

## Why Tests?

Let's start by talking about what test suites can do for you (as a developer).

Reference: Kat Busch. "A beginner's guide to automated testing."

<https://hackernoon.com/treat-yourself-e55a7c522f71>

Back to the situation in Lecture 1. You write some new code, and want some assurance that it works. Passing code review at Dropbox at the time required tests along with the code.

"Lo and behold, I soon needed to fix a small bug." But, of course, it's easy to introduce even more bugs when fixing something. Fortunately, she had some tests.

I ran the tests. Within a few seconds, I knew that everything still worked! Not just a single code path (as in a manual test), but all code paths for which I'd written tests! It was magical. It was so much faster than my manual testing. And I knew I didn't forget to test any edge cases, since they were all still covered in the automated tests.

Not writing tests is incurring technical debt. You'll pay for it later, when you have to maintain the code. Having tests allows you to move faster later, without worrying about breaking your code.

If your code is still in the codebase a year (or five) after you've committed it and there are no tests for it, bugs will creep in and nobody will notice for a long time.

Writing tests is like eating your vegetables. It'll enable your code to go big and strong.

**If it matters that the code works you should write a test for it.** There is no other way you can guarantee it will work.

(We'll look at other ways in this course, but tests are the state of the industry.)

## Exploratory Testing

Exploratory testing is usually (but not always) carried out by dedicated testers. In that sense, it's somewhat different from the other testing activities in this course, which are more developer-focussed—our usual goal is learning, as developers, how to deploy better automated test suites for our software. Hallway usability testing, though, is an application of exploratory testing. Furthermore, the dedicated QA function is important, and we should learn about how it works.

**Resources.** James Bach has an introduction to exploratory testing:

- <https://www.satisfice.com/exploratory-testing>

There is an exhaustive set of notes on exploratory testing by Cem Kaner:

- <https://www.kaner.com/pdfs/QAExploring.pdf>

“Exploratory testing is simultaneous learning, test design, and test execution.”

Contrast this to scripted testing: test design happens ahead of time and then test execution happens (repeatedly) throughout the product’s development cycle. When we think of dedicated QA teams, we think they are manually executing scripted tests. In 2025, that is not an effective use of staff.

There is a continuum between scripted testing and exploratory testing. Good exploratory testing may use prepared scripts for certain tasks.

**Scenarios where Exploratory Testing Excels.** (from Bach’s article)

- providing rapid feedback on new product/feature;
- learning product quickly;
- diversifying testing beyond scripts;
- finding single most important bug in shortest time;
- independent investigation of another tester’s work;
- investigating and isolating a particular defect;
- investigate status of a particular risk to evaluate need for scripted tests.

**Exploratory Testing Process.** Exploratory testing should not be randomly bumbling around (we can call that “ad hoc testing”)—the random approach finds bugs but isn’t the most efficient at giving you an idea of how well the software works.

- Start with a charter for your testing activity, e.g. “Explore and analyze the product elements of the software.” These charters should be somewhat ambiguous.
- Decide what area of the software to test.
- Design a test (informally).
- Execute the test; log bugs.
- Repeat.

Exploratory testing shouldn’t produce an exhaustive set of notes. Good testers will be able to reproduce the bugs that they encounter during their testing from brief notes. Taking full notes takes too long.

The output from exploratory testing is at least a set of bug reports. It may also include test notes, which include overall impressions and a summary of the test strategy/thought process. Artifacts such as test data or test materials are also both inputs and outputs from exploratory testing.

**Primary vs contributing tasks.** One way to classify tasks that software can do (or, in other words, its features) is *primary* vs *contributing*. A *primary* task is core functionality of the system; it's something that you would say "You Had One Job!" about. As examples, text editors must be able to load text files, add text, and save the text files. On the other hand, *contributing* tasks are secondary. A macro system for a text editor would be a contributing task. Being able to read email in your editor is definitely a contributing task. Sometimes it's not black-and-white. Spell-check can go either way.

## In-class exercise: Exploratory testing of WaterlooWorks.

We will try out exploratory testing with WaterlooWorks. I believe that all non-exchange students here should have access to the system, although there may be no jobs visible right now.

The charter will be "Explore the overall functionality of WaterlooWorks". Summarize in one or two sentences what the purpose of WaterlooWorks is. Identify the tasks that WaterlooWorks should be able to do and classify them as primary or contributing. Identify areas of potential instability. Test each function and record results (bugs).

Of course, don't do things that have actual effects. Usually, testers would have access to a development server and could test those areas more aggressively. But we are working with production systems here.

## Regression Testing

Regression testing refers to any software testing that uncovers errors by retesting the modified program (Wikipedia). This form of testing often refers to comprehensive sets of test cases to detect regressions:

- of bug fixes that a developer has proposed.
- of related and unrelated other features that have been added.

Regression testing usually refers to system level (integration level) testing that runs the entire process.

### Attributes of Regression Tests

Regression tests usually have the following attributes:

- **Automated:** no real reason to have manual regression tests.
- **Appropriately Sized:** too small and bugs will be missed. Too large and they will take a long time to run. Optimally, we want to run tests continuously.
- **Up-to-date:** ensure that tests are valid for the version of program being tested.

## Automating Regression Tests

Regression tests often have a low yield in terms of finding bugs (and are boring to run). Automation is key.

### Input

If the input is from a file, regression tests are easy to run (but should still be automatically triggered on a regular interval). There may still be a problem with validating output. We can also create special mocks that can take input from a file or other sources (e.g. scripting engines).

For UIs, the standard approach is to capture and replay events. This approach can be fragile! For example, tests may fail based on window placement or whitespace. For web applications, there is capture and replay for HTTP using Selenium (see L02-extra notes). Mozilla has a project named Marionette<sup>1</sup> that is used to test Firefox and Thunderbird; it is like Selenium but also works on Chrome elements.

### Output

Verifying output can be hard!! Problems can arise from issues such as resolution, whitespace, window placement etc.

### Mozilla Case Study<sup>2</sup>

The case study presents an approach to testing Gecko based applications. Gecko is the layout engine for Mozilla applications like Firefox and Thunderbird.

In the past, a frame tree with coordinates for all UI elements was created and manual testers performed tests. Not optimal! The new approach was to capture screenshots after test cases and compare them with the expected screenshot. There were a few problems with the approach due to nondeterminism:

- Animated images: no way to ensure tests would function with animated images.
- Font Hinting: the same character would appear slightly differently after each run of the application.
- Other bugs: resolution problems, minor changes in layout etc.

The problem was “really hellish” and the partial solution was to enable logging in the application. The logs would essentially be compared with expected logs. This became very ugly given 1300 or so test cases; distributing across different computers helped.

## Industrial Best Practices

Some more references:

- Gerard Meszaros. *xUnit Test Patterns: Refactoring Test Code*.
- Kent Beck. *Test Driven Development: By Example*.

---

<sup>1</sup><https://developer.mozilla.org/en-US/docs/Mozilla/QA/Marionette>

<sup>2</sup>[http://robert.ocallahan.org/2005/03/visual-regression-tests\\_04.html](http://robert.ocallahan.org/2005/03/visual-regression-tests_04.html)

- Roy Oshero. *The Art of Unit Testing: with examples in C#*.
- **Unit Tests:** Each class has an associated unit test. If you change a class, you must also modify the unit test.
- **Code Reviews:** Each branch in the central version control system has owners. In order to commit code to that branch, you must have your code reviewed and approved by one of the owners. This ensures code quality.
- **Continuous Builds:** There is often a machine that continuously checks out and tests the latest code. All unit and regression tests are run and the status is made public to the team. The status contains information about the whether the code was built successfully, whether all unit and regression tests passed, and a list of the last few commits that were made to the branch. This ensures developers try to submit good code since if you break something, everyone knows about it :)
- **One-button Deploy:** If all tests have passed, one should be able to deploy to production with one command.
- **Back Button:** Systems should be designed so that it's possible to roll back changes.

## Unit tests

Unit tests are more low-level than integration tests and focus on one particular “class, module, or function”. They should execute quickly. Sometimes you need to create fake inputs (or mocks) for unit tests; we’ll talk about that too. You should generally not use an entire real input for a unit test.

There are many unit testing frameworks (e.g. JUnit, NUnit), though I don’t know of any dominant ones for C or C++. Here’s an NUnit test I found on the Internet<sup>3</sup>.

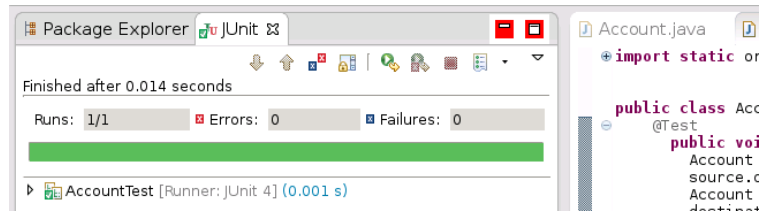
```
1 [Test]
2 public void GetMinimum_UnsortedIntegerArray_ReturnsSmallestValue()
3 {
4     var unsortedArray = new int[] {7,4,9,2,5}; // Arrange
5     var minimum = Statistics.GetMinimum(unsortedArray); // Act
6     Assert.AreEqual(2, minimum); // Assert
7 }
```

Note that tests have three phases: arrange, act, and assert. I’ll say a bit more about good test design, and could say even more, but I don’t think I will this term.

You’ll find that writing tests as you go makes your interfaces better and makes your code more testable. If you find yourself writing something hard to test, you’ll notice it early on when there’s still time to improve the design.

**Goal.** Well-designed tests are *self-checking*. That means that if the test runs with no errors and no failures (and hence produces a green bar in your IDE), we know that the test was successful.

<sup>3</sup><https://dzone.com/articles/the-anatomy-of-good-unit-testing>



Writing self-checking tests means that the tests automatically report the status of the code. This enables a “keep the bar green” coding style. Implications: 1) you can worry less about introducing bugs (but still take ordinary care); and 2) the tests help document your system’s specs.

**How to write self-checking tests.** One might think:

“Isn’t it just calling asserts?”

Sadly, no. That’s not enough.

Two questions about actually deploying asserts:

- Q: what for?  
A: check method call results
- Q: where?  
A: usually after calling SUT (System Under Test)

**Counter example.** Here’s some example code.

```
1 public class Counter {  
2     int count;  
3  
4     public int getCount() { return count; }  
5     public void addToCount(int n) { count += n; }  
6 }
```

We can test it with the following JUnit test.

```

1 // java -cp /usr/share/java/junit4.jar:. \
2 //   org.junit.runner.JUnitCore CounterTest
3 import static org.junit.Assert.*;
4 import org.junit.Test;
5
6 public class CounterTest {
7     @org.junit.Test
8     public void add10() {
9         Counter c = new Counter(); // arrange
10        c.addToCount(10); // act
11        // after calling SUT, read off results
12        int count = c.getCount();
13        assertEquals(10, count); // assert
14    }
15 }

```

What kind of test is this? Let's consider two kinds of tests: state-based tests vs. behaviour-based tests.

- **State:** e.g. object field values. Verify by calling accessor methods.
- **Behaviour:** which calls SUT makes. Verify by inserting observation points, monitoring interactions.

Does Counter Test verify state or behaviour?

**Flight example.** Here's more example code.

```

1 // Meszaros, p. 471
2 // not self-checking
3 public void testRemoveFlightLogging_NSC() {
4     // arrange:
5     FlightDto expectedFlightDto=createRegisteredFlight();
6     FlightMgmtFacade facade=new FlightMgmtFacadeImpl();
7     // act:
8     facade.removeFlight(expectedFlightDto.getFlightNo());
9     // assert:
10    // have not found a way to verify the outcome yet
11    // Log contains record of Flight removal
12 }

```

## Implementing State Verification

We can verify state:

```
1 // Meszaros, p. 471
2 // extended state specification
3 public void testRemoveFlightLogging_NSC() {
4     // arrange:
5     FlightDto expectedFlightDto=createRegisteredFlight();
6     FlightMgmtFacade=new FlightMgmtFacadeImpl();
7     // act:
8     facade.removeFlight(expectedFlightDto.getFlightNo());
9     // assert:
10    assertFalse("flight still exists after removed",
11                facade.flightExists(expectedFlightDto,
12                                     getFlightNo()));
13 }
```

Note that we are exercising the SUT, verifying state, and checking return values.

In state-based tests, we inspect only outputs, and only call methods from SUT. We do not instrument the SUT. We do not check interactions.

You have two options for verifying state:

1. procedural (bunch of asserts); or,
2. via expected objects (won't talk about them this year).

Returning to the flight example:

- We do check that the flight got removed.
- We don't check that the removal got logged.
- Hard to check state and observe logging.
- Solution: Spy on SUT behaviour.



## Implementing Procedural Behaviour Verification

Or, we can implement behaviour verification. This is one way to do so, behaviourally:

```
1 // Meszaros, p. 472
2 // procedural behaviour verification
3 public void testRemoveFlightLogging_PBV() {
4     // arrange:
5     FlightDto expectedFlightDto=createRegisteredFlight();
6     FlightMgmtFacade=new FlightMgmtFacadeImpl();
7     AuditLogSpy logSpy = new AuditLogSpy();
8     facade.setAuditLog(logSpy);
9     // act:
10    facade.removeFlight(expectedFlightDto.getFlightNo());
11    // assert:
12    assertEquals("number of calls",
13                1, logSpy.getNumberOfCalls());
14    // ...
15    assertEquals("detail",
16                expectedFlightDto.getFlightNumber(),
17                logSpy.getDetail());
18 }
```

As an alternative, we can use a mock object framework (e.g. JMock) to define expected behaviour.

**Idea.** Observe calls to the logger, make sure right calls happen.

## Assertions

We build tests using assertions. In JUnit, there are three basic built-in choices:

1. `assertTrue(aBooleanExpression)`
2. `assertEquals(expected, actual)`
3. `assertEquals(expected, actual, tolerance)`

(There are others too, but let's start with these.)

`assertTrue` is more flexible, since you can write anything with a boolean value. However, it can give hard-to-diagnose error messages—you need try harder when using it if you want good tests.

**Using Assertions.** Why use assertions? Assertions are good for:

- checking all the things that should be true (more = better);
- serving as documentation: when system in state  $S_1$ , and I do  $X$ , assert that the result should be  $R$ , and that system should be in  $S_2$ .
- allowing failure diagnosis (include assertion messages!)

There are alternatives to using assertions. For instance, one can also do external result verification:

- write output to files; and
- use diff (or your own custom diff) to compare expected and actual output.

The twist is that the expected result is then not visible when looking at test's source code. (What's a good workaround?)

**Verifying Behaviour.** The key is to observe actions (calls) of the SUT. Some options:

- procedural behaviour verification (the challenge in that case: recording and verifying behaviour); or
- expected behaviour specification (capturing the outbound calls of the SUT).

## Representation invariants

We said that assertions are good for checking all the things that should be true. In particular, they can help with checking the integrity of complex data structures; in such a case, we can aim to write out all the things that should be true of the data structure. Let's look at examples from the Fuzzing chapter of the Fuzzing Book<sup>4</sup>. The example isn't really a complex data structure, but it does have some expected properties. Your typical CS341 data structure is going to have more complex properties, but they are harder to illustrate.

```
1 airport_codes: dict[str, str] = {
2     "YVR": "Vancouver",
3     "JFK": "New York-JFK",
4     "CDG": "Paris-Charles de Gaulle",
5     "CAI": "Cairo",
6     "LED": "St. Petersburg",
7     "PEK": "Beijing",
8     "HND": "Tokyo-Haneda",
9     "AKL": "Auckland"
10 }
```

*# plus many more*

We can introduce a *representation invariant* for this dict:

```
1 def code_repOK(code: str) -> bool:
2     assert len(code) == 3, "Airport code must have three characters: "
3         + repr(code)
4     for c in code:
5         assert c.isalpha(), "Non-letter in airport code: " + repr(
6             code)
7         assert c.isupper(), "Lowercase letter in airport code: " +
8             repr(code)
9     return True
```

---

<sup>4</sup><https://www.fuzzingbook.org/html/Fuzzer.html>

and we can check it:

```
1 assert code_repOK("SEA")
```

which quietly does not fail. You can also invoke `code_repOK` on all of the codes in `airport_codes`, and in fact, you can make a function to do so:

```
1 def airport_codes_repOK():
2     for code in airport_codes:
3         assert code_repOK(code)
4     return True
```

If this is an abstract data type, you would usually use a function to add to it:

```
1 def add_new_airport(code: str, city: str) -> None:
2     assert code_repOK(code)
3     assert airport_codes_repOK()
4     airport_codes[code] = city
5     assert airport_codes_repOK()
```

which checks the data structure before-and-after as well as the input. Might be slow.

Note that this checks properties that are *specific to your data structure*. I call them domain-specific properties. It's not just "program doesn't crash", it's something that encodes information about how your program is supposed to behave. It's a matter of taste to encode the right things.

Here's a collection of checks for a red-black tree:

```
1 class RedBlackTree:
2     def repOK(self):
3         assert self.rootHasNoParent()
4         assert self.rootIsBlack()
5         assert self.rootNodesHaveOnlyBlackChildren()
6         assert self.treeIsAcyclic()
7         assert self.parentsAreConsistent()
8         return True
9
10    def rootIsBlack(self):
11        if self.parent is None:
12            assert self.color == BLACK
13        return True
```

You might (automatically) run these checks every time you add or remove from the red-black tree, though an acyclicity check might be expensive and best to disable in production. This kind of check does work well with automatically-generated inputs from fuzzing, which we'll discuss in a few weeks.

## Test Doubles

Mock objects are a particular kind of test double. We need test doubles because objects collaborate with other objects, but we only want to test one object at a time. Meszaros categorizes test doubles as follows:

- dummy objects: these are not actually test doubles; they don't do anything, but just take up space in parameter lists. Are like `null`, but get past nullness checks in code.
- fake objects: have actual behaviour (which is correct), but somehow unsuitable for use in production; typical example is an in-memory database.
- stubs: produce canned answers in response to interactions from the class under test.
- mocks: like stubs, also produce canned answers. Difference: mock objects also check that the class under test makes the appropriate calls.
- spies: usually wraps the real object (instead of the mock, which stubs it), and records interactions for later verification.

Shorter reference about test doubles: [martinfowler.com/articles/mocksArentStubs.html](http://martinfowler.com/articles/mocksArentStubs.html)

## Mock Objects

Before we talk about mock objects, let's look at a stub. Imagine that you have a service that sends out emails. You don't actually want to send out emails while you're testing. So here's a class that pretends to send out emails.

```
1 public class MailServiceStub implements MailService {
2     private List<Message> messages = new ArrayList<Message>();
3     public void send (Message msg) {
4         messages.add(msg);
5     }
6     public int numberSent() {
7         return messages.size();
8     }
9 }
```

This stub permits *state verification*, as seen in the following assert in a test:

```
1 assertEquals(1, mailer.numberSent());
```

This is state verification because it's checking the contents of memory (which should reflect interactions that have happened in the past). One could also check the recipients, contents of messages, etc.

**jMock example.** Instead of state verification, we can also do behaviour verification. This is jMock syntax.

```

1 class OrderInteractionTester... {
2     public void testOrderSendsMailIfUnfilled() {
3         Order order = new Order(TALISKER, 51);
4         Mock warehouse = mock(Warehouse.class);
5         Mock mailer = mock(MailService.class);
6         order.setMailer((MailService) mailer.proxy());
7
8         mailer.expects(once()).method("send");
9         warehouse.expects(once()).method("hasInventory")
10             .withAnyArguments()
11             .will(returnValue(false));
12
13         order.fill((Warehouse) warehouse.proxy());
14     }
15 }

```

The calls to `mock()` create mock objects which have the appropriate type. If you are using the objects as simple dummy objects, calling `mock()` and `proxy()` is enough. Note that we have a real `Order` object but we're giving it the fake proxy objects, as created by the `Mock`'s `proxy()` methods.

We also specify the expected behaviour of the `mailer` and the `warehouse`. The test case is saying that the mailer ought to have `send()` called on it once, and that the warehouse ought to have `hasInventory()` called; that method should return `false()`.

**EasyMock example.** Different mock object libraries have different syntax. Here's another example, this time for EasyMock.

```

1 @RunWith(EasyMockRunner.class)
2 public class ExampleTest {
3
4     @TestSubject
5     private ClassUnderTest classUnderTest = new ClassUnderTest();
6
7     @Mock // creates a mock object
8     private Collaborator mock;
9
10    @Test
11    public void testRemoveNonExistingDocument() {
12        replay(mock);
13        classUnderTest.removeDocument("Does not exist");
14    }
15 }

```

Here we are testing the `ClassUnderTest` and creating a mock object of `Collaborator` type. EasyMock 2.3 reads the `@Mock` annotation and automatically fills in a mock object of the appropriate type. In our test case, we call `replay(mock)` to indicate that we are no longer recording expectations, but are instead starting the test case itself. In the above code, there are currently no expectations.

Let's add some expectations.

```

1  @Test
2  public void testAddDocument() {
3      // ** recording phase **
4      // expect document addition
5      mock.documentAdded("Document");
6      // expect to be asked to vote for document removal, and vote for it
7      expect(mock.voteForRemoval("Document"))
8          .andReturn((byte) 42);
9      // expect document removal
10     mock.documentRemoved("Document");
11     replay(mock);
12     // ** replaying phase ** we expect the recorded actions to happen
13     classUnderTest.addDocument("New Document", new byte[0]);
14     // check that the behaviour actually happened:
15     verify(mock);
16 }

```

Here we record the fact that the mock should be called with `documentAdded` and a parameter "New Document". We also record that the mock's `voteForRemoval` method should be called, and when that happens, it should return value 42. Finally, we add a call `verify()` to let EasyMock know that we're done and that it can go ahead and check that the expected behaviour actually happened.

## Flaky Tests

The second test engineering topic I want to talk about today is flaky tests. Flaky tests are those that sometimes fail (nondeterministically). Flakiness is not something you want in your test cases. (I have heard one defense of a flaky test: it lets you know that the system has the potential to actually work.) In general, flaky tests don't play well with the expectation that your test suite passes 100%.

Reference:

Qingzhou Luo, Farah Hariri, Lamyaa Eloussi, Darko Marinov. "An Empirical Analysis of Flaky Tests". In Proceedings of Foundations of Software Engineering '14.

**Dealing with flaky tests.** Companies with large test suites have found mitigations for the flaky test suite problem. One can label known-flaky tests as flaky and automatically re-run them to see if they eventually pass. One can also ignore or remove flaky tests. But this is unsatisfactory: it takes a long time to re-run failing tests.

**Causes of flakiness.** Luo et al studied 201 fixes to flaky tests in open-source projects. They found that the three most common causes of fixable flaky tests were:

1. improper waits for asynchronous responses;
2. concurrency; and
3. test order dependency.

The problem that caused flakiness for asynchronous waits was that there was typically a `sleep()`

call which didn't wait long enough for the action (perhaps a network call) to finish. The best practice is to use some sort of `wait()` call to wait for the result instead of hardcoding a sleep time.

Concurrency problems were what one might expect. The problem could either be in the system under test or in the test itself. Problems included data races, atomicity violations, and deadlocks; the solutions were the proper use of concurrency primitives (e.g. locks) as seen in CS343.

Test order dependency problems arose when some tests expected other tests to have already executed (and left a side effect like a file in the filesystem). They came up especially in the transition from Java 6 to Java 7 because that transition changed the (not-guaranteed) test execution order. The solution is to remove the dependency.

## **Not unit tests**

Many interesting things happen when different units interact. Integration tests verify end-to-end functionality. They're slower, flakier, and harder both to write and use. Focus on unit tests while developing code. But you eventually need integration tests to make sure the user sees the right thing.