

Test Suite Quality Metrics



One-Slide Summary

- **Test suite quality metrics** help us decide which suite to use. **Line coverage**, the fraction of lines visited when running a suite, is simple but gives limited confidence. **Branch coverage**, which requires both true and false values for conditionals, is richer (incorporating data values indirectly).
- **Mutation analysis** measures the fraction of **seeded defects** detected by a suite; it is expensive but effective.
- **Beta** and **A/B testing** involve real users and their experiences.

The Story So Far ...

- **Testing** is the most common **dynamic** technique for software quality assurance.
- Testing is **very expensive** (e.g., 35% of total IT spending). [Capgemini World Quality Report. 2015]
- Not testing, or testing badly, is **even more expensive**: [Minimizing code defects to improve software quality and lower development costs. IBM 2008] (Remember the IBM Story?)

Design and architecture	Implementation	Integration testing	Customer beta test	Postproduct release
1X*	5X	10X	15X	30X

*X is a normalized unit of cost and can be expressed in terms of person-hours, dollars, etc.

Source: National Institute of Standards and Technology (NIST)†

By catching defects as early as possible in the development cycle, you can significantly reduce your development costs.

Story Time

- Abboty Labs (St. Jude Medical) makes pacemakers
- In 2016, 465,000 of them were discovered to have security vulnerabilities

“The wireless protocol used for communication amongst St. Jude Medical cardiac devices has serious security vulnerabilities that make it possible to convert Merlin@home devices into **weapons capable of disabling therapeutic care and delivering shocks to patients at distances of 10 feet**, a range that could be extended using off-the-shelf parts to modify Merlin@home units.”

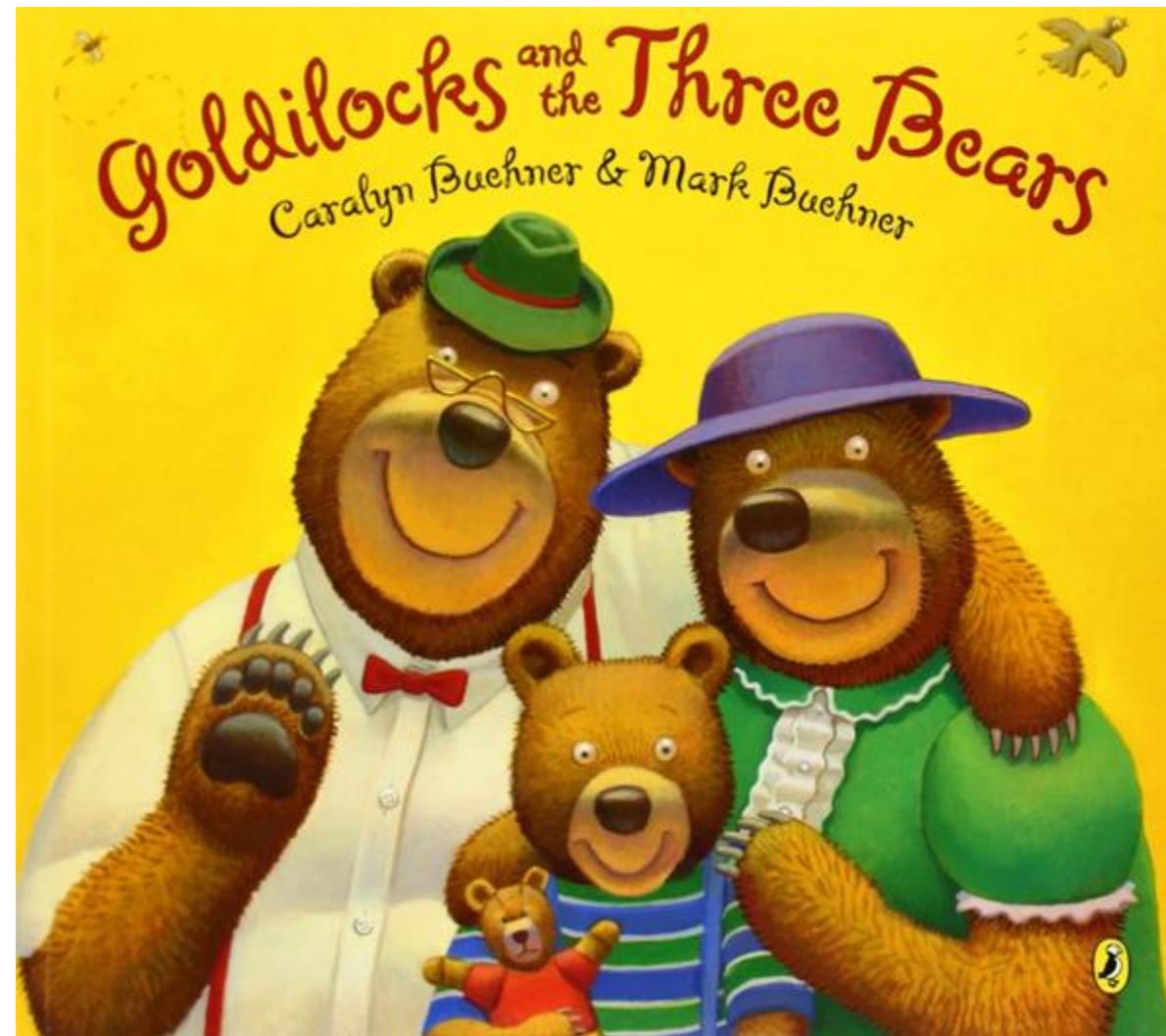


Turtles All The Way Down

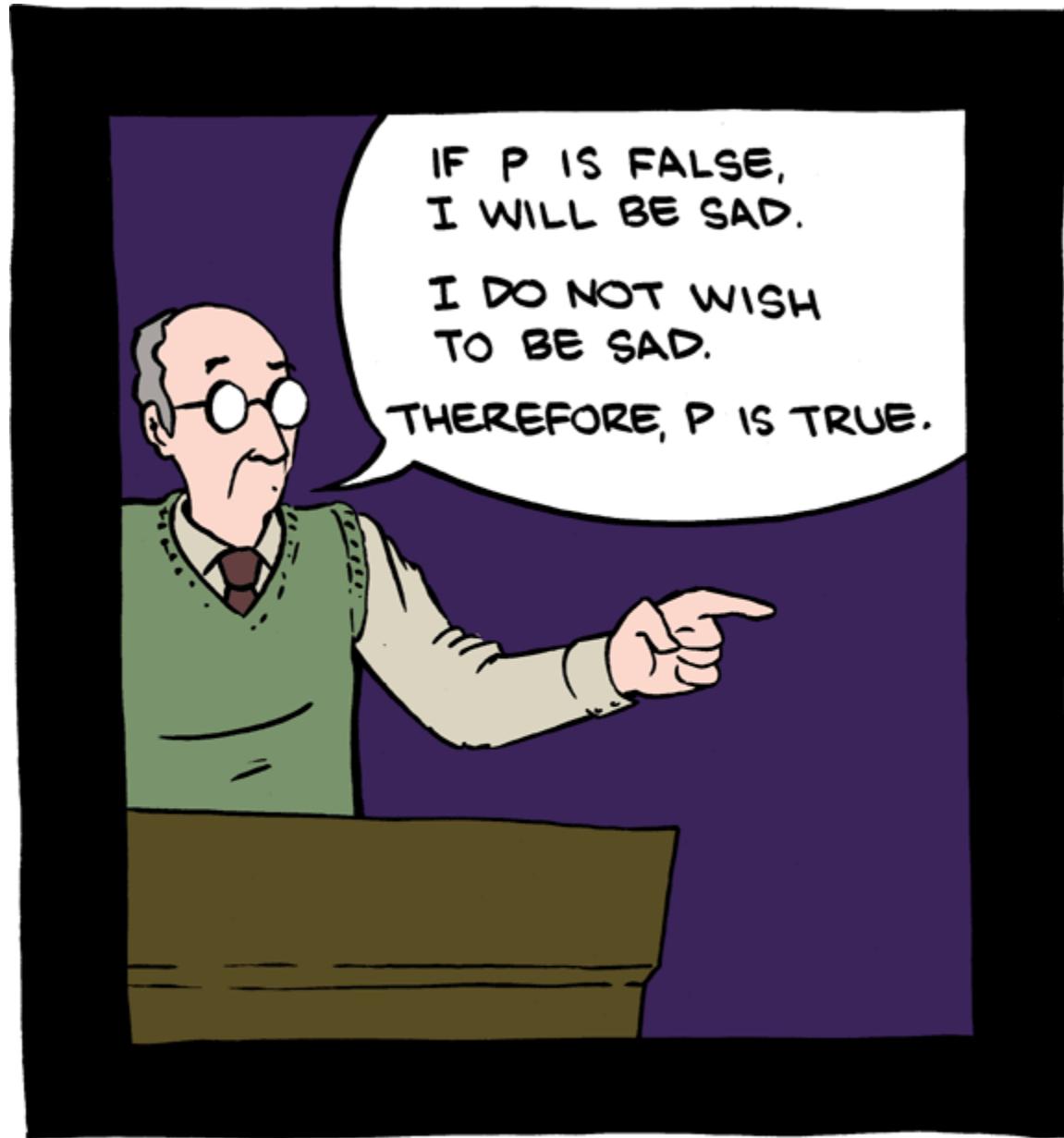
- “The “fix” is not a surgical replacement pacemaker, but a firmware update that takes about three minutes to complete and carries a “very low risk of update malfunction;” a very small percentage of people might experience a “complete loss of device functionality” during the firmware update. The patch covers St. Jude Medical’s pacemakers: Accent, Anthem, Accent MRI, Accent ST, Assurity and Allure.”
- <https://www.csoonline.com/article/3222068/hacking/465000-abbott-pacemakers-vulnerable-to-hacking-need-a-firmware-fix.html>

Guiding Narrative

- How should we think about testing?
- Lens of Logic
- Lens of Statistics
- Lens of Adversity



Lens of Logic



There. Now you can skip 99% of philosophical debates.

The Motivation

- If testing is our best way to gain confidence in the quality of software, but testing is expensive, how can we ensure that we are testing in an **effective** manner?
- Informal Desideratum: The program passes the tests **if and only if** it does all the right things and none of the wrong things.
- Pass all tests → program adheres to requirements
- Each failing test → program behaves incorrectly

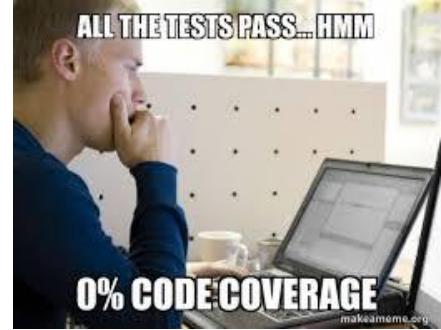
Intuition (Gedankenexperiment)



- Suppose you were writing a `sqrt` program and one of the requirements was that it should abort gracefully on negative inputs.
- Suppose further that your test suite does not include any negative inputs.
- Can we conclude that passing all of the tests implies adhering to all of the requirements?

```
def safe_sqrt(x):  
    if x < 0:  
        raise ValueError("Cannot do sqrt of negative number")  
    return math.sqrt(x)
```

Coverage



- We desire all of the requirements to be covered (“checked”) by the test suite.
- For our purposes, ***X coverage*** is the degree to which *X* is executed/exercised by the test suite.
- Examples:
 - **Code coverage** is the degree to which the source code is executed by the test suite.
 - **Statement coverage** is the fraction of source statements that are executed by the test suite.

Do Tests Cover All Requirements?

- In ideal world we would have **traceability** between requirements and test cases
- That is, each test case would have an annotation like “**a program that passes this test satisfies requirement X**” or “**passing this test gives confidence that a program adheres to requirement Y**”
- Outside of certain industries (e.g., Aerospace), such **formal traceability is rare**
- e.g., <https://en.wikipedia.org/wiki/DO-178C>

An Approximation

- We will cover requirements and their elicitation later in this course (mid-semester)
- But suppose for now you don't have formal traceability to your requirements
- So **testing that the program does all and only the good things that it is required to do is not possible** (or not feasible)
- Analogy: “Lie of Omission”

While at a lake, you see your friend is drowning.

Are you obligated to rescue them?

Aside: Ethics

- It is very tempting to say “yes, you are morally **obligated** to save your friend” (many would agree!)
- However, it can be surprisingly difficult to make a consistent moral system that requires particular positive actions, as opposed to just forbidding negative actions
- cf. “Thou shalt not kill” (Old Testament) or “An it harm none, do what ye will” (Wiccan Rede) or “Everything which is not forbidden is allowed” (English law), etc.
- For more information, take a class on *Ethics* (normative ethics) from the Philosophy department

Don't Do Bad Things

(How does good/bad intuition relate to testing?)

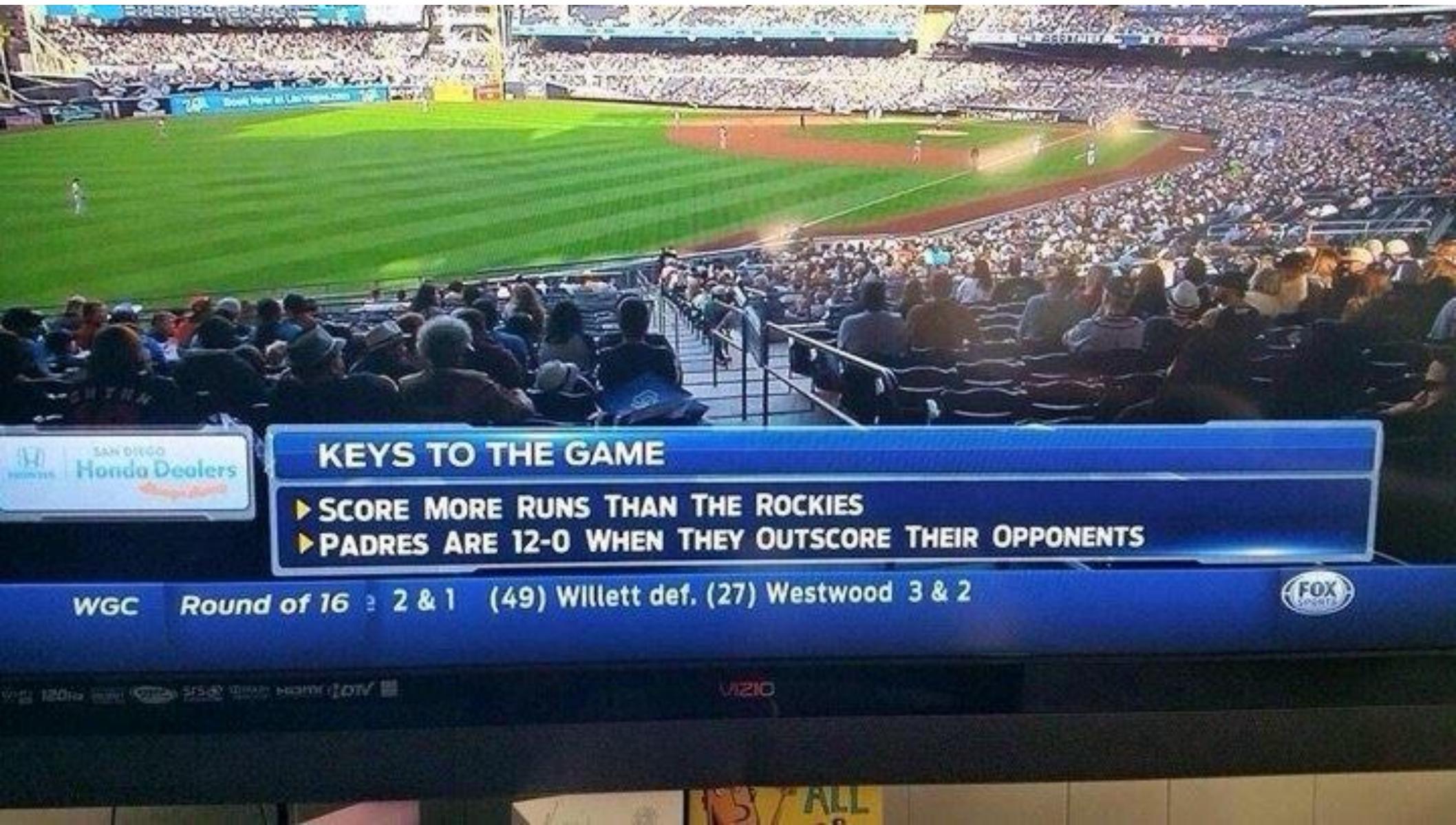


- We can at least test that the program does *not* do certain bad things
- e.g., “don't segfault”, “don't send my password to Microsoft”, “on this one particular input, don't get the wrong answer”
- Note that “I never do bad things” is not the same as “I always/eventually do good things”
- For more information, take a class on *Modal Logic* or read about *Liveness* vs. *Safety* properties

Testing to Find Bugs

- So now we want to test to gain confidence that the program does not do “bad things”
- That is, that the program does not have bugs
- Key Logical Observation: If we **never test** line X then testing **cannot rule out** the presence of a bug on line X
- (You could read line X, but this lecture is on testing. Later this semester: code review.)

If this seems “too obvious” so far,
just wait ...



$$P \rightarrow Q$$

“No test covers X → may have bug in X”

- Note that you could test line X and still have a bug on line X:

```
foo(a,b) { return a/b; }
```

```
test: foo(6,2)
```

- But testing X gives us some small but non-zero confidence in the correctness of X

“All Other Things Being Equal”

- If test A visits lines 1 and 2
- And test B visits lines 1, 2, 3 and 4
- Then, **all other things being equal**, we prefer test B
 - Test A gives some confidence about 1 and 2 and no confidence (no information) about 3 and 4
 - Test B gives some confidence about 1, 2, 3 and 4
- If the confidence/info gained per tested line is $c > 0$, test A gives us $2c + 0$ and test B gives us $4c$.
- Because $c > 0$, we have $4c > 2c$. So $B > A$.

Simplifying Assumptions

- Assumption 1. We gain the same amount of confidence (or information) for each visited line.
- Assumption 2. The amount of confidence (or information) we gain per visited line is positive.
- Assumption 3. ...



Line Coverage: A Test Suite Quality Metric

- A **test suite quality metric** or **test suite adequacy criterion** assesses the quality of a test suite (with respect to an external notion of utility) and allows test suites to be compared.
- **Line (or statement) coverage** is a test suite quality metric: it is the number of unique lines (statements) visited (exercised) by the program when running the test suite.
 - (Informally: visiting more lines is better because you have no information about un-visited lines.)

Using Line Coverage

(What can we do with it?)

- Given two test suites that both run within your resource budget (“AOTBE”, etc.), if we can only run one, we prefer the test suite with higher line coverage
- Thus **coverage is a metric that allows us to compare two test suites** and pick the “better” one (remember we managers want us to get better!)
- We use this **information to guide decision-making in a software process** (“how should we do testing?”)

Collecting Line Coverage

(How do we measure it?)

- At its simplest, this is just print-statement debugging
- Put a print statement before every line of the program
- Run all the tests, collect all the printed information, remove duplicates, count
- Practical concern:
 - the **observer effect** (from physics) is the fact that simply observing a situation or phenomenon necessarily changes that phenomenon.



Coverage Instrumentation

(How do we measure it?)

- **Coverage instrumentation** modifies a program to record coverage information in a way that minimizes the observer effect.
 - This can be done at the source or binary level.
- Don't actually print to stdout/stderr
- Don't slow things down too much
 - Pre-check before printing a duplicate?
- Don't introduce infinite loops
 - Instrument “print” with a call to “print”?

Good News: “Solved” Problem

- This is a well-studied problem and many push-button solutions exist for various forms of coverage
 - Either built in to your IDE or as external tools
- You will use three in the Homework
 - Python's coverage, gcc's gcov, Java's cobertura
- For more information on how to write one yourself, take a (graduate?) PL or Compilers class.

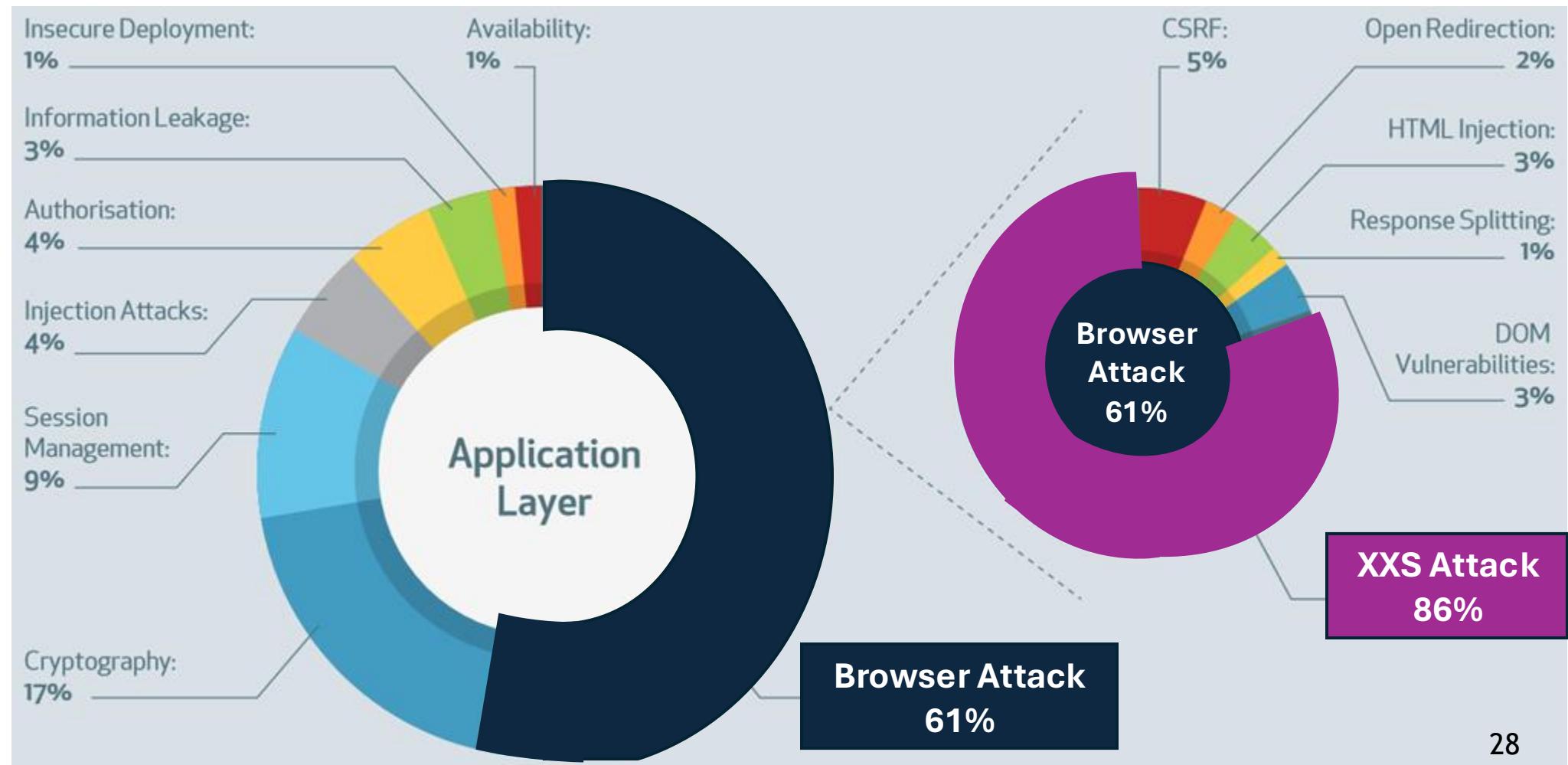
Problems with Line Coverage

- What could go wrong with line coverage?
- Can you think of situations with 100% line coverage where the program might still have bugs?



Example Where Statement Coverage is Inadequate

- Cross-site scripting (XXS) attacks: [2016 Vulnerability Statistics Report, edgescan]



Example Where Statement Coverage is Inadequate

• Cross-site scripting (XXS) attacks: [2016 Vulnerability Statistics Report, edgescan]



Testing Data Values Matter Too!

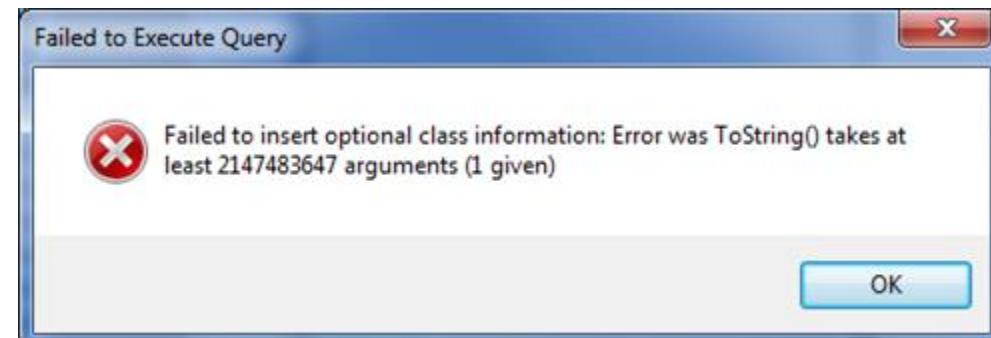
Data Values and Implicit Control Flow

```
return a/b      -----> if (b != 0)  
                return a/b  
                else  
                ABORT
```

```
print ptr->fld → if (ptr != NULL)  
    print ptr->fld  
else  
    ABORT
```

Intuition

- Many interesting **data** values cause implicit or explicit changes of **control**
- That is, they cause different branches of conditionals to execute
- Informally, the problem of ensuring that we cover interesting data values may *reduce to* the problem of ensuring that we cover all branches of conditionals



Branch Coverage

- **Branch coverage** is a test suite quality metric that counts the total number of conditional branches exercised by that test suite (i.e., if→true and if→false are counted separately)
- Note that branch coverage can **subsume** line coverage:

```
def foo(a):  
    if a > 5:  
        print("x")  
    print("y")
```

Test Suite {foo(7)} has 100% line coverage but 50% branch coverage.

Test Suite {foo(7), foo(0)} has 100% line and 100% branch coverage.

Branch vs. Line

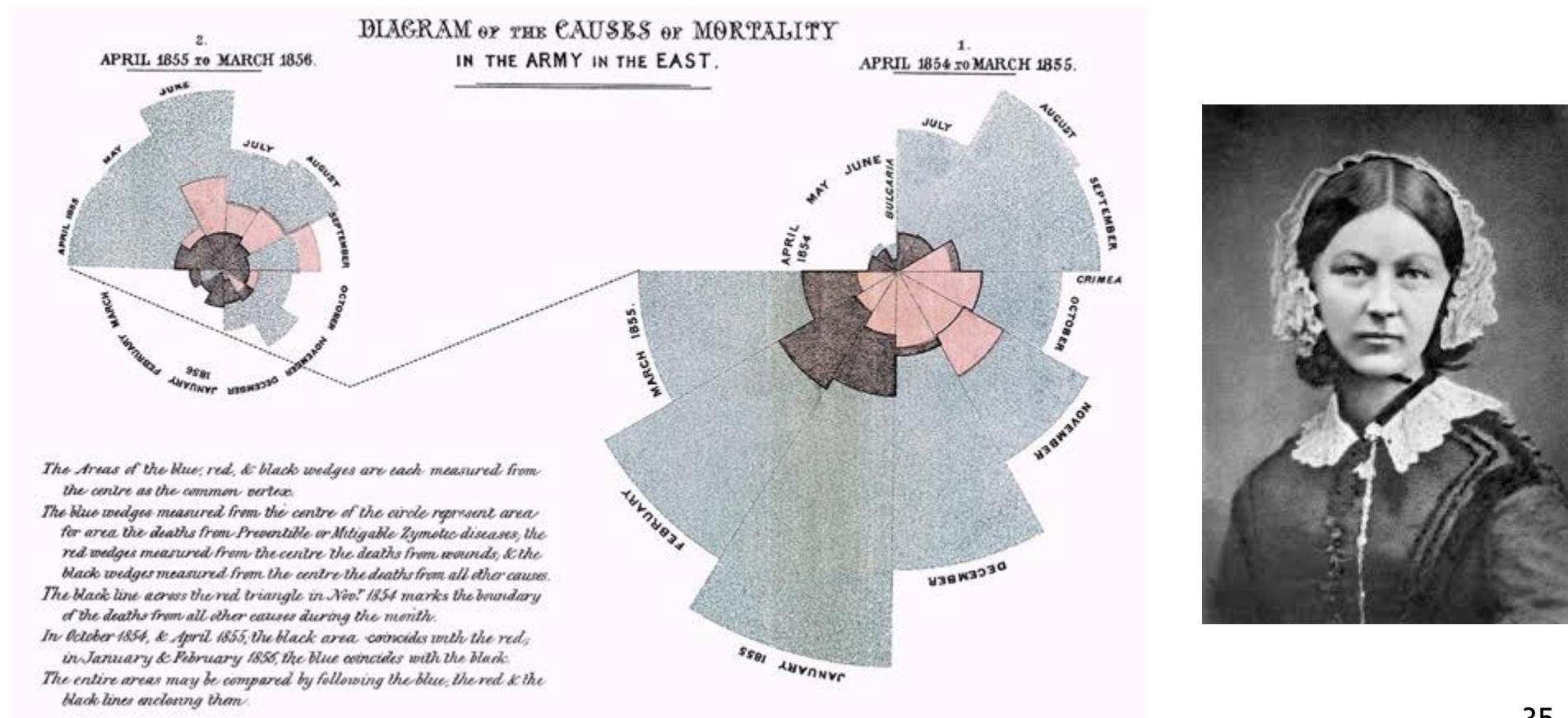
- Branch coverage typically gives us **more confidence** than line coverage
- Typically, 100% branch coverage implies 100% line coverage
- However, branch coverage is “more expensive” in the sense that it is harder for a test suite to have high branch coverage than to have high line coverage
 - Note: quality isn't really “more expensive”, you were just fooling yourself before by thinking line coverage was OK. Being correct is expensive.

Other Flavors

- **Function Coverage**: what fraction of functions have been called?
- **Condition Coverage**: what fraction of boolean subexpressions have been evaluated to both true and also (e.g., on another run) to false?
 - Comparing this to branch coverage is a not-uncommon test question ...
- **Modified Condition / Decision Coverage (MCDC)**: function coverage + branch coverage (this is a simplification)
 - Used in mission critical (e.g., avionics) software

Trivia: Statistics

- This English social reformer and statistician (among other activities, ~1850) was a pioneer in the use of infographics: the effective graphical presentation of statistical data.*



Trivia: Geography

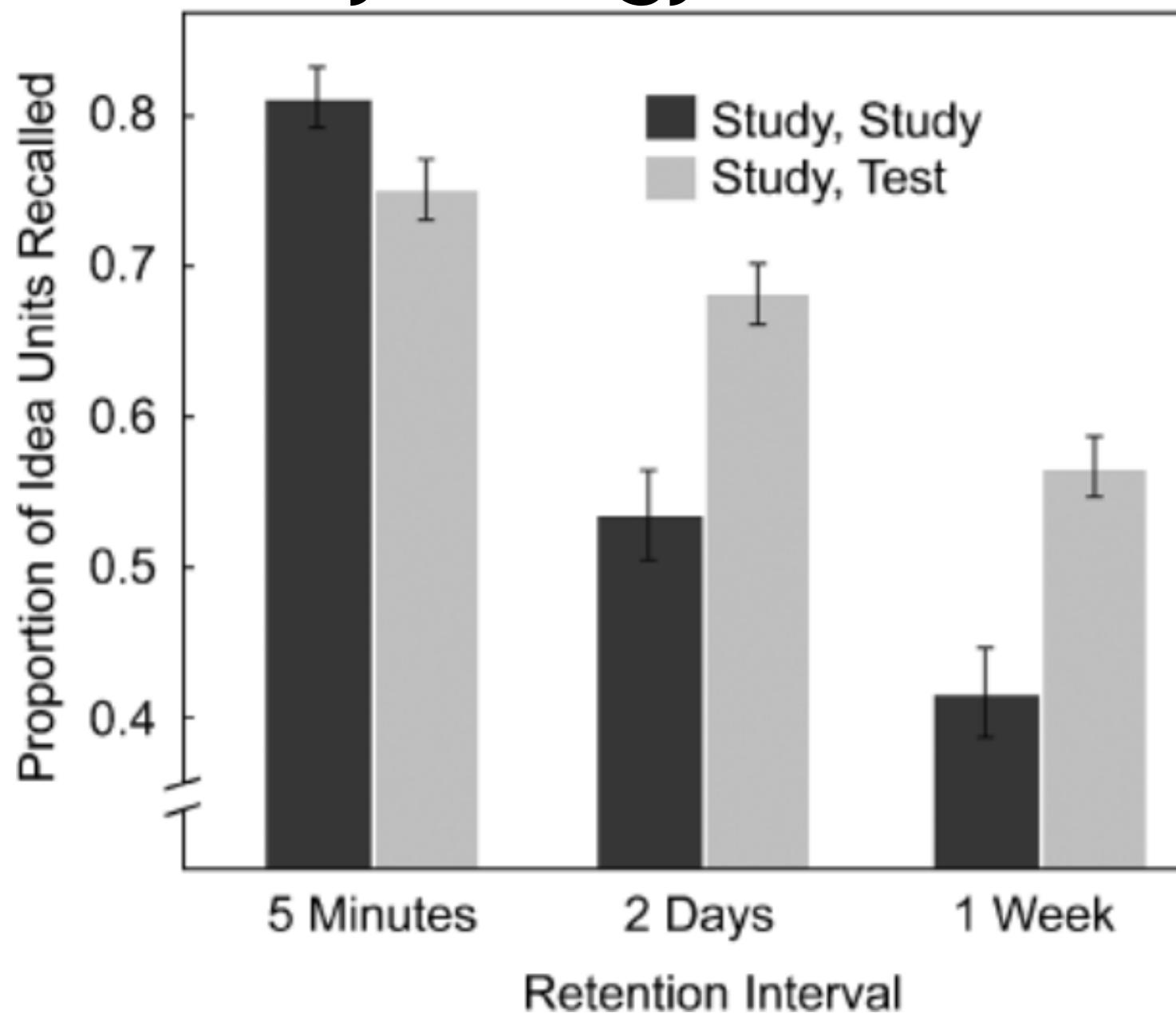
- *This is the capital of, and most populous city in, Malaysia. Its name means “muddy confluence” in Malay. It was captured by Japan in 1942 and gained independence from British rule in 1957. Islam and Buddhism are the majority religions.*



Psychology: Recall

- 120 students (age 18 to 24) were asked to study prose passages (e.g., 300 words on “Sea Otters”) and also do math problems
- Group 1: Read for 7m, math for 2m, re-read for 7m, math for 5m
- Group 2: Read for 7m, math for 2m, test for 10m, math for 5m
- Both groups: later → test for 10 minutes
- Which group did better? By how much?

Psychology: Recall



Psychology: Testing Effect

- The **testing effect**: long-term memory is **increased** when some of the learning period is devoted to **retrieving** the to-be-remembered information through testing with feedback.
- “They found that re-studying or re-reading memorized information had no effect, but trying to recall the information had an effect.”
- Implication for SE: Code comprehension.

[Roediger, H. L.; Karpicke, J. D. (2006). "Test-Enhanced Learning: Taking Memory Tests Improves Long-Term Retention". *Psychological Science*. 17 (3): 249-255.]

Lens of Statistics



Alternate View

- The bugs experienced by **users** are the ones that matter.
- Dually, bugs never experienced by users do not matter.



Positive User View



- Suppose you are writing a point-of-sale cashier application that makes change for a dollar (\$1). Given any price between 1 and 100 cents, you must indicate the coins to give out as change.
 - e.g., Price is 23 cents, User Gives you a \$1 → return 3 quarters and 2 pennies (77 cents)
- In this scenario, you can **exhaustively** test all 100 inputs that will occur to real users in the real world
- In some sense, it does not matter if that is 100% statement or code coverage (e.g., dead code)

Negative User View

- Suppose users will only ever cause lines 1, 2 and 3 of your program to be executed
- Then you do not need to test line 4
 - Even if it has a bug, users will **never** encounter that bug
- Note “will” → this either requires a prediction of the future or a finite input domain

Testing as Sampling

- If user-experienced bugs are the ones that matter, testing should be devoted to **sampling** those inputs that users will provide
- Two views:
 - Sample what users do **most commonly**
 - Sample what causes the **most harm** if users do it
- Compare:
 - $\text{Risk} = (\text{Prob. of Event}) * (\text{Damage if Event Occurs})$

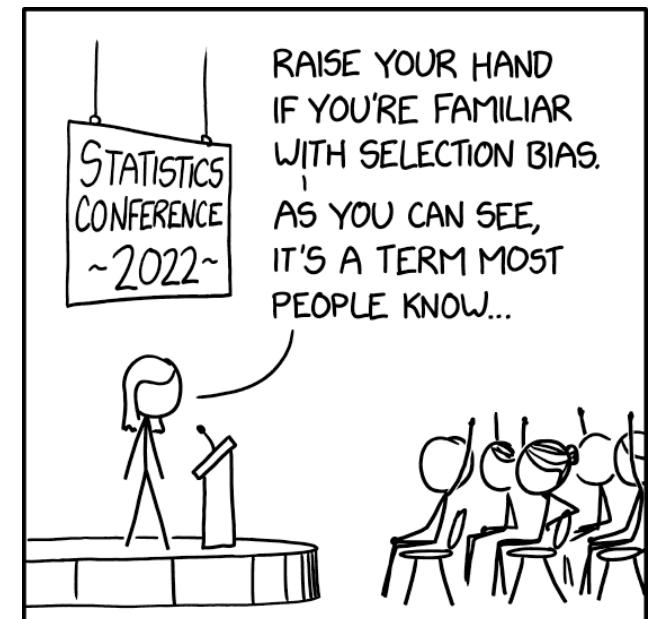
Sampling Error

- In statistics, **sampling error** occurs when the statistical characteristics of a population are estimated from a subset, or sample, of that population.
 - “Our test suite is a sample of inputs that could occur in the real world. Our program behaves well on our test suite.” → later → “Our program behaves badly on some other untested real input. Sampling error!”
- Testing gives confidence the same way sampling (or polling) gives confidence.



Sampling Bias

- In statistics, **sampling bias** is a bias in which a sample is collected in such a way that some members of the intended population are less likely to be included than others.
- Suppose you are conducting a poll to see who will win the next election, but you only poll republicans.
- Suppose you are creating tests to see if your program will crash, but you only poll nice, small, inputs.



Solution?

- There are a number of well-established sampling techniques in the field of statistics to help address such biases
- They often require knowing something about the **distribution** of the full population from which you want to sample a subpopulation
- The basic problem in SE is that the underlying distribution of real user inputs **is not known**

Beta Testing

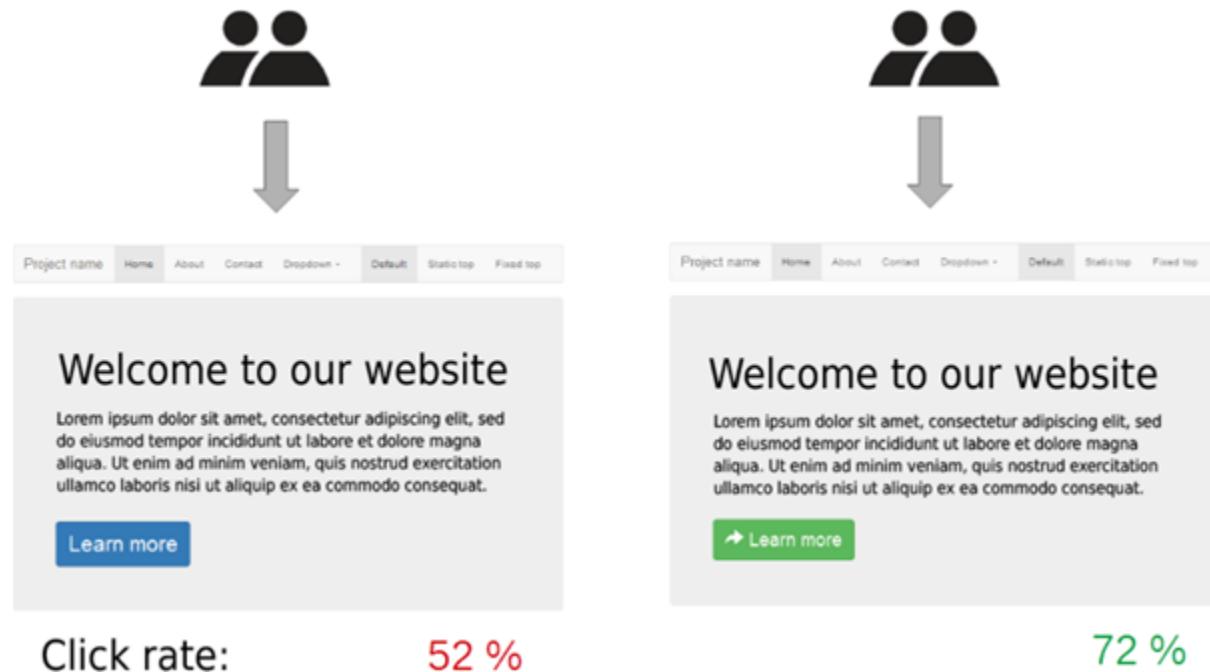
(What can we do? Ask Users!)

- **Alpha testing** is testing done by developers.
- **Beta testing** is testing done by external users (often using a special beta version of the program).
 - Beta testing can be viewed as directly sampling the space of user inputs



A/B Testing

- **A/B testing** involves **two variants** of your **software**, A and B, which **differ only in one feature**. Different users are shown different variants and responses are recorded. It is an instance of two-sample statistical hypothesis testing.



Likely or Damaging?

- Recall two guiding approaches:
 - (1) Sample what users will do **most commonly**
 - (2) Sample what will cause the **most harm**
- The former (1) is sometimes called **workload generation**

Common for databases, webservers, etc. (Youtube?)
- The latter (2) often relates to **computer security**
 - Exploit generation, penetration testing, etc.
 - cf. AFL in Homework 2

Non-Security Damage

- For Amazon (etc.), “damaging” is “customer does not complete the purchase”

- Cascading Stylesheet Error.** An error in loading the stylesheet between the *current* and *next* pages.
- Code on the Screen.** Any error that results in programming language code appear on screen, including any error referring to a line number (with the exception of visible HTML code).
- Other Error/Error Message.** Either any error message, or any error that cannot be classified in any other category.
- Form Error.** Missing, malformed, or extra buttons, form fields, drop-down menus, etc, including incorrectly validating forms.
- Missing Information.** Any part of a webpage that is missing, not including images.
- Wrong Page/No Redirect.** An unexpected page is loaded.
- Authentication.** Any errors that occur during login.
- Permission.** Any errors occurring with respect to user permissions in an application, such as access being incorrectly denied to a user.

Feature	Correlation	F	Pr(> F)
Code on the Screen	+	19.47	0.00
Cosmetic	-	13.23	0.00
Database	+	12.36	0.00
Authentication	+	6.99	0.01
Functional Display	-	6.00	0.01
Other Error	+	4.40	0.03

[Dobolyi et al. Modeling Consumer-Perceived Web Application Fault Severities for Testing. ISSTA 2010.]

Lens of Adversity



Finding Bugs



- Suppose you wanted to evaluate the quality of two truffle-sniffing pigs (or bomb-sniffing dogs)
- You might hide some truffles and see how many each pig finds (etc.)
- The pig that finds more of the hidden truffles in your backyard is assumed to find more real truffles in the wild
- Suppose you wanted to evaluate the quality of two bug-finding test suites ...

Mutation Testing

- **Mutation testing** (or **mutation analysis**) is a test suite adequacy metric in which the quality of a test suite is related to the number of intentionally-added defects it finds.
- Informally: “You claim your test suite is really great at finding security bugs? Well, I'll just **intentionally add a bug** to my source code and see if your test suite finds it!”



How Close is this to Real Life?

- In the truffle-pig example, if every truffle I hide in my back yard is next to a smelly red flower, a pig that finds them all may not actually do well in the real world
 - The truffle placements I made up were **not indicative** of real-world truffles
 - Similarly, if I add a bunch of defects to my software that are not at all the sort of defects real humans would make, then mutation testing is uninformative

Defect Seeding

- **Defect seeding** is the process of intentionally introducing a defect into a program. The defect introduced is typically intentionally similar to defects introduced by real developers. The seeding is typically done by changing the source code.
- For mutation testing, defect seeding is typically done **automatically** (given a model of what human bugs look like)
 - You will do this in Homework 3

Mutation Operators

(How do we **seed defects**?)

- A **mutation operator** systematically changes a program point. In mutation testing, the mutation operators are modeled on historical human defects. Example mutations:

`if (a < b)` → `if (a <= b)`

`if (a == b)` → `if (a != b)`

`a = b + c` → `a = b - c`

`f(); g();` → `g(); f();`

`x = y;` → `x = z;`

Mutant



- A **mutant** (or **variant**) is a version of the original program produced by **applying one or more mutation operators** to one or more program locations.
- The **order of a mutant** is the number of mutation operators applied.

// original

```
if (a < b) :  
    x = a + b  
print(x)
```

// 2nd-order mutant

```
if (a <= b) :  
    x = a - b  
print(x)
```

Competent Programmers

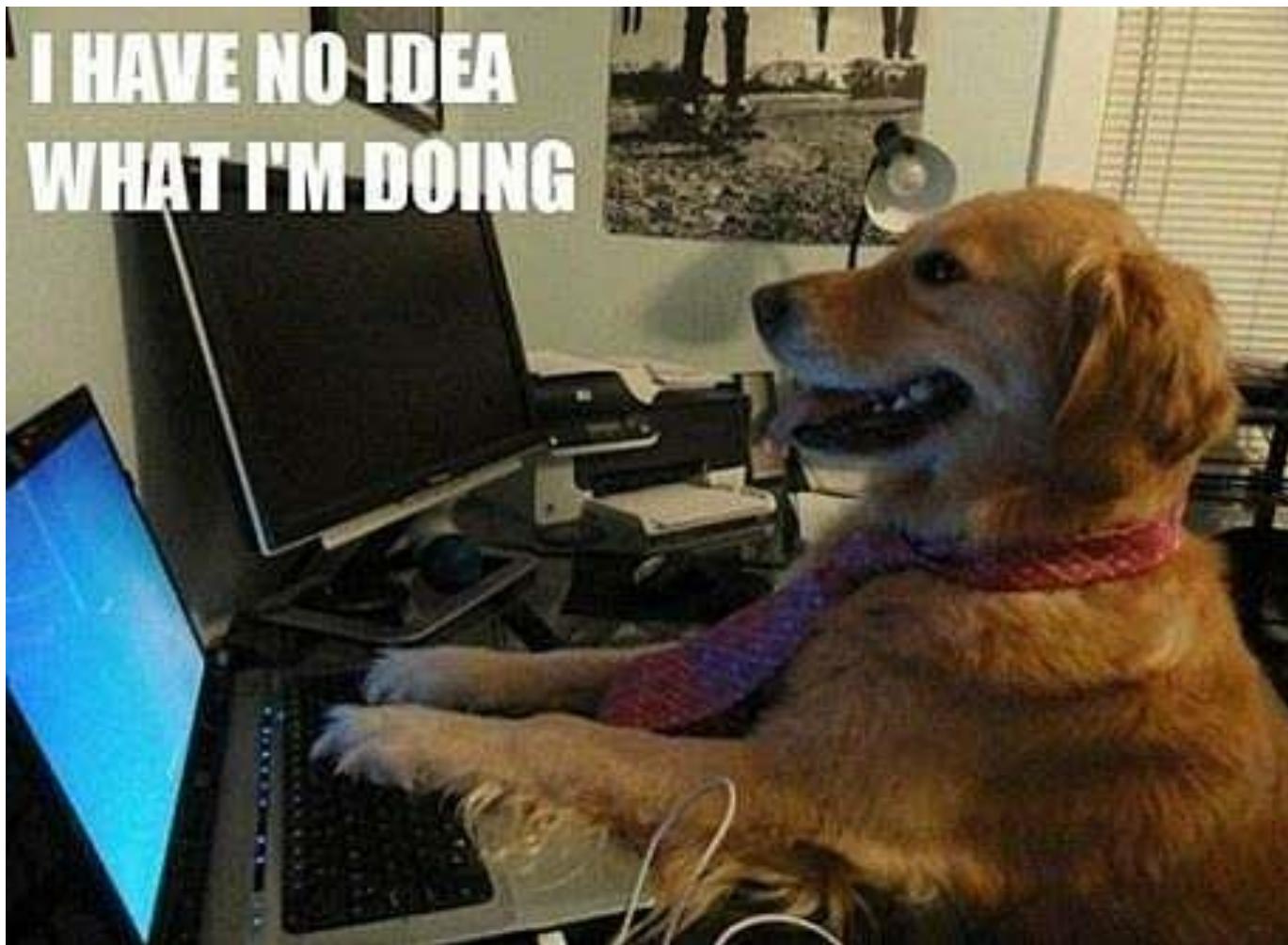
- The **competent programmer hypothesis** holds that program faults are syntactically small and can be corrected with a few keystrokes.
- Programmers write programs that are largely correct. Thus, the **mutants simulate the likely effect of real faults**.
- Therefore, **if the test suite is good at catching the artificial mutants, it will also be good at catching the unknown but real faults in the program**.

Do Humans Really Make Simple Mistakes?



Competent?

- Is the competent programmer hypothesis true?



Competent?

- Is the competent programmer hypothesis true?
- Yes and no.
- It is certainly true that humans often make simple typos (e.g., + to -).
- But it is also true that some bugs are more **complex** than that.

Coupling Effect

- The **coupling effect hypothesis** holds that complex faults are “coupled” to simple faults in such a way that a test suite that detects all simple faults in a program will detect a high percentage of the complex faults.
- Is it true?
 - Tests that detect simple mutants were **also** able to detect over 99% of second- and third-order mutants **historically** [A. J. Offutt. Investigations of the software testing coupling effect. ACM Trans. Softw. Eng. Methodol., 1(1):5-20, Jan. 1992.]

Mutation Testing

- A test suite is said to **kill** (or **detect**, or **reveal**) a mutant if the mutant fails a test that the original passes.
- **Mutation testing** (or **mutation analysis**) of a test suite proceeds by making a number of mutants and measuring the fraction of them killed by that test suite. This fraction is called the **mutation adequacy score** (or **mutation score**).
- A test suite with a higher score is better.

The wording can be tricky, I know ...



Mutation Analysis: Pros and Cons

- Has the potential to **subsume** other test suite adequacy criteria (it can be very good)
- Which mutation operators do you use?
- Where do you apply them? How often do you apply them?
 - Typically done at random, but how?
- It is **very expensive**. If you make 1,000 mutants, you must now run your test suite 1,000 times!
 - We started by saying testing (1x) was expensive!



Equivalent Mutant Problem

- Suppose you have “ $x = a + b; y = c + d;$ ” and you swap those two statements.
- The resulting program is a mutant, but it is **semantically equivalent** to the original.
 - So it will pass and fail all of the tests that the original passes and fails.
- So it will dilute the mutation score (make worse)
- Detecting **equivalent mutants** is a big deal. How hard is it?

Equivalent Mutant Problem

- Detecting equivalent mutants is a big deal.
How hard is it?
- It is **undecidable!**
 - By direct reduction to the **halting problem**, or by **Rice's Theorem**

```
def foo():          # foo halts if and only if
    if p1() == p2(): # p1 is equivalent to p2
        return 0

foo()
```

Questions?

- Lens of Logic: “no visit X → no find bug in X”
 - Leads to statement and branch coverage.
- Lens of Statistics: “sample the inputs the users will make”
 - Leads to beta testing, A/B testing.
- Lens of Adversity: “poke realistic holes in the program and see if you find them”
 - Leads to mutation testing.
- ***Don't neglect HW 1 components!***