

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; Hellemans, 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.

Acknowledgements

Thank you, Dr. Steve Drekić 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 Figures

Chapter 1

Introduction

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 [?].

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 $\gamma_N \in \Gamma_N$,*

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

de Finetti's Theorem states that exchangeability implies the collection to be conditionionally independent (in a Markovian sense) and identically distributed [?]. Moreover, for partition $N = \bigcup_{i \leq k} N_i$, exchangeability over partition $\{N_i\}_{i \leq k}$ such that

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

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

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 [?]. 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.

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, *Threshold(R, d)* has risen to prominence as a means to balance workload in queueing systems.

Definition 4 (Workload and System Load).

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 $\#$ which counts the jobs waiting in a queue at some particular time.

System Load refers to the amount of work remaining in the entire system,

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

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, ψ will henceforth refer to this more general case. As such, ψ will henceforth refer to this more general case.

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

Figure 1.1: Commutativity of Limits

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 ?? describes the behaviour which would be expected in an ideal system wherein both a mean field and asymptotic independence can be achieved [?]. In particular, P describes a fixed “equilibrium” point of the system, a state of the system (in terms of queue-counts) which is consistently held once reached, giving the collection a distribution of δ_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}^ψ 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(\mathbf{R}, \mathbf{d}), denoted $\mathcal{D}_{\text{THRESH}(\mathbf{R}, \mathbf{d}), \mathbf{Z}}$, 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 queues.
2. If $i = 0$, place the original arrival in a queue from the d chosen at random.

\mathbf{Z} refers to any imposed job cancellation cost (e.g., an added temporary workload).

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 $N \longrightarrow \infty$, the system $\mathcal{D}_{THRESH(R,d),Z}$ becomes (ψ, Γ) -exchangeable. As $t \longrightarrow \infty$, 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 [?]. Most importantly, assuming we delegate a probability space $(\Omega, \mathcal{F}_t, \mathcal{F}, P)$ the measurable flow $\{\theta_t\}$ which is P/θ_t invariant (i.e., 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

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

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 a series of random variables defining their behaviour within the system [?]. In general, for the systems we shall consider, we consider the marked processes σ_n as being the required service of arrival n and the marks T_n as its time of arrival. The tuple*

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

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

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

In a simple $G/M/c$ queue, as discussed before, the double $(T, \sigma)_n$ would be sufficient for reducing the problem into a deterministic one. 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 [?] to prove the conjecture in the case of Join-The-Shortest-Queue(d) systems, although with an additional job dependency term.

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, $\sup_{n \in \chi} z_n(t) := \nu(t)$.

“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 by \tilde{X}_t^*) and is such that

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

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

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

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 (i.e., 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, Definition ?? can be represented more succinctly in a hypergraph form, allowing one to represent the cavity process studied in [?].

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 , $j < d$ 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 by \mathcal{G}_t . See Fig ?? for a visual representation.

When one creates the infinite graph iteratively, it is meant that one could view the graph described in Definition ?? 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 [?], we can fully describe the marked process for queues with general service times.

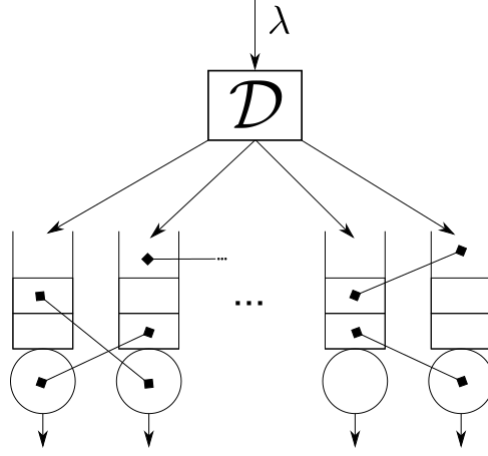


Figure 2.1: A graphical representation of queues with replicated jobs as hyperedges.

Definition 10 (State Space With Redundancy).

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

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

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

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

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 [?]. In order to visualize and analyse the large class of queueing systems within this paradigm [?, ?], 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 time 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 as demonstrated by Figure ??.

Take Figures ?? and ?? for example, which permits one to simulate a Redundancy-2 queueing system with 100 queues in parallel for 1000 units of time while returning the total queue counts over time (which are updated upon a change in queue count).

Remarkably, the simulation itself is performed speedily on consumer hardware despite the size of the system as demonstrated in Figure ??.


```
#!/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.

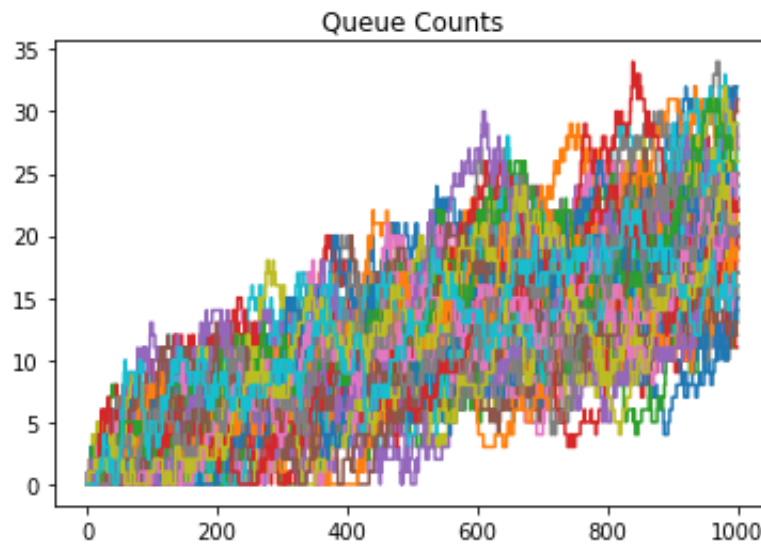


Figure 3.2: Plot using totals of Figure ??

```

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

```

Figure 3.3: Runtime Statistics

Altogether, this makes the package easy to parallelize with and thus to compare systems of different sizes and with large running-times. While currently not implemented in any development branch of *ParallelQueue*, the base Python package *multiprocessing* is used in throughout this paper when simulating for the same system across parameters. In general, the main caveat when processing many models is that the storage of the simulation results can quickly begin to consume storage; when processing many models, therefore ensure that they are saved (e.g., using *pickle*) and removed from the local environment when doing analysis.

In terms of development, the models implemented in the *base_models* module use the framework established in ???. That is, modelling redundancy, a hyperedge of sorts is generated whence the dispatcher receives a job to be cloned. This hyperedge then exists for the duration of time for which the replica class is in the system and is defined in such a way that *Monitor* class objects can interact with them in order to acquire data. In Python, such a data structure can be implemented rather easily by employing the *Dict* type which defines a keyed set of values. By keying based on the job arrivals (before cloning), a unique set of marks can be retrieved for the set by simply using the *Dict* object as a reference.

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 Figure ?? 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

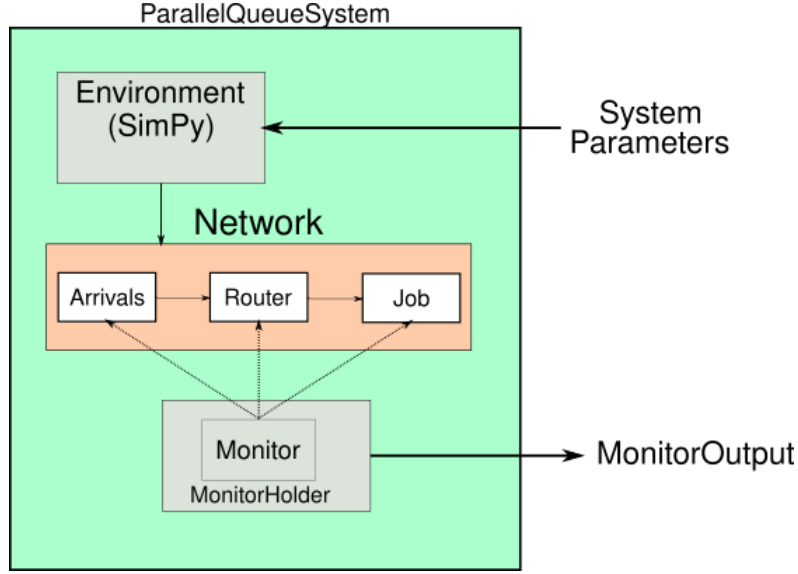


Figure 3.4: Overview of the ParallelQueue API

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

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

As [?] show, a Redundancy(d) system is asymptotically stable if and only if $\rho < 1$. Given that, in premise, Threshold(r, d) 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 ???. 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 $\mathbf{X}^{(N)}$ being a system of queues such that $\rho < 1$,

$$\mathcal{D}_{\text{Thresh}(r,d)}(X^{(N)}) \stackrel{r \rightarrow \infty}{\rightsquigarrow} \mathcal{D}_{\text{Red}(d)}(X^{(N)}).$$

Proof:

Take \leq_{st} to mean that [?]

$$\mathbf{X}_1^{(N)} \leq_{st} \mathbf{X}_2^{(N)} \iff \#X_{1,i} \leq \#X_{2,i} \quad \forall i \in [N] \quad P - a.s.$$

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 [?] . For $\mathbf{X}_1^{(N)}$ under $\text{Thresh}(r, d)$, await the first arrival, denoted T_{ξ_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 [?] 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_{ξ_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 [?]. 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 [?] 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 if and only if there exists process G_t where $\forall T \subset \mathbb{R}, \gamma \in (X, \Gamma)$

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

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

where $G_t^{(T)}$ refers to G_t restricted to T and d is a premetric which can be generated by means of finding the canonical premetric of the process.

Figure 3.6: Threshold(2,2) for Varying N and ρ

In essence, this is formulation provides a means of reducing time-exchangeability – a problem concerning the probabilistic behaviour of a group-theoretic operator – into a problem of analysis. While not necessarily implying exchangeability, local exchangeability provides a sufficient condition upon fulfilment of an additional criterion for full exchangeability: $d(t, t') \xrightarrow{|t-t'| \rightarrow 0} O(|t - t'|^{1+a})$ for $a > 1$.

In terms of $N \rightarrow \infty$, obviously it is important to consider ρ . For extremely low values of ρ , 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 ???. For this image, each value of ρ , $N \in \{2, 12, \dots 42\}$ is tested and prescribed a bar representing that $N \times \rho$ combination's $E(T)$ until $t = 1000$ (each combination being run 30 times with independent seeds). Going forward, this suggests that we can examine rather high levels of ρ 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 whether replicas are or are not being considered [?, ?].

First, we look at Redundancy(2) under $\rho = 2.5$ for varying levels of N . Under the conjecture, we would expect convergence in t whence the ECDFs over t no longer change. As Figure ??? shows, this indeed seems to be the case when evaluating the ECDF 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). Figure ??? can be produced by plotting histograms at different event times $\tilde{\tau}$.

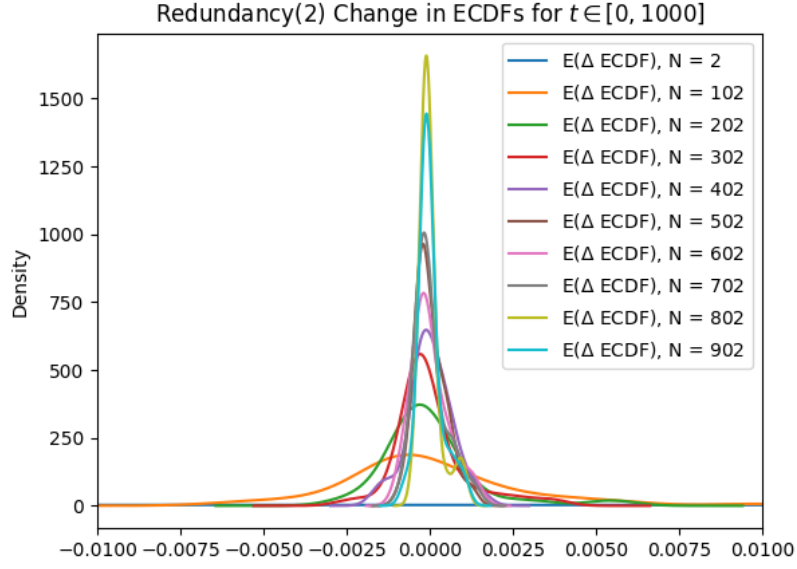


Figure 3.7: Redundancy(2) ECDFs for Varying N

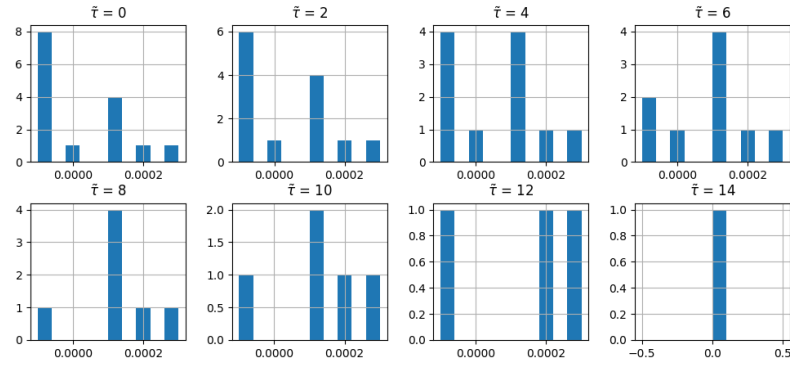


Figure 3.8: Redundancy(2) Histograms of ECDF Changes for Varying $\tilde{\tau}$ at $N = 1000$

APPENDICES