

Availability of Globally Distributed Nodes: An Empirical Evaluation*

Timo Warns, Christian Storm, Wilhelm Hasselbring
Carl von Ossietzky University of Oldenburg
Department of Computer Science
26111 Oldenburg, Germany

{timo.warns, christian.storm, hasselbring}@informatik.uni-oldenburg.de

Abstract

Dependability models of distributed systems are often parameterised by the failure characteristics of the nodes that form a system. For realistic results, these parameters must be estimated accurately, for example, based on evaluations of real-world systems. We empirically evaluate over 400 globally distributed nodes of the PlanetLab research cluster and estimate the popular parameters mean-time-to-failure, mean-time-to-repair, availability, and failure correlation coefficients. We fit the resulting empirical distributions by simple theoretical distributions and find that the mean-time-to-failure, the availability, and the failure correlation coefficient correlate with the geographical distance between nodes.

1. Introduction

Popular parameters of dependability models for distributed systems are the *availability*, the *time-to-failure* (TTF), and the *time-to-repair* (TTR) of the individual nodes. Intuitively, the availability of a node gives the fraction of time that the node operates as expected. The TTF and TTR give the distributions of how long a node operates as expected and how long it takes to restore this mode of operation after the occurrence of a failure. For realistic results, such parameters must be estimated accurately, for example, by empirical evaluations of real-world systems.

For reducing complexity, dependability models abstract away details of real-world systems. For example, the TTF and TTR are often assumed to be exponentially distributed and summarised by their means, the *mean-time-to-failure* (MTTF) and *mean-time-to-repair* (MTTR). Further simplifications are that failures are independent and that each node has the same MTTF and MTTR. While such abstractions make models

more tractable, some accuracy is lost and results become less realistic. For higher accuracy, less abstract models are necessary that, for example, consider that the MTTF and the MTTR are not single-valued but follow certain distributions. For tractable models, conceptually simple distributions are desirable raising the challenge to find well-suited simple distributions.

If an assessment of a model reveals that the underlying assumptions of the model are not accurate, new models with better abstractions must be found. For example, the interest in models for dependent failures raised when empirical evaluations indicated that failures are not independent in many distributed systems [1, 3, 9]. Such models often rely on the failure correlation coefficients as input parameters [2, 9] or explicitly model the causes of dependent failures. For example, the workload of nodes and shared resources were identified as causes for dependent failures and added as model parameters. Such parameters may be abstract in the sense that a single abstract cause models several real-world causes of dependent failures.

In this paper, we make the following contributions. We empirically estimate the popular model parameters MTTF, MTTR, availability, and failure correlation coefficient for a large number of globally distributed nodes. We fit the resulting empirical distributions against simple theoretical distributions and assess their goodness-of-fit. We confirm previous findings that the TTF and TTR distributions are unlikely to be exponentially distributed and find more likely simple theoretical distributions. We also confirm that failures are not independent and find that the geographical distance of nodes is a good candidate for an abstract cause of dependent failures.

Our results are based on data provided by the CoMon [7] monitoring system from the Network Systems Group of the Princeton University. CoMon continuously monitors all nodes running in PlanetLab [8], a global cluster of currently more than 800 nodes that are distributed over more than 35 countries.

*This work is supported by the German Research Foundation (DFG), grant GRK 1076/1.

2. Related Work

Numerous empirical studies have estimated basic failure characteristics such as the MTTF and the MTTR of nodes in distributed systems [1, 3, 6, 12]. For example, Long et al. [6] evaluated about 1,100 Internet hosts and found that, on average, a node has an MTTF of 29.39 days and an MTTR of 3.88 days. Our estimations of these characteristics are not conceptually new, but are based on recent data, which is required for dependability models of current systems. Furthermore, we evaluate data from an observation period of two years, which is significantly longer than the observation periods of previous studies. Longer periods of observation bear the advantage of increased accuracy as, generally, estimations become better with more samples. Additionally, we go beyond estimating failure characteristics and fit the empirical distributions to simple theoretical ones.

Different empirical studies have assessed the dependence of failures [1, 3, 9, 12]. For example, Amir and Wool [1] and Bakkaoglu et al. [2] studied globally distributed nodes and find that the correlation of their failures is significant, which raised the interest in how to model dependent failures [2, 9] and how to tolerate them [1, 10, 11]. Many dependent failure models rely on the failure correlation coefficient of nodes as input parameter or explicitly model the causes for dependent failures. In this paper, we estimate the failure correlation coefficient parameter, identify a new abstract cause of dependent failures, the geographical distance of nodes, and quantify its impact.

Our evaluations are based on monitoring data from a single observer node and, therefore, yield failure characteristics as perceived by the observer node. Such an approach does not mask failures of the communication infrastructure, but attributes them to the nodes that are affected. A similar approach has been taken, for example, by Bakkaoglu et al. [2], who monitored about 100 web servers by a central node that periodically retrieved web pages from each server. Other approaches as, for example, followed by Long et al. [6] rely on distributed monitoring systems, which mask some communication failures. Their estimations for a node are, therefore, more close to the real values, but possibly differ significantly from what is perceived by a node in the system. Which approach is favourable depends on the dependability model, whose parameters are estimated.

3. Interpretation of Monitoring Entries

PlanetLab is a globally distributed research cluster that is used as a testbed for deploying and evaluating large-scale distributed systems. Its nodes are

monitored by the CoMon [7] monitoring system. Each node runs a *node-centric daemon* that measures node-specific metrics such as the current CPU and memory utilisation. A central node, the *data collector*, queries all node-centric daemons once every 5 minutes in parallel. Most daemons respond within a second or do not respond at all; successful queries with response times of more than twenty seconds are rare [7].

The data collector records each query to a node-centric daemon by a monitoring entry that contains information about the query itself (e.g., the start time of the query and the name of the queried node) as well as the results from the node-centric daemon. If a query fails, the data collector records an entry only containing information about the failed query (e.g., the start time of the query and the name of the queried node).

As CoMon is a non-trivial system, it is not completely free of faults. Unfortunately, some faults manifest themselves in faulty monitoring entries. We do not try to repair the faulty entries, but reject them from further evaluation, because many of them are too crippled to deduce valid information even by manual inspection. We evaluate the monitoring entries from the period of July 2005 to July 2007 with overall 143,750,690 entries of which 19,860 entries are faulty (i.e., $< 0.014\%$).

Based on the monitoring entries, we consider a node as available if the data collector is able to successfully query the node's node-centric daemon. More precisely, we interpolate between query times and consider a node *available at a time t* iff (if and only if) the most recent query of the node before t was successful. We say that a node *fails at time t* iff the most recent query before t was successful, but the query at time t failed. Analogously, a node is *repaired at time t* iff the most recent query before t failed, but the query at time t was successful. As the PlanetLab cluster is not a completely static system, we additionally need to decide when a node joins and leaves the system. We say that a node *n joined the system at time t* iff n was queried for the first time at t . It *left the system at time t* iff n was queried the last time at t .

4. Estimation of Dependability Metrics

With the interpretation of the monitoring entries, we determine the state of each queried node over time. Figure 1 shows the number of nodes being monitored and being available over time. Overall, 979 nodes have been queried at least once by the data collector during the observation period. The number of nodes being monitored grows from 556 at observation start time to 794 nodes at the end of the observation period. More precisely, 423 nodes joined the system after the obser-

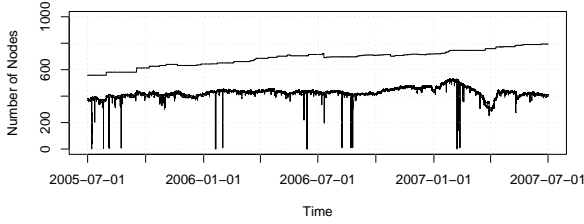


Figure 1: Number of monitored nodes (top) and available nodes (bottom) over time.

vation started, 185 left before it ended, and 90 nodes joined the system after start and left before the end. On average, more than one third of all nodes monitored in a query period were not available in the query period. For some query periods, no node was available at all. These “anomalies” are probably caused by failures of the communication infrastructure near the data collector. Nevertheless, we attribute these failures to the monitored nodes to obtain the dependability characteristics as perceived by the data collector.

For further evaluation, we exclude all nodes that were not queried for the whole observation period. This restriction allows us to base each estimation on the data for a complete two-years period and, therefore, to improve the accuracy of the estimations. It also eases the evaluation of relationships among failure characteristics. For example, computing failure correlation coefficients between nodes is only reasonable for periods in which both nodes are queried. For the 461 permanently queried nodes, 123,306 failures and 123,193 repairs have been observed. Hence, on average, each node failed 267.5 time and was repaired 267.2 times. For 15 nodes, we cannot derive any sensible failure or repair data as these were not queried successfully at all. We additionally exclude these 15 nodes from further evaluation and only consider the 446 nodes that were queried for the whole observation period and were queried successfully at least once.

The estimation of the *MTTF* of a node is based on all observed periods between a failure and a subsequent repair of the node. We additionally assume that a node is repaired (fails) at observation start if the node is (not) available at observation start. We point estimate the actual *MTTF* of a node by the sample mean $\hat{\Theta}$ of the observed periods. As the point estimates rarely coincide with the actual values, we assess the quality of the estimation by the confidence interval meaning: we compute ϵ such that the actual *MTTF* value is in the interval $(\hat{\Theta} - \epsilon, \hat{\Theta} + \epsilon)$ with a probability of 0.95.

Figure 2a gives the cumulative distribution function of the estimated *MTTF* values and a summary of the results. The lowest *MTTF* found for a node is ca. 5 minutes, which is approximately the smallest

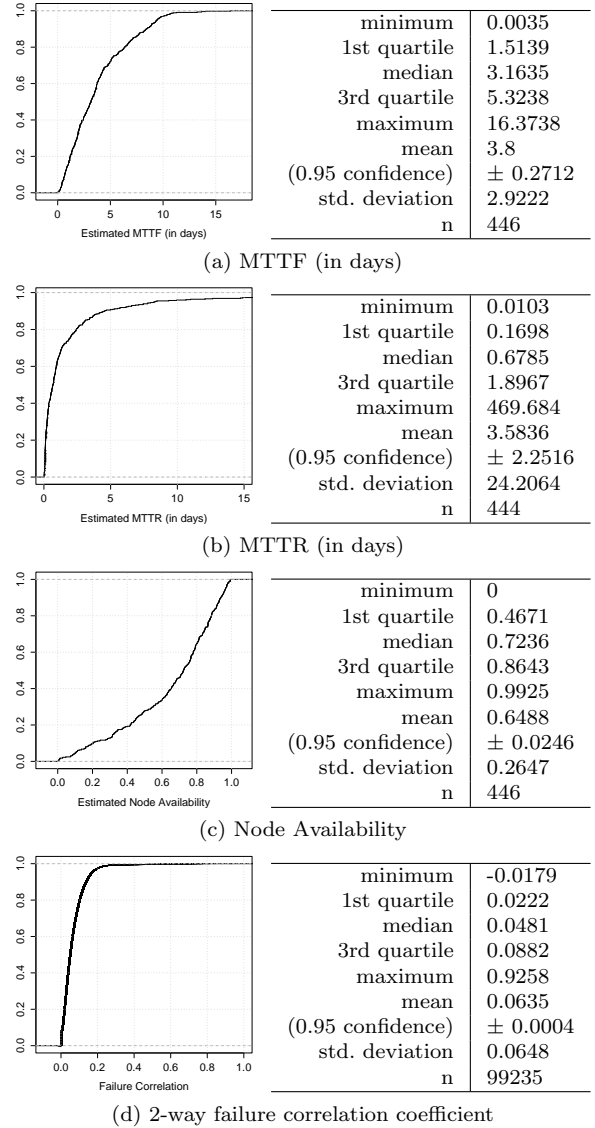


Figure 2: Estimated parameters for nodes

value that can be observed due to the query frequency. The highest *MTTF* found is 16.37 days. On average, a node has an *MTTF* of 3.8 days, but half of the nodes have an *MTTF* below 3.16 days. Overall, the values are rather widespread with a standard deviation of 2.9 days. Computing the confidence intervals reveals that, although we have a long observation period of two years and a high query frequency of $1/5 \text{ min}^{-1}$, $\hat{\Theta}$ is a rather rough estimation of the actual value: on average, ϵ is 1.42 days. In particular, ϵ grows with an increasing *MTTF*. For example, ϵ is 11.2 days for the node with the maximum *MTTF* of 16.4 days.

The *MTTR* values are estimated analogously to the *MTTF* values. As summarised in Fig. 2b, the lowest *MTTR* found is 15 minutes and the average one 3.58 days, but half of the nodes have an *MTTR* below 16.3

hours. The highest MTTR found is 469.7 days, i.e., for at least one node, it took the data collector over a year between two successful queries. Clearly, such extreme values can be considered as outliers. When omitting the top 5 % of the values, the mean shrinks to 1.21 days, the median to 14.4 hours, and the standard deviation to 1.61 days. Note that the MTTR estimations are only based on the data of 444 instead of 446 nodes as two nodes failed once and were not repaired at all during the observation period. The estimation of MTTF values is even less accurate than the estimation of the MTTR values: on average, ϵ is 4.1 days.

The (limiting) *availability* A of a node can be computed from its MTTF and MTTR by $A = MTTF / (MTTF + MTTR)$. However, as some nodes were not repaired during the observation period and due to the accuracy issues of the MTTF and MTTR estimations, we take a different approach and determine the estimated availability \hat{A} of a node by the expectation value of the random variable that denotes whether a node is available over time. The results are summarised in Fig. 2c. The node with the minimum availability is hardly available at all ($\hat{A} < 0.00005$); the most available one has $\hat{A} = 0.9925$; that is, the most available node operated correctly 362 days per year. Just like the MTTF and the MTTR, the availability values are widespread: with a standard deviation of 0.26, a node has $\hat{A} = 0.65$ on average, but half of the nodes have an availability above 0.72. Compared to the estimation of the MTTF and the MTTR, the accuracy of the estimated availability is better: on average, ϵ is below 0.0017 and even the maximum ϵ is below 0.0022.

With an increasing MTTF, a node generally becomes more available as the node is able to continuously provide service for longer periods without a failure. Availability also improves with a decreasing MTTR as the node then sooner continues to provide service after a failure. Hence, a node with high MTTF and low MTTR has a high availability. But does the reverse also hold, has a node with a high availability a high MTTF and a low MTTR? Evaluating the relationship among availability, MTTF, and MTTR reveals that availability correlates with the MTTF and the MTTR with coefficients of 0.27 and -0.25 , respectively. Interestingly, the MTTF and the MTTR only correlate with a coefficient of -0.04 , i.e., good MTTF and good MTTR values are rather independent.

A common assumption of many dependability models is that failures of different nodes are independent. However, there has been empirical evidence that this assumption does not hold in many real-world systems [1–3, 9, 12]. We confirm these results by computing the (Pearson) *correlation coefficient*, which also is an input

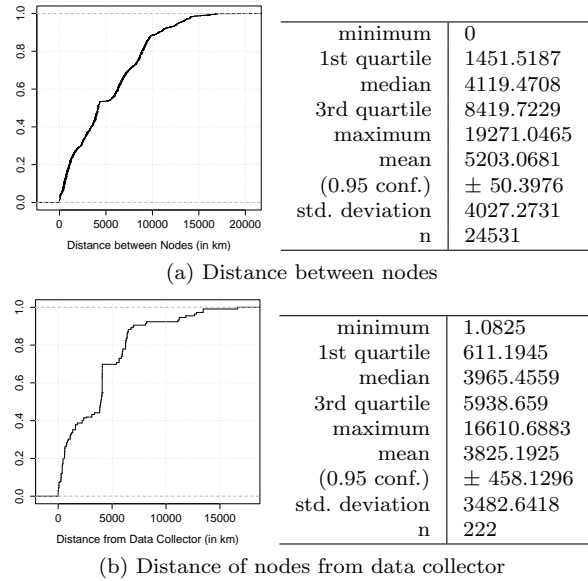


Figure 3: Geographical distance

parameter for many dependent failure models. More precisely, we compute the coefficient for random variables that are defined for each node as being 1 for a query period if the node fails in the query period and 0 otherwise. If failures of different nodes are independent, then the correlation coefficient is approximately 0. If two nodes have a high failure correlation coefficient, the assumption that their failures are independent can be safely rejected. Figure 2d shows the correlation coefficients for each pair of different nodes. The average correlation coefficient is 0.06 with an approximately equal standard deviation of 0.06. One quarter of the node pairs have a coefficient above 0.09. The maximum coefficient was above 0.92.

When investigating high failure correlations more closely, we found that nodes with high correlation coefficients are generally physically co-located (judging from their domain names). This result motivated us to investigate the relationship between dependability metrics and geographical distance. With the geo-lookup service <http://hostip.info/>, we were able to determine the geographical location of 222 of the evaluated nodes. Figure 3a shows the results of computing the distance between pairs of nodes. On average, two nodes have a distance of 5,203 km; half of the nodes have a distance less than 4,119 km. The maximal distance is 19,271 km, which approximately equals the largest distance possible on Earth. We found that the geographical distance between a pair of nodes and their failure correlation coefficient correlates with a coefficient of -0.23 (i.e., this correlation has approximately the same absolute coefficient as the correlation between availability and $MTTF / MTTR$).

While PlanetLab is a globally distributed system, CoMon operates in a centralised manner with one data collector statically deployed at a single node. The more distant a node is from the collector, the more likely querying the node fails as, generally, more communication infrastructure is involved to query more distant nodes. We quantify this effect by relating the dependability metrics of a node to its geographical distance from the data collector. Figure 3b summarises the distances between the monitored nodes and the data collector. The maximal distance of a node from the collector is 16,610 km, the average one 3,825 km. The MTTF values for the nodes and their distance to the data collector considerably correlate with a coefficient of -0.2 . The MTTR is only weakly correlated with a coefficient of -0.01 . The availability exhibits a higher correlation with a coefficient of -0.27 .

5. Fitting to Theoretical Distributions

Empirical distributions, as presented so far, are rather inconvenient for dependability models. Simple theoretical distributions (e.g., the exponential distribution) are more easy to handle and often allow simplifications resulting in more tractable models. We try to find well-fitting simple theoretical distributions and perform goodness-of-fit tests by computing the Kolmogorov-Smirnov (KS) statistic against the empirical distributions. We assume that a theoretical distribution with a lower KS statistic fits better than another distribution with a higher KS statistic.

For the empirical distributions of the MTTF, MTTR, node availability, and failure correlation coefficient values, we estimate the parameters for different simple distribution families (e.g., Exponential, Weibull, Log-Normal, Gaussian, and Gamma) by *maximum likelihood estimations* summarised in Tab. 1. The MTTF and availability values are best approximated by Weibull distributions, while MTTR and failure correlation values are best fitted by Log-Normal Distributions. While the KS statistic is low for the MTTF, MTTR, and availability distributions, it is rather high for the failure correlation coefficient. Unfortunately, the failure correlation coefficient is not well approximated by a simple theoretical distribution such that dependent failure models must either rely on more complex distributions or accept potential accuracy issues.

The TTF and the TTR of an individual node are often assumed to be exponentially distributed. We try to verify this assumption by a goodness-of-fit test for each node. More precisely, we conduct a KS test for the TTF and TTR of each node and the exponential distribution with unknown mean, a level of significance $\alpha = 0.01$, and critical values taken from Lilliefors [4].

When conducting the tests, the hypothesis of exponentially distributed TTF is rejected for 419 out of 446 nodes. For the TTR, the hypothesis is rejected for 442 out of 444 nodes. When fitting the empirical distributions of TTF and TTR to different distribution families, we find that their distributions are best approximated by distributions of the log-normal distribution family with an average KS statistic of 0.1489 for the TTF and 0.27 for the TTR. Although the fit by log-normal distributions is better than the fit by exponential distributions, the assumption that the TTF / TTR are log-normal distributed must unfortunately be rejected for most nodes. Overall, we find that TTF / TTR are not well-approximated by distributions from a single simple distribution family, but either require distributions from a single complex distribution family or from different families. For example, as the assumption of exponentially distributed TTF was accepted for 27 nodes, there is no pressing need to find another distribution family for these nodes.

6. Conclusion

A model that relies on inaccurate assumptions or parameters is unlikely to yield realistic results. Empirical evaluations are essential for dependability models of distributed systems as they allow to accurately assess a model's parameters and assumptions. In this paper, we empirically evaluated popular model assumptions and parameters for more than 400 globally distributed nodes based on long-term monitoring data as perceived by a central observer node. As one result, we found that the nodes are very heterogeneous w.r.t. their failure characteristics (e.g., the standard deviations are relatively high for all characteristics). This indicates that dependability models should, for example, not assume the same MTTF for all nodes, but model the MTTF value by a distribution function. We tried to fit the identified empirical distributions by simple theoretical ones that are favourable for tractable models.

The evaluated nodes have surprisingly bad failure characteristics with, for example, an availability of 0.65 on average. In contrast, Long et al. [6] conducted a related study in 1995 and found significantly better failure characteristics with, for example, an average availability of 0.88. The differences are partially explained by the different approach of monitoring. While Long et al. relied on decentralised monitoring, which masks failures of the communication infrastructure, we relied on centralised monitoring and attributed all communication failures to the affected nodes, which seemed to have a significant impact. This indicates that such failures should be considered in future dependability models and empirical evaluations.

Table 1: Summary of distribution fitting

Metric	Distr. Family	Parameter	KS-Statistic
MTTF	Weibull	shape = 1.2864, scale = 4.1027	0.0323
MTTR	Log-normal	mean = -0.4822, std. dev. = 1.6366 (log scale)	0.0406
Availability	Weibull	shape = 1.2976, scale = 0.3797	0.0425
Fail. Cor. Coeff.	Log-Normal	mean = 0.0599, std. dev. = 0.0559 (log scale)	0.1321

Another likely reason for the bad failure characteristics is that PlanetLab is a research cluster with low requirements on high availability or reliability. For example, Peterson et al. [8] document that PlanetLab weakly isolates resources for the sake of efficiency such that projects with low-quality implementations may cause node failures. As the individual nodes are of low criticality, their repair often has no priority. This raises the question of how representative our evaluation is for other kinds of systems. On the one hand, PlanetLab is likely to have worse failure characteristics than, for example, clusters that are critical for commercial revenue. On the other hand, it probably has better characteristics than, for example, many peer-to-peer networks with high churn-rates. Unfortunately, these differences can hardly be quantified without further empirical evaluations of other systems.

While our assessment of independent failures and exponentially distributed TTF and TTR confirm previous results, the finding of correlation between failure characteristics and geographical distance in large-scale distributed systems has not received attention so far. Our evaluation reveals that the distance of a node to the observer node correlates with the MTTF and the availability of the node, but hardly with its MTTR. Likely reasons are that more communication infrastructure is involved to communicate with a more distant nodes. This means that there are more possible sources of failure, which overall decrease the time-to-failure and availability. As each component is restored independently, this hardly affects the repair process.

Likewise, the failure correlation coefficient of node pairs correlates with the distance between the nodes. A possible reason is that co-located nodes are connected to the Internet via the same communication infrastructure and administrated by the same personnel. If the communication infrastructure at the physical location of co-located nodes fails or an administrator mis-configures co-located nodes at the same time, the nodes jointly fail. While correlations between distance and failure characteristics have been considered for embedded systems [5], they seem to also occur in large-scale distributed systems. Future research on dependability models of such systems may address the question of how to consider geographical distance for more accurate models.

Acknowledgements We would like to thank the people behind CoMon, KyoungSoo Park and Vivek Pai, for making CoMon’s monitoring data available to the public and, thereby, enabling this study.

References

- [1] Y. Amir and A. Wool. Evaluating quorum systems over the internet. In *Proc. 26th Int. Fault-Tolerant Comp. Symp.*, pages 26–35. IEEE, 1996.
- [2] M. Bakkaloglu, J. J. Wylie, C. Wang, and G. R. Ganger. Modeling correlated failures in survivable storage systems. In *Proc. 2002 Int. Conf. on Dependable Systems and Networks*. IEEE, 2002.
- [3] W. J. Bolosky, J. R. Douceur, D. Ely, and M. Theimer. Feasibility of a serverless distributed file system deployed on an existing set of desktop PCs. In *Proc. ACM SIGMETRICS '00*, pages 34–43. ACM, 2000.
- [4] H. W. Lilliefors. On the Kolmogorov-Smirnov test for the exponential distribution with mean unknown. *Journal of the American Statistical Association*, 64(325):387–389, Mar. 1969.
- [5] P. Limbourg, H.-D. Kochs, K. Echtele, and I. Eusgeld. Reliability prediction in systems with correlated component failures. In *Proc. 20th Int. Conf. on Architecture of Comp. Systems*, pages 55 – 62. VDE, 2007.
- [6] D. Long, A. Muir, and R. Golding. A longitudinal survey of internet host reliability. In *Proc. 14th Symp. on Reliable Distr. Systems*, pages 2–9. IEEE, 1995.
- [7] K. Park and V. S. Pai. CoMon: A mostly-scalable monitoring system for PlanetLab. *ACM SIGOPS Operating Systems Review*, 40(1):65–74, 2006.
- [8] L. Peterson, A. Bavier, M. Fluczynski, and S. Muir. Experiences building planetlab. In *Proc. 7th USENIX Symp. on Operating Systems Design and Implementation*, pages 351 – 366. USENIX Association, 2006.
- [9] D. Tang and R. K. Iyer. Analysis and modeling of correlated failures in multicomputer systems. *IEEE Trans. on Computers*, 41(5):567–577, 1992.
- [10] T. Warns, F. Freiling, and W. Hasselbring. Solving consensus using structural failure models. In *Proc. 25th Symp. on Reliable Distr. Systems*, pages 212–221. IEEE, 2006.
- [11] H. Weatherspoon, T. Moscovitz, and J. Kubiawicz. Introspective failure analysis: Avoiding correlated failures in peer-to-peer systems. In *Proc. 21st Symp. on Reliable Distr. Systems*, pages 362–369. IEEE, 2002.
- [12] P. Yalagandula, S. Nath, H. Yu, P. B. Gibbons, and S. Sesh. Beyond availability: Towards a deeper understanding of machine failure characteristics in large distributed systems. In *Proc. of the 1st Workshop on Real, Large Distr. Systems*. USENIX, 2004.