

Routing and Staffing when Servers are Strategic

Ragavendran Gopalakrishnan

Xerox Research Centre India, Bangalore, Karnataka 560103,
Ragavendran.Gopalakrishnan@xerox.com

Sherwin Doroudi

Tepper School of Business, Carnegie Mellon University, Pittsburgh, PA 15213,
sdoroudi@andrew.cmu.edu

Amy R. Ward

Marshall School of Business, University of Southern California, Los Angeles, CA 90089,
amyward@marshall.usc.edu

Adam Wierman

Department of Computing and Mathematical Sciences, California Institute of Technology, Pasadena, CA 91125,
adamw@caltech.edu

Traditionally, research focusing on the design of routing and staffing policies for service systems has modeled servers as having fixed (possibly heterogeneous) service rates. However, service systems are generally staffed by people. Furthermore, people respond to workload incentives; that is, how hard a person works can depend both on how much work there is, and how the work is divided between the people responsible for it. In a service system, the routing and staffing policies control such workload incentives; and so the rate servers work will be impacted by the system’s routing and staffing policies. This observation has consequences when modeling service system performance, and our objective in this paper is to investigate those consequences.

We do this in the context of the $M/M/N$ queue, which is the canonical model for large service systems. First, we present a model for “strategic” servers that choose their service rate in order to maximize a trade-off between an “effort cost”, which captures the idea that servers exert more effort when working at a faster rate, and a “value of idleness”, which assumes that servers value having idle time. Next, we characterize the symmetric Nash equilibrium service rate under any routing policy that routes based on the server idle time (such as the longest idle server first policy). We find that the system must operate in a quality-driven regime, in which servers have idle time, in order for an equilibrium to exist. The implication is that to have an equilibrium solution the staffing must have a first-order term that strictly exceeds that of the common square-root staffing policy. Then, within the class of policies that admit an equilibrium, we (asymptotically) solve the problem of minimizing the total cost, when there are linear staffing costs and linear waiting costs. Finally, we end by exploring the question of whether routing policies that are based on the service rate, instead of the server idle time, can improve system performance.

Key words: service systems; staffing; routing; scheduling; routing; strategic servers

Subject classifications: Primary: Queues: applications, limit theorems; secondary: Games/group decisions: noncooperative

1. Introduction. There is a broad and deep literature studying the scheduling and staffing of service systems that bridges operations research, applied probability, and computer science. This literature has had, and is continuing to have, a significant practical impact on the design of call centers (see, for example, the survey papers [18] and [1]), health care systems (see, for example, the recent book [29]), and large-scale computing systems (see, for example, the recent book [26]), among other areas. Traditionally, this literature on scheduling and staffing has modeled the servers of the system as having fixed (possibly heterogeneous) service rates and then, given these rates, scheduling and staffing policies are proposed and analyzed. However, in reality, when the servers are *people*, the rate a server chooses to work can be, and often is, impacted by the scheduling and staffing policies used by the system.

For example, if requests are always scheduled to the “fastest” server whenever that server is available, then this server may have the incentive to slow her rate to avoid being overloaded with

work. Similarly, if extra staff is always assigned to the division of a service system that is the busiest, then servers may have the incentive to reduce their service rates in order to ensure their division is assigned the extra staff. The previous two examples are simplistic; however, strategic behavior has been observed in practice in service systems. For example, empirical data from call centers shows many calls that last near 0 seconds [18]. This strategic behavior of the servers allowed them to obtain “rest breaks” by hanging up on customers – a rather dramatic means of avoiding being overloaded with work. For another example, academics are often guilty of strategic behavior when reviewing for journals. It is rare for reviews to be submitted before an assigned deadline since, if someone is known for reviewing papers very quickly, then they are likely to be assigned more reviews by the editor.

Clearly, the strategic behavior illustrated by the preceding examples can have a significant impact on the performance provided by a service system. One could implement a staffing or scheduling policy that is provably optimal under classical scheduling models, where servers are nonstrategic, and end up with far from optimal system performance as a result of undesirable strategic incentives created by the policy. Consequently, it is crucial for service systems to be designed in a manner that provides the proper incentives for such “strategic servers”.

In practice, there are two approaches used for creating the proper incentives for strategic servers: one can either provide structured bonuses for employees depending on their job performance (performance-based payments) or one can provide incentives in how scheduling and staffing is performed that reward good job performance (incentive-aware scheduling). While there has been considerable research on how to design performance-based payments in the operations management and economics communities; the incentives created by scheduling and staffing policies are much less understood. In particular, *the goal of this paper is to initiate the study of incentive-aware scheduling and staffing policies for strategic servers.*

The design of incentive-aware scheduling and staffing policies is important for a wide variety of service systems. In particular, in many systems performance-based payments such as bonuses are simply not possible, e.g., in service systems staffed by volunteers such as academic reviewing. Furthermore, many service systems do not use performance-based compensation schemes; for example, the 2005 benchmark survey on call center agent compensation in the U.S. shows that a large fraction of call centers pay a fixed hourly wage (and have no performance-based compensation) [3].

Even when performance-based payments are possible, the incentives created by scheduling and staffing policies impact the performance of the service system, and thus impact the success of performance-based payments. Further, since incentive-aware scheduling and staffing does not involve monetary payments (beyond a fixed employee salary), it may be less expensive to provide incentives through scheduling and staffing than through monetary bonuses. Additionally, providing incentives through scheduling and staffing eliminates many concerns about “unfairness” that stem from differential payments to employees.

Of course, the discussion above assumes that the incentives created by scheduling and staffing can be significant enough to impact the behavior. A priori it is not clear if they are, since simply changing the scheduling and staffing policies may not provide strong enough incentives to strategic servers to significantly change service rates, and thus system performance. It is exactly this uncertainty that motivates the current paper, which seeks to understand the impact of the incentives created by scheduling and staffing, and then to design incentive-aware staffing and scheduling policies that provide near-optimal system performance without the use of monetary incentives.

1.1. Contributions of this paper. This paper makes three main contributions. We introduce a new model for the strategic behavior of servers in large service systems and, additionally, we initiate the study of staffing and routing in the context of strategic servers. Each of these contributions is described in the following.

Modeling Strategic Servers (Sections 2 and 3): The essential first step for an analysis of strategic servers is a model for server behavior that is simple enough to be analytically tractable and yet rich enough to capture the salient influences on how each server may choose her service rate. Our model is motivated by work in labor economics that identifies two main factors that impact the utility of agents: effort cost and idleness. More specifically, it is common in labor economics to model agents as having some “effort cost” function that models the decrease in utility which comes from an increase in effort [12]. Additionally, it is a frequent empirical observation that agents in service systems engage in strategic behavior to increase the amount of idle time they have [18]. The key feature of the form of the utility we propose in Section 2 is that it captures the inherent trade-off between idleness and effort. In particular, a faster service rate would mean quicker completion of jobs and might result in a higher idle time, but it would also result in a higher effort cost.

In Section 3 of this paper, we apply our model in the context of a $M/M/N$ system, analyzing the first order condition, and provide a necessary and sufficient condition for a solution to the first order condition to be a symmetric equilibrium service rate (Theorem 4). In addition, we discuss the existence of solutions to the first order condition, and provide a sufficient condition for a unique solution (Theorem 5). These results are necessary in order to study staffing and routing decisions, as we do in Sections 4 and 5; however, it is important to note that the model is applicable more generally as well.

Staffing Strategic Servers (Section 4): The second piece of the paper studies the impact strategic servers have on staffing policies in multi-server service systems. The decision of a staffing level for a service system has a crucial impact on the performance of the system. As such, there is a large literature focusing on this question in the classical, nonstrategic, setting, and the optimal policy is well understood. In particular, the number of servers that must be staffed to ensure stability in a conventional $M/M/N$ queue with arrival rate λ and fixed service rate μ should be strictly larger than the offered load, λ/μ . However, when there are linear staffing and waiting costs, the economically optimal number of servers to staff is more. Specifically, the optimal policy employs the square root of the offered load more servers [8]. This results in efficient operation, because the system loading factor $\lambda/(N\mu)$ is close to one; and maintains quality of service, because the customer wait times are small (on the order of $1/\sqrt{\lambda}$). Thus, this is often referred to as the Quality and Efficiency Driven (QED) regime or as square-root staffing.

Our contribution in this paper is to initiate the study of staffing strategic servers. In the presence of strategic servers, the offered load depends on the arrival rate, the staffing, and the routing, through the servers’ choice of their service rate. We show that an equilibrium service rate exists only if the number of servers staffed is order λ more than the aforementioned square-root staffing (Theorem 7). In particular, the system must operate in a quality-driven regime, in which the servers have idle time, instead of the quality-and-efficiency driven regime that arises under square-root staffing, in which servers do not have idle time. Then, within the set of policies that admit an equilibrium service rate, we (asymptotically) solve the problem of minimizing the total cost, when there are linear staffing costs and linear waiting costs (Theorem 8).

Routing to Strategic Servers (Section 5): The final piece of this paper studies the impact of strategic servers on the design of scheduling policies in multi-server service systems. When servers are not strategic, how to schedule (dispatch) jobs to servers in multi-server systems is well understood. In particular, the most commonly proposed policies for this setting include Fastest Server First (FSF), which dispatches arriving jobs to the idle server with the fastest service rate; Longest Idle Server First (LISF), which dispatches jobs to the server that has been idle for the longest period of time; and Random, which dispatches the job to each idle server with equal probability. When strategic servers are not considered, FSF is the natural choice for reducing the mean response time (though it is not optimal in general [16, 35]). However, in the context of strategic servers the story changes. In particular, we prove that FSF has no symmetric equilibria when strategic servers

are considered, even when there are just two servers. Further, we prove that LISF, a commonly suggested policy for call centers due to its fairness properties, has the same, unique, symmetric equilibrium as random dispatching. In fact, we prove that there is a large *policy-space collapse* – all routing policies that are idle-time-order-based are equivalent in a very strong sense (Theorem 9).

With this in mind, one might suggest that Slowest Server First (SSF) would be a good dispatch policy, since it could incentivize servers to work fast; however, we prove that, like FSF, SSF has no symmetric equilibria (Theorem 10). However, by “softening” SSF’s bias toward slow servers, we are able to identify policies that are guaranteed to have a unique symmetric equilibrium and provide mean response times that are smaller than that under LISF and Random (Theorem 11).

A key message provided by the results described above is that scheduling policies must carefully balance two conflicting goals in the presence of strategic servers: making efficient use of the service capacity (e.g., by sending work to fast servers) while still incentivizing servers to work fast (e.g., by sending work to slow servers). While these two goals are inherently in conflict, our results show that it is possible to balance them in a way that provides improved performance over Random.

1.2. Related work. As we have already described, the question of how to route and staff in many-server systems when servers have fixed, nonstrategic, service rates is well-studied. In general, this is a very difficult question, because the routing depends on the staffing and vice versa. However, when all the servers serve at the same rate, the routing question is moot. Then, [8] shows that square-root staffing, first introduced in [17] and later formalized in [25], is economically optimal when both staffing and waiting costs are linear. Furthermore, square root staffing is remarkably robust: there is theoretical support for why it works so well for systems of moderate size [30], and it continues to be economically optimal both when abandonment is added to the $M/M/N$ model [19] and when there is uncertainty in the arrival rate [33]. Hence, to study the joint routing and staffing question for more complex systems, that include heterogeneous servers that serve at different rates and heterogeneous customers, many authors have assumed square root staffing and show how to optimize the routing for various objective functions (see, for example, [4, 23, 6, 40, 41]). In relation to this body of work, this paper shows that scheduling and routing results for classical many-server systems that assume fixed service rates must be revisited when servers exhibit strategic behavior. This is because they may not admit a symmetric equilibrium service rate in the case of square-root staffing (see Section 4) or be feasible in the case of Fastest Server First routing (see Section 5).

Importantly, the Fastest Server First routing policy mentioned earlier has already been recognized to be potentially problematic because it may be perceived as “unfair”. The issue from an operational standpoint is that there is strong indication in the human resource management literature that the perception of fairness affects employee performance [15, 14]. This has motivated the analysis of “fair” routing policies that, for example, equalize the cumulative server idleness [7, 38], and the desire to find an optimal “fair” routing policy [5, 42]. Another approach is to formulate a model in which the servers choose their service rate in order to balance their desire for idle time (which is obtained by working faster) and the exertion required to serve faster. This leads to a non-cooperative game for a $M/M/N$ queue in which the servers act as strategic players that selfishly maximize their utility.

Finally, the literature that is, perhaps, most closely related to the current paper is the literature on queueing games, which is surveyed in [28]. The bulk of this literature focuses on the impact of customers acting strategically (e.g., deciding whether to join and which queue to join) on queueing performance. Still, there is a body of work within this literature that considers settings where servers can choose their service rate, e.g., [31, 21, 10, 11]. However, in all of the aforementioned papers, there are two servers that derive utility from some monetary compensation per job or per unit of service that they provide, and there are no staffing decisions. In contrast, our work considers systems with more than two servers, and considers servers that derive utility from idle time (and

have a cost of effort). The idea that servers value idle time is most similar to the setting in [20], but that paper restricts its analysis to a two server model. Perhaps the closest previous work to the current paper in analysis spirit is [2], which characterizes approximate equilibria in a market with many servers that compete on price and service level. However, this is similar in theme to [31, 10] in the sense that they consider servers as competing firms in a market. This contrasts with the current paper, where our focus is on competition between servers *within the same firm*.

2. A model for strategic servers. The objective of this paper is to initiate an investigation into the effects of strategic servers on classical management decisions in service systems, e.g., staffing and routing. We start by, in this section, describing formally our model for the behavior of a strategic server.

The term “strategic server” could be interpreted in many ways depending on the server’s goal. Thus, the key feature of the model is the utility function for a strategic server. Our motivation comes from a service system staffed by people who are paid a fixed wage, independent of performance. In such settings, one may expect two key factors to have a first-order impact on the experience of the servers: the amount of effort they put forth and the amount of idle time they have.

Thus, a first-order model for the utility of a strategic server is to linearly combine the cost of effort with the idle time of the server. This gives the following form for the utility of server i in a service system with N servers:

$$U_i(\boldsymbol{\mu}) = I_i(\boldsymbol{\mu}) - c(\mu_i), \quad i \in \{1, \dots, N\}, \quad (1)$$

where $\boldsymbol{\mu}$ is a vector of the rate of work chosen by each server (i.e., the service rate vector), $I_i(\boldsymbol{\mu})$ is the time-average idle time experienced by server i given the service rate vector $\boldsymbol{\mu}$, and $c(\mu_i)$ is the effort cost of server i . We take c to be an increasing, convex function which is the same for all servers. We assume that the strategic behavior of servers (choosing a utility-maximizing service rate) is independent of the state of the system and that the server has complete information about the steady state properties of the system when choosing a rate, i.e., they know the arrival rate, scheduling policy, staffing policy, etc., and thus can optimize $U_i(\boldsymbol{\mu})$.

The key feature of the form of the utility in (1) is that it captures the inherent trade-off between idleness and effort. The idleness, and hence the utility, is a steady state quantity. In particular, a faster service rate would mean quicker completion of jobs and might result in higher idle time in steady state, but it would also result in a higher effort cost. This trade-off then creates a difficult challenge for staffing and routing in a service system. To increase throughput and decrease response times, one would like to route requests to the fastest servers, but by doing so the utility of servers decreases, making it less desirable to maintain a fast service rate. Our model should be interpreted as providing insight into the *systemic* incentives created by scheduling and staffing policies rather than the *transitive* incentives created by the stochastic behavior of the system.

Our focus in this paper will be to explore the consequences of strategic servers for staffing and routing in large service systems, specifically, in the $M/M/N$ setting. However, the model is generic and can be studied in non-queueing contexts as well.

To quickly illustrate the issues created by strategic servers, a useful example to consider is that of a $M/M/1$ queue with a strategic server.

EXAMPLE 1 (THE $M/M/1$ QUEUE WITH A STRATEGIC SERVER). In a classic $M/M/1$ system, jobs arrive at rate λ into a queue with an infinite buffer, where they wait to obtain service from a single server having fixed service rate μ . When the server is strategic, instead of serving at a fixed rate μ , the server chooses her service rate $\mu > \lambda$ in order to maximize the utility in (1). To understand what service rate will emerge, recall that in a $M/M/1$ queue with $\mu > \lambda$ the steady

state fraction of time that the server is idle is given by $I(\mu) = 1 - \frac{\lambda}{\mu}$. Substituting this expression into (1) means that the utility of the server is given by the following concave function:

$$U(\mu) = 1 - \frac{\lambda}{\mu} - c(\mu).$$

We now have two possible scenarios. First, suppose that $c'(\lambda) < 1/\lambda$, so that the cost function does not increase too fast. Then, $U(\mu)$ attains a maximum in (λ, ∞) at a unique point μ^* , which is the optimal (utility maximizing) operating point for the strategic server. Thus, a stable operating point emerges, and the performance of this operating point can be derived explicitly when a specific form of a cost function is considered.

On the other hand, if $c'(\lambda) \geq 1/\lambda$, then $U(\mu)$ is strictly decreasing in (λ, ∞) and hence does not attain a maximum in this interval. We interpret this case to mean that the server's inherent skill level (as indicated by the cost function) is such that the server must work extremely hard just to stabilize the system, and therefore should not have been hired in the first place.

For example, consider the class of cost functions $c(\mu) = c_E \mu^p$. If $c(\lambda) < \frac{1}{p}$, then μ^* solves $\mu^* c(\mu^*) = \frac{\lambda}{p}$, which gives $\mu^* = \left(\frac{\lambda}{c_E p}\right)^{\frac{1}{p+1}} > \lambda$. On the other hand, if $c(\lambda) \geq \frac{1}{p}$, then $U(\mu)$ is strictly decreasing in (λ, ∞) and hence does not attain a maximum in this interval.

Before moving on to the analysis of the $M/M/N$ model with strategic servers, it is important to point out that the model we study focuses on a linear trade-off between idleness and effort. There are certainly many generalizations that are interesting to study in future work. One particularly interesting generalization would be to consider a concave (and increasing) function of idle time in the utility function, since it is natural that the gain from improving idle time from 10% to 20% would be larger than the gain from improving idle time from 80% to 90%. A preliminary analysis highlights that the results in this paper would not qualitatively change in this context.¹

3. The $M/M/N$ queue with strategic servers. Our focus in this paper is on the staffing and routing decisions in large service systems, and so we adopt a classical model of this setting, the $M/M/N$, and adjust it by considering strategic servers, as described in Section 2. The analysis of staffing and routing policies is addressed in Sections 4 and 5, but before moving to such questions, we start by formally introducing the $M/M/N$ model, and performing some preliminary analysis that is useful both in the context of staffing and routing.

3.1. Model and notation. In a $M/M/N$ queue, customers arrive to a service system having N servers according to a Poisson process with rate λ . Delayed customers (those that arrive to find all servers busy) are served according to the First In First Out (FIFO) discipline. Each server is fully capable of handling any customer's service requirements. The time required to serve each customer is independent and exponential, and has a mean of one time unit when the server works at rate one. However, each server strategically chooses her service rate to maximize her own (steady state) utility, and so it is not a priori clear what the system service rates will be.

¹ Specifically, if $g(I_i(\boldsymbol{\mu}))$ replaces $I_i(\boldsymbol{\mu})$ in (1), all the results in Section 3 characterizing equilibria service rates are maintained so long as $g''' < 0$, except for Theorem 5, whose sufficient condition would have to be adjusted to accommodate g . In addition, our results could be made stronger depending on the specific form of g . For example, if g is such that $\lim_{\mu_i \rightarrow \underline{\mu}_i} U_i(\boldsymbol{\mu}) = -\infty$, then, a preliminary analysis reveals that it would not be necessary to impose the stability constraint $\mu_i > \lambda/N$ exogenously. Moreover, every solution to the symmetric first order condition (9) would be a symmetric equilibrium (i.e., the sufficient condition of Theorem 4 as generalized for this case by Footnote 2 would automatically be satisfied).

In this setting, the utility functions that the servers seek to maximize are given by

$$U_i(\boldsymbol{\mu}; \lambda, N, R) = I_i(\boldsymbol{\mu}; \lambda, N, R) - c(\mu_i), \quad i \in \{1, \dots, N\}, \quad (2)$$

where $\boldsymbol{\mu}$ is the vector of service rates, λ is the arrival rate, N is the number of servers (staffing level), and R is the routing policy. $I_i(\boldsymbol{\mu}; \lambda, N, R)$ is the steady state fraction of time that server i is idle. $c(\mu)$ is an increasing, convex function with $c'''(\mu) \geq 0$, that represents the server effort cost.

Note that, as compared with (1), we have emphasized the dependence on the arrival rate λ , staffing level N , and routing policy of the system, R . In the remainder of this article, we expose or suppress the dependence on these additional parameters as relevant to the discussion. In particular, note that the idle time fraction I_i (and hence, the utility function U_i) in (2) depends on how arriving customers are routed to the individual servers.

There are a variety of routing policies that are feasible for the system manager. In general, the system manager may use information about the order in which the servers became idle, the rates at which servers have been working, etc. This leads to the possibility of using simple policies such as Random, which chooses an idle server to route to uniformly at random, as well as more complex policies such as Longest/Shortest Idle Server First (LISF/SISF) and Fastest/Slowest Server First (FSF/SSF). We study the impact of this decision in detail in Section 5.

Given the routing policy chosen by the system manager and the form of the server utilities in (2), the situation that emerges is a competition among the servers for the system idle time. In particular, the routing policy yields a division of idle time among the servers, and both the division and the amount of idle time will depend on the service rates chosen by the servers.

As a result, the servers can be modeled as strategic players in a noncooperative game, and thus the operating point of the system is naturally modeled as an equilibrium of this game. In particular, a Nash equilibrium of this game is a set of service rates $\boldsymbol{\mu}^*$, such that,

$$U_i(\mu_i^*, \boldsymbol{\mu}_{-i}^*; R) = \max_{\mu_i > \frac{\lambda}{N}} U_i(\mu_i, \boldsymbol{\mu}_{-i}^*; R), \quad (3)$$

where $\boldsymbol{\mu}_{-i}^* = (\mu_1^*, \dots, \mu_{i-1}^*, \mu_{i+1}^*, \dots, \mu_N^*)$ denotes the vector of service rates of all the servers except server i . Note that we exogenously impose the (symmetric) constraint that each server must work at a rate strictly greater than $\frac{\lambda}{N}$ in order to define a product action space that ensures the stability of the system.² Such a constraint is necessary to allow steady state analysis, and does not eliminate any feasible symmetric equilibria. We treat this bound as exogenously fixed, however in some situations a system manager may wish to impose quality standards on servers, which would correspond to imposing a larger lower bound (likely with correspondingly larger payments for servers). Investigating the impact of such quality standards is an interesting topic for future work.

Our focus in this paper is on symmetric Nash equilibria. With a slight abuse of notation, we say that $\boldsymbol{\mu}^*$ is a symmetric Nash equilibrium if $\boldsymbol{\mu}^* = (\mu^*, \dots, \mu^*)$ is a Nash equilibrium (solves (3)). Throughout, the term ‘‘equilibrium service rate’’ means a symmetric Nash equilibrium service rate.

We focus on symmetric Nash equilibria for two reasons. First, because the agents we model intrinsically have the same skill level (as quantified by the effort cost functions), a symmetric

² One can imagine that servers, despite being strategic, would endogenously stabilize the system. To test this, one could study a related game where the action sets of the servers are $(0, \infty)$. Then, the definition of the idle time $I_i(\boldsymbol{\mu})$ must be extended into the range of $\boldsymbol{\mu}$ for which the system is overloaded; a natural way to do so is to define it to be zero in this range, which would ensure continuity at $\boldsymbol{\mu}$ for which the system is critically loaded. However, it is not differentiable there, which necessitates a careful piecewise analysis. A preliminary analysis indicates that in this scenario, no $\boldsymbol{\mu} \in (0, \frac{\lambda}{N}]$ can ever be a symmetric equilibrium, and then, the necessary and sufficient condition of Theorem 4 would become $U(\boldsymbol{\mu}^*, \boldsymbol{\mu}^*) \geq \lim_{\mu_1 \rightarrow 0^+} U(\mu_1, \boldsymbol{\mu}^*)$, which is more demanding than (10) (e.g., it imposes a finite upper bound on μ^*), but not so much so that it disrupts the staffing results that rely on this theorem (e.g., Lemma 1 still holds).

equilibrium corresponds to a fair outcome. As we have already discussed, this sort of fairness is often crucial in service organizations [15, 14, 5]. A second reason for focusing on symmetric equilibria is that analyzing symmetric equilibria is already technically challenging, and it is not clear how to approach asymmetric equilibria in the contexts that we consider. Note that we do not rule out the existence of asymmetric equilibria; in fact, they likely exist, and it would be interesting to study whether they lead to better or worse system performance than their symmetric counterparts.

3.2. The $M/M/N$ queue with strategic servers and Random routing. Before analyzing staffing and routing in detail, we first study the $M/M/N$ queue with strategic servers and Random routing. We focus on Random routing first because it is, perhaps, the most commonly studied policy in the classical literature on nonstrategic servers. Further, this importance is magnified by a new “policy-space collapse” result included in Section 5.1.1, which shows that all idle-time-order-based routing policies (e.g., LISF and SISF) have equivalent steady state behavior, and thus have the same steady state behavior as Random routing. We stress that this result stands on its own in the classical, nonstrategic setting of a $M/M/N$ queue with heterogeneous service rates, but is also crucial to analyze routing to strategic servers (Section 5).

The key goal in analyzing a queueing system with strategic servers is to understand the equilibria service rates, i.e., show conditions that guarantee their existence and characterize the equilibria when they exist. Theorems 4 and 5 of Section 3.2.2 summarize these results for the $M/M/N$ queue with Random routing. However, in order to obtain such results we must first characterize the idle time in a $M/M/N$ system in order to be able to understand the “best responses” for servers, and thus analyze their equilibrium behavior. Such an analysis is the focus of Section 3.2.1.

3.2.1. The idle time of a tagged server. In order to characterize the equilibria service rates, a key first step is to understand the idle time of a $M/M/N$ queue. This is, of course, a well-studied model, and so one might expect to be able to use off-the-shelf results. While this is true when the servers are homogeneous (i.e., all the server rates are the same), for heterogeneous systems, closed form expressions are challenging to obtain in general, and the resulting forms are quite complicated [22].

To characterize equilibria, we do need to understand the idle time of heterogeneous $M/M/N$ queues. However, due to our focus on symmetric equilibria, we only need to understand a particular, mild, form of heterogeneity. In particular, we need only understand the best response function for a “deviating server” when all other servers have the same service rate. Given this limited form of heterogeneity, the form of the idle time function simplifies, but still remains quite complicated, as the following theorem shows.

THEOREM 1. *Consider a heterogenous $M/M/N$ system with Random routing and arrival rate $\lambda > 0$, where $N - 1$ servers operate at rate $\mu > \frac{\lambda}{N}$, and a tagged server operates at rate $\mu_1 > \underline{\mu}_1 = (\lambda - (N - 1)\mu)^+$. The steady state probability that the tagged server is idle is given by:*

$$I(\mu_1, \mu; \lambda, N) = \left(1 - \frac{\rho}{N}\right) \left(1 - \frac{\rho}{N} \left(1 - \frac{\mu}{\mu_1}\right) \left(1 + \frac{C(N, \rho)}{N - \left(\rho + 1 - \frac{\mu_1}{\mu}\right)}\right)\right)^{-1}, \quad (4)$$

where $\rho = \frac{\lambda}{\mu}$, and $C(N, \rho)$ denotes the Erlang C formula, given by:

$$C(N, \rho) = \frac{\frac{\rho^N}{N!} \frac{N}{N - \rho}}{\sum_{j=0}^{N-1} \frac{\rho^j}{j!} + \frac{\rho^N}{N!} \frac{N}{N - \rho}}.$$

In order to understand this idle time function more, we derive expressions for the first two derivatives of I with respect to μ_1 in the following theorem. These results are crucial to the analysis of equilibrium behavior.

THEOREM 2. *The first two partial derivatives of I with respect to μ_1 are given by*

$$\frac{\partial I}{\partial \mu_1} = \frac{I^2}{\mu_1^2} \frac{\lambda}{N-\rho} \left(1 + \frac{C(N, \rho)}{N - (\rho + 1 - \frac{\mu_1}{\mu})} + \left(1 - \frac{\mu_1}{\mu} \right) \frac{\mu_1}{\mu} \frac{C(N, \rho)}{\left(N - (\rho + 1 - \frac{\mu_1}{\mu}) \right)^2} \right) \quad (5)$$

$$\frac{\partial^2 I}{\partial \mu_1^2} = -\frac{2I^3}{\mu_1^3} \frac{\lambda}{N-\rho} \left(\left(1 - \frac{\rho C(N, \rho)}{\left(N - (\rho + 1 - \frac{\mu_1}{\mu}) \right)^2} \right) \left(1 + \frac{C(N, \rho)}{N - (\rho + 1 - \frac{\mu_1}{\mu})} \right) + \left(N - \left(1 - \frac{\mu_1}{\mu} \right)^2 \right) \frac{\mu_1}{\mu} \frac{C(N, \rho)}{\left(N - (\rho + 1 - \frac{\mu_1}{\mu}) \right)^3} \right) \quad (6)$$

Importantly, it can be shown that the right hand side of (5) is always positive, and therefore, the idle time is increasing in the service rate μ_1 , as expected. However, it is not clear through inspection of (6) whether the second derivative is positive or negative. Our next theorem characterizes the second derivative, showing that the idle time could be convex at $\mu_1 = \underline{\mu}_1$ to begin with, but if so, then as μ_1 increases, it steadily becomes less convex, and is eventually concave. This behavior adds considerable complication to the equilibrium analysis.

THEOREM 3. *The second derivative of the idle time satisfies the following properties:*

- (a) *There exists a threshold $\mu_1^\dagger \in [\underline{\mu}_1, \infty)$ such that $\frac{\partial^2 I}{\partial \mu_1^2} > 0$ for $\underline{\mu}_1 < \mu_1 < \mu_1^\dagger$, and $\frac{\partial^2 I}{\partial \mu_1^2} < 0$ for $\mu_1^\dagger < \mu_1 < \infty$.*
- (b) *$\frac{\partial^2 I}{\partial \mu_1^2} > 0 \Rightarrow \frac{\partial^3 I}{\partial \mu_1^3} < 0$.*

We remark that it is possible that the threshold μ_1^\dagger could be greater than $\frac{\lambda}{N}$, so, restricting the service rate of server 1 to be greater than $\frac{\lambda}{N}$ does not necessarily simplify the analysis.

3.2.2. Symmetric equilibrium analysis for a finite system. The properties of the idle time function derived in the previous section provide the key tools we need to characterize the symmetric equilibria service rates under Random routing for a $M/M/N$ system.

To characterize the symmetric equilibria, we consider the utility of a tagged server, without loss of generality, server 1, under the mildly heterogeneous setup of Theorem 1. We denote it by

$$U(\mu_1, \mu; \lambda, N) = I(\mu_1, \mu; \lambda, N) - c(\mu_1) \quad (7)$$

For a symmetric equilibrium in $(\frac{\lambda}{N}, \infty)$, we explore the first order and second order conditions for U as a function of μ_1 to have a maximum in $(\underline{\mu}_1, \infty)$.

The first order condition for an interior local maximum at μ_1 is given by:

$$\frac{\partial U}{\partial \mu_1} = 0 \quad \Rightarrow \quad \frac{\partial I}{\partial \mu_1} = c'(\mu_1) \quad (8)$$

Since we are interested in a symmetric equilibrium, we analyze the symmetric first order condition, obtained by plugging in $\mu_1 = \mu$ in (8):

$$\frac{\partial U}{\partial \mu_1} \Big|_{\mu_1=\mu} = 0 \quad \Rightarrow \quad \frac{\lambda}{N^2 \mu^2} \left(N - \frac{\lambda}{\mu} + C \left(N, \frac{\lambda}{\mu} \right) \right) = c'(\mu) \quad (9)$$

Now, suppose that $\mu^* > \frac{\lambda}{N}$ satisfies the symmetric first order condition (9). Then, $\mu_1 = \mu^*$ is a stationary point of $U(\mu_1, \mu^*)$. It follows then, that μ^* will be a symmetric equilibrium for the servers (satisfying (3)) if and only if $U(\mu_1, \mu^*)$ attains a global maximum at $\mu_1 = \mu^*$ in the interval $(\frac{\lambda}{N}, \infty)$. While an obvious necessary condition for this is that $U(\mu^*, \mu^*) \geq U(\frac{\lambda}{N}, \mu^*)$, we show, perhaps surprisingly, that it is also sufficient, in the following theorem.

THEOREM 4. $\mu^* > \frac{\lambda}{N}$ is a symmetric equilibrium if and only if it satisfies the symmetric first order condition (9), and the inequality $U(\mu^*, \mu^*) \geq U(\frac{\lambda}{N}, \mu^*)$, i.e.,

$$c(\mu) \leq c\left(\frac{\lambda}{N}\right) + \left(1 - \frac{\rho}{N}\right) \left(1 + \left(1 - \frac{\rho}{N} + \frac{C(N, \rho)}{N-1}\right)^{-1}\right)^{-1}. \quad (10)$$

Finally, we need to understand when the symmetric first order condition (9) admits a feasible solution $\mu^* > \frac{\lambda}{N}$. Towards that, we present sufficient conditions for at least one feasible solution, as well as for a *unique* feasible solution.

THEOREM 5. If $c'(\frac{\lambda}{N}) < \frac{1}{\lambda}$, then the symmetric first order condition (9) has at least one solution for μ in $(\frac{\lambda}{N}, \infty)$. In addition, if $2\frac{\lambda}{N}c'(\frac{\lambda}{N}) + (\frac{\lambda}{N})^2 c''(\frac{\lambda}{N}) \geq 1$, then the symmetric first order condition (9) has a unique solution for μ in $(\frac{\lambda}{N}, \infty)$.

In the numerical results that follow, we see instances of zero, one, and two equilibria.³ Interestingly, when more than one equilibrium exists, the equilibrium with the largest service rate, which leads to best system performance, also leads to highest server utility, and hence is also most preferred by the servers, as the following theorem shows.

THEOREM 6. If the symmetric first order condition (9) has two solutions, say μ_1^* and μ_2^* , with $\mu_1^* > \mu_2^* > \frac{\lambda}{N}$, then $U(\mu_1^*, \mu_1^*) > U(\mu_2^*, \mu_2^*)$.

3.3. Numerical examples. Because of the complexity of the expression for the equilibrium service rate(s) given by the first order condition (9) and the possibility of multiple equilibria, we discuss a few numerical examples here in order to provide intuition. In addition, we point out some interesting characteristics that emerge as a consequence of strategic server behavior.

We present two sets of graphs below: one that varies the arrival rate λ while holding the staffing level fixed at $N = 20$ (Figure 1), and one that varies the staffing level N while holding the arrival rate fixed at $\lambda = 2$ (Figure 2). In each set, we plot the following two equilibrium quantities: (a) service rates, and (b) mean steady state waiting times. Note that the graphs in Figure 2 only show data points corresponding to integer values of N ; the thin line through these points is only meant as a visual tool that helps bring out the pattern. Each of the four graphs shows data for three different effort cost functions: $c(\mu) = \mu$, $c(\mu) = \mu^2$, and $c(\mu) = \mu^3$, which are depicted in red, blue, and green respectively. The data points in Figure 2 marked \times and \diamond correspond to special staffing levels $N^{ao,2}$ and $N^{opt,2}$ respectively, which are introduced later, in Section 4.

The first observation we make is that there are at most two equilibria. Further, for large enough values of the minimum service rate $\frac{\lambda}{N}$, there is no equilibrium. (In Figure 1(a) where N is fixed, this happens for large λ , and in Figure 2(a) where λ is fixed, this happens for small N .) On the other hand, when the minimum service rate $\frac{\lambda}{N}$ is small enough, there is a unique equilibrium; for this range, even if the symmetric first order condition (9) has another solution greater than $\frac{\lambda}{N}$, it fails to satisfy (10). If an intermediate value of $\frac{\lambda}{N}$ is small enough for (9) to have two feasible solutions, but not too small so that both solutions satisfy (10), then there are two equilibria.

The second observation we make is that the two equilibria have very different behaviors. As illustrated in Figure 1(a), the larger equilibrium service rate first increases and then decreases while

³ In general, the symmetric first order condition (9) can be rewritten as

$$\mu^2 c'(\mu) + \frac{\lambda}{N^2} (\rho - C(N, \rho)) - \frac{\lambda}{N} = 0.$$

Note that, when the term $\rho - C(N, \rho)$ is convex in μ , it follows that the left hand side of the above equation is also convex in μ , which implies that there are at most two symmetric equilibria.

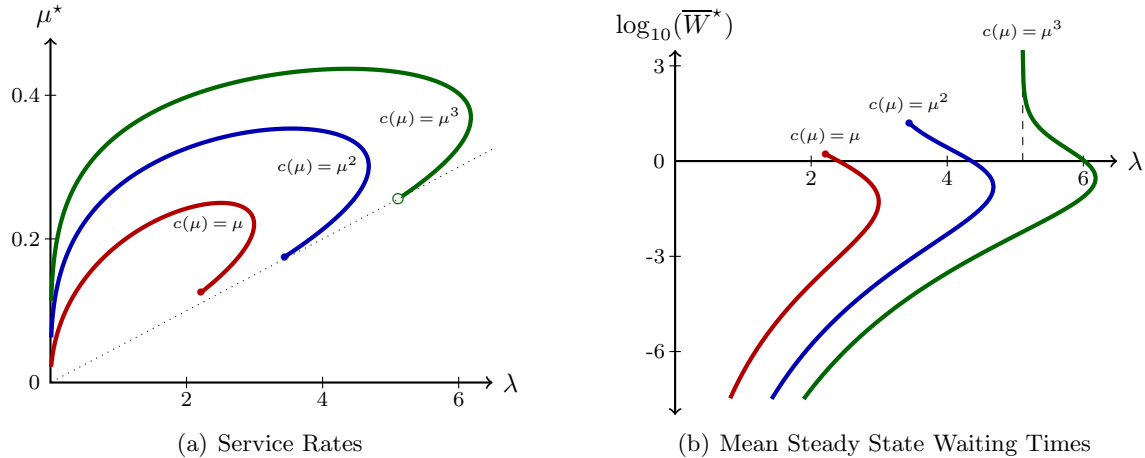


FIGURE 1. Equilibrium behavior as a function of the arrival rate when the staffing level is fixed at $N = 20$, for three different effort cost functions: linear, quadratic, and cubic. The dotted line in (a) is $\mu = \lambda/N = \lambda/20$.

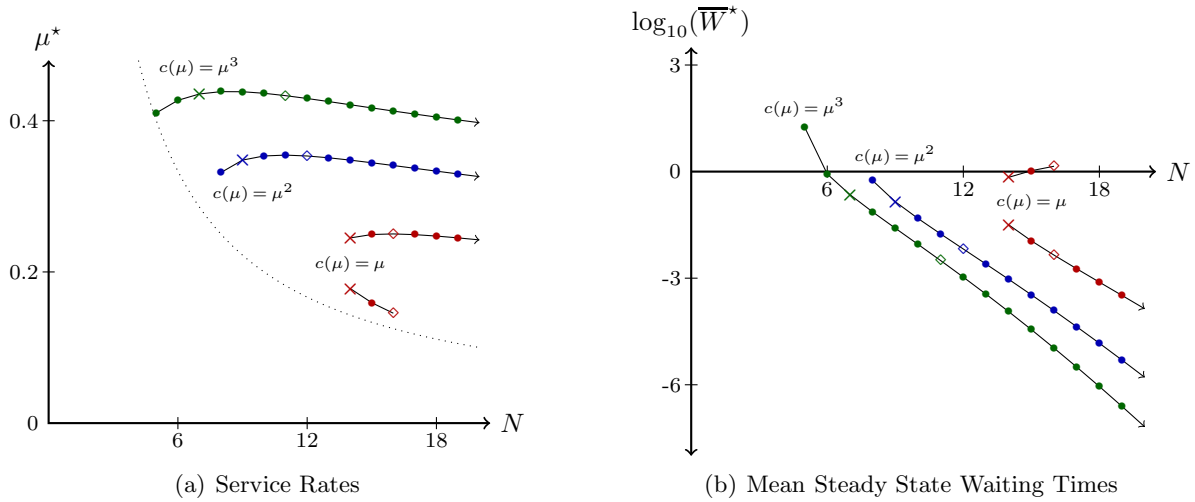


FIGURE 2. Equilibrium behavior as a function of the staffing level when the arrival rate is fixed at $\lambda = 2$, for three different effort cost functions: linear, quadratic, and cubic. The dotted curve in (a) is $\mu = \lambda/N = 2/N$. The data points marked \times and \diamond correspond to $N^{ao,2}$ and $N^{opt,2}$ respectively.

the corresponding mean steady state waiting time in Figure 1(b) steadily increases. In contrast, as the smaller equilibrium service rate increases, the corresponding mean steady state waiting time decreases. The relationship between the equilibrium service rates and waiting times is similarly inconsistent in Figure 2. This behavior is not consistent with results from classical, nonstrategic models, and could serve as a starting point to explaining empiric observations that are also not consistent with classical, nonstrategic models. For example, the non-monotonicity of service rate in workload is consistent with behavior observed in a hospital setting in [32].

4. Staffing strategic servers. One of the most studied questions for the design of service systems is staffing. Specifically, how many servers should be used for a given arrival rate. In the classical, nonstrategic setting, this question is well understood. In particular, as mentioned in the introduction, square-root staffing is known to be optimal when there are linear staffing and waiting costs [8].

In contrast, there is no previous work studying staffing in the context of strategic servers. *The goal of this section is to initiate the study of the impact that strategic servers have on staffing.* To get a feeling for the issues involved, consider a system with arrival rate λ and two possible staffing

policies: $N_1 = \lambda$ and $N_2 = 2\lambda$, where N_i is the number of servers staffed under policy i given arrival rate λ . Under N_1 , if the servers work at any rate slightly larger than 1, then they will have almost no idle time, and so they will have incentive to work harder. However, if servers are added, so that the provisioning is as in N_2 , then servers will have plentiful idle time when working at rate 1, and thus not have incentive to work harder. Thus, the staffing level has a fundamental impact on the incentives of the servers.

The above highlights that one should expect significant differences in staffing when strategic servers are considered. In particular, the key issue is that the staffing level itself creates incentives for the servers to speed up or slow down, because it influences the balance between effort and idle time. Thus, the policies that are optimal in the nonstrategic setting are likely suboptimal in the strategic setting, and vice versa.

The goal of the analysis in this section is to find the staffing level that minimizes costs when the system manager incurs linear staffing and waiting costs, within the class of policies that admit a symmetric equilibrium service rate. However, the analysis in the previous section highlights that determining the exact optimal policy is difficult, since we only have an implicit characterization of the symmetric equilibrium service rate in (9). As a result, we focus our attention on the setting where λ is large, and look for an asymptotically optimal policy.

As expected, the asymptotically optimal staffing policy we design for the case of strategic servers differs considerably from the optimal policies in the nonstrategic setting. In particular, in order for a symmetric equilibrium service rate to exist, the staffing level must be *order λ larger* than the optimal staffing in the classical, nonstrategic setting. Then, the system operates in a *quality-driven (QD)* regime instead of the *quality-and-efficiency-driven (QED)* regime that results from square-root staffing. This is intuitive given that the servers value their idle time, and in the QD regime they have idle time but in the QED regime their idle time is negligible.

The remainder of this section is organized as follows. We first introduce the cost structure and define asymptotic optimality in Section 4.1. Then, in Section 4.2, we provide a simple approximation of a symmetric equilibrium service rate and an asymptotically optimal staffing policy. Finally, in Section 4.3, we compare our asymptotically optimal staffing policy for strategic servers with the square-root staffing policy that is asymptotically optimal in the nonstrategic setting.

4.1. Preliminaries. Our focus in this section is on a $M/M/N$ queue with strategic servers, as introduced in Section 3. We assume Random routing throughout this section. It follows that our results hold for any “idle-time-order-based” routing policy (as explained in the beginning of Section 3.2 and validated by Theorem 9). The cost structure we assume is consistent with the one in [8], under which square-root staffing is asymptotically optimal when servers are not strategic. In their cost structure, there are linear staffing and waiting costs. One difference in our setting is that the equilibrium service rate may not be unique. In light of Theorem 6, we focus on the largest symmetric equilibrium service rate, and assume \bar{W}^* denotes the mean steady state waiting time in a $M/M/N$ queue with arrival rate λ , and strategic servers that serve at the largest symmetric equilibrium service rate (when there is more than one equilibrium).⁴ Then, the total system cost is

$$C^*(N, \lambda) = c_s N + \bar{w} \lambda \bar{W}^*.$$

The \star superscript indicates that the mean steady state waiting time, and hence, the cost function, depends on the (largest) symmetric equilibrium service rate μ^* , which in turn depends on N and λ .

⁴Note that the staffing policy we derive in this section (Theorem 8) will be asymptotically optimal regardless of which equilibrium service rate the servers choose.

The function $C^*(N, \lambda)$ is well-defined only if a symmetric equilibrium service rate, under which the system is stable, exists. Furthermore, we would like to rule out having an unboundedly large symmetric equilibrium service rate because then the server utility (1) will be large and negative – and it is hard to imagine servers wanting to participate in such a game.

DEFINITION 1. A staffing policy N^λ is *admissible* if the following two properties hold:

- (i) There exists a symmetric equilibrium $\mu^{*,\lambda}$ under which the system is stable ($\lambda < \mu^{*,\lambda} N^\lambda$) for all large enough λ .
 - (ii) There exists a sequence of symmetric equilibria $\{\mu^{*,\lambda}, \lambda > 0\}$ for which $\limsup_{\lambda \rightarrow \infty} \mu^{*,\lambda} < \infty$.
- If the requirement (ii) in the above definition is not satisfied, then the server utility will approach $-\infty$ as the service rates become unboundedly large. The servers will not want to participate in such a game. As long as the requirement (ii) is satisfied, we can assume the server payment is sufficient to ensure that the servers have positive utility.

We let Π denote the set of admissible staffing policies. We would like to solve for

$$N^{opt,\lambda} = \arg \min_{N \in \Pi} C^*(N, \lambda). \quad (11)$$

However, given the difficulty of deriving $N^{opt,\lambda}$ directly, we instead characterize the first order growth term of $N^{opt,\lambda}$ in terms of λ . To do this, we consider a sequence of systems, indexed by the arrival rate λ , and let λ become large.

Our convention when we wish to refer to any process or quantity associated with the system having arrival rate λ is to superscript the appropriate symbol by λ . In particular, N^λ denotes the staffing level in the system having arrival rate λ , and $\mu^{*,\lambda}$ denotes an equilibrium service rate (assuming existence) in the system with arrival rate λ and staffing level N^λ . We assume $\overline{W}^{*,\lambda}$ equals the mean steady state waiting time in a $M/M/N^\lambda$ queue with arrival rate λ when the servers work at the largest equilibrium service rate. The associated cost is

$$C^{*,\lambda}(N^\lambda) = c_S N^\lambda + \overline{w} \lambda \overline{W}^{*,\lambda}. \quad (12)$$

Given this setup, we would like to find an admissible staffing policy N^λ that has close to the minimum cost $C^{*,\lambda}(N^{opt,\lambda})$.

DEFINITION 2. A staffing policy N^λ is *asymptotically optimal* if it is admissible ($N^\lambda \in \Pi$) and

$$\lim_{\lambda \rightarrow \infty} \frac{C^{*,\lambda}(N^\lambda)}{C^{*,\lambda}(N^{opt,\lambda})} = 1.$$

In what follows, we use the o and ω notations to denote the limiting behavior of functions. Formally, for any two real-valued functions $f(x), g(x)$ that take nonzero values for sufficiently large x , we say that $f(x) = o(g(x))$ (equivalently, $g(x) = \omega(f(x))$) if $\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = 0$. In other words, f is dominated by g asymptotically (equivalently, g dominates f asymptotically).

4.2. An asymptotically optimal staffing policy. The class of policies we study are those that staff independently of the equilibrium service rates, which are endogenously determined according to the analysis in Section 3.2. More specifically, these are policies that choose N^λ purely as a function of λ . Initially, it is unclear what functional form an asymptotically optimal staffing policy can take in the strategic server setting. Thus, to begin, it is important to rule out policies that *cannot* be asymptotically optimal. The following proposition does this, and highlights that asymptotically optimal policies must be asymptotically linear in λ .

PROPOSITION 1. *Suppose $N^\lambda = f(\lambda) + o(f(\lambda))$ for some function f . If either $f(\lambda) = o(\lambda)$ or $f(\lambda) = \omega(\lambda)$, then the staffing policy N^λ cannot be asymptotically optimal.*

Intuitively, if $f(\lambda) = o(\lambda)$, understaffing forces the servers to work too hard, their service rates growing unboundedly (and hence their utilities approaching $-\infty$) as λ becomes large. On the other hand, the servers may prefer to have $f(\lambda) = \omega(\lambda)$ because the overstaffing allows them to be lazier; however, the overstaffing is too expensive for the system manager.

Proposition 1 implies that to find a staffing policy that is asymptotically optimal, we need only search within the class of policies that have the following form:

$$N^\lambda = \frac{1}{a}\lambda + o(\lambda), \text{ for } a \in (0, \infty). \quad (13)$$

However, before we can search for the cost-minimizing a , we must ensure that the staffing (13) guarantees the existence of a symmetric equilibria $\mu^{*,\lambda}$ for all large enough λ . It turns out that this is only true when a satisfies certain conditions. After providing these conditions (see Theorem 7 in the following), we evaluate the cost function as λ becomes large to find the a^* (defined in (17)) under which (13) is an asymptotically optimal staffing policy (see Theorem 8).

Equilibrium characterization. The challenge in characterizing equilibria comes from the complexity of the first order condition derived in Section 3. This complexity drives our focus on the large λ regime.

The first order condition for a symmetric equilibrium (9) is equivalently written as

$$\frac{\lambda}{N^\lambda} \left(\mu \left(1 + \frac{C(N^\lambda, \lambda/\mu)}{N^\lambda} \right) - \frac{\lambda}{N^\lambda} \right) = \mu^3 c'(\mu). \quad (14)$$

Under the staffing policy (13), when the limit $\lambda \rightarrow \infty$ is taken, this becomes

$$a(\mu - a) = \mu^3 c'(\mu). \quad (15)$$

Since $\mu^3 c'(\mu) > 0$, it follows that any solution μ has $\mu > a$. Therefore, under the optimistic assumption that a symmetric equilibrium solution $\mu^{*,\lambda}$ converging to the aforementioned solution μ exists, it follows that

$$\lambda/\mu^{*,\lambda} < \lambda/a$$

for all large enough λ . In words, the presence of strategic servers that value their idle time forces the system manager to staff order λ more servers than the offered load $\lambda/\mu^{*,\lambda}$. In particular, since the growth rate of N^λ is λ/a , *the system will operate in the quality-driven regime.*

The properties of the equation (15) are easier to see when it is rewritten as

$$\frac{1}{a} - \frac{1}{\mu} = \frac{\mu^2}{a^2} c'(\mu). \quad (16)$$

Note that the left-hand side of (16) is a concave function that increases from $-\infty$ to $1/a$ and the right-hand side is a convex function that increases from 0 to ∞ . These functions either cross at exactly two points, at exactly one point, or never intersect, depending on a . That information then can be used to show that the first order condition (14) has either two or zero solutions, depending on the value of a in the staffing policy (13).

THEOREM 7. *The following holds for all large enough λ .*

- (i) *Suppose $a > 0$ is such that there exists $\mu_2 > \mu_1 > 0$ that solve (16). Then, there exist exactly two solutions that solve (14).*
- (ii) *Suppose $a > 0$ is such that there exists exactly one $\mu_1 > 0$ that solves (16).*
 - (a) *Suppose $N^\lambda - \frac{\lambda}{a} \geq 0$. Then, there exists exactly two solutions that solve (14).*
 - (b) *Otherwise, if $N^\lambda - \frac{\lambda}{a} < -3$, then there does not exist a solution μ^λ to (14).*

Furthermore, for any $\epsilon > 0$, if μ^λ solves (14), then $|\mu^\lambda - \mu| < \epsilon$ for some μ that solves (16).

We are not sure if there is 0, 1, or 2 solutions in the case of $N^\lambda - \frac{1}{a}\lambda \in [-3, 0)$; however, given that we are focusing on a large λ asymptotic regime, the range $[-3, 0)$ is vanishingly small.

Moving forward, once the existence of a solution to the first order condition (14) is established, to conclude that solution is a symmetric equilibrium service rate also requires verifying the condition (10) in Theorem 4. This can be done for any staffing policy (13) under which the system operates in the quality driven-regime.

LEMMA 1. *For any staffing policy N^λ and associated μ^λ that satisfies the first order condition (14), if*

$$\liminf_{\lambda \rightarrow \infty} \frac{N^\lambda \mu^\lambda}{\lambda} = d > 1 \text{ and } \limsup_{\lambda \rightarrow \infty} \mu^\lambda < \infty,$$

then $\mu^{*,\lambda} = \mu^\lambda$ is a symmetric equilibrium for all large enough λ .

Under the conditions for the existence of a solution to the first order condition (14) in Theorem 7, it is also true that the conditions of Lemma 1 are satisfied. In particular, there exists a bounded sequence $\{\mu^\lambda\}$ having

$$\liminf_{\lambda \rightarrow \infty} \frac{N^\lambda \mu^\lambda}{\lambda} = \liminf_{\lambda \rightarrow \infty} \frac{\mu^\lambda}{a} + \mu^\lambda \frac{o(\lambda)}{\lambda} > 1.$$

This then guarantees that, for all large enough λ , there exists a solution $\mu^{*,\lambda}$ to (14) that is a symmetric equilibrium, under the conditions of Theorem 7.

There are either two symmetric equilibria for each λ or 0, because from Theorem 7 there are either two or zero solutions to the first order condition (14). These two symmetric equilibria will be close when there exists exactly one μ that solves (16); however, they may not be close when there exist two μ that solve (16). We show in the following that this does not affect what staffing policy should be asymptotically optimal.

Optimal staffing. Given the characterization of symmetric equilibria under a staffing policy (13), we can now move to the task of optimizing the staffing level, i.e., optimizing a . The first step is to characterize the associated cost, which is done in the following proposition.

PROPOSITION 2. *Suppose $a > 0$ is such that there exists $\mu > 0$ that solves (16). Then, under the staffing policy (13),*

$$\frac{C^{*,\lambda}(N^\lambda)}{\lambda} \rightarrow \frac{1}{a}c_S, \text{ as } \lambda \rightarrow \infty.$$

Proposition 2 implies that to minimize costs within the class of staffing policies that satisfy (13), the maximum a under which there exists at least one solution to (16) should be used. That is, we should choose a to be

$$a^* := \sup \mathcal{A}, \text{ where } \mathcal{A} := \{a > 0 : \text{there exists at least one solution } \mu > 0 \text{ to (16)}\}. \quad (17)$$

LEMMA 2. *$a^* \in \mathcal{A}$ is finite.*

Importantly, this a^* is not only optimal among the class of staffing policies that satisfy (13), it is asymptotically optimal among all admissible staffing policies. In particular, the following theorem shows that as λ becomes unboundedly large, no other admissible staffing policy can asymptotically achieve strictly lower cost than the one in (13) with $a = a^*$.

THEOREM 8. *If $N^{a^o,\lambda}$ satisfies (13) with $a = a^*$, then $N^{a^o,\lambda}$ is admissible and asymptotically optimal. Furthermore,*

$$\lim_{\lambda \rightarrow \infty} \frac{C^{*,\lambda}(N^{a^o,\lambda})}{\lambda} = \lim_{\lambda \rightarrow \infty} \frac{C^{*,\lambda}(N^{opt,\lambda})}{\lambda} = c_S \frac{1}{a^*}.$$

Note that an inspection of the proof of Theorem 8 shows that it holds regardless of which equilibrium service rate is used to define $\overline{W}^{\star,\lambda}$. Hence, even though we have defined $\overline{W}^{\star,\lambda}$ to be the mean steady state waiting time when the servers serve at the largest equilibrium service rate, this is not necessary. The staffing policy $N^{ao,\lambda}$ in Theorem 8 will be asymptotically optimal regardless of which equilibrium service rate the servers choose.

Though the above theorem characterizes an asymptotically optimal staffing level, because the definition of a^* is implicit, it is difficult to develop intuition. To highlight the structure more clearly, the following lemma characterizes a^* for a specific class of effort cost functions.

LEMMA 3. *Suppose $c(\mu) = c_E \mu^p$ for some $c_E \in [1, \infty)$ and $p \geq 1$. Then,*

$$a^* = \left[\frac{(p+1)}{(p+2)} \left(\frac{1}{c_E p(p+2)} \right)^{\frac{1}{p+1}} \right]^{(p+1)/p} < \mu^* = \left(\frac{p+1}{c_E p(p+2)^2} \right)^{\frac{1}{p}} < 1,$$

and a^* and μ^* are both increasing in p . Furthermore,

$$\text{if } a \begin{cases} < \\ > \\ = \end{cases} a^*, \text{ then } \begin{cases} \text{there are at least 2 non-negative solutions to (16)} \\ \text{there is no non-negative solution to (16)} \\ \text{there is exactly one solution to (16)} \end{cases}.$$

There are several interesting relationships between the effort cost function and the staffing level that follow from Lemma 3. First, for fixed p ,

$$a^*(p) \downarrow 0 \text{ as } c_E \rightarrow \infty.$$

In words, the system manager must staff more and more servers as effort becomes more costly. Second, for fixed c_E , since $a^*(p)$ is increasing in p , the system manager can staff less servers when the cost function becomes “more convex”. The lower staffing level forces the servers to work at a higher service rate since $\mu^*(p)$ is also increasing in p . We will revisit this idea that convexity is helpful to the system manager in the next section.

4.3. Contrasting staffing policies for strategic and nonstrategic servers. One of the most crucial observations that the previous section makes about the impact of strategic servers on staffing is that the strategic behavior leads the system to a quality-driven regime. In this section, we explore this issue in more detail, by comparing to the optimal staffing rule that arises when servers are not strategic, and then attempting to implement that staffing rule.

Nonstrategic servers. Recall that, for the conventional $M/M/N$ queue (without strategic servers), square-root staffing minimizes costs as λ becomes large (see equation (1), Proposition 6.2, and Example 6.3 in [8]). So, we can define

$$C_\mu^\lambda(N) = c_S N + \overline{w} \lambda \overline{W}_\mu^\lambda$$

to be the cost associated with staffing N nonstrategic servers that work at the fixed service rate μ . Further,

$$N_\mu^{opt,\lambda} = \arg \min_{N > \frac{\lambda}{\mu}} C_\mu^\lambda(N)$$

is the staffing level that minimizes expected cost when the system arrival rate is λ and the service rate is fixed to be μ . So, the staffing rule

$$N_\mu^{BMR,\lambda} = \frac{\lambda}{\mu} + y^* \sqrt{\frac{\lambda}{\mu}} \tag{18}$$

is asymptotically optimal in the sense that

$$\lim_{\lambda \rightarrow \infty} \frac{C_\mu^\lambda(N_\mu^{BMR,\lambda})}{C_\mu^\lambda(N_\mu^{opt,\lambda})} = 1.$$

Here, $y^* := \arg \min_{y>0} \left\{ c_S y + \frac{w\alpha(y)}{y} \right\}$, where $\alpha(y) = \left(1 + \frac{y}{h(-y)} \right)^{-1}$ with $h(\cdot)$ being the hazard rate function of the standard normal distribution, namely, $h(x) := \frac{\phi(x)}{1-\Phi(x)}$ with $\phi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$ and $\Phi(x) = \int_{-\infty}^x \phi(t) dt$. The staffing rule (18) is the famous square-root safety staffing rule.

Contrasting strategic and nonstrategic servers. In order to compare the case of strategic servers to the case of nonstrategic servers, it is natural to fix μ in (18) to the limiting service rate that results from using the optimal staffing rule $N^{ao,\lambda}$ defined in Theorem 8. We see that $N^{ao,\lambda}$ staffs order λ more servers than $N_{\mu^*}^{BMR,\lambda}$, where μ^* solves (16) for $a = a^*$, because any solution to (16) has $a > \mu$. When the effort cost function is $c(\mu) = c_E \mu^p$ for $p \geq 1$, we know from Lemma 3 and Theorem 7 (since the a^* is unique) that

$$\mu^{*,\lambda} \rightarrow \mu^* \text{ as } \lambda \rightarrow \infty,$$

where μ^* is as given in Lemma 3. Then, the difference in the staffing levels is

$$N^{ao,\lambda} - N_{\mu^*}^{BMR,\lambda} = \left(\frac{1}{a^*} - \frac{1}{\mu^*} \right) \lambda + o(\lambda) = \frac{1}{a^*} \left(\frac{1}{p+2} \right) \lambda + o(\lambda).$$

Since $a^* = a^*(p)$ is increasing in p from Lemma 3, we see that the difference $N^{ao,\lambda} - N_{\mu^*}^{BMR,\lambda}$ decreases to 0 as the cost function becomes “more convex”. This is consistent with our observation at the end of the previous subsection that convexity is helpful to the system manager.

It is natural to wonder if a system manager can force the servers to work harder by adopting the staffing policy suggested by the analysis of nonstrategic servers, i.e.,

$$N^{*,BMR,\lambda} = \frac{\lambda}{\mu^{*,\lambda}} + y^* \sqrt{\frac{\lambda}{\mu^{*,\lambda}}}. \quad (19)$$

The interpretation of this staffing rule requires care, because the offered load $\lambda/\mu^{*,\lambda}$ is itself a function of the staffing level (and the arrival rate) through the equilibrium service rate $\mu^{*,\lambda}$. The superscript \star emphasizes this dependence.

The first question concerns whether or not the staffing policy (19) is even possible in practice, because the staffing level depends on the equilibrium service rate and vice versa. More specifically, for a given staffing level, the servers relatively quickly arrive at an equilibrium service rate. Then, when system demand grows, the system manager increases the staffing, and the servers again arrive at an equilibrium service rate. In other words, there are two games, one played on a faster time scale (that is the servers settling to an equilibrium service rate), and one played on a slower time scale (that is the servers responding to added capacity).

To analyze the staffing policy (19), note that the first order condition for a symmetric equilibrium (9) is equivalently written as

$$\frac{\lambda/\mu}{(N^{*,BMR,\lambda})^2} \left(N^{*,BMR,\lambda} - \frac{\lambda}{\mu} + C \left(N^{*,BMR,\lambda}, \frac{\lambda}{\mu} \right) \right) = \mu c'(\mu).$$

Then, if μ^λ is a solution to the first order condition under the staffing $N^{*,BMR,\lambda}$ from (19), from substituting $N^{*,BMR,\lambda}$ into the above expression, μ^λ must satisfy

$$\frac{\lambda/\mu^\lambda}{\left(\lambda/\mu^\lambda + y^* \sqrt{\lambda/\mu^\lambda}\right)^2} \left(y^* \sqrt{\lambda/\mu^\lambda} + C \left(\lambda/\mu^\lambda + y^* \sqrt{\lambda/\mu^\lambda}, \frac{\lambda}{\mu^\lambda} \right) \right) = \mu^\lambda c'(\mu^\lambda).$$

As λ becomes large, since $C \left(\lambda/\mu + y \sqrt{\lambda/\mu}, \frac{\lambda}{\mu} \right)$ is bounded above by 1, the left-hand side of the above expression has limit 0. Furthermore, the right-hand side of the above equation is non-negative and increasing as a function of μ . Hence any sequence of solutions μ^λ to the first order condition has the limiting behavior

$$\mu^\lambda \rightarrow 0, \text{ as } \lambda \rightarrow \infty,$$

which cannot be a symmetric equilibrium service rate because we require the servers to work fast enough to stabilize the system.

One possibility is to expand the definition of an equilibrium service rate in (1) to allow the servers to work exactly at the lower bound λ/N . In fact, the system manager may now be tempted to push the servers to work even faster. However, faster service cannot be mandated for free – there must be a trade-off; for example, the service quality may suffer or the salaries should be higher.

4.4. Numerical examples. In order to understand how well our asymptotically optimal staffing policy $N^{ao,\lambda}$ performs in comparison with the optimal policy $N^{opt,\lambda}$ for finite λ , and how fast the corresponding system cost converges to the optimal cost, we present some results from numerical analysis in this section.

We consider two staffing policies: (i) $N^{opt,\lambda}$ (defined in (11)), and (ii) $N^{ao,\lambda}$ (defined in Theorem 8 and (17)) where we ignore the $o(\lambda)$ term of (13). For each, we first round up the staffing level if necessary, and then plot the following two equilibrium quantities as a function of the arrival rate λ : (a) service rates $\mu^{*,\lambda}$ (if there is more than one, we pick the largest), and (b) normalized costs $C^{*,\lambda}/\lambda$. We calculate $N^{opt,\lambda}$ numerically, by iterating over the staffing levels that admit equilibria (and we choose the lowest cost when there are multiple equilibria). These plots are shown in Figure 3 for three different effort cost functions: $c(\mu) = \mu$, $c(\mu) = \mu^2$, and $c(\mu) = \mu^3$, which are depicted in red, blue, and green respectively. For each color, the curve with the darker shade corresponds to $N^{opt,\lambda}$ and the curve with the lighter shade corresponds to $N^{ao,\lambda}$. The horizontal dashed lines correspond to the limiting values as $\lambda \rightarrow \infty$.

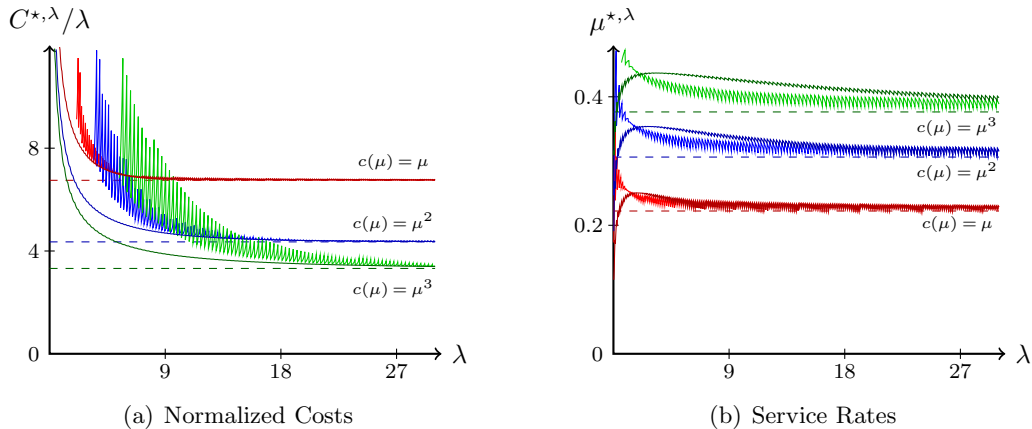


FIGURE 3. Equilibrium behavior as a function of the arrival rate for the optimal and asymptotically optimal staffing policies, for three different effort cost functions: linear, quadratic, and cubic.

An immediate first observation is the jaggedness of the curves, which is a direct result of the discreteness of the staffing levels $N^{opt,\lambda}$ and $N^{ao,\lambda}$. In particular, as the arrival rate λ increases, the equilibrium service rate $\mu^{*,\lambda}$ decreases (respectively, the equilibrium normalized cost $C^{*,\lambda}/\lambda$ increases) smoothly until the staffing policy adds an additional server, which causes a sharp increase (respectively, decrease). The jaggedness is especially pronounced for smaller λ , resulting in a complex pre-limit behavior that necessitates asymptotic analysis in order to obtain analytic results.

However, despite the jaggedness, the plots illustrate clearly that both the equilibrium service rates and normalized costs of the optimal policy $N^{ao,\lambda}$ converge quickly to those of the optimal policy $N^{opt,\lambda}$, highlighting that our asymptotic results are predictive at realistically sized systems.

5. Routing to strategic servers. Thus far we have focused our discussion on staffing, assuming that jobs are routed randomly to servers when there is a choice. Of course, the decision of how to route jobs to servers is another crucial aspect of the design of service systems. As such, the analysis of routing policies has received considerable attention in the queueing literature, when servers are not strategic. In this section, we begin to investigate the impact of strategic servers on the design of routing policies.

In the classical literature studying routing when servers are nonstrategic, a wide variety of policies have been considered. These include “rate-based policies” such as Fastest Server First (FSF) and Slowest Server First (SSF); as well as “idle-time-order-based policies” such as Longest Idle Server First (LISF) and Shortest Idle Server First (SISF). Among these routing policies, FSF is a natural choice to minimize the mean response time (although, as noted in the Introduction, it is not optimal in general). This leads to the question: how does FSF perform when servers are strategic? In particular, does it perform better than the Random routing that we have so far studied?

Before studying optimal routing to improve performance, we must first answer the following even more fundamental question: what routing policies admit symmetric equilibria? This is a very challenging goal, as can be seen by the complexity of the analysis for the $M/M/N$ under Random routing. This section provides a first step towards that goal.

The results in this section focus on two broad classes of routing policies *idle-time-order-based policies* and *rate-based policies*, which are introduced in turn in the following.

5.1. Idle-time-order-based policies. Informally, idle-time-order-based policies are those routing policies that use only the rank ordering of when servers last became idle in order to determine how to route incoming jobs. To describe the class of idle-time-order-based policies precisely, let $\mathcal{I}(t)$ be the set of servers idle at time $t > 0$, and, when $\mathcal{I}(t) \neq \emptyset$, let $\mathbf{s}(t) = (s_1, \dots, s_{|\mathcal{I}(t)|})$ denote the ordered vector of idle servers at time t , where server s_j became idle before server s_k whenever $j < k$. For $n \geq 1$, let $\mathcal{P}_n = \Delta(\{1, \dots, n\})$ denote the set of all probability distributions over the set $\{1, \dots, n\}$. An idle-time-order-based routing policy is defined by a collection of probability distributions $\mathbf{p} = \{p^S\}_{S \in 2^{\{1,2,\dots,N\}} \setminus \emptyset}$, such that $p^S \in \mathcal{P}_{|S|}$, for all $S \in 2^{\{1,2,\dots,N\}} \setminus \emptyset$. Under this policy, at time t , the next job in queue is assigned to idle server s_j with probability $p^{\mathcal{I}(t)}(j)$. Examples of idle-time-order-based routing policies are as follows.

1. *Random.* An arriving customer that finds more than one server idle is equally likely to be routed to any of those servers. Then, $p^S = (1/|S|, \dots, 1/|S|)$ for all $S \in 2^{\{1,2,\dots,N\}} \setminus \emptyset$.
2. *Weighted Random.* Each such arriving customer is routed to one of the idle servers with probabilities that may depend on the order in which the servers became idle. For example, if

$$p^S(j) = \frac{|S| + 1 - j}{\sum_{n=1}^{|S|} n}, \quad j \in S, \text{ for } s_j \in S, \text{ for all } S \in 2^{\{1,2,\dots,N\}} \setminus \emptyset,$$

then the probabilities are decreasing according to the order in which the servers became idle. Note that $\sum_j p^S(j) = \frac{|S|(|S|+1) - \frac{1}{2}|S|(|S|+1)}{\frac{1}{2}|S|(|S|+1)} = 1$.

3. *Longest Idle Server First (Shortest Idle Server First)*. Each such arriving customer is routed to the server that has idled the longest (idled the shortest). Then, $p^S = (1, 0, \dots, 0)$ ($p^S = (0, \dots, 0, 1)$) for all $S \subseteq \{1, 2, \dots, N\}$.

5.1.1. Policy-space collapse. Surprisingly, it turns out that all idle-time-order-based policies are “equivalent” in a very strong sense — they all lead to the same steady state probabilities, resulting in a remarkable *policy-space collapse* result, which we discuss in the following.

Fix R to be some idle-time-order-based routing policy, defined through the collection of probability distributions $\mathbf{p} = \{p^S\}_{\emptyset \neq S \subseteq \{1, 2, \dots, N\}}$. The states of the associated continuous time Markov chain are defined as follows:

- State B is the state where all servers are busy, but there are no jobs waiting in the queue.
- State $\mathbf{s} = (s_1, s_2, \dots, s_{|\mathcal{I}|})$ is the ordered vector of idle servers \mathcal{I} . When $\mathcal{I} = \emptyset$, we identify the empty vector \mathbf{s} with state B .
- State m ($m \geq 0$) is the state where all servers are busy and there are m jobs waiting in the queue (i.e., there are $N + m$ jobs in the system). We identify state 0 with state B .

When all servers are busy, there is no routing, and so the system behaves exactly as a $M/M/1$ queue with arrival rate λ and service rate $\mu_1 + \dots + \mu_N$. Then, from the local balance equations, the associated steady state probabilities π_B and π_m for $m = 0, 1, 2, \dots$, must satisfy

$$\pi_m = (\lambda/\mu)^m \pi_B \text{ where } \mu = \sum_{j=1}^N \mu_j. \quad (20)$$

One can anticipate that the remaining steady state probabilities satisfy

$$\pi_{\mathbf{s}} = \pi_B \prod_{s \in \mathcal{I}} \frac{\mu_s}{\lambda} \quad \text{for all } \mathbf{s} = (s_1, s_2, \dots, s_{|\mathcal{I}|}) \text{ with } |\mathcal{I}| > 0, \quad (21)$$

and the following theorem verifies this by establishing that the detailed balance equations are satisfied.

THEOREM 9. *All idle-time-order-based policies have the steady state probabilities that are uniquely determined by (20)-(21), together with the normalization constraint that their sum is one.*

Theorem 9 is remarkable because there is no dependence on the collection of probability distributions \mathbf{p} that define R . Therefore, it follows that all idle-time-order-based routing policies result in the same steady state probabilities. Note that, concurrently, a similar result has been discovered independently in the context of loss systems [24].

In relation to our server game, it follows from Theorem 9 that all idle-time-order-based policies have the same equilibrium behavior as Random. This is because an equilibrium service rate depends on the routing policy through the server idle time vector $(I_1(\boldsymbol{\mu}; R), \dots, I_N(\boldsymbol{\mu}; R))$, which can be found from the steady state probabilities in (20)-(21). As a consequence, if there exists (does not exist) an equilibrium service rate under Random, then there exists (does not exist) an equilibrium service rate under any idle-time-order-based policy. In summary, it is not possible to achieve better performance than under Random by employing any idle-time-order-based policy.

5.2. Rate-based policies. Informally, a rate-based policy is one that makes routing decisions using only information about the rates of the servers. As before, let $\mathcal{I}(t)$ denote the set of idle servers at time t . In a rate-based routing policy, jobs are assigned to idle servers only based on their service rates. We consider a parameterized class of rate-based routing policies that we term

r -routing policies ($r \in \mathbb{R}$). Under an r -routing policy, at time t , the next job in queue is assigned to idle server $i \in \mathcal{I}(t)$ with probability

$$p_i(\boldsymbol{\mu}, t; r) = \frac{\mu_i^r}{\sum_{j \in \mathcal{I}(t)} \mu_j^r}$$

Notice that for special values of the parameter r , we recover well-known policies. For example, setting $r = 0$ results in Random; as $r \rightarrow \infty$, it approaches FSF; and as $r \rightarrow -\infty$, it approaches SSF.

In order to understand the performance of rate-based policies, the first step is to perform an equilibrium analysis, i.e., we need to understand what the steady state idle times look like under any r -routing policy. The following proposition provides us with the required expressions.

PROPOSITION 3. *Consider a heterogeneous $M/M/2$ system under an r -routing policy, with arrival rate $\lambda > 0$ and servers 1 and 2 operating at rates μ_1 and μ_2 respectively. The steady state probability that server 1 is idle is given by:*

$$I_1^r(\mu_1, \mu_2) = \frac{\mu_1(\mu_1 + \mu_2 - \lambda) \left[(\lambda + \mu_2)^2 + \mu_1\mu_2 + \frac{\mu_2^r}{\mu_1^r + \mu_2^r}(\lambda\mu_1 + \lambda\mu_2) \right]}{\mu_1\mu_2(\mu_1 + \mu_2)^2 + (\lambda\mu_1 + \lambda\mu_2) \left[\mu_1^2 + 2\mu_1\mu_2 - \frac{\mu_1^r}{\mu_1^r + \mu_2^r}(\mu_1^2 - \mu_2^2) \right] + (\lambda\mu_1)^2 + (\lambda\mu_2)^2},$$

and the steady state probability that server 2 is idle is given by $I_2^r(\mu_1, \mu_2) = I_1^r(\mu_2, \mu_1)$.

Note that we restrict ourselves to a 2-server system for this analysis. This is due to the fact that there are no closed form expressions known for the resulting Markov chains for systems with more than 3 servers. It may be possible to extend these results to 3 servers using results from [37]; but, the expressions are intimidating, to say the least. However, the analysis for two servers is already enough to highlight important structure about the impact of strategic servers on policy design.

In particular, our first result concerns the FSF and SSF routing policies, which can be obtained in the limit when $r \rightarrow \infty$ and $r \rightarrow -\infty$ respectively. Recall that FSF is asymptotically optimal in the nonstrategic setting. Intuitively, however, it penalizes the servers that work the fastest by sending them more and more jobs. In a strategic setting, this might incentivize servers to decrease their service rate, which is not good for the performance of the system. One may wonder if by doing the opposite, that is, using the SSF policy, servers can be incentivized to increase their service rate. However, our next theorem (Theorem 10) shows that neither of these policies is useful if we are interested in symmetric equilibria.

Recall that our model for strategic servers already assumes an increasing, convex effort cost function with $c'''(\mu) \geq 0$. For the rest of this section, in addition, we assume that $c'(\frac{\lambda}{2}) < \frac{1}{\lambda}$. (Recall that this is identical to the sufficient condition $c'(\frac{\lambda}{N}) < \frac{1}{\lambda}$ which we introduced in Section 3.2, on substituting $N = 2$.)⁵

THEOREM 10. *Consider a $M/M/2$ queue with strategic servers. Then, FSF and SSF do not admit a symmetric equilibrium.*

Moving beyond FSF and SSF, we continue our equilibrium analysis (for a finite r) by using the first order conditions to show that whenever an r -routing policy admits a symmetric equilibrium, it is unique. Furthermore, we provide an expression for the corresponding symmetric equilibrium service rate in terms of r , which brings out a useful monotonicity property.

⁵ The sufficient condition $c'(\frac{\lambda}{2}) < \frac{1}{\lambda}$ might seem rather strong, but it can be shown that it is necessary for the symmetric first order condition to have a unique solution. This is because, if $c'(\frac{\lambda}{2}) > \frac{1}{\lambda}$, then the function $\varphi(\mu)$, defined in (22), ceases to be monotonic, and as a result, for any given r , the first order condition $\varphi(\mu) = r$ could have more than one solution.

THEOREM 11. *Consider a $M/M/2$ queue with strategic servers. Then, any r -routing policy that admits a symmetric equilibrium, admits a unique symmetric equilibrium, given by $\mu^* = \varphi^{-1}(r)$, where $\varphi: (\frac{\lambda}{2}, \infty) \rightarrow \mathbb{R}$ is the function defined by*

$$\varphi(\mu) = \frac{4(\lambda + \mu)}{\lambda(\lambda - 2\mu)} (\mu(\lambda + 2\mu)c'(\mu) - \lambda). \quad (22)$$

Furthermore, among all such policies, μ^ is decreasing in r , and therefore, $\mathbb{E}[T]$, the mean response time (a.k.a. sojourn time) at symmetric equilibrium is increasing in r .*

In light of the inverse relationship between r and μ^* that is established by this theorem, the system manager would ideally choose the smallest r such that the corresponding r -routing policy admits a symmetric equilibrium, which is in line with the intuition that a bias towards SSF (the limiting r -routing policy as $r \rightarrow -\infty$) incentivizes servers to work harder. However, there is a hard limit on how small an r can be chosen (concurrently, how large an equilibrium service rate μ^* can be achieved) so that there exists a symmetric equilibrium, as evidenced by our next theorem.

THEOREM 12. *Consider a $M/M/2$ queue with strategic servers. Then, there exists $\bar{\mu}, \underline{r} \in \mathbb{R}$, with $\underline{r} = \varphi(\bar{\mu})$, such that no service rate $\mu > \bar{\mu}$ can be a symmetric equilibrium under any r -routing policy, and no r -routing policy with $r < \underline{r}$ admits a symmetric equilibrium.*

The proof of this theorem is constructive and we do exhibit an \underline{r} , however, it is not clear whether this is tight, that is, whether there exists a symmetric equilibrium for all r -routing policies with $r \geq \underline{r}$. We provide a partial answer to this question of what r -routing policies do admit symmetric equilibria in the following theorem.

THEOREM 13. *Consider a $M/M/2$ queue with strategic servers. Then, there exists a unique symmetric equilibrium under any r -routing policy with $r \in \{-2, -1, 0, 1\}$.*

Notice that we show equilibrium existence for four integral values of r . It is challenging to show that all r -routing policies in the interval $[-2, 1]$ admit a symmetric equilibrium. This theorem provides an upper bound on the \underline{r} of the previous theorem, that is, $\underline{r} \leq -2$. Therefore, if the specific cost function c is unknown, then the system manager can guarantee better performance than Random ($r = 0$), by setting $r = -2$. If the specific cost function is known, the system manager may be able to employ a lower r to obtain even better performance. For example, consider a 2-server system with $\lambda = 1/4$ and one of three different effort cost functions: $c(\mu) = \mu$, $c(\mu) = \mu^2$, and $c(\mu) = \mu^3$. Figure 4 shows the corresponding equilibrium mean response times (in red, blue, and green, respectively). It is worth noting that the more convex the effort cost function, larger the range of r (and smaller the minimum value of r) for which a symmetric equilibrium exists.

6. Concluding remarks. The rate at which each server works in a service system has important consequences for service system design. However, traditional models of large service systems do not capture the fact that human servers respond to incentives created by scheduling and staffing policies, because traditional models assume each server works at a given fixed service rate. In this paper, we initiate the study of a class of strategic servers that seek to optimize a utility function which values idle time and includes an effort cost.

Our focus is on the analysis of staffing and routing policies for a $M/M/N$ queue with strategic servers, and our results highlight that strategic servers have a dramatic impact on the optimal policies in both cases. In particular, policies that are optimal in the classical, nonstrategic setting can perform quite poorly when servers act strategically.

For example, a consequence of the strategic server behavior is that the cost-minimizing staffing level is order λ larger than square-root staffing, the cost minimizing staffing level for systems with fixed service rate. In particular, any system with strategic servers operates in the quality-driven

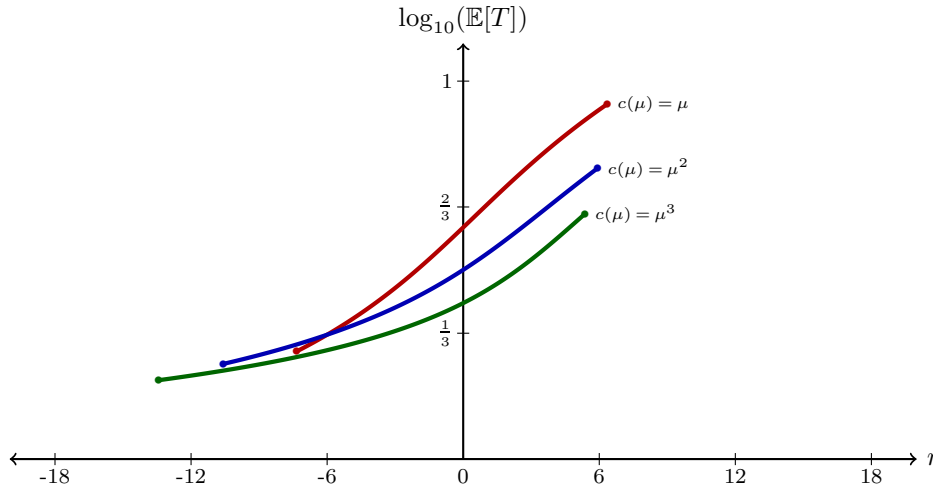


FIGURE 4. Equilibrium mean response time (a.k.a. sojourn time) as a function of the policy parameter, r , when the arrival rate is $\lambda = \frac{1}{4}$, for three different effort cost functions: linear, quadratic, and cubic.

regime at equilibrium (as opposed to the quality-and-efficiency-driven regime that arises under square-root staffing), in which the servers all enjoy non-negligible idle time.

The intuitive reason square-root staffing is not feasible in the context of strategic servers is that the servers do not value their idleness enough in comparison to their effort cost. This causes the servers to work too slowly, making idle time scarce. In the economics literature [9, 36], it is common to assume that scarce goods are more highly valued. If we assume that the servers valued their idle time more heavily as the idle time becomes scarcer, then the servers would work faster in order to make sure they achieved some. This suggests the following interesting direction for future research: what is the relationship between the assumed value of idle time in (2) and the resulting cost minimizing staffing policy? Another situation in which servers may not care about idle time becoming scarce is when their compensation depends on their service volume (which is increasing in their service rate). Then, it is reasonable to expect the servers prefer to have negligible idle time. It would be interesting to be able to identify a class of compensation schemes under which that is the case.

The aforementioned two future research directions become even more interesting when the class of routing policies is expanded to include rate-based policies. This paper solves the *joint* routing and staffing problem within the class of idle-time-order-based policies. Section 5 suggests that by expanding the class of routing policies to also include rate-based policies we should be able to achieve better system performance (although it is clear that the analysis becomes much more difficult). The richer question also aspires to understand the relationship between the server idle time value, the compensation scheme, the (potentially) rate-based routing policy, and the number of strategic servers to staff.

Finally, it is important to note that we have focused on symmetric equilibrium service rates. We have not proven that asymmetric equilibria do not exist. Thus, it is natural to wonder if there are routing and staffing policies that result in an asymmetric equilibrium. Potentially, there could be one group of servers that have low effort costs but negligible idle time and another group of servers that enjoy plentiful idle time but have high effort costs. The question of asymmetric equilibria becomes even more interesting when the servers have different utility functions. For example, more experienced servers likely have lower effort costs than new hires. Also, different servers can value their idle time differently. How do we design routing and staffing policies that are respectful of such considerations?

Acknowledgments. This work was supported by NSF grant #CCF-1101470, AFOSR grant #FA9550-12-1-0359, and ONR grant #N00014-09-1-0751.

References

- [1] Aksin, Z., M. Armony, V. Mehrotra. 2007. The modern call-center: A multi-disciplinary perspective on operations management research. *Prod. Oper. Manag.* **16**(6) 665–688.
- [2] Allon, G., I. Gurvich. 2010. Pricing and dimensioning competing large-scale service providers. *M&SOM* **12**(3) 449–469.
- [3] Anton, J. 2005. One-minute survey report #488: Agent compensation & advancement. Document Tracking Number SRV488-080305.
- [4] Armony, M. 2005. Dynamic routing in large-scale service systems with heterogeneous servers. *Queueing Syst. Theory Appl.* **51**(3-4) 287–329.
- [5] Armony, M., A. R. Ward. 2010. Fair dynamic routing in large-scale heterogeneous-server systems. *Oper. Res.* **58** 624–637.
- [6] Atar, R. 2005. Scheduling control for queueing systems with many servers: Asymptotic optimality in heavy traffic. *Ann. Appl. Probab.* **15**(4) 2606–2650.
- [7] Atar, R., Y. Y. Shaki, A. Shwartz. 2011. A blind policy for equalizing cumulative idleness. *Queueing Syst.* **67**(4) 275–293.
- [8] Borst, S., A. Mandelbaum, M. I. Reiman. 2004. Dimensioning large call centers. *Oper. Res.* **52**(1) 17–34.
- [9] Brock, T. C. 1968. Implications of commodity theory for value change. *Psychological Foundations of Attitudes* 243–275.
- [10] Cachon, G. P., P. T. Harker. 2002. Competition and outsourcing with scale economies. *Manage. Sci.* **48**(10) 1314–1333.
- [11] Cachon, G. P., F. Zhang. 2007. Obtaining fast service in a queueing system via performance-based allocation of demand. *Manage. Sci.* **53**(3) 408–420.
- [12] Cahuc, P., A. Zylberberg. 2004. *Labor Economics*. MIT Press.
- [13] Cheng, S.-F., D. M. Reeves, Y. Vorobeychik, M. P. Wellman. 2004. Notes on equilibria in symmetric games. *International Workshop On Game Theoretic And Decision Theoretic Agents (GTDT)*. 71–78.
- [14] Cohen-Charash, Y., P. E. Spector. 2001. The role of justice in organizations: A meta-analysis. *Organ. Behav. and Hum. Dec.* **86**(2) 278–321.
- [15] Colquitt, J. A., D. E. Conlon, M. J. Wesson, C. O. L. H. Porter, K. Y. Ng. 2001. Justice at the millennium: A meta-analytic review of 25 years of organizational justice research. *J. Appl. Psychol.* **86**(3) 425–445.
- [16] de Véricourt, F., Y.-P. Zhou. 2005. Managing response time in a call-routing problem with service failure. *Oper. Res.* **53**(6) 968–981.
- [17] Erlang, A. K. 1948. On the rational determination of the number of circuits. E. Brockmeyer, H. L. Halstrom, A. Jensen, eds., *The Life and Works of A. K. Erlang*. The Copenhagen Telephone Company, 216–221.
- [18] Gans, N., G. Koole, A. Mandelbaum. 2003. Telephone call centers: Tutorial, review, and research prospects. *M&SOM* **5**(2) 79–141.
- [19] Garnett, O., A. Mandelbaum, M. Reiman. 2002. Designing a call center with impatient customers. *M&SOM* **4**(3) 208–227.
- [20] Geng, X., W. T. Huh, M. Nagarajan. 2013. Strategic and fair routing policies in a decentralized service system. Working paper.
- [21] Gilbert, S. M., Z. K. Weng. 1998. Incentive effects favor nonconsolidating queues in a service system: The principal-agent perspective. *Manage. Sci.* **44**(12) 1662–1669.
- [22] Gumbel, H. 1960. Waiting lines with heterogeneous servers. *Oper. Res.* **8**(4) 504–511.
- [23] Gurvich, I., W. Whitt. 2007. Scheduling flexible servers with convex delay costs in many-server service systems. *M&SOM* **11**(2) 237–253.

- [24] Haji, B., S. M. Ross. 2013. A queueing loss model with heterogenous skill based servers under idle time ordering policies. Working paper.
- [25] Halfin, S., W. Whitt. 1981. Heavy-traffic limits for queues with many exponential servers. *Oper. Res.* **29**(3) 567–588.
- [26] Harchol-Balter, M. 2013. *Performance Modeling and Design of Computer Systems: Queueing Theory in Action*. Cambridge University Press.
- [27] Harel, A. 1988. Sharp bounds and simple approximations for the Erlang delay and loss formulas. *Manage. Sci.* **34**(8) 959–972.
- [28] Hassin, R., M. Haviv. 2003. *To Queue or Not to Queue: Equilibrium Behavior in Queueing Systems*. Kluwer.
- [29] Hopp, W., W. Lovejoy. 2013. *Hospital Operations: Principles of High Efficiency Health Care*. Financial Times Press.
- [30] Janssen, A. J. E. M., J. S.H. van Leeuwen, B. Zwart. 2011. Refining square-root safety staffing by expanding Erlang C. *Oper. Res.* **59**(6) 1512–1522.
- [31] Kalai, E., M. I. Kamien, M. Rubinovitch. 1992. Optimal service speeds in a competitive environment. *Manage. Sci.* **38**(8) 1154–1163.
- [32] Kc, D. S., C. Terwiesch. 2009. Impact of workload on service time and patient safety: An econometric analysis of hospital operations. *Manage. Sci.* **55**(9) 1486–1498.
- [33] Kocaga, L., M. Armony, A. R. Ward. 2013. Staffing call centers with uncertain arrival rates and co-sourcing. Working paper.
- [34] Krishnamoorthi, B. 1963. On Poisson queue with two heterogeneous servers. *Oper. Res.* **11**(3) 321–330.
- [35] Lin, W., P. Kumar. 1984. Optimal control of a queueing system with two heterogeneous servers. *IEEE Trans. Autom. Contr.* **29**(8) 696–703.
- [36] Lynn, M. 1991. Scarcity effects on value: A quantitative review of the commodity theory literature. *Psychol. Market.* **8**(1) 43–57.
- [37] Mokaddis, G. S., C. H. Matta, M. M. El Genaidy. 1998. On Poisson queue with three heterogeneous servers. *International Journal of Information and Management Sciences* **9** 53–60.
- [38] Reed, J., Y. Shaki. 2013. A fair policy for the G/GI/N queue with multiple server pools. Preprint.
- [39] Saaty, T. L. 1960. Time-dependent solution of the many-server Poisson queue. *Oper. Res.* **8**(6) 755–772.
- [40] Tezcan, T. 2008. Optimal control of distributed parallel server systems under the Halfin and Whitt regime. *Math. Oper. Res.* **33** 51–90.
- [41] Tezcan, T., J. Dai. 2010. Dynamic control of N-systems with many servers: Asymptotic optimality of a static priority policy in heavy traffic. *Oper. Res.* **58**(1) 94–110.
- [42] Ward, A. R., M. Armony. 2013. Blind fair routing in large-scale service systems with heterogeneous customers and servers. *Oper. Res.* **61** 228–243.
- [43] Whitt, W. 2002. IEOR 6707: Advanced Topics in Queueing Theory: Focus on Customer Contact Centers. Homework 1e Solutions, see <http://www.columbia.edu/~ww2040/ErlangBandCFormulas.pdf>.

Routing and Staffing when Servers are Strategic: Technical Appendix

Ragavendran Gopalakrishnan, Sherwin Doroudi, Amy R. Ward, and Adam Wierman

In this technical appendix, we provide proofs for the results stated in the main body of the manuscript titled: ‘‘Routing and Staffing when Servers are Strategic’’. The proofs of these results are in the order in which they appear in the main body.

PROOFS FROM SECTION 3

Proof of Theorem 1. The starting point of this proof is the expression for the steady state probabilities of a general heterogeneous $M/M/N$ system with Random routing, which was derived in [22]. Before stating this more general result, we first set up the required notation. Let $\mu_1, \mu_2, \dots, \mu_N$ denote the service rates of the N servers, and let $\rho_j = \frac{\lambda}{\mu_j}$, $1 \leq j \leq N$. We assume that $\sum_{j=1}^N \rho_j^{-1} > 1$ for stability. Let (a_1, a_2, \dots, a_k) denote the state of the system when there are k jobs in the system ($0 < k < N$) and the busy servers are $\{a_1, a_2, \dots, a_k\}$, where $1 \leq a_1 < a_2 < \dots < a_k \leq N$. Let $P(a_1, a_2, \dots, a_k)$ denote the steady state probability of the system being in state (a_1, a_2, \dots, a_k) . Also, let P_k denote the steady state probability of k jobs in the system. Then,

$$P(a_1, a_2, \dots, a_k) = \frac{(N-k)! P_0 \rho_{a_1} \rho_{a_2} \cdots \rho_{a_k}}{N!}, \quad (\text{EC.1})$$

where P_0 , the steady state probability that the system is empty, is given by:

$$P_0 = \frac{N! C_N^N}{D_N}, \quad (\text{EC.2})$$

where, for $1 \leq j \leq N$,

$$\begin{aligned} C_j^N &= \text{sum of combinations of } j \rho_i^{-1} \text{ values from } N \rho_i^{-1} \text{ values} \\ &= \sum_{a_1=1}^{N-j+1} \sum_{a_2=a_1+1}^{N-j+2} \cdots \sum_{a_{j-1}=a_{j-2}+1}^{N-j+j-1} \sum_{a_j=a_{j-1}+1}^N \rho_{a_1}^{-1} \rho_{a_2}^{-1} \cdots \rho_{a_j}^{-1}, \end{aligned} \quad (\text{EC.3})$$

and

$$D_N = \sum_{j=1}^N j! C_j^N + \frac{C_1^N}{C_1^N - 1}. \quad (\text{EC.4})$$

Note that,

$$C_N^N = \prod_{i=1}^N \rho_i^{-1} \quad \text{and} \quad C_1^N = \sum_{i=1}^N \rho_i^{-1}.$$

Also, by convention, we write $C_0^N = 1$. The steady state probability that a tagged server, say server 1, is idle is obtained by summing up the steady state probabilities of every state in which server 1 is idle:

$$I(\mu_1, \mu_2, \dots, \mu_N; \lambda, N) = P_0 + \sum_{k=1}^{N-1} \sum_{2 \leq a_1 \leq \dots \leq a_k \leq N} P(a_1, a_2, \dots, a_k) \quad (\text{EC.5})$$

We now simplify the expressions above for our special system where the tagged server works at a rate μ_1 and all other servers work at rate μ . Without loss of generality, we pick server 1 to be the

tagged server, and we set $\mu_2 = \mu_3 = \dots = \mu_N = \mu$, and therefore, $\rho_2 = \rho_3 = \dots = \rho_N = \rho = \frac{\lambda}{\mu}$. Then, (EC.1) simplifies to:

$$P(a_1, a_2, \dots, a_k) = \frac{(N-k)! P_0 \rho^k}{N!}, 2 \leq a_1 \leq \dots \leq a_k \leq N \quad (\text{EC.6})$$

In order to simplify (EC.3), we observe that

$$C_j^N = \rho_1^{-1} C_{j-1}^{N-1} + C_j^{N-1}$$

where the terms C_{j-1}^{N-1} and C_j^{N-1} are obtained by applying (EC.3) to a homogeneous $M/M/(N-1)$ system with arrival rate λ and all servers operating at rate μ . This results in:

$$C_j^N = \frac{\rho}{N\rho_1} \left(j \binom{N}{j} \rho^{-j} \right) + \frac{1}{N} \left((N-j) \binom{N}{j} \rho^{-j} \right) \quad (\text{EC.7})$$

The corresponding special cases are given by: $C_0^N = 1$, $C_1^N = \rho_1^{-1} + (N-1)\rho$, and $C_N^N = \frac{\rho}{\rho_1} \rho^{-N}$. We then simplify (EC.4) by substituting for C_j^N from (EC.7), to obtain:

$$\begin{aligned} D_N &= \left(\frac{N!}{\rho^N} \left(\frac{\rho}{\rho_1} + \frac{\rho}{N} \left(1 - \frac{\rho}{\rho_1} \right) \right) \sum_{j=0}^{N-1} \frac{\rho^j}{j!} + \frac{\rho}{\rho_1} - 1 \right) + \left(1 + \frac{1}{C_1^N - 1} \right) \\ &= \frac{\rho}{\rho_1} \left(\frac{N!}{\rho^N} \left(1 - \frac{\rho}{N} \left(1 - \frac{\rho_1}{\rho} \right) \right) \sum_{j=0}^{N-1} \frac{\rho^j}{j!} + 1 + \frac{\rho_1}{\rho} \frac{\rho}{N - \left(\rho + 1 - \frac{\rho}{\rho_1} \right)} \right) \end{aligned} \quad (\text{EC.8})$$

Next, we simplify (EC.2) by substituting for D_N from (EC.8), to obtain:

$$P_0 = \left(\left(1 - \frac{\rho}{N} \left(1 - \frac{\rho_1}{\rho} \right) \right) \sum_{j=0}^{N-1} \frac{\rho^j}{j!} + \frac{\rho^N}{N!} \left(1 + \frac{\rho_1}{N - \left(\rho + 1 - \frac{\rho}{\rho_1} \right)} \right) \right)^{-1}$$

To express P_0 in terms of $C(N, \rho)$, the Erlang C formula, we add and subtract the term $\frac{N}{N-\rho} \frac{\rho^N}{N!}$ within, to obtain:

$$P_0 = \left(\left(1 - \frac{\rho}{N} \left(1 - \frac{\rho_1}{\rho} \right) \right) \left(\sum_{j=0}^{N-1} \frac{\rho^j}{j!} + \frac{N}{N-\rho} \frac{\rho^N}{N!} \right) + \frac{\rho^N}{N!} \left(1 + \frac{\rho_1}{N - \left(\rho + 1 - \frac{\rho}{\rho_1} \right)} - \frac{N}{N-\rho} \left(1 - \frac{\rho}{N} \left(1 - \frac{\rho_1}{\rho} \right) \right) \right) \right)^{-1}$$

which reduces to:

$$\begin{aligned} P_0 &= \left(\left(1 - \frac{\rho}{N} \left(1 - \frac{\rho_1}{\rho} \right) \right) \left(\sum_{j=0}^{N-1} \frac{\rho^j}{j!} + \frac{N}{N-\rho} \frac{\rho^N}{N!} \right) - \frac{\rho}{N} \left(1 - \frac{\rho_1}{\rho} \right) \frac{\frac{N}{N-\rho} \frac{\rho^N}{N!}}{N - \left(\rho + 1 - \frac{\rho}{\rho_1} \right)} \right)^{-1} \\ &= \left(\sum_{j=0}^{N-1} \frac{\rho^j}{j!} + \frac{N}{N-\rho} \frac{\rho^N}{N!} \right)^{-1} \left(1 - \frac{\rho}{N} \left(1 - \frac{\rho_1}{\rho} \right) \left(1 + \frac{C(N, \rho)}{N - \left(\rho + 1 - \frac{\rho}{\rho_1} \right)} \right) \right)^{-1} \end{aligned} \quad (\text{EC.9})$$

Finally, (EC.5) simplifies to:

$$I(\mu_1, \mu, \mu, \dots, \mu; \lambda, N) = P_0 + \sum_{k=1}^{N-1} \binom{N-1}{k} P(2, 3, \dots, k+1)$$

Substituting for P_0 from (EC.9) and $P(2, 3, \dots, k+1)$ from (EC.6), we get:

$$\begin{aligned}
I(\mu_1, \mu; \lambda, N) &= P_0 + \sum_{k=1}^{N-1} \binom{N-1}{k} \left(\frac{(N-k)! P_0 \rho^k}{N!} \right) \\
&= \left(1 - \frac{\rho}{N} \right) \left(\sum_{k=0}^{N-1} \frac{\rho^k}{k!} + \frac{N}{N-\rho} \frac{\rho^N}{N!} \right) P_0 \\
&= \left(1 - \frac{\rho}{N} \right) \left(1 - \frac{\rho}{N} \left(1 - \frac{\rho_1}{\rho} \right) \left(1 + \frac{C(N, \rho)}{N - \left(\rho + 1 - \frac{\rho}{\rho_1} \right)} \right) \right)^{-1} \\
&= \left(1 - \frac{\rho}{N} \right) \left(1 - \frac{\rho}{N} \left(1 - \frac{\mu}{\mu_1} \right) \left(1 + \frac{C(N, \rho)}{N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right)} \right) \right)^{-1},
\end{aligned}$$

as desired. ■

Proof of Theorem 2. We start with the expression for I from (4), and take its first partial derivative with respect to μ_1 :

$$\begin{aligned}
\frac{\partial I}{\partial \mu_1} &= - \left(1 - \frac{\rho}{N} \right) \left(1 - \frac{\rho}{N} \left(1 - \frac{\mu}{\mu_1} \right) \left(1 + \frac{C(N, \rho)}{N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right)} \right) \right)^{-2} \frac{\partial}{\partial \mu_1} \left(1 - \frac{\rho}{N} \left(1 - \frac{\mu}{\mu_1} \right) \left(1 + \frac{C(N, \rho)}{N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right)} \right) \right) \\
&= - \frac{N}{N-\rho} I^2 \frac{\partial}{\partial \mu_1} \left(1 - \frac{\rho}{N} \left(1 - \frac{\mu}{\mu_1} \right) \left(1 + \frac{C(N, \rho)}{N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right)} \right) \right) = \frac{\rho}{N-\rho} I^2 \frac{\partial}{\partial \mu_1} \left(\left(1 - \frac{\mu}{\mu_1} \right) \left(1 + \frac{C(N, \rho)}{N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right)} \right) \right)
\end{aligned}$$

Applying the product rule, and simplifying the expression, we get (5). Next, for convenience, we rewrite (5) as:

$$\frac{N-\rho}{\lambda} \frac{\partial I}{\partial \mu_1} = \frac{I^2}{\mu_1^2} \left(1 + \frac{C(N, \rho)}{N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right)} + \left(1 - \frac{\mu_1}{\mu} \right) \frac{\mu_1}{\mu} \frac{C(N, \rho)}{\left(N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right) \right)^2} \right) \quad (\text{EC.10})$$

Differentiating this equation once more with respect to μ_1 by applying the product rule, we get:

$$\begin{aligned}
\frac{N-\rho}{\lambda} \frac{\partial^2 I}{\partial \mu_1^2} &= \left(\frac{2I}{\mu_1^2} \frac{\partial I}{\partial \mu_1} - \frac{2I^2}{\mu_1^3} \right) \left(1 + \frac{C(N, \rho)}{N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right)} + \left(1 - \frac{\mu_1}{\mu} \right) \frac{\mu_1}{\mu} \frac{C(N, \rho)}{\left(N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right) \right)^2} \right) \\
&\quad + \frac{I^2}{\mu_1^2} \frac{\partial}{\partial \mu_1} \left(1 + \frac{C(N, \rho)}{N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right)} + \left(1 - \frac{\mu_1}{\mu} \right) \frac{\mu_1}{\mu} \frac{C(N, \rho)}{\left(N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right) \right)^2} \right) \\
&= \left(\frac{2I}{\mu_1^2} \frac{\partial I}{\partial \mu_1} - \frac{2I^2}{\mu_1^3} \right) \frac{\mu_1^2}{I^2} \frac{N-\rho}{\lambda} \frac{\partial I}{\partial \mu_1} + \frac{I^2}{\mu_1^2} \frac{\partial}{\partial \mu_1} \left(1 + \frac{C(N, \rho)}{N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right)} + \left(1 - \frac{\mu_1}{\mu} \right) \frac{\mu_1}{\mu} \frac{C(N, \rho)}{\left(N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right) \right)^2} \right)
\end{aligned}$$

Applying the product rule for the second term, and simplifying the expression, we get:

$$\frac{\partial^2 I}{\partial \mu_1^2} = \frac{2}{I} \left(\frac{\partial I}{\partial \mu_1} \right)^2 - \frac{2}{\mu_1} \left(\frac{\partial I}{\partial \mu_1} \right) - \frac{2I^2}{\mu_1 \mu^2} \frac{\lambda}{N-\rho} \frac{C(N, \rho)}{\left(N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right) \right)^2} \left(1 + \left(1 - \frac{\mu_1}{\mu} \right) \frac{1}{N - \left(\rho + 1 - \frac{\mu_1}{\mu} \right)} \right)$$

The expression in (6) is then obtained by substituting for $\frac{\partial I}{\partial \mu_1}$ from (5), and carefully going through some incredibly messy (but straightforward) algebra. ■

Proof of Theorem 3. In order to prove this theorem, we make the transformation

$$t = \rho + 1 - \frac{\mu_1}{\mu} \quad (\text{EC.11})$$

For example, when $\mu_1 = \underline{\mu}_1 = (\lambda - (N-1)\mu)^+$, $t = \bar{t} = \min(\rho + 1, N)$. Using this transformation, the $\frac{I}{\mu_1}$ term that appears in the beginning of the expression for the second derivative of the idle time (6) can be written in terms of t as follows.

$$\frac{I}{\mu_1} = \frac{(N-\rho)(N-t)}{\mu g(t)}$$

where

$$g(t) = N(N-t)(\rho+1-t) - \rho(\rho-t)(N-t + C(N, \rho))$$

Note that $g(t) > 0$, since $I > 0$, $N > \rho$, and from stability, $N > t$. Substituting this in (6), and using (EC.11) to complete the transformation, we get the following expression for the second derivative of the idle time in terms of t .

$$\frac{\partial^2 I}{\partial \mu_1^2} = H(t) = -\frac{2\lambda(N-\rho)^2 f(t)}{\mu^3 g^3(t)}$$

where we use the notation $g^3(t)$ to denote $(g(t))^3$, and

$$f(t) = \left((N-t)^2 - \rho C(N, \rho) \right) (N-t + C(N, \rho)) + \left(N - (\rho-t)^2 \right) (\rho+1-t) C(N, \rho)$$

In order to prove the theorem, we now need to show that

- (a) There exists a threshold $t^\dagger \in (-\infty, \bar{t}]$ such that $H(t) < 0$ for $-\infty < t < t^\dagger$, and $H(t) > 0$ for $t^\dagger < t < \bar{t}$.
- (b) $H(t) > 0 \Rightarrow H'(t) > 0$.

To show these statements, we prove the following three properties of f and g .

- $f(t)$ is a decreasing function of t .
- $g(t)$ is a decreasing function of t .
- $f(0) > 0$.

In what follows, for convenience, we denote $C(N, \rho)$ simply by C . Differentiating $f(t)$, we get

$$\begin{aligned} f'(t) &= -((N-t)^2 - \rho C) - 2(N-t)(N-t+C) - (N - (\rho-t)^2)C + 2(\rho-t)(\rho+1-t)C \\ &= -3((N-t)^2 + (-(\rho-t)^2 + (N-\rho))C) \\ &= -3((N-t)^2(1-C) + ((N-t)^2 - (\rho-t)^2) + (N-\rho))C \\ &= -3((N-t)^2(1-C) + ((N-t+\rho-t)(N-\rho) + (N-\rho))C) \\ &= -3((N-t)^2(1-C) + (N-t+\rho+1-t)(N-\rho)C) \\ &< 0 \end{aligned}$$

The last step follows by noting that $N-t > 0$, $\rho+1-t \geq 0$, $N-\rho > 0$, and $0 < C(N, \rho) < 1$ when $0 < \rho < N$. This shows that $f(t)$ is a decreasing function of t . Next, differentiating $g(t)$, we get

$$\begin{aligned} g'(t) &= -N(N-t) - N(\rho+1-t) + \rho(\rho-t) + \rho(N-t+C) \\ &= -N(N-t+\rho+1-t) + \rho(\rho+1-t) + \rho(N-t) - \rho(1-C) \\ &= -(N-\rho)(N-t+\rho+1-t) - \rho(1-C) \\ &< 0 \end{aligned}$$

The last step follows by noting that $N - t > 0$, $\rho + 1 - t \geq 0$, $N - \rho > 0$, and $0 < C(N, \rho) < 1$ when $0 < \rho < N$. This shows that $g(t)$ is a decreasing function of t . Finally, evaluating $f(0)$, we get

$$\begin{aligned} f(0) &= (N^2 - \rho C)(N + C) + (N - \rho^2)(\rho + 1)C \\ &= N^3 - \rho^3 C + N^2 C - \rho^2 C + NC - \rho C^2 \\ &= (N^3 - \rho^3) + \rho^3(1 - C) + (N^2 - \rho^2)C + (N - \rho)C + \rho C(1 - C) \\ &> 0 \end{aligned}$$

The last step follows by noting that $N - \rho > 0$, and $0 < C(N, \rho) < 1$ when $0 < \rho < N$.

We are now ready to prove the statements (a-b).

- (a) First, note that because $f(t)$ is decreasing and $f(0) > 0$, there exists a threshold $t^\dagger \in (0, \bar{t}]$ such that $f(t) > 0$ for $-\infty < t < t^\dagger$, and $f(t) < 0$ for $t^\dagger < t < \bar{t}$. (Note that if $f(\bar{t}) > 0$, then we let $t^\dagger = \bar{t}$ so that $f(t) < 0$ in an empty interval.) Next, since $g(t) > 0$ for all $t \in (-\infty, \bar{t}]$, the sign of $H(t)$ is simply the opposite of the sign of $f(t)$. Statement (a) now follows directly.
- (b) Statement (b) is equivalent to showing that $f(t) < 0 \Rightarrow H'(t) > 0$. Differentiating $H(t)$, we get

$$\begin{aligned} H'(t) &= -\frac{2\lambda(N - \rho)^2}{\mu^3} \left(\frac{g^3(t)f'(t) - 3f(t)g^2(t)g'(t)}{g^6(t)} \right) \\ &= -\frac{2\lambda(N - \rho)^2}{\mu^3} \left(\frac{g(t)f'(t) - 3f(t)g'(t)}{g^4(t)} \right) \end{aligned}$$

Since $g(t) > 0$, $f'(t) < 0$, and $g'(t) < 0$, it follows that $H'(t) > 0$ whenever $f(t) < 0$. This concludes the proof. ■

Proof of Theorem 4. The “only if” direction is straightforward. Briefly, it follows from the fact that, by definition, any symmetric equilibrium $\mu^* > \frac{\lambda}{N}$ must be an interior global maximizer of $U(\mu_1, \mu^*)$ in the interval $\mu_1 \in (\frac{\lambda}{N}, \infty)$.

The “if” direction requires more care. We first show that the utility function $U(\mu_1, \mu^*)$ inherits the properties of the idle time function $I(\mu_1, \mu^*)$ as laid out in Theorem 3, and then consider the two cases when it is either increasing or decreasing at $\mu_1 = \frac{\lambda}{N}$.

Recall that $U(\mu_1, \mu^*) = I(\mu_1, \mu^*) - c(\mu_1)$. Let $\mu_1^\dagger \in [\underline{\mu}_1, \infty)$ be the threshold of Theorem 3. We subdivide the interval $(\underline{\mu}_1, \infty)$ as follows, in order to analyze $U(\mu_1, \mu^*)$.

- Consider the interval $(\underline{\mu}_1, \mu_1^\dagger)$, where, from Theorem 3, we know that $I'''(\mu_1, \mu^*) < 0$. Therefore, $U'''(\mu_1, \mu^*) = I'''(\mu_1, \mu^*) - c'''(\mu_1) < 0$. This means that $U''(\mu_1, \mu^*)$ is decreasing in this interval. (Note that this interval could be empty, i.e., it is possible that $\mu_1^\dagger = \underline{\mu}_1$.)
- Consider the interval (μ_1^\dagger, ∞) , where, from Theorem 3, we know that $I''(\mu_1, \mu^*) < 0$. Therefore, $U''(\mu_1, \mu^*) = I''(\mu_1, \mu^*) - c''(\mu_1) < 0$. This means that $U(\mu_1, \mu^*)$ is concave in this interval.

Thus, the utility function $U(\mu_1, \mu^*)$, like the idle time function $I(\mu_1, \mu^*)$, may start out as a convex function at $\mu_1 = \underline{\mu}_1$, but it eventually becomes concave, and stays concave thereafter. Moreover, because the cost function c is increasing and convex, $\lim_{\mu_1 \rightarrow \infty} U(\mu_1, \mu^*) = -\infty$, which implies that $U(\mu_1, \mu^*)$ must eventually be *decreasing* concave.

We now consider two possibilities for the behavior of $U(\mu_1, \mu^*)$ in the interval $(\frac{\lambda}{N}, \infty)$:

Case (I): $U(\mu_1, \mu^*)$ is increasing at $\mu_1 = \frac{\lambda}{N}$. If $\mu_1^\dagger > \frac{\lambda}{N}$ (see Figure 1(a)), $U(\mu_1, \mu^*)$ would start out being increasing convex, reach a rising point of inflection at $\mu_1 = \mu_1^\dagger$, and then become increasing concave. (Otherwise, if $\mu_1^\dagger \leq \frac{\lambda}{N}$, $U(\mu_1, \mu^*)$ would just be increasing concave to begin with.) It would then go on to attain a (global) maximum, and finally become decreasing concave. This means that the unique stationary point of $U(\mu_1, \mu^*)$ in this interval must be at this (interior) global maximum. Since $U'(\mu^*, \mu^*) = 0$ (from the symmetric first order condition (9)), $\mu_1 = \mu^*$ must be the global maximizer of the utility function $U(\mu_1, \mu^*)$, and hence a symmetric equilibrium.

Case (II): $U(\mu_1, \mu^*)$ is decreasing at $\mu_1 = \frac{\lambda}{N}$. Because $U(\mu^*, \mu^*) \geq U(\frac{\lambda}{N}, \mu^*)$, $U(\mu_1, \mu^*)$ must eventually increase to a value at or above $U(\frac{\lambda}{N}, \mu^*)$, which means it must start out being decreasing *convex* (see Figure 1(b)), attain a minimum, then become increasing convex. It would then follow the same pattern as in the previous case, i.e., reach a rising point of inflection at $\mu_1 = \mu_1^\dagger$, and then become increasing concave, go on to attain a (global) maximum, and finally become decreasing concave. This means that it admits two stationary points – a minimum and a maximum. Since $U'(\mu^*, \mu^*) = 0$ (from the symmetric first order condition (9)) and $U(\mu^*, \mu^*) \geq U(\frac{\lambda}{N}, \mu^*)$, $\mu_1 = \mu^*$ must be the (global) maximizer, and hence a symmetric equilibrium.

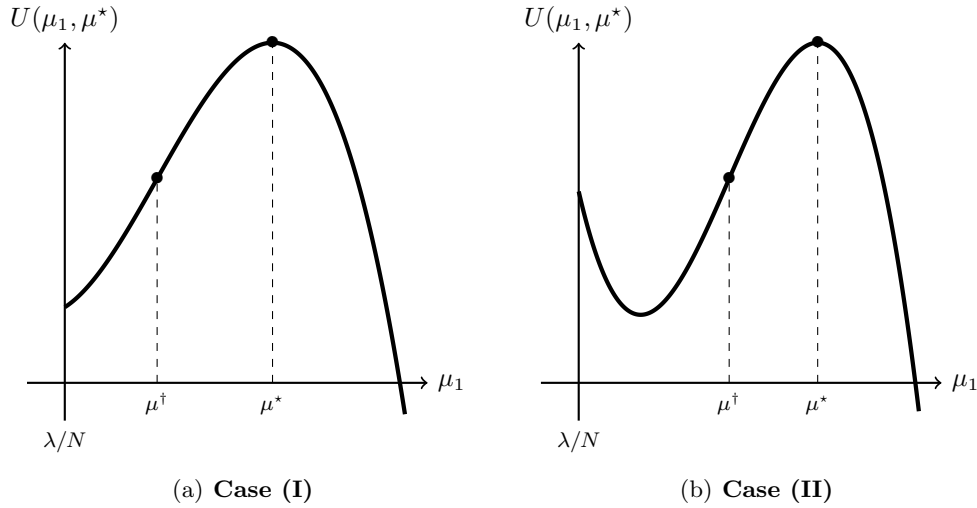


FIGURE EC.1. The graphic depiction of the proof of Theorem 4.

Note that if $U(\mu_1, \mu^*)$ is stationary at $\mu_1 = \frac{\lambda}{N}$, it could either start out increasing or decreasing in the interval $(\frac{\lambda}{N}, \infty)$, and one of the two cases discussed above would apply accordingly.

Finally, to conclude the proof, note that (10) is equivalent to the inequality $U(\mu^*, \mu^*) \geq U(\frac{\lambda}{N}, \mu^*)$, obtained by plugging in and evaluating the utilities using (7) and (4). This completes the proof. ■

Proof of Theorem 5. The symmetric first order condition (9) can be rewritten as

$$C\left(N, \frac{\lambda}{\mu}\right) = \mu^2 c'(\mu) \frac{N^2}{\lambda} + \frac{\lambda}{\mu} - N$$

It suffices to show that if $\lambda c'(\frac{\lambda}{N}) < 1$, then the left hand side and the right hand side intersect at least once in $(\frac{\lambda}{N}, \infty)$. We first observe that the left hand side, the Erlang-C function, is shown to be convex and increasing in $\rho = \frac{\lambda}{\mu}$ (pages 8 and 11 of [43]). This means that it is decreasing and convex in μ . Moreover, $C(N, N) = 1$ and $C(N, 0) = 0$, which means that the left hand side decreases from 1 to 0 in a convex fashion as μ runs from $\frac{\lambda}{N}$ to ∞ . The right hand side is clearly convex in μ , and is equal to $\lambda c'(\frac{\lambda}{N})$ when $\mu = \frac{\lambda}{N}$, and approaches ∞ as μ approaches ∞ . Therefore, if $\lambda c'(\frac{\lambda}{N}) < 1$, then the two curves must intersect at least once in $(\frac{\lambda}{N}, \infty)$.

Next, it is sufficient to show that if $2\frac{\lambda}{N}c'(\frac{\lambda}{N}) + (\frac{\lambda}{N})^2 c''(\frac{\lambda}{N}) \geq 1$, then the right hand side is non-decreasing in μ . In order to do so, it suffices to show that

$$\begin{aligned} & \frac{\partial}{\partial \mu} \left(\mu^2 c'(\mu) \frac{N^2}{\lambda} + \frac{\lambda}{\mu} - N \right) \geq 0 \\ \Leftrightarrow & (\mu^2 c''(\mu) + 2\mu c'(\mu)) \frac{N^2}{\lambda} - \frac{\lambda}{\mu^2} \geq 0 \\ \Leftrightarrow & (\mu^2 c''(\mu) + 2\mu c'(\mu)) \left(\frac{N\mu}{\lambda} \right)^2 \geq 1 \end{aligned}$$

The left hand side is a non-decreasing function of μ , therefore, in the interval $(\frac{\lambda}{N}, \infty)$, we have

$$(\mu^2 c''(\mu) + 2\mu c'(\mu)) \left(\frac{N\mu}{\lambda} \right)^2 \geq \left(\frac{\lambda}{N} \right)^2 c''\left(\frac{\lambda}{N}\right) + 2\frac{\lambda}{N} c'\left(\frac{\lambda}{N}\right) \geq 1.$$

This completes the proof. ■

Proof of Theorem 6. The utility at any symmetric point can be evaluated as $U(\mu, \mu) = 1 - \frac{\lambda}{N\mu} - c(\mu)$, using (7) and (4). Therefore, it follows that showing $U(\mu_1^*, \mu_1^*) > U(\mu_2^*, \mu_2^*)$ is equivalent to showing that

$$\frac{c(\mu_1^*) - c(\mu_2^*)}{\mu_1^* - \mu_2^*} < \frac{\lambda}{N\mu_1^* \mu_2^*}.$$

The function c is convex by assumption. It follows that

$$c(\mu_1^*) - c(\mu_2^*) \leq (\mu_1^* - \mu_2^*) c'(\mu_1^*). \quad (\text{EC.12})$$

Therefore, rearranging and substituting for $c'(\mu_1^*)$ from the symmetric first order condition (9),

$$\frac{c(\mu_1^*) - c(\mu_2^*)}{\mu_1^* - \mu_2^*} \leq \frac{\lambda}{N^2(\mu_1^*)^2} \left(N - \frac{\lambda}{\mu_1^*} + C\left(N, \frac{\lambda}{\mu_1^*}\right) \right).$$

It has been shown (page 14 of [43], and [27]) that $C\left(N, \frac{\lambda}{\mu}\right) < \frac{\lambda}{N\mu}$. Using this,

$$\frac{c(\mu_1^*) - c(\mu_2^*)}{\mu_1^* - \mu_2^*} < \frac{\lambda}{N^2(\mu_1^*)^2} \left(N - \frac{\lambda}{\mu_1^*} \left(1 - \frac{1}{N} \right) \right) < \frac{\lambda}{N^2(\mu_1^*)^2} (N) = \frac{\lambda}{N(\mu_1^*)^2} < \frac{\lambda}{N\mu_1^* \mu_2^*}.$$

This completes the proof. ■

PROOFS FROM SECTION 4

Proof of Proposition 1. We first observe that if $f(\lambda) = \omega(\lambda)$, then

$$\frac{C^{*,\lambda}(N^\lambda)}{\lambda} \geq c_s \frac{N^\lambda}{\lambda} \rightarrow \infty \text{ as } \lambda \rightarrow \infty.$$

Since Proposition 2 evidences a staffing policy under which $C^{*,\lambda}(N^\lambda)/\lambda$ has a finite limit, having $f(\lambda) = \omega(\lambda)$ cannot result in an asymptotically optimal staffing policy.

Next, we consider the case $f(\lambda) = o(\lambda)$. In this case, $\lambda/N^\lambda \rightarrow \infty$. Since any symmetric equilibrium must have $\mu^{*,\lambda} > \lambda/N^\lambda$ from (3), it follows that if there exists a sequence of symmetric equilibria $\{\mu^{*,\lambda}\}$, then $\mu^{*,\lambda} \rightarrow \infty$ as $\lambda \rightarrow \infty$. We conclude that such a staffing policy cannot be admissible. ■

Proof of Theorem 7. We can rewrite (16) as

$$f(\mu) = g(\mu)$$

where

$$f(\mu) = \frac{1}{a} \text{ and } g(\mu) = \frac{\mu^2}{a^2} c'(\mu) + \frac{1}{\mu}.$$

The two cases of interest (i) and (ii) are as shown in Figure EC.2. Our strategy for the proof is to rewrite (14) in terms of functions f^λ and g^λ that are in some sense close to f and g . Then, in case (i), the fact that $g(\mu)$ lies below $f(\mu)$ for $\mu \in [\mu_1, \mu_2]$ implies that f^λ and g^λ intersect twice. The case (ii) is more delicate, because the sign of $o(\lambda)$ determines if the functions f^λ and g^λ will cross at least twice or not at all. (We remark that it will become clear in that part of the proof where the condition $o(\lambda) < -3$ is needed.)

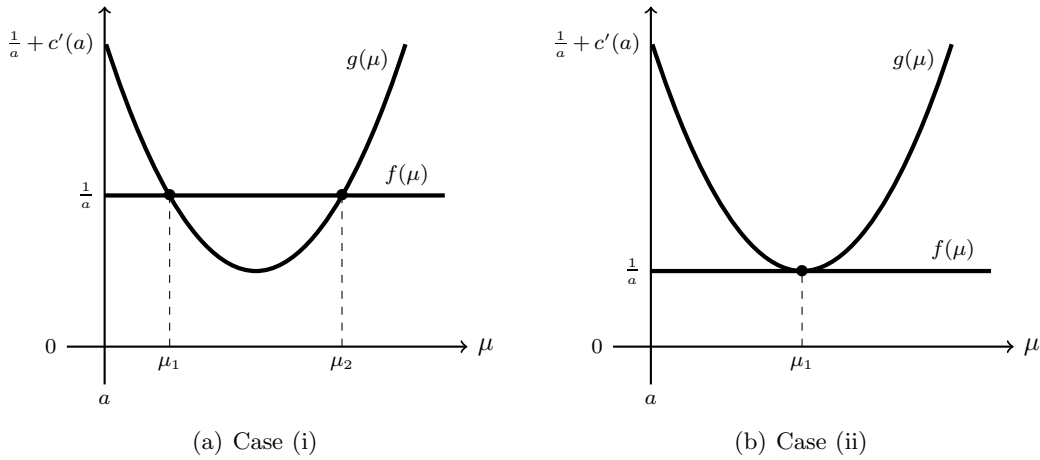


FIGURE EC.2. The limiting first order condition (16).

The first step is to rewrite (14) as

$$f^\lambda(\mu) = g^\lambda(\mu)$$

where

$$f^\lambda(\mu) = \frac{1}{\lambda} C\left(N^\lambda, \frac{\lambda}{\mu}\right) + \frac{N^\lambda}{\lambda}$$

$$g^\lambda(\mu) = \mu^2 c'(\mu) \left(\frac{N^\lambda}{\lambda}\right)^2 + \frac{1}{\mu}.$$

The function g^λ converges uniformly on compact sets to g since for any $\bar{\mu} > 0$, substituting for N^λ in (13) shows that

$$\sup_{\mu \in [0, \bar{\mu}]} |g^\lambda(\mu) - g(\mu)| \leq \bar{\mu}^2 c'(\bar{\mu}) \left(\frac{2}{a} \left| \frac{o(\lambda)}{\lambda} \right| + \left(\frac{o(\lambda)}{\lambda} \right)^2 \right) \rightarrow 0, \quad (\text{EC.13})$$

as $\lambda \rightarrow \infty$. Next, recall $C(N, \rho) \leq 1$ whenever $\rho/N < 1$. Since

$$|f^\lambda(\mu) - f(\mu)| \leq \frac{1}{\lambda} C\left(\frac{1}{a}\lambda + o(\lambda), \frac{\lambda}{\mu}\right) + \left| \frac{o(\lambda)}{\lambda} \right| \quad (\text{EC.14})$$

and $C(\lambda/a + o(\lambda), \lambda/\mu) \leq 1$ for all $\mu > a$ for all large enough λ , the function f^λ converges uniformly to f on any compact set $[a + \epsilon, \bar{\mu}]$ with $\bar{\mu} > a + \epsilon$ and ϵ arbitrarily small. The reason we need only consider compact sets having lower bound $a + \epsilon$ is that it is straightforward to see any solution to (16) has $\mu > a$. It is also helpful to note that g^λ is convex in μ because

$$\frac{d^2}{d\mu^2} g^\lambda(\mu) = (2c'(\mu) + 4\mu c''(\mu) + \mu^2 c'''(\mu)) + 2\frac{1}{\mu^3} > 0 \text{ for } \mu \in (0, \infty),$$

and f^λ is convex decreasing in μ because $C(N, \rho)$ is convex increasing in ρ (pages 8 and 11 of [43]).

We prove (i) and then (ii).

Proof of (i): There exists $\mu_m \in (\mu_1, \mu_2)$ for which $f(\mu_m) > g(\mu_m)$. Then, it follows from (EC.13) and (EC.14) that $f^\lambda(\mu_m) > g^\lambda(\mu_m)$ for all large enough λ . Also,

$$\lim_{\mu \rightarrow \infty} f^\lambda(\mu) = \frac{1}{a} < \lim_{\mu \rightarrow \infty} g^\lambda(\mu) = \infty.$$

and

$$\lim_{\mu \downarrow \lambda/N^\lambda} f^\lambda(\mu) = \frac{1}{\lambda} + \frac{N^\lambda}{\lambda} < g^\lambda\left(\frac{\lambda}{N^\lambda}\right) = c'\left(\frac{\lambda}{N^\lambda}\right) + \frac{N^\lambda}{\lambda}$$

for all large enough λ , where the inequality follows because c is strictly increasing. Since f^λ is convex decreasing and g^λ is convex, we conclude that there exists exactly two solutions to (14).

Proof of (ii): We prove part (a) and then part (b). Recall that μ_1 is the only $\mu > 0$ for which $f(\mu_1) = g(\mu_1)$.

Proof of (ii)(a): For part (a), it is enough to show that for all large enough λ ,

$$f^\lambda(\mu_1) - g^\lambda(\mu_1) > 0. \tag{EC.15}$$

The remainder of the argument follows as in the proof of part (i).

From the definition of f^λ and g^λ in the second paragraph of this proof, and substituting for N^λ ,

$$\begin{aligned} & f^\lambda(\mu_1) - g^\lambda(\mu_1) \\ &= \frac{1}{a} - \frac{1}{\mu_1} - \left(\frac{\mu_1}{a}\right)^2 c'(\mu_1) + \frac{1}{\lambda} C\left(\frac{1}{a}\lambda + o(\lambda), \frac{\lambda}{\mu_1}\right) + \frac{o(\lambda)}{\lambda} \left(1 - \frac{2}{a}(\mu_1)^2 c'(\mu_1)\right) - (\mu_1)^2 c'(\mu_1) \left(\frac{o(\lambda)}{\lambda}\right)^2. \end{aligned}$$

It follows from $f(\mu_1) = g(\mu_1)$ that $1/a - 1/\mu_1 - (\mu_1/a)^2 c'(\mu_1) = 0$, and so, also noting that $C(\lambda/a + o(\lambda), \lambda/\mu_1) > 0$,

$$f^\lambda(\mu_1) - g^\lambda(\mu_1) > \frac{o(\lambda)}{\lambda} \left(1 - \frac{2}{a}(\mu_1)^2 c'(\mu_1)\right) - (\mu_1)^2 c'(\mu_1) \left(\frac{o(\lambda)}{\lambda}\right)^2. \tag{EC.16}$$

Again using the fact that $f(\mu_1) = g(\mu_1)$,

$$1 - \frac{2}{a}(\mu_1)^2 c'(\mu_1) = 1 - 2a \left(\frac{1}{a} - \frac{1}{\mu_1}\right) = -1 + 2\frac{a}{\mu_1}.$$

Then, the term multiplying $o(\lambda)/\lambda$ in (EC.16) is positive if

$$-1 + 2\frac{a}{\mu_1} > 0, \tag{EC.17}$$

which implies (EC.15) holds for all large enough λ .

To see (EC.17), and so complete the proof of part (ii)(a), note that since μ_1 solves (16), and the left-hand side of (16) is convex increasing while the right-hand side is concave increasing, μ_1 also solves the result of differentiating (16), which is

$$\frac{1}{\mu_1^2} = \frac{1}{a^2} (\mu_1^2 c''(\mu_1) + 2\mu_1 c'(\mu_1)).$$

Algebra shows that

$$\frac{1}{\mu_1} - 2 \left(\frac{\mu_1^2}{a^2} c'(\mu_1) \right) = \frac{\mu_1^3}{a^2} c''(\mu_1).$$

We next use (16) to substitute for $\frac{\mu_1^2}{a^2} c'(\mu_1)$ to find

$$\frac{3}{\mu_1} - \frac{2}{a} = \frac{\mu_1^3}{a^2} c''(\mu_1).$$

Since c is convex,

$$\frac{3}{\mu_1} - \frac{2}{a} \geq 0,$$

and so $1.5a \geq \mu_1$, from which (EC.17) follows.

Proof of (ii)(b): Let $\mu^\lambda \in (0, \infty)$ be the minimizer of the function g^λ . The minimizer exists because g^λ is convex and

$$\frac{d}{d\mu} g^\lambda(\mu) = [\mu c'(\mu) + \mu^2 c''(\mu)] \left(\frac{N^\lambda}{\lambda} \right)^2 - \frac{1}{\mu^2},$$

which is negative for all small enough μ , and positive for all large enough μ . It is sufficient to show that for all large enough λ

$$g^\lambda(\mu) - f^\lambda(\mu) > 0 \text{ for all } \mu \in [a, \mu^\lambda]. \quad (\text{EC.18})$$

This is because for all $\mu > \mu^\lambda$, g^λ is increasing and f^λ is decreasing.

Suppose we can establish that for all large enough λ

$$\frac{1}{\mu} \geq \frac{2}{3a} - \frac{\epsilon^\lambda}{2a}, \text{ for all } \mu \in [a, \mu^\lambda], \quad (\text{EC.19})$$

where ϵ^λ satisfies $\epsilon^\lambda \rightarrow 0$ as $\lambda \rightarrow \infty$. Since $g(\mu) \geq f(\mu)$ for all μ , it follows that

$$g^\lambda(\mu) = \mu^2 c'(\mu) \left(\frac{N^\lambda}{\lambda} \right)^2 + \frac{1}{\mu} \geq \left(a - \frac{a^2}{\mu} \right) \left(\frac{N^\lambda}{\lambda} \right)^2 + \frac{1}{\mu}.$$

Substituting for N^λ and algebra shows that

$$\left(a - \frac{a^2}{\mu} \right) \left(\frac{N^\lambda}{\lambda} \right)^2 + \frac{1}{\mu} = \frac{1}{a} + \left(\frac{o(\lambda)}{\lambda} \right) 2 \left(1 - \frac{a}{\mu} \right) + a \left(1 - \frac{a}{\mu} \right) \left(\frac{o(\lambda)}{\lambda} \right)^2.$$

Then, from the definition of f^λ and the above lower bound on g^λ , also using the fact that the assumption $N^\lambda - \lambda/a < 0$ implies the term $o(\lambda)$ is negative,

$$g^\lambda(\mu) - f^\lambda(\mu) \geq \left| \frac{o(\lambda)}{\lambda} \right| \left(\frac{2a}{\mu} - 1 \right) - \frac{1}{\lambda} C \left(N^\lambda, \frac{\lambda}{\mu} \right) + a \left(1 - \frac{a}{\mu} \right) \left(\frac{o(\lambda)}{\lambda} \right)^2.$$

Since $-C(N^\lambda, \lambda/\mu) > -1$ and $1/a - 1/\mu > 0$ from (16) implies $1 - a/\mu > 0$,

$$g^\lambda(\mu) - f^\lambda(\mu) \geq \frac{1}{\lambda} \left(|o(\lambda)| \left(\frac{2a}{\mu} - 1 \right) - 1 \right).$$

Next, from (EC.19),

$$\frac{2a}{\mu} \geq \frac{4}{3} - \epsilon^\lambda,$$

and so

$$g^\lambda(\mu) - f^\lambda(\mu) \geq \frac{1}{\lambda} \left(|o(\lambda)| \left(\frac{1}{3} - \epsilon^\lambda \right) - 1 \right).$$

The fact that $o(\lambda) > 3$ and $\epsilon^\lambda \rightarrow 0$ then implies that for all large enough λ (EC.18) is satisfied.

Finally, to complete the proof, we show that (EC.19) holds. First note that μ^λ as the minimizer of g^λ satisfies

$$(2\mu c'(\mu) + \mu^2 c''(\mu)) \left(\frac{N^\lambda}{\lambda} \right)^2 - \frac{1}{\mu^2} = 0,$$

and that solution is unique and continuous in λ . Hence $\mu^\lambda \rightarrow \mu_1$ as $\lambda \rightarrow \infty$. Then,

$$g^\lambda(\mu^\lambda) \rightarrow g(\mu_1) = \frac{1}{a} \text{ as } \lambda \rightarrow \infty.$$

Furthermore, $g^\lambda(\mu^\lambda)$ approaches $g(\mu_1)$ from above; i.e.,

$$g^\lambda(\mu^\lambda) \downarrow \frac{1}{a} \text{ as } \lambda \rightarrow \infty,$$

because, recalling that the term $o(\lambda)$ is negative,

$$g^\lambda(\mu) = \mu^2 c'(\mu) \left(\frac{1}{a} - \frac{o(\lambda)}{\lambda} \right)^2 + \frac{1}{\mu} > g(\mu) \text{ for all } \mu > 0.$$

Therefore, there exists $\epsilon^\lambda \rightarrow 0$ such that

$$g^\lambda(\mu^\lambda) = \frac{1}{a} - \frac{3}{4} \frac{\epsilon^\lambda}{a},$$

where the $3/(4a)$ multiplier of ϵ^λ is chosen for convenience when obtaining the bound in the previous paragraph. Finally,

$$\frac{1}{\mu} \geq \frac{1}{\mu^\lambda}$$

means that (EC.19) follows if

$$\frac{1}{\mu^\lambda} \geq \frac{2}{3} g^\lambda(\mu^\lambda) = \frac{2}{3} \frac{1}{a} - \frac{1}{2} \frac{\epsilon^\lambda}{a}.$$

To see the above display is valid, note that μ^λ solves

$$(g^\lambda(\mu))' = 0,$$

which from algebra is equivalent to

$$2g^\lambda(\mu^\lambda) - \frac{3}{\mu^\lambda} + (\mu^\lambda)^3 c''(\mu^\lambda) \left(\frac{N^\lambda}{\lambda} \right)^2 = 0.$$

Hence

$$2g^\lambda(\mu^\lambda) - \frac{3}{\mu^\lambda} \leq 0,$$

as required. ■

Proof of Lemma 1. It is enough to show the inequality (10) of Theorem 4 holds. The function c is convex by assumption. It follows that

$$c(\mu^\lambda) - c\left(\frac{\lambda}{N^\lambda}\right) \leq \left(\mu^\lambda - \frac{\lambda}{N^\lambda}\right) c'(\mu^\lambda). \quad (\text{EC.20})$$

Plugging in for $c'(\mu^\lambda)$ from the symmetric first order condition (14) yields (after algebra)

$$\left(\mu^\lambda - \frac{\lambda}{N^\lambda}\right) c'(\mu^\lambda) = \frac{\lambda}{\mu^\lambda N^\lambda} \left(1 - \frac{\lambda}{\mu^\lambda N^\lambda}\right) \left(1 - \frac{\lambda}{\mu^\lambda N^\lambda} + \frac{C(N^\lambda, \lambda/\mu^\lambda)}{N^\lambda}\right).$$

Hence, in order to show the inequality (10) is true, also substituting for $\rho^\lambda = \lambda/\mu^\lambda$, it is enough to verify that

$$\begin{aligned} \frac{\lambda}{\mu^\lambda N^\lambda} \left(1 - \frac{\lambda}{\mu^\lambda N^\lambda}\right) \left(1 - \frac{\lambda}{\mu^\lambda N^\lambda} + \frac{C(N^\lambda, \lambda/\mu^\lambda)}{N^\lambda}\right) &\leq \left(1 - \frac{\lambda}{\mu^\lambda N^\lambda}\right) \left(1 + \left(1 - \frac{\lambda}{\mu^\lambda N^\lambda} + \frac{C(N^\lambda, \lambda/\mu^\lambda)}{N-1}\right)^{-1}\right)^{-1} \\ \iff 1 + \frac{1}{1 - \frac{\lambda}{\mu^\lambda N^\lambda} + \frac{C(N^\lambda, \lambda/\mu^\lambda)}{N^\lambda - 1}} &\leq \frac{1}{\left(\frac{\lambda}{\mu^\lambda N^\lambda}\right) \left(1 - \frac{\lambda}{\mu^\lambda N^\lambda} + \frac{C(N^\lambda, \lambda/\mu^\lambda)}{N^\lambda}\right)}. \end{aligned}$$

Since $N^\lambda - 1 < N^\lambda$, it is enough to show that

$$\begin{aligned} 1 + \frac{1}{1 - \frac{\lambda}{\mu^\lambda N^\lambda} + \frac{C(N^\lambda, \lambda/\mu^\lambda)}{N^\lambda}} &\leq \frac{1}{\left(\frac{\lambda}{\mu^\lambda N^\lambda}\right) \left(1 - \frac{\lambda}{\mu^\lambda N^\lambda} + \frac{C(N^\lambda, \lambda/\mu^\lambda)}{N^\lambda}\right)} \\ \iff 1 &\leq \frac{\left(1 - \frac{\lambda}{\mu^\lambda N^\lambda}\right)}{\frac{\lambda}{\mu^\lambda N^\lambda} \left(1 - \frac{\lambda}{\mu^\lambda N^\lambda} + \frac{C(N^\lambda, \lambda/\mu^\lambda)}{N^\lambda}\right)} \\ \iff \frac{\lambda}{\mu^\lambda N^\lambda} \left(1 - \frac{\lambda}{\mu^\lambda N^\lambda}\right) + \frac{\lambda}{\mu^\lambda N^\lambda} \frac{C(N^\lambda, \lambda/\mu^\lambda)}{N^\lambda} &\leq \left(1 - \frac{\lambda}{\mu^\lambda N^\lambda}\right) \\ \iff \frac{C(N^\lambda, \lambda/\mu^\lambda)}{\lambda/\mu^\lambda} &\leq \left(\frac{N^\lambda \mu^\lambda}{\lambda} - 1\right)^2. \end{aligned}$$

Since $N^\lambda \mu^\lambda / \lambda \rightarrow d > 1$ by assumption, the limit of the right-hand side of the above expression is positive, and, since $C(N^\lambda, \lambda/\mu^\lambda) \leq 1$, the limit of the left-hand side of the above expression is 0. We conclude that for all large enough λ , the above inequality is valid. \blacksquare

Proof of Proposition 2. Let

$$\underline{\mu}^{*,\lambda} = \frac{1}{2} \arg \min\{\mu > 0 : (16) \text{ holds}\}.$$

Next, recalling that $\underline{\mu}^{*,\lambda} > a$, also let

$$\underline{\mu} = \underline{\mu}^{*,\lambda} - \frac{1}{2} (\underline{\mu}^{*,\lambda} - a) > a,$$

so that the system is stable if all servers were to work at rate $\underline{\mu}$ ($\lambda < \underline{\mu} N^\lambda$ for all large enough λ .) It follows from lemma 7 that, for all large enough λ , any μ^λ that satisfies the first order condition (14) also satisfies $\mu^\lambda > \underline{\mu}$. Hence any symmetric equilibrium $\mu^{*,\lambda}$ must also satisfy $\mu^{*,\lambda} > \underline{\mu}$ for all large enough λ , and so

$$\overline{W}^{*,\lambda} < \overline{W}_{\underline{\mu}}^\lambda.$$

Therefore, also using the fact that $\overline{W}^{*,\lambda} > 0$, it follows that

$$c_S \frac{N^\lambda}{\lambda} < \frac{C^{*,\lambda}(N^\lambda)}{\lambda} = c_S \frac{N^\lambda}{\lambda} + \overline{w} \overline{W}^{*,\lambda} < c_S \frac{N^\lambda}{\lambda} + \overline{w} \overline{W}_\mu^\lambda.$$

Then, since $N^\lambda/\lambda \rightarrow 1/a$ as $\lambda \rightarrow \infty$ from (13), it is sufficient to show

$$\overline{W}_\mu^\lambda \rightarrow 0 \text{ as } \lambda \rightarrow \infty.$$

This follows from substituting the staffing $N^\lambda = \lambda/a + o(\lambda)$ in (13) into the well-known formula for the steady state mean waiting time in a $M/M/N^\lambda$ queue with arrival rate λ and service rate $\underline{\mu}$ as follows

$$\begin{aligned} \overline{W}_\mu^\lambda &= \frac{1}{\lambda} \frac{\lambda/\underline{\mu}}{N^\lambda - \lambda/\underline{\mu}} C\left(N^\lambda, \frac{\lambda}{\underline{\mu}}\right) \\ &= \frac{1/\underline{\mu}}{(1/a - 1/\underline{\mu})\lambda + o(\lambda)} C\left(N^\lambda, \frac{\lambda}{\underline{\mu}}\right) \\ &\rightarrow 0, \text{ as } \lambda \rightarrow \infty, \end{aligned}$$

since $C(N^\lambda, \lambda/\underline{\mu}) \in [0, 1]$ for all λ . ■

Proof of Lemma 2. It follows from the equation

$$a(\mu - a) = \mu^3 c'(\mu)$$

that

$$a = \frac{\mu}{2} \left(1 \pm \sqrt{1 - 4\mu c'(\mu)} \right).$$

The condition $4\mu c'(\mu) \leq 1$ is required to ensure that there is a real-valued solution for a . Hence

$$\mathcal{A} = \left\{ \frac{\mu}{2} \left(1 \pm \sqrt{1 - 4\mu c'(\mu)} \right) : 0 \leq 4\mu c'(\mu) \leq 1 \right\}.$$

Since $c'(\mu)$ is well-behaved, this implies that \mathcal{A} is compact, and, in particular, closed. We conclude that $a^* = \sup \mathcal{A} \in \mathcal{A}$, which implies that a^* is finite. ■

Proof of Theorem 8. It follows from Proposition 1 that

$$0 \leq \liminf_{\lambda \rightarrow \infty} \frac{N^{opt,\lambda}}{\lambda} \leq \limsup_{\lambda \rightarrow \infty} \frac{N^{opt,\lambda}}{\lambda} < \infty,$$

because any staffing policy that is not asymptotically optimal also is not optimal for each λ . Consider any subsequence λ' on which either $\liminf_{\lambda \rightarrow \infty} N^{opt,\lambda}/\lambda$ or $\limsup_{\lambda \rightarrow \infty} N^{opt,\lambda}/\lambda$ is attained, and suppose that

$$\frac{N^{opt,\lambda'}}{\lambda'} \rightarrow \frac{1}{a} \text{ as } \lambda \rightarrow \infty, \text{ where } a \in [0, \infty). \quad (\text{EC.21})$$

The definition of asymptotic optimality requires that for each λ' , there exists a symmetric equilibrium service rate $\mu^{*,\lambda'}$. As in the proof of Lemma 1, it is enough to consider sequences $\{\mu^\lambda\}$ that satisfy the first order condition (14). Then, any sequence of solutions $\{\mu^{\lambda'}\}$ to (14) must have $\mu^{\lambda'} \rightarrow \mu$ as $\lambda' \rightarrow \infty$ for μ that satisfies (16), given a in (EC.21). In summary, the choice of a in (EC.21) is constrained by the requirement that a symmetric equilibrium service rate must exist.

Given that there exists at least one symmetric equilibrium service rate for all large enough λ' , it follows in a manner very similar to the proof of Proposition 2 that

$$\overline{W}^{*,\lambda'} \rightarrow 0 \text{ as } \lambda' \rightarrow \infty,$$

even though when there are multiple equilibria we may not be able to guarantee which symmetric equilibrium $\mu^{*,\lambda'}$ the servers choose for each λ' . We conclude that

$$\frac{C^{*,\lambda'}(N^{opt,\lambda'})}{\lambda'} = c_S \frac{N^{opt,\lambda'}}{\lambda'} + \overline{w} \overline{W}^{*,\lambda'} \rightarrow c_S \frac{1}{a}, \text{ as } \lambda' \rightarrow \infty. \quad (\text{EC.22})$$

We argue by contradiction that a in (EC.22) must equal a^* . Suppose not. Then, since

$$\frac{C^{*,\lambda}(N^{ao,\lambda})}{\lambda} \rightarrow c_S \frac{1}{a^*} \text{ as } \lambda \rightarrow \infty$$

by Proposition 2 (and so the above limit is true on any subsequence), and $a^* > a$ by its definition, it follows that

$$C^{*,\lambda'}(N^{ao,\lambda'}) < C^{*,\lambda'}(N^{opt,\lambda'}) \text{ for all large enough } \lambda'.$$

The above inequality contradicts the definition of $N^{opt,\lambda'}$.

The previous argument did not depend on if λ' was the subsequence on which $\liminf_{\lambda \rightarrow \infty} N^{opt,\lambda}/\lambda$ or $\limsup_{\lambda \rightarrow \infty} N^{opt,\lambda}/\lambda$ was attained. Hence

$$\lim_{\lambda \rightarrow \infty} \frac{N^{opt,\lambda}}{\lambda} = \frac{1}{a^*},$$

and, furthermore,

$$\lim_{\lambda \rightarrow \infty} \frac{C^{*,\lambda}(N^{opt,\lambda})}{\lambda} = c_S \frac{1}{a^*}.$$

Since also

$$\lim_{\lambda \rightarrow \infty} \frac{C^{*,\lambda}(N^{ao,\lambda})}{\lambda} = c_S \frac{1}{a^*},$$

the proof is complete. ■

Proof of Lemma 3. We first observe that (16) is equivalently written as:

$$0 = c_E p \mu^{p+2} - a \mu + a^2.$$

The function

$$f(\mu) = c_E p \mu^{p+2} - a \mu + a^2$$

attains its minimum value in $(0, \infty)$ at

$$\underline{\mu} = \left(\frac{a}{c_E p (p+2)} \right)^{1/(p+1)}.$$

The function f is convex in $(0, \infty)$ because $f''(\mu) > 0$ for all $\mu \in (0, \infty)$ and so $\underline{\mu}$ is the unique minimum. It follows that

$$\text{if } f(\underline{\mu}) \begin{cases} < \\ > \\ = \end{cases} 0, \text{ then } \begin{cases} \text{there are at least 2 non-negative solutions to (16)} \\ \text{there is no non-negative solution to (16)} \\ \text{there is exactly one solution to (16)} \end{cases}.$$

Since

$$f(\underline{\mu}) = a^{\frac{p+2}{p+1}} \left(a^{2-\frac{p+2}{p+1}} - \Delta \right)$$

for

$$\Delta := \left(\frac{1}{c_E p(p+1)} \right)^{\frac{1}{p+1}} \left(1 - \left(\frac{1}{c_E p} \right)^{p+1} \left(\frac{1}{p+2} \right)^{p+2} \right) > 0,$$

it follows that

$$\text{if } a^{\frac{p}{p+1}} - \Delta \begin{cases} < \\ > \\ = \end{cases} 0, \text{ then } \begin{cases} \text{there are at least 2 non-negative solutions to (16)} \\ \text{there is no non-negative solution to (16)} \\ \text{there is exactly one solution to (16)} \end{cases}.$$

The expression for Δ can be simplified so that

$$\Delta = \frac{(p+1)}{(p+2)} \left(\frac{1}{c_E p(p+2)} \right)^{\frac{1}{p+1}}.$$

Then, a^* follows by noting that $a^* = \Delta^{(p+1)/p}$ and μ^* follows by noting that $\mu^* = \underline{\mu}$ and then substituting for a^* .

To complete the proof, we must show that a^* and μ^* are both increasing in p . This is because we have already observed that any solution to (16) has $a < \mu$, and the fact that $\mu < 1$ follows directly from the expression for $\underline{\mu}$. We first show a^* is increasing in p , and then argue that this implies μ^* is increasing in p .

To see that a^* is increasing in p , we take the derivative of $\log a^*(p)$ and show that this is positive. Since

$$\begin{aligned} \log a^*(p) &= \log(p+1) - \log(p+2) + \frac{1}{p} \log(p+1) \\ &\quad - \frac{1}{p} \log c_E - \frac{1}{p} \log p - \frac{2}{p} \log(2+p), \end{aligned}$$

it follows that

$$\begin{aligned} (\log a^*(p))' &= \frac{1}{p+1} - \frac{1}{p+2} + \left(\frac{p/(p+1) - \log(p+1)}{p^2} \right) + \frac{1}{p^2} \log c_E \\ &\quad - \frac{\frac{p}{p} - \log(p)}{p^2} - 2 \left(\frac{\frac{p}{p+2} - \log(p+2)}{p^2} \right). \end{aligned}$$

After much simplification, we have

$$(\log a^*(p))' = \frac{1}{p^2} \log c_E + \frac{1}{p^2} \left(\log \left(\frac{p(p+2)^2}{p+1} \right) - \frac{p^2+p+4}{(p+1)(p+2)} \right).$$

Hence it is enough to show that

$$\Delta(p) = \log \left(\frac{p(p+2)^2}{p+1} \right) - \frac{p^2+p+4}{(p+1)(p+2)} \geq 0, \text{ for } p \geq 1.$$

This follows because the first term is increasing in p , and has a value that exceeds 1 when $p = 1$; on the other hand, the second term has a value that is strictly below 1 for all $p \geq 1$.

Finally, it remains to argue that μ^* is increasing in p . At the value $\mu = \mu^*$

$$g(\mu) = \mu^3 c'(\mu) - a\mu + a^2 = 0.$$

At the unique point where the minimum is attained, it is also true that

$$g'(\mu) = \mu^3 c''(\mu) + 3\mu^2 c'(\mu) - a = 0.$$

Since $\mu^3 c''(\mu) + 3\mu^2 c'(\mu)$ is an increasing function of μ , it follows that if a increases, then μ must increase. ■

PROOFS FROM SECTION 5

Proof of Theorem 9. It is sufficient to verify the detailed balance equations. For reference, it is helpful to refer to Figure EC.3, which depicts the relevant portion of the Markov chain. We require the following additional notation. For all $\mathcal{I} \subseteq \{1, 2, \dots, N\}$, all states $\mathbf{s} = (s_1, s_2, \dots, s_{|\mathcal{I}|})$, all servers $s' \in \{1, 2, \dots, N\} \setminus \mathcal{I}$, and integers $j \in \{1, 2, \dots, |\mathcal{I}| + 1\}$, we define the state $\mathbf{s}[s', j]$ by

$$\mathbf{s}[s', j] \equiv (s_1, s_2, \dots, s_{j-1}, s', s_j, \dots, s_{|\mathcal{I}|}).$$

We first observe that:

$$\begin{aligned} \text{Rate into state } \mathbf{s} \text{ due to an arrival} &= \lambda \sum_{s' \notin \mathcal{I}} \sum_{j=1}^{|\mathcal{I}|+1} \pi_{\mathbf{s}[s', j]} p^{\mathcal{I} \cup \{s'\}}(j) \\ &= \lambda \sum_{s' \notin \mathcal{I}} \sum_{j=0}^{|\mathcal{I}|} \frac{\mu_{s'} \pi_B}{\lambda} \prod_{s \in \mathcal{I}} \left(\frac{\mu_s}{\lambda} \right) p^{\mathcal{I} \cup \{s'\}}(j) \\ &= \sum_{s' \notin \mathcal{I}} \mu_{s'} \pi_B \prod_{s \in \mathcal{I}} \frac{\mu_s}{\lambda} = \sum_{s' \notin \mathcal{I}} \mu_{s'} \pi_{\mathbf{s}} \\ &= \text{Rate out of state } \mathbf{s} \text{ due to a departure.} \end{aligned}$$

Then, to complete the proof, we next observe that for each $s' \notin \mathcal{I}$:

$$\begin{aligned} \text{Rate into state } \mathbf{s} \text{ due to a departure} &= \mu_{s_{|\mathcal{I}|}} \pi_{(s_1, s_2, \dots, s_{|\mathcal{I}|-1})} \\ &= \mu_{s_{|\mathcal{I}|}} \pi_B \prod_{s \in \mathcal{I} \setminus \{s_{|\mathcal{I}|\}} \} \frac{\mu_s}{\lambda} \end{aligned}$$

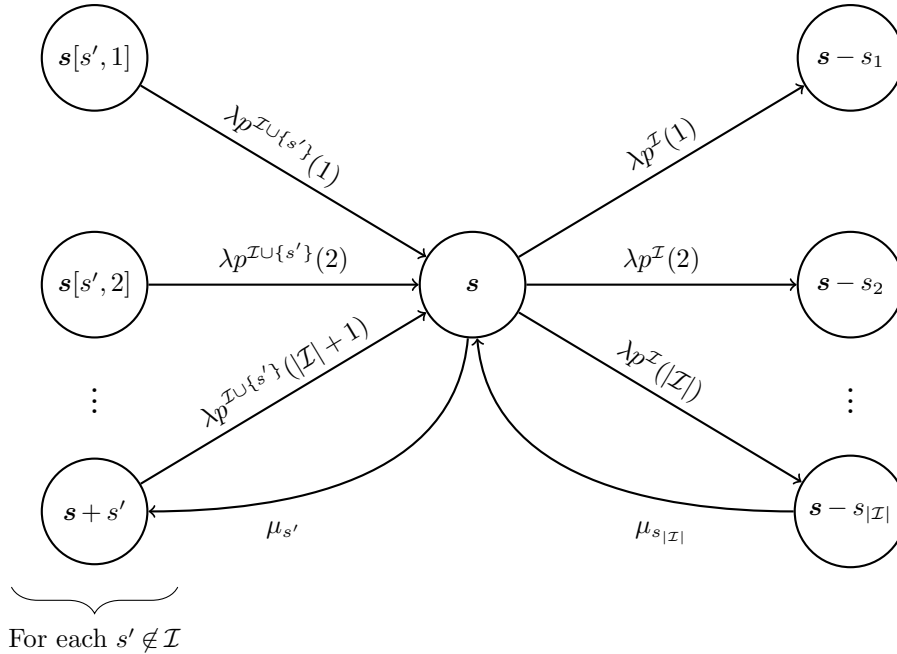
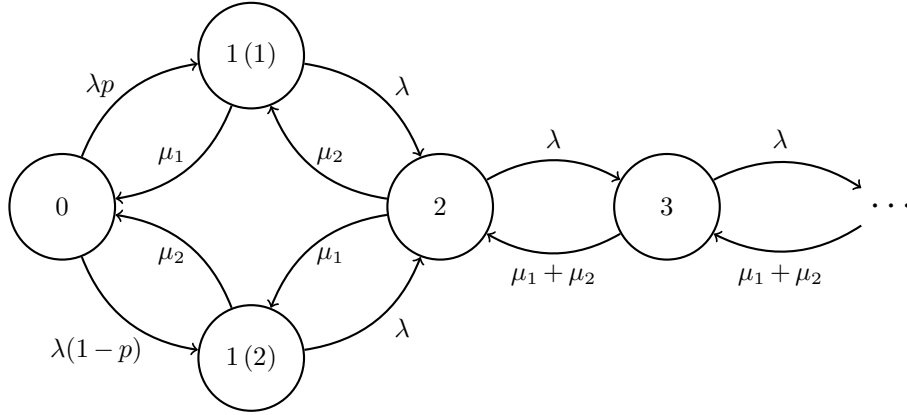


FIGURE EC.3. Snippet of the Markov chain showing the rates into and out of state $\mathbf{s} = (s_1, \dots, s_{|\mathcal{I}|})$. For convenience, we use $\mathbf{s} - s_j$ to denote the state $(s_1, s_2, \dots, s_{j-1}, s_{j+1}, \dots, s_{|\mathcal{I}|})$ and $\mathbf{s} + s'$ to denote the state $\mathbf{s}[s', |\mathcal{I}| + 1] = (s_1, s_2, \dots, s_{|\mathcal{I}|}, s')$.

FIGURE EC.4. The $M/M/2$ Markov chain with probabilistic routing

$$\begin{aligned}
 &= \lambda \pi_B \prod_{s \in \mathcal{I}} \frac{\mu_s}{\lambda} = \lambda \pi_s \\
 &= \text{Rate out of state } s \text{ due to an arrival.}
 \end{aligned}$$

■

Proof of Proposition 3. In order to derive the steady state probability that a server is idle, we first solve for the steady state probabilities of the $M/M/2$ system (with arrival rate λ and service rates μ_1 and μ_2 respectively) under an arbitrary probabilistic routing policy where a job that arrives to find an empty system is routed to server 1 with probability p and server 2 with probability $1 - p$. Then, for an r -routing policy, we simply substitute $p = \frac{\mu_1^r}{\mu_1^r + \mu_2^r}$.

It should be noted that this analysis (and more) for 2 servers has been carried out by [34]. Prior to that, [39] carried out a partial analysis (by analyzing an r -routing policy with $r = 1$). However, we rederive the expressions using our notation for clarity.

The dynamics of this system can be represented by a continuous time Markov chain shown in Figure EC.4 whose state space is simply given by the number of jobs in the system, except when there is just a single job in the system, in which case the state variable also includes information about which of the two servers is serving that job. This system is stable when $\mu_1 + \mu_2 > \lambda$ and we denote the steady state probabilities as follows:

- π_0 is the steady state probability that the system is empty.
- $\pi_1^{(j)}$ is the steady state probability that there is one job in the system, served by server j .
- For all $k \geq 2$, π_k is the steady state probability that there are k jobs in the system.

We can write down the balance equations of the Markov chain as follows:

$$\begin{aligned}
 \lambda \pi_0 &= \mu_1 \pi_1^{(1)} + \mu_2 \pi_1^{(2)} \\
 (\lambda + \mu_1) \pi_1^{(1)} &= \lambda p \pi_0 + \mu_2 \pi_2 \\
 (\lambda + \mu_2) \pi_1^{(2)} &= \lambda (1 - p) \pi_0 + \mu_1 \pi_2 \\
 (\lambda + \mu_1 + \mu_2) \pi_2 &= \lambda \pi_1^{(1)} + \lambda \pi_1^{(2)} + (\mu_1 + \mu_2) \pi_3 \\
 \forall k \geq 3: \quad (\lambda + \mu_1 + \mu_2) \pi_k &= \lambda \pi_{k-1} + (\mu_1 + \mu_2) \pi_{k+1},
 \end{aligned}$$

yielding the following solution to the steady state probabilities:

$$\pi_0 = \frac{\mu_1 \mu_2 (\mu_1 + \mu_2 - \lambda) (\mu_1 + \mu_2 + 2\lambda)}{\mu_1 \mu_2 (\mu_1 + \mu_2)^2 + \lambda (\mu_1 + \mu_2) (\mu_2^2 + 2\mu_1 \mu_2 + (1 - p) (\mu_1^2 - \mu_2^2)) + \lambda^2 (\mu_1^2 + \mu_2^2)} \quad (\text{EC.23})$$

$$\begin{aligned}\pi_1^{(1)} &= \frac{\lambda(\lambda + p(\mu_1 + \mu_2))\pi_0}{\mu_1(\mu_1 + \mu_2 + 2\lambda)} \\ \pi_1^{(2)} &= \frac{\lambda(\lambda + (1-p)(\mu_1 + \mu_2))\pi_0}{\mu_2(\mu_1 + \mu_2 + 2\lambda)}.\end{aligned}$$

Consequently, the steady state probability that server 1 is idle is given by

$$I_1(\mu_1, \mu_2; p) = \pi_0 + \pi_1^{(2)} = \left(1 + \frac{\lambda(\lambda + (1-p)(\mu_1 + \mu_2))}{\mu_2(\mu_1 + \mu_2 + 2\lambda)}\right) \pi_0.$$

Substituting for π_0 , we obtain

$$I_1(\mu_1, \mu_2; p) = \frac{\mu_1(\mu_1 + \mu_2 - \lambda) [(\lambda + \mu_2)^2 + \mu_1\mu_2 + (1-p)\lambda(\mu_1 + \mu_2)]}{\mu_1\mu_2(\mu_1 + \mu_2)^2 + \lambda(\mu_1 + \mu_2) [\mu_2^2 + 2\mu_1\mu_2 + (1-p)(\mu_1^2 - \mu_2^2)] + \lambda^2(\mu_1^2 + \mu_2^2)}. \quad (\text{EC.24})$$

Finally, for an r -routing policy, we let $p = \frac{\mu_1^r}{\mu_1^r + \mu_2^r}$ to obtain:

$$\begin{aligned}I_1^r(\mu_1, \mu_2) &= I_1(\mu_1, \mu_2; p = \frac{\mu_1^r}{\mu_1^r + \mu_2^r}) \\ &= \frac{\mu_1(\mu_1 + \mu_2 - \lambda) \left[(\lambda + \mu_2)^2 + \mu_1\mu_2 + \frac{\mu_2^r}{\mu_1^r + \mu_2^r} \lambda(\mu_1 + \mu_2) \right]}{\mu_1\mu_2(\mu_1 + \mu_2)^2 + \lambda(\mu_1 + \mu_2) \left[\mu_2^2 + 2\mu_1\mu_2 + \frac{\mu_2^r}{\mu_1^r + \mu_2^r} (\mu_1^2 - \mu_2^2) \right] + \lambda^2(\mu_1^2 + \mu_2^2)}.\end{aligned}$$

By symmetry of the r -routing policy, it can be verified that $I_2^r(\mu_1, \mu_2) = I_1^r(\mu_2, \mu_1)$, completing the proof. \blacksquare

Proof of Theorem 10. We first highlight that when all servers operate at the same rate $\mu \in (\frac{\lambda}{N}, \infty)$, both FSF and SSF are equivalent to Random routing. Henceforth, we refer to such a configuration as a symmetric operating point μ . In order to prove that there does not exist a symmetric equilibrium under either FSF or SSF, we show that at any symmetric operating point μ , any one server can attain a strictly higher utility by unilaterally setting her service rate to be slightly lower (in the case of FSF) or slightly higher (in the case of SSF) than μ .

We borrow some notation from the proof of Proposition 3 where we derived the expressions for the steady state probability that a server is idle when there are only 2 servers under any probabilistic policy, parameterized by a number $p \in [0, 1]$ which denotes the probability that a job arriving to an empty system is routed to server 1. Recall that $I_1(\mu_1, \mu_2; p)$ denotes the steady state probability that server 1 is idle under such a probabilistic policy, and the corresponding utility function for server 1 is $U_1(\mu_1, \mu_2; p) = I_1(\mu_1, \mu_2; p) - c(\mu_1)$. Then, by definition, the utility function for server 1 under FSF is given by:

$$U_1^{FSF}(\mu_1, \mu_2) = \begin{cases} U_1(\mu_1, \mu_2; p = 0) & , \mu_1 < \mu_2 \\ U_1(\mu_1, \mu_2; p = \frac{1}{2}) & , \mu_1 = \mu_2 \\ U_1(\mu_1, \mu_2; p = 1) & , \mu_1 > \mu_2. \end{cases}$$

Similarly, under SSF, we have:

$$U_1^{SSF}(\mu_1, \mu_2) = \begin{cases} U_1(\mu_1, \mu_2; p = 1) & , \mu_1 < \mu_2 \\ U_1(\mu_1, \mu_2; p = \frac{1}{2}) & , \mu_1 = \mu_2 \\ U_1(\mu_1, \mu_2; p = 0) & , \mu_1 > \mu_2. \end{cases}$$

Note that while the utility function under any probabilistic routing policy is continuous everywhere, the utility function under FSF or SSF is discontinuous at symmetric operating points. This

discontinuity turns out to be the crucial tool in the proof. Let the two servers be operating at a symmetric operating point μ . Then, it is sufficient to show that there exists $0 < \delta < \mu - \frac{\lambda}{2}$ such that

$$U_1^{FSF}(\mu - \delta, \mu) - U_1^{FSF}(\mu, \mu) > 0, \quad (\text{EC.25})$$

and

$$U_1^{SSF}(\mu + \delta, \mu) - U_1^{SSF}(\mu, \mu) > 0. \quad (\text{EC.26})$$

We show (EC.25), and (EC.26) follows from a similar argument. Note that

$$\begin{aligned} U_1^{FSF}(\mu - \delta, \mu) - U_1^{FSF}(\mu, \mu) &= U_1(\mu - \delta, \mu; p = 0) - U_1\left(\mu, \mu; p = \frac{1}{2}\right) \\ &= (U_1(\mu - \delta, \mu; p = 0) - U_1(\mu, \mu; p = 0)) \\ &\quad + \left(U_1(\mu, \mu; p = 0) - U_1\left(\mu, \mu; p = \frac{1}{2}\right) \right) \end{aligned}$$

Since the first difference, $U_1(\mu - \delta, \mu; p = 0) - U_1(\mu, \mu; p = 0)$, is zero when $\delta = 0$, and is continuous in δ , it is sufficient to show that the second difference, $U_1(\mu, \mu; p = 0) - U_1(\mu, \mu; p = \frac{1}{2})$, is strictly positive:

$$\begin{aligned} U_1(\mu, \mu; p = 0) - U_1\left(\mu, \mu; p = \frac{1}{2}\right) &= I_1(\mu, \mu; p = 0) - I_1\left(\mu, \mu; p = \frac{1}{2}\right) \\ &= \frac{\lambda(2\mu - \lambda)}{(\mu + \lambda)(2\mu + \lambda)} > 0 \quad (\text{using (EC.24)}). \end{aligned}$$

This completes the proof. ■

Proof of Theorem 11. The proof of this theorem consists of two parts. First, we show that under any r -routing policy, any symmetric equilibrium $\mu^* \in (\frac{\lambda}{2}, \infty)$ must satisfy the equation $\varphi(\mu^*) = r$. This is a direct consequence of the necessary first order condition for the utility function of server 1 to attain an interior maximum at μ^* . The second part of the proof involves using the condition $c'(\frac{\lambda}{2}) < \frac{1}{\lambda}$ to show that φ is a *strictly decreasing bijection* onto \mathbb{R} , which would lead to the following implications:

- φ is invertible; therefore, if an r -routing policy admits a symmetric equilibrium, it is unique, and is given by $\mu^* = \varphi^{-1}(r)$.
- $\varphi^{-1}(r)$ is strictly decreasing in r ; therefore, so is the unique symmetric equilibrium (if it exists). Since the mean response time $\mathbb{E}[T]$ is inversely related to the service rate, this establishes that $\mathbb{E}[T]$ at symmetric equilibrium (across r -routing policies that admit one) is increasing in r .

We begin with the first order condition for an interior maximum. The utility function of server 1 under an r -routing policy, from (2), is given by

$$U_1^r(\mu_1, \mu_2) = I_1^r(\mu_1, \mu_2) - c(\mu_1)$$

For $\mu^* \in (\lambda/2, \infty)$ to be a symmetric equilibrium, the function $U_1^r(\mu_1, \mu^*)$ must attain a global maximum at $\mu_1 = \mu^*$. The corresponding first order condition is then given by:

$$\left. \frac{\partial I_1^r}{\partial \mu_1}(\mu_1, \mu^*) \right|_{\mu_1 = \mu^*} = c'(\mu^*), \quad (\text{EC.27})$$

where I_1^r is given by Proposition 3. The partial derivative of the idle time can be computed and the left hand side of the above equation evaluates to

$$\left. \frac{\partial I_1^r}{\partial \mu_1}(\mu_1, \mu^*) \right|_{\mu_1 = \mu^*} = \frac{\lambda(4\lambda + 4\mu^* + \lambda r - 2\mu^* r)}{4\mu^*(\lambda + \mu^*)(\lambda + 2\mu^*)}. \quad (\text{EC.28})$$

Substituting in (EC.27) and rearranging the terms, we obtain:

$$\frac{4(\lambda + \mu^*)}{\lambda(\lambda - 2\mu^*)} (\mu^*(\lambda + 2\mu^*)c'(\mu^*) - \lambda) = r.$$

The left hand side is equal to $\varphi(\mu^*)$, thus yielding the necessary condition $\varphi(\mu^*) = r$.

Next, we proceed to show that if $c'(\frac{\lambda}{2}) < \frac{1}{\lambda}$, then φ is a strictly decreasing bijection onto \mathbb{R} . Note that the function

$$\varphi(\mu) = \frac{4(\lambda + \mu)}{\lambda(\lambda - 2\mu)} (\mu(\lambda + 2\mu)c'(\mu) - \lambda)$$

is clearly a continuous function in $(\frac{\lambda}{2}, \infty)$. In addition, it is a surjection onto \mathbb{R} , as evidenced by the facts that $\varphi(\mu) \rightarrow -\infty$ as $\mu \rightarrow \infty$ and $\varphi(\mu) \rightarrow \infty$ as $\mu \rightarrow \frac{\lambda}{2}+$ (using $c'(\frac{\lambda}{2}) < \frac{1}{\lambda}$).

To complete the proof, it is sufficient to show that $\varphi'(\mu) < 0$ for all $\mu \in (\frac{\lambda}{2}, \infty)$. First, observe that

$$\varphi'(\mu) = \frac{4\psi(\mu)}{\lambda(\lambda - 2\mu)^2},$$

where

$$\psi(\mu) = \mu(\lambda + \mu)(\lambda^2 - 4\mu^2)c''(\mu) + (\lambda^3 + 6\lambda^2\mu - 8\mu^3)c'(\mu) - 3\lambda^2.$$

Since $c'(\frac{\lambda}{2}) < \frac{1}{\lambda}$, as $\mu \rightarrow \frac{\lambda}{2}+$, $\psi(\mu) < 0$. Moreover, since $c'''(\mu) > 0$, for all $\mu > \frac{\lambda}{2}$, we have

$$\psi'(\mu) = -4\mu(\lambda + \mu) \left(\mu^2 - \left(\frac{\lambda}{2} \right)^2 \right) c'''(\mu) - 4 \left(\mu - \frac{\lambda}{2} \right) (\lambda^2 + 6\lambda\mu + 6\mu^2)c''(\mu) - 24 \left(\mu^2 - \left(\frac{\lambda}{2} \right)^2 \right) c'(\mu) < 0.$$

It follows that $\psi(\mu) < 0$ for all $\mu > \frac{\lambda}{2}$. Since $\varphi'(\mu)$ has the same sign as $\psi(\mu)$, we conclude that $\varphi'(\mu) < 0$, as desired. \blacksquare

Proof of Theorem 12. From Theorem 11, we know that if a symmetric equilibrium exists, then it is unique, and is given by $\mu^* = \varphi^{-1}(r)$, where φ establishes a one-to-one correspondence between r and μ^* (μ^* is strictly decreasing in r and vice versa). Therefore, it is enough to show that there exists a finite upper bound $\bar{\mu} > \frac{\lambda}{2}$ such that no service rate $\mu > \bar{\mu}$ can be a symmetric equilibrium under *any* r -routing policy. It would then automatically follow that for $\underline{r} = \varphi(\bar{\mu})$, no r -routing policy with $r \leq \underline{r}$ admits a symmetric equilibrium. We prove this by exhibiting a $\bar{\mu}$ and showing that if $\mu \geq \bar{\mu}$, then the utility function of server 1, $U_1^r(\mu_1, \mu)$, cannot attain a global maximum at $\mu_1 = \mu$ for any $r \in \mathbb{R}$.

We begin by establishing a lower bound for the maximum utility $U_1^r(\mu_1, \mu)$ that server 1 can obtain under any r -routing policy:

$$\max_{\mu_1 > \frac{\lambda}{2}} U_1^r(\mu_1, \mu) \geq U_1^r\left(\frac{\lambda}{2}, \mu\right) = I_1^r\left(\frac{\lambda}{2}, \mu\right) - c\left(\frac{\lambda}{2}\right) \geq -c\left(\frac{\lambda}{2}\right) = U_1^r\left(\frac{\lambda}{2}, \frac{\lambda}{2}\right). \quad (\text{EC.29})$$

By definition, if μ^* is a symmetric equilibrium under any r -routing policy, then the utility function of server 1, $U_1^r(\mu_1, \mu^*)$, is maximized at $\mu_1 = \mu^*$, and hence, using (EC.29), we have

$$U_1^r(\mu^*, \mu^*) \geq U_1^r\left(\frac{\lambda}{2}, \frac{\lambda}{2}\right). \quad (\text{EC.30})$$

Next, we establish some properties on $U_1^r(\mu, \mu)$ that help us translate this necessary condition for a symmetric equilibrium into an upper bound on any symmetric equilibrium service rate. We have,

$$U_1^r(\mu, \mu) = 1 - \frac{\lambda}{2\mu} - c(\mu),$$

which has the following properties:

- Since $c'(\frac{\lambda}{2}) < \frac{1}{\lambda}$, $U_1^r(\mu, \mu)$, as a function of μ , is strictly increasing at $\mu = \frac{\lambda}{2}$.
- $U_1^r(\mu, \mu)$ is a concave function of μ .

This means that $U_1^r(\mu, \mu)$ is strictly increasing at $\mu = \frac{\lambda}{2}$, attains a maximum at the unique $\mu_{\dagger} > \frac{\lambda}{2}$ that solves the first order condition $\mu_{\dagger}^2 c'(\mu_{\dagger}) = \frac{\lambda}{2}$, and then decreases forever. This shape of the curve $U_1^r(\mu, \mu)$ implies that there must exist a unique $\bar{\mu} > \mu_{\dagger}$, such that $U_1^r(\bar{\mu}, \bar{\mu}) = U_1^r(\frac{\lambda}{2}, \frac{\lambda}{2})$.

Since $U_1^r(\mu, \mu)$ is a strictly decreasing function for $\mu > \mu_{\dagger}$, it follows that if $\mu^* > \bar{\mu}$, then, $U_1^r(\mu^*, \mu^*) < U_1^r(\bar{\mu}, \bar{\mu}) = U_1^r(\frac{\lambda}{2}, \frac{\lambda}{2})$, contradicting the necessary condition (EC.30). This establishes the required upper bound $\bar{\mu}$ on any symmetric equilibrium service rate, completing the proof. ■

Proof of Theorem 13. A useful tool for proving this theorem is Theorem 3 from [13], whose statement we have adapted to our model:

THEOREM EC.1. *A symmetric game with a nonempty, convex, and compact strategy space, and utility functions that are continuous and quasiconcave has a symmetric (pure-strategy) equilibrium.*

We begin by verifying that our 2-server game meets the qualifying conditions of Theorem EC.1:

- *Symmetry:* First, all servers have the same strategy space of service rates, namely, $(\frac{\lambda}{2}, \infty)$. Moreover, since an r -routing policy is symmetric and all servers have the same cost function, their utility functions are symmetric as well. Hence, our 2-server game is indeed symmetric.
- *Strategy space:* The strategy space $(\frac{\lambda}{2}, \infty)$ is nonempty and convex, but not compact, as required by Theorem EC.1. Hence, for the time being, we modify the strategy space to be $[\frac{\lambda}{2}, \bar{\mu} + 1]$ so that it is compact, where $\bar{\mu}$ is the upper bound on any symmetric equilibrium, established in Theorem 12, and deal with the implications of this modification later.
- *Utility function:* $U_1^r(\mu_1, \mu_2)$ is clearly continuous. From Mathematica, it can be verified that the idle time function $I_1^r(\mu_1, \mu_2)$ is concave in μ_1 for $r \in \{-2, -1, 0, 1\}$, and since the cost function is convex, this means the utility functions are also concave. (Unfortunately, we could not get Mathematica to verify concavity for non-integral values of r , though we strongly suspect that it is so for the entire interval $[-2, 1]$.)

Therefore, we can apply Theorem EC.1 to infer that an r -routing policy with $r \in \{-2, -1, 0, 1\}$ admits a symmetric equilibrium in $[\frac{\lambda}{2}, \bar{\mu} + 1]$. We now show that the boundaries cannot be symmetric equilibria. We already know from Theorem 12 that $\bar{\mu} + 1$ cannot be a symmetric equilibrium. (We could have chosen to close the interval at any $\mu > \bar{\mu}$. The choice $\bar{\mu} + 1$ was arbitrary.) To see that $\frac{\lambda}{2}$ cannot be a symmetric equilibrium, observe that $c'(\frac{\lambda}{2}) < \frac{1}{\lambda}$ implies that $U_1^r(\mu_1, \frac{\lambda}{2})$ is increasing at $\mu_1 = \frac{\lambda}{2}$ (using the derivative of the idle time computed in (EC.28)), and hence server 1 would have an incentive to deviate. Therefore, any symmetric equilibrium must be an interior point, and from Theorem 11, such an equilibrium must be unique. This completes the proof. ■