

Last time: performance (QoS - quality of service)

This class: correctness

Correctness is always important (compilers, query optimizers, etc.)

have program Q and convert to Q'

↑
org

↑
modified

prove they're
equivalent
 $Q \equiv Q'$

(reducible to
halting problem,
but at least if
test limit times,
can say they're
behaviorally
equivalent)

what if there isn't 2
programs to check
equivalence? what if only P?

need complete spec
(usually written in logic)

strong

2 kinds of correctness:

① total: assumes there is a spec
and you try to show P
is behaviorally equivalent to
spec
 $P \equiv \text{spec}$ ($P \leq \text{spec}$)

$$\forall s : O(s) \equiv O(P)$$

↑
Output

② partial:

ex: program never derefs null

(any program that does deref null is wrong)

ex: no uninitialized vars

ex: use after free (C/C++ problem)

ex: buffer overflow

ex: exceptions (e.g. file not found error)

ex: segv

rules
that
should
be true

We assume that any correct program
does not break these rules

implicit spec for programs

ex: if (x == 0) {

: } assignment, so will always fall through
linters (pattern checkers) would report
bad programming practices

disadv of linters?

precision
inverse of fp
report \Rightarrow bug

recall
bug \Rightarrow report

- 1000s of warning messages,
referred to as false positives

- also suffers from false negatives
(if bug \Rightarrow not part of pattern)

perfect recall:

for (i=0; i < n; i++) {

if there's a
bug, def report! }

cout << i << "no bugs" << endl

perfect precision;

cout << "no bugs" << endl

precision-recall
tradeoff

ppl now prefer
high precision in general

The Spec Problem

- writing spec is hard
- Spec is a program just in another lang spec is long
+ complicated
too
- is the spec correct? we don't know...
- who guards the guardians?
- Program is the spec (there's no real spec)

↳ paper: Bugs as Deviant Behavior

programs act weird sometimes means probably a bug

static v. dynamic analysis



Coverity: static analysis tool

collect factual info by running program in abstract interpretation
(execute symbolically)

- Code has infeasible paths in program, but static analysis does not know about these
- State Space explosion: static analysis needs to control this

ex: $x \geq 12$
 $x \leq 130$ } many states between start + end
losing precision

also lose precision w/ loops

ex: while (true) {
 x = x + 1
}

all we know is $x \geq x'$
long x

False pos example:

index ≥ 0 index $\leq 1 \text{ mil}$

a = new int [1000]

a[index] buffer
overflow?

Coverity ranks errors!

- ranked on state space, precision, etc

- better than giving someone 1mil errors
and saying good luck

most analyses
v.
many analyses

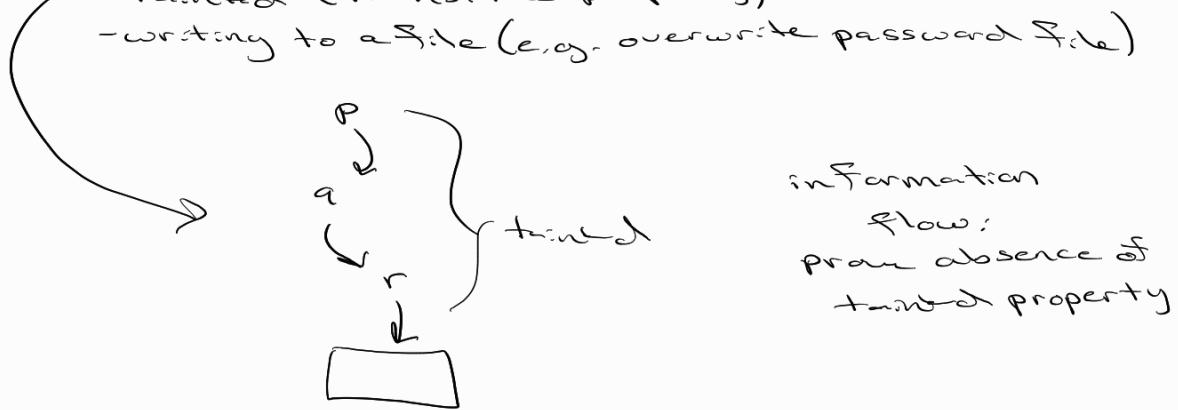
some other static analysis
tools focus on security
over correctness:

- exploitable bug $\Rightarrow \$\$$

- security-based analysis

- classic bugs:

- use after free
- tainted (transitive property)
- writing to a file (e.g. overwrite password file)



Dynamic Analysis:

- Valgrind

- shadow memory: version of the mem (identical or compact representation)

- uninitialized values → 0 or 1
- malloc + free mem → M or F

- precision super high

- if no error found, then particular execution path doesn't have error
 - generally do not generalise (static can)
 - not much recall

	Dynamic	Static
P	✓	✗
R	✗	✓

- valgrind uses memcheck, cache grnd, helgrind (race detection)

- Google has ASan, UBSan, + TSan

Rather than ASan (overhead), typically use logging + telemetry (Watson)

↳ usually get stacktraces (not helpful)
wish we had MRE (minimized reproducible example?)
PII

GORR (rules for managing + collecting PII)

↳ makes complicated to find bugs

So nothing's finding all bugs...

Do better testing

If ppl don't write specs, will they write tests?

- unit tests
- assertions
- regression tests

} add to continuous integration

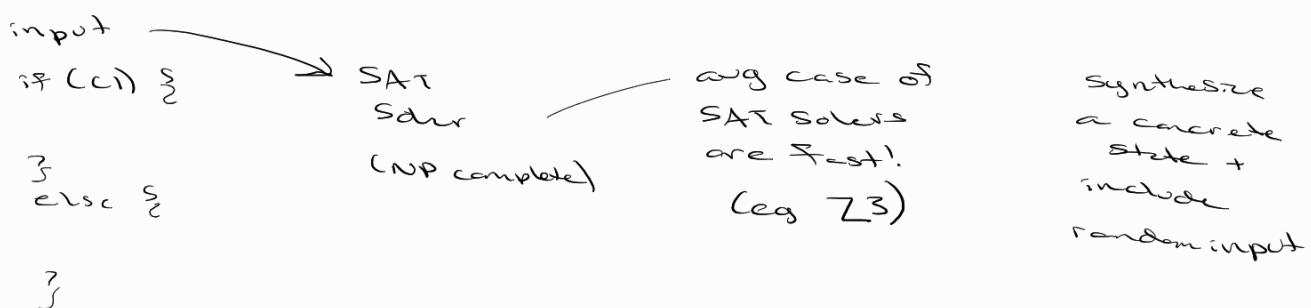
As # of tests grows, cost grows (time b/w commits grows)

Also more complicated tests!

systems != dev env

can't anticipate all

- integration tests
 - end to end tests
- Ppl hate writing tests, and written tests aren't great (corner case?)
- developers v. testers (moral hazard)
 - ↓
 - if testers are true, devs will write worst code (this actually was true?)
 - coverage: need tests to increase coverage
 - line
 - branch
 - full path
 - flaky tests: heisenbug in tests (says test failed)
 - fuzzing:
 - discovered by accident (line noise ex)
 - all Unix tools crashed on random inputs
 - no one anticipated randomness!
 - easy to generate randomness, but not random structured inputs
 - we now have great fuzzers (ex. JSONs)
 - these have found tons of bugs
 - 2 ways to work:
 - through static analysis
 - DART: concrete + symbolic analysis



AFL is another fuzzer

- American Fuzzy Lop
- one of the most effective fuzzers
- mutates inputs to see if any further cause
- uses genetic algorithms (really a giant hash but works great!)
- parallelizable, very quick
- very successful @ finding security bugs

Thursday Exam!

Last wk is for final project (Juan taking attendance)

