



## Going Beyond Metadata

Why We need To Think of Adapting Static Analysis in  
Dependency Tools

Joseph Hejderup

A11102 808852

NATL INST OF STANDARDS & TECH. N.I.C.



A11102808852

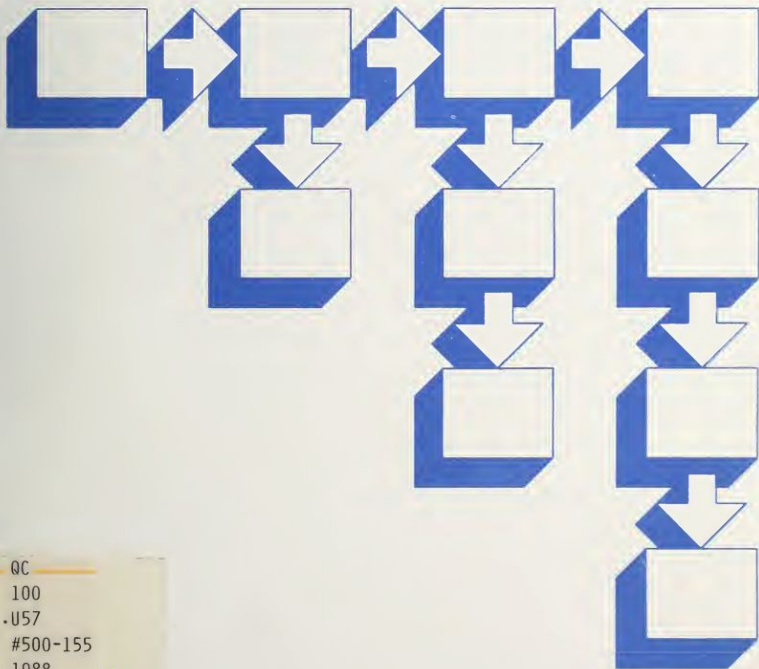
Wong, William/Management guide to software reuse

QC100 .U57 NO.500-155 1986 V19 C.1 NBS-P

NBS Special Publication 500-155

## Management Guide to Software Reuse

William Wong



# Back to 1980

Let's first  
understand the ideas  
behind Software  
Reuse

Wong, William. *A management overview of software reuse*. US Department of Commerce, National Bureau of Standards, 1986.

QC

100

.U57

#500-155

1986

C.2





# // Software Reuse: Ideas

## 5.1 Productivity

Reusing well-designed, well-developed, and well-documented software improves productivity and reduces software development time, costs, and risks.

## 5.2 Quality

Improvements in the quality of software developed from well-designed, well-tested, and well-documented reusable software components



# // Software Reuse: Implementation

Gateway to  
thousands of  
libraries and  
frameworks

```
jhejderup — -bash — 68x20
[Josephs-MBP:~ jhejderup$ npm help

Usage: npm <command>

where <command> is one of:
  access, adduser, bin, bugs, c, cache, completion, config,
  ddp, dedupe, deprecate, dist-tag, docs, doctor, edit,
  explore, get, help, help-search, i, init, install,
  install-test, it, link, list, ln, login, logout, ls,
  outdated, owner, pack, ping, prefix, prune, publish, rb,
  rebuild, repo, restart, root, run, run-script, s, se,
  search, set, shrinkwrap, star, stars, start, stop, t, team,
  test, tst, un, uninstall, unpublish, unstar, up, update, v,
  version, view, whoami

npm <cmd> -h      quick help on <cmd>
npm -l           display full usage info
npm help <term>  search for help on <term>
npm help npm     involved overview
```

Dependency  
management with  
automated conflict  
resolution

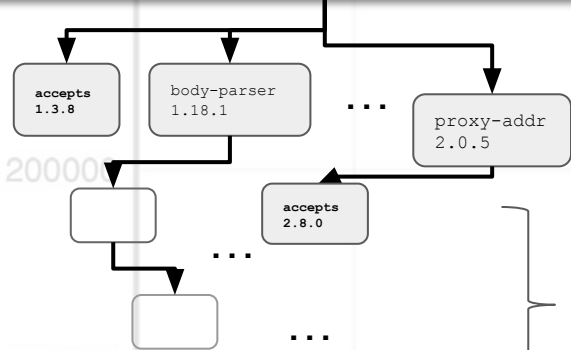
Centralized distribution

Package Managers

# // Dependency Management

## Problem #1: Local Level

```
Compiling async-executor v1.4.1
Compiling rand_chacha v0.2.2
Building [=====>] 194/307
```



Transitive  
deps

```
[package]
name = "webframework"
version = "4.16.3"
```

[dependencies]

accepts = "~1.3.5"

body-parser = "1.18.2"

depd = "~1.1.2"

encodeurl = "1.0.2"

escape-html = "~1.0.3"

etag = "1.8.1"

proxy-addr = "~2.0.3"

TODAY

...3 DAYS LATER

1.3.7

1.3.12

1.1.3

1.1.3

1.0.6

1.0.9

2.0.5

2.0.5

Complex & Large Package Compositions

Temporal Properties



# // Dependency Management

Problem #2: Global/Repository Level

**QZ** Quartz

## How one programmer broke the internet by deleting a tiny piece of code

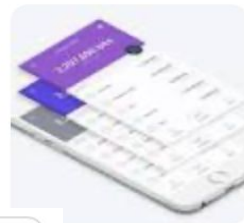
Lots of npm packages  
how this tiny bit of code  
Mar 27, 2016

```
module.exports = lattpad;
function lattpad(str, len, old) {
  str = String(str);
  var i = 0;
  while (i < str.length) {
    len = len - str.length;
    str = str + len;
  }
}
```

**ZD** ZDNet

## Hacker backdoors popular JavaScript library to steal Bitcoin funds

The library loading the malicious code is named Event-Stream, a JavaScript



**W** WIRED

## A Second SolarWinds Hack Deepens Third-Party Software Fears

It's been more than two months since revelations that alleged Russia-backed hackers broke into the IT management firm SolarWinds and used ...  
5 days ago





# // Dependency Management

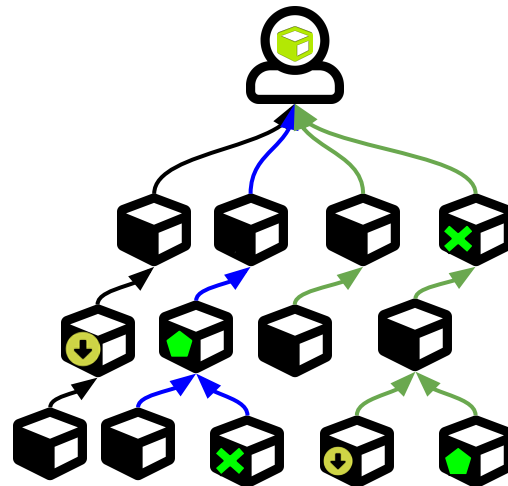
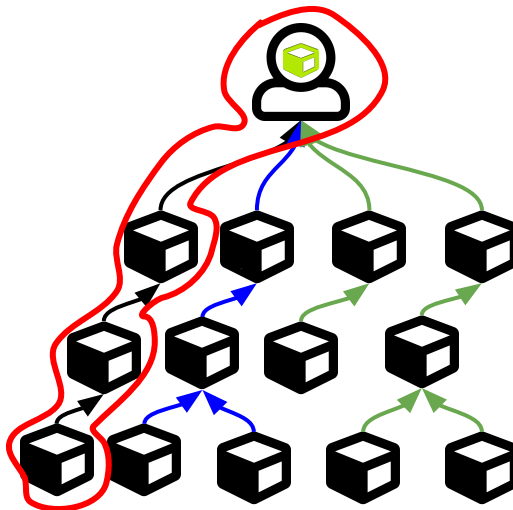
Temporal Properties => Version Pinning/Lock files

Everything Else => Dependency Analyzers/Bots/Plugins

- \* alert fatigue
- \* **actionability?**
- \* **precision?**



- \* vulnerabilities
- \* updates
- \* audit
- \* quality
- \* deprecation
- \* ...more to come!





# // Revisit

## 5.1 Productivity

Reusing well-designed, well-developed, and well-documented software improves productivity and reduces software development time, costs, and risks.







# // Classic Alert Fatigue? Nope

We need first address the quality of analyzers!

```
[package]
name = "webframework"
version = "4.16.3"
```

```
[dependencies]
accepts = "~1.3.5"
body-parser = "1.18.2"
depd = "~1.1.2"
encodeurl = "1.0.2"
escape-html = "~1.0.3"
etag = "1.8.1"
proxy-addr = "~2.0.3"
```

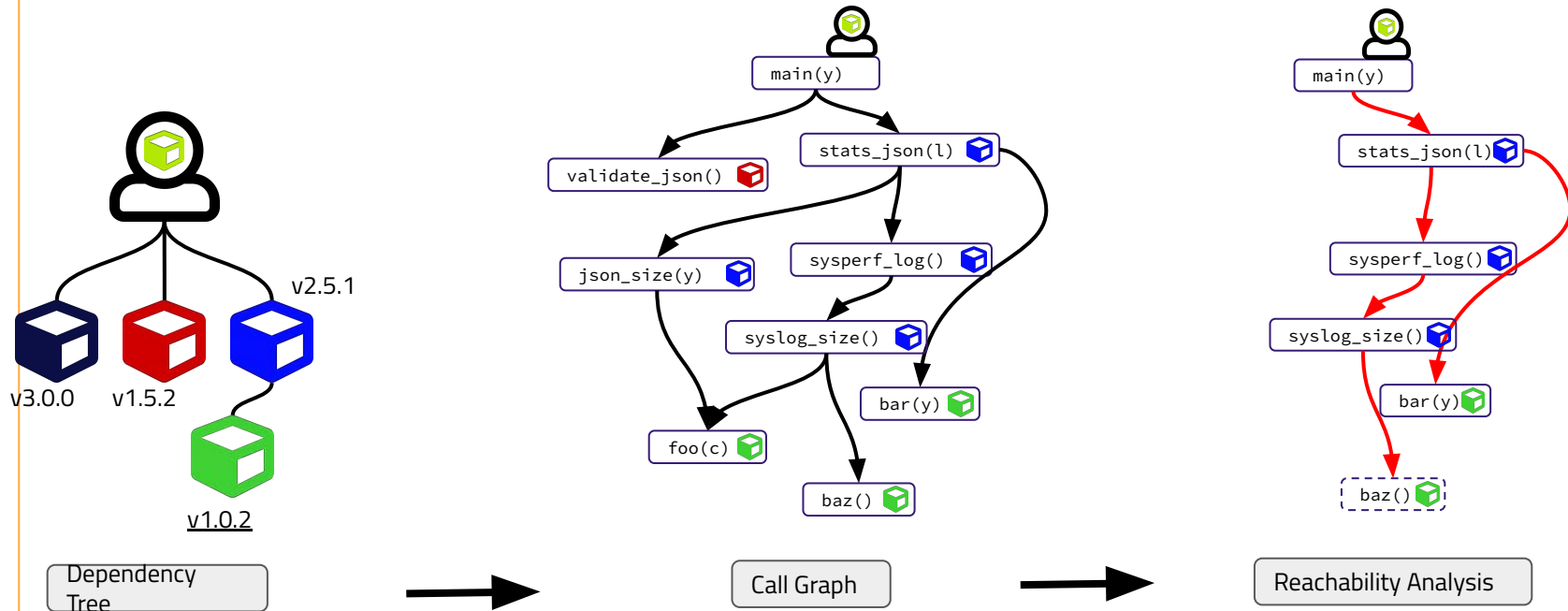
Rule #1: Metadata is not source code!

Declaration  $\neq$  usage



# // Classic Alert Fatigue? Nope

We need first address the quality of analyzers!



Rule #2: Make Code first-class citizens



# // Program Analysis

great, got it but expensive and not scalable, dependencies are many, right?



ktrianta/rust-callgraphs

35,896 Packages

208,023 Releases



Repository index

23,767 Packages

142,301 Releases



call graphs

10



days

67%

(actually 80%)

Rule #3: Aim for Light-Weight Analysis Techniques!



# // Program Analysis

- \* “Overkill for me to add Program Analysis”
- \* “What about Python/JavaScript?”
- \* “Program Analysis suffers from False Negatives, my security customers won’t be happy about it”





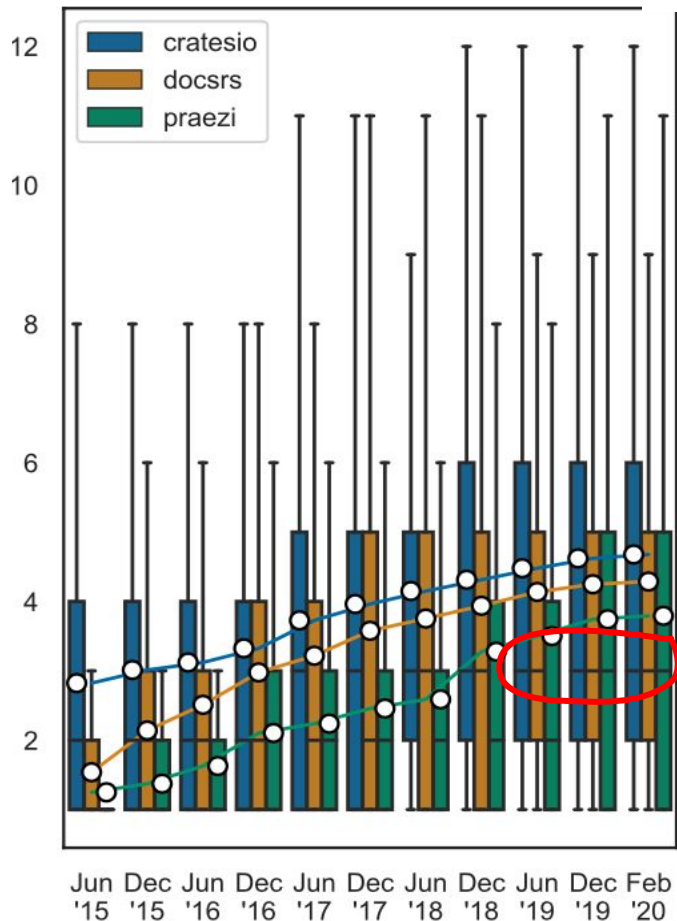
# // Program Analysis

Let's settle the question through some research!

RQ: What is the difference in the number of reported dependencies between traditional metadata-based approaches vs program analysis approaches?



# // Number of Direct Dependencies

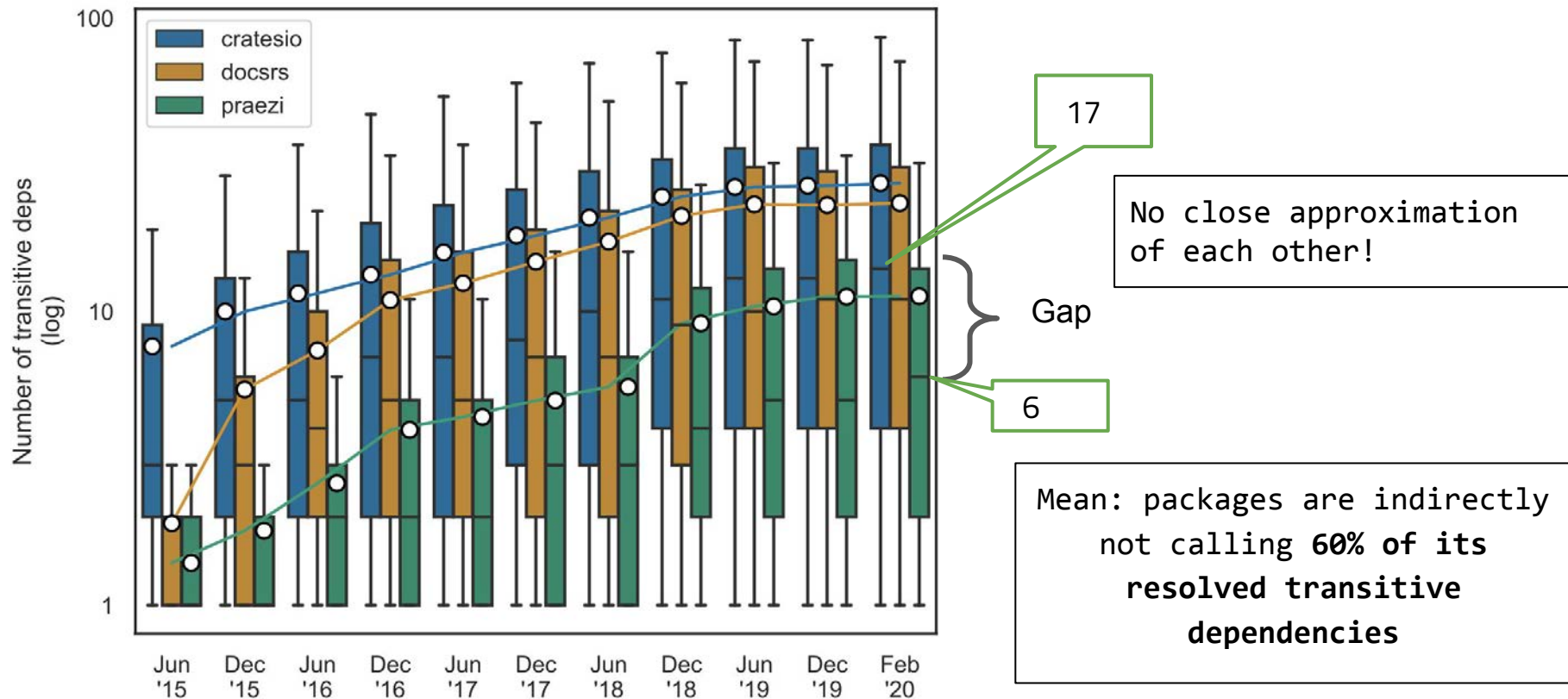


The median is similar; the mean relatively when comparing the three representations.

A metadata-based direct deps slightly over-approx. a static analysis inferred direct deps!



# // Number of Transitive Dependencies





// Why is there such a large difference for transitive depz?

diff metadata-based call-based dependency tree  
in 34 randomly selected cases





// Why is there such a large difference for transitive depz?

Diff #1

- \* 3 x No import statements
- \* 4 x Import statements but no usage

Not totally unexpected but also not that common!



// Why is there such a large difference for transitive depz?

Diff #2

- \* 1 x conditional compilation (e.g., `[#cfg(...)]`)
- \* 2 x derive macro libraries (code generation)
- \* 1 x test dependency (in wrong section!)

Not all dependencies are runtime libraries!

Optional dependencies surface and enabled in code; how do we handle them?



// Why is there such a large difference for transitive depz?

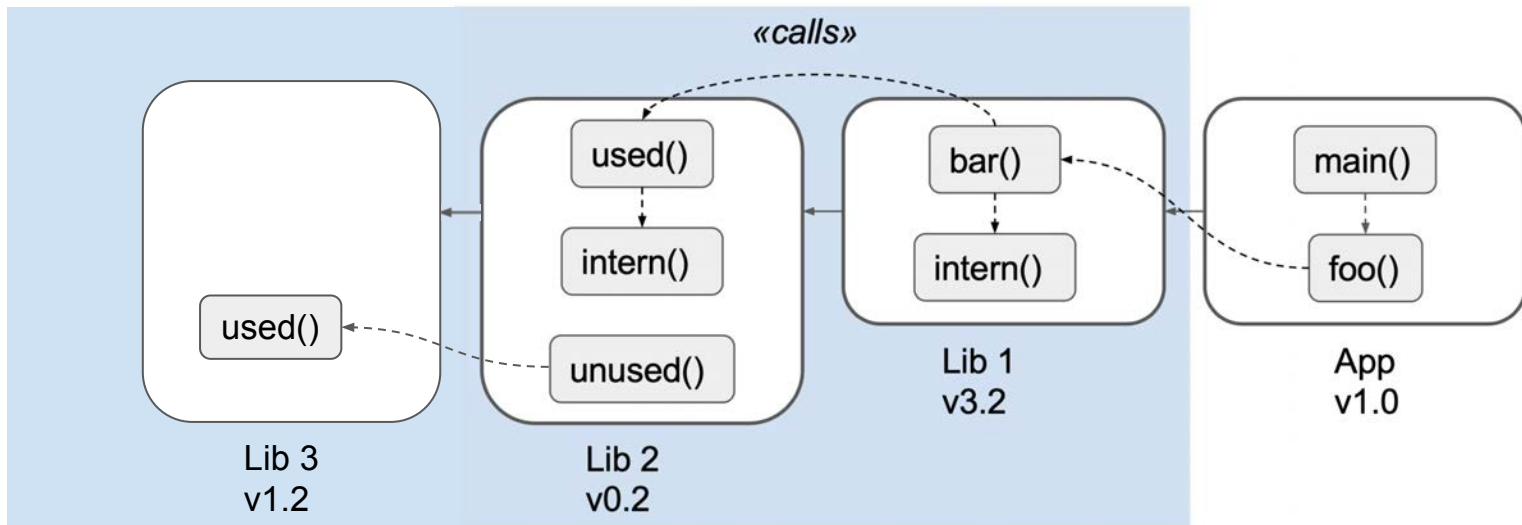
Diff #3

\* 16 x non-reachable transitive dependencies!

What do you mean?

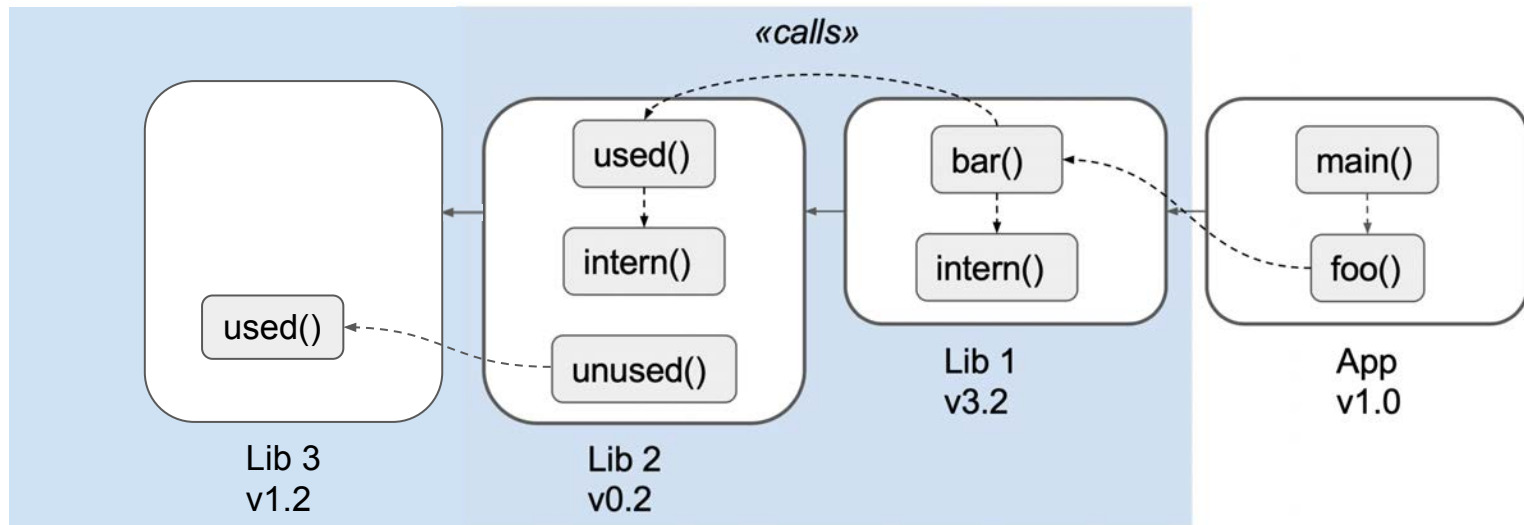


// How many dependencies is App v1.0 using?





# // Context Matters!



We make the general assumption that we use

- **ALL APIs** of all direct dependencies,
- and **then ALL APIs** of transitive dependencies



# // Practical Implications - Trade-off

Analysis Scope - Direct Dependencies

- \* Declared dependency closely estimates a utilized dependency.
- \* Metadata > Static Analysis

Pros:

- No need to implement program analysis
- Higher recall for security/soundness-sensitive applications

Cons:

- Insensitive to simple things like no import statement & specifics of APIs being utilized
- Cannot eliminate dependencies that are solely doing code-generation (i.e, no actual runtime dependency)



# // Practical Implications - Trade-off

Analysis Scope - Transitive Dependencies

\* Static Analysis > Metadata

Pros:

- Capture the actual usage, reducing false positives.
- Higher actionability

Cons:

- False Negatives => May not be ideal for security applications.
- Highly dynamic libraries are challenging to analyze



# // Program Analysis: Precision & Recall

- \* Recall

- Coverage of language features
- Lost Precision, still better than metadata

- \* Precision

- Algorithms like Points-to algorithm not that scalable

- \* General Implications

- Scope of analysis: Project + its dependency tree
- Package repositories are not homogeneous collection of libraries





**Joseph Hejderup**  
joseph@endor.ai