

Setup Times in Multiserver Systems

Jalani K. Williams

CMU-CS-24-104

May 2024

Computer Science Department
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213

Thesis Committee:

Weina Wang, Chair

Mor Harchol-Balter

Alan Scheller-Wolf

Jamol J. Pender (Cornell)

William A. Massey (Princeton)

*Submitted in partial fulfillment of the requirements
for the degree of Doctor of Philosophy in Computer Science.*

Copyright © 2024 Jalani K. Williams

This research was supported in part by NSF CMMI-1938909, NSF CSR-1763701, NSF CNS-200773, and NSF ECCS-2145713, in addition to the Gates Millennium and GEM Fellowships.

The views and conclusions in this document are those of the author and do not necessarily reflect the opinions of the NSF or any other sponsors.

Keywords: queueing, multiserver systems, setup times, exceptional first service

To all those I hold dear, and have held me in kind.

Abstract

In many systems, servers do not turn on instantly; instead, a *setup time* must pass before a server can begin work. These “setup times” can wreak havoc on a system’s queueing; this is especially true in modern systems, where servers are regularly turned on and off as a way to reduce operating costs (energy, labor, CO_2 emissions, etc.). To design modern systems which are both efficient *and* performant, we need to understand how setup times affect queues.

Unfortunately, despite successes in understanding setup in the single server setting, setup in the multiserver setting remains poorly understood. To circumvent the main difficulty in analyzing multiserver setup, all existing results assume that setup times are memoryless, i.e. distributed Exponentially. However, in most practical settings, setup times are close to Deterministic, and the widely used Exponential-setup assumption leads to unrealistic model behavior and a dramatic underestimation of the true harm caused by setup times.

This thesis represents a comprehensive characterization of the average waiting time in a multiserver system with *Deterministic* setup times, the $M/M/k/Setup$ -Deterministic. In particular, we derive multiplicatively-tight lower and upper bounds on the average waiting time, demonstrating that **setup times, along with their distributions, can not be ignored; setup times can cause profound increases in waiting time, especially when the distribution of setup time has low variability.** Our bounds are the first closed-form bounds on waiting time in *any* many-server system with setup times, including the extensively-studied Exponential setup system. Furthermore, we use our bounds to derive a highly-accurate approximation, which we evaluate in a variety of settings. These results are made possible via our new method for bounding the expectation of a random time integral, called the Method of Intervening Stopping Times or MIST.

Acknowledgments

I have far too many people to thank and far too limited a memory to even attempt to explicitly thank everyone who deserves it on this tiny page. That said, I want to thank my partner May, my family, my friends, and all of the collaborators, colleagues, students, and administrative staff that I've the pleasure of interacting with over the past six years.

It's my impression that completing a PhD is never an easy thing; even so, it's clear to me that, without all your support, guidance, understanding, and occasional prodding, I would never have been up to the task. I hope, sincerely, that all of you have enjoyed your time with me to even a fraction of the degree that I have enjoyed my time with you. I will treasure this time always.

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 1 |
| 1.1 | Brief problem description | 1 |
| 1.2 | What makes the setup effect difficult to understand? | 2 |
| 1.2.1 | Setup effect can be indirect | 2 |
| 1.2.2 | Setup effect is even harder to understand in the multiserver setting | 3 |
| 1.3 | State of the Art | 3 |
| 1.4 | Our Contributions | 5 |
| 1.4.1 | Discussion of Bounds. | 5 |
| 1.4.2 | Discussion of Approximation. | 6 |
| 1.4.3 | Impact: How our results transform capacity provisioning | 7 |
| 1.5 | Outline | 9 |
| 2 | Prior Work | 10 |
| 2.1 | Systems without Simultaneous Setup | 10 |
| 2.1.1 | The $M/G/1/Setup$ | 10 |
| 2.1.2 | $M/M/k$ and $M/G/k$ with staggered setup | 10 |
| 2.2 | The $M/M/k/Setup$ -Exponential | 10 |
| 2.2.1 | $M/M/k/Setup$ -Exponential, Approximations | 10 |
| 2.2.2 | $M/M/k/Setup$ -Exponential, Exact Analysis | 11 |
| 2.2.3 | Distributed Setting | 11 |
| 2.3 | Scheduling with Setup | 12 |
| 2.4 | Prior Work on Deterministic Setup Times | 12 |
| 3 | Model | 14 |
| 3.1 | Detailed Model Description | 14 |
| 3.2 | Construction | 16 |
| 4 | Key Ideas and Techniques | 17 |
| 4.1 | Initial Steps: Applying the Renewal-Reward Theorem | 17 |
| 4.1.1 | Reduction to analyzing $\mathbb{E}[N(\infty) - R]$ | 17 |
| 4.1.2 | Applying Renewal-Reward. | 17 |
| 4.1.3 | Key Idea: Define renewals around stability. | 18 |
| 4.2 | The Method of Intervening Stopping Times (MIST) | 19 |
| 4.2.1 | Why we need it | 19 |

| | | |
|----------|--|-----------|
| 4.2.2 | <i>What it does</i> | 20 |
| 4.2.3 | <i>IST Lemma: Statement and Proof</i> | 20 |
| 5 | <i>The Lower Bounds</i> | 23 |
| 5.1 | <i>Why we need a lower bound</i> | 23 |
| 5.2 | <i>The First Lower Bound</i> | 23 |
| 5.3 | <i>The New Lower Bound</i> | 24 |
| 5.3.1 | <i>The New Lower Bound: Theorem Statement</i> | 24 |
| 5.3.2 | <i>The New Lower Bound: Proof Outline.</i> | 24 |
| 5.3.3 | <i>Proof of Lemma 5.1: Lower Bound on Cycle Integral.</i> | 25 |
| 5.3.4 | <i>Lower Bound on Integral until $\tau_L + \beta$.</i> | 26 |
| 5.3.5 | <i>Lower Bound on Integral after $\tau_L + \beta$.</i> | 28 |
| 5.3.6 | <i>Proof of Lemma 5.2: Upper Bound on the Accumulation Time $\mathbb{E}[T_A]$.</i> | 30 |
| 5.3.7 | <i>Proof of Lemma 5.3: Upper Bound on the Remaining Cycle Time $\mathbb{E}[X - T_A]$.</i> | 33 |
| 5.4 | <i>The Lower Bounds: Review of Findings</i> | 37 |
| 6 | <i>The Upper Bound</i> | 38 |
| 6.1 | <i>Why we need an upper bound.</i> | 38 |
| 6.2 | <i>The Upper Bound</i> | 38 |
| 6.2.1 | <i>Proof of Lemma 6.1, Upper Bound on Integral Over Accumulation Period</i> | 40 |
| 6.2.2 | <i>Proof of Lemma 6.2, Upper Bound on Integral Over Draining Period.</i> | 46 |
| 6.2.3 | <i>Proof of Lemma 6.3: Lower Bound on the Cycle Length.</i> | 51 |
| 6.3 | <i>The Upper Bound: Review of Findings</i> | 52 |
| 7 | <i>The Approximation</i> | 53 |
| 7.1 | <i>Why we need an approximation</i> | 53 |
| 7.2 | <i>The approximation</i> | 53 |
| 7.3 | <i>Justification</i> | 54 |
| 7.3.1 | <i>Justification of Numerator</i> | 54 |
| 7.3.2 | <i>Justification of Denominator</i> | 55 |
| 8 | <i>Conclusion</i> | 56 |
| 8.1 | <i>Summary and Takeaways</i> | 56 |
| 8.2 | <i>Broader Impacts</i> | 57 |
| 8.2.1 | <i>Computer Science</i> | 57 |
| 8.2.2 | <i>Operations/Management</i> | 57 |
| 8.2.3 | <i>Healthcare</i> | 57 |
| 8.3 | <i>Open Problems</i> | 58 |
| 8.3.1 | <i>Standby States</i> | 58 |
| 8.3.2 | <i>Analyzing Tail Performance in the M/M/k/Setup.</i> | 58 |
| A | <i>Miscellaneous Claims</i> | 59 |
| A.1 | <i>Proof of Multiplicative Tightness.</i> | 59 |
| A.1.1 | <i>Proof for Lower Bound.</i> | 59 |

| | | |
|-------|---|----|
| A.1.2 | <i>Proof for Upper Bound.</i> | 60 |
| A.2 | <i>Construction and Coupling Claims.</i> | 61 |
| A.2.1 | <i>Construction</i> | 61 |
| A.2.2 | <i>Three Coupling Claims</i> | 62 |
| A.2.3 | <i>Proof of Claim A.2, the Coupling Integral Bound.</i> | 63 |
| A.2.4 | <i>Proof of Claim A.3, the Coupling Probability Bound.</i> | 65 |
| A.2.5 | <i>Proof of Claim A.4: Bound on Expected Value After Coupling.</i> | 66 |
| A.3 | <i>Hitting Time Bounds.</i> | 68 |
| A.3.1 | <i>Proof of Claim A.5, Discrete-Time Hitting Time Tail Bound.</i> | 68 |
| A.3.2 | <i>Proof of Claim A.6, Continuous-Time Hitting Time Tail Bound.</i> | 69 |
| A.3.3 | <i>Proof of Claim A.7, Bound on Expected Length of Stopped Random Walk.</i> | 71 |
| A.3.4 | <i>Proof of Claim A.8, Bound on the Expected Hitting Time in the M/M/∞.</i> | 72 |
| A.4 | <i>Helper Claims.</i> | 73 |
| A.4.1 | <i>Proof of Claim A.9, the Busy Period Integral Bound.</i> | 73 |
| A.4.2 | <i>Proof of Claim 6.10, the Wait Busy Claim.</i> | 74 |
| A.4.3 | <i>Proof of Claim 6.3</i> | 77 |
| A.4.4 | <i>Proof of Claim A.10.</i> | 80 |
| A.4.5 | <i>Proof.</i> | 80 |
| A.4.6 | <i>Proof of (5.3): Lower Bound on $\mathbb{E}[L]$, Expected Value of First Long Epoch Index.</i> | 80 |
| A.4.7 | <i>Proof of Claim 6.11.</i> | 83 |

| | |
|----------------------------|-----------|
| <i>Bibliography</i> | 87 |
|----------------------------|-----------|

List of Figures

- 1.1 An example of an M/M/k/Setup-Deterministic queue with $k = 4$. Jobs (blue rectangles) enter into a central queue, where they wait in FCFS order until they are served by one of k servers (white circles with black outline containing 1 of 3 elements). Servers can be *off* (red “X”), *on* (blue rectangle), or *in setup* (green hourglass). 2
- 1.2 A comic illustrating of the indirect nature of setup (Section 1.2.1). In the first panel, the queue starts empty. Then, a job arrives triggering the setup of the server. While the initial job is waiting for the server to turn, more jobs arrive, all observing the server in setup —the main cause of their delay. However, jobs which arrive after the server is ready *still* experience additional delay —but do not directly observe *why* they are delayed. 3
- 1.3 Simulation results for the M/M/k/Setup-Deterministic, M/M/k/Setup-Exponential, M/M/k (no setup), varying the number of servers k and keeping fixed the service rate $\mu = 1$, the setup time $\beta = 1000$, and the load $\rho = 0.5$. Note the high separation between all three models. 4
- 1.4 Our results along with simulation data for the M/M/k/Setup-Deterministic and M/M/k/Setup-Exponential, varying the number of servers k while keeping the mean service time $\frac{1}{\mu} = 1$ ms, the mean setup time $\beta = 1000$ ms, and the load $\rho = 0.5$ fixed. (a) A comparison of our results to the true average waiting time in the M/M/k/Setup-Deterministic. Our results behave like the true average waiting time, while the Exponential model behaves differently. (b) A plot showing the increase in average waiting time as one moves from moderate variance (M/M/k/Setup-Exponential) to zero variance (M/M/k/Setup-Deterministic); we use an Erlang distribution to interpolate between the two settings. 6
- 1.5 A provisioning example. 7
- 1.6 A provisioning example highlighting the differences between the Deterministic and Exponential models. To achieve a target waiting time of 20 ms, our approximation correctly predicts it will take $k \approx 2000$ servers, while the Exponential model predicts that only $k \approx 50$ servers should suffice. See Figure 1.7 for a more comprehensive evaluation. 7

| | | |
|-----|--|----|
| 1.7 | Some examples demonstrating the excellent accuracy of our simple approximation (1) to the average waiting time in the M/M/k/Setup-Deterministic. For each of these 9 plots, we plot the behavior of the average waiting time as one varies the load ρ from 0 to 1, holding fixed the total number of servers k as well as the setup time β . In each row, we hold the number of servers k constant while testing increasing values of the setup times β . In each column, we hold the setup time β constant while increasing the number of servers. We also include, as a reference, a dotted line illustrating the point at which the offered load $R \triangleq k\rho = 2$ | 8 |
| 3.1 | An example of M/M/k/Setup-Deterministic with $k = 4$. The state pictured has $Z(t) = 2$ busy servers, which means there are 2 jobs in service. There is $Q(t) = 1$ job in queue, and thus there are $N(t) = Z(t) + Q(t) = 3$ jobs in system. . . . | 15 |
| 4.1 | A depiction of the decomposition of a renewal cycle into an accumulating phase and draining phase, described in Section 4.1. (a) In the M/M/1/Setup, the canonical renewal cycle is split into three parts: 1) the system is empty until a job arrives; 2) jobs accumulate in the queue while the server sets up; and 3) the system's single server turns on, starting a busy period. (b) Likewise, for the M/M/k/Setup, our renewal cycle splits into two parts: 1) during the accumulating phase, the departure rate $\mu Z(t) \leq \mu R = k\lambda$, so that the system is transiently unstable and a queue <i>accumulates</i> ; and 2) during the draining phase, the departure rate $\mu Z(t) > k\lambda$, so that the queue <i>drains</i> | 18 |
| 4.2 | A depiction of our decomposition of a renewal cycle into an accumulating phase and draining phase, described in Section 4.1. During the accumulating phase, the departure rate $\mu Z(t) \leq \mu R = k\lambda$, so that the system is transiently unstable and a queue <i>accumulates</i> . During the draining phase, the departure rate $\mu Z(t) > k\lambda$, so that the queue <i>drains</i> | 19 |
| 6.1 | A depiction of the up-crossings and down-crossings defined in Section 6.2.1. In this example, we see that the number of up-crossings in epoch 3 is $n_e^{(3)} = 2$ and that, in this case, epoch 3 ends when epoch 4 begins (i.e. at time τ_4). | 44 |

Chapter 1

Introduction

What are setup times?

In many systems, servers do not turn on instantly; instead, a significant *setup time* must pass before a server can begin work [1, 25]. For example, when photocopying a document, the first person to use the photocopier must wait for the machine to warm up; when using cloud computing resources, one must wait for a virtual machine (VM) to boot up before the VM can be used [17, 34]; when hosting applications online, application replicas must take time to boot up before they can begin fulfilling requests [14]. When studying the effect of setup times on a system's queueing behavior, we model the effect of setting up via an abstract *setup time* [1].

Why do setup times matter?

Setup times can have a significant impact on a system's queueing behavior, especially in modern systems. This is because, in modern computing systems, 1) servers are regularly turned on and off and 2) setup times are much, much larger than service times. Because servers don't turn on instantly, jobs in a system with setup times end up delayed compared to their no-setup counterparts. Since setup times are sometimes more than 1000 times larger than the average job size (e.g. a VM's average boot time of 10 s compared to an average job size of 10 ms [17]), this additional delay can be significant.

Nevertheless, many systems still regularly turn their servers off and on. Why? Because by doing so, one can save a considerable amount on operating costs, e.g. energy, money, CO_2 emissions, etc. That said, this cost-saving measure is only a viable option if the additional delay caused by setup times is not too large. Therefore, if we want to design systems which are simultaneously efficient *and* performant, we need a good understanding of how setup times affect queueing performance.

1.1 Brief problem description

In this work, we study the effect of setup times on the average waiting time in the $M/M/k/Setup$, a simple variation on the classic $M/M/k$ queue.

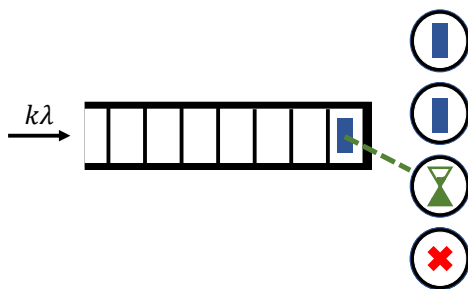


Figure 1.1: An example of an $M/M/k/\text{Setup}$ -Deterministic queue with $k = 4$. Jobs (blue rectangles) enter into a central queue, where they wait in FCFS order until they are served by one of k servers (white circles with black outline containing 1 of 3 elements). Servers can be *off* (red “X”), *on* (blue rectangle), or *in setup* (green hourglass).

Job Dynamics. Outside of its setup dynamics, the $M/M/k/\text{Setup}$ behaves essentially identically to the usual $M/M/k$. Jobs arrive in a Poisson process to a central queue, where they wait in First-Come-First-Served order until they are served by one of k servers. The job at the head of the queue enters service whenever either 1) a server finishes setting up and turns on or 2) a server finishes its current job. Once in service, the job stays in service for an i.i.d. Exponential amount of time, after which the job departs.

Setup Dynamics. To complete our description of the $M/M/k/\text{Setup}$, it suffices to describe how the system controls the setup process. Servers can be in one of three states: *off*, *on*, or *in setup*. Servers turn *off* whenever they finish their current job and there are no jobs waiting in the queue. Servers turn *on* when they have remained *in setup* for a full *setup time*; in general, this setup time is some i.i.d. *random variable* which is sampled at the moment that setup is first initiated. Servers initiate and cancel *setup* based on job arrivals and departures, respectively. In particular, if a job arrives to the system and sees *off* servers, it initiates setup at one of these *off* servers (we assume every server is identical). Likewise, if the number of *in setup* servers ever exceeds the number of jobs waiting in the queue, then the system turns off the server that has been *in setup* for the least amount of time. Using the $M/M/k/\text{Setup}$, we can now study the complex interactions which arise from simultaneous setup.

1.2 What makes the setup effect difficult to understand?

1.2.1 Setup effect can be indirect

Example: the $M/M/1/\text{Setup}$. Note, though, that the additional delay caused by setup does not always manifest in an obvious way; we illustrate this in Figure 1.2. For example, consider a simple single server queue with setup times, the $M/M/1/\text{Setup}$. The job which is the first to arrive to an empty system triggers the *off* server’s setup, and must wait a full setup time before it is served. Likewise, every job that arrives while the server is still setting up must wait in line behind the setup-triggering job, and so also partially observes the server setting up. However,



Figure 1.2: A comic illustrating of the indirect nature of setup (Section 1.2.1). In the first panel, the queue starts empty. Then, a job arrives triggering the setup of the server. While the initial job is waiting for the server to turn, more jobs arrive, all observing the server in setup—the main cause of their delay. However, jobs which arrive after the server is ready *still* experience additional delay—but do not directly observe *why* they are delayed.

even after it turns on, the fact that the server was off for a while has a lasting impression on the length of the queue; the queue is longer than it would otherwise be. Thus, setup can even affect the delay of jobs that never actually observe a server in setup.

The effect of setup on queueing behavior is made even more complex when the setup time is allowed to be a random variable, sampled independently every time setup is initiated, and when the length of each job follows an arbitrary distribution. This more complex model is called the $M/G/1/Setup$. Further, despite its apparent complexity, the full waiting time distribution of the $M/G/1/Setup$ was completely characterized in 1964 by [36].

1.2.2 Setup effect is *even harder* to understand in the multiserver setting

Unfortunately, the effect of setup on delay is even harder to understand when multiple servers can set up at the same time. Recall that, in the single server setting, the server’s setup process always completes once initiated, and there is always at least one job to work on once the server turns on. This implies that, although setup has a complex effect on job delay, the server’s behavior itself is quite simple: it first initiates setup; then, once setup completes, the server begins working; then, after finishing all the work in the system, the server turns off. On the other hand, when multiple servers can *simultaneously* be *in setup*, their server states begin to interact.

In particular, via the speed of their processing, the *busy* servers indirectly control the setup behavior of the *not-busy* servers. For example, if server A is *on* while server B is setting up, then server A might finish all the work in the queue before server B even has a chance to turn *on*. As such, in the multiserver setting, it can sometimes make sense to *cancel* a server’s setup process; a situation which would *never* occur in the single server setting. Of course, the opposite can also happen: if the busy servers are working much more slowly than expected, then the queue might grow large enough that we begin setting up a server that would otherwise be left *off*. This interaction between departure behavior and setup behavior is exactly what makes the setup effect so much harder to understand in the multiserver setting.

1.3 State of the Art

Everyone uses the Exponential model. Despite continued academic interest, our understanding of the $M/M/k/Setup$ is still extremely limited. Perhaps the most significant limitation is that

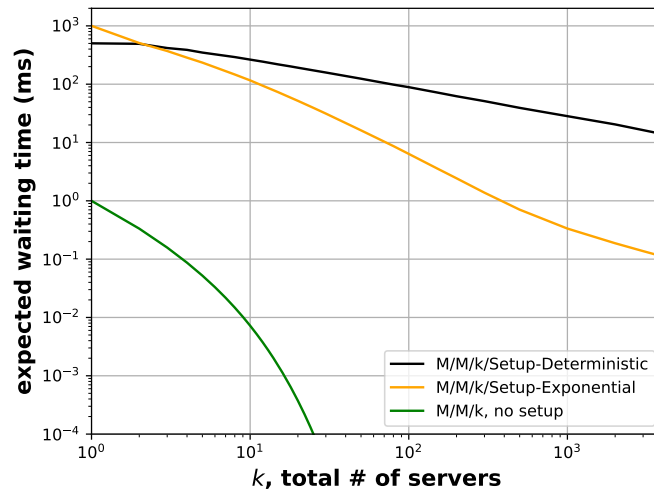


Figure 1.3: Simulation results for the M/M/k/Setup-Deterministic, M/M/k/Setup-Exponential, M/M/k (no setup), varying the number of servers k and keeping fixed the service rate $\mu = 1$, the setup time $\beta = 1000$, and the load $\rho = 0.5$. Note the high separation between all three models.

all state-of-the-art research [15, 32, 33] assumes that setup times are Exponentially distributed. This limitation has major consequences for the utility of their work.

The Exponential model is unrealistic. We give two reasons why this limitation is so significant. First, the “Exponential setup” assumption leads to extremely unrealistic behavior in some situations. To illustrate where the breakdown in realism happens, consider a scenario where only a single server is setting up and compare it to a scenario where 100 servers begin setup at the same time. In the Exponential setup model, the 100-server system receives its first server on average 100x faster than the single-server system receives its first (and only) server. This is not a quirk of our specific example: in the Exponential model, the *longer* the system’s queue is, the *more rapidly* the system’s servers turn on to help *drain* that queue. In a sense, the Exponential system can rapidly “react” to increases in queue length.

The Exponential model underestimates waiting. This unrealistic “reactivity” phenomenon causes a further, more concerning, problem: in real systems, the Exponential model dramatically underestimates how much waiting actually occurs. To be more precise, in modern systems: 1) average setup times are often larger than average job sizes by two or three orders of magnitude [14, 17, 26, 27, 28]; and 2) as noted in the paper that *introduces* the Exponential model [15], setup times are actually closer to *Deterministic*. When these two criteria are satisfied, as observed in Figure 1.3, the true waiting time is often orders of magnitude larger than what the Exponential model predicts. Accordingly, in many practical studies of setup [14, 15, 20, 22], setup times are assumed to be Deterministic, e.g. servers take exactly 2 minutes to set up. However, despite its apparent limitations, the Exponential setup model remains the *de facto* choice for theoretical analysis, since it allows for the application of a number of existing theoretical techniques.

Challenges of the Deterministic model. Although modeling setup times as Deterministic might be more realistic, it also comes with a set of unique theoretical challenges. In the Deterministic case, one must, even in simulation, track the individual remaining setup time of *every* server that is currently setting up. In contrast, because the Exponential distribution is *memoryless*, in the Exponential case it suffices to track only the *total number* of servers setting up instead, greatly simplifying the system state. Moreover, the Exponential setup model’s simplified state forms a Continuous-Time Markov Chain, a well-studied class of stochastic processes for which a number of techniques have been developed. For the Deterministic setup model, no such techniques exist.

1.4 Our Contributions

In this thesis, I demonstrate that **setup times, along with their distributions, can not be ignored; setup times can cause profound increases in waiting time, especially when the distribution of setup time has low variability.** (See Figures 1.4a and 1.4b for an illustration.)

Summary. The contributions of this thesis are both theoretical and practical. On the theoretical side, in Chapters 5 and 6, respectively, we derive the first lower and upper bounds on the average waiting time in the $M/M/k/\text{Setup-Deterministic}$. Notably, these results are the first closed-form bounds on the average waiting time in *any* $M/M/k/\text{Setup}$ system, including the extensively-studied Exponential setup system. We obtain these bounds via a new technique for bounding random time integrals called MIST, described in Chapter 4. On the practical side, in Chapter 7, we then show how to take the components of our upper and lower bounds, and combine them to make a highly accurate approximation; an example of the approximation alongside the bounds is shown in Figure 1.7.

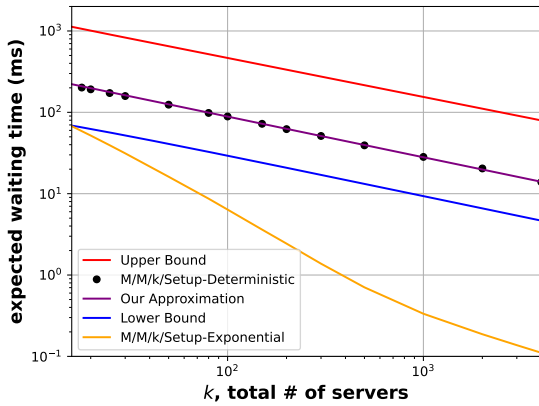
Assumptions. In our theorems, in order to more simply characterize the system’s behavior, we make two assumptions. First, we assume that setup times are large compared to service time, i.e. the average setup time $\beta \geq 100\frac{1}{\mu}$; this is often satisfied in practice [14, 17, 26, 27, 28]. Second, we assume that, on average, the system utilizes at least 100 servers, i.e. the offered load $R \triangleq k\rho \geq 100$. Note also that, while these conditions are often satisfied in practice, the specific 100 value stated here is not strictly necessary for our analysis to be accurate, as evidenced by the apparent accuracy of Approximation 1, shown in Figure 1.7.

1.4.1 Discussion of Bounds.

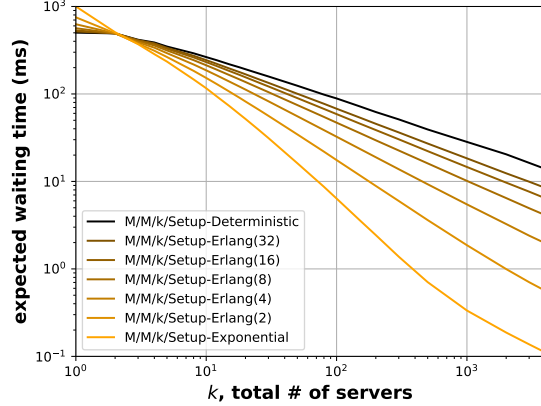
After simplification, our bounds state

$$\mathbb{E}[Q(\infty)] =_c \mu\beta\sqrt{R} + \frac{1}{1-\rho};$$

where the operator $=_c$ denotes equality up to multiplicative constants; we show this explicitly in Section A.1. This simplified characterization of the queue length gives us insight into how the



(a) An illustration of our results.



(b) Variance-dependent waiting behavior.

Figure 1.4: Our results along with simulation data for the $M/M/k/Setup$ -Deterministic and $M/M/k/Setup$ -Exponential, varying the number of servers k while keeping the mean service time $\frac{1}{\mu} = 1$ ms, the mean setup time $\beta = 1000$ ms, and the load $\rho = 0.5$ fixed. (a) A comparison of our results to the true average waiting time in the $M/M/k/Setup$ -Deterministic. Our results behave like the true average waiting time, while the Exponential model behaves differently. (b) A plot showing the increase in average waiting time as one moves from moderate variance ($M/M/k/Setup$ -Exponential) to zero variance ($M/M/k/Setup$ -Deterministic); we use an Erlang distribution to interpolate between the two settings.

system parameters govern the system’s queueing behavior. The first term, $\mu\beta\sqrt{R}$, scales linearly with the setup time β and captures the effect of turning servers *off* and *on*. The second term, $\frac{1}{1-\rho}$, dominates the first term only when the load ρ is high enough, say in the renowned super-Halfin-Whitt regime [16, 24] where $\rho = 1 - \gamma k^{-\alpha}$ with $0 < \gamma < 1$ and $\alpha > 0.5$. In this case, it recovers the $\frac{1}{1-\rho}$ scaling seen in the $M/M/k$ without setup times.

1.4.2 Discussion of Approximation.

Moreover, our analyses suggest a simple approximation for the waiting time. Taking $C_{\text{apx}} \triangleq \sqrt{\frac{\pi}{2}}$,

$$\mathbb{E}[T_Q] \approx \frac{C_{\text{apx}}}{2} \frac{\beta}{\sqrt{R}} + \frac{C_{\text{apx}}\sqrt{R}}{k(1-\rho) + C_{\text{apx}}\sqrt{R}} \left(\frac{1}{k\lambda} \left[\frac{1}{1-\rho} + \frac{1}{2} \right] \right).$$

In simulations, we find this approximation to be highly accurate across a variety of parameter settings (Figure 1.7), so long as the offered load $R > 2$. This is despite the fact that, in our analysis, we assume that the offered load $R \gg \sqrt{R}$ and often make considerable use of this fact. In that sense, it is a testament to the strength of our approach that our resulting approximation remains accurate all the way up to an offered load of $R = 2$. Furthermore, while offered loads smaller than this are of limited interest in practical settings, when the offered load is that small, we seem to recover the single-server behavior observed in [36].

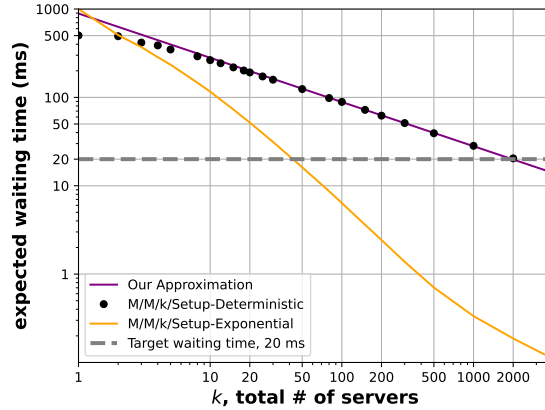


Figure 1.5: A provisioning example.

Figure 1.6: A provisioning example highlighting the differences between the Deterministic and Exponential models. To achieve a target waiting time of 20 ms, our approximation correctly predicts it will take $k \approx 2000$ servers, while the Exponential model predicts that only $k \approx 50$ servers should suffice. See Figure 1.7 for a more comprehensive evaluation.

1.4.3 Impact: How our results transform capacity provisioning

A common but complex problem which arises in many areas is that of designing a system such that the average waiting time of a customer is below some target waiting time. Historically, we understand this problem well for systems without setup times, e.g. there’s a straightforward formula for the average waiting time in the M/M/k without setup. Unfortunately, our understanding of this problem is quite poor for more modern systems, since their average waiting times are affected by setup times. As mentioned before, previous results on understanding the relationship between setup times and the average waiting time leave much to be desired. Our new results expand on the state-of-the-art Exponential model in two important ways: 1) predicting the waiting time is much less computationally-intensive, and 2) the prediction is of much higher quality.

Easier predictions. Compared to the Exponential model, our new Deterministic approximation greatly simplifies the design process. In particular, when predicting the average waiting time in the Exponential model using the state-of-the-art method from [15], one must solve a system of $O(k^2)$ quadratic equations to find the average waiting time $\mathbb{E}[T_Q]$. Two practical issues arise from this fact. First, the equations change depending on the number of servers k , meaning that the computation must be repeated every time one wishes to test a new number of servers. Second, the opacity of the process makes it difficult to get intuition about how the average waiting time changes as one alters the system parameters. In contrast, Approximation 1 is a relatively simple function of the relevant parameters. The simplicity of our approximation has, likewise, two benefits: 1) computing the waiting time becomes easy, and 2) our approximation’s form makes it clear how and why the waiting time behaves the way it does.

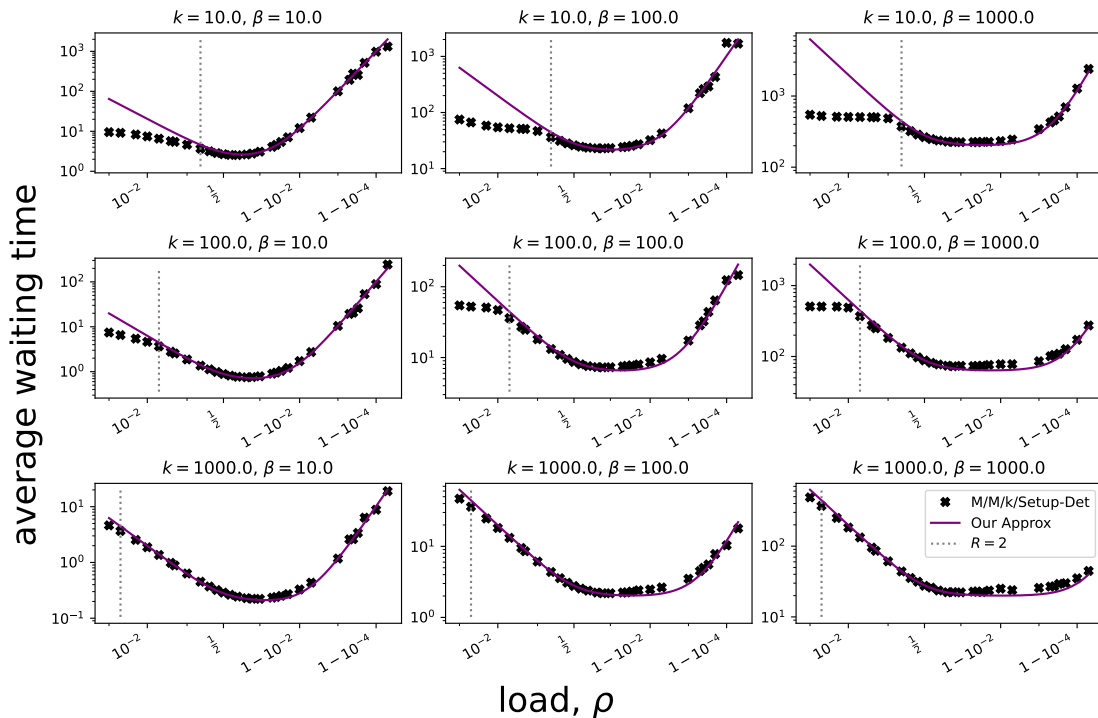


Figure 1.7: Some examples demonstrating the excellent accuracy of our simple approximation (1) to the average waiting time in the M/M/k/Setup-Deterministic. For each of these 9 plots, we plot the behavior of the average waiting time as one varies the load ρ from 0 to 1, holding fixed the total number of servers k as well as the setup time β . In each row, we hold the number of servers k constant while testing increasing values of the setup times β . In each column, we hold the setup time β constant while increasing the number of servers. We also include, as a reference, a dotted line illustrating the point at which the offered load $R \triangleq k\rho = 2$.

Higher quality predictions. Moreover, when compared to the predictions of the Exponential model, the predictions we obtain using our Deterministic approximation are of a much higher quality. This difference in quality is perhaps best illustrated by looking at a simple example. In Figure 1.5, we compare the prediction from the Exponential model to the prediction from our approximation, plotting how the predicted average waiting time changes as one increases the number of servers k while fixing the load $\rho = 0.5$, the average setup time $\beta = 1000$ ms, and the average service time $\frac{1}{\mu} = 1$ ms (note that these are typical relative values in many applications [27]). Our goal is to determine how large the number of servers k needs to be before we reach the target waiting time $T^{\text{target}} = 20$ ms. In both models, the average waiting time decreases as the system gets larger. However, the Exponential model predicts that the average waiting time will be small enough once $k = 50$. On the other hand, as captured by our approximation, the Deterministic setup system will only reach the target waiting time once the number of servers $k \approx 2000$ —a full 40 times larger than what the Exponential system predicts! At even a modest number of servers, the Exponential system underestimates the waiting time by orders of magnitude.

1.5 Outline

Chapter 2: Prior Work. In Chapter 2, we begin our study of setup by reviewing some related work. We start by discussing the single server setting, then we move through the history of the study of setup times up to the state of the art. For each result we review, we compare and contrast their work with the main results developed in this thesis.

Chapter 3: Model. In Chapter 3, we give a more detailed description of our model. Besides reviewing the brief description we gave in this chapter, we also describe our notation and give a construction of our processes of interest using Poisson processes.

Chapter 4 : Key Ideas and Techniques. Next, in Chapter 4, we describe the key ideas and techniques of this thesis. In particular, we introduce the Method of Intervening Stopping Times; the *MIST* method. We describe the MIST method by first describing its general function, then stating its associated formal definition, then proving a key lemma which allows it to be generally applied.

Chapter 5: The Lower Bounds. In Chapter 5, we describe our first two main results (Theorems 5.1 and 5.2), both lower bounds on the average queue length in the $M/M/k/\text{Setup-Deterministic}$. We begin by describing in greater detail why a lower bound is needed, then proceed by stating both bounds and proving the stronger one.

Chapter 6: The Upper Bound. In Chapter 6, we describe our final main result (Theorem 6.1), an upper bound on the average queue length in the $M/M/k/\text{Setup-Deterministic}$. As we did in Chapter 5 with the lower bounds, we first describe why we need this upper bound. Afterwards, we give its proof.

Chapter 7: The Approximation. After proving these results, in Chapter 7, we develop an approximation to the average waiting time in the $M/M/k/\text{Setup-Deterministic}$. As before, we first describe why such an approximation is needed. Afterwards, we explicitly state the approximation formula and describe how to derive the approximation from the bounds in Chapters 5 and 6.

Chapter 8: Conclusion. Finally, in Chapter 8, we summarize the main results of this thesis, discuss some possible applications, and describe a few open problems.

Chapter 2

Prior Work

2.1 Systems without Simultaneous Setup

2.1.1 The M/G/1/Setup

The best-understood case is the single-server case. The foremost result on this model is the result of [36]; the author considers a generalization of the M/G/1 queue where, if a customer arrives while the server is idle, then they have a different service distribution than if they arrive while the server is busy. By observing that the system state at customer departure times forms a discrete Markov chain, then analyzing that embedded chain, Welch characterizes the steady-state and transient distributions of the queue length; via distributional Little's Law, this gives the same result for delay and response time. This important result has been extended in a variety of different directions, by adjusting the service discipline or arrival process[3, 4, 18].

2.1.2 M/M/k and M/G/k with staggered setup

The easiest case of multiserver systems with setup times involves the *staggered setup* model, where at most one server can be in setup at a time, greatly simplifying the analysis. In [2], the authors obtain an expression for the steady-state distribution of queue length for the system when setup times are Exponential, using the method of difference equations. In [11], the authors simplify the solution of the staggered M/M/k with exponential setup times considerably, and prove a decomposition result for mean delay. In [9], the decomposition result is generalized to a hyperexponential job size distribution, and shown to hold approximately for a general job size distribution.

2.2 The M/M/k/Setup-Exponential

2.2.1 M/M/k/Setup-Exponential, Approximations

All previous theoretical results that investigate an M/M/k/Setup system assume Exponential setup times. We first highlight the state-of-the-art papers concerning approximating the M/M/k/Setup-

Exponential. In particular, we highlight the work in [32] and [11]. Gandhi et al. [11] seek useful intuitive approximations to the M/M/k/Setup-Exponential system. Their approximations stem from an exact analysis of the M/M/ ∞ /Setup-Exponential system, which they then modify in various ways to capture the finite server case. The approximations in [11] work well, except when both load and setup times are moderately high ($\rho > 0.5$ and $\frac{\mu}{\alpha} > 10$).

Pender and Phung-Duc [32] consider a generalization of the M/M/k/Setup-Exponential model which includes non-stationary arrival rate and customer abandonment. Within this model, they derive a mean field approximation for the system dynamics, which they prove converges as the number of servers, k , approaches infinity.

Unlike our work, neither Pender and Phung-Duc [32] nor Gandhi et al. [11] provide explicit bounds on the delay. The approximations themselves are also not stated as an explicit function of the system parameters. Finally, neither considers Deterministic setup times.

2.2.2 M/M/k/Setup-Exponential, Exact Analysis

There are only a few results that deal with the exact analysis of the M/M/k with Exponential setup times. The most well-known are [15] and [33].

In a followup to the approximation work done in [11], in [15] the authors give the first exact analysis of the M/M/k/Setup-Exponential system. To do this, they develop the *Recursive Renewal Reward (RRR)* technique, which allows them to analyze 2-dimensional Markov chains of a certain structure. They apply this technique to the M/M/k/Setup-Exponential system, and thus provide a method for computing the time-average value of *any* function of the system state; applying this method to the correct function gives the mean and Laplace transform for the number of jobs in queue.

In [33], the author rederives the exact solutions for the queue length obtained in [15] using two different methods: an analysis using generating functions, and an analysis applying the matrix analytic method after casting the system as a quasi-birth-death process. Although these techniques appear different, the author highlights some core correspondences between them, and also between these methods and the RRR technique of [15].

Despite the fact that [15] and [33] represented the first breakthrough in our understanding of the setup effect in 50 years, their results are limited in two significant ways. First, instead of a closed-form formula for the average waiting time, the authors only derived an algorithm for computing the average waiting time. This algorithm is useful in the sense that it bypasses the need to simulate the system, but unfortunately fails to give intuition about wait times scale with system parameters. Moreover, like all of the works mentioned in this section, their work assumes that setup times are distributed Exponentially, which turns out to severely limit the utility of their results; see Section 1.3 for a detailed discussion.

2.2.3 Distributed Setting

There has also been some work on servers with setup times outside of the centralized queue setting. In [30], the authors consider a queueing system which functions much like the M/M/k/Setup-Exponential, except, instead of a central queue, each server has its own queue, and there is a central dispatcher which routes arriving jobs to one of these queues. In this model, they describe

a token-based load balancing and scaling scheme called TABS, and prove that its performance (as $k \rightarrow \infty$) is asymptotically optimal. In particular, they show that the relative energy wastage and the mean delay both go to 0 under their scheme, by analyzing an appropriate fluid limit. In a followup paper, [29], the authors consider the performance of TABS in the infinite-buffer case. They give two results. First, they show that, somewhat counterintuitively, there exist parameter settings under which the TABS scheme is unstable. Second, they show that, in spite of this finite instability, for sufficiently large k , the system under TABS is stable. Moreover, its performance continues to be asymptotically optimal. In our opinion, it is best to think of these results as complementary to the body of work on the M/M/k/Setup, as one typically does when comparing distributed queueing work to centralized queueing work. Although all papers discussed deal with setup times in some capacity, the nature of the questions being asked and answered in [30] and [29] are very different from the central-queue-oriented work we discuss.

2.3 Scheduling with Setup

Gittins in the G/G/k/Setup. In [19], the authors consider a very general queueing model, the G/G/k/Setup, and show that the scheduling performance of the Gittins policy is near-optimal in this setting. In particular, they explicitly bound the deviation from optimality of the average waiting time under the Gittins policy, showing that this “suboptimality loss” is uniformly bounded at all loads. They thus conclude that the Gittins policy is heavy-traffic optimal.

Their work differs from ours in two important ways. First, they investigate a model of setup where the setup process is never cancelled. While this may be accurate in certain situations, a reasonable amount of complexity in our problem stems from the fact that the setup process can be cancelled. Second, they are mainly concerned with bounding the performance of a scheduling policy as compared to the optimal scheduling policy. By contrast, our principle results are concerned with directly characterizing the average waiting time. In that sense, [19] serves as an interesting study whose results are somewhat orthogonal to ours.

2.4 Prior Work on Deterministic Setup Times

M/G/2/Setup-Deterministic, with dispatching In the control literature, deterministic setup times have been incorporated into models in order to enhance realism. Hyytiä et al. [20] consider a dispatching version of the M/G/2/Setup-Deterministic model, and attempt to build near-optimal policies for the joint control of setup initiation and the dispatching of jobs. We hope that our analysis here could open the door to more fine-grained stochastic analysis of such control policies.

M/M/k/Setup-Deterministic, simulation only The only work we have found which discusses the M/M/k/Setup-Deterministic model explicitly is a simulation-based thesis by Kara [22]. They observe that the mean delay in the M/M/k/Setup-Deterministic is consistently larger than that of the M/M/k/Setup-Exponential, and, as the mean setup time $\frac{1}{\alpha}$ increases, the relative increase in mean delay between the M/M/k/Setup-Deterministic and the M/M/k/Setup-Exponential also increases. We corroborate and expand on their results in Section 1.3.

Algorithms for reducing the effect of setup times on delay and energy usage Setup times are both a problem from a delay perspective and also from an energy perspective (servers utilize peak power while in setup [14]). One can of course avoid setup times altogether by always leaving servers on, but this results in wasted energy as well, since a server which is on, but idle, utilizes 60-70% of peak energy [14]. To manage power efficiently, several algorithms have been developed to reduce the costly effects of setup times. One idea is *DelayedOff*, whereby a one waits some time before turning off a server, so as to avoid a future setup time [10, 11, 14, 32]. When using *DelayedOff*, the choice of which idle server to route a job to now matters. One idea is routing jobs to the *Most Recently Busy server (MRB)*, so as to minimize the size of the pool of servers that are turning on and off [10]. Similar to MRB is the idea of creating a rank ordering of all servers and always sending each job to the *lowest-numbered server in the rank* [14]. The goal of all such algorithms is to minimize the Energy-Response-time-Product (ERP) [10], maximize the Normalized-Performance-Per-Watt (NPPW) [8], or minimize energy given a fixed tail cutoff for response time [14]. Other ideas for minimizing delay and energy involve utilizing sleep states in servers, which require more power than being off, but have a lower setup time [12, 13].

Chapter 3

Model

In Chapter 3, we discuss our model of interest, the M/M/k/Setup-Deterministic. We begin the chapter by going through a detailed model description, then discuss how to construct the relevant stochastic processes via Poisson processes.

3.1 Detailed Model Description

The system behavior, excluding setup. As in the typical M/M/k queue, jobs arrive in a Poisson process of rate $k\lambda$ into a FCFS queue where jobs wait to be served at one of k servers. The job at the head of queue enters service whenever a server frees up, either from a job completing service or from a server finishing set up. Once a job enters service, it remains in service for $\text{Exp}(\mu)$ time before departing. We assume all the servers have identical service and setup properties. As such, we can assign each server an index from 1 to k , and without loss of generality assume that departures always occur at the busy server with the highest index; i.e., we re-index the servers when a job departs so the server with the newly departed job has the highest index among the busy servers. From here, we define the quantity $Z(t)$ to be the number of busy servers (or jobs in service) at time t , the quantity $Q(t)$ to be the number of jobs waiting in the queue at time t , and the quantity $N(t) = Q(t) + Z(t)$ to be the total number of jobs in our system. Excluding the setup dynamics, one sees that, as promised, our model behaves identically to the M/M/k queue.

The setup dynamics. From here, it suffices to describe precisely how servers will be turned *on* and *off*. We assume that each server is always in one of three states: *on*, *off*, or *in setup*. A given server remains *on* only as long as that server remains busy. In other words, a server turns *off* when it finishes its current job and the queue is empty. On the other hand, server i begins *setup* when a job arrives to the system and there are only $i - 1$ jobs in the system. Server i remains in setup until one of two events occurs: either 1) some fixed quantity β time has passed, or 2) there are fewer than i jobs in the system; accordingly, we refer to β as the *setup time* of a server. In the first case, if β time has passed without $N(t)$ dipping below i , then server i has completed its setup and begins working on the job at the head of the queue. In the second case, if the number of jobs $N(t)$ dips below i before server i completes setup, then the setup is canceled and server i turns *off*. We use $Y_i(t)$ to denote the detailed state of server i at time t . If server i is *off*, we set

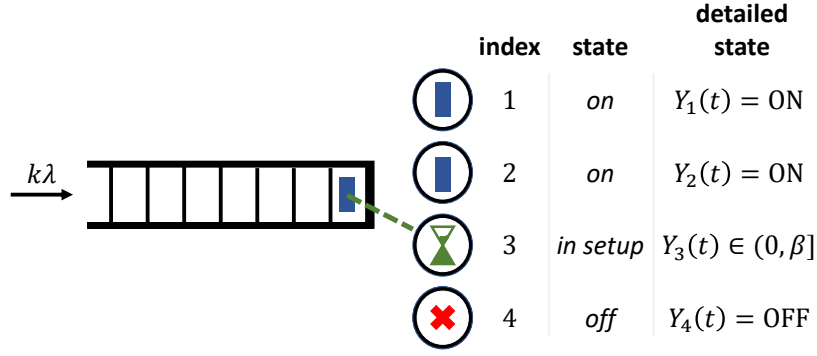


Figure 3.1: An example of $M/M/k/\text{Setup-Deterministic}$ with $k = 4$. The state pictured has $Z(t) = 2$ busy servers, which means there are 2 jobs in service. There is $Q(t) = 1$ job in queue, and thus there are $N(t) = Z(t) + Q(t) = 3$ jobs in system.

$Y_i(t) = \text{OFF}$; if server i is *on*, we set $Y_i(t) = \text{ON}$; if server i is *in setup*, we let $Y_i(t)$ denote the remaining amount of time until server i would finish setup, if left uninterrupted. To be precise, $Y_i(t)$ is set to β when server i first initiates *setup*, and this value decreases at rate 1 until either setup completes or setup is canceled. For convenience, we assume, without loss of generality, that $\text{ON} < s < \text{OFF}$ for every possible remaining setup time $s \in (0, \beta]$; this ensures that the detailed state $Y_i(t)$ is non-decreasing in i . As a shorthand, we use $\mathbf{Y} = (Y_1(t), Y_2(t), \dots, Y_k(t))$ to denote the vector of detailed server states.

A state descriptor. Accordingly, a Markovian state descriptor for our system at time t is $S(t) \triangleq (N(t), \mathbf{Y}(t))$. Note that, since one can recover the number of jobs in service $Z(t)$ from the detailed server states $\mathbf{Y}(t)$, one could also choose the state to be $(Q(t), \mathbf{Y}(t))$. Either suffices in providing a complete description of the forward dynamics of the system. Furthermore, when discussing the steady-state distribution of, say, the number of jobs $N(t)$, we use the notation $N(\infty)$.

Some important constants. We define some system parameters which are critical to system behavior. We use $\rho \triangleq \frac{\lambda}{\mu}$ to refer to the load of our system, i.e., the time-average utilization of an average server. We call the offered load $R \triangleq k\rho$; this is the time-average *number* of busy servers in our system. To enforce stability, we require that $\rho < 1$. As discussed previously, the symbol β refers to the fixed (Deterministic) setup time of a server.

Busy period notation. Our results can be stated more concisely with two quantities related to a busy period of an $M/M/1$ queue. We give the notation below. We use $T^{\text{busy}}(n, j)$ to denote the expectation of the random *length* of an $M/M/1$ busy period with arrival rate $k\lambda$, service rate $k\lambda + \mu j$, and which starts with n jobs in the system. Likewise, we use $I^{\text{busy}}(n, j)$ to denote expectation of the random *time integral of the number of jobs* within the $M/M/1$ over the same

period. Explicitly, we have

$$T^{\text{busy}}(n, j) = \frac{n}{\mu j} \quad (3.1)$$

and

$$I^{\text{busy}}(n, j) = \frac{n}{\mu j} \left[\frac{n+1}{2} + \frac{R}{j} + 1 \right]. \quad (3.2)$$

3.2 Construction

We now discuss how we formally construct this system using Poisson processes; being explicit here will prove useful when we make coupling arguments in the future.

The arrival and departure processes. We take the number of jobs that have arrived at time t to be $\Pi_A(t)$, where Π_A is a Poisson process of rate $k\lambda$. In a slight abuse of notation, we let $\Pi_A([a, b])$ denote the number of arrivals that occur in the interval $[a, b]$; we apply the same extension to all other counting processes mentioned here. We set the potential departure process of, say, server i to be $\Pi_i(t)$, where Π_i is a Poisson process of rate μ . A potential departure from server i only “counts” if server i is busy when that potential departure occurs, i.e., if the number of busy servers $Z(t) \geq i$ at the time. Thus, the total number of departures from our system by time t is

$$D(t) \triangleq \sum_{i=1}^k \int_0^t \mathbf{1}\{Z(s) \geq i\} d\Pi_i(s),$$

where these integrals are with respect to the Π_i 's as counting processes.

The number of busy servers $Z(t)$. To find the number of busy servers $Z(t)$, one could count the number of setup completion events that have occurred so far and the number of server shutoffs that have occurred so far; this description is a bit difficult to work with. Alternatively, one can see from the initial description of setup dynamics that server i is *on* at time t if and only if the total number of jobs $N(s) \geq i$ for all $s \in [t - \beta, t]$, where one should recall that β is the setup time. An easier description of $Z(t)$ follows:

$$Z(t) = \min \left(k, \min_{s \in [t-\beta, t]} N(s) \right).$$

A departure operator. We can extend our departure process $D(t)$ to a departure operator $\mathcal{D}[f(s)](\mathcal{I})$ which takes a function $f(s) \in \{0, 1, \dots, k\}$ defined on some interval \mathcal{I} and computes the number of departures that would occur in that interval provided that the number of busy servers $Z(s) = f(s)$, i.e.

$$\mathcal{D}[f(s)]((a, b)) \triangleq \sum_{i=1}^k \int_a^b \mathbf{1}\{f(s) \geq i\} d\Pi_i(s).$$

Note that we can write the total number of departures using our newly-defined operator as $D(t) = \mathcal{D}[Z(s)]([0, t])$.

Chapter 4

Key Ideas and Techniques

We now describe our approach to analyzing the average waiting time in the M/M/k/Setup-Deterministic, using the upper bound, Theorem 6.1, as a case study. To begin, we go through the first few steps in our proof, leading us to our first technical challenge: defining a renewal cycle so that we can apply the Renewal-Reward Theorem. After explaining how to choose the renewal cycle, we then address our second technical challenge: analyzing the time integrals resulting from our application of the Renewal-Reward Theorem. To solve this challenge, we introduce a method we call the *Method of Intervening Stopping Times* (MIST).

4.1 Initial Steps: Applying the Renewal-Reward Theorem

4.1.1 Reduction to analyzing $\mathbb{E}[N(\infty) - R]$.

We begin by applying the Renewal Reward theorem. Although it is tempting to apply the theorem directly to the queue length $\mathbb{E}[Q(\infty)]$, it simplifies the analysis considerably if one analyzes the number of jobs $\mathbb{E}[N(\infty) - R]$ instead. To justify this, note that, in steady-state, the number of busy servers $\mathbb{E}[Z(\infty)] = R$, and that the total number of jobs in system $N(t)$ satisfies $N(t) = Q(t) + Z(t)$. It follows that the average queue length

$$\mathbb{E}[Q(\infty)] = \mathbb{E}[N(\infty) - Z(\infty)] = \mathbb{E}[N(\infty)] - R = \mathbb{E}[N(\infty) - R].$$

4.1.2 Applying Renewal-Reward.

Applying the Renewal Reward Theorem, for any renewal cycle,

$$\mathbb{E}[N(\infty) - R] = \frac{\mathbb{E}\left[\int_{\text{cycle}} [N(t) - R] dt\right]}{\mathbb{E}[\text{cycle length}]} = \frac{\mathbb{E}\left[\int_0^X [N(t) - R] dt\right]}{\mathbb{E}\left[\int_0^X 1 dt\right]}, \quad (4.1)$$

where we have set time 0 to be an arbitrary renewal point and time X to be the renewal point which immediately follows it. To upper bound (4.1), it suffices to upper-bound the right side's numerator and lower-bound the right side's denominator. These bounds constitute our three main

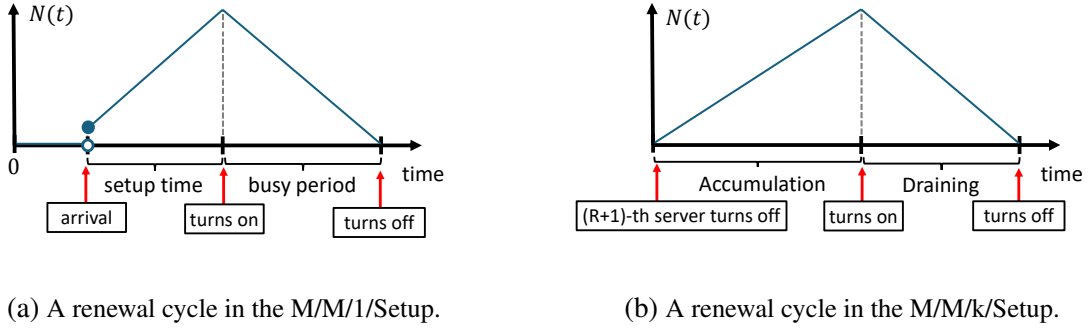


Figure 4.1: A depiction of the decomposition of a renewal cycle into an accumulating phase and draining phase, described in Section 4.1. (a) In the M/M/1/Setup, the canonical renewal cycle is split into three parts: 1) the system is empty until a job arrives; 2) jobs accumulate in the queue while the server sets up; and 3) the system’s single server turns on, starting a busy period. (b) Likewise, for the M/M/k/Setup, our renewal cycle splits into two parts: 1) during the accumulating phase, the departure rate $\mu Z(t) \leq \mu R = k\lambda$, so that the system is transiently unstable and a queue *accumulates*; and 2) during the draining phase, the departure rate $\mu Z(t) > k\lambda$, so that the queue *drains*.

lemmas; two lemmas for the numerator (which is harder to bound) and one for the denominator. But before we can state these main lemmas, we must first address our first technical challenge: defining the collection of renewal points \mathcal{X} .

4.1.3 Key Idea: Define renewals around stability.

Drawing insight from the M/M/1/Setup. To explain our choice, we draw insight from the renewal-reward analysis of the M/M/1/Setup queue. In that setting, the canonical renewal points are those moments when the system empties. By choosing this renewal point, it becomes possible to break up our renewal cycle into distinct, easy-to-analyze phases; see Figure 4.1a for an illustration. In particular: 1) the system is empty until a job arrives; 2) jobs accumulate in the queue while the server sets up; and 3) the system’s single server turns on, at which point the queue drains until empty (i.e. the system enters an M/M/1 busy period). These three phases are much simpler to analyze, since we understand the transient behavior of constant rate Poisson processes and M/M/1 queues very well.

Two possible interpretations. When choosing a renewal point for the M/M/k/Setup, there are two natural interpretations of the M/M/1/Setup renewal point, but only one of these interpretations leads to the same natural decomposition. In particular, one could interpret the M/M/1 renewal point as occurring either 1) when the system empties, i.e. when the number of jobs $N(t)$ goes from 0 to 1, or 2) when the system moves from transient stability to transient instability, i.e. when the departure rate $\mu Z(t)$ goes from above the arrival rate to below it. In the M/M/1/Setup queue, these points are equivalent. However, when one moves to the M/M/k/Setup, only the

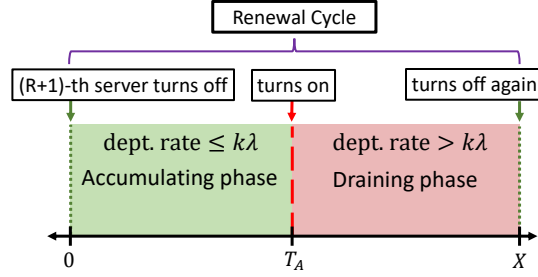


Figure 4.2: A depiction of our decomposition of a renewal cycle into an accumulating phase and draining phase, described in Section 4.1. During the accumulating phase, the departure rate $\mu Z(t) \leq \mu R = k\lambda$, so that the system is transiently unstable and a queue *accumulates*. During the draining phase, the departure rate $\mu Z(t) > k\lambda$, so that the queue *drains*.

second interpretation gives us the same decomposition; in our analysis, we define our collection of renewal points $\mathcal{X} = \{t : Z(t^-) = R + 1, Z(t) = R\}$ as those moments when the $(R + 1)$ -th server is turned off.

A natural decomposition. By defining our collection of renewal points \mathcal{X} as those moments when the $(R + 1)$ -th server is turned off, we naturally split a renewal cycle into two parts: the first part of the cycle *before* the $(R + 1)$ -th server turns on, and the second part of the cycle *after* it turns on; we depict this decomposition in Figure 4.2. During the beginning of a cycle (when $Z(t) \leq R$), the departure rate $\mu Z(t)$ is smaller than the arrival rate, making the system behave, in a transient sense, like a critically- or over-loaded queue. On the other hand, after the $(R + 1)$ -th server turns on and until the renewal cycle is over, the departure rate $\mu Z(t)$ is guaranteed to be strictly greater than the arrival rate, making the queue, on the whole, drain over time. This observation turns out to be hugely useful in our analysis. As such, we have special names for each of these special times: we call the time before the $(R + 1)$ -th server turns on the *accumulation period*, we call the moment when the $(R + 1)$ -th server turns on the *accumulation time* T_A , and we call the period from time T_A until the cycle ends the *draining period*. At times, we will also refer to these periods as *phases*.

4.2 The Method of Intervening Stopping Times (MIST)

4.2.1 Why we need it

Our second technical challenge is to actually bound the integrals we obtained from applying the Renewal-Reward Theorem, i.e. the integrals in (4.1). To do so, we have developed a general technique we call the Method of Intervening Stopping Times, or MIST. The basic function of this lemma is to bound the expected time integral between two random events in some Markov system, an initial event and a final event, a problem arises often in the study of stochastic systems.

4.2.2 What it does

The basic idea of this lemma is to break up our random time interval of interest into a *random number* of smaller, more manageable pieces. We do this by defining *intervening events*, moments where something special happens to the system state that gives us an opportunity to characterize the system's behavior. From there, we can define a “small piece” of time as the time in between these intervening events. For example, in this work, it can often be useful to analyze the system around time points where the number of jobs $N(t)$ gets large. Because we work in a system with setup times, if we have a lot of jobs for a long enough period of time, then, by the end of that long period of time, we can guarantee that a lot of servers are turned on.

By performing this decomposition of the integral into smaller pieces, we reduce our initial bounding problem to showing two facts:

- First, we must show that the time integral of these smaller pieces is not too big; in this thesis, we typically use martingale arguments combined with worst-case coupling arguments to prove this fact.
- Second, we must show that not too many of these these smaller pieces actually occur. For this “not too many” condition, it's particularly helpful if we can show that, if the i -th intervening event has occurred, then the $(i + 1)$ -th event has at most a constant probability of occurring.

By formalizing our notion of events using stopping times and applying some ideas from Wald's equation, we obtain the Intervening Stopping Time Lemma, Lemma 4.1, which we now state and prove.

4.2.3 IST Lemma: Statement and Proof

Lemma 4.1 (Intervening Stopping Time Lemma). *Given a starting stopping time T_0 , an ending stopping time P , and a collection of intervening stopping times $(T_i : i \in \mathbb{Z}_+)$, define the random variable F to be such that $T_F \leq P < T_{F+1}$. Now, given some time-varying random variable $Y_t \geq 0$ which is a function of the underlying Markov state of the system $\mathcal{S}(t)$, suppose that:*

1. $\mathbb{E} \left[\int_{T_0}^{\min(T_1, P)} Y_t dt \middle| \mathcal{F}_{T_0} \right] \leq G_0(\mathcal{S}(T_0))$,
2. $\mathbb{E} \left[\int_{T_i}^{\min(T_{i+1}, P)} Y_t dt \middle| \mathcal{F}_{T_i}, F \geq i \right] \leq G_i + B \cdot \mathbb{E} [\min(T_{i+1}, P) - T_i | \mathcal{F}_{T_i}, F \geq i]$,
3. and $\Pr(F \geq i | \mathcal{F}_{T_i}, F \geq i - 1) \leq 1 - p_i$,

where G_0 is also some function of the system state, and the G_i 's, the p_i 's, and B are all constants (possibly depending on system parameters).

Then,

$$\mathbb{E} \left[\int_{T_0}^P Y_t dt \right] \leq \mathbb{E} [G_0(\mathcal{S}(T_0))] + \Pr(F > 0) \sum_{j=1}^{\infty} G_j \prod_{i=2}^j (1 - p_i) + B \cdot \mathbb{E} [P - T_0].$$

Proof.

We begin with a manipulation of the integral, finding

$$\begin{aligned} \int_{T_0}^P Y_t dt &= \int_{T_0}^{\min(T_1, P)} Y_t dt + \sum_{i=1}^{\infty} \int_{\min(T_i, P)}^{\min(T_{i+1}, P)} Y_t dt \\ &= \int_{T_0}^{\min(T_1, P)} Y_t dt + \sum_{i=1}^{\infty} \mathbf{1}_{T_i < P} \int_{T_i}^{\min(T_{i+1}, P)} Y_t dt. \end{aligned}$$

Applying linearity of expectation and the tower property, we find that

$$\begin{aligned} &\mathbb{E} \left[\int_{T_0}^P Y_t dt \right] \\ &= \mathbb{E} \left[\mathbb{E} \left[\int_{T_0}^{\min(T_1, P)} Y_t dt \middle| \mathcal{F}_{T_0} \right] \right] + \sum_{i=1}^{\infty} \mathbb{E} \left[\mathbb{E} \left[\mathbf{1}_{T_i < P} \int_{T_i}^{\min(T_{i+1}, P)} Y_t dt \middle| \mathcal{F}_{T_i} \right] \right] \\ &= \mathbb{E} \left[\mathbb{E} \left[\int_{T_0}^{\min(T_1, P)} Y_t dt \middle| \mathcal{F}_{T_0} \right] \right] + \sum_{i=1}^{\infty} \mathbb{E} \left[\mathbf{1}_{T_i < P} \mathbb{E} \left[\int_{T_i}^{\min(T_{i+1}, P)} Y_t dt \middle| \mathcal{F}_{T_i} \right] \right]. \end{aligned}$$

Noting that the event $\{T_i < P\} = \{F \geq i\}$, we have

$$\begin{aligned} &= \mathbb{E} \left[\mathbb{E} \left[\int_{T_0}^{\min(T_1, P)} Y_t dt \middle| \mathcal{F}_{T_0} \right] \right] + \sum_{i=1}^{\infty} \mathbb{E} \left[\mathbf{1}_{F \geq i} \mathbb{E} \left[\int_{T_i}^{\min(T_{i+1}, P)} Y_t dt \middle| \mathcal{F}_{T_i} \right] \right] \\ &\leq \mathbb{E} [G_0 (S(T_0))] + \sum_{i=1}^{\infty} \mathbb{E} [\mathbf{1}_{F \geq i} (G_i + B \cdot \mathbb{E} [\min(T_{i+1}, P) - T_i | \mathcal{S}(T_i), F \geq i])] \\ &= \mathbb{E} [G_0 (S(T_0))] + B \cdot \mathbb{E} [P - T_0] + \sum_{i=1}^{\infty} G_i \Pr (F \geq i). \end{aligned}$$

Applying our final assumption to bound $\Pr (F \geq i)$,

$$\begin{aligned} \Pr (F \geq i) &= \Pr (F > 0) \prod_{j=2}^i \Pr (F \geq j | F \geq j - 1) \\ &= \Pr (F > 0) \prod_{j=2}^i \mathbb{E} [\Pr (F \geq j | F \geq j - 1, \mathcal{F}_{T_{j-1}})] \\ &\leq \Pr (F > 0) \prod_{j=2}^i \mathbb{E} [1 - p_j] \\ &= \Pr (F > 0) \prod_{j=2}^i (1 - p_j). \end{aligned}$$

Applying this final result, we find

$$\mathbb{E} \left[\int_{T_0}^P Y_t dt \right] \leq \mathbb{E} [G_0 (\mathcal{S}(T_0))] + \Pr (F > 0) \sum_{j=1}^{\infty} G_j \prod_{i=2}^j (1 - p_i) + B \cdot \mathbb{E} [P - T_0],$$

as desired. □

With Lemma 4.1 proven, we are ready to apply the MIST method to obtain our main results.

Chapter 5

The Lower Bounds

In this chapter, we discuss two of our results, our lower bounds on the average waiting time in the M/M/k/Setup-Deterministic. First, we discuss why these lower bounds are needed, then state both bounds, then prove the stronger and more recent bound.

5.1 Why we need a lower bound

From a provisioning standpoint, a lower bound tells us what system parameters are **necessary** to achieve a certain average waiting time. Accordingly, we now discuss our two lower bounds. The first lower bound that we present, Theorem 5.1 (the main result of [38]), is the first-ever result bounding the average waiting time in the M/M/k/Setup-Deterministic. Notably, it is also the first closed-form result bounding the average waiting time in any M/M/k/Setup system. The second lower bound that we present, Theorem 5.2 (one of the two main results in [37]), is an improvement of Theorem 5.1. The improved lower bound now applies to systems with an arbitrarily large numbers of servers k , removes an unnecessary and restrictive condition on the system parameters, and also has a far simpler proof. We state both theorems, but only prove the improved theorem.

5.2 The First Lower Bound

We now state the first lower bound for the average queue length in the M/M/k/Setup-Deterministic, from [38].

Theorem 5.1 (First Lower Bound On Average Queue Length). *In an M/M/k/Setup-Deterministic system with load ρ , setup time $\beta \geq 1000\frac{1}{\mu}$, and offered load $R \triangleq k\rho \geq 128$, if the setup time $\beta \geq \frac{1}{\mu} \log^2(k\rho)$, then the average queue length $\mathbb{E}[Q(\infty)]$ is lower bounded by*

$$\mathbb{E}[Q(\infty)] \geq \frac{\frac{1}{2}\beta^2 \frac{\mu\sqrt{R}}{2} + I^{\text{busy}} \left(\left[(\mu\beta - 1) \frac{\sqrt{R}}{2} - k(1 - \rho) \right]^+, k - R \right)}{C_1^{(\text{old})} \left(3\beta + \frac{1}{\mu} \right) + \beta + C_2^{(\text{old})} \frac{\mu\beta\sqrt{R}}{\mu k(1-\rho)} + C_3^{(\text{old})} \frac{1}{\mu} \log \left(C_2^{(\text{old})} \mu\beta\sqrt{R} \right)},$$

where $C_1^{(\text{old})}$, $C_2^{(\text{old})}$, and $C_3^{(\text{old})}$ are absolute constants.

5.3 The New Lower Bound

5.3.1 The New Lower Bound: Theorem Statement

After tightening and clarifying our techniques into the MIST method of Chapter 4, we obtained the following lower bound; its proof follows.

Theorem 5.2 (Improved Lower Bound on Average Queue Length). *For an $M/M/k$ /Setup-Deterministic with an offered load $R \triangleq k\rho \geq 100$ and a setup time $\beta \geq 100\frac{1}{\mu}$, the expected number of jobs in queue in steady state is lower-bounded as*

$$\mathbb{E}[Q(\infty)] \geq \frac{L_1\beta^2\sqrt{R} + I^{\text{busy}} \left(\left[L_1\beta\sqrt{R} - (k - R) \right]^+, k - R \right)}{2.08\beta + \frac{1}{\mu} \frac{F_1\beta\sqrt{R}}{k-R} + \frac{1}{\mu} \frac{3}{2} \ln(\beta) + \frac{1}{\mu} \ln(F_1 D_1) + \frac{2}{\mu} + \frac{1}{\mu} \left[D_2 + \frac{D_3}{\sqrt{R}} \right] \max \left(\frac{1}{D_1\sqrt{\mu\beta}}, \frac{1}{\sqrt{R}} \right)},$$

where $L_1, F_1, D_1, D_2,$ and D_3 are constants independent of system parameters.

5.3.2 The New Lower Bound: Proof Outline.

Basic Structure. We prove Theorem 5.2 via the MIST method. As noted in Chapter 4, we begin by applying the Renewal-Reward theorem to the queue length $Q(t)$, defining our renewal points as those points in time where the $(R + 1)$ -th server turns off. Defining time 0 to be one of these points, and defining the cycle time $X \triangleq \min \{t > 0 : Z(t^-) = R + 1, Z(t) = R\}$ as the next point, this gives

$$\mathbb{E}[Q(\infty)] = \frac{\mathbb{E} \left[\int_0^X Q(t) dt \right]}{\mathbb{E}[X]}.$$

To obtain our lower bound, it suffices to lower bound the numerator and upper bound the denominator of this fraction, i.e. lower bound $\mathbb{E} \left[\int_0^X Q(t) dt \right]$ and upper bound $\mathbb{E}[X]$. The time integral lower bound is handled by Lemma 5.1, which we state at the end of this section. The cycle length upper bound is split into two separate lemmas: Lemma 5.2 upper bounds the length of the cycle's "first part" and Lemma 5.3 bounds the length of its "second part".

Decomposition into phases. However, before we state or prove these lemmas, we first discuss the decomposition of the renewal cycle $[0, X)$ into two parts; one might think of this as a "miniature" application of the MIST method. We begin by noting that the end of the renewal cycle is moment when the $(R + 1)$ -th server turns off. Since the $(R + 1)$ -th server is off at the start of a renewal period, we can break the renewal cycle into two phases based on whether the $(R + 1)$ -th server has turned on yet. Formally, we define the accumulation time $T_A \triangleq \min \{t > 0 : Z(t) = R + 1\}$ as the first moment that the $(R + 1)$ -th server turns on. From here, we can focus separately on the accumulation phase, from time 0 to time T_A , and the draining phase, from time T_A to time X .

With this decomposition, we can now state our main lemmas. Their proofs follow in sequence afterwards.

Lemma 5.1 (Lower bound on Cycle Integral). *Define busy period integral $I^{\text{busy}}(x, z)$ as*

$$I^{\text{busy}}(x, z) \triangleq \frac{x}{\mu z} \left[\frac{x+1}{2} + \frac{1}{1 - \frac{k\lambda}{k\lambda + \mu z}} \right] = \frac{x}{\mu z} \left[\frac{x+1}{2} + \frac{R}{z} \right].$$

For the time integral of the queue length $Q(t)$ over a renewal cycle, we have

$$\mathbb{E} \left[\int_0^X Q(t) dt \right] \geq \frac{1}{2} \mu \beta^2 L_1 \sqrt{R} + I^{\text{busy}} \left(\left[\mu \beta L_1 \sqrt{R} - (k - R) \right]^+, k - R \right).$$

Lemma 5.2 (Upper bound on Accumulation Phase Length). *Recall that*

$$T_A \triangleq \min \{ t > 0 : Z(t) \geq R + 1 \}$$

is the amount of time until the $(R + 1)$ -th server turns on. Then we can bound the expectation $\mathbb{E}[T_A]$ by

$$\mathbb{E}[T_A] \leq e^{\frac{1}{24R\beta}} \left(\sqrt{1 + \frac{1}{2R\beta}} \right) \left[1 + \frac{e^{\frac{1}{12R}}}{\sqrt{2\mu\beta}} \right] \beta \leq 1.08 * \beta.$$

Lemma 5.3 (Upper bound on Draining Phase Length). *Recall that the accumulation time*

$$T_A \triangleq \min \{ t > 0 : Z(t) \geq R + 1 \}$$

is the amount of time until the $(R + 1)$ -th server turns on and the cycle time X is the moment when it turns off. Then, one can bound $\mathbb{E}[X - T_A]$ by

$$\mathbb{E}[X - T_A] \leq \beta + \frac{1}{\mu} \frac{F_1 \beta \sqrt{R}}{k - R} + \frac{1}{\mu} \frac{3}{2} \ln(\beta) + \frac{1}{\mu} \ln(F_1 D_1) + \frac{2}{\mu} + \left[D_2 + \frac{D_3}{\sqrt{R}} \right] \max \left(\frac{1}{D_1 \sqrt{\mu\beta}}, \frac{1}{\sqrt{R}} \right),$$

where $F_1, D_1, D_2,$ and D_3 are constants not depending on system parameters.

5.3.3 Proof of Lemma 5.1: Lower Bound on Cycle Integral.

Lemma 5.1 Proof Outline

Basic Strategy. First, we split the first phase $[0, T_A)$ into epochs, where epoch i begins when the number of busy servers $Z(t)$ first drops to $R - i$, and an epoch ends either when the next epoch starts or when the first phase ends. Our goal will be to analyze a specific ‘‘significant’’ epoch. In particular, we say that an epoch is long if it lasts for longer than a setup time β . Because the accumulation phase ends when the $(R + 1)$ -th server turns on, at least one epoch must be long. We use L to denote the index of the *first* long epoch. From here, we argue via a martingale/coupling argument that the expected time integral over the first β time in epoch L is at least $\frac{1}{2} \beta^2 \mathbb{E}[L]$. To bound the integral afterwards, we couple the behavior of the total number of jobs $N(t)$ to the queue length in an M/M/1 queue with arrival rate $k\lambda$ and departure rate $k\mu$.

Formalization. Define the stopping time $\tau_i \triangleq \min \{t \geq 0 : N(t) \leq R - u\}$ as the beginning of epoch i . We say that the epoch *occurs* is $\tau_i < T_A$, and define the end of epoch i as $\gamma_i \triangleq \min(\tau_{i+1}, T_A)$ the moment when either epoch $i + 1$ begins or when the first phase ends. If epoch i occurs, we say it is long if $\gamma_i - \tau_i \geq \beta$. Let $L \triangleq \min \{i \in \mathbb{N} : \gamma_i - \tau_i \geq \beta\}$ be the index of the first long epoch. It suffices to show two claims; we state and prove them in sequence.

5.3.4 Lower Bound on Integral until $\tau_L + \beta$.

We show the following claim.

Claim 5.1. *Let L be the index of the first long epoch. Then,*

$$\mathbb{E} \left[\int_0^{\tau_L + \beta} Q(t) dt \right] \geq \frac{1}{2} \mu \beta^2 L_1 \sqrt{R}, \quad (5.1)$$

where L_1 is some absolute constant.

Claim 5.1 Proof Strategy. First, we show that the initial integral is bounded by

$$\mathbb{E} \left[\int_0^{\tau_L + \beta} Q(t) dt \right] \geq \frac{1}{2} \mu \beta^2 \mathbb{E}[L]. \quad (5.2)$$

Afterwards, we give a bound on $\mathbb{E}[L]$, showing that

$$\mathbb{E}[L] \geq L_1 \sqrt{R}. \quad (5.3)$$

Proof of (5.2), Bound in terms of $\mathbb{E}[L]$.

To show (5.2), we first condition on whether $L \geq i$, giving

$$\begin{aligned} \mathbb{E} \left[\int_0^{\tau_L + \beta} Q(t) dt \right] &= \sum_{i=0}^{\infty} \mathbb{E} \left[\int_{\tau_i}^{\min(\tau_i + \beta, \tau_{i+1})} Q(t) dt \mathbf{1}_{L \geq i} \right] \\ &= \sum_{i=0}^{\infty} \mathbb{E} \left[\int_{\tau_i}^{\min(\tau_i + \beta, \tau_{i+1})} Q(t) dt \middle| \mathcal{F}_{\tau_i} \right] \Pr(L \geq i). \end{aligned}$$

To further develop this conditional expectation, we note that during the interval $[\tau_i, \min(\tau_i + \beta, \tau_{i+1}))$, the system must have exactly $Z(t) = R - i$ busy servers running, meaning that $Q(t) = N(t) - (R - i)$. Defining a coupled process $\tilde{Q}(t)$ as

$$\tilde{Q}(t) = A(\tau_i, t) - \mathcal{D}[R - i](\tau_i, t),$$

we see that $Q(t)$ and $\tilde{Q}(t)$ coincide during the interval in question. Moreover, one can redefine the stopping time $\gamma = \tau_{i+1}$ as $\min \{t > \tau_i : \tilde{Q}(t) = -1\}$. Noting that $Q(\min(\gamma, t)) = -1$ for

any time $t > \gamma$, we find that

$$\begin{aligned}
\int_{\tau_i}^{\min(\tau_i+\beta, \tau_{i+1})} Q(t) dt &= \int_{\tau_i}^{\min(\tau_i+\beta, \tau_{i+1})} \tilde{Q}(t) dt \\
&= \int_{\tau_i}^{\min(\tau_i+\beta, \tau_{i+1})} \tilde{Q}(\min(t, \tau_{i+1})) dt + \int_{\min(\tau_i+\beta, \tau_{i+1})}^{\tau_i+\beta} \left(\tilde{Q}(\min(t, \tau_{i+1})) + 1 \right) dt \\
&= \int_{\tau_i}^{\tau_i+\beta} \tilde{Q}(\min(t, \tau_{i+1})) dt + [\beta - \min(\beta, \tau_{i+1} - \tau_i)] \\
&\geq \int_{\tau_i}^{\tau_i+\beta} \tilde{Q}(\min(t, \tau_{i+1})) dt.
\end{aligned}$$

Taking the conditional expectation at time τ_i , we find

$$\mathbb{E} \left[\int_{\tau_i}^{\tau_i+\beta} \tilde{Q}(\min(t, \tau_{i+1})) dt \middle| \mathcal{F}_{\tau_i} \right] = \int_{\tau_i}^{\tau_i+\beta} \mathbb{E} \left[\tilde{Q}(\min(t, \tau_{i+1})) \middle| \mathcal{F}_{\tau_i} \right] dt.$$

Noting that $V_L(t) = \tilde{Q}(t) - \mu i [t - \tau_i]$ is a martingale, and that $\min(t, \tau_{i+1})$ is an almost-surely bounded stopping time, we have that

$$\begin{aligned}
\tilde{Q}(\tau_i) &= V_L(\tau_i) = 0 \\
&= \mathbb{E} [V_L(\min(t, \tau_{i+1})) | \tau_i] \\
&= \mathbb{E} \left[\tilde{Q}(\min(t, \tau_{i+1})) \middle| \mathcal{F}_{\tau_i} \right] - \mu i \mathbb{E} [\min(t, \tau_{i+1}) | \mathcal{F}_{\tau_i}].
\end{aligned}$$

Since

$$\mathbb{E} [\min(t, \tau_{i+1}) | \mathcal{F}_{\tau_i}] \geq t \cdot \Pr(\tau_{i+1} - \tau_i \geq t) \geq t \cdot \Pr(\tau_{i+1} - \tau_i \geq \beta) = t \Pr(L = i | L \geq i),$$

we have

$$\begin{aligned}
\Pr(L \geq i) \mathbb{E} \left[\int_{\tau_i}^{\min(\tau_i+\beta, \tau_{i+1})} Q(t) dt \middle| L \geq i \right] &\geq \Pr(L \geq i) \mathbb{E} \left[\int_{\tau_i}^{\tau_i+\beta} \tilde{Q}(\min(t, \tau_{i+1})) dt \middle| L \geq i \right] \\
&\geq \int_{\tau_i}^{\tau_i+\beta} \mu i t \Pr(L = i) \\
&= \mu \frac{\beta^2}{2} i \Pr(L = i)
\end{aligned}$$

Summing across all i , we obtain (5.2).

Proof Sketch for (5.3), bound on $\mathbb{E}[L]$.

We defer the full proof of this to Section A.4.6, and for now give a proof sketch.

We prove (5.3) by first showing that

$$\Pr(L > j | L \geq j) \geq \left(1 - \frac{j}{R}\right) \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right),$$

where $b_1 = \frac{2}{\sqrt{\pi}}$. Next, we show that this implies that, for any $\delta \in (0, 1)$ and any $j < \delta R$,

$$\Pr(L > j) \geq \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right)^{j+1} e^{-\frac{j(j+1)}{2R} \frac{1}{1-\delta}}.$$

From here, we use the sum of tails formula $\mathbb{E}[L] = \sum_{j=0}^{\infty} \Pr(L > j)$ to show

$$\mathbb{E}[L] \geq \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right) \left(\left[\sqrt{\frac{\pi}{2}}(1-\delta) - \frac{1.15(1-\delta)}{\sqrt{\mu\beta}} \right] \sqrt{R} - \frac{1}{2} - \frac{2(1-\delta)}{\delta} e^{-R \frac{\delta^2}{1-\delta}} \right).$$

Choosing $\delta = \frac{2}{\sqrt{R}}$ then noting that $\mu\beta \geq 100$ and $R \geq 100$ gives the result.

5.3.5 Lower Bound on Integral after $\tau_L + \beta$.

To finish our lower bound on the integral, we now show the following claim.

Claim 5.2. *Let L be the index of the first long epoch. Then,*

$$\mathbb{E} \left[\int_{\tau_L + \beta}^X Q(t) dt \right] \geq I^{\text{busy}} \left(\left[\mu\beta L_1 \sqrt{R} - (k - R) \right]^+, k - R \right) \quad (5.4)$$

where L_1 is some absolute constant.

Claim 5.2: Proof Strategy. First, we show that the remaining integral is bounded by

$$\mathbb{E} \left[\int_{\tau_L + \beta}^X Q(t) dt \right] \geq I^{\text{busy}} \left([\mathbb{E}[N(\tau_L + \beta)] - k]^+, k - R \right). \quad (5.5)$$

Then, we use martingales again to show that

$$\mathbb{E}[N(\tau_L + \beta)] \geq R + \mu\beta \mathbb{E}[L]. \quad (5.6)$$

Applying (5.3), our bound on $\mathbb{E}[L]$, we obtain the result.

Proof of (5.5), Bound in terms of $\mathbb{E}[N(\cdot)]$.

To prove (5.5), we make a simple coupling argument. Let $\eta_k \triangleq \min \{t \geq \tau_L + \beta : N(t) \leq k\}$. Since the draining phase starts at $T_A \geq \tau_L + \beta$ and the end of the cycle $X = \min \{t \geq T_A : N(t) \leq R\}$, we know that $X \geq \eta_k$. Moreover, we know the number of busy servers $Z(t) \leq k$; it follows by Claim A.1 that we can define $\tilde{N}(t)$ as

$$\tilde{N}(t) \triangleq N(\tau_L + \beta) + A(\tau_L + \beta, t) - \mathcal{D}[k]((\tau_L + \beta, t))$$

and have $\tilde{N}(t) \leq N(t)$ for any $t > \tau_L + \beta$. Even further, we can defined a coupled hitting time $\tilde{\eta}_k \triangleq \min \left\{ t > \tau_L + \beta : \tilde{N}(t) \leq k \right\}$ which must happen before η_k . In other words,

$$\begin{aligned} \int_{\tau_L + \beta}^X Q(t) dt &\geq \int_{\tau_L + \beta}^{\eta_k} Q(t) dt \\ &\geq \int_{\tau_L + \beta}^{\eta_k} [N(t) - k] dt \\ &\geq \int_{\tau_L + \beta}^{\tilde{\eta}_k} [\tilde{N}(t) - k] dt. \end{aligned}$$

This final term is just the time integral of the number of jobs in system over a M/M/1 busy period started by $[N(\tau_L + \beta) - k]^+$ jobs, where jobs arrive at rate $k\lambda$ and depart at rate $k\mu$. Accordingly, we have

$$\begin{aligned} \mathbb{E} \left[\int_{\tau_L + \beta}^X Q(t) dt \right] &\geq \mathbb{E} [I^{\text{busy}} ([N(\tau_L + \beta) - k]^+, k - R)] \\ &\geq I^{\text{busy}} (\mathbb{E} [N(\tau_L + \beta) - k]^+, k - R) \\ &\geq I^{\text{busy}} (\mathbb{E} [[N(\tau_L + \beta) - R - (k - R)]^+], k - R), \end{aligned}$$

where in the last two lines we have applied Jensen's inequality. \square

Proof of (5.6), Bound on $\mathbb{E} [N(\cdot)]$.

To bound $\mathbb{E} [N(\tau_L + \beta)]$, we condition on the value of L , then make a martingale argument.

$$\begin{aligned} \mathbb{E} [N(\tau_L + \beta)] &= \sum_{i=0}^R \mathbb{E} [N(\tau_i + \beta) \mathbf{1}_{L=i}] \\ &= \sum_{i=0}^R \mathbb{E} [N(\tau_i + \beta) \mathbf{1}_{L=i, L \geq i}] \\ &\geq \sum_{i=0}^R \Pr(L \geq i) \mathbb{E} [N(\tau_i + \beta) \mathbf{1}_{L=i} | L \geq i]. \end{aligned}$$

Continuing with this conditional expectation,

$$\begin{aligned} &\mathbb{E} [N(\tau_i + \beta) \mathbf{1}_{L=i} | \mathcal{F}_{\tau_i}] \\ &= \mathbb{E} [N(\tau_i + \beta) \mathbf{1}_{\tau_i + \beta < \tau_{i+1}} | \mathcal{F}_{\tau_i}] \\ &= \mathbb{E} [[N(\tau_i + \beta) - (R - (i + 1))] \mathbf{1}_{\tau_i + \beta < \tau_{i+1}} | \mathcal{F}_{\tau_i}] + (R - i - 1) \Pr(\tau_i + \beta < \tau_{i+1} | \mathcal{F}_{\tau_i}) \\ &= \mathbb{E} [N(\min(\tau_i + \beta, \tau_{i+1})) - (R - (i + 1)) | \mathcal{F}_{\tau_i}] + (R - i - 1) \Pr(\tau_i + \beta < \tau_{i+1} | \mathcal{F}_{\tau_i}) \\ &= 1 + \mu i \mathbb{E} [\min(\beta, \tau_{i+1} - \tau_i)] + (R - i - 1) \Pr(\tau_i + \beta < \tau_{i+1}) \\ &\geq 1 + \mu i \beta \Pr(\tau_{i+1} - \tau_i \geq \beta) + (R - i - 1) \Pr(\tau_i + \beta < \tau_{i+1}) \\ &= 1 + \mu i \beta \Pr(L = i | L \geq i) + (R - i - 1) \Pr(L = i | L \geq i). \end{aligned} \tag{5.7}$$

Summing across i , we find

$$\begin{aligned}\mathbb{E}[N(\tau_L + \beta)] &= \sum_{i=0}^R \Pr(L \geq i) \quad ((5.7)) \\ &= (1 + \mathbb{E}[L]) + (\mu\beta\mathbb{E}[L]) + (R - \mathbb{E}[L] - 1) \\ &= R + \mu\beta\mathbb{E}[L],\end{aligned}$$

as desired. \square

Combining Claims 5.1 and 5.2, we obtain a lower bound on $\mathbb{E}\left[\int_0^X Q(t)dt\right]$, proving Lemma 5.1.

5.3.6 Proof of Lemma 5.2: Upper Bound on the Accumulation Time $\mathbb{E}[T_A]$.

Defining a coupling. To prove Lemma 5.2, we first note that, during the accumulation phase, we have two bounds on the number of busy servers $Z(t)$: it must be less than the total number of jobs $N(t)$ and it must be less than R ; the former because every busy server must be working on a job, and the latter because otherwise the accumulation phase would be over. Thus, we can define a coupled $M/M/R$ system for which the number of jobs $\tilde{N}(t)$ in the coupled system is always at least the number of jobs $N(t)$ in the original system.

How we use the coupling. To use this coupled process to bound $\mathbb{E}[T_A]$, recall that the accumulation point T_A is the first time the $(R+1)$ -th server turns on. Accordingly, one can also think of this as the first time that there has been at least $R+1$ jobs in the system for β time. Thus, if we define a coupled accumulation point $\tilde{T}_A \triangleq \min\left\{t \geq \beta : \min_{s \in [t-\beta, t]} \tilde{N}(t) \geq R+1\right\}$, then we know $\tilde{T}_A \geq T_A$. In other words, it suffices to bound $\mathbb{E}[\tilde{T}_A]$.

General Strategy. We bound $\mathbb{E}[\tilde{T}_A]$ using the MIST method of Lemma 4.1. As such, we define a few stopping times, then list the preconditions/claims that we will satisfy to complete the proof of Lemma 5.2.

Definition of γ and α . Let the initial cycle-downcrossing occur at $\alpha_0 \triangleq 0$ and iteratively define the upcrossings γ and downcrossings α as

$$\gamma_i \triangleq \min\left\{t \geq \alpha_i : \tilde{N}(t) \geq R+1\right\}$$

and

$$\alpha_{i+1} \triangleq \min\left\{t \geq \gamma_i : \tilde{N}(t) \geq R+1\right\}.$$

Application of Lemma 4.1, the IST Lemma. Applying Lemma 4.1 using $0 = \alpha_0$ as our starting point, the coupled accumulation point \tilde{T}_A as our ending point, our test function as $Y_t = 1$, and the cycle-upcrossings (α_i) as our intervening stopping times, we now must prove that

$$\mathbb{E} [\gamma_i - \alpha_i | n_\alpha \geq i] \leq \frac{1}{\mu} e^{\frac{1}{12R}} \sqrt{1 + \frac{1}{R} \frac{\sqrt{2\pi}}{\sqrt{R}}} \leq \frac{c_3}{\mu\sqrt{R}}, \quad (5.8)$$

$$\mathbb{E} \left[\min \left(\tilde{T}_A, \alpha_{i+1} \right) - \gamma_i \mid n_\alpha \geq i \right] \leq b_1 \sqrt{\frac{\beta}{\mu R}} + \frac{6}{\mu R}, \quad (5.9)$$

and

$$\Pr (n_\alpha \geq i + 1 | n_\alpha \geq i) \leq 1 - \frac{b_1}{\sqrt{2}} e^{-\frac{1}{3(\mu^2 R \beta - 1)}} \frac{1}{\sqrt{\mu^2 R \beta + 2}} \leq 1 - \frac{b_1 c_4}{\sqrt{\mu R \beta}}, \quad (5.10)$$

where $b_1 \triangleq \sqrt{\frac{2}{\pi}}$, $c_3 = 1.001\sqrt{2\pi}$, and $c_4 = 0.499$.

Completion of Proof, assuming (5.8), (5.9), and (5.10). Applying Lemma 4.1, one finds that

$$\begin{aligned} \mathbb{E} [\tilde{T}_A] &= \sum_{i=0}^{\infty} \mathbb{E} \left[\min \left(\tilde{T}_A, \alpha_{i+1} \right) - \alpha_i \mid n_\alpha \geq i \right] \Pr (n_\alpha \geq i) \\ &\leq \left[\frac{c_3}{\mu\sqrt{R}} + \frac{b_1\sqrt{\beta}}{\sqrt{\mu R}} + \frac{6}{\mu R} \right] \frac{\sqrt{\mu R \beta}}{b_1 c_4} \\ &= \frac{1}{\mu} \left[\frac{c_3}{b_1 c_4} \sqrt{\mu \beta} + \frac{1}{c_4} \beta + \frac{6}{b_1 c_4} \sqrt{\frac{\mu \beta}{R}} \right]. \end{aligned}$$

Proof of (5.8): Upper bound on initial up-crossing time.

To prove (5.8), we note that, since our coupled system is an $M/M/R$, the expected time $\mathbb{E} [\gamma_i - \alpha_i | n_\alpha \geq i]$ is simply the expected passage time from state R to $(R + 1)$ in an $M/M/R$ (and equivalently an

$M/M/R/(R+1)$, an $M/M/R$ which can contain only $R+1$ jobs. Solving, one finds that

$$\begin{aligned}
\mathbb{E} [T_{R \rightarrow (R+1)}] &\leq \mathbb{E} [T_{(R+1) \rightarrow (R+1)}] \\
&= \frac{1}{\mu(R+1)} \frac{1}{\pi_{R+1}} \\
&= \frac{1}{\mu(R+1)} \frac{\sum_{i=0}^{R+1} \frac{R^i}{i!}}{\frac{R^{R+1}}{(R+1)!}} \\
&\leq \frac{1}{\mu(R+1)} e^R \frac{(R+1)!}{R^{R+1}} \\
&\leq \frac{1}{\mu(R+1)} e^R \frac{e^{\frac{1}{12(R+1)}} \sqrt{2\pi(R+1)} (R+1)^{R+1} e^{-(R+1)}}{R^{R+1}} \\
&= e^{\frac{1}{12R}} \frac{1}{\mu} \sqrt{2\pi} \frac{\sqrt{R+1}}{R} \left(1 + \frac{1}{R}\right)^R e^{-1} \\
&\leq \frac{1}{\mu} e^{\frac{1}{12R}} \sqrt{1 + \frac{1}{R}} \frac{\sqrt{2\pi}}{\sqrt{R}} \\
&\leq \frac{1}{\mu} 1.006 \frac{\sqrt{2\pi}}{\sqrt{R}} \\
&\triangleq \frac{c_3}{\mu\sqrt{R}}
\end{aligned}$$

Proof of (5.9): Bound on time between up-crossings.

To bound the expected time $\mathbb{E} \left[\min \left(\tilde{T}_A, \alpha_{i+1} \right) - \gamma_i \mid n_\alpha \geq i \right]$, we first note that, if $\gamma_i + \beta \leq \alpha_{i+1}$, then $\tilde{T}_A = \gamma_i + \beta$. Likewise, if $\gamma_i + \beta > \alpha_{i+1}$, then $\tilde{T}_A > \alpha_{i+1}$. It follows that, given that $n_\alpha \geq i$, the time $\min \left(\tilde{T}_A, \alpha_{i+1} \right) = \min \left(\beta + \gamma_i, \alpha_{i+1} \right)$. Thus, we have that

$$\mathbb{E} \left[\min \left(\tilde{T}_A, \alpha_{i+1} \right) - \gamma_i \mid n_\alpha \geq i \right] = \mathbb{E} \left[\min \left(\beta, \alpha_{i+1} - \gamma_i \right) \mid n_\alpha \geq i \right] = \int_0^\beta \Pr \left(\alpha_{i+1} - \gamma_i > s \mid n_\alpha \geq i \right) ds.$$

We continue by bounding this tail probability. To begin, note that, while $\tilde{N}(t)$ stays above $R+1$, the dynamics of \tilde{N} are precisely that of a critically-loaded $M/M/1$ queue with arrival rate and departure rate equal to $k\lambda$. The tail probability we are interested in bounding is precisely the probability that a busy period (started with 1 job) in such a system lasts longer than s time. Applying Claim A.6, one finds that, for any $t \geq \frac{3}{\mu 2R}$,

$$\Pr \left(\alpha_{i+1} - \gamma_i > s \mid n_\alpha \geq i \right) \leq b_1 \left(\frac{1}{\sqrt{\mu 2R s}} + \frac{b_2}{(\mu 2R s)^{3/2}} \right).$$

Integrating, we find that

$$\begin{aligned}
\int_0^\beta \Pr(\alpha_{i+1} - \gamma_i > s | n_\alpha \geq i) \, ds &\leq \frac{3}{\mu 2R} + \frac{b_1}{\sqrt{2}} \int_{\frac{3}{\mu 2R}}^\beta \frac{1}{\sqrt{\mu 2R s}} + \frac{b_2}{(\mu 2R s)^{3/2}} \, ds \\
&\leq \frac{3}{\mu 2R} + \frac{b_1}{\sqrt{2}} \left[\sqrt{\frac{2\beta}{\mu R}} + b_2 \sqrt{\frac{2}{3}} \frac{1}{\mu R} \right] \\
&= b_1 \sqrt{\frac{\beta}{\mu R}} + \left(\frac{3}{2} + (b_1 + 2.5) \sqrt{\frac{2}{3}} \right) \frac{1}{\mu R} \\
&\leq \frac{2}{\sqrt{\pi}} \sqrt{\frac{\beta}{\mu R}} + \frac{6}{\mu R}.
\end{aligned}$$

□

Proof of (5.10): Bound on probability of another γ up-crossing.

To prove (5.10), it suffices to note that, upon conditioning on the filtration at γ_i , the probability $\Pr(n_\alpha \geq i + 1 | n_\alpha \geq i)$ is simply the probability that a busy period in a critically-loaded M/M/1, with arrival and departure rate equal to μR , ends before β time has passed. Applying Claim A.6, one finds that this is

$$\Pr(n_\alpha \geq i + 1 | n_\alpha \geq i) \geq 1 - \frac{b_1}{\sqrt{2}} e^{-\frac{1}{3(\mu 2R\beta - 1)}} \frac{1}{\sqrt{\mu 2R\beta + 2}}.$$

□

5.3.7 Proof of Lemma 5.3: Upper Bound on the Remaining Cycle Time $\mathbb{E}[X - T_A]$.

We now prove the upper bound on $\mathbb{E}[X - T_A]$. We make use of the “wait-busy” idea from Section 6.2.2 as well as our main tool, Lemma 4.1. As such, we begin by defining some stopping times.

Definition of $v_i^{(\text{down})}$ and $v_i^{(\text{up})}$. Recall that the draining phase begins at time T_A . Let $M_L \triangleq \min\left(k - R, \max\left(\frac{\sqrt{R}}{D_1 \sqrt{\beta}}, 1\right)\right)$ be a specially-set analysis threshold. Let the stopping time $v_1^{(\text{down})} \triangleq \min\{t \geq T_A : N(t) < R + M_L\}$ be the first time the number of jobs $N(t)$ drops below $R + M_L$, and recursively define

$$v_i^{(\text{up})} \triangleq \min\left\{t \geq v_i^{(\text{down})} : N(t) \geq R + M_L\right\}$$

and

$$v_{i+1}^{(\text{down})} \triangleq \min\left\{t \geq v_i^{(\text{up})} : N(t) < R + M_L\right\}.$$

Specification Step. Now, we apply Lemma 4.1 using the accumulation point T_A as our initial point, the cycle end X as our ending point, the constant function $Y_t = 1$ as our test function, and the draining-downcrossing points $(v_i^{(\text{down})})$ as our intervening points; we use n_ζ to count the number of intervening points. To complete the proof, we must show that the following claims:

$$\mathbb{E} \left[v_1^{(\text{down})} - T_A \right] \leq \beta + \frac{1}{\mu} \frac{F_1 \beta \sqrt{R}}{k - R} + \frac{1}{\mu} \frac{3}{2} \ln(\beta) + \frac{1}{\mu} \ln(F_1 D_1), \quad (5.11)$$

$$\mathbb{E} \left[\min \left(X, v_{i+1}^{(\text{down})} \right) - v_i^{(\text{down})} \mid n_\zeta \geq i \right] \leq \frac{D_2}{\mu \sqrt{R}} + \frac{D_3}{\mu R} + \frac{2}{\mu M_L} \quad (5.12)$$

$$\Pr(n_\zeta \geq i + 1 \mid n_\zeta \geq i) \leq \frac{1}{M_L}. \quad (5.13)$$

Completion of Proof assuming (5.11), (5.12), and (5.13). Before proving the claims, we now prove the lemma. It suffices to give a bound on $\mathbb{E} \left[X - v_1^{(\text{down})} \right]$; applying Lemma 4.1 gives

$$\begin{aligned} \mathbb{E} \left[X - v_1^{(\text{down})} \right] &\leq M_L \left[\frac{D_2}{\mu \sqrt{R}} + \frac{D_3}{\mu R} + \frac{2}{\mu M_L} \right] \\ &= \frac{2}{\mu} + \left[D_2 + \frac{D_3}{\sqrt{R}} \right] \max \left(\frac{1}{D_1 \sqrt{\mu \beta}}, \frac{1}{\sqrt{R}} \right). \end{aligned}$$

Proof of (5.11): Upper bound on time until first downward visit.

To bound $\mathbb{E} \left[v_1^{(\text{down})} - T_A \right]$, we make a coupling argument then apply basic results on M/M/1 busy periods. Moreover, instead of proving (5.11) directly, we first show a more general claim.

Claim 5.3. For $M_L \leq j \leq N(T_A) - R$, define η_j as the first time after T_A that $N(t) \leq R + j$. Note that this means that $\eta_{N(T_A)} = T_A$ and $\eta_{M_L} = v_1^{(\text{down})}$. Then we have the following bound:

$$\mathbb{E} \left[\eta_{M_L} - \eta_j \mid \mathcal{F}_{\eta_j} \right] \leq Y_{R+j}(\eta_j) + \frac{1}{\mu} \sum_{i=M_L}^j \frac{1}{\min(i, k - R)}.$$

Afterwards, we complete the proof by noting that $[N(T_A) - k]^+ \leq [N(T_A) - R]$, taking expectations, applying Jensen's inequality to the minimum function and the $\ln(\cdot)$ (which is concave), using the bound on $\mathbb{E} [N(T_A) - R]$ from Claim 6.11, then letting $h = M_L$.

Proof of Claim 5.3. We prove Claim 5.3 by induction. In the base case, suppose that $j = M_L + 1$. Note that at time η_{M_L+1} , the numbers of jobs $N(\eta_{M_L+1}) = R + M_L + 1$ and the remaining time until the $(R + M_L + 1)$ -th server turns on is $Y_{R+M_L+1}(\eta_{M_L+1})$. As such, we can simply wait until either that server turns on, in which case we can analyze the system as an M/M/1 busy period with departure rate $\mu \min(R + M_L + 1, k)$, or the number of jobs $N(t)$ drops below $R + M_L + 1$ on its own. In other words, (using j here to save space)

$$\mathbb{E} \left[\eta_{j-1} - \eta_j \mid \mathcal{F}_{\eta_j} \right] \leq Y_{R+j}(\eta_j) + \frac{\mathbb{E} \left[[N(\eta_j + Y_{R+j}(\eta_j)) - (R + (j - 1))] \mathbf{1}_{\eta_{j-1} > \eta_j + Y_{R+j}(\eta_j)} \mid \mathcal{F}_{\eta_j} \right]}{\mu \min(j, k - R)}.$$

Now, we reframe the expectation as an expectation up to a stopping time. We note that, if $\eta_{j-1} > \eta_j + Y_{R+j}(\eta_j)$, then we have that

$$N(\eta_j + Y_{R+j}(\eta_j)) = N(\min(\eta_j + Y_{R+j}(\eta_j), \eta_{j-1})).$$

Likewise, if $\eta_{j-1} \leq \eta_j + Y_{R+j}(\eta_j)$, then

$$R + j - 1 = N(\eta_{j-1}) = N(\min(\eta_j + Y_{R+j}(\eta_j), \eta_{j-1})).$$

Using this and applying a simple coupling argument, one sees that

$$\begin{aligned} & \mathbb{E} \left[[N(\eta_j + Y_{R+j}(\eta_j)) - (R + (j - 1))] \mathbf{1}_{\eta_{j-1} > \eta_j + Y_{R+j}(\eta_j)} \middle| \mathcal{F}_{\eta_j} \right] \\ &= \mathbb{E} \left[N(\min(\eta_j + Y_{R+j}(\eta_j), \eta_{j-1})) - (R + j - 1) \middle| \mathcal{F}_{\eta_j} \right] \\ &\leq N(\eta_j) - (R + j - 1) = 1. \end{aligned}$$

Thus, we find that

$$\mathbb{E} [\eta_{j-1} - \eta_j \middle| \mathcal{F}_{\eta_j}] \leq Y_{R+j}(\eta_j) + \frac{1}{\mu \min(j, k - R)}.$$

Inductive case. The inductive case proceeds in much the same way, except now, if $N(t)$ does drop below $R + j$ “early”, then we can factor in the time that has elapsed in the value of $Y_{R+j}(\eta_j)$. In particular, note that, since the $(R + j)$ -th server would have already turned on,

$$\mathbb{E} [\eta_{M_L} - \eta_j \middle| \mathcal{F}_{\eta_j}] \mathbf{1}_{\eta_j \geq \eta_{j+1} + Y_{R+j+1}(\eta_{j+1})} \leq \frac{1}{\mu} \sum_{i=M_L}^j \frac{1}{\mu \min(i, k - R)} \mathbf{1}_{\eta_j \geq \eta_{j+1} + Y_{R+j+1}(j+1)}.$$

It follows that

$$\mathbb{E} [\eta_{M_L} - \eta_j \middle| \mathcal{F}_{\eta_j}] \leq Y_{R+j}(\eta_j) \mathbf{1}_{\eta_j < \eta_{j+1} + Y_{R+j+1}(\eta_{j+1})} + \frac{1}{\mu} \sum_{i=M_L}^j \frac{1}{\mu \min(i, k - R)}.$$

Now, we note that

$$\begin{aligned} Y_{R+j}(\eta_j) \mathbf{1}_{\eta_j < \eta_{j+1} + Y_{R+j+1}(\eta_{j+1})} &= [Y_{R+j}(\eta_j) + \eta_j - \eta_j] \mathbf{1}_{\eta_j < \eta_{j+1} + Y_{R+j+1}(\eta_{j+1})} \\ &= [Y_{R+j}(\eta_{j+1}) + \eta_{j+1} - \eta_j] \mathbf{1}_{\eta_j < \eta_{j+1} + Y_{R+j+1}(\eta_{j+1})} \\ &\leq [Y_{R+j+1}(\eta_{j+1}) + \eta_{j+1} - \eta_j] \mathbf{1}_{\eta_j < \eta_{j+1} + Y_{R+j+1}(\eta_{j+1})} \\ &= [Y_{R+j+1}(\eta_{j+1}) + \eta_{j+1} - \eta_j]^+, \end{aligned}$$

so that we find

$$\mathbb{E} [\eta_{M_L} - \eta_j \middle| \mathcal{F}_{\eta_j}] \leq [Y_{R+j+1}(\eta_{j+1}) + \eta_{j+1} - \eta_j]^+ + \frac{1}{\mu} \sum_{i=M_L}^j \frac{1}{\mu \min(i, k - R)}.$$

Finally, we note that

$$\mathbb{E} [\eta_j - \eta_{j+1} \middle| \mathcal{F}_{\eta_{j+1}}] \leq \mathbb{E} [\min(\eta_j - \eta_{j+1}, Y_{R+j+1}(\eta_{j+1})) \middle| \mathcal{F}_{\eta_{j+1}}] + \frac{1}{\mu \min(j + 1, k - R)}.$$

Summing these final two expressions gives the inductive result, proving Claim 5.3.

Using Claim 5.3. Thus, we obtain that, using H_i to denote the i -th harmonic number,

$$\begin{aligned}\mathbb{E} \left[v_1^{(\text{down})} - T_A \middle| \mathcal{F}_{T_A} \right] &\leq \beta + \frac{1}{\mu} \frac{[N(T_A) - k]^+}{k - R} + \frac{1}{\mu} [H_{\min(N(T_A) - R, k - R)} - H_{M_L}] \\ &\leq \beta + \frac{1}{\mu} \frac{[N(T_A) - R]^+}{k - R} + \frac{1}{\mu} \ln \left(\frac{\min(N(T_A) - R, k - R)}{M_L} \right).\end{aligned}$$

Taking expectations and applying Jensen's inequality twice, we find

$$\begin{aligned}\mathbb{E} \left[v_1^{(\text{down})} - T_A \middle| \mathcal{F}_{T_A} \right] &\leq \beta + \frac{1}{\mu} \frac{F_1 \mu \beta \sqrt{R}}{k - R} + \frac{1}{\mu} \ln \left(\frac{F_1 \mu \beta \sqrt{R}}{M_L} \right) \\ &\leq \beta + \frac{1}{\mu} \frac{F_1 \mu \beta \sqrt{R}}{k - R} + \frac{1}{\mu} \ln \left(\frac{\min(F_1 \mu \beta \sqrt{R}, k - R)}{\min\left(\max\left(1, \frac{\sqrt{R}}{D_1 \sqrt{\beta}}\right), k - R\right)} \right) \\ &\leq \beta + \frac{1}{\mu} \frac{F_1 \mu \beta \sqrt{R}}{k - R} + \frac{1}{\mu} \ln(F_1 D_1 \beta^{3/2})\end{aligned}$$

Proof of (5.12): Upper Bound on Time between Consecutive Downward Visits.

To bound the expectation $\mathbb{E} \left[\min(v_{i+1}^{(\text{down})}, X) - v_i^{(\text{down})} \middle| \mathcal{F}_{v_i^{(\text{down})}} \right]$, we split the interval into two parts, $\left[\min(v_i^{(\text{up})}, X) - v_i^{(\text{down})} \right]$ and $\left[v_{i+1}^{(\text{down})} - v_i^{(\text{down})} \right]$.

To bound the expectation of the first quantity, it suffices to note that, if we couple the system to an $M/M/\infty$, the coupled number of jobs $\tilde{N}(t)$ will reach $R + M_L$ only after the original system. Using Claim A.8 to bound this passage time, we thus know that

$$\begin{aligned}\mathbb{E} \left[\min(v_i^{(\text{up})}, X) - v_i^{(\text{down})} \middle| \mathcal{F}_{v_i^{(\text{down})}} \right] &\leq \mathbb{E} \left[\min\left(T_{(R+M_L-1) \rightarrow (R+M_L)}^{M/M/\infty} + v_i^{(\text{down})}, X\right) - v_i^{(\text{down})} \middle| \mathcal{F}_{v_i^{(\text{down})}} \right] \\ &\leq \mathbb{E} \left[T_{(R+M_L-1) \rightarrow (R+M_L)}^{M/M/\infty} \right] \\ &\leq \frac{D_2}{\sqrt{R}}.\end{aligned}$$

To bound the expectation of the second quantity, we provide two bounds. First, we again make use of the “wait-busy” idea; as we argued in the proof of (5.11),

$$\mathbb{E} \left[v_{i+1}^{(\text{down})} - v_i^{(\text{up})} \middle| \mathcal{F}_{v_i^{(\text{up})}} \right] \leq \mathbb{E} \left[\min(v_{i+1}^{(\text{down})} - v_i^{(\text{up})}, \beta) \middle| \mathcal{F}_{v_i^{(\text{up})}} \right] + \frac{1}{\mu M_L}.$$

From here, we note, by coupling to an $M/M/1$ with arrival rate and departure rate both equal to $k\lambda$, we can bound $\mathbb{E} \left[\min(v_{i+1}^{(\text{down})} - v_i^{(\text{up})}, \beta) \middle| v_i^{(\text{up})} < X \right]$ by the expected minimum between β and the length of a single-job busy period in that system. Applying Claim A.7, we can complete the proof, finding that

$$\mathbb{E} \left[\min(v_{i+1}^{(\text{down})} - v_i^{(\text{up})}, \beta) \middle| \mathcal{F}_{v_i^{(\text{up})}} \right] \leq D_1 \frac{\sqrt{\beta}}{\sqrt{\mu R}} + \frac{6}{\mu R}.$$

For the second bound, we simply note that, during the draining phase, the number of busy servers $Z(t) \geq R + 1$. It follows from a simple coupling argument that

$$\mathbb{E} \left[\min \left(v_{i+1}^{(\text{down})} - v_i^{(\text{up})}, \beta \right) \middle| \mathcal{F}_{v_i^{(\text{up})}} \right] \leq \frac{1}{\mu}.$$

Combining the bounds pessimistically, we find that

$$\begin{aligned} \mathbb{E} \left[\min \left(v_{i+1}^{(\text{down})}, X \right) - v_i^{(\text{down})} \middle| \mathcal{F}_{v_i^{(\text{down})}} \right] &\leq \frac{D_2}{\mu\sqrt{R}} + \frac{D_3}{\mu R} + \frac{1}{\mu M_L} + \min \left(D_1 \frac{\sqrt{\beta}}{\sqrt{\mu R}}, \frac{1}{\mu} \right) \\ &\leq \frac{D_2}{\mu\sqrt{R}} + \frac{D_3}{\mu R} + \frac{2}{\mu M_L} \end{aligned}$$

Proof of (5.13): Upper Bound on Probability of Another Downward Visit.

To bound the probability of an additional downcrossing, we again make a coupling argument. In particular, we couple again to the system which only has R servers busy, which gives an upper bound on the number of jobs in the system $N(t)$. If, in our coupled system, we reach $\tilde{N}(t) = R + M_L$ before we reach $\tilde{N}(t) = R$, then another upcrossing *must* have previously occurred in the original system, and thus another downcrossing must also occur. But, of course, we know classically that the probability that this happens is just $\frac{1}{M_L}$; this is precisely what is asserted by (5.13). \square

5.4 The Lower Bounds: Review of Findings

In this chapter, we proved two lower bounds on the average waiting time in the M/M/k/Setup-Deterministic. The first lower bound, Theorem 5.1, was the first-ever explicit result for the average waiting time in this model. The second lower bound, Theorem 5.2, is a considerable strengthening of Theorem 5.1, and also was far easier to prove once we made use of the MIST method.

Chapter 6

The Upper Bound

In this chapter, we present our upper bound on the average waiting time in the M/M/k/Setup-Deterministic.

6.1 Why we need an upper bound.

From a provisioning standpoint, an upper bound tells us what system parameters **sufficient** to achieve a certain average waiting time. By combining this bound with our lower bound, we find out what is necessary and sufficient for good performance. Theoretically-speaking, having the two bounds allows us to fully characterize how the average waiting time in the M/M/k/Setup-Deterministic scales with its system parameters, modulo some constant multiplicative factors.

6.2 The Upper Bound

We now state and prove the upper bound.

Theorem 6.1 (Upper Bound on Average Queue Length). *For an M/M/k/Setup-Deterministic with an offered load $R \triangleq k\rho \geq 100$ and a setup time $\beta \geq 1000\frac{1}{\mu}$, the expected number of jobs in queue in steady state is upper-bounded as*

$$\mathbb{E}[Q(\infty)] \leq A_1\sqrt{\mu\beta R} + A_2\frac{R}{M} + \frac{A_3\beta^2\mu\sqrt{R} + I^{\text{busy}}(B_5\sqrt{\mu\beta R} + B_6\mu\beta\sqrt{R}, M) + A_4I^{\text{busy}}(M, M)}{\beta + T^{\text{busy}}(D_1\beta\mu\sqrt{R}, k - R)},$$

where $A_1, A_2, A_3, A_4, B_5, B_6,$ and D_1 are constants independent of system parameters, and

$$M \triangleq \min(C_1\sqrt{\mu\beta R}, k - R)$$

for some constant C_1 independent of system parameters.

We now describe the full proof of Theorem 6.1. As discussed in Chapter 4, it suffices to prove the three following lemmas.

Lemma 6.1 (Accumulation Period Upper Bound). *Suppose the system begins at time 0 with R jobs in service and no jobs in the queue (and thus no servers in setup), and define the accumulation time*

$$T_A \triangleq \min \{t \geq 0 : Z(t) = R + 1\}$$

to be the moment the $(R + 1)$ -th server turns on.

Then,

$$\mathbb{E} \left[\int_0^{T_A} [N(t) - R] dt \right] \leq B_1 \sqrt{\mu\beta R} \cdot \mathbb{E}[T_A] + B_2 \beta^2 \mu \sqrt{R},$$

where $B_1 = 3.6$ and $B_2 = 1.04$.

Lemma 6.2 (Draining Period Upper Bound). *Recall that accumulation time T_A is the first (and only) time the $(R + 1)$ -th server turns on during a renewal cycle, and that the next renewal point $X = \min \{t > T_A : Z(t) = R\}$ is simply the next time the $(R + 1)$ -th server turns off. Then,*

$$\begin{aligned} & \mathbb{E} \left[\int_{T_A}^X [N(t) - R]^+ dt \right] \\ & \leq \left(B_5 \sqrt{\mu\beta R} + B_6 \mu \beta \sqrt{R} \right) \cdot \beta + I^{\text{busy}} \left(B_5 \sqrt{\mu\beta R} + B_6 \mu \beta \sqrt{R}, M \right) + B_7 \frac{\beta R}{M} \\ & + \left[2M + 2 \frac{R}{M} \right] \mathbb{E}[X - T_A] + \frac{1}{1 - p_2} I^{\text{busy}}(M, M), \end{aligned}$$

where all of these quantities are defined in Chapter 3 and Section A.4.7.

Lemma 6.3 (Cycle Length Lower Bound). *Suppose the system begins at time 0 with R jobs in service and no jobs in the queue (and thus no servers in setup), and let*

$$X \triangleq \min \{t > 0 : Z(t^-) = R + 1, Z(t) = R\}$$

be the next time the $(R + 1)$ -th server turns off.

Then,

$$\mathbb{E}[X] \geq \beta + T^{\text{busy}} \left(D_1 \beta \mu \sqrt{R}, k - R \right),$$

where D_1 is a constant independent of system parameters.

After proving these lemmas, the result follows by a bit of algebra. First, note that, by sum-

ming the two integral bounds, one obtains

$$\begin{aligned}
& \mathbb{E} \left[\int_0^X [N(t) - R] dt \right] \\
& \leq B_1 \sqrt{\mu\beta R} \cdot \mathbb{E}[T_A] + B_2 \beta^2 \mu \sqrt{R} + \left(B_5 \sqrt{\mu\beta R} + B_6 \mu \beta \sqrt{R} \right) \cdot \beta \\
& \quad + I^{\text{busy}} \left(B_5 \sqrt{\mu\beta R} + B_6 \mu \beta \sqrt{R}, M \right) + B_7 \frac{\beta R}{M} + \left[2M + 2 \frac{R}{M} \right] \mathbb{E}[X - T_A] + \frac{1}{1-p_2} I^{\text{busy}}(M, M) \\
& \leq B_1 \sqrt{\mu\beta R} \cdot \mathbb{E}[T_A] + B_2 \beta^2 \mu \sqrt{R} + \left(B_5 \sqrt{\mu\beta R} \right) \mathbb{E}[T_A] + B_6 \mu \beta^2 \sqrt{R} \\
& \quad + I^{\text{busy}} \left(B_5 \sqrt{\mu\beta R} + B_6 \mu \beta \sqrt{R}, M \right) + B_7 \frac{R}{M} \mathbb{E}[T_A] + \left[2M + 2 \frac{R}{M} \right] \mathbb{E}[X - T_A] + \frac{1}{1-p_2} I^{\text{busy}}(M, M) \\
& \leq \max \left(2M + 2 \frac{R}{M}, (B_1 + B_5) \sqrt{\mu\beta R} + B_7 \frac{R}{M} \right) \mathbb{E}[X] + (B_2 + B_6) \beta^2 \mu \sqrt{R} \\
& \quad + I^{\text{busy}} \left(B_5 \sqrt{\mu\beta R} + B_6 \mu \beta \sqrt{R}, M \right) + \frac{1}{1-p_2} I^{\text{busy}}(M, M) \\
& = \left(A_1 \sqrt{\mu\beta R} + A_2 \frac{R}{M} \right) \mathbb{E}[X] + A_3 \beta^2 \mu \sqrt{R} + I^{\text{busy}} \left(B_5 \sqrt{\mu\beta R} + B_6 \mu \beta \sqrt{R}, M \right) + A_4 I^{\text{busy}}(M, M),
\end{aligned}$$

where we have taken the constant $A_1 \triangleq \max(B_1 + B_5, C_3)$, the constant $A_2 \triangleq B_2 + B_3$, the constant $A_3 \triangleq B_2 + B_6$, and the constant $A_4 = \frac{1}{1-p_2}$. Upon dividing the reward integral by the cycle length, we obtain that

$$\begin{aligned}
\mathbb{E}[Q(\infty)] &= \frac{\mathbb{E} \left[\int_0^X [N(t) - R] dt \right]}{\mathbb{E}[X]} \\
&\leq \frac{\left(A_1 \sqrt{\mu\beta R} + A_2 \frac{R}{M} \right) \mathbb{E}[X] + A_3 \beta^2 \mu \sqrt{R} + I^{\text{busy}} \left(B_5 \sqrt{\mu\beta R} + B_6 \mu \beta \sqrt{R}, M \right) + A_4 I^{\text{busy}}(M, M)}{\mathbb{E}[X]} \\
&= A_1 \sqrt{\mu\beta R} + A_2 \frac{R}{M} + \frac{A_3 \beta^2 \mu \sqrt{R} + I^{\text{busy}} \left(B_5 \sqrt{\mu\beta R} + B_6 \mu \beta \sqrt{R}, M \right) + A_4 I^{\text{busy}}(M, M)}{\mathbb{E}[X]} \\
&\leq A_1 \sqrt{\mu\beta R} + A_2 \frac{R}{M} + \frac{A_3 \beta^2 \mu \sqrt{R} + I^{\text{busy}} \left(B_5 \sqrt{\mu\beta R} + B_6 \mu \beta \sqrt{R}, M \right) + A_4 I^{\text{busy}}(M, M)}{\beta + T^{\text{busy}} \left(D_1 \beta \mu \sqrt{R}, k - R \right)},
\end{aligned}$$

which is the upper bound stated in Theorem 6.1.

6.2.1 Proof of Lemma 6.1, Upper Bound on Integral Over Accumulation Period

We prove this result via two applications of the Intervening Stopping Time Lemma, Lemma 4.1. To apply this decomposition lemma, there are two broad steps. First, we must specify a starting time (T_0), an ending time (P), a series of intervening stopping times (T_i), the process (Y_t), and

an counting variable (F). Second, we must prove that the three preconditions of the lemma hold, given these specifications.

First application of Lemma 4.1, at the epoch level.

Definition of (τ_j) . We define the sequence of stopping times $(\tau_j : j = 0, 1, \dots, R)$ as $\tau_j \triangleq \min \{t > 0 : N(t) \leq R - j\}$, i.e., τ_j is the first time there are only $R - j$ jobs within the system. Note that, by definition, $\tau_0 = 0$. We call the period $\left[\tau_j, \min(\tau_{j+1}, T_A) \right)$ the j -th epoch, and say epoch j occurs whenever $\tau_j < T_A$. We then let n_e denote the number of epochs which occur in a given renewal cycle.

Specification step. Since we are interested in bounding $\mathbb{E} \left[\int_0^{T_A} [N(t) - R] dt \right]$, we let our starting stopping time be $T_0 = 0$, our ending stopping time be $P = T_A$, our intervening stopping times be $T_j = \tau_j$, the process of interest $Y_t = N(t) - R$ and our counting variable be $F = n_e$. Let the quantity

$$p_{\text{rise}}^{(j)} \triangleq \Pr \left(\max_{t \in [\tau_j, \min(\tau_{j+1}, T_A)]} N(t) \geq R + C_3 \sqrt{\mu\beta R} \mid n_e \geq j \right) \quad (6.1)$$

be the probability that the total number of jobs $N(t)$ exceeds $R + C_3 \sqrt{\mu\beta R}$ during epoch j .

Required claims. From here, we can apply Lemma 4.1 after showing the following claims:

Claim 6.1 (Upper Bound on the Probability of Another Epoch). *Recall that the total number of epochs $n_e \triangleq \max \{j \in \mathbb{Z}_+ : \tau_j < T_A\}$. Then, taking $C_4 = 0.98$, we have $\Pr(n_e \geq j + 1 \mid n_e \geq j) \leq 1 - C_4 p_{\text{rise}}^{(j)}$.*

Claim 6.2 (Upper Bound on the Integral Over an Epoch). *Let $\tau_j \triangleq \min \{t \geq 0 : N(t) \leq R - j\}$, $T_A \triangleq \min \{t \geq 0 : Z(t) = R + 1\}$, and let $n_e \triangleq \max \{i \in \mathbb{Z}_+ : \tau_i < T_A\}$. Then,*

$$\mathbb{E} \left[\int_{\tau_j}^{\min(\tau_{j+1}, T_A)} [N(t) - R] dt \mid n_e \geq j \right] \leq B_1 \sqrt{\mu\beta R} \cdot \mathbb{E} [\min(\tau_{j+1}, T_A) - \tau_j \mid n_e \geq j] + C_2 \beta^2 \mu_j p_{\text{rise}}^{(j)},$$

where $B_1 = 3.6$ and $C_2 = \frac{1}{2 \cdot 0.98} > 0.511$.

Proof of Lemma 6.1 assuming Claims 6.1 and 6.2.

Before going further, we show how to complete the proof of Lemma 6.1, assuming the two prior claims. Applying Lemma 4.1, we find that

$$\begin{aligned} \mathbb{E} \left[\int_0^{T_A} [N(t) - R] dt \right] &\leq B_1 \sqrt{\mu\beta R} \cdot \mathbb{E}[T_A] + C_2 \beta^2 \mu \sum_{j=1}^R j p_{\text{rise}}^{(j)} \prod_{i=1}^{j-1} \left(1 - C_4 p_{\text{rise}}^{(j)}\right) \\ &\leq B_1 \sqrt{\mu\beta R} \cdot \mathbb{E}[T_A] + \frac{C_2}{C_4} \beta^2 \mu \left[\sum_{j=1}^R j C_4 p_{\text{rise}}^{(j)} \prod_{i=1}^{j-1} \left(1 - C_4 p_{\text{rise}}^{(j)}\right) \right] \\ &= B_1 \sqrt{\mu\beta R} \cdot \mathbb{E}[T_A] + \frac{C_2}{C_4} \beta^2 \mu \left[\sum_{j=1}^R \prod_{i=1}^j \left(1 - C_4 p_{\text{rise}}^{(j)}\right) \right], \end{aligned}$$

where we have used the ‘‘expectation as a sum of tails’’ trick. We now apply the following claim:

Claim 6.3 (Bound on the Probability of an Up-crossing $p_{\text{rise}}^{(j)}$). *Let $p_{\text{rise}}^{(j)}$ be the probability that the total number of jobs $N(t)$ exceeds $R + C_3 \sqrt{\mu\beta R}$ during epoch j defined in (6.1). Then, for any epoch $j \geq A_5 \sqrt{R}$, we have $p_{\text{rise}}^{(j)} \geq 0.99 \frac{A_5}{\sqrt{R}}$.*

Continuation: Proof of Lemma 6.1 assuming Claims 6.1, 6.2, and 6.3.

We defer the proof of Claim 6.3 to Section A.4.3. Applying the claim’s result, we find that

$$\sum_{j=1}^R \prod_{i=1}^j \left(1 - C_4 p_{\text{rise}}^{(j)}\right) \leq \sum_{j=1}^R \left(1 - \frac{0.99 C_4 A_5}{\sqrt{R}}\right)^{\lceil j - A_5 \sqrt{R} \rceil^+} \leq \sum_{j=1}^{\infty} \left(1 - \frac{0.99 C_4 A_5}{\sqrt{R}}\right)^{\lceil j - A_5 \sqrt{R} \rceil^+}. \quad (6.2)$$

Bounding this as a Geometric sum, we obtain

$$(6.2) = A_5 \sqrt{R} + \sum_{j=0}^{\infty} \left(1 - \frac{0.99 C_4 A_5}{\sqrt{R}}\right)^j = A_5 \sqrt{R} + \frac{1}{0.99 C_4 A_5} \sqrt{R}.$$

Returning to our original inequality, we obtain that

$$\mathbb{E} \left[\int_0^{T_A} [N(t) - R] dt \right] \leq B_1 \sqrt{\mu\beta R} \cdot \mathbb{E}[T_A] + \frac{C_2}{C_4} \left(A_5 + \frac{1}{0.99 C_4 A_5} \right) \beta^2 \mu \sqrt{R}.$$

Noting that $A_5 = 1$ and taking $B_2 \triangleq 1.04 > \frac{C_2}{C_4} \left(A_5 + \frac{1}{0.99 C_4 A_5} \right)$, we finish the proof of Lemma 6.1. \square

Proof of Claim 6.1, Upper Bound on Probability of Another Epoch.

Rewriting the claim. And so, assuming the preconditions of Lemma 4.1 (Claims 6.1 and 6.2) as well as the helper claim 6.3, we have proven Lemma 6.1. We thus begin proving Claim 6.1. We begin by rewriting the probability of another epoch occurring as

$$\Pr(n_e \geq j + 1 | n_e \geq j) = 1 - \Pr(n_e = j | n_e \geq j) = 1 - \Pr(T_A < \tau_{j+1} | n_e \geq j).$$

It thus suffices to show a bound on the probability that the accumulation phase ends in epoch j :

$$\Pr(T_A < \tau_{j+1} | n_e \geq j) \geq C_4 p_{\text{rise}}^{(j)}. \quad (6.3)$$

Lower bound based on up-crossing and down-crossing times. To show (6.3), we analyze a particular sequence of events which results in the accumulation phase ending in the current epoch, i.e. $T_A < \tau_{j+1}$. Specifically, we define the up-crossing time $u = \min \{t > \tau_j : N(t) \geq R + C_3 \sqrt{\mu\beta R}\}$ and the down-crossing time $d = \min \{t > u : N(t) \leq R\}$. We consider the event where (1) the up-crossing occurs during the accumulation phase ($u < T_A$) and (2) the accumulation phase ends before the next down-crossing occurs ($d > T_A$). Symbolically, we have (at the end, recalling that $p_{\text{rise}}^{(j)}$ is the probability of an up-crossing occurs)

$$(6.3) \geq \Pr(u < T_A < d | n_e \geq j) = \Pr(d > T_A | u < T_A) \Pr(u < T_A) = \Pr(d > T_A | u < T_A) p_{\text{rise}}^{(j)}.$$

Development of conditional probability. To bound the conditional probability $\Pr(d > T_A | u < T_A)$, we condition on the filtration at time u , then make a coupling argument. To begin, note that, if the number of jobs $N(t)$ does not fall to R before the $(R+1)$ -th server finishes setting up, then the accumulation time T_A occurs exactly when the $(R+1)$ -th server finishes, i.e. the accumulation time $T_A = u + Y_{R+1}(u)$. Furthermore, the number of busy servers $Z(t) \leq R$ at any time during the accumulation phase $t < T_A$. Applying a basic coupling argument (Claim A.1), we have a lower bound on $N(t)$ in the coupled process

$$\tilde{N}(t) \triangleq N(u) + \Pi_A((u, t]) - \mathcal{D}[R]((u, t]),$$

for any time $t \in [u, T_A]$. Let $\tilde{d} \triangleq \min \{t > u : \tilde{N}(t) \geq R\}$ be the analogous down-crossing time in the coupled system. Since the coupled $\tilde{N}(t)$ is a lower bound, the coupled down-crossing time $\tilde{d} \leq d$. Thus,

$$\Pr(d > T_A | \mathcal{F}_u, u < T_A) = \Pr(d > u + Y_{R+1}(u) | \mathcal{F}_u, u < T_A) \geq \Pr(\tilde{d} > u + Y_{R+1}(u) | \mathcal{F}_u, u < T_A). \quad (6.4)$$

Analyzing the coupled probability. Continuing, the probability that $\{\tilde{d} \geq \ell\}$ is decreasing in ℓ . Thus,

$$(6.4) \geq \Pr(\tilde{d} > u + \beta | \mathcal{F}_u, u < T_A) = \Pr(\tilde{d} - u > \beta) \geq 1 - 2\Phi\left(-\frac{C_3}{\sqrt{2}}\right) - \frac{2}{3\sqrt{\mu\beta R}} \geq 0.98,$$

where in the final inequalities we have applied both the down-crossing probability bound of Claim A.3 and our assumptions. Taking $C_4 \triangleq 0.98$, we have the inter-epoch probability bound of Claim 6.1. \square

Proof of Claim 6.2, Upper Bound on the Integral Over an Epoch.

We now prove Claim 6.2, the upper bound on the time integral over an epoch. We do this via another application of Lemma 4.1 —first specifying the intervening times, then completing the proof, then proving that the preconditions hold.

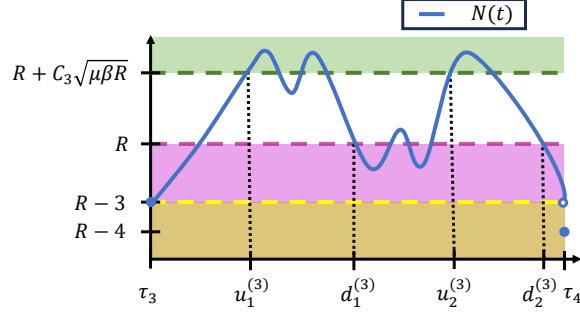


Figure 6.1: A depiction of the up-crossings and down-crossings defined in Section 6.2.1. In this example, we see that the number of up-crossings in epoch 3 is $n_e^{(3)} = 2$ and that, in this case, epoch 3 ends when epoch 4 begins (i.e. at time τ_4).

Definition of up-crossings and down-crossings. Let the 0-th down-crossing time in epoch j occur at time τ_j , i.e. let $d_0^{(j)} \triangleq \tau_j$. Next, define the first up-crossing in epoch j as the first time during epoch j that the total number of jobs $N(t)$ exceeds $R + C_3\sqrt{\mu\beta R}$, i.e.

$$u_1^{(j)} \triangleq \min \left\{ t > \tau_j : N(t) \geq R + C_3\sqrt{\mu\beta R} \right\}.$$

From here, define i -th down-crossing in epoch j and the $i + 1$ -th up-crossing in epoch j as

$$d_i^{(j)} \triangleq \min \left\{ t \geq u_i^{(j)} : N(t) \leq R \right\}$$

and

$$u_{i+1}^{(j)} \triangleq \min \left\{ t \geq d_i^{(j)} : N(t) \geq R + C_3\sqrt{\mu\beta R} \right\},$$

respectively; we visualize these definitions in Figure 6.1. We say the i -th up-crossing occurs if $u_i^{(j)} < \min(T_A, \tau_{j+1})$ and let $n_u \triangleq \max \left\{ i \geq 0 : u_i^{(j)} < \min(T_A, \tau_{j+1}) \right\}$ be the random number of up-crossings which occur in epoch j . We call the interval $\left[d_i^{(j)}, \min(u_i^{(j)}, \min(T_A, \tau_{j+1})) \right)$ the i -th rise, and the interval $\left[u_i^{(j)}, \min(d_i^{(j)}, \min(T_A, \tau_{j+1})) \right)$ the i -th fall. Note that, if the i -th up-crossing occurs, then, by definition, $d_i < \tau_{j+1}$; this means that the i -th fall can always be written as $\left[u_i, \min(T_A, d_i) \right)$. For readability, we fix our epoch of interest and freely omit the superscript j on our up-crossings and down-crossings.

Specification step. With up-crossings and down-crossings defined, we are now ready to specify our application of the IST Lemma, Lemma 4.1. We define our starting time as $T_0 = \tau_j = d_0$, our ending time as $P = \min(T_A, \tau_{j+1})$, our intervening sequence as $(u_i)_{i=1}^\infty$, and our counting variable as $F = n_u$.

Required Claims. From here, in order to apply Lemma 4.1, we must show the following three claims:

Claim 6.4 (Upper Bound on Integral Until First Up-crossing). *Taking $B_1 = 3.6$, the integral until u_1 is*

$$\mathbb{E} \left[\int_{d_0}^{\min(u_1, \min(T_A, \tau_{j+1}))} [N(t) - R] dt \middle| n_u \geq j \right] \leq B_1 \sqrt{\mu\beta R} \cdot \mathbb{E} [\min(u_1, \min(T_A, \tau_{j+1})) - \tau_j | n_u \geq j].$$

Claim 6.5 (Upper Bound on Integral Between Up-crossings). *The integral between up-crossings u_i is*

$$\begin{aligned} \mathbb{E} \left[\int_{u_i}^{\min(u_{i+1}, \min(T_A, \tau_{j+1}))} [N(t) - R] dt \middle| n_u \geq i \right] &\leq B_1 \sqrt{\mu\beta R} \cdot \mathbb{E} [\min(u_{i+1}, \min(T_A, \tau_{j+1})) - u_i | n_u \geq i] \\ &\quad + \frac{1}{2} \beta^2 \mu j. \end{aligned}$$

Claim 6.6 (Upper Bound on Probability of Another Up-crossing). *Recall that $p_{\text{rise}}^{(j)}$ is the probability that the number of jobs $(N(t) \geq C_3 \sqrt{\mu\beta R})$ at some point during epoch j , given that epoch j occurs. Then, $\Pr(n_u > 0) = p_{\text{rise}}^{(j)}$, and, for all counts $i \geq 1$ and $p_2 = 0.98$, we have $\Pr(n_u \geq i + 1 | n_u \geq i) \leq 0.02 = 1 - p_2$.*

Proof of Claim 6.2, assuming Claims 6.4, 6.5, and 6.6.

Once again, before we move on to proving these claims, we show that they indeed suffice to prove Claim 6.2. By Lemma 4.1, taking $C_2 \triangleq \frac{0.5}{p_2}$,

$$\begin{aligned} \mathbb{E} \left[\int_{\tau_j}^{\min(T_A, \tau_{j+1})} [N(t) - R] dt \middle| n_e \geq j \right] &\leq B_1 \sqrt{\mu\beta R} \cdot \mathbb{E} [\min(T_A, \tau_{j+1}) - \tau_j | n_e \geq j] \\ &\quad + p_{\text{rise}}^{(j)} 0.5 \beta^2 \mu j \sum_{i=1}^{\infty} (1 - p_2)^{i-1} \\ &= B_1 \sqrt{\mu\beta R} \cdot \mathbb{E} [\min(T_A, \tau_{j+1}) - \tau_j | n_e \geq j] + p_{\text{rise}}^{(j)} \frac{0.5}{p_2} \beta^2 \mu j. \quad \square \end{aligned}$$

Proofs of Claims 6.4, 6.5, and 6.6.

All that remains to be proven are our three aforementioned claims.

Proof of Claim 6.4: Upper Bound on Integral until First Up-crossing. Proving this claim is quite simple. In fact, we now prove a far more general claim, that the integral from a down-crossing to the next up-crossing

$$\int_{d_i}^{\min(u_i, \min(T_A, \tau_{j+1}))} [N(t) - R] dt \leq C_3 \sqrt{\mu\beta R} \cdot [\min(T_A, \tau_{j+1}) - d_i]. \quad (6.5)$$

To see this, note that, at any point between a down-crossing and up-crossing, the total number of jobs $N(t)$ must be strictly less than $R + C_3\sqrt{\mu\beta R}$. Apply this to the 0-th down-crossing and we have the claim. \square

Proof of Claim 6.5: Upper Bound on Integral Between Up-crossings. This proof is a bit more involved. We separate the interval $\left[u_i, \min(u_{i+1}, \min(T_A, \tau_{j+1})) \right)$ into the i -th fall and the i -th rise, as discussed previously. For the rising portion, we can simply apply the simple bound from (6.5). For the falling portion, we apply the integral coupling claim, Claim A.2. In particular, note that $Z(t) \geq R - j$ until time τ_{j+1} and that the interval $[u_i, \min(d_i, T_A))$ is equivalent to the interval $[u_i, \min(d_i, u_i + Y_{R+1}(u_i))]$. Applying Claim A.2,

$$\begin{aligned} \mathbb{E} \left[\int_{u_i}^{\min(d_i, T_A)} [N(t) - R] dt \middle| S(u_i) \right] &\leq \frac{1}{2} \beta^2 \mu j + [N(u_i) - R] \cdot Y_{R+1}(u_i) \\ &= \frac{1}{2} \beta^2 \mu j + C_3 \sqrt{\mu\beta R} \cdot Y_{R+1}(u_i). \end{aligned}$$

By our analysis in the proof of Claim 6.1, we note that the remaining setup time $Y_{R+1}(u_i) \leq \min(d_i, T_A)$ with probability at least $p_2 \leq \min_{\mathcal{F}_{u_i}} \Pr(d_i < T_A | \mathcal{F}_{u_i}, n_u \geq i)$. By Markov's inequality,

$$Y_{R+1}(u_i) \leq \frac{1}{p_2} \mathbb{E}[\min(d_i, T_A) - u_i | S(u_i)].$$

Combining our bounds on the rises and falls and taking $B_1 = \frac{C_3}{p_2}$, we have Claim 6.5. \square

Proof of Claim 6.6: Upper Bound on Probability of Another Up-Crossing. We now proceed to our final claim, concerning the up-crossing probabilities $\Pr(n_u > 0)$ and $\Pr(n_u \geq i + 1 | n_u \geq i)$. To begin, we first note that, since the first up-crossing occurs precisely at the moment that $N(t)$ exceeds $R + C_3\sqrt{\mu\beta R}$ during epoch j , one has $p_{\text{rise}}^{(j)} = \Pr(u_1 < \min(T_A, \tau_{j+1})) = \Pr(n_u > 0)$.

To prove the second part of the claim, first observe that an $(i + 1)$ -th up-crossing can only occur if an i -th down-crossing occurs, i.e. $\Pr(n_u \geq i + 1 | n_u \geq i) \leq \Pr(d_i < T_A | n_u \geq i)$.

To bound this conditional probability, we can apply a previous result. Recall the proof of the inter-epoch probability bound (Claim 6.1). In (6.3), we have already argued a bound on the conditional probability that the *first* down-crossing occurs, *in a state-independent manner*. The bound derived there thus *also applies here*:

$$\Pr(d_i < T_A | n_u \geq i) = 1 - \Pr(d_i > T_A | n_u \geq i) \leq 1 - C_4.$$

Taking $p_2 \triangleq C_4 = 0.98$, we have bounded the probability of another up-crossing (Claim 6.5). \square

6.2.2 Proof of Lemma 6.2, Upper Bound on Integral Over Draining Period.

To prove this lemma, we again make use of Lemma 4.1. We proceed through the usual two-step process, first defining the stopping time sequence we will analyze over, then proving the preconditions of the lemma.

Definition of the upward visits $v_i^{(\text{up})}$ and downward visits $v_i^{(\text{down})}$. Recall that the draining phase begins at time T_A . Let $M_B \triangleq \min(k - R, \sqrt{R})$ be a specially-set analysis threshold. Let the stopping time $v_1^{(\text{down})} \triangleq \min\{t \geq T_A : N(t) < R + M_B\}$ be the first time the number of jobs $N(t)$ drops below $R + M_B$, and recursively define

$$v_i^{(\text{up})} \triangleq \min\{t \geq v_i^{(\text{down})} : N(t) \geq R + M_B\}$$

and

$$v_{i+1}^{(\text{down})} \triangleq \min\{t \geq v_i^{(\text{up})} : N(t) < R + M_B\}.$$

Application of Lemma 4.1.

Applying Lemma 4.1, we take our initial stopping time to be the accumulation time T_A , our final stopping time to be the end of the renewal cycle X , our intervening stopping times to be the downward visits $v_i^{(\text{down})}$, and our counting index to be n_b .

Required claims. To apply Lemma 4.1, we need to show is the usual three claims: a bound on the initial integral, a bound on the continuing integral, and a bound on the probability.

Claim 6.7 (Upper Bound on Integral Until First Downward Visit). *Let $g(x, y, z) \triangleq x \frac{1}{2\mu z} + y \left[\frac{R}{\mu z^2} + \frac{3}{2\mu z} \right]$. Then, one can bound the integral immediately after time T_A with*

$$\begin{aligned} \mathbb{E} \left[\int_{T_A}^{v_1^{(\text{down})}} [N(t) - R] dt \right] &\leq \left[\beta + \frac{1}{\mu} \right] \left[3\mu\beta\sqrt{R} + \max\left(\sqrt{R}, \frac{\rho}{1-\rho}\right) \right] + \frac{2}{\mu} \ln \left(3 \frac{\mu\beta}{\sqrt{R}} \max\left(\sqrt{R}, \frac{\rho}{1-\rho}\right) \right) \\ &\quad + g\left((9(\mu\beta)^2 R, 3\mu\beta\sqrt{R}, k(1-\rho))\right) + \beta \frac{\rho}{1-\rho}. \end{aligned}$$

Claim 6.8 (Upper Bound on Integral Between Downward Visits). *One can bound the integral between consecutive downward visits by*

$$\mathbb{E} \left[\int_{v_i^{(\text{down})}}^{v_{i+1}^{(\text{down})}} [N(t) - R] dt \middle| \mathcal{F}_{v_i^{(\text{down})}} \right] \leq \frac{1}{\mu M_B} \left[\max\left(\sqrt{R}, \frac{\rho}{1-\rho}\right) + 14 + \mu\beta\sqrt{R} + 2b_1\sqrt{\mu\beta R} + b_2\sqrt{R} \right]$$

Claim 6.9 (Upper Bound on Probability of Another Downward Visit). *One can bound the probability of another downward visit occurring by*

$$\Pr(n_b \geq i + 1 | n_b \geq i) \leq \frac{1}{M_B}. \quad (6.6)$$

Proof of Lemma 6.2 assuming Claims 6.7, 6.8, and 6.9.

Simplifying the first bound further,

$$\begin{aligned}
\mathbb{E} \left[\int_{T_A}^{v_1^{(\text{down})}} [N(t) - R] dt \right] &\leq \left[\beta + \frac{3}{\mu} \right] \left[2.9\mu\beta\sqrt{R} + \max \left(\sqrt{R}, \frac{\rho}{1-\rho} \right) \right] \\
&\quad + g \left(9(\mu\beta)^2 R, 3\mu\beta\sqrt{R}, k(1-\rho) \right) + \beta \frac{\rho}{1-\rho} \\
&\leq 1.03\beta \left[2.91\mu\beta\sqrt{R} + \frac{\rho}{1-\rho} \right] + g \left(9(\mu\beta)^2 R, 3\mu\beta\sqrt{R}, k(1-\rho) \right) + \beta \frac{\rho}{1-\rho} \\
&\leq 3\mu\beta^2\sqrt{R} + 2.03\beta \frac{\rho}{1-\rho} + g \left(9(\mu\beta)^2 R, 3\mu\beta\sqrt{R}, k(1-\rho) \right),
\end{aligned}$$

where we have bounded $\ln(x)/x < 0.1$ for values $x > 100$, noted that $\max(x, y) \leq x + y$, and done some upwards rounding. Simplifying the ‘‘continuing integral’’ term, we have

$$\begin{aligned}
\mathbb{E} \left[\int_{v_1^{(\text{down})}}^X [N(t) - R] dt \right] &\leq \frac{1}{\text{Claim 6.9}} \cdot [\text{Claim 6.8}] \\
&= \frac{1}{\mu} \left[\max \left(\sqrt{R}, \frac{\rho}{1-\rho} \right) + 14 + \mu\beta\sqrt{R} + 2b_1\sqrt{\mu\beta R} + b_2\sqrt{R} \right] \\
&\leq \frac{1}{\mu} \left[\mu\beta\sqrt{R} \left(1 + \frac{2b_1}{\sqrt{\mu\beta}} + \frac{b_2}{\beta} + \frac{14}{\mu\beta\sqrt{R}} \right) + \frac{\rho}{1-\rho} \right] \\
&\leq \frac{1}{\mu} \left[1.6\mu\beta\sqrt{R} + \frac{\rho}{1-\rho} \right].
\end{aligned}$$

Combining these two pieces, we obtain, as desired,

$$\mathbb{E} \left[\int_{T_A}^X [N(t) - R] dt \right] \leq 3.01\mu\beta^2\sqrt{R} + 2.04\beta \frac{\rho}{1-\rho} + g \left(9(\mu\beta)^2 R, 3\mu\beta\sqrt{R}, k(1-\rho) \right). \square$$

Precursor: The ‘‘Wait-Busy’’ idea.

As such, to complete our proof it suffices to show Claims 6.7, 6.8, and 6.9. To prove these claims, we make heavy use of the following idea.

Claim 6.10 (Wait Busy Claim). *Let τ be some stopping time, let the number of jobs $N(\tau) = R + h$, and define $\text{ns}(h) \triangleq \min \{h, k(1-\rho)\}$. Let the down-crossing $d_{\text{gen}} \triangleq \min \{t > 0 : N(\tau + t) = R + h - 1\}$. If $Z(\tau) \geq R$, then*

$$\mathbb{E} \left[\int_{\tau}^{\tau+d_{\text{gen}}} [N(t) - (R + h - 1)] dt \middle| \mathcal{F}_{\tau} \right] \leq Y_{R+\text{ns}(h)}(\tau) + g \left(1 + 2\mu R \mathbb{E} \left[\min \left(Y_{R+\text{ns}(h)}(\tau), d_{\text{gen}} \right) \right], 1, \mu \text{ns}(h) \right), \tag{6.7}$$

where the function $g(x, y, z) \triangleq x \frac{1}{2\mu z} + y \left[\frac{R}{\mu z^2} + \frac{3}{2\mu z} \right]$.

Furthermore,

$$\begin{aligned} \mathbb{E} \left[\int_{T_A}^{v_1^{(\text{down})}} [N(t) - R] dt \middle| \mathcal{F}_{T_A} \right] &\leq \left[\beta + \frac{1}{\mu} \right] \left[\mathbb{E} [N(T_A) - R] + \frac{R}{M_B} \right] + \frac{2}{\mu} \ln \left(\frac{\mathbb{E} [N(T_A) - R]}{M_B} \right) \\ &+ g \left(\mathbb{E} [[N(T_A) - R]^2] + 2\mathbb{E} [N(T_A) - R], \mathbb{E} [N(T_A) - R], k(1 - \rho) \right). \end{aligned} \quad (6.8)$$

Intuition. We defer the proof of Claim 6.10 until Section 6.2.2. For now, we give some brief intuition for how the bound is derived and how we use it in our proof. Essentially, we can consider performing the following procedure at time τ : First, watch the system for β time. If the number of jobs ever dips below $R + h$ during this watching period, we can end our integral immediately. If the number of jobs $N(t)$ never dips below $R + h$ during this watching period, then we know for sure that we have at least $\min(R + h, k)$ servers on at time $\tau + Y_{\min(R+h, k)}(\tau)$, since we have continually had at least $R + h$ servers either busy or setting up during that period. Moreover, since we only turn off servers when there isn't work for them to do, those servers will *stay* on until the number of jobs $N(t)$ dips below $R + h$; in other words, they will stay on until time d_{gen} . The proof of the claim follows along essentially the same lines, formalizing things and performing computations using coupling and martingales.

Proof of Claim 6.7: Bound on Integral Until First Visit.

We return to proving our claims. The proof of this claim is simple; it is essentially rolled into the proof of Claim 6.10. From here, it suffices to apply the following claim, substituting in and simplifying constants:

Claim 6.11 (Upper Bound on $\mathbb{E} [N(T_A)]$). *Recall that $T_A \triangleq \min \{t > 0 : Z(t) = R + 1\}$. Then,*

$$\mathbb{E} [N(T_A) - R] \leq F_1 \mu \beta \sqrt{R} \left(1 + \frac{F_2}{\sqrt{\mu \beta}} \right) \leq 2.9 \mu \beta \sqrt{R}$$

and

$$\mathbb{E} [(N(T_A) - R)^2] \leq F_1^2 (\mu \beta)^2 R \left(1 + \frac{F_2}{\sqrt{\mu \beta}} \right)^2 + 2\mu \beta R \leq 8.4 (\mu \beta)^2 R + 2\mu \beta R$$

where $F_1 = 2.12$ and $F_2 = 3.645$.

Proof of Claim 6.8: Bound on Integral Between Visits.

To prove Claim 6.8, we break the integral into two parts: from the down-crossing $v_i^{(\text{down})}$ to the up-crossing $v_i^{(\text{up})}$, and vice-versa.

First part: from $v_i^{(\text{down})}$ to $v_i^{(\text{up})}$. To bound the integral from the down-crossing to the next up-crossing, we first make the simple observation that

$$\int_{v_i^{(\text{down})}}^{\min(v_i^{(\text{up})}, X)} [N(t) - R] dt \leq \left[\min(v_i^{(\text{up})}, X) - v_i^{(\text{down})} \right] \cdot M_B, \quad (6.9)$$

since $v_i^{(\text{up})}$ is the next time $N(t) \geq R + M_B$. To bound $\mathbb{E} \left[\min \left(v_i^{(\text{up})}, X \right) - v_i^{(\text{down})} \middle| \mathcal{F}_{v_i^{(\text{down})}} \right]$, we couple the system to an infinite-server $M/M/\infty$ queue at time $v_i^{(\text{down})}$, and note that the coupled up-crossing time

$$\tilde{T}_{(R+M_B-1) \rightarrow (R+M_B)} + v_i^{(\text{down})} \geq v_i^{(\text{up})} \geq \min \left(v_i^{(\text{up})}, X \right).$$

Since $M_B \leq \sqrt{R}$, from standard results on the $M/M/\infty$ (reproduced in Section A.3.4), we have that

$$M_B \mathbb{E} \left[\min \left(v_i^{(\text{up})}, X \right) - v_i^{(\text{down})} \middle| \mathcal{F}_{v_i^{(\text{down})}} \right] \leq M_B \frac{b_2}{\mu \sqrt{R}} = \frac{1}{\mu M_B} b_2 \frac{M_B^2}{\sqrt{R}} \leq \frac{1}{\mu M_B} b_2 \sqrt{R}.$$

Second part: from $v_i^{(\text{up})}$ to $v_{i+1}^{(\text{down})}$. From $v_i^{(\text{up})}$ onwards, we use the “wait-busy” bound. Applying the “Wait Busy” Claim (Claim 6.10) with $h = M_B$, we obtain

$$\begin{aligned} \mathbb{E} \left[\int_{v_i^{(\text{up})}}^{v_{i+1}^{(\text{down})}} [N(t) - R] dt \middle| \mathcal{F}_{v_i^{(\text{up})}} \right] &\leq \beta + g \left(1 + 2\mu R \mathbb{E} \left[\min \left(\beta, v_{i+1}^{(\text{down})} - v_i^{(\text{up})} \middle| \mathcal{F}_{v_i^{(\text{up})}} \right) \right], 1, M_B \right) \\ &\quad + M_B \mathbb{E} \left[\min \left(\beta, v_{i+1}^{(\text{down})} - v_i^{(\text{up})} \middle| \mathcal{F}_{v_i^{(\text{up})}} \right) \right]. \end{aligned}$$

To bound the conditional expectation of $\mathbb{E} \left[\min \left(\beta, v_{i+1}^{(\text{down})} - v_i^{(\text{up})} \right) \middle| \mathcal{F}_{v_i^{(\text{up})}} \right]$, we make our usual coupling argument. Define a coupled system $M/M/1$ queue with departure rate μR , and let d_{gen}^{\sim} be the length of its busy period. It suffices to bound $\mathbb{E} \left[\min \left(d_{\text{gen}}^{\sim}, \beta \right) \right]$, for the coupled relative down-crossing time d_{gen}^{\sim} . From standard results on simple random walks (Claim A.7), we have $\mathbb{E} \left[\min \left(d_{\text{gen}}^{\sim}, \beta \right) \right] \leq \frac{b_1 \sqrt{\mu \beta}}{\mu \sqrt{R}} + \frac{6}{\mu R}$, giving

$$\begin{aligned} \mathbb{E} \left[\int_{v_i^{(\text{up})}}^{v_{i+1}^{(\text{down})}} [N(t) - R] dt \middle| \mathcal{F}_{v_i^{(\text{up})}} \right] &\leq \frac{1}{\mu M_B} \left[\frac{R}{M_B} + 2 \right] + \beta + \left[\frac{b_1 \sqrt{\mu \beta}}{\mu \sqrt{R}} + \frac{6}{\mu R} \right] \left(M_B + \frac{R}{M_B} \right) \\ &= \frac{1}{\mu M_B} \left[\frac{R}{M_B} + 2 + \mu \beta M_B + (M_B^2 + R) \left[\frac{b_1 \sqrt{\mu \beta}}{\sqrt{R}} + \frac{6}{R} \right] \right] \\ &\leq \frac{1}{\mu M_B} \left[\max \left(\sqrt{R}, \frac{\rho}{1 - \rho} \right) + 14 + \mu \beta \sqrt{R} + 2b_1 \sqrt{\mu \beta R} \right], \end{aligned}$$

where in the last line we used $M_B \leq \sqrt{R}$; combining these two parts gives Claim 6.8. \square

Proof of Claim 6.9, Upper Bound on the Probability of Another Visit.

To see (6.6), we first note that, if there is another upcrossing, then there must be another down-crossing. As such, it suffices to upper bound $\Pr \left(v_i^{(\text{up})} < X \middle| v_i^{(\text{down})} \right)$. To do this, we note that the number of busy servers $Z(t) \geq R$. From Claim A.1, it thus suffices to bound the corresponding probability in the coupled system with exactly R busy servers. But this is simply the probability that a simple random walk started at $W(0) = M_B - 1$ hits $W(t) = M_B$ before it hits $W(t) = 0$. Classically, this probability is $\frac{1}{M_B}$; this proves the claim, and thus Lemma 6.2. \square

6.2.3 Proof of Lemma 6.3: Lower Bound on the Cycle Length.

Preliminaries. The proof of this lemma is much simpler than the others. Before describing our strategy, we first state some preliminaries. Recall the definition of the start of the j -th epoch $\tau_j \triangleq \min \{t \geq 0 : N(t) \leq R - j\}$, that we call the period $\left[\tau_j, \min(\tau_{j+1}, T_A) \right)$ the j -th epoch, and we say epoch j occurs if $\tau_j < T_A$. Now, say epoch j is *long* if it lasts longer than a setup time β ; note that such an epoch must exist, since servers can only turn on during long epochs, and a server must turn on before the accumulation phase ends at time T_A . Let $L \triangleq \min \{j \in \{0, 1, 2, \dots, R\} : \min(\tau_{j+1}, T_A) - \tau_j > \beta\}$ be the index of the *first* long epoch. Note that, although the random time τ_L is *not* a stopping time (we do not know how long an epoch will last when the epoch starts), the first moment we can identify epoch L , the random time $\tau_L + \beta$, is a stopping time. Moreover, we know that $\tau_L + \beta < T_A$. From here, one sees that $\mathbb{E}[X] = \mathbb{E}[\tau_L + \beta] + \mathbb{E}[X - (\tau_L + \beta)] \geq \beta + \mathbb{E}[X - (\tau_L + \beta)]$. To complete the proof, it suffices to show

$$\mathbb{E}[X - (\tau_L + \beta)] \geq \frac{L_1 \mu \beta \sqrt{R}}{\mu k (1 - \rho)}. \quad (6.10)$$

Proof of (6.10): Lower Bound on the Remaining Cycle Length.

To show (6.10), we first show we can bound an analogous quantity in a coupled process, then appeal to standard results on the M/M/1 queue.

Defining the coupled process $\tilde{N}(t)$. To define the coupled process, note that the number of busy servers $Z(t) \leq k$. It follows from Claim A.1 that, for any time $t \geq \tau_L + \beta$, the coupled process

$$\tilde{N}(t) \triangleq N(\tau_L + \beta) + A((\tau_L + \beta, t]) - \mathcal{D}[k]((\tau_L + \beta, t])$$

satisfies $\tilde{N}(t) \leq N(t)$.

Using the coupled process to bound $\mathbb{E}[X - \tau_L + \beta]$. We now use this process to bound $\mathbb{E}[X - \tau_L + \beta]$. Recall that the end of the renewal cycle $X \triangleq \min \{t > 0 : Z(t^-) = R + 1, Z(t) = R\}$ occurs when the $(R + 1)$ -th server turns off. It is useful to view X in a different way: since the accumulation time T_A is the moment when the $(R + 1)$ -th server turns on, we also know that the end of the renewal cycle $X = \min \{t > T_A : N(t) \leq R\}$ is the first moment after time T_A that the number of jobs $N(t) \leq R$. Furthermore, since the time $\tau_L + \beta$ happens *before* any server could possibly turn on, the time $\tau_L + \beta < T_A$. Denoting the end of the coupled renewal as \tilde{X} as the first moment the coupled process $\tilde{N}(t) \leq R$, we have

$$\tilde{X} \triangleq \min \left\{ t > \tau_L + \beta : \tilde{N}(t) \leq R \right\} \leq \min \{ t > \tau_L + \beta : N(t) \leq R \} \leq \min \{ t > T_A : N(t) \leq R \} = X.$$

Bounding the end of coupled renewal $\mathbb{E}[X - \tau_L + \beta]$. To bound the quantity $\mathbb{E}[X - \tau_L + \beta]$, we condition on the filtration at time $\tau_L + \beta$ and use standard results on the M/M/1 busy period. Note that, since the departure rate of the coupled system is fixed at μk , the period $\left[\tilde{X} - \tau_L + \beta \right]$

is precisely the length of an M/M/1 busy period with 1) arrival rate $k\lambda$, 2) departure rate $k\mu$, and 3) started by $[N(\tau_L + \beta) - R]^+$ jobs. It follows that $\mathbb{E}[X - (\tau_L + \beta) | \mathcal{F}_{\tau_L + \beta}] \geq \frac{[N(\tau_L + \beta) - R]^+}{\mu k(1 - \rho)}$. Taking expectations, applying Jensen's and results from the lower bound ((5.6) and (5.3)), we obtain, proving (6.10), Lemma 6.3, and Theorem 6.1 simultaneously,

$$\mathbb{E}[X - (\tau_L + \beta)] \geq \frac{[\mathbb{E}[N(\tau_L + \beta) - R]^+]}{\mu k(1 - \rho)} \geq \frac{\mu\beta\mathbb{E}[L]}{\mu k(1 - \rho)} \geq \frac{L_1\mu\beta\sqrt{R}}{\mu k(1 - \rho)}. \square$$

6.3 The Upper Bound: Review of Findings

In this chapter, we proved an upper bound on the average waiting time in the M/M/k/Setup-Deterministic. We proved this bound via a number of applications of the MIST Lemma, Lemma 4.1. In fact, to bound the accumulation phase integral, we needed to use the MIST Lemma in a nested way: First, we used it to break the accumulation phase into epochs, and then we used it to break each epoch into "rises" and "falls," periods of time punctuated by up-crossings and down-crossings. Compared to the lower bound of Chapter 5, the upper bound proven here truly highlights the utility of the MIST method.

Chapter 7

The Approximation

In Chapter 7, we present our approximation for the average waiting time in the M/M/k/Setup-Deterministic. We begin by discussing why we need such an approximation, then state the approximation, then give a short justification for its form.

7.1 Why we need an approximation

Despite our success in analyzing the M/M/k/Setup-Deterministic, our upper and lower bounds *alone* are not suitable for practical use in predicting the value of the average waiting time $\mathbb{E}[T_Q]$.

There are two reasons for this. First, although we can prove that our bounds are both within a constant factor of the true waiting time (Theorem A.1), it's not *a priori* obvious whether the true value of $\mathbb{E}[T_Q]$ will get closer to one bound or the other as we vary the system parameters. Although the true value does not *seem* to ever get closer to a particular bound (and so we could conceivably just scale our lower bound to serve as a predictor), it would be better to have a more concrete theoretical justification for our prediction.

The second reason why our bounds are unsuitable for practical use is their complexity. Although both the upper and lower bounds are far more straightforward to compute than, for example, the average waiting time in the M/M/k/Setup-Exponential, both bounds incorporate a large number of terms and are thus somewhat difficult to reason about on the fly. As such, it would be better to have a predictor which incorporates only a few, easy-to-remember terms.

7.2 The approximation

To this end, we introduce the following approximation; the justification for the approximation follows. An empirical evaluation of this approximation can be found in Figure 1.7; it is extremely accurate.

Approximation 1 (Approximation to the average queue length.). *In the M/M/k/Setup-Deterministic,*

for offered loads $R \triangleq k\rho > 2$,

$$\mathbb{E}[Q(\infty)] \approx Q_{\text{apx}} \triangleq \frac{\frac{1}{2}\beta^2 C_{\text{apx}} \sqrt{R} + \frac{\beta C_{\text{apx}} \sqrt{R}}{\mu k(1-\rho)} \left[\frac{\beta C_{\text{apx}} \sqrt{R} + 1}{2} + \frac{1}{1-\rho} \right]}{\beta + \frac{\beta C_{\text{apx}} \sqrt{R}}{\mu k(1-\rho)}}, \quad (7.1)$$

where $C_{\text{apx}} \triangleq \sqrt{\frac{\pi}{2}}$.

7.3 Justification

We arrive at this bound via a straightforward combination of our results from Chapters 5 and 6, along with a few modifications. We follow our renewal-reward analysis, separately approximating the expected time integral over our renewal cycle and the expected length of that renewal cycle, the numerator and denominator of 7.1, respectively.

7.3.1 Justification of Numerator

We first approximate the numerator of our expression, the expected time integral over our chosen renewal cycle. We begin by recalling the lower bound on the time integral, Lemma 5.1, which states

$$\mathbb{E} \left[\int_0^X Q(t) dt \right] \geq L_1 \beta^2 \sqrt{R} + I^{\text{busy}} \left(\left[L_1 \beta \sqrt{R} - (k - R) \right]^+, k - R \right),$$

where

$$I^{\text{busy}}(x, z) \triangleq \frac{x}{\mu z} \left[\frac{x+1}{2} + \frac{1}{1 - \frac{k\lambda}{k\lambda + \mu z}} \right]$$

represents the time integral of the queue length a certain M/M/1 queue over a busy period started by x jobs.

To obtain the appropriate constant C_{apx} , we next note that, although our theorem states L_1 as an absolute constant, as the setup time β and the offered load R grow, the best possible constant will become $C_{\text{apx}} = \sqrt{\frac{\pi}{2}}$. Under the hood, this convergence stems from the fact that

$$\sum_{j=1}^R \prod_{i=1}^j \left(1 - \frac{j}{R} \right) \approx \int_0^\infty e^{-\frac{j^2}{2R}} dj = \frac{1}{2} \sqrt{2\pi R};$$

see the proof of Lemma 5.1 for more details.

To complete the bound, it suffices to remove the subtraction of $(k - R)$ in the busy period term, which we anticipate is an artifact of our analysis. Removing it, we obtain the desired approximation

$$\mathbb{E} \left[\int_0^X Q(t) dt \right] \approx \frac{1}{2} \beta^2 C_{\text{apx}} \sqrt{R} + \frac{\beta C_{\text{apx}} \sqrt{R}}{\mu k(1-\rho)} \left[\frac{\beta C_{\text{apx}} \sqrt{R} + 1}{2} + \frac{1}{1-\rho} \right]. \quad (7.2)$$

7.3.2 Justification of Denominator

We next approximate the denominator of our expression, the expected length of our chosen renewal cycle. To do so, we again make use of the lower bound on the expected cycle length $\mathbb{E}[X]$ from Lemma 6.3, which states

$$\mathbb{E}[X] \geq \beta + \frac{L_1 \beta \sqrt{R}}{\mu k (1 - \rho)}.$$

By making the same convergence argument for L_1 , i.e. that $L_1 \rightarrow C_{\text{apx}}$ for large setup times β and large offered loads R , we obtain the denominator, completing both parts of our bound.

Chapter 8

Conclusion

In this chapter, we summarize the thesis, discuss some broader impacts of this thesis, and state some related open problems.

8.1 Summary and Takeaways

In this thesis, we studied the effect of setup times on the queueing behavior of multiserver systems. In particular, we studied how the average waiting time $\mathbb{E}[T_Q]$ in the M/M/k/Setup depends on the system parameters like the number of servers k , the average setup time β , and the load ρ . In Chapter 1, we first noted that the fundamental difficulty in analyzing setup in multiserver systems was the fact that multiple servers can be *in setup* at the same time. We then noted that all prior theoretical work made the simplifying assumption that setup times were distributed i.i.d. Exponential, even though, practically-speaking, setup times are much closer to *Deterministic*; see Chapter 2 for more details. Furthermore, we found in simulation that this distributional assumption has a large impact on the behavior of the system: systems with Deterministic setup times have very different behavior from systems with Exponential setup times.

Accordingly, we narrowed our focus to studying the average waiting time in the M/M/k/Setup-Deterministic (defined in Chapter 3), deriving the first-ever lower and upper bounds on this quantity in Chapters 5 and 6, respectively. Next, in Chapter 7, we described how to take the tightest parts of our bounds and combine them to make an approximation which is extremely accurate. Finally, in this chapter, we summarize our results and state the practical takeaways of our work:

- that the *average waiting time in the M/M/k/Setup-Exponential is drastically smaller* than the average waiting time in the corresponding M/M/k/Setup-Deterministic (Section 1.3);
- that *our approximation is highly accurate* in predicting the average waiting time in the M/M/k/Setup-Deterministic (Figure 1.7);
- and that the simplicity and accuracy of *our approximation radically simplifies capacity provisioning* for dynamically-scaled systems (Section 1.4.3).

8.2 Broader Impacts

This thesis has the potential to impact a large number of different fields, since setup times arise in so many different settings.

8.2.1 Computer Science

In computer science, setup times arise most directly when performing dynamic-scaling in the cloud. There, booting up another container (or virtual machine) might take a few seconds while the actual runtime of a specific task might only take a few milliseconds. Because we do not fully understand how to manage systems with setup times, these servers could be burning much more energy than necessary; when Google introduced Autopilot[34], they were able to cut resource waste in half, from 46% to 23%. Moreover, the energy that these datacenters waste does not just affect these companies' profitability —it also affects our climate via CO_2 emissions and increased demand on the energy grid. Given that one percent of *all power globally* is spent running these datacenters, if we can save two or three percent more energy in their operations, that would be a significant gain for the entire world.

8.2.2 Operations/Management

From an operations/management perspective, the effect of setup times is well-illustrated in employee turnover. When hiring, it might take months to fully onboard a new team member, whereas a typical task might be completed in a day; on the other hand, many employees can be laid off more-or-less instantly. The way in which a firm goes about hiring people, migrating them between different teams, and deciding to lay them off is a great example of the human side of dynamic scaling. Effective management is timelessly relevant, and a setup-time-oriented perspective could provide insights and tools in the same vein as the Pollaczek–Khinchine formula or the Erlang-C model.

8.2.3 Healthcare

Setup times also occur in the medical setting, e.g. when managing on-call doctors. Because patient need (i.e. service demand) is unpredictable, some doctors are often kept “on-call” for up to 36 hours at a time. While on-call, although a physician may not always have work to do, if their service is requested, then they are expected to respond within, say, 30 minutes (which includes travel time to the hospital, if required). For context, most requests can be handled in a very short amount of time, e.g. under a minute. Because these physicians must stay ready-to-respond for multiple days, the current on-call system can lead to extreme sleep deprivation and, accordingly, a poor standard of care for patients. Along the lines of this thesis, further research on dynamically allocating physicians might someday lead us to a new, more sustainable on-call system, with both better care quality and better physician well-being.

8.3 Open Problems

8.3.1 Standby States

Within this thesis, we assume that servers have two persistent states: *on* and *off*; for some systems, this assumption turns out to be wrong. In reality, many servers possess intermediate *standby* states. A server on *standby* takes a shorter amount of time to get ready than a server that is completely *off*, but it also burns more energy. Since the setup process itself takes energy, by using these *standby* states cleverly, we might be able to both improve performance and improve energy efficiency within these systems. Given this, one might ask: **“When should we put a server on *standby* versus turning it completely *off*? What are the benefits of using the *standby* states?”**

8.3.2 Analyzing Tail Performance in the M/M/k/Setup.

Another important open problem lies in analyzing tail performance in the M/M/k/Setup. For context, when customers purchase cloud hosting, an ubiquitous component in their purchase agreements is some kind of “tail/deadline constraint” on their job delay. For example, the agreement will stipulate that “95% of submitted jobs must complete service within one second of their arrival,” with some sort of financial penalty if this constraint is not honored.

Tail constraints in queueing pose a number of technical challenges. In even the single-server case, we do not yet understand how to schedule jobs to optimally meet these constraints. In the multiserver case, though, we have another perspective from which we can analyze the problem: that of *dynamic-scaling*. Instead of thinking about how to schedule these tail-constrained jobs, we can instead think about how we can dynamically-scale our system to ensure these tail constraints are met. This scaling perspective provides a natural way of thinking about the different costs involved. With enough servers, we should be able to ensure that our tail-constraint is met. As such, we can now ask: **“How and when should a system use additional servers to satisfy a given tail constraint?”**

The above question is challenging, and worth considering even in systems without setup times. However, as we have made clear throughout this thesis, setup times often have an enormous impact on the queueing behavior of a dynamically-scaled system. Although there exists extensive study of the performance of dynamic staffing [6, 31], especially in the time-varying arrival rate case [7, 21], much of that work has yet to be extended to the setup time case. As such, we should also ask a more fundamental question: **“How does setup time impact the *distribution* of waiting time in the M/M/k/Setup?”**

Appendix A

Miscellaneous Claims

A.1 Proof of Multiplicative Tightness.

We now show that the upper and lower bounds of Theorems 6.1 and 5.2, respectively, differ by at most a multiplicative factor.

Theorem A.1. *The bounds of Theorems 6.1 and 5.2 lie within a constant multiplicative factor of each other. In particular, using $=_c$ to denote equivalence modulo a multiplicative constant,*

$$\mathbb{E}[Q(\infty)] =_c \beta\sqrt{R} + \frac{1}{1-\rho}. \quad (\text{A.1})$$

A.1.1 Proof for Lower Bound.

We prove Theorem A.1 in two parts, showing equivalence for the lower bound, then for the upper bound. For the lower bound, we first discard all constants and a number of terms in the denominator, since $\beta > \frac{1}{\mu}$ by assumption. Doing so, we obtain

$$\mathbb{E}[Q(\infty)] \geq_c \frac{\mu\beta^2\sqrt{R} + \frac{[L_1\beta\sqrt{R} - k(1-\rho)]^+}{\mu k(1-\rho)}}{\beta + \frac{\beta\sqrt{R}}{\mu k(1-\rho)}} \left[[L_1\beta\sqrt{R} - k(1-\rho)]^+ + \frac{1}{1-\rho} \right], \quad (\text{A.2})$$

where \geq_c denotes that the inequality holds up to an (unspecified) constant factor.

Replacing the $[\cdot]^+$ term.

Now, we show that the $[\beta\sqrt{R} - k(1-\rho)]^+$ term can be replaced by the term $\mu\beta\sqrt{R}$, while losing only a constant factor; this turns out to be the difficult part. We approach this by casing on whether the positive term $\frac{1}{2}L_1\mu\beta\sqrt{R} \geq k(1-\rho)$.

First case. If we have $\frac{1}{2}L_1\mu\beta\sqrt{R} \geq k(1-\rho)$, then $[\mu L_1\beta\sqrt{R} - k(1-\rho)]^+ \geq \frac{1}{2}L_1\mu\beta\sqrt{R} =_c \mu\beta\sqrt{R}$.

Second case. In the second case, assume that $\frac{1}{2}L_1\beta\sqrt{R} < k(1 - \rho)$. In this case, even if we increase the value of the numerator by replacing the $[\cdot]^+$ term, the relevant term in the numerator becomes

$$\frac{\beta\sqrt{R}}{k(1 - \rho)} \left[\mu\beta\sqrt{R} + \frac{1}{1 - \rho} \right] = \mu\beta^2\sqrt{R} \cdot \frac{\sqrt{R}}{k(1 - \rho)} \cdot \left[1 + \frac{1}{\mu\beta\sqrt{R}} \right] =_c \mu\beta^2\sqrt{R} \cdot \frac{\sqrt{R}}{k(1 - \rho)} \leq_c \mu\beta^2\sqrt{R} \cdot \frac{1}{\beta}$$

where in the second equality we have used that $\sqrt{R} \geq 1$ and $\mu\beta \geq 1$, and in the final inequality we have used our case assumption. From here, it's clear that one can replace the term $[L_1\mu\beta\sqrt{R} - k(1 - \rho)]^+$ with the term $\mu\beta\sqrt{R}$ without altering the scaling behavior of numerator. In other words,

$$\mathbb{E}[Q(\infty)] \geq_c \frac{\mu\beta^2\sqrt{R} + \frac{\mu\beta\sqrt{R}}{\mu k(1-\rho)} \left[\beta\sqrt{R} + \frac{1}{1-\rho} \right]}{\beta + \frac{\beta\sqrt{R}}{\mu k(1-\rho)}} =_c \mu\beta\sqrt{R} + \frac{\sqrt{R}}{k(1 - \rho) + \sqrt{R}} \frac{1}{1 - \rho}. \quad (\text{A.3})$$

Bounding the final term.

We now show equivalence for this final term, i.e. that

$$\mu\beta\sqrt{R} + \frac{\sqrt{R}}{k(1 - \rho) + \sqrt{R}} \frac{1}{1 - \rho} =_c \mu\beta\sqrt{R} + \frac{1}{1 - \rho}. \quad (\text{A.4})$$

To do so, we bound the rightmost term in (A.3). Note that, since $\frac{\sqrt{R}}{k - R + \sqrt{R}} \leq \frac{R}{k - R + R} = \rho$, in order for this term to have an appreciable effect on the scaling, we must have that $\frac{\rho}{1 - \rho} \geq_c \mu\beta\sqrt{R}$, or, phrased more usefully, we must have $\sqrt{R} \geq_c \mu\beta k(1 - \rho)$. But even in this case, we can bound the factor in the rightmost term of (A.3) with $\frac{\sqrt{R}}{k(1 - \rho) + \sqrt{R}} \geq_c \frac{\beta k(1 - \rho)}{k(1 - \rho) + \mu\beta k(1 - \rho)} = \frac{\mu\beta}{1 + \mu\beta} =_c 1$; the multiplicative equivalence (A.4) follows.

A.1.2 Proof for Upper Bound.

Initial Steps.

The proof for the upper bound follows along the same lines. First, note that the terms outside of the fraction are $\frac{R}{M} \leq \frac{R}{k(1 - \rho)} = \frac{\rho}{1 - \rho}$ and $\sqrt{\mu\beta R} \ll \mu\beta\sqrt{R}$. Discarding the lower order terms and constants, we obtain

$$\mathbb{E}[Q(\infty)] \leq_c \mu\beta\sqrt{R} + \frac{\rho}{1 - \rho} + \frac{\mu\beta^2\sqrt{R} + g\left(9(\mu\beta)^2 R, 3\mu\beta\sqrt{R}, k(1 - \rho)\right)}{\beta + \frac{\mu\beta\sqrt{R}}{\mu k(1 - \rho)}}, \quad (\text{A.5})$$

where one should recall that $M =_c \min(k(1 - \rho), \sqrt{\mu\beta R})$. For the terms in the fraction, the denominator of the upper bound is already up-to-constants-equivalent to the denominator of

(A.2). It thus suffices to show that the numerator of the upper bound aligns with the numerator of (A.3). However, note that, by definition, the function

$$g\left(9(\mu\beta)^2R, 3(\mu\beta)\sqrt{R}, k(1-\rho)\right) \triangleq \frac{3\mu\beta\sqrt{R}}{k(1-\rho)} \left[\frac{3\mu\beta\sqrt{R}+1}{2} + \frac{1}{1-\rho} \right];$$

thus, the terms are clearly equivalent up to scaling. \square

A.2 Construction and Coupling Claims.

A.2.1 Construction

We now discuss how we formally construct this system using Poisson processes; being explicit here will prove useful when we make coupling arguments in the future.

The arrival and departure processes. We take the number of jobs that have arrived at time t to be $\Pi_A(t)$, where Π_A is a Poisson process of rate $k\lambda$. In a slight abuse of notation, we let $\Pi_A([a, b])$ denote the number of arrivals that occur in the interval $[a, b]$; we apply the same extension to all other counting processes mentioned here. We set the potential departure process of, say, server i to be $\Pi_i(t)$, where Π_i is a Poisson process of rate μ . A potential departure from server i only “counts” if server i is busy when that potential departure occurs, i.e., if the number of busy servers $Z(t) \geq i$ at the time. Thus, the total number of departures from our system by time t is, taking integrals with respect to the Poisson processes Π_i as counting processes,

$$D(t) \triangleq \sum_{i=1}^k \int_0^t \mathbf{1}\{Z(s) \geq i\} d\Pi_i(s).$$

The number of busy servers $Z(t)$. To find the number of busy servers $Z(t)$, one could count the number of setup completion events that have occurred so far and the number of server shutoffs that have occurred so far; this description is a bit difficult to work with. Alternatively, one can see from the initial description of setup dynamics that server i is *on* at time t if and only if the total number of jobs $N(s) \geq i$ for all $s \in [t - \beta, t]$, where one should recall that β is the setup time. An easier description of $Z(t)$ follows:

$$Z(t) = \min\left(k, \min_{s \in [t-\beta, t]} N(s)\right).$$

A departure operator. We can extend our departure process $D(t)$ to a departure operator $\mathcal{D}[f(s)](\mathcal{I})$ which takes a function $f(s) \in \{0, 1, \dots, k\}$ defined on some interval \mathcal{I} and computes the number of departures that would occur in that interval provided that the number of busy servers $Z(s) = f(s)$, i.e.

$$\mathcal{D}[f(s)]((a, b)) \triangleq \sum_{i=1}^k \int_a^b \mathbf{1}\{f(s) \geq i\} d\Pi_i(s).$$

Note that the total number of departures can now be written as $D(t) = \mathcal{D}[Z(s)]([0, t])$.

A.2.2 Three Coupling Claims

We now describe three useful claims applied throughout the proof. The first, we will state and prove immediately. For the latter two, we first give a high-level explanation, then state and prove them.

Basic coupling claim: Maintaining an initial relation.

Claim A.1 (Basic Coupling). *Suppose that we have two processes N_1 and N_2 with an initial relation $N_1(a) \leq N_2(a)$, where the behavior of each process is governed, for all times s from a up to some stopping time τ , by the equation*

$$N_j(s) \triangleq N_j(a) + \Pi_A((a, s]) - \mathcal{D}[Z_j(x)]((a, s]), \text{ for } j \in \{1, 2\}.$$

Furthermore, suppose that the first system's number of busy servers $Z_1(s) \geq Z_2(s)$ for all times $s \in [a, \tau]$. Then, for all $s \in [a, \tau]$, the relation is maintained, i.e. $N_1(s) \leq N_2(s)$.

Proof. We show equivalently that $N_2(s) - N_1(s) \geq 0$. Applying the definitions of N_1 and N_2 ,

$$\begin{aligned} N_2(s) - N_1(s) &= N_2(a) - N_1(a) + [\mathcal{D}[Z_1(x)]((a, s]) - \mathcal{D}[Z_2(x)]((a, s])] \\ &\geq [\mathcal{D}[Z_1(x)]((a, s]) - \mathcal{D}[Z_2(x)]((a, s))] \\ &= \sum_{i=1}^k \int_a^s \mathbf{1}\{Z_1(x) \geq i\} d\Pi_i(x) - \sum_{i=1}^k \int_a^s \mathbf{1}\{Z_2(x) \geq i\} d\Pi_i(x) \\ &= \sum_{i=1}^k \int_a^s \left[\mathbf{1}\{Z_1(x) \geq i\} - \mathbf{1}\{Z_2(x) \geq i\} \right] d\Pi_i(x). \end{aligned}$$

Since $Z_1(x) \geq Z_2(x)$, the integrand $\left[\mathbf{1}\{Z_1(x) \geq i\} - \mathbf{1}\{Z_2(x) \geq i\} \right] \geq 0$; the claim follows.

□

Statement and proof of remaining coupling claims.

High-level explanation. This claim leads nicely into a couple more claims. Both are concerned with bounding a quantity involving a general “down-crossing” time. In particular, our analysis will begin at a stopping time τ and will “end” at the down-crossing time d_{gen} , where $d_{\text{gen}} \triangleq \min\{t \geq 0 : N(t + \tau) \leq h\}$ is the length of time it takes for the number of jobs $N(t)$ to become lower than some given threshold h . The first claim, Claim A.2, uses a coupling argument to bound the expected integral of $N(t)$ from some arbitrary time τ until $N(t)$ drops below some pre-defined threshold h , provided that one has a lower bound on the number of busy servers $Z(t)$ over that period. The second claim, Claim A.3, uses a related argument to bound the probability that $N(t)$ drops below some threshold h within some amount of time ℓ , given that one has bounds on $Z(t)$ over the relevant period. We defer the proof of these claims to Sections A.2.3 and A.2.4.

Claim A.2 (Coupling Integral Bound). *Let τ be some stopping time and d_{gen} be the next down-crossing as described in Section A.2.2. Suppose that, at time τ , we have a lower bound on the number of busy servers over a period, i.e. we know that the number of busy servers $Z(t) \geq R - j$, for all $t \in [\tau, \tau + \min(\ell, d_{\text{gen}})]$ and for some non-negative j . Then we have the following bound on the integral over this time period:*

$$\mathbb{E} \left[\int_{\tau}^{\tau + \min(d_{\text{gen}}, \ell)} [N(t) - h] dt \middle| \mathcal{F}_{\tau} \right] \leq \ell \cdot [N(\tau) - h]^+ + \frac{1}{2} \mu j \ell^2.$$

Claim A.3 (Coupling Probability Bound). *Let τ be some stopping time and d_{gen} be the next down-crossing as described in Section A.2.2. We consider two cases.*

*In the first case, suppose that we have a **lower** bound on the number of busy servers $Z(t)$ over some length ℓ interval starting at time τ , i.e. the busy servers $Z(t) \geq R - j$, for all $t \in [\tau, \tau + \min(\ell, d_{\text{gen}})]$ and for some non-negative j . Then, we can bound the threshold-crossing probability by*

$$\Pr(d_{\text{gen}} < \ell | \mathcal{F}_{\tau}) \geq 2\Phi \left(- \left[\frac{N(\tau) - h + \mu j \ell}{\sqrt{\ell(2k\lambda - \mu j)}} \right] \right) - \frac{2}{3\sqrt{\ell(2k\lambda - \mu j)}}.$$

In particular, if $N(\tau) - h = c_1 \sqrt{\mu\beta R}$, then the probability $\Pr(d_{\text{gen}} < \ell | \mathcal{F}_{\tau}) \geq 2\Phi \left(-\frac{c_1}{\sqrt{2}} \right) - \frac{1}{100}$.

*In the second case, suppose that we instead have the **upper** bound on $Z(t) \leq R$ during this interval instead. Then,*

$$\Pr(d_{\text{gen}} < \ell | \mathcal{F}_{\tau}) \leq 2\Phi \left(- \left[\frac{N(\tau) - h}{\sqrt{2\ell k \lambda}} \right] \right) - \frac{2}{3\sqrt{2k\lambda\ell}}.$$

As before, if $N(\tau) - h = c_2 \sqrt{\mu\beta R}$, then the probability $\Pr(d_{\text{gen}} < \ell | \mathcal{F}_{\tau}) \leq 2\Phi \left(-\frac{c_2}{\sqrt{2}} \right) + \frac{1}{100}$.

A.2.3 Proof of Claim A.2, the Coupling Integral Bound.

Proof. We prove this claim in three parts. First, we construct a coupled process $\tilde{N}(t) \geq N(t)$ on the interval of interest. Then, we give an upper bound on $\mathbb{E} \left[\int_{\tau}^{\tau + \min(\ell, d_{\text{gen}})} \tilde{N}(t) dt \middle| \mathcal{F}_{\tau} \right]$. Define $\tilde{N}(t)$ as

$$\tilde{N}(t) \triangleq N(\tau) + A(\tau, t) - \mathcal{D}[R - j]((\tau, t)).$$

Then, by Claim A.1, we have that

$$\tilde{N}(t) \geq N(t).$$

on the interval of interest. To develop the integral, we first move the minimum from the bounds of integration into the integrand. In particular, we note that the quantity $N(d_{\text{gen}}) - h = 0$, and thus, for any $t > \tau + d_{\text{gen}}$, the quantity $N(\min(\tau + d_{\text{gen}}, t)) - h = 0$. On the other hand, for any

$t < \tau + d_{\text{gen}}$, the quantity $N(\min(\tau + d_{\text{gen}}, t)) = N(t)$. It follows that

$$\begin{aligned}
\int_{\tau}^{\tau + \min(\ell, d_{\text{gen}})} [N(t) - h] dt &= \int_{\tau}^{\tau + \min(\ell, d_{\text{gen}})} [N(\min(t, \tau + d_{\text{gen}})) - h] dt \\
&= \int_{\tau}^{\tau + \min(\ell, d_{\text{gen}})} [N(\min(t, \tau + d_{\text{gen}})) - h] dt \\
&\quad + \int_{\tau + \min(\ell, d_{\text{gen}})}^{\tau + \ell} [N(\min(t, \tau + d_{\text{gen}})) - h] dt \\
&= \int_{\tau}^{\tau + \ell} [N(\min(t, \tau + d_{\text{gen}})) - h] dt \\
&\leq \int_{\tau}^{\tau + \ell} [\tilde{N}(\min(t, \tau + d_{\text{gen}})) - h] dt.
\end{aligned}$$

Defining $d_{\text{gen}}^{\tilde{}} \triangleq \min\{t > 0 : \tilde{N}(\tau + t) \leq h\}$, since $\tilde{N}(t) \geq N(t)$, we know both that $d_{\text{gen}}^{\tilde{}} \geq d_{\text{gen}}$ and that, for any $t \in [\tau + d_{\text{gen}}, \tau + d_{\text{gen}}^{\tilde{}}]$,

$$\tilde{N}(t) - h \geq 0.$$

Moreover, the process $V(t)$ defined as

$$V(t) \triangleq \tilde{N}(t) - \mu j t$$

is a martingale. Thus, we have

$$\int_{\tau}^{\tau + \ell} [\tilde{N}(\min(t, \tau + d_{\text{gen}})) - h] dt \leq \int_{\tau}^{\tau + \ell} [\tilde{N}(\min(t, \tau + d_{\text{gen}}^{\tilde{}})) - h] dt.$$

Taking the expectation, we find that

$$\mathbb{E} \left[\int_{\tau}^{\tau + \ell} [\tilde{N}(\min(t, \tau + d_{\text{gen}}^{\tilde{}})) - h] dt \middle| \mathcal{F}_{\tau} \right] \tag{A.6}$$

$$\begin{aligned}
&= \int_{\tau}^{\tau + \ell} \mathbb{E} [\tilde{N}(\min(t, \tau + d_{\text{gen}}^{\tilde{}})) - h \middle| \mathcal{F}_{\tau}] dt \\
&= \int_{\tau}^{\tau + \ell} \mathbb{E} [V(\min(t, \tau + d_{\text{gen}}^{\tilde{}})) + \mu j (\min(\tau + d_{\text{gen}}^{\tilde{}}, t)) - h \middle| \mathcal{F}_{\tau}] dt \\
&= \int_{\tau}^{\tau + \ell} \mathbb{E} [V(\tau) + \mu j (\min(\tau + d_{\text{gen}}^{\tilde{}}, t)) - h \middle| \mathcal{F}_{\tau}] dt \tag{A.7} \\
&\leq \int_{\tau}^{\tau + \ell} \mathbb{E} [V(\tau) + \mu j t - h \middle| \mathcal{F}_{\tau}] dt \\
&= \int_{\tau}^{\tau + \ell} \mathbb{E} [\tilde{N}(\tau) - \mu j \tau + \mu j t - h \middle| \mathcal{F}_{\tau}] dt \\
&= [\tilde{N}(\tau) - h] \ell + \frac{1}{2} \mu j \ell^2,
\end{aligned}$$

where (A.7) is an application of Doob's Optimal Stopping Theorem. \square

A.2.4 Proof of Claim A.3, the Coupling Probability Bound.

Proof. We prove this result in three parts. First, we use Claim A.1 to construct a process $\tilde{N}(t) \geq N(t)$ on the interval of interest. Afterwards, we analyze the down-crossing probability of this coupled process. In particular, we use a reflection argument to show that

$$\Pr(d_{\text{gen}} < \ell) \geq 2 \Pr\left(\tilde{N}(\tau + \ell) \leq h\right),$$

then use a Berry-Esseen bound to bound this final probability. In what follows, we focus on the lower-bound; the upper bound follows in precisely the same way.

To construct our coupled process, we note that, by assumption, the number of busy servers $Z(t) \geq R - j$ for any $t \in [\tau, \tau + \min(\ell, d_{\text{gen}})]$. Thus, by Claim A.1, the process $\tilde{N}(t)$ defined as

$$\tilde{N}(t) \triangleq N(\tau) + A(\tau, \tau + t) + \mathcal{D}[R - j](\tau, \tau + t)$$

is an upper bound for $N(t + \tau)$, i.e.

$$\tilde{N}(t) \geq N(\tau + t)$$

for any $t \in [0, \min(\ell, d_{\text{gen}})]$. By definition, we have that

$$\begin{aligned} \Pr(d_{\text{gen}} < \ell) &= \Pr\left(\inf_{t \in [0, \ell]} N(\tau + t) \leq h\right) \\ &\geq \Pr\left(\inf_{t \in [0, \ell]} \tilde{N}(t) \leq h\right). \end{aligned}$$

From a reflection argument, since \tilde{N} is upwards-biased,

$$\begin{aligned} \Pr\left(\inf_{t \in [0, \ell]} \tilde{N}(t) \leq h\right) &= \Pr\left(\inf_{t \in [0, \ell]} \tilde{N}(t) \leq h, \tilde{N}(\ell) < h\right) + \Pr\left(\inf_{t \in [0, \ell]} \tilde{N}(t) \leq h, \tilde{N}(\ell) \geq h\right) \\ &\geq 2 \Pr\left(\inf_{t \in [0, \ell]} \tilde{N}(t) \leq h, \tilde{N}(\ell) < h\right) \\ &= 2 \Pr\left(\tilde{N}(\ell) < h\right). \end{aligned}$$

Let $\sigma \triangleq \sqrt{\ell(2k\lambda - \mu j)}$. We now apply Now, assume that, for any x ,

$$\left|\Pr\left(\tilde{N}(\ell) < \tilde{N}(0) + \mu j \ell + x\sigma\right) - \Phi\left(\frac{x}{\sigma}\right)\right| \leq \frac{0.3328}{\sigma}, \quad (\text{A.8})$$

we have

$$\begin{aligned} \Pr\left(\tilde{N}(\ell) < h\right) &= \Pr\left(\tilde{N}(\ell) < \tilde{N}(0) + \mu j \ell + \frac{h - \mu j \ell - \tilde{N}(0)}{\sigma} \cdot \sigma\right) \\ &\geq \Phi\left(\frac{h - \mu j \ell - \tilde{N}(0)}{\sigma}\right) - \frac{1}{3\sigma} \\ &= \Phi\left(-\frac{[N(\tau) - h + \mu j \ell]}{\sigma}\right) - \frac{1}{3\sigma}. \end{aligned}$$

Putting this all together, we find

$$\Pr(d_{\text{gen}} < \ell | \mathcal{F}_\tau) \geq 2\Phi\left(-\frac{[N(\tau) - h + \mu j \ell]}{\sigma}\right) - \frac{2}{3\sigma}.$$

From here, then, it suffices to show (A.8). To begin, note that, if we choose some arbitrarily large n and define

$$X_i \triangleq \Pi'_i\left(\frac{k\lambda\ell}{n}\right) - \Pi''_i\left(\frac{\mu(R-j)\ell}{n}\right) - \frac{\mu j \ell}{n},$$

where each $\Pi(y)$ is an independent Poisson random variable with mean y , then

$$\tilde{N}(\ell) =_d \sum_{i=1}^n X_i + \mu j \ell + \tilde{N}(0).$$

To compute the moments of X_i , note that one can define centered Poisson random variables $A_i = \Pi\left(\frac{k\lambda\ell}{n}\right) - \frac{k\lambda\ell}{n}$ and $B_i = \Pi\left(\frac{\mu(R-j)\ell}{n}\right) - \frac{\mu(R-j)\ell}{n}$, and then take $X_i = A_i - B_i$. Doing this, one finds that

$$\mathbb{E}[X_i^2] = \mathbb{E}[(A_i - B_i)^2] = \frac{k\lambda\ell}{n} + \frac{\mu(R-j)\ell}{n} = \frac{\mu(2R-j)\ell}{n}$$

and, using the triangle inequality, that

$$\mathbb{E}[|X_i|^3] = \mathbb{E}[|A_i - B_i|^3] \leq \mathbb{E}[|A_i|^3] + \mathbb{E}[|B_i|^3] = \frac{\mu(2R-j)\ell}{n} + o\left(\frac{1}{n^2}\right).$$

We now apply the main result of [35]. Let $\sigma_n \triangleq \sqrt{\mathbb{E}[X_i^2]} = \sqrt{\frac{\mu(2R-j)\ell}{n}} = \frac{\sigma}{\sqrt{n}}$ and note that $\rho_n = \mathbb{E}[|X_i|^3] < \sigma_n + o\left(\frac{1}{n^2}\right)$ (from [5]). Then, noting that $\rho_n \geq 1.286\sigma_n^3$ for sufficiently large n , we have

$$\max_x \left| \Pr\left(\frac{\sum X_i}{\sqrt{n}\sigma_n} < x\right) - \Phi(x) \right| \leq \frac{0.3328\rho_n + 0.429\sigma_n^3}{\sigma_n^3\sqrt{n}} = \frac{0.3328}{\sqrt{\mu(2R-j)\ell}} + o\left(\frac{1}{n}\right).$$

Now noting that

$$\frac{\sum_{i=1}^n X_i}{\sqrt{n}\sigma} = \frac{\tilde{N}(\ell) - \tilde{N}(0) - \mu j \ell}{\sigma}$$

and taking $n \rightarrow \infty$, we have our result. \square \square

A.2.5 Proof of Claim A.4: Bound on Expected Value After Coupling.

Claim A.4 (Bound on Expected Value after Coupling.). *Let τ be some stopping time and d_{gen} be the next down-crossing as described in Section A.2.2. Suppose that we have a **lower** bound on the number of busy servers $Z(t)$ over some length ℓ interval starting at time τ , i.e. the busy servers $Z(t) \geq R - j$, for all $t \in [\tau, \tau + \min(\ell, d_{\text{gen}})]$ and for some non-negative j . Then, bounding the first moment,*

$$\mathbb{E}[N(\tau + \ell) - h] \mathbf{1}_{d_{\text{gen}} > \ell} | \mathcal{F}_\tau \leq [N(\tau) - h] + \mu j \ell, \quad (\text{A.9})$$

and, bounding the second moment,

$$\mathbb{E}[N(\tau + \ell) - h] \mathbf{1}_{d_{\text{gen}} \geq \ell} \leq [N(\tau) - h + \mu j \ell]^2 + 2\mu R \ell. \quad (\text{A.10})$$

Proof.

The proof is essentially an application of Doob's Optional Stopping Theorem to an appropriately selected martingale. To begin, we define a coupled process $\tilde{N}(t)$ with

$$\tilde{N}(t - \tau) \triangleq N(\tau) + A[\tau, t] - \mathcal{D}[R - j](\tau, t);$$

by Claim A.1, we know that $\tilde{N}(t - \tau) \geq N(t)$ for any $t \in [\tau, \tau + \min(d_{\text{gen}}, \ell)]$, and that the coupled hitting time $d_{\text{gen}}^{\tilde{}} \triangleq \min\{t > 0 : \tilde{N}(t) \leq h\}$ can not be smaller than the original hitting time d_{gen} . It follows that

$$N(\tau + \ell) \mathbf{1}_{d_{\text{gen}} > \ell} \leq \tilde{N}(\ell) \mathbf{1}_{d_{\text{gen}}^{\tilde{}} > \ell}.$$

Thus, we bound coupled versions of (A.9) and (A.10).

Construction of martingales. We now construct our martingales and set up the language of optional stopping. Note that, for any process $\tilde{N}(t)$ with independent, stationary increments, both functions V_1 and V_2 , defined as

$$V_1(t) \triangleq [\tilde{N}(t) - h] - \mathbb{E}[\tilde{N}(t) - \tilde{N}(0)]$$

and

$$\begin{aligned} V_2(t) &\triangleq [\tilde{N}(t) - h - \mathbb{E}[\tilde{N}(t) - \tilde{N}(0)]]^2 - \mathbb{E}[[\tilde{N}(t) - h - \mathbb{E}[\tilde{N}(t) - \tilde{N}(0)]]^2] \\ &= (\tilde{N}(t) - h - \mu j t)^2 - \mu(2R - j)t \end{aligned}$$

are martingales [23]. Moreover, one has that

$$\begin{aligned} [\tilde{N}(\ell) - h] \mathbf{1}_{d_{\text{gen}}^{\tilde{}} > \ell} &= [\tilde{N}(\min(d_{\text{gen}}^{\tilde{}}, \ell)) - h] \mathbf{1}_{d_{\text{gen}}^{\tilde{}} > \ell} \\ &= [\tilde{N}(\min(d_{\text{gen}}^{\tilde{}}, \ell)) - h] \mathbf{1}_{d_{\text{gen}}^{\tilde{}} > \ell} + [\tilde{N}(\min(d_{\text{gen}}^{\tilde{}}, \ell)) - h] \mathbf{1}_{\ell \leq d_{\text{gen}}^{\tilde{}}} \\ &= [\tilde{N}(\min(d_{\text{gen}}^{\tilde{}}, \ell)) - h]. \end{aligned}$$

Proof of (A.9). Combining these facts allows us to prove our desired result. Applying Doob's Optional Stopping Theorem along with our previous deductions, we obtain

$$\begin{aligned} \mathbb{E}[[N(\tau + \ell) - h] \mathbf{1}_{d_{\text{gen}} > \ell} | \mathcal{F}_\tau] &\leq \mathbb{E}[[\tilde{N}(\ell) - h] \mathbf{1}_{d_{\text{gen}}^{\tilde{}} > \ell}] \\ &= \mathbb{E}[\tilde{N}(\min(d_{\text{gen}}^{\tilde{}}, \ell)) - h] \\ &= \mathbb{E}[V_1(\min(d_{\text{gen}}^{\tilde{}}, \ell))] + \mu j \mathbb{E}[\min(d_{\text{gen}}^{\tilde{}}, \ell)] \\ &= \mathbb{E}[V_1(0)] + \mu j \mathbb{E}[\min(d_{\text{gen}}^{\tilde{}}, \ell)] \\ &= [\tilde{N}(0) - h] + \mu j \mathbb{E}[\min(d_{\text{gen}}^{\tilde{}}, \ell)] \\ &\leq [\tilde{N}(0) - h] + \mu j \ell \\ &= [N(\tau) - h] + \mu j \ell. \end{aligned}$$

Proof of (A.10). To do the same for the squared martingale $V_2(t)$, we must first note, via some algebra, that

$$\left(\tilde{N}(t) - h\right)^2 = V_2(t) + \left(\tilde{N}(t) - h\right) \mu j t - \mu j^2 t^2 + \mu (2R - j) t.$$

Now, applying the same deductions we made previously,

$$\begin{aligned} & \mathbb{E} \left[[N(\tau + \ell) - h] \mathbf{1}_{d_{\text{gen}} > \ell} \middle| \mathcal{F}_\tau \right] \\ & \leq \mathbb{E} \left[\left[\tilde{N}(\ell) - h \right]^2 \mathbf{1}_{\tilde{d}_{\text{gen}} > \ell} \right] \\ & = \mathbb{E} \left[\left(\tilde{N} \left(\min \left(\tilde{d}_{\text{gen}}, \ell \right) \right) - h \right)^2 \right] \\ & = \mathbb{E} \left[V_2 \left(\min \left(\tilde{d}_{\text{gen}}, \ell \right) \right) \right] + \mathbb{E} \left[\left(\tilde{N} \left(\min \left(\tilde{d}_{\text{gen}}, \ell \right) \right) - h \right) \mu j \min \left(\tilde{d}_{\text{gen}}, \ell \right) \right] - \mu j^2 \left(\min \left(\tilde{d}_{\text{gen}}, \ell \right) \right)^2 \\ & \quad + \mu (2R - j) \mathbb{E} \left[\min \left(\tilde{d}_{\text{gen}}, \ell \right) \right] \\ & \leq \mathbb{E} \left[V_2 \left(\min \left(\tilde{d}_{\text{gen}}, \ell \right) \right) \right] + \mathbb{E} \left[\left(\tilde{N} \left(\min \left(\tilde{d}_{\text{gen}}, \ell \right) \right) - h \right) \mu j \ell + \mu (2R) \ell \right] \\ & \leq \mathbb{E} [V_2(0)] + \left[\tilde{N}(0) - h + \mu j \ell \right] \mu j \ell + \mu (2R) \ell \\ & = \left[\tilde{N}(0) - h \right]^2 + \left[\tilde{N}(0) - h \right] \mu j \ell + (\mu j \ell)^2 + \mu (2R) \ell \\ & = \left[\tilde{N}(0) - h + \mu j \ell \right]^2 - \left[\tilde{N}(0) - h \right] \mu j \ell + \mu 2R \ell \\ & \leq \left[\tilde{N}(0) - h + \mu j \ell \right]^2 + 2\mu R \ell \\ & = [N(\tau) - h + \mu j \ell]^2 + 2\mu R \ell. \end{aligned}$$

A.3 Hitting Time Bounds.

A.3.1 Proof of Claim A.5, Discrete-Time Hitting Time Tail Bound.

Claim A.5 (Discrete-Time Hitting Time Tail Bound). *Suppose one has an upwards-biased discrete random walk $V(t)$ where in each step*

$$\Pr(V(t+1) = V(t) + 1 | \mathcal{F}_t) = p = 1 - q,$$

where $p \geq \frac{1}{2} \geq q$. Suppose that $V(0) = 1$ and let the hitting time $\gamma \triangleq \min \{t \in \mathcal{N} : V(t) = 0\}$ be the first timestep where the walk $V(t) = 0$. Then, for $n \geq 1$,

$$\Pr(\gamma \geq 2m + 1) \leq \frac{1}{\sqrt{\pi}} \frac{2q}{\sqrt{m}} \left(1 + \frac{1}{2(m+1)} \right).$$

Moreover, if $p = q = \frac{1}{2}$, then

$$\Pr(\gamma \geq 2m + 1) \geq \frac{1}{\sqrt{\pi}} e^{-\frac{1}{6m}} \frac{1}{\sqrt{m+1}}.$$

Proof

We first note, as in [38], that by a counting argument $\Pr(\gamma = 2\ell + 1) = q(qp)^\ell C_\ell$, where $C_\ell \triangleq \frac{1}{\ell+1} \frac{(2\ell)!}{\ell!\ell!}$ is the ℓ -th Catalan number; note that γ can not be even, since the number of downward steps must exceed the number of upward steps by exactly 1.

We proceed by bounding the Catalan numbers using Stirling's approximation. For $m = 0$, then $\Pr(\gamma \geq 1) = \Pr(\gamma \geq 2) = p$, i.e. the probability that the first step is an upward step. For $m \geq 1$, applying Stirling's approximation and simplifying gives

$$e^{-\frac{1}{6\ell}} \frac{1}{\sqrt{\pi\ell(\ell+1)}} q(4pq)^\ell \leq \Pr(\gamma = 2\ell + 1) \leq \frac{1}{\sqrt{\pi\ell(\ell+1)}} q(4pq)^\ell.$$

Lower bound. Since we are interested in the lower bound only when $q = p = \frac{1}{2}$, we obtain that

$$\begin{aligned} \Pr(\gamma \geq 2m + 1) &\geq \frac{1}{\sqrt{\pi i}} \frac{1}{2} \sum_{\ell=m}^{\infty} \frac{e^{-\frac{1}{6\ell}}}{\sqrt{\ell(\ell+1)}} \\ &\geq \frac{1}{\sqrt{\pi}} e^{-\frac{1}{6m}} \frac{1}{2} \sum_{\ell=m}^{\infty} \frac{1}{\sqrt{\ell(\ell+1)}} \\ &\geq \frac{1}{\sqrt{\pi}} e^{-\frac{1}{6m}} \frac{1}{2} \int_m^{\infty} \frac{1}{(\ell+1)^{3/2}} d\ell \\ &= \frac{1}{\sqrt{\pi}} e^{-\frac{1}{6m}} \frac{1}{\sqrt{m+1}}. \end{aligned}$$

Upper bound. Noting that $4pq \leq 1$, we have likewise that

$$\begin{aligned} \Pr(\gamma \geq 2m + 1) &\leq \frac{1}{\sqrt{\pi}} q \sum_{\ell=m}^{\infty} \frac{1}{\sqrt{\ell(\ell+1)}} \\ &\leq \frac{1}{\sqrt{\pi}} q \frac{1}{\sqrt{m(m+1)}} + \int_m^{\infty} \frac{1}{\ell^{3/2}} d\ell \\ &= \frac{1}{\sqrt{\pi}} q \frac{2}{\sqrt{m}} \left(1 + \frac{1}{2(m+1)} \right). \end{aligned}$$

A.3.2 Proof of Claim A.6, Continuous-Time Hitting Time Tail Bound.

We further extend this discrete-time bound into a continuous-time bound.

Claim A.6 (Continuous-Time Hitting Time Tail Bound). *Suppose one has an Poisson arrival process $Y_A(t)$ of rate $k\lambda$ and a Poisson departure process $Y_D(t)$ of rate $\mu(R - j)$, for some integer $j \geq 0$. Let the continuous random walk $X_c(t) = Y_A(t) - Y_D(t)$, with $X_c(0) = 1$, and define $\gamma_c \triangleq \min\{t > 0 : X_c(t) = 0\}$. Let $\nu = (2R - j)\mu t$. For any $\nu \geq 3$, we have*

$$\Pr(\gamma_c \geq t) \leq \frac{b_1}{\sqrt{2}} \left(\frac{1}{\sqrt{\nu}} + \frac{b_2}{\nu^{3/2}} \right)$$

where $b_1 = \sqrt{\frac{2}{\pi}}$ and $b_2 = 1 + \frac{2.5}{b_1\sqrt{2}}$.

Moreover, if $j = 0$, then

$$\Pr(\gamma_c \geq t) \geq \frac{b_1}{\sqrt{2}} e^{-\frac{1}{3(\nu-1)}} \frac{1}{\sqrt{\nu+2}}.$$

Proof of Upper Bound.

To prove this claim, we first condition on the value of $Y_T = Y_A(t) + Y_D(t)$, the total number of Poisson events during the interval $[0, t]$, then relate that to the same question in a discrete-time random walk, a la Claim A.5. Note that $Y_T \sim \text{Poisson}(\nu)$, and thus

$$\begin{aligned} \Pr(\gamma_c \geq t) &= \Pr(\gamma \geq Y_T) \\ &= \sum_{j=0}^{\infty} e^{-\nu} \frac{\nu^j}{j!} \Pr(\gamma \geq j) \\ &= e^{-\nu} + 2p\nu e^{-\nu} + \sum_{j=3}^{\infty} e^{-\nu} \frac{\nu^j}{j!} \Pr(\gamma \geq j + \mathbf{1}_{j \text{ is even}}) \\ &= e^{-\nu} + 2p\nu e^{-\nu} + \sum_{j=0}^{\infty} e^{-\nu} \frac{\nu^j}{j!} \Pr\left(\gamma \geq 2 \left(\frac{j + \mathbf{1}_{j \text{ is even}} - 1}{2}\right) + 1\right). \end{aligned}$$

Applying the discrete upper bound to the sum, we obtain

$$\begin{aligned} &\sum_{j=3}^{\infty} e^{-\nu} \frac{\nu^j}{j!} \Pr\left(\gamma \geq 2 \left(\frac{j + \mathbf{1}_{j \text{ is even}} - 1}{2}\right) + 1\right) \\ &\leq b_1 \sqrt{2q} \sum_{j=3}^{\infty} e^{-\nu} \frac{\nu^j}{j!} \frac{1}{\sqrt{j + \mathbf{1}_{j \text{ is even}} - 1}} \left(1 + \frac{1}{(j + \mathbf{1}_{j \text{ is even}} + 1)}\right) \\ &= b_1 \sqrt{2q} \frac{1}{\nu} \sum_{j=3}^{\infty} e^{-\nu} \frac{\nu^{(j+1)}}{(j+1)!} \frac{1}{\sqrt{j + \mathbf{1}_{j \text{ is even}} - 1}} \left(j + 1 + \frac{j+1}{j + \mathbf{1}_{j \text{ is even}} + 1}\right) \\ &\leq b_1 \sqrt{2q} \frac{1}{\nu} \sum_{j=3}^{\infty} e^{-\nu} \frac{\nu^{(j+1)}}{(j+1)!} \frac{j+2}{\sqrt{j + \mathbf{1}_{j \text{ is even}} - 1}} \\ &\leq b_1 \sqrt{2q} \frac{1}{\nu} \sum_{j=3}^{\infty} e^{-\nu} \frac{\nu^{(j+1)}}{(j+1)!} \frac{j + \mathbf{1}_{j \text{ is even}} - 1 + 3}{\sqrt{j + \mathbf{1}_{j \text{ is even}} - 1}} \\ &= b_1 \sqrt{2q} \frac{1}{\nu} \sum_{j=3}^{\infty} e^{-\nu} \frac{\nu^{(j+1)}}{(j+1)!} \left(\sqrt{j + \mathbf{1}_{j \text{ is even}} - 1} + \frac{3}{\sqrt{j + \mathbf{1}_{j \text{ is even}} - 1}}\right). \end{aligned}$$

From here, we note that the function $f(x) = \sqrt{x} + \frac{3}{\sqrt{x}}$ is both increasing and concave for all $x \geq 3$. After increasing the argument and applying Jensen's inequality, we find that

$$\begin{aligned} &\leq b_1 \sqrt{2} q \frac{1}{\nu} \sum_{j=3}^{\infty} e^{-\nu} \frac{\nu^{(j+1)}}{(j+1)!} \left(\sqrt{j+1} + \frac{3}{\sqrt{j+1}} \right) \\ &\leq b_1 \sqrt{2} q \frac{1}{\nu} \left(\sqrt{\nu} + \frac{3}{\sqrt{\nu}} \right), \end{aligned}$$

where in the final line we have used that the function $f(x)$ is increasing in x for any $x \geq 3$, and that $\mathbb{E}[Y_T \mathbf{1}_{Y_T \geq 4}] \geq \nu - 3 \geq 3$. Thus, we have that

$$\begin{aligned} \Pr(\gamma_c \geq t) &\leq (3\nu) e^{-\nu} + 2q \sqrt{\frac{2}{\pi}} \left(\frac{1}{\sqrt{\nu}} + \frac{1}{\nu^{3/2}} \right) \\ &\leq \frac{2.5}{\nu^{3/2}} + 2q \sqrt{\frac{2}{\pi}} \left(\frac{1}{\sqrt{\nu}} + \frac{1}{\nu^{3/2}} \right) \end{aligned}$$

Proof of Lower Bound.

We approach the initial stages of the proof in the precisely the same way, obtaining

$$\begin{aligned} \Pr(\gamma_c \geq t) &= \Pr(\gamma \geq Y_T) \\ &= e^{-\nu} + 2p\nu e^{-\nu} + \sum_{j=3}^{\infty} e^{-\nu} \frac{\nu^j}{j!} \Pr\left(\gamma \geq 2 \left(\frac{j + \mathbf{1}_{j \text{ is even}} - 1}{2} \right) + 1\right) \\ &\geq \sum_{j=3}^{\infty} e^{-\nu} \frac{\nu^j}{j!} b_1 e^{-\frac{1}{3(j + \mathbf{1}_{j \text{ is even}} - 1)}} \frac{q\sqrt{2}}{\sqrt{(j + \mathbf{1}_{j \text{ is even}} + 1)}} \\ &\geq \sum_{j=3}^{\infty} e^{-\nu} \frac{\nu^j}{j!} b_1 e^{-\frac{1}{3(j-1)}} \frac{q\sqrt{2}}{\sqrt{(j+2)}}. \end{aligned}$$

Applying Jensen's inequality, we obtain

$$\geq b_1 q \sqrt{2} e^{-\frac{1}{3(\nu-1)}} \frac{1}{\sqrt{\nu+2}}.$$

□

A.3.3 Proof of Claim A.7, Bound on Expected Length of Stopped Random Walk.

Claim A.7 (Bound on Expected Length of Stopped Random Walk). *Suppose we have a critically loaded M/M/1 queue with arrival rate and departure rate both equal to $k\lambda$, with offered load*

$R > 100$ and setup time $\beta > 100$. Suppose also that at time 0, a job arrives. Let τ be the length of the busy period which follows. Then,

$$\mathbb{E} [\min (\beta, \tau)] \leq b_1 \frac{\sqrt{\beta}}{\sqrt{\mu R}} + \frac{6}{\mu R}.$$

Proof.

From Claim A.6, the continuous-time random walk hitting time bound, we have that

$$\Pr (\tau \geq t) \leq \frac{b_1}{\sqrt{2}} \left(\frac{1}{\sqrt{\nu}} + \frac{b_2}{\nu^{3/2}} \right), \quad (\text{A.11})$$

where $\nu = 2\mu R t$ and we require that $\nu \geq 3$. By integrating this bound (using a bound of 1 wherever this bound doesn't apply), we obtain

$$\begin{aligned} \mathbb{E} [\min (\beta, \tau)] &= \int_0^\beta \Pr (\tau > t) dt \leq \frac{3}{2\mu R} + \int_{\frac{3}{2\mu R}}^\beta \frac{b_1}{\sqrt{2}} \left(\frac{1}{\sqrt{2\mu R t}} + \frac{b_2}{(2\mu R)^{3/2}} \right) dt \\ &\leq \frac{3}{2\mu R} + b_1 \frac{\sqrt{\beta}}{\sqrt{\mu R}} + \frac{b_1 b_2}{4\mu R} \left[2\sqrt{\frac{2}{3}} \right] \leq b_1 \frac{\sqrt{\beta}}{\sqrt{\mu R}} + \frac{6}{\mu R}. \square \end{aligned}$$

A.3.4 Proof of Claim A.8, Bound on the Expected Hitting Time in the M/M/ ∞ .

Claim A.8 (M/M/ ∞ Passage Time Bound). *Given an M/M/ ∞ queue, let $T_{x \rightarrow y}$ denote the random amount of time taken to go from state x to state y . Suppose this system has an arrival rate of $k\lambda$ and a per-server departure rate of μ . Let $R \triangleq k\frac{\lambda}{\mu}$. Then, for any h such that $1 \leq h \leq \sqrt{R}$,*

$$\mathbb{E} [T_{(R+h-1) \rightarrow (R+h)}] \leq \frac{\sqrt{2\pi}}{\mu\sqrt{R}} \left(1 + \frac{h}{R} \right)^{h-\frac{1}{2}} e^{\frac{1}{12R}} \leq D_2 \frac{\sqrt{\pi}}{\mu\sqrt{R}}.$$

Proof.

The proof here is quite simple. First, we note that the passage time in the M/M/ ∞ from state $(R + h - 1)$ to state $(R + h)$ is exactly the passage time from those states in the M/M/ $(R +$

$h)/(R + h)$. This new system has a nice product form, so that

$$\begin{aligned}
\mathbb{E} [T_{(R+h-1) \rightarrow (R+h)}] &\leq \mathbb{E} [T_{(R+h) \rightarrow (R+h)}] \\
&= \frac{1}{\mu(R+h)} \frac{1}{\pi_{R+h}} \\
&= \frac{1}{\mu(R+h)} \frac{\sum_{i=0}^{R+h} \frac{R^i}{i!}}{\frac{R^{R+h}}{(R+h)!}} \\
&\leq \frac{1}{\mu(R+h)} e^R \frac{(R+h)!}{R^{R+h}} \\
&\leq \frac{1}{\mu(R+h)} e^R \frac{e^{\frac{1}{12(R+h)}} \sqrt{2\pi(R+h)} (R+h)^{R+h} e^{-(R+h)}}{R^{R+h}} \\
&\leq e^{\frac{1}{12R}} \frac{1}{\mu} \frac{1}{\mu \sqrt{R+h}} \sqrt{2\pi} \left(1 + \frac{h}{R}\right)^{R+h} e^{-h} \\
&\leq \frac{\sqrt{2\pi}}{\mu \sqrt{R}} \left(1 + \frac{h}{R}\right)^{h-\frac{1}{2}} e^{\frac{1}{12R}} \\
&\leq \frac{1}{\mu} \frac{\sqrt{2\pi}}{\sqrt{R}} e^{\frac{h^2}{R}} e^{\frac{1}{12R}} \\
&\leq \frac{7}{\mu \sqrt{R}},
\end{aligned}$$

where we have made extensive use of Stirling's approximation and the bound $(1+x) \leq e^x$.

A.4 Helper Claims.

A.4.1 Proof of Claim A.9, the Busy Period Integral Bound.

Claim A.9 (Busy Period Integral Bound). *Suppose that, at time τ , we can guarantee that $N(\tau) \geq Z(\tau) \geq R+j$. Let $\eta_i \triangleq \min \{t > 0 : N(t) \leq R+i\}$, for $i \in \{j, j+1, \dots, [N(\tau) - R]\}$. Then,*

$$\mathbb{E} \left[\int_{\tau}^{\eta_j} [N(t) - R] dt \middle| \mathcal{F}_{\tau} \right] \leq (N(\tau) - (R+j)) \left[\frac{3}{2\mu j} + \frac{1}{\mu} + \frac{R}{\mu j^2} \right] + \frac{(N(\tau) - (R+j))^2}{2\mu j} \triangleq I^{\text{busy}}([N(\tau) - R]),$$

Proof. We prove this claim via an appeal to conventional M/M/1 busy period analysis. In particular, we first note that

$$\int_{\tau}^{\eta_j} [N(t) - R] dt = \sum_{i=j+1}^{N(\tau)-R} \int_{\eta_i}^{\eta_{i-1}} [N(t) - R] dt,$$

meaning we need only bound the integrals between the η_i 's. To bound that process, we define a coupled process $\tilde{N}(t)$ and bound the integrals over that process.

To do so, note that, until time η_j , the number of busy servers $Z(t) \geq R + j$. By Claim A.1, we can define, for each index i , the i -th coupled process $\tilde{N}_i(t)$ as

$$\tilde{N}_i(t) = N(\eta_{i+1}) + A(\eta_{i+1}, t) - \mathcal{D}[R + j](\eta_{i+1}, t),$$

and have $\tilde{N}(t) \geq N(t)$ on the interval $[\eta_{i+1}, \eta_i]$. Furthermore, we can extend our integral of interest from the interval $[\eta_{i+1}, \eta_i)$ to the interval $[\eta_{i+1}, \tilde{\eta}_i)$, where $\tilde{\eta}_i \triangleq \min\{t > 0 : N(t) \leq R + i\}$. Now, we note that

$$\mathbb{E} \left[\int_{\eta_{i+1}}^{\tilde{\eta}_i} [\tilde{N}_i(t) - R] dt \middle| \mathcal{F}_\tau \right] = \mathbb{E} \left[\int_{\eta_{i+1}}^{\tilde{\eta}_i} [\tilde{N}_i(t) - (R + i)] dt \middle| \mathcal{F}_\tau \right] + i \mathbb{E}[\eta_{i+1} - \tilde{\eta}_i | \mathcal{F}_\tau].$$

The first term on the right is simply the expected time integral of the number of jobs in an M/M/1 queue over a busy period, with arrival rate $k\lambda$ and departure rate $\mu(R + j)$. The second term is simply the quantity i multiplied by the expected length of that M/M/1 busy period. Let $\rho_j = \frac{k\lambda}{\mu(R+j)}$. Then, from standard results on the M/M/1 busy period,

$$\mathbb{E} \left[\int_{\eta_{i+1}}^{\tilde{\eta}_i} [\tilde{N}_i(t) - (R + i)] dt \middle| \mathcal{F}_\tau \right] = \frac{1}{\mu j} \left[\frac{1}{1 - \rho_j} \right] = \frac{1}{\mu j} \left[\frac{R}{j} + 1 \right] = \frac{1}{\mu j} + \frac{R}{\mu j^2}.$$

Summing over all values of i , we obtain

$$\begin{aligned} & \mathbb{E} \left[\int_{\tau}^{\eta_j} [N(t) - R] dt \middle| \mathcal{F}_\tau \right] \\ & \leq \sum_{i=j+1}^{N(\tau)-R} \mathbb{E} \left[\int_{\eta_i}^{\tilde{\eta}_{i-1}} [\tilde{N}_i(t) - R] dt \middle| \mathcal{F}_\tau \right] \\ & = \sum_{i=j+1}^{N(\tau)-R} \left[\frac{1}{\mu j} + \frac{R}{\mu j^2} \right] + i \frac{1}{\mu j} \\ & = (N(\tau) - (R + j)) \left[\frac{1}{\mu j} + \frac{R}{\mu j^2} \right] + (N(\tau) - (R + j)) \frac{1}{\mu} + \frac{1}{\mu j} \left[\frac{(N(\tau) - (R + j))(N(\tau) - (R + j) + 1)}{2} \right] \\ & = (N(\tau) - (R + j)) \left[\frac{3}{2\mu j} + \frac{1}{\mu} + \frac{R}{\mu j^2} \right] + \frac{(N(\tau) - (R + j))^2}{2\mu j}, \end{aligned}$$

as desired. \square \square

A.4.2 Proof of Claim 6.10, the Wait Busy Claim.

Claim 6.10 (Wait Busy Claim). *Let τ be some stopping time, let the number of jobs $N(\tau) = R + h$, and define $\text{ns}(h) \triangleq \min\{h, k(1 - \rho)\}$. Let the down-crossing $d_{\text{gen}} \triangleq \min\{t > 0 : N(\tau + t) = R + h - 1\}$. If $Z(\tau) \geq R$, then*

$$\mathbb{E} \left[\int_{\tau}^{\tau + d_{\text{gen}}} [N(t) - (R + h - 1)] dt \middle| \mathcal{F}_\tau \right] \leq Y_{R+\text{ns}(h)}(\tau) + g \left(1 + 2\mu R \mathbb{E} \left[\min(Y_{R+\text{ns}(h)}(\tau), d_{\text{gen}}) \right], 1, \mu \text{ns}(h) \right), \quad (6.7)$$

where the function $g(x, y, z) \triangleq x \frac{1}{2\mu z} + y \left[\frac{R}{\mu z^2} + \frac{3}{2\mu z} \right]$.

Furthermore,

$$\begin{aligned} \mathbb{E} \left[\int_{T_A}^{v_1^{(\text{down})}} [N(t) - R] dt \middle| \mathcal{F}_{T_A} \right] &\leq \left[\beta + \frac{1}{\mu} \right] \left[\mathbb{E} [N(T_A) - R] + \frac{R}{M_B} \right] + \frac{2}{\mu} \ln \left(\frac{\mathbb{E} [N(T_A) - R]}{M_B} \right) \\ &+ g \left(\mathbb{E} [[N(T_A) - R]^2] + 2\mathbb{E} [N(T_A) - R], \mathbb{E} [N(T_A) - R], k(1 - \rho) \right). \end{aligned} \quad (6.8)$$

Proof of (6.7).

We prove the two parts of Claim 6.10 separately; we first show (6.7) by applying coupling, martingales, and busy period analysis. First, note that, if the down-crossing at $\tau + d_{\text{gen}}$ does not occur by time $Y_{R+\text{ns}(h)}(\tau)$, then the system must have at least $(R + \text{ns}(h))$ servers at its disposal afterwards. (Note that the $\text{ns}(\cdot)$ function here is just to account for the case where you have more jobs than servers.) Accordingly, we split our analysis into two parts.

First portion. For the first portion, since the number of busy servers $Z(t) \geq R$, by coupling our system to a critically-loaded M/M/1 using Claim A.2, we have

$$\mathbb{E} \left[\int_{\tau}^{\tau + \min(d_{\text{gen}}, Y_{R+\text{ns}(h)}(\tau))} [N(t) - (R + h - 1)] dt \middle| \mathcal{F}_{\tau} \right] \leq Y_{R+\text{ns}(h)}(\tau).$$

Second portion. For the second portion, since, at that point the number of busy servers $Z(t) \geq R + \text{ns}(h)$, we can apply a stronger bound. In particular, at time $(\tau + Y_{R+\text{ns}(h)}(\tau))$, we can couple to an accordingly-stronger M/M/1 with the same number of jobs. From basic busy period analysis and Claim A.1, this tells us that, letting the adjusted number of jobs $N_{\text{adj}}(t) \triangleq N(\tau + t) - (R + h - 1)$ and the important remaining setup time $Y_{\text{imp}} \triangleq Y_{R+\text{ns}(h)}(\tau)$ as a shorthand,

$$\mathbb{E} \left[\mathbf{1}_{Y_{\text{imp}} < d_{\text{gen}}} \int_{Y_{\text{imp}}}^{d_{\text{gen}}} [N_{\text{adj}}(t)] dt \middle| \mathcal{F}_{\tau + Y_{\text{imp}}} \right] \leq \mathbf{1}_{Y_{\text{imp}} < d_{\text{gen}}} g(N_{\text{adj}}(Y_{\text{imp}})^2, N_{\text{adj}}(Y_{\text{imp}}), \text{ns}(h)),$$

since the function $g(x^2, x, z)$ describes the integral of the number of jobs over a busy period started by x jobs in an M/M/1 with arrival rate $k\lambda$ and departure rate $\mu(R + z)$. Note that $\mathbf{1}_{Y_{\text{imp}} < d_{\text{gen}}} [N(\tau + Y_{\text{imp}}) - (R + h - 1)] = [N(\tau + \min(d_{\text{gen}}, Y_{\text{imp}})) - (R + h - 1)]$, since $N(\tau + d_{\text{gen}}) = R + h - 1$ by definition. Coupling our system to the critically-loaded M/M/1 \tilde{N} system as before, then taking expectations and applying Doob's Optional Stopping Theorem, we obtain, as desired,

$$\mathbb{E} \left[\mathbf{1}_{Y_{R+\text{ns}(h)}(\tau) < d_{\text{gen}}} \int_{Y_{R+\text{ns}(h)}(\tau)}^{d_{\text{gen}}} N_{\text{adj}}(t) dt \middle| \mathcal{F}_{\tau + Y_{R+\text{ns}(h)}(\tau)} \right] \leq g(1 + 2\mu R \mathbb{E} [\min(Y_{R+\text{ns}(h)}(\tau), d_{\text{gen}})], 1, \text{ns}(h)).$$

Combining these two terms, we obtain (6.7).

Proof of (6.8)

. We now apply (6.7) to prove (6.8).

Decomposition in terms of η_i 's. We first fix the filtration/state at the accumulation time T_A , then define the hitting times $\eta_i \triangleq \min \{t > T_A : N(t) \leq R + i\}$ for a set number of jobs $i \in [M_B, N(T_A)]$; we accordingly omit the filtration at time T_A in our expectations. Note that the down-crossing $v_1^{(\text{down})} = \eta_{M_B}$, by this definition. From there, it's clear that

$$\int_{T_A}^{v_1^{(\text{down})}} [N(t) - R] dt = \sum_{i=M_B}^{N(T_A)-R-1} \int_{\eta_{i+1}}^{\eta_i} [N(t) - R] dt.$$

For each of these terms, we can separate $[N(t) - R] = [N(t) - (R + i)] + i$ and apply (6.7) to find

$$\mathbb{E} \left[\int_{\eta_{i+1}}^{\eta_i} [N(t) - R] dt \right] \leq \mathbb{E} [Y_{R+\text{ns}(h)}(\eta_{i+1})] + i \mathbb{E} [\min(Y_{R+\text{ns}(i+1)}(\eta_{i+1}), \eta_i - \eta_{i+1})] \quad (\text{A.12})$$

$$+ g \left([1 + 2\mu R \mathbb{E} [\min(Y_{R+\text{ns}(i+1)}(\eta_{i+1}), \eta_i - \eta_{i+1})]] , 1, \mu \text{ns}(i+1) \right) \quad (\text{A.13})$$

$$+ i \cdot \frac{1}{\mu \text{ns}(i+1)}. \quad (\text{A.14})$$

We analyze each of these terms separately.

Bound on (A.12), the remaining setup time portion. From here, it suffices to note that the sum

$$\sum_{i=1}^{N(T_A)-R} i \min(Y_{R+\text{ns}(h)}(\tau), \eta_i - \eta_{i+1}) + \sum_{i=1}^{N(T_A)-R} \mathbb{E} [Y_{R+\text{ns}(i+1)}(\eta_{i+1})] \leq \beta [N(T_A) - R]; \quad (\text{A.15})$$

actually, the statement is true without expectations. To see this, we make an interchange of summation argument. First note that, if the $(R+i)$ -th server becomes busy, then, by the monotonicity of server states, all servers of index smaller than $(R+i)$ must also be busy; in other words, the remaining setup time $Y_{R+\text{ns}(i+1)}(\eta_{i+1}) = 0$ for all $i < s$. To use this, we let s be the largest index for which $Y_{R+\text{ns}(s)}(\eta_s) < (\eta_{s-1} - \eta_s)$. Note also that η_i is the first time after T_A that the number of jobs $N(t) \leq R + i$ (and so must have been continuously decreasing from time T_A); it follows that the remaining setup time $Y_{R+\text{ns}(i)}(\eta_i) \leq [\beta - (\eta_i - T_A)]^+$. Breaking things down further,

$$\sum_{i=s}^{N(T_A)-R} i \min(Y_{R+\text{ns}(i+1)}(\eta_{i+1}), (\eta_i - \eta_{i+1})) = \sum_{i=s}^{N(T_A)-R} i (\eta_i - \eta_{i+1}) + s Y_{R+\text{ns}(s)}(\eta_s).$$

From here, by an interchange of summation argument,

$$\sum_{i=s}^{N(T_A)-R} i (\eta_i - \eta_{i+1}) = \sum_{i=s}^{N(T_A)-R} \sum_{j=1}^i (\eta_i - \eta_{i+1}) = \sum_{j=s}^{N(T_A)-R} \sum_{i=j}^{N(T_A)-R} (\eta_i - \eta_{i+1}) = \sum_{j=1}^{N(T_A)-R} \eta_s - T_A,$$

where, by definition, the (relative to T_A) hitting time $(\eta_s - T_A) \leq \beta - Y_{R+\text{ns}(s)}(\eta_s)$; using this,(A.15) follows.

Bound on (A.13), the busy period integral portion. Applying similar reason to the sum of the first terms in g , and using the independence of g in its first and second arguments (i.e. that $g(x, y, z) = f_1(x, z) + f_2(y, z)$ for two functions f_1 and f_2 linear in their first argument),

$$\sum_{i=M_B}^{N(T_A)-R} \text{(A.13)} \leq g(2\mu R\beta, 0, M_B) + g([N(T_A) - k]^+, [N(T_A) - k]^+, k(1 - \rho)) + \sum_{i=M_B}^{\min(N(T_A)-R-1, k(1-\rho))} g(1, 1, i)$$

This last term above can be bounded by replacing it with an integral, which gives

$$\begin{aligned} \sum_{i=M_B}^{\min(N(T_A)-R, k(1-\rho))} g(1, 1, i) &\leq \int_{M_B}^{\min(N(T_A)-R, k(1-\rho))} \frac{2}{\mu i} + \frac{R}{\mu i^2} \mathbf{d}i \\ &\leq \frac{2}{\mu} \ln \left(\frac{\min(N(T_A) - R, k(1 - \rho))}{M_B} \right) + \frac{R}{\mu} \left[\frac{1}{M_B} \right] \\ &\leq \frac{2}{\mu} \ln \left(\frac{N(T_A) - R}{M_B} \right) + \frac{R}{\mu M_B}. \end{aligned}$$

Bound on (A.14), the busy period length portion. Using the definition of the function g , we also find that $\sum_{i=M_B}^{N(T_A)-R} \frac{i}{\mu \text{ns}(i)} \leq g\left(\left([N(T_A) - k]^+\right)^2 + [N(T_A) - k]^+, 0, k(1 - \rho)\right) + \frac{1}{\mu} \min(k(1 - \rho), N(T_A) - R, \cdot)$

Combining the terms to bound (6.8), the integral from T_A to $v_1^{(\text{down})}$. Combining terms, noting that $N(T_A) - k \geq N(T_A) - R \geq 0$, and applying Jensen's inequality to the $\ln(\cdot)$ term, we find that

$$\begin{aligned} (6.8) &\leq \beta \mathbb{E}[N(T_A) - R] + g(2\mu R\beta, 0, M_B) + \mathbb{E}\left[g\left([N(T_A) - k]^+, [N(T_A) - k]^+, k(1 - \rho)\right)\right] \\ &\quad + \mathbb{E}\left[\frac{2}{\mu} \ln \left(\frac{N(T_A) - R}{M_B} \right)\right] + \frac{1}{\mu} \frac{R}{M_B} \\ &\quad + \mathbb{E}\left[g\left(\left([N(T_A) - k]^+\right)^2 + [N(T_A) - k]^+, 0, k(1 - \rho)\right)\right] + \frac{1}{\mu} \mathbb{E}[\min(k(1 - \rho), N(T_A) - R)] \\ &\leq \left[\beta + \frac{1}{\mu}\right] \left[\mathbb{E}[N(T_A) - R] + \frac{R}{M_B}\right] + \frac{2}{\mu} \ln \left(\frac{\mathbb{E}[N(T_A) - R]}{M_B} \right) \\ &\quad + g\left(\mathbb{E}\left[[N(T_A) - R]^2\right] + 2\mathbb{E}[N(T_A) - R], \mathbb{E}[N(T_A) - R], k(1 - \rho)\right). \square \end{aligned}$$

A.4.3 Proof of Claim 6.3

We now prove Claim 6.3, restated here.

Claim 6.3 (Bound on the Probability of an Up-crossing $p_{\text{rise}}^{(j)}$). *Let $p_{\text{rise}}^{(j)}$ be the probability that the total number of jobs $N(t)$ exceeds $R + C_3\sqrt{\mu\beta R}$ during epoch j defined in (6.1). Then, for any epoch $j \geq A_5\sqrt{R}$, we have $p_{\text{rise}}^{(j)} \geq 0.99\frac{A_5}{\sqrt{R}}$.*

We show a more general claim: that, for $j \geq A_5\sqrt{R}$,

$$p_{\text{rise}}^{(j)} \geq 0.99\frac{j}{R}. \quad (\text{A.16})$$

Proof of (A.16): Lower Bound on $p_{\text{rise}}^{(j)}$.

We begin with a simple probability manipulation:

$$\begin{aligned} p_{\text{rise}}^{(j)} &\triangleq \Pr\left(N(t) \geq C_3\sqrt{\mu\beta R} \text{ at some point during epoch } j \middle| \mathcal{F}_{\tau_j}\right) \\ &\geq \Pr\left(N(t) \geq C_3\sqrt{\mu\beta R} \text{ during the interval } [\tau_j, \min(\tau_j + \beta, \tau_{j+1})] \middle| \mathcal{F}_{\tau_j}\right). \end{aligned}$$

From here, we make with a useful observation: since there are no server *in setup* at the beginning of an epoch (as we have just turned off a server), no servers can complete setup in the first β time of an epoch. Thus, the number of busy servers $Z(t) \leq R - j$ during this time, and, by Claim A.1, the coupled process

$$\tilde{N}(t) \triangleq N(\tau_j) + A(\tau_j, t) - \mathcal{D}[R - j](\tau_j, t)$$

must be a lower bound on $N(t)$, during the interval $[\tau_j, \tau_j + \beta]$. Moreover, the number of busy servers $Z(t)$ can not be smaller than $R - j$ until the beginning of epoch $j + 1$ either. Thus, we find that the behavior of $N(t)$ corresponds *exactly* with the behavior of $\tilde{N}(t)$ during the interval $[\tau_j, \min(\tau_{j+1}, \tau_j + \beta)]$.

We now use this coupled process to analyze our original probability. Define the up-crossing time τ_{up} as

$$\tau_{\text{up}} \triangleq \min\left\{t > 0 : \tilde{N}(\tau_j + t) \geq R + C_3\sqrt{\mu\beta R}\right\}.$$

Likewise, define the down-crossing time τ_{down} as

$$\tau_{\text{down}} \triangleq \min\left\{t > 0 : \tilde{N}(\tau_j + t) \leq R - (j + 1)\right\}.$$

It follows that

$$\begin{aligned} &\Pr\left(\text{reach } N(t) \leq R - (j + 1) \text{ during the interval } [\tau_j, \min(\tau_j + \beta, \tau_{j+1})] \middle| \mathcal{F}_{\tau_j}\right) \\ &= \Pr\left(\text{reach } \tilde{N}(t - \tau_j) \leq R - (j + 1) \text{ during the interval } [\tau_j, \min(\tau_j + \beta, \tau_{j+1})] \middle| \mathcal{F}_{\tau_j}\right) \\ &= \Pr(\tau_{\text{up}} \leq \beta, \tau_{\text{up}} < \tau_{\text{down}}) \\ &= \Pr(\tau_{\text{up}} \leq \beta) - \Pr(\tau_{\text{up}} \leq \beta, \tau_{\text{up}} \geq \tau_{\text{down}}) \\ &= \Pr(\tau_{\text{up}} \leq \beta) - \Pr(\tau_{\text{up}} \leq \beta | \tau_{\text{up}} \geq \tau_{\text{down}}) \Pr(\tau_{\text{up}} \geq \tau_{\text{down}}). \end{aligned}$$

We now observe that

$$\Pr(\tau_{\text{up}} \leq \beta | \tau_{\text{up}} \geq \tau_{\text{down}}) \leq \Pr(\tau_{\text{up}} \leq \beta), \quad (\text{A.17})$$

since the process has farther to go, less time to do so, and the process's behavior is translation-invariant (this last point is why we needed to analyze the coupled process instead).

Continuing from where we left off, we find that

$$\begin{aligned}
p_{\text{rise}}^{(j)} &= \Pr(\tau_{\text{up}} \leq \beta) - \Pr(\tau_{\text{up}} \leq \beta | \tau_{\text{up}} \geq \tau_{\text{down}}) \Pr(\tau_{\text{up}} \geq \tau_{\text{down}}) \\
&\geq \Pr(\tau_{\text{up}} \leq \beta) - \Pr(\tau_{\text{up}} \leq \beta) \Pr(\tau_{\text{up}} \geq \tau_{\text{down}}) \\
&= \Pr(\tau_{\text{down}} > \tau_{\text{up}}) \Pr(\tau_{\text{up}} \leq \beta) \\
&\geq \Pr(\tau_{\text{down}} > \infty) \Pr(\tau_{\text{up}} \leq \beta) \\
&= \frac{j}{R} \Pr(\tau_{\text{up}} \leq \beta),
\end{aligned}$$

where the last equality is a classical result on upwards-biased discrete random walks (one can think of \tilde{N} as a discrete random walk driven by a Poisson process of rate $(k\lambda + \mu(R - j))$, where the probability that \tilde{N} increases at a Poisson event is $\frac{k\lambda}{k\lambda + \mu(R - j)} = \frac{R}{2R - j}$).

From here, it suffices to lower bound $\Pr(\tau_{\text{up}} \leq \beta)$. To begin, note

$$\begin{aligned}
\Pr(\tau_{\text{up}} \leq \beta) &= \Pr\left(\sup_{t \in [0, \beta]} \tilde{N}(t) \geq R + C_3 \sqrt{\mu\beta R}\right) \\
&\geq \Pr\left(\tilde{N}(\beta) \geq R + C_3 \sqrt{\mu\beta R}\right) \\
&= \Pr\left(A(\tau_j, \tau_j + \beta) - \mathcal{D}[R - j](\tau_j, \tau_j + \beta) \geq j + C_3 \sqrt{\mu\beta R}\right).
\end{aligned}$$

Noting that the number of arrivals $A(\tau_j, \tau_j + \beta)$ and the number of departures $\mathcal{D}[R - j](\tau_j, \tau_j + \beta)$ are independent Poisson r.v.'s, we can apply the Berry-Esseen bound of Claim A.10 to find

$$\begin{aligned}
&= 1 - \Phi\left(\frac{\mu\beta j - j - C_3 \sqrt{\mu\beta R}}{\sqrt{\mu\beta(2R - j)}}\right) - \frac{1}{3\sqrt{\mu\beta(2R - j)}} \\
&\geq 1 - \Phi\left(\frac{0.99\mu\beta j - C_3 \sqrt{\mu\beta R}}{\sqrt{2\mu\beta R}}\right) - \frac{1}{3\sqrt{\mu\beta R}} \\
&= 1 - \Phi\left(-0.99\frac{j}{\sqrt{R}}\sqrt{\mu\beta} + \frac{C_3}{\sqrt{2}}\right) - \frac{1}{3\sqrt{\mu\beta R}} \\
&\geq 1 - \Phi\left(-9.9A_5 + \frac{C_3}{\sqrt{2}}\right) - \frac{1}{300}.
\end{aligned}$$

To complete the proof, we set the constant A_5 such that the final probability is ≥ 0.99 . In particular, we need

$$\Phi\left(-9.9A_5 + \frac{C_3}{\sqrt{2}}\right) \leq \frac{2}{300},$$

which is achieved when $A_5 > \frac{C_3}{9.9\sqrt{2}} + 0.25$; choosing $A_5 = 1$ gives the result. \square

A.4.4 Proof of Claim A.10.

Claim A.10 (Berry-Esseen bound for the Skellam distribution). *Given two independent random variables $Y_1 \sim \text{Poisson}(\mu_1)$ and $Y_2 \sim \text{Poisson}(\mu_2)$, as well as a constant C with $\mu_1 > \mu_2 + C$, one has*

$$\Pr(Y_1 - Y_2 \geq C) \geq 1 - \Phi\left(-\left[\frac{\mu_1 - \mu_2 - C}{\mu_1 + \mu_2}\right]\right) - \frac{1}{3\sqrt{\mu_1 + \mu_2}}.$$

A.4.5 Proof.

This follows directly from the Poisson Berry-Esseen bound of [5], applied twice; first approximating Y_1 then approximating Y_2 . \square

A.4.6 Proof of (5.3): Lower Bound on $\mathbb{E}[L]$, Expected Value of First Long Epoch Index.

We prove this result by first showing that

$$\Pr(L > j | L \geq j) \geq \left(1 - \frac{j}{R}\right) \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right), \quad (\text{A.18})$$

where $b_1 = \frac{2}{\sqrt{\pi}}$. Next, we show that this implies that for any $\delta \in (0, 1)$ and any $j < \delta R$,

$$\Pr(L > j) \geq \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right)^{j+1} e^{-\frac{j(j+1)}{2R} \frac{1}{1-\delta}}. \quad (\text{A.19})$$

From here, we use the sum of tails formula $\mathbb{E}[L] = \sum_{j=0}^{\infty} \Pr(L > j)$ to show

$$\mathbb{E}[L] \geq \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right) \left(\left[\sqrt{\frac{\pi}{2}}(1-\delta) - \frac{1.15(1-\delta)}{\sqrt{\mu\beta}}\right] \sqrt{R} - \frac{1}{2} - \frac{2(1-\delta)}{\delta} e^{-R \frac{\delta^2}{1-\delta}}\right).$$

Choosing $\delta = \frac{2}{\sqrt{R}}$ then noting that $\mu\beta \geq 100$ and $R \geq 100$ gives the result.

Proof of (A.18): Lower Bound on Probability that Current Epoch is Short.

Recall that an epoch j is *long* if $\tau_{j+1} - \tau_j > \beta$, that L is the index of the first long epoch, and that, if $L \geq j$, then we learn that $L \geq j$ precisely at time τ_j , i.e. when epoch j begins. Moreover, since the system is Markovian, the behavior of the system from τ_j onwards is completely independent of what happened previously. Thus,

$$\Pr(L > j | L \geq j) = \Pr(L > j | \mathcal{F}_{\tau_j}, L \geq j) = \Pr(\tau_{j+1} - \tau_j \leq \beta | \mathcal{F}_{\tau_j}, L \geq j) = \Pr(\tau_{j+1} - \tau_j \leq \beta).$$

From here, we note that the random time $\tau_{j+1} - \tau_j$ is a stopping time; a hitting time, to be exact. Moreover, since the number of servers $Z(t)$ can not increase before time $\tau_j + \beta$ and can not decrease until τ_{j+1} , we have that the coupled process $\tilde{N}(t)$ defined as

$$\tilde{N}(t - \tau_j) \triangleq 1 + A(\tau_j, t) - \mathcal{D}[R - j](t, \tau_j)$$

is in correspondence with $N(t)$; in particular,

$$N(t) = \tilde{N}(t - \tau_j) + R - j - 1$$

for any time $t \in [\tau_j, \min(\tau_j + \beta, \tau_{j+1})]$. If we define the coupled hitting time $\gamma_c \triangleq \min\{t > 0 : \tilde{N}(t) \leq 0\}$, then we also have that the hitting time $\gamma_c = \tau_{j+1} - \tau_j$, whenever the event $\{\tau_{j+1} - \tau_j \leq \beta\}$ occurs. From here, we can apply Claim A.6 to find that

$$\Pr(\gamma_c \leq \beta) \geq \left(1 - \frac{j}{R}\right) \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right), \text{ as desired.}$$

Proof of (A.19).

Having shown the above bound on the conditional extension of the tail, we note that, for $j \leq \delta R$,

$$\begin{aligned} \Pr(L \geq j + 1) &= \Pr(L \geq j + 1 | L \geq j) \Pr(L \geq j | L \geq j - 1) \cdots \Pr(L \geq 1) \\ &\geq \prod_{i=0}^j \left(1 - \frac{i}{R}\right) \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right) \\ &\geq \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right)^{j+1} \prod_{i=0}^j e^{-\frac{i}{R-i}} \\ &\geq \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right)^{j+1} e^{-\sum_{i=0}^j \frac{i}{R-i}} \\ &= \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right)^{j+1} e^{-\sum_{i=0}^j \frac{i}{R-j}} \\ &= \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right)^{j+1} e^{-\frac{j(j+1)}{2R} \frac{1}{1-\delta}}, \end{aligned}$$

as desired. □

Proof of (A.4.6): Final Bound on $\mathbb{E}[L]$ using Gaussian Integral.

We now complete the proof. Let $a \triangleq \frac{1}{2R} \frac{1}{1-\delta}$ and $\psi \triangleq -\ln\left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right)$ as a shorthand. Then we can rewrite (A.19) as

$$\Pr(L \geq j + 1) \geq e^{-aj^2 - (\psi+a)j - \psi}.$$

Now, using the sum-of-tails formula for expectations, we find that

$$\begin{aligned}
\mathbb{E}[L] &= \sum_{j=0}^{R-1} \Pr(L \geq j+1) \\
&\geq \sum_{j=0}^{\delta R-1} \Pr(L \geq j+1) \\
&\geq \sum_{j=0}^{\delta R-1} e^{-aj^2 - (\psi+a)j - \psi} \\
&\geq \int_0^{\delta R} e^{-aj^2 - (\psi+a)j - \psi} \mathbf{d}j \\
&= \int_0^{\delta R} e^{-a(j^2 + (\frac{\psi}{a}+1)j) - \psi} \mathbf{d}j \\
&= \int_0^{\delta R} e^{-a(j + \frac{1}{2}(\frac{\psi}{a}+1))^2 + \frac{a}{4}(\frac{\psi}{a}+1)^2 - \psi} \mathbf{d}j \\
&= e^{\frac{a}{4}(\frac{\psi}{a}+1)^2 - \psi} \int_0^{\delta R} e^{-a(j + \frac{1}{2}(\frac{\psi}{a}+1))^2} \mathbf{d}j.
\end{aligned}$$

Evaluating the integral further, we find that

$$\begin{aligned}
\int_0^{\delta R} e^{-a(j + \frac{1}{2}(\frac{\psi}{a}+1))^2} \mathbf{d}j &= \int_{\frac{1}{2}(\frac{\psi}{a}+1)}^{\delta R + \frac{1}{2}(\frac{\psi}{a}+1)} e^{-aj^2} \mathbf{d}j \\
&= \int_0^{\infty} e^{-aj^2} \mathbf{d}j - \int_0^{\frac{1}{2}(\frac{\psi}{a}+1)} e^{-aj^2} \mathbf{d}j - \int_{\delta R + \frac{1}{2}(\frac{\psi}{a}+1)}^{\infty} e^{-aj^2} \mathbf{d}j.
\end{aligned}$$

We now bound each of these integrals in turn. First, we know classically that

$$\int_0^{\infty} e^{-aj^2} \mathbf{d}j = \frac{1}{2} \sqrt{\frac{\pi}{a}} = \sqrt{\frac{\pi}{2}} \cdot \sqrt{1-\delta} \sqrt{R} \geq \sqrt{\frac{\pi}{2}} \cdot (1-\delta) \sqrt{R}$$

Next, we note that, since the integrand is ≤ 1 ,

$$\begin{aligned}
\int_0^{\frac{1}{2}(\frac{\psi}{a}+1)} e^{-aj^2} \mathbf{d}j &\leq \frac{1}{2} \left(\frac{\psi}{a} + 1 \right) \\
&= \frac{1}{2} \left(2R(1-\delta) \ln \left(\frac{1}{1 - \frac{b_1}{\sqrt{\mu\beta R}}} \right) \right) + \frac{1}{2} \\
&\leq R(1-\delta) \frac{b_1}{\sqrt{\mu\beta R}} \frac{1}{1 - \frac{b_1}{\sqrt{\mu\beta R}}} + \frac{1}{2} \\
&\leq \left(\frac{1}{\sqrt{\beta}} \right) (1-\delta) \cdot \frac{100 \cdot b_1}{100 - b_1} \sqrt{R} + \frac{1}{2} \\
&\leq \frac{1.15(1-\delta)}{\sqrt{\mu\beta}} \sqrt{R} + \frac{1}{2}.
\end{aligned}$$

Finally, we have that,

$$\int_{\delta R + \frac{1}{2}(\frac{\psi}{a} + 1)}^{\infty} e^{-aj^2} \mathbf{d}j \leq \int_{\delta R}^{\infty} e^{-aj^2} \mathbf{d}j \leq \int_{\delta R}^{\infty} e^{-a\delta R j} \mathbf{d}j = \frac{1}{a\delta R} e^{-a\delta^2 R^2} = \frac{2(1-\delta)}{\delta} e^{-R \frac{\delta^2}{1-\delta}}.$$

To complete the proof, we note that $e^{-\psi} = \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right)$, thus

$$\begin{aligned} \mathbb{E}[L] &\geq e^{\frac{a}{4}(\frac{\psi}{a} + 1)^2 - \psi} \int_0^{\delta R} e^{-a(j + \frac{1}{2}(\frac{\psi}{a} + 1))^2} \mathbf{d}j \\ &\geq e^{-\psi} \int_0^{\delta R} e^{-a(j + \frac{1}{2}(\frac{\psi}{a} + 1))^2} \mathbf{d}j \\ &\geq \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right) \left[\sqrt{\frac{\pi}{2}}(1-\delta) - \frac{1.15(1-\delta)}{\sqrt{\mu\beta}} \right] \sqrt{R} - \frac{1}{2} - \frac{2(1-\delta)}{\delta} e^{-R \frac{\delta^2}{1-\delta}} \\ &= \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right) \left[\sqrt{\frac{\pi}{2}}(1-\delta) - \frac{1.15(1-\delta)}{\sqrt{\mu\beta}} \right] \sqrt{R} - \frac{1}{2} - \frac{2(1-\delta)}{\delta} e^{-R \frac{\delta^2}{1-\delta}} \\ &= \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right) \left(\left[\sqrt{\frac{\pi}{2}}(1-\delta) - \frac{1.15(1-\delta)}{\sqrt{\mu\beta}} \right] \sqrt{R} - \frac{1}{2} - \frac{2(1-\delta)}{\delta} e^{-R \frac{\delta^2}{1-\delta}} \right). \end{aligned}$$

From here, we could choose δ to maximize our lower bound further based on system parameters, but a simple choice is $\delta = \frac{2}{\sqrt{R}}$. This gives

$$\begin{aligned} \mathbb{E}[L] &\geq \left(1 - \frac{b_1}{\sqrt{\mu\beta R}}\right) \left[\left(1 - \frac{2}{\sqrt{R}}\right) \left(\sqrt{\frac{\pi}{2}} - \frac{1.15}{\sqrt{\mu\beta}} - 2e^{-4} \right) - \frac{1}{2\sqrt{R}} \right] \sqrt{R} \\ &\geq \frac{2}{3} \sqrt{\frac{\pi}{2}} \sqrt{R}, \end{aligned}$$

as desired.

A.4.7 Proof of Claim 6.11.

Claim 6.11 (Upper Bound on $\mathbb{E}[N(T_A)]$). *Recall that $T_A \triangleq \min\{t > 0 : Z(t) = R + 1\}$. Then,*

$$\mathbb{E}[N(T_A) - R] \leq F_1 \mu \beta \sqrt{R} \left(1 + \frac{F_2}{\sqrt{\mu\beta}}\right) \leq 2.9 \mu \beta \sqrt{R}$$

and

$$\mathbb{E}[(N(T_A) - R)^2] \leq F_1^2 (\mu\beta)^2 R \left(1 + \frac{F_2}{\sqrt{\mu\beta}}\right)^2 + 2\mu\beta R \leq 8.4 (\mu\beta)^2 R + 2\mu\beta R$$

where $F_1 = 2.12$ and $F_2 = 3.645$.

Proof.

The beginning of the proof will be the same for both of these inequalities. Using the up-crossing and down-crossing decomposition of Section 6.2.1, we know that time T_A occurs either during a *rise* or during a *fall*. Since the number of jobs $N(t) \leq R + C_3\sqrt{\mu\beta R}$ during a rise,

$$[N(T_A) - R] \mathbf{1}_{T_A \text{ during a rise}} \leq C_3\sqrt{\mu\beta R} \mathbf{1}_{T_A \text{ during a rise}}.$$

If T_A occurs during a fall, we need a more nuanced bound. Writing out the event $\{T_A \text{ during a fall}\}$ in terms of disjoint events, we find

$$\{T_A \text{ during a fall}\} = \bigcup_{j=0}^R \bigcup_{i=1}^{\infty} \left\{ u_i^{(j)} \leq T_A < d_i^{(j)} \right\},$$

so that, for $c \in \{1, 2\}$,

$$\begin{aligned} \mathbb{E} \left[[N(T_A) - R]^c \mathbf{1}_{T_A \text{ during a fall}} \right] &= \sum_{j=0}^{R-1} \sum_{i=1}^{\infty} \mathbb{E} \left[[N(T_A) - R]^c \mathbf{1}_{u_i^{(j)} \leq T_A < d_i^{(j)}} \right] \\ &= \sum_{j=0}^{R-1} \sum_{i=1}^{\infty} \mathbb{E} \left[[N(T_A) - R]^c \mathbf{1}_{T_A < d_i^{(j)}} \middle| \mathcal{F}_{u_i^{(j)}}, n_u^{(j)} \geq i \right] \Pr(n_u^j \geq i). \end{aligned}$$

To bound this conditional expectation, we apply Claim A.4. Notice that $N(u_i^{(j)}) - R = C_3\sqrt{\mu\beta R}$, the $(R+1)$ -th server starts up at time $T_A = u_i^{(j)} + Y_{R+1}(u_i^{(j)})$ if $T_A < d_i^{(j)}$, the time $d_i^{(j)}$ is a hitting time, and that $Z(t) \geq R - j$ until time $\tau_{j+1} \geq d_i^{(j)}$. It follows that

$$\mathbb{E} \left[[N(T_A) - R] \mathbf{1}_{T_A < d_i^{(j)}} \middle| \mathcal{F}_{u_i^{(j)}}, n_u^{(j)} \geq i \right] \leq C_3\sqrt{\mu\beta R} + \mu j Y_{R+1}(u_i^{(j)}) \leq C_3\sqrt{\mu\beta R} + \mu j \beta,$$

and that

$$\begin{aligned} \mathbb{E} \left[[N(T_A) - R]^2 \mathbf{1}_{T_A < d_i^{(j)}} \middle| \mathcal{F}_{u_i^{(j)}}, n_u^{(j)} \geq i \right] &\leq \left(C_3\sqrt{\mu\beta R} + \mu j Y_{R+1}(u_i^{(j)}) \right)^2 + \mu 2R Y_{R+1}(u_i^{(j)}) \\ &\leq \left(C_3\sqrt{\mu\beta R} + \mu j \beta \right)^2 + \mu 2R \beta \\ &= C_3^2 \mu \beta R + 2C_3\sqrt{\mu\beta R} \mu \beta j + (\mu \beta)^2 j^2 + 2\mu R \beta. \end{aligned}$$

It now suffices to bound $\sum_j \sum_i j^c \Pr(n_u^{(j)} \geq i)$, where $c \in \{0, 1, 2\}$. We do this via the same method used in Section 6.2.1:

$$\begin{aligned} \sum_{j=1}^R \sum_{i=1}^{\infty} j^c \Pr(n_u^{(j)} \geq i) &\leq \sum_{j=1}^R \sum_{i=1}^{\infty} j^c \Pr(n_e \geq j) p_{\text{rise}}^{(j)} (1 - p_2)^{i-1} \\ &= \frac{1}{p_2} \sum_{j=1}^R j^c p_{\text{rise}}^{(j)} \Pr(n_e \geq j) \\ &\leq \frac{1}{C_4 p_2} \sum_{j=1}^R j^c C_4 p_{\text{rise}}^{(j)} \prod_{\ell=0}^{j-1} \left(1 - C_4 p_{\text{rise}}^{(\ell)} \right). \end{aligned}$$

This is simply the expectation of a time-varying geometric random variable G , with $\Pr(G = j | G \geq j) = C_4 p_{\text{rise}}^{(j)}$. It follows that if one lower-bounds $p_{\text{rise}}^{(j)}$, then an upper bound on the desired expectation is obtained. Applying Claim 6.3, we note that we are essentially bounding G using $Y \sim \text{Geometric}\left(\frac{0.99C_4A_5}{\sqrt{R}}\right)$ and saying $Y + A_5\sqrt{R}$ stochastically-dominates G . It follows that

$$\mathbb{E}[G] \leq A_5\sqrt{R} + \frac{1}{0.99C_4A_5}\sqrt{R}$$

and that, for any b ,

$$\begin{aligned} \mathbb{E}[(G + b)^2] &\leq \mathbb{E}[(Y + A_5\sqrt{R} + b)^2] = \mathbb{E}[Y^2] + 2(A_5\sqrt{R} + b)\mathbb{E}[Y] + (A_5\sqrt{R} + b)^2 \\ &= 2\mathbb{E}[Y]^2 - \mathbb{E}[Y] + 2(A_5\sqrt{R} + b)\mathbb{E}[Y] + (A_5\sqrt{R} + b)^2 \\ &\leq (\mathbb{E}[Y] + A_5\sqrt{R} + b)^2 + \mathbb{E}[Y]^2. \end{aligned}$$

Defining $B_5 \triangleq \frac{C_3}{C_4p_2}$, $B_6 \triangleq \frac{1}{C_4p_2} \left(\frac{1}{0.99C_4A_5} + A_5 \right)$, and $B_7 \triangleq \frac{1}{2C_4p_2} \left[\frac{1}{(0.99C_4A_5)^2} + 2 \right]$, it follows that

$$\begin{aligned} &\mathbb{E} \left[[N(T_A) - R]^2 \mathbf{1}_{T_A \text{ during a fall}} \right] \\ &\leq \frac{1}{C_4p_2} \sum_{j=1}^R C_4p_{\text{rise}}^{(j)} \left[\left(C_3\sqrt{\mu\beta R} + \mu j\beta \right)^2 + \mu 2R\beta \right] \prod_{\ell=0}^{j-1} \left(1 - C_4p_{\text{rise}}^{(\ell)} \right) \\ &= \frac{1}{C_4p_2} \mathbb{E} \left[\left(C_3\sqrt{\mu\beta R} + \mu G\beta \right)^2 + \mu 2R\beta \right] \\ &\leq \frac{1}{C_4p_2} \mathbb{E} \left[\left(C_3\sqrt{\mu\beta R} + \mu\beta Y + \mu\beta A_5\sqrt{R} \right)^2 + \mu 2R\beta \right] \\ &= \frac{1}{C_4p_2} \left[\left(C_3\sqrt{\mu\beta R} + \mu\beta \frac{1}{0.99C_4A_5}\sqrt{R} + \mu\beta A_5\sqrt{R} \right)^2 + \mu\beta \frac{1}{(0.99C_4A_5)^2}R + 2\mu\beta R \right] \\ &\leq \frac{1}{C_4^2p_2^2} \left[\left(C_3\sqrt{\mu\beta R} + \mu\beta \frac{1}{0.99C_4A_5}\sqrt{R} + \mu\beta A_5\sqrt{R} \right)^2 \right] + \frac{1}{C_4p_2} \left[\frac{1}{(0.99C_4A_5)^2} + 2 \right] \mu\beta R \\ &= \left(B_5\sqrt{\mu\beta R} + B_6\mu\beta\sqrt{R} \right)^2 + 2B_7\mu\beta R. \end{aligned}$$

and that

$$\begin{aligned}
& \mathbb{E} \left[[N(T_A) - R] \mathbf{1}_{T_A \text{ during a fall}} \right] \\
& \leq \frac{1}{C_4 p_2} \sum_{j=1}^R C_4 p_{\text{rise}}^{(j)} \left[C_3 \sqrt{\mu \beta R} + \mu j \beta \right] \prod_{\ell=0}^{j-1} \left(1 - C_4 p_{\text{rise}}^{(\ell)} \right) \\
& = \frac{1}{C_4 p_2} \mathbb{E} \left[C_3 \sqrt{\mu \beta R} + \mu \beta G \right] \\
& \leq \frac{1}{C_4 p_2} \left[C_3 \sqrt{\mu \beta R} + \mu \beta \left(A_5 \sqrt{R} + \frac{1}{0.99 C_4 A_5} \sqrt{R} \right) \right] \\
& = \left(B_5 \sqrt{\mu \beta R} + B_6 \mu \beta \sqrt{R} \right).
\end{aligned}$$

Defining $F_1 \triangleq B_6$ and $F_2 \triangleq \frac{(B_5 + C_3)}{B_6}$, it follows that

$$\mathbb{E} [N(T_A) - R] \leq \left((B_5 + C_3) \sqrt{\mu \beta R} + B_6 \mu \beta \sqrt{R} \right) = F_1 \mu \beta \sqrt{R} \left(1 + \frac{F_2}{\sqrt{\mu \beta}} \right)$$

and that

$$\mathbb{E} \left[[N(T_A) - R]^2 \right] \leq \left(B_5 \sqrt{\mu \beta R} + B_6 \mu \beta \sqrt{R} \right)^2 + 2B_7 \mu \beta R \leq F_1 (\mu \beta)^2 R \left(1 + \frac{F_2}{\sqrt{\mu \beta}} \right)^2 + 2\mu \beta R \square$$

Bibliography

- [1] Ali Allahverdi and HM Soroush. The significance of reducing setup times/setup costs. *European journal of operational research*, 187(3):978–984, 2008. 1
- [2] Jesus R Artalejo, Antonis Economou, and Maria Jesus Lopez-Herrero. Analysis of a Multiserver Queue with Setup Times. *Queueing Syst.*, 51(1):53–76, 2005. 2.1.2
- [3] Wolfgang Bischof. Analysis of M/G/1-Queues with Setup Times and Vacations under Six Different Service Disciplines. *Queueing Syst.*, 39(4):265–301, 2001. 2.1.1
- [4] Gautam Choudhury. On a batch arrival Poisson queue with a random setup time and vacation period. *Comp. & Oper. Res.*, 25(12):1013–1026, 1998. 2.1.1
- [5] John D. Cook. Details for error bound on normal approximation to the Poisson distribution. https://www.johndcook.com/blog/berry_esseen_poisson/, 2024. Accessed: 01/16/2024. A.2.4, A.4.5
- [6] Andrew Daw, Robert C Hampshire, and Jamol Pender. How to staff when customers arrive in batches. *arXiv preprint arXiv:1907.12650*, 2019. 8.3.2
- [7] Zohar Feldman, Avishai Mandelbaum, William A Massey, and Ward Whitt. Staffing of time-varying queues to achieve time-stable performance. *Management Science*, 54(2):324–338, 2008. 8.3.2
- [8] Anshul Gandhi and Mor Harchol-Balter. How Data Center Size Impacts the Effectiveness of Dynamic Power Management. In *Proc. Ann. Allerton Conf. Communication, Control and Computing*, pages 1164–1169, Urbana-Champaign, IL, September 2011. 2.4
- [9] Anshul Gandhi and Mor Harchol-Balter. M/G/k with staggered setup. *Oper. Res. Lett.*, 41(4):317–320, 2013. 2.1.2
- [10] Anshul Gandhi, Varun Gupta, Mor Harchol-Balter, and Michael Kozuch. Optimality analysis of energy-performance trade-off for server farm management. In *Proc. Int. Symp. Computer Performance, Modeling, Measurements and Evaluation (IFIP Performance)*, Namur, Belgium, November 2010. 2.4
- [11] Anshul Gandhi, Mor Harchol-Balter, and Ivo Adan. Server farms with setup costs. *Performance Evaluation*, 67(11):1123–1138, 2010. 2.1.2, 2.2.1, 2.2.2, 2.4
- [12] Anshul Gandhi, Mor Harchol-Balter, and Mike Kozuch. The case for sleep states in servers. In *SOSP Workshop on Power-Aware Computing and Systems (HotPower)*, pages 1–5, Cascais, Portugal, October 2011. 2.4
- [13] Anshul Gandhi, Mor Harchol-Balter, and Mike Kozuch. Are sleep states effective in data

- centers? In *Int. Conf. Green Computing (IGCC)*, pages 1–10, San Jose, CA, 2012. 2.4
- [14] Anshul Gandhi, Mor Harchol-Balter, Ram Raghunathan, and Michael A Kozuch. AutoScale: Dynamic, Robust Capacity Management for Multi-Tier Data Centers. *ACM Trans. Comput. Syst.*, 30(4):1–26, 2012. 1, 1.3, 1.4, 2.4
- [15] Anshul Gandhi, Sherwin Doroudi, Mor Harchol-Balter, and Alan Scheller-Wolf. Exact analysis of the M/M/k/setup class of Markov chains via Recursive Renewal Reward. In *Queueing Syst.*, pages 153–166, 2013. 1.3, 1.3, 1.4.3, 2.2.2
- [16] Shlomo Halfin and Ward Whitt. Heavy-traffic limits for queues with many exponential servers. *Oper. Res.*, 29(3):567–588, 1981. 1.4.1
- [17] Jianwei Hao, Ting Jiang, Wei Wang, and In Kee Kim. An empirical analysis of vm startup times in public iaas clouds. In *2021 IEEE 14th International Conference on Cloud Computing (CLOUD)*, pages 398–403. IEEE, 2021. 1, 1, 1.3, 1.4
- [18] Qi-Ming He and E Jewkes. Flow time in the MAP/G/1 queue with customer batching and setup times. *Stochastic Models*, 11(4):691–711, 1995. 2.1.1
- [19] Yige Hong and Ziv Scully. Performance of the gittins policy in the g/g/1 and g/g/k, with and without setup times. *ACM SIGMETRICS Performance Evaluation Review*, 51(2):33–35, 2023. 2.3
- [20] Esa Hyytiä, Douglas Down, Pasi Lassila, and Samuli Aalto. Dynamic Control of Running Servers. In *Int. Conf. Measurement, Modelling and Evaluation of Comput. Systems*, pages 127–141, Erlangen, Germany, 2018. Springer. 1.3, 2.4
- [21] Otis B Jennings, Avishai Mandelbaum, William A Massey, and Ward Whitt. Server staffing to meet time-varying demand. *Management Science*, 42(10):1383–1394, 1996. 8.3.2
- [22] Aytaç Kara. Energy Consumption in Data Centers with Deterministic Setup Times. Master’s thesis, Middle East Technical University, 2017. 1.3, 2.4
- [23] Ioannis Karatzas and Steven Shreve. *Brownian motion and stochastic calculus*, volume 113. Springer Science & Business Media, 2012. A.2.5
- [24] Xin Liu and Lei Ying. Universal scaling of distributed queues under load balancing in the super-Halfin-Whitt regime. *IEEE/ACM Trans. Netw.*, 30(1):190–201, 2022. 1.4.1
- [25] Francesco Longo, Letizia Nicoletti, and Antonio Padovano. Smart operators in industry 4.0: A human-centered approach to enhance operators’ capabilities and competencies within the new smart factory context. *Computers & industrial engineering*, 113:144–159, 2017. 1
- [26] Sumit Maheshwari, Dipankar Raychaudhuri, Ivan Seskar, and Francesco Bronzino. Scalability and performance evaluation of edge cloud systems for latency constrained applications. In *2018 IEEE/ACM Symposium on Edge Computing (SEC)*, pages 286–299. IEEE, 2018. 1.3, 1.4
- [27] Ming Mao and Marty Humphrey. A Performance Study on the VM Startup Time in the Cloud. In *IEEE Int. Conf. Cloud Computing (CLOUD)*, pages 423–430, Honolulu, HI, 2012. 1.3, 1.4, 1.4.3
- [28] Jeffrey C Mogul and Ramana Rao Kompella. Inferring the network latency requirements

- of cloud tenants. In *15th Workshop on Hot Topics in Operating Systems (HotOS XV)*, 2015. 1.3, 1.4
- [29] Debankur Mukherjee and Alexander Stolyar. Join Idle Queue with Service Elasticity: Large-Scale Asymptotics of a Nonmonotone System. *Stoch. Syst.*, 9(4):338–358, 2019. 2.2.3
- [30] Debankur Mukherjee, Souvik Dhara, Sem C Borst, and Johan SH van Leeuwen. Optimal Service Elasticity in Large-Scale Distributed Systems. *Proc. ACM SIGMETRICS Int. Conf. Measurement and Modeling of Computer Systems*, 1:1–28, 2017. 2.2.3
- [31] Jerome Niyirora and Jamol Pender. Optimal staffing in nonstationary service centers with constraints. *Naval Research Logistics (NRL)*, 63(8):615–630, 2016. 8.3.2
- [32] Jamol Pender and Tuan Phung-Duc. A law of large numbers for m/m/c/delayoff-setup queues with nonstationary arrivals. In *Int. Conf. on Analytical and Stochastic Modeling Techniques and Applications*, pages 253–268, Cardiff, UK, 2016. Springer. 1.3, 2.2.1, 2.4
- [33] Tuan Phung-Duc. Exact solutions for M/M/c/setup queues. *Telecommun. Syst.*, 64(2): 309–324, 2017. 1.3, 2.2.2
- [34] Krzysztof Rządca, Paweł Findeisen, Jacek Swiderski, Przemysław Zych, Przemysław Broniek, Jarek Kusmirek, Paweł Nowak, Beata Strack, Piotr Witusowski, Steven Hand, et al. Autopilot: workload autoscaling at Google. In *Proc. European Conf. Computer Systems (EuroSys)*, pages 1–16, Heraklion, Crete, Greece, 2020. 1, 8.2.1
- [35] Irina Shevtsova. On the absolute constants in the Berry-Esseen type inequalities for identically distributed summands, 2011. A.2.4
- [36] Peter D Welch. On a Generalized M/G/1 Queuing Process in Which the First Customer of Each Busy Period Receives Exceptional Service. *Oper. Res.*, 12(5):736–752, 1964. 1.2.1, 1.4.2, 2.1.1
- [37] Jalani Williams, Weina Wang, and Mor Harchol-Balter. Average waiting time in the M/M/k/Setup with Deterministic Setup Times. *Oper. Res.*, 2024, in preparation. 5.1
- [38] Jalani K. Williams, Mor Harchol-Balter, and Weina Wang. The M/M/k with Deterministic Setup Times. *Proc. ACM SIGMETRICS Int. Conf. Measurement and Modeling of Computer Systems*, 6(3), dec 2022. doi: 10.1145/3570617. 5.1, 5.2, A.3.1