



Breaking the Availability Barrier I  
Survivable Systems for Enterprise Computing

Volume 1 of 3 Volume Series  
Dr. Bruce Holenstein, Dr. Bill Highleyman, and Paul J.  
Holenstein

© 2007 Gravic, Inc. All rights reserved.

No part of this book may be reproduced, stored in a retrieval system, or transmitted by any means, electronic, mechanical, photocopying, recording, or otherwise, without written permission from the authors.

ISBN: 978-1-4107-9231-0 (e-book)

ISBN: 978-1-4107-9232-7 (Paperback)

ISBN: 978-1-4107-9233-4 (Dust Jacket)

Library of Congress Control Number: 2003099708

All products mentioned in this book are trademarks of their respective owners. The information in this book is provided on an as-is informational basis. The authors, owners, and publisher disclaim liability for any errors or omissions. The reader accepts all risks associated with the use of the contents of this book.

### **About the Authors:**

The authors of the book, Dr. Bill Highleyman, Paul J. Holenstein, and Dr. Bruce Holenstein, have a combined experience of over 90 years in the implementation of fault-tolerant, highly available computing systems. This experience ranges from the early days of custom redundant systems to today's fault-tolerant offerings from HP (NonStop) and Stratus.

Click for Book Order Information: [Breaking the Availability Barrier](#)  
or visit Amazon.com or Authorhouse.com to purchase.

## Chapter 5 - The Facts of Life

In the first four chapters of this book, we focused on some basic concepts and applied these concepts to the development of architectures that offered significantly enhanced availability. The model that we used was that of a system that comprised multiple identical subsystems. Of these subsystems,  $s$  were spares; and the system could tolerate the failure of any  $s$  subsystems. However, the failure of  $s+1$  subsystems might cause the failure of the system. In the event of a system failure caused by  $s+1$  subsystem failures, the system was immediately restored to service upon the repair of one of the failed subsystems.

This model is accurate for replicated and split systems. For instance, if a split system comprises two nodes, and should both nodes fail, then service is restored as soon as one of the failed nodes is repaired and brought back into service (assuming that no database recovery is required before the system can be used).

However, within a failed node, things are not so simple. Once the requisite number of subsystems have been repaired and are operational, the node often must be recovered before it can be returned to service. This may require a variety of actions taking several hours.

In this chapter, we take a look at the impact of system recovery which may have to follow system repair. But first, let us review what we have done to date.

## A Review of Availability

In our previous chapters, we developed the general availability relation

$$A \approx 1 - f(1 - a)^{s+1} \approx 1 - f \left( \frac{mtr}{mtbf} \right)^{s+1} \quad (5-1)$$

where

- $A$  is the availability of a system comprising similar redundant subsystems.
- $a$  is the availability of a subsystem.
- $s$  is the number of spares provided in the system. That is, any  $s$  subsystems may fail; and the system will continue to operate.
- $f$  is the number of failure modes, or the number of ways that  $s+1$  subsystems will fail in such a way that a system outage will result.
- $mtr$  is the mean time to repair a subsystem.
- $mtbf$  is the mean time before failure for a subsystem.

Availability is the probability that the system will be operational and is therefore

$$A = \frac{MTBF}{MTBF + MTR} \quad (5-2)$$

where

- $MTBF$  is the mean time before failure for the system.
- $MTR$  is the mean time to repair the system (or, more accurately, the mean time to *restore* the system to service)

We discussed the characterization of availability in terms of 9s; e.g., an availability of .9999 is called four 9s. We showed that adding

a backup doubles the nines, which is the basic power of fault-tolerant systems. We discussed the relative efficiencies of methods to keep replicated databases in exact synchronization. We also showed that splitting a system into several independent nodes using data replication can provide substantial improvement in availability at little or no additional cost.

However, systems of multiple redundant processors suffer from a plurality of failure modes which can reduce system availability by one or two 9s. We showed that by careful allocation of processes to processors, we can substantially reduce this impact.

We so far have applied these concepts to hardware failures as if hardware failures were the predominant cause of system outages. This has been the mechanism that has allowed us to demonstrate the application of these concepts. However, this is not real life. Recent experience indicates that only a very small portion of system outages are caused by dual hardware failures. The rest are caused by complex interactions between hardware failures, software failures, operational errors, and environmental faults.

Adding to the rules which we derived in our earlier chapters, we have

**Rule 18:** *Redundant hardware systems have an availability of five to six 9s. Software and people reduce this to four 9s or less.*

The work to date is no waste of time. The concepts are valid. But in order to apply these concepts to real-life availability, we have to understand better what is going on.

## **Why Do Computers Stop?**

Jim Gray is unquestionably one of the key contributors to fault-tolerant computing. In 1985, he published “*Why Do Computers Stop*

*Dr. Bill Highleyman, Paul J. Holenstein, and Dr. Bruce Holenstein*  
*and What Can We Do About It,*<sup>1</sup> a defining paper on the real causes of system outages. This paper formed the basis for his chapter on “*Software Fault Tolerance*” in his subsequent book Transaction Processing: Concepts and Techniques,<sup>2</sup> which sets forth the basic principles underlying fault-tolerant systems.

Gray’s footprints can be found all throughout this chapter. Rather than citing each reference, let us suggest that you read his paper, which can be found at <http://www.cs.berkeley.edu/~yelick/294-f00/papers/Gray85.txt>.

Though Gray’s paper was published in 1985, it remains strikingly applicable today. Little has changed except that hardware has become somewhat more reliable, as he predicted. We summarize his observations next with some updating comments.

Gray studied 166 unscheduled system outages reported to Tandem Computers over a seven-month period. The outages covered over 2000 systems. Note that this equates to a system availability of .99993, assuming an MTR of four hours (see later). It confirmed Tandem’s claim to availabilities of about four 9s.

About one-third of the reported outages related to “infant mortality” problems, defined as recurring problems which were later fixed. If these were removed, only 107 system outages were reported during this period. The outages were characterized as follows (the Standish results are described later):

---

<sup>1</sup> Gray, J.; “*Why Do Computers Stop and What Can We Do About It?*” 5<sup>th</sup> Symposium on Reliability in Distributed Software and Database Systems; 1986.

<sup>2</sup> Gray, J.; et al.; Transaction Processing: Concepts and Techniques, Morgan Kaufmann; 1993.

	<b>Gray 1985</b>	<b>Standish 2002</b>
People	42%	38%
Software	25%	28%
Hardware	18%	17%
Environment (power, a/c, etc.)	14%	17%
Unknown	1%	-

**Contributors to Tandem NonStop Outages**  
**Table 5-1**

There are some caveats that must be made about these observations. First, they probably don't include many outages caused by application software faults or environmental problems. Customers don't generally report these. Second, they probably don't include all operator errors. Operators don't always report their goofs. Finally, they don't include scheduled down time.

Gray points out in his 1985 paper that "...hardware will be even more reliable due to better design, increased levels of integration, and reduced numbers of connectors." In fact, this forecast has turned out to be quite correct. Experience now indicates that outages due to dual hardware failures represent less than 5% of all outages.

Gray goes on to say, "...the trend for software and system administration is not so positive. Systems are getting more complex." A recent study by Standish Group<sup>26</sup> supports this observation. Clearly, not much has changed since 1985. Citing a NonStop server reliability of .9998 (classic Tandem systems are now known as NonStop servers following the acquisition of Tandem by HP), Standish's measurements of the causes of system outages are also shown in Table 5-1. They are amazingly consistent with Gray's findings. Interestingly, the Standish study also indicates that network failures happen more often than server system failures; and outages caused by applications swamp server system failures by almost four to one.

---

<sup>26</sup> Standish Group, VirtualBEACON, Issue 244; September, 2002

*Dr. Bill Highleyman, Paul J. Holenstein, and Dr. Bruce Holenstein*

Note that about 40% of all outages seem to have been caused by human error. However, this isn't as bad as it may look. Forty-five outages per year over 2000 systems represent one human error per system every 44 years. Don't we all wish that we were that accurate?

An interesting insight into operator errors has been nicely phrased by Wendy Bartlett of HP. It recognizes the stress factor which accompanies an unexpected failure:

**Rule 19: (Bartlett's Law)** *When things go wrong, people get stupider.*

Because of this human characteristic, it is very important to have a recurrent training program to ensure that operations staff can react properly in an emergency and that the established procedures are correct. Ideally, failure situations should be simulated periodically so that the operations staff can practice proper recovery. The more automatic their reactions become, the less stupid they will appear when the heat is on. This leads to the following corollary to Rule 19:

**Rule 20:** *Conduct periodic simulated failures to keep the operations staff trained and to ensure that recovery procedures are current.*

Today, the numbers may be a little different. In fact, they may be different each year because the sample size is too small. However, all indications are that the story these numbers tell is the same:

**Rule 21:** *System outages are predominantly caused by human and software errors.*

**Rule 22:** *Seldom does a recovery entail hardware repair. It entails a reload of the system.*

Let us consider further these Rules 21 and 22.

## **Some Definitions**

Let us first define some terms:

A *fault* is a lurking incorrectness waiting to strike. It may be a hardware or software design error, a hardware component failure, a software coding error, or even a bit of human ignorance (such as an operator's confusion over the effects of a given command).

A *failure* is the exercise of a fault. Failures in themselves do not cause outages. Fault-tolerant systems are designed to survive any single failure.

A *failover* is a recovery from a failure by switching to a backup component.

A *failover fault* is a failure of a failover.

An *outage* is a denial of some service to part or all of the user community. Outages can range from unacceptable response times to total system unavailability.

A *trigger* is an initial failure that begins an event sequence that leads to a subsequent failure from which the system cannot recover. The result is a system outage.

A *repair* is the return to service of a component which has experienced a failure.

A *recovery* is the return to service of a system which has experienced an outage.

Note that the *repair* of a failed subsystem does not necessarily result in the *recovery* of its system.

## **Triggered Outages**

Our availability analysis so far has assumed (in today's fault-tolerant environment) that an outage is caused by two failures *and* that these two failures are independent (i.e., one failure does not induce the other). Although simplistic, this view of things has led us to several important concepts in availability.

We now see that dual failures which cause outages are hardly independent at all. Rather, one random failure often leads to a directly related second failure, which causes the outage. For example,

- a disk unit fails, and the good disk is erroneously pulled for replacement.
- a critical process aborts, and a bug in its backup checkpointing procedures causes the backup to fail as well.
- An error (a fault) in the operations manual is followed by the system operator and results in an outage.

In general, a failure occurs and must be corrected but otherwise does not seriously affect system operation. But then that failure triggers another failure and a system outage results.

This leads to another interesting observation. You may have noted that new systems seem to be less reliable than established systems. A system seems to "burn in" with time. How can this be?

New systems are subject to continuous change. Functional errors are corrected, bugs are worked out, and enhancements are made. Each change carries with it the potential for further errors which may act as outage triggers. As the system matures, changes become less frequent; and system reliability improves. As Carl Neihaus of HP has said,

**Rule 23: (Neihaus' Law) - *Change causes outages.***

## The Impact of Failover Faults

In Chapter 1 of this series, we showed that the probability of an outage for a system configured with one spare and with randomly distributed processes is

$$F \approx \frac{n(n-1)}{2}(1-a)^2 \quad (5-3)$$

where

- $F$  is the probability of a system outage.
- $n$  is the number of subsystems in the system.
- $a$  is the probability of a subsystem failure.

That is, a subsystem will fail with a probability of  $(1-a)$ . A pair of subsystems will fail with a probability of  $(1-a)^2$ , and there are  $n(n-1)/2$  ways in which two subsystems might fail.

This relation assumes that the outage is caused by the independent failures of two subsystems, whether those failures are caused by hardware or by software. Implicit in Chapters 1 and 2 was the assumption that critical process pairs were trusted and did not fail. A system outage was caused by a dual hardware failure, which took down a critical process pair or which denied a critical process pair access to needed data. We now relax that assumption and allow a subsystem to fail due either to a hardware failure or to a software failure. Furthermore, we realize that the *failover* mechanism which should recover from that failure by invoking the hardware or software backup may fail itself and thereby create a *failover fault*.

Thus, we now know that not all outages are caused by dual independent failures. With a probability of  $p$ , a *single* subsystem failure, whether it be due to hardware or software, will experience a failover fault that leads to a system outage. Only  $(1-p)$  of all outages are caused by dual subsystem failures, where

$p$  is the probability that a failover attempt will fail (a failover fault).

*Dr. Bill Highleyman, Paul J. Holenstein, and Dr. Bruce Holenstein*

We now have two failure modes for a system:

- two subsystems have failed.
- one subsystem has failed, and the failover has failed.

Furthermore, there are two components to system restoration – subsystem repair and system recovery. Let

- $r$  be the mean time to repair a subsystem.
- $R$  be the mean time to recover a system.

Also, remember that the average time to return the system to service when there are two failed subsystems is  $r/2$  (see Rule 7 in Chapter 1).

Let us consider each of the two failure modes listed above.

(1) Two subsystems have failed.

The probability that two subsystems will fail is given by Equation (5-3).  $(1-p)$  of all outages will be caused by a dual subsystem failure. Furthermore, the effective system repair time is the time to repair one of the subsystems,  $r/2$ , plus the time to recover the system,  $R$ , giving a total repair time of  $r/2+R$  rather than just  $r/2$ . Since down time is proportional to repair time, the failure probability for this mode is increased by a factor of  $(r/2+R)/r/2$ :

outage probability due to dual failures

$$\approx \frac{r/2+R}{r/2} (1-p) \frac{n(n-1)}{2} (1-a)^2 \quad (5-4a)$$

(2) One subsystem has failed, and the failover has failed

The probability that a particular subsystem will fail is  $(1-a)$ . There are  $n$  ways in which a system can experience a single subsystem failure. Thus, the probability that one subsystem in

the system will fail is  $n(1-a)$ . Of all outages,  $p$  are caused by single subsystem failures followed by a failover fault.

Should an outage be caused by a failover fault, there is no need for a subsystem repair since no more than one subsystem is down. Thus, subsystem repair time,  $r$ , has been replaced by system recovery time,  $R$ . Therefore, the probability of failure for this mode is modified by a factor of  $R/r$ :

$$\begin{aligned} &\text{outage probability due to failover fault} \\ &\approx \frac{R}{r}pn(1-a) \end{aligned} \quad (5-4b)$$

Recognizing that the system failure probability  $F$  is the sum of the above two probabilities, and assuming that  $p$  is very much less than 1, then

$$F \approx \frac{\frac{r}{2} + R}{\frac{r}{2}} \frac{n(n-1)}{2} (1-a)^2 + \frac{R}{r}pn(1-a) \quad (5-4c)$$

These relationships are good approximations describing the impact of failover faults on system availability. They are proven formally in Appendix 3, “*Failover Faults.*” The following observations are made about these relationships:

### ***A Better Value for Subsystem Availability***

So far, we have assumed that system availability  $A$  is about four 9s, given a subsystem availability  $a$  of .995. Now we can deduce a better value for  $a$  from Equation (5-4c).

Let us consider an 8-processor system ( $n = 8$ ) with an availability  $A$  of four 9s ( $F = .0001$ ). We assume a repair time  $r$  of 24 hours, a recovery time  $R$  of 4 hours, and a failover fault probability of 1%. Solving Equation (5-4c) for  $a$ , we find that a better value for subsystem availability  $a$  is .9986, more than three times better than we previously assumed. This is because we now recognize that some

*Dr. Bill Highleyman, Paul J. Holenstein, and Dr. Bruce Holenstein*

outages are caused by failover faults and not by dual subsystem failures.

### ***Effect of Failover Faults on System Availability***

We can calculate from Equations (5-4) that a 1% failover fault rate causes 20% of all system down time under the previous example. Note from Equation (5-4b) that this is directly affected by system recovery time  $R$ . To the extent that we can reduce recovery time, we can minimize the effects of failover faults.

### ***Effect of Failover Faults on Effective Subsystem Availability***

A further insight into the impact of failover faults on system availability is gained as follows.

We can rewrite Equation (5-4c) as

$$F \approx \frac{\frac{r}{2} + R}{\frac{r}{2}} \frac{n(n-1)}{2} (1-a)(1-a') \quad (5-5)$$

where

$$a' = a - \frac{R}{\frac{r}{2} + R} \frac{p}{n-1} \quad (5-6)$$

Note that  $a'$  is less than  $a$ , being reduced by a subtractive term. Thus, comparing Equation (5-5) to Equation (5-4a) (and ignoring the  $1-p$  term as being very close to one), we can make an interesting interpretation. The failure of the first subsystem will occur with a probability of  $(1-a)$  as expected. *However, once one subsystem has failed, the system then behaves as if it comprises  $n-1$  remaining subsystems with decreased availability  $a'$ .*

A simple example will serve to illustrate this. Let us consider a 4-processor node ( $n = 4$ ) comprising subsystems with an availability,  $a$ , of .9986 (as calculated above). Subsystem repair time,  $r$ , is 24 hours and system recovery time,  $R$ , is four hours. If the probability of a

failover fault,  $p$ , is 1%, then the effective subsystem availability,  $a'$ , following a single subsystem failure, is reduced from .9986 to .9978. Failure probability has increased from .0014 to .0022. Under the above parameter assumptions, a 1% chance of a failover fault makes the system 60% less reliable following a single subsystem failure!

This leads to the following rule:

**Rule 24:** *Following the failure of one subsystem, failover faults cause the system to behave as if it comprises  $n-1$  remaining subsystems with decreased availability.*

Note also that as recovery time  $R$  decreases, Equation (5-6) shows that the reduced subsystem availability  $a'$  improves and approaches the subsystem availability  $a$ . Thus, reducing recovery time reduces directly the impact of failover faults on system availability.

### ***Effect of Failover Faults on System Splitting***

In Chapter 2, we showed that splitting a system into  $k$  nodes improved system reliability by at least a factor of  $k$ :

$$\text{Split System Reliability Improvement} = k \frac{n-1}{n-k} > k \quad (5-7)$$

where

- $k$  is the number of nodes into which the system is split.
- $n$  is the number of processors in the system.

Thus, Equation (5-7) predicts that the reliability improvement achieved by system splitting is always greater than  $k$ . However, such a gain is compromised when we have the possibility of failover faults. As shown in Chapter 11, “*Failover Faults*,” this relation then can be shown to become

$$\text{Split System Reliability Improvement} = k \frac{(n-1) + x}{(n-k) + kx} \quad (5-8)$$

*Dr. Bill Highleyman, Paul J. Holenstein, and Dr. Bruce Holenstein*

where we can think of  $x$  as an exasperation factor. If  $x$  is very small, then Equation (5-8) approaches Equation (5-7); and we have the split system advantage we are seeking. However, if  $x$  is very large, then Equation (5-8) approaches one; and the reliability advantage of system splitting disappears.

$x$  is given by

$$x = \frac{R}{\frac{r}{2} + R} \frac{2p}{1-a} \quad (5-9)$$

Note that as recovery time  $R$  becomes smaller,  $x$  becomes less significant; and we recover our availability advantage provided by system splitting. Thus,

**Rule 25:** *The possibility of failover faults erodes the availability advantages of system splitting (see Rule 9).*

As an example, from Equation (5-7) we expect that splitting a 16-processor system into four 4-processor nodes will give us a reliability advantage of a factor of 5. But suppose that we have a failover fault probability,  $p$ , of 1%. Furthermore, assume that subsystem availability,  $a$ , is .9986, that subsystem repair time,  $r$ , is 24 hours, and that system recovery time,  $R$ , is four hours. In this case, from Equations (5-8) and (5-9), the availability advantage of system splitting decreases from a factor of 5 to a factor of 2.8.

## **The Golden Rule – Reduce Recovery Time**

We have seen the importance of minimizing recovery time to improve system availability. A further insight can be gained by noting that about 20% of CPU halts are caused by hardware failures, and about 80% are caused by software faults or human errors. Thus, only about 4% of system outages (20% x 20%) are caused by dual hardware failures.

The remaining 96% of outages are caused by no more than one hardware failure combined with a software fault or a human error.

These outages do not require a repair to return them to service. They only require a recovery:

**Rule 22 (restated):** *A system outage usually does not require a repair of any kind. Rather, it entails a recovery of the system.*

We now have seen that recovery time is a predominant factor in system availability:

- The time to recover from most outages is recovery time rather than repair time. Therefore, any reduction in recovery time is *directly* reflected in system availability.
- Minimizing recovery time helps negate the effect of failover faults on system failure rate.
- Minimizing recovery time helps maximize the availability advantages of system splitting.

This leads to what might be considered the golden rule of availability:

**Rule 26:** (The Golden Rule) - *Design your systems for fast recovery to maximize availability, to reduce the effect of failover faults, and to take full advantage of system splitting.*

What can we do to reduce recovery time?<sup>27</sup> The recovery process is quite complex. It may entail

- realizing that a problem has occurred.
- diagnosing the cause of the problem.
- deciding what to do.
- obtaining permission to follow a course of recovery action.
- collecting diagnostic data (such as a processor memory dump).

---

<sup>27</sup> For a description of on-going research into ways to minimize recovery time via Recovery-Oriented Computing (ROC), see A. Fox, D. Patterson, “*Self-Repairing Computers*,” Scientific American; June, 2003.

*Dr. Bill Highleyman, Paul J. Holenstein, and Dr. Bruce Holenstein*

- cold-loading the system.
- restarting software subsystems.
- restarting the applications.
- restarting the network.
- perhaps recovering the database.

A typical recovery procedure takes anywhere from an hour to several hours. Four hours is a pretty good guess at an average recovery time.

As can be seen from the above list, recovery time isn't minimized by simply building our applications for quick recovery, although this, of course, is very important. Efficient recovery also requires

- good operator training.
- efficient decision making by the entire management team.
- well-documented recovery procedures for
  - o the system.
  - o the application.
  - o the database.
  - o the network.

**Rule 27:** *Rapid recovery of a system outage is not simply a matter of command line entries. It is an entire business process.*

## **The Importance of Restore Time**

Let us define *restore time* as the time required to return a system to service. Restore time may entail a hardware repair and probably requires a system recovery. The *MTR* term in Equation (5-2) should really stand for *mean time to restore*.

The importance of restore time (which can also be called down time) so far has only been considered as a factor in availability. To the extent that one can reduce restore time, one can increase system availability.

### *Breaking the Availability Barrier*

However, restore time has a far greater importance in its own right. It may be that the average cost per down time hour is a user's measure of the value of reliability. But this cost is also a function of the length of a down time interval. As down time grows longer, the outage cost may increase. Customer annoyance may give way to customer anger, then lost sales, then lost customers.

Down time is even more profound for safety-critical functions. For instance, 911 operators will probably feel that a five-second outage is an annoyance. But a five-minute outage can mean a death due to cardiac arrest or a building being burned to the ground.

Fault-tolerant systems are full of outages that are seen not as outages at all but as very brief periods of perhaps degraded response times. These are the normal outages caused by subsystem failures that are recovered in seconds. For instance, in many fault-tolerant systems, recovery is accomplished by failing over to a backup process or by resubmitting a transaction to a surviving server process. It is the severe outages requiring a system recovery that currently limit fault-tolerant systems to four 9s.

Since down time is predominantly recovery time for fault-tolerant systems, we perhaps can add a nine to the current fault-tolerant availability if we pay more attention to designing systems and business processes for fast recovery. An excellent example of this is a major corporation which runs a variety of applications on several HP NonStop servers. Every one of their applications is backed up by another system that either normally runs other applications or is a dedicated backup. Recovery time for many of their replicated applications is about two minutes. A typical NonStop system has an availability of four 9s and an average recovery time of four hours. By cutting the recovery time from four hours to two minutes – over a 100:1 improvement in restore time – this corporation has added two 9s to its availability and thus has created an availability in excess of six 9s.

*Dr. Bill Highleyman, Paul J. Holenstein, and Dr. Bruce Holenstein*

## **What's Next?**

We have developed in the last four chapters several concepts regarding availability and the related parameters of mean time before failure and mean time to repair. We have shown that splitting a system into several independent nodes can significantly improve the reliability of a system, and have discussed how these nodes can remain synchronized. We have discussed the impact of system recovery on these parameters.

In Chapter 6, “*RPO and RTO*,” we look at two important objectives for high availability systems and discuss how various replication architectures map into these objectives.