# Mutation Testing

Stuart Anderson

School of informatics

**School of informatics**

# Overview

*Mutation testing is a structural testing method, i.e. we use the structure of the code to guide the test process.*

We cover the following aspects of Mutation Testing:

- What is a mutation?

- What is mutation testing?

- When should we use mutation testing?

- Mutations

- Examples

- Mutation testing tools

# What is a mutation?

- A mutation is a small change in a program.

- Such small changes are intended to model low level defects that arise in the process of coding systems.

- Ideally mutations should model low-level defect creation.

# What is Mutation Testing?

- Mutation testing is a structural testing method aimed at assessing/improving the **adequacy** of test suites, and estimating the number of faults present in systems under test.
- The process, given program $P$ and test suite $T$, is as follows:
  - We systematically apply mutations to the program $P$ to obtain a sequence $P_1, P_2, \dots P_n$ of mutants of $P$. Each mutant is derived by applying a single mutation operation to $P$.
  - We run the test suite $T$ on each of the mutants, $T$ is said to **kill** mutant $P_j$ if it detects an error.
  - If we kill $k$ out of $n$ mutants the adequacy of $T$ is measured by the quotient $k/n$. $T$ is **mutation adequate** if $k = n$.
- One of the benefits of the approach is that it can be almost completely automated.

# When should we use mutation testing?

- Structural test suites are directed at identifying defects in the code. One goal of mutation testing is to assess or improve the efficacy of test suites in discovering defects.

- When we are carrying out structural testing we are worried about defects remaining in the code. Often we are keen to measure the **Residual Defect Density (RDD)** in the program P under test.

- The Residual Defect Density is usually measured in defects per thousand lines of code.

- Advocates of mutation testing argue that it can provide us with an estimate of the RDD of a program P that has satisfied all the tests in a test suite T.

# Using Mutation Testing to Estimate the RDD

We want to estimate the RDD of Program $P$ given that it has satisfied all the tests in test suite $T$. We follow the procedure:

- Suppose we have an estimate $r$ of the RDD of programs produced by our development process before they are subject to test (this could be gathered using production data and field experience, or it could be based on the number of faults our tests have already detected).
- Generate $n$ mutants of the program $P$.
- Test each mutant with the test suite $T$.
- Find the number, $k$, of mutants that are killed by $T$. To yield a non-zero RDD we need to test enough mutants to ensure that $0 < k < n$.
- Use $r . (n-k)/k$ as the estimate for the RDD of the tested program.
- $k/n$ is a measure of the adequacy of $T$ in finding defects in $P$.

# Assumptions

The validity of this rests on many assumptions:

- That mutations are a good model for defects.

- That defects are usually independent

- That the construction of $T$ is not influenced by knowledge of the mutation process (i.e. we do not use knowledge of the mutation process to build tests that are better at finding defects generated by mutations than normal defects).

- If we are interested in making confident estimates of very low RDDs we will need very large numbers of mutants.

- For example, if our development process left us with 10 defects per kLoc before test and we want to be confident our RDD after test is lower that 0.1 per kLoC then we need to test many mutants to be confident of such an RDD estimate.

# An Approach to Mutation

- Ideally we need systematically to apply mutations to the program under test. This would involve some criterion of applicability.

- Usually we consider mutation operators in the form of rules that match a context and create some systematic mutation of the context to create a mutant.

- The simple approach to coverage is to consider all possible mutants but this may create a very large number of mutants (in the case of estimating RDDs the value and confidence required of the estimated RDD would control the number of mutants to be generated).

- Mutation testing is best supported by tools because of the potentially very large numbers of mutations to be generated during testing.

# Kinds of Mutation

- **Value Mutations:** these mutations involve changing the values of constants or parameters (by adding or subtracting values etc), e.g. loop bounds – being one out on the start or finish is a very common error.

- **Decision Mutations:** this involves modifying conditions to reflect potential slips and errors in the coding of conditions in programs, e.g. a typical mutation might be replacing a $>$ by a $<$ in a comparison.

- **Statement Mutations:** these might involve deleting certain lines to reflect omissions in coding or swapping the order of lines of code. There are other operations, e.g. changing operations in arithmetic expressions. A typical omission might be to omit the increment on some variable in a while loop.

A wide range of mutation operators is possible...

School of **informatics**

# Offutt's Mutations for Inter-Class Testing

| Language Feature | Operator | Description |
|---|---|---|
| Access Control | AMC | Access modifier change |
| Inheritance | IHD | Hiding variable deletion |
| | IHI | Hiding variable insertion |
| | IOD | Overriding method deletion |
| | IOP | overriding method calling position change |
| | IOR | Overriding method rename |
| | ISK | *super* keyword deletion |
| | IPC | Explicit call of a parent's constructor deletion |
| Polymorphism | PNC | *new* method call with child class type |
| | PMD | Instance variable declaration with parent class type |
| | PPD | Parameter variable declaration with child class type |
| | PRV | Reference assignment with other comparable type |
| Overloading | OMR | Overloading method contents change |
| | OMD | Overloading method deletion |
| | OAO | Argument order change |
| | OAN | Argument number change |
| Java-Specific Features | JTD | *this* keyword deletion |
| | JSC | *static* modifier change |
| | JID | Member variable initialization deletion |
| | JDC | Java-supported default constructor creation |
| Common Programming Mistakes | EOA | Reference assignment and content assignment replacement |
| | EOC | Reference comparison and content comparison replacement |
| | EAM | Accessor method change |
| | EMM | Modifier method change |

School of
**informatics**

# Value Mutation

- Here we attempt to change values to reflect errors in reasoning about programs.

- Typical examples are:

  - Changing values to one larger or smaller (or similar for real numbers).
  - Swapping values in initialisations.

- The commonest approach is to change constants by one in an attempt to generate a one-off error (particularly common in accessing arrays).

- Coverage criterion: Here we might want to perturb all constants in the program or unit at least once or twice.

Example 11

# Value Mutation

```
public int Segment(int t[], int l, int u){
    // Assumes t is in ascending order, and l<u,
    // counts the length of the segment
    // of t with each element l<t[i]<u
    int k = 0;

    for(int i=0; i<t.length && t[i]<u; i++){
        if(t[i]  ){
            k++;
        }
    }
    return(k);
}
```

Mutating to k=1 causes miscounting

Here we might mutate the code to read i=1, a test that would kill this would have t length 1 and have l < t[0] < u, then the program would fail to count t[0] and return 0 rather than 1 as a result

# Decision Mutation

- Here again we design the mutations to model failures in reasoning about conditions in programs. As before this is a very limited model of programming error really modelling slips in coding rather than a design error.
- Typical examples are:
  - Modelling "one-off" errors by changing $<$ to $<=$ or vice versa (this is common in checking loop bounds).
  - Modelling confusion about larger and smaller, so changing $>$ to $<$ or vice versa.
  - Getting parenthesisation wrong in logical expressions e.g. mistaking precedence between $\&\&$ and $||$
- Coverage Criterion: We might consider one mutation for each condition in the program. Alternatively we might consider mutating all relational operators (and logical operators e.g. replacing $||$ by $\&\&$ and vice versa)

Example 13

School of **informatics**

# Decision Mutation

```
public int Segment(int t[], int l, int u){
    // Assumes t is in ascending order, and l<u,
    // counts the length of the segment
    // of t with each element l<t[i]<u
    int k = 0;

    for(int i=0; i<t.length && t[i]<u; i++){
        if(t[i]>l){
            k++;
        }
    }
    return(k);
}
```

Mutating to t[i]>u will cause miscounting

We can model "one-off" errors in the loop bound by changing this condition to i<=t.length - provided array bounds are checked exactly this will provoke an error on every execution.

# Statement Mutation

- Here the goal is primarily to model editing slips at the line level – these typically arise when the developer is cutting and pasting code. The result is usually omitted or duplicated code. In general we might consider arbitrary deletions and permutations of the code.

- Typical examples include:
  - Deleting a line of code
  - Duplicating a line of code
  - Permuting the order of statements.

- Coverage Criterion: We might consider applying this procedure to each statement in the program (or all blocks of code up to and including a given small number of lines).

Example                                                                                     15

School of
**informatics**

# Statement Mutation

```
public int Segment(int t[], int l, int u){
    // Assumes t is in ascending order, and l<u,
    // counts the length of the segment
    // of t with each element l<t[i]<u
    int k = 0;

    for(int i=0; i<t.length && t[i]<u; i++){
        if(t[i]>l){
            k++;
        }
    }
    return(k)
}
```

Here we might consider deleting this statement (then count would be zero for all inputs, we might also duplicate this line in which case all counts would be doubled.

# Observations

- Mutations model low level errors in the mechanical production process. Modelling design errors is much harder because they involve large numbers of coordinated changes throughout the program.

- Ensuring test sets satisfy coverage criteria are often enough to ensure they kill mutants (because mutants often do not "make sense" and so provoke a failure if they are ever executed).

- Black-box test sets are poorer at killing mutants – we'd expect this because black-box tests are driven more by the operational profile than by the need to cover statements.

- We could see mutation testing as a way of forcing more diversity on the development of test sets if we use a black-box approach as our primary test development approach.

# Concepts from the literature

- **Syntactic vs semantic size of a mutant** – the size the source change a mutant involves, versus the size of its effect on program behaviour. It has been hypothesised that mutation operators which produce semantically small faults are better (because semantically large faults will be caught by most tests). Justification for elimination of certain types of mutation.

- **Competent programmer hypothesis** – the program under test is "close to" the correct program. So exploring the space of small mutations will lead us to that program.

- **Coupling effect hypothesis** — tests for detecting simpler faults will be sufficient also for detecting more complex faults. So even though many faults are a product of logical errors with wide consequences in the code, small mutants will lead to recognition of these faults.
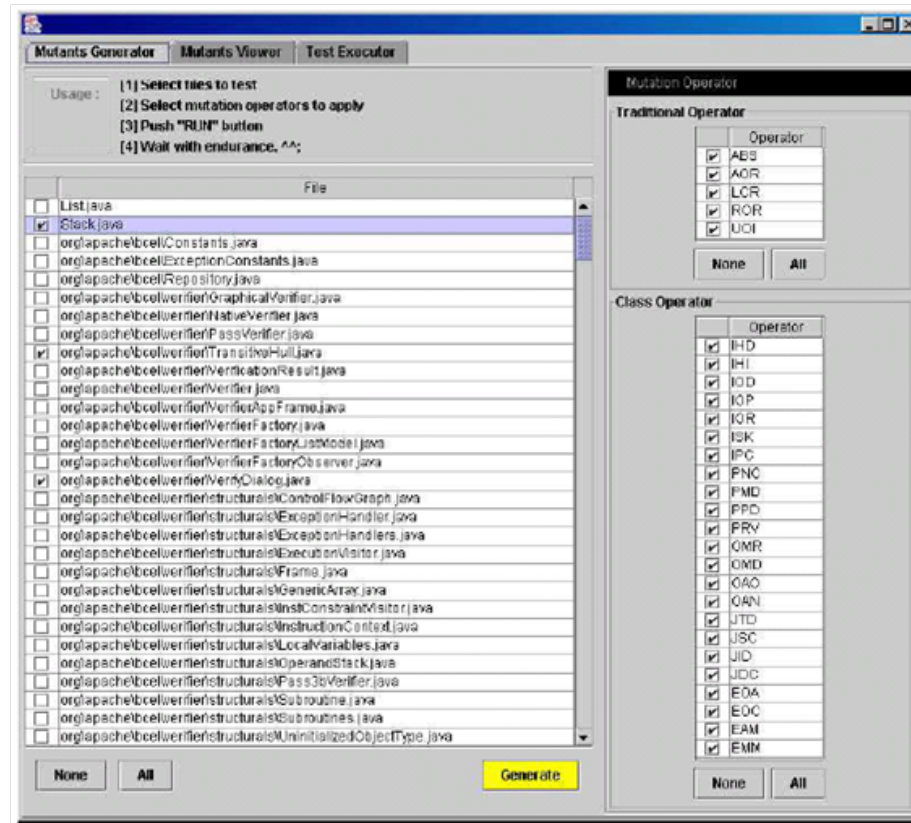
# Mutation Testing Tools

- There is a range of possible mutation tools. Recently Offutt and others have created MuJava, a tool for creating Java mutants.
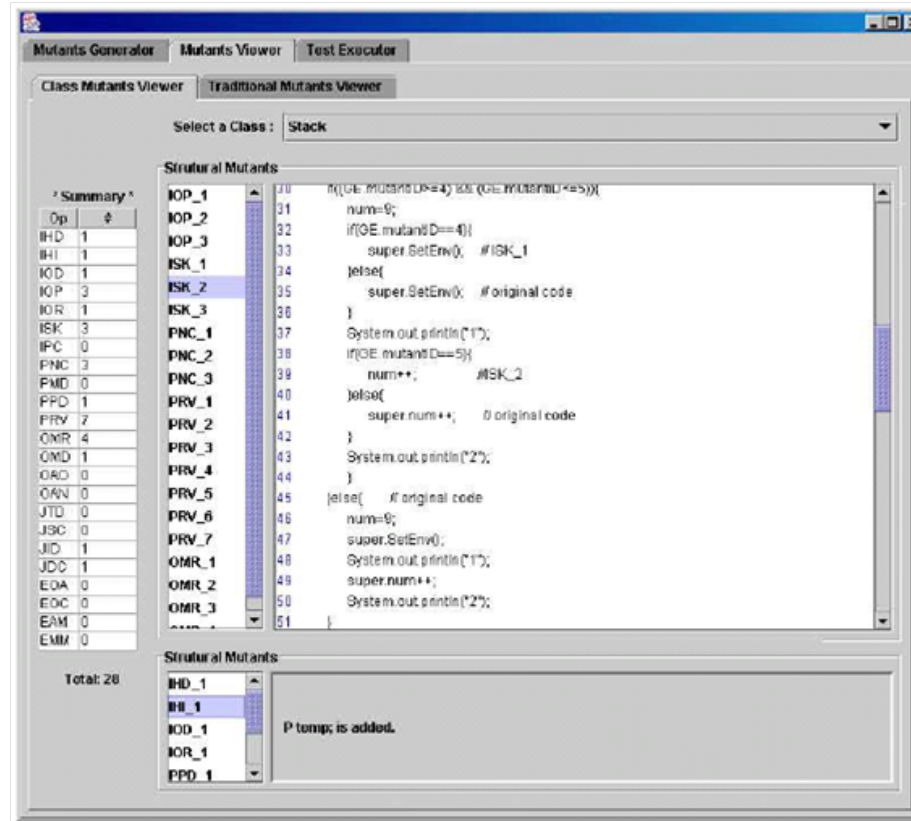
  *MuJava: An Automated Class Mutation System*, Yu-Seung Ma, Jeff Offutt and Yong-Rae Kwon. Journal of Software Testing, Verification and Reliability, 15(2):97-133, June 2005. http://dx.doi.org/10.1002/stvr.v15:2

- Their system is designed specifically to include a range of mutation operators that target OO languages in particular.

- They incorporate an efficient version of generating a "metamutant" that is capable of behaving like all mutants of the program (using Java reflection to instantiate operators at execution time).
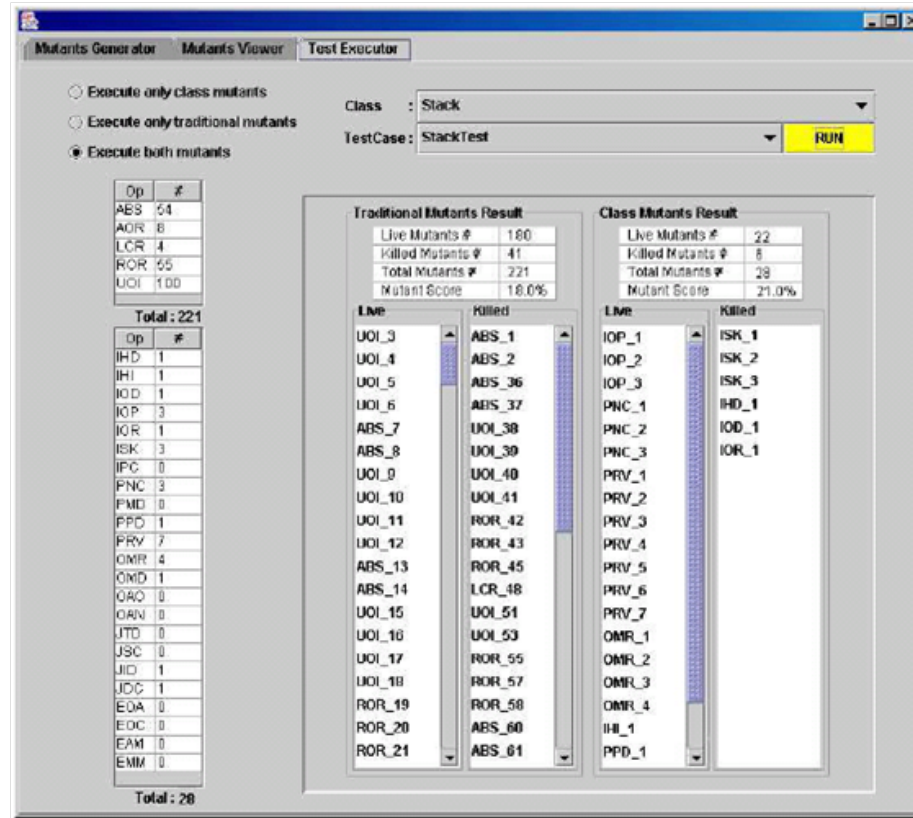
# Mutant Generation Interface

# Mutant Analysis Interface

# Test Execution Interface

School of **informatics**

# Summary

- Mutation testing can be a useful addition to the test process.

- It is laborious and really requires tool assistance if it is to be cost-effective.

- Improving Residual Defect Density estimates requires very large numbers of mutants if we are to have confidence in the results.

- Object Orientation has a wide range of structural and operational mutants that are specific to objects.

- Tools like mu-Java use features of Java to enable the efficient generation and test of mutants.

- Even with efficient techniques execution times can be very slow for large numbers of mutants.

School of **informatics**

# Required Readings

- **Textbook (Pezzè and Young):** Chapter 16, Fault-Based Testing

- MuJava: An Automated Class Mutation System, Yu-Seung Ma, Jeff Offutt and Yong-Rae Kwon. Journal of Software Testing, Verification and Reliability, 15(2):97-133, June 2005. http://dx.doi.org/10.1002/stvr.v15:2

- A. Jefferson Offutt, Ammei Lee, Gregg Rothermel, Roland H. Untch, and Christian Zapf. 1996. An experimental determination of sufficient mutant operators. ACM Trans. Softw. Eng. Methodol. 5, 2 (April 1996), 99-118. http://dx.doi.org/10.1145/227607.227610