

Automatically Detecting Equivalent Mutants and Infeasible Paths

A. JEFFERSON OFFUTT¹* AND JIE PAN[†]

¹ISSE Department, George Mason University, Fairfax, VA 22030, U.S.A.

²Template Software, 45365 Vintage Park Plaza, Suite 100, Pulles, VA 20166, U.S.A.

SUMMARY

Mutation testing is a technique for testing software units that has great potential for improving the quality of testing, and thereby increasing the ability to assure the high reliability of critical software. It will be shown that recent advances in mutation research have brought a practical mutation testing system closer to reality. One recent advance is a partial solution to the problem of automatically detecting equivalent mutant programs. Equivalent mutants are currently detected by hand, which makes it very expensive and time-consuming. The problem of detecting equivalent mutants is a specific instance of a more general problem, commonly called the *feasible path problem*, which says that for certain structural testing criteria some of the test requirements are infeasible in the sense that the semantics of the program imply that no test case satisfies the test requirements. Equivalent mutants, unreachable statements in path testing techniques, and infeasible DU-pairs in data flow testing are all instances of the feasible path problem. This paper presents a technique that uses mathematical constraints, originally developed for test data generation, to detect some equivalent mutants and infeasible paths automatically. © 1997 by John Wiley & Sons, Ltd.

Softw. Test. Verif. Reliab., 7, 165–192 (1997)

No. of Figures: 8 No. of Tables: 7 No. of References: 24

KEY WORDS software testing; mutation testing; constraints; feasible path analysis

1. INTRODUCTION

Mutation testing is a technique, originally proposed by DeMillo *et al.* (1978) and Hamlet (1977), that requires testers to create test data that cause a finite, well-specified set of faults to result in failure. The testers do 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 the absence of those types of faults) or cause the test program to fail (detecting a fault). The technique thus serves two goals: it provides a test adequacy criterion, and leads to the detection of faults in the program being tested.

Unit level testing techniques such as mutation hold great promise for improving the quality of software. Although most of these techniques are currently so expensive that

* Correspondence to: A. Jefferson Offutt, ISSE Department, George Mason University, Fairfax, VA 22030, U.S.A.

practical use is rare, recent advances in mutation research have brought a practical mutation testing system closer to reality. Commercial tools (e.g. PiSCES, Voas *et al.*, 1991) have recently become available, although they are still not widely used in industry. A major problem with mutation testing is that it is too expensive to apply; one aspect of the cost is identifying programs that are intended to be faulty, but in fact are not. These programs are said to be *equivalent*. Previous papers have presented ways to speed up mutation testing (Howden, 1982; Offutt *et al.*, 1996; DeMillo and Offutt, 1991; Untch *et al.*, 1993), including a technique for partially solving the problem of automatically detecting equivalent mutants that was based on compiler optimization techniques (Offutt and Craft, 1994). Equivalent mutants are currently detected by hand, which makes it very expensive and time-consuming. This paper presents a new technique for partial detection of equivalent mutants based on constraint-based testing that gives much better results than the compiler optimization techniques.

The rest of this paper presents a method for detecting equivalent mutants. The remainder of the introduction provides background material on mutation testing, presents the problem of recognizing equivalent mutants, and discusses previous work on this problem. Section 2 describes constraint-based testing, then provides full details, including algorithms, on how to use constraints to detect infeasible constraint systems, and therefore equivalent mutants. Next, a proof-of-concept tool that implements the equivalent mutant detection strategies is discussed; then empirical results from using the tool are presented. Conclusions and several specific suggestions for improving the technique and the tool are given.

1.1. Mutation testing overview

Mutation testing helps testers create test data by interacting with them to strengthen the quality of the test data. Faults are introduced into programs by creating many versions of the software, each containing one fault. Test cases are used to execute these faulty programs with the goal of causing each faulty program to produce incorrect output (i.e. to *fail*). Hence the term mutation; faulty programs are *mutants* of the original, and a mutant is *killed* when it fails. When this happens, the mutant is considered *dead* and no longer needs to remain in the testing process because the faults represented by that mutant have been detected.

Figure 1 shows a short function and five mutations; the original program is shown to the left, and the mutant programs are represented on the right by the lines preceded by the Δ symbol. Note that each mutated statement represents the change that is made to create one mutant program. A mutation *operator* is a rule that is applied to a program to create mutants. The Mothra mutation system (DeMillo *et al.*, 1988) uses 22 mutation operators to test Fortran 77 programs, as shown in Table I. The Mothra operators replace each operand by all other syntactically legal operands ($\Delta 1$, $\Delta 3$ and $\Delta 5$ in Figure 1), modify expressions by replacing operators and inserting new operators ($\Delta 2$), and modify entire statements ($\Delta 4$).

The mutation testing process begins with an automated mutation system creating the mutants of a test program. Test cases are then added, either manually or automatically, to the mutation system and the user checks the output of the program for each test case to see if it is correct. If incorrect, a fault has been found and the program must be modified and the process restarted. If the output is correct, that test case is executed against each live mutant. If the output of a mutant differs from that of the original

Original Program		With Embedded Mutants
<pre> FUNCTION Min (I, J : Integer) RETURN Integer IS MinVal : Integer; BEGIN MinVal := I; IF (J < I) THEN MinVal := J; END IF; RETURN (MinVal); End Min; </pre>		<pre> FUNCTION Min (I, J : Integer) RETURN Integer IS MinVal : Integer; BEGIN MinVal := I; MinVal := J; IF (J < I) THEN IF (J > I) THEN IF (J < MinVal) THEN MinVal := J; TRAP; MinVal := I; END IF; RETURN (MinVal); End Min; </pre>
	Δ1	
	Δ2	
	Δ3	
	Δ4	
	Δ5	

Figure 1. Function MIN.

Table I. Mothra mutation operators for Fortran 77

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

program, it is assumed to be incorrect, the mutant is killed, and the test case is kept as an 'effective' test case.

In Figure 1, note that the mutant $\Delta 3$ replaces the variable l with $MinVal$. In the previous statement, the value of l was assigned to $MinVal$. At this point in the program, l and $MinVal$ will *always* have the same value, thus $\Delta 3$ represents an equivalent mutant.

After each mutant has been executed with each test case, each remaining mutant falls into one of two categories. One, the mutant is killable, but the set of test cases is insufficient to kill it. In this case, a new test case needs to be created. Two, the mutant is functionally *equivalent* to the original program. An equivalent mutant will always produce the same output as the original program, so no test case can kill it. Once identified as equivalent, there is no need for the mutant to remain in the system for further consideration.

1.2. Equivalent mutant problem

This paper focuses on the problem of detecting equivalent mutants, which has been an obstacle to the practical application of mutation. Without detecting all the equivalent mutants, the mutation score will never be 100%. Thus, the tester will not have complete confidence in the program and the test data. Worse, the tester will be left wondering whether the remaining mutants are equivalent or if the test set is insufficient. Detecting equivalent mutants by hand is very time-consuming, which contributes to making the cost of mutation testing prohibitively high.

1.2.1. Distribution of equivalent mutants among mutant types

Equivalent mutants are not evenly distributed among the mutant types; they tend to cluster among only a few types. Table II summarizes statistics from the programs used in

Table II. Equivalent mutants among mutant types (programs used in this paper)

Mutant type	% of equivalent	% of all mutants
ABS	47.19	4.30
ACR	14.10	1.28
SCR	7.05	0.64
UOI	6.04	0.55
SRC	4.89	0.45
SVR	4.46	0.41
ROR	3.60	0.33
SDL	2.16	0.20
CRP	1.58	0.14
AAR	1.44	0.13
RSR	1.44	0.13
LCR	1.15	0.10
ASR	1.15	0.10
CSR	1.01	0.09
SAR	1.01	0.09
All others	1.73	0.16
Total	100.00	9.10

Section 3.2 of this paper. The first column in the table gives the mutant operator type and the second column gives the percentage of the total number of equivalent mutants represented by each type. The third column gives the percentage of all mutants that are equivalent of that type—in the programs studied, 9.1% of the mutants were equivalent.

Offutt and Craft (1994) and Budd (1980) give similar statistics of the distribution of equivalent mutants among mutant types. Both sources indicate one very interesting fact. One mutant type, *absolute value insertion* (ABS), has many more equivalent mutants than any other mutant type. The ABS mutant operator inserts three unary operators before each expression—ABS computes the absolute value of the expression, NEGABS computes the negative of the absolute value, and ZPUSH kills the mutant if the expression is zero, otherwise it does nothing (this forces the tester to cause each expression to have the value zero, a common testing heuristic). The fact that equivalent mutants are clustered among the ABS mutants is used in the techniques presented in this paper.

1.2.2. Do procedures for automatically detecting equivalence exist?

Budd and Angluin (1982) examine the relationships between equivalence and test data generation. They prove that if there is a computable procedure for checking if two programs are equivalent, there is also a computable procedure for generating adequate test data for a program and vice versa. They also show that, in general, neither of these computable procedures exists. Thus, there cannot be a complete algorithmic solution to the equivalence problem. That is, detecting equivalence either between two arbitrary programs or two mutants is an undecidable problem.

Fortunately, the equivalent mutant problem has one advantage over the general equivalence problem. Specifically, it is not necessary to determine the equivalence of arbitrary pairs of programs. Because of the definitions of mutation operators, mutants are syntactically very much like their original program. Although Budd and Angluin also prove that this problem is undecidable, it has been suggested that for many specific cases, equivalence can be decided (Acree, 1980; Offutt and Craft, 1994). This paper reports results from techniques and heuristics that automatically detect many of the equivalent mutants.

1.2.3. Previous work in detecting equivalent mutants

It is obvious that automatically detecting equivalent mutants can save much time and energy for the testers, but Acree (1980) found that it could also prevent people from making errors in marking equivalent mutants. In a study of 50 mutants, half of which were equivalent, Acree found that people judged mutant equivalence correctly only about 80% of the time. The people marked equivalent mutants non-equivalent (type 2 errors) 12% of the time and non-equivalent mutants equivalent (type 1 errors) 8% of the time. Because type 2 errors can be corrected during later testing, a conservative approach would be to try to eliminate type 1 errors, even if that means allowing a few type 2 errors. If an automated tool is used, the effect of the types of errors depends on how the tool is used. Type 1 errors result in the mutation score being overestimated, and type 2 errors result in the mutation score being underestimated. If all mutants that are *not* determined to be equivalent are examined by hand, then the remaining equivalent mutants can be determined by hand. If the remaining mutants are not examined by hand (as suggested in this paper) then type 1 errors will result in an overestimate of the number of equivalent mutants, hence an overestimate of the mutation score. Likewise, type 2 errors will result

in an underestimate of the mutation score. Thus, type 1 errors could be viewed as being less conservative. The advantage of using automated techniques to detect equivalent mutants is that they can be more conservative by avoiding as many type 1 errors as possible.

Baldwin and Sayward (1979) proposed using compiler optimization techniques to detect equivalent mutants. The key intuition behind this approach is that code optimization transformations produce equivalent programs, and some equivalent mutants are, in some sense, either optimizations or de-optimizations of the original program. So when a mutant satisfies a code optimization rule, algorithms can detect that it is equivalent. Baldwin and Sayward describe six types of compiler optimization techniques that can be used to detect equivalent mutants.

Offutt and Craft (1994) designed algorithms for these six techniques, and developed and implemented Equalizer, an automated tool for detecting equivalent mutants. They found that Equalizer detected an average of 10% of the equivalent mutants for 15 programs.

In his Ph.D. dissertation, Offutt (1988) presented a technique for using mathematical constraints for testing, which is called 'constraint-based testing' (CBT). How constraints can be used to generate test cases to satisfy mutation testing was presented by DeMillo and Offutt (1991). The dissertation also suggested using CBT to detect equivalent mutants, but did not offer details for how to do it. This paper develops the idea, presents specific strategies and algorithms for detecting equivalent mutants, and presents results from an implementation of these algorithms.

1.3. Feasible path problem

For most structural testing criteria, some of the test requirements are infeasible in the sense that the semantics of the program imply that no test case exists that meets the test requirements. Goldberg *et al.* (1994) defined the problem as follows: 'given a description of a set of control flow paths through a procedure, *feasible path analysis (FPA)* determines if there is input data that causes execution to flow down some path in the collection'. This is generalized here to the *feasible test problem (FTP)*: given a requirement for a test case, the feasible test problem is to determine if there is input data that can satisfy the requirement. Feasible path analysis is one instance of this problem. Mutation provides another—the FTP is to decide whether the mutant is equivalent or killable.

The problem also arises in: (1) statement coverage, where the FTP is whether the statement is reachable; (2) branch coverage techniques (Myers, 1979), where the FTP is whether a branch, predicate, condition or combination of conditions can be covered; and (3) data flow testing (Frankl and Weyuker, 1988), where the FTP is whether a definition-clear subpath can be found between a DU-pair. Unfortunately, this problem is formally undecidable (Goldberg *et al.*, 1994; DeMillo and Offutt, 1991) for all variations mentioned here; thus approximations must be used.

The most recent work on the feasible test problem is by Goldberg *et al.* (1994) and Jasper *et al.* (1994). Goldberg *et al.* presented a structural testing system that uses a theorem prover to attempt to decide if test requirements can be satisfied. Jasper *et al.* also used theorem proving techniques, attempted to determine the feasibility of expressions containing references to arrays, and applied their techniques to the modified condition decision coverage criterion.

This paper focuses on the FTP within the context of mutation testing, and uses a very

different approach to determine feasibility. Specifically, a heuristic-based set of transformations is applied to mathematical constraints to look for infeasibility. These transformations, and their application to infeasibility, are formally proved in the following sections.

2. USING CONSTRAINTS TO DETECT EQUIVALENT MUTANTS

In this paper, a *constraint* is a mathematical algebraic expression that restricts the input space of the program to be the portion of the input domain that satisfies a certain goal. As a simple example, $(x > 0)$ describes the portion of the input domain where x is positive. More complicated constraints can be used to describe higher-level goals, such as that an array must be sorted or that a shape represented by a set of points must be rectangular.

CBT (DeMillo and Offutt, 1991) uses constraints for automatic test data generation. In CBT, a constraint represents the conditions under which a mutant will die. The technique in this paper uses the fact that if a test case kills the mutant, the constraint system will be satisfied by that test case. If the constraint system cannot be true, then there is no test case that can kill the mutant and the mutant is equivalent. The general approach to using constraints to detect equivalent mutants is to look for infeasibility in constraint systems.

This section gives a detailed discussion of how constraints can be used to detect equivalent mutants, and how the procedure works. Since constraint-based testing was previously used for test data generation, this technique is introduced with constraint-based test data generation. Then proofs are given for three theorems that show how CBT can be used to detect equivalent mutants.

2.1. The CBT technique

CBT uses the concepts of mutation analysis to create test data automatically. This test data is designed specifically to kill the mutants of the test program. Such test data can be used to kill mutants within a mutation system. For each mutant, a test case is said to be *effective* if it causes the mutant to fail, and a test case is *ineffective* if it does not. For the following, P is a program, M is a mutant of P on statement S , and t is a test case for P . The *state* of a program is, as usual, the values of all data items and the program counter. To kill M , t has to have the following three broad characteristics.

- (1) **Reachability.** Since a mutant is represented as a syntactic change to an executable statement, and the other statements in the mutated program are syntactically equal to the statements in the original program, a minimum requirement for a test case to kill the mutant is that it executes the mutated statement. This characteristic is called the *reachability condition*.

If t cannot reach S , it is guaranteed that t will not kill M (DeMillo and Offutt, 1991).

- (2) **Necessity.** For a test case to kill a mutant, it must be able to cause the mutant to have an incorrect state if it reaches the mutated statement. This characteristic is called the *necessity condition*.

For t to kill M , it is *necessary* that if S is reached, the state of M immediately following some execution of S must be different from the state of P at the same point (DeMillo and Offutt, 1991).

To see why, note that M is syntactically equal to P except for the mutated statement S . Thus, if the states of the two programs do not differ after some execution of S , they will never differ. If S is in a loop, necessity requires that the state be incorrect after any given execution of the mutated statement. Of course, the incorrect state must propagate through the program to the last execution of the statement, but that requirement is part of sufficiency.

- (3) **Sufficiency.** The *sufficiency condition* is that the final state of M differs from that of P (DeMillo and Offutt, 1991).

Note that the necessity condition specifically does not include the reachability condition. That is, a test case can satisfy the necessity condition even if it does not actually execute the statement. Although this may seem counter-intuitive, the orthogonality makes the necessity constraints easier to create by automated tools.

Let D represent the entire domain of all test cases t for P . In light of the three conditions above, D can be divided in several ways for each mutant:

- (a) $D = D_r \cup D_{\bar{r}}$, where D_r is the portion of D that will satisfy the *reachability condition* C_r for a given mutant and $D_{\bar{r}}$ is the portion of D that will not.
- (b) $D = D_n \cup D_{\bar{n}}$, where D_n is the portion of D that will satisfy the *necessity condition* C_n for a given mutant and $D_{\bar{n}}$ is the portion of D that will not.
- (c) $D = D_s \cup D_{\bar{s}}$, where D_s is the portion of D that will satisfy the *sufficiency condition* C_s for a given mutant and $D_{\bar{s}}$ is the portion of D that will not.

The letter C stands for a condition. The subscripts r, n and s are taken from the terms reachability, necessity and sufficiency. From the definitions of the above three conditions and sub-domains of P , the following facts are observed.

Fact 1

t is an effective test case that will kill M if and only if $t \in D_s$ for M .

Fact 2

If t is an effective test case that will kill M then $t \in D_r \cap D_n$.

Fact 3

$$D_s \subseteq D_r \cap D_n$$

A test case from the intersection of D_r and D_n will reach the mutated statement and cause an incorrect intermediate state to be created. But this incorrect state may not always propagate through the execution to result in incorrect output. Morell (1990) discussed various reasons for this, including masking of data values and error correction. Thus the set of test cases that cause incorrect output is actually a subset of $D_r \cap D_n$.

Unfortunately, finding t such that $t \in D_r$ is an undecidable problem. This is because determining whether t executes S is reducible to the halting problem. Thus, a weaker condition is defined. C_R is defined such that if S is executed, then C_R is true. D_R is the

domain that contains all inputs t that satisfy C_R . Of course, a trivial solution to C_R is $C_R = \text{true}$, but in practice the aim is to find the strongest reachability condition possible. Since $C_r \Rightarrow C_R$, the following fact is clear.

Fact 4

$$D_r \subseteq D_R$$

Figure 2 is a Venn diagram that graphically shows these domains and their relationships for one mutant. The tests that kill the mutant (D_s) are located in the intersection of the tests that reach the mutant (D_R) and the tests that satisfy necessity (D_n). The above facts and this diagram are used in later proofs.

CBT uses a *path expression* to describe the reachability condition (the weaker condition), C_R , for a statement. A path expression for a statement S in a program P is an algebraic expression that describes a condition on test cases that will be true when P reaches S . Path expressions usually describe multiple paths to S by using a disjunctive formula, where each clause represents a separate path. Path expressions are automatically derived from the program by extracting the predicate expressions on the program's control flow graph. There are a variety of ways to do this, and a variety of analysis techniques (most of which are based on symbolic evaluation) are used to refine or simplify the expressions. As should be expected, more effort in the analysis usually leads to more successful results.

Since mutation operators represent syntactic changes, a test case that satisfies the necessity condition must ensure that the syntactic change effected by the mutation results in an incorrect state for the program. CBT uses a *necessity constraint* to describe this necessity condition C_n . That is, if a test case that satisfies the necessity condition will cause the mutated statement to be reached, the state immediately following some execution of the mutated statement will be incorrect.

As an example, Figure 3 shows the function Mid with a mutant that replaces the relational operator $<$ with \leq on statement 5. The path expression for statement 5 is taken from the true branch of statement 2, and the false branch of statement 3, giving $(Y < Z) \wedge (X \geq Y)$. Note that in this case, $C_R = C_r$. The necessity constraint for the mutant requires that the two predicates evaluate to different results, i.e. $((X < Z) \neq (X \leq Z))$. CBT tries to generate a test case t by satisfying a constraint system, which can be either a reachability constraint (a path expression), a necessity constraint, or a conjunction of a path expression and a necessity constraint. t must be in the union of D_R and D_n , i.e. $t \in D_R \cup D_n$. Of course, CBT is not always successful, even if the constraint system is satisfiable.

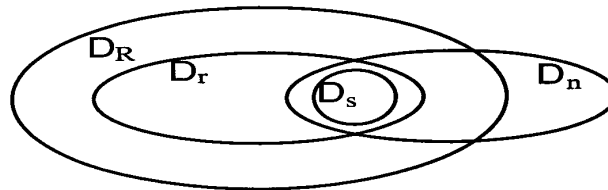


Figure 2. Input domain subsets.

```

      FUNCTION Mid (X, Y, Z : Integer)
      RETURN Integer IS
      MidVal : Integer;
      BEGIN
1       MidVal := Z;
2       IF (Y < Z) THEN
3         IF (X < Y) THEN
4           MidVal := Y;
5         ELSE IF (X < Z) THEN
      Δ    ELSE IF (X <= Z) THEN
6           MidVal := X;
7         END IF;
8       ELSE
9         IF (X > Y) THEN
10          MidVal := Y;
11        ELSE IF (X > Z) THEN
12          MidVal := X;
13        END IF;
14      END IF;
15      RETURN (MidVal);
16    END IF;

```

Figure 3. The function Mid.

2.2. Constraints and detecting equivalent mutants

In Section 1.1, an equivalent mutant was said to have the same functional behaviour as the original program. However, the term ‘the same functional behaviour’ was not formally defined. Now, equivalent mutants are formally defined in terms of inputs and outputs.

Definition

Let P be a program, M a mutated program of P , and $P(t)$ and $M(t)$ be the outputs of P and M on t . Then M is an *equivalent mutant program* of P if and only if $P(t) = M(t)$ for all t , $t \in D$.

The above definition says that if a mutant is functionally equivalent to the original program, it is impossible to find any test data to kill the mutant, i.e.

$$\neg(\exists t|_{t \in D} \bullet P(t) \neq M(t)) \Leftrightarrow \forall t|_{t \in D} \bullet P(t) = M(t)$$

To support efforts in automatically detecting equivalent mutants, the following three theorems are stated and proved. These proofs are based on the definition of an equivalent mutant, descriptions of input domains, and the facts given in Section 2.1. These theorems are not deep; the essential concepts are actually encoded in the definitions.

Theorem_1

Let D_r be the domain in which test cases satisfy the reachability condition (C_r) for a mutant M . If C_r is infeasible (D_r is empty) then M is equivalent. That is, $D_r = \emptyset \Rightarrow M$ is equivalent.

Proof of Theorem_1

- (1) M is equivalent $\Leftrightarrow D_s = \emptyset$ —Definition, Fact 1
- (2) $D_r \cap D_n \supseteq D_s$ —Fact 3
- (3) $D_r = \emptyset \Rightarrow D_s = \emptyset$ —Rules of sets, 2
- (4) $D_r = \emptyset \Rightarrow M$ is equivalent —Substitution of 1 in 3

Theorem_2

Let D_n be the domain in which test cases satisfy the necessity condition (C_n) for a mutant M . If C_n is infeasible (D_n is empty) then M is equivalent. That is, $D_n = \emptyset \Rightarrow M$ is equivalent.

Proof of Theorem_2

- (1) M is equivalent $\Leftrightarrow D_s = \emptyset$ —Definition, Fact 1
- (2) $D_r \cap D_n \supseteq D_s$ —Fact 3
- (3) $D_n = \emptyset \Rightarrow D_s = \emptyset$ —Rules of sets, 2
- (4) $D_n = \emptyset \Rightarrow M$ is equivalent —Substitution of 1 in 3

Theorem_3

Let D_r be the domain in which test cases satisfy the reachability condition (C_r) for a mutant M , and let D_n be the domain in which test cases satisfy the necessity condition (C_n) for M . If $C_r \wedge C_n$ is infeasible ($D_r \cap D_n$ is empty) then M is equivalent.

Proof of Theorem_3

- (1) M is equivalent $\Leftrightarrow D_s = \emptyset$ —Definition, Fact 1
- (2) $D_r \cap D_n \supseteq D_s$ —Fact 3
- (3) $D_r \cap D_n = \emptyset \Rightarrow M$ is equivalent —Substitution of 1 in 2

If the weaker reachability condition (C_R) is used instead of the reachability condition (C_r), since $C_r \Rightarrow C_R$, the following two claims can easily be derived.

- (a) $D_R = \emptyset \Rightarrow M$ is equivalent.
- (b) $D_R \cap D_n = \emptyset \Rightarrow M$ is equivalent.

Proof of the Derivations

- (1) $D_r = \emptyset \Rightarrow M$ is equivalent —Theorem_1
- (2) $D_R \supseteq D_r$ —Fact 4
- (3) $D_R = \emptyset \Rightarrow D_r = \emptyset$ —Rules of sets, 2
- (4) $D_R = \emptyset \Rightarrow M$ is equivalent —Transition of implication, 1, 3
- (5) $D_r \cap D_n = \emptyset \Rightarrow M$ is equivalent —Theorem_3
- (6) $D_R \cap D_n = \emptyset \Rightarrow D_r \cap D_n = \emptyset$ —Rules of sets, 2

(7) $D_R \cap D_n = \emptyset \Rightarrow M$ is equivalent—Transition of implication, 5, 6

Section 2.1 mentions that CBT uses *path expression constraint systems* to represent reachability conditions (the weaker conditions) and *necessity constraint systems* to represent necessity conditions. So the following statements are true.

- (a) If a path expression constraint system (C_R) for a statement modified by a mutant M is infeasible, then the set of test cases (D_R) that can kill M is empty—implying that M can never be killed. So M is equivalent.
- (b) If a necessity constraint system (C_n) for a mutant M is infeasible, then the set of test cases (D_n) that can kill M is empty—implying that M can never be killed. So M is equivalent.
- (c) If a constraint system that is a conjunction of a path expression constraint system and a necessity constraint system ($C_R \wedge C_n$) is infeasible, then the set of test cases ($D_R \cap D_n$) that can kill M is empty—implying that M can never be killed. So M is equivalent.

To decide if a constraint system is infeasible, there must be a contradiction in the constraint system itself. For example, the constraint system $(x > 0) \wedge (x < 0)$ is a contradiction, because x can never be assigned a value that is both greater than 0 and less than 0. If a mutant M has the above constraint as a path expression associated with it, then M is equivalent.

So far, the problem of detecting equivalent mutants has been translated to the problem of recognizing contradictions in mathematical constraint systems. Test data generation uses constraints to generate test cases, while equivalent mutant detection uses constraints to detect equivalent mutants.

2.3. Constraint representation

Before describing how to use constraints to detect equivalent mutants, it is appropriate to describe how constraints are represented in CBT. The basic component of a constraint is an algebraic expression composed of variables, parentheses and programming language operators. Expressions are taken directly from the test program and come from the right-hand sides of assignment statements, predicates within decision statements, etc. A constraint is a pair of algebraic expressions related by one of the conditional operators $\{>, \geq, <, \leq, =, \neq\}$. Constraints evaluate to one of the Boolean values TRUE or FALSE and can be modified by the negation operator NOT (\neg). A *clause* is a list of constraints connected by the logical operators AND (\wedge) and OR (\vee). A *conjunctive clause* uses only the logical AND. All constraints are kept in *disjunctive normal form* (DNF), which is a list of conjunctive clauses connected by logical ORs. For example, $(x > 0)$ represents a constraint, $(x > 0) \wedge (y < 0)$ is a conjunctive clause, and $((x > 0) \wedge (y < 0)) \vee (z = 0)$ is a disjunctive formula.

In CBT, a DNF formula is referred to as a *constraint system*. DNF is used during two steps. In constraint generation, each conjunctive clause within a path expression represents a unique path to a statement. During constraint satisfaction, only one conjunctive clause needs to be satisfied. Constraints are originally created using the variables that occur in

the program text. Unfortunately, this includes variables that are ‘internal’ to the program, i.e. variables that are not given values as part of the test case. For test case generation, symbolic evaluation (King, 1976; Offutt, 1991) is used to rewrite the variables to be in terms of input variables.

2.4. Strategies for detecting equivalent mutants

The general strategy for detecting equivalent mutants is to find contradictions in the constraint systems. Because recognizing infeasible constraint systems is generally undecidable, the problem cannot be completely solved algorithmically. However, equivalent mutants are currently detected by hand, which means that partial solutions are potentially valuable.

One approach for detecting infeasible constraint systems is to apply off-the-shelf general theorem provers. Although this approach was considered, it was rejected for two reasons. First, a general purpose theorem prover would provide much more than is needed. Second, the difficulty of integrating such a system with already-existing software did not seem to be worth the effort for this research. Of course, if these results were to be used in a production system, a special purpose theorem prover might be worth the effort and probably should be used.

This general strategy is applied through a collection of special case analysis techniques and heuristics to recognize infeasible constraint systems. The focus is on cases that occur the most in mutation, based on the observation that mutants differ from their original programs in small, well defined ways. Section 3.2 presents empirical results that measure how well these strategies work.

This paper describes and evaluates three broad strategies that attempt to recognize infeasible constraint systems. These are: negation; constraint splitting; and constants comparison.

2.4.1. Negation

Definition 1

Constraint C1 is the *negation* of C2 if and only if the domains they describe:

- (a) are non-overlapping; and
- (b) cover the entire domain of the variables in C1 and C2.

To recognize infeasible constraint systems, this paper concentrates on constraint systems that are non-overlapping but not necessarily domain covering. The notion of partial negation is used to loosen the restriction of covering the entire domain in negation.

Definition 2

Constraint C1 is a *partial negation* of C2 if and only if the domains they describe:

- (a) are non-overlapping; and

- (b) do not cover the entire domain.

Definition 3

Two constraints are *semantically equal* if they describe the same domain.

Definition 4

Two constraints are *syntactically equal* if they describe the same domain and also have the same string of symbols.

Clearly, two syntactically equal constraints are also semantically equal.

Examples

- (1) Let A be the constraint $x > 1$ and B be the constraint $x \leq 1$. Then A is the *negation* of B , and B is the *negation* of A (negation is commutative). Both constraints cannot be satisfied at the same time (their domains are non-overlapping), but the domain of x that makes the two constraints TRUE covers the entire domain of x .
- (2) Let A be the constraint $x > 1$ and B be the constraint $x < 1$. Then A is a *partial negation* of B , and B is a *partial negation* of A . Both constraints cannot be satisfied at the same time, but the domain of x that makes the two constraints TRUE does not cover the entire domain of x .
- (3) Suppose constraint A is $x > 0$ and constraint B is $x > 0$. Then A and B are syntactically equal. Thus, A and B are semantically equal.
- (4) Let x be an integer variable, A be the constraint $x > 0$ and B be the constraint $x \geq 1$. A and B are not syntactically equal, but they are semantically equal.

The negation strategy is the basic technique used to recognize infeasible constraints. Given two constraints, *negation* or *partial negation* is first used to rewrite one of the constraints, then these two constraints are compared. If they are syntactically equal, the constraints conflict, and the constraint system is infeasible. So, a mutant with this infeasible constraint system is equivalent.

For example, assume two constraints A and B , where A is $(x + y) > z$ and B is $(x + y) \leq z$. The negation of A is $(x + y) \leq z$, denoted A' . Since A' and B are syntactically equal, A and B conflict. The negation algorithm is given in Figure 4; it uses the negation and partial negation functions defined in Table III.

2.4.2. Constraint splitting

Constraint splitting is also used to recognize infeasible constraints. A commonly occurring case is a necessity constraint such as $(x + y) > 0$, together with a path expression such as $(x < 0) \wedge (y < 0)$. The negation strategy cannot recognize that these conflict.

To detect such conflicts, a strategy called *constraint splitting* is used. Given two constraints (say C and D), two new constraints (say A and B) are generated such that $C \Rightarrow A \vee B$. Then A and B are compared with D . The following proves that if both A

```

algorithm:    Negation (A, B)
precondition: A and B are properly initialized constraints.
postcondition: Returns conflict if A and B conflict, no-conflict otherwise.

BEGIN
  -- Use Table 3 (Negation table) to negate A.
  neg-A = Negate(A)
  IF (neg-A syntactically equals B)
    RETURN conflict
  ELSE
    IF (the relation operator in A is one of {>, <, =})
      -- Use Table 3 (Negation table) to negate B.
      partneg-A = PartialNegate1(A)
      IF (partneg1-A syntactically equals B)
        RETURN conflict
      ELSE
        partneg2-A = PartialNegate2(A)
        IF (partneg2-A syntactically equals B)
          RETURN conflict
        ELSE
          RETURN no-conflict
        END IF
      END IF
    END IF
  END IF
END Negation

```

Figure 4. The negation algorithm—decides if two constraints conflict with each other.

Table III. Negation and partial negation

Constraint C	Negation of C	Partial negation of C	
		Partial negation1	Partial negation2
expr1 > expr2	expr1 ≤ expr2	expr1 < expr2	expr1 = expr2
expr1 ≥ expr2	expr1 < expr2	—	—
expr1 < expr2	expr1 ≥ expr2	expr1 > expr2	expr1 = expr2
expr1 ≤ expr2	expr1 > expr2	—	—
expr1 = expr2	expr1 ≠ expr2	expr1 > expr2	expr1 < expr2
expr1 ≠ expr2	expr1 = expr2	—	—
True	False	—	—
False	True	—	—

and B conflict with D , then C conflicts with D , thus showing the correctness of the constraint splitting strategy.

$$\begin{aligned}
C &\Rightarrow A \vee B \\
\Leftrightarrow \neg C \vee (A \vee B) &\quad \text{—implication} \\
\Leftrightarrow (A \vee B) \vee \neg C &\quad \text{—commutativity} \\
\Leftrightarrow \neg\neg(A \vee B) \vee \neg C &\quad \text{—logical negation} \\
\Leftrightarrow \neg(\neg A \wedge \neg B) \vee \neg C &\quad \text{—DeMorgan's law} \\
\Leftrightarrow \neg A \wedge \neg B \Rightarrow \neg C &\quad \text{—implication}
\end{aligned}$$

By showing that A and B conflict with D , i.e. $\neg A \wedge \neg B \wedge D$, and using the above proof, $\neg A \wedge \neg B \Rightarrow \neg C$, one is sure that C conflicts with D , i.e. $\neg C \wedge D$. The following proves this.

From $\neg A \wedge \neg B \wedge D$, $\neg A \wedge \neg B \Rightarrow \neg C$ Infer $\neg C \wedge D$

- (1) $\neg A \wedge \neg B \wedge D$ —premise
- (2) $\neg A \wedge \neg B$ —property of And, 1
- (3) $\neg A \wedge \neg B \Rightarrow \neg C$ —premise
- (4) $\neg C$ —implication eliminating, 2, 3
- (5) D —property of And, 1
- (6) $\neg C \wedge D$ —property of And, 4, 5

For the *constraint splitting* strategy, the cases shown in Table IV are analysed. Note that, for most of these, A and B are weaker than C , but it is usually easier to decide if A or B conflicts with D . The algorithm for performing constraint splitting is given in Figure 5.

2.4.3. Constants comparison

A third strategy uses a property that is common in constraints created for test case generation. The property is *grounding*, where both constraints have the format $(V \text{ rop } K)$, where V is a variable, rop is a relational operator and K is a constant. In addition, the variables in both constraints must be the same.

Table IV. Constraint splitting cases analysis

Original constraint		New constraint1		New constraint2
$(x + y) > 0$	\Rightarrow	$x > 0$	\vee	$y > 0$
$(x + y) \geq 0$	\Rightarrow	$x \geq 0$	\vee	$y \geq 0$
$(x + y) < 0$	\Rightarrow	$x < 0$	\vee	$y < 0$
$(x + y) \leq 0$	\Rightarrow	$x \leq 0$	\vee	$y \leq 0$
$(x + y) = 0$	\Rightarrow	$x \leq 0$	\vee	$y \leq 0$
$(x + y) \neq 0$	\Rightarrow	$x \neq -y$		
$(x - y) > 0$	\Rightarrow	$x > 0$	\vee	$y < 0$
$(x - y) \geq 0$	\Rightarrow	$x \geq 0$	\vee	$y \leq 0$
$(x - y) < 0$	\Rightarrow	$x < 0$	\vee	$y > 0$
$(x - y) \leq 0$	\Rightarrow	$x \leq 0$	\vee	$y \geq 0$
$(x - y) = 0$	\Rightarrow	$x \leq 0$	\vee	$y \geq 0$
$(x - y) \neq 0$	\Rightarrow	$x \neq y$		
$(x \times y) > 0$	\Rightarrow	$x > 0 \wedge y > 0$	\vee	$x < 0 \wedge y < 0$
$(x \times y) \geq 0$	\Rightarrow	$x \geq 0 \wedge y \geq 0$	\vee	$x \leq 0 \wedge y \leq 0$
$(x \times y) < 0$	\Rightarrow	$x > 0 \wedge y < 0$	\vee	$x < 0 \wedge y > 0$
$(x \times y) \leq 0$	\Rightarrow	$x \geq 0 \wedge y \leq 0$	\vee	$x \leq 0 \wedge y \geq 0$
$(x \times y) = 0$	\Rightarrow	$x = 0$	\vee	$y = 0$
$(x \times y) \neq 0$	\Rightarrow	$x \neq 0 \wedge y \neq 0$		
$(x \div y) > 0$	\Rightarrow	$x > 0 \wedge y > 0$	\vee	$x < 0 \wedge y < 0$
$(x \div y) \geq 0$	\Rightarrow	$x \geq 0 \wedge y > 0$	\vee	$x \leq 0 \wedge y < 0$
$(x \div y) < 0$	\Rightarrow	$x > 0 \wedge y < 0$	\vee	$x < 0 \wedge y > 0$
$(x \div y) \leq 0$	\Rightarrow	$x \geq 0 \wedge y < 0$	\vee	$x \leq 0 \wedge y > 0$
$(x \div y) = 0$	\Rightarrow	$x = 0$		
$(x \div y) \neq 0$	\Rightarrow	$x \neq 0$		

```

algorithm:      SplitConstraint (NecConst, PEConst)
precondition:  NecConst and PEConst are properly initialized constraints.
postcondition: Returns conflict if NecConst and PEConst conflict, no-conflict otherwise.

BEGIN
  -- V1 and V2 are variables, K is a constant, aop is an arithmetic operator,
  -- and rop is a relational operator.
  IF (the format of NecConst is not ((V1 aop V2) rop K))
    RETURN no-conflict
  ELSE
    -- Use Table 4 (Constraint Splitting table) to split NecConst.
    A = NewConstraint1(NecConst)
    B = NewConstraint2(NecConst)
  END IF
  IF (Negation (A, PEConst)) AND (Negation (B, PEConst))
    RETURN conflict
  ELSE
    IF (CompareConstraints (A, PEConst)) AND (CompareConstraints (B, PEConst))
      RETURN conflict
    ELSE
      RETURN no-conflict
    END IF
  END IF
END SplitConstraint

```

Figure 5. The SplitConstraint algorithm—splits a necessity constraint into two parts and decides if they both conflict with the PE.

Let A be the constraint $(X \text{ rop1 } K1)$, and B be the constraint $(X \text{ rop2 } K2)$. By evaluating the two constants and relational operators, it can often be decided whether A conflicts with B . This strategy is called *constants comparison*. Table V shows the cases analysed for the *constants comparison* strategy. The first and the second columns are the two given constraints. The third column is a predicate on constants $k1$ and $k2$, which is used to decide whether the two constraints conflict. Note that it does not always have such a predicate available. The last column is the conclusion of whether the two constraints conflict. The word ‘pred’ stands for ‘the predicate holds’; the predicate is in column 3. The letter ‘T’ stands for ‘True’, which means the two given constraints conflict; the letter ‘F’ stands for ‘False’, which means the two given constraints do not conflict.

Given the two constraints $(x > 1)$ and $(x < 0)$, if the negation strategy is used, the first constraint $(x > 1)$ can be partially negated or negated to be $(x < 1)$ or $(x \leq 1)$, but neither $(x < 1)$ nor $(x \leq 1)$ is syntactically equal to $(x < 0)$, so it cannot be determined from this that they conflict. Yet constants comparison can be used to show that they conflict.

To expand the use of constants comparison, if a constraint has the format $(V \text{ aop } K1) \text{ rop } K2$, it is rewritten as $V \text{ rop } (K2 \overline{\text{aop}} K1)$, such that $(V \text{ aop } K1) \text{ rop } K2 \Leftrightarrow V \text{ rop } (K2 \overline{\text{aop}} K1)$, where V is a variable, aop is an arithmetic operator, $\overline{\text{aop}}$ is the mathematical inverse operation of aop , rop is a relational operator, and $K1$ and $K2$ are constants. The algorithm for performing constants comparison is given in Figure 6.

Table V. Constants comparison cases analysis

Constraint A	Constraint B	Predicate (pred)	Conclusion
$x > k1$	$x > k2$	—	F
$x > k1$	$x \geq k2$	—	F
$x > k1$	$x < k2$	$k1 \geq k2 - 1$	if pred T, else F
$x > k1$	$x \leq k2$	$k1 \geq k2$	if pred T, else F
$x > k1$	$x = k2$	$k1 \geq k2$	if pred T, else F
$x > k1$	$x \neq k2$	—	F
$x \geq k1$	$x > k2$	—	F
$x \geq k1$	$x \geq k2$	—	F
$x \geq k1$	$x < k2$	$k1 \geq k2$	if pred T, else F
$x \geq k1$	$x \leq k2$	$k1 > k2$	if pred T, else F
$x \geq k1$	$x = k2$	$k1 > k2$	if pred T, else F
$x \geq k1$	$x \neq k2$	—	F
$x < k1$	$x > k2$	$k1 \leq k2 + 1$	if pred T, else F
$x < k1$	$x \geq k2$	$k1 \leq k2$	if pred T, else F
$x < k1$	$x < k2$	—	F
$x < k1$	$x \leq k2$	—	F
$x < k1$	$x = k2$	$k1 \leq k2$	if pred T, else F
$x < k1$	$x \neq k2$	—	F
$x \leq k1$	$x > k2$	$k1 \leq k2$	if pred T, else F
$x \leq k1$	$x \geq k2$	$k1 < k2$	if pred T, else F
$x \leq k1$	$x < k2$	—	F
$x \leq k1$	$x \leq k2$	—	F
$x \leq k1$	$x = k2$	$k1 < k2$	if pred T, else F
$x \leq k1$	$x \neq k2$	—	F
$x = k1$	$x > k2$	$k1 \leq k2$	if pred T, else F
$x = k1$	$x \geq k2$	$k1 < k2$	if pred T, else F
$x = k1$	$x < k2$	$k1 \geq k2$	if pred T, else F
$x = k1$	$x \leq k2$	$k1 > k2$	if pred T, else F
$x = k1$	$x = k2$	$k1 \neq k2$	if pred T, else F
$x = k1$	$x \neq k2$	$k1 = k2$	if pred T, else F
$x \neq k1$	$x > k2$	—	F
$x \neq k1$	$x \geq k2$	—	F
$x \neq k1$	$x < k2$	—	F
$x \neq k1$	$x \leq k2$	—	F
$x \neq k1$	$x = k2$	$k1 = k2$	if pred T, else F
$x \neq k1$	$x \neq k2$	—	F

3. EMPIRICAL EVALUATION

The mathematical basis for these strategies allows it to be proved that they can successfully detect at least some equivalent mutants. But this leaves the question of how well they will work—how many equivalent mutants they will detect. This section describes a proof-of-concept tool that was built to demonstrate that the technique could be applied in practice, and then describes some empirical results from using the tool.

```

algorithm:    CompareConstants (A, B)
precondition: A and B are properly initialized constraints.
postcondition: Returns conflict if A and B conflict, no-conflict otherwise.

BEGIN
  -- V is a variable, K, K1, and K2 are constants,
  -- rop is a relational operator, and aop is an arithmetic operator.
  IF (the format of A is (V rop K))
    keep the format the same
  ELSE IF (the format of A is (K rop V))
    modify the format to (V  $\overline{rop}$  K)
  ELSE IF (the format of A is ((V aop K1) rop K2))
    modify the format to (V rop (K2  $\overline{aop}$  K1))
  ELSE IF (the format of A is (K1 rop (V aop K2)))
    modify the format to (V  $\overline{rop}$  (K1  $\overline{aop}$  K2))
  ELSE
    RETURN no-conflict
  END IF
  IF (the format of B is (V rop K))
    keep the format the same
  ELSE IF (the format of B is (K rop V))
    modify the format to (V  $\overline{rop}$  K)
  ELSE IF (the format of B is ((V aop K1) rop K2))
    modify the format to (V rop (K2  $\overline{aop}$  K1))
  ELSE IF (the format of B is (K1 rop (V aop K2)))
    modify the format to (V  $\overline{rop}$  (K1  $\overline{aop}$  K2))
  ELSE
    RETURN no-conflict
  END IF
  IF (V in A and B are not the same)
    RETURN no-conflict
  END IF
  -- Use Table 5 (Constants Comparison table) to check A and B.
  IF (ConstantsComparison(A, B) == True)
    RETURN conflict
  ELSE
    RETURN no-conflict
  END IF
END CompareConstants

```

Figure 6. The CompareConstants algorithm—compares two constraints using the principle of grounding to see if they conflict.

3.1. A constraint-based equivalence detection tool

The proof-of-concept tool, Equivalencer, is integrated with Godzilla, a test data generator within the Mothra mutation tool set. Although the CBT technique is language independent, Mothra works with Fortran 77 programs. Equivalencer was implemented in the programming language C on a Sun Sparc classic workstation running the SunOS 4.1.3 operating system. Equivalencer contains more than 2000 lines of executable source code and also uses several Mothra and Godzilla library packages.

Figure 7 graphically shows the general execution flow of Equivalencer. For each mutant created for a program, each of the three equivalence detection strategies described in Section 2.4 is applied in turn.

This section first describes assertion constraints, which affect the design and implementation, and how they are inserted into a program under test. Then the architectural design for the tool Equivalencer is given.

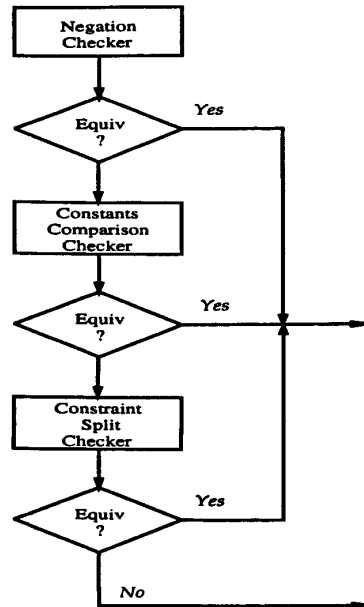


Figure 7. Infeasible constraint checker.

3.1.1. Inserting assertion constraints

In Equivalencer, assertion constraints are used to help detect equivalent mutants. Assertions are constraints that a user inserts into a test program to restrict the input domain of some variables manually. They could be preconditions to the program, predicates on specific statements, or predicates that apply to an entire function or program.

There are two kinds of variables in a program or a function: parameter variables and internal variables. The assertions on the parameter variables can usually only be derived by hand (semantically). These are often part of the specifications, i.e. preconditions. The assertions on the internal variables can sometimes be derived automatically (syntactically). Three kinds of assertions in the programs are used in this paper.

- (1) Assertions on parameter variables: these assertions usually encode preconditions.
- (2) Assertions on internal variables that could be derived automatically: there are known techniques for doing this, such as slicing (Weiser, 1984) and control flow analysis (Fischer and Leblanc, 1988). Because the focus of this research was on detecting equivalent mutants, these assertions were generated by hand and inserted as assertion constraints.
- (3) Assertions on internal variables that could not easily be derived automatically.

When Godzilla generates a constraint that contains arrays, it takes a conservative approach that does not provide array index expressions associated with array references. For example, a constraint system such as $A(i) \geq 0 \wedge A(j) < 0$ will be generated as $A() \geq 0 \wedge A() < 0$. Depending on the negation strategy, this constraint system could be recognized as having a contradiction in it, which is incorrect. To avoid this, Equivalencer

skips checking constraints with arrays, except in the case described below. This is a severe limitation on the experimental proof-of-concept tool that should be addressed in a practical implementation.

Constraints with arrays are used to recognize contradiction only if this constraint is an assertion constraint. This is because it is assured that an assertion constraint involving an array will apply to every element in the array. For example, an assertion constraint is $A() \geq 0$, which means every element in array $A()$ is greater than or equal to 0. If there is a necessity constraint, say $A() < 0$, it is definite that this assertion conflicts with this necessity constraint. To take advantage of this, a routine was added to check whether assertion constraints conflict with other constraints, such as necessity and path expression constraints. In this routine, Equivalencer works on array constraints. This routine is referred to as *array-extension*.

3.1.2. Architectural design

Figure 8 shows the execution flow of Equivalencer. The first step is the initialization. Equivalencer opens all the files that are needed and brings the data into memory. In the second step, Equivalencer consults the *Failure Information* offered by Godzilla. In simple cases, the test data generator is able to recognize equivalent mutants. If the information says the mutant being checked is equivalent, it outputs an equivalent message and exits, otherwise Equivalencer goes to the next step. Next Equivalencer gets the *Path Expression* (**pe**) (from Godzilla) and *Assertion Constraints* (**assertion**), and combines them to form

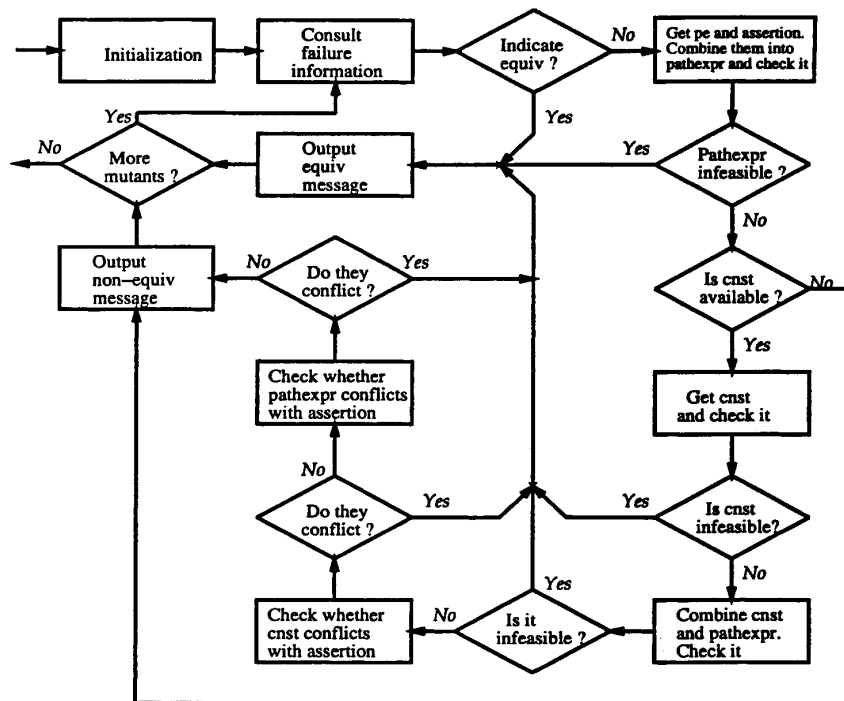


Figure 8. Equivalencer execution diagram.

a completed path expression (**pathexpr**). Then it checks for infeasible constraints in **pathexpr** by using the negation, constants comparison and constraint splitting strategies. If the constraints are infeasible, it outputs the equivalent message, otherwise it goes to the next step. Equivalencer gets the *Necessity Constraint* (**cnst**) if it is available, otherwise it outputs a message indicating that fact. If **cnst** is available, then Equivalencer checks whether **cnst** is infeasible by using the negation, constraint splitting and constants comparison strategies. If it is infeasible, Equivalencer outputs an equivalence message. If it is not, Equivalencer goes to the next step. After combining **cnst** and **pathexpr**, Equivalencer checks for a contradiction in the combination by using the three strategies. If a contradiction is found, it outputs an equivalence message. If it is not found, the array-extension checking routine is applied to assertion constraints against **cnst** and against **pathexpr**. If a conflict is found, an equivalence message is output. If none of these steps find that the mutant is equivalent, Equivalencer outputs a message indicating that fact. After outputting the messages, Equivalencer checks whether there are more mutants. If not, it is finished and it terminates.

As a proof-of-concept prototype, efficiency was not a high priority when Equivalencer was built. The tool compares each **pathexpr** from each statement with each **cnst** that is defined at that statement. So the running time is constant in the number of necessity constraints, which is on the order of the number of mutants. As shown in a previous paper (Offutt *et al.*, 1996), the number of mutants is $O(S^2)$, where S is the number of statements in the program. Because the focus of this study was a proof-of-concept, execution time information was not precisely monitored, but when compared with hand analysis of mutants, the execution time of this tool was trivial.

3.2. Empirical results

A controlled study was carried out using Equivalencer to determine equivalent mutants in 11 Fortran 77 programs that cover a range of applications. These programs range in size from about 11 to 30 executable statements and have from about 180 to 3000 mutants. Because the speed of these equivalence procedures is infinitesimal when compared to the speed of mutation testing, the size of the programs is relatively unimportant. However, each program had to be analysed by hand to determine the true number of equivalent mutants. Equivalencer's effectiveness has been compared based on the percentage of equivalent mutants that it detected. The equivalent mutants for these programs have been analysed by several researchers.

3.2.1. Equivalence detection

The procedure was relatively straightforward.

- Step 1.* The lists of equivalent mutants were gathered from previous research.
- Step 2.* For each program, Mothra was used to generate all mutants, and Godzilla was used to generate the path expression (reachability) constraints and necessity constraints.
- Step 3.* Assertion constraints were generated for programs that included hand-inserted assertions. These assertions were created by the first author.
- Step 4.* Equivalencer was run to see how many equivalent mutants could be detected.

Step 5. The mutants that were marked equivalent were compared with the hand-generated lists.

The data are displayed in Table VI. The number of executable statements, mutants, equivalent mutants, equivalent mutants detected and the percentage of equivalent mutants detected for each program is shown. The bottom line shows the total number of statements, mutants and equivalent mutants in the programs, and the average number detected. The average detected is calculated from the totals, not the percentage detected for each program (i.e. it is not a weighted average).

3.2.2. Discussion

Although Equivalencer was able to detect roughly half of the equivalent mutants on average, the percentage found for each program varied greatly. This is due to a variety of factors, perhaps the most important being limitations of the tools on which Equivalencer is based.

Godzilla treats arrays as single data items, and references to array elements are treated as references to the entire array. For this reason, many constraints involving array references could not be used to determine equivalence. An example of this is the statement $A(5) = 0$. From this definition, the fact that the fifth element of A is set to zero can be determined, and a later mutant that could only be killed if $A(5)$ was not 0 is equivalent. If elements of an array could be treated as individual data items, the constraint-based technique could be used to detect more equivalent mutants. This problem is particularly important for the Cal program, for which only 15% of the equivalent mutants were detected. Recent results with a 'dynamic-domain reduction' approach to generating tests (Offutt *et al.*, 1994) indicate that arrays can be handled within constraints as individual elements—perhaps allowing many more equivalent mutants to be detected.

Godzilla associates every variable in each constraint with a statement number. Usually, the statement number for a variable in a constraint is the number of the statement associated with the constraint. For example, a constraint $A < 0$ in statement 5 will be expressed as $A.5 < 0$. But if a variable's value has not changed from the previous

Table VI. Equivalent mutants detected

Program	Statements	Mutants	Equivalent	Equiv. detected	Percent. detected
Bsearch	20	299	27	19	70.37
Bub	11	338	35	24	68.57
Cal	29	3010	236	37	15.67
Euclid	11	196	24	18	75.00
Find	28	1022	75	63	84.00
Insert	14	460	46	32	69.57
Mid	16	183	13	3	23.08
Pat	17	513	61	29	47.54
Quad	10	359	31	4	12.90
Trityp	28	951	109	80	73.39
Warshall	11	305	35	22	62.86
Total/Avg	185	7636	695	331	47.63

statement, the same statement number for the variable might be used. For example, using the constraint in the example above, if A has not changed from statement 4, it might be shown as $A.4 < 0$ on statement 5. This affects both the handling of assertions and symbolic evaluation.

Assertions in Godzilla are handled in one of two ways. *Global assertions* are applied to all statements in the program, and *local assertions* are applied only to the statement immediately following the assertion. Global assertions are used as constraints that carry the same statement number (usually 1) to each statement. Unfortunately, Godzilla is limited to treating all assertions of a program as either global or local, thus it was necessary to gather data using global assertions separately from the data using local assertions. The assertions increased the overall average percentage detected by almost 5%.

Godzilla's symbolic evaluator presented a more formidable problem. The symbolic evaluator rewrites constraints to be solely in terms of the input parameters by symbolically evaluating the program. Also, Godzilla rewrites the statement numbers in the constraints. Initially, the statement number for each variable is the statement where the constraint appears. Godzilla rewrites a variable's statement number to a previous statement number if the variable has not been assigned a new value from the previous statement. This is called *constraints propagation* because the statement numbers are 'propagated' through the constraints. Constraints propagation, especially in path expression constraints, increases Equivalencer's detection ability. Unfortunately, Godzilla has a design flaw that causes it to work incorrectly on constraints involving loops. If a constraint rewritten by Godzilla is in a loop, it is correct for the first iteration, but usually not correct in the second and subsequent iterations. This caused no problems for the original purpose of test data generation, but made equivalent mutant detection more difficult. Since Equivalencer detects equivalent mutants using constraints, Equivalencer will detect incorrect equivalent mutants if the constraints are wrong. Unfortunately, fixing the design flaw for Godzilla was impractical, thus the mutants that were incorrectly identified as equivalent were removed by hand.

A previous automatic equivalent mutant detector was presented by Offutt and Craft (1994). It used compiler optimization techniques to identify equivalent mutants of the original program. When run on the same programs, the constraint-based technique detected almost five times more equivalent mutants than the compiler optimization techniques.

A test suite of eleven Fortran 77 programs was used for these studies. It cannot be claimed that these programs represented a statistically valid sample of programs. There is no generally accepted way to choose such a sample of programs, and this paper does not attempt to solve that problem. However, the programs (functions) were taken from the literature and chosen to represent different types of problems to exercise the equivalent mutant detection capabilities in as wide a manner as possible.

3.2.3. Feasible path results

Although the focus in this research was on equivalent mutants, the same techniques can be applied to the problem of detecting infeasible paths. Specifically, a simple corollary to Theorem 1 is that if the reachability condition for a statement is infeasible, the statement is unreachable. Thus, a preliminary evaluation of using the constraint-based technique to recognize unreachable statements was made. To do this, programs were artificially constructed so that some of their statements were not reachable, and the

Equivalencer was used to try to detect the unreachability. This was done by applying the mutation operator that ensures statement coverage; the mutants are killed by a test case if and only if the test case reaches the statement. If Equivalencer found that the mutant is equivalent, that means the statement is unreachable.

Artificially creating programs with unreachable statements implies a severe limitation and potential bias, so these results should not be generalized without further research. On the other hand, the preliminary results in Table VII are encouraging. Of the 9 programs, there were 14 unreachable paths, and the tool was able to find 10 of those. Many of the equivalent mutants that Equivalencer was not able to detect were equivalent because of computation that occurred after the mutated statement was executed; since this is not a problem for unreachable paths, it is not surprising to get better results for unreachable paths than for equivalent mutants.

4. SUGGESTIONS FOR IMPROVEMENT AND FUTURE WORK

There are several improvements that should be used in a practical tool. These are divided into two categories: (1) recommendations for improving the software; and (2) recommendations for improving the technique.

4.1. Improving the software

The Equivalencer tool relied very heavily on the pre-existing Godzilla test data generator. Godzilla implements symbolic evaluation as a separate step from infeasible constraint recognition. Although symbolic evaluation helps Equivalencer detect equivalent mutants by propagating the constraints, it increases the difficulty of detection by throwing away considerable information that Equivalencer needs, such as all references to internal variables. Also, Godzilla does not propagate the assertion constraints. In future systems, the symbolic evaluation should be merged with infeasible constraint recognition. This should allow for more detection of equivalent mutants, as well as being more efficient and flexible.

Table VII. Unreachable statements detected

Program	Unreachable	Detected	Percentage detected
Prog 1	2	1	50.00
Prog 2	1	0	0.00
Prog 3	1	1	100.00
Prog 4	1	1	100.00
Prog 5	1	1	100.00
Prog 6	1	1	100.00
Prog 7	2	2	100.00
Prog 8	3	3	100.00
Prog 9	2	0	0.00
Total/Avg	14	10	71.43

4.2. Improving the detection techniques

Three ways are recommended to improve the techniques. One is another strategy that could recognize infeasible constraints, another is to have better constraints, and the third is to analyse the execution after the mutated statement.

Three strategies for recognizing infeasible constraints were presented in Section 2.4. Another potential strategy is the following. Assume a constraint X . Suppose other constraints can be found, say $C1, C2, \dots, Cn$, such that $X \wedge C1 \wedge C2 \wedge \dots \wedge Cn \Rightarrow Y$, and $C1, C2, \dots, Cn$ are TRUE. If it can be proved that $\neg Y$ holds, then it can be said that $\neg X$ holds, which means that X is an infeasible constraint. An analysis of the programs studied uncovered only one equivalent mutant that could be detected using this strategy. Thus, this strategy was not implemented, but it is possible that it could detect more equivalent mutants in other programs.

One limitation has to do with array constraints generation. Currently, Godzilla generates array constraints without indexes, which is a safe approach when indexes are variables. But programs often contain array constraints with indexes that are constants. Analysis of the program Cal shows that if array constraints with constant-indexes could be used, Equivalencer would detect 69 more equivalent mutants, which would increase the detection percentage from 15.67% to 44.92% in the program Cal.

Another thought on improving the constraints is that of using humans to help with difficult constraints. One can imagine a system that interacts with the tester to get help with difficult constraints. Although this would require more work from the tester than automatic detection, it is still less than requiring hand recognition of equivalent mutants. That is, it is easier for a human to recognize infeasible constraints than equivalent mutants.

This research only analysed the execution up to and including the mutated statement. If execution after the mutated statement could be analysed, then many more equivalent mutants could be recognized. Analysis of the execution after the mutated statements requires analysis of sufficiency conditions, which are not generated by Godzilla, but no reason is seen why this could not be done.

5. CONCLUSIONS

This paper presents a partial solution to the problem of equivalent mutant detection, and to the feasible path problem, and shows the specific relationship between them. The solution is proposed, specific algorithms are presented, and a proof-of-concept experimental tool is described. The results show that the approach is an effective partial solution to this problem. It is also shown that this technique can be applied to the feasible path problem, possibly giving better results than with equivalent mutants. This general technique is generalizable to all instances of what is defined as the feasible test problem, and can thus be used to support branch coverage techniques and data flow testing.

By using the CBT technique, Equivalencer is able to detect a significant percentage for most programs automatically, although it is not possible to detect all equivalent mutants. In the experiments, the detection rate is over 60% for 7 of the 11 programs and the average detection rate over all programs is over 45% (see Table VI in Section 3.2). With appropriate extensions, the detection rates could be even higher. Also the detection time is reasonably fast, even in this implementation, which was not optimized for speed. Compared with running every test case against those equivalent mutants, the time saving

is large. Compared with detecting those equivalent mutants by hand, the time saving is significant.

Previous research (Offutt and Craft, 1994) reported results from using compiler optimization techniques for this problem. The ideas in this paper are derived from that work, and it is clear that CBT is a far more effective approach.

This research is part of a long-term project to provide practical, powerful automated test environments to testers, so that highly assured software can be produced at reasonable cost. A system is envisaged that provides almost complete automation to the tester. This type of system would allow a programmer to submit a software module, and after a few minutes of computation, respond with a set of test cases that are assured of providing the software with a very effective test, and a set of outputs that can be examined to find failures in the software. Furthermore, these input-output pairs can be used as a basis for debugging when failures are found.

Acknowledgement

Partially supported by the National Science Foundation under Grant CCR-93-11967.

References

- Acree, A. T. (1980) 'On mutation', Ph.D. thesis, Georgia Institute of Technology, Atlanta, Georgia, U.S.A.
- Budd, T. A. and Angluin, D. (1982) 'Two notions of correctness and their relation to testing', *Acta Informatica*, **18**(1), 31–45.
- Baldwin, D. and Sayward, F. (1979) 'Heuristics for determining equivalence of program mutations', Research report 276, Department of Computer Science, Yale University, New Haven, Connecticut, U.S.A.
- Budd, T. A. (1980) 'Mutation analysis of program test data', Ph.D. thesis, Yale University, New Haven, Connecticut, U.S.A.
- 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', *Proceedings of the Second Workshop on Software Testing, Verification and Analysis*, Banff, Alberta, Canada, July, IEEE Computer Society Press, Los Alamitos, California, U.S.A., pp. 142–151.
- DeMillo, R. A., Lipton, R. J. and Sayward, F. G. (1978) 'Hints on test data selection: Help for the practicing programmer', *IEEE Computer*, **11**(4), 34–41.
- DeMillo, R. A. and Offutt, A. J. (1991) 'Constraint-based automatic test data generation', *IEEE Transactions on Software Engineering*, **17**(9), 900–910.
- Fischer, C. N. and Leblanc, R. J. (1988) *Crafting a Compiler*, Benjamin/Cummings, Menlo Park, California, U.S.A.
- Frankl, P. G. and Weyuker, E. J. (1988) 'An applicable family of data flow testing criteria', *IEEE Transactions on Software Engineering*, **14**(10), 1483–1498.
- Goldberg, A., Wang, T. C. and Zimmerman, D. (1994) 'Applications of feasible path analysis to program testing', *Proceedings of the 1994 International Symposium on Software Testing and Analysis*, Seattle, Washington, U.S.A., August, ACM Press, New York, U.S.A., pp. 80–94.
- Hamlet, R. G. (1977) 'Testing programs with the aid of a compiler', *IEEE Transactions on Software Engineering*, **3**(4), 279–290.
- Howden, W. E. (1982) 'Weak mutation testing and completeness of test sets', *IEEE Transactions on Software Engineering*, **8**(4), 371–379.
- Jasper, R., Brennan, M., Williamson, K., Currier, B. and Zimmerman, D. (1994) 'Test data generation and feasible path analysis', *Proceedings of the 1994 International Symposium on Software Testing and Analysis*, Seattle, Washington, U.S.A., August, ACM Press, New York, U.S.A., pp. 95–107.
- King, J. C. (1976) 'Symbolic execution and program testing', *Communications of the ACM*, **19**(7), 385–394.

-
- Morell, L. J. (1990) 'A theory of fault-based testing', *IEEE Transactions on Software Engineering*, **16**(8), 844–857.
- Myers, G. (1979) *The Art of Software Testing*, Wiley, New York, U.S.A.
- Offutt, A. J. and Craft, W. M. (1994) 'Using compiler optimization techniques to detect equivalent mutants', *Software Testing, Verification and Reliability*, **4**(3), 131–154.
- Offutt, A. J. (1988) 'Automatic test data generation', Ph.D. thesis, Georgia Institute of Technology, Atlanta, Georgia, U.S.A., Technical Report GIT-ICS 88/28.
- Offutt, A. J. (1991) 'An integrated automatic test data generation system', *Journal of Systems Integration*, **1**(3), 391–409.
- Offutt, J., Jin, Z. and Pan, J. (1994) 'The dynamic domain reduction approach for test data generation: design and algorithms', Department of Information and Software Systems Engineering, George Mason University, Fairfax, Virginia, U.S.A., September, Technical Report ISSE-TR-94-110.
- Offutt, A. J., Lee, A., Rothermel, G., Untch, R. and Zapf, C. (1996) 'An experimental determination of sufficient mutation operators', *ACM Transactions on Software Engineering Methodology*, **5**(2), 99–118.
- Untch, R., Offutt, A. J. and Harrold, M. J. (1993) 'Mutation analysis using program schemata', *Proceedings of the 1993 International Symposium on Software Testing and Analysis*, Cambridge, Massachusetts, U.S.A., June, ACM Press, New York, U.S.A., pp. 139–148.
- Voas, J. M., Morell, L. and Miller, K. W. (1991) 'Predicting where faults can hide from testing', *IEEE Software*, **8**(2), 41–58.
- Weiser, M. (1984) 'Program slicing', *IEEE Transactions on Software Engineering*, **10**(4), 352–357.