

Asymptotic Performance of Energy-Aware Multiserver Queueing Systems with Setup Times

Maccio, Vincent J.
macciiov@mcmaster.ca

Down, Douglas G.
downd@mcmaster.ca

October 20, 2016

Abstract

An intuitive solution to address the immense energy demands of datacentres is to turn servers off to incur less costs. However, when to turn a specific server off, and when to then turn that server back on, are far from trivial questions. As such, many different authors have modeled this problem as an $M/M/C$ queue where each server can be turned on, with an exponentially distributed setup time, or turned off instantaneously. Due to the complexity of the model analysis, authors often examine a specific policy. Moreover, different authors examine different policies under different cost functions. This in turn causes difficulties when making statements or drawing conclusions regarding competing policies. Therefore, we analyse this well established model under the asymptotic regime where the number of servers approaches infinity, i.e. $C \rightarrow \infty$, while the load remains fixed, i.e. $0 < \lambda/(C\mu) < 1$, and show that not only are many of the policies in the literature equivalent under this regime, but they are also optimal under any cost function which is non-decreasing in the expected energy cost and response time.

1 Introduction

Over the past several years, energy concerns in datacentres have driven an interest in queueing systems where individual servers can be turned on to improve performance, and turned off to save on costs. This interest has led to different authors studying the same, or similar, queueing models. However, due to the complexity of the problem, i.e. the choice of cost function, policy implemented, model details, etc., different conclusions can be drawn from similar underlying problems. One consequence of the variety in the problems studied and the corresponding variety of insights, is that it is difficult to confidently draw conclusions which are overarching across the problem domain. To address this issue this work presents a result under an asymptotic regime, which states that when the system parameters are appropriately scaled up, a general class of policies is optimal, under all reasonable cost functions.

To the best of our knowledge, Chen et al. [3] and Sledger et al. [19] were the first to apply queueing models in the context of energy-aware datacentres. This work introduced a popular queueing model which extends the traditional $M/M/C$ queue where each of the C servers can be switched on after a setup delay (to improve performance) and instantly switched off (to increase efficiency). The question remained however, when should servers be turned on, and when should they be turned off. This question was and currently remains a topic of interest. Gandhi et al. [4–7] produced a body of work examining this model under the *staggered setup* policy, where the number of jobs in the system is equal to the number of servers on and in setup when possible, and servers turn off when idle. Furthermore, they also studied the *delayed off* policy, which extends the staggered setup policy, by allowing an idle server to wait an exponentially distributed period of time before it turns off. Mitrani [15] studied this model where a reserved set of servers were brought into setup when the number of jobs in the system exceeds a threshold, and then shuts those servers off once the number of jobs drops below another threshold. This policy was further studied in [9]. Xu and Tian [21] examined the model where e servers will turn off when d servers idle.

From here more sophisticated policies began to emerge. Specifically, the use of a threshold decision variable regarding the number of jobs in the system was often employed. Maccio and Down [13] derived several structural properties pertaining to the optimal policy to allow for more confidence in existing policies,

as well as to create guidelines when determining a new policy to study. Kuehn and Mashly [11] looked at the system which turns on servers if some threshold k is met, and turns them off when idle, under the constraint of a finite buffer. Maccio and Down [13] studied two policies both of which allow for a static provisioning of servers (servers which are always on), and the remaining servers turn off the moment they idle. The first of the policies they examined was the *bulk setup* policy, which begins the setup process of all available servers simultaneously, cancelling the remaining setups once one server has turned on. The second, *staggered threshold*, is in the same vein as [11]. Informally, a dedicated number of servers always remain on, and of the remaining servers, the i th server will begin its setup process once ik jobs have accumulated in the queue. Phung-Duc [16] has also inspected this model where he analysed the policies studied by Gandhi et al., and Phung-Duc and Kawanishi [17] examined the *s-staggered setup* policy which turns on s servers for each job waiting in the queue. Moreover, Ren et al. [18] also studied the staggered threshold policy under a slight modification, in the context of virtual machines. This model also has applications to or is studied in problems which arise in other fields such as manufacturing, logistics, and vacation models [2, 14, 20].

As mentioned previously and as one may infer from the literature review, an existing problem is that analysing a specific policy, or family of policies, is far from trivial. As such, the examination and analysis of a single policy could span an entire work. Moreover, when evaluating one of these policies it may do well for a specific cost function, but poorly under another. Therefore, saying one policy is strictly better than another can be a difficult claim to justify. With the goal in mind for one to be able to make broad claims across a large set of policies and cost functions, this work makes the following contributions:

1. It is shown that as the number of servers and arrival rate are appropriately scaled to infinity such that the load on the system remains fixed on some value less than one, then there exists an infinite set of policies such that any policy belonging to this set will simultaneously minimize both the expected response time and expected energy costs, and therefore will be optimal under all cost functions which are non-decreasing in those metrics.
2. This set of policies is formally described, and observed to include many of the policies currently examined in the literature and discussed previously in this section.
3. Numerical experiments are conducted to determine how quickly the asymptotic behaviours of these system are reached, and it is shown that particular choices of policy parameters may be used to induce faster convergence to the minimum values.

2 Model

The model under study is an $M/M/C$ queue where each server can be switched on and off, and where turn-offs are instantaneous, but turn-ons take an exponentially distributed setup time. This is described formally as follows. Jobs arrive to a central queue following a Poisson process with rate λ , are processed on a first come first served basis, and have processing times (job sizes) which are exponentially distributed with rate μ . Furthermore, there are C homogeneous servers present, each of which can be in one of four energy states: *off*, *setup*, *idle*, or *busy*. For ease of exposition this work often refers to a server being busy, idle, off, or in setup as shorthand for a server being in the corresponding energy state. Regarding definitions and transitions, a server is *idle* if and only if it is on and not processing a job. Moreover, a server can only begin serving a job if it is currently *idle*, in which case the server becomes *busy*. At any time, a server can be switched *off*. Regarding the process of turning a server on, an *off* server can transition to *setup*. Once in *setup*, the server will remain there for a time exponentially distributed with rate γ , after which the server will become *idle*, from which it may become *busy* instantaneously. A system which meets the criteria of the above model is said to be an *energy-aware system*.

This work refers to an energy-aware system as a four-tuple $(C, \lambda, \mu, \gamma)$, where C is the number of servers, and λ , μ , and γ are the arrival, processing, and setup rates, respectively. The system load ρ is defined as $\rho = \lambda/(C\mu)$. Moreover, the well known $M/M/C$ queue is referred to by a three-tuple (C, λ, μ) , with the traditional interpretation of those parameters. As such, one may view an energy-aware system as an extension of an $M/M/C$ queue. To fully understand how a specific energy-aware system behaves, one must also know when, or how it is determined when, each server is turned on and off. Such a description of the server behaviour is referred to as a policy. In this work a specific policy is denoted by π . Some examples of

a policy π are: the number of servers which are on or in setup equals the number of jobs in the system, turn on all servers once there are k jobs in the system and turn them off when they idle, keep all servers on all the time, etc. It is natural to want to compare such policies against each other. That is, to determine if one policy is better than another.

In order to compare different policies, one must first have metrics to evaluate. This work examines the trade-off between efficacy and efficiency. The expected response time, denoted by $\mathbb{E}[R]$ is employed to evaluate efficacy, while the expected energy cost, denoted by $\mathbb{E}[E]$, is employed to evaluate efficiency. The expected response time is the expected amount of time a job spends in the system, from arrival to departure. The expected energy cost takes a little more care to define. Each of the energy states (*off*, *idle*, *busy*, and *setup*) have a corresponding energy consumption rate. Let these rates be denoted by E_{Off} , E_{Idle} , E_{Busy} , and E_{Setup} , respectively. Furthermore, let the random variables C_{Off} , C_{Idle} , C_{Busy} , and C_{Setup} denote the number of servers which are off, idle, busy, or in setup, respectively. Then

$$\mathbb{E}[E] = E_{\text{Off}}\mathbb{E}[C_{\text{Off}}] + E_{\text{Idle}}\mathbb{E}[C_{\text{Idle}}] + E_{\text{Busy}}\mathbb{E}[C_{\text{Busy}}] + E_{\text{Setup}}\mathbb{E}[C_{\text{Setup}}]. \quad (1)$$

Without loss of generality, it is assumed that $E_{\text{Busy}} = 1$, and the remaining rates are appropriately normalized. Furthermore, it is also assumed that $E_{\text{Idle}} < E_{\text{Setup}}$, $E_{\text{Idle}} < E_{\text{Busy}}$, and $E_{\text{Off}} = 0$, although the latter could be relaxed to account for lower energy consumption states where the server cannot process jobs, e.g. sleep states. We will see that it is often more instructive to look at the expected energy cost of the system on a per job basis. This viewpoint is extended easily from the previous notation. Letting E^J denote the energy cost incurred by job J , it is not difficult to see $\mathbb{E}[E] = \lambda\mathbb{E}[E^J]$. Therefore, the expected energy cost incurred by a single job is simply a normalized version of the expected total energy cost. However, it will be seen later in this section that looking at $\mathbb{E}[E^J]$ can grant a greater understanding of the system behaviour. This viewpoint is examined in greater detail in Section 4.

With the cost metrics defined, one can then create a cost function dependent on these metrics and begin to compare policies. The number of cost functions which can be defined is infinite, but common cost functions do arise in the literature, e.g. $\mathbb{E}[R]\mathbb{E}[E]$, and $\mathbb{E}[R] + \beta\mathbb{E}[E]$. Due to the diversity of the set of possible cost functions, conclusions which span across many cost functions are often difficult to make, and moreover, niche behaviours are often easy to invoke by tweaking parameters within a cost function, such as β in $\mathbb{E}[R] + \beta\mathbb{E}[E]$. For example, one could imagine a system with a small number of servers where the policy which keeps all servers on would be optimal when β is small, since this minimizes the expected response time. But when β is large, the optimal policy for the same system may be one where servers are kept off for long periods of time while they wait for a large number of jobs to accumulate before expending the energy to turn on. As such, this work strives to draw conclusions applicable to large sets of cost functions, and in fact does so for all *well-formed cost functions*.

Definition 1. Well-Formed Cost Function: A cost function $\mathcal{C}(\cdot)$ is a well-formed cost function if it is non-decreasing in, dependent on, and only dependent on, the expected response time, i.e. $\mathbb{E}[R]$, and the expected energy costs, i.e. $\mathbb{E}[E]$.

As stated previously, the point of these cost functions is to allow the comparison of policies applied to the same energy-aware system. Similar to how this work strives to make conclusions across a large set of cost functions, it also strives to make conclusions across large sets of policies. To this point, two important sets of policies are defined below.

Definition 2. Class A Policy: A policy is said to be a Class A policy if the following conditions are met:

1. Server setups are invoked following a threshold scheme.
2. A server will never turn off if there is a job which it could be processing.

To elaborate on the definition of a Class A policy, a threshold scheme pertaining to server turn-ons implies that each server i , $0 \leq i \leq C$, has a corresponding threshold value $k_{i,j}$, $0 \leq j < i$, such that while there are j servers currently on, if the number of jobs in the system is greater than or equal to $k_{i,j}$ and server i is currently off, then server i begins its setup process, and if the number of jobs in the system is less than $k_{i,j}$ and server i is in setup, then it is switched off. It is worth noting that due to the homogeneity of the servers, from the subscript i and j one can infer how many servers are also in setup. As an example, if the

system is turning on its third server while one server is already on, then the second server must also be in setup. We would argue that the properties of Class A policies are intuitively appealing, and moreover, from [12] it is known that optimal policies are contained within the set of Class A policies.

Before the second class of policies is given, another definition must first be introduced. Because this work examines energy-aware systems as $C \rightarrow \infty$ it may be the case that while a policy gives a criterion to turn on some server s , the probability of this criterion being met approaches 0. Therefore to reason about these cases, and others, the following framework is introduced. Let $X_{\mathcal{E}}(s, t)$ be an indicator variable such that

$$X_{\mathcal{E}}(s, t) = \begin{cases} 1, & \text{if server } s \text{ is in energy state } \mathcal{E} \text{ at time } t, \\ 0, & \text{otherwise,} \end{cases}$$

where $\mathcal{E} \in \{\text{off}, \text{setup}, \text{idle}, \text{busy}\}$. Then it is said that s is an always \mathcal{E} server if and only if as $t \rightarrow \infty$, $P(X_{\mathcal{E}}(s, t) = 1) \rightarrow 1$. As an example, if a server s has a criterion which turns it on and it is known that the server will always eventually turn off, but the probability that the turn on criterion is met approaches 0 as $t \rightarrow \infty$, s would be called an always off server, since as $t \rightarrow \infty$, $P(X_{\text{off}}(s, t) = 1) \rightarrow 1$. With these notions in mind the second class of policies is defined as follows.

Definition 3. Class B Policy: A policy is said to be a Class B policy if the following conditions are met:

1. It is a Class A policy.
2. There exists an $\alpha < 1$ such that the number of always idle servers is less than $(1 - \rho)C^\alpha$.
3. If a server s turns off when there are n_1 servers on and n_2 jobs in the system, then while there are at least n_1 servers on, s will not begin its setup until there are at least $n_2 + 1$ jobs in the system.

The second condition for Class B policies states that the number of servers which are always idle cannot be on the same order as the total number of servers. The third condition protects against known suboptimal behaviour, see Theorem 3 of [12]. That is, it is never the case that when a server switches off it immediately begins its setup process. It is worth noting that most policies studied in the literature are Class B policies, i.e. the policies of focus in [4–7, 11, 13, 16–18, 21] are all Class B policies. The sets of Class A and Class B policies are denoted by Π_A and Π_B respectively, and furthermore, for a specific policy π , $\mathbb{E}[R^\pi]$ and $\mathbb{E}[E^\pi]$ denote the expected response time and expected energy costs under policy π , respectively.

3 Main Results

This work examines the behaviour of energy-aware systems under a fixed-load, many-server asymptotic regime. That is, for an energy-aware system $S = (C, \lambda, \mu, \gamma)$, the metrics $\mathbb{E}[R]$ and $\mathbb{E}[E]$ are evaluated as $C \rightarrow \infty$ while $\rho = \lambda/(C\mu)$ is held constant. Formally, we consider a sequence of energy-aware systems, indexed by n , where each system has a fixed ρ , and as $n \rightarrow \infty$, $C \rightarrow \infty$. All proofs of the results are given in Section 4.

Theorem 1. All policies in Π_A are asymptotically optimal with regards to expected response time. In other words, given an energy-aware system, for any $\pi_a \in \Pi_A$, as $\lambda, C \rightarrow \infty$ and $\lambda/\mu C$ is fixed to be ρ , where $0 < \rho < 1$, $\mathbb{E}[R^{\pi_a}] \rightarrow 1/\mu$.

While perhaps surprising at first, such a result becomes intuitive as one considers the details of the system behaviour. Informally, there is a significant proportion of jobs which are served immediately on arrival, and therefore a significant proportion of jobs have a response time equal to their service time. And while it is true that some jobs will have to wait to be served, whether it be for a server to complete a job or finish a setup, the number of these jobs turns out to be negligible under the asymptotic regime. It is worth noting that Theorem 1 would not necessarily hold for policies which turned servers off while there are waiting jobs that they could process. On the other hand, belonging to Π_A is not necessary for minimizing the expected response time. With optimal policies now known for $\mathbb{E}[R]$, our focus shifts to the second cost metric, $\mathbb{E}[E]$.

Theorem 2. All policies in Π_B are asymptotically optimal with regards to expected energy cost. In other words, given an energy-aware system, for any $\pi_b \in \Pi_B$, as $\lambda, C \rightarrow \infty$ and $\lambda/\mu C$ is fixed to be ρ , where $0 < \rho < 1$, $\mathbb{E}[E^{\pi_b}]/\lambda \rightarrow \mathbb{E}[E^{J, \pi_b}] \rightarrow E_{\text{Busy}}/\mu$.

Corollary 1. *All Class B policies are asymptotically optimal under any well-formed cost function.*

The optimality result for the expected energy cost is arguably more surprising than the result for the expected response time. One may have the intuition that some of these policies would regularly have an infinite number of servers in a specific energy state, such as setup, which would in turn incur an infinite amount more cost than some other policy. For example, one policy may regularly have an infinite number of servers in setup, while another may instead have an infinite number of servers idle. Considering these two policies, it may be fair to think one would incur infinitely more cost than the other, which makes the notion of these policies being equivalent under all cost functions difficult to wrap one’s mind around. While this line of thinking is not contradicted by Theorem 2, it does say that focusing on details of servers which do spend a certain amount of time in setup or idling is a misleading way to think about these systems. In other words, under the asymptotic regime, servers which find themselves in setup or idle are negligible when compared to those which spend all of their time busy. In fact, as will be seen in the Sections 3.1 and 4, reasoning about the energy cost from the perspective of the servers, as was seen in (1), results in complexities which are washed away when the energy costs are viewed from the perspective of the jobs, i.e. $\mathbb{E}[E] = \lambda\mathbb{E}[E^J]$. Once these observations are made, these seemingly complex systems become simple to reason about.

The most significant implication of Theorems 1 and 2 is that under the asymptotic regime the trade-off between $\mathbb{E}[R]$ and $\mathbb{E}[E]$ is in fact not a trade-off at all. That is, not only are both cost metrics minimized across a large set of policies, but over all well-formed cost functions. This is a powerful result, since if a system is *close* to this asymptotic regime, then a manager can confidently employ a Class B policy knowing that it will be reasonably close to optimal. Of course this begs the question, what does it mean for a system to be *close* to the system of study? This work addresses this question by numerically inspecting energy-aware systems with a finite number of servers C , to see how quickly the cost metrics approach their bounds described in Theorems 1 and 2.

3.1 Numerical Experiments

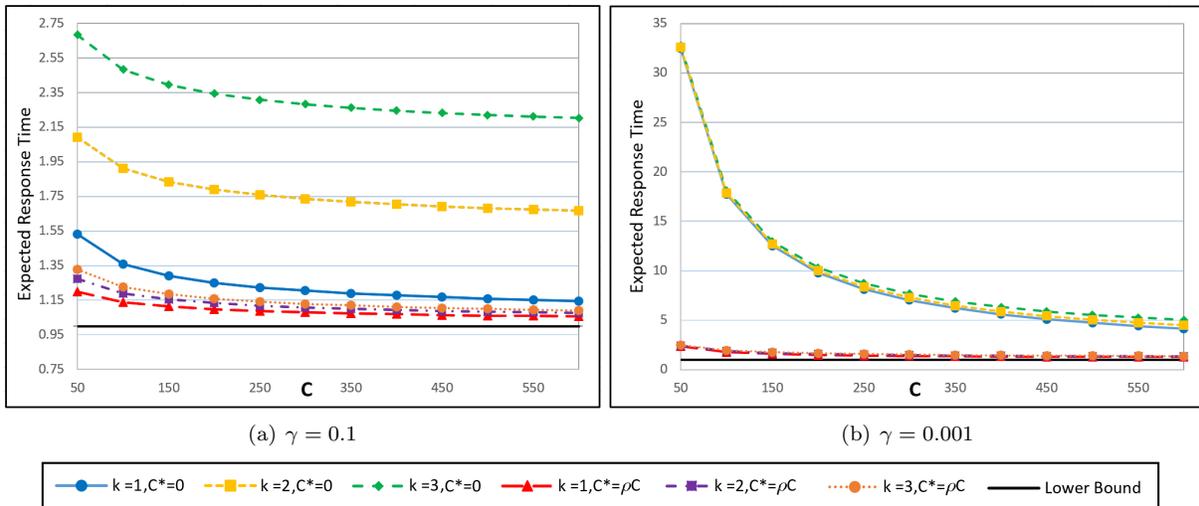


Figure 1: Expected response time vs C for $\lambda = C/2, \mu = 1$

All numerical experiments presented here are done for an energy-aware system employing a *staggered threshold* policy, with a specific instantiation of its decision variables k , and C^* . A brief description of this policy is as follows. Regardless of the system state, C^* of the servers always remain on, the remaining $(C - C^*)$ servers will turn off the moment they idle, and the number of servers in setup is determined by the threshold value k , which equals $\{[\{j - C^*\}^+ / k] - i\}^+$, where $C^* + i$ is the number of servers currently on, and j is the number of jobs in the system. Informally, the greater the value of k , the more jobs are required to accumulate before a server will begin its setup. A more in depth examination and analysis of this policy can be found in [13]. It is worth noting that by appropriately choosing the decision variables, other studied

policies can be instantiated. As an example, letting $k = 1$ and $C^* = 0$ results in the staggered setup policy mentioned in Section 1.

Letting the system state space be (i, j) , where i is the number of servers on (idle or busy), and j is the number of jobs in the system, allows for the underlying CTMC to be expressed as a quasi birth-death process. This can be analysed using any of a number of well understood methods. We chose to analyse the CTMCs using the RRR technique described in [4]. Therefore, all numerical results are exact and were evaluated using standard Matlab libraries. The source code for the numerical analysis can be found at [1]. The purpose of these numerical experiments is firstly to ensure that exact analysis agrees with our results pertaining to the system under the asymptotic regime, and secondly to examine how quickly the system approaches the corresponding optimal behaviour as the parameters are appropriately scaled up.

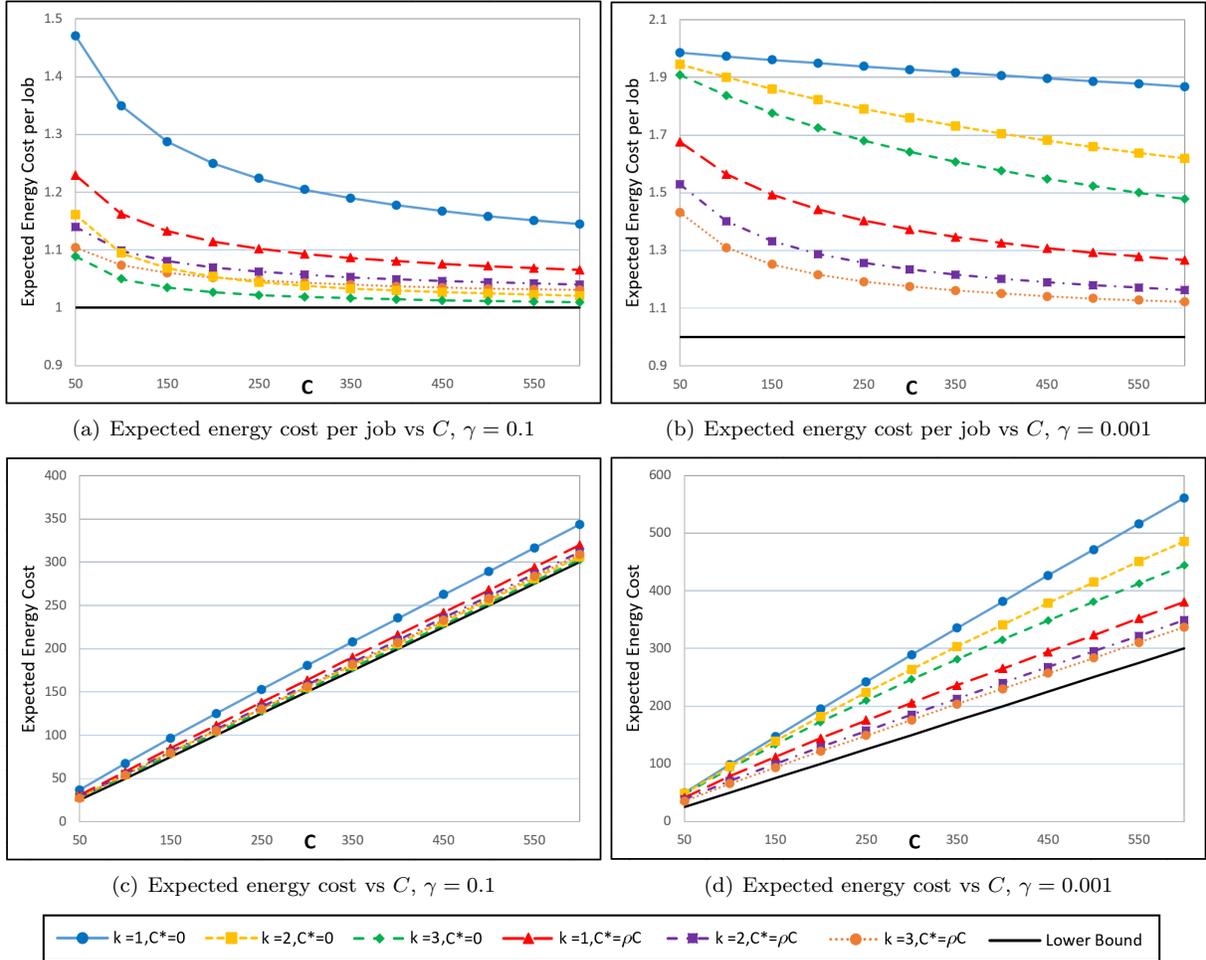


Figure 2: Expected energy cost and Expected energy cost per job vs C for $\lambda = C/2$, $\mu = 1$

Figure 1 shows the behaviour of $\mathbb{E}[R]$ as the system is scaled up. A preliminary observation is that for the curves where $C^* = 0$, the corresponding value of $\mathbb{E}[R]$ can be far from optimal even for large values of C . As an example, the curve where $k = 3$ and $C^* = 0$ in Figure 1-(a) is more than double that of the optimal value even for the largest values of C analysed. And perhaps even worse than that, the curves where $C^* = 0$ have extremely slow convergence rates. On the other hand, one may also note that when $C^* = \rho C$, $\mathbb{E}[R]$ becomes reasonably close to its optimal value relatively quickly. This effect is accentuated further in 1-(b), where the setup times are large. Here, all curves which share the same choice of C^* are visually grouped together, and moreover, when $C^* = 0$ the expected response time can be far from optimal even for larger values of C . But the curves which have a number of servers which are forced to be on, i.e. $C^* = \rho C$, gets much closer to the minimum value. In other words, the convergence rate is sensitive to the choice of C^* ,

while relatively insensitive to the threshold value k , especially when setup times are large. And forcing the system to have a number of servers always on, where the number always on equals ρC , causes the system to more quickly approach its minimum expected response time.

The appealing choice of forcing $\lambda/\mu = \rho C$ servers to always remain on is interesting, since as will be seen in Section 4, the number of servers which are always busy approaches $\lambda/\mu = \rho C$ under the asymptotic regime. In other words, when the system parameters are finite, setting $C^* = \rho C$ forces the system to behave in a manner in which it is known to behave under the asymptotic regime. As such, it is intuitive that when the system is constrained to invoke certain asymptotic behaviour, i.e. $C^* = \rho C$, the corresponding values of $\mathbb{E}[R]$ are closer to values which would be seen under the asymptotic regime.

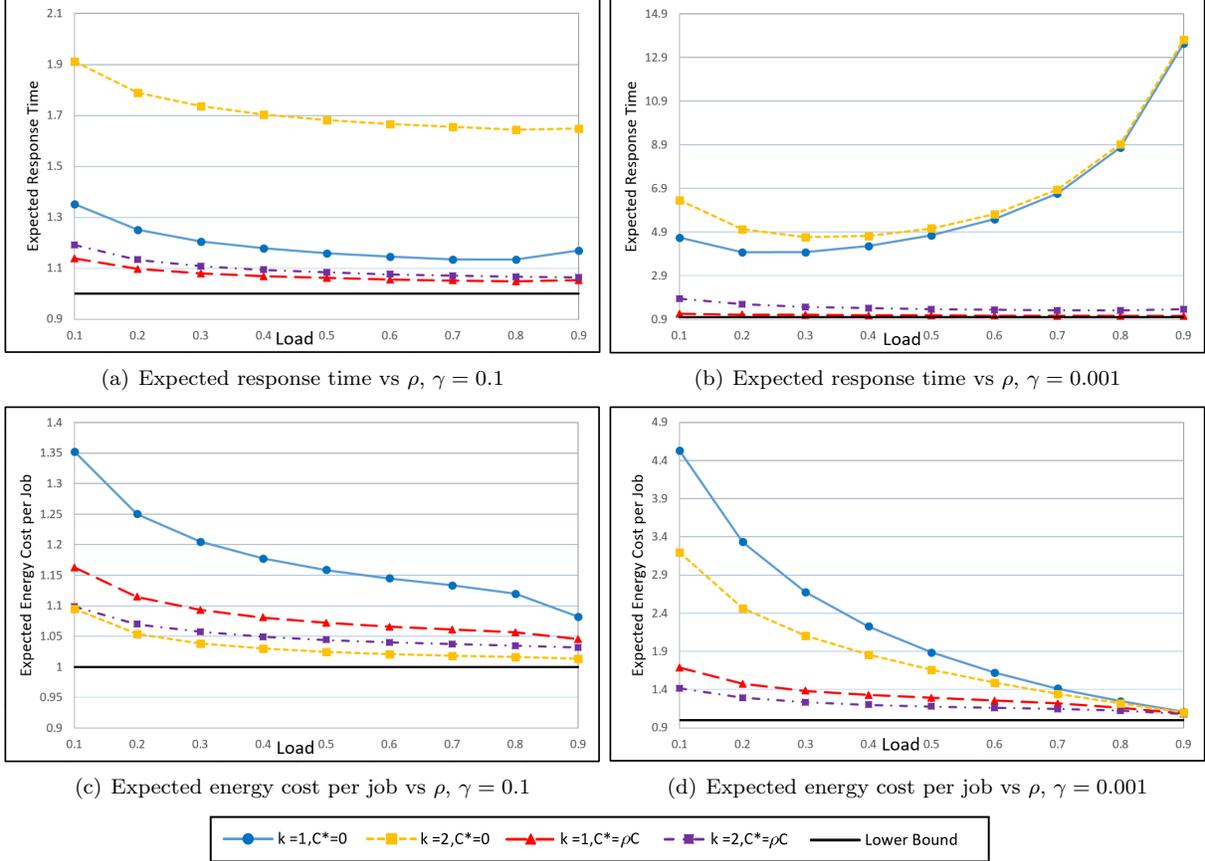


Figure 3: Expected response time and Expected energy cost per job vs ρ for $C = 500$, $\mu = 1$

Shifting focus to the expected energy costs per job and Figures 2 (a) and (b), a similar trend regarding the choice of C^* is seen. That is, when C^* is forced to take on the value it practically approaches under the asymptotic regime, $\mathbb{E}[E^J]$ gets closer to its optimal value. Here plotting the expected energy cost on a per job basis allows one to clearly see how the system is becoming more efficient as the parameters are scaled up. This is in contrast to viewing the total expected energy cost as seen in Figure 2 (c) and (d). Here one in fact sees that the curves are diverging from their lower bound. At first this may seem to be a direct contradiction of Theorem 2 due to the observed divergence, but this is not the case. While it is true that the difference between the expected energy costs and the optimal value is getting larger, it is important to remember that the lower bound itself is also growing with the system parameters, specifically λ . In other words, as $C \rightarrow \infty$ the difference between the energy costs and lower bound is growing to infinity, but is growing at a sub-linear rate. At the same time, the lower bound is growing to infinity at a linear rate. Therefore, the difference between the expected energy cost and the lower bound, while infinite, becomes negligible when compared to the total cost. One of the key insights that looking at the per job expectations of the energy cost grants under this scaling, is that a divergence from the lower bound of the total expected energy cost does not

imply suboptimality, specifically under the asymptotic regime.

Another aspect of these systems which warrants attention, is how sensitive the convergence rate is to the load. This is seen in Figure 3. Examining the load's effect on the expected response time, one can observe that when the setup times are relatively short, the convergence rate is relatively insensitive to choice of load. Furthermore, the most sensitive parts of the curves are when the load is light or heavy, especially in Figure 3 (b) where the setup times are longer. This makes some intuitive sense, since when the loads are light or heavy the system is more likely to exhibit behaviour that is not described by the asymptotic regime. When the load is light, the system has a significant chance to be empty, and in turn has a significant chance to have the minimum number of servers on. When jobs arrive it begins to overcompensate with more setups than are needed and the servers begin to thrash. On the other hand when the system load is high there is a significant chance that there will be more than C jobs in the system. So even if all servers were on, jobs would still have to wait. Having many servers regularly thrashing, or having more jobs in the system than servers are two characteristics which are never exhibited by a system under the asymptotic regime. Therefore, it is intuitive that a system under light or heavy loads would be slower to exhibit asymptotic behaviour than that of a system with a medium load.

With this sensitivity in mind, one can still clearly note that the previous observation regarding having ρC servers always on induces the asymptotic behaviour to occur sooner. The only curve not to agree with this notion is the case where $k = 2$ and $C^* = 0$ in Figure 3 (c). In this case the system is approaching the minimum value $\mathbb{E}[E^J]$ slightly quicker than the curves where $C^* = \rho C$. This is a product of the servers thrashing, causing most jobs to see the system when many other jobs are present and therefore little energy is wasted, but is only achieved with a large sacrifice to $\mathbb{E}[R]$, and therefore this configuration would not be suggested.

4 Proofs

Here the proofs of Theorems 1 and 2 are presented in detail. Before this can be done however, some preliminary material, which is used throughout the proofs, must first be defined and presented. Consider the following three sequences of systems with $0 < \rho < 1$:

1. Let S_1 be a sequence of n energy-aware queueing systems, where the i th energy-aware queueing system is given by $S_{1,i} = (C_{1,i}, \lambda_{1,i}, \mu_{1,i}, \gamma_{1,i})$, which employ some policy $\pi_i \in \Pi_A$, where $(\forall i \text{ s.t. } 0 < i \leq n : \mu_{1,i} = 1, C_{1,i} = i, \text{ and } \lambda_{1,i}/C_{1,i} = \rho)$, and as $n \rightarrow \infty, \lambda_{1,n} \rightarrow \infty$.
2. Let S_2 be a sequence of n $M/M/C$ queues, where the i th $M/M/C$ queue is denoted by $S_{2,i} = (C_{2,i}, \lambda_{2,i}, \mu_{2,i})$, where $(\forall i \text{ s.t. } 0 < i \leq n : \mu_{2,i} = 1, C_{2,i} = i, \text{ and } \lambda_{2,i}/C_{2,i} = \rho)$, and as $n \rightarrow \infty, \lambda_{2,n} \rightarrow \infty$.
3. Let S_3 be a sequence of n $M/M/C$ queues, where the i th $M/M/C$ queue is denoted by $S_{3,i} = (C_{3,i}, \lambda_{3,i}, \mu_{3,i})$, where $(\forall i \text{ s.t. } 0 < i \leq n : \mu_{3,i} = 1 \text{ and } C_{3,i} = \lambda_{3,i} + (\lambda_{3,i})^{0.5+\epsilon})$, where $0 < \epsilon < 0.5$, and as $n \rightarrow \infty, \lambda_{3,n} \rightarrow \infty$.

In order to compare and reason about these systems, let it also hold that: $(\forall i \text{ s.t. } 0 < i \leq n : \lambda_{1,i} = \lambda_{2,i} = \lambda_{3,i} = \lambda_i)$. Note that imposing such a constraint implies that $(\forall i \text{ s.t. } 0 < i \leq n : C_{1,i} = C_{2,i} = C_i)$. In other words, the arrival rate of the systems across all three sequences are equal, and the total number of servers in the systems of sequences 1 and 2 are also equal.

Let $B_{1,i}, B_{2,i}$, and $B_{3,i}$ denote the number of always busy servers in the i th system of sequence S_1, S_2 , and S_3 respectively. This work now provides a Lemma which is key to understanding the behaviours of these energy-aware systems under the asymptotic regime.

Lemma 1. *Given a sequence of energy-aware systems where each system employs a policy $\pi \in \Pi_A$, C_n, λ_n , and μ_n denote the number of servers, arrival rate, and processing rate of the n th system respectively, and as $n \rightarrow \infty, \lambda_n, C_n \rightarrow \infty$ while $\lambda_n/(\mu_n C_n)$ is fixed to some ρ , where $0 < \rho < 1$, it holds that for the number of always busy servers in the n th system denoted by B_n ,*

$$\lim_{n \rightarrow \infty} \frac{B_n}{\lambda} = \frac{1}{\mu}.$$

Proof. Without loss of generality one can set $\mu = 1$. Therefore, to show Lemma 1, it is equivalent to prove

$$\lim_{n \rightarrow \infty} \frac{B_{1,n}}{\lambda_n} = 1.$$

This is done via a sample path argument regarding the sequences S_1 and S_2 . Consider the systems $S_{1,n}$ and $S_{2,n}$ as $n \rightarrow \infty$. At any point in time the number of servers currently available in $S_{1,n}$ is less than or equal to the number of servers available in $S_{2,n}$. This follows from the fact that $S_{1,n}$ may have some of its C_n servers off or in setup, while $S_{2,n}$ has C_n servers on at all times. Therefore, taking the same arrival stream and job sizes for both systems, the number of jobs in $S_{2,n}$ is less than or equal to the number of jobs in $S_{1,n}$. Therefore, if s is busy in $S_{2,n}$, then s has enough workload to also be busy in $S_{1,n}$, but may not be busy due to it being switched off or in setup. Therefore, if s is an always busy server in $S_{2,n}$, then s has enough workload to be always busy in $S_{1,n}$. However, as $S_{1,n}$ is employing a policy from Π_A , specifically, from the second condition in the definition of Class A policies it is known a server will never turn off if there is work to do and from the first condition a server will eventually turn on from the threshold scheme, it follows that almost surely the servers which can be always busy, will be always busy. That is to say, if s is an always busy server in $S_{2,n}$, then s is an always busy server in $S_{1,n}$. Therefore,

$$\lim_{n \rightarrow \infty} B_{1,n} \geq \lim_{n \rightarrow \infty} B_{2,n}. \quad (2)$$

Furthermore, it is known that

$$\lim_{n \rightarrow \infty} \frac{B_{2,n}}{\lambda_n} = 1.$$

This is shown via the following argument. Let $N_{2,i}(t)$ denote the number of jobs in the system $S_{2,i}$ at time t , and let a corresponding diffusion scaling be denoted by $\hat{N}_{2,n}(t)$, where

$$\hat{N}_{2,n}(t) = \frac{N_{2,n}(t) - n\rho}{\sqrt{n\rho}}.$$

From Theorem 4.1 in [10] it can be seen that as $n \rightarrow \infty$, $\hat{N}_{2,n}(t)$ weakly converges to an Ornstein-Uhlenbeck process, and therefore, at each time point t as $n \rightarrow \infty$, $\hat{N}_{2,n}(t)$ is normally distributed. After some elementary algebra,

$$N_{2,n}(t) = \hat{N}_{2,n}(t)\sqrt{\lambda_n} + \lambda_n \quad \Rightarrow \quad \frac{N_{2,n}(t)}{\lambda_n} = \frac{\hat{N}_{2,n}(t)}{\sqrt{\lambda_n}} + 1.$$

Because as $n \rightarrow \infty$, $\hat{N}_{2,n}(t)$ is normally distributed with finite mean and variance, $\lim_{n \rightarrow \infty} \hat{N}_{2,n}(t)/\sqrt{\lambda_n} = 0$. This immediately implies,

$$\frac{N_{2,n}(t)}{\lambda_n} = 1.$$

Moreover, one can say that as $n \rightarrow \infty$ if there are almost surely at least x jobs in the system at all time points t , then as $n \rightarrow \infty$ there are at least x always busy servers at all time points t . Therefore,

$$\lim_{n \rightarrow \infty} \frac{N_n(t)}{\lambda_n} = 1 \quad \Rightarrow \quad \lim_{n \rightarrow \infty} \frac{B_{2,n}}{\lambda_n} \geq 1.$$

After realizing this property of S_2 , one can begin examining the implication on the behaviour of S_1 . Specifically from (2) it is known,

$$\lim_{n \rightarrow \infty} \frac{B_{2,n}}{\lambda_n} \geq 1 \quad \Rightarrow \quad \lim_{n \rightarrow \infty} \frac{B_{1,n}}{\lambda_n} \geq 1.$$

Hence all that remains to show Lemma 1 is to prove

$$\lim_{n \rightarrow \infty} \frac{B_{1,n}}{\lambda_n} \leq 1.$$

This follows immediately however, after the observation that

$$\lim_{n \rightarrow \infty} \frac{B_{1,n}}{\lambda_n} > 1$$

is true only if the arrival rate of the system is greater than λ_n , which is a direct contradiction to the system definition. Therefore,

$$\lim_{n \rightarrow \infty} \frac{B_{1,n}}{\lambda_n} \geq 1, \text{ alongside } \lim_{n \rightarrow \infty} \frac{B_{1,n}}{\lambda_n} \leq 1$$

implies,

$$\lim_{n \rightarrow \infty} \frac{B_{1,n}}{\lambda_n} = 1.$$

□

4.1 Proof of Theorem 1

For simplicity of navigation, Theorem 1 is restated.

Theorem 1. *All policies in Π_A are asymptotically optimal with regards to expected response time. In other words, given an energy-aware system, for any $\pi_a \in \Pi_A$, as $\lambda, C \rightarrow \infty$ and $\lambda/\mu C$ is fixed to be ρ , where $0 < \rho < 1$, $\mathbb{E}[R^{\pi_a}] \rightarrow 1/\mu$.*

A high-level description of the proof is as follows. It is determined that if a job J is served by an always busy server in the system $S_{3,n}$ as $n \rightarrow \infty$, then the expected response time of job J approaches its expected service time. With this in mind, the system $S_{1,n}$ is compared to $S_{3,n}$ as $n \rightarrow \infty$. It becomes clear that the expected response time of these systems is dominated by jobs which are served by always busy servers. Moreover, the limits of the expected response time of these two systems are equal. Therefore, the expected response time of $S_{1,n}$ approaches the expected service time, as $n \rightarrow \infty$.

Proof. As with the proof of Lemma 1, without loss of generality one can set $\mu = 1$. Therefore, to prove the theorem it is enough to show that $\lim_{n \rightarrow \infty} \mathbb{E}[R_{1,n}] = 1$, where $\mathbb{E}[R_{i,j}]$ denotes the expected response time corresponding to the system at the j th index of the i th sequence of systems. The equality of this limit is shown via a sample path argument regarding S_1 and S_3 . Before this argument is made however, some properties of S_3 must be shown.

Consider the system $S_{3,n}$. From the square root staffing rule, Theorem 15.2 of [8], it is known that for a system with C servers, where $C = \lambda_n + c\sqrt{\lambda_n}$, the probability of queueing gets arbitrarily close to 0 as c gets large. Keeping in mind the definition of S_3 , specifically that $C_{3,i} = \lambda_i + (\lambda_i)^{0.5+\epsilon}$, it can be said that for all finite c ,

$$\lim_{n \rightarrow \infty} \frac{c\sqrt{\lambda_n}}{\lambda_n^{0.5+\epsilon}} = \frac{c}{\lambda_n^\epsilon} = 0,$$

and therefore in the system $S_{3,n}$ as $n \rightarrow \infty$, $P(N_{3,n} > C) \rightarrow 0$, where $N_{3,n}$ is a random variable denoting the number of jobs in $S_{3,n}$, which implies,

$$\lim_{n \rightarrow \infty} \mathbb{E}[R_{3,n}] = \frac{1}{\mu_{3,n}} = 1. \quad (3)$$

Moreover, identical to the reasoning presented in Lemma 1, one can conclude

$$\lim_{n \rightarrow \infty} \frac{B_{3,n}}{\lambda_n} = 1. \quad (4)$$

From here the server pool is split into two distinct conceptual sets. The first set consists of the always busy servers, and the second set consists of all remaining servers, i.e. the servers which spend some of their time idle. Let these two sets of servers be denoted by $A_{3,i}$ and $I_{3,i}$, respectively. Then the expected response time can be expressed as a sum of terms,

$$E[R_{3,n}] = P_{A_{3,n}} \mathbb{E}[R_{3,n}^A] + P_{I_{3,n}} \mathbb{E}[R_{3,n}^I], \quad (5)$$

where $P_{A_{3,n}}$ and $P_{I_{3,n}}$ denote the probability of a job being served from set $A_{3,n}$ or $I_{3,n}$, respectively, and $E[R_{3,n}^A]$ and $E[R_{3,n}^I]$ denote the expected response times given that a job is served by a server from set $A_{3,n}$ or $I_{3,n}$, respectively. Due to the homogeneity of the servers, it is assumed that if there are c servers currently

on and $i \leq c$ jobs in the system, then servers one through i will be the servers which are busy. From this it can be noted that always busy servers, i.e. servers from set $A_{3,n}$, have serving priority over the others when jobs arrive. That is, it is known that the probability of a particular job being served by an always busy server is greater than or equal to choosing it randomly and uniformly from the entire pool. In other words,

$$P_{A_{3,n}} \geq \frac{|A_{3,n}|}{C_{3,n}} = \frac{|A_{3,n}|}{\lambda_n + \lambda_n^{0.5+\epsilon}}.$$

Furthermore, from the definition of $A_{3,n}$ it is known that $|A_{3,n}| = B_{3,n}$, which implies,

$$\lim_{n \rightarrow \infty} P_{A_{3,n}} \geq \lim_{n \rightarrow \infty} \frac{B_{3,n}}{\lambda_n + \lambda_n^{0.5+\epsilon}} = \lim_{n \rightarrow \infty} \frac{B_{3,n}}{\lambda_n} = 1.$$

Noting $P_{I_{3,n}} = 1 - P_{A_{3,n}}$ alongside (3) and (5) it becomes clear that,

$$\lim_{n \rightarrow \infty} E[R_{3,n}] = \lim_{n \rightarrow \infty} \mathbb{E}[R_{3,n}^A] = 1.$$

With the limit of $\mathbb{E}[R_{3,n}^A]$ explicitly solved, the proof proceeds with a sample path argument involving S_1 and S_3 . As $A_{3,n}$ and $I_{3,n}$ denote the sets of always busy and sometimes idle servers respectively for $S_{3,n}$, let $A_{1,n}$ and $I_{1,n}$ denote the corresponding always busy and sometimes idle sets for $S_{1,n}$.

It is worth noting that while calling $I_{3,n}$ a set of sometimes idle servers is apt, applying the same notion to $I_{1,n}$ is somewhat inappropriate. $I_{1,n}$ is the set of servers which are not always busy. These servers potentially could spend zero time idling, i.e. they immediately switch off when they complete a job and the queue is empty. However, for the purposes of the proof it is only required that these servers spend some portion of time not busy.

Continuing with the sample path argument, consider systems $S_{1,n}$ and $S_{3,n}$, as $n \rightarrow \infty$, with identical arrival processes. Jobs are viewed as being marked to be served by a server belonging to the sets $A_{1,n}$ and $A_{3,n}$, or $I_{1,n}$ and $I_{3,n}$, without loss of generality. If a job arrives to $S_{3,n}$ and is served from the set $A_{3,n}$, then the corresponding job in $S_{1,n}$ will almost surely be served from the set $A_{1,n}$. This is the case by noting from (4) and Lemma 1 that

$$\lim_{n \rightarrow \infty} \frac{B_{1,n}}{B_{3,n}} = 1. \quad (6)$$

Moreover, if a job arrives to $S_{3,n}$ and is served from the set $I_{3,n}$, then the corresponding job in $S_{1,n}$ will almost surely be served from the set $I_{1,n}$. Firstly, it must be seen that $I_{1,n}$ has the capacity to handle the load which $I_{3,n}$ has the capacity to handle. This follows immediately by noting that

$$\lim_{n \rightarrow \infty} \frac{|I_{3,n}|}{\lambda_n} = 0 \text{ while } \lim_{n \rightarrow \infty} \frac{|I_{1,n}|}{C} = (1 - \rho),$$

which implies,

$$\lim_{n \rightarrow \infty} \frac{|I_{3,n}|}{|I_{1,n}|} = 0.$$

As was done with $S_{3,n}$, one can decompose the expected response time of $S_{1,n}$ into two distinct components as follows,

$$E[R_{1,n}] = P_{A_{1,n}} \mathbb{E}[R_{1,n}^A] + P_{I_{1,n}} \mathbb{E}[R_{1,n}^I].$$

From here two important observations are made. Firstly because $S_{1,n}$ is a stable system $\mathbb{E}[R_{1,n}^I]$ is known to be finite. Secondly, because $S_{1,n}$ and $S_{3,n}$ have identical arrival processes, it can be said

$$\lim_{n \rightarrow \infty} P_{A_{1,n}} = \lim_{n \rightarrow \infty} P_{A_{3,n}} = 1$$

which alongside (6) implies,

$$\lim_{n \rightarrow \infty} \mathbb{E}[R_{1,n}^A] = \lim_{n \rightarrow \infty} \mathbb{E}[R_{3,n}^A] = 1.$$

Therefore,

$$\lim_{n \rightarrow \infty} E[R_{1,n}] = \lim_{n \rightarrow \infty} \mathbb{E}[R_{1,n}^A] = 1.$$

□

4.2 Proof of Theorem 2

For simplicity of navigation, Theorem 2 is restated.

Theorem 2. *All policies in Π_B are asymptotically optimal with regards to expected energy cost. In other words, given an energy-aware system, for any $\pi_b \in \Pi_B$, as $\lambda, C \rightarrow \infty$ and $\lambda/\mu C$ is fixed to be ρ , where $0 < \rho < 1$, $\mathbb{E}[E^{\pi_b}]/\lambda \rightarrow \mathbb{E}[E^{J, \pi_b}] \rightarrow E_{Busy}/\mu$.*

A high-level description of the proof is as follows. The notion of the energy cost contributed by a single job is examined. Specifically, the energy cost of a single job is examined based on the criterion of it being served by an always busy server or not. Lemma 2 gives an exact value of the energy cost of a single job assuming it was served by an always busy server. Moreover, this is a minimum value. On the other hand, Lemma 3 shows if a job is not served by an always busy server, the energy cost is finite. From there, similar to the procedure in the proof of Theorem 1, it becomes clear that the total expected energy cost is dominated by jobs which are served by always busy servers, and therefore is minimized.

Definition 4. Energy Cost of a Job: *Let E^J be a random variable which denotes the energy cost contributed by a randomly chosen job J . There are four contributing factors to consider when determining E^J for some job J .*

1. Each job J is responsible for the energy required to process it.
2. If a job J is the first job which server s serves after completing its setup process, then it is said that J is responsible for contributing the entire cost of the setup process of s , as well as any idling costs of s until the next time s is switched off.
3. If a job J is responsible for causing a server setup which is canceled due to that job entering service before the setup process completes, then J is also responsible for the energy cost incurred by that setup.
4. The idling cost of servers which never turn off is divided evenly among all jobs that pass through the system. Furthermore, a special case is added to this factor, which is if some of the servers are always busy servers, i.e. some of the system parameters are approaching infinity, then the aforementioned idling cost is evenly distributed only across jobs which are served by always busy servers, rather than across all jobs which pass through the system.

As mentioned in Section 2 and from the definition of $[E^J]$, it is clear that

$$E[E] = \lambda \mathbb{E}[E^J]. \quad (7)$$

Therefore, when $\mathbb{E}[E^J]$ is minimized, $\mathbb{E}[E]$ is minimized.

Lemma 2. *In an energy-aware queueing system employing a Class B policy, if a job J is served by an always busy server, then $\mathbb{E}[E^J] = E_{Busy}/\mu$.*

Proof. The proof of this Lemma is argued after several key observations. To prove the Lemma, it is equivalent to show that if J is served by an always busy server, then $\lim_{n \rightarrow \infty} \mathbb{E}[E_{1,n}^J] = E_{Busy}$. Moreover, for there to be any always busy servers present in the system, it is required that $n \rightarrow \infty$. Therefore, it will be argued that

$$\lim_{n \rightarrow \infty} \mathbb{E}[E_{1,n}^J \mid J \text{ was served by an always busy server}] = E_{Busy}.$$

To show the above equality, the four contributing factors to $\mathbb{E}[E_{1,n}^J]$, from Definition 4, are addressed individually and summed.

1. It is known that each job J will eventually be served, by an always busy server or otherwise, and therefore J incurs an expected cost of $E_{Busy}/\mu_{1,i} = E_{Busy}$, for all i .
2. If it is known that J is served by an always busy server, then it is trivially known that it is not the first job to be processed after a server completes its setup process. Therefore it can be said that J incurs no cost from this contributing factor.

3. The third contributing factor takes a little more care, but can be shown to also incur no cost via a contradiction argument. Assume that some job J is responsible for causing a server to start its setup process. If this is the case however, it is known that such a server will almost surely never cancel its setup nor turn off from the structural property, using Theorem 3 of [12], which Class B policies adhere to. Therefore, it cannot be the case that J caused a setup which was then canceled. Moreover, if J did cause some server s to start a setup process, then s would be idle an infinite amount of time before turning off, a direct violation of Class B policies. Taking this all into account, it is clear that if J is served by an always busy server, then it incurs no cost from starting an eventually canceled setup.
4. Because $n \rightarrow \infty$ the special case of this factor is invoked i.e. jobs which are served by always busy servers are responsible for the full cost of the always idle servers. From the definition of Class B policies it is known that the number of always idle servers is less than $(1 - \lambda_{1,n}/C_{1,n})C_{1,n}^\alpha$, where $0 \leq \alpha < 1$. Therefore, costs from these idle servers are incurred at some rate less than $(1 - \rho)C_{1,n}^\alpha E_{\text{Idle}}$. Moreover, from the proof of Theorem 1, it is known that as $n \rightarrow \infty$, the probability of being served by an always busy server approaches 1. It then follows that the rate at which jobs are served by always busy servers approaches the arrival rate, $\lambda_{1,n}$, as $n \rightarrow \infty$. Therefore, letting the expected contributing cost from these idle servers per job be denoted by $\mathbb{E}[E^{J,I}]$ and from the definition of S_1 knowing that $\lambda_{1,n} = C_{1,n}/\rho$:

$$\mathbb{E}[E_{1,n}^{J,I}] = \lim_{n \rightarrow \infty} \frac{(1 - \rho)C_{1,n}^\alpha E_{\text{Idle}}}{\lambda_{1,n}} = \lim_{n \rightarrow \infty} \rho(1 - \rho)E_{\text{Idle}} \frac{C_{1,n}^\alpha}{C_{1,n}} = \lim_{n \rightarrow \infty} \rho(1 - \rho)E_{\text{Idle}} \frac{1}{C_{1,n}^{(1-\alpha)}} = 0.$$

Therefore, the only contributing factor to $\mathbb{E}[E_{1,n}^J]$ given that J has been served by an always busy server, is the cost of processing it, which implies

$$\lim_{n \rightarrow \infty} \mathbb{E}[E_{1,n}^J \mid J \text{ is served by an always busy server}] = E_{\text{Busy}}.$$

□

Lemma 3. *In an energy-aware queueing system employing a Class B policy, if a job J is not served by an always busy server, then $\mathbb{E}[E^J]$ is finite.*

Proof. Similar to the proof of Lemma 2, this proof iterates through the contributing factors of Definition 4 to show that the expectation of each factor is finite, and therefore, the expectation of their sum is finite, i.e. $\mathbb{E}[E^J]$ is finite.

1. The cost directly incurred from processing a job is trivially known to be finite, since it is directly proportional to the service time of the job.
2. If it is known that a job J has been served by a server s which regularly completes its setup process, i.e. as $t \rightarrow \infty$, $0 < \mathbb{E}[X_{\text{Setup}}(s, t)] < 1$, two cases must be considered due to Definition 3. The first and simpler case being that J is not the first job to be served by s , after s completes a setup process. Here no costs will be incurred and are trivially finite. The second and more interesting case is that J is the first job to be served by s following a setup process. Here J is responsible for the setup costs incurred by s alongside the idling costs incurred by s until the next time it shuts off. From the definition of Class B policies it is immediately known that the cost incurred by s from idling is finite. So all that remains is to show the setup cost associated with s is also finite. While at first glance this may seem trivial since it is known that the expected setup time of a server is finite, this is not quite the case. The policy could be such that when a threshold is reached not only does the setup process of one server begin, but many, potentially even all servers begin their setup process, and the next one to complete its setup is chosen as the server to use, while the remainder of the servers terminate or cancel their setups. If s was the only server used in the setup process, then the expected cost is simply E_{Setup}/γ , which is finite. But since the setup times are exponentially distributed, for all m , where m is the number of servers used in the setup process of s , the cost incurred is expected to be $mE_{\text{Setup}}/m\gamma = E_{\text{Setup}}/\gamma$, which is finite. Therefore, the expected energy cost incurred from the contributing factor of being the first job to be served after a completed setup process is finite.

3. The last contributing factor to consider is the case of job J causing servers to start their setup processes, but have them be canceled before they complete. This is similar to the previous case of the server(s) completing their setup process. If it is known that a server's setup process is interrupted before it completes, due to the underlying exponential distribution, it can be said that the expected cost incurred is less than E_{Setup}/γ . Moreover, if it is known that the setup processes of m servers are interrupted before any of them complete, due to the underlying exponential distribution it can be said that the total cost incurred is less than $mE_{\text{Setup}}/m\gamma = E_{\text{Setup}}/\gamma$. Therefore, the cost incurred by J from causing server setups which are eventually interrupted is finite.
4. Because there are always busy servers present in the system, from the definition of E^J it is trivially known that zero costs are incurred by this factor, which is finite.

All of the contributing terms of $\mathbb{E}[E^J]$ are finite, therefore $\mathbb{E}[E^J]$ is also finite. \square

With the proof of the Lemmas complete, this work proceeds with the proof of Theorem 2.

Proof. To prove the theorem, it is equivalent to show $\lim_{n \rightarrow \infty} \mathbb{E}[E_{1,n}]/\lambda_n = E_{\text{Busy}}$ under the assumption that the systems of S_1 are employing Class B policies. From the definition of $\mathbb{E}[E^J]$, it is known $\mathbb{E}[E] = \lambda \mathbb{E}[E^J]$. Therefore,

$$\lim_{n \rightarrow \infty} \mathbb{E}[E_{1,n}]/\lambda_n = \lim_{n \rightarrow \infty} \mathbb{E}[E_{1,n}^J].$$

Furthermore, using the notation introduced in the proof of Theorem 1, it is also known that,

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{E}[E_{1,n}^J] &= \lim_{n \rightarrow \infty} P_{A_{1,n}}[E_{1,n}^J \mid J \text{ is served from } A_{1,n}] \\ &\quad + P_{I_{1,n}}[E_{1,n}^J \mid J \text{ is served from } I_{1,n}]. \end{aligned}$$

Leveraging past equalities allows one to simplify the above equation. From Lemma 2 it is known

$$\lim_{n \rightarrow \infty} [E_{1,n}^J \mid J \text{ is served from } A_{1,n}] = E_{\text{Busy}}.$$

From Lemma 2, $[E_{1,n}^J \mid J \text{ is served from } I_{1,n}]$ is finite, i.e. for some $L_n > 0$

$$\lim_{n \rightarrow \infty} [E_{1,n}^J \mid J \text{ is served from } I_{1,n}] = L_n$$

From the proof of Theorem 1, it is known,

$$\lim_{n \rightarrow \infty} P_{A_{1,n}} = 1 \quad \text{and} \quad \lim_{n \rightarrow \infty} P_{I_{1,n}} = 0.$$

Therefore,

$$\lim_{n \rightarrow \infty} \mathbb{E}[E_{1,n}^J] = E_{\text{Busy}}. \tag{8}$$

It is worth noting that E_{Busy} is a lower bound for $\mathbb{E}[E_{1,n}^J]$ since for the system to be stable the job must be processed. In other words, as $n \rightarrow \infty$, $\mathbb{E}[E_{1,n}^J]$ approaches its minimum value. Returning to equality (7),

$$\lim_{n \rightarrow \infty} \mathbb{E}[E_{1,n}]/\lambda_n = \lim_{n \rightarrow \infty} \mathbb{E}[E_{1,n}^J],$$

it becomes clear from (8) that,

$$\lim_{n \rightarrow \infty} \mathbb{E}[E_{1,n}]/\lambda_n = E_{\text{Busy}}.$$

Moreover, in system $S_{1,i}$, $\lambda_i E_{\text{Busy}}$ is a trivial lower bound for the expected energy cost, implying that as $n \rightarrow \infty$, $\mathbb{E}[E_{1,n}]$ is minimized. That is, the policy which $S_{1,n}$ is employing, which is a Class B policy, minimizes the expected energy cost as $n \rightarrow \infty$. \square

5 Conclusions

This work examined an established multiserver queueing model, where each server can be turned on, which takes an exponentially distributed amount of time, to improve performance, or turned off instantly to save on costs. How these servers are turned on and off define a policy, and this policy is evaluated under a cost function which takes performance and energy costs into account. A problem which arises is which policy should be employed for a given cost function, and furthermore, whether this policy will be robust enough to be a reasonable choice when evaluated under other cost functions. While the problem seems to be complex since the sets of potential policies and cost functions are infinite, as the number of servers $C \rightarrow \infty$ with a fixed system load $\rho = \lambda/(C\mu)$, a large set of policies become equivalent, i.e. Class B policies. Therefore, under this asymptotic regime the choice of which Class B policy to employ is irrelevant. Furthermore, not only are these policies equivalent, but they also simultaneously minimize $\mathbb{E}[R]$ and $\mathbb{E}[E]$. It then follows that all Class B policies will be optimal under all well-formed cost functions.

This work then numerically evaluated the staggered threshold policy to inspect how quickly these asymptotic behaviours are seen under finite parameters. These numerical results suggested that a finer-grained differentiation of Class B policies can be given such that the resulting subsets exhibit faster or slower convergence rates. Specifically, it was determined that if a static provisioning of $C\rho$ servers is enforced within the policy, then as C increases the metrics approach their corresponding minimum values more quickly than if no static provisioning was present.

Looking forward to the future of this work, an issue to address is how well these policies do under a time varying arrival rate. It is our intuition that if the arrival rate varied on a relatively long time scale, with regards to other system parameters such as the expected setup time, then the results given here may be reasonable to apply. However, questions remain which require a formal treatment. Can some of the asymptotic results be extended for a subset of Class B policies? If so, what new criteria must these policies adhere to? If not, what complication is the limiting factor in the analysis? While these questions are certainly deserving of attention, the theorems presented here regarding the optimality of all Class B policies under all well-formed cost functions allows one to confidently make overarching statements and conclusions across the problem domain.

6 Acknowledgments

This research was funded by the Natural Sciences and Engineering Research Council of Canada.

References

- [1] Source code. <http://www.cas.mcmaster.ca/~macciiov/publications.html>. Accessed: 2016-10-10.
- [2] A. Allahverdi, C. Ng, T. Cheng, and M. Y. Kovalyovc. A survey of scheduling problems with setup times or costs. *European Journal of Operational Research*, 187(3):985–1032, 2008.
- [3] Y. Chen, A. Das, W. Qin, A. Sivasubramaniam, Q. Wang, and N. Gautam. Managing server energy and operational costs in hosting centers. *SIGMETRICS Performance Evaluation Review*, 33(1):303–314, June 2005.
- [4] A. Gandhi, S. Doroudi, M. Harchol-Balter, and A. Scheller-Wolf. Exact analysis of the M/M/k/setup class of Markov chains via recursive renewal reward. In *ACM SIGMETRICS Performance Evaluation Review*, pages 153–166, 2013.
- [5] A. Gandhi, V. Gupta, M. Harchol-Balter, and M. A. Kozuch. Optimality analysis of energy-performance trade-off for server farm management. *Performance Evaluation*, 67(11):1155–1171, 2010.
- [6] A. Gandhi and M. Harchol-Balter. M/M/k with exponential setup. Technical report, Carnegie Mellon University, 2010.
- [7] A. Gandhi, M. Harchol-Balter, and I. Adan. Server farms with setup costs. *Performance Evaluation*, 67(11):1123–1138, 2010.

- [8] M. Harchol-Balter. *Performance Modeling and Design of Computer Systems: Queueing Theory in Action*. Cambridge University Press, 2013.
- [9] J. Hu and T. Phung-Duc. Power consumption analysis for data centers with independent setup times and threshold controls. In *AIP*, 2015.
- [10] D. L. Iglehart. Limiting diffusion approximations for the many server queue and the repairman problem. *Journal of Applied Probability*, 2(2):429–441, 1965.
- [11] P. J. Kuehn and M. E. Mashaly. Automatic energy efficiency management of data center resources by load-dependent server activation and sleep modes. *Ad Hoc Networks*, 25(2):497–504, 2015.
- [12] V. J. Maccio and D. G. Down. On optimal control for energy-aware queueing systems. In *27th International Teletraffic Congress (ITC 27)*, pages 98–106, 2015.
- [13] V. J. Maccio and D. G. Down. Exact analysis of energy-aware multiserver queueing systems with setup times. Technical report, McMaster University, 2016.
- [14] M. J. Magazine. On optimal control of multi-channel service systems. *Naval Research Logistics Quarterly*, 18(2):429–441, 1971.
- [15] I. Mitrani. Managing performance and power consumption in a server farm. *Annals of Operations Research*, 202(1):121–134, 2013.
- [16] T. Phung-Duc. Exact solutions for M/M/c/setup queues. *Telecommunication Systems*, pages 1–16, 2016.
- [17] T. Phung-Duc and K. Kawanishi. Energy-aware data centers with s -staggered setup and abandonment. In *International Conference on Analytical and Stochastic Modeling Techniques and Applications*, pages 269–283. Springer, 2016.
- [18] Y. Ren, T. Phung-Duc, Z. W. Yu, and J. C. Chen. Design and analysis of dynamic auto scaling algorithm (DASA) for 5G mobile networks. *arXiv preprint arXiv:1604.05803*, 2016.
- [19] J. Slegers, N. Thomas, and I. Mitrani. Dynamic server allocation for power and performance. In *SPEC International Workshop on Performance Evaluation: Metrics, Models and Benchmarks*, pages 247–261, 2008.
- [20] N. Tian and Z. G. Zhang. A two threshold vacation policy in multiserver queueing systems. *European Journal of Operational Research*, 168(1):153–163, 2006.
- [21] X. Xu and N. Tian. The M/M/c queue with (e, d) setup time. *Journal of Systems Science and Complexity*, 21(3):446–455, 2008.