5}} - (N+K)^4 - 4\right) \\ &\geq t^{\frac{4}{5}} - 1. \end{aligned}$$

Lemma 8 then gives for  $N_k(t)$ :

$$\begin{aligned} \mathbb{P}\left(N_k(t) \leq (1 - \delta) \mathbb{E}[N_k(t)]\right) &\leq \exp\left(-\frac{\delta^2 \mathbb{E}[N_k(t)]}{2}\right) \\ \mathbb{P}\left(N_k(t) \leq (1 - \delta) \lambda_i \left(t^{\frac{4}{5}} - 1\right)\right) &\leq \exp\left(-\frac{\lambda_i \delta^2 \left(t^{\frac{4}{5}} - 1\right)}{2}\right). \end{aligned}$$

Which leads for  $\delta = \frac{1}{2}$  to

$$\mathbb{P}\left(N_k(t) \leq \frac{\lambda_i}{2} \left(t^{\frac{4}{5}} - 1\right)\right) \leq \exp\left(-\lambda_i \frac{t^{\frac{4}{5}} - 1}{8}\right). \quad (19)$$

The number of exploratory pulls and the observations on the server  $k$  are independent. Thanks to this, Hoeffding's inequality can be directly used as follows

$$\mathbb{P}\left(|\hat{\mu}_k^i(t) - \mu_k| \geq \varepsilon \mid N_k(t)\right) \leq 2 \exp(-2N_k(t)\varepsilon^2).$$

Using Equation (19) now gives the first concentration inequality for  $\varepsilon \leq \frac{1}{4} \leq \frac{1}{2\sqrt{2}}$ :

$$\begin{aligned} \mathbb{P}\left(|\hat{\mu}_k^i(t) - \mu_k| \geq \varepsilon\right) &\leq 2 \exp\left(-\lambda_i \left(t^{\frac{4}{5}} - 1\right) \varepsilon^2\right) + \exp\left(-\lambda_i \frac{t^{\frac{4}{5}} - 1}{8}\right) \\ &\leq 3 \exp\left(-\lambda_i \left(t^{\frac{4}{5}} - 1\right) \varepsilon^2\right). \end{aligned}$$



**Concentration for  $\hat{\lambda}$ .** Consider agent  $i$  in the following. First show a concentration inequality for  $\tilde{\mu}$ . Denote by  $N(t)$  the total number of exploratory pulls on servers proceeded by player  $i$  at round  $t$ , i.e.,  $N(t) = \sum_{k=1}^K N_k(t)$ . Similarly to Equation (19), it can be shown that

$$\mathbb{P}\left(N(t) \leq \lambda_i K \frac{t^{\frac{4}{5}} - 1}{2}\right) \leq \exp\left(-\lambda_i K \frac{t^{\frac{4}{5}} - 1}{8}\right).$$

Lemma 8 then gives for  $\delta \in (0, 1)$ :

$$\begin{aligned} \mathbb{P}(|\tilde{\mu} - \bar{\mu}| \geq \delta \bar{\mu}) &\leq 2 \exp\left(-\lambda_i K \delta^2 \bar{\mu} \frac{t^{\frac{4}{5}} - 1}{8}\right) + \exp\left(-\lambda_i K \frac{t^{\frac{4}{5}} - 1}{8}\right) \\ &\leq 3 \exp\left(-\lambda_i K \delta^2 \bar{\mu} \frac{t^{\frac{4}{5}} - 1}{8}\right). \end{aligned}$$

Note that  $|\tilde{\mu} - \bar{\mu}| \leq \delta \bar{\mu}$  implies  $|\frac{1}{\tilde{\mu}} - \frac{1}{\bar{\mu}}| \leq \frac{\delta}{(1-\delta)\bar{\mu}}$ . So this gives the following inequality:

$$\mathbb{P}\left(\left|\frac{1}{\tilde{\mu}} - \frac{1}{\bar{\mu}}\right| \geq \frac{\delta}{(1-\delta)\bar{\mu}}\right) \leq 3 \exp\left(-\lambda_i K \delta^2 \bar{\mu} \frac{t^{\frac{4}{5}} - 1}{8}\right). \quad (20)$$

A concentration bound on  $\hat{S}_j^i$  can be shown similarly for any  $\delta \in (0, 1)$

$$\mathbb{P}\left(\left|\hat{S}_j^i(t) - \left(1 - \frac{\lambda_j}{2}\right)\bar{\mu}\right| \geq \delta \bar{\mu}\right) \leq 3 \exp\left(-\lambda_i K \delta^2 \bar{\mu} \frac{t^{\frac{4}{5}} - 1}{8}\right). \quad (21)$$

Now recall that the estimate of  $\lambda_j$  is defined by  $\hat{\lambda}_j = 2 - \frac{2\hat{S}_j^i}{\bar{\mu}}$ . We then have the following identity:

$$\hat{\lambda}_j - \lambda_j = 2\left(\frac{1}{\tilde{\mu}} - \frac{1}{\bar{\mu}}\right)\hat{S}_j^i + \frac{2}{\bar{\mu}}\left(\left(1 - \frac{\lambda_j}{2}\right)\bar{\mu} - \hat{S}_j^i\right).$$

Since  $\hat{S}_j^i \in [0, 1]$ , it yields for  $\varepsilon \leq \frac{1}{4}$  and  $x \in (0, 1)$ :

$$\begin{aligned} \mathbb{P}\left(\left|\hat{\lambda}_j^i(t) - \lambda_j\right| \geq \varepsilon\right) &\leq \mathbb{P}\left(\left|2\left(\frac{1}{\tilde{\mu}} - \frac{1}{\bar{\mu}}\right)\hat{S}_j^i\right| \geq x\varepsilon \text{ or } \left|\frac{2}{\bar{\mu}}\left(\left(1 - \frac{\lambda_j}{2}\right)\bar{\mu} - \hat{S}_j^i\right)\right| \geq (1-x)\varepsilon\right) \\ &\leq \mathbb{P}\left(\left|\frac{1}{\tilde{\mu}} - \frac{1}{\bar{\mu}}\right| \geq \frac{x\varepsilon}{2\hat{S}_j^i} \mid \hat{S}_j^i \leq \left(1 + \frac{1-x}{8}\right)\bar{\mu}\right) + \mathbb{P}\left(\left|\hat{S}_j^i(t) - \left(1 - \frac{\lambda_j}{2}\right)\bar{\mu}\right| \geq \frac{(1-x)\bar{\mu}\varepsilon}{2}\right) \\ &\leq \mathbb{P}\left(\left|\frac{1}{\tilde{\mu}} - \frac{1}{\bar{\mu}}\right| \geq \frac{\delta}{(1-\delta)\bar{\mu}} \text{ for } \delta = \frac{4x\varepsilon}{9}\right) + \mathbb{P}\left(\left|\hat{S}_j^i(t) - \left(1 - \frac{\lambda_j}{2}\right)\bar{\mu}\right| \geq \frac{(1-x)\bar{\mu}\varepsilon}{2}\right) \end{aligned}$$

Taking  $x = \frac{9}{17}$  leads to  $\frac{4x}{9} = \frac{1-x}{2}$  and thus, using Equations (20) and (21):

$$\begin{aligned} \mathbb{P}\left(\left|\hat{\lambda}_j^i(t) - \lambda_j\right| \geq \varepsilon\right) &\leq 6 \exp\left(-\lambda_i K \left(\frac{8}{17}\right)^2 \bar{\mu} \frac{t^{\frac{4}{5}} - 1}{32} \varepsilon^2\right) \\ &\leq 6 \exp\left(-\lambda_i K \bar{\mu} \frac{t^{\frac{4}{5}} - 1}{145} \varepsilon^2\right) \end{aligned}$$

□

In the following, let  $c_1 = 0.1$  and  $c_2 = \frac{4}{(1-2c_1)/\sqrt{e}-2c_1} \approx 14$ . For a problem instance, let the good event  $\mathcal{E}_t$  at time  $t$  be defined as

$$\mathcal{E}_t := \left\{ \|\hat{\lambda}^i - \lambda, \hat{\mu}^i - \mu\|_\infty \leq \frac{0.1\Delta^2}{2c_2 2^{2K^2} KN}, \forall i \in [N] \right\}.$$

As  $\Delta$  is smaller than 1, the right hand term in the definition of  $\mathcal{E}_t$  is smaller than  $c_1\Delta$ . Thanks to Lemmas 1 and 3,  $\mathcal{E}_t$  then guarantees that any player will collide with another player with probability at most  $0.1\Delta$ , i.e.,  $\forall i \in [N]$ ,

$$\mathbb{P}_{\omega \sim \mathcal{U}(0,1)} \left( \exists j \in [N], \hat{A}_t^i(\omega) \neq \hat{A}_t^j(\omega) \mid \mathcal{E}_t \right) \leq 0.1\Delta.$$

Moreover, thanks to Lemma 1, under  $\mathcal{E}_t$ ,

$$\lambda_i \leq \sum_{k=1}^K \hat{P}_{i,k} \mu_k - \left( \frac{1-2c_1}{\sqrt{e}} - 2c_1 \right) \Delta.$$

These two last inequalities lead to the following lemma.

**Lemma 10.** For  $t \geq \frac{2^5 K^5}{0.08^5 \Delta^5}$ , denote by  $\mathcal{H}_t$  the history of observations up to round  $t$ . Then

$$\mathbb{E} [S_t^i \mid \mathcal{E}_t, \mathcal{H}_t] \geq \lambda_i + 0.1\Delta.$$

*Proof.* This is a direct consequence of the following decomposition:

$$\begin{aligned} \mathbb{E} [S_t^i \mid \mathcal{E}_t, \mathcal{H}_t] &\geq \overbrace{(1 - (N+K)t^{-\frac{1}{5}})}^{\text{proba to exploit}} \left( \overbrace{\hat{P}_{i,k} \mu_k}^{\text{proba to clear}} - \overbrace{\mathbb{P}(\exists j \in [N], \hat{A}_t^i(\omega) \neq \hat{A}_t^j(\omega) \mid \mathcal{E}_t)}^{\text{collision proba}} \right) \\ &\geq (1 - (N+K)t^{-\frac{1}{5}}) \left( \lambda_i + \left( \frac{1-2c_1}{\sqrt{e}} - 2c_1 \right) \Delta - 0.1\Delta \right) \\ &\geq (1 - (N+K)t^{-\frac{1}{5}}) (\lambda_i + 0.18\Delta). \end{aligned}$$

The last inequality is given by  $c_1 = 0.1$  and it leads to

$$\mathbb{E} [S_t^i \mid \mathcal{E}_t, \mathcal{H}_t] \geq \lambda_i + 0.18\Delta - (N+K)t^{-\frac{1}{5}}.$$

For  $t \geq \frac{2^5 K^5}{0.08^5 \Delta^5}$ , the last term is smaller than  $0.08\Delta$ , giving Lemma 10.  $\square$

Define the stopping time

$$\tau := \min \left\{ t \geq \frac{2^5 K^5}{0.08^5 \Delta^5} \mid \forall t \geq s, \mathcal{E}_s \text{ holds} \right\}. \quad (22)$$

**Lemma 11.** For any integer  $r \geq 1$ ,

$$\mathbb{E}[\tau^r] = \mathcal{O} \left( KN \left( \frac{N^{\frac{5}{2}} K^{\frac{5}{2}} 2^{5K^2}}{(\min(1, K\bar{\mu})\lambda)^{\frac{5}{4}} \Delta^5} \right)^r \right),$$

where the  $\mathcal{O}$  notation hides constant factors that only depend on  $r$ .

*Proof.* Define for this proof  $t_0 = \left\lceil \frac{2^5 K^5}{0.08^5 \Delta^5} \right\rceil$ . By definition, if  $\mathcal{E}_t$  does not hold for  $t > t_0$ , then  $\tau \geq t$ . As a consequence, for any  $t > t_0$  and thanks to Lemma 9:

$$\begin{aligned} \mathbb{P}(\tau \geq t) &\leq \mathbb{P}(\neg \mathcal{E}_t) \\ &\leq (3eKN + 6eN^2) \exp \left( -ct^{\frac{4}{5}} \right), \end{aligned}$$

where  $c = c_0 \frac{\min(1, K\bar{\mu})\lambda\Delta^4}{N^2 K^2 2^{4K^2}}$  for some universal constant  $c_0 \leq 1$ .

We can now bound the moments of  $\tau$ :

$$\begin{aligned}\mathbb{E}[\tau^r] &= r \int_0^\infty t^{r-1} \mathbb{P}(\tau \geq t) dt \\ &\leq t_0^r + (3eKN + 6eN^2)r \int_0^\infty t^{r-1} e^{-ct^{\frac{4}{5}}} dt.\end{aligned}$$

Using the change of variable  $u = ct^{\frac{4}{5}}$ , it can be shown that

$$\int_0^\infty t^{r-1} e^{-ct^{\frac{4}{5}}} dt = \frac{5}{4} c^{-\frac{5r}{4}} \Gamma\left(\frac{5r}{4}\right),$$

where  $\Gamma$  denotes the Gamma function. It finally allows to conclude:

$$\begin{aligned}\mathbb{E}[\tau^r] &= \mathcal{O}\left(\frac{K^{5r}}{\Delta^{5r}} + KNc^{-\frac{5r}{4}}\right) \\ &= \mathcal{O}\left(KN \left(\frac{N^{\frac{5}{2}} K^{\frac{5}{2}} 2^{5K^2}}{(\min(1, K\bar{\mu})\lambda)^{\frac{5}{4}} \Delta^5}\right)^r\right).\end{aligned}$$

□

Let  $X_t$  be a random walk biased towards 0 with the following transition probabilities:

$$\begin{aligned}\mathbb{P}(X_{t+1} = X_t + 1) &= p, \quad \mathbb{P}(X_{t+1} = X_t - 1 | X_t > 0) = q, \\ \mathbb{P}(X_{t+1} = X_t | X_t > 0) &= 1 - p - q, \quad \mathbb{P}(X_{t+1} = X_t | X_t = 0) = 1 - p,\end{aligned}\tag{23}$$

and  $X_0 = 0$ .

**Lemma 12.** *The non-asymptotic moments of the random walk defined by Equation (23) are bounded. For any  $t > 0, r > 0$ :*

$$\mathbb{E}[(X_t)^r] \leq \frac{r!}{(\ln(q/p))^r}.$$

*Proof:* Let  $\pi$  be the stationary distribution of the random walk. It verifies the following system of equations:

$$\begin{cases} \pi(z) = p\pi(z-1) + q\pi(z+1) + (1-p-q)\pi(z), \quad \forall z > 0 \\ \pi(0) = (1-p)\pi(0) + q\pi(1) \\ \sum \pi(z) = 1 \end{cases}$$

which gives:

$$\pi(z) = \frac{q-p}{q} \left(\frac{p}{q}\right)^z.$$

Equivalently,  $\pi(z) = \mathbb{P}(\lfloor Y \rfloor = z)$  with  $Y$  an exponential random variable of parameter  $\ln(q/p)$ . This gives:

$$\mathbb{E}_{X \sim \pi} [(X)^r] \leq \frac{r!}{(\ln(q/p))^r}.$$

Let  $\tilde{X}_t$  be the random walk with the same transition probabilities as  $X_t$  and  $\tilde{X}_0 \sim \pi$ . For any  $t > 0$ ,  $\tilde{X}_t \sim \pi$ . Moreover, for any  $t > 0$ ,  $\tilde{X}_t$  stochastically dominates  $X_t$ , which terminates the proof. □

*Proof of Theorem 5.* For  $\tau$  the stopping time defined by Equation (22), Lemma 11 bounds its moments as follows

$$\mathbb{E}[\tau^r] = \mathcal{O}\left(KN \left(\frac{N^{\frac{5}{2}} K^{\frac{5}{2}} 2^{5K^2}}{(\min(1, K\bar{\mu})\lambda)^{\frac{5}{4}} \Delta^5}\right)^r\right).$$

Let

$$p_i = \lambda_i(1 - \lambda_i - 0.1\Delta) \text{ and } q_i = (\lambda_i + 0.1\Delta)(1 - \lambda_i).$$

Let  $X_t^i$  be the random walk biased towards 0 with parameters  $p_i$  and  $q_i$ , with  $X_t^i = 0$  for any  $t \leq 0$ . According to Lemma 10, past time  $\tau$ ,  $Q_t^i$  is stochastically dominated by the random process  $\tau + X_{t-\tau}^i$ . Thus, for any  $t > 0$ , for any  $r > 0$

$$\begin{aligned}
\mathbb{E}[(Q_t^i)^r] &\leq \max(1, 2^{r-1}) (\mathbb{E}[\tau^r] + \mathbb{E}[(X_{t-\tau}^i)^r]) \\
&= \mathcal{O} \left( KN \left( \frac{N^{\frac{5}{2}} K^{\frac{5}{2}} 2^{5K^2}}{(\min(1, K\bar{\mu})\lambda)^{\frac{5}{4}} \Delta^5} \right)^r + \frac{1}{\ln(q_i/p_i)^r} \right) \\
&= \mathcal{O} \left( KN \left( \frac{N^{\frac{5}{2}} K^{\frac{5}{2}} 2^{5K^2}}{(\min(1, K\bar{\mu})\lambda)^{\frac{5}{4}} \Delta^5} \right)^r + \Delta^{-r} \right) \\
&= \mathcal{O} \left( KN \left( \frac{N^{\frac{5}{2}} K^{\frac{5}{2}} 2^{5K^2}}{(\min(1, K\bar{\mu})\lambda)^{\frac{5}{4}} \Delta^5} \right)^r \right).
\end{aligned}$$

□