



HAL
open science

The Impact of Consistency on System Latency in Fault Tolerant Internet Computing

Olga Tarasyuk, Anatoliy Gorbenko, Alexander Romanovsky, Vyacheslav Kharchenko, Vitalii Ruban

► **To cite this version:**

Olga Tarasyuk, Anatoliy Gorbenko, Alexander Romanovsky, Vyacheslav Kharchenko, Vitalii Ruban. The Impact of Consistency on System Latency in Fault Tolerant Internet Computing. 15th IFIP International Conference on Distributed Applications and Interoperable Systems (DAIS), Jun 2015, Grenoble, France. pp.179-192, 10.1007/978-3-319-19129-4_15 . hal-01775037

HAL Id: hal-01775037

<https://inria.hal.science/hal-01775037>

Submitted on 24 Apr 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

The Impact of Consistency on System Latency in Fault Tolerant Internet Computing

Olga Tarasyuk¹, Anatoliy Gorbenko¹, Alexander Romanovsky²,
Vyacheslav Kharchenko¹, Vitalii Ruban¹

¹Department of Computer Systems and Networks
National Aerospace University, Kharkiv, Ukraine
{O.Tarasyuk, A.Gorbenko}@csn.khai.edu, V.Kharchenko@khai.edu
²School of Computing Science, Newcastle University, Newcastle upon Tyne, UK
Alexander.Romanovsky@ncl.ac.uk

Abstract. The paper discusses our practical experience and theoretical results in investigating the impact of consistency on latency in distributed fault tolerant systems built over the Internet. Trade-offs between consistency, availability and latency are examined, as well as the role of the application timeout as the main determinant of the interplay between system availability and performance. The paper presents experimental results of measuring response time for replicated service-oriented systems that provide different consistency levels: ONE, ALL and QUORUM. These results clearly show that improvements in system consistency increase system latency. A set of novel analytical models is proposed that would enable quantified response time prediction depending on the level of consistency provided by a replicated system.

Keywords. Internet computing, fault-tolerance, consistency, latency, response time, modelling

1 Introduction

Distributed computing has become an industrial trend, indispensable in dealing with enormous data growth. High availability requirements for many modern Internet applications require the use of system redundancy and data replication. Basic fault tolerant solutions such as N -modular, hot- and cold-spare redundancy usually assume a synchronous communication between replicas, which means that every message is delivered within a fixed and known amount of time [1]. This is a reasonable simplification for the local-area systems whose components are compactly located, for instance, within a single data centre.

This assumption does not appear to be relevant, however, for the wide-area systems, in which replicas are deployed over the Internet and their updates cannot be propagated immediately, which makes it difficult to guarantee consistency.

The Internet and, more generally, the wide-area networked systems are characterized by a high level of uncertainty, which makes it hard to guarantee that a client will receive a response from the service within a finite time. It has been previously shown that there is a significant uncertainty of response time in service-oriented systems invoked over the Internet [2–4]. Besides, our experience and other studies [4–7] show that failures are a regular occurrence on the Internet, clouds and in scale-out data centre networks. When developers apply replication and other fault tolerant techniques in the Internet- and cloud-based systems, they need to understand the time overheads and be concerned about delays and their uncertainty.

In this paper we examine, both in experimental and theoretical terms, how different fault-tolerance solutions [8] implemented over the Internet affect system latency depending on the level of consistency provided. The paper discusses the trade-offs between consistency, availability and latency. Although these relations have been identified by the CAP theorem in qualitative terms [9, 10], it is still necessary to quantify how different fault-tolerant techniques affect system latency depending on the consistency level. The main contributions of the paper are probabilistic models that can predict the system response time depending on the chosen fault-tolerance technique and/or the selected consistency level, with the probabilistic behaviour of replicas as an input parameter.

The rest of the paper is organized as follows. In Section 2 we discuss the impact of the CAP theorem [9, 10] on distributed fault-tolerant systems and examine the trade-offs between system consistency, availability and latency. Section 3 summarises the results of experimental response time measurements for testbed fault-tolerant systems that have three replicas distributed over the Internet and support different consistency levels. The probabilistic models introduced in Section 4 define the relation between system response time and the consistency level provided. Section 5 evaluates the accuracy of the proposed analytical models by applying them in practice and comparing their results with our experimental data. Finally, some practical lessons learnt from our experimental and theoretical work are summarised in Section 6.

2 Understanding Trade-offs Between Consistency, Availability and Latency in Distributed Fault-Tolerant Systems

The CAP conjecture [9], which first appeared in 1998-1999, defines a trade-off between system availability, consistency and partition tolerance, stating that only two of the three properties can be preserved in distributed replicated systems at the same time. Gilbert and Lynch [10] view the CAP theorem as a particular case of a more general trade-off between consistency and availability in unreliable distributed systems which assume that updates are eventually propagated.

System partitioning, availability and latency are tightly connected. A replicated fault-tolerant system becomes partitioned when one of its parts does not respond due to arbitrary message loss, delay or replica failure, resulting in a timeout. System availability can be interpreted as a probability that each client request eventually receives a response.

In many real systems, however, a response that is too late (i.e. beyond the application timeout) is treated as a failure. High latency is an undesirable effect for many interactive web applications. In [13] the authors showed that if a response time increases by as little as 100 ms, it dramatically reduces the probability of the customer continuing to use the system.

Failure to receive responses from some of the replicas within the specified timeout causes partitioning of the replicated system. Thus, partitioning can be considered as a bound on the replica's response time. A slow network connection, a slow-responding replica or the wrong timeout settings can lead to an erroneous decision that the system has become partitioned. When the system detects a partition, it has to decide whether to return a possibly inconsistent response to a client or to send an exception message in reply, which undermines system availability.

The designers of the distributed fault-tolerant systems cannot prevent partitions which happen due to network failures, message losses, hacker attacks and components crashes and, hence, have to choose between availability and consistency. One of these two properties has to be sacrificed. If system developers decide to forfeit consistency they can also improve the system response time by returning the fastest response to the client without waiting for other replica responses until the timeout, though this would increase the probability of providing inconsistent results. Besides, timeout settings are also important. If the timeout is lower than the typical response time, a system is likely to enter the partition mode more often [11].

It is important to remember that none of these three properties is binary. For example, modern distributed database systems, e.g. Cassandra [14], can provide a discrete set of different consistency levels for each particular read or write request. The response time can theoretically vary between zero and infinity, although in practice it ranges between a minimal affordable time higher than zero and the application timeout. Availability varies between 0% and 100% as usual.

The architects of modern distributed database management systems and large-scale web applications such as Facebook, Twitter, etc. often decide to relax consistency requirements by introducing asynchronous data updates in order to achieve higher system availability and allow a longer response time. Yet the most promising approach is to balance these properties. For instance, the Cassandra NoSQL DDBS introduces a tunable replication factor and an adjustable consistency model so that a customer can choose a particular level of consistency to fit with the desired system latency.

The CAP theorem helps the developers to understand the system trade-offs between consistency and availability/latency [12]. Yet even though this theorem strongly suggests that better consistency undermines system availability and latency, developers do not have quantitative models to help them to estimate the system response time for the chosen consistency level and to achieve a precise trade-off between them.

Our interpretation of the CAP theorem and the trade-offs resulting from the CAP is depicted in Fig. 1. The application timeout can be considered as a bound between system availability and performance (in term of latency or response time) [15]. Thus, system designers should be able to set up timeouts according to the desired system response time, also keeping in mind the choice between consistency and availability.

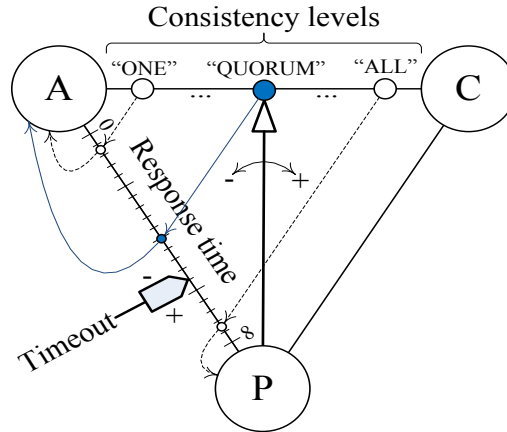


Fig. 1. The CAP trade-offs.

In the following sections we discuss our practical experience on measuring latency of fault-tolerant service-oriented system depending on the provided consistency level and also introduce analytical models predicting system response time.

3 Experimental Investigation of the CAP Impact on Fault-Tolerant Service-Oriented Systems

3.1 Description of the Testbed Architecture

To investigate the CAP impact on fault-tolerant distributed systems we developed a testbed service-oriented system composed out of the three replicated web services (see Fig. 2). This is a typical setup employed in many fault-tolerant solutions.

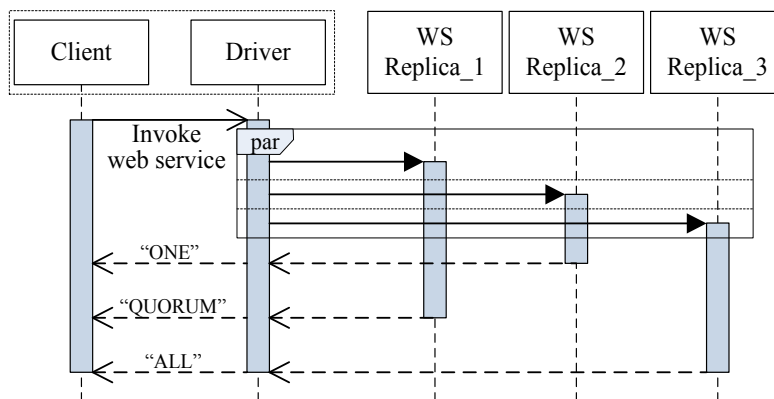


Fig. 2. Fault-tolerant service-oriented system.

A testbed web service was written in Java and its replicas uploaded to Amazon Elastic Beanstalk and were deployed in the three different location domains: (i) US West (Oregon); (ii) South America (Sao Paulo) and Asia Pacific (Tokyo). Each web service replica performs a heavy-computational arithmetic calculation such as finding the n digit of Pi when n is a large number and returns the result to the driver. The driver is responsible for invoking each of the replicated web services, waiting for the web services to complete their execution and return response, and, finally, implementing a particular fault-tolerant scheme upon the obtained results.

AWS SDK for Java was used to connect web service replicas on Amazon EC2 from clients (driver) programming code that helps to take the complexity out of coding by providing Java APIs for AWS services.

In our study we investigated the three basic fault-tolerant patterns for web services [16] corresponding to different consistency levels (ONE, ALL, QUORUM). In all cases the driver simultaneously forwards client's request to all replicated web services. The consistency level determines the number of replicas which must return a response to the driver before it sends an adjudicated result to the client application:

- ONE (*hot-spare redundancy*) – when the FASTEST response is received the driver forwards it to the client. This is the weakest consistency level though it guarantees the minimal latency;
- ALL (*N-modular redundancy*) – the driver must wait until ALL replicas return their responses. In this case the response time is constrained by the slowest replica though the strongest consistency is provided;
- QUORUM – the driver must wait for the responses from a QUORUM of replica web services. It provides a compromise between the ONE and ALL options trading off latency versus consistency. The quorum is calculated as: $(\text{amount_of_replicas} / 2) + 1$, rounded down to an integer value. As far as in our experiments we use the replication factor of 3, the quorum is 2.

The driver also implements a timeout mechanism aimed to protect clients from endless waiting in case of network or web-services failures or cloud outages.

3.2 Response Time Measurement

The driver was implemented as part of the Java client software. The client software was run at a host in the Newcastle University (UK) corporate network. It invoked replica web services several thousand times in a loop using the driver as a proxy.

For the particular client's request we measured the response time of the each web service replica and also times when the driver produces responses corresponding to different consistency levels. The delay induced by the driver itself was negligible in our experiments.

The measurement results obtained for the first 100 invocations are presented in Figs. 3 and 4. Table 1 summarizes basic statistical characteristics of the measured data whereas probability density series (*pds*) of system and replicas response times are depicted in Figs. 5 and 6.

As expected, when the system is configured to provide consistency level ONE its latency in average is less than the average response time of the fastest replica. Average system latency in case it provides consistency level ALL is larger than the average response time of the slowest replica. System latency associated with consistency level QUORUM is in the middle.

However, our main observation is that it is hardly possible to make an accurate prediction of the average system latency corresponding to the certain consistency level when the only common statistical measures of replicas response time (i.e. minimal, maximal and average estimates and standard deviation) are known.

This finding resulting from our massive experiments and also confirmed by other researches [17] show that it is extremely difficult to predict the timing characteristics of various types of wide-area distributed systems, including fault-tolerant SOAs, distributed databases and file systems (e.g. Cassandra, GFS, HDFS), parallel processing systems (e.g. Hadoop Map-Reduce). The dynamic and changing nature of timing characteristics of such systems can be better captured by employing probability density functions.

In the next section we propose a probabilistic modelling approach that addresses this problem. It relies on using probability density functions (PDF) of replica response times to predict system latency at different consistency levels.

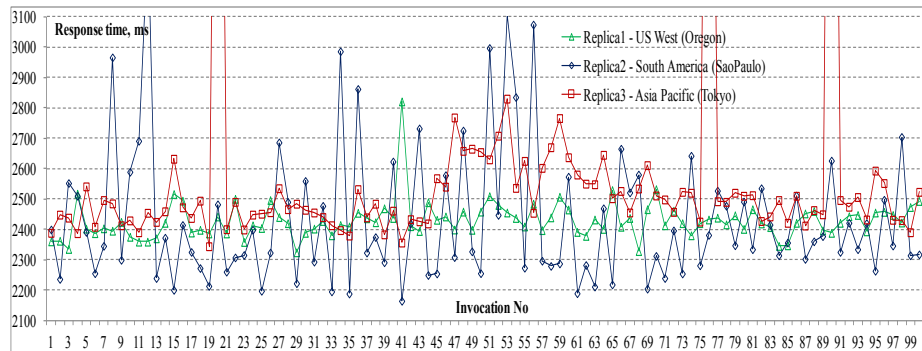


Fig. 3. Response time of different web service replicas.

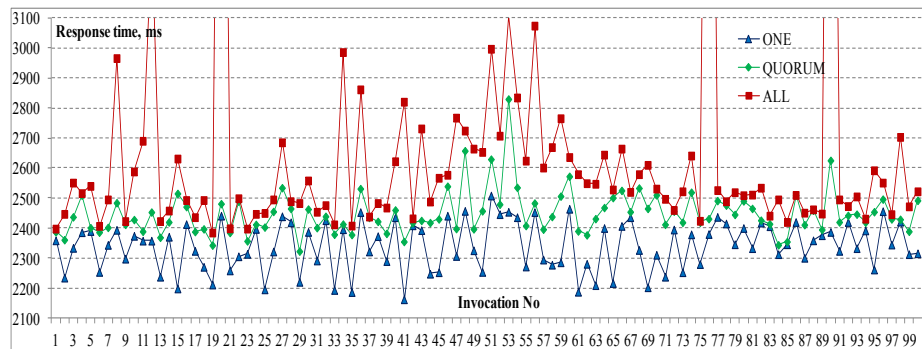


Fig. 4. System response time corresponding to different consistency levels.

Table 1. Response time statistics.

Response Time, ms	Replica1	Replica2	Replica3	System consistency level		
	(Oregon)	(Sao Paulo)	(Tokyo)	ONE	QUORUM	ALL
Minimal	2324	2164	2344	2164	2324	2386
Average	2428	2434	2588	2342	2449	2660
Maximal	2821	3371	5573	2509	2830	5573
Std. deviation	60	228	522	80	72	529

4 Probabilistic Models of System Response Time for Different Consistency Levels

We propose a set of probabilistic models that allow us to build a combined probability density function of system response time by taking into account provided consistency level and incorporating response time probability density functions for each replica.

When the system is configured to provide consistency level ALL, the probability of returning response to the client at time t is equal to the probability that one of the replicas (e.g. the first one) returns its response exactly at time t , i.e. $g_1(t)$ while two other replicas return their responses not later than t (by time t), i.e. $\int_0^t g_2(t) = G_2(t)$ and $\int_0^t g_3(t) = G_3(t)$.

So far as we have three replicas, all three possible combinations have to be accounted. As a result, the probability density function of the system response time for consistency level ALL can be defined as following:

$$f_{ALL}(t) = g_1(t)G_2(t)G_3(t) + g_2(t)G_1(t)G_3(t) + g_3(t)G_1(t)G_2(t). \quad (1)$$

where $g_1(t)$, $g_2(t)$ and $g_3(t)$ – are response time probability density functions of the first, second and third replicas respectively; $G_1(t)$, $G_2(t)$ and $G_3(t)$ – are response time cumulative distribution functions of the first, second and third replicas respectively.

When the system is configured to provide consistency level ONE, the probability of returning a response to the client at time t is equal to the probability that if only one of the replicas (e.g. the first one) returns its response exactly at time t , i.e. $g_1(t)$, while two other replicas return their responses at the same time or later on, i.e. $\int_t^\infty g_2(t) = 1 - G_2(t)$ and $\int_t^\infty g_3(t) = 1 - G_3(t)$.

Keeping in mind three possible combinations we can deduce the probability density function of the system response time for consistency level ALL as:

$$f_{ONE}(t) = g_1(t)(1 - G_2(t))(1 - G_3(t)) + g_2(t)(1 - G_1(t))(1 - G_3(t)) + g_3(t)(1 - G_1(t))(1 - G_2(t)). \quad (2)$$

Deducing the response time probability density function for the QUORUM consistency level is based on a combination of the previous two cases.

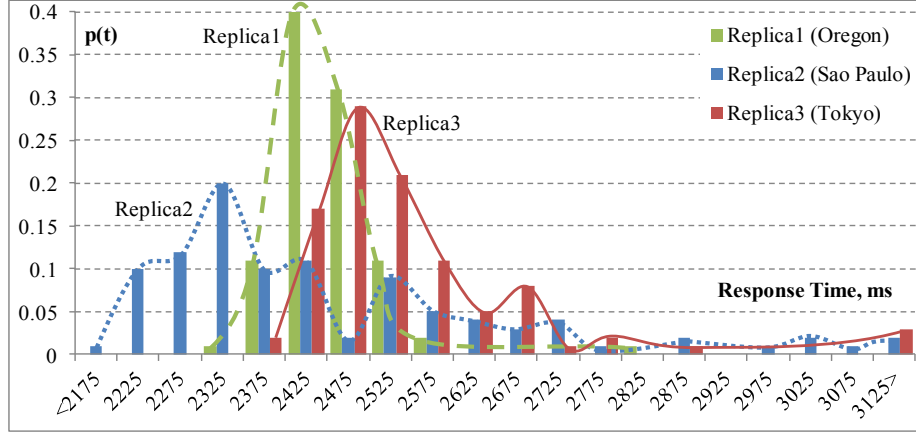


Fig. 5. Probability density series of replicas response times.

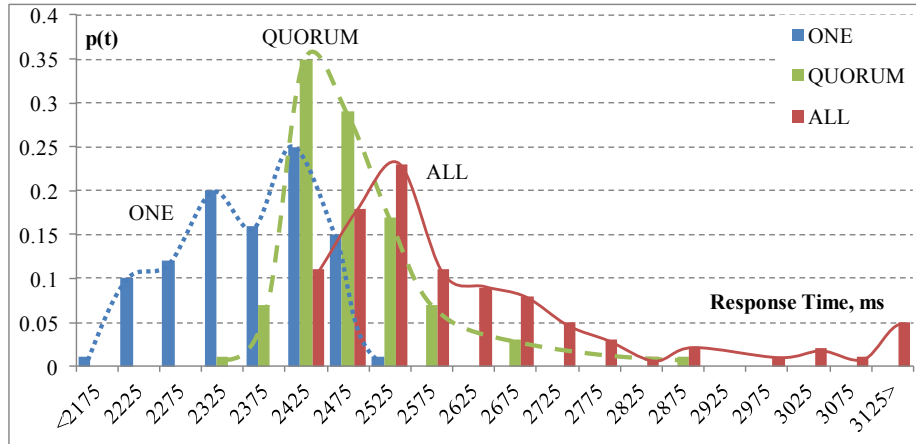


Fig. 6. Probability density series of system response time for different consistency levels.

The probability of returning response to the client at time t is equal to the probability that one of the replicas returns its response exactly at time t ; one of the two remained replicas returns its response by time t and another one responds at time t or later on. Taking into account all possible combinations the probability density function of the system response time for consistency level QUORUM can be deduced as:

$$\begin{aligned}
 f_{QUORUM}(t) = & (g_1(t)G_2(t) + g_2(t)G_1(t))(1 - G_3(t)) + \\
 & + (g_1(t)G_3(t) + g_3(t)G_1(t))(1 - G_2(t)) + \\
 & + (g_2(t)G_3(t) + g_3(t)G_2(t))(1 - G_1(t)).
 \end{aligned} \tag{3}$$

Using similar reasoning it is possible to deduce response time probability density functions of a system composed of n replicas:

$$f_{ALL}(t) = \sum_{i=1}^n \left(\frac{g_i(t)}{G_i(t)} \cdot \prod_{j=1}^n G_j(t) \right). \quad (4)$$

$$f_{ONE}(t) = \sum_{i=1}^n \left(\frac{g_i(t)}{1 - G_i(t)} \cdot \prod_{j=1}^n (1 - G_j(t)) \right). \quad (5)$$

It is extremely hard to build a general form of the probability density function of the system response time for consistency level QUORUM. However, the general reasoning is as following. The composed probability density function should be presented as a sum of m items, where m is a number of k -combinations of n (k is a number of replicas constituting a quorum). Each of the m items is a product of two factors. The first one defines the probability that a particular combination of k replicas return their responses by time t . Another factor defines the probability that the remaining $(n-k)$ replicas return their responses after t .

5 Models Validity

In this section we check the validity and accuracy of the proposed models by comparing their prediction with the experimental data presented in Section 3. This check includes the following four steps:

- finding out theoretical distribution laws that accurately approximate the measured replica response times;
- applying the proposed mathematical models (1), (2) and (3) to deduce probability density functions of the system response time for different consistency levels;
- estimating replica and system average response times using the theoretical probability distribution functions;
- comparing the theoretical and experimental values of replica and system average response times.

5.1 Finding Theoretical Distribution Laws of Replica Response Times

Theoretical distribution laws approximating replica response times can be found in a way described in [2]. It is based on performing a series of hypotheses checks in the Matlab numeric computing environment. The techniques of hypothesis testing consist of the two basic procedures. First, the values of distribution parameters are estimated by analysing an experimental sample. Second, the null hypothesis that experimental data has a particular distribution with certain parameters should be tested.

To perform hypothesis testing itself we used the `kstest` function: `[h, p] = kstest(t, cdf)`, conducting the Kolmogorov-Smirnov test to compare the distribution of t with the hypothesized distribution defined by matrix *cdf*.

The null hypothesis for the Kolmogorov-Smirnov test is that t has a distribution defined by *cdf*. The alternative hypothesis is that x does not have that distribution.

Result h is equal to '1' if we can reject the hypothesis, or '0' if we cannot. The function also returns the p -value which is the probability that x does not contradict the null hypothesis. We reject the hypothesis if the test is significant at the 5% level (if p -value is less than 0.05). The p -value returned by `kstest` was used to estimate the goodness-of-fit of the hypothesis. As a result of hypothesis testing we found out that the *Weibull* distribution fits well the response time of the first (Oregon) and the third (Tokyo) replicas. The response time of the second replica (Sao Paulo) can be accurately approximated by the *Gamma* distribution.

5.2 Deducing Probability Density Functions of the System Response Time

Mathcad has been used at the second stage of our investigation to deduce theoretical distributions of system response times for different consistency levels. It also allows to estimate average system latency and to plot probability density functions. Mathcad worksheet is shown in Fig. 7. It includes seven modelling steps.

At the 1st step we define abscissa axis t and its dimension in milliseconds. Secondly, we set up parameters of replicas response time distribution functions estimated in Matlab and also their shifts on the abscissa axis (i.e. minimal response time values).

At the 3rd and 4th steps the replica response time probability density functions $g_1(t)$, $g_2(t)$, $g_3(t)$ and the corresponding cumulative distribution functions $G_1(t)$, $G_2(t)$, $G_3(t)$ are defined using Mathcad library functions `dweibull` and `dgamma`.

At the 5th step we define probability density functions of the system response time corresponding to different consistency levels by combining replicas *pdf* and *cdf* according to the proposed equations (1), (2) and (3).

1	$t := 2000, 2010.. 3000$		
2	$a1 := 113.3578$ $b1 := 2.3041$ $min1 := 2324$	$a2 := 1.5952$ $b2 := 164.1599$ $min2 := 2164$	$a3 := 176.8796$ $b3 := 1.7467$ $min3 := 2344$
3	$g1(t) := \frac{1}{a1} \cdot dweibull\left[\frac{(t - min1)}{a1}, b1\right]$	$g2(t) := \frac{1}{b2} \cdot dgamma\left[\frac{(t - min2)}{b2}, a2\right]$	$g3(t) := \frac{1}{a3} \cdot dweibull\left[\frac{(t - min3)}{a3}, b3\right]$
4	$G1(t) := \int_0^t g1(t) dt$	$G2(t) := \int_0^t g2(t) dt$	$G3(t) := \int_0^t g3(t) dt$
5	$fALL(t) := g1(t) \cdot G2(t) \cdot G3(t) + g2(t) \cdot G1(t) \cdot G3(t) + g3(t) \cdot G1(t) \cdot G2(t)$ $fONE(t) := g1(t) \cdot (1 - G2(t)) \cdot (1 - G3(t)) + g2(t) \cdot (1 - G1(t)) \cdot (1 - G3(t)) + g3(t) \cdot (1 - G1(t)) \cdot (1 - G2(t))$ $fQUORUM(t) := (g1(t) \cdot G2(t) + g2(t) \cdot G1(t)) \cdot (1 - G3(t)) + (g1(t) \cdot G3(t) + g3(t) \cdot G1(t)) \cdot (1 - G2(t)) + (g2(t) \cdot G3(t) + g3(t) \cdot G2(t)) \cdot (1 - G1(t))$		
6	$\int_0^{10000} t \cdot g1(t) dt = 2.424 \times 10^3$	$\int_0^{10000} t \cdot g2(t) dt = 2.426 \times 10^3$	$\int_0^{10000} t \cdot g3(t) dt = 2.502 \times 10^3$
7	$\int_0^{10000} t \cdot fALL(t) dt = 2.567 \times 10^3$	$\int_0^{10000} t \cdot fONE(t) dt = 2.341 \times 10^3$	$\int_0^{10000} t \cdot fQUORUM(t) dt = 2.444 \times 10^3$

Fig. 7. Mathcad's worksheet.

Probability distribution functions of replicas and system response times are shown in Figs. 8 and 9. The bulk of the values of probability density function $f_{ALL}(t)$ is shifted to the right on the abscissa axis as it was expected. The shapes of the $f_{ONE}(t)$ and $f_{QUORUM}(t)$ probability density functions are also in line with the reasonable expectations and experimentally obtained probability density series (see Fig. 6).

Finally, at steps 6 and 7 we estimate the system and replicas average response time by integrating their theoretical probability distribution functions.

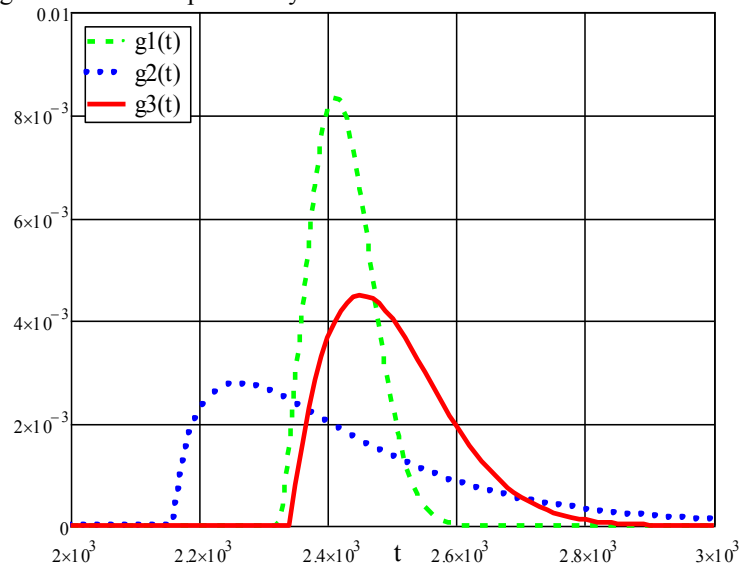


Fig. 8. Probability density functions of replicas response times.

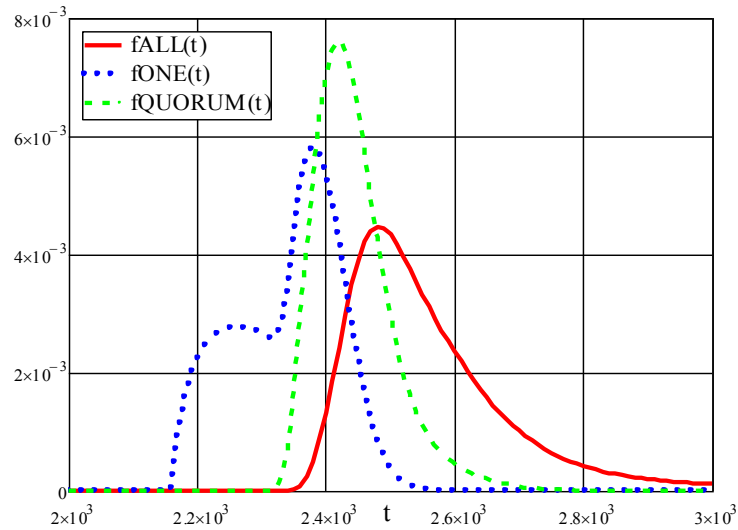


Fig. 9. Probability density functions of system response time for different consistency levels.

5.3 Accuracy of Mathematical Modelling

Table 2 shows the deviation between the average values of the system and replicas response time estimated practically (see Table 1) and theoretically with the help of the obtained probability distribution functions. These results confirm the significant closeness between actual and modelled timing characteristics. To be sure that not only the average value can be accurately predicted we compare theoretical system probability density functions (see Fig. 9) and practically obtained probability density series (Fig. 6). With this purpose we estimated experimental and theoretical probabilities that system latency at different consistency levels is less than the specified time.

Table 2. Accuracy of mathematical modelling.

	Replica1 (Oregon)	Replica2 (Sao Paulo)	Replica3 (Tokyo)	System consistency level		
				ONE	QUORUM	ALL
Approximating theoretical distributions and their parameters						
distribution	Weibull	Gamma	Weibull			
alpha	113.3578	1.5952	176.8796			
beta	2.3041	164.1599	1.7467			
x-shift	2324	2164	2344			
Average response time, ms						
measured	2428	2434	2588	2342	2449	2660
modelled	2424	2426	2502	2341	2444	2567
Deviation, %	0.18	0.34	3.32	0.03	0.19	3.51

Table 3. Deviation between theoretical system *pdf* and *pds* obtained experimentally.

Time, ms	Probability that system latency is less than the specified time								
	ONE			QUORUM			ALL		
	pds	pdf	dev.,%	pds	pdf	dev.,%	pds	pdf	dev.,%
2175	0.01	0.009	10.00	0	0	-	0	0	-
2225	0.11	0.116	5.45	0	0	-	0	0	-
2275	0.23	0.252	9.57	0	0	-	0	0	-
2325	0.43	0.385	10.47	0.01	0	-	0	0	-
2375	0.59	0.596	1.02	0.08	0.097	21.25	0	0.003	-
2425	0.84	0.858	2.14	0.43	0.434	0.93	0.11	0.073	33.64
2475	0.99	0.975	1.52	0.72	0.752	4.44	0.29	0.263	9.31
2525	1	0.998	0.20	0.89	0.903	1.46	0.52	0.476	8.46
2575	1	1	0	0.96	0.961	0.10	0.63	0.643	2.06
2625	1	1	0	0.96	0.984	2.50	0.72	0.761	5.69
2675	1	1	0	0.99	0.994	0.40	0.8	0.841	5.13
2725	1	1	0	0.99	0.998	0.81	0.85	0.892	4.94
2775	1	1	0	0.99	0.999	0.91	0.88	0.924	5.00
2825	1	1	0	0.99	1	1.01	0.89	0.945	6.18
2875	1	1	0	1	1	0	0.91	0.959	5.38
2925	1	1	0	1	1	0	0.91	0.969	6.48
2975	1	1	0	1	1	0	0.92	0.977	6.20
3025	1	1	0	1	1	0	0.94	0.982	4.47
3075	1	1	0	1	1	0	0.95	0.987	3.89
Average deviation, %			2.12			2.25			7.12

The results of this comparison (see Table 3) show a close approximation of the experimental data by the proposed analytical models, especially for the consistency levels ONE and QUORUM. The probabilistic model of the system response time for consistency level ALL gives slightly optimistic prediction, though the average deviation from the experimental data is only 7% – that is close enough.

6 Conclusion and Lessons Learnt

When employing fault-tolerance techniques over the Internet and clouds, engineers need to deal with delays, their uncertainty, timeouts, adjudication of asynchronous replies from replicas, and other specific issues involved in global distributed systems. The overall aim of this work was to study the impact of consistency on system latency in fault tolerant Internet computing.

Our experimental results clearly show that improving system consistency makes system latency worse. This finding confirms one of the generally accepted qualitative implications of the CAP theorem [9, 10]. However, so far system developers have not had any mathematical tools to help them to accurately predict the response time of large-scale replicated systems. While estimating the system worst-case execution time remains common practice for many applications (e.g. embedded computer systems, server fault-tolerance solutions, like STRATUS, etc.), this is no longer a viable solution for the wide-area service-oriented systems in which components can be distributed all over the Internet. In our previous works [2, 3] we demonstrated that extreme unpredictable delays exceeding the value of ten average response times can happen in such systems quite often. In this paper we have proposed a set of novel analytical models providing a *quantitative basis* for the system response time prediction depending on the consistency level provided for (or requested by) clients. The models allow us to derive the probability distribution function of the system response time which corresponds to a particular consistency level (ONE, ALL or QUORUM) by incorporating the probability density functions of the replica response times.

The validity of the proposed models has been verified against the experimental data reported in Section 3. It has been demonstrated that the proposed models ensure a significant level of accuracy in the system average response time prediction, especially in case of ONE and QUORUM consistency levels. The proposed models provide a mathematical basis for predicting latency of distributed fault and intrusion-tolerance techniques operating over the Internet. The models take into account the probabilistic uncertainty of replicas' response time and the required consistency level.

The practical application of our work is in allowing practitioners to predict system performance, and in offering them crucial support for the optimal timeout setup and for understanding the trade-off between system consistency and latency. Trading off system consistency against latency requires the knowledge of probability density functions (and parameter values) that accurately approximate replicas' response time. These probabilistic characteristics, which can be obtained by testing or during the trial usage, will need to be corrected at run-time or at tune-time to improve prediction accuracy. It would be possible to replace the response time probability density functions in the proposed models with probability density series. This would make it easier to use the models in practice.

7 Acknowledgements

We are grateful to Aad van Moorsel for his feedback on the earlier version of this work and Batyrkhan Omarov for his help with running some of the experiments. Alexander Romanovsky is partially supported by the EPSRC TRAMS-2 platform grant.

8 References

1. Lee, P.A., Anderson, T.: *Fault Tolerance. Principles and Practice*. Springer-Verlag (1990)
2. Gorbenko, A. et al.: Real Distribution of Response Time Instability in Service-Oriented Architecture. In: 29th IEEE Int'l Symp. Reliable Distributed Systems, 92–99 (2010)
3. Gorbenko, A. et al.: Exploring Uncertainty of Delays as a Factor in End-to-End Cloud Response Time. In: 9th European Dependable Computing Conference, 185–190 (2012)
4. Omar Bakr, Idit Keidar: Evaluating the running time of a communication round over the internet. In: 21th Ann. ACM Symposium on Principles of Distributed Computing (PODC'00), 243–252 (2002)
5. Chen, Y. et al.: Measuring and Dealing with the Uncertainty of the SOA Solutions. In: *Performance and Dependability in Service Computing: Concepts, Techniques and Research Directions*, V. Cardellini et al., eds., IGI Global, 265–294 (2011)
6. Potharaju, R., Jain, N.: When the Network Crumbles: An Empirical Study of Cloud Network Failures and their Impact on Services. In: 4th ACM Symposium on Cloud Computing (SoCC), (2013)
7. Scott, C., Choffnes, D.R., Cunha, I., etc. LIFEGUARD: practical repair of persistent route failures. In: ACM SIGCOMM 2012 conference on Applications, technologies, architectures, and protocols for computer communication, 395–406 (2012)
8. Avizienis, A., Laprie, J.-C., Randell, B., Landwehr, C.: Basic Concepts and Taxonomy of Dependable and Secure Computing. *IEEE Trans. on Dependable and Secure Computing*, 1(1), 11–33 (2004)
9. Brewer, E.: Towards Robust Distributed Systems. In: 19th Ann. ACM Symposium on Principles of Distributed Computing (PODC'00), 7–10 (2000)
10. Gilbert, S., Lynch, N.: Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services. *ACM SIGACT News*. 33(2), 51–59 (2002)
11. Gorbenko, A., Romanovsky, A., Tarasyuk, O., Kharchenko, V. Dependability of Service-Oriented Computing: Time-Probabilistic Failure Modelling. In Avgeriou, Paris (Ed.). *Software Engineering for Resilient Systems*, LNCS, Vol. 7527/2012, 121–133 (2012)
12. Abadi, D.J.: Consistency Tradeoffs in Modern Distributed Database System Design. *IEEE Computer*. 45(2), 37–42 (2012)
13. Brutlag, J.: Speed Matters for Google Web Search. Google, http://services.google.com/fh/files/blogs/google_delayexp.pdf (2009)
14. Lakshman, A., Malik, P.: Cassandra: a decentralized structured storage system. *ACM SIGOPS Operating Systems Review*. 44(2), 35–40 (2010)
15. Gorbenko, A., Romanovsky, A.: Time-Outing Internet Services. *IEEE Security & Privacy*. 11(2), 68–71 (2013)
16. Gorbenko, A., Kharchenko, V., Romanovsky, A.: Using Inherent Service Redundancy and Diversity to Ensure Web Services Dependability. In: *Methods, Models and Tools for Fault Tolerance*, LNCS 5454, Springer, 324–341 (2009)
17. Rao, J., Shekita, E.J., Tata, S.: Using Paxos to Build a Scalable, Consistent, and Highly Available Datastore. In: *VLDB Endowment*, 243–254 (2011)