
CSE 331

Software Design & Implementation

James Wilcox

Autumn 2021

Testing

How do we ensure correctness?

Best practice: use three techniques

1. **Tools**

- e.g., type checking, @Override, libraries, etc.

2. **Inspection**

- think through your code carefully
- have another person review your code

3. **Testing**

- usually >50% of the work in building software

Each removes $\sim 2/3$ of bugs. Together >97%

What can you learn from testing?

“Program testing can be used to show the presence of bugs, but never to show their absence!”

Edsger Dijkstra

Notes on Structured Programming,
1970



Testing is *essential* but it is insufficient by itself

Only **reasoning** can prove there are no bugs. Yet...

How do we ensure correctness?



“Beware of bugs in the above code;
I have only proved it correct, not tried it.”
-Donald Knuth, 1977

Trying it is a surprisingly useful way to find mistakes!

No **single activity** or approach can guarantee correctness

We need tools **and** inspection **and** testing to ensure correctness

Why you will care about testing

In all likelihood, you will be expected to **test your own code**

- Industry-wide trend toward developers doing more testing
 - 20 years ago we had large test teams
 - now, test teams are small to nonexistent
- Reasons for this change:
 1. easy to update products after shipping (users are testers)
 2. often lowered quality expectations (startups, games)
 - some larger companies want to be more like startups

This has positive and negative effects...

It's hard to test your own code

Your **psychology** is fighting against you:

- confirmation bias
 - tendency to avoid evidence that you're wrong
- operant conditioning
 - programmers get cookies when the code works
 - testers get cookies when the code breaks

You can avoid some effects of confirmation bias by

writing most of your tests before the code

Not much you can do about operant conditioning

Kinds of testing

- Testing field has terminology for different kinds of tests
 - we won't discuss all the kinds and terms
- Here are three orthogonal dimensions [so 8 varieties total]:
 - *unit* testing versus *system/integration/end-to-end* testing
 - one module's functionality versus pieces fitting together
 - *opaque* testing versus *transparent* testing
 - did you look at the code before writing the test?
 - *specification* testing versus *implementation* testing
 - test only behavior guaranteed by specification or other behavior expected for the implementation

Unit Testing

- A unit test focuses on one class / module (or even less)
 - could write a unit test for a single method
- Tests a single unit in isolation from all others
- Integration tests verify that the modules fit together properly
 - usually don't want these until the units are well tested
 - i.e., unit tests come first

How is testing done?

Write the test

- 1) Choose input / configuration
- 2) Define the expected outcome

Run the test

- 3) Run with input and record the actual outcome
- 4) Compare *actual* outcome to *expected* outcome

What's So Hard About Testing?

“Just try it and see if it works...”

```
// requires:  $1 \leq x, y, z \leq 100,000$   
// returns: computes some  $f(x, y, z)$   
int func1(int x, int y, int z) {...}
```

Exhaustive testing would require 1 quadrillion cases!

- impractical even for this trivially small problem

Key problem: choosing test suite

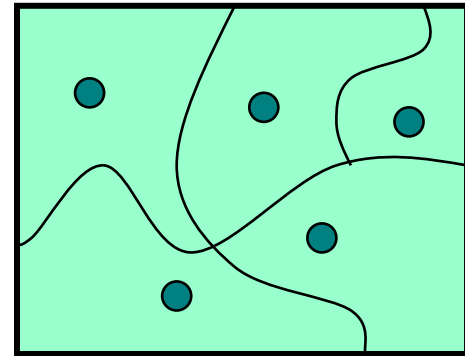
- Large/diverse enough to provide a useful amount of validation
- (Small enough to write/run in reasonable amount of time.)
 - less important... very few projects have too many tests

Approach: Partition the Input Space

Ideal test suite:

Identify sets with “same behavior”
(actual and expected)

Test **at least** one input from each set
(we call this set a *subdomain*)



Two problems:

1. Notion of **same behavior** is subtle
 - Naive approach: **execution equivalence**
 - Better approach: **revealing subdomains**
2. Discovering the sets requires perfect knowledge
 - If we had it, we wouldn't need to test
 - Use heuristics to approximate cheaply

Naive Approach: Execution Equivalence

```
// returns:  x < 0      => returns -x
//           otherwise => returns  x
int abs(int x) {
    if (x < 0) return -x;
    else      return  x;
}
```

All $x < 0$ are **execution equivalent**:

- Program takes same sequence of steps for any $x < 0$

All $x \geq 0$ are execution equivalent

Suggests that $\{-3, 3\}$, for example, is a good test suite

Execution Equivalence Can Be Wrong

```
// returns:  x < 0      => returns -x
//           otherwise => returns  x

int abs(int x) {
    if (x < -2) return -x;
    else       return  x;
}
```

{-3, 3} does not reveal the error!

Two possible executions: $x < -2$ and $x \geq -2$

Three possible behaviors:

- $x < -2$ OK, $x = -2$ or $x = -1$ (BAD)
- $x \geq 0$ OK

Revealing Subdomains

- A *subdomain* is a subset of possible inputs
- A subdomain is *revealing* for error E if either:
 - *every* input in that subdomain triggers error E , *or*
 - *no* input in that subdomain triggers error E
- Need test at least one input from a revealing subdomain to find bug
 - if you test one input from every revealing subdomain for E , you are guaranteed to find the bug
- The trick is to *guess* revealing subdomains for **the errors present**
 - even though your reasoning says your code is correct, make educated guesses where the bugs might be

Testing Heuristics

- Testing is *essential* but difficult
 - want set of tests likely to reveal the bugs present
 - but we don't know where the bugs are
- Our approach:
 - split the input space into enough subsets (subdomains) such that inputs in each one are likely all correct or incorrect
 - can then take just one example from each subdomain
- Some heuristics are useful for choosing subdomains...

Heuristics for Designing Test Suites

A good heuristic gives:

- for all errors in some class of errors E : high probability that some subdomain is revealing for E
- not an *absurdly* large number of subdomains

Different heuristics target different classes of errors

- in practice, combine multiple heuristics
 - (we will see several)
- a way to think about and communicate your test choices

Specification Testing

Heuristic: Explore alternate cases in the specification

Procedure is **opaque**: specification visible, internals hidden

Example

```
// returns: a > b => returns a  
//          a < b => returns b  
//          a = b => returns a  
int max(int a, int b) {...}
```

3 cases lead to 3 tests

$(4, 3) \Rightarrow 4$ (*i.e. any input in the subdomain $a > b$*)

$(3, 4) \Rightarrow 3$ (*i.e. any input in the subdomain $a < b$*)

$(3, 3) \Rightarrow 3$ (*i.e. any input in the subdomain $a = b$*)

Specification Testing: Advantages

Process is not influenced by component being tested

- avoids psychological biases we discussed earlier
- can only do this for your own code if you **write tests first**

Robust with respect to changes in implementation

- test data need not be changed when code is changed

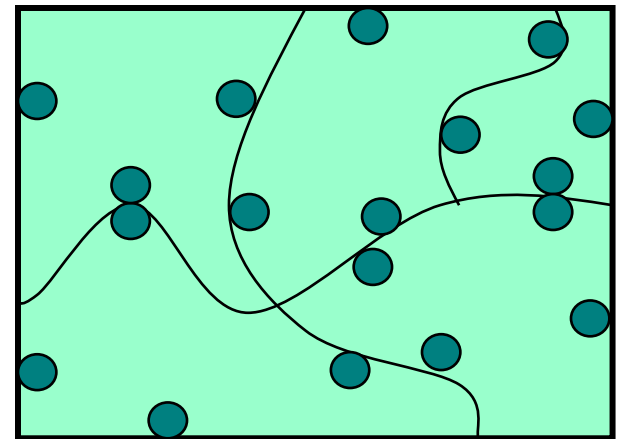
Allows others to test the code (rare nowadays)

Heuristic: Boundary Cases

Create tests at the boundaries between subdomains

Edges of the “main” subdomains have a high probability of revealing errors

- e.g., off-by-one bugs



Include one example **on each side** of the boundary

Also want to test the side edges of the subdomains...

Heuristic: Special Cases

Arithmetic

- zero
- smallest/largest values
- arithmetic overflow

Objects

- null
- zero elements
- same object passed as multiple arguments (aliasing)

All of these are common cases where bugs lurk

- you'll find more as you encounter more bugs

Heuristic: Transparent testing

Focus on features not described by specification

- control-flow details (e.g., conditions of “if” statements in code)
- performance optimizations
- alternate algorithms for different cases

Example: **abs** from before

- had different behavior > 2 and ≤ 2

Transparent Example

There are some subdomains that opaque testing won't catch:

```
boolean[] primeTable = new boolean[CACHE_SIZE];

boolean isPrime(int x) {
    if (x > CACHE_SIZE) {
        for (int i=2; i <= x/2; i++) {
            if (x % i == 0)
                return false;
        }
        return true;
    } else {
        return primeTable[x];
    }
}
```

Transparent Testing: [Dis]Advantages

- Finds an important class of boundaries
 - yields useful test cases
- Consider `CACHE_SIZE` in `isPrime` example
 - important tests `CACHE_SIZE-1`, `CACHE_SIZE`, `CACHE_SIZE+1`
 - if `CACHE_SIZE` is mutable, may need to test with different `CACHE_SIZE` values

Disadvantage:

- buggy code tricks you into thinking it's right once you look at it
 - (confirmation bias)
- can end up with tests having same bugs as implementation
- so also write tests **before** looking at the code

Pragmatics: Regression Testing

- Whenever you find a bug
 - store the input that elicited that bug, plus the correct output
 - add these to the test suite
 - verify that the test suite **fails**
 - fix the bug
 - verify the fix
- Ensures that your fix solves the problem
 - don't add a test that succeeded to begin with!
 - another reason to try to write tests before coding
- Protects against reversions that reintroduce bug
 - it happened at least once, and it might happen again (especially when trying to change the code in the future)

Summary of Heuristics

- Split subdomains on boundaries appearing in the specification
- Split subdomains on boundaries appearing in the implementation
- Test boundaries that commonly lead to errors
- Test special cases like nulls, empty arrays, 0, etc.
- Test any cases that caused bugs before (to avoid regression)

On the other hand, don't confuse *volume* with *quality* of tests

- look for revealing subdomains
- want tests in every revealing subdomain not **just** lots of tests