

# Mean Field Behaviour of Job Redundancy Queueing Models

by

Aaron Janeiro Stone

I understand that my thesis may be made electronically available to the public.

## Abstract

Technological advancement in cloud computing has resulted in the viability of a new class of routing algorithm, the so-called redundancy models which are able to replicate jobs for processing on different servers. An asymptotic amount of queue-endowed servers could in this way employ their plentiful, otherwise idle servers towards processing. Research on the behaviour of such systems, however, has only conjectured the asymptotic independence of queues (Gardner, Harchol-Balter, & Scheller-Wolf, 2016; Hellermann, Bodas, & Van Houdt, 2019). To this end, by modelling the process as a time-indexed family of hypergraphs wherein hyperedges represent cloned dependents, along with the notion of exchangeable random groups derived by Austin (2008), we are able to demonstrate the conjectured behaviour in a simulation setting.

Is it 2016  
or 2017?  
Please check and  
be consistent in  
your document  
and references  
section.

et al.

## **Acknowledgements**

Thank you, Dr. Steve Drekic and Dr. Tim Hellemans, for your guidance as I delved into an area once foreign to me.

## Dedication

Thank you, Fatima and my family, for always being there for me in such a fast-moving world.

*The whole entire world is a very narrow bridge; the main thing is to have no fear at all.*

- Nachman

# Table of Contents

List of Tables	vii
List of Figures	viii
1 Introduction	1
1.1 Background . . . . .	1
2 Model Specification	5
3 Simulation Study	9
3.1 ParallelQueue Package . . . . .	9
3.2 Results . . . . .	11
References	17
APPENDICES	19

Is there a reason why this is labelled  
as a subsection? Are you planning on  
adding another subsection? If not, I would  
then suggest you simply remove the subsection  
"1.1 Background" vi

completely and just have the "Introduction"  
heading of  
Section 1.

## List of Tables

What is happening here?

Do you plan on having any  
tables included in your essay?

If not, then please <sup>remove</sup>✓ this from here  
as well as any reference to  
it in the Table of Contents.

# List of Figures

1.1	Commutativity of Limits . . . . .	3
2.1	A graphical representation of queues with replicated jobs as hyperedges. . . . .	8
3.1	Python code using ParallelQueue to simulate a Redundancy-2 System. . . . .	10
3.2	Plot using totals of Figure 3.1 . . . . .	10
3.3	Runtime Statistics . . . . .	11
3.4	Overview of the ParallelQueue API . . . . .	12
3.5	Comparisons of Systems: Values are averaged over 30 independent iterations each, running for $t = 1000$ . . . . .	13
3.6	Threshold(2,2) for Varying $N$ and $\rho$ . . . . .	15
3.7	Redundancy(2) ECDFs for Varying $N$ . . . . .	16
3.8	Redundancy(2) Histograms of ECDF Changes for Varying $\tilde{\tau}$ at $N = 1000$ .	16

# Chapter 1

## Introduction

### 1.1 Background

Enamoured by physicists for its ability to turn probabilistic behaviour into matters of determinism, Mean Field Theory (MFT) also has a place in the study of queueing systems as the number of queues become asymptotic.

**Definition 1** (Mean Field).

Over a time-filtered probability space  $(\Omega, \mathcal{F}_t, \mathcal{F}, P)$ , for  $N \in \mathbb{N}$ , a mean field describes the behaviour of any set of stochastic random variables

$$\mathbf{X}^{(N)}(t) := \{X_i(t)\}_{i \leq N}$$

that turns deterministic in law as  $N \rightarrow \infty$ , irrespective of if the finite- $N$  or finite- $t$  cases resulted in this set of "bodies" being dependent [Mukhopadhyay & Mazumdar, 2015].

Next, it is important to define exchangeability, a condition which provides useful properties for MFT.

**Definition 2** (Exchangeability). A collection of random variables  $\mathbf{X}^{(N)}$  is exchangeable if for any  $N$ -permutation  $\gamma_N \in \Gamma_N$ ,

$$\text{Law}(\gamma_N \mathbf{X}^{(N)}) = \text{Law}(\mathbf{X}^{(N)}).$$

de Finetti's Theorem states that exchangeability implies the collection to be conditionally independent (in a Markovian sense) and identically

distributed [Austin, 2015]. Moreover, for partition  $N = \bigcup_{i \leq k} N_i$ , exchangeability over partition  $\{N_i\}_{i \leq k}$  such that

$$\text{Law}(\mathbf{X}^{(\gamma\{N_i\}_{i \leq k})}) = \text{Law}(\mathbf{X}^{\{N_i\}_{i \leq k}})$$

ITALICIZE

is known as  $(N, \Gamma)$  exchangeability (with  $\gamma \in \Gamma$ ) [Austin, 2008].

With multiprocessing being employed at its current scale in server farms, society is indeed approaching a time wherein the asymptotic behaviour of parallel queueing systems can be considered realistically. Lately, interest has been given to queueing systems which employ job redundancy in order to lower total processing time (Ayesta et al., 2018). That is, routing policies which take advantage of scenarios wherein a surplus of queues and/or servers are available, replicating each job such that it might be completed faster should it happen to make its way through a less-busy queue than the original.

The letter "n" is missing

**Definition 3** (Job Redundancy).

A scheduler,  $\mathcal{D}$ , follows a job redundancy policy if it systematically clones arriving jobs and removes all clones upon (or after a delay following) completion of any one clone.

Prior studies have pointed towards certain redundancy policies as being inefficient or even unrealistic due to over-relying on cloning. Take for example *Redundancy-d*, wherein  $d$  servers are chosen per arrival; each is given a clone and upon completion of any one clone, all others are removed immediately without any cost. As such, implementing a threshold on when to clone a job becomes useful for budgeting cancellation costs. In particular,  $\text{Threshold}(R, d)$  has risen to prominence as a means to balance workload in queueing systems.

and System Load

**Definition 4** (Workload).

**Workload** refers to the total amount of work remaining (in time) for a queue. In trivial cases not involving enqueued bodies potentially leaving, this would merely be the sum of individual jobs' service times. With  $w_j^{(i)}(t)$  denoting the (random) service time of the  $j$ th job in queue  $X_i$  at time  $t$ , the workload of a queue would be:

$$W_i(t) = \sum_{j \leq \#X_i(t)} w_j^{(i)}(t) \quad (1.1)$$

for counting measure  $\#$ , given that  $\#(X_i^{-1}) : t \mapsto z$ , (or that a queue has countable jobs). **System Load** refers to the amount of work remaining in the entire system,

$$W(t) = \sum_{i \in \psi} W_i(t)$$

This underlined text needs to be better expressed/clarified.  
Can you please improve the wording/meaning?

namely

$$\begin{array}{ccc} \mathbf{X}^{(N)}(t) & \xrightarrow{N \rightarrow \infty} & \mathbf{X}(t) \\ \downarrow t \rightarrow \infty & & \downarrow t \rightarrow \infty \\ \mathbf{X}^{(N)}(\infty) & \xrightarrow{N \rightarrow \infty} & \pi \end{array}$$

Figure 1.1: Commutativity of Limits

where  $\psi \subseteq \mathbb{N}$ , given that we will be considering the case of systems operating in finite time as the number of queues grows indefinitely. As such,  $\psi$  will refer to this more-general case from henceforth.

As an example, a service time in a  $G/M/c$  system will be drawn from an exponential distribution. In this simple case, the  $m$ th arriving job can be given the "marks"  $(T_m, S_m) \equiv (T, S)_m$  where  $S_m \stackrel{IID}{\sim} \text{EXP}(\lambda)$  and  $T_m$  is the time of arrival; conditioning on the process  $(T, S)_m$ ,  $W_i(t)$  turns into a matter of merely adding up the enqueued service times and that remaining of the currently serviced job, which is conveniently memoryless.

Altogether, Figure 1.1 describes the behaviour which would be expected in an ideal system wherein both a mean field and asymptotic independence can be achieved [Mukhopadhyay and Mazumdar, 2015]. In particular,  $P$  describes a fixed point, giving the collection a distribution of  $\delta_P$ . In terms of the predictability and stability of a system, needless to say, this would be a "gold-standard"; one would be interested in their ability to achieve such a system in practice.

For the sake of brevity, the notation of  $[n] := \{i \in \mathbb{N} | i \leq n\}$  will be used, along with the understanding of  $\mathbf{X}^{[n]} \equiv \mathbf{X}^{(n)}$  for maximal element  $n$ . Moreover, accepting this set-index notation,  $\mathbf{X}^\psi$  will refer to the general case of  $[n]$ , given we will also consider  $[n] \xrightarrow{n \rightarrow \infty} \mathbb{N}$ .

**Definition 5** (Threshold( $R, d$ )).

Threshold( $R, d$ ), denoted  $D_{\text{THRESH}(R,d), Z}$  selects  $d$  queues upon a job arrival. Next:

1. For  $i \leq d$  queues which have workload less than or equal to  $R$ , place copies in these

$i$

2. If  $i = 0$ , place the original arrival in a queue from the  $d$  chosen at random.

queues

$Z$  refers to any imposed job cancellation cost (e.g., an added temporary workload).

This underlined text needs to be better expressed/clarified.  
It reads confusing as is. Can you please improve the wording/meaning here?

One important question, however, is yet to be answered. It is unknown whether or not there exists sufficient arrival rate or service rate parameters such that a mean field will be observed for particular values of  $R$  or  $d$  in the threshold model. This leads us to the following conjecture for which this paper aims to prove.

**Conjecture 1.**

*As  $N \rightarrow \infty$ , the system  $\mathcal{D}_{THRESH(R,d),Z}$  becomes  $(\psi, \Gamma)$ -exchangeable. As  $t \rightarrow \infty$ , the system becomes deterministic.*

## Chapter 2

# Model Specification

In order to model such a system, we incorporate the notation most frequently seen in the study of palm calculus [Bacelli and Brémaud, 2003]. Most importantly, assuming we delegate a probability space  $(\Omega, \mathcal{F}_t, \mathcal{F}, P)$  the measurable flow  $\{\theta_t\}$  which is  $P/\theta_t$  invariant (Ergodic), such that  $\theta_t M = M$  for some  $M \in \mathcal{F} \Rightarrow P(M) \in \{0, 1\}$ , then the arrival process  $A$  can be associated with the flow. In a system with Markovian arrivals, the counting measure associated with the flow would necessarily be Poisson; other so-called counting processes wherein inter-arrivals are determined by random or even deterministic periods of time can be used. Inter-arrival times (for arrival  $n$  occurring at  $T_n$ ), regardless of the generating process, shall be denoted

$$\tau_n = T_{n+1} - T_n, n \in \mathbb{Z}.$$

Intensities (even if non-Poisson) will be denoted  $\lambda = E(A((0, 1]))$  for arrival process  $A$ , being interpreted as the intensity of a process moving a state (i.e., counting an additional element) within unit time.

**Definition 6** (Marked Process). For each job which enters the system, they can be "marked" by series of random variables defining their behaviour within the system [Bacelli and Brémaud, 2003]. In general, for the systems we shall consider, we consider the mark processes  $\sigma_n$  as being the required service of arrival  $n$  and the marks  $T_n$  as its time of arrival. The tuple

$$(T, \sigma)_n, n \in \mathbb{Z}$$

will therefore be used to mark the  $n$ th job. We shall extend this notation, however, to include sets; for set of jobs  $\eta$ , wherein each element has their own arrival time:

$$(\sigma, T)_\eta = \{(T, \sigma)_n\}_{n \in \eta}.$$

marked process

5

Shouldn't this be  
 $(T, \sigma)_n$   
here?

Change to  
comma

Should  
this  
be  
 $\mathbb{Z}^+$   
here?

This underlined text is confusing  
and needs to be better  
expressed/clarified.  
Can you please improve it here?

queue

For the marked process to be useful, we would like to be able to condition onto it in order to ease the problem of determining workload into one of considering (conditionally) non-random terms. In a simple  $G/M/c$ , as mentioned before, the double  $(T, \sigma)_n$  would be sufficient. In our case, given that jobs can find their status in the system tied to the behaviour of their replicas, we must append a marking to track this phenomenon. As a matter of fact, we instead move to marking the queues as was done by [Bramson et al., 2012] to prove the conjecture in the case of Join-The-Shortest-Queue (d) systems, although with an additional job dependency term.

ITALICIZE

put  
in  
math  
font

### Definition 7 (Job Dependency Graph: Finite Case).

The **Job Dependency Graph**,  $G_t = (V, E)_t$  is held constant between the events of arrivals and job completions, where an edge is drawn between two nodes if and only if they are job-dependent.

1. Movement in a queue requires appropriate redrawing of graph.
2.  $G_t$  is stochastic with law in  $\Pr(\Omega, \mathcal{F}_t)$ .
3. Maximal system queue size is bounded,  $\sup_{n \in \chi} z_n(t) := \nu(t)$ .

66

"Appropriate redrawing" in this case means

1. The graph is redrawn to reflect movement within each queue (moving due to job completions or arrivals).
2.  $G_t$  can depend only on  $G_s, s < t$  and other current values of  $X_t^*$  (denote these other values  $\tilde{X}_t^*$ ) and is such that

$$P(G_t | \tilde{X}_t^*, \{G_a\}_{a \in S}) = P(G_t | \tilde{X}_t^*, G_{\max(S)})$$

by

insert  
comma  
here

for any set  $S$  such that  $S \subset [0, t]$ .

### Definition 8 (Job Dependency Matrix: Finite Case).

The **Job Dependency Matrix** is an adjacency matrix (non-unique but one-to-one) for  $G_t$ . Specifically,

$$\rho(t) \equiv \rho(G_t) = \left[ \begin{array}{c|c|c|c} B_{1,1} & B_{1,2} & \dots & B_{1,\nu} \\ \hline \vdots & \vdots & \vdots & \vdots \\ \hline B_{\nu,1} & B_{\nu,2} & \dots & B_{\nu,\nu} \end{array} \right] \quad (2.1)$$

is a blocked matrix such that sub matrix

$$B_{i,j} = [b_{q,m}]_{q,m \leq N} = \begin{cases} 1, & \text{queue } q \text{ in row } i \text{ is connected in } G_t \\ & \text{to queue } m \text{ in row } j \\ 0, & \text{otherwise} \end{cases}$$

i.e.,

where connections between jobs occur if and only if the completion of one job implies the removal of all other connected jobs from the system.

While complicated in the above form, this graph merely draws an edge between jobs of separate queues while tracking their current place in their respective queues. Because at most one event can happen in infinitesimal time (an arrival or departure), the assumptions merely state that between any two events this graph shall not need to be redrawn. Extending this notion, one can iterate such a procedure indefinitely, viewing the resultant infinite graph as one where jobs are joined by hyperedges if and only if there exists an element of a set of replicas in each of the queues. In other words, one can represent queues which are connected by means of having replicas of the same jobs enqueued within them by viewing each as a node connected with a hyperedge; we will call queues related by such a hyperedge members of a *replica class*. Thus, 7 can be represented more succinctly in a hypergraph form, allowing one to 10 represent the cavity process studied in [Hellemans et al., 2019].

**Definition 9.** (Hypergraph Representation of the Dependency Matrix) By collapsing edges as jobs and vertices as servers, all connected subgraphs of  $G_t$  can be thought of as incident nodes with edges  $i_1, \dots, i_{j,d}$  representing connected jobs. This gives us now only as many nodes as there are currently servers. Such a representation for  $G_t$  will be denoted  $\tilde{G}_t$ . See Fig 2.1 for a visual representation.

$i_1, \dots, i_j, j < d,$

When one creates the infinite graph iteratively, it is meant that one could view the graph described in 7 as an embedding into a larger graph by merely considering more queues or a larger maximal queue size at any finite time  $t$ , corresponding to the adding of an additional column or row of vertices, respectively. Building a metric on an infinite graph space simplifies the matter of quantifying convergence in terms of  $N$  significantly because all possible finite embeddings can be expressed in the same space as  $N \rightarrow \infty$ . Thus, let us consider

$$\mathbb{E} := \{\text{locally finite graphs}\}.$$

Now, extending the notation of [Bramson et al., 2012], we can fully describe the marked process for queues with general service times.

ITALICIZE

Do you mean  
Definition 7  
here?  
If so, you should  
then write  
"Definition 7"  
and not just "7."

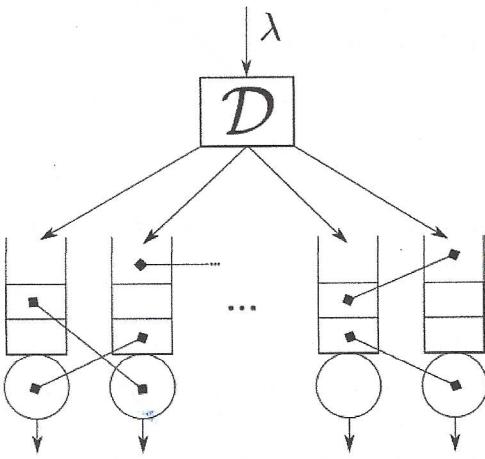


Figure 2.1: A graphical representation of queues with replicated jobs as hyperedges.

**Definition 10** (State Space With Redundancy).

For  $X_n^*(t)$  corresponding to the  $n$ th queue out of  $N$ ,

$$X_n^*(t) \in (\mathbb{N}, (\mathbb{R}^+)^3, \rho(\mathbb{E})) := \mathcal{E}_n^{(\psi)}$$

$$X^*(t) \in \{\mathcal{E}_n^{(\mathbb{N})}\}_{n \leq N} := \mathcal{E}^{(\psi)}$$

Corresponding to

1.  $z_n(t) \in \mathbb{N}$  for queue size
2.  $w_n(t) \in \mathbb{R}^+$  for workload
3.  $\ell_n(t) \in \mathbb{R}^+$  amount of service time spent on job currently in server
4.  $v_n(t) \in \mathbb{R}^+$  time remaining for job in server
5.  $A \in \rho(\mathbb{E})$  is a representation of the graph

possible queues

This definition is expressed very awkwardly.  
It needs improvement and better clarification.  
Please look after improving this.

# Chapter 3

## Simulation Study

### 3.1 ParallelQueue Package

In order to generate and study parallel queueing processes, few trivial options currently exist. Moreover, while there exist some discrete event simulation (DES) frameworks which indeed focus on queueing networks, they currently tend not to permit the simultaneous study of asynchronous, redundancy-based schemes [Palmer et al., 2019]. In order to visualize and analyse the large class of queueing systems within this paradigm [Shneer & Stolyar, 2020, Cruise et al., 2020] I introduced a novel module for Python which is currently available on PyPi: ParallelQueue, extending the DES package SimPy.

The package currently allows for the studying of parallel systems with or without redundancy as well as with the option of allowing thresholds to be implemented in either case. Moreover, the package allows one to specify any inter-arrival and service distribution as well as their own Monitors, being a class which can gather data from the ongoing simulation to be distributed back to the user upon the completion of a simulation. In particular, the Monitors are currently configured to collect data upon arrival, routing, and job completion Fig 3.3.

Take Figures 3.1 and 3.2, for example, which permits one to simulate a Redundancy-2 queueing system with 100 queues in parallel for 10000 units of time while returning the total queue counts over time (which are updated upon a change in queue count).

Remarkably, the simulation itself was fast despite the size of the system as in Fig 3.3.

3.1 and 3.2,

9

Don't you mean  
1000 here?

Insert  
Comma

demonstrated  
in  
Figure 3.3.

ITALICIZE

time

(as demonstrated by Fig 3.4).

```

#!/usr/bin/python3
from parallelqueue.base_models import RedundancyQueueSystem
from parallelqueue.monitors import TimeQueueSize
import random

sim = RedundancyQueueSystem(
    maxTime=1000.0, parallelism=100, seed=1234,
    d=2, Arrival=random.expovariate,
    AArgs=9, Service=random.expovariate,
    SArgs=0.08, Monitors=[TimeQueueSize])
# Note RedundancyQueueSystem is a ParallelQueueSystem wrapper
sim.RunSim()
totals = sim.MonitorOutput["TimeQueueSize"]

```

Figure 3.1: Python code using ParallelQueue to simulate a Redundancy-2 System.

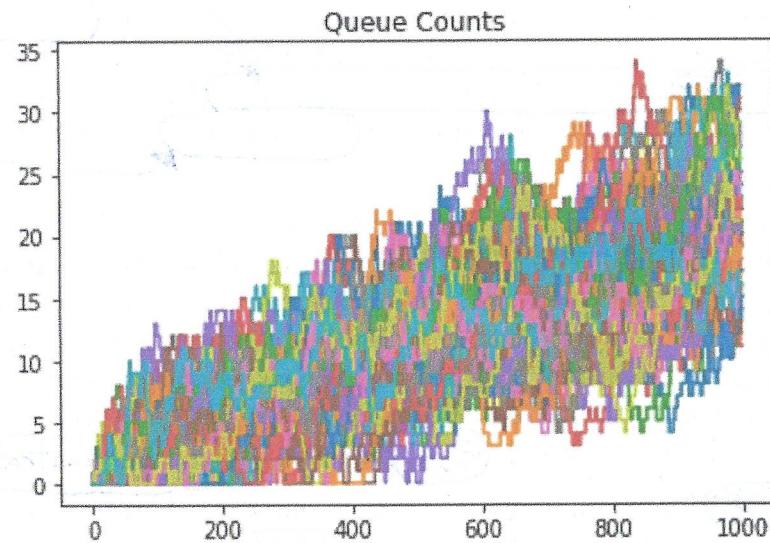
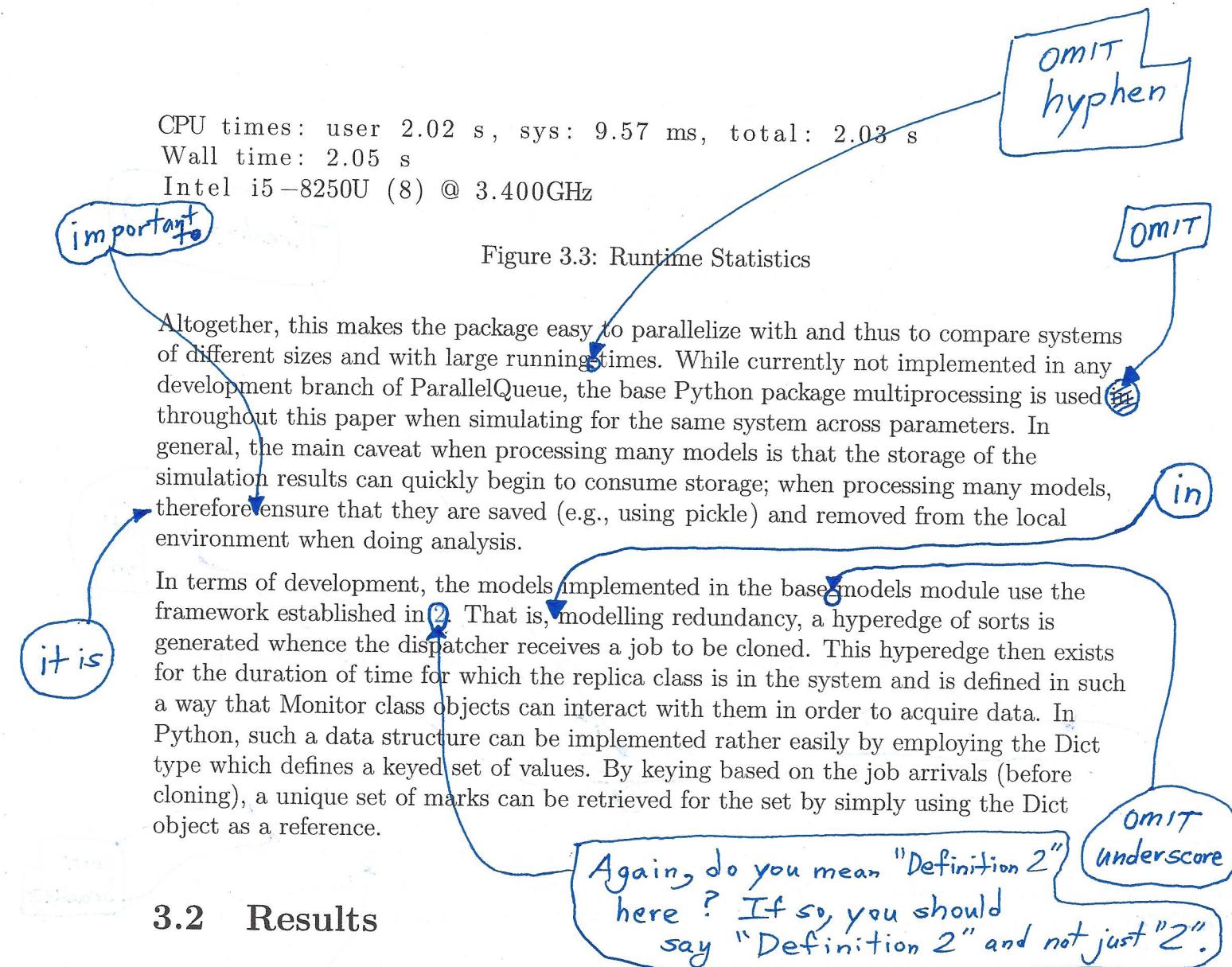


Figure 3.2: Plot using totals of Figure 3.1

CPU times: user 2.02 s, sys: 9.57 ms, total: 2.03 s  
Wall time: 2.05 s  
Intel i5-8250U (8) @ 3.400GHz



## 3.2 Results

First, we examine each model in terms of their respective performance in  $E(T)$ , the expected time each job spends in the system. As Fig 3.5 shows, for a load  $\rho \triangleq \frac{\lambda}{\mu} = 0.5$  by taking  $\mu = 1, \lambda = 0.5$  (we will assume  $\mu \equiv 1$  for the rest of the simulations), Redundancy(2) and Threshold(2,2) policies are rather alike with low loads as  $N \rightarrow \infty$ . This is to be expected, of course, given that even such a low threshold is unlikely to be exceeded with the processors acting faster than arrivals on average. Ignoring cancellation costs, this clearly demonstrates how utilizing otherwise dormant queues comes to benefit the system's performance. Note that the figure is generated with the *same* seed generation scheme for 30 different seeds per iteration, making the overlap a product of

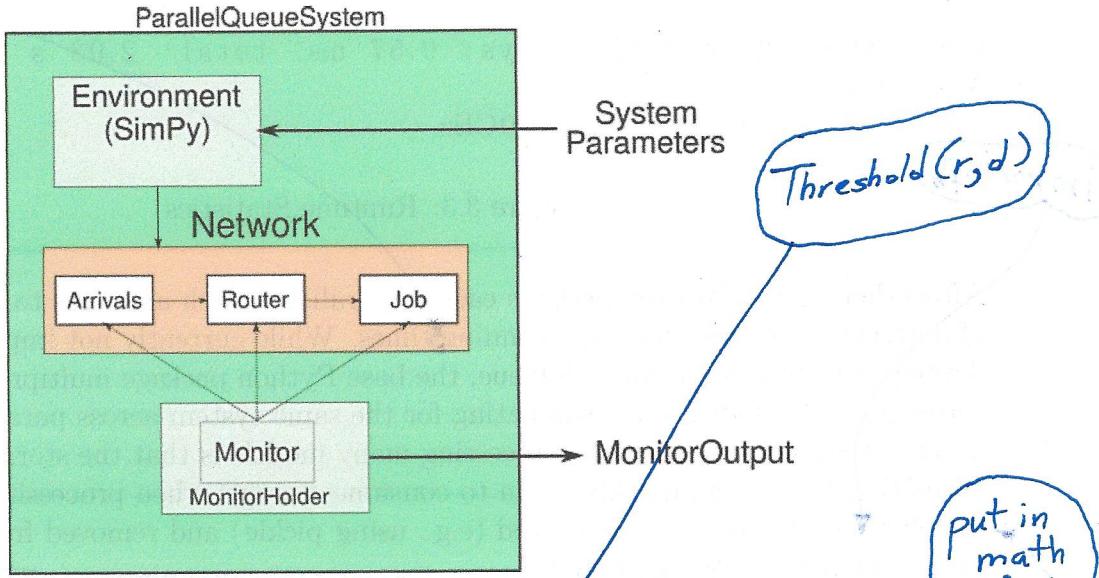


Figure 3.4: Overview of the ParallelQueue API

the two accessing random numbers from the system at the same points in time (in their respective simulations).

As [Gardner et al., 2017] show, a Redundancy- $d$  system is asymptotically stable if and only if  $\rho < 1$ . Given that, in premise, Threshold models are more or less a superclass of Join-The-Shortest queue and Redundancy- $d$  (trivially with rising  $r$  implying no threshold exists and thus copies should always be made as in the case of Redundancy- $(d)$ ), it is perhaps most interesting to examine if Threshold models are better able to handle high-load environments. Proving the superclass property in terms of Redundancy is relatively easy and is done in Lemma 1. By contrast, after merely setting  $r \equiv 0$ , we get  $\mathcal{D}_{\text{Thresh}(0,d)} \stackrel{d}{=} \mathcal{D}_{\text{JSQ}(d)}$  by definition.

**Lemma 1.** For  $X^{(n)}$  being a system of queues such that  $\rho < 1$ ,

$$\mathcal{D}_{\text{Thresh}(r,d)}(X^{(n)}) \xrightarrow{r \rightarrow \infty} \mathcal{D}_{\text{Red}(d)}(X^{(n)})$$

*Proof:*

Take  $\leq_{st}$  to mean that [Bramson et al., 2012]

$$X_1^{(N)} \leq_{st} X_2^{(N)} \iff \#X_{1,i} \leq \#X_{2,i} \quad \forall i \in [N] \quad P - \text{as}$$

as or  
a.s.?  
You use both  
acronyms. Choose  
one AND be consistent.  
place period here

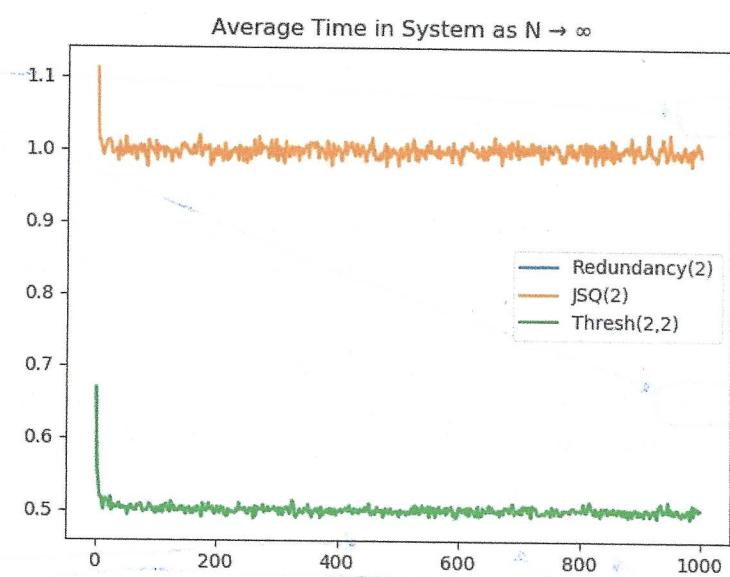


Figure 3.5: Comparisons of Systems: Values are averaged over 30 independent iterations each, running for  $t = 1000$ )

This proof is "jammed packed" with tight spacing, making it difficult for the reader to understand everything being introduced and discussed. I encourage you to better "organize" this proof for improved readability.

To show the sequence to be convergent, we will construct a coupling in such a manner that  $\forall r \in \mathbb{R}^+$ , the system is bounded by another (Baccelli & Brémaud, 2003). For  $\mathbf{X}_1^{(N)}$  under  $\text{Thresh}(r, d)$ , await the first arrival, denoted  $T_{\xi_1}$ , where a selection set,  $\nu \subset [N]$  is prescribed such that  $\exists \hat{X} \in \{X_i\}_{i \in \nu}$  such that  $\#(\hat{X}) > r$ . Thus, for  $t \in [0, T_{\xi_1}]$ ,  $\mathcal{D}_{\text{Thresh}(r, d)}(X^{(n)}) = \mathcal{D}_{\text{Red}(d)}(X^{(n)})$ . For clarity, let us now consider  $\mathbf{Y}^{(N)}$  to be a copy of  $\mathbf{X}^{(N)}$  such that they are independent, identical in distribution and in terms of the marks of arrival process and job-size draws along with the queues parsed (as was similarly done in (Bramson et al., 2012) Lemma 4.1), effectively, the only difference being left between these copies is  $r$  changing which queues receive clones (and thus, too, the mark of queue-dependency).  $\mathcal{D}_{\text{Red}(d)}(X^{(n)}) = \mathcal{D}_{\text{Thresh}(r, d)}(Y^{(n)})$ ; clearly, at time  $T_{\xi_1}$ ,  $\mathbf{X}^{(N)} \leq_{st} \mathbf{Y}^{(N)}$  due to jobs only being added for  $\hat{Y}_i$  (defined analogously to  $\hat{X}_i$ ).

Now, assume  $\rho < 1$ , giving us  $\#Y_i < \infty \quad \forall t \in \mathbb{R}^+$  with probability 1

(Gardner et al., 2017). As such, we have that

$\forall r \exists \xi_1(r) | P_r(\mathbf{X}^{(N)}(t) = \mathbf{Y}^{(N)}(t) | t \in [0, T_{\xi_1}]) = 1$  such that  $\xi_1(r)$  is monotone increasing in  $r$  and where  $P_r(A) = P(\mathcal{D}_{M(r)}(A))$  for routing algorithm of  $A$  being  $M$ . Letting  $r \rightarrow \infty \Rightarrow \xi_1 \rightarrow \infty$ , we then have  $\mathbf{X}^{(N)}(t) = \mathbf{Y}^{(N)}(t) \quad \forall t \in \mathbb{R}^+$  in terms of distribution (i.e., as the result holds  $\forall t \in \mathbb{R}^+$  as  $r \rightarrow \infty$ ), implying the required weak convergence for  $\mathcal{D}$  in law over system  $\mathbf{X}^{(N)}(t)$ .  $\square$

In order to evaluate the results of these simulations directly, the work of

(Campbell et al., 2020) provides statistical notions of asymptotic exchangeability in the de Finetti sense by means of quantifying local empirical measure sequences.

**Definition 11** (Local exchangeability).  $\mathbb{X}$  is a locally exchangeable process there exists process  $G_t$  where  $\forall (T \subset \mathbb{R}, \gamma \in (X, \Gamma))$ ,

$$P\left(\bigcap_{t \in T} \{X_t \in A\} | G_t^{(T)}\right) = \prod_{t \in T} G_t \quad P\text{-a.s.}$$

$$\sup_{\omega} E|G_t^{(T)}(\omega) - G_{\gamma t}^{(T)}| \leq \sum_{t \in T} d(t, \gamma(t))$$

Omit these brackets

Insert comma here

where  $G_t^{(T)}$  refers to  $G_t$  restricted to  $T$  and  $d$  is a premetric.

In terms of  $N \rightarrow \infty$ , obviously it is important to consider  $\rho$ . For extremely low values of  $\rho$ , we expect the probability of the threshold being breached at a high  $N$  at any time to approach zero. This is examined, for example, in Figure 3.6. For this image, each value of  $\rho$ ,  $N \in \{2, 12, \dots, 42\}$  is tested and prescribed a bar representing that  $N \times \rho$  combination's

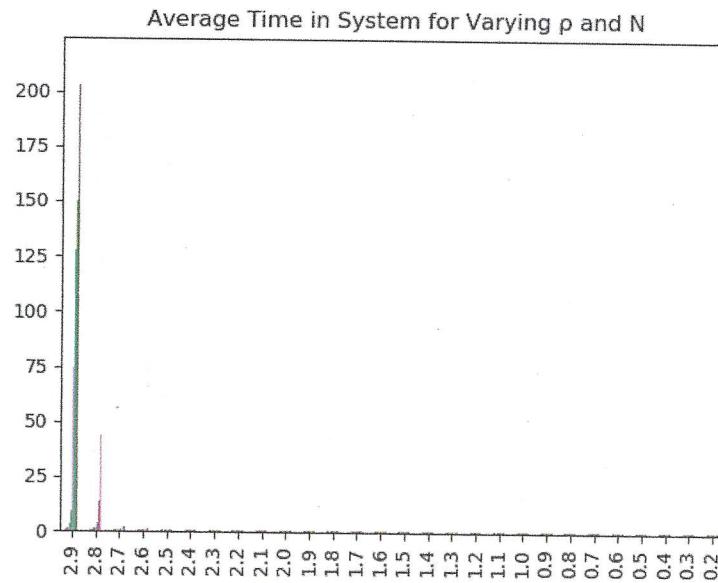
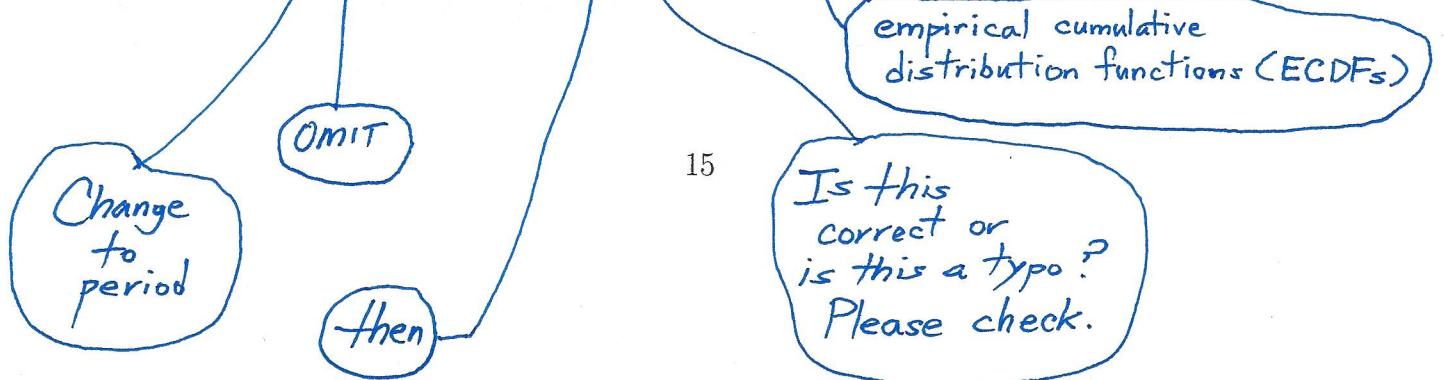


Figure 3.6: Threshold(2,2) for Varying  $N$  and  $\rho$

$E(T)$  until  $t = 1000$  (each combination being run 30 times with independent seeds). Going forward, this suggests that we can examine rather high levels of  $\rho$  to better contrast the difference between redundancy and non-redundancy models. Moreover, as we are evaluating the efficacy of this model for any  $d$  choices, we place particular emphasis on  $d \equiv 2$  given that returns from increasing  $d$  tends to decrease regardless of replicas are or are not being considered [Gardner et al., 2017, Mitzenmacher, 2001].

First, we look at Redundancy(2) under  $\rho = 5$  for varying levels of  $N$ . Under the conjecture, we would expect convergence in  $t$  whence the empirical CDFs over  $t$  no longer change. As Figure 3.7 shows, this indeed seems to be the case when evaluating the CDFs at each time point wherein an event occurs for values up to 10 (which is never surpassed across the simulations). That is, a point mass at 0 would indicate that the random variable seldom changes - a feat which seems to occur in this instance. To visualize the effect of then taking  $t \rightarrow \infty$ , let  $\tilde{\tau}$  be the times at which an arrival or exit occurs (after the first arrival), then Figure 3.8 can be produced by plotting histograms at different event times  $\tilde{\tau}$ .



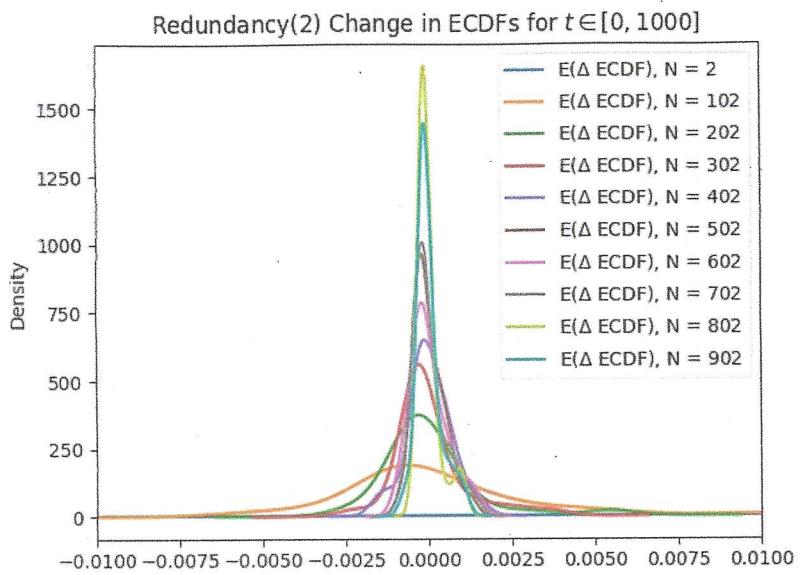


Figure 3.7: Redundancy(2) ECDFs for Varying  $N$

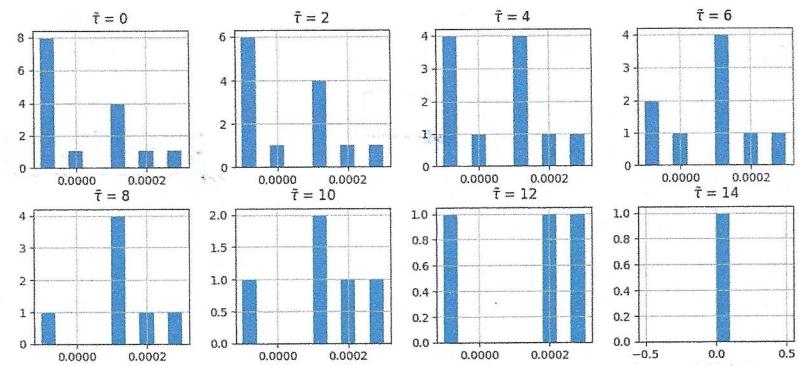


Figure 3.8: Redundancy(2) Histograms of ECDF Changes for Varying  $\tilde{\tau}$  at  $N = 1000$

\* Please look at each of these references and DOUBLE CHECK that you have them written accurately and correctly.  
Please include the NAME OF THE JOURNAL as many of your citations below fail to include this information.

## References

- [Austin, 2008] Austin, T. (2008). On exchangeable random variables and the statistics of large graphs and hypergraphs. 5(0):80–145.
- [Austin, 2015] Austin, T. (2015). Exchangeable random measures. 51(3):842–861.
- [Ayesta et al., 2018] Ayesta, U., Bodas, T., and Verloop, I. (2018). On a unifying product form framework for redundancy models. 127-128:93–119.
- [Bacelli and Brémaud, 2003] Bacelli, F. and Brémaud, P. (2003). *Elements of Queueing Theory*. Springer Berlin Heidelberg.
- [Bramson et al., 2012] Bramson, M., Lu, Y., and Prabhakar, B. (2012). Asymptotic independence of queues under randomized load balancing. 71(3):247–292.
- [Campbell et al., 2020] Campbell, T., Syed, S., Yang, C.-Y., Jordan, M. I., and Broderick, T. (2020). Local exchangeability.
- [Cruise et al., 2020] Cruise, J., Jonckheere, M., and Shneer, S. (2020). Stability of JSQ in queues with general server-job class compatibilities.
- [Gardner et al., 2017] Gardner, K., Harchol-Balter, M., Scheller-Wolf, A., Velednitsky, M., and Zbarsky, S. (2017). Redundancy-d: The power of d choices for redundancy. 65(4):1078–1094.
- [Hellemans et al., 2019] Hellemans, T., Bodas, T., and Van Houdt, B. (2019). Performance analysis of workload dependent load balancing policies. *Proc. ACM Meas. Anal. Comput. Syst.*, 3(2).
- [Mitzenmacher, 2001] Mitzenmacher, M. (2001). The power of two choices in randomized load balancing. *IEEE Trans. Parallel Distrib. Syst.*, 12(10):1094–1104.

- [Mukhopadhyay and Mazumdar, 2015] Mukhopadhyay, A. and Mazumdar, R. R. (2015). Analysis of randomized join-the-shortest-queue (JSQ) schemes in large heterogeneous processor sharing systems.
- [Palmer et al., 2019] Palmer, G. I., Knight, V. A., Harper, P. R., and Hawa, A. L. (2019). Ciw: An open-source discrete event simulation library.
- [Shneer and Stolyar, 2020] Shneer, S. and Stolyar, A. (2020). Large-scale parallel server system with multi-component jobs.

## APPENDICES

\*What intends to  
go here?

Are you planning on  
using this section?

If not, please remove  
it as well as any  
reference to it in the  
Table of Contents.

