

Softwaretechnik / Software-Engineering

Lecture 16: Testing & Review

2015-07-13

Prof. Dr. Andreas Podtiski, Dr. Bernd Westphal
Albert-Ludwigs-Universität Freiburg, Germany

Contents of the Block "Quality Assurance"

(i) Introduction and Vocabulary	L 1: 20.4, Mo
• correctness illustrated	L 2: 27.4, Mo
• correctness: success or failure	L 3: 30.4, Do
• three basic approaches	L 4: 7.5, Mo
(ii) Formal Verification	L 5: 11.5, Mo
• Hoare calculus	L 6: 14.5, Do
• Verifying C Compiler (VCC)	L 7: 21.5, Do
• over- / under-approximations	L 8: 28.5, Mo
(iii) Systematic Tests	L 9: 4.6, Do
• systematic test vs. experiment	L 10: 13.6, Mo
• classification of test procedures	L 11: 22.6, Mo
• global test: coverage measure	L 12: 29.6, Do
(iv) Runtime Verification	L 13: 2.7, Do
(v) Review	L 14: 2.7, Do
(vi) Concluding Discussion	L 15: 6.7, Mo
• Dependability	L 16: 13.7, Mo
	L 17: 16.7, Do
	L 18: 23.7, Do

2, 48

The Verifying C Compiler

- The Verifying C Compiler (VCC) basically implements Hoare-style reasoning.
- Special syntax:**
 - #include <vec.h>
 - (requires p) — pre-condition, p is a C expression
 - (ensures q) — post-condition, q is a C expression
 - (invariant $expr$) — loop invariant, $expr$ is a C expression
 - (assert p) — intermediate invariant, p is a C expression
 - (writes lx) — VCC considers concurrent C programs, we need to declare for each procedure which global variables it is allowed to write to (also checked by VCC)
- Special expressions:**
 - thread_local(lv) — no other thread writes to variable v (in pre-conditions)
 - volatile — the value of v when procedure was called (useful for post-conditions)
 - retval — return value of procedure (useful for post-conditions)

4, 68

VCC

- The Verifying C Compiler (VCC) basically implements Hoare-style reasoning.
- Special syntax:**
 - #include <vec.h>
 - (requires p) — pre-condition, p is a C expression
 - (ensures q) — post-condition, q is a C expression
 - (invariant $expr$) — loop invariant, $expr$ is a C expression
 - (assert p) — intermediate invariant, p is a C expression
 - (writes lx) — VCC considers concurrent C programs, we need to declare for each procedure which global variables it is allowed to write to (also checked by VCC)
- Special expressions:**
 - thread_local(lv) — no other thread writes to variable v (in pre-conditions)
 - volatile — the value of v when procedure was called (useful for post-conditions)
 - retval — return value of procedure (useful for post-conditions)

16 - 2015-07-13 - Svcc -

5, 68

Contents & Goals

- Last Lecture:**
 - Completed the block "Architecture & Design"
- This Lecture:**
 - Educational Objectives:** Capabilities for following tasks/questions:
 - What can we conclude from the outcome of tests like VCC?
 - What is an example for not a test, non-systematic test, systematic test?
 - Given a test case and a software, is the outcome successful or unsuccessful?
 - How many test cases are necessary for exhaustive testing of a given software?
 - Content:**
 - The Verifying C Compiler (VCC)
 - Systematic test, test case, test suite
 - Testing notions
 - Coverage measures

16 - 2015-07-13 - Sptelm -

3, 68

VCC Syntax Example

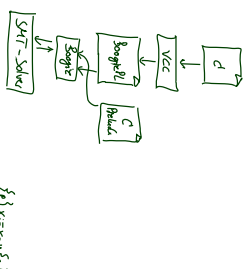
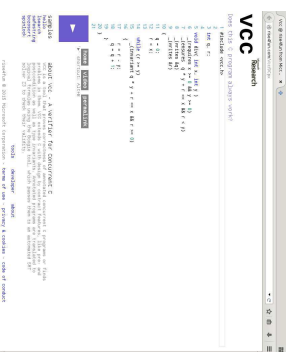
```

1 #include <vec.h>
2 int g;
3 void div( int x, int y )
4 {
5     requires x >= 0 && y >= 0
6     {
7         invariant q + x + y + r == x && r < y
8         -writes &q
9     }
10    q = 0;
11    while ( x >= y )
12    {
13        invariant q + x + y + r == x && r >= 0
14        r = x;
15        x = x - y;
16        q = q + 1;
17    }
18 }
19 DIV ≡ q := 0; r := x; while ( x >= y ) do r := r - y; q := q + 1 do
20 { x >= 0 & y >= 0 } DIV { q + y + r == x & r < y }

```

16 - 2015-07-13 - Svcc -

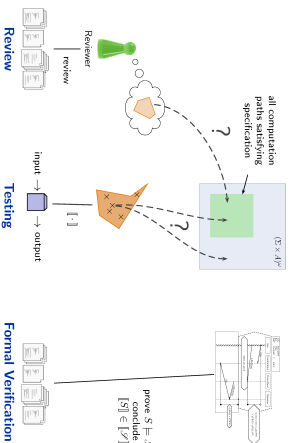
6, 68



- For the exercises, we use VCC only for sequential, single-thread programs.
- VCC checks a number of **implicit assertions**:
 - no arithmetic overflow in expressions (according to C-standard),
 - array-out-of-bounds access,
 - NULL-pointer dereference,
 - and many more.
- VCC also supports:
 - concurrency**: different threads may write to shared global variables. VCC can check whether concurrent access to shared variables is properly managed;
 - data structure invariants**: we may declare invariants that have to hold for, e.g., records (e.g. the length field *l* is always equal to the length of the string field *str*); those invariants may **temporarily** be violated when updating the data structure.
 - and much more.
- Verification does not always succeed!**
 - The backend SMT-solver may not be able to discharge proof-obligations (in particular non-linear multiplication and division are challenging);
 - In many cases, we need to provide **loop invariants** manually.



- VCC says: "verification succeeded"
 - We can **only conclude** that the tool — under its platform assumptions (32bit), etc. — thinks "that it can prove $\models (p) \text{ DIV } (q)$ ". Can be due to an error in the tool!
 - Yes we can ask for a printout of the proof and check it manually (lucky people in practice) or with other tools like interactive theorem provers.
 - Note:** $\models (false) \text{ } f \text{ } (q)$ **always holds**
 - so a mistake in writing down the pre-condition can provide a **false negative**.
- VCC says: "verification failed"
 - One case: "timeout" etc. — completely inconclusive outcome.
 - The tool **does not provide counter-examples** in the form of a computation path. It (only) gives **hints on input values** satisfying *p* and causing a violation of *q*. May be a **false negative** if these inputs are actually never used. Make pre-condition *p* stronger, and try again.



Quotes On Testing

"Testing is the execution of a program with the goal to discover errors." ¹
G.J. Myers, 1979

"Testing is the demonstration of a program or system with the goal to show that it does what it is supposed to do." ²
W. Hetzel, 1984

¹
"Software can be used to show the presence of bugs, but never to show their absence!" ²
E.W. Dijkstra, 1970

Rule-of-thumb: (fairly systematic) tests discover half of all errors.
(Ludwing and Lichten, 2013)

13/69

Tests vs. Systematic Tests

Test — (one or multiple) execution(s) of a program on a computer with the goal to find errors.
(Ludwing and Lichten, 2013)

(Our) **Synonyms**: Experiment, Rumpöbchen.

Not (even) a test (in the sense of this weak definition):

- any inspection of the program,
- demo of the program,
- analysis by software-tools, e.g. for values of metrics,
- investigation of the program with a debugger.

Systematic Test — a test with

- (environment) conditions are defined or precisely documented,
- inputs have been chosen systematically,
- results documented and assessed according to criteria that have been fixed before.
(Ludwing and Lichten, 2013)

In the following: **experiment** := test — test := systematic test.

16

14/69

More Formally: Test Case

- A **test case** T is a set of pairs $\{(h_1, Sd1), \dots\}$ consisting of
 - a (description of a) finite **input** sequence h_1 (activities different in T),
 - a (description of a) finite set of **expected** computation path $Sd1$.

Examples

- $T_1 = \{(\text{FILLUP, C50: water button on}) \mid \text{insert C50 coin (at any time), expect water button is enabled (same time later)}\}$ (shorthand notation)
- $T_2 = \{(o_0^0 \xrightarrow{-1} o_1^1; o_1^1 \xrightarrow{-2} o_2^2) \mid o_1^1(x) = T \wedge \sigma_1(y) = 49\}$ (input T , expect output 49, don't care for other variables' values; shorthand notation: $T; 49$)
- $T_3 = \{(o_0^0 \xrightarrow{-5} o_1^1; o_1^1 \xrightarrow{-4} o_2^2) \mid o_1^1 = 0(x := T); o_2^2 = 0(\sigma_1 = 0)y := 49\}$ (each and every variable value at start and at end fixed) ¹
Random selection

16

15/69

Test Case Execution, Test Suite

- An **execution** of test case T for software S is a computation path of S
 $\pi = (o_0^0 \xrightarrow{-1} o_1^1 \xrightarrow{-2} o_2^2 \dots o_{i-1}^{i-1} \xrightarrow{-i} o_i^i \xrightarrow{-i+1} o_{i+1}^{i+1} \dots o_{n-1}^{n-1} \xrightarrow{-n} o_n^n)$, where $o_0^0 \xrightarrow{-1} o_1^1 \xrightarrow{-2} o_2^2 \dots = h_1$ for some i in T .

- The **test case execution** is called

- **successful** (or **positive**) if it discovered an error, i.e. if $\pi \notin Sd1$.
(Alternative test item failed to pass test; containing: "test failed")
- **unsuccessful** (or **negative**) if it did not discover an error, i.e. if $\pi \in Sd1$.
(Alternative test item passed test; okay; "test passed")

Note: if input sequence not adhered to, or power outage, etc., it is not a test execution.

- A **test suite** is a set of test cases.
Execution, **positive**, and **negative** are listed canonically

16/69

16

16/69

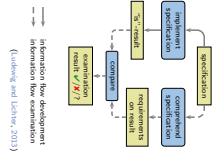
The Outcome of Systematic Tests Depends on...

- **inputs**:
 - the input vector of the test case (of course), possibly with timing constraints,
 - other interaction, e.g. from network,
 - initial memory content,
 - etc.
- **(environmental) conditions**
 - any aspects which could have an effect on the outcome of the test such as
 - which program (version) is tested? built with which compiler, linker, etc.?
 - test host (OS architecture, memory size, connected device (configuration), etc.)
 - which other software (in which version, configuration) is involved?
 - who tested when?
 - etc.
 - ... so strictly speaking **all** of them need to be specified within (or as an extension to) h_1 .
 - **In practice**, this is hardly possible — but one wants to specify as much as possible in order to achieve **reproducibility**.
- **One approach**:
 - have a fixed build environment, a fixed test host which does not do any other jobs, etc.

16

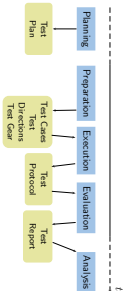
17/69

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



Recall:

	checking procedure	product error
time to market	long	short
artefact has error	low	high
	time to market	time to market



- **Test Case:**
 - **test driver**— A software module used to invoke a module under test and, often, provide test inputs, control and monitor execution and report test results.
 - **Synonym:** test harness.
- **stub(1)** A skeletal or 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 statements substituting for the body of a software module that is or will be defined elsewhere.
- **IEEE 60112 (1990)**
- **hardware-in-the-loop, software-in-the-loop:** the final implementation is running on (prototype) hardware, other system component are simulated by a separate computer.

- How are the test cases **chosen**?
- Considering the structure of the test item (glass-box or structure test).
- Considering only the specification (black-box or function test).
- How much **effort** is put into testing?
- **execution trial** — does the program run at all?
- **throw-away-test** — invent input and judges output on-the-fly,
- **systematic test** — somebody (not author) derives test cases, defines input, test, documents test execution.
- In the long run, **systematic tests** are more **economic**.
- **Complexity** of the test item:
 - **unit test** — a single program unit is tested (function, sub-routine, method, class, etc.)
 - **module test** — a component is tested.
 - **integration test** — the interplay between components is tested.
 - **system test** — tests whole system.

- Which **property** is tested?
- **function test** — functionality as specified by the requirements documents.
- **installation test** — is it possible to install the software with the provided documentation and tools?
- **recommissioning test** — is it possible to bring the system back to operation after operation was stopped?
- **availability test** — does the system run for the required amount of time without issues, **load** and **stress test** — does the system behave as required under high or highest load? ... under overload?
- **regression test** — does the new version of the software behave like the old one on inputs where no behaviour change is expected?
- **response time** . **minimal hardware (software) requirements**, etc.
- Which roles are **involved** in testing?
- **only** the developer, or selected (potential) customers (**alpha** and **beta** test).
- **acceptance test** — the customer tests whether the system (or parts of it, at milestones) test whether the system is acceptable.

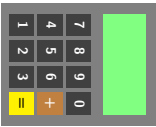


- **Requirement:**
 - If the display shows x , $+$, and y , then after pressing **=**.
 - the sum of x and y is displayed if $x + y$ has at most 8 digits.
 - otherwise **“E”** is displayed.



- **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.

Testing the Pocket Calculator



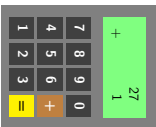
Test some representatives of "equivalence classes":

- $n + 1$, n small.
- $n + m$, n small, m small (for non error).
- $n + m$, n big, m big (for non error).
- $n + m$, n huge, m small (for error).
- ...

e.g. $27 + 1$
e.g. $13 + 27$
e.g. $12345 + 678$
e.g. $99999999 + 1$

23.08

Testing the Pocket Calculator



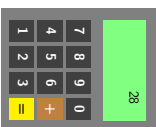
Test some representatives of "equivalence class":

- $n + 1$, n small.
- $n + m$, n small, m small (for non error).
- $n + m$, n big, m big (for non error).
- $n + m$, n huge, m small (for error).
- ...

e.g. $27 + 1$
e.g. $13 + 27$
e.g. $12345 + 678$
e.g. $99999999 + 1$

23.08

Testing the Pocket Calculator



Test some representatives of "equivalence class":

- $n + 1$, n small.
- $n + m$, n small, m small (for non error).
- $n + m$, n big, m big (for non error).
- $n + m$, n huge, m small (for error).
- ...

e.g. $27 + 1$
e.g. $13 + 27$
e.g. $12345 + 678$
e.g. $99999999 + 1$

23.08

Testing the Pocket Calculator



Test some representatives of "equivalence classes":

- $n + 1$, n small.
- $n + m$, n small, m small (for non error).
- $n + m$, n big, m big (for non error).
- $n + m$, n huge, m small (for error).
- ...

e.g. $27 + 1$
e.g. $13 + 27$
e.g. $12345 + 678$
e.g. $99999999 + 1$

23.08

Testing the Pocket Calculator



Test some representatives of "equivalence class":

- $n + 1$, n small.
- $n + m$, n small, m small (for non error).
- $n + m$, n big, m big (for non error).
- $n + m$, n huge, m small (for error).
- ...

e.g. $27 + 1$
e.g. $13 + 27$
e.g. $12345 + 678$
e.g. $99999999 + 1$

23.08

Testing the Pocket Calculator



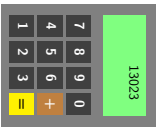
Test some representatives of "equivalence class":

- $n + 1$, n small.
- $n + m$, n small, m small (for non error).
- $n + m$, n big, m big (for non error).
- $n + m$, n huge, m small (for error).
- ...

e.g. $27 + 1$
e.g. $13 + 27$
e.g. $12345 + 678$
e.g. $99999999 + 1$

23.08

Testing the Pocket Calculator



Test some representatives of "equivalence classes":

- $n + 1, n$ small.
- $n + m, n$ small, m small (for non error).
- $n + m, n$ big, m big (for non error).
- $n + m, n$ huge, m small (for error).
- ...

e.g. $27 + 1$
e.g. $13 + 27$
e.g. $12345 + 678$
e.g. $99999999 + 1$

23/68

Testing the Pocket Calculator



Test some representatives of "equivalence class":

- $n + 1, n$ small.
- $n + m, n$ small, m small (for non error).
- $n + m, n$ big, m big (for non error).
- $n + m, n$ huge, m small (for error).
- ...

e.g. $27 + 1$
e.g. $13 + 27$
e.g. $12345 + 678$
e.g. $99999999 + 1$

23/68

Testing the Pocket Calculator



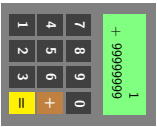
Test some representatives of "equivalence classes":

- $n + 1, n$ small.
- $n + m, n$ small, m small (for non error).
- $n + m, n$ big, m big (for non error).
- $n + m, n$ huge, m small (for error).
- ...

e.g. $27 + 1$
e.g. $13 + 27$
e.g. $12345 + 678$
e.g. $99999999 + 1$

23/68

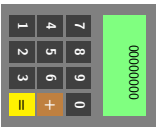
Testing the Pocket Calculator: One More Try



Oops...

24/68

Testing the Pocket Calculator: One More Try



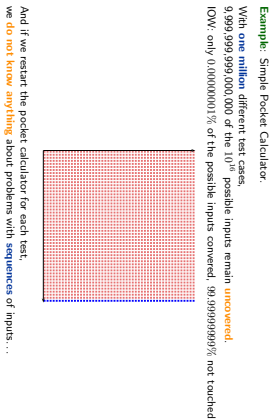
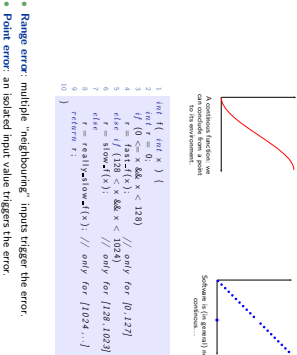
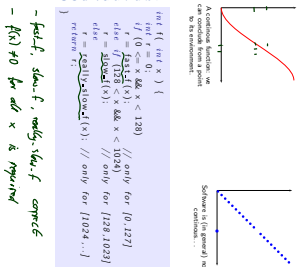
Oops...

24/68

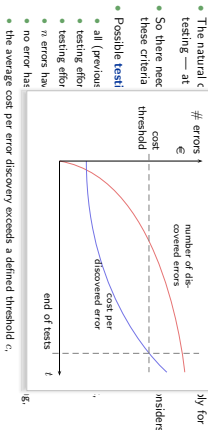
Behind the Scenes: Test "99999999 + 1" Failed Because...

```
1 int add( int x, int y )
2 {
3     if (y == 1) // be fast
4         return x + 1;
5     int r = x + y;
6     if (r > 99999999)
7         r = -1;
8     return r;
9 }
```

25/68



- The natural criterion "when **everything has been done**" does not apply for testing — at least not for testing pocket calculators.
- So there need to be defined **criteria** to stop testing; project planning considers these criteria and experience with them.
- Possible **testing is done** criteria:
 - all (previously) specified test cases have been executed with negative result.
 - testing effort sums up to x hours (days, weeks).
 - testing effort sums up to y (any other useful unit).
 - n errors have been discovered.
 - no error has been discovered during the last z hours (days, weeks) of testing.
 - the average cost per error discovery exceeds a defined threshold c .



- The natural criterion “**when everything has been done**” does not apply for testing — at least not for testing pocket calculators.
- So there need to be defined **criteria** to stop testing: project planning considers these criteria and experience with them.
- Possible **testing is done** criteria:
 - all (previously) specified test cases have been executed with negative result.
 - testing effort sums up to x hours (days, weeks).
 - testing effort sums up to y (any other useful unit).
 - n errors have been discovered.
 - no error has been discovered during the last z hours (days, weeks) of testing.
 - the average cost per error discovery exceeds a defined threshold c .
- Values for x , y , n , z , c are fixed based on experience, estimation, budget, etc..
- **Of course:** not all equally reasonable or compatible with each testing approach.

Choosing Test Cases

Choosing Test Cases

- A test case is a **good test case** if discovers with high probability an unknown error.
- An **ideal** test case should be
- **representative**, i.e. represent a whole class of inputs,
 - **error sensitive**, i.e. has high probability to detect an error,
 - **of low redundancy**, i.e. it does not test what other test cases also test.
- The wish for representative test cases is **particularly problematic**:
- Recall **point errors** (pocket calculator, fast/slow $f \dots$).

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

- “Acceptable” equivalence classes: Based on requirement specification, e.g.
- valid and invalid inputs (to check whether input validation works),
 - different classes of inputs considered in the requirements, e.g. “buy water”, “buy soft-drink”, “buy tea” vs. “buy beverage”.

Lion and Error Hunting

“The *šim* who is hunting lions, should know how a lion looks like. He/she should also know where the lion likes to stay, which traces the lion leaves behind, and which sounds the lion makes.”
(Ludewig and Lüthke, 2013)

Hunting errors in software is (basically) the same.

Some traditional popular belief on software error habitat:

- Software errors — in contrast to lions — (seem to) enjoy
 - range boundaries, e.g.
 - 0, 1, 27 if software works on inputs from $[0, 27]$,
 - -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,

Where Do We Get The “*Soil*“-Values From?

- In an **ideal world**, all test cases are pairs $(In, Soil)$ with proper “soil“-values. As, for example, defined by the formal requirements specification.
- **Advantage:** we can mechanically, objectively check for positive/negative.

• 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 an **oracle**.

- I prefer not to call decisions based on **formally defined** test cases “oracle” ... :-)

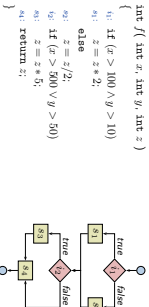
Glass-Box Testing: Coverage

Class-Box Testing: Coverage

- **Coverage** is a property of **test cases** and **test suite**.
- **Reach**: An execution of test case $T = (I_n, SdI)$ for software S is a computation path $\left(\sigma_0^1 \right)_{\sigma_1^1} \dots \left(\sigma_i^1 \right)_{\sigma_{i+1}^1} \dots \left(\sigma_n^1 \right)_{\sigma_{n+1}^1}$ where $\sigma_0^1, \sigma_1^1, \sigma_2^1, \dots, \sigma_n^1, \sigma_{n+1}^1 = I_n$.
- Let S be a **program** (or **model**) consisting of **statements** S_{stmt} , **conditions** S_{con} , and a **control flow graph** (V, E) (as defined by the programming language).
- Assume that each state σ gives information on statements, conditions, and control flow graph edges **which were executed** right before obtaining σ .
 $den : \Sigma \rightarrow 2^{S_{stmt}}, \quad edg : \Sigma \rightarrow 2^E$
 $den : \Sigma \rightarrow 2^{S_{con}}, \quad cond : \Sigma \rightarrow 2^{con}$
- T achieves $p\%$ **statement coverage** if and only if $p = \frac{|\bigcup_{\sigma \in den(C_T)} S_{stmt}|}{|S_{stmt}|}, |S_{stmt}| \neq 0$.
- T achieves $p\%$ **branch coverage** if and only if $p = \frac{|\bigcup_{\sigma \in den(C_T)} E|}{|E|}, |E| \neq 0$.
- **Define**: $p = 100$ for empty program.
- Statement/branch coverage canonically extends to test suite T .

34.68

Coverage Example

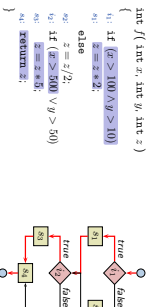


- Requirement: $\{true\} f \{true\}$ (no abnormal termination)

16 - 2015-07-13 - Score -

35.68

Coverage Example

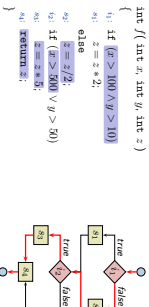


- Requirement: $\{true\} f \{true\}$ (no abnormal termination)

16 - 2015-07-13 - Score -

35.68

Coverage Example



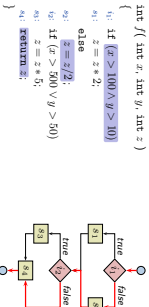
- Requirement: $\{true\} f \{true\}$ (no abnormal termination)

x, y, z	$v1/f$	$v1/f$	$a1$	$v2/f$	$v2/f$	$a3$	$a3$	$stmt$	$cond$	$term$	$\%$	$v2/\%$
501,11,0	✓	✓	✓	✓	✓	✓	✓	75	50	25		
501,0,0												

16 - 2015-07-13 - Score -

35.68

Coverage Example



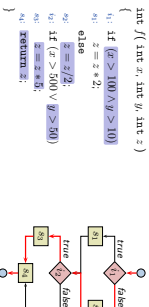
- Requirement: $\{true\} f \{true\}$ (no abnormal termination)

x, y, z	$v1/f$	$v1/f$	$a1$	$v2/f$	$v2/f$	$a3$	$a3$	$stmt$	$cond$	$term$	$\%$	$v2/\%$
501,11,0	✓	✓	✓	✓	✓	✓	✓	75	50	25		
501,0,0												
0,0,0								100	75	25		

16 - 2015-07-13 - Score -

35.68

Coverage Example



- Requirement: $\{true\} f \{true\}$ (no abnormal termination)

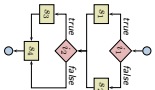
x, y, z	$v1/f$	$v1/f$	$a1$	$v2/f$	$v2/f$	$a3$	$a3$	$stmt$	$cond$	$term$	$\%$	$v2/\%$
501,11,0	✓	✓	✓	✓	✓	✓	✓	75	50	25		
501,0,0								100	75	25		
0,0,0								100	100	75		
0,51,0												

16 - 2015-07-13 - Score -

35.68

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)

π_0, z	s_1/f	s_1/f	s_2	b_2/f	$1/z/f$	c_1	c_2	c_3	s_{eff}	% stem	% dg	$t_2/\%$
50, 1.0	✓	✓	✓	✓	✓	✓	✓	✓	75	50	75	50
50, 0.0	✓	✓	✓	✓	✓	✓	✓	✓	100	75	75	75
0.51, 0	✓	✓	✓	✓	✓	✓	✓	✓	100	100	100	100
0.51, 0	✓	✓	✓	✓	✓	✓	✓	✓	100	100	100	100

35/6

Conclusions from Coverage Measures

- Assume, we are testing property $\varphi: (v_1 \neq v_2) \wedge \{v_1\}$ (maybe just $v_1 = \text{true}$ with $\{v_1\}$).
- assume our test suite T achieved 100% statement-, branch / term coverage.
- What does this tell us about φ ? 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 , e.g. differing in variable / valuation).
- There is no unreachable statement!

- 100% **branch (term)** coverage

- "there is no single branch (term) which necessarily causes violations of φ "
- IOW: "for each condition (term), there is one computation path satisfying φ where the condition (term) evaluates to true/false"
- "there is no unused condition (term)"

Not more (\rightarrow exercises)!

That's **something**, but not as much as "100 %" may sound. . . .

38/e

Term Coverage

- Consider the statement

$$\text{if } (A \wedge B \vee (C \wedge D)) \vee E \text{ then } \dots$$

$$\text{if } (A \wedge B \vee (C \wedge D)) \vee E \text{ then } \dots$$
 - A, \dots, E are minimal boolean terms, e.g. $x > 0$, but not $a < b$.
 - Branch coverage is easy: use $(A = 0, \dots, E = 0)$ and $(A = 0, \dots, E = 1)$.
 - Additional goal:** check whether there are useless terms, or terms causing abnormal program termination.
 - Term coverage** (for an expression $x \cdot y \cdot z$)
 - Let $\beta: \{A_1, \dots, A_n\} \rightarrow B$ be a valuation of the terms.
 - Term A_i is **effective** in β if x is *exp* if and only if
- | A | B | C | D |
|---|---|---|---|
| 1 | 1 | 1 | 1 |
| 1 | 1 | 1 | 0 |
| 1 | 1 | 0 | 1 |
| 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 |
| 0 | 1 | 0 | 1 |
| 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 1 |
| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 0 |

A	B	C	D	E	%
1	1	0	0	0	20
1	0	0	1	0	50
1	0	1	1	0	70
0	0	1	0	1	80

$$\beta(A_i) = b \text{ and } \llbracket \text{expr} \rrbracket(\beta(A_i/\text{true})) \neq \llbracket \text{expr} \rrbracket(\beta(A_i/\text{false}))$$

- $\Xi \subseteq (\{A_1, \dots, A_n\} \rightarrow B)$ achieves $p\%$ **term coverage** if and only if
$$p = \frac{|\{A_i^b \mid \exists \beta \in \Xi \bullet A_i \text{ is } b\text{-effective in } \beta\}|}{\gamma_n}.$$

$$\frac{2n}{\left| \left\{ \left(x_1, \dots, x_n \right) \in \mathbb{R}^n : x_1 + \dots + x_n = 0 \right\} \right|}$$

- 16 - 2015.07.13 - Score -

36/6

Unreachable Code

```

int f( int x, int y, int z )
{
    if (x != y)
    {
        z = y/0;
        if (x = x * z / 0 = 27)
            z = z * 2;
        return z;
    }
}

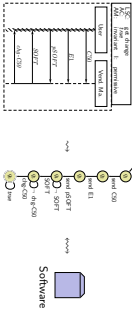
```

- Statement s_1 is **never executed** ($x \neq x \iff \text{false}$), thus 100% coverage **not achievable**.
- Is statement s_1 an **error** anyway...?
- Term $y/0$ is never evaluated either (short-circuit evaluation)

37/6

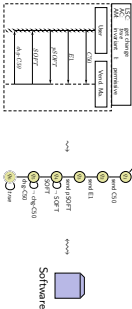
Model-Based Testing

40/6



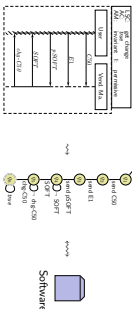
- If the LSC has designated **environment instance** lines, we can distinguish:
 - messages expected to originate **from** the environment (driver role),
 - messages expected addressed to the environment (monitor role).

42.09



- If the LSC has designated **environment instance** lines, we can distinguish:
 - messages expected to originate **from** the environment (driver role),
 - messages expected addressed to the environment (monitor role).
- Adjust the TBA-construction algorithm to construct a **test driver & monitor** and have it (possibly with some **glue logic** in the middle) interact with the software (or a model or it).
- **Test passed** (i.e., test unsuccessful) if and only if TBA state q_0 is reached.

42.09



- If the LSC has designated **environment instance** lines, we can distinguish:
 - messages expected to originate **from** the environment (driver role),
 - messages expected addressed to the environment (monitor role).
- Adjust the TBA-construction algorithm to construct a **test driver & monitor** and have it (possibly with some **glue logic** in the middle) interact with the software (or a model or it).
- **Test passed** (i.e., test unsuccessful) if and only if TBA state q_0 is reached.
- We may need to **refine** the LSC by adding an activation condition, or communication which drives the system under test into the desired start state.

42.09

Another Approach: Statistical Tests

One proposal to deal with the **uncertainty of tests**, and to **avoid bias** (people tend to choose expected inputs): **classical statistical testing**.

Statistical Testing

43.09

44.09

Another Approach: Statistical Tests

One proposal to deal with the **uncertainty of tests**, and to **avoid bias** (people tend to choose expected inputs): **classical statistical testing**.

- Randomly choose and apply test cases T_1, \dots, T_n .
 - **if an error is found**: good, we certainly know there is an error,
 - **if no error is found**: "program is not correct" with a certain confidence interval. refine hypothesis "program is not correct" with a certain confidence interval.
- Needs stochastic assumptions on error distribution and truly random test cases.
(Confidence interval may get large — reflecting the low information tests give.)

44.09

One proposal to deal with the **uncertainty of tests**, and to **avoid bias** (people tend to choose expected inputs): classical **statistical testing**.

- Randomly choose and apply test cases T_1, \dots, T_n .
- if an error is found**: good, we certainly know there is an error.
- if no error is found**:

refuse hypothesis: "program is not correct" with a certain confidence interval.

Needs stochastic assumptions on error distribution and truly random test cases.

(Confidence interval may get large — reflecting the low information tests give.)

(Ludewig and Lichter, 2013) name the following objections against statistical testing:

One proposal to deal with the **uncertainty of tests** and to **avoid bias** (people tend to choose expected inputs): classical **statistical testing**.

- Randomly choose and apply test cases T_1, \dots, T_n .
- if an error is found**: good, we certainly know there is an error.
- if no error is found**:

refuse hypothesis: "program is not correct" with a certain confidence interval.

Needs stochastic assumptions on error distribution and truly random test cases.

(Confidence interval may get large — reflecting the low information tests give.)

(Ludewig and Lichter, 2013) name the following objections against statistical testing:

- In particular for interactive software, the primary goal is often that the "typical user" does not experience failures. Statistical testing (in general) may also cover a bit of "untypical user behaviour", unless **user models** are used.

One proposal to deal with the **uncertainty of tests**, and to **avoid bias** (people tend to choose expected inputs): classical **statistical testing**.

- Randomly choose and apply test cases T_1, \dots, T_n .
- if an error is found**: good, we certainly know there is an error.
- if no error is found**:

refuse hypothesis: "program is not correct" with a certain confidence interval.

Needs stochastic assumptions on error distribution and truly random test cases.

(Confidence interval may get large — reflecting the low information tests give.)

(Ludewig and Lichter, 2013) name the following objections against statistical testing:

- In particular for interactive software, the primary goal is often that the "typical user" does not experience failures. Statistical testing (in general) may also cover a bit of "untypical user behaviour", unless **user models** are used.
- Statistical testing needs a method to compute "soft" values for the randomly chosen inputs, that is easy for "does not crash" but can be difficult in general.

One proposal to deal with the **uncertainty of tests**, and to **avoid bias** (people tend to choose expected inputs): classical **statistical testing**.

- Randomly choose and apply test cases T_1, \dots, T_n .
- if an error is found**: good, we certainly know there is an error.
- if no error is found**:

refuse hypothesis: "program is not correct" with a certain confidence interval.

Needs stochastic assumptions on error distribution and truly random test cases.

(Confidence interval may get large — reflecting the low information tests give.)

(Ludewig and Lichter, 2013) name the following objections against statistical testing:

- In particular for interactive software, the primary goal is often that the "typical user" does not experience failures. Statistical testing (in general) may also cover a lot of "untypical user behaviour", unless **user models** are used.
- Statistical testing needs a method to compute "soft" values for the randomly chosen inputs, that is easy for "does not crash" but can be difficult in general.
- There is a high risk for not finding point or small-range errors which do live in their "natural habitat" as expected by testers.

Findings in the literature can at best be called **inconclusive**.

- A low profile approach¹ when a formal (requirements) specification is not available, not even "agile-style" in form of test cases
- whenever a feature² is considered finished,

(i) make up **inputs** for (at least one) test case,

(ii) create **script** which **runs** the program on these inputs,

(iii) carefully **examine** the outputs for whether they are acceptable,

(iv) **if no** repair,

(v) **if yes**: define the observed output as "soft",

(vi) extend **script** to compare ist/soil and add to test suite.

1: best for pipe/filter style software, where comparing output with "soil" is trivial.

2: if test case creation is postponed too long, chances are high that there will not be any test cases at all. **Experiment**: "too long" is very short.

3: error handling is also a feature.

Advantages of testing (in particular over inspection):

Discussion

Advantages of testing (in particular over inspection):

- Testing is a **“human”** checking procedure: “everybody can test”.

Discussion

Advantages of testing (in particular over inspection):

- Testing is a **“human”** checking procedure: “everybody can test”.
- The systematic test is reproducible and **objective** (if the start configuration is reproducible and the test environment deterministic).

Discussion

Advantages of testing (in particular over inspection):

- Testing is a **“human”** checking procedure: “everybody can test”.
- The systematic test is reproducible and **objective** (if the start configuration is reproducible and the test environment deterministic).
- Invested effort can be **re-used**: properly prepared and documented tests can be re-executed with low effort, in particular fully automatic tests; important in maintenance.

Discussion

Advantages of testing (in particular over inspection):

- Testing is a **“human”** checking procedure: “everybody can test”.
- The systematic test is reproducible and **objective** (if the start configuration is reproducible and the test environment deterministic).
- Invested effort can be **re-used**: properly prepared and documented tests can be re-executed with low effort, in particular fully automatic tests; important in maintenance.
- The **test environment** is (implicitly) subject of testing; errors in additional components and tools may show up.

Discussion

Advantages of testing (in particular over inspection):

- Testing is a **“human”** checking procedure: “everybody can test”.
- The systematic test is reproducible and **objective** (if the start configuration is reproducible and the test environment deterministic).
- Invested effort can be **re-used**: properly prepared and documented tests can be re-executed with low effort, in particular fully automatic tests; important in maintenance.
- The **test environment** is (implicitly) subject of testing; errors in additional components and tools may show up.
- System behaviour (efficiency, usability) becomes **visible**, even if not explicitly subject of a test.

Discussion

Advantages of testing (in particular over inspection):

- Testing is a **“human”** checking procedure: “everybody can test”.
- The systematic test is reproducible and **objective** (if the start configuration is reproducible and the test environment deterministic).
- Invested effort can be **re-used**: properly prepared and documented tests can be re-executed with low effort, in particular fully automatic tests; important in maintenance.
- The **test environment** is (implicitly) subject of testing; errors in additional components and tools may show up.
- System behaviour (efficiency, usability) becomes **visible**, even if not explicitly subject of a test.

Disadvantages:

Advantages of testing (in particular over inspection):

- Testing is a **"natural"** checking procedure: "everybody can test".
- The systematic test is **reproducible** and **objective** (if the start configuration is reproducible and the test environment deterministic).
- Invested effort can be **re-used**: properly prepared and documented tests can be re-executed with low effort, in particular fully automatic tests; important in maintenance.
- The **test environment** is (implicitly) subject of testing; errors in additional components and tools may show up.
- System behaviour (efficiency, usability) becomes **visible**, even if not explicitly subject of a test.

Disadvantages:

- A proof of correctness is practically **impossible**, tests are seldomly **exhaustive**.

Advantages of testing (in particular over inspection):

- Testing is a **"natural"** checking procedure: "everybody can test".
- The systematic test is **reproducible** and **objective** (if the start configuration is reproducible and the test environment deterministic).
- Invested effort can be **re-used**: properly prepared and documented tests can be re-executed with low effort, in particular fully automatic tests; important in maintenance.
- The **test environment** is (implicitly) subject of testing; errors in additional components and tools may show up.
- System behaviour (efficiency, usability) becomes **visible**, even if not explicitly subject of a test.

Disadvantages:

- A proof of correctness is practically **impossible**, tests are seldomly **exhaustive**.
- It can be **extremely hard** to provide environment conditions like interrupts or critical timings (two buttons pressed at the same time").

Advantages of testing (in particular over inspection):

- Testing is a **"natural"** checking procedure: "everybody can test".
- The systematic test is **reproducible** and **objective** (if the start configuration is reproducible and the test environment deterministic).
- Invested effort can be **re-used**: properly prepared and documented tests can be re-executed with low effort, in particular fully automatic tests; important in maintenance.
- The **test environment** is (implicitly) subject of testing; errors in additional components and tools may show up.
- System behaviour (efficiency, usability) becomes **visible**, even if not explicitly subject of a test.

Disadvantages:

- A proof of correctness is practically **impossible**, tests are seldomly **exhaustive**.
- It can be **extremely hard** to provide environment conditions like interrupts or critical timings (two buttons pressed at the same time").
- Other **properties** of the implementation (like readability, maintainability) are not subject of the tests (but, e.g., of reviews).

Advantages of testing (in particular over inspection):

- Testing is a **"natural"** checking procedure: "everybody can test".
- The systematic test is **reproducible** and **objective** (if the start configuration is reproducible and the test environment deterministic).
- Invested effort can be **re-used**: properly prepared and documented tests can be re-executed with low effort, in particular fully automatic tests; important in maintenance.
- The **test environment** is (implicitly) subject of testing; errors in additional components and tools may show up.
- System behaviour (efficiency, usability) becomes **visible**, even if not explicitly subject of a test.

Disadvantages:

- A proof of correctness is practically **impossible**, tests are seldomly **exhaustive**.
 - It can be **extremely hard** to provide environment conditions like interrupts or critical timings (two buttons pressed at the same time").
 - Other **properties** of the implementation (like readability, maintainability) are not subject of the tests (but, e.g., of reviews).
 - Tests tend to focus only on the code, **other artefacts** (documentation, etc.) are hard to test. (Some say, developers tend to focus (too much) on coding, anyway).
- Recall: some agile methods turn this into a feature: there's only requirements, tests, and code, the positive result may be false, caused by flawed test gear.

Advantages of testing (in particular over inspection):

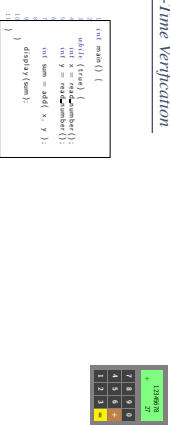
- Testing is a **"natural"** checking procedure: "everybody can test".
- The systematic test is **reproducible** and **objective** (if the start configuration is reproducible and the test environment deterministic).
- Invested effort can be **re-used**: properly prepared and documented tests can be re-executed with low effort, in particular fully automatic tests; important in maintenance.
- The **test environment** is (implicitly) subject of testing; errors in additional components and tools may show up.
- System behaviour (efficiency, usability) becomes **visible**, even if not explicitly subject of a test.

Disadvantages:

- A proof of correctness is practically **impossible**, tests are seldomly **exhaustive**.
 - It can be **extremely hard** to provide environment conditions like interrupts or critical timings (two buttons pressed at the same time").
 - Other **properties** of the implementation (like readability, maintainability) are not subject of the tests (but, e.g., of reviews).
 - Tests tend to focus only on the code, **other artefacts** (documentation, etc.) are hard to test. (Some say, developers tend to focus (too much) on coding, anyway).
- Recall: some agile methods turn this into a feature: there's only requirements, tests, and code, the positive result may be false, caused by flawed test gear.

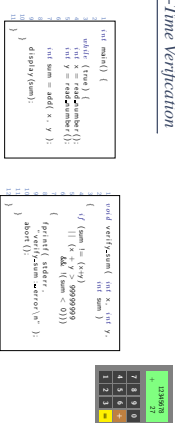
Run-Time Verification

Run-Time Verification



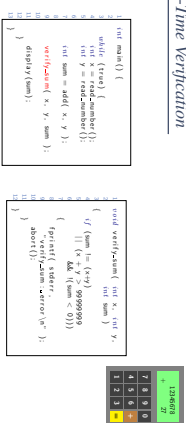
- If we have an **implementation** for checking whether an output is correct wrt. a given input (according to requirements),
 - we can just embed this **implementation** into the actual software, and
 - thereby **check satisfaction** of the requirement during **each run**.
- **run-time verification**

Run-Time Verification



- If we have an **implementation** for checking whether an output is correct wrt. a given input (according to requirements),
 - we can just embed this **implementation** into the actual software, and
 - thereby **check satisfaction** of the requirement during **each run**.
- **run-time verification**

Run-Time Verification



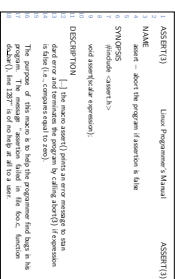
- If we have an **implementation** for checking whether an output is correct wrt. a given input (according to requirements),
 - we can just embed this **implementation** into the actual software, and
 - thereby **check satisfaction** of the requirement during **each run**.
- **run-time verification**

Simplest Case: Assertions

- Maybe the simplest instance of runtime verification: **Assertions**.
- Available in standard libraries of many programming languages, e.g. C:

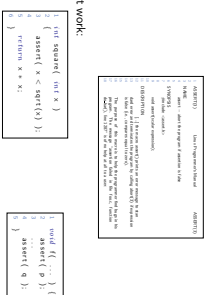
Simplest Case: Assertions

- Maybe the simplest instance of runtime verification: **Assertions**.
- Available in standard libraries of many programming languages, e.g. C:

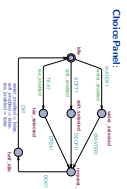


Simplest Case: Assertions

- Maybe the simplest instance of runtime verification: **Assertions**.
- Available in standard libraries of many programming languages, e.g. C:

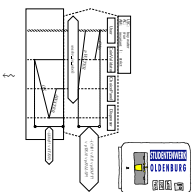
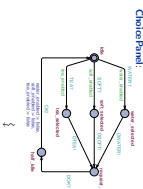


More Complex Case: LSC Observer

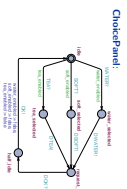


50/65

More Complex Case: LSC Observer

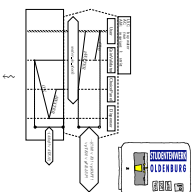
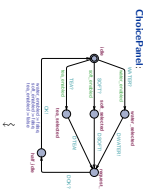


More Complex Case: LSC Observer

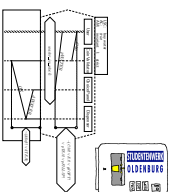
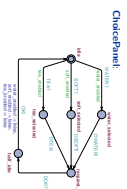


50/65

More Complex Case: LSC Observer

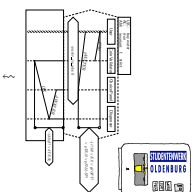
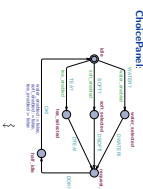


More Complex Case: LSC Observer



50/65

More Complex Case: ISC Observer

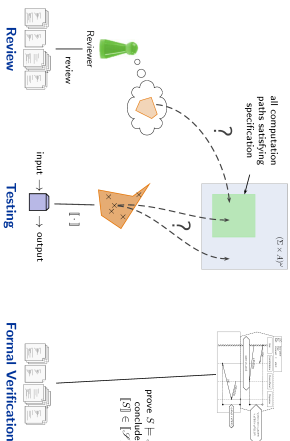


- **Experience:**
During development, **assertions** for pre/post conditions and intermediate invariants are an extremely powerful tool with **very good gain/error ratio** (low effort, high gain).
 - Effectively work as safe-guard against unexpected use of functions and regression, e.g. during later maintenance or efficiency improvement.
 - Can serve as **formal** (support of) documentation.
 - User reader: at this point in the program, I expect this condition to hold, because... "

- **Experience:**
During development, **assertions** for pre/post conditions and intermediate invariants are an extremely powerful tool with **very good gain/error ratio** (low effort, high gain).
 - Effectively work as safe-guard against unexpected use of functions and regression, e.g. during later maintenance or efficiency improvement.
 - Can serve as **formal** (support of) documentation.
 - User reader: at this point in the program, I expect this condition to hold, because... "
- **Usually:**
Development version **with** (cf. assert(3)) / release version **without** run-time verification.
If run-time verification enabled in release version,
 - software should terminate as gracefully as possible (e.g. try to save data),
 - save information from assertion failure if possible.

- **Experience:**
During development, **assertions** for pre/post conditions and intermediate invariants are an extremely powerful tool with **very good gain/error ratio** (low effort, high gain).
 - Effectively work as safe-guard against unexpected use of functions and regression, e.g. during later maintenance or efficiency improvement.
 - Can serve as **formal** (support of) documentation.
 - User reader: at this point in the program, I expect this condition to hold, because... "
- **Usually:**
Development version **with** (cf. assert(3)) / release version **without** run-time verification.
If run-time verification enabled in release version,
 - software should terminate as gracefully as possible (e.g. try to save data),
 - save information from assertion failure if possible.
- Run-time verification can be arbitrarily complicated and complex, e.g., construction of observers for LSCs or temporal logic, e.g., expensive checking of data, etc.

- **Experience:**
During development, **assertions** for pre/post conditions and intermediate invariants are an extremely powerful tool with **very good gain/error ratio** (low effort, high gain).
 - Effectively work as safe-guard against unexpected use of functions and regression, e.g. during later maintenance or efficiency improvement.
 - Can serve as **formal** (support of) documentation.
 - User reader: at this point in the program, I expect this condition to hold, because... "
- **Usually:**
Development version **with** (cf. assert(3)) / release version **without** run-time verification.
If run-time verification enabled in release version,
 - software should terminate as gracefully as possible (e.g. try to save data),
 - save information from assertion failure if possible.
- Run-time verification can be arbitrarily complicated and complex, e.g., construction of observers for LSCs or temporal logic, e.g., expensive checking of data, etc.
- **Drawback:** development and release software have different computation paths — with bad luck, the software only behaves well **because** of the run-time verification code. -

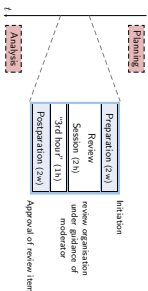


- **Review Item:** can be every closed, human-readable part of software (document, module, test data, installation manual, etc.)
Social aspect: it is an **artifact** which is examined, not the **human** (who created it).

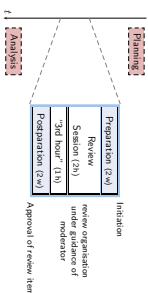
- **Review Item:** can be every closed, human-readable part of software (document, module, test data, installation manual, etc.)
Social aspect: it is an **artifact** which is examined, not the **human** (who created it).
- **Input to Review Session:**
 - the review item, and reference documents which enable an assessment (requirements specification, guidelines (e.g. coding conventions), catalogue of questions ("all variables initialised?"), etc.)

- **Review Item:** can be every closed, human-readable part of software (document, module, test data, installation manual, etc.)
Social aspect: it is an **artifact** which is examined, not the **human** (who created it).
- **Input to Review Session:**
 - the review item, and reference documents which enable an assessment (requirements specification, guidelines (e.g. coding conventions), catalogue of questions ("all variables initialised?"), etc.)
- **Roles:**
 - Moderator:** leads session, responsible for properly conducted procedure.
 - Author:** (representative of the) creator(s) of the artifact under review, is present to listen to the discussions, can answer questions, does not speak up if not asked.
 - Reviewer(s):** person who is able to judge the artifact under review; maybe different reviewers for different aspects (programming, tool usage, etc.), at least experienced in detecting inconsistencies or incompleteness.
 - Transcript Writer:** keeps minutes of review session, can be assumed by author.
- The review team consists of everybody but the author(s).

Review Procedure



Review Procedure



- review triggered e.g. by submission to revision control system;
- moderator invites (include review item in invitation), states review missions;
- **preparation:** reviewers investigate review item;
- **review session:** reviewers report, evaluate and document issues; solve open questions;
- **"3rd hour":** time for informal chat, reviewers may state proposals for solutions or improvements;
- **postparation:** wrap-up; responsibility of author(s);
- reviewers re-assess reviewed review item (until approval)
- **planning:** reviews need time in project plan; **analysis:** improve development and review process.

Review Rules (Indewig and Löhner, 2013)

- moderator organises, invites to, conducts review,**
- the review session is limited to 2 hours** — (if needed: more sessions)
- moderator may terminate review if condition not possible** (inputs, preparation, or people missing).
- the review item is under review, not the author(s).**
 - reviewers choose their wording accordingly;
 - authors neither defend themselves nor the review item;
 - roles are not mixed up; the moderator does not act as reviewer;
- style issues** (outside fixed conventions) are not discussed.
- the review team is not supposed to develop solutions;** issues are not noted in form of tasks for the author(s).
- each reviewer gets the opportunity to present her/his findings appropriately;**
- reviewers need to reach consensus on issues; consensus is noted down;**
- issues are classified as:** **critical** (review unusable for purpose), **major** (usability severely affected), **minor** (usability hardly affected), **good** (no problem).
- review team decides: accept without changes, accept with changes, do not accept.**
- protocol** is signed by all participants.

- **Review**
 - **Design and Code Inspection** (Fagan, 1976, 1989)
 - deluxe variant of review,
 - approx. 50% more time, approx. 50% more faults found.
- 57/68

- **Review**
 - **Design and Code Inspection** (Fagan, 1976, 1989)
 - deluxe variant of review,
 - approx. 50% more time, approx. 50% more faults found.
- 57/68

- **Structured Walkthrough**
 - simple variant of review: developer moderates walkthrough-session, presents artifact, reviewer pose (prepared or spontaneous) questions, issues are noted down, review is moderated and documented (preparation (do reviewers see the artifact before the session?) less effort, less effective.
 - **Review**
 - **Design and Code Inspection** (Fagan, 1976, 1989)
 - deluxe variant of review,
 - approx. 50% more time, approx. 50% more faults found.
- 57/68

- **Comment** ('Stellungnahme')
 - colleague(s) of developer read artifacts,
 - developer considers feedback,
 - **advantage:** low organisational effort; **disadvantages:** choice of colleagues may be biased; no protocol; consideration of comments at discretion of developer.
 - **Structured Walkthrough**
 - simple variant of review: developer moderates walkthrough-session, presents artifact, reviewer pose (prepared or spontaneous) questions, issues are noted down, review is moderated and documented (preparation (do reviewers see the artifact before the session?) less effort, less effective.
 - **disadvantages:** unclear responsibilities; "salaman"-author may trick reviewers
 - **Review**
 - **Design and Code Inspection** (Fagan, 1976, 1989)
 - deluxe variant of review,
 - approx. 50% more time, approx. 50% more faults found.
- 57/68

- **Careful Reading** ('Durchsicht')
 - done by developer,
 - recommendation: "way from screen" (use print-out or different device and situation)
 - **Comment** ('Stellungnahme')
 - colleague(s) of developer read artifacts,
 - developer considers feedback,
 - **advantage:** low organisational effort; **disadvantages:** choice of colleagues may be biased; no protocol; consideration of comments at discretion of developer.
 - **Structured Walkthrough**
 - simple variant of review: developer moderates walkthrough-session, presents artifact, reviewer pose (prepared or spontaneous) questions, issues are noted down, review is moderated and documented (preparation (do reviewers see the artifact before the session?) less effort, less effective.
 - **disadvantages:** unclear responsibilities; "salaman"-author may trick reviewers
 - **Review**
 - **Design and Code Inspection** (Fagan, 1976, 1989)
 - deluxe variant of review,
 - approx. 50% more time, approx. 50% more faults found.
- 57/68

- **Careful Reading** ('Durchsicht')
 - done by developer,
 - recommendation: "way from screen" (use print-out or different device and situation)
 - **Comment** ('Stellungnahme')
 - colleague(s) of developer read artifacts,
 - developer considers feedback,
 - **advantage:** low organisational effort; **disadvantages:** choice of colleagues may be biased; no protocol; consideration of comments at discretion of developer.
 - **Structured Walkthrough**
 - simple variant of review: developer moderates walkthrough-session, presents artifact, reviewer pose (prepared or spontaneous) questions, issues are noted down, review is moderated and documented (preparation (do reviewers see the artifact before the session?) less effort, less effective.
 - **disadvantages:** unclear responsibilities; "salaman"-author may trick reviewers
 - **Review**
 - **Design and Code Inspection** (Fagan, 1976, 1989)
 - deluxe variant of review,
 - approx. 50% more time, approx. 50% more faults found.
- 57/68



Concluding Discussion

Techniques Revisited

	auto- matic	prove "can run"	toolchain considered	exhaustive	prove correct	partial results	entry cost
Test	(✓)						(✓)
Runtime- Verification							
Review							
Static Checking							
Verification							

- Strengths:**
- fully automatic (yet not easy for CLI programs);
 - negative test proves "program not completely broken"; "can run" (or positive scenarios);
 - final product is examined, thus toolchain and platform considered;
 - one can stop at any time and take partial results;
 - few, simple test cases are usually easy to obtain;
 - provides reproducible counter-examples (good starting point for repair).
- Weaknesses:**
- often produces *justly incomplete* (no proof of correctness);
 - creating test cases for complex functions (or complex conditions) can be difficult;
 - maintaining many, complex test cases be challenging;
 - executing many tests may need substantial time (but can be run in parallel).

Techniques Revisited

	auto- matic	prove "can run"	toolchain considered	exhaustive	prove correct	partial results	entry cost
Test	(✓)						
Runtime- Verification	(✓)	(✓)		(✗)	✗	✓	(✓)
Review							
Static Checking	✗	✗	✗	(✓)	(✓)	✓	(✓)
Verification							

- Strengths:**
- human readers can *understand* the code, may spot point errors;
 - reported to be highly effective;
 - one can stop at any time and take partial results;
 - intermediate entry costs, good effort/effect ratio achievable.
- Weaknesses:**
- no tool support;
 - no results on actual execution, toolchain not reviewed;
 - human readers may *overlook* errors, usually not aiming at proofs;
 - does (in general) not provide counter-examples, developers may deny existence of error.

Techniques Revisited

	auto- matic	prove "can run"	toolchain considered	exhaustive	prove correct	partial results	entry cost
Test	(✓)	(✓)					
Runtime- Verification	(✓)	(✓)					
Review				(✗)	✗	✓	(✓)
Static Checking							
Verification							

- Strengths:**
- highly automatic (once changes are in place);
 - provides reproducible counter-examples (good starting point for repair);
 - (finally) final product is examined, thus toolchain and platform considered;
 - one can stop at any time and take partial results;
 - assert-statements have a very good effort/effect ratio.
- Weaknesses:**
- may negatively affect performance;
 - code is changed, program may only run *because* of the observers;
 - completeness depends on usage, may also be vastly incomplete, so no correctness proofs;
 - constructing observers for complex properties may be difficult, one needs to learn how to construct observers.

Techniques Revisited

	auto- matic	prove "can run"	toolchain considered	exhaustive	prove correct	partial results	entry cost
Test	(✓)	(✓)					
Runtime- Verification	(✓)	(✓)		(✗)	✗	✓	(✓)
Review							
Static Checking	✓	(✗)	✗	(✓)	(✓)	✓	(✓)
Verification	(✓)	✗	✗	✓	✓	(✗)	✗

- Strengths:**
- some tool support available (few commercial tools);
 - reported to be highly effective on large language semantics, platform, etc.);
 - can provide correctness proofs;
 - can be more efficient than other techniques.
- Weaknesses:**
- no results on actual execution, toolchain not reviewed;
 - not many intermediate results, "half of a proof" may not allow any useful conclusions;
 - proving things is difficult, failing to find a proof does not allow any useful conclusion;
 - false negatives (broken program "proved" correct) hard to detect.

Concluding Recommendations

- Not having at least one (systematic) test for each feature is **(grossly?) negligent**. IOW, without at least one test for each feature, **it is not software engineering**.

General Guidelines: Do's and Don'ts

- Do not use special **examination versions** for examination. (Test-cases, stubs, etc. can be used; yet may have errors which may undermine results.)

General Guidelines: Do's and Don'ts

- Do not use special **examination versions** for examination. (Test-cases, stubs, etc. can be used; yet may have errors which may undermine results.)
 - Do not **stop examination** when first error is detected.
- Clear:** Examination can (and should) be aborted if the examined program is not executable at all.

General Guidelines: Do's and Don'ts

- Do not use special **examination versions** for examination. (Test-cases, stubs, etc. can be used; yet may have errors which may undermine results.)
 - Do not **stop examination** when first error is detected.
- Clear:** Examination can (and should) be aborted if the examined program is not executable at all.
- Do not **modify** the artefact under examination during examinatio.
 - changes/corrections during examination in the end unclear **what exactly** has been examined ("moving target"). (results need to be uniquely traceable to one artefact version.)
 - fundamental flaws sometimes easier to detect with a **complete picture** of unsuccessful/successful tests.
 - **changes are particularly error-prone**, should not happen "en passant" in examination.
 - fixing flaws during examination may cause them to go uncounted in the **statistics** (which we need for all kinds of estimation).
 - roles **developer** and **examiner** are different anyway; an **examiner** fixing flaws would violate the role assignment.

General Guidelines: Do's and Don'ts

- Do not use special **examination versions** for examination. (Test-cases, stubs, etc. can be used; yet may have errors which may undermine results.)
 - Do not **stop examination** when first error is detected.
- Clear:** Examination can (and should) be aborted if the examined program is not executable at all.
- Do not **modify** the artefact under examination during examinatio.
 - changes/corrections during examination in the end unclear **what exactly** has been examined ("moving target"). (results need to be uniquely traceable to one artefact version.)
 - fundamental flaws sometimes easier to detect with a **complete picture** of unsuccessful/successful tests.
 - **changes are particularly error-prone**, should not happen "en passant" in examination.
 - fixing flaws during examination may cause them to go uncounted in the **statistics** (which we need for all kinds of estimation).
 - roles **developer** and **examiner** are different anyway; an **examiner** fixing flaws would violate the role assignment.
- In particular: Do not switch (fine grained) between **examination and debugging**.

So All Hope is Lost...?

- Seems like computer systems more or less inevitably have errors.

So All Hope is Lost...?

- Seems like computer systems more or less inevitably have errors.
- So **why** does my (heavily computerised) Airbus fly at all?

62.66

So All Hope is Lost...?

- Seems like computer systems more or less inevitably have errors.
- So **why** does my (heavily computerised) Airbus fly at all?
 - Firstly, aerospace software maybe has the lowest error rate of all softwares due to **very careful development**, **very thorough analysis** (e.g. fault tree analysis), and **strong regulatory obligations** ("no proof of correctness, no take-off").

62.66

So All Hope is Lost...?

- Seems like computer systems more or less inevitably have errors.
- So **why** does my (heavily computerised) Airbus fly at all?
 - Firstly, aerospace software maybe has the lowest error rate of all softwares due to **very careful development**, **very thorough analysis** (e.g. fault tree analysis), and **strong regulatory obligations** ("no proof of correctness, no take-off").
 - **Plus**: classical engineering wisdom for high reliability: **Redundancy**.
- Highly-critical components may be present 3-times redundant, developed by 3 different teams, compiled by 3 different complers, running on 3 different platforms, ...

62.66

So All Hope is Lost...?

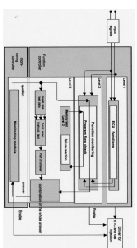
- Seems like computer systems more or less inevitably have errors.
- So **why** does my (heavily computerised) Airbus fly at all?
 - Firstly, aerospace software maybe has the lowest error rate of all softwares due to **very careful development**, **very thorough analysis** (e.g. fault tree analysis), and **strong regulatory obligations** ("no proof of correctness, no take-off").
 - **Plus**: classical engineering wisdom for high reliability: **Redundancy**.
- Highly-critical components may be present 3-time redundant, developed by 3 different teams, compiled by 3 different complers, running on 3 different platforms, ...
- And **why** does my (heavily computerised) car, infusion pump, etc. **not do harm**?

62.66

So All Hope is Lost...?

- Seems like computer systems more or less inevitably have errors.
- So **why** does my (heavily computerised) Airbus fly at all?
 - Firstly, aerospace software maybe has the lowest error rate of all softwares due to **very careful development**, **very thorough analysis** (e.g. fault tree analysis), and **strong regulatory obligations** ("no proof of correctness, no take-off").
 - **Plus**: classical engineering wisdom for high reliability: **Redundancy**.
- Highly-critical components may be present 3-times redundant, developed by 3 different teams, compiled by 3 different complers, running on 3 different platforms, ...
- And **why** does my (heavily computerised) car, infusion pump, etc. **not do harm**?
- Again, classical engineering wisdom for high reliability: **Run-time monitoring**.

62.66



Proposai: Dependability Cases (Jackson, 2009)

- A **dependable** system is one you can **depend** on — that is, you can place your trust in it.

63.66

- A **dependable** system is one you can **depend** on — that is, you can place your trust in it.
- **Proposed Approach:**

- A **dependable** system is one you can **depend** on — that is, you can place your trust in it.
- **Proposed Approach:**
- Identify the **critical requirements**, and determine what **level of confidence** is needed. Most systems do also have **non-critical** requirements.

- A **dependable** system is one you can **depend** on — that is, you can place your trust in it.
- **Proposed Approach:**
- Identify the **critical requirements**, and determine what **level of confidence** is needed. Most systems do also have **non-critical** requirements.
- Construct a **dependability case**:
- an argument, that the software, in concert with other components, establishes the critical properties.

- A **dependable** system is one you can **depend** on — that is, you can place your trust in it.
- **Proposed Approach:**
- Identify the **critical requirements**, and determine what **level of confidence** is needed. Most systems do also have **non-critical** requirements.
- Construct a **dependability case**:
- an argument, that the software, in concert with other components, establishes the critical properties.
- The case should be
- **auditable**: can (easily) be evaluated by third-party certifier.
- **complete**: no holes in the argument, any assumptions that are not justified should be noted (e.g. assumptions on compiler, on protocol obeyed by users, etc.)
- **sound**: e.g. should not claim full correctness [...] based on nonexhaustive testing; should not make unwarranted assumptions on independence of component failures, etc.

- A **dependable** system is one you can **depend** on — that is, you can place your trust in it.
- **Proposed Approach:**
- Identify the **critical requirements**, and determine what **level of confidence** is needed. Most systems do also have **non-critical** requirements.
- Construct a **dependability case**:
- an argument, that the software, in concert with other components, establishes the critical properties.
- The case should be
- **auditable**: can (easily) be evaluated by third-party certifier.
- **complete**: no holes in the argument, any assumptions that are not justified should be noted (e.g. assumptions on compiler, on protocol obeyed by users, etc.)
- **sound**: e.g. should not claim full correctness [...] based on nonexhaustive testing; should not make unwarranted assumptions on independence of component failures, etc.
- **IDW**: "Developers [should] express the critical properties and make an explicit argument that the system satisfies them." (As opposed to, e.g. requiring term coverage (which is usually not exhaustive), or requiring **only** coding conventions and procedure models, which may support, but do not prove dependability)

References

References

- Fagan, M. (1976). Design and code inspections to reduce errors in program development. *IBM Systems Journal*, 15(3)182-211.
- Fagan, M. (1986). Advances in software inspections. *IEEE Transactions On Software Engineering*, 12(7)744-751.
- IEEE (1990). *IEEE Standard Glossary of Software Engineering Terminology*. Std 610.1:1990.
- Jackson, D. (2009). A direct path to dependable software. *Comm. ACM*, 52(4).
- Lettrani, M. and Kiese, J. (2001). Scenario-based monitoring and testing of real-time UML models. In Gagola, M. and Kobryn, C., editors, *UML*, number 2185 in Lecture Notes in Computer Science, pages 317-328. Springer-Verlag.
- Ludewig, J. and Lichte, H. (2013). *Software Engineering*. dpunkt.verlag, 3. edition.