

Software Testing

[Carnegie Mellon University](#)
18-849b Dependable Embedded Systems
Spring 1999
Authors: [Jiantao Pan](#)
jpan@cmu.edu

Abstract:

Software testing is any activity aimed at evaluating an attribute or capability of a program or system and determining that it meets its required results. [\[Hetzel88\]](#) Although crucial to software quality and widely deployed by programmers and testers, software testing still remains an art, due to limited understanding of the principles of software. The difficulty in software testing stems from the complexity of software: we can not completely test a program with moderate complexity. Testing is more than just debugging. The purpose of testing can be quality assurance, verification and validation, or reliability estimation. Testing can be used as a generic metric as well. Correctness testing and reliability testing are two major areas of testing. Software testing is a trade-off between budget, time and quality.

Contents:

- [Introduction](#)
 - [Key Concepts](#)
 - [Taxonomy](#)

 - [Testing automation](#)
 - [When to stop testing?](#)
 - [Alternatives to testing](#)

 - [Available tools, techniques, and metrics](#)
 - [Relationship to other topics](#)
 - [Conclusions](#)
 - [Annotated Reference List & Further Reading](#)
-

Introduction

Software Testing is the process of executing a program or system with the intent of finding errors. [\[Myers79\]](#) Or, it involves any activity aimed at evaluating an attribute or capability of a program or system and determining that it meets its required results. [\[Hetzel88\]](#) Software is not unlike other physical processes where inputs are received and outputs are produced. Where software differs is in the manner in which it fails. Most physical systems fail in a fixed (and reasonably small) set of ways. By contrast, software can fail in many bizarre ways. Detecting all of the different failure modes for software is generally infeasible. [\[Rstcorp\]](#)

Unlike most physical systems, most of the defects in software are design errors, not manufacturing

defects. Software does not suffer from corrosion, wear-and-tear -- generally it will not change until upgrades, or until obsolescence. So once the software is shipped, the design defects -- or bugs -- will be buried in and remain latent until activation.

Software bugs will almost always exist in any software module with moderate size: not because programmers are careless or irresponsible, but because the complexity of software is generally intractable -- and humans have only limited ability to manage complexity. It is also true that for any complex systems, design defects can never be completely ruled out.

Discovering the design defects in software, is equally difficult, for the same reason of complexity. Because software and any digital systems are not continuous, testing boundary values are not sufficient to guarantee correctness. All the possible values need to be tested and verified, but complete testing is infeasible. Exhaustively testing a simple program to add only two integer inputs of 32-bits (yielding 2^{64} distinct test cases) would take hundreds of years, even if tests were performed at a rate of thousands per second. Obviously, for a realistic software module, the complexity can be far beyond the example mentioned here. If inputs from the real world are involved, the problem will get worse, because timing and unpredictable environmental effects and human interactions are all possible input parameters under consideration.

A further complication has to do with the dynamic nature of programs. If a failure occurs during preliminary testing and the code is changed, the software may now work for a test case that it didn't work for previously. But its behavior on pre-error test cases that it passed before can no longer be guaranteed. To account for this possibility, testing should be restarted. The expense of doing this is often prohibitive. [\[Rstcorp\]](#)

An interesting analogy parallels the difficulty in software testing with the pesticide, known as the Pesticide Paradox [\[Beizer90\]](#): *Every method you use to prevent or find bugs leaves a residue of subtler bugs against which those methods are ineffectual.* But this alone will not guarantee to make the software better, because the Complexity Barrier [\[Beizer90\]](#) principle states: *Software complexity (and therefore that of bugs) grows to the limits of our ability to manage that complexity.* By eliminating the (previous) easy bugs you allowed another escalation of features and complexity, but his time you have subtler bugs to face, just to retain the reliability you had before. Society seems to be unwilling to limit complexity because we all want that extra bell, whistle, and feature interaction. Thus, our users always push us to the complexity barrier and how close we can approach that barrier is largely determined by the strength of the techniques we can wield against ever more complex and subtle bugs. [\[Beizer90\]](#)

Regardless of the limitations, testing is an integral part in software development. It is broadly deployed in every phase in the software development cycle. Typically, more than 50% percent of the development time is spent in testing. Testing is usually performed for the following purposes:

- **To improve quality.**

As computers and software are used in critical applications, the outcome of a bug can be severe. Bugs can cause huge losses. Bugs in critical systems have caused airplane crashes, allowed space shuttle missions to go awry, halted trading on the stock market, and worse. Bugs can kill. Bugs can cause disasters. The so-called year 2000 (Y2K) bug has given birth to a cottage industry of consultants and programming tools dedicated to making sure the modern world doesn't come to a screeching halt on the first day of the next century. [\[Bugs\]](#) In a computerized embedded world, the quality and reliability of software is a matter of life and death.

Quality means the conformance to the specified design requirement. Being correct, the minimum requirement of quality, means performing as required under specified circumstances. Debugging, a narrow view of software testing, is performed heavily to find out design defects by the programmer.

The imperfection of human nature makes it almost impossible to make a moderately complex program correct the first time. Finding the problems and get them fixed [Kaner93], is the purpose of debugging in programming phase.

- **For Verification & Validation (V&V)**

Just as topic [Verification and Validation](#) indicated, another important purpose of testing is verification and validation (V&V). Testing can serve as metrics. It is heavily used as a tool in the V&V process. Testers can make claims based on interpretations of the testing results, which either the product works under certain situations, or it does not work. We can also compare the quality among different products under the same specification, based on results from the same test.

We can not test quality directly, but we can test related factors to make quality visible. Quality has three sets of factors -- functionality, engineering, and adaptability. These three sets of factors can be thought of as dimensions in the software quality space. Each dimension may be broken down into its component factors and considerations at successively lower levels of detail. Table 1 illustrates some of the most frequently cited quality considerations.

Functionality (exterior quality)	Engineering (interior quality)	Adaptability (future quality)
Correctness	Efficiency	Flexibility
Reliability	Testability	Reusability
Usability	Documentation	Maintainability
Integrity	Structure	

Table 1. Typical Software Quality Factors [Hetzel88]

Good testing provides measures for all relevant factors. The importance of any particular factor varies from application to application. Any system where human lives are at stake must place extreme emphasis on reliability and integrity. In the typical business system usability and maintainability are the key factors, while for a one-time scientific program neither may be significant. Our testing, to be fully effective, must be geared to measuring each relevant factor and thus forcing quality to become tangible and visible. [Hetzel88]

Tests with the purpose of validating the product works are named clean tests, or positive tests. The drawbacks are that it can only validate that the software works for the specified test cases. A finite number of tests can not validate that the software works for all situations. On the contrary, only one failed test is sufficient enough to show that the software does not work. Dirty tests, or negative tests, refers to the tests aiming at breaking the software, or showing that it does not work. A piece of software must have sufficient exception handling capabilities to survive a significant level of dirty tests.

A testable design is a design that can be easily validated, falsified and maintained. Because testing is a rigorous effort and requires significant time and cost, design for testability is also an important design rule for software development.

- **For reliability estimation [Kaner93] [Lyu95]**

Software reliability has important relations with many aspects of software, including the structure, and the amount of testing it has been subjected to. Based on an operational profile (an estimate of the relative frequency of use of various inputs to the program [Lyu95]), testing can serve as a statistical

sampling method to gain failure data for reliability estimation.

Software testing is not mature. It still remains an art, because we still cannot make it a science. We are still using the same testing techniques invented 20-30 years ago, some of which are crafted methods or heuristics rather than good engineering methods. Software testing can be costly, but not testing software is even more expensive, especially in places that human lives are at stake. Solving the software-testing problem is no easier than solving the Turing halting problem. We can never be sure that a piece of software is correct. We can never be sure that the specifications are correct. No verification system can verify every correct program. We can never be certain that a verification system is correct either.

Key Concepts

Taxonomy

There is a plethora of testing methods and testing techniques, serving multiple purposes in different life cycle phases. Classified by purpose, software testing can be divided into: correctness testing, performance testing, reliability testing and security testing. Classified by life-cycle phase, software testing can be classified into the following categories: requirements phase testing, design phase testing, program phase testing, evaluating test results, installation phase testing, acceptance testing and maintenance testing. By scope, software testing can be categorized as follows: unit testing, component testing, integration testing, and system testing.

Correctness testing

Correctness is the minimum requirement of software, the essential purpose of testing. Correctness testing will need some type of oracle, to tell the right behavior from the wrong one. The tester may or may not know the inside details of the software module under test, e.g. control flow, data flow, etc. Therefore, either a white-box point of view or black-box point of view can be taken in testing software. We must note that the black-box and white-box ideas are not limited in correctness testing only.

- **Black-box testing**

The black-box approach is a testing method in which test data are derived from the specified functional requirements without regard to the final program structure. [Perry90] It is also termed data-driven, input/output driven [Myers79], or requirements-based [Hetzel88] testing. Because only the functionality of the software module is of concern, black-box testing also mainly refers to functional testing -- a testing method emphasized on executing the functions and examination of their input and output data. [Howden87] The tester treats the software under test as a black box -- only the inputs, outputs and specification are visible, and the functionality is determined by observing the outputs to corresponding inputs. In testing, various inputs are exercised and the outputs are compared against specification to validate the correctness. All test cases are derived from the specification. No implementation details of the code are considered.

It is obvious that the more we have covered in the input space, the more problems we will find and therefore we will be more confident about the quality of the software. Ideally we would be tempted to exhaustively test the input space. But as stated above, exhaustively testing the combinations of valid inputs will be impossible for most of the programs, let alone considering invalid inputs, timing, sequence, and resource variables. Combinatorial explosion is the major roadblock in functional testing. To make things worse, we can never be sure whether the specification is either correct or

complete. Due to limitations of the language used in the specifications (usually natural language), ambiguity is often inevitable. Even if we use some type of formal or restricted language, we may still fail to write down all the possible cases in the specification. Sometimes, the specification itself becomes an intractable problem: it is not possible to specify precisely every situation that can be encountered using limited words. And people can seldom specify clearly what they want -- they usually can tell whether a prototype is, or is not, what they want after they have been finished. Specification problems contributes approximately 30 percent of all bugs in software. [Beizer95]

The research in black-box testing mainly focuses on how to maximize the effectiveness of testing with minimum cost, usually the number of test cases. It is not possible to exhaust the input space, but it is possible to exhaustively test a subset of the input space. Partitioning is one of the common techniques. If we have partitioned the input space and assume all the input values in a partition is equivalent, then we only need to test one representative value in each partition to sufficiently cover the whole input space. Domain testing [Beizer95] partitions the input domain into regions, and consider the input values in each domain an equivalent class. Domains can be exhaustively tested and covered by selecting a representative value(s) in each domain. Boundary values are of special interest. Experience shows that test cases that explore boundary conditions have a higher payoff than test cases that do not. Boundary value analysis [Myers79] requires one or more boundary values selected as representative test cases. The difficulties with domain testing are that incorrect domain definitions in the specification can not be efficiently discovered.

Good partitioning requires knowledge of the software structure. A good testing plan will not only contain black-box testing, but also white-box approaches, and combinations of the two.

- **White-box testing**

Contrary to black-box testing, software is viewed as a white-box, or glass-box in white-box testing, as the structure and flow of the software under test are visible to the tester. Testing plans are made according to the details of the software implementation, such as programming language, logic, and styles. Test cases are derived from the program structure. White-box testing is also called glass-box testing, logic-driven testing [Myers79] or design-based testing [Hetzel88].

There are many techniques available in white-box testing, because the problem of intractability is eased by specific knowledge and attention on the structure of the software under test. The intention of exhausting some aspect of the software is still strong in white-box testing, and some degree of exhaustion can be achieved, such as executing each line of code at least once (statement coverage), traverse every branch statements (branch coverage), or cover all the possible combinations of true and false condition predicates (Multiple condition coverage). [Parrington89]

Control-flow testing, loop testing, and data-flow testing, all maps the corresponding flow structure of the software into a directed graph. Test cases are carefully selected based on the criterion that all the nodes or paths are covered or traversed at least once. By doing so we may discover unnecessary "dead" code -- code that is of no use, or never get executed at all, which can not be discovered by functional testing.

In mutation testing, the original program code is perturbed and many mutated programs are created, each contains one fault. Each faulty version of the program is called a mutant. Test data are selected based on the effectiveness of failing the mutants. The more mutants a test case can kill, the better the test case is considered. The problem with mutation testing is that it is too computationally expensive to use. The boundary between black-box approach and white-box approach is not clear-cut. Many testing strategies mentioned above, may not be safely classified into black-box testing or white-box testing. It is also true for transaction-flow testing, syntax testing, finite-state testing, and many other testing strategies not discussed in this text. One reason is that all the above techniques will need some knowledge of the specification of the software under test. Another reason is that the idea of

specification itself is broad -- it may contain any requirement including the structure, programming language, and programming style as part of the specification content.

We may be reluctant to consider random testing as a testing technique. The test case selection is simple and straightforward: they are randomly chosen. Study in [Duran84] indicates that random testing is more cost effective for many programs. Some very subtle errors can be discovered with low cost. And it is also not inferior in coverage than other carefully designed testing techniques. One can also obtain reliability estimate using random testing results based on operational profiles. Effectively combining random testing with other testing techniques may yield more powerful and cost-effective testing strategies.

Performance testing

Not all software systems have specifications on performance explicitly. But every system will have implicit performance requirements. The software should not take infinite time or infinite resource to execute. "Performance bugs" sometimes are used to refer to those design problems in software that cause the system performance to degrade.

Performance has always been a great concern and a driving force of computer evolution. Performance evaluation of a software system usually includes: resource usage, throughput, stimulus-response time and queue lengths detailing the average or maximum number of tasks waiting to be serviced by selected resources. Typical resources that need to be considered include network bandwidth requirements, CPU cycles, disk space, disk access operations, and memory usage [Smith90]. The goal of performance testing can be performance bottleneck identification, performance comparison and evaluation, etc. The typical method of doing performance testing is using a benchmark -- a program, workload or trace designed to be representative of the typical system usage. [Vokolos98]

Reliability testing

Software reliability refers to the probability of failure-free operation of a system. It is related to many aspects of software, including the testing process. Directly estimating software reliability by quantifying its related factors can be difficult. Testing is an effective sampling method to measure software reliability. Guided by the operational profile, software testing (usually black-box testing) can be used to obtain failure data, and an estimation model can be further used to analyze the data to estimate the present reliability and predict future reliability. Therefore, based on the estimation, the developers can decide whether to release the software, and the users can decide whether to adopt and use the software. Risk of using software can also be assessed based on reliability information. [Hamlet94] advocates that the primary goal of testing should be to measure the dependability of tested software.

There is agreement on the intuitive meaning of dependable software: it does not fail in unexpected or catastrophic ways. [Hamlet94] Robustness testing and stress testing are variances of reliability testing based on this simple criterion.

The robustness of a software component is the degree to which it can function correctly in the presence of exceptional inputs or stressful environmental conditions. [IEEE90] Robustness testing differs with correctness testing in the sense that the functional correctness of the software is not of concern. It only watches for robustness problems such as machine crashes, process hangs or abnormal termination. The oracle is relatively simple, therefore robustness testing can be made more portable and scalable than correctness testing. This research has drawn more and more interests recently, most of which uses commercial operating systems as their target, such as the work in [Koopman97] [Kropp98] [Ghosh98] [Devale99] [Koopman99].

Stress testing, or load testing, is often used to test the whole system rather than the software alone. In such tests the software or system are exercised with or beyond the specified limits. Typical stress includes resource exhaustion, bursts of activities, and sustained high loads.

Security testing

Software quality, reliability and security are tightly coupled. Flaws in software can be exploited by intruders to open security holes. With the development of the Internet, software security problems are becoming even more severe.

Many critical software applications and services have integrated security measures against malicious attacks. The purpose of security testing of these systems include identifying and removing software flaws that may potentially lead to security violations, and validating the effectiveness of security measures. Simulated security attacks can be performed to find vulnerabilities.

Testing automation

Software testing can be very costly. Automation is a good way to cut down time and cost. Software testing tools and techniques usually suffer from a lack of generic applicability and scalability. The reason is straight-forward. In order to automate the process, we have to have some ways to generate oracles from the specification, and generate test cases to test the target software against the oracles to decide their correctness. Today we still don't have a full-scale system that has achieved this goal. In general, significant amount of human intervention is still needed in testing. The degree of automation remains at the automated test script level.

The problem is lessened in reliability testing and performance testing. In robustness testing, the simple specification and oracle: doesn't crash, doesn't hang suffices. Similar simple metrics can also be used in stress testing.

When to stop testing?

Testing is potentially endless. We can not test till all the defects are unearthed and removed -- it is simply impossible. At some point, we have to stop testing and ship the software. The question is when.

Realistically, testing is a trade-off between budget, time and quality. It is driven by profit models. The pessimistic, and unfortunately most often used approach is to stop testing whenever some, or any of the allocated resources -- time, budget, or test cases -- are exhausted. The optimistic stopping rule is to stop testing when either reliability meets the requirement, or the benefit from continuing testing cannot justify the testing cost. [Yang95] This will usually require the use of reliability models to evaluate and predict reliability of the software under test. Each evaluation requires repeated running of the following cycle: failure data gathering -- modeling -- prediction. This method does not fit well for ultra-dependable systems, however, because the real field failure data will take too long to accumulate.

Alternatives to testing

Software testing is more and more considered a problematic method toward better quality. Using testing to locate and correct software defects can be an endless process. Bugs cannot be completely ruled out. Just as the complexity barrier indicates: chances are testing and fixing problems may not necessarily improve the quality and reliability of the software. Sometimes fixing a problem may introduce much more severe problems into the system, happened after bug fixes, such as the telephone outage in California and eastern seaboard in 1991. The disaster happened after changing 3

lines of code in the signaling system.

In a narrower view, many testing techniques may have flaws. Coverage testing, for example. Is code coverage, branch coverage in testing really related to software quality? There is no definite proof. As early as in [Myers79], the so-called "human testing" -- including inspections, walkthroughs, reviews -- are suggested as possible alternatives to traditional testing methods. [Hamlet94] advocates inspection as a cost-effect alternative to unit testing. The experimental results in [Basili85] suggests that code reading by stepwise abstraction is at least as effective as on-line functional and structural testing in terms of number and cost of faults observed.

Using formal methods to "prove" the correctness of software is also an attracting research direction. But this method can not surmount the complexity barrier either. For relatively simple software, this method works well. It does not scale well to those complex, full-fledged large software systems, which are more error-prone.

In a broader view, we may start to question the utmost purpose of testing. Why do we need more effective testing methods anyway, since finding defects and removing them does not necessarily lead to better quality. An analogy of the problem is like the car manufacturing process. In the craftsmanship epoch, we make cars and hack away the problems and defects. But such methods were washed away by the tide of pipelined manufacturing and good quality engineering process, which makes the car defect-free in the manufacturing phase. This indicates that engineering the design process (such as clean-room software engineering) to make the product have less defects may be more effective than engineering the testing process. Testing is used solely for quality monitoring and management, or, "design for testability". This is the leap for software from craftsmanship to engineering.

Available tools, techniques, and metrics

There are an abundance of software testing tools exist. The correctness testing tools are often specialized to certain systems and have limited ability and generality. Robustness and stress testing tools are more likely to be made generic.

Mothora [DeMillo91] is an automated mutation testing tool-set developed at Purdue University. Using *Mothora*, the tester can create and execute test cases, measure test case adequacy, determine input-output correctness, locate and remove faults or bugs, and control and document the test.

NuMega's Boundschecker [NuMega99] *Rational's Purify* [Rational99]. They are run-time checking and debugging aids. They can both check and protect against memory leaks and pointer problems.

Ballista COTS Software Robustness Testing Harness [Ballista99]. The *Ballista* testing harness is an full-scale automated robustness testing tool. The first version supports testing up to 233 POSIX function calls in UNIX operating systems. The second version also supports testing of user functions provided that the data types are recognized by the testing server. The *Ballista* testing harness gives quantitative measures of robustness comparisons across operating systems. The goal is to automatically test and harden Commercial Off-The-Shelf (COTS) software against robustness failures.

Relationship to other topics

Software testing is an integrated part in software development. It is directly related to software quality. It has many subtle relations to the topics that software, software quality, software reliability and system reliability are involved.

Related topics

- [Software reliability](#) Software testing is closely related to software reliability. Software reliability can be augmented by testing. Also testing can be served as a metric for software reliability.
- [Fault injection](#) Fault injection can be considered a special way of testing. Fault injection and testing are usually combined and performed to validate the reliability of critical fault-tolerant software and hardware.
- [Verification, validation and certification](#) The purpose of software testing is not only for revealing bugs and eliminate them. It is also a tool for verification, validation and certification.

Conclusions

- Software testing is an art. Most of the testing methods and practices are not very different from 20 years ago. It is nowhere near maturity, although there are many tools and techniques available to use. Good testing also requires a tester's creativity, experience and intuition, together with proper techniques.
- Testing is more than just debugging. Testing is not only used to locate defects and correct them. It is also used in validation, verification process, and reliability measurement.
- Testing is expensive. Automation is a good way to cut down cost and time. Testing efficiency and effectiveness is the criteria for coverage-based testing techniques.
- Complete testing is infeasible. Complexity is the root of the problem. At some point, software testing has to be stopped and product has to be shipped. The stopping time can be decided by the trade-off of time and budget. Or if the reliability estimate of the software product meets requirement.
- Testing may not be the most effective method to improve software quality. Alternative methods, such as inspection, and clean-room engineering, may be even better.

Annotated References

- [Ballista99] <http://www.cs.cmu.edu/afs/cs/project/edrc-ballista/www/>

Ballista COTS Software Robustness Testing Harness homepage.

- [Basili85] Victor R. Basili, Richard W. Selby, Jr. "Comparing the Effectiveness of Software Testing Strategies",

Technical Report, Department of Computer Science, University of Maryland, College Park, 1985.

- [Beizer90] Boris Beizer, Software Testing Techniques. Second edition. 1990
- A very comprehensive book on the testing techniques. Many testing techniques are

enumerated and discussed in detail. Domain testing, data-flow testing, transactin-flow testing, syntax testing, logic-based testing, etc.

- [Beizer95] Beizer, Boris, *Black-box Testing: techniques for functional testing of software and systems*. Publication info: New York : Wiley, c1995. ISBN: 0471120944 Physical description: xxv, 294 p.: ill. ; 23 cm.

This book is a comprehensive introduction to various methods of testing, using intuitive examples. Complete coverage of all important testing techniques, and up-to-date. The focus is black-box/functional testing. The author is an internationally known software consultant with almost four decades of experience in the computer industry.

- [Duran84] Joe W. Duran, Simeon C. Ntafos, "An Evaluation of Random Testing", IEEE Transactions on Software Engineering, Vol. SE-10, No. 4, July 1984, pp438-443.

This research investigates the effectiveness of random testing, in comparison of other testing methods. It convinces us that random testing method is much more powerful than it appears and should be deployed more.

- [Hetzel88] Hetzel, William C., *The Complete Guide to Software Testing, 2nd ed.* Publication info: Wellesley, Mass. : QED Information Sciences, 1988. ISBN: 0894352423. Physical description: ix, 280 p. : ill ; 24 cm.

This book is a good guide to software testing. But it may not be as complete a guide as it was 10 years ago.

- [Howden87] William E. Howden. Functional program Testing and Analysis. McGraw-Hill, 1987.

In-depth discussion about functional testing throughout all product life-cycle.

- [IEEE90] IEEE Standard Glossary of Software Engineering Terminology (IEEE Std 610.12-1990), IEEE Computer Soc., Dec. 10, 1990.

Terminology standard.

- [Kaner93] Cem Kaner, Testing Computer Software. 1993.
- Discusses the purpose and techniques for software testing.
- [Lyu95] Michael R. Lyu , Handbook of Software Reliability Engineering. McGraw-Hill publishing, 1995, ISBN 0-07-039400-8

See topic [Software Reliability](#) reference.

- [Myers79] Myers, Glenford J., *The art of software testing*, Publication info: New York : Wiley, c1979. ISBN: 0471043281 Physical description: xi, 177 p. : ill. ; 24 cm.

Classical book on software testing and still has its influence today. Software testing still remains an art, since the first day. We still don't know how to make it a science.

- [NuMega99] <http://www.numega.com/devcenter/bc.shtml>
- Introduction of the tools BoundsChecker by NuMega
- [Parrington89] Norman Parrington and Marc Roper, *Understanding Software Testing*, Published by John Willey & Sons, 1989. ISBN:0-7458-0533-7; 0-470-21462-7
- This book surveys current software methods and techniques available. The authors show how, when carried out effectively in the commercial sphere, software testing can result in higher quality software products, more satisfied users and lower costs, which in turn leads to more accurate and reliable results.
- [Rational99] http://www.rational.com/products/purify_unix/index.jtmpl
- Introduces the Purify tool.
- [Rstcorp] http://www.rstcorp.com/definitions/software_testing.html
- A definition of and introduction to software testing.
- [PERRY90] A standard for testing application software, William E. Perry, 1990
- This book summarizes and standardizes many testing techniques.
- [Musa97] Software-reliability-engineered testing practice (tutorial); John D. Musa; Proceedings of the 1997 international conference on Software engineering , 1997, Pages 628 - 629

This 2-page short tutorial gets us started with an introduction of the basic ideas underlying SRE testing of software. Definitions of reliability, failure, fault, severity class, operation, operational profile, feature testing, load testing and regression testing are included. Four steps, developing operational profiles, preparing for testing, executing tests and interpreting failure data are introduced. This article is a good introduction to testing-for-reliability.

- [Vokolos98] Performance testing of software systems; Filippos I. Vokolos, and Elaine J. Weyuker; Proceedings of the first international workshop on Software and performance , 1998, Pages 80 - 87

Performance testing papers are very rare in literature, although performance is what most people care about. This paper gives a discussion on approaches to software performance testing. A case study describing the experience of using these approaches for testing the performance of a system used as a gateway in a large industrial client/server transaction processing application is presented. It is not a very representative paper, though.

- [Bugs] <http://www.cnet.com/Content/Features/Dlife/Bugs/?dd>
- Interesting readings about bugs' life. "Bugs -- How they breed and the damage they do" and "10 great bugs of history", etc
- [DeMillo91] Progress toward automated software testing; Richard A. DeMillo; Proceedings of the 13th international conference on Software engineering , 1991, Page 180

This paper mainly discusses their toolset Mothra, integrated environment for automated software

validation. Mothra uses program mutation to drive the tests. Few human intervention needed. Automated test generation.

- [Hamlet94] Dick Hamlet; Foundations of software testing: dependability theory; Proceedings of the second ACM SIGSOFT symposium on Foundations of software engineering , 1994, Pages 128 - 139
- The point "four epochs of testing" is interesting.
- [Koopman97] Philip Koopman, John Sung, Christopher Dingman, Daniel Siewiorek, Ted Marz. Comparing Operating Systems Using Robustness Benchmarks. 16th IEEE Symposium on Reliable Distributed Systems, Durham, NC, October 22-24, 1997, pp.72-79.
- Precursor work of Ballista project. Tests selected POSIX functions for robustness.
- [Koopman99] Philip Koopman, John DeVale. Comparing the Robustness of POSIX Operating Systems. Proceedings of FTCS'99, 15-18 June 1999, Madison, Wisconsin.
- Ballista paper: multi-version comparison method to find silent and hindering failures.
- [Kropp98] Kropp, N. P.; Koopman, P. J.; Siewiorek, D. P. Automated robustness testing of off-the-shelf software components. Twenty-eighth Annual International Symposium on Fault-Tolerant Computing (Cat. No.98CB36224)
- Ballista paper. Full scale, automated robustness test harness. Can test up to 233 POSIX function calls. Tested 15 operating systems from 10 vendors. Between 42% and 63% of components tested had robustness problems, with a normalized failure rate ranging from 10% to 23% of tests conducted. Robustness testing could be used by developers to measure and improve robustness, or by consumers to compare the robustness of competing COTS component libraries.
- [Ghosh98] Testing the Robustness of Windows NT Software
- A simple heuristic method to test the robustness of some NT and GNU library.
- [DeVale99] John DeVale, Philip Koopman & David Guttendorf. The Ballista Software Robustness Testing Service. Proceedings of TCS'99, Washington DC.
- Describes the Ballista testing web server and client architecture.
- [Koehnemann93] Harry Koehnemann, and Timothy Lindquist; Towards target-level testing and debugging tools for embedded software; Conference proceedings on TRI-Ada '93 , 1993, Page 288
- Argues that current debugging method is not efficient for embedded software and propose an improved method.
- [Smith90] Smith, C. U. Performance Engineering of Software Systems. Addison-Wesley, 1990.
- A software performance testing paper.
- [Yang95] Yang, M.C.K.; Chao, A. Reliability-estimation and stopping-rules for software

testing, based on repeated appearances of bugs; IEEE Transactions on Reliability, vol.44, no.2, p. 315-21, 1995

When to stop testing? The choice is usually based on one of two decision criteria: (1) when the reliability has reached a given threshold, and (2) when the gain in reliability cannot justify the testing cost. Various stopping rules and software reliability models are compared by their ability to deal with these two criteria. Two new stopping rules, initiated by theoretical study of the optimal stopping rule based on cost, are more stable than other rules for a large variety of bug structures.

Future Reading

- Software testing slides of RST corporation by Jeff Voas:
<http://www.rstcorp.com/presentations/tampa98/>
- Says testing and debugging can worsen reliability. 15% chance
- [Pham95] Software Reliability and Testing, pp29
- A collection of papers on software reliability testing.
- <http://www.cs.jmu.edu/users/foxcj/cs555/Unit12/Testing/index.htm> slide show of software testing.
- comprehensive annotated www testing resource list with http links
- <http://www.eg3.com/softd/index.htm>
- Embedded s/w net resources

Comprehensive annotated www testing resource list:

- <http://www.thegrid.net/tech/softd/softtest.htm>
- Software testing and software quality are major issues. Here is the best the net has to offer on them. keywords include: testing, software testing, software reliability, software verification. This page discusses key issues such as testing, software testing, software reliability, software verification.. testing, software testing, software reliability, software verification. are discussed on this page.
- <http://www.io.com/~wazmo/qa.html>
- <http://isse.gmu.edu/faculty/ofut/rsrch/mut.html>
- Mutation testing page and list of papers

Security Links:

- <http://cs-www.ncsl.nist.gov/>
- Computer Security Resource Clearinghouse, National Institute of Standards and Technology
- <http://www.sni.net/~deckm/>

- Clean-room software engineering, inc.

Fun stuff:

- <http://www.mtsu.edu/cgi-bin/users/storm/researchers.pl>
- Meet the people: Who's Who in Software Testing Research
- <http://www.bugnet.com/>
- WWW resource for bug fixes, patches, and news.
- <http://hongkong1.cnet.com/Briefs/Guidebook/Bugs2/ss05i.html>
- 10 greatest bugs in history.
So says Ivars Peterson, author of Fatal Defect: [Chasing Killer Computer Bugs](#). As he concludes in his book, "The fact that we can never be sure that a computer system will function flawlessly constitutes a fatal defect. It limits what we can hope to achieve by using computers as our servants and surrogates. As computer-controlled systems become more complex and thoroughly entwined in the fabric of our lives, their potential for costly, life-threatening failures keeps growing."

[Index of other topics](#)

[Home page](#)
