# Department of

# Computer Science

## The Relationship Between Test Coverage and Reliability

Yashwant K. Malaiya, Naixin Li, Jim Bieman,
Rick Karcich, Bob Skibbe

# Colorado State University

# The Relationship Between Test Coverage and Reliability

Yashwant K. Malaiya*

Naixin Li

Jim Bieman†

Computer Science Dept.

Colorado State University

Fort Collins, CO 80523

malaiya@cs.colostate.edu

Rick Karcich

Bob Skibbe

StorageTek

2270 South 88th Street

Louisville, CO 80028-2286

(303) 673-6223

Rick_Karcich@stortek.com

March 15, 1994

### Abstract

In this paper, we model the relation between testing effort, coverage and reliability. We present a logarithmic model that relates testing effort to test coverage (block, branch, c-use or p-use). The model is based on the hypothesis that the enumerables (like branches or blocks) for any coverage measure have different detectability, just like defects have different detectability. This model allows us to relate a test coverage measure directly with defect coverage. Data sets for programs with real defects are used to validate the model. The results are consistent with the known inclusion relationships among block, branch and p-use coverage measures. We show how defect density controls *time to next failure*.

The model can eliminate the variables like test application strategy from consideration. It is suitable for high reliability applications where automatic (or manual) test generation is used to cover enumerables which have not yet been tested.

# 1 Introduction

Developers can increase the reliability of software systems by measuring reliability as early as possible during development. Early indications of reliability problems allow developers to correct errors and make process adjustments.

---

Reliability can be first measured as soon as running code exists. To quantify reliability during testing, the code (or portion of code) is executed using inputs randomly selected following an operational distribution. Then, appropriate reliability models can be used to predict the amount of effort required to satisfy product reliability requirements.

The needs of early reliability measurement and modeling unfortunately are not met by common testing practices. The focus of testing is on finding defects, and defects can be often found much faster by non-random methods [bei90]. Testing is directed towards inputs and program components where errors are more likely. For example, testing may be conducted to insure that particular portions of the program and/or boundary cases are covered.

Models that can measure and predict reliability based on the status of non-random testing are clearly needed. Reliability models will be affected by:

- the testing strategy: Test coverage may be based on the functional specification (black-box), or it may be based on internal program structure (white-box). Strategies can vary in their ability to find defects.

- the relationship between calendar time and execution time: The testing process can be accelerated through the possibly parallel, intensive execution of tests at a faster rate that would occur during operational use. However, testing, in some environments, might occur at a slower rate than normal operational system use.

- the testing of rarely executed modules: Such modules include exception handling or error recovery routines. These modules rarely run, and are notoriously difficult to test. Yet, they are critical components of a system that must be highly reliable. Only by forcing the coverage of such critical components, can reliability be predicted at very high levels.

Intuition and empirical evidence suggests that test coverage must be related to reliability. Yet, the connection between structure based measurements, like test coverage, and reliability is still not well understood.

Ramsey and Basili [ram85] experimented with different permutations of the same test set and collected data relating the number of tests to statement coverage growth. A variety of models were attempted to fit the data. The best fit was obtained using the Goel and Okumoto's exponential model (GO model). Ramsey and Basili also noticed that the faults revealed in a procedure are independent of the number of times the procedure is exercised. Their results support the view that structural (procedure) coverage may be used as an indicator of testing thoroughness. However, they did not model the relation between test coverage and software reliability.

Dalal et al [dhk93] also examined the correlation between test coverage and the error removal rate. They give a scatter plot of the number of faults detected during system testing versus the block coverages achieved during unit testing for 28 program modules. The plot clearly shows that modules covered more thoroughly during unit testing are much less likely to contain errors.

Vouk [vou92] found that the relation between structural coverage and fault coverage is a variant of the Rayleigh distribution. He assumed that the fault detection rate during testing is proportional to the number of faults present in the software and test coverage values including block, branch, data-flow, and functional group coverage. Vouk's experimental results, however, support the use of a more general Weibull distribution. Using the Rayleigh model, Vouk computed that, in terms of error removal capability, the relative power of the coverage measures block:p-use:DUD-chains is 1:2:6.

Chen et al [chm92] incorporate structural coverage into traditional time-based software reliability models (SRMs). Their model only includes test cases that increase coverage. The included test effort data is used to fit existing time-based models. Thus, they avoided the overestimation from traditional time-based SRMs due to the saturation effect of testing strategies. They do not relate test coverage directly to the error removal process as we do here.

Assuming random testing, Piwowarski, Ohba and Caruso [poc93] analyzed block coverage growth during function test, and derived an exponential model relating the number of tests to block coverage Their model is equivalent to the GO model attempted in [ram85]. They also derived an exponential model relating the covering frequency to the error removal ratio. However, the utility of the model relies on prior knowledge of the error distribution over different functional groups in a product.

Frankl and Weiss [fra93] compared the fault exposing capability of branch coverage and data flow coverage criteria. They found that for 4 out of 7 programs, the effectiveness of a test in exposing an error is positively correlated with the two coverage measures. They also observed complex relationships between test coverage growth and the probability of exposing an error for a test set. Since the 7 programs they used are very small and they only considered subtle errors, the result can not be extrapolated to practical software. They did not model the relation between test coverage and fault coverage.

The Leone test coverage model given in [neu93] is a weighted average of four different coverage metrics achieved during test phases: lines of executable code, independent test paths, functions/requirements, and hazard test cases. The weighted average is used as an indicator of software reliability. The model assumes that full coverage of all four metrics implies that the software tested is 100% reliable. In reality, such software may have some remaining faults. A similar approach, but with different coverage metrics, was taken to provide a test quality report [pos93].

In this paper, we explore the connection between test coverage and reliability. We develop a model that relates test coverage to defect coverage. With this model and with knowledge of the operational profile, we can predict reliability from test coverage measures.

## 2    Coverage of Enumerables

The concept of test coverage is applicable for both hardware and software. In hardware, coverage is measured in terms of the number of possible faults covered. For example, each

node in a digital system can possibly be stuck-at 0 or stuck-at 1. A stuck-at test coverage of 80% means that the tests applied would have detected any one of the 80% faults covered.

Test coverage in software is measured in terms of structural or data-flow units that have been exercised. These units can be statements (or blocks), branches, etc. as defined below:

- Statement (or block) coverage: the fraction of the total number of statements that have been executed by the test data.

- Branch (or decision) coverage: the fraction of the total number of branches that have been executed by the test data.
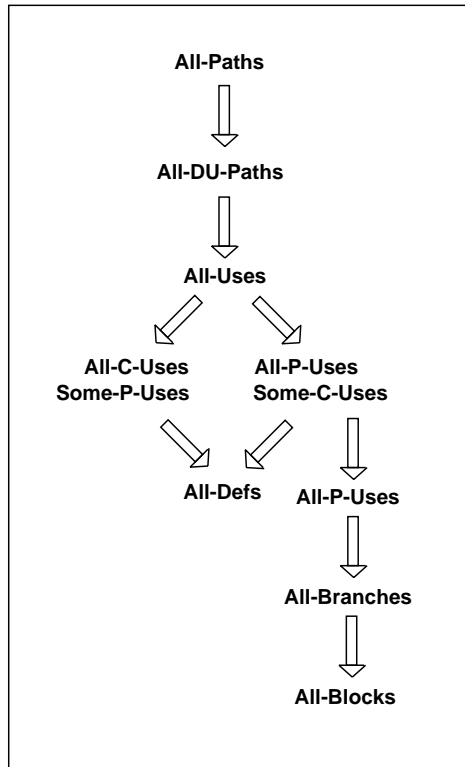


Figure 1: The subsumption relationships of different coverage criteria [weu84]

- C-use coverage: the fraction of the total number of c-uses that have been covered by one c-use path during testing. A c-use path is a path through a program from each point where the value of a variable is modified to each computation use or c-use (without the variable being modified along the path).

- P-use coverage: the fraction of the total number of p-uses that have been covered by one p-use path during testing. A p-use path is a path from each point where the value of a variable is modified to each p-use, a use in a predicate or decision (again, without modifications to the variable along the path).

When such a unit is exercised, it is possible that an associated fault is detected. Counting the number of units gives us a measure of the extent of sampling. The defect coverage in

software can be defined in an analogous manner, it is the fraction of actual defects initially present, that would be detected by a given test set.

Table 1: The complexity (test length) for achieving different coverage criteria [weu93]

| Coverage Criterion | Upper bound | Observed |
|---|---|---|
| All-Blocks | $d + 1$ | |
| All-Branches | $d + 1$ | |
| All-P-Uses | $\frac{1}{4}(d^2 + 4d + 3)$ | $0.38d + 3.17$ |
| All-Defs | $m + (i \times n)$ | |
| All-P-Uses/Some-C-Uses | $\frac{1}{4}(d^2 + 4d + 3)$ | |
| All-C-Uses/Some-P-Uses | $\frac{1}{4}(t^2 + 4d + 3)$ | $0.36d + 2.82$ |
| All-Uses | $\frac{1}{4}(d^2 + 4d + 3)$ | $0.39d + 3.76$ |
| All-DU-Paths | $2^d$ | $0.49d + 4.03$ $2^d$ |
| All-Paths | $\infty$ | |

$n$: variables, $m$: assignments, $i$: input statements, $d$: binary decisions.

Figure 1 is taken from [weu84], which shows the relative strength of some well-known criteria. If there is a directed path from criteria $A$ to criteria $B$, then test sets that meet criteria $A$ (complete coverage) are guaranteed to satisfy criteria $B$. Table 1 shows the upper bounds on the test length [weu84] to satisfy these criteria and the observed complexities [weu93] for some of the criteria. Such knowledge of complexities can be very useful for testers for selecting appropriate test coverage criteria. The upper bound for all-du-paths was reached in one subroutine out of 143 considered by Bieman and Schultz [bisc89,bisc92].

In order to keep the following discussion general, we will use the term *enumerable*. For branch coverage, the enumerables are branches, for defect coverage the enumerables are defects and so on. In this paper, the term enumerable-type implies one of these: defects, blocks, branches, c-uses and p-uses. We will use superscript $i$, $i = 0$ to 4, to indicate one of the five types in this sequence: 0: defects, 1: blocks, 2:branches, 3: c-uses, 4: p-uses.

# 3    Detectability Profiles of Enumerables

The coverage achieved by a set of tests, under a testing strategy, depends not only on the number of tests applied (or, equivalently, the testing time) but also on the distribution of *testability* values of the enumerables. A statement which is reached more easily, is more testable. Such statements are likely to get covered (i.e. exercised at lease once) with only a small number of tests. Testability also depends on the likelihood that a fault that is reached actually causes a failure [voas92]. On the other hand a statement which gets executed in rare situations has low testability. It may not get exercised by most of the tests which would

normally be applied. As testing progresses, the distribution of testability values will shift. The easy-to-test enumerables will get covered early during testing, and are thus removed from consideration. The enumerables remaining to be covered include a larger fraction of hard-to-test enumerables. Thus the growth of coverage will be slow.

**Definition**: Detectability of an enumerable $d_l^j$ is the probability that the $l$-th enumerable of type $j$ will be exercised by a randomly chosen test.

The distribution of detectability values in the system under test is given by the *detectability profile*. The detectability profile concept was introduced by Malaiya and Yang [mal84] and has been used to characterize testing of hardware [wag87] as well as software [mvs93]. A continuous version of the detectability profile was defined by Seth, Agrawal and Farhat [set90]. For convenience, we use the normalized detectability profile (NDP) as defined below.

**Definition**: The discrete NDP for the system under test is given by the vector,

$$P^j = \{p_{d1}^j, p_{d2}^j, ..., p_{di}^j, ..., p_{du}^j\} \qquad\qquad d_{i-1}^j < d_i^j < d_{i+1}^j \qquad\qquad (1)$$

where $p_{di}^j$ is the fraction of all enumerables of type $j$ which have detectability exactly equal to $di$. Thus $p_{0.3}^{branch}$ represents the fraction of all branches which have detectability equal to 0.3. In Equation 1, $du$ refers to the detectability value of unity (1), which is the highest detectability value possible. Notice that $\sum_{di=0}^1 p_{di}^j = 1$ since all fractions added will be unity.

Notice that a detectability value of 0 is possible, since a branch might be infeasible, or a defect might not be testable because of redundancy in implementation. Detectability profiles of several digital circuits [mal84, wag87] and software systems [tra92, dun86] have been compiled by researchers.

If the number of enumerables is very large, the discrete NDP above can be approximated by a continuous function defined below.

**Definition**: The continuous NDP, for the system under test is given by the function $p^j(x), 0 \le x \le 1$

$$p^j(x)dx = \frac{nr\_enumerables^j(x, x+dx)}{all\_enumerables^j} \qquad\qquad (2)$$

where $nr\_enumerables^j(x, x+dx)$ stands for all enumerables of type $j$ with detectability values between $x$ and $x + dx$.

Notice that $\int_0^1 p^j(x)dx = 1$, just like the discrete NDP case.


# 4    A one-parameter Model

The detectability profile gives the probability of each enumerable getting exercised. Hence it can be used to calculate expected coverage when a given number of tests have been applied. In this section, we will assume that testing is random, i.e. any single test is selected randomly with replacement. Malaiya and Yang [mal84], and Wagner et al [wag87] have shown that

the expected coverage of the enumerables of type $j$ is given by

$$C^j(n) = 1 - \sum_{i=1}^{n}(1 - d_i^j)^n p_i^j \tag{3}$$

provided testing is random. The same result can be obtained for continuous NDP [set90]

$$C^j(n) = 1 - \int_0^1 (1 - x)^n p(x)dx \tag{4}$$

In practice, testing is more likely to be pseudo-random, when a test will not be repeated. In this cases random testing can be considered to be an approximation. This approximation can be fairly good, except when close to 100% coverage has been achieved.

Equations 3 and 4 give expected coverage. In a specific case, the coverage obtained can be different. If the number of vectors applied is large, then the central limit theorem suggests that results obtained should be close to these given by Equations 3 and 4.

The use of Equations 3 and 4 requires the knowledge of detectability profiles. Obtaining exact detectability profiles requires a lot of computation. Discrete detectability profiles have been calculated for several small and large combinational circuits. Continuous detectability profiles for some benchmark circuits have been estimated [set90]. However software systems are generally much more complex.
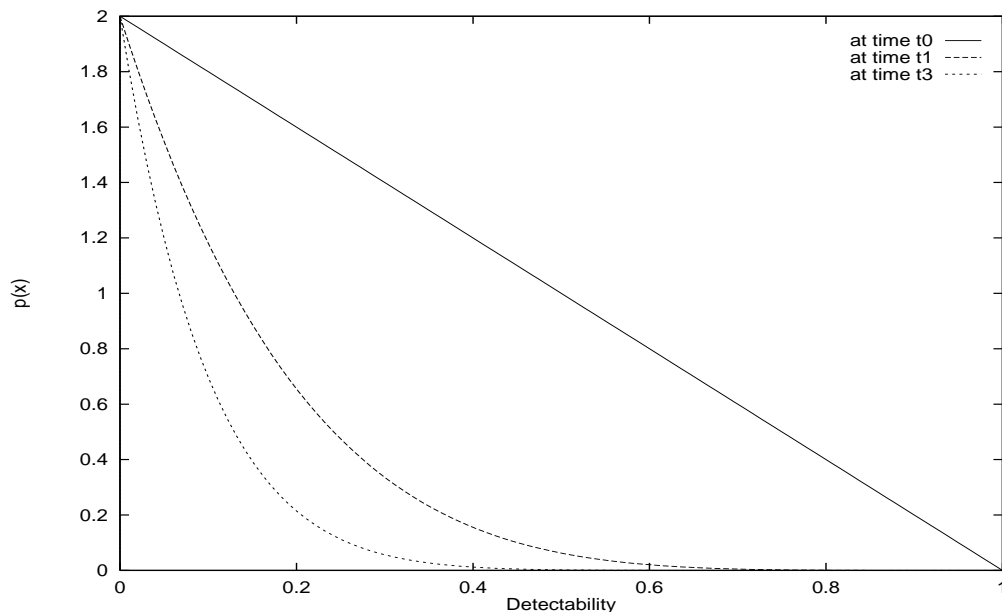


Figure 2: Distribution of detectability

Fortunately, it is possible to obtain reasonable approximation for the detectability profiles. When one test is applied, the probability that an enumerable with detectability $d_i^j$ will not be covered is $(1 - d_i^j)$. The probability that an enumerable will not be covered by $n$ tests and thus remain a part of profile is $(1 - d_i^j)^n$. Thus if the initial profile was given by Equation 1, the profile after having applied $n$ tests, will be given by

7

$$P_n^j = \{p_{d1}^j(1 - d1^j)^n, p_{d2}^j(1 - d2^j)^n, ...\}$$

Equivalently the continuous profile is given by

$$p_n^j(x) = p_n(x)(1 - x)^n$$

The above equations state that the enumerables with high testability are likely to get covered earlier. This would cause the profile to "erode" as testing progresses, as shown in Figure 1. The enumerables with low testability will get removed at a much lower rate, and thus will soon dominate. Thus during much of the testing, the shape of the profile will appear like the bottom curve in Figure 1, regardless of the initial profile.

The available results for hardware components suggest that the initial detectability profile may be of the form

$$p^j(x) = (mj + 1)(1 - x)^{mj} \tag{5}$$

where $mj$ is a parameter. The factors $(mj + 1)$ ensures that the area under the initial profile curve is unity. The significance of the parameter $mj$ can be seen by substituting the right hand side of Equation 5 in Equation 4. We get

$$C^j(n) = 1 - (mj + 1)\int_0^1 (1 - x)^{mj+n}dx = 1 - \frac{mj + 1}{mj + n + 1} = \frac{n}{mj + n + 1} \tag{6}$$

The curve given by Equation 6 has the general shape found with experimental data. However it does not provide a good fit. One problem is that Equation 6 includes only a single parameter which can be adjusted for fitting. We can assume a more general initial detectability in Equation 5, involving two parameters, but even that will not be accurate, as we discuss in the next section. The approach considered next, yields a much better model.

# 5   A New Logarithmic Coverage Model

Random testing is rarely done in practice. Randomness implies that a new test is selected regardless of the tests that have been applied thus far, and that tests are selected based only on the operational distribution. In actual practice, a test case is selected in order to exercise a functionality or enumerable that has remained untested so far. This process makes actual testing more directed and hence more efficient than random testing.

Malaiya, von Mayrhauser and Srimani [mvs92] show that this non-random process leads to a defect finding behavior described by the logarithmic growth model [mus87]. Their analysis gives an interpretation for the model parameters. The coverage growth of an enumerable-type depends on the detectability profile of the type and the test selection strategy. If the defect coverage growth in practice is described by the logarithmic model, it is likely that the

coverage growth for other enumerable-types may also be logarithmic. We thus suggest the following model.

$$C^i(t) = \frac{1}{N^i}\beta_0^i \ln(1 + \beta_1^i t), \qquad\qquad C^i(t) \leq 1 \qquad\qquad (7)$$

where $N^i$ is the total number of enumerables of type $i$, $\beta_0^i$ and $\beta_1^i$ are model parameters. If a single application of a test, assuming it is a constant, takes $T_s$ seconds, then the time $t$, needed to apply $n$ tests is $nT_s$. Substituting in 7,

$$C^i(n) = \frac{\beta_0^i}{N^i} \ln(1 + \beta_1^i T_s n)$$

Defining $b_0^i$ as $(\frac{\beta_0^i}{N^i})$ and $b_1^i$ as $(\beta_1^i T_s)$, we can rewrite the above as,

$$C^i(n) = b_0^i \ln(1 + b_1^i n), \qquad\qquad C^i(n) \leq 1 \qquad\qquad (8)$$

Notice that when $C^i = 1$, there are no more additional enumerables of that type to be found. With non-random testing assumption, it takes a finite, although possibly large, number of tests to achieve 100% coverage of the feasible enumerables.

For defects ($i = 0$), the parameters $\beta_0^0$ and $\beta_1^1$ have the following interpretation [mvs93].

$$\beta_0^0 = \frac{K^0(0)N^0(0)}{a^0 T_L} \qquad\qquad (9)$$

and

$$\beta_1^0 = a^0 \qquad\qquad (10)$$

where $K^0$ is the exposure ratio, $T_L$ is the linear execution time and $a^0$ is a parameter that describes the variation in the exposure ratio.

Equation 8 relates coverage $C^i$ with the number of tests applied. We can use it to obtain an expression giving defect coverage $C^0$ in terms of one of the coverage metrics $C^i$, $i = 1$ to 4. Using Equation 8, we can solve for $n$,

$$n = \frac{1}{b_1^i}[exp(\frac{C^i}{b_0^i}) - 1], \qquad\qquad i = 1 \quad to \quad 4$$

Substituting this for $C^0$, again using Equation 8,

$$C^0 = b_0^0 \ln[1 + \frac{b_1^0}{b_1^i}(exp(\frac{C^i}{b_0^i}) - 1)], \qquad\qquad i = 1 \quad to \quad 4$$

Defining $a_0^i = b_0^0$, $a_1^i = \frac{b_1^0}{b_1^i}$ and $a_2^i = \frac{1}{b_0^i}$, we can write the above using three parameters as,

$$C^0 = a_0^i \ln[1 + a_1^i(exp(a_2^i C^i) - 1)] \qquad\qquad i = 1 \quad to \quad 4 \qquad\qquad (11)$$

Equation 11 gives us a convenient three-parameter model for defect coverage in terms of a measurable test coverage metric. Notice that Equation 11 is applicable for only $C^0 \leq 1$. It is possible to approximate Equation 11 by a linear relation, but it would be valid for only a small range.

# 6    Analysis of Data

Here we will evaluate the proposed model, as given by Equations 8 and 11 using four data sets. The first data set, DS1, is from a multiple-version automatic airplane landing system [lyu93]. The twelve versions have a total of 30,694 lines. The data used here is for integration and acceptance test phases, where 66 defects were found. One additional defect was found during operational testing. The second data set, DS2, is from a NASA supported project implementing sensor management in inertial navigation system [vou92]. For this program, 1196 test cases were applied and 9 defects were detected. The third data set, DS3, is for a simple program used to illustrate test coverage measures [agr93]. The fourth data set, DS4, is from an evolving software system containing a large number of modules.

Table 2: Summary table for DS1 (total 21,000 tests applied)

|             | Blocks i=1 | Decisions i=2 | c-uses i=3 | p-uses i=4 | Defects i=0 |
|-------------|------------|---------------|------------|------------|-------------|
| Total enums | 6977       | 3524          | 8851       | 4910       | 67          |
| Final cov.  | 91.8%      | 83.9%         | 91.7%      | 73.5%      | 98.4%       |
| $b_0^i$     | 0.031      | 0.049         | 0.036      | 0.041      | 0.184       |
| $b_1^i$     | 2E+8       | 1234          | 3.4E+6     | 2439       | 0.01        |
| LSE         | 5.7E-4     | 3.5E-5        | 5.8E-4     | 8.1E-5     | 7.3E-7      |

The first data set and the results from it are summarized in Table 2. The first row gives the total number of enumerables for all versions. The second row gives the average coverage when 21,000 tests had been applied. The values of the estimate parameters $b_0^i$ and $b_1^i$ and the least square error are given in the rows below. The model given by Equation 8 fits the data very well. The data shows that $C^1 > C^2 > C^4$. This relationship is expected. Complete decision coverage implies complete block coverage, and complete p-uses coverage implies complete decision coverage [bei90, fra88, n88]. The c-uses coverage has no such relation relative to the other metrics. Indeed the data shows that while $C^3 < C^1$ at the beginning of testing, near the end of testing $C^3$ is almost equal to $C^1$.

Table 3 summarizes the result for DS2. Nine faults were revealed by application of 1196 tests; we assume that one fault (i.e. 10%) is still undetected. In spite of the small number of faults, the models given in Equations 8 and 11 still fit the data very well.

10

Table 3: Summary table for DS2 (total 1196 tests applied)

| | Blocks i=1 | Branches i=2 | c-uses i=3 | p-uses i=4 | Defects i=0 |
|---|---|---|---|---|---|
| Final cov. | 89% | 84% | 76% | 61% | 90% |
| $b_0^i$ | 0.032 | 0.060 | 0.034 | 0.039 | 0.166 |
| $b_1^i$ | 2E+8 | 870 | 3E+7 | 2500 | 0.11 |
| LSE | 0.02 | 6.2E-4 | 3.5E-3 | 4.9E-3 | 0.025 |
| $a_0^i$ | 1.31 | 0.46 | 0.23 | 0.29 | |
| $a_1^i$ | 1.8E-3 | 4.6E-3 | 9.11E-7 | 5.2E-3 | |
| $a_2^i$ | 6.95 | 3.84 | 23.12 | 13.46 | |
| LSE | 0.017 | 0.018 | 0.041 | 0.025 | |

One important observation can be made from Table 3. We can use the number of tests applied and parameters $b_0^0$ and $b_1^0$ (equation 8) to estimate the defect coverage. Alternatively we can use branch coverage and the parameters $a_0^2$, $a_1^2$ and $a_2^2$ (equation 11) to estimate the defect density. The second approach provides a slightly better fit and it is unaffected by the test selection strategy.
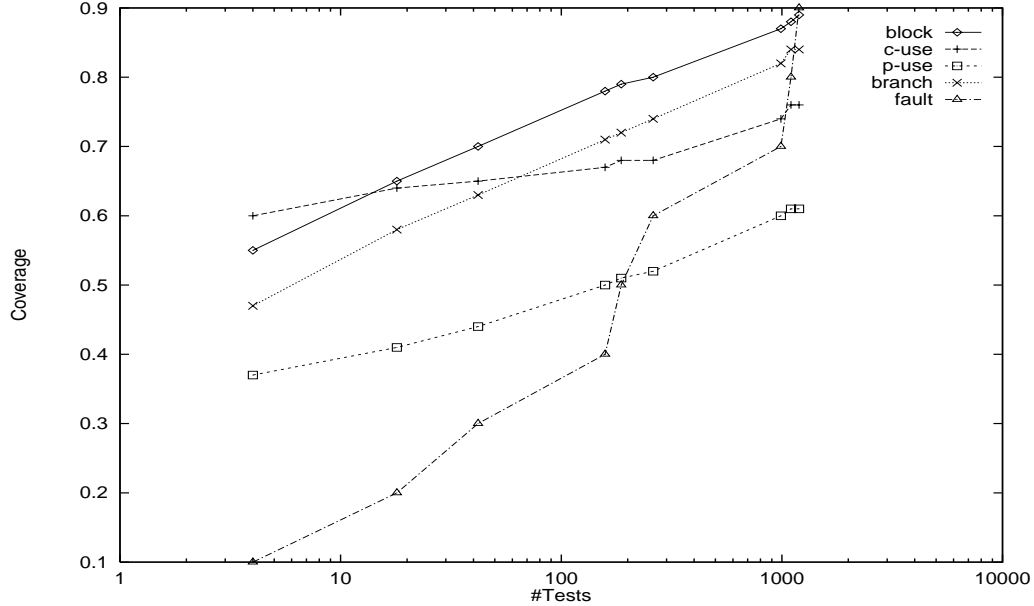


Figure 3: Growth of coverage measures with testing

As shown in Figure 3, we can see that the curves for $C^1$, $C^2$, and $C^4$ have the same shapes, and we again observe that $C^1 > C^2 > C^4$. Again the c-uses coverage does not have the same relation with other three. With the logarithmic x-axis (number of tests) $C^1$, $C^2$, $C^3$, $C^4$ appear as straight lines. This can be explained by examining equation 8. If $b_1^i n \gg 1$, is true for $i = 1, 2, 3, 4$, then we can write,
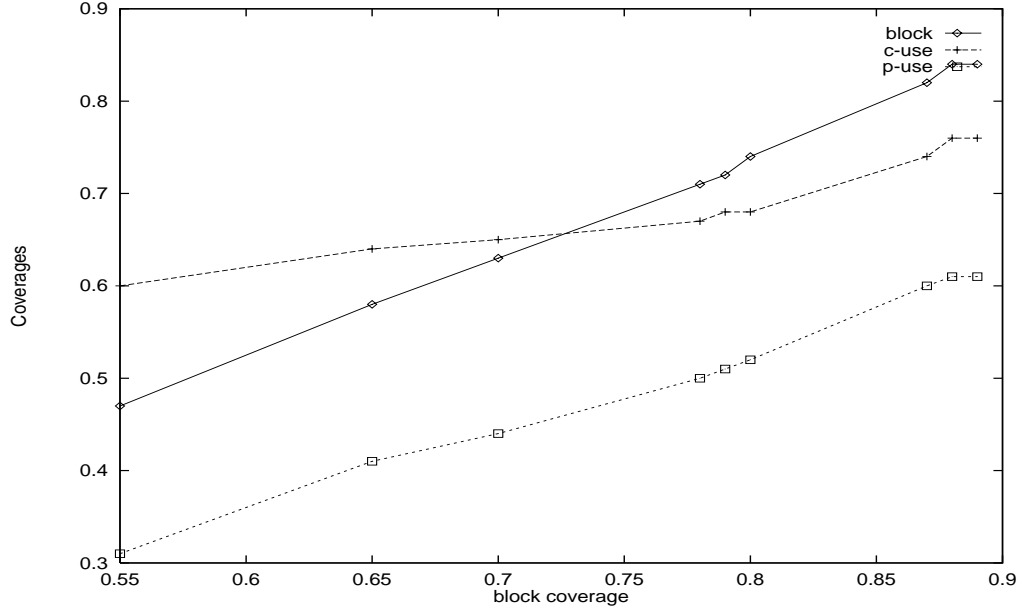
$$C^i(n) = b_0^i \ln(b_1^i n)$$

Figure 4: Scatter plot of C2, C3 and C4 against C1

or

$$C^i(n) = b_0^i \ln(b_1) + b_0^i \ln(n)$$

This gives us the linear curves of Figure 3. Notice that since $b_1^0 < 1$, the approximation above is not applicable for $C^0(n)$ as observed in Figure 3. Figure 4 shows the correlation of other test coverage measures $C^2$, $C^3$ and $C^4$ with block coverage $C^1$. As we would expect, branch coverage, and to a lesser extent p-use coverage, are both strongly correlated with block coverage. The correlation with c-use coverage is weaker. Figure 5 shows actual and computed values for fault coverage. The computed values have been obtained using branch coverage and Equation 11. Notice that at 50% branch coverage, the fault coverage is still quite low (about 10%), however with only 84% branch coverage, 90% fault coverage is obtained. The branch coverage shows saturation at about 84%. This may provide an explanation for why it is often considered quite adequate to achieve 80% branch coverage [gra92].

Figure 6 gives a scatter plot of computed values of defect coverage against actual values. The computed values are obtained using the number of tests and Equation 8 (traditional reliability growth modeling), and using test coverage measures $C^1$, $C^2$, $C^3$ and $C^4$ using Equation 11. The calculated values are all quite close, showing that coverage based modeling can replace time-based modeling.

Table 4 presents similar results for a very small illustrative program. No defects were involved. However this again illustrates applicability of our modeling scheme. We again notice that $C^1 \geq C^2 \geq C^4$. The c-use coverage again behaves differently.

In evolving programs, significant changes are being made while testing is in progress. Because new modules are being added new defects as well as non-covered enumerables are
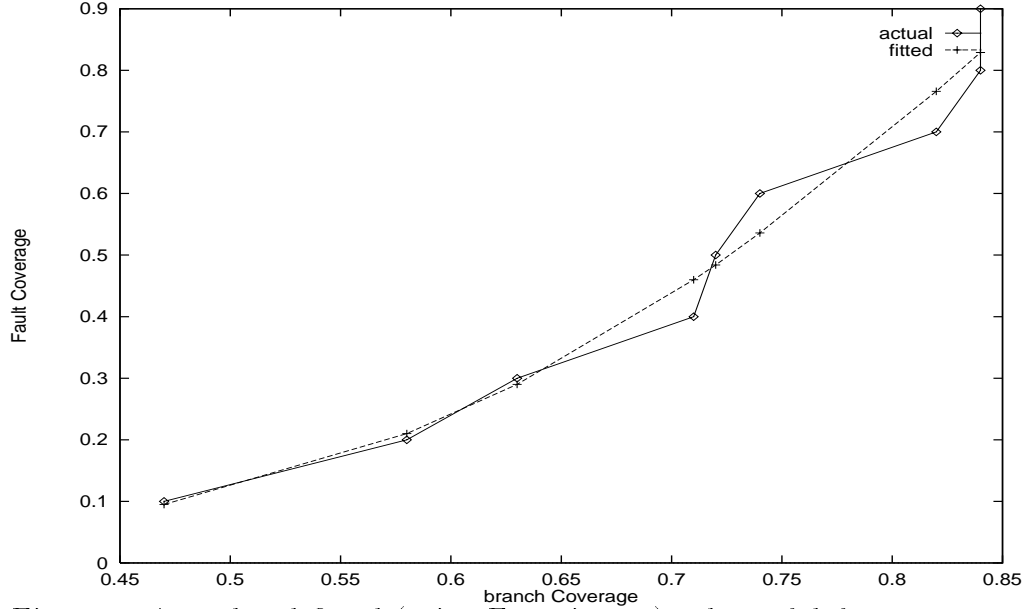
12

Figure 5: Actual and fitted (using Equation 11) values of defect coverage

Table 4: Summary table for DS3 (total 16 tests applied)

|  | Blocks i=1 | Branches i=2 | c-uses i=3 | p-uses i=4 |
|---|---|---|---|---|
| Total enums | 12 | 10 | 10 | 26 |
| Final cov. | 100% | 100% | 100% | 93% |
| $b_0^i$ | 0.06 | 0.11 | 0.11 | 0.12 |
| $b_1^i$ | 8.4E5 | 769 | 561 | 162 |
| LSE | 2E-5 | 1.7E-3 | 1.6E-3 | 2.3E-3 |

also being added. The coverage obtained by a test set can actually go down in some cased. From DS4, as shown in Figure 7, we see that the linear correlation between coverage measures can still be applicable. The part of the data used here covers an intermediate phase of the process. The analysis of evolving programs is however more complex and is the subject of future research.

# 7   Model Parameters

Researchers have noticed that the logarithmic model works best among other two-parameter models [mkv92], however interpretation of its parameters has been difficult. One interpretation is given by Malaiya et al [mvs92], described by Equations 9 and 10. We can argue that the same interpretation may be applicable for enumerables other than defects. The first parameter of Equation 8 then is,
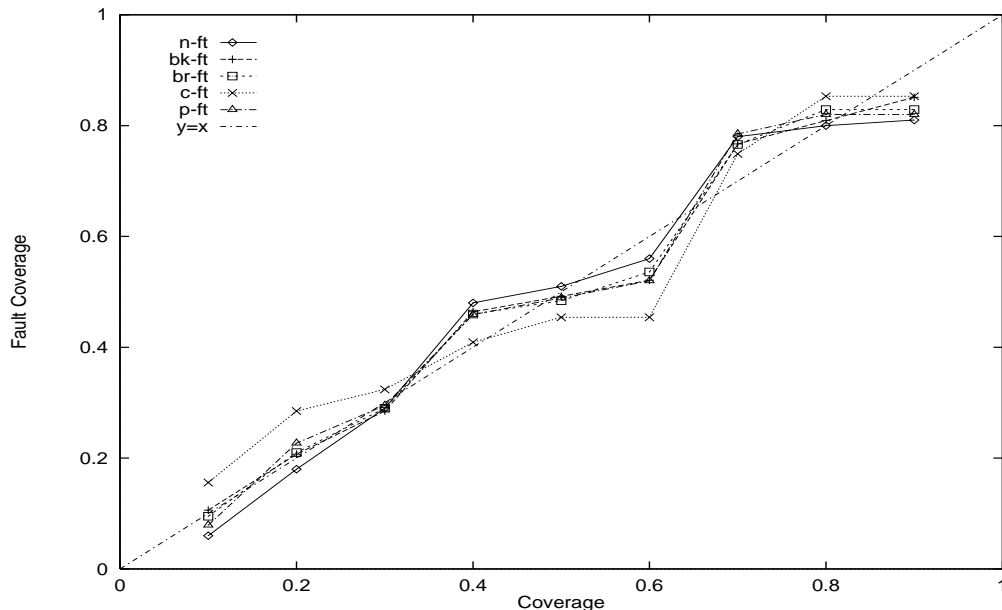
13

Figure 6: Actual defect coverage vs computed values

$$b_0^i = \frac{K^i(0)N^i}{a^i T_L N^i} = \frac{K^i(0)}{a^i T_L} \tag{12}$$

Notice that the linear execution time is given by the number of lines of code multiplied by the average execution time of each line. An empirical method to estimate the initial fault exposure ratio $K^0(0)$ has been suggested by Li and Malaiya [lim93]. Estimation of $a^i$ remains an open problem. The second parameter is given by,

$$b^i = a^0 T_s \tag{13}$$

The single test execution time $T_s$ depends on the program size and its structure. The product $b_0^i b_1^i$ then should be independent of the program size.

The parameters $a_0^i$, $a_1^i$, and $a_2^i$ are defined in terms of $b_0^i$ and $b_1^i$ above. When this definition for $a_0^i$, $a_1^i$, and $a_2^i$ is as an initial estimate for numerically fitting Equation 11, the initial estimate itself provides a least-square fit. If the initial estimates are significantly different, then the least square fit may yield somewhat different parameter values.

A-priori estimation of model parameters remains a partly unsolved problem. Currently we must rely on curve fitting based approaches.

# 8 Defect density and reliability

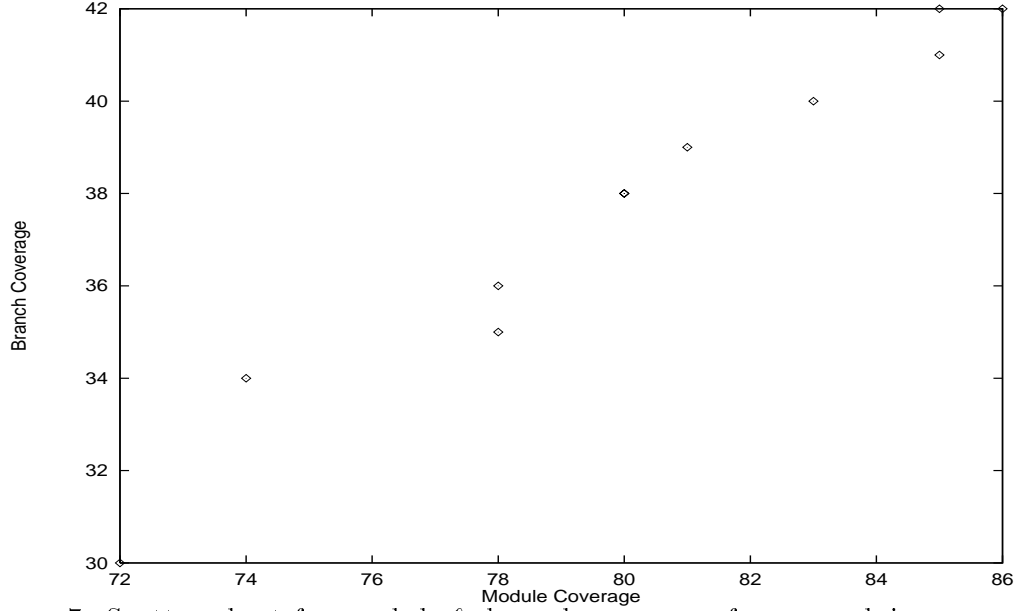Since the failure intensity is proportional to the number of defects, we have [mvs92, mvs93],

Figure 7: Scatter chart for module & branch coverages for an evolving program

$$\lambda = \frac{K}{T_L} N$$

Where $K$ is the overall value of fault exposure ratio.

Let $N_0$ be the total number of faults initially present in the program and there is no new fault introduced during testing process. Then $N$ can be computed as:

$$N = N_0(1 - C^0)$$

Substituting $C^0$ using Equation 11,

$$N = N_0(1 - a_0^i \ln[1 + a_1^i(exp(a_2^i C^i) - 1)])$$

Hence, the expected duration between successive failures can be obtained as

$$\frac{1}{\lambda} = \frac{T_L}{K} \frac{1}{N_0(1 - a_0^i \ln[1 + a_1^i(exp(a_2^i C^i) - 1)])} \tag{14}$$

The Equation 14 can be used for operational period also, provided the appropriate value for fault exposure ratio is used. In can be estimated by testing using the operational input profile. Alternatively, appropriate compression factor should be used.

# 9 Conclusions and Discussions

In this paper we have presented a modeling scheme that relates defect density to measurable coverage metrics. The scheme is based on the observation that defects have some detectability distribution like other coverage enumerables, and the same model may govern them. There are two advantages of using logarithmic model to describe test effort and enumerables covered.

1. The logarithmic model has been found to be superior to other models for predicting the number of defects.

2. The logarithmic model can take into account the fact that 100% coverage can be achieved in finite time. For high reliability applications, it is quite possible that 100% block coverage is achieved, but the reliability requirements will require further testing. A more strict coverage measure, branch or p-use coverage then can be used to estimate the defect density.

The data sets used suggest that the model works well. The results are consistent with the theoretical results obtained in the past about coverage inclusion relationships. There is a need to collect additional data sets with larger number of points. The model presented here is simple and easily explained, and is thus suitable for industrial use.

The model, as given by Equation 11, can be used in two different ways. *Extrapolation* requires collecting data for part of the testing process, which is then used for estimating the applicable parameter values. These are used for making projections for planning the rest of the test effort. *A priori* parameter estimation requires empirical estimation of parameters even before testing begins. We have some observations on what factors control the parameter values. Further work is needed to fully develop these techniques. This would include a careful study of enumerable exposure ratios.

As we have shown, any test coverage measure can be used to estimate the defect density, by using applicable parameter values. This raise an important question. Should several coverage measures be used or just one ? Which individual measure (or selected set) would provide the best estimate ? We have the following considerations.

1. We need further studies to determine which coverage measure would provide the best project about the number of defects. For DS2, we have observed that block coverage $C^1$ provides the best and c-uses coverage $C^3$ the worst fit. This however may be true for only specific data sets, or for specific coverage/defect density ranges. Since the different coverage measures can be strongly correlated, as we have seen, perhaps most of them may work equally well in many situations.

2. For every high reliability, we may need a "scale" that works in that region. If the requirements are such that 100% block coverage is not enough, branch or p-use coverage may be more appropriate. From prevailing practices in industry today, it appears that branch coverage may be an adequate measure in many cases, since about 80% branch

coverage often produces acceptable results [gra92]. However for testing of individual modules or for highly reliable software p-use may be a better measure.

3. Several researchers have suggested use of some form of a weighted risk measure [agr93, pos93, neu93], where a weighted average is computed using some coverage measures. The weights are chosen on the basis of relative significance of each measure. As we have seen the structural coverage measures tend to be strongly correlated, and thus a weighted average may not provide more information than a single measure. We need to identify measures with weaker correlation. It is possible that other types of coverage measures like functional coverage measures may be suitable for this purpose. Further study is needed to evaluate correlation among different types of coverage measures, and how they can be optimally combined.

The results presented here can serve as a basis for further data collection and analysis. We need to examine the behavior at different fault densities, specially at very low defect densities (for highly reliable applications). We also need to validate the model for different testing strategies on the modeling scheme and the parameter values. In general deterministic (coverage driven) is more efficient than true random testing. Testing using special values or use of equivalence partitioning can significantly compress the test time. Since test coverage measures provide direct sampling of the state of the software, we expect the model Equation 11 to hold because time is eliminated as a variable. Additional data will allow us to validate and refine the modelling scheme presented here. In addition we need to develop schemes for evolving programs where new faults and other non-covered enumerables are being added.

# 10    Acknowledgement

# References

[agr93]    H. Agrawal, J.R. Horgan, E.W. Krauser, S.A. London, "A testing-Based model and Risk Browser for C" *Proc. Int. Conf.on Rel., Qual.Control and Risk Asses.*, Oct. 1993, pp 1-7.

[bei90]    B. Beizer, *Software Testing Techniques*, Van Nostrand Reinhold, 1990, pp. 74-75, 161-171.

[bisc89]    J.M. Bieman and J.L. Schultz, "Estimating the Number of Test Cases Required to Satisfy the All-du-paths Testing Criterion," *Proc. ACM SIGSOFT, in Software Engineering Notes*, Dec. 1989, pp.179-186.

[bisc92]    J.M. Bieman and J.L. Schultz, "An Empirical Evaluation (and Specification) of the All-du-paths Testing Criterion," *Software Engineering Journal*, Jan. 1992, pp. 43-51.

[chm92]    M.H. Chen, J.R. Horgan, A.P. Mathur and V.J. Rego, "A time/structure based model for estimating software reliability," *SERC-TR-117-P*, Purdue University, Dec. 1992.

[dhk93]    S.R. Dalal, J.R. Horgan and J.R. Kettenring, "Reliable Software and Communications: Software Quality, Reliability and Safety," *Proc. 15th Int. Conf. Software Engineering*, May 1993, pp. 425-435

[dun86]    J.R. Dunham, "Experiments in software reliability: Life Critical Applications," *IEEE Trans. Soft. Eng.*, Jan. 1986, pp. 110-123.

[fra88]    P.G. Frankl and E.J. Wayuker, "An Applicable Family of Data Flow Testing Criteria," *IEEE Trans. Soft. Eng.*, Oct. 1988, pp. 1483-1498.

[fra93]    P.G.Frankl and N.Weiss, "An Experimental Comparison of the Effectiveness of Branch Testing and Data Flow Testing," *IEEE Trans. Soft. Eng.*, Aug. 1993, pp. 774-787.

[gra92]    R.E. Grady, *Practical Software Metrics for Project Management and Process improvement*, PTR prentice-Hall, 1992, pp. 58-60.

[hec94]    H. Hecht and P. Crane, "Rare Conditions and Their Effect on Software Failures," Proceedings of *Ann. reliability and maintainability Symp.*, pp. 334-337, Jan. 1994.

[lim93]    N. Li and Y. K. Malaiya, "Fault Exposure Ratio and Reliability Estimation," *Proc. 3rd Workshop on Issues in Software Reliability*, November 1993, pp. 6.3.1-6.3.18.

[limi93]   N. Li and Y.K. Malaiya, "Enhancing Acuracy of Software reliability Prediction" IEEE Int. Symp. on Software Reliability Engineering, 1993.

[lyu93]    M.R. Lyu, J.R. Horgan and S. London, "A coverage Analysis Tool for the Effectiveness of Software Testing" IEEE Int. Symp. on Software Reliability Engineering, 1993, pp. 25-34.

[mal84]    Y.K. Malaiya and S. Yang, "The Coverage Problem for Random Testing," Proceedings of *Int. Test Conference*, pp. 237-242, October 1984.

[mkv92]    Y. K. Malaiya, N. Karunanithi and P. Verma, "Predictability of Software Reliability Models," *IEEE Trans. Reliability*, Dec. 1992, pp. 539-546.

[mus87]    J.D. Musa, A Iannino, K. Okumoto, *Software Reliability, Measurement, Prediction, Application*, McGraw-Hill, 1987.

[mvs92]    Y. K. Malaiya, A. von Mayrhauser and P. Srimani, "The Nature of Fault Exposure Ratio," *Proc. IEEE Int. Symp. Soft. Rel. Eng.*, Oct. 1992, pp. 23-32.

[mvs93]    Y. K. Malaiya, A. von Mayrhauser and P. Srimani, "An examination of Fault Exposure Ratio," to appear in *IEEE Trans. Software Engineering*, 1993.

[n88]      S.C. Ntafos, "A comparision of some structural testing strategies" *IEEE Trans. Software Engineering*, June 1988, pp.868-874.

[neu93]    A.M. Neufelder, *Ensuring Software Reliability*, Marcel Dekker Inc., 1993, pp. 137-140.

[poc93]    P. Piwowarski, M. Ohba and J. Caruso, "Coverage measurement experience during function test," *Proc. 15th Int. Conf. Software Engineering*, May 1993, pp. 287-300

[pos93]    R.M. Poston, "The Power of Simple Software Testing Metrics", Software Testing Times, Vol. 3, No. 1993.

[ram85]    J. Ramsey and V.R.Basili, "Analyzing the Test Process Using Structural Coverage", *Proc. 8th Int. Conf. on Software Engineering*, August 1985, pp. 306-312.

[set90]    S.C. Seth, V.D. Agrawal and H. Farhat, "A Statistical Theory of Digital Circuit Testability," *IEEE Trans. Comp.*, April, 1990, pp. 582-586.

[tra92]    M. Trachtenberg, "Why Failure Rates observe Zipf's Law in Operational Software," *IEEE Trans. Reliability*, Sept. 1992, pp. 386-389.

[voas92]   J. Voas and K. Miller, "Improving the Software Development Process Using Testability Research," Proc. Int. Symp. on Software Reliability Engineering, 1992, pp. 114-121.

[vou92]    M.A. Vouk "Using Reliability Models During Testing With Non-operational Profiles," *Proc. 2nd Bellcore/Purdue workshop on issues in Software Reliability Estimation,* Oct. 1992, pp. 103-111

[wag87]    K. Wagnor, C. Chin and E. McCluskey, "Pseudorandom Testing," *IEEE Trans. Comput.*, Vol. C-36, pp. 332-343, March 1987.

[weu84]    E.J. Weyuker, "More Experience with Data Flow Testing", *IEEE Trans. Software Engineering*, September 1993, pp. 912-919.

[weu93]    E.J. Weyuker, "An Empirical Study of the Complexity of Data Flow Testing," *2nd Workshop on Software Testing, Verification, and Analysis*, July 1988.