

Kill All Mutants!

(Intro to Mutation Testing)

by Dave Aronson
T.Rex-2025@Codosaur.us



Codosaur.us

@davearonson

(Blank slide so I can flip to a new one to start my timer, ignore this.)

CURRENT TIME: 37-38, want 35-40 (slot is 45, w/ 5-10 for Q&A), can slow or ad-lib a LITTLE

Kill All Mutants!

(Intro to Mutation Testing)

by Dave Aronson
T.Rex-2025@Codosaur.us



Codosaur.us

@davearonson

Bonjour encore, CERN!



(Hello again, CERN!)

Codosaur.us

Image: standard emoji

@davearonson

Je me suis présenté ce matin,



(I introduced myself this morning,)

Codosaur.us

Image: me speaking at JSConf Hawai'i 2020

@davearonson

et maintenant je suis ici



(and now I'm here)

Codosaur.us

Image: https://commons.wikimedia.org/wiki/File:Clock_01-15.svg

@davearonson

pour vous enseigner



(to teach you)

à tuer les mutants!



(to kill mutants!)

Codosaur.us

Image: <https://pixabay.com/vectors/turtle-tortoise-cartoon-animal-152079/>

@davearonson

(Mutants, pas moutons.)



(Mutants, not sheep.)

Mais encore . . .



(But again . . .)

je le ferai en anglais.



(I will do it in English.)

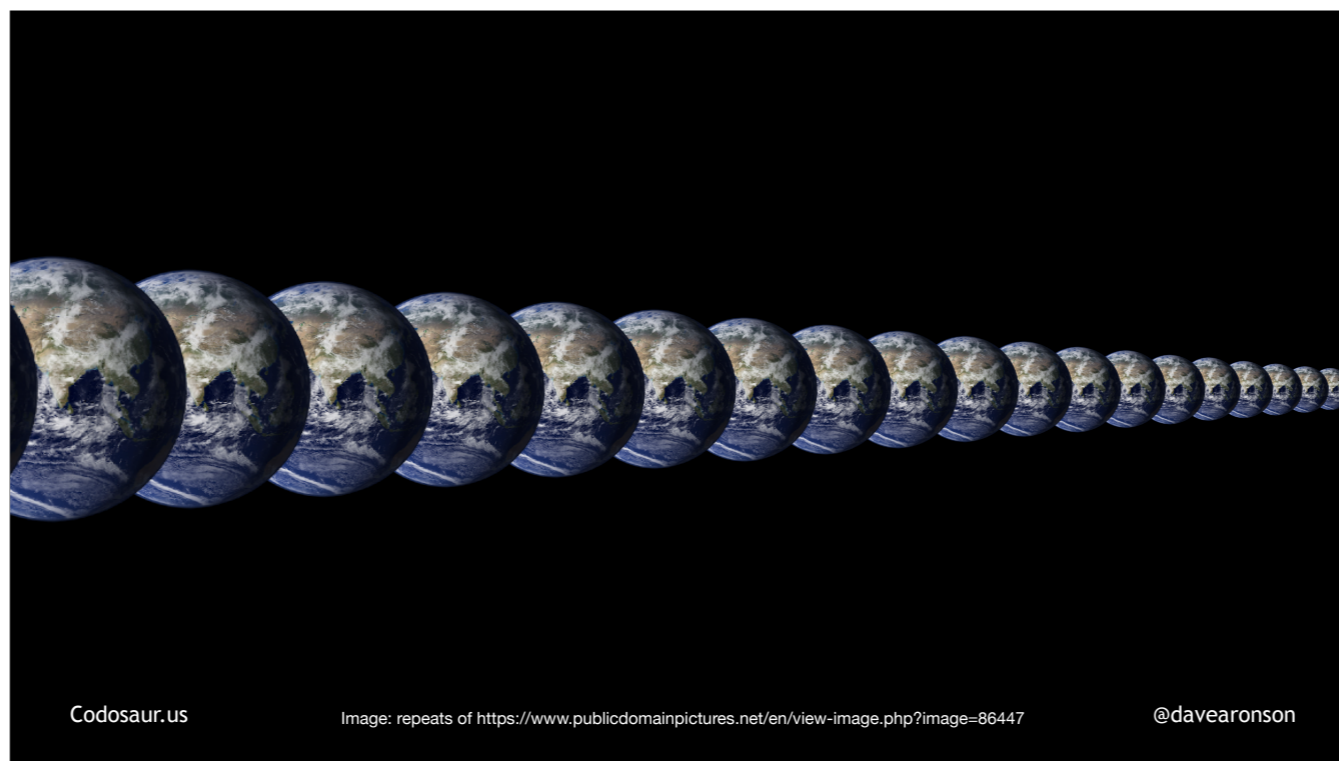
Codosaur.us

Image: standard emoji

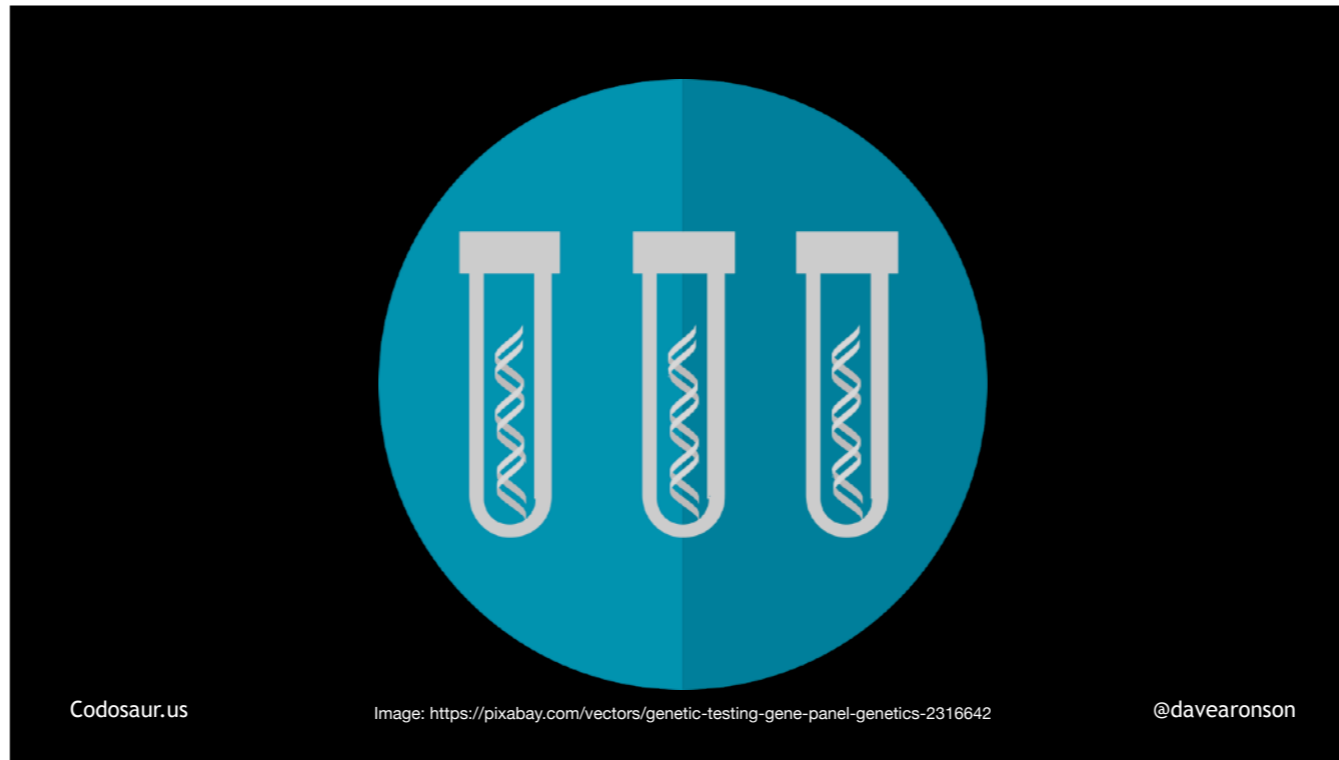
@davearonson

(PAUSE!)

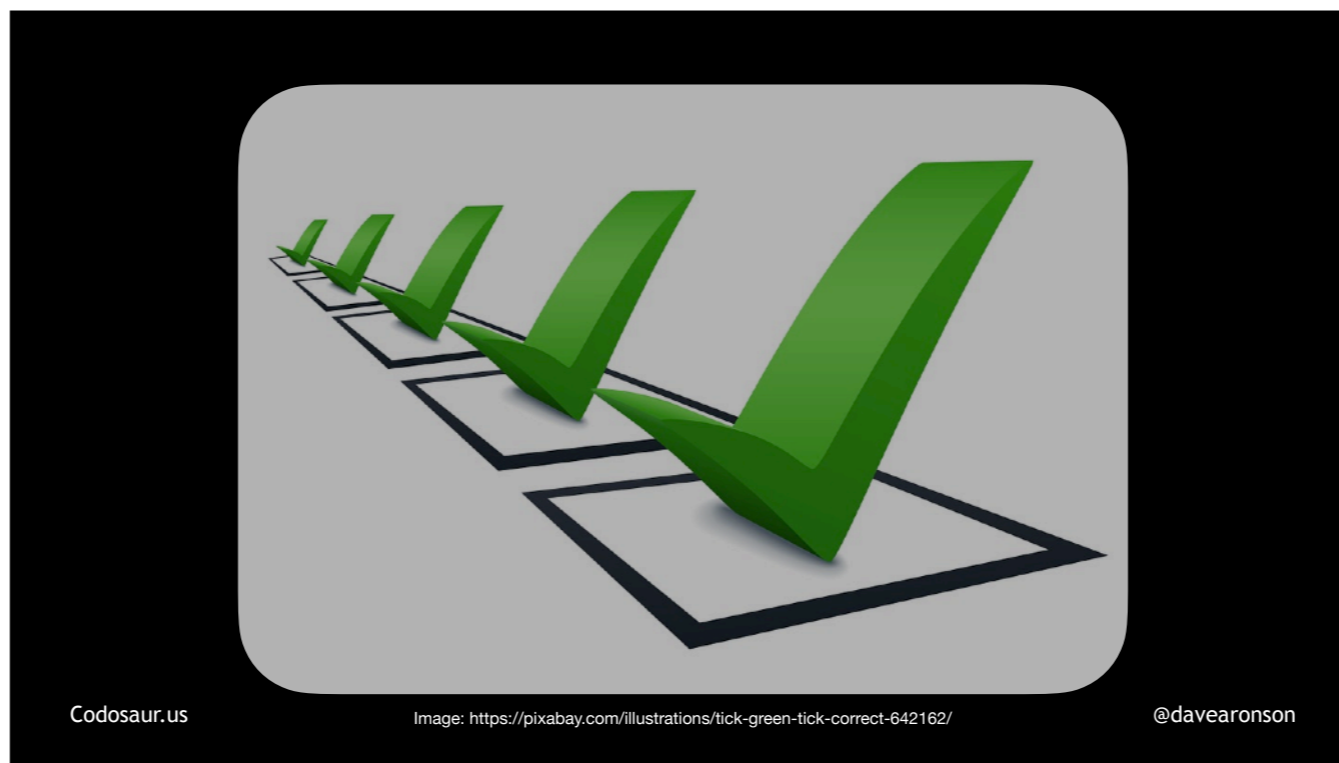
So, what on . . .



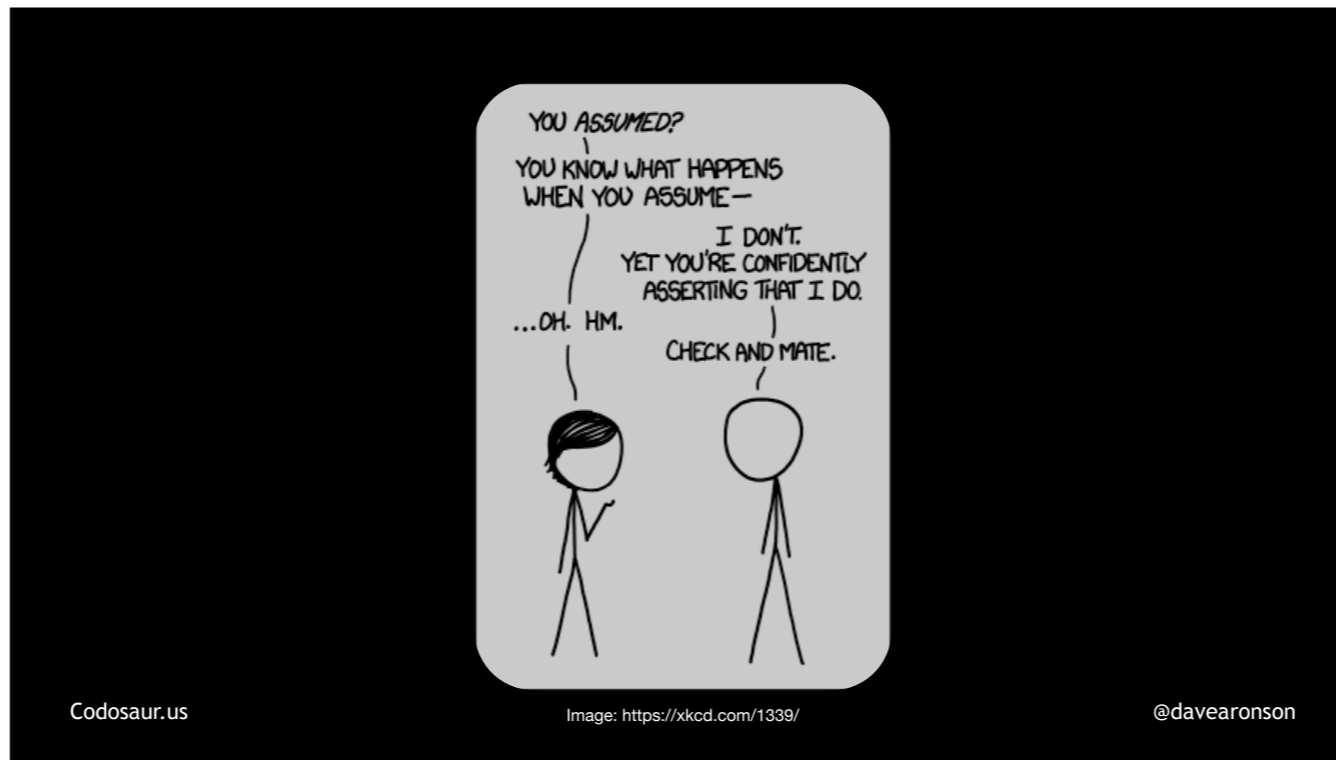
. . . Infinite Earths, makes . . .



. . . mutation testing different from all the *other* software testing techniques? The main difference is that most of the others are about . . .



. . . checking whether our code is correct. But mutation testing . . .



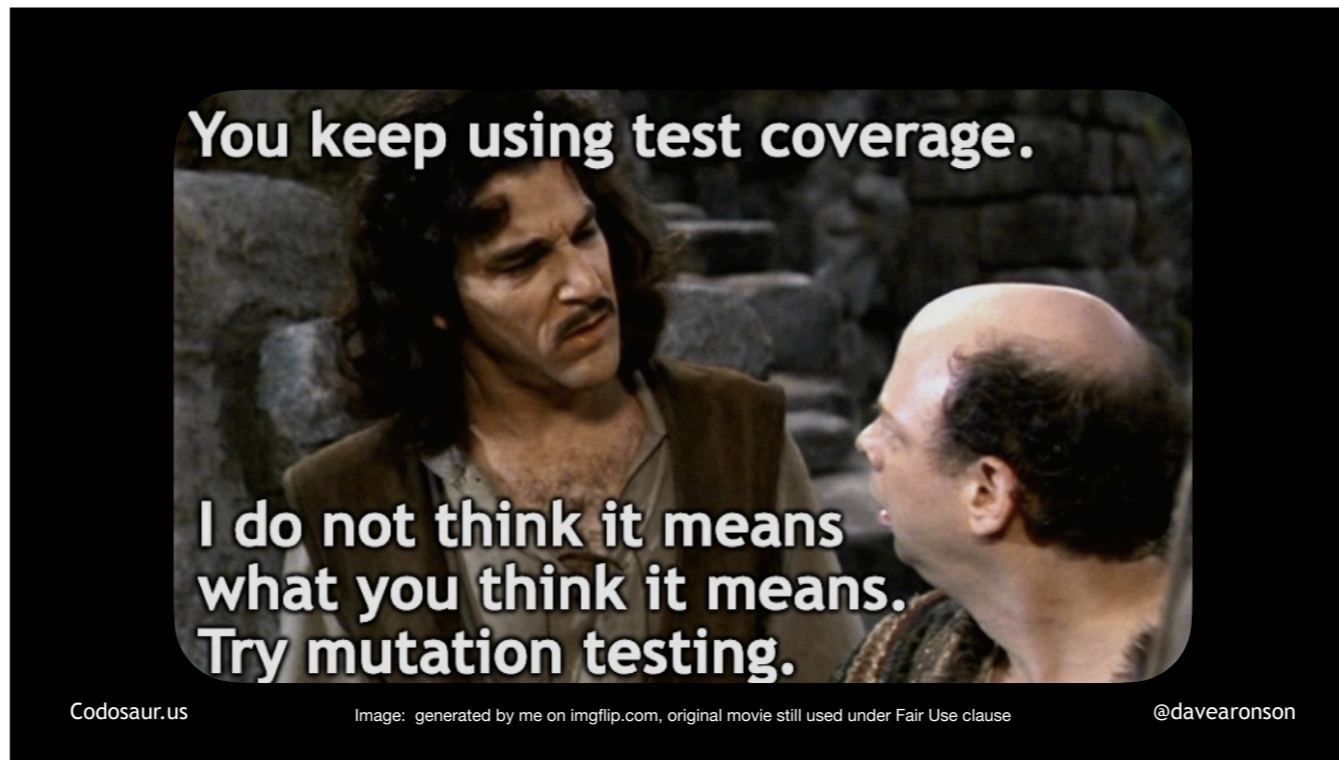
. . . *assumes* that our code is correct, at least in the sense of passing its tests. Instead, mutation testing checks for two *other* qualities. In a typical codebase, I think the more *important* one is that our test suite is . . .

```
"use strict";
```

Codosaur.us

@davearonson

. . . *strict*. Now you may be thinking, “Isn’t that what test coverage is for? If we have 100% test coverage, doesn’t that mean all our code is fully tested?”



No. The *only* thing that test coverage tells us is that at least one test *ran* . . .

```
class Conway:
    ALIVE = "*"
    DEAD = " "

    @classmethod
    def next_state(cls, cur_state, neighbors):
        if cur_state == cls.ALIVE:
            r = cls.ALIVE if neighbors in [2,3] else cls.DEAD
        else:
            r = cls.ALIVE if neighbors == 3 else cls.DEAD
        return r

    def another_func:
        # whatever
```

Codosaur.us

@davearonson

. . . the code it claims is “covered”, and *no* tests ran the rest of it. It tells us NOTHING about whether the *correctness* of the code *made any difference* to whether any test passed, let alone any *specific* test. And that’s what “tested” really *means*, right?

By way of illustration, let’s look at . . .

```
test_live_with_3_survives:  
    expected = Conway.ALIVE  
    actual = next_state(Conway.ALIVE, 3)  
    assert actual == expected
```

Codosaur.us

@davearonson

... a test, and ...

```
test_live_with_3_survives:  
  expected = Conway.ALIVE  
  actual = next_state(Conway.ALIVE, 3)  
  # assert actual == expected
```

Codosaur.us

@davearonson

. . . comment out the assertion. Assertions might be commented out, removed, or not even written in the first place, for any number of reasons, usually not very good, but it's often done anyway. Let's have a show of hands, who's seen that? I'm not asking who's *done* it, so no need to be ashamed! Anyway, our test still *runs* the function, so the lines *show as covered*, but we're not *asserting* anything, so the code is *obviously* not tested. This is just *one of many* ways that coverage can be misleading.

So how *can* we tell if the code really is tested? As you may have guessed, that . . . is where mutation testing comes in.

To check that our test suite is *strict*, a mutation testing tool will try to . . .



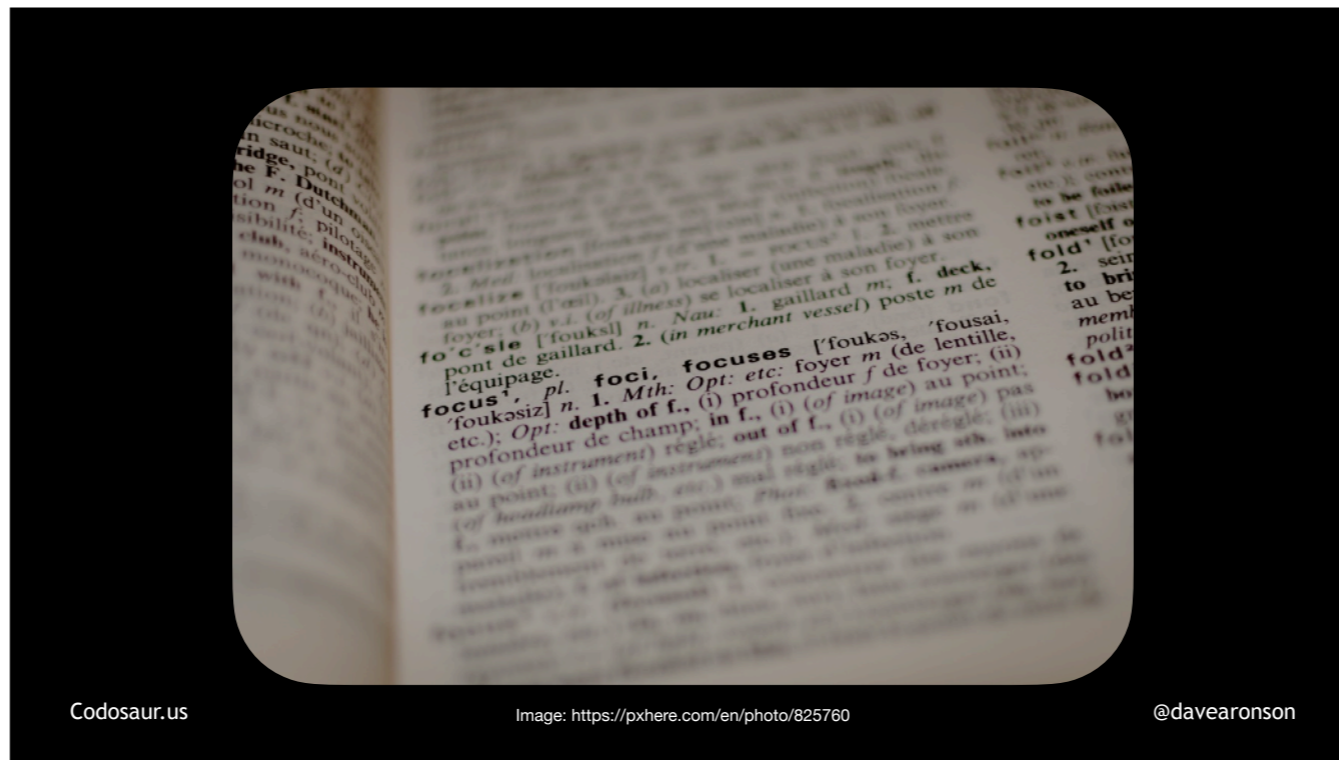
Codosaur.us

Image: https://commons.wikimedia.org/wiki/File:Mind_the_gap_2.JPG

@davearonson

. . . find the gaps in our test suite, that let our code get away with unintended behavior. Once we find gaps, we can close them by either adding tests or improving existing tests. Lack of strictness comes mainly from *lack* of tests, or poorly *written* tests.

The other thing mutation testing checks is that our code is . . .



. . . *meaningful*, so that any semantic change to the code will produce a noticeable change in its behavior. Lack of *meaning* comes mainly from code being unreachable, redundant, or otherwise just not having any real effect. When we find "meaningless" code, we can figure out *why* it's meaningless, then make it meaningful, if that fits our intent, but the usual fix is just to remove it.

Mutation testing . . .



. . . puts these two together, by checking that every change to the code, that the tool knows how to do, does indeed make a noticeable change to its behavior, *and* that the test suite is indeed strict enough that at least one test will notice that change, and fail.

That's the positive side, but there are some drawbacks. As . . .



**Fred Brooks, author of
"No Silver Bullet –
Essence and Accident in
Software Engineering"
(1986 paper)**

Codosaur.us

Image: https://commons.wikimedia.org/wiki/File:Frederick_Brooks_IMG_2279.jpg

@davearonson

. . . Fred Brooks told us back in 1986, there's no . . .

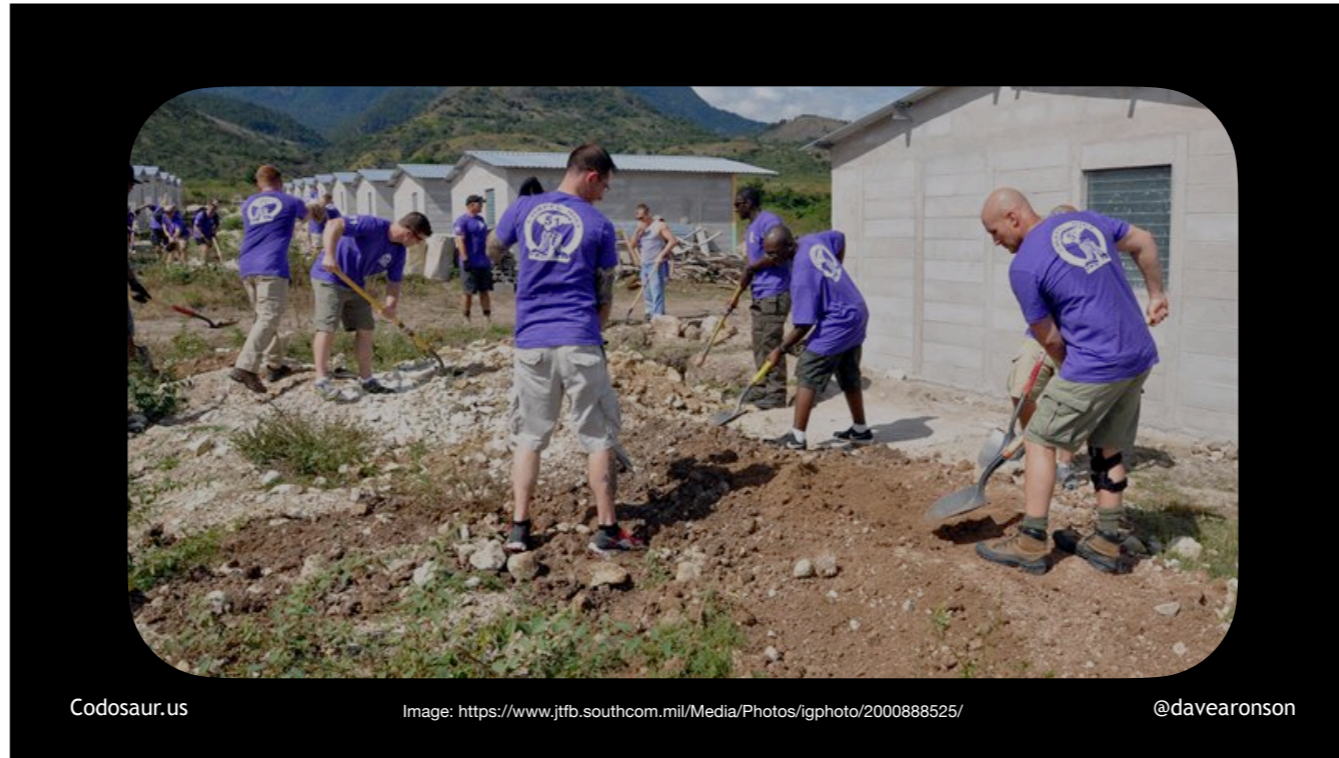


. . . silver bullet! Besides, those are for killing . . .



. . . werewolves, not mutants!

The first drawback is that it's . . .

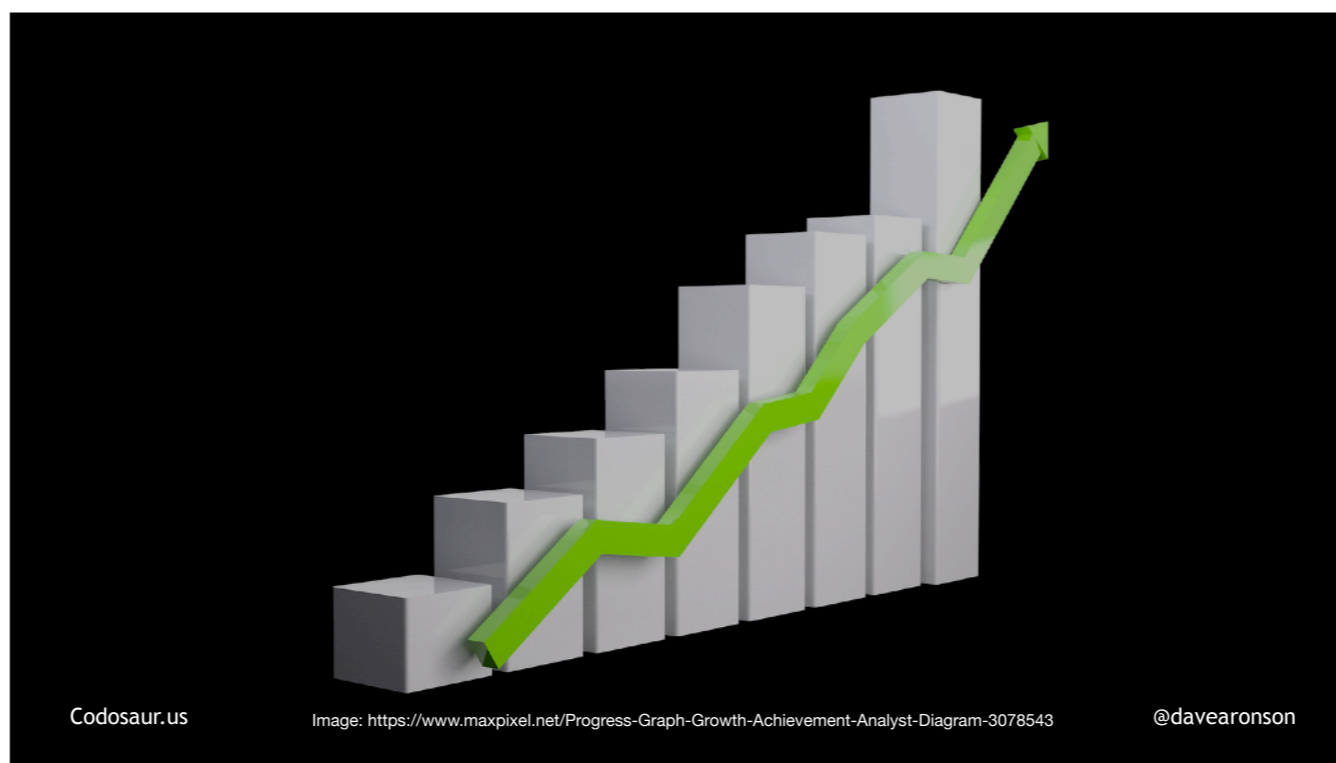


Codosaur.us

Image: <https://www.jtfb.southcom.mil/Media/Photos/igphoto/2000888525/>

@davearonson

. . . hard labor for the CPU, and therefore usually rather sllloooooow. We certainly won't want to mutation-test our entire codebase every time we save a file! Maybe over a lunch break for a smallish system, or a weekend for a large one. Fortunately, most tools let us just check specific functions, classes, files, and so on. Also, they often include some kind of . . .



. . . incremental mode, so that we can test only the changes since the last mutation test, or the last git commit, or the changes from the main branch, or some such milestone. With such filtering, maybe we can test just the relevant changes on each save, or at least over a much shorter break.

Another drawback is that it's often . . .

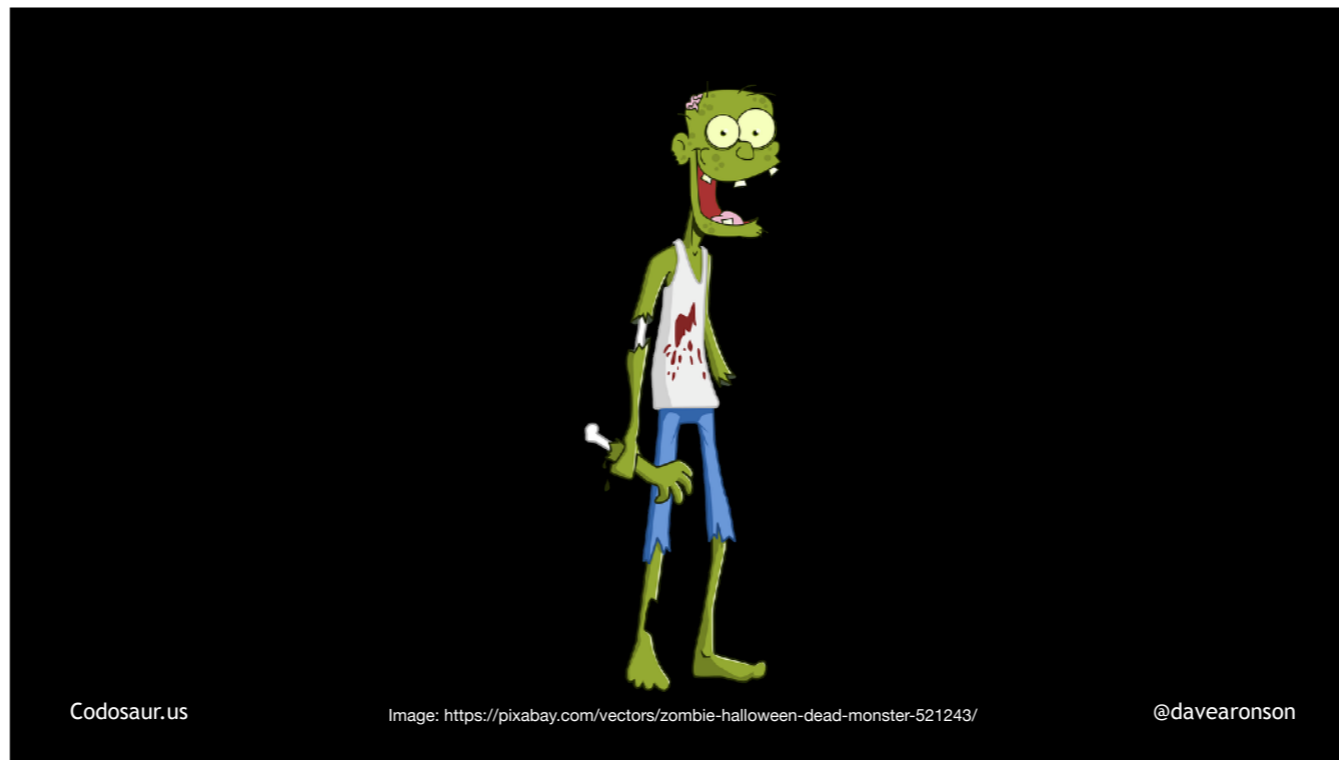


Codosaur.us

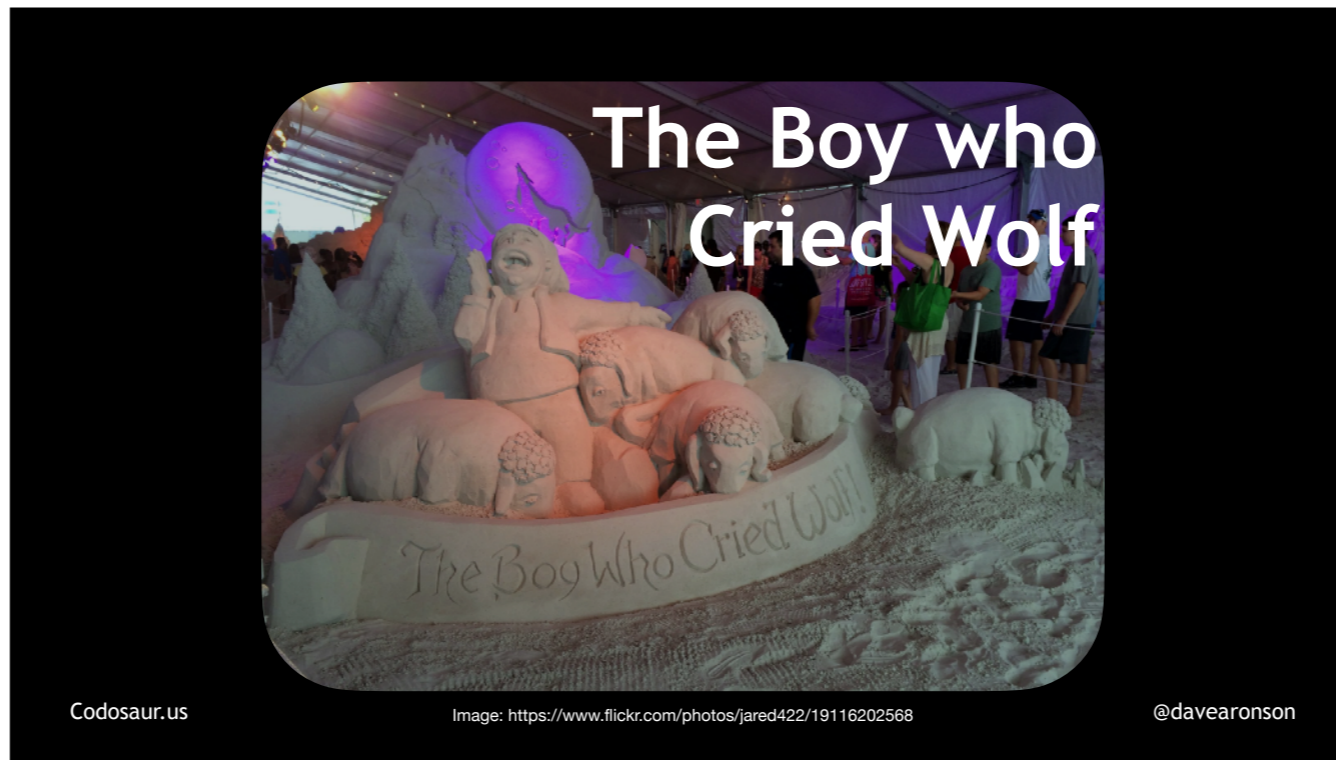
Image: <https://pxhere.com/en/photo/717939>

@davearonson

. . . not at all clear what to do about the results! It tells us that some particular change to the code made no difference to the test results, but what does *that* even mean? It takes a lot of interpretation to figure out what a mutant is trying to tell us. Their accent is verrah strayinge, and they're almost as incoherent as . . .



. . . zombies, but with a much bigger vocabulary, so they're not always on about braaaaaains. They're *usually* trying to tell us that our code is meaningless, or our tests are lax, or both, but it can be very hard to figure out exactly *how!* Even worse, sometimes it's a . . .



Codosaur.us

Image: <https://www.flickr.com/photos/jared422/19116202568>

@davearonson

. . . false alarm, because the mutation didn't make a test fail, but it didn't make any actual difference in the first place. It can still take quite a lot of time and effort to figure *that* out.

And even if a mutation *does* make a difference, most programs have quite a lot of code that we just . . .



. . . *shouldn't bother* to test, like debugging traces! Fortunately, most tools have ways to tell them "don't bother mutating this line", or even this whole function, class, file, or whatever . . . but that's usually with comments, which can clutter up the code, and make it less readable.

Now that we've seen some of the pros and cons, how does mutation testing work, unlike this guy? It . . .

Point Mutations

DNA ——— **mRNA**

Normal
DNA: GUUCGAUUGA / CAAGCTAACT
mRNA: GUUCGAUUGA

Missense
DNA: GUUCGUUGA / CAAGCAACT
mRNA: GUUCGUUGA

Frameshift insertion
DNA: GUUCGGAUUGA / CAAGGCTAACT
mRNA: GUUCGGAUUGA

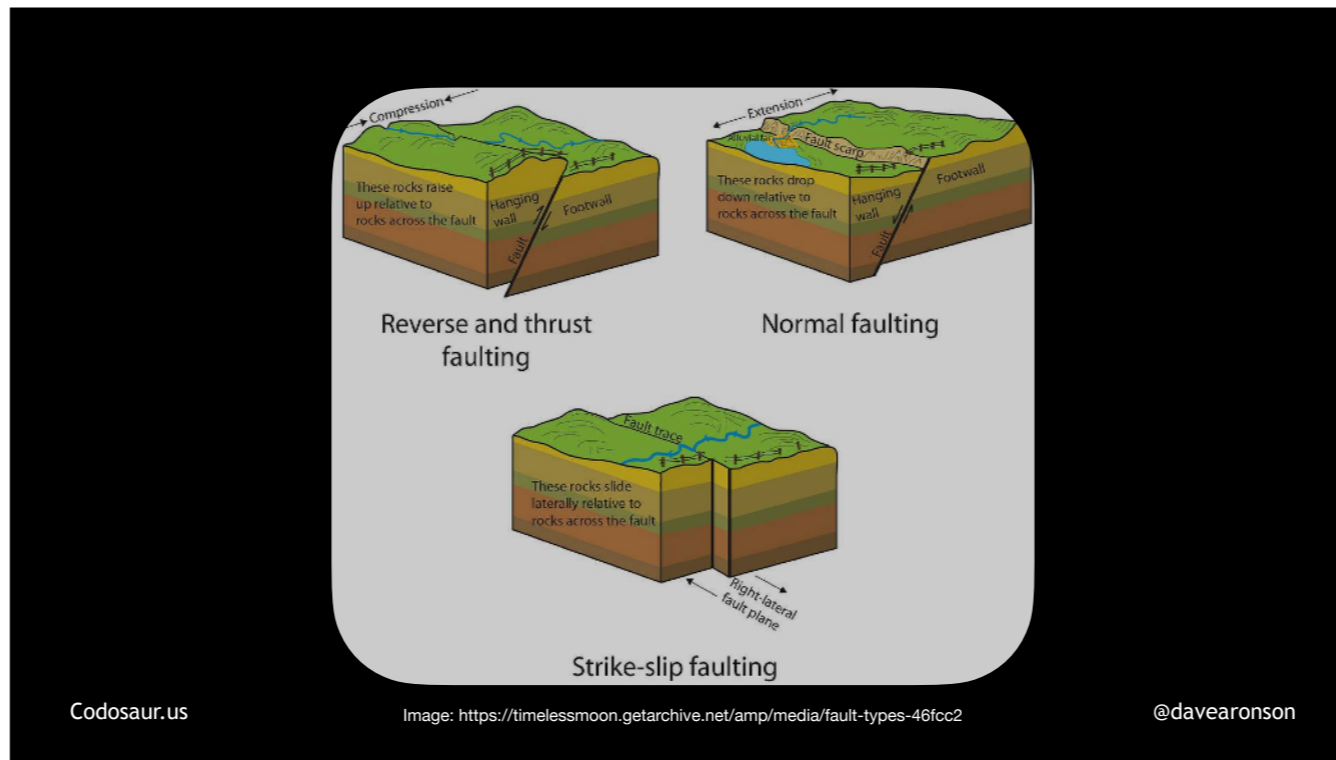
Frameshift deletion
DNA: GUUCUUGA / CAAGAACT
mRNA: GUUCUUGA

Nonsense
DNA: GUUUGG / CAATCG
mRNA: GUUUGG (STOP)

NATIONAL CANCER INSTITUTE

Codosaur.us Image: https://commons.wikimedia.org/wiki/File:Thrust_with_fault_propagation_fold.svg @davearonson

. . . *mutates* copies of our code, hence the name. It does this to create test failures, also known as . . .



. . . faults. So, mutation testing can be categorized as a “*fault-based*” testing technique. And this means that it is related to something you might already be familiar with:



. . . Chaos Monkey, from Netflix. Just like Chaos Monkey uses faults to help Netflix discover flaws in their error recovery, mutation testing uses faults to help us discover certain flaws in our tests and our code. But the way mutation testing does it, is sort of . . .



Codosaur.us

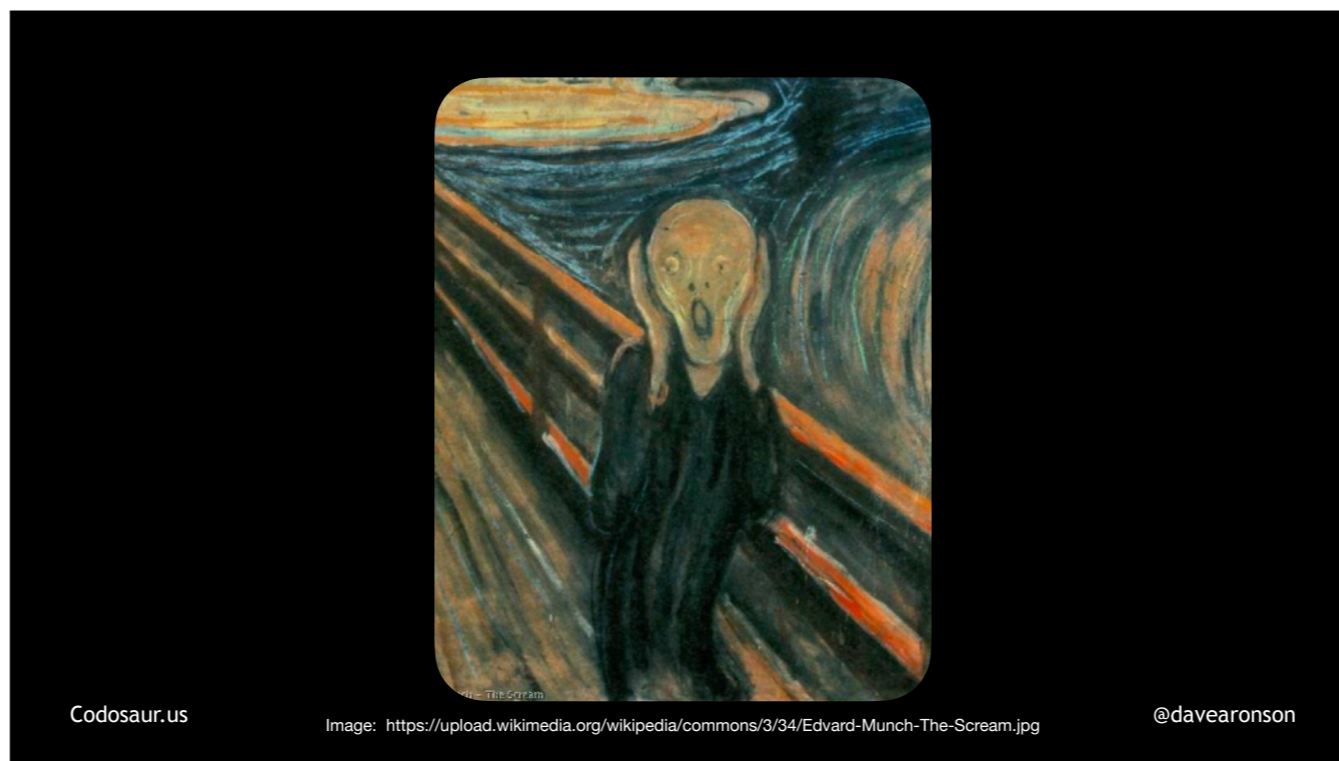
Image: <https://github.com/Netflix/chaosmonkey/raw/master/docs/logo.png>
(used for educational Fair Use purposes)

@davearonson

. . . upside down from what Chaos Monkey does. Chaos Monkey is best known for . . .



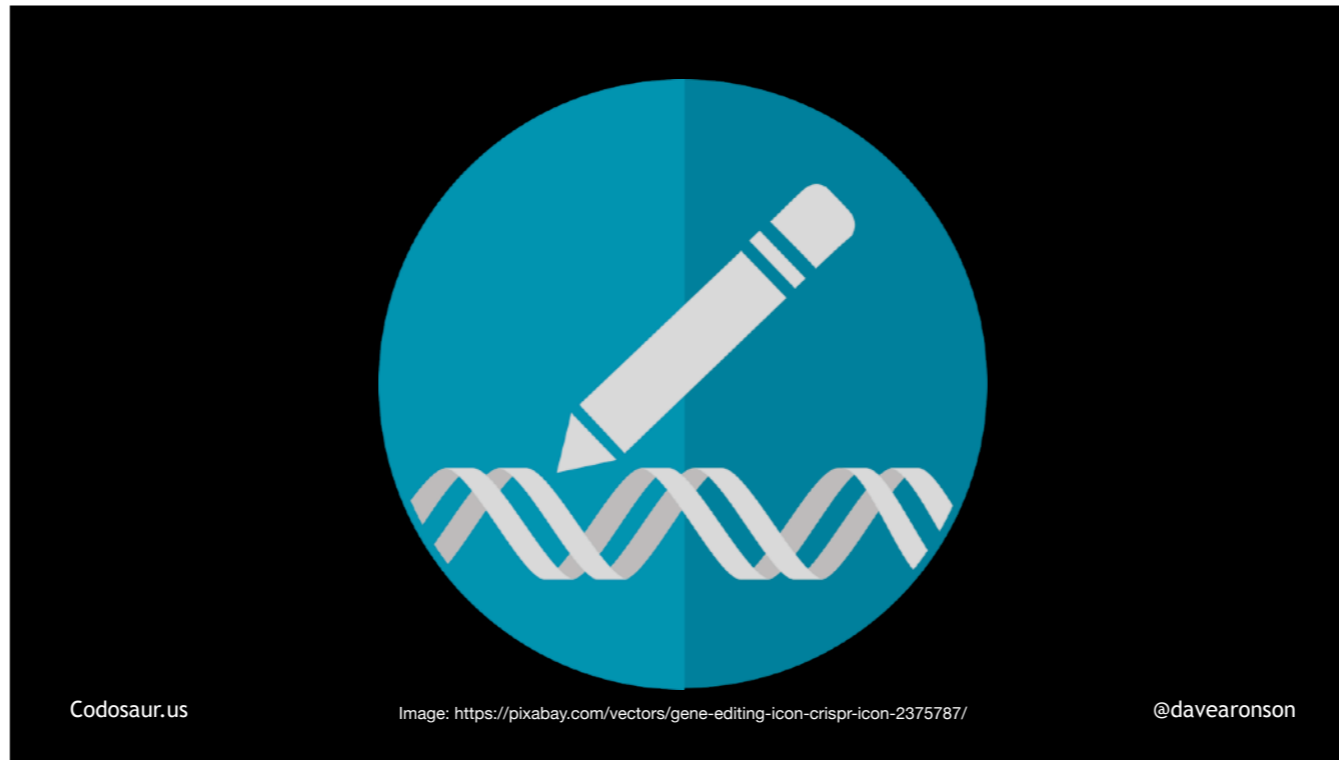
. . . injecting faults into Netflix's production network. (QUICK-CUT TO NEXT SLIDE!)



If all still goes well, in the sense that Netflix's customers don't notice, and their metrics still look good, then Netflix knows that their error recovery is working fine. Mutation testing, however, injects *semantic* . . .



. . . *changes*, not necessarily *problems*. We *hope* each of these changes will create faults, but that depends on the test suite. It injects them *into* . . .



. . . copies of our code, not our actual network, and it does this in our . . .



. . . *test* environment, not production. (Whew!) And if everything still goes well, *in the sense that* . . .

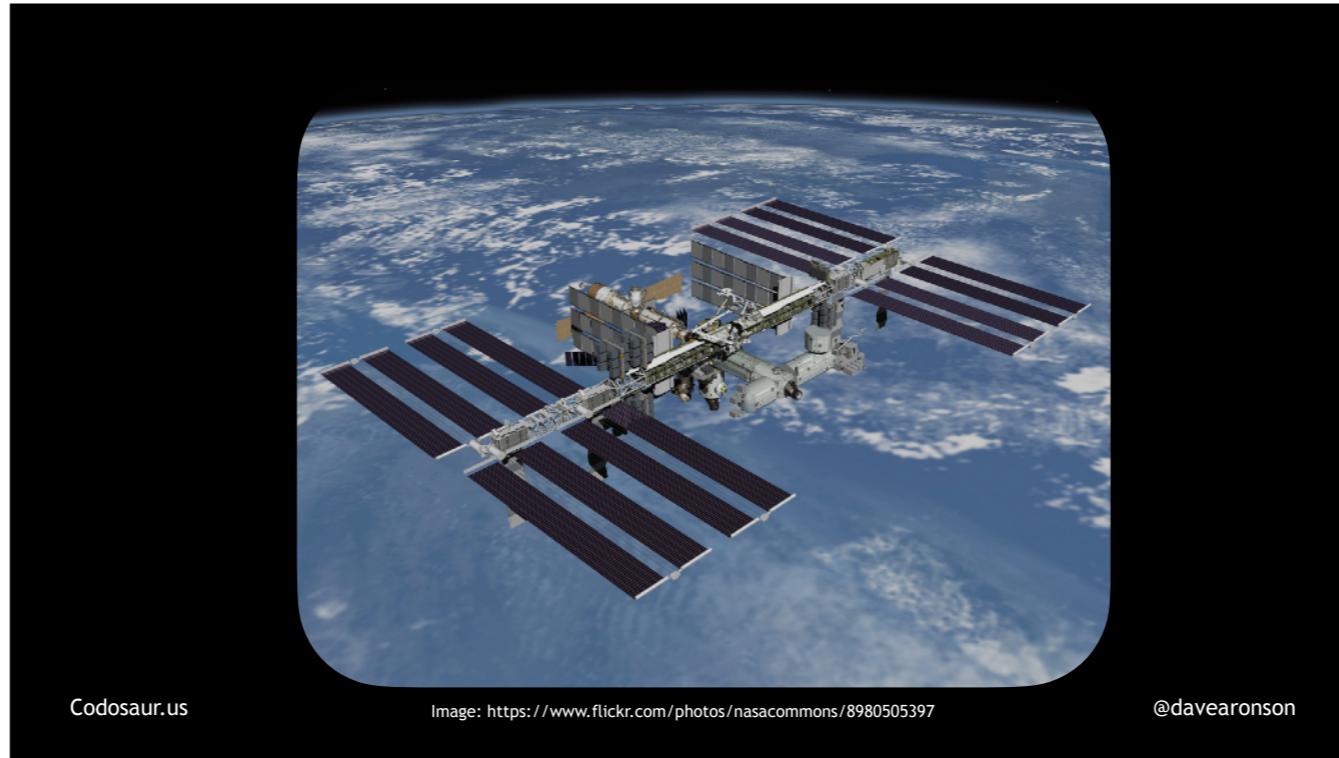


. . . fuzzing, a security penetration technique involving throwing lots of random data at an application. Mutation testing is like fuzzing our *code* rather than the *data*, plus it's . . .



. . . not random. These tools have a set of mutations they know how to do. The smarter ones can use the results of simpler mutations, to know they don't need to bother with more complex ones, so they may sometimes do different things and therefore *look* random, but they're not.

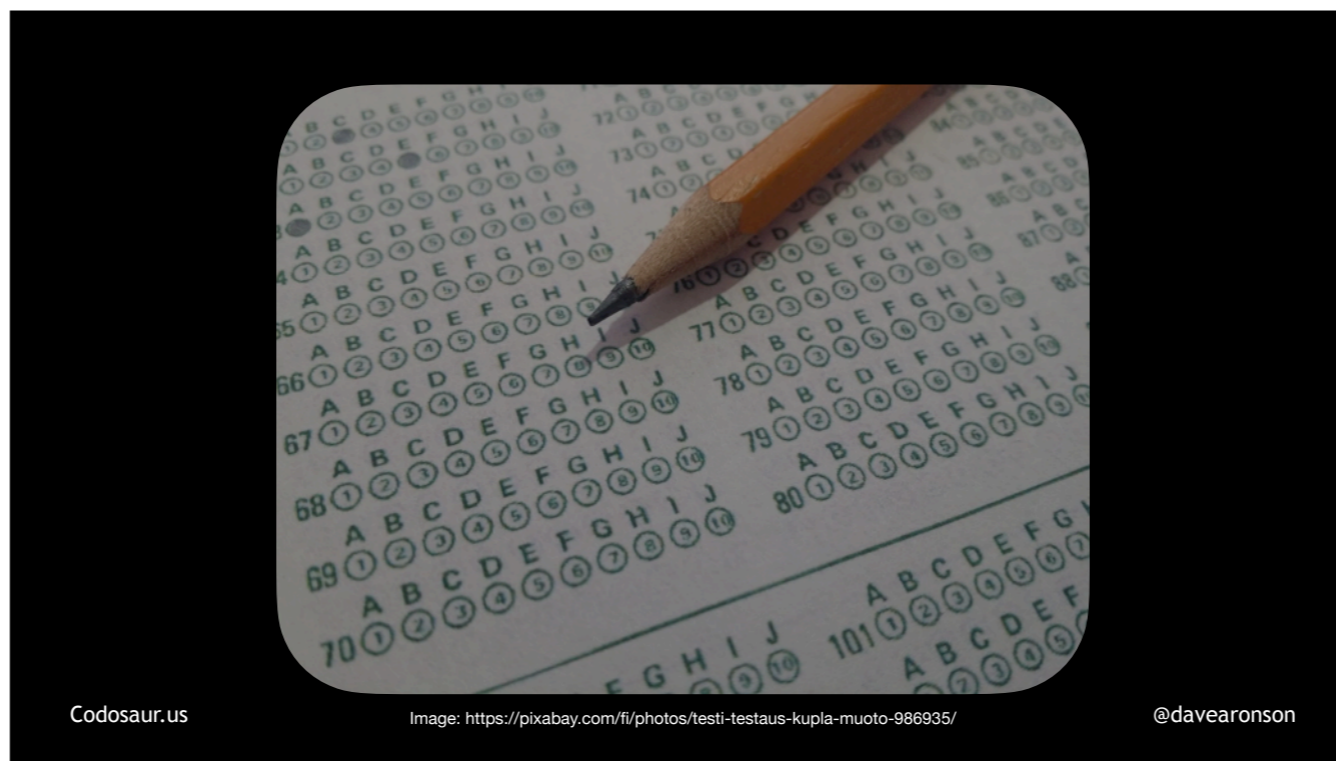
But enough about differences. What does mutation testing *do*, and how? Let's start with . . .



. . . a high-level view. First, our chosen tool . . .



. . . breaks our code apart into pieces to test. Usually, these are our functions -- or methods if we're *object-oriented*, but I'm just going to say functions. Then, for each function, it finds . . .



Codosaur.us

Image: <https://pixabay.com/fi/photos/testi-testaus-kupla-muoto-986935/>

@davearonson

. . . the *tests* that cover that function. If it can't find any, most will just skip this function, ideally warning us so we know we need to add or annotate some tests. Some, though, will use the entire test suite, which of course is horribly inefficient. If we're not skipping this function, then next the tool . . .



. . . makes mutants from that function. To do that, it looks closely at it to see how it can be changed. For each tiny little way the tool sees to change it, the tool makes . . .



. . . one mutant, with *that one mutation*.

Once our tool is done creating all the mutants it can for a given function, it iterates over . . .



Codosaur.us

Image: <https://www.flickr.com/photos/39160147@N03/15074089655>

@davearonson

. . . that list. And now we get to the heart of the concept.

Mutating function whatever, at something.py:42											
Test #	1	2	3	4	5	6	7	8	9	10	Result
Mutant #											
1	✓	✓	✓	✓	🕒						In Progress
2											To Do
3											To Do
4											To Do
5											To Do

Codosaur.us

@davearonson

This chart represents the progress of our tool. The tools generally don't give us quite all this information, let alone so neatly organized, but it's a conceptual model I use to help illustrate the point.

For each . . .

Mutating function whatever, at something.py:42

Test #	1	2	3	4	5	6	7	8	9	10	Result
Mutant #											
1	✓	✓	✓	✓	🕒						In Progress
2											To Do
3											To Do
4											To Do
5											To Do

. . . mutant, derived from . . .

Mutating function whatever, at something.py:42

Test #	1	2	3	4	5	6	7	8	9	10	Result
Mutant #											
1	✓	✓	✓	✓	🕒						In Progress
2											To Do
3											To Do
4											To Do
5											To Do

... a given function, the tool runs the function's ...

Mutating function whatever, at something.py:42

Test #	1	2	3	4	5	6	7	8	9	10	Result
Mutant #											
1	✓	✓	✓	✓	🕒						In Progress
2											To Do
3											To Do
4											To Do
5											To Do

. . . tests, but it runs them . . .

Mutating function whatever, at something.py:42

Test #	1	2	3	4	5	6	7	8	9	10	Result
Mutant #											
1	✓	✓	✓	✓	🕒						In Progress
2											To Do
3											To Do
4											To Do
5											To Do

. . . using the *current mutant* in place of the original function.

(PAUSE) If any test . . .

Mutating function whatever, at something.py:42

Test #	1	2	3	4	5	6	7	8	9	10	Result
Mutant #											
1	✓	✓	✓	✓	✗						In Progress
2											To Do
3											To Do
4											To Do
5											To Do

... *fails*, this is called ...



. . . “killing the mutant”, and it’s a . . .



. . . *good* thing. It means that our code is *meaningful* enough that the change that the tool made, to *create* this mutant, made a noticeable difference in the function's behavior, *and* that at least one test *noticed* that difference, and failed. Then, the tool will . . .

Mutating function whatever, at something.py:42

Test #	1	2	3	4	5	6	7	8	9	10	Result	
Mutant #												
1	✓	✓	✓	✓	✗						Killed	
2											To Do	
3											To Do	
4											To Do	
5											To Do	

. . . mark that mutant killed, . . .

Mutating function whatever, at something.py:42

Test #	1	2	3	4	5	6	7	8	9	10	Result	
Mutant #												
1	✓	✓	✓	✓	✗	-	-	-	-	-	Killed	
2											To Do	
3											To Do	
4											To Do	
5											To Do	

. . . stop running any more tests against it, and . . .

Mutating function whatever, at something.py:42											
Test #	1	2	3	4	5	6	7	8	9	10	Result
Mutant #											
1	✓	✓	✓	✓	✗	-	-	-	-	-	Killed
2	🕒										In Progress
3											To Do
4											To Do
5											To Do

Codosaur.us

@davearonson

. . . move on to the next one. Once a mutant has made *one* test fail, we don't care how many more it *could* make fail. Like so much in computers, we only care about ones and zeroes.

On the other claw, if a mutant . . .

Mutating function whatever, at something.py:42											
Test #	1	2	3	4	5	6	7	8	9	10	Result
Mutant #											
1	✓	✓	✓	✓	✗	-	-	-	-	-	Killed
2	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	In Progress
3											To Do
4											To Do
5											To Do

... lets all the tests pass, then the mutant is said to have ...

Mutating function whatever, at something.py:42											
Test #	1	2	3	4	5	6	7	8	9	10	Result
Mutant #											
1	✓	✓	✓	✓	✗	-	-	-	-	-	Killed
2	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	Survived!
3											To Do
4											To Do
5											To Do

... *survived*. That means that the mutant has the ...



. . . superpower of mimicry, skilled enough to *fool our tests!* This usually means that our code is meaningless, or our tests are lax, or both — and now it's up to us to figure out how.

Now let's peel back one . . .



. . . layer of the onion, and look at some *technical details* of how this works. First, our tool parses . . .

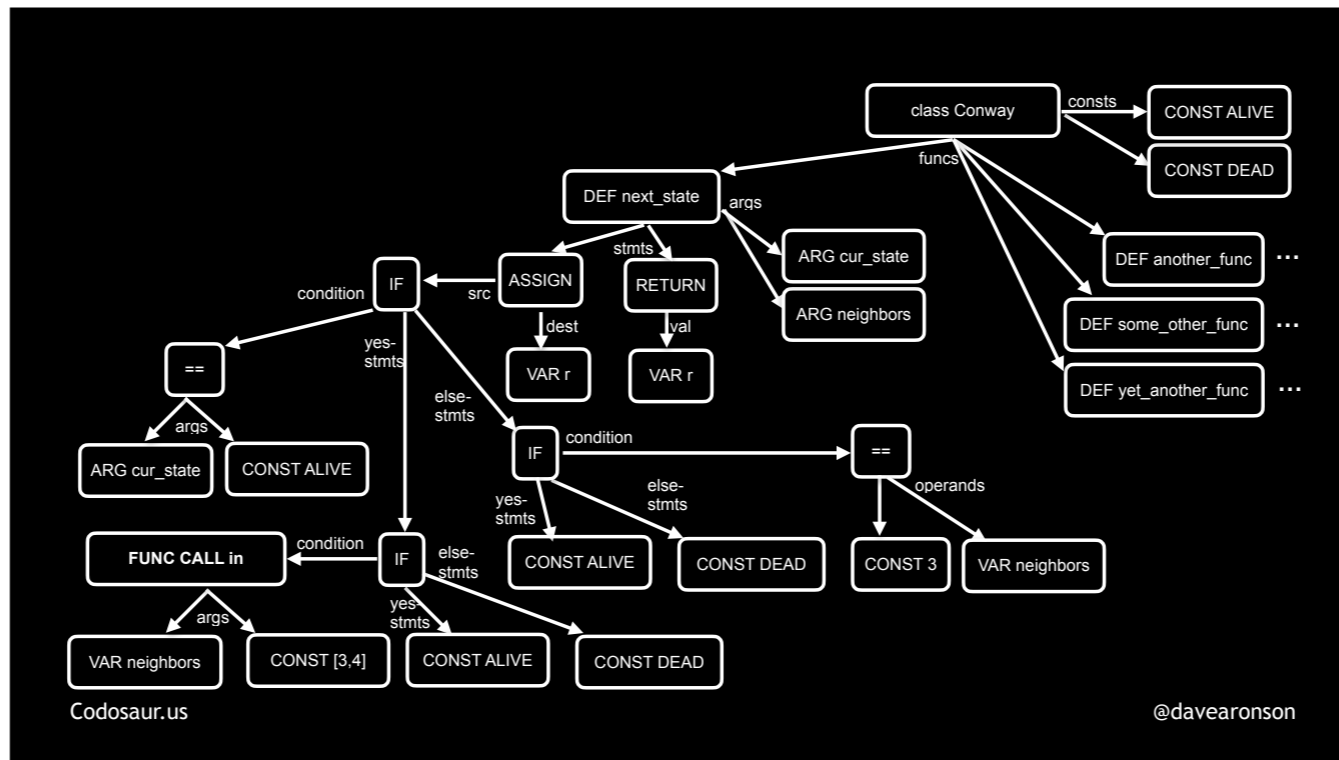
```
class Conway:
    ALIVE = "*"
    DEAD = " "

    @classmethod
    def next_state(cls, cur_state, neighbors):
        if cur_state == cls.ALIVE:
            result = cls.ALIVE if neighbors in [2,3] else cls.DEAD
        else:
            result = cls.ALIVE if neighbors == 3 else cls.DEAD
        return result

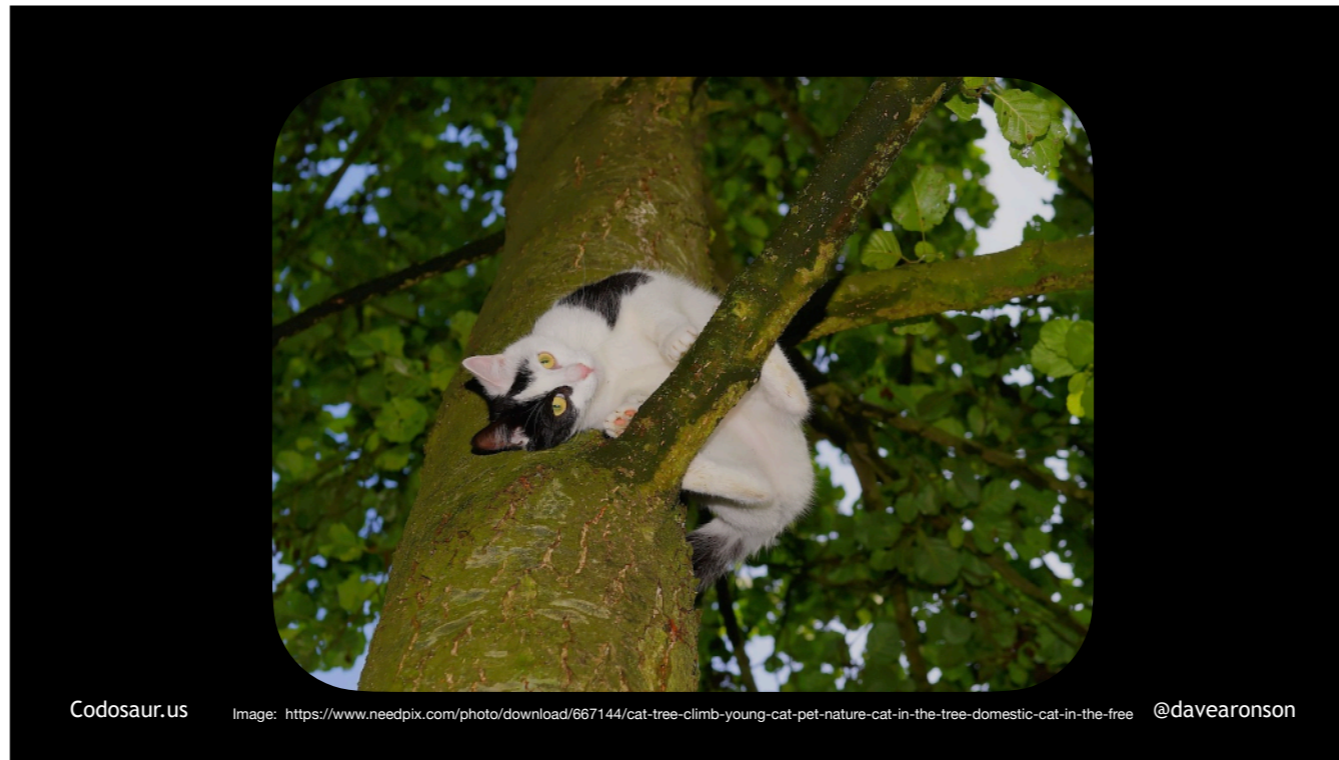
    def another_func:
        # whatever
    def some_other_func:
        # whatever
    def yet_another_func:
        # whatever
```

Codosaur.us @davearonson

. . . our code, usually into an Abstract Syntax Tree. Some work on bytecode, or even raw source code, but most of the tools that are actually ready for real use, use an AST, so let's just go with that approach. So, this code (which you don't really need to read and understand) becomes . . .



. . . this Abstract Syntax Tree (which you also don't really need to understand, except to notice that it includes some function definitions, rooted at the . . .



. . . traverses the tree, looking for sub-trees, or branches if you will, that represent each function. After finding *them*, it looks for each one's *tests*. That usually relies mainly on us developers, either . . .

```
@mumu tests-for foo
```

```
def test_foo_turns_3_into_6:  
  foo(3).must_equal 6
```

```
def test_foo_turns_4_into_10:  
  foo(4).must_equal 10
```

Codosaur.us

@davearonson

... annotating our tests, or following some kind of ...

```
def test_foo_turns_3_into_6:  
    foo(3).must_equal 6
```

```
def test_foo_turns_4_into_10:  
    foo(4).must_equal 10
```

Codosaur.us

@davearonson

... naming convention. These manual techniques are often supplemented and sometimes even replaced by ...

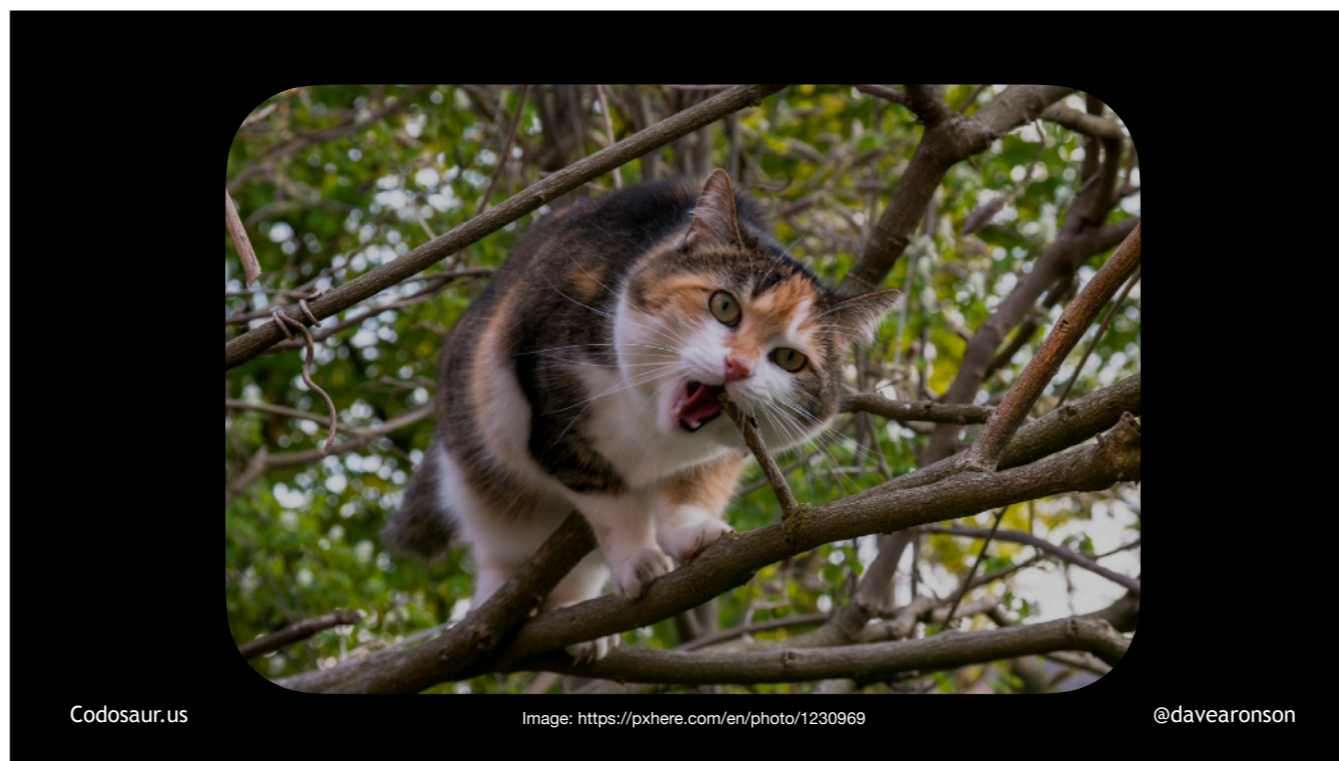
```
def test_foo_turns_3_into_6:  
    foo(3).must_equal 6
```

```
def test_foo_turns_4_into_10:  
    foo(4).must_equal 10
```

Codosaur.us

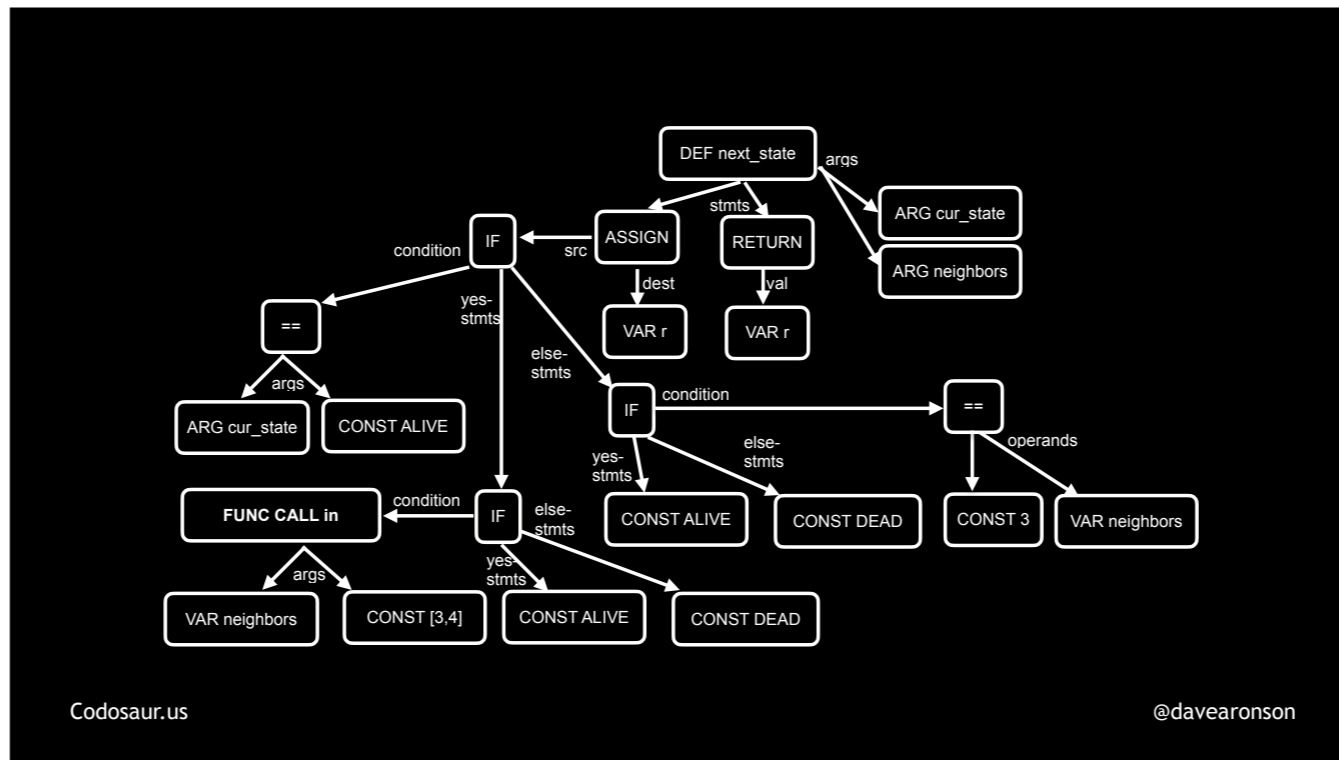
@davearonson

. . . the tool looking at what tests call what functions, out of our own codebase. Anyway, assuming it's not skipping this function due to not *finding* any tests, next it makes the mutants. To make mutants *from* an AST subtree, it . . .

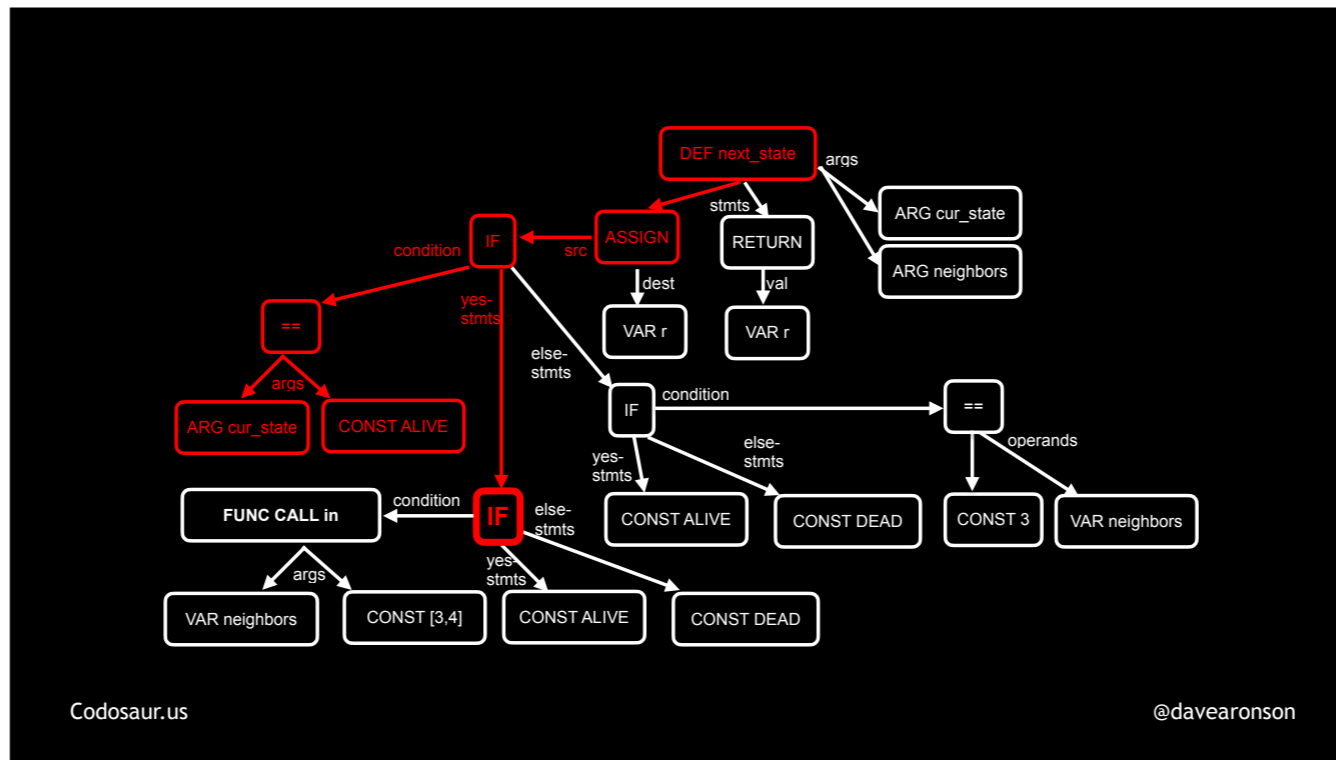


. . . traverses that subtree, just like it did to the whole thing. But now, instead of looking for even *smaller* subtrees it can *extract*, like twigs or something, it's looking for *nodes* where it can *change* something. Each time it finds one, then for each way it can change that node, it makes one copy of the function's AST subtree, with that one node changed, in that one way.

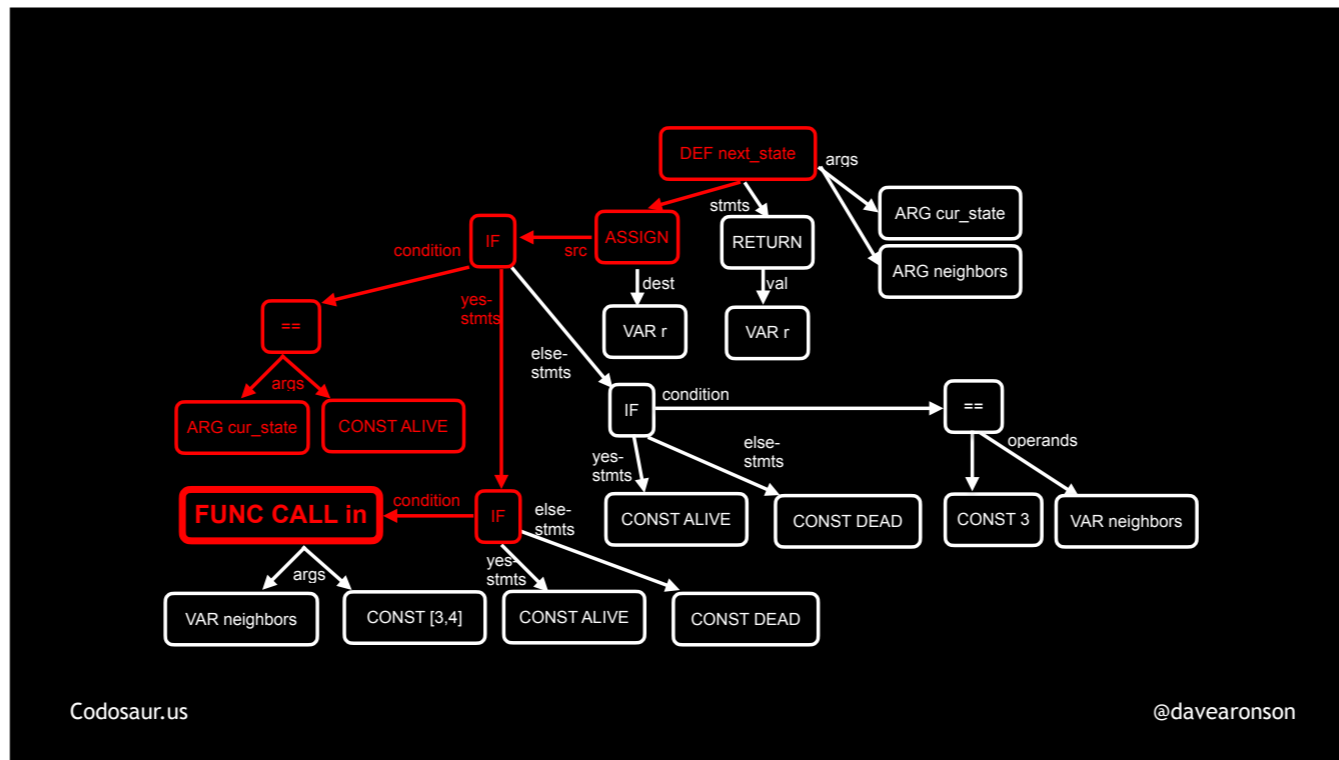
For instance, suppose our tool has started traversing . . .



. . . the function subtree from that AST I showed earlier, and has gotten down to . . .



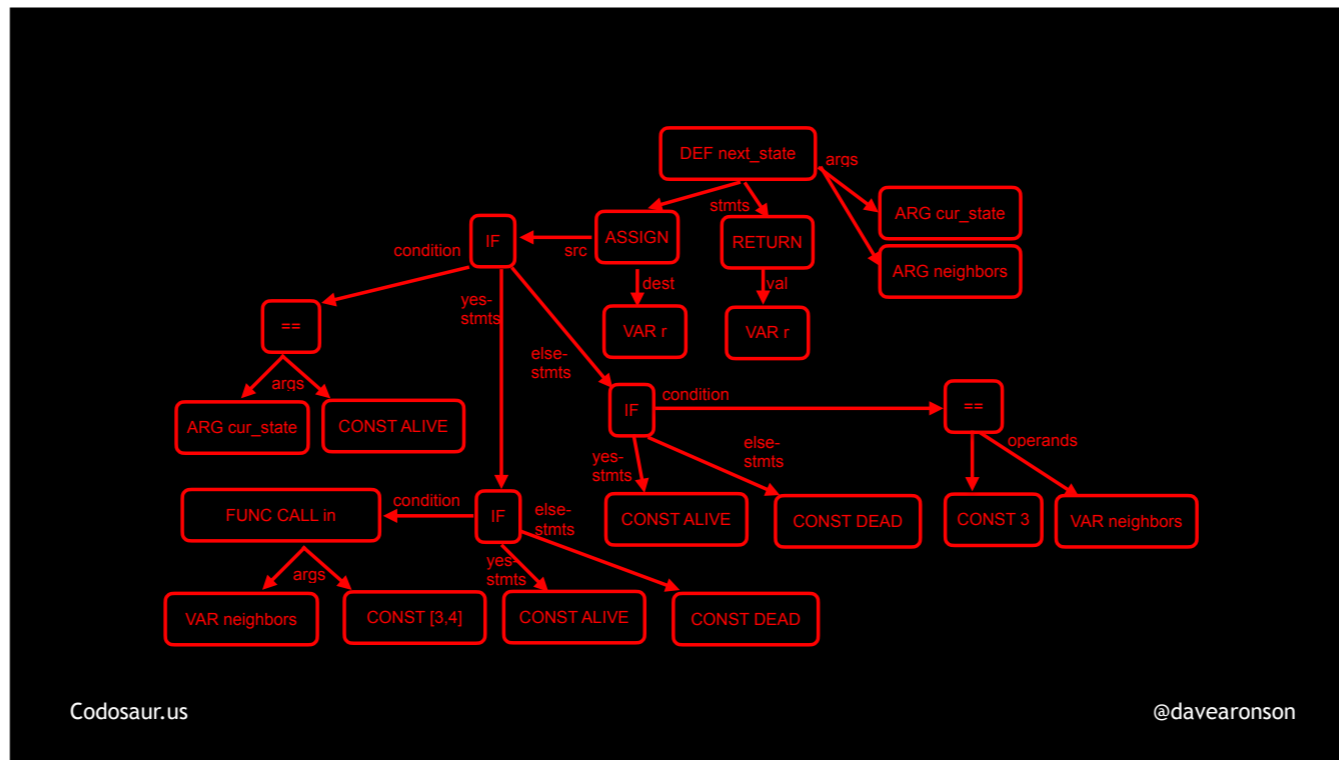
... this if statement. For each way the tool could change that node, it would make a fresh copy, of this whole subtree, with only that one node changed, in that one way. In other words, it would create a mutant with that mutation. After it's done making as many mutants as it can by mutating *that* node, it would continue traversing the subtree, to ...



Codosaur.us

@davearonson

... the *next* node. Again, for each way it could change *that* node, it would make a copy of this whole subtree, with only that mutation. And so on, until it has ...



... traversed the entire subtree.

Now, I've been talking a lot about changing things, so what kind of changes are we talking about? There are quite a lot!

`x + y` could become: `x - y`
`x * y`
`x / y`
`x ** y`

`x || y` could become: `x && y`
`x ^ y`

`x | y` could become: `x & y`
`x ^ y`

Maybe even swap *between sets!*

Codosaur.us

@davearonson

It could change a mathematical, logical, or bitwise operator from one to another. When the language allows, it could even cross these categories. For instance, in many languages, we can treat anything as a boolean, and all kinds of integers as bitfields, so `x TIMES y` could become `x AND y`, or `x BITWISE-EXCLUSIVE-OR y`, and so on.

$x - y$ could *also* become $y - x$

x / y could *also* become y / x

$x ** y$ could *also* become $y ** x$

" x " + " y " could *also* become " y " + " x "

When the *order* of operands matters, it could *swap* them.

`x < y`

could become:

`x <= y`

`x == y`

`x != y`

`x >= y`

`x > y`

It could change a *comparison* from one to another.

x

could become:

$-x$

$!x$

$\sim x$

. . . or vice-versa!

It could insert *or remove* a mathematical, logical, or bitwise *negation*.

```
a = foo(x)  
b = bar(y)
```

could become:

```
a = foo(x)
```

or

```
b = bar(y)
```

Codosaur.us

@davearonson

It can remove an entire *statement* or *expression*.

```
if x == y:  
    foo(z)
```

could become:

```
foo(z)
```

It can remove an if-condition, so that something that might be done or not, is always done.

```
while x == y:  
    foo(z)
```

could become:

```
foo(z)
```

It can remove a looping condition, so something that might be done, skipped, or done multiple times, is always done once.

```
def f(x, y): # lots of code here
could become:
def f(x, y): return 0
def f(x, y): return :math.max_int
def f(x, y): return "a string"
def f(x, y): return nil
def f(x, y): return x # or y
def f(x, y): fail("kaboom")
def f(x, y): # nothing
etc.
```

Codosaur.us

@davearonson

It could replace a function's *entire contents* with returning a constant, or any of the arguments, or raising an error, or nothing at all.

```
42      43      "42"      math.min_int
could    41      [42]      math.max_int
become:  -42     {42}     math.min_float
         1      []      math.max_float
         0      ()     math.infinity
         -1     {}     math.epsilon
         42.1   None
         41.9
```

Codosaur.us

@davearonson

It could change a value to some other value, such as changing 42 to any of these, and many more but I had to stop somewhere. It could even change it to something of a different and possibly incompatible type, such as changing a number into a, if I may quote . . .



. . . Gollum, “string, or nothing!”

There are *many* many more types of changes, but I trust you get the idea!

From here on, there are no more low-level details I want to add, so let’s *finally* walk through some *examples!* We’ll start with an easy one. Suppose we have a function . . .

```
def power(x, y):  
    x ** y
```

Codosaur.us

@davearonson

. . . like so. Never mind *why*, it just makes a good simple example, so let's just roll with it.

Think about what a mutant made from this might *return*, since that's what our tests would probably be looking at. It sure doesn't look like it has side effects.

Mainly, such a mutant could return results such as . . .

```
x + y
x - y
x * y
x / y
y ** x
x
y
0
1
-1
0.1
-0.1
math.min_int
math.max_int
math.max_float
math.min_float
math.infinity
math.epsilon
raise(DeliberateError)
"some random string"
[]
()
{}
None
and many more
```

Codosaur.us

@davearonson

. . . any of *these* expressions or constants, and, again, many more but I had to stop somewhere.

Now suppose we had only one test . . .

```
assert power(2, 2) == 4
```

Codosaur.us

@davearonson

. . . like so. This is a rather poor test, and I think at least one reason why is clear to most of us, but even so, *most* of those mutants on the previous slide *would get killed* by this test, the ones shown . . .

<code>x + y</code>	<code>math.min_int</code>
<code>x - y</code>	<code>math.max_int</code>
<code>x * y</code>	<code>math.max_float</code>
<code>x / y</code>	<code>math.min_float</code>
<code>y ** x</code>	<code>math.infinity</code>
<code>x</code>	<code>math.epsilon</code>
<code>y</code>	<code>raise(DeliberateError)</code>
<code>0</code>	<code>"some-random-string"</code>
<code>1</code>	<code>[]</code>
<code>-1</code>	<code>()</code>
<code>0.1</code>	<code>{}</code>
<code>-0.1</code>	<code>None</code>
	and many more

Codosaur.us

@davearonson

... here in crossed-out green. The ones returning constants, are very unlikely to match. There's no particular reason a tool would put a 4 there, as opposed to zero, 1, -1, minimum and maximum signed and unsigned integers and floats, and other such significant numbers. Changing the exponentiation into subtracting one argument from the other gets us zero, dividing them gets us one, returning either one alone gets us two, and the mismatched types and deliberate errors will at *least* make the test not pass. But ...

```
x + y
x - y
x * y
x / y
y ** x
y
0
1
-1
0.1
-0.1
math.min_int
math.max_int
math.max_float
math.min_float
math.infinity
math.epsilon
raise(DeliberateError)-
"some-random-string"
[]
()
{}
None
and many more
```

Codosaur.us @davearonson

. . . addition, multiplication, and exponentiation in the reverse order, all get us the correct answer. Mutants based on *these* mutations will therefore "survive" our test.

So how do we see that happening? When we run our tool, it gives us a report, that looks roughly like . . .

```
function "power" (demo.py:42)  
has 4 surviving mutants:
```

```
42 - def power(x, y):  
42 + def power(y, x):
```

```
43 -     x ** y  
43 +     x + y
```

```
43 -     x ** y  
43 +     x * y
```

```
43 -     x ** y  
43 +     y ** x
```

Codosaur.us

@davearonson

... this. The format will vary greatly depending on exactly which tool we use, but *semantically*, the information should be the same. And that is that if we changed ...

```
function "power" (demo.py:42)  
has 4 surviving mutants:
```

```
42 - def power(x, y):  
42 + def power(y, x):
```

```
43 -     x ** y  
43 +     x + y
```

```
43 -     x ** y  
43 +     x * y
```

```
43 -     x ** y  
43 +     y ** x
```

Codosaur.us

@davearonson

. . . the function called power, in . . .

```
function "power" (demo.py:42)  
has 4 surviving mutants:
```

```
42 - def power(x, y):  
42 + def power(y, x):
```

```
43 -     x ** y  
43 +     x + y
```

```
43 -     x ** y  
43 +     x * y
```

```
43 -     x ** y  
43 +     y ** x
```

Codosaur.us

@davearonson

... file demo.py, at line 42 ...

```
function power (demo py:42)  
has 4 surviving mutants:
```

```
42 - def power(x, y):  
42 + def power(y, x):
```

```
43 - x ** y  
43 + x + y
```

```
43 - x ** y  
43 + x * y
```

```
43 - x ** y  
43 + y ** x
```

Codosaur.us

@davearonson

... in any of four ways, then all its tests still pass.

And, that those four ways are: ...

```
function "power" (demo.py:42)  
has 4 surviving mutants:
```

```
42 - def power(x, y):  
42 + def power(y, x):
```

```
43 - x ** y  
43 + x + y
```

```
43 - x ** y  
43 + x * y
```

```
43 - x ** y  
43 + y ** x
```

Codosaur.us

@davearonson

. . . to change the function declaration to swap the arguments, or . . .

```
function "power" (demo.py:42)  
has 4 surviving mutants:
```

```
42 - def power(x, y):  
42 + def power(y, x):
```

```
43 - x ** y  
43 + x + y
```

```
43 - x ** y  
43 + x * y
```

```
43 - x ** y  
43 + y ** x
```

Codosaur.us

@davearonson

. . . change the function body to change the exponentiation into addition or multiplication, or . . .

```
function "power" (demo.py:42)
has 4 surviving mutants:
```

```
42 - def power(x, y):
42 + def power(y, x):
```

```
43 -     x ** y
43 +     x + y
```

```
43 -     x ** y
43 +     x * y
```

```
43 -     x ** y
43 +     y ** x
```

Codosaur.us

@davearonson

... to change the body to swap the exponentiation's operands.

So what is ...

```
function "power" (demo.py:42)
has 4 surviving mutants:
```

```
42 - def power(x, y):
42 + def power(y, x):
```

```
43 -     x ** y
43 +     x + y
```

```
43 -     x ** y
43 +     x * y
```

```
43 -     x ** y
43 +     y ** x
```

Codosaur.us

@davearonson

. . . this set of surviving mutants trying to tell us? We can tell from a glance at . . .

```
def power(x, y):  
    x ** y
```

Codosaur.us

@davearonson

. . . our code, that it's probably not trying to tell us about redundant or unreachable code. The body is just one line, so that sort of problem is extremely unlikely. So it's probably a test gap! The question now boils down to, how are . . .

```
function "power" (demo.py:42)
has 4 surviving mutants:
```

```
42 - def power(x, y):
42 + def power(y, x):
```

```
43 -     x ** y
43 +     x + y
```

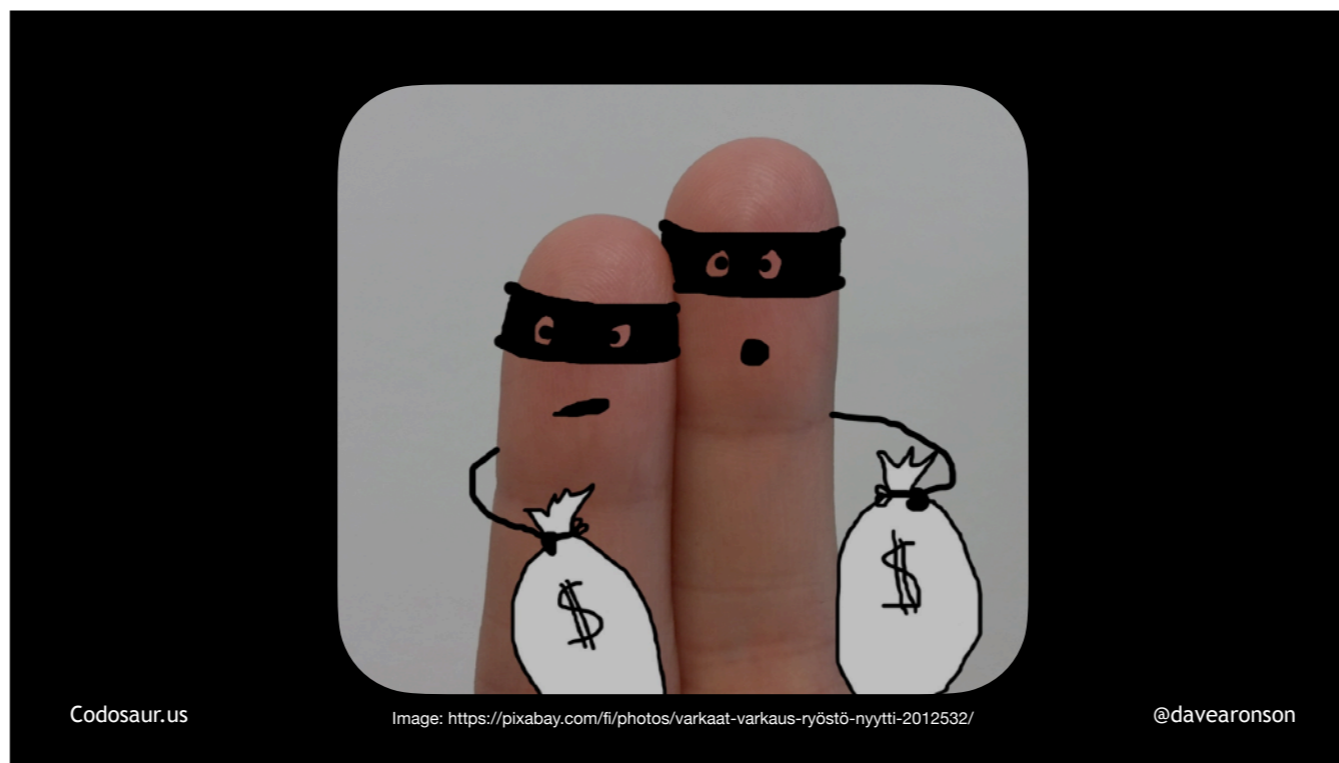
```
43 -     x ** y
43 +     x * y
```

```
43 -     x ** y
43 +     y ** x
```

Codosaur.us

@davearonson

... these mutants surviving? Are they ...



. . . pulling heists? Are they getting free room and board at the . . .



. . . Xavier Institute? Or what?

The usual answer is that . . .

```
mutant_power(x, y)
==
original_power(x, y)
```

Codosaur.us

@davearonson

. . . they return the same result as the original function. Or they have the same side effect — whatever our tests are looking at. To determine how that *happens*, it helps to take a closer look at it *along with* a test it passes. Let's start with . . .

the change:

```
43 - x ** y
```

```
43 + x + y
```

our test:

```
assert power(2, 2) == 4
```

Codosaur.us

@davearonson

... the "plus" mutant. Looking at the change, together with our test, makes it clear that this one survives because ...



. . . two *plus* two equals two *to* the two. (T-tu, t-tu, t-tu-t-tu-tu) (And so does two *times* two, but he's in the background, we can save him for later.)

So how can we *kill* . . .

the change:

```
43 - x ** y  
43 + x + y
```

our test:

```
assert power(2, 2) == 4
```

Codosaur.us

@davearonson

. . . this mutant, in other words, make at least one test fail when run against it, that would pass when run against the original code? To do that, we need to make at least one test use inputs such that *x plus y* is different from *x to the y*. For instance, we could add a test or change our existing test to . . .

```
assert power(2, 4) == 16
```

Codosaur.us

@davearonson

. . . assert that two to the *fourth* power is *sixteen*. All the mutants that our original test killed, *this would still kill*. But in addition, two *plus* four is six, not *sixteen*, so **this kills the plus mutant**. (See how that works?)

Better yet, two *times* four is eight, which is *also* not sixteen! We devs should certainly know our powers of two at least *that* well! So, this kills the "times" mutant as well. Killing one mutant often kills many other mutants of the same function.

But . . .



. . . the pair of argument-swapping mutants survive! What, how can that be? It's because . . .

$$4^{**} 2 == 16$$

$$2^{**} 4 == 16$$

Codosaur.us

@davearonson

. . . four squared is the same as two to the fourth, they're both sixteen. But that's not a big deal, we can . . .



. . . attack these mutants separately, no need to kill all the mutants in one shot and be some kind of superhero about it. To kill *them*, again, we can either add a test, or adjust an existing test, to something like . . .

```
assert power(2, 3) == 8
```

Codosaur.us

@davearonson

. . . this, asserting that two to the *third* power is *eight*. Three squared is nine, not eight, so **this kills the argument-swapping mutants**. Better yet, two *plus* three is five, two *times* three is six, and both of those are, guess what: not eight! So the "plus" and "times" mutants *stay* dead, and we don't get any . . .



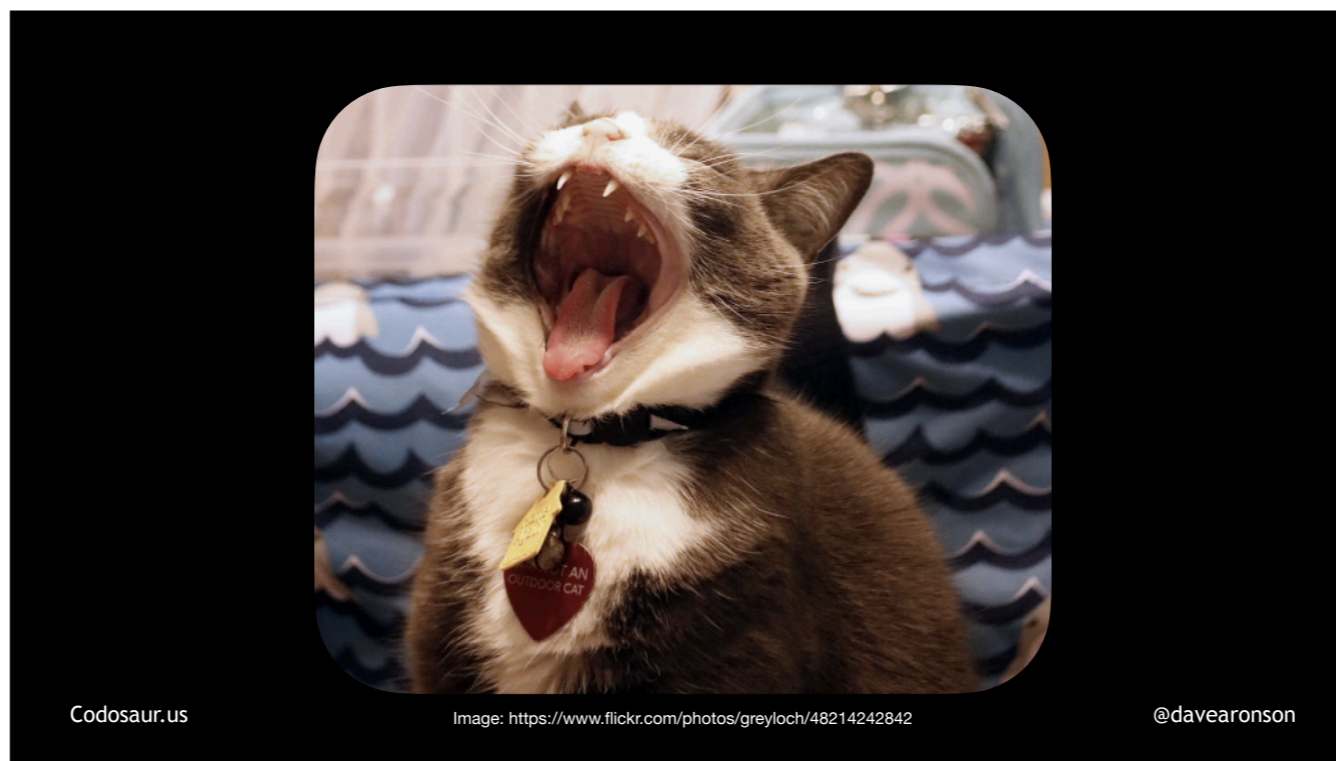
. . . zombie mutants wandering around, even if . . .

```
assert power(2, 3) == 8
```

Codosaur.us

@davearonson

. . . this were still our one and only test. (PAUSE!) With these inputs, the correct operation is the only simple common one that yields the correct answer. This isn't the *only* solution, though; even if we stuck to single digits, there are *lots* of ways to skin . . .



Codosaur.us

Image: <https://www.flickr.com/photos/greyloch/48214242842>

@davearonson

. . . *that* flerken!

This may make mutation testing sound simple, but this was a downright trivial example. So let's look at a more *complex* one!

Suppose we have a function to send a message, . . .

```
def send_message(buf, len):  
    sent = 0  
    while sent < len:  
        sent += send_bytes(buf, sent,  
                             len - sent)  
    return sent
```

Codosaur.us

@davearonson

. . . like so. This function, `send_message`, uses `send_bytes` to send as many bytes as `send_bytes` *could* send, like a woodchuck, looping to pick up where it left off, until the message is all sent. This is a very common pattern in communication software.

A mutation testing tool could make lots of mutants from this, but one of particular interest, would be . . .

```
def send_message(buf, len):  
    sent = 0  
    while sent < len:  
        sent += send_bytes(buf, sent,  
                            len - sent)  
    return sent
```

Codosaur.us

@davearonson

. . . this, an example of removing a looping control.

Now suppose that this mutant does indeed survive our test suite, which consists mainly of . . .

```
assert send_message(msg, size) == size
```

Codosaur.us

@davearonson

. . . *this*. (PAUSE!) There's a bit more that I'm not going to show you *quite* yet, dealing with setting the size and actually creating the message. But even without seeing that test code, what does the survival of that non-looping mutant tell us? (PAUSE!)

If a mutant that only goes through . . .

```
def send_message(buf, len):  
    sent = 0  
    while sent < len:  
        sent += send_bytes(buf, sent,  
                             len - sent)  
    return sent
```

Codosaur.us

@davearonson

. . . that loop once, acts the same as our normal code, as far as our tests can tell, that means that our *tests* are only making our *normal* code go through that loop once. So, what does *that* mean? (PAUSE!) By the way, you'll find that interpreting mutants often involves a lot of asking yourself "so, *what does that mean*", often deeply recursively!

In this case, it means that we're not testing sending a message larger than `send_bytes` can handle in one chunk! There are many ways that can happen, but we're only going to look at two possibilities. The most *likely* is that we should have, but simply *forgot*, or didn't *bother*, to test with a big enough message. For instance, . . .

```
in module Network:
```

```
max_chunk_size = 10_000
```

```
in test_send_message:
```

```
msg = "foo"
```

```
size = length(msg)
```

```
# other setup, like stubbing send_bytes
```

```
assert send_message(msg, size) == size
```

Codosaur.us

@davearonson

. . . suppose our maximum chunk size, what `send_bytes` can handle in one chunk, is 10,000 bytes. But . . .

```
in module Network:
max_chunk_size = 10_000

in test_send_message:
msg = "foo"
size = length(msg)
# other setup, like stubbing send_bytes
assert send_message(msg, size) == size
```

Codosaur.us

@davearonson

. . . we're only testing with an itty-bitty *three* byte message. (PAUSE!)

The obvious fix is to deliberately use a message larger than our maximum chunk size. With this kind of message, we can easily construct one, as shown . . .

```
in module Network:
```

```
max_chunk_size = 10_000
```

```
in test_send_message:
```

```
size = network.max_chunk_size + 1
```

```
msg = "x" * size
```

```
# other setup, like stubbing send_bytes
```

```
assert send_message(msg, size) == size
```

Codosaur.us

@davearonson

... here. (PAUSE!) We just take the maximum size, add some, and construct that big a message.

But now let's look at another possible cause and solution. Maybe we *did* test with the *largest* permissible message, out of a set of predefined messages, or at least message *sizes*. For instance, ...

```
in module Message:
```

```
SmallMsgSize = 1_000
```

```
LargeMsgSize = 5_000 # the largest
```

```
in test_send_message:
```

```
size = Message.LargeMsgSize
```

```
msg = Message.make_msg("a" * size)
```

```
# other setup, like stubbing send_bytes
```

```
assert send_message(msg, size) == size
```

Codosaur.us

@davearonson

. . . here we have Small and Large message sizes, we test with a Large, and yet, this mutant survives! In other words, we're still sending the whole message in one chunk. What could possibly be wrong with that? It sounds like a *good* thing to me! What is this mutant trying to tell us in this case? (PAUSE!)

In *this* scenario, it's trying to tell us that a version of `send_message` with the looping removed will do the job just fine. If we remove the looping, we wind up with . . .

```
def send_message(buf, len):  
    sent = 0  
    sent += send_bytes(buf, sent,  
                        len - sent)  
    return sent
```

Codosaur.us

@davearonson

. . . this. Now some other stuff is *clearly* redundant, because we only needed it to support the looping. If we *also* remove *that*, then it boils down to . . .

```
def send_message(buf, len):  
    return send_bytes(buf, 0, len)
```

Codosaur.us

@davearonson

. . . this. (PAUSE!) Now the message is clear: the *entire* send_message *function* may well be *redundant*, so we can just use send_bytes *directly!* In real-world code, though, it might not be, because there may be some . . .

```
def send_message(buf, len):
    # logging?
    res = send_bytes(buf, 0, len)
    # hi-lvl error handling?
    # other record-keeping?
```

Codosaur.us

@davearonson

logging, error handling, and so on, needed in `send_message`, that we can't shove down the stack into `send_bytes`, but at the very least, the *looping* was redundant. Fortunately, when it's this kind of problem, the usual solution is clear and easy, just rip out the extra junk that the mutant doesn't have. This will also make our code more *maintainable*, by getting rid of useless cruft that just gets in the way of understanding it.

Now that we've seen a few different examples, of spotting both bad tests and redundant code, I'd like to address some . . .



. . . occasionally asked questions — mutation testing is still rare enough that there are very few *frequently* asked questions. First, this all sounds pretty weird, deliberately making tests fail, to prove that the code succeeds! Where did this whole . . .



. . . bizarre idea come from? Mutation testing has a surprisingly . . .



Codosaur.us

Image: <https://www.flickr.com/photos/brickset/33236853148>

@davearonson

. . . long history -- at least in the context of computers. It was first proposed in 1971, in Richard Lipton's term paper titled "Fault Diagnosis of Computer Programs", at Carnegie-Mellon University. The first *tool* didn't appear until 1980, as part of Timothy Budd's PhD work at Yale. However, it wasn't *practical* on typical developer-grade computers, until the early 2000s, with significant advances in CPU *speed*, multi-*core* CPUs, larger and cheaper memory, and so on. But now, it's practical even on fairly low-end modern systems, like this 2020 M1 MacBook Air.

That leads us to the next question: *why* is it so CPU- and memory-intensive? To answer that, we need do some math, but don't worry, it's pretty basic. Suppose our functions have, on average, . . .

10 lines

Codosaur.us

@davearonson

. . . about ten lines each. And each line has about . . .

x **10 lines**
5 mutation points

Codosaur.us

@davearonson

. . . five places where it can be mutated, to any of about . . .

10 lines
x 5 mutation points
x 20 alternatives

. . . twenty alternatives. That works out to about . . .

$$\begin{array}{r} 10 \text{ lines} \\ \times 5 \text{ mutation points} \\ \times 20 \text{ alternatives} \\ \hline = 1000 \text{ mutants/function!} \end{array}$$

Codosaur.us

@davearonson

. . . a thousand mutants for each function! And for each one, we'll have to run somewhere between one test, if we're lucky and kill it on the first try, all the way up to *all* of that function's tests, if we kill it on the last try, or worse yet, it survives.

Suppose we wind up running just . . .

10 lines
x 5 mutation points
x 20 alternatives

= 1000 mutants/function!
x 20 % of the tests, each

Codosaur.us

@davearonson

. . . one *fifth* of the tests for each mutant, on average. Since we start with a thousand mutants, that's still . . .

10 lines
x 5 mutation points
x 20 alternatives

= 1000 mutants/function!
x 20 % of the tests, each

= 200 x as many test runs!

Codosaur.us

@davearonson

. . . *two hundred times* the test runs for that function, compared to regular testing. If our test suite normally takes a zippy ten seconds, then with these assumptions, mutation testing will take about *two thousand* seconds. That might not sound like much, because I'm saying "seconds", but it's over *half an hour!* I don't want to sit and wait for that!

But there is some . . .



Codosaur.us

Image: <https://www.flickr.com/photos/akuchling/50310316>

@davearonson

. . . good news! Over the past decade or so, there has been a lot of research on trimming down the number of mutants, mainly by weeding out ones that are semantically equivalent to the original code, redundant with other mutants, or trivial in various ways such as creating an obvious uncaught error condition. Such things have reduced the mutant horde down to about one third! But even with that rare level of success, it's still . . .



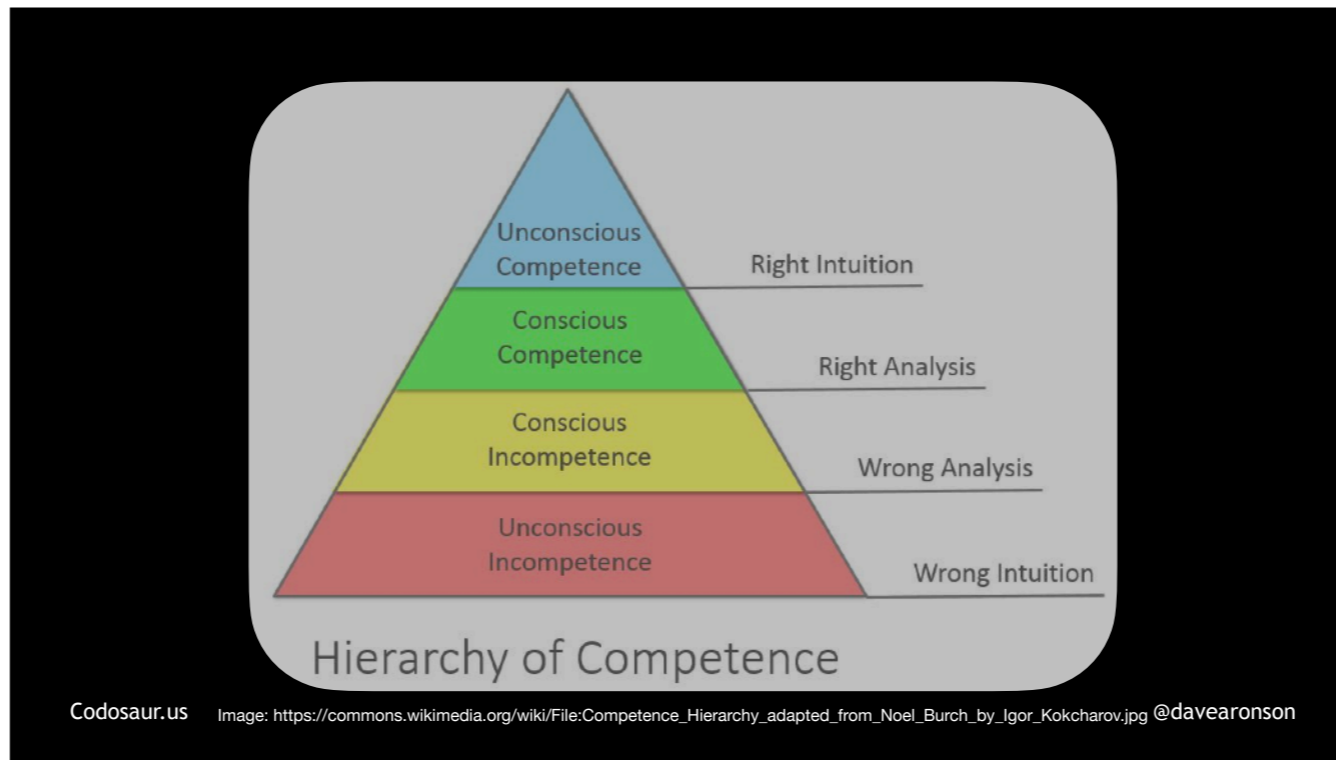
. . . no silver bullet, as this takes lots of CPU time itself -- and there are still quite a lot of mutants left to deal with.

The next question is, when making each mutant, why change it in only . . .



. . . one way?

There are multiple reasons. First off, the main theoretical underpinning of mutation testing is . . .



. . . the Competent Programmer hypothesis. Let's give that a quick check. Raise your hand if you're competent! (PAUSE!) Okay, looks like most of us. The rest of you, you probably really are competent, so you might want to read up on Impostor Syndrome.

Long story short, the Competent Programmer Hypothesis is the idea that we generally have a pretty good clue what we're doing, and when we make a mistake, it's usually a single small mistake, like adding when we should subtract, or comparing using "less than or equal" when we mean "strictly less than", or greater than, or whatever. Does this kind of simple substitution sound familiar? It's exactly what a mutation testing tool does! So we can think of mutation testing as sort of a "did you mean" function, like how Google suggests a different search if ours didn't have many hits.

There are also practical considerations. For one, it helps us poor humans . . .



Codosaur.us

Image: <https://pixnio.com/objects/camera/camera-focus-nostalgia-snapshot-lens-photography-old-mechanism>

@davearonson

. . . FOCUS! It's much easier to tell what a surviving mutant is trying to tell us, if we're only talking about one thing in the first place. Another reason is that multiple changes may . . .



. . . balance each other out, leading to more false alarms. Lastly, allowing multiple mutations would create a combinatorial . . .



. . . explosion of mutants, with the tool making multiple *orders of magnitude* more mutants per function, which would make it even *more* CPU- and memory-intensive. Never mind *running* the *tests*, just *creating* the *mutants* could get to be a huge workload! But we can avoid this huge workload, *and* the increased false alarms, *and* the lack of focus, if we just . . .



. . . limit it to one mutation per mutant.


Lastly, a more practical question: where should we fit this into . . .

- Claim Ticket and Make Branch
- Write Tests
- Write Code
- Lint?
- Refactor?
- Create Pull Request
- Get PR Approved
- Merge PR and Delete Branch
- Go Back, Jack, Do It Again

Codosaur.us

@davearonson


. . . our development process? Mainly, I think it belongs at *least* . . .

- Claim Ticket and Make Branch
- Write Tests
- Write Code
- Lint?
- Refactor?
- Create Pull Request
-  - Get PR Approved
- Merge PR and Delete Branch
- Go Back, Jack, Do It Again

Codosaur.us

@davearonson

. . . here, as part of the requirements for a Pull Request (or whatever your process uses) to be approved. You can set some standards for what you're willing to tolerate, such as no surviving mutants on new code and no increase of them on old code. Ideally this would be automated, as part of a CI pipeline, started automatically when the PR is created, stopping the job if the requirements aren't met. That said, I personally would also do it in my *own* work as part of . . .

- Claim Ticket and Make Branch
- Write Tests
- Write Code
-  - Lint?
- Refactor?
- Create Pull Request
- Get PR Approved
- Merge PR and Delete Branch
- Go Back, Jack, Do It Again

Codosaur.us

@davearonson

. . . the Linting step, where I apply all sorts of quality checking tools, to make sure my code is as good as possible, before making anyone else bother with it.

If you'd like to try mutation testing for yourself . . .

Alloy:	MuAlloy
Android:	mdroid+
C:	mutate.py, SRCIROR
C/C++:	accmut, dextool, MART, MuCPP, Mutate++, mutate_cpp, SRCIROR
C#/.NET/Mono:	nester, NinjaTurtles, Stryker.NET, Testura.Mutation, VisualMutator
Clojure:	mutant
Crystal:	crytic
Dart:	mutation_test
Elixir:	darwin, exavier, exmen, mutation, Muzak [Pro]
Erlang:	mu2
Etherium:	vertigo
FORTRAN-77:	Mothra (written in mid 1980s!)
Go:	go-mutesting, gremlins, ooze
Haskell:	fitspec, muCheck
Java:	jumble, major, metamutator, muJava, pit/pitest, and many more
JavaScript:	stryker, grunt-mutation-testing
Pharo:	MUTALK
PHP:	infection, humbug
PL/SQL:	MuPLSQL
Python:	cosmic-ray, mutmut, mutpy, pester, xmutant
Ruby:	mutant, mutest , heckle
Rust:	mutagen
Scala:	scalamu, stryker4s
Smalltalk:	mutalk
Solidity:	RegularMutator
SQL:	SQLMutation
Swift:	muter
Anything on LLVM:	llvm-mutate, mull
Tool to make more:	Wodel-Test (https://gomezabajo.github.io/Wodel/Wodel-Test/)

Codosaur.us

@davearonson

. . . here is a list of tools for some popular languages and platforms, and some others, I doubt many of you are doing FORTRAN-77 these days! The tools I *know* are outdated, are crossed out, but there may be others; I don't know or follow quite *all* of these languages and platforms. Don't worry about taking pictures, the last slide has the URL for the whole deck.

To summarize at last, mutation testing is a powerful technique to . . .

😊 Checks that code is meaningful

Codosaur.us

@davearonson

. . . help ensure that our code is meaningful and . . .

😊 Checks that code is meaningful

😊 Checks that tests are strict

. . . our tests are strict. It's . . .

- 😊 Checks that code is meaningful
- 😊 Checks that tests are strict
- 😊 Easy to get started with

Codosaur.us

@davearonson

easy to get started with, in terms of setting up most of the tools and annotating our tests if needed (which may be *tedious* and *time-consuming* but at least it's *easy*), but it's . . .

😊 Checks that code is meaningful

😊 Checks that tests are strict

😊 Easy to get started with

😞 **Difficult to interpret results**

. . . not so easy to interpret the results, nor is it . . .

😊 Checks that code is meaningful

😊 Checks that tests are strict

😊 Easy to get started with

😞 Difficult to interpret results

😞 Hard labor on the CPU

Codosaur.us

@davearonson

. . . easy on the CPU.

Even if these drawbacks mean it might not be a good fit for our current projects, I still think it's just . . .

- 😊 Checks that code is meaningful
- 😊 Checks that tests are strict
- 😊 Easy to get started with
- 😞 Difficult to interpret results
- 😞 Hard labor on the CPU
- 😎 Fascinating concept! 😎

Codosaur.us

@davearonson

. . . a really cool idea . . . in a geeky kind of way.

If you have any questions, . . .



T.Rex-2025@Codosaur.us
twitter.com/DaveAronson
linkedin.com/in/DaveAronson

Slides and FULL SCRIPT:
[Codosaur.us/reds/mutants-vdcern-25-slides](https://codosaur.us/reds/mutants-vdcern-25-slides)

Codosaur.us

@davearonson

. . . I'll take them now, and if you think of anything later, there's my contact info, plus the URL for the slides, complete with a full script, which I've *mostly* stuck to. Any questions?

AFTER APPLAUSE: When you get back to work, remember to use mutation testing to make sure your coverage report means what you think it means! Please remember to vote about the talks, and if you don't give me the top rating for whatever reason, please tell me what should be improved, and ideally how.