

Topic Area: Code Quality Assurance: Content

VL15	• Introduction and Vocabulary
VL16	• Limits of Software Testing
VL17	• Glass-Box Testing
VL18	• Statement-, branch-, term-, coverage
VL19	• Other Approaches
VL20	• Model-based testing
VL21	• Runtime verification
VL22	• Software quality assurance
VL23	• In a larger scope
VL24	• Program Verification
VL25	• partial and total correctness,
VL26	• Proof System PD
VL27	• Review

Content

- Limits of Software Testing
 - Software examination paths
 - Is exhaustive testing feasible?
 - Range- vs. point errors
 - When To Stop Testing?
 - Choosing Test Cases
 - Requirements on test cases
 - The natural habitat of many errors
 - Test Oracle
- Glass-Box Testing
 - Statement coverage
 - Branch and term coverage
 - Conclusions from coverage measures
- Model-Based Testing
- Testing in the Development Process

Testing

Recall: Test Case

Definition. A test case T is a pair (I_n, S_{out}) consisting of

- a description I_n of sets of finite input sequences,
- a description S_{out} of expected outcomes,
- and an interpretation $[\cdot]$ of these descriptions.

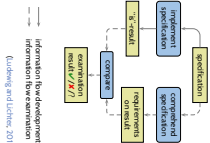
- A test execution π is $(\pi^1, \dots, \pi^n) \in \Sigma_n$ for some $n \in \mathbb{N}_0$ is called
 - successful (or positive) if it discovered an error, i.e. if $\pi \notin [S_{out}]$.
 - unsuccessful (or negative) if it did not discover an error, i.e. if $\pi \in [S_{out}]$.
 - (Alternative: test item S failed to pass test: confusing "test failed")
 - (Alternative: test item S passed test: okay "test passed")

Software Examination (in Particular: Testing)

- In each examination, there are two paths from the specification to results
 - the production path using model, source code, executable, etc.), and
 - the examination path (using requirements specifications).
- A check can only discover errors on **exactly one** of the paths.
- If difference detected: examination result is positive.
- What is not on the paths, is not checked: crucial: **specification** and comparison.

Recall:

	checking production	checking specification
spec	no error	error
no spec	new error	bug in prod



- Limits of Software Testing
 - Software examination paths is exhaustive testing feasible?
 - Range- vs. point errors
- When To Stop Testing?
- Choosing Test Cases
 - Requirements on test cases
 - The natural habitat of many errors
- Test Oracle
- Glass-Box Testing
 - Statement coverage
 - Branch and term coverage
 - Conclusions from coverage measures
- Model-Based Testing
- Testing in the Development Process

The Crux of Software Testing

Why Can't We Show The Absence of Errors (in General)?

Recall:

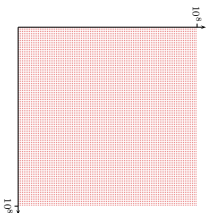
"Software testing can be used to show the presence of bugs, but never to show their absence!"
(E.W. Dijkstra, 1970)

- Consider a simple pocket calculator for adding 8-digit decimals:
- Requirement: If the display shows x , $+$, and y , then after pressing $=$
 - the sum of x and y is displayed if $x + y$ has at most 8 digits
 - otherwise $-E-$ is displayed
 - With 8 digits, both x and y range over $[0, 10^8 - 1]$.
 - Thus there are 10^{16} possible input pairs (x, y) to be considered for exhaustive testing!
 - And if we reset the pocket calculator for each test, we **do not know anything** about problems with sequences of inputs..
(local variables may not be re-initialized properly, for example)



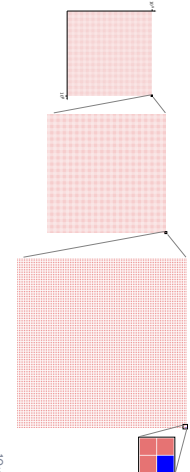
Observation: Software Usually Has Many Inputs

- Example: Simple Pocket Calculator:
 - With ten thousand different test cases (that's a lot!), 9,999,999,999,990,000 of the 10^{16} possible inputs remain **uncovered**.
 - In other words:
 - Only 0.0000000001% of the possible inputs are covered; 99,999,999,999% not touched in diagrams (red: untested; blue: tested)



Observation: Software Usually Has Many Inputs

- Example: Simple Pocket Calculator:
 - With ten thousand different test cases (that's a lot!), 9,999,999,999,990,000 of the 10^{16} possible inputs remain **uncovered**.
 - In other words:
 - Only 0.0000000001% of the possible inputs are covered; 99,999,999,999% not touched in diagrams (red: untested; blue: tested)



Conclusion

- Question 1:
 - If we **cannot consider all** test cases, are there **clever choices** of test cases?

More Observations

- Software is (in general) **not continuous**.
 - Consider a continuous function, e.g. For sufficiently small ϵ -environments of an input, the outputs differ only by a small amount δ .



- For software, adjacent inputs **may yield arbitrarily distinct** output values.

Vocabulary

- **Range error**: multiple "neighbouring" inputs trigger the error.
- **Point error**: an isolated input value triggers the error.



- For Software, we can (in general) **not conclude from some values to others**.

- For example, if a bridge endures a single car of 1000kg, we strongly expect the bridge to hold a single person of 100kg.
- If the pocket calculator is correct for $7245678 \cdot 27$, we can (in general) not expect anything of the other numbers.

12/44

Conclusion Cont'd

- **Question 1:** If we **cannot consider all** test cases, are there **clever choices** of test cases?
- **Question 2:** If we **cannot conclude** from few test cases to all inputs, when should we **stop testing**?

13/44

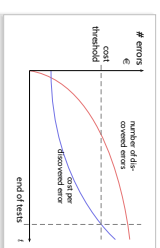
When To Stop Testing?

- There need to be defined **criteria** for when to stop testing. project planning should consider these criteria and previous experience!
- Possible "Testing is done" criteria:
 - all (previously) specified test cases have been executed with negative result.
 - **Special case**: All test cases resulting from a certain strategy (like maximal statement coverage) (\rightarrow if a minute) have been executed)
 - testing effort time sums up to x (hours, days, weeks).
 - testing effort sums up to y (any other useful unit).
 - x errors have been discovered.
 - **no error** has been discovered **during** the last z hours/days, weeks of testing.
- Values for x , y , z , x are fixed based on experience, estimation, budget, etc.
- **Of course**, not all criteria are equally reasonable or compatible with each testing approach

15/44

Another Criterion

- Another possible "Testing is done" criterion:
 - The average cost per error discovery exceeds a defined threshold c .



Value for c is again fixed based on experience, estimation, budget, etc.

16/44

Content

- Limits of Software Testing
 - Software examination paths
 - Is exhaustive testing feasible?
 - Range- vs. point errors
- **When To Stop Testing?**
 - **Choosing Test Cases**
 - Requirements on test cases
 - The natural habitat of many errors
 - Test Oracle
- **Glass-Box Testing**
 - Statement coverage
 - Branch and term coverage
 - Conclusions from coverage measures
- **Model-Based Testing**
- **Testing in the Development Process**

17/44

When To Stop Testing?

18/44

Choosing Test Cases

Choosing Test Cases

A test case is a **good test case** if it discovers – with high probability – an **unknown error**.
An **ideal test case** (I_n, Sol) would be

- of **low redundancy**, i.e. it does not test what other test cases also test.
- **error sensitive**, i.e. has high probability to detect an error.
(probability should at least be greater than 0)
- **representative**, i.e. represent a whole class of inputs.
(i.e. software S passes (I_n, Sol) if and only if S behaves well for all I' from the class)

The wish for representative test cases is **particularly problematic**:
(Recall: **pothole errors**)

- In general, we **do not know** which inputs lie in an equivalence class wrt. a certain error.
- **assuming** we know the equivalence classes.

Still, it is perfectly reasonable to test representatives of equivalence classes induced by the specification, e.g.

- valid and invalid inputs to check whether input validation works as all.
- different classes of inputs considered in the requirements, like "buy water", "buy soft-drink", "buy tea" vs. "buy beverage".

[Recall: strive to have at least one test case per failure.]

Who is hunting lions should know how a lion looks like
and how it sounds. You can't tell a lion apart from a lioness
and which sounds the lion makes. (Ludwig and Lorenz, 2013)



Lion and Error Hunting

Hunting errors in software is (basically) the same.
Some traditional popular belief on software error hunting:

- Software errors (in contrast to lions?) **(seem to) enjoy**
- **large boundaries**, e.g.
 - 0, 1, 27 if software works on inputs from $[0, 2^31]$.
 - -1, 28 for error handling.
 - $-2^{31} - 1, 2^{31}$ on 32-bit architectures.
 - boundaries of arrays (first, last element).
 - boundaries of loops (first, last iteration).
- **special cases** of the problem (empty list, use-case without actor...).
- special cases of the programming language semantics.
- **complex implementations**.

→ **Good idea**: for each test case, note down why it has been chosen.
For example: "demonstrate that error handling is not completely broken".

Where Do We Get The "Sol"-Values From?

Recall: A test case is a pair (I_n, Sol) with proper expected for "sol" values.

- In an **ideal world**, all "sol"-values are defined by the (formal) requirements specification and effectively **pre-computable**.
- In the **this world**,
 - the formal requirements specification may only **reflectively** describe acceptable results without giving a **procedure** to compute the results.
 - there may not be a formal requirements specification, e.g.
 - "the game objects should be rendered properly".
 - "the compiler must translate the program correctly".
 - "the notification message should appear on a proper screen position".
 - "the data must be available for at least 10 days".
 - etc.
 - Then: need another instance to decide whether the observation is acceptable.
- The testing community prefers to call **any instance** which decides whether results are acceptable a **(test) oracle**.
(I'd prefer **not to call** automatic derivation of "sol"-values from a **formal specification** an "oracle"... :)) (Openion or agency considered to provide wise and insightful,] prophetic predictions or pre-cognition of the future, inspired by the gods," says Wikipedia)

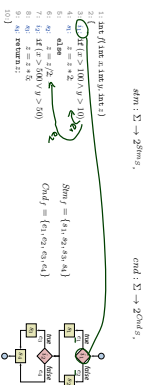
Glass-Box Testing: Coverage

Statements and Branches by Example

Definition 5.6 (run). A *finite description* S of a (possibly infinite) set $[S]$ of (finite or infinite) *computation paths* of the form $\sigma_0 \xrightarrow{a_0} \sigma_1 \xrightarrow{a_1} \sigma_2 \dots$ where

- $\sigma_i \in \Sigma$, $i \in \mathbb{N}_0$, is called *state* (or *configuration*), and
- $a_i \in A$, $i \in \mathbb{N}_0$, is called *action* (or *event*).

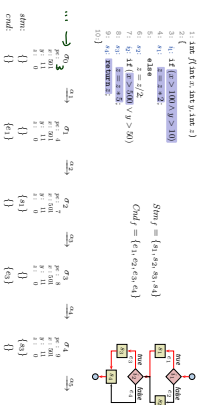
- In the following, we assume that
 - S has a control flow graph $(V, E)_S$, and statements $Sins \subseteq V$ and branches $Outs \subseteq E$.
 - each state σ gives information on statements and control flow graph branch edges which were executed right before obtaining σ .



24/4

Statements and Branches by Example

- In the following, we assume that
 - S has a **control flow graph** $\langle V, E \rangle_S$, and **statements** $Stms \subseteq V$ and **branches** $Brnds \subseteq E$
 - each state σ gives information on statements and control flow graph branch edges which were executed right before obtaining σ :



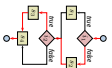
244

Coverage Example

```

int f(int x, int y, int z)
{
    if (x > 100 & y > 10)
        z = z + 2;
    else
        z = z/2;
    if (x > 500 & y > 50)
        z = z + 5;
    return z;
}

```



- **Requirement:** $\{true\} f \{true\}$ (no abnormal termination), i.e. $Soll = \Sigma^* \cup \Sigma^\omega$.

h	test 1: safe coverage											
	t_1/t_1	t_1/t_1	s_1	s_2	t_2/t_1	t_2	c_2	s_3	s_4	% 75 50 25	% 75 50 25	% 75 50 25
20 ± 2	✓	✓			✓	✓				100	100	100
10 ± 1.0	✓				✓	✓				100	75	25
5 ± 0.5	✓				✓					100	100	75
3 ± 1.0	✓				✓					100	100	100

26/4

Term Coverage

- Consider the statement

$$\text{if } \underbrace{(A \wedge (B \vee (C \wedge D))) \vee E}_{\text{exp}} \text{ then } \dots$$
 where A, \dots, E are **minimal** boolean terms, e.g. $x > 0$, but not $a \vee b$.
 Branch coverage is easy in this case:
 Use tr_1 such that $(A = 0, \dots, E = 0)$, and tr_2 such that $(A = 0, \dots, E = 1)$

- **Additional goal:**

or terms causing abnormal program te

- **Term Coverage** (for an expression $expr$)

- Let $\beta : \{A_1, \dots, A_n\} \rightarrow B$ be a valuation of the terms

- Term A_i is **b-effective** in β for *expr* if and only if

$$\beta(A_i) = b \text{ and } [\text{exp}^r](\beta(A_i/\text{true})) \neq [\text{exp}^r](\beta(A_i/\text{false}))$$

	A	B	C	D	E	b	%
β_1	1	1	0	0	1		20
β_2	1	0	1	0	0		50
β_3	1	0	1	0	1		70
β_4	0	0	1	0	1		80

- $\Xi \subseteq (\{A_1, \dots, A_n\} \rightarrow B)$ achieves p % **term coverage** if and only if

$$p = \frac{|(A_i^b \mid \exists \beta \in \Xi \bullet A_i \text{ is } b\text{-effective in } \beta)|}{2^n}$$

27/4

Glass-Box Testing: Coverage

- Coverage is a property of test cases and test suites.
- Execution π of test case T achieves $p\%$ statement coverage if and only if

$$p = cov_{stm}(\pi) := \frac{|\bigcup_{i \in N_0} stm(\sigma_i)|}{|Stm_S|}, |Stm_S| \neq 0.$$

Test case T achieves p % statement coverage if and only if $p = \min_{\pi \in \text{execution of } T} \text{cov}_{stm}(\pi)$

- Execution π of T achieves $p\%$ branch coverage if and only if

$$p = \text{cov}_{\text{end}}(\pi) := \frac{|\cup_{i \in N_0} \text{end}(\sigma_i)|}{|\text{Cnd } s|}, |\text{Cnd } s| \neq 0$$

Test case T achieves p % branch coverage if and only if $p = \min_{\pi \text{ execution of } T} cov_{cond}(\pi)$

- **Define:** $p = 100$ for empty program.

- Statement/branch coverage canonically extends to test suite $T = \{T_1, \dots, T_n\}$
e.g. given executions π_1, \dots, π_n , T achieves

$$p = \frac{|\bigcup_{1 \leq j \leq n} \bigcup_{i \in N_0} stm(\pi_j^i)|}{|Stm_S|}, |Stm_S| \neq 0, \text{ statement coverage}$$

25/4

Unreachable Code

```
int f(int x, int y, int z)
{
    i1: if (x ≠ x)
    s1:  z = y/0;
    i2:  if (x = x ∨ z/0 = 27
    s2:  z = z * 2;
    s3:  return z;
}
```

- Statement s_1 is **never executed** because $x \neq x \iff \text{false}$
thus 100 % statement-/branch-/term-coverage is **not achievable**.
- Assume evaluating $s_2()$ causes (undesired) **abnormal program termination**
Is statement s_3 an **error** in the program...?
- Term $s_3()$ in t_2 also looks critical...
- In programming languages with short-circuit evaluation, it is never evaluated)

28/4

- Assume, test suite T tests software S for the following property φ :
 - **pre-condition** p , **post-condition** q ,
and S passes (!) T , and the execution achieves 100 % statement / branch / term coverage.
 - **What does this tell us** about S ? Or: what can we conclude from coverage measures?
 - 100 % statement coverage:
 - There is no statement, which necessarily violates φ .
 - Still, there may be many, many computation paths which violate φ , and which just have not been touched by T .
 - There is no unreachable statement!
 - 100 % branch (term) coverage:
 - There is no single branch (term) which necessarily **causes** violations of φ .
 - In other words: "for each condition (term), there is one computation path satisfying φ where the condition (term) evaluates to true and one for false".
 - There is no unused condition (term)!
- Not more** (\leftarrow exercises)!
That's definitely **something**, but not as much as "100 %" may sound like. .

29/44

- (Seems that) DO-178B, "Software Considerations in Airborne Systems and Equipment Certification", (which deals with the safety of software used in certain airborne systems) requires that certain **coverage measures** are reached, in particular something similar to **term coverage** (MC/DC coverage) (Next to development process requirements, reviews, unit testing, etc.)
- If not required, ask: what is the effort / gain ratio?
(Average effort to detect an error: term coverage needs high effort)
- Currently, the standard moves towards accepting certain **qualification** or **static analysis** tools to support (or even replace) some testing obligations.

30/44

- **Limits of Software Testing**
 - Software examination paths
 - be exhaustive testing feasible?
 - Range- vs. point errors
 - **When To Stop Testing?**
 - Choosing Test Cases
 - Requirements on test cases
 - The natural habitat of many errors
 - Test Oracle
 - **Glass-Box Testing**
 - Statement coverage
 - Branch and term coverage
 - Conclusions from coverage measures
 - **Model-Based Testing**
 - Testing in the Development Process

31/44

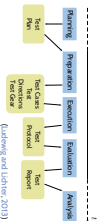
- Classical **statistical testing** is another approach to deal with
- in practice not exhaustively testable **huge input space**,
 - **tester bias**:
(People tend to choose "good-will" inputs and disregard some cases; recall: the developer is not a good tester)
- Procedure:**
- Randomly (!) choose test cases T_1, \dots, T_n for test suite T .
 - Execute test suite T .
 - If an error is found:
 - **good**: we certainly know there is an error,
 - if no error is found
 - **refute hypothesis**: "program is **not correct**" with a certain significance n_{mean} .
(Significance n_{mean} may be unsatisfactory with small test suites)
(And: Needs statistical assumptions on error distribution and truly random test cases.)

33/44

- (Ludewig and Ucker, 2013) name the following objections against statistical testing
- In particular for **interactive software**, the primary requirement is often **no failures are experienced by the "typical user"**.
Statistical testing in general may also cover a lot of "untypical user behaviours" unless sophisticated user models are used.
 - Statistical testing needs a method to compute "self"-values for the randomly chosen inputs.
That is easy for requirement "does not crash" but can be difficult in general.
If they live in their "natural habitat", carefully crafted test cases would probably uncover them.
 - **There is a high risk for not finding point or small-range errors**.
- Findings in the literature can at best be called **inconclusive**.

34/44

32/44



- **Test Case:** may need to be developed in the project!

test driver – A software module used to invoke a module under test and, often, to provide the test data and monitor execution, and report test results.
Synchron. test harness IEEE 61012 (1990)

stub – A module for special purpose implementation of a software module, used to develop or test a module that calls or is otherwise dependent on it.
(2) A computer program statement substituting for the body of a software module that is or will be defined elsewhere. IEEE 61012 (1990)

- **Roles:** **tester** and **developer** should be different persons!

- **A checker can only discover errors on exactly one path.**
- **Software testing is challenging** because
 - typically, huge input space
 - software is non-continuous
- Define criteria for “testing done” (like coverage or cost per error)
- There is a **test amount of iterations** on how to choose test cases.

A good starting point:

- at least one test case per feature,
- corner-cases, external values,
- error handling, etc.

- **Class-box testing**

- considers the control flow graph,
- defines coverage measures.

- **Other approaches:**

- statistical testing, model-based testing

Continuous integration → V4, 18

- **Proces:** tester and developer should be different persons.

References

References

IEEE (1990). *IEEE Standard Glossary of Software Engineering Terminology*. Std 61012-1990.

Lethal, M. and Voss, J. (2001). Scenario-based monitoring and testing of real-time UML models. In *ICSE*. IEEE Press, 2001.

Lecture Notes in Computer Science, pages 317–328. Springer-Verlag.

Ludewig, J. and Uicker, H. (2013). *Software Engineering*. dpunkt.verlag, 3. edition.