

Preliminaries

The past 50 years have seen computers move from being expensive computational engines used by government and big corporations to becoming an everyday commodity, deeply embedded in practically every aspect of our lives. Not only are computers visible everywhere, in desktops, laptops, and PDAs, it is also a commonplace that they are *invisible* everywhere, as vital components of cars, home appliances, medical equipment, aircraft, industrial plants, and power generation and distribution systems. Computer systems underpin most of the world's financial systems: given current transaction volumes, trading in the stock, bond, and currency markets would be unthinkable without them. Our increasing willingness, as a society, to place computers in life-critical and wealth-critical applications is largely driven by the increasing possibilities that computers offer. And yet, as we depend more and more on computers to carry out all of these vital actions, we are—implicitly or explicitly—gambling our lives and property on computers doing their jobs properly.

Computers (hardware plus software) are quite likely the most complex systems ever created by human beings. The complexity of computer hardware is still increasing as designers attempt to exploit the higher transistor density that new generations of technology make available to them. Computer software is far more complex still, and with that complexity comes an increased propensity to failure. It is probably fair to say that there is not a single large piece of software or hardware today that is free of bugs. Even the space shuttle, with software that was developed and tested using some of the best and most advanced techniques known to engineering, is now known to have flown with bugs that had the potential to cause catastrophe.

Computer scientists and engineers have responded to the challenge of designing complex systems with a variety of tools and techniques to reduce the number of faults in the systems they build. However, that is not enough: we need to build systems that will acknowledge the existence of faults as a fact of life, and incorpo-

rate techniques to tolerate these faults while still delivering an acceptable level of service. The resulting field of *fault tolerance* is the subject of this book.

1.1 Fault Classification

In everyday language, the terms *fault*, *failure*, and *error* are used interchangeably. In fault-tolerant computing parlance, however, they have distinctive meanings. A *fault* (or *failure*) can be either a hardware defect or a software/programming mistake (bug). In contrast, an *error* is a manifestation of the fault/failure/bug.

As an example, consider an adder circuit, with an output line stuck at 1; it always carries the value 1 independently of the values of the input operands. This is a fault, but not (yet) an error. This fault causes an error when the adder is used and the result on that line is supposed to have been a 0, rather than a 1. A similar distinction exists between programming mistakes and execution errors. Consider, for example, a subroutine that is supposed to compute $\sin(x)$ but owing to a programming mistake calculates the absolute value of $\sin(x)$ instead. This mistake will result in an execution error only if that particular subroutine is used and the correct result is negative.

Both faults and errors can spread through the system. For example, if a chip shorts out power to ground, it may cause nearby chips to fail as well. Errors can spread because the output of one unit is used as input by other units. To return to our previous examples, the erroneous results of either the faulty adder or the $\sin(x)$ subroutine can be fed into further calculations, thus propagating the error.

To limit such contagion, designers incorporate *containment zones* into systems. These are barriers that reduce the chance that a fault or error in one zone will propagate to another. For example, a fault-containment zone can be created by ensuring that the maximum possible voltage swings in one zone are insulated from the other zones, and by providing an independent power supply to each zone. In other words, the designer tries to electrically isolate one zone from another. An error-containment zone can be created, as we will see in some detail later on, by using redundant units/programs and voting on their output.

Hardware faults can be classified according to several aspects. Regarding their duration, hardware faults can be classified into *permanent*, *transient*, or *intermittent*. A *permanent fault* is just that: it reflects the permanent going out of commission of a component. As an example of a permanent fault think of a burned-out lightbulb. A *transient fault* is one that causes a component to malfunction for some time; it goes away after that time and the functionality of the component is fully restored. As an example, think of a random noise interference during a telephone conversation. Another example is a memory cell with contents that are changed spuriously due to some electromagnetic interference. The cell itself is undamaged: it is just that its contents are wrong for the time being, and overwriting the memory cell will make the fault go away. An *intermittent fault* never quite goes away entirely; it oscillates between being quiescent and active. When the fault is quiescent, the

component functions normally; when the fault is active, the component malfunctions. An example for an intermittent fault is a loose electrical connection.

Another classification of hardware faults is into *benign* and *malicious* faults. A fault that just causes a unit to go dead is called *benign*. Such faults are the easiest to deal with. Far more insidious are the faults that cause a unit to produce reasonable-looking, but incorrect, output, or that make a component “act maliciously” and send differently valued outputs to different receivers. Think of an altitude sensor in an airplane that reports a 1000-foot altitude to one unit and a 8000-foot altitude to another unit. These are called *malicious* (or Byzantine) faults.

1.2 Types of Redundancy

All of fault tolerance is an exercise in exploiting and managing *redundancy*. Redundancy is the property of having more of a resource than is minimally necessary to do the job at hand. As failures happen, redundancy is exploited to mask or otherwise work around these failures, thus maintaining the desired level of functionality.

There are four forms of redundancy that we will study: hardware, software, information, and time. Hardware faults are usually dealt with by using hardware, information, or time redundancy, whereas software faults are protected against by software redundancy.

Hardware redundancy is provided by incorporating extra hardware into the design to either detect or override the effects of a failed component. For example, instead of having a single processor, we can use two or three processors, each performing the same function. By having two processors, we can detect the failure of a single processor; by having three, we can use the majority output to override the wrong output of a single faulty processor. This is an example of *static hardware redundancy*, the main objective of which is the immediate masking of a failure. A different form of hardware redundancy is *dynamic redundancy*, where spare components are activated upon the failure of a currently active component. A combination of static and dynamic redundancy techniques is also possible, leading to *hybrid hardware redundancy*.

Hardware redundancy can thus range from a simple duplication to complicated structures that switch in spare units when active ones become faulty. These forms of hardware redundancy incur high overheads, and their use is therefore normally reserved for critical systems where such overheads can be justified. In particular, substantial amounts of redundancy are required to protect against malicious faults.

The best-known form of information redundancy is error detection and correction coding. Here, extra bits (called check bits) are added to the original data bits so that an error in the data bits can be detected or even corrected. The resulting error-detecting and error-correcting codes are widely used today in memory units

and various storage devices to protect against benign failures. Note that these error codes (like any other form of information redundancy) require extra hardware to process the redundant data (the check bits).

Error-detecting and error-correcting codes are also used to protect data communicated over noisy channels, which are channels that are subject to many transient failures. These channels can be either the communication links among widely separated processors (e.g., the Internet) or among locally connected processors that form a local network. If the code used for data communication is capable of only detecting the faults that have occurred (but not correcting them), we can retransmit as necessary, thus employing time redundancy.

In addition to transient data communication failures due to noise, local and wide-area networks may experience permanent link failures. These failures may disconnect one or more existing communication paths, resulting in a longer communication delay between certain nodes in the network, a lower data bandwidth between certain node pairs, or even a complete disconnection of certain nodes from the rest of the network. Redundant communication links (i.e., hardware redundancy) can alleviate most of these problems.

Computing nodes can also exploit time redundancy through re-execution of the same program on the same hardware. As before, time redundancy is effective mainly against *transient* faults. Because the majority of hardware faults are transient, it is unlikely that the separate executions will experience the same fault. Time redundancy can thus be used to detect transient faults in situations in which such faults may otherwise go undetected. Time redundancy can also be used when other means for detecting errors are in place and the system is capable of recovering from the effects of the fault and repeating the computation. Compared with the other forms of redundancy, time redundancy has much lower hardware and software overhead but incurs a high performance penalty.

Software redundancy is used mainly against software failures. It is a reasonable guess that *every* large piece of software that has ever been produced has contained faults (bugs). Dealing with such faults can be expensive: one way is to independently produce two or more versions of that software (preferably by disjoint teams of programmers) in the hope that the different versions will not fail on the same input. The secondary version(s) can be based on simpler and less accurate algorithms (and, consequently, less likely to have faults) to be used only upon the failure of the primary software to produce acceptable results. Just as for hardware redundancy, the multiple versions of the program can be executed either concurrently (requiring redundant hardware as well) or sequentially (requiring extra time, i.e., time redundancy) upon a failure detection.

1.3 Basic Measures of Fault Tolerance

Because fault tolerance is about making machines more dependable, it is important to have proper measures (yardsticks) by which to gauge such dependability. In this section, we will examine some of these yardsticks and their application.

A measure is a mathematical abstraction that expresses some relevant facet of the performance of its object. By its very nature, a measure only captures some subset of the properties of an object. The trick in defining a suitable measure is to keep this subset large enough so that behaviors of interest to the user are captured, and yet not so large that the measure loses focus.

1.3.1 Traditional Measures

We first describe the traditional measures of dependability of a single computer. These metrics have been around for a long time and measure very basic attributes of the system. Two of these measures are *reliability* and *availability*.

The conventional definition of reliability, denoted by $R(t)$, is the probability (as a function of the time t) that the system has been up continuously in the time interval $[0, t]$. This measure is suitable for applications in which even a momentary disruption can prove costly. One example is computers that control physical processes such as an aircraft, for which failure would result in catastrophe.

Closely related to reliability are the *Mean Time to Failure*, denoted by MTTF, and the *Mean Time Between Failures*, MTBF. The first is the average time the system operates until a failure occurs, whereas the second is the average time between two consecutive failures. The difference between the two is due to the time needed to repair the system following the first failure. Denoting the *Mean Time to Repair* by MTTR, we obtain

$$\text{MTBF} = \text{MTTF} + \text{MTTR}$$

Availability, denoted by $A(t)$, is the average fraction of time over the interval $[0, t]$ that the system is up. This measure is appropriate for applications in which continuous performance is not vital but where it would be expensive to have the system down for a significant amount of time. An airline reservation system needs to be highly available, because downtime can put off customers and lose sales; however, an occasional (very) short-duration failure can be well tolerated.

The long-term availability, denoted by A , is defined as

$$A = \lim_{t \rightarrow \infty} A(t)$$

A can be interpreted as the probability that the system will be up at some random point in time, and is meaningful only in systems that include repair of faulty components. The long-term availability can be calculated from MTTF, MTBF, and MTTR as follows:

$$A = \frac{\text{MTTF}}{\text{MTBF}} = \frac{\text{MTTF}}{\text{MTTF} + \text{MTTR}}$$

A related measure, *Point Availability*, denoted by $A_p(t)$, is the probability that the system is up at the particular time instant t .

It is possible for a low-reliability system to have high availability: consider a system that fails every hour on the average but comes back up after only a second.

Such a system has an MTBF of just 1 hour and, consequently, a low reliability; however, its availability is high: $A = 3599/3600 = 0.99972$.

These definitions assume, of course, that we have a state in which the system can be said to be “up” and another in which it is not. For simple components, this is a good assumption. For example, a lightbulb is either good or burned out. A wire is either connected or has a break in it. However, for even simple systems, such an assumption can be very limiting. For example, consider a processor that has one of its several hundreds of millions of gates stuck at logic value 0. In other words, the output of this logic gate is always 0, regardless of the input. Suppose the rest of the processor is functional, and that this failed logic gate only affects the output of the processor about once in every 25,000 hours of use. For example, a particular gate in the divide unit when being faulty may result in a wrong quotient if the divisor is within a certain subset of values. Clearly, the processor is not fault-free, but would one define it as “down”?

The same remarks apply with even greater force to systems that degrade gracefully. By this, we mean systems with various levels of functionality. Initially, with all of its components operational, the system is at its highest level of functionality. As these components fail, the system degrades from one level of functionality to the next. Beyond a certain point, the system is unable to produce anything of use and fails completely. As with the previous example, the system has multiple “up” states. Is it said to fail when it degrades from full to partial functionality? Or when it fails to produce any useful output at all? Or when its functionality falls below a certain threshold? If the last, what is this threshold, and how is it determined?

We can therefore see that traditional reliability and availability are very limited in what they can express. There are obvious extensions to these measures. For example, we may consider the average computational capacity of a system with n processors. Let c_i denote the computational capacity of a system with i operational processors. This can be a simple linear function of the number of processors, $c_i = ic_1$, or a more complex function of i , depending on the ability of the application to utilize i processors. The *Average Computational Capacity* of the system can then be defined as $\sum_{i=1}^n c_i P_i(t)$, where $P_i(t)$ is the probability that exactly i processors are operational at time t . In contrast, the reliability of the system at time t will be

$$R(t) = \sum_{i=m}^n P_i(t)$$

where m is the minimum number of processors necessary for proper operation of the system.

1.3.2 Network Measures

In addition to the general system measures previously discussed, there are also more specialized measures, focusing on the network that connects the processors together. The simplest of these are classical node and line *connectivity*, which are

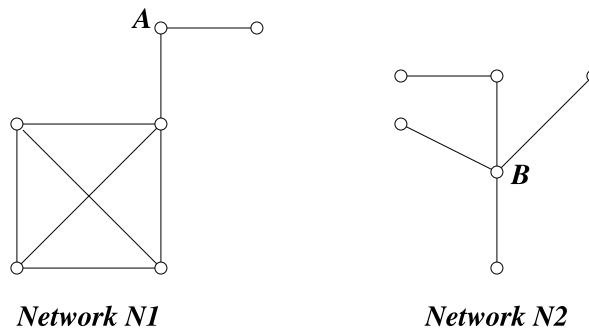


FIGURE 1.1 Inadequacy of classical connectivity.

defined as the minimum number of nodes and lines, respectively, that have to fail before the network becomes disconnected. This gives a rough indication of how vulnerable a network is to disconnection: for example, a network that can be disconnected by the failure of just one (critically positioned) node is potentially more vulnerable than another that requires at least four nodes to fail before it becomes disconnected.

Classical connectivity is a very basic measure of network reliability. Like reliability, it distinguishes between only two network states: connected and disconnected. It says nothing about how the network degrades as nodes fail before, or after, becoming disconnected. Consider the two networks shown in Figure 1.1. Both networks have the same classical node connectivity of 1. However, in a real sense, network *N1* is much more “connected” than *N2*. The probability that *N2* splinters into small pieces is greater than that for *N1*.

To express this type of “connectivity robustness,” we can use additional measures. Two such measures are the average node-pair distance, and the network diameter (the maximum node-pair distance), both calculated given the probability of node and/or link failure. Such network measures, together with the traditional measures listed above, allow us to gauge the dependability of various networked systems that consist of computing nodes connected through a network of communication links.

1.4 Outline of This Book

The next chapter is devoted to hardware fault tolerance. This is the most established topic within fault-tolerant computing, and many of the basic principles and techniques that have been developed for it have been extended to other forms of fault tolerance. Techniques to evaluate the reliability and availability of fault-tolerant systems are introduced here, including the use of Markov models.

Next, several variations of information redundancy are covered, starting with the most widely used error detecting and correcting codes. Then, other forms of

information redundancy are discussed, including storage redundancy (RAID systems), data replication in distributed systems, and, finally, the algorithm-based fault-tolerance technique that tolerates data errors in array computations using some error-detecting and error-correcting codes.

Many computing systems nowadays consist of multiple networked processors that are subject to interconnection link failures, in addition to the already-discussed single node/processor failures. We, therefore, present in this book suitable fault tolerance techniques for these networks and analysis methods to determine which network topologies are more robust.

Software mistakes/bugs are, in practice, unavoidable, and consequently, some level of software fault tolerance is a must. This can be as simple as acceptance tests to check the reasonableness of the results before using them, or as complex as running two or more versions of the software (sequentially or in parallel). Programs also tend to have their state deteriorate after running for long periods of time and eventually crash. This situation can be avoided by periodically restarting the program, a process called rejuvenation. Unlike hardware faults, software faults are very hard to model. Still, a few such models have been developed and several of them are described.

Hardware fault-tolerance techniques can be quite costly to implement. In applications in which a complete and immediate masking of the effect of hardware faults (especially, of transient nature) is not necessary, checkpointing is an inexpensive alternative. For programs that run for a long time and for which re-execution upon a failure might be too costly, the program state can be saved (once or periodically) during the execution. Upon a failure occurrence, the system can roll back the program to the most recent checkpoint and resume its execution from that point. Various checkpointing techniques are presented and analyzed in the book.

Case studies illustrating the use of many of the fault-tolerance techniques described previously are presented, including Tandem, Stratus, Cassini, and microprocessors from IBM and Intel.

Two fault-tolerance topics that are rapidly increasing in practical importance, namely, defect tolerant VLSI design and fault tolerance in cryptographic devices are discussed. The increasing complexity of VLSI chip design has resulted in a situation in which manufacturing defects are unavoidable. If nothing is done to remedy this situation, the expected *yield* (the fraction of manufactured chips which are operational) will be very low. Thus, techniques to reduce the sensitivity of VLSI chips to defects have been developed, some of which are very similar to the hardware redundancy schemes.

For cryptographic devices, the need for fault tolerance is two-fold. Not only is it crucial that such devices (e.g., smart cards) operate in a fault-free manner in whatever environment they are used, but more importantly, they must stay secure. Fault-injection-based attacks on cryptographic devices have become the simplest and fastest way to extract the secret key from the device. Thus, the incorporation of fault tolerance is a must in order to keep cryptographic devices secure.

An important part of the design and evaluation process of a fault-tolerant system is to demonstrate that the system does indeed function at the advertised level of reliability. Often the designed system is too complex to develop analytical expressions of its reliability. If a prototype of the system has already been constructed, then fault-injection experiments can be performed and certain dependability attributes measured. If, however, as is very common, a prototype does not yet exist, statistical simulation must be used. Simulation programs for complex systems must be carefully designed to produce accurate results. We discuss the principles that should be followed when preparing a simulation program, and show how simulation results can be analyzed to infer system reliability.

1.5 Further Reading

Several textbooks and reference books on the topic of fault tolerance have been published in the past. See, for example, [2,4,5,9,10,13–16]. Journals have published several special issues on fault-tolerant computing, e.g., [7,8]. The major conference in the field is the *Conference on Dependable Systems and Networks* (DSN) [3]; this is a successor to the *Fault-Tolerant Computing Symposium* (FTCS).

The concept of computing being invisible everywhere appeared in [19], in the context of *pervasive computing*, that is, computing that pervades everyday living, without being obtrusive.

The definitions of the basic terms and measures appear in most of the textbooks mentioned above and in several probability and statistics books. For example, see [18]. Our definitions of fault and error are slightly different from those used in some of the references. A generally used definition of an error is that it is that part of the system state that leads to system failure. Strictly interpreted, this only applies to a system with state, i.e., with memory. We use the more encompassing definition of anything that can be construed as a manifestation of a fault. This wider interpretation allows purely combinational circuits, which are stateless, to generate errors.

One measure of dependability that we did not describe in the text is to consider everything from the perspective of the application. This approach was taken to define the measure known as *performability*. The application is used to define “accomplishment levels” L_1, L_2, \dots, L_n . Each of these represents a level of quality of service delivered by the application. For example, L_i may be defined as follows: “There are i system crashes during the time period $[0, T]$.” Now, the performance of the computer affects this quality (if it did not, by definition, it would have nothing to do with the application!). The approach taken by performability is to link the performance of the computer to the accomplishment level that this enables. Performability is then a vector, $(P(L_1), P(L_2), \dots, P(L_n))$, where $P(L_i)$ is the probability that the computer functions well enough to permit the application to reach up to accomplishment level L_i . For more on performability, see [6,11,12].

References

- [1] A. Avizienis and J. Laprie, "Dependable Computing: From Concepts to Design Diversity," *Proceedings of the IEEE*, Vol. 74, pp. 629–638, May 1986.
- [2] W. R. Dunn, *Practical Design of Safety-Critical Computer Systems*, Reliability Press, 2002.
- [3] *Dependable Systems and Networks (DSN) Conference*, <http://www.dsn.org>.
- [4] C. E. Ebeling, *An Introduction to Reliability and Maintainability Engineering*, McGraw-Hill, 1997.
- [5] J.-C. Geffroy and G. Motet, *Design of Dependable Computing Systems*, Kluwer Academic Publishers, 2002.
- [6] M.-C. Hsueh, R. K. Iyer, and K. S. Trivedi, "Performability Modeling Based on Real Data: A Case Study," *IEEE Transactions on Computers*, Vol. 37, pp. 478–484, April 1988.
- [7] *IEEE Computer*, Vol. 23, No. 5, July 1990. [Special issue on fault-tolerant systems]
- [8] *IEEE Transactions on Computers*, Vol. 41, February 1992; Vol. 47, April 1998; and Vol. 51, February 2002. [Special issues on fault-tolerant systems]
- [9] P. Jalote, *Fault Tolerance in Distributed Systems*, PTR Prentice Hall, 1994.
- [10] B. W. Johnson, *Design and Analysis of Fault-Tolerant Digital Systems*, Addison-Wesley, 1989.
- [11] J. F. Meyer, "On Evaluating the Performability of Degradable Computing Systems," *IEEE Transactions on Computers*, Vol. 29, pp. 720–731, August 1980.
- [12] J. F. Meyer, D. G. Furchtgott, and L. T. Wu, "Performability Evaluation of the SIFT Computer," *IEEE Transactions on Computers*, Vol. 29, pp. 501–509, June 1980.
- [13] D. K. Pradhan (Ed.), *Fault Tolerant Computer System Design*, Prentice Hall, 1996.
- [14] L. L. Pullum, *Software Fault Tolerance Techniques and Implementation*, Artech House, 2001.
- [15] D. P. Siewiorek and R. S. Swarz, *Reliable Computer Systems: Design and Evaluation*, A. K. Peters, 1998.
- [16] M. L. Shooman, *Reliability of Computer Systems and Networks: Fault Tolerance, Analysis, and Design*, Wiley-Interscience, 2001.
- [17] A. K. Somani and N. H. Vaidya, "Understanding Fault-tolerance and Reliability," *IEEE Computer*, Vol. 30, pp. 45–50, April 1997.
- [18] K. S. Trivedi, *Probability and Statistics with Reliability, Queuing, and Computer Science Applications*, John Wiley, 2002.
- [19] M. Weiser, "The Computer for the Twenty-first Century," *Scientific American*, pp. 94–104, September 1991. Available at: <http://www.ubiq.com/hypertext/weiser/SciAmDraft3.html>.