**RESEARCH ARTICLE**

WILEY

# The price of anarchy in loss systems

## Shoshana Anily[1] | Moshe Haviv[2,3]

[1]Coller School of Management, Tel Aviv University, Tel Aviv, Israel

[2]School of Data Science, The Chinese University of Hong Kong, Shenzhen, China

[3]Department of Statistics and the Center for the Study of Rationality, Hebrew University of Jerusalem, Jerusalem, Israel

**Correspondence**
Shoshana Anily, Coller School of Management, Tel Aviv University, Tel Aviv 69978, Israel.
Email: anily@tauex.tau.ac.il

**Abstract**

Assume a multi-server memoryless loss system. Each server is associated with a service rate and a value of service. Customers from a common Poisson arrival process are routed to the servers in an unobservable way, where the goal is to maximize the long-run expected reward per customer (which is the service value times the probability that the customer is not blocked). We first solve this problem under two criteria: social optimization and Nash equilibrium. Our main result is that the price of anarchy, defined as the ratio between the expected gain under the two criteria, is bounded by 2. We also show, via examples, that this bound is tight for any number of servers.

**KEYWORDS**

loss systems, price of anarchy, routing games, symmetric Nash equilibrium, unobservable queues

## 1 | INTRODUCTION

Customers seek service due to its value. Yet, their net utility might be a function not only of their own actions, but also of the decisions made by the other customers. Consider, for example, customers that choose a local service provider, like a plumber, a handyman, or a technician, by using internet sites that rate service providers in each category of profession. Some of the internet sites, like the Angies List that operates in many big cities in the US (its link is https://www.angieslist.com), base their ratings on customers' reviews. For this reason, customers are asked to rate the service provider upon service completion and specify their satisfaction with respect to the service's quality, price, and possibly other characteristics such as the service provider's availability, promptness, and so forth. These reviews are analyzed according to a number of criteria and summarized into a paragraph that is posted on the internet together with the average rating of that service provider, and the number of reviews it is based on. These ratings are common knowledge, and customers, who seek service and want to maximize the value of service, choose a service provider based on their needs and the posted ratings. Even if the posted ratings reflected the real value of service obtained by each customer, yet, the selection process of a service provider by customers could not be said to be optimal, as servers with high ratings are overloaded with requests for

service, whereas others, whose ratings are somewhat lower, are not called by that many customers. The routing of customers to servers could be optimized in terms of customers' overall utility, if a central controller that had full information on the servers' service rates and their ratings, but had no information on which servers are currently idle, was the one to direct customers to the servers. A natural question to ask here is how much the expected utility of the customers could be improved, in comparison with the case where individuals make the choices by themselves, if a central controller was to direct the customers to the servers. In this article, we pose this question for a variant of the above-described problem, where waiting for a busy server is not an option, as service is needed immediately. More specifically, we assume that if the first server that a customer calls is busy, then the customer quits the system as she turns immediately to a different service system that guarantees urgent service provision.

In the model dealt with here, each customer chooses a server from a collection of servers, without being able to observe if the selected server is busy or idle. The only pieces of information on the servers that are common knowledge among all customers are the arrival rate, and the service's value, and service's rate of each server. Customers, upon arrival, choose a server, and if the selected server is busy serving some other customer, they quit the system immediately. Such systems are called *loss systems*.

*Loss systems* of the type described above, have been considered by Erlang since the second decade of the 20th century, in the context of telephony operators. Such models have been reinforced during the current Covid-19 pandemic, in view of various restrictions that were issued regarding social distancing, while waiting for a service. These restrictions impose a capacity limit on the number of customers waiting for service in a queue inside service or retailer facilities. In many parts of the world, including Canada and Israel, when retailers are allowed to open their premises, they must obey restrictions on social distancing and, in particular, on the maximum number of customers that are allowed to be at the same time inside commercial establishments, which is a function of their floor area. For example, in the province of Quebec in Canada, one customer can be admitted per 20 m² (see https://www.quebec.ca/en/health/health-issues/a-z/2019-coronavirus/restriction-on-the-number-of-customers-admitted-to-commercial-enterprises). This restriction has implied that many small street-front boutiques, vintage, jewelry, and accessory shops whose floor size is about 20 m² are limited to a reception capacity of just one customer at a time. It is interesting to note that one of the implications of the pandemic and the capacity restrictions is that many clients, especially among elderly ones, prefer shopping in small shops rather than in malls, in order to reduce the risk of catching the virus while shopping. A direct effect of such severe capacity restrictions in small shops, is a low departure rate which prolongs the waiting time of potential customers, implying that many of them give up and leave, especially if the weather is bad.

When analyzing decision problems in queueing systems, there are two common criteria, self-optimization, and social-optimization. Under the former, it is selfish individuals who decide which server to try getting service from. Since their utility is not a function only of their choice but also of that made by others, they face a symmetric non-cooperative game and the solution one looks for is that of symmetric Nash equilibrium. In our stochastic setting we, in fact, look for a strategy which is used by all, then under the resulting steady-state conditions, self-optimization leads a singled out customer to follow this strategy as well.

Nash equilibrium is usually not optimal from a social (overall) point of view. In fact, in a symmetric situation, there usually exists a socially optimal strategy, which might be administrated by a central planner and is obeyed by all players. In our setting, a utility of a symmetric strategy is defined as the average social gain per unit of time if this strategy is used indefinitely by all (where social gain is defined as the sum of gains across participants). The ratio between the optimal social utility and the one achieved by the Nash equilibrium strategy, is referred to as the *price of anarchy* (PoA), as it states the social loss due to lack of coordination among the players.

To be specific, we consider a model where customers that are generated by a Poisson arrival process, seek service from one out of a number of servers. Service time is exponential with a server-dependent rate. Likewise, the value of service is server-dependent. The service rate and the value of service of each server, as well as the arrival rate, are common knowledge overall customers. We assume a loss model so a customer who selects a busy server leaves empty-handed. A symmetric profile assigns a routing strategy which states a probability distribution over the servers. Both the equilibrium profile (which turns out to be unique) and the socially optimal profile (where uniqueness is not an issue), come with an expected individual utility, and the inverse of the ratio between these two values, is the PoA.

Given a model, the PoA is parameter-dependent. In our model, it is a function of the arrival rate, the servers' service rates and the servers' valuations, and (indirectly) the number of servers. A commonly asked question is whether it is possible to bound the PoA of a given model by a function of a subset of the parameters defining the model. In the extreme and most looked-at case, a constant upper bound on the PoA is found, namely, a bound that is free of all model's parameters. Such a constant bound may be finitely or infinitely large. In either case, proving that the bound is tight is desirable. Proving tightness means finding an instance, which achieves the bound, or, more generally, proving that there exists a sequence of instances whose associated PoAs converge to the bound. The main result presented in this article is that the PoA of the model presented above is bounded by 2, and that this bound is tight for any number of servers.

The rest of this article is organized as follows. A literature review on the PoA, in the context of routing games, is described in Section 2. Section 3 states the formal model and introduces the required notation. Section 4 develops the equilibrium routing profile and the social value associated with it. Section 5 does the same but now for the socially optimal routing profile. Section 6 compares the results of the previous two sections. In Section 7, the tight PoA of 2 is developed. Section 8 concludes the paper. Three long proofs are relegated in an Appendix.

## 2 | LITERATURE REVIEW

Routing decisions are common in transportation systems, network flow models, and queueing theory. The area of customers decision-making in queues commenced with the pioneering work of Naor (1969). Usually, a number of individuals, which can be discrete or nonatomic, need to select a route, a path, or a server from a set of available options, each of which is associated with a cost. This cost is a function of both the choice to be made and the number of customers that have made the same choice. The terminology and the concept of PoA were first introduced two decades ago in Koutsoupias and Papadimitriou (1999), to measure the inefficiency of equilibria.

A routing model of nonatomic participants that need to travel from some origin to a certain destination along the

edges of a congested network, where the cost of using an edge is a monotone function of its usage level, is analyzed in Roughgarden (2003) and Roughgarden and Tardos (2002). For an extension of these papers, see Roughgarden and Tardos (2004). In Roughgarden (2003) and Roughgarden and Tardos (2002), the cost along each edge represents latency as a function of the edge's congestion. The PoA of minimizing the total latency for affine cost functions was proved in Roughgarden and Tardos (2002) to be tightly bounded by 4/3. In Roughgarden (2003), it is proved that the PoA under a wide class of edge latency functions, is achieved by the simplest networks having a single commodity and parallel links. Moreover, if the cost functions are polynomials with nonnegative coefficients, and their maximal degree is $p$, then the PoA equals $\left[1 - p(p+1)^{-(p+1)/p}\right]^{-1}$, which is asymptotically $\Theta\left(\frac{p}{\ln p}\right)$ as $p \to \infty$, that is, it depends only on the largest degree of the polynomials. For the case of atomic congestion games with latency functions that are polynomial with positive coefficients, and of degree at most $p$, a tight bound on the PoA, one that is a function of $p$, and improves upon the bound given in Awerbuch et al. (2005), is derived in Aland et al. (2006). See also Gkatzelis et al. (2016). For games with a concave cost function (as in our case) but with a finite number of players and unnecessarily limited to symmetric strategies (under both the equilibrium and the optimization criteria), it was shown in Vetta (2002) that the PoA is bounded by 2. The issue of tightness of the bound was not discussed there.

In the context of queues, the pioneering paper is that of Bell and Stidham Jr. (1982), and, in fact, its model is the closest to ours. As in our model, there is a Poisson arrival process of customers where each needs service from one of a number of exponential servers that are characterized by their own service rate and their own waiting line. In Bell and Stidham Jr. (1982), which deals with an unobservable parallel $M/M/1$ queues, homogenous customers select a server while being unable to observe which servers are idle and how many customers are waiting in front of each busy server. The common goal of the customers is to minimize their own mean waiting time based on the servers' service rates. Both the equilibrium and the socially optimal profiles are found in Bell and Stidham Jr. (1982). See also Hassin and Haviv (2003, pp. 62–64). In Haviv and Roughgarden (2007), it is shown that the PoA in this model is tightly bounded by the number of servers. The PoA is unbounded for models similar to Haviv and Roughgarden (2007), but with servers-dependent waiting costs and service distributions that are not necessarily exponential, see Altman et al. (2011) and Ayesta et al. (2010). For another queueing routing problem, where service is granted on a relative priority basis, see Oz et al. (2017). In Gilboa-Freedman et al. (2014), it is shown that in the observable version of the $M/M/1$ queue, the PoA is bounded by 2 as long as the arrival rate is smaller than the service rate. See also Hassin and Snitkovsky (2017) where the dilemma whether to check

(upon a fee) a loss system, prior to joining a regular queue, is dealt with. For more on the PoA for various queueing models, see Hassin (2016) and the references cited therein.

A similar unobservable routing problem of parallel production loss systems, controlled by a central planner, is considered in Anily and Haviv (2017). More specifically, each of $n$ parallel production loss systems, called machines, is characterized by its exponential service rate and its arrival rate. A central planner is allowed to outsource some of the production at a constant cost per unit, as well as reroute the remaining units among the machines, in order to minimize the total cost that consists of the outsourcing cost and the cost of lost units. The problem is formulated as a cooperative game where the players are the machines. The authors show that the game can be reduced to a market game, implying that the competitive equilibrium price cost allocation is at its core. In the problem considered in this article, there is a single stream of arrivals, and outsourcing is not allowed.

## 3 | THE MODEL

Consider a set of $n$ servers where server $i \in N = \{1, \ldots, n\}$ is associated with a general service time distribution with a mean service rate of $\mu_i > 0$. The servers serve a common Poisson arrival process with a mean arrival rate of $\lambda$. The servers have no buffers, and so that a customer that arrives at a server when the server is busy quits the system without being served. Note that due to the insensitivity property of the $M/G/1/1$ model (see, e.g., Haviv, 2013 p. 166), our results hold for any service distribution, where the mean service time at server $i \in N$ is denoted by $(\mu_i)^{-1}$. Our point of departure is that if the arrival rate to a server is $x$ and its service rate is $y$, then, in steady-state, the probability that the server is idle is $y/(x+y)$. Moreover, by the PASTA property, this probability applies also to the state of the server as faced by any new arrival (see, e.g., Haviv, 2013, pp. 133 and 120).

A customer that is served by server $i$ gets a reward of $\alpha_i > 0$, for any $i = 1, \ldots, n$. We assume that the servers are indexed in a strictly decreasing[1] order of their reward values. Furthermore, without loss of generality, the rewards are assumed to be scaled so that the largest reward, namely $\alpha_1$ equals 1, that is, $1 = \alpha_1 > \alpha_2 > \cdots > \alpha_n > 0$. Let $\alpha_{n+1} = 0$ and denote $\mu_{(i)} = \sum_{j=1}^{i} \mu_j$ for $i = 1, \ldots, n$.

Next, we derive the equilibrium and the socially optimal profiles. In both cases, we determine which servers are open, and the rate of arrivals that are routed to each one of them. No customers are routed to servers that are not open. Note that neither the customers nor the controller can see when making a decision, which servers are idle and which are not.

---

[1] We will see below that the assumption of strictly decreasing rewards is without loss of generality relative to the assumption of weakly decreasing rewards.

## 4 | THE EQUILIBRIUM SOLUTION

In this section, we consider the equilibrium solution. Note that by an equilibrium strategy, we refer to a routing strategy of customers to the servers under steady-state conditions, such that if followed by all customers, then an individual customer cannot do better but also follow this strategy. It is not necessarily the case that this is the best option for an individual customer under steady-state conditions (and usually it is not). We present below a mixed strategy that assigns probabilities $(p_i)_{i=1}^n$ to servers, $\sum_{i=1}^n p_i = 1$, where $p_i \geq 0$ is the probability of being routed to server $i$, $1 \leq i \leq n$. A server is said to be *open* if $p_i > 0$. Note that a necessary and sufficient condition for an equilibrium is that any individual customer is indifferent among all open servers, and joining any of the non-open servers is not better than joining any of the open ones.

Our main goal, in this section, is to derive the, in fact, unique equilibrium. In particular, we show that in equilibrium there exists a server, to be denoted by $i^e$, such that only servers indexed by $i$, $1 \leq i \leq i^e$, are open. This is expected, as a low-rewarding server, even if idle with probability one, might be less appealing than another sufficiently high-rewarding server, even when the latter comes with a low probability of being free. Below, we state explicitly the value of $i^e$ and the corresponding routing probabilities. In addition, we derive the individual mean reward resulting from this equilibrium behavior.

**Theorem 1** *In equilibrium, the set of open servers is* $\{1, \ldots, i^e\}$, *where*

$$i^e = \min \left\{ i \in N : \quad \alpha_{i+1} < \frac{\sum_{j=1}^i \mu_j \alpha_j}{\mu_{(i)} + \lambda} \right\}. \tag{1}$$

*Denote by $p_i^e$ the routing probability to open server $i$, $1 \leq i \leq i^e$. Then,*

$$p_i^e = \frac{\mu_i}{\lambda} \left( \frac{\alpha_i}{R^e} - 1 \right), \quad 1 \leq i \leq i^e, \tag{2}$$

*where*

$$R^e = \frac{\sum_{j=1}^{i^e} \mu_j \alpha_j}{\mu_{(i^e)} + \lambda}. \tag{3}$$

*Moreover, $R^e$ is the expected utility of a random customer in equilibrium.*

*Finally, let $\pi_i^e$ be the probability that server $i$, $1 \leq i \leq i^e$, is idle in equilibrium. Then,*

$$\pi_i^e = \frac{\mu_i}{\lambda p_i^e + \mu_i} = \frac{R^e}{\alpha_i}, \quad 1 \leq i \leq i^e. \tag{4}$$

*Proof* First, observe that if, in equilibrium, server $j$ is open, then so is any server whose reward is at least as large as $\alpha_j$, as otherwise there would exist a server $i \in \{1, \ldots, j-1\}$, who is idle and whose reward satisfies $\alpha_i > \alpha_j$. Such a server would be more appealing to some customers than server $j$, implying a migration of some customers to server $i$. Hence, in equilibrium, server $i$ is open for all $1 \leq i \leq i^e$ for some (to be determined) $i^e$, $1 \leq i^e \leq n$. Second, observe that if, in equilibrium, servers $i$ and $j$ are open, then

$$\alpha_i \frac{\mu_i}{\lambda p_i + \mu_i} = \alpha_j \frac{\mu_j}{\lambda p_j + \mu_j}. \tag{5}$$

This is the case since $\mu_k / (\mu_k + \lambda p_k)$ is the probability that server $k$ is idle, $1 \leq k \leq i^e$. Thus, the common value in (5) is the individual expected reward for joining any open server. This observation coupled with the condition $\sum_{i=1}^{i^e} p_i^e = 1$ leads, after some algebra, to (2), (3), and (4). Finally, the fraction on the right-hand side of inequality (1) is the commonly expected reward in equilibrium for the customers who join an open server, which applies to all customers, (see (3)). The first server that is left closed (if such a server exists) is the one who obeys the inequality in (1); that is, the reward to a customer that is routed to this server, even if the server is free, is strictly lower than the expected reward obtained from being routed to the open servers. ∎

*Remark* 1 Inspecting (1) and (3) indicates why assuming strictly decreasing rewards come without loss of generality: had $\alpha_{i+1}$ been equal to $\alpha_i$, one could merge servers $i$ and $i+1$ into a single server whose service rate equals $\mu_i + \mu_{i+1}$ and whose service value coincides with their common service value. Moreover, if one of them is open, so is the other.

The next question is how $i^e$ varies with $\lambda$. It is intuitively clear that the higher $\lambda$ is, the more customers consider migrating from the low-indexed servers to the less busy ones. Hence, the higher $\lambda$ is, the more servers are open in equilibrium. The following theorem states this result in a precise way.

**Corollary 1** *Let $\Lambda^e(m)$ be the set of all values of $\lambda > 0$ where $i^e = m$, $1 \leq m \leq n$. Then, $\lambda \in \Lambda^e(m)$, $1 \leq m \leq n$, where*

$$\Lambda^e(1) \equiv \left( 0, \mu_1 \left( \frac{\alpha_1}{\alpha_2} - 1 \right) \right] \tag{6}$$

$$\Lambda^e(m) \equiv \left( \sum_{i=1}^{m-1} \mu_i \left( \frac{\alpha_i}{\alpha_m} - 1 \right), \sum_{i=1}^m \mu_i \left( \frac{\alpha_i}{\alpha_{m+1}} - 1 \right) \right],$$
$$1 < m < n, \tag{7}$$

$$\Lambda^e(n) \equiv \left( \sum_{i=1}^{n-1} \mu_i \left( \frac{\alpha_i}{\alpha_n} - 1 \right), \infty \right). \tag{8}$$

*Proof* The proof follows immediately from the definition of $i^e$ as given in (1). ∎

Let $\lambda_{\max}^e(m)$, $1 \leq m \leq n-1$, be the largest value of $\lambda$ for which $i^e = m$. Its actual value can be read from (6), (7), and (8). Server 1 is the only open server if and only if $\lambda \in \Lambda^e(1) = (0, \lambda_{\max}^e(1)]$. The set of open servers is $\{1, \ldots, m\}$, for $2 \leq m \leq n-1$, if and only if $\lambda \in \Lambda^e(m) = (\lambda_{\max}^e(m-1), \lambda_{\max}^e(m)]$ and all servers are open if and only if $\lambda > \lambda_{\max}^e(n-1)$. One can check that for the arrival rates $\lambda = \lambda_{\max}^e(m)$, $1 \leq m \leq n-1$, the routing probability $p_m^e$ is strictly positive whereas the routing probability $p_e^{m+1}$ is zero, but any infinitesimal increase of $\lambda$ necessitates the opening of server $m + 1$.

Next, we explain why the equilibrium solution is socially suboptimal: in equilibrium, customers over utilize the high-rewarding servers, implying that it would have been socially better had some of them migrated to those with lower rewards. Their loss would be more than compensated by the gain obtained by those who did not migrate.

**Theorem 2** *The marginal social reward contribution per customer of server i in the equilibrium solution, strictly increases with i, $1 \leq i \leq i^e$.*

*Proof* The expected social reward per customer, given some routing probabilities $p_i$, $1 \leq i \leq n$, equals

$$\sum_{i=1}^{n} \frac{\mu_i \alpha_i p_i}{\lambda p_i + \mu_i}.$$

Its derivative with respect to $p_i$ equals

$$\frac{\alpha_i \mu_i^2}{(\lambda p_i + \mu_i)^2}, \quad 1 \leq i \leq n.$$

Using the value for $p_i^e$, $1 \leq i \leq i^e$, as it appears in (2), coupled with some algebra, leads to the fact that the value of the derivative of the expected social reward per customer for $p_i = p_i^e$ is equal to

$$\frac{(R^e)^2}{\alpha_i}, \quad 1 \leq i \leq i^e,$$

which is indeed increasing with $i$, $1 \leq i \leq i^e$. ∎

## 5 | SOCIAL OPTIMIZATION

In equilibrium, $R^e$ is the utility of an individual customer (see (3)), implying that $\lambda R^e$ is the corresponding utility of the society per unit of time. Yet, a central planner who dictates the routing probabilities might achieve a better social utility. The problem that the central planner faces is

$$\max_{p_1, \ldots, p_n} \sum_{i=1}^{n} \frac{\mu_i \alpha_i p_i}{\lambda p_i + \mu_i} \qquad (9)$$

$$\text{s.t.} \quad \sum_{i=1}^{n} p_i = 1$$

$$p_i \geq 0, \quad 1 \leq i \leq n$$

and, in particular, to derive $p_i^s$, $1 \leq i \leq n$, which forms the optimal solution to this constrained optimization problem.

As in the equilibrium case, it is possible that some of the servers are open and some are not (but the set of open servers does not necessarily coincide with the corresponding one in the equilibrium solution). Also, in social optimization, as expected, the open servers are the lower indexed ones. The details of the optimal social utility are given in the next theorem.

**Theorem 3** *The set of open servers in social optimization is $\{1, \ldots, i^s\}$, where*

$$i^s = \min \left\{ i \in N : \quad \sqrt{\alpha_{i+1}} < \frac{\sum_{j=1}^{i} \mu_j \sqrt{\alpha_j}}{\mu_{(i)} + \lambda} \right\}. \qquad (10)$$

*The socially optimal routing probabilities to the open servers are*

$$p_i^s = \frac{\mu_i}{\lambda} \left( \sqrt{\frac{\alpha_i}{\Theta}} - 1 \right), \quad 1 \leq i \leq i^s, \qquad (11)$$

*where $\Theta$ is the value of the Lagrange multiplier of the equality constraint $\sum_{i=1}^{n} p_i = 1$, which is equal to*

$$\Theta = \left( \frac{\sum_{j=1}^{i^s} \mu_j \sqrt{\alpha_j}}{\mu_{(i^s)} + \lambda} \right)^2. \qquad (12)$$

*The optimal expected reward per customer, to be denoted by $R^s$, equals*

$$R^s = \sum_{i=1}^{i^s} \frac{\mu_i \alpha_i p_i^s}{\mu_i + \lambda p_i^s} = \frac{1}{\lambda} \left( \sum_{j=1}^{i^s} \mu_j \alpha_j - \frac{\left( \sum_{j=1}^{i^s} \mu_j \sqrt{\alpha_j} \right)^2}{\mu_{(i^s)} + \lambda} \right). \qquad (13)$$

*Finally, denote by $\pi_i^s$ the probability that server i is idle, $1 \leq i \leq i^s$, when the socially optimal routing is used. Then,*

$$\pi_i^s = \frac{\mu_i}{\lambda p_i^s + \mu_i} = \sqrt{\Theta \alpha_i} = \frac{\sum_{j=1}^{i^s} \mu_j \sqrt{\alpha_j}}{\mu_{(i^s)} + \lambda} \cdot \frac{1}{\sqrt{\alpha_i}}, \quad 1 \leq i \leq i^s. \qquad (14)$$

*Proof* The observation made in the equilibrium criterion that if server $j$ is open, then so is server $i$, for $1 \leq i \leq j \leq n$, holds here too. Yet, the sets of open servers in the two solutions, do not necessarily coincide. Thus, let the set of open servers be $\{1, 2, \ldots, i^s\}$ for some (to be determined) $i^s$. Inspecting the objective function (9) indicates that this function is separable in its decision variables, $p_i$, $1 \leq i \leq n$. Hence, the derivatives of the summands with respect to the corresponding variables should coincide for all variables which come with an interior solution. The derivative of the $i$-th function with respect

to $p_i$ equals

$$\alpha_i \frac{\mu_i^2}{(\lambda p_i + \mu_i)^2}, \quad 1 \le i \le i^s.$$

Hence, the condition that replaces condition (5), for $1 \le i \le i^s$, is

$$\sqrt{\alpha_i} \frac{\mu_i}{\lambda p_i + \mu_i} = \sqrt{\alpha_j} \frac{\mu_j}{\lambda p_j + \mu_j}. \quad (15)$$

From now on, the proof follows the proof of Theorem 1, where the term $\alpha_i$ is now replaced by $\sqrt{\alpha_i}$. As for the value of the Lagrange multiplier, note that it equals the common derivative value with respect to the optimal routing probabilities of the open servers. ∎

*Remark* 2    By inspecting the right-hand side of (13), it is possible to see that $\sum_{j=1}^{i^s} \mu_j \alpha_j$ is the social gain had servers $\{1, \dots, i^s\}$ worked in a non-stop manner. Obviously, this is an upper bound on the actual term and the second term adjusts to the correct value.

*Remark* 3    Note that Remark 1 holds for the social optimization criterion as well.

Similarly to the analysis of the equilibrium solution, we show next how $i^s$ varies with $\lambda$. It is intuitively clear that the higher $\lambda$ is, the more the central planner wants customers to migrate from low-indexed servers to idle ones. Hence, the higher $\lambda$ is, the more servers are open in the optimal solution. The following theorem states this result quantitatively.

**Corollary 2**    *Let $\Lambda^s(m)$ be the set of all values of $\lambda$ where $i^s = m, 1 \le m \le n$. Then, $\lambda \in \Lambda^s(m)$, $1 \le m \le n$, if and only if*

$$\Lambda^s(1) \equiv \left( 0, \mu_1 \left( \sqrt{\frac{\alpha_1}{\alpha_2}} - 1 \right) \right] \quad (16)$$

$$\Lambda^s(m) \equiv \left( \sum_{i=1}^{m-1} \mu_i \left( \sqrt{\frac{\alpha_i}{\alpha_m}} - 1 \right), \ \sum_{i=1}^{m} \mu_i \left( \sqrt{\frac{\alpha_i}{\alpha_{m+1}}} - 1 \right) \right]$$
$$1 < m < n, \quad (17)$$

$$\Lambda^s(n) \equiv \left( \sum_{i=1}^{n-1} \mu_i \left( \sqrt{\frac{\alpha_i}{\alpha_n}} - 1 \right), \infty \right) \quad (18)$$

*Proof*    The proof follows immediately from the definition of $i^s$ as given in (10). ∎

Let $\lambda_{\max}^s(m)$, $1 \le m \le n-1$, be the largest value of $\lambda$ for which $i^s = m$. Its actual value can be read from (16) and (17). Similarly to the equilibrium case, see Section 4, server 1 is the only open server if and only if $\lambda \in \Lambda^s(1) = \left( 0, \lambda_{\max}^s(1) \right]$. The set of open servers is $\{1, \dots, m\}$, for $2 \le m \le n-1$,

if and only if $\lambda \in \Lambda^s(m) = \left( \lambda_{\max}^s(m-1), \lambda_{\max}^s(m) \right]$ and all servers are open if and only if $\lambda > \lambda_{\max}^s(n-1)$. In addition, for any arrival rate $\lambda = \lambda_{\max}^s(m)$, $1 \le m \le n-1$, the routing probability $p_m^s$ is strictly positive, whereas the routing probability $p_{m+1}^s$ is zero, but any infinitesimal increase of the arrival rate $\lambda$, beyond $\lambda = \lambda_{\max}^s(m)$, necessitates the opening of server $m + 1$. As in Section 4, the range of arrival rates for the case where all $n$ servers are open is not bounded from above.

## 6 | COMPARISON OF THE TWO SOLUTIONS

By definition, the individual reward in equilibrium, $R^e$, is bounded from above by the corresponding reward under the socially optimal routing, $R^s$. In this section, we compare some other properties of the two solutions, while in the next section we bound from above the ratio $\frac{R^s}{R^e}$, which is the PoA that we are looking for. Recall that the largest arrival rate for which server $m + 1$, $1 \le m \le n-1$, is still closed, equals, in equilibrium, $\lambda_{\max}^e(m) = \sum_{i=1}^{m} \mu_i \left( \frac{\alpha_i}{\alpha_{m+1}} - 1 \right)$, see (7), and, under the social optimization criterion, $\lambda_{\max}^s(m) = \sum_{i=1}^{m} \mu_i \left( \sqrt{\frac{\alpha_i}{\alpha_{m+1}}} - 1 \right)$, see (17).

Next, we generalize the definition of $\Theta$, the Lagrange multiplier of the solution of $R^s$, see (12), and define

$$\Theta_i = \left( \frac{\sum_{j=1}^{i} \mu_j \sqrt{\alpha_j}}{\mu_{(i)} + \lambda} \right)^2 \quad \text{for any } i, \ 1 \le i \le n. \quad (19)$$

In particular, $\Theta_1 = \left( \frac{\mu_1}{\lambda + \mu_1} \right)^2$, and $\Theta = \Theta_{i^s}$. The next proposition states a joint property of $R^e$ and $\Theta$.

**Proposition 1**    *For any instance:*

1. *The sequence $\frac{\sum_{j=1}^{i} \mu_j \alpha_j}{\mu_{(i)} + \lambda}$ is nondecreasing in $i$ for $1 \le i \le i^e$, implying that $R^e = \max_{1 \le i \le i^e} \frac{\sum_{j=1}^{i} \mu_j \alpha_j}{\mu_{(i)} + \lambda}$.*
2. *The sequence $\Theta_i$, (see(19)), is nondecreasing in $i$ for $1 \le i \le i^s$, implying that $\Theta = \max_{1 \le i \le i^s} \Theta_i$.*

*Proof*    Note that for any $i$, $1 \le i \le n$, $\frac{\sum_{j=1}^{i} \mu_j \alpha_j}{\mu_{(i)} + \lambda}$ can be seen as the weighted average of the following decreasing sequence of $i + 1$ numbers, whose first element is 1 and the last element is zero: $1 = \alpha_1 > \alpha_2 > \cdots > \alpha_i > a_{n+1} = 0$. Similarly, $\sqrt{\Theta_i}$ is the weighted average of the following decreasing sequence of $i+1$ numbers, whose first element is 1 and the last element is zero: $1 = \sqrt{\alpha_1} > \sqrt{\alpha_2} > \cdots > \sqrt{\alpha_i} > \sqrt{\alpha_{n+1}} = 0$. The weight of the $j$th element, $1 \le j \le i$, in any of these two sequences is proportional to $\mu_j$, and the weight of the last element in the sequences,

namely 0, is proportional to $\lambda$. The computation process of $R^e$ (respectively, $\Theta$), starts with $i = 1$, that is, the weighted average of the sequence of two elements, namely $\frac{\mu_1 \alpha_1}{\mu_{(1)} + \lambda}$ $\left( \text{respectively,} \right.$

$\left. \frac{\mu_1 \sqrt{\alpha_1}}{\mu_{(1)} + \lambda} \right)$. If the next element that is considered for entering the weighted average is at least as large as the current weighted average, then the new element enters and the new weighted average is at least as large as the former one. The process for computing $R^e$ ends after inserting $\alpha_{i^e}$ into the corresponding former weighted average, see (1) and (3). Similarly, for $\Theta$, the process ends after inserting $\sqrt{\alpha_{i^s}}$ into the corresponding former weighted average, see (10) and (12). ∎

As can be seen from (1) and (10), except for the fact that $i^e$ and $i^s$ may have different values, the routing probabilities to the open servers in the two cases, see (2) together with (3), and (11) together with (12), have similar forms: under equilibrium, the probabilities $p_i^e$ for $1 \leq i \leq i^e$, are linear increasing in $\frac{\mu_i \alpha_i}{\sum_{j=1}^{i^e} \mu_j \alpha_j}$, where in social optimization, the probabilities $p_i^s$ for $1 \leq i \leq i^s$, are linear increasing in $\frac{\mu_i \sqrt{\alpha_i}}{\sum_{j=1}^{i^s} \mu_j \sqrt{\alpha_j}}$. Also, in both solutions, when the issue is whether the next server should be open or not, what matters is only its reward (regardless of its service rate). The next proposition compares some of the properties of the two solutions. In particular, we show that the number of open servers in equilibrium does not exceed the number of open severs under social optimality. This can be explained by the fact that individuals decisions ignore the externalities they inflict on others, making them joining servers who are too congested from the social point of view. A similar situation exists in the case where unlimited queues are formed in front of the servers, as shown in Haviv and Roughgarden (2007).

**Proposition 2** *1. The number of open servers under the social optimization criterion is at least as large as the number of open servers in equilibrium, i.e., $i^e \leq i^s$.*
*2. $\alpha_{i^e+1} < R^e \leq \alpha_{i^e}$.*
*3. $\alpha_{i^s+1} < \Theta \leq \alpha_{i^s}$.*
*4. If $i^e = i^s$ then $\Theta \leq R^e$.*

*Proof*

1. In view of the fact that the sequence of rewards is strictly decreasing, the inequalities, $\lambda_{\max}^e(m) > \lambda_{\max}^s(m)$ for $1 \leq m \leq n - 1$, hold. Thus, if $i^s = m$ for $1 \leq m < n$, then $\lambda \leq \lambda_{\max}^s(m) < \lambda_{\max}^e(m)$, implying that $i^e \leq m$.

2. As explained in the proof of Proposition 1, $R^e = \frac{\sum_{i=1}^{i^e} \mu_i \alpha_i}{\mu_{(i^e)} + \lambda}$ is the weighted average of the

strictly decreasing sequence $\alpha_1, \dots, a_{i^e}, 0$, where the weight of the reward $\alpha_j$ is proportional to $\mu_j$ for $j = 1, \dots, i^e$, and the weight of the last term in the sequence, namely 0, is proportional to $\lambda$. By the definition of $i^e$, see (1), $\alpha_{i^e} \geq \frac{\sum_{i=1}^{(i^e-1)} \mu_i \alpha_i}{\mu_{(i^e-1)} + \lambda}$, implying that if $\alpha_{i^e}$ is added to the sequence $\alpha_1, \dots, \alpha_{i^e-1}, 0$, the weighted average of the sequence does not decrease, and in fact, is less than or equal to $\alpha_{i^e}$, that is, $\frac{\sum_{i=1}^{i^e} \mu_i \alpha_i}{\mu_{(i^e)} + \lambda} \leq \alpha_{i^e}$. Finally, the definition of $\alpha_{i^e+1}$, see (1), implies the reverse of this inequality.

3. By similar arguments to the ones in the previous item, one can prove, by (10), (12), and Proposition 1, that $\sqrt{\alpha_{i^s+1}} < \sqrt{\Theta} \leq \sqrt{\alpha_{i^s}}$, concluding the proof.

4. The last item follows from the well-known *weighted power-mean inequality,* which can be proved by the Cauchy–Schwarz inequality: $\Theta$ is a weighted mean of power 0.5 and $R^e$ is a weighted mean of power 1, and the inequality states that the weighted power mean increases with the power. ∎

Next we show that customers behave greedily in equilibrium, at least in comparison to the social optimization case, that is, the equilibrium routing probabilities to the low-indexed servers are greater than the corresponding routing probabilities under social optimization.

**Proposition 3** *There exists an integer $k$, $1 \leq k \leq i^e$ such that $p_i^e \geq p_i^s$ for $i = 1, \dots, k$, and $p_i^e < p_i^s$ for $i = k+1, \dots, i^e$.*

*Proof* The inequality $i^e \leq i^s$ implies that $\sum_{i=1}^{i^e} p_i^e = 1$ and $\sum_{i=1}^{i^e} p_i^s \leq 1$, and, in particular, that $\sum_{i=1}^{i^e} \left( p_i^e - p_i^s \right) \geq 0$. Consider the sequence $p_i^e - p_i^s$ for $i = 1, \dots, i^e$:

$$p_i^e - p_i^s = \frac{\mu_i}{\lambda} \left( \frac{\mu_{(i^e)} + \lambda}{\sum_{j=1}^{i^e} \mu_j \alpha_j} \alpha_i - 1 \right)$$

$$- \frac{\mu_i}{\lambda} \left( \frac{\mu_{(i^s)} + \lambda}{\sum_{j=1}^{i^s} \mu_j \sqrt{\alpha_j}} \sqrt{\alpha_i} - 1 \right)$$

$$= \frac{\mu_i \alpha_i \left( \lambda + \mu_{(i^e)} \right) \sum_{j=1}^{i^s} \mu_j \sqrt{\alpha_j} - \mu_i \sqrt{\alpha_i} \left( \lambda + \mu_{(i^s)} \right) \sum_{j=1}^{i^e} \mu_j \alpha_j}{\lambda \sum_{j=1}^{i^e} \mu_j \alpha_j \sum_{j=1}^{i^s} \mu_j \sqrt{\alpha_j}}$$

$$= \mu_i \sqrt{\alpha_i} \left( \frac{\sqrt{\alpha_i} \left( \lambda + \mu_{(i^e)} \right) \sum_{j=1}^{i^s} \mu_j \sqrt{\alpha_j} - \left( \lambda + \mu_{(i^s)} \right) \sum_{j=1}^{i^e} \mu_j \alpha_j}{\lambda \sum_{j=1}^{i^e} \mu_j \alpha_j \sum_{j=1}^{i^s} \mu_j \sqrt{\alpha_j}} \right).$$

The sign of the above difference for $i = 1, \dots, i^e$, is determined by the sign of the

numerator of the fraction in parentheses in the last expression, which is a decreasing function in $i \in \{1, \dots, i^e\}$, as the sequence of rewards is decreasing. As $\sum_{i=1}^{i^e} p_i^e = 1$, $\sum_{i=1}^{i^e} p_i^s \leq 1$, and $\sum_{i=1}^{i^e} \left( p_i^e - p_i^s \right) \geq 0$, there exists an index $k \geq 1$ that satisfies $p_i^e - p_i^s \geq 0$ for $i = 1, \dots, k$, and $p_i^e - p_i^s < 0$ for $i = k + 1, \dots, i^e$. ∎

Under social optimization, the marginal social reward is identical across all open servers, and its value $\Theta$ is strictly higher than any of the square root rewards of the closed servers. This is certainly not the case in equilibrium, as implied by Theorem 2. Taking the equilibrium routing probabilities as a point of departure, it is clear that social optimization requires a migration from the low-indexed servers to some of the high-indexed ones, maybe even leading to the opening of some additional servers.

For the sake of the next proposition, denote the probability that a random customer be served by $TP$. Note that $\lambda * TP$ is the throughput of the system. Let $TP^s$ and $TP^e$ be the corresponding values in social optimum and equilibrium, respectively.

In the next proposition, we prove that the system throughput is lower in equilibrium than in social optimum.

> **Proposition 4** *For any instance, $TP^s \geq TP^e$. Furthermore, a strict inequality, $TP^s > TP^e$, holds for all instances, except for instances where $i^e = i^s = 1$..*
>
> *Proof* See Appendix. ∎

## 7 | THE PRICE OF ANARCHY

The PoA is the ratio between the expected gain under the socially optimal routing (13) and the corresponding value under the equilibrium routing (3). In other words, the PoA measures social inefficiency due to both the lack of coordination between the individuals who act selfishly and the assumed result of their egocentric behavior, namely, the use of Nash equilibrium.

Comparing the equilibrium condition (5) with the social optimization condition (15), it is possible to see that the two strategies coincide in the case where the rewards are identical across all servers. In particular, the PoA then equals 1. It is interesting to observe that this is the case in spite of the fact that the service rates can be different.

In general, the ultimate goal once both the socially optimal and the Nash equilibrium solutions are derived, is to bound the PoA as tightly as possible. Generally, the PoA is a function of all the parameters of the model, which in our case include (a) the number of servers $n$, (b) the arrival rate $\lambda$, (c) the individual service rates, $\mu_j, j = 1, \dots, n$, and (d) the servers' rewards $\alpha_j$, for $j = 1, \dots, n$. A bound on the PoA is

usually a function of a proper subset of these parameters (and hence it is applicable to all possible values of the complementary subset of parameters). Clearly, a tight bound is preferable, where by "tight" we mean that there exists an instance for which the bound is achievable. Constant bounds, which are parameter-free, are, in particular, of special interest. In some problems, it is possible to show that no finite bound exists. The objective of this article is, in fact, to characterize the PoA as much as possible.

The next proposition provides a bound on the PoA. It does not lead to the tightest bound (as we show shortly), but nevertheless, it does shed some light on the PoA.

**Proposition 5**

$$PoA = \frac{R^s}{R^e} \leq \frac{\sum_{j=1}^{i^s} \mu_j \sqrt{\alpha_j}}{\sum_{j=1}^{i^s} \mu_j \alpha_j}. \tag{20}$$

*Proof* It is easy to check that had the rewards been replaced by their square roots, the resulting equilibrium routing probabilities would have coincided with the original socially optimal ones. The reward in equilibrium, under the modified data, would have been equal to

$$\frac{\sum_{j=1}^{i^s} \mu_j \sqrt{\alpha_j}}{\mu_{(i^s)} + \lambda}. \tag{21}$$

Thus, (i) the equilibrium optimal routing probabilities for a sequence of rewards $\left( \sqrt{\alpha_j} \right)_{j=1}^n$ coincide with the optimal routing probabilities that maximize the social reward $R^s$ for the sequence of rewards $\left( \alpha_j \right)_{j=1}^n$, and (ii) $\sqrt{\alpha_j} \geq \alpha_j$ for $1 \leq j \leq n$, as the sequence $\left( \alpha_j \right)_{j=1}^n$ is bounded from above by 1. These observations imply that

$$R^s \leq \frac{\sum_{j=1}^{i^s} \mu_j \sqrt{\alpha_j}}{\mu_{(i^s)} + \lambda}. \tag{22}$$

As $i^e$ is chosen so that $\frac{\sum_{j=1}^k \mu_j \alpha_j}{\mu_{(k)} + \lambda}$ is maximized over $k \in \{1, \dots, i^e\}$, see Proposition 1, the inequality $R^e \geq \frac{\sum_{j=1}^{i^s} \mu_j \alpha_j}{\mu_{(i^s)} + \lambda}$ follows. ∎

Unfortunately, the upper bound on the PoA given in (20) can be made infinitely large while actually, the PoA itself is relatively small. This possibility is exemplified next.

> **Example 1** *Consider an instance with $n = 2$, $\lambda = 1$, $\alpha_1 = 1$, $\alpha_2 = \varepsilon < 1$, $\mu_1 = \varepsilon$, and $\mu_2 = 1$. As $\sqrt{\alpha_2} > \alpha_2 = \varepsilon > \frac{\varepsilon}{1+\varepsilon}$, the two servers are open in equilibrium and under social optimization. In order to compute the PoA we need to compare $R^e = \frac{\varepsilon + \varepsilon}{1 + \varepsilon + 1} = \frac{2\varepsilon}{2 + \varepsilon}$, with $R^s = \varepsilon + \varepsilon - \frac{(\varepsilon + \sqrt{\varepsilon})^2}{2 + \varepsilon}$. Thus, $PoA = \frac{R^s}{R^e} = 1.5 + 0.5\varepsilon - \sqrt{\varepsilon}$, implying that $\lim_{\varepsilon \to 0} PoA = 1.5$. However,*

$PoA \leq \lim_{\varepsilon \to 0} \frac{\varepsilon + \sqrt{\varepsilon}}{2\varepsilon} = \lim_{\varepsilon \to 0} 0.5 \left(1 + \varepsilon^{-0.5}\right) = \infty$, *by (20)*.

In the next lemma, we prove that for case of $n = 2$, though the bound on the PoA as given in Proposition 5, can be made infinitely large, as shown by Example 1, the PoA is quite small. The lemma bounds the PoA by a function of the service rates of the two servers, namely $\mu_1$ and $\mu_2$, and is independent of the rewards. Further, the lemma shows that the bound on the PoA is tightly bounded by two. Indeed, the assumption $1 = \alpha_1 > \alpha_2$ on the reward values, explains why the bound on the PoA as given in the following lemma, is not symmetric in the two service rates $\mu_1$ and $\mu_2$.

**Lemma 1** *For the case where $n = 2$, the PoA is bounded by*

$$\frac{1 + \mu_1}{1 + \mu_1 + \mu_2} \left(1 + 2\mu_2 - 2\mu_2 \frac{\sqrt{\mu_1}}{\sqrt{1 + \mu_1}}\right). \quad (23)$$

*This value itself is tightly bounded by 2.*

*Proof* See Appendix. ∎

The next theorem generalizes Lemma 1, and is the main contribution of this article.

**Theorem 4** *The PoA of the loss system is bounded by 2. More specifically, for the case where $i^e = i^s = 1$ the PoA = 1. For any number of open servers $i^s \geq 2$, the PoA is smaller than or equal to 2. Finally, for any $i^s \geq 2$, this bound is tight.*

*Proof* See Appendix. ∎

As shown above, the PoA is much less sensitive to the number of servers in the parallel multi-server loss system with rewards than in the classic multi-server queueing model. This can be explained by the fact that in the standard M/M/1 model with unlimited capacity, the individual cost is unbounded, while in the loss system the individual cost is bounded by $\alpha_1 = 1$, that is, the loss from not getting service from the best server.

## 8 | CONCLUSIONS

As discussed in Section 1, many service systems have been greatly affected by the recent pandemic due to severe restrictions on social distancing. As a result, service systems with short buffers have become common in many parts of the world and we expect that this will also impact the research on such systems. This article contributes to the currently well-established area of research on the PoA, in particular to models in which decision makers need to select a route from a number of feasible ones in a given network. In most of the

decision models studied, the cost usually comes in the form of latency. We looked at a multi-server model in which the reward is due to the value of service, if received. As we noted, at the beginning of Section 7, there is no loss of efficiency due to selfish behavior in the case of homogeneity of service valuations across servers and hence we considered the case where servers vary with the value of service associated with them. We solved both equilibrium and socially optimal routing problems, compared the two resulting routing profiles, and showed that the PoA is tightly bounded by 2, regardless of the number of servers. As we have shown in the analysis, the behavior of customers in equilibrium is more affected by the rewards, that is, customers tend to overvalue servers with high rewards and undervalue the other servers, than is the case when a central planner routes the customers according to the optimal social solution. This result should be compared with the corresponding multi-server routing problem with homogeneous waiting costs, in which the PoA is tightly bounded by the number of servers, but the PoA turns out to be unbounded when the homogeneity assumption is removed (see Altman et al., 2011). As we noted in our literature review, the value of 2 for the PoA appears in other queueing models. A possible explanation for this discrepancy is that although in both cases the costs or rewards are convex functions of the arrival rates, in our loss model, the reward is bounded. Note that the model considered in this article can be classified as an unobservable decision problem as customers are routed to a server without an earlier inspection if the server is busy or idle. As for the observable version of our problem, it is clear that the greedy policy of joining the most rewarding idle server, is an equilibrium. Yet, this is not necessarily the case when social optimization is of concern and hence the PoA is greater than 1. A central controller will probably take into account the number of free servers and their parameters, that is, their service rate and their value of service, while determining which free server will be assigned the next customer. We leave this open question for future research for the PoA of the observable model under the decision if to join or not to join the standard single server Markovian queue (Naor's model; see Gilboa-Freedman et al., 2014).

Recall that we investigated here the PoA for loss systems, where no buffers exist, and thus waiting in queue for service is impossible. However, if we kept the same model, but assumed a finite, positive-sized buffer for waiting customers before each server, then the main concern for customers would have been their mean waiting time, which could replace the valuation parameters in our loss system. A natural question to be asked in this context is, how does the PoA of the steady-state mean waiting time depend on the buffers' size. In the extreme case where the buffers are unlimited in size, as in the case mentioned above, it has been proved in Haviv and Roughgarden (2007) that the PoA, which is the ratio between the equilibrium and the socially optimal mean waiting times, is bounded by the number of servers, and that this bound is tight. Intuitively, it seems that the smaller the buffers' size

$k \geq 1$ is, the smaller the variation within the system of parallel $M/M/1/k$ servers in steady-state is, and hence the smaller the PoA is. In the extreme case analyzed here, where no buffers exist, that is, $k = 1$, the PoA equals 2. It is an open question whether the minimum PoA value over $k \geq 1$ of $M/M/1/k$ systems is indeed 2, consistent with the PoA value that we obtained for loss systems ($k = 1$).

## ACKNOWLEDGMENTS

## DATA AVAILABILITY STATEMENT

Data sharing not applicable to this article as no datasets were generated or analysed during the current study.

## ORCID

*Shoshana Anily* https://orcid.org/0000-0001-8527-1659

## REFERENCES

Aland, S., Dumrauf, D., Gairing, M., Monien, B., & Schoppman, F. (2006). Exact price of anarchy for polynomial congestion games. *Annual symposium on theoretical aspects of computer science (STACS)*, 218–229.

Altman, A., Ayesta, U., & Prabhu, B. J. (2011). Load balancing in processor sharing systems. *Telecommunication Systems*, 47(1–2), 35–48.

Anily, S., & Haviv, M. (2017). Line balancing in parallel $M/M/1$ and loss systems as cooperative games. *Production and Operations Management*, 26(8), 1568–1584.

Awerbuch, B., Azar, Y., & Epstein, A.. (2005). The price of routing unsplittable flow. In *Proceedings of the 37th Annual ACM Symposium on Theory of Computing (STOC)*, 57–66.

Ayesta, U., Brun, O., & Prabhu, B. J. (2010). Price of anarchy in non-cooperative load balancing. In *2010 Proceedings of the IEEE Transaction on INFOCOM, San Diego*, 1–5.

Bell, C. E., & Stidham, S., Jr. (1982). Individual versus social optimization in the allocation of customers to alternative servers. *Management Science*, 29(7), 831–839.

Gilboa-Freedman, G., Hassin, R., & Kerner, Y. (2014). The price of anarchy in the Markovian single server queue. *IEEE Transactions on Automatic Control*, 59(2), 455–459.

Gkatzelis, V., Kollias, K., & Roughgarden, T. (2016). Optimal cost-sharing in general resource games. *Operations Research*, 64(6), 1230–1238.

Hassin, R. (2016). Rational Queueing. CRC Press.

Hassin, R., & Haviv, M. (2003). To queue or not to queue: Equilibrium behavior in queues. Kluwer.

Hassin, R., & Snitkovsky, R. (2017). Strategic customer behavior in queueing system with a loss subsystem. *Queueing Systems: Theory and Applications*, 86(3), 361–387.

Haviv, M. (2013). Queues: A course in queueing theory. Springer.

Haviv, M., & Roughgarden, T. (2007). The price of anarchy in an exponential multi-server. *Operations Research Letters*, 35(4), 421–426.

Koutsoupias, E., & Papadimitriou, C. H. (1999). Worst-case equilibria. In *Proceedings of the 16th Symposium on Theoretical Aspects of Computer Science. LNCS*, 1563, 404–413.

Naor, P. (1969). The regulation of queue size by levying tolls. *Econometrica*, 37, 15–24.

Oz, B., Haviv, M., & Puterman, M. L. (2017). The advantage of relative priority regimes in multi-class multi-server queueing systems with strategic routing. *Operations Research Letters*, 45, 498–502.

Roughgarden, T. (2003). The price of anarchy is independent of the network topology. *Journal of Computer and System Sciences*, 67(2), 341–364.

Roughgarden, T., & Tardos, E. (2002). How bad is selfish routing? *Journal of the Association for Computing Machinery*, 49(2), 236–259.

Roughgarden, T., & Tardos, E. (2004). Bounding the inefficiency of equilibria in nonatomic congestion games. *Games and Economic Behavior*, 47(2), 389–403.

Vetta, A. (2002). Nash equilibria in competitive societies with applications to facility location, traffic routing and auctions. In *Proceedings of the 43rd Annual Symposium on the Foundations of Computer Science (FOCS'02)*, 416–425.

## APPENDIX: PROOFS

### Proof for Proposition 4

The success probability in social optimum is given by $TP^s = \sum_{i=1}^{i^s} p_i^s \pi_i^s$, where $p_i^s$ is given in (11) and $\pi_i^s$ is given in (14). Similarly, the equilibrium is given by $TP^e = \sum_{i=1}^{i^e} p_i^e \pi_i^e$, where $p_i^e$ is given in (2) and $\pi_i^e$ is given in (4). By simple algebra, we get that

$$TP^s = \sum_{i=1}^{i^s} \frac{\mu_i}{\lambda} \left( 1 - \sqrt{\frac{\theta}{\alpha_i}} \right)$$

and,

$$TP^e = \sum_{i=1}^{i^e} \frac{\mu_i}{\lambda} \left( 1 - \frac{R^e}{\alpha_i} \right).$$

Thus,

$$TP^s - TP^e = \sum_{i=1}^{i^e} \frac{\mu_i}{\lambda} \left( \frac{R^e}{\alpha_i} - \sqrt{\frac{\theta}{\alpha_i}} \right) + \sum_{i=i^e+1}^{i^s} \frac{\mu_i}{\lambda} \left( 1 - \sqrt{\frac{\theta}{\alpha_i}} \right)$$

$$= \sum_{i=1}^{i^e} \frac{\mu_i}{\lambda \alpha_i} \left( R^e - \theta \sqrt{\alpha_i} \right) + \sum_{i=i^e+1}^{i^s} \frac{\mu_i}{\lambda} \left( 1 - \sqrt{\frac{\theta}{\alpha_i}} \right) \quad \text{(A1)}$$

We consider first the case where $i^e = i^s$, where the second summand in (A1) vanishes. The fact that the rewards are bounded from above by 1, with Proposition 2 part 4, imply that for any $i = 1 \ldots i^e$, $\theta \sqrt{\alpha_i} \le \theta \le R^e$, implying that

$$TP^s - TP^e = \sum_{i=1}^{i^e} \frac{\mu_i}{\lambda \alpha_i} \left( R^e - \theta \sqrt{\alpha_i} \right) \ge 0$$

The only case where the above difference equals 0, is when $i^s = i^e = 1$, that is, the case where the equilibrium solution coincides with the social optimum.

Next, we consider the case where $i^s > i^e$, and, in particular, $i^s > 1$. The first summand of (A1) is strictly positive as for any $i$, $1 \le i \le i^e$, $\theta \alpha_i \le \theta \le \alpha_{i^s} \le \alpha_{i^e+1} < R^e$, where the first inequality follows from the fact that the rewards are bounded from above by 1 and $\alpha_1 = 1$, the second inequality follows from Part 3 of Proposition 2, the third weak inequality follows from the fact that we consider here the case where $i^s > i^e$, and the last inequality follows from part 2 of Proposition 2. In addition, we assume now that $i^s > 1$ and since the rewards are strictly increasing, we conclude that the first summand of (A1) is strictly positive.

The second summand of (A1) is nonnegative, as by the third part of Proposition 2, $\theta \le \alpha_{i^s} < \alpha_{i^s-1} < \cdots < \alpha_{i^e+1}$, and equality to zero of this summand may hold only if $i^s = i^e + 1$.

Thus, we conclude that except for the case $i^e = i^s = 1$, where $TP^s = TP^e$, (A1) is strictly positive, implying that $TP^s > TP^e$.

**Proof for Lemma 1**

In order to ease notation assume without loss of generality that $\lambda = 1$, on top of our assumptions that $\alpha_1 = 1$ and that $\alpha_2 < 1$. In particular, $\alpha_1$ belongs to the unit interval and we denote the resulting PoA by $PoA(\alpha_2)$. We consider below three exhaustive sub-intervals that may contain $\alpha_2$:

1. If

$$\alpha_2 < \left( \frac{\mu_1}{1 + \mu_1} \right)^2,$$

then $i^s = i^e = 1$, implying that in this range for $\alpha_2$, $PoA(\alpha_2) = 1$.

2. If

$$\left( \frac{\mu_1}{1 + \mu_1} \right)^2 \le \alpha_2 < \frac{\mu_1}{1 + \mu_1},$$

then, $i^e = 1$ and $i^s = 2$. Therefore,

$$PoA(\alpha_2) = \frac{1 + \mu_1}{1 + \mu_1 + \mu_2}$$
$$\times \left( 1 + \mu_2 + \mu_2 \left( \frac{1 + \mu_1}{\mu_1} \alpha_2 - 2\sqrt{\alpha_2} \right) \right).$$

Since $PoA(\alpha_2)$ is a monotone increasing function of $\alpha_2$ in this interval, we obtain an upper bound on $PoA(\alpha_2)$ by replacing $\alpha_2$ by the right end-point of

this interval, namely by $\frac{\mu_1}{1+\mu_1}$, implying that

$$PoA \le \frac{1 + \mu_1}{1 + \mu_1 + \mu_2}$$
$$\times \left( 1 + 2\mu_2 - 2\mu_2 \frac{\sqrt{\mu_1}}{\sqrt{1 + \mu_1}} \right).$$
$$(A2)$$

3. If $\mu_1 / (1 + \mu_1) \le \alpha_2 < 1$, then $i^e = i^s = 2$ and

$$PoA(\alpha_2) = (1 + \mu_1 + \mu_2)$$
$$- \left( \mu_1 + \mu_2 \sqrt{\alpha_2} \right)^2 / (\mu_1 + \mu_2 \alpha_2).$$
$$(A3)$$

Since, in this interval, the function $PoA(\alpha_2)$ is monotone decreasing in $\alpha_2$, it obtains its maximum at the left end-point of this interval, namely, at $\alpha_2 = \mu_1 / (1 + \mu_1)$, which coincides with the maximal point of the previous interval.

Thus, the bound given in (A2) holds in this interval too, and, in fact, as is shown next, it is tight at $\mu_1 / (1 + \mu_1)$ :

$$PoA \left( \frac{\mu_1}{1 + \mu_1} \right) = \frac{1 + \mu_1}{1 + \mu_1 + \mu_2}$$
$$\times \left( 1 + 2\mu_2 - 2\mu_2 \frac{\sqrt{\mu_1}}{\sqrt{1 + \mu_1}} \right). \quad (A4)$$

This concludes the proof of the first part of the Lemma.

Next, we show that (23) is bounded by 2. We commence by substituting $\mu_1 + 1$ by $\mu_{1'}$, implying that $\mu_{1'} > 1$. Then,

$$PoA \left( \frac{\mu_1}{1 + \mu_1} \right) = \frac{\mu_1'}{\mu_1' + \mu_2} \left( 1 + 2\mu_2 - 2\mu_2 \frac{\sqrt{\mu_1' - 1}}{\sqrt{\mu_1'}} \right)$$
$$= \frac{\mu_1'}{\mu_1' + \mu_2} (1 + 2\mu_2) - 2\mu_2 \frac{\sqrt{\mu_1'}\sqrt{\mu_1' - 1}}{\mu_1' + \mu_2}$$
$$< \frac{\mu_1'}{\mu_1' + \mu_2} (1 + 2\mu_2) - 2\mu_2 \frac{(\mu_1' - 1)}{\mu_1' + \mu_2}$$
$$= \frac{\mu_1' + 2\mu_2}{\mu_1' + \mu_2} < 2 \quad (A5)$$

Thus, the PoA for the case of two servers is strictly smaller than 2. Next, we show that this bound is tight by using again (23), while decreasing $\mu_1$ to an infinitesimal value and increasing $\mu_2$ to infinity:

$$\lim_{\mu_2 \to \infty} \lim_{\mu_1 \to 0} \frac{1 + \mu_1}{1 + \mu_1 + \mu_2} \left( 1 + 2\mu_2 - 2\mu_2 \frac{\sqrt{\mu_1}}{\sqrt{1 + \mu_1}} \right)$$
$$= \lim_{\mu_2 \to \infty} \frac{1}{1 + \mu_2} (1 + 2\mu_2) = 2.$$

**Proof for Theorem 4**

The fact that when $i^s = 1$ the PoA equals 1, is trivial, so we proceed to the case where we prove that PoA = 2 for any $i^s \ge$

2. We start by showing that for any number $i^e$, where $1 \leq i^e < i^s$, of open servers in equilibrium, the PoA is bounded by 2 and that this bound is tight. We then complete the proof by showing that the same is the case if $i^s = i^e$.

For ease of notation, we let again, without loss of generality, that $\lambda = 1$, as in Lemma 1. Denote $i^e + 1$ by $k$. In view of the fact that the sequence $(\alpha_i)_{i=1}^n$ is strictly decreasing, and by Proposition 1, $\alpha_k < \frac{\sum_{j=1}^{k-1} \mu_j \alpha_j}{\mu_{(k-1)}+1}$, see (1). Also, $\alpha_i \geq \Theta_{i-1}$ for $i = 2, \ldots, i^s$, see (10), where the definition of $\Theta_i$ is given in (19). Finally, $\Theta_{i^s} \leq \alpha_{i^s}$ which follows from the argument that $\Theta_{i^s-1} \leq \alpha_{i^s}$, see (10), and the fact that $\Theta_{i^s}$ is a weighted average of $\Theta_{i^s-1}$ and $\alpha_{i^s}$. In summary,

$$0 < \alpha_{i^s+1} < \Theta_{i^s} \leq \alpha_{i^s} < \cdots < \alpha_k < \frac{\sum_{j=1}^{k-1} \mu_j \alpha_j}{\mu_{(k-1)} + 1}$$
$$= R^e \leq \alpha_{k-1} < \cdots < \alpha_1 = 1 \qquad (A6)$$

Similarly to the proof for the case where $i^s = 2$, we look for the lim sup of the PoA as a function of the rewards of the servers that are open in social optimization but are closed in equilibrium. We start by proving that under the inequalities of (A6), the assumption that $i^e = k - 1 < i^s$, implies that

$$PoA = \limsup_{\alpha_{i^e+1}, \ldots, \alpha_{i^s}} \frac{R^s}{R^e} = \frac{\sum_{j=1}^{i^s} \mu_j \alpha_j - \frac{\left(\sum_{j=1}^{i^s} \mu_j \sqrt{\alpha_j}\right)^2}{\mu_{(i^s)}+1}}{\frac{\sum_{j=1}^{k-1} \mu_j \alpha_j}{\mu_{(k-1)}+1}} \leq 2$$

For this sake, note that $R^s$ is increasing in any of the rewards $\alpha_k, \ldots, \alpha_{i^s}$. We indeed increase them as much as possible under the constraints of (A6), that is, while the index $i^e$, and therefore the value of $R^e$, are kept unaltered. Towards this end, redefine the sequence of rewards $\alpha_k, \alpha_{k+1}, \ldots, i^s$ recursively as follows: for a sufficiently small $\varepsilon > 0$, start with $\alpha_k = \frac{\sum_{j=1}^{k-1} \mu_j \alpha_j}{\mu_{(k-1)}+1} - \varepsilon = R^e - \varepsilon$ and then let $\alpha_i = \alpha_{i-1} - \varepsilon = \frac{\sum_{j=1}^{k-1} \mu_j \alpha_j}{\mu_{(k-1)}+1} - (i-k+1)\varepsilon = R^e - (i-k+1)\varepsilon$ for $i = k+1, \ldots, i^s$. In addition, we substitute $\mu_1 + 1$ by $\mu_1' > 1$ and, accordingly, we let $\mu_{(j)}' = \mu_1' + \sum_{i=2}^j \mu_i = \mu_{(j)} + 1$ for any $j$, $1 \leq j \leq i^s$. Thus, by the help of (13) and (3), we get

$$\lim_{\varepsilon \to 0} R^s = \sum_{j=1}^{k-1} \mu_j \alpha_j + R^e \sum_{j=k}^{i^s} \mu_j - \frac{\left(\sum_{j=1}^{k-1} \mu_j \sqrt{\alpha_j} + \sqrt{R^e} \sum_{i=k}^{i^s} \mu_i\right)^2}{\mu_{(i^s)}'}$$

$$= R^e \mu_{(k-1)}' + R^e \sum_{j=k}^{i^s} \mu_j$$

$$- R^e \frac{\left(\left(\sqrt{R^e}\right)^{-1} \sum_{j=1}^{k-1} \mu_j \sqrt{\alpha_j} + \sum_{i=k}^{i^s} \mu_i\right)^2}{\mu_{(i^s)}'}$$

$$= R^e \left(\mu_{(i^s)}' - \frac{\left(\sqrt{\mu_{(k-1)}'} \cdot \frac{\sum_{j=1}^{k-1} \mu_j \sqrt{\alpha_j}}{\sqrt{\sum_{j=1}^{k-1} \mu_j \alpha_j}} + \mu_{(i^s)}' - \mu_{(k-1)}'\right)^2}{\mu_{(i^s)}'}\right)$$

$$= \frac{R^e}{\mu_{(i^s)}'} \left(\mu_{(i^s)}'^2 - \left(\sqrt{\mu_{(k-1)}'} \cdot \frac{\sum_{j=1}^{k-1} \mu_j \sqrt{\alpha_j}}{\sqrt{\sum_{j=1}^{k-1} \mu_j \alpha_j}} + \mu_{(i^s)}' - \mu_{(k-1)}'\right)^2\right)$$
$$(A7)$$

Next, by using the equation $a^2 - b^2 = (a+b)(a-b)$, we get that the above equals

$$\frac{R^s}{R^e} = \frac{1}{\mu_{(i^s)}'} \left(2\mu_{(i^s)}' - \mu_{(k-1)}' + \sqrt{\mu_{(k-1)}'} \frac{\sum_{j=1}^{k-1} \mu_j \sqrt{\alpha_j}}{\sqrt{\sum_{j=1}^{k-1} \mu_j \alpha_j}}\right)$$

$$\times \left(\mu_{(k-1)}' - \sqrt{\mu_{(k-1)}'} \frac{\sum_{j=1}^{k-1} \mu_j \sqrt{\alpha_j}}{\sqrt{\sum_{j=1}^{k-1} \mu_j \alpha_j}}\right)$$

In addition, note that $R^s$ is an increasing function of $\mu_{i^s}$ while $R^e$ is not a function of it, implying that the above expression is increasing in $\mu_{i^s}$. Further note that when one increases $\mu_{i^s}$ to infinity, the same is the effect on $\mu_{(i^s)}'$. This, in turn, makes the above ratio converge to 2, implying that the limit of the ratio (A7) when $\mu_{i^s}$ goes to infinity is 2. This concludes the proof that $PoA = 2$ for the case where $i^e < i^s$.

It remains to prove that this bound holds also for the case where $i^s = i^e$. Without loss of generality, assume that $i^e = i^s = n$, that is, all servers are open both in equilibrium and in social optimization. For that to happen, it must hold that $\alpha_n \geq \frac{\sum_{i=1}^{n-1} \mu_i \alpha_i}{\mu_{(n-1)}+1}$, see (1). Next, we show that the PoA as a function of $\alpha_n$, in the interval

$$\alpha_n \in \left[\frac{\sum_{i=1}^{n-1} \mu_i \alpha_i}{\mu_{(n-1)}+1}, \alpha_{n-1}\right) \qquad (A8)$$

is maximized at the left end-point of the interval, namely at $\alpha_n = \frac{\sum_{i=1}^{n-1} \mu_i \alpha_i}{\mu_{(n-1)}+1}$, which we denote by $R_{(n-1)}^e$, as it is the equilibrium reward per customer when the first $n-1$ servers are open.

Towards that end, let $PoA(\alpha_n)$ denote the PoA as a function of $\alpha_n$ in the above interval, while all the other parameters are fixed, so in particular, $i^e = i^s = n$. Using (3) and (13) for this special case, it can be shown, after some algebra, that for $\alpha_n$ in the interval (A8)

$$PoA(\alpha_n) = (\mu_{(n)} + 1) - \left(\sum_{i=1}^n \mu_i \sqrt{\alpha_i}\right)^2 / \sum_{i=1}^n \mu_i \alpha_i. \qquad (A9)$$

The sign of the partial derivative of (A9) with respect to $\alpha_n$ is

$$\text{sign}\left(\frac{\partial PoA(\alpha_n)}{\partial \alpha_n}\right)$$

$$= -\text{sign}\left(2\left(\sum_{i=1}^{n}\mu_i\sqrt{\alpha_i}\right)\frac{\mu_n}{2\sqrt{\alpha_n}}\sum_{i=1}^{n}\mu_i\alpha_i - \mu_n\left(\sum_{i=1}^{n}\mu_i\sqrt{\alpha_i}\right)^2\right)$$

$$= \text{sign}\left(\frac{\mu_n}{\sqrt{\alpha_n}}\sum_{i=1}^{n}\mu_i\alpha_i - \mu_n\sum_{i=1}^{n}\mu_i\sqrt{\alpha_i}\right)$$

$$= -\text{sign}\sum_{i=1}^{n}\mu_i\left(\alpha_i/\sqrt{\alpha_n} - \sqrt{\alpha_i}\right) < 0.$$

Thus, the function $PoA(\alpha_n)$ is decreasing in $\alpha_n$, implying that its maximum is obtained at the break-point $\alpha_n = \frac{\sum_{i=1}^{n-1}\mu_i\alpha_i}{\mu_{(n-1)}+\lambda}$, where the equilibrium solution does not route any customers to server $n$,. This implies that, effectively, the worst PoA for $i^s = n$ in that interval is obtained when $i^e = n-1$, and the rest of the proof follows from the former part of the proof that the PoA is bounded by 2 for $i^e < i^n$ and that this bound is tight.