

An Experimental Determination of Sufficient Mutant Operators

A. JEFFERSON OFFUTT and AMMEI LEE

George Mason University

and

GREGG ROTHERMEL

Clemson University

and

ROLAND H. UNTCH

Middle Tennessee State University

and

CHRISTIAN ZAPF

Siemens Corporation

Mutation testing is a technique for unit-testing software that, although powerful, is computationally expensive. The principal expense of mutation is that many variants of the test program, called mutants, must be repeatedly executed. This article quantifies the expense of mutation in terms of the number of mutants that are created, then proposes and evaluates a technique that reduces the number of mutants by an order of magnitude. Selective mutation reduces the cost of mutation testing by reducing the number of mutants. This article reports experimental results that compare selective mutation testing with standard, or nonselective, mutation testing, and results that quantify the savings achieved by selective mutation testing. The results support the hypothesis that selective mutation is almost as strong as nonselective mutation; in experimental trials selective mutation provides almost the same coverage as nonselective mutation, with a four-fold or more reduction in the number of mutants.

Categories and Subject Descriptors: D.2.5 [Software Engineering]: Testing and Debugging—*test data generators*

General Terms: Experimentation, Measurement, Reliability

The first author was partially supported by the National Science Foundation under grant CCR-93-11967.

Authors' addresses: A. J. Offutt and A. Lee, Department of Information and Software Systems Engineering, 4A4, George Mason University, Fairfax, VA 22030; email: {ofut, alee}@issee.gmu.edu; G. Rothermel, Department of Computer Science, Clemson University, Clemson, SC 29634-1906; email: grother@cs.clemson.edu; R. H. Untch, Department of Computer Science, Box 48, Middle Tennessee State University, Murfreesboro, TN 37132; email: untch@knuth.mtsu.edu; C. Zapf, Siemens AG—Medical Group, BNES 13, Henkestrasse 127, 91052 Erlangen, Germany; email: zapf@erlh.siemens.de.

Permission to make digital/hard copy of part or all of this work for personal or classroom use is granted without fee provided that the copies are not made or distributed for profit or commercial advantage, the copyright notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or a fee.

© 1996 ACM 1049-331X/96/0400-0099 \$03.50

Additional Key Words and Phrases: Fault-based testing, mutation, software testing, unit testing

1. INTRODUCTION

Mutation testing is a technique, originally proposed by DeMillo et al. [1978] and Hamlet [1977], that requires a person testing a program to create test data that causes a finite, well-specified set of faults to result in failure. The tester does this by finding test cases that cause faulty versions of the program to fail. These test cases will then either result in correct output from the test program (demonstrating its quality) or cause the test program to fail (detecting a fault). The technique thus serves two goals: it provides a test adequacy criterion, and it leads to detection of faults in the program being tested.

Unit-level-testing techniques such as mutation hold great promise for improving the quality of software. Unfortunately, most of these techniques are currently so expensive that we cannot afford to use them. In a previous paper [Offutt et al. 1993], we gave a preliminary definition and evaluation of a mutation approximation method called *selective mutation*. There, we reported results that indicate savings of over 50% of the execution cost of mutation, with negligible loss of effectiveness by one measure. Since then, we have extended selective mutation and have gotten results that, on average, show savings of over 75% of the execution cost on our experimental programs. More importantly, from a theoretical point of view, our analysis in Section 4 shows that selective mutation reduces the number of mutants generated by mutation systems from quadratic in the number of variable references to linear.

In the following sections, we describe mutation testing, analyze its cost, and define the method of selective mutation. We then present experimental results to validate the effectiveness of selective mutation, present data that quantify the savings of selective mutation, and discuss possible future directions.

1.1 Mutation-Testing Overview

Mutation testing helps users create test data. During mutation testing, faults are introduced into a program by creating many versions of the program, each of which contains one fault. Test data are used to execute these faulty programs with the goal of causing each faulty program to fail. Hence we use the term mutation; faulty programs are *mutants* of the original, and a mutant is *killed* when a test case causes it to fail. When this happens, the mutant is considered *dead* and no longer needs to remain in the testing process, since the faults represented by that mutant have been detected and, more importantly, since the mutant has satisfied its requirement of identifying a useful test case.

```

      FUNCTION Min (I,J)
1      Min = I
      Δ Min = J
2      IF (J .LT. I) Min = J
      Δ IF (J .GT. I) Min = J
      Δ IF (J .LT. Min) Min = J
3      RETURN

```

Fig. 1. Function Min.

Figure 1 contains a small Fortran function with three mutated lines (preceded by the Δ symbol). Note that each mutated statement represents a separate program. The most recent mutation system, Mothra [DeMillo et al. 1988; King and Offutt 1991], uses 22 mutant operators to test Fortran-77 programs. The *coupling effect*, which says that complex faults are coupled to simple faults in such a way that a test data set that detects all simple faults in a program will detect most complex faults [DeMillo et al. 1978], allows us to limit mutation operators to those that make simple changes. The coupling effect has been supported experimentally [Offutt 1992] in a study that compared test sets generated for mutants that involved changes in two places with test sets generated for single-change mutants. Morell [1984] also analyzed the coupling effect from a theoretical perspective. The strongest support is by Wah [1995], who proved probabilistically that the coupling effect holds under certain restrictions (for finite bijective functions); his dissertation removes the restrictions.

The mutation-testing process begins with the construction of mutants of a test program. The user then adds test cases (generated manually or automatically) to the mutation system and checks the output of the program on each test case to see if it is correct. If the output is incorrect a fault has been found, and the program must be modified and the process restarted. If the output is correct, the test case is executed against each live mutant. If the output of a mutant differs from that of the original program on the same test case, the mutant's output is assumed to be incorrect, and it is killed.

After each new test case has been executed against each live mutant, each remaining mutant falls into one of two categories. One, the mutant is functionally *equivalent* to the original program. An equivalent mutant always produces the same output as the original program, so no test case can kill it. Two, the mutant is killable, but the set of test cases is insufficient to kill it. In this case, new test cases need to be created, and the process iterates until the test set is strong enough to satisfy the tester.

The *mutation score* for a set of test data is the *percentage of nonequivalent mutants killed by that data*. A test data set is said to be *mutation adequate* if its mutation score is 100%.

1.2 The Cost of Mutation Testing

The major computational cost of mutation testing is incurred when running the mutant programs against the test cases. Budd [1980] analyzed the

Table I. Mothra Mutant Operators

Mutation Operator	Description
AAR	array reference for array reference replacement
ABS	absolute value insertion
ACR	array reference for constant replacement
AOR	arithmetic operator replacement
ASR	array reference for scalar variable replacement
CAR	constant for array reference replacement
CNR	comparable array name replacement
CRP	constant replacement
CSR	constant for scalar variable replacement
DER	DO statement end replacement
DSA	DATA statement alterations
GLR	GOTO label replacement
LCR	logical connector replacement
ROR	relational operator replacement
RSR	RETURN statement replacement
SAN	statement analysis
SAR	scalar variable for array reference replacement
SCR	scalar for constant replacement
SDL	statement deletion
SRC	source constant replacement
SVR	scalar variable replacement
UOI	unary operator insertion

number of mutants generated for a program and found it to be roughly proportional to the product of the number of data references times the number of data objects. Typically, this is a large number for even small program units. For example, 44 mutants are generated for the function *Min* shown in Figure 1. Since each mutant must be executed against at least one, and potentially many, test cases, mutation testing requires large amounts of computation. We explore this cost in more detail in Section 4.

2. SELECTIVE MUTATION TESTING

One way to reduce the cost of mutation testing is to reduce the number of mutant programs created, using an approximation approach originally suggested by Mathur [1991]. To understand his proposal, we must first discuss the method by which mutant programs are generated.

The syntactic modifications responsible for mutant programs are determined by a set of *mutant operators*. This set is determined by the language of the program being tested and the mutation system used for testing. Mutant operators are created with one of two goals: to induce simple syntax changes based on errors that programmers typically make (such as using the wrong variable name) or to force common testing goals (such as executing each branch). The Mothra mutant operators represent more than 10 years of refinement through several mutation systems. Table I lists the operators; their detailed descriptions are elsewhere [King and Offutt 1991]. Each of the 22 mutant operators is represented by a three-letter acronym.

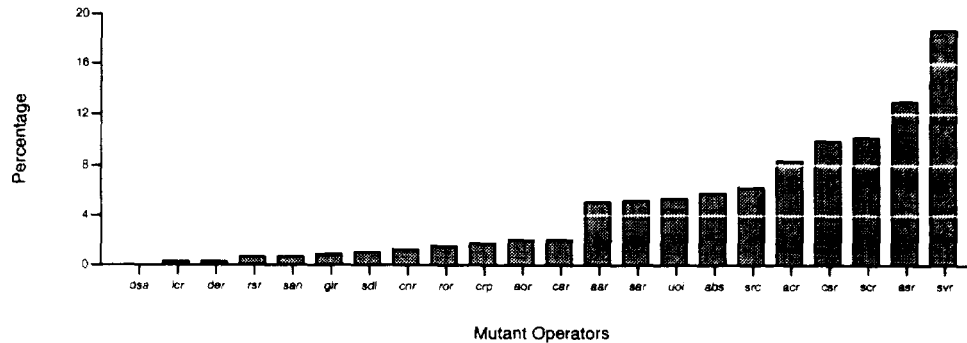


Fig. 2. Distribution of mutants by operator.

For example, the “array reference for array reference replacement” (AAR) mutant operator causes each array reference in a program to be replaced by each other distinct array reference in the program. Budd’s analysis of the number of mutants, discussed in Section 1.2, is largely based on the fact that the SVR mutant operator is the dominant operator.

Since mutant operators generate different numbers of mutant programs, Mathur [1991] proposed *selective mutation*¹ as mutation *without* the operators that create the most mutants. We show the distribution of mutants by operator for a set of 28 programs in Figure 2. Mathur suggested omitting the SVR and ASR operators, which are the two operators that result in the most mutations. In our previous paper [Offutt et al. 1993], we called this two-selective mutation, and we extended the theory to *N-selective mutation*, which omits the *N* most prevalent operators. We presented results based on two-selective, four-selective, and six-selective mutation.

Since the results from six-selective mutation were so positive, we have extended selective mutation to try to experimentally determine whether further reductions in the number of operators can yield effective mutation testing. The mutant operators used by Mothra can be divided into three general categories based on the syntactic elements that they modify. *Replacement-of-operand* operators replace each operand in a program with each other legal operand. Listed in Table I, the operators AAR, ACR, ASR, CAR, CNR, CRP, CSR, SAR, SCR, SRC, and SVR perform operand replacement. *Expression modification* operators (ABS, AOR, LCR, ROR, UOI) modify expressions by replacing operators and inserting new operators. *Statement modification* operators (DER, DSA, GLR, RSR, SAN, SDL) modify entire statements. If (1) the number of constants and variables in a program is *Vals*, (2) the number of value references is *Refs*, and (3) the number of source executable lines is *Lines*, then operand replacement operators result in $O(Vals * Refs)$ mutants; expression modification operators result in $O(Refs)$ mutants; and statement modification operators result in $O(Lines)$ mutants. We define *expression/statement-selective muta-*

¹In his paper, Mathur used the term “constrained mutation.” With his agreement, we have chosen the new term “selective mutation.”

tion (*ES-selective*)² to be mutation with only the expression and statement operators (not using the 11 operand replacement mutant operators), *replacement/statement-selective mutation* (*RS-selective*) to be mutation with only the replacement and statement mutant operators, and *replacement/expression-selective mutation* (*RE-selective*) to be mutation with only the replacement and expression mutant operators. Finally, in Section 3.3, we give results using *expression-selective mutation* (*E-selective*), which is mutation with only the five expression modification mutant operators.

3. EXPERIMENTATION WITH SELECTIVE MUTATION

We evaluated selective mutation empirically. Two relationships that have been defined elsewhere are ProbBetter and ProbSubsumes. Weyuker et al. [1991] suggest a relationship called ProbBetter:

A testing criterion C_1 is ProbBetter than C_2 for a program P if a randomly selected test set T that satisfies C_1 is more “likely” to detect a failure than a randomly selected test set that satisfies C_2 .

Mathur and Wong [Mathur and Wong 1994; Wong 1993] suggest a different relationship called ProbSubsumes:

A testing criterion C_1 ProbSubsumes C_2 for a program P if a test set T that is adequate with respect to C_1 is “likely” to be adequate with respect to C_2 . If C_1 ProbSubsumes C_2 , C_1 is said to be at least as more “difficult” to satisfy than C_2 .

The ProbBetter relation is defined with respect to the fault detection capability of test sets, whereas the ProbSubsumes relation is defined with respect to the difficulty of satisfying one criterion in terms of another. Both are probabilistic relations between two testing criteria and are defined in terms of specific programs. Although this means that it is difficult to draw general conclusions from any one study, as the number and variety of programs studied increases, our confidence in the validity of a ProbSubsumes or a ProbBetter relationship with a larger set of programs also increases.

While it can easily be seen that nonselective mutation subsumes selective mutation, we hypothesize that selective mutation ProbSubsumes nonselective mutation. While it seems unlikely that selective mutation will consistently ProbSubsume nonselective, a fundamental aspect of software testing is that we are usually satisfied with partial results, and people who have experimented with mutation have observed that test sets with mutation scores exceeding 95% are difficult to create and believe such test sets are effective [Budd et al. 1980; Choi and Mathur 1993; DeMillo and Offutt 1993; Frankl et al. 1994]. Thus, we refine our hypothesis to be: test sets that kill all mutants under selective mutation will have a high coverage

²In our previous paper [Offutt et al. 1993], we defined selective in terms of the operators excluded; here we define selective in terms of the operators chosen.

Table II. Experimental Programs

Program	Description	Statements	Mutants	Equivalent
Banker	Deadlock avoidance algorithm	48	2765	43
Bub	Bubble sort on an integer array	11	338	35
Cal	Days between two dates	29	3010	236
Euclid	Greatest common divisor (Euclid's)	11	196	24
Find	Partitions an array	28	1022	75
Insert	Insertion sort on an integer array	14	460	46
Mid	Median of three integers	16	183	13
Quad	Real roots of quadratic equation	10	359	31
Trityp	Classifies triangle types	28	951	109
Warshall	Transitive closure of a matrix	11	305	35

level (high mutation score) under nonselective mutation. Specifically, we hypothesize that ES-selective mutation, RS-selective mutation, and RE-selective mutation will lead to test sets that are over 95% mutation adequate.

To evaluate our hypothesis, we compared selective mutation with nonselective mutation. To do this, we created test data sets that were selective mutation adequate (achieving mutation scores of 100% when run against the selective mutants), and we measured the mutation adequacy of those test data sets. This section discusses the procedure used by, and the results of, those experiments.

3.1 Experimental Procedure

Ten Fortran-77 program units that cover a range of applications were chosen for the experiments. These programs ranged in size from 10 to 48 executable statements and had from 183 to 3010 mutants. The programs are described in Table II.

We began by experimenting with ES-selective mutation. For each program, we first created the ES-selective mutants for the program (leaving out the AAR, ACR, ASR, CAR, CNR, CRP, CSR, SAR, SCR, SRC, and SVR operators). We then eliminated all equivalent mutants (previously determined by hand).

Next, we used the automatic test data generator Godzilla [DeMillo and Offutt 1991] to generate test cases to kill as many nonequivalent mutants as possible, and we generated additional test cases by hand when necessary. Godzilla attempts to generate a test case to kill a mutant by analyzing the properties a test case must have to kill the mutant and by representing these properties in the form of algebraic systems of constraints. The constraint system for a mutant represents (1) the paths that can be taken to reach the mutated statement and (2) what must be true at some execution of that statement for the mutant to cause an incorrect program state. Thus, the constraints are taken from the conditionals on the control flow graph of the program and the expression that is mutated. After creating these constraints, one for each mutant, Godzilla tries to satisfy each constraint system using a sophisticated heuristic procedure, which is

fully described elsewhere [DeMillo and Offutt 1991; Offutt 1991]. The actual values chosen are within the space described by the constraint system, but are chosen arbitrarily from within that region.

Godzilla has very little user input; it primarily uses the test program and information about the mutants. The user only provides two inputs: a domain of values for the input parameters and a limited set of assertions that can be inserted into the program. We used those assertions in six of the programs to further refine the input domains for certain variables. For example, the assertions encoded relationships among input variables, such as that an input length must be less than or equal to the declared length of the array and that stated requirements such as that certain input variables must be greater than zero. Godzilla does nothing except try to generate a test case that will kill each mutant—it has no other properties that are likely to bias the process. That is, based on a mutant, it creates a test requirement, formulates that as a constraint system, and tries to satisfy the constraint system to produce values.

To eliminate any possible bias that could be introduced by a specific test set, we generated five separate sets of selective mutation-adequate test sets for each program and level of selective mutation. All our results are based on average scores for five test sets. The variances in both test set size and number of mutants killed for the 200 test sets were all quite low.

There were 10 programs, and we generated five sets of test cases per program. We tried four different variations of selective mutation test sets. Thus we had a total of 200 distinct sets of test data. However, Godzilla cannot always kill all mutants, so a few test cases had to be created by hand. These were created by the authors. The process was to look at each remaining live mutant and try to handcraft a test case to kill it. Although the actual number varied, the overwhelming number of test cases were generated by Godzilla. The 20 sets for Banker required the most hand-generated test sets—between 1 and 5. Banker also required the most test cases—between 50 and 105. The average generated by hand was 1%, and in no case was more than 10% generated by hand.

Because of the overlap in requirements for killing mutants, many mutants are often killed by one test case, and many test cases do not kill any mutants. To be consistent with other empirical-testing researchers, we define a test case to be *ineffective* if it does not kill any mutants, and we eliminate ineffective test cases before running them on nonselective mutant sets.

Next, for each program, we created all *nonselective* mutants (using *all* mutant operators) and eliminated all equivalent mutants. We then ran each set of selective mutation-adequate test cases on the nonselective mutants and computed the mutation score for these sets. This final score measures the effectiveness of selective mutation-adequate test sets on sets of nonselective mutants, and thus provides an estimate for how close selective mutation comes to ProbSuming nonselective mutation.

Table III. Nonselective Mutation Scores of ES-Selective Mutation-Adequate Sets

Program	Test Cases	Number of Live Mutants	Mutation Score
Banker	62.8	9.8	99.64
Bub	4.6	0.4	99.87
Cal	25.4	1.8	99.93
Euclid	3.0	1.2	99.30
Find	13.2	1.4	99.85
Insert	2.6	0.2	99.95
Mid	24.6	0.0	100.00
Quad	8.0	2.8	99.15
Trityp	40.8	6.2	99.26
Warshall	3.6	4.2	98.45
Average	13.9	2.0	99.54

3.2 Experimental Results for Selective Mutation

Mutation scores for each program, averaged over the five test sets, are shown in Tables III (ES-selective), IV (RS-selective), and V (RE-selective). These tables show that test sets that were 100% adequate for selective mutation were all almost 100% adequate for nonselective mutation. In fact, the mutation scores were above 98% for all but one program (Mid under RS-selective mutation) and above 99% for all programs under ES-Selective and RE-Selective mutation. The columns are the numbers of effective test cases.

Tables VI–VIII show the savings obtained by selective mutation in terms of the number of mutants. The columns “Percentage Saved” were computed by subtracting the number of selective mutants from the number of nonselective mutants and dividing the difference by the number of nonselective mutants (the percentage of mutants that did not have to be generated and killed with selective mutation). Selective sets save from 6% to 72% of the mutants. Not surprisingly, the biggest gain is from ES-selective, which eliminates the operators that result in $O(Vals * Refs)$ mutants.

Varying the number of test cases used on each program had little effect on results. The variances in size and number of mutants killed are very small. For example, the program with the largest variation in test set size for ES-selective mutation-adequate test sets was Find, which had from 10 to 17 test cases. The largest set had a (nonselective) mutation score that was only 0.11% larger than the smallest set—the scores ranged from 99.89 to 100. In fact, the largest test set killed only one more mutant (of the 1022 that existed) than the smallest set. This was typical of the programs we looked at. In other words, even the smallest selective mutation-adequate test sets killed almost as many mutants as the largest such sets.

3.3 Combining ES- and RE-Selective Mutation

Since the mutation scores in Table IV are lower than the scores in Tables III and V, we wondered if the expression operators are the most useful

Table IV. Nonselective Mutation Scores of RS-Selective Mutation-Adequate Sets

Program	Test Cases	Number of Live Mutants	Mutation Score
Banker	90.0	9.2	99.66
Bub	3.0	6.0	98.02
Cal	38.6	6.0	99.78
Euclid	2.6	0.6	99.65
Find	13.4	8.8	99.07
Insert	1.8	5.6	98.65
Mid	6.0	32.8	80.71
Quad	7.0	1.8	99.45
Trityp	31.6	15.6	98.15
Warshall	4.2	0.0	100.00
Average	12.0	8.6	97.31

Table V. Nonselective Mutation Scores of RE-Selective Mutation-Adequate Sets

Program	Test Cases	Number of Live Mutants	Mutation Score
Banker	105.4	2.8	99.90
Bub	6.8	0.0	100.00
Cal	39.8	0.0	100.00
Euclid	3.6	0.0	100.00
Find	15.8	1.4	99.85
Insert	4.0	0.0	100.00
Mid	25.0	0.0	100.00
Quad	13.0	0.0	100.00
Trityp	47.0	0.0	100.00
Warshall	6.2	0.0	100.00
Average	17.91	0.2	99.97

Table VI. Savings Obtained by ES-Selective Mutation

Program	Nonselective Mutants	Selective Mutants	Percentage Saved
Banker	2765	629	77.25
Bub	338	147	56.51
Cal	3009	288	90.43
Euclid	195	115	41.03
Find	1022	422	58.71
Insert	460	197	57.17
Mid	183	133	27.32
Quad	359	190	47.08
Trityp	951	494	48.05
Warshall	305	115	62.30
Total	9587	2730	71.52

operators. Thus, we extended our experimentation to measure the effectiveness of combining ES- and RE-selective mutation. *Expression-selective* (*E-selective*) mutation omits both the replacement and statement operators, leaving only five Mothra operators. Table IX gives the mutation scores for

Table VII. Savings Obtained by RS-Selective Mutation

Program	Nonselective Mutants	Selective Mutants	Percentage Saved
Banker	2765	2281	17.50
Bub	338	221	34.62
Cal	3009	2772	7.88
Euclid	195	106	45.64
Find	1022	765	25.15
Insert	460	316	31.30
Mid	183	68	62.84
Quad	359	185	48.47
Trityp	951	505	46.90
Warshall	305	217	28.85
Total	9587	7436	22.44

Table VIII. Savings Obtained by RE-Selective Mutation

Program	Nonselective Mutants	Selective Mutants	Percentage Saved
Banker	2765	2620	5.24
Bub	338	308	8.88
Cal	3009	2958	1.69
Euclid	195	169	13.33
Find	1022	857	16.14
Insert	460	407	11.52
Mid	183	165	9.84
Quad	359	343	4.46
Trityp	951	903	5.05
Warshall	305	278	8.85
Total	9587	9008	6.04

Table IX. Nonselective Mutation Scores of E-Selective Mutation-Adequate Sets

Program	Test Cases	Number of Live Mutants	Mutation Score
Banker	50.6	11.6	99.57
Bub	5.6	0.2	99.93
Cal	17.6	10.4	99.63
Euclid	3.0	1.2	99.30
Find	10.8	2.8	99.70
Insert	4.4	0.2	99.95
Mid	25.0	0.0	100.00
Quad	8.8	3.0	99.09
Trityp	42.0	5.4	99.36
Warshall	2.8	3.6	98.67
Average	13.3	2.9	99.51

E-selective mutation. All scores are above 98.5%, with an average over our 10 programs of 99.5%. This indicates that the five operators ABS, AOR, LCR, ROR, and UOI are sufficient to achieve almost full mutation coverage.

Table X. Savings Obtained by E-Selective Mutation

Program	Nonselective Mutants	Selective Mutants	Percentage Saved
Banker	2765	484	82.50
Bub	338	117	65.38
Cal	3009	237	92.12
Euclid	195	89	54.36
Find	1022	257	74.85
Insert	460	144	68.70
Mid	183	115	37.16
Quad	359	174	51.53
Trityp	951	446	53.10
Warshall	305	88	71.15
Total	9587	2151	77.56

Table X gives the savings for E-selective mutation. These range from 37% to 92%, for an average of 78%. At the two extremes, Mid is a very small routine that has a relatively large number of branches and very little computation. Cal, on the other hand, contains very few branches and is mostly a straight-line series of simple computations. Since the five operators in E-selective mutation modify logical and arithmetic expressions, programs with a large number of decisions and arithmetic operations (such as Mid) can be expected to result in less savings than programs with few decisions and few arithmetic operations (such as Cal).

3.4 Comparing E-Selective Mutation with Randomly Generated Tests

As a control on our method for generating test data sets, we decided to evaluate E-selective mutation using randomly generated test data. This is an attempt to eliminate the possibility of a bias being introduced by our method of generating test data. We followed the same procedure as before, except that instead of using Godzilla, we developed random-number generators for our programs (based on the Unix `random()` function). We generated inputs to the program in groups of 200, ran them against the E-selective mutants, and then eliminated all ineffective test cases. We repeated this process until all nonequivalent mutants were killed, or until we determined that further runs were not going to kill any more mutants.

The results from this control experiment are in Table XI. As can be seen, the results are almost identical to those in Table IX. Thus we see no evidence that the method of test data generation has an impact on our results.

The biggest difference between Table IX and Table XI is that random generation required many more test cases. While this is not a subject of this research, it points out a major problem of automatically generating test data using purely random techniques—it is expensive. Additionally, the test case sizes reported in Tables IX and XI are of the *effective* test cases only. The number of test cases generated for mutation was bound by the number of mutants being targeted (this depends on which version of

Table XI. Nonselective Mutation Scores of E-Selective Randomly Generated Adequate Sets

Program	Test Cases	Number of Live Mutants	Mutation Score
Banker	61.4	10.4	99.62
Bub	4.0	0.0	100.00
Cal	19.2	9.2	99.67
Euclid	4.0	1.0	99.42
Find	15.8	0.8	99.92
Insert	8.4	0.4	99.90
Mid	26.2	0.0	100.00
Quad	6.0	0.0	100.00
Trityp	63.4	4.2	99.50
Warshall	3.2	1.2	98.56
Average	21.2	2.7	99.76

selective mutation is used, but Table II contains the total number of nonselective mutants). For all programs, we generated many more test cases using random-number generation than constraint-based generation—the lowest was 800 for Mid, and the highest was 10,000 for Trityp. Additionally, we had to generate significantly more test cases by hand when using the random-number generators, greatly increasing the cost of test data generation.

3.5 Further Reductions in the Selective Set

The E-selective mutant operators form a coherent set of operators that result in far fewer mutants than the full set of mutants (in the next section, we show that E-selective mutation is linear in the number of variable references, rather than quadratic). But that does not mean that all of the E-selective operators are necessary for mutation testing. To determine whether further mutant operators can be eliminated, more data were gathered for each E-selective mutant operator.

For each program, test sets were generated that killed all the mutants generated by each of the five E-selective mutant operators. These test sets were then run against all mutants. For example, the ABS mutants were created, and test sets were generated to kill all of the ABS mutants. These test sets were generated using the same procedure as before; we generated five independent test sets, using Godzilla as far as possible, and when Godzilla did not ensure complete coverage, a small number of test cases were generated by hand.

The test sets generated were then run against all mutants, as before. This gave us, for each E-selective mutant operator, a list of all nonselective mutants that were killed by the test set based on that operator. We then asked the question, if only four of the five E-selective mutants are used, how many nonselective mutants will be killed? We computed this by merging four of the five lists (eliminating duplicates). The results of this are shown in Table XII. We show the results of eliminating each mutant operator in turn.

Table XII. Nonselective Mutation Scores Using Four of Five E-Selective Mutation-Adequate Sets

Program	With All Five	Without UOI	Without ROR	Without LCR	Without AOR	Without ABS
Banker	99.57	98.57	99.57	99.57	91.57	91.57
Bub	99.93	98.93	99.93	99.93	99.93	97.93
Cal	99.63	95.63	99.63	99.63	99.63	99.63
Euclid	99.00	97.30	99.30	99.30	99.30	99.30
Find	99.30	98.70	99.70	99.70	99.70	99.70
Insert	99.75	98.95	99.95	99.95	99.95	97.95
Mid	99.90	98.00	100.00	100.00	100.00	94.00
Quad	100.00	99.09	99.09	99.09	99.09	97.09
Trityp	99.36	98.36	98.36	99.36	99.36	98.36
Warshall	99.67	97.67	98.67	98.67	98.67	98.67

In the table, the first column represents the nonselective mutation scores achieved by all five test sets, and the subsequent columns represent the nonselective mutation scores achieved by test sets generated *without* using that particular mutant operator. For example, the second column contains the mutation scores achieved by test data sets generated using the ABS, AOR, LCR, and ROR operators, but not UOI.

Unfortunately, the software testing literature offers no clear evidence that 100% coverage provides better testing than coverage at a lower level. Therefore we must draw conclusions from this data in a conservative way. Suppose we consider drawing the line at 99%. In Table XII, there are several scores below 99%. In the column "Without UOI," the mutation score for Cal is 95.63%. Thus, we cannot conclude that UOI should be eliminated. Likewise, the score for Banker is 91.57% in both the "Without AOR" and "Without ABS" columns.

All the scores for LCR, however, are above 99%, except Warshall, for which the score is 98.67%. Thus we could conclude that the LCR operator contributed less to mutation testing of the subject programs than the other operators. If we follow our "99%" rule, we could round the score for Warshall up and eliminate the LCR operator. However, there is good reason not to take this step. The LCR operator replaces each logical connector (AND and OR) with other logical connectors. These operators are somewhat uncommon; in particular, of the 10 programs studied, eight had no LCR mutants. It would not be safe to draw strong conclusions from such a small number of mutants; therefore we do not recommend eliminating the LCR mutant.

The lowest score for ROR is 98.36%, which is only slightly below 99%, so an argument could also be made that it contributes less to mutation testing than the other operators. The ROR operator replaces each relational operator by other relational operators and TRUE and FALSE. Although it seems safer to eliminate this operator than any of the others, the ROR operator is precisely the operator that ensures branch coverage [Offutt and Voas 1996]; therefore, a strong argument can be made for keeping it. In the

absence of further evidence, we conclude that all five E-selective mutant operators should be used.

4. QUANTIFYING THE SAVINGS OF SELECTIVE MUTATION

This section tries to justify why an approximation method such as selective mutation is needed and why selective mutation is valuable. A question we might like to ask is: which operators generate mutants that are most useful from a testing perspective? Unfortunately, we do not presently have a context for answering, or even for precisely posing, this question. This article really asks the twin questions: (a) which mutant operators generate the most mutants? and (b) what happens if we do not consider these operators? The data in this section answer question (a) and give some justification that the answer to (b) is worth investigating.

It is clearly desirable to have an equation, or model, that relates the number of mutants generated for a given program to lexical characteristics of that program. Various models have been suggested for nonselective mutation. Budd [1980] claimed that the number of mutants for a program is $O(Vals * Refs)$, where *Vals* is the number of data objects, and *Refs* is the number of data references. Acree et al. [1979] claimed that the number of mutants is $O(Lines * Refs)$ —assuming that the number of data objects in a program is proportional to the number of lines. They went on to say that this is actually $O(Lines * Lines)$ for most programs.

We have checked these claims statistically on a sample of 96 Fortran-77 programs. The programs vary in size from 5 to 735 lines of code. Using Mothra, the mutants of each program were created. The number of mutants created for an individual program ranged from a low of 76 to a high of 3,911,460. For the entire sample, 11,180,104 mutants were generated. For each claim, the corresponding regression model was fitted to the data using the method of least squares. The SAS statistical package [Freund and Littell 1981] was used for all our analysis.

Simple linear regression models with one explanatory, or independent, variable x_i can be written as:

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i = 1, 2, \dots, n. \quad (1)$$

And multiple linear regression models with p explanatory variables x_{ij} can be expressed as:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i, \quad i = 1, 2, \dots, n. \quad (2)$$

The parameters, or β_i values, represent the constants of the equations. Since we are only concerned with the limits of the equations (the Big-O forms), the β values have little significance. However, we do supply these values for the interested readers. The random variables $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ are errors that create scatter around the linear relationships.

For each model, the coefficient of determination was calculated. The *coefficient of determination*, or R^2 value, for a model is defined as the

proportion of the variability in the predicted, or dependent, variable that can be accounted for by the explanatory variables of the model [Ott 1988]. The R^2 value provides a summary measure of how well the regression equation fits the data. For example, an R^2 value of 0.95 for a particular regression model means that the explanatory variables explain 95% of the variability in the Y values. Also, the *observed significance level*, or ρ value, of each β parameter was calculated; levels of 0.05 and less were deemed statistically significant.

We found the most accurate predictors of the number of mutants to be the number of lines (*Lines*), the number of variables and constant values (*Vals*), and the number of value references (*Refs*). First, we statistically test Acree's contention that the number of mutants is $O(\text{Lines} * \text{Lines})$ for most programs, obtaining the following model:

$$\text{Mutants} = \beta_0 + \beta_1 * \text{Lines} + \beta_2 * \text{Lines} * \text{Lines}. \quad (3)$$

The β values are $\beta_0 = -37961$, $\beta_1 = 1521$, and $\beta_2 = 3$. For nonselective mutation, we find that (1) there is only a 0.94% probability that the *Lines* term does not contribute and (2) a 1.2% probability that the *Lines * Lines* term does not contribute (the ρ values). Thus, we conclude that both terms are significant. The R^2 value is 0.488.

For E-selective mutation, we find a 0.01% probability that the *Lines* term does not contribute, but a 96.68% probability that the *Lines * Lines* term does not contribute. Thus, we conclude that a more appropriate model for E-selective mutation is:

$$E - \text{SelectiveMutants} = \beta_0 + \beta_1 * \text{Lines}. \quad (4)$$

The β values are $\beta_0 = -320$ and $\beta_1 = 37$.

On the other hand, the R^2 values for the two models are 0.488 and 0.508, respectively, meaning that only about half of the variation in the number of mutants is predictable from lines of code as an explanatory variable. Thus, we conclude that Acree et al. [1979] were mistaken and that the number of lines in a program is not a valid predictor for the number of mutants.

Since Budd suggested that the number of mutants is based on the number of values times the number of value references, we obtained the following model:

$$\text{Mutants} = \beta_0 + \beta_1 * \text{Lines} + \beta_2 * \text{Vals} * \text{Refs}. \quad (5)$$

The β values are $\beta_0 = 23790$, $\beta_1 = -946$, and $\beta_2 = 1$. The R^2 value for this model was 0.988 for nonselective mutation. The ρ values were 0.0001 for each term.

For E-selective mutation, we obtained a different model:

$$E - \text{SelectiveMutants} = \beta_0 + \beta_1 * \text{Lines} + \beta_2 * \text{Refs}. \quad (6)$$

The β values are $\beta_0 = -25$, $\beta_1 = -24$, and $\beta_2 = 12$. The R^2 value for this model was 0.98, and the ρ value was 0.0001 for each term. Thus, we

conclude that E-selective mutation eliminates the *Vals* factor from the number of mutants. This is as we expected, because the major category of mutants we are eliminating is the category that replaces variable references by all data values. We also tried to use the number of language operators (as well as many other variables), but did not find a reasonable correlation.

Finally, we looked at the simple ratio of mutants to E-selective mutants (M/E) for our sample programs. On average, that ratio was approximately 17, which indicates that E-selective gives us an order-of-magnitude improvement. Unfortunately, the standard deviation is 21.1, indicating that this ratio is very program dependent. For example, in one program we had a ratio of only 1.59, and in another of 129.8. Although we tried to relate this ratio to the programs, the closest we could find is that the ratio roughly goes up as the number of mutants goes up (as can be seen in Tables VI–VIII and X), but we do not feel comfortable drawing any general conclusions from this.

Another way of reporting the improvement is to compare all mutants for the entire 96 programs. For standard mutation, we had to consider all 11,180,104 mutants. For E-selective, we only had to consider 231,972 mutants, a ratio of 48.2.

5. CONCLUSIONS AND FUTURE WORK

This article presents results which indicate that previous mutation implementations are much more expensive than necessary and which support our hypothesis that selective mutation is effective and efficient. Specifically, our results indicate that the mutant operators that replace all operands with all syntactically legal operands add very little to the effectiveness of mutation testing. Additionally, the mutant operators that modify entire statements add very little. Our data indicate that of the 22 mutant operators used by Mothra, five of these suffice to effectively implement mutation testing. Section 4 shows that standard mutation takes time quadratic in the number of data references (actually, the number of data objects times the number of data references), whereas selective mutation can be effectively implemented in time linear in the number of data references. This is a major step forward in the practical application of mutation.

The five sufficient operators are ABS, which forces each arithmetic expression to take on a zero value, a positive value, and a negative value; AOR, which replaces each arithmetic operator with every syntactically legal operator; LCR, which replaces each logical connector (AND and OR) with several kinds of logical connectors; ROR, which replaces relational operators with other relational operators; and UOI, which inserts unary operators in front of expressions. It is interesting to note that this set includes the operators that are required to satisfy branch and extended branch coverage [Offutt and Voas 1996], leading us to believe that extended branch coverage is in some sense a major part of mutation.

This research is unusual in that it takes a different direction from previous research in mutation testing. Previous researchers have tried to

enhance the strength of mutation by adding more mutation operators, taking the philosophy “more is better.” On the contrary, we find that “less is more,” or at least, that “less is nearly as good.” Previous research focused on improving the strength of mutation testing, but without a clear metric for strength. We focus on efficiency and find that strength is (by our measure) only marginally compromised.

The results reported here, especially when combined with the techniques of weak mutation [Offutt and Lee 1994] and schema-based mutation [Untch et al. 1993], lead us to hope that it may soon be possible to make mutation testing a practical reality in the form of tools that can be used by software developers. Of course, we cannot make strong claims about the statistical significance of our results. We have no basis for claiming that our programs are a representative sample or that the test sets are representative. We tried to generate test cases in a realistic way—using tool support as far as possible and hand generation thereafter. But the data certainly indicate that the idea is sound, and the analysis shows that the technique has great potential.

Our experiment has two assumptions that warrant further investigation. One is our measurement of selective effectiveness. We decided whether selective mutation was effective by computing the nonselective mutation scores of the test sets. Other measurements could be the relative abilities of the test sets to detect actual faults in the program, or other testing criteria. Another assumption is that our results will scale up to larger programs. The data in Section 2 indicate that larger programs have relatively more replacement operator mutants; thus the savings from selective mutation could be relatively greater.

These observations prompt one final question. If selective mutation involves the useful elimination of certain mutant operators, what is it about those operators that makes them easily killed by test cases that are adequate for other operators? Except for a very few special cases [Offutt and Voas 1996], there has been no analytical work on the subsumption of some mutant operators by others. Investigation of this question might shed light on the coupling effect and on the nature of the sensitivity of errors to tests.

One possible reason for the E-selective adequate test sets to also kill R-selective and S-selective mutants is that most statements contain operators or references for which E-selective mutants are generated. That is, the E-selective mutant operators modify most statements. The five operators in the E-selective set modify every operator, every variable except variables on the left side of assignments, and all constants except TRUE and FALSE. In the programs studied in Section 3, we only found two kinds of statements that are not mutated by the E-selective set. Statements that assign a boolean constant (of which there is one, “X = FALSE”) and GOTO statements are not mutated. Therefore, the E-selective mutant operators require the tester to execute every statement and test most statements in several ways.

5.1 Fault Size

Another possible explanation for the value of selective mutation can be found in the relationship between semantic and syntactic faults. Budd and

Angluin [1982] present a theoretical model of mutation in which mutants are considered to be in the “neighborhood” of the original program. Our perception of mutation may be subtly different depending on whether we think of the neighborhood of a program as being close in the syntactic sense or in the semantic sense.

In another paper [Offutt and Hayes 1996], we define the *syntactic size* of a fault to be the fewest number of statements or tokens that need to be changed to get a correct program. Mutations have traditionally modified a single statement or token; therefore mutants are very small syntactically. We define the *semantic size* of a fault to be the relative size of the input domain for which the program is incorrect. Obviously, a mutant has both a syntactic and semantic aspect, but the sizes do not always correspond. A small syntactic change can result in a very large semantic fault—that is, it can affect almost every input. Also, a major syntactic rewrite of P can result in a very small semantic change.

Mutation systems create mutants by making syntactic changes. Some of these mutants are semantically small, even size zero (equivalent mutants), and some are “trivial”—that is, the corresponding semantic fault is very large, and the mutant dies very quickly. We think it is likely that the most interesting mutants are those that are small semantically.

We suggest that the underlying goal of selective mutation is to try to only use operators that tend to produce mutants that have semantically small faults. If this model is correct, then selective mutants should contain a very high percentage of equivalent mutants. It turns out that the five operators in the selective set account for 57% of the equivalent mutants for the programs studied in Section 3.

ACKNOWLEDGMENTS

We would like to thank Aditya Mathur for initially suggesting the notion of selective mutation. We would also like to thank one of the reviewers for suggesting that E-selective mutation works because the E-selective operators modify most statements in most programs.

REFERENCES

- ACREE, A. T., BUDD, T. A., DEMILLO, R. A., LIPTON, R. J., AND SAYWARD, F. G. 1979. Mutation analysis. Tech. Rep. GIT-ICS-79/08, School of Information and Computer Science, Georgia Inst. of Technology, Atlanta, Ga. Sept.
- BUDD, T. A. 1980. Mutation analysis of program test data. Ph.D. thesis, Yale Univ., New Haven, Conn.
- BUDD, T. A. AND ANGLUIN, D. 1982. Two notions of correctness and their relation to testing. *Acta Informatica* 18, 1 (Nov.), 31–45.
- BUDD, T., DEMILLO, R. A., LIPTON, R. J., AND SAYWARD, F. G. 1980. Theoretical and empirical studies on using program mutation to test the functional correctness of programs. In *Proceedings of the 1980 ACM Principles of Programming Languages Symposium*. ACM, New York, 220–233.
- CHOI, B. AND MATHUR, A. P. 1993. High performance mutation testing. *J. Syst. Softw.* 20, 2 (Feb.), 135–152.

- DEMILLO, R. A. AND OFFUTT, A. J. 1991. Constraint-based automatic test data generation. *IEEE Trans. Softw. Eng.* 17, 9 (Sept.), 900-910.
- DEMILLO, R. A. AND OFFUTT, A. J. 1993. Experimental results from an automatic test case generator. *ACM Trans. Softw. Eng. Meth.* 2, 2 (Apr.), 109-127.
- DEMILLO, R. A., GUINDI, D. S., KING, K. N., MCCracken, W. M., AND OFFUTT, A. J. 1988. An extended overview of the Mothra software testing environment. In *Proceedings of the 2nd Workshop on Software Testing, Verification, and Analysis* (Banff, Alberta, Canada, July). IEEE Computer Society Press, Los Alamitos, Calif., 142-151.
- DEMILLO, R. A., LIPTON, R. J., AND SAYWARD, F. G. 1978. Hints on test data selection: Help for the practicing programmer. *IEEE Comput.* 11, 4 (Apr.), 34-41.
- FREUND, R. J. AND LITTELL, R. C. 1981. *SAS for Linear Models: A Guide to the ANOVA and GLM Procedures*. SAS Institute, Inc., Cary, N. Carol.
- FRANKL, P. G., WEISS, S. N., AND HU, C. 1994. All-uses versus mutation testing: An experimental comparison of effectiveness. Tech. Rep. PUCS-100-94, Dept. of Computer Science, Polytechnic Univ., Brooklyn, N.Y. Feb.
- HAMLET, R. G. 1977. Testing programs with the aid of a compiler. *IEEE Trans. Softw. Eng.* 3, 4 (July).
- HOW TAI WAH, K. S. 1995. Fault coupling in finite bijective functions. *J. Softw. Testing Verif. Reliab.* 5, 1 (Mar.), 3-47.
- KING, K. N. AND OFFUTT, A. J. 1991. A Fortran language system for mutation-based software testing. *Softw. Pract. Exper.* 21, 7 (July), 685-718.
- MATHUR, A. P. 1991. Performance, effectiveness, and reliability issues in software testing. In *Proceedings of the 15th Annual International Computer Software and Applications Conference* (Tokyo, Japan, Sept.). 604-605.
- MATHUR, A. P. AND WONG, W. E. 1994. An empirical comparison of data flow and mutation-based test adequacy criteria. *J. Softw. Testing Verif. Reliab.* 4, 1 (Mar.), 9-31.
- MORELL, L. J. 1984. A theory of error-based testing. Ph.D. thesis, Tech. Rep. TR-1395, Univ. of Maryland, College Park, Md.
- OFFUTT, A. J. 1991. An integrated automatic test data generation system. *J. Syst. Integration* 1, 3 (Nov.), 391-409.
- OFFUTT, A. J. 1992. Investigations of the software testing coupling effect. *ACM Trans. Softw. Eng. Methodol.* 1, 1 (Jan.), 3-18.
- OFFUTT, A. J. AND HAYES, J. H. 1996. A semantic model of program faults. In *Proceedings of the 1996 International Symposium on Software Testing, Verification, and Analysis* (San Diego, Calif., Jan.). ACM Press, New York, 195-200.
- OFFUTT, A. J. AND LEE, S. D. 1994. An empirical evaluation of weak mutation. *IEEE Trans. Softw. Eng.* 20, 5 (May), 337-344.
- OFFUTT, A. J. AND VOAS, J. M. 1996. Subsumption of condition coverage techniques by mutation testing. Tech. Rep. ISSE-TR-96-100, Dept. of Information and Software Systems Engineering, George Mason Univ., Fairfax, Va.
- OFFUTT, A. J., ROTHERMEL, G. AND ZAPF, C. 1993. An experimental evaluation of selective mutation. In *Proceedings of the 15th International Conference on Software Engineering* (Baltimore, Md., May). IEEE Computer Society Press, Los Alamitos, Calif., 100-107.
- OTT, L. 1988. *An Introduction to Statistical Methods and Data Analysis*, 3rd ed. PWS-Kent, Boston, Mass.
- UNTCH, R., OFFUTT, A. J., AND HARROLD, M. J. 1993. Mutation analysis using program schemata. In *Proceedings of the 1993 International Symposium on Software Testing, Analysis, and Verification* (Cambridge, Mass., June). ACM, New York, 139-148.
- WONG, W. E. 1993. On mutation and data flow. Ph.D. thesis, Tech. Rep. SERC-TR-149-P, Software Engineering Research Center, Purdue Univ., W. Lafayette, Ind.
- WEYUKER, E. J., WEISS, S. N., AND HAMLET, R. G. 1991. Comparison of program testing strategies. In *Proceedings of the 4th Symposium on Software Testing, Analysis, and Verification* (Victoria, British Columbia, Canada, Oct.). ACM, New York, 1-10.

Received January 1994; revised November 1995; accepted January 1996