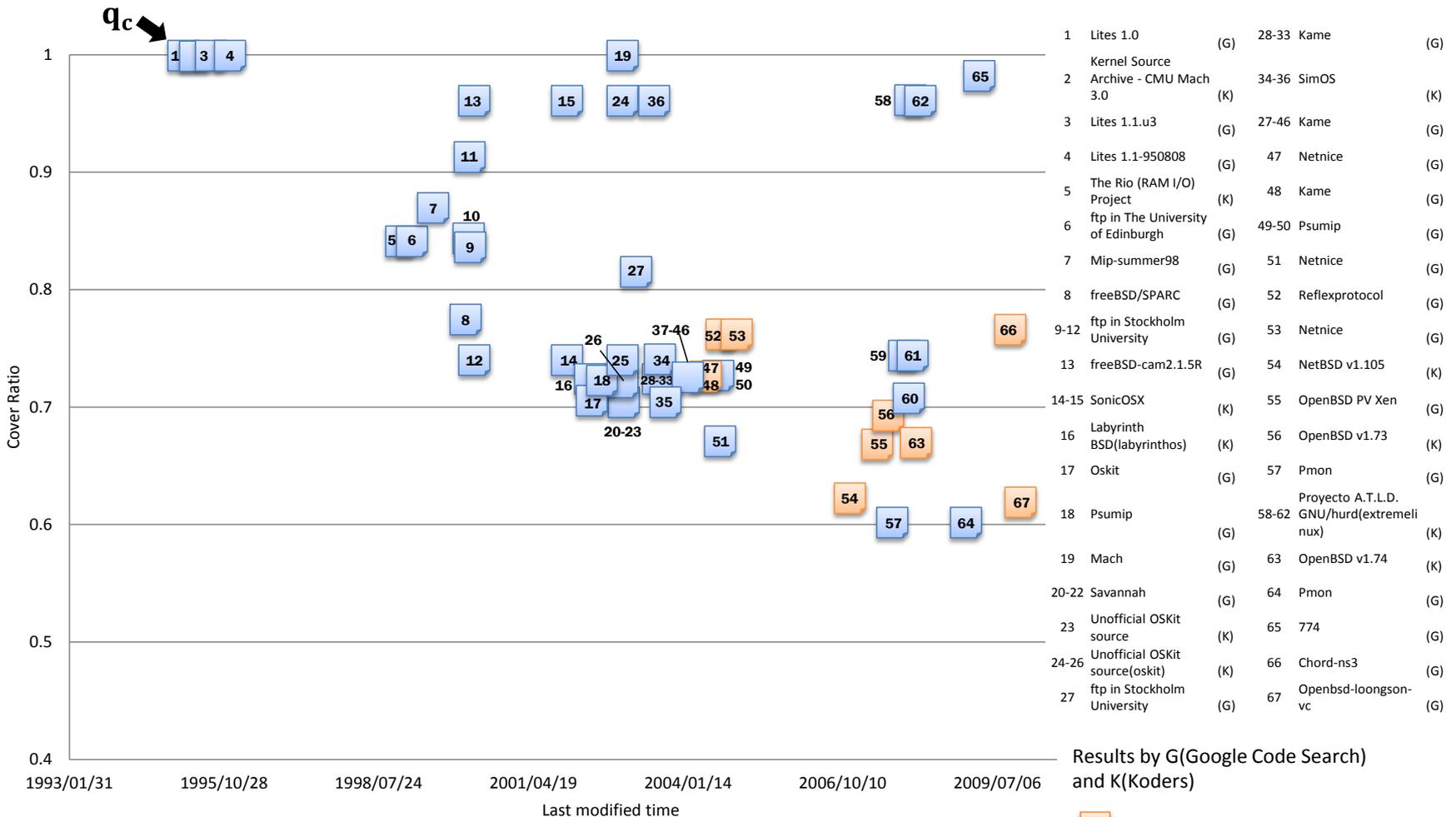


Challenges in Mining Whole Software Universe

Katsuro Inoue
Osaka University



Analyzing Evolution of kern_malloc



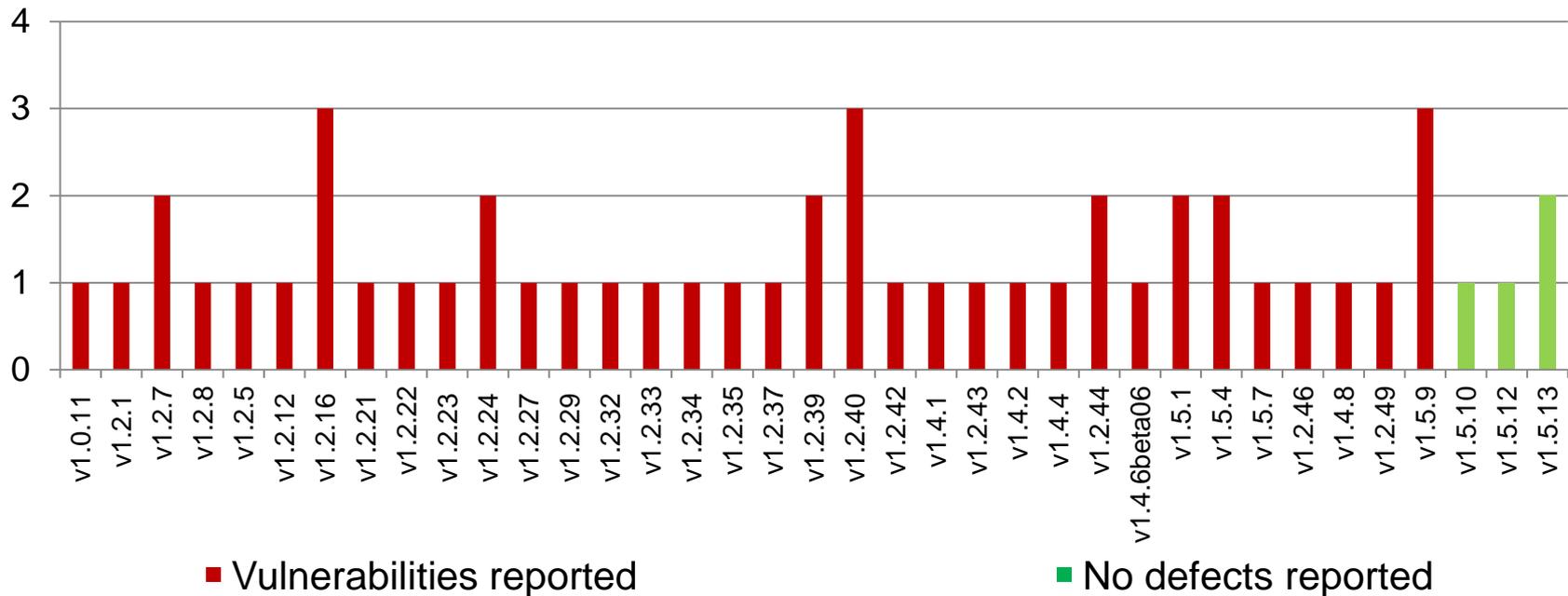
Results by G(Google Code Search) and K(Koders)

- : File in New BSD License
- : File in Original BSD License



Analyzing Reuse of Outdated Libraries

Vulnerability of 50 OSS Projects Using `libpng`



Result from Google Code Search and Koders



Experience and Concern

Mining source code repositories, e.g.,
SourceForge, Github, Open Hub, Google Code, Marven, ...
(BlackDuck)



- Outcomes heavily depend on repository contents
- Aren't we mining a small world?
- There may be many other source code contents in the universe







SOURCEforge

fedora™



codeBeamer
JavaForge



Fusion Forge

GITGO

CloudForge
by CollabNet



berliOS
Entdecke Open Source



GitLab.com



our
project.org



SEUL.org

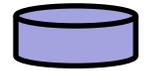


sourceforge

CodePlex



Google
Developers



Alioth.debian.org

OW2



Bitbucket



assembla

freepository.

GITORIOUS

Tigris.org
Open Source Software Engineering Tools

launchpad

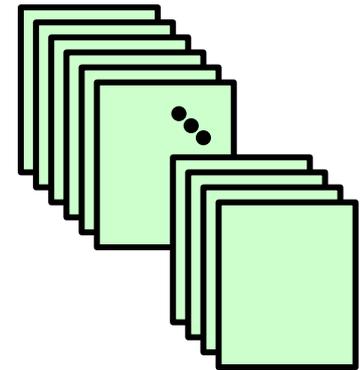
Whole Software Universe U

- Whole Software Universe

$U \equiv$ *Collection of All Software*

Developed by Human in the Past

- Open source software
- Personally-developed software
- Proprietary software
- ... any others
- P : Set of all meaningful software
(a countable infinite set)
- $U \subseteq P$



Questions for U

A) How do we get U ?



B) What do we mine from U ?

C) How do we mine U ?

D) Why do we mine U ?



A) How Do We Get U ?

- No one knows actual U
- So we would collect many repositories, and construct a subset $U' \subseteq U$
- U' should be as large as possible, of course
- U' should reflect characteristics of U
- **Challenges**
 - Collecting and unifying different repositories into U'
 - Duplication, coherence, ...
 - Performance and capacity for U'
 - Updating and maintaining U'



B) What Do We Mine from U ?

Examples

- Simple metrics of U over history
 - Size $|U|_{t1}, |U|_{t2}, \dots$
 - Language usage
 - ...
- Density of U with respect to P
- History and evolution of code c in U
 - Origin version of c
 - Closely related code c' (clone, variation, family, ...)
 - Future prediction for c



C) How Do We Mine U (U')?

1. Direct mining

- Good model
- Powerful machine

2. Indirect mining

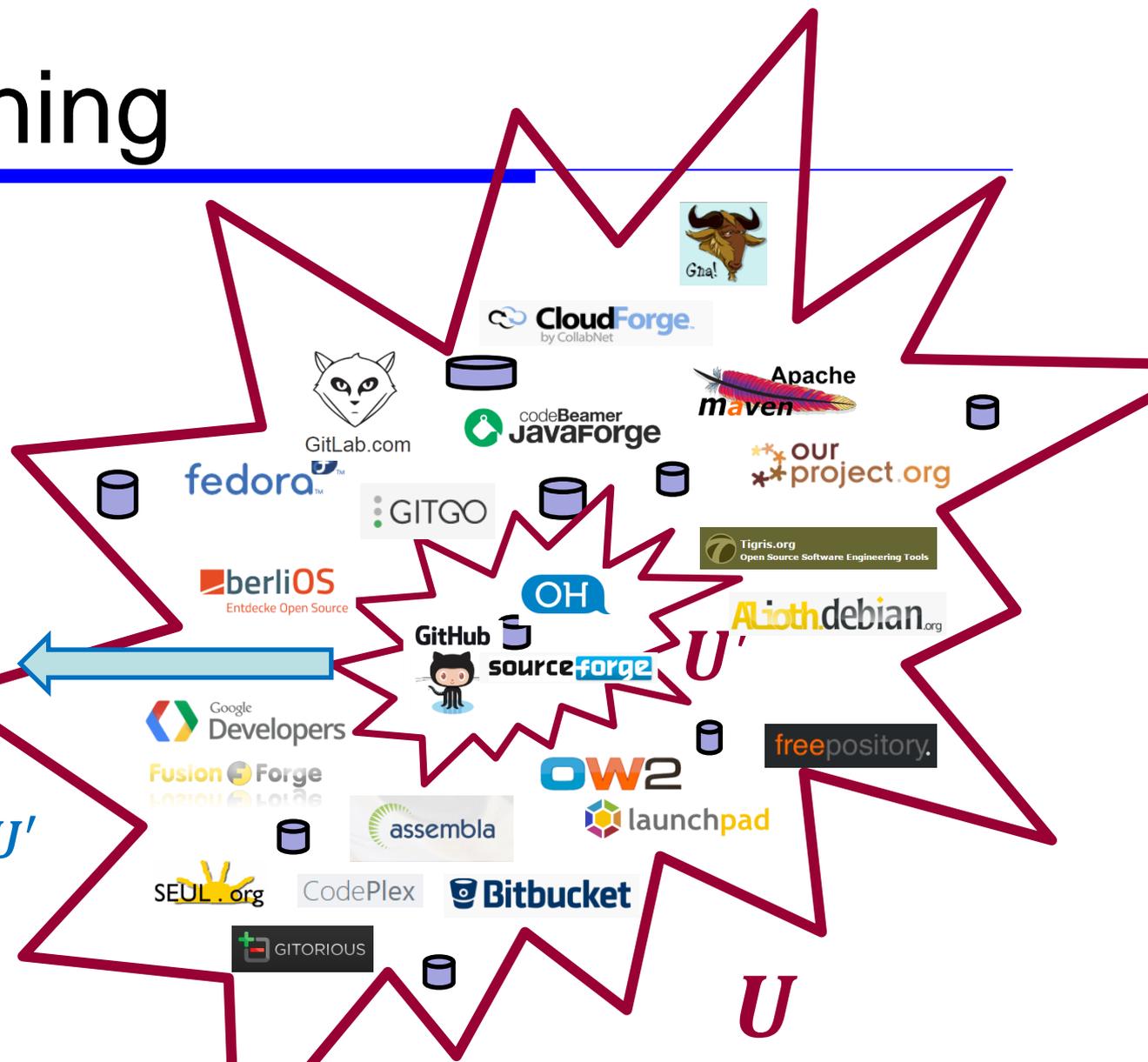
- Use external services
- Reconstruct mining result from those external services



Direct Mining



Copy of U'



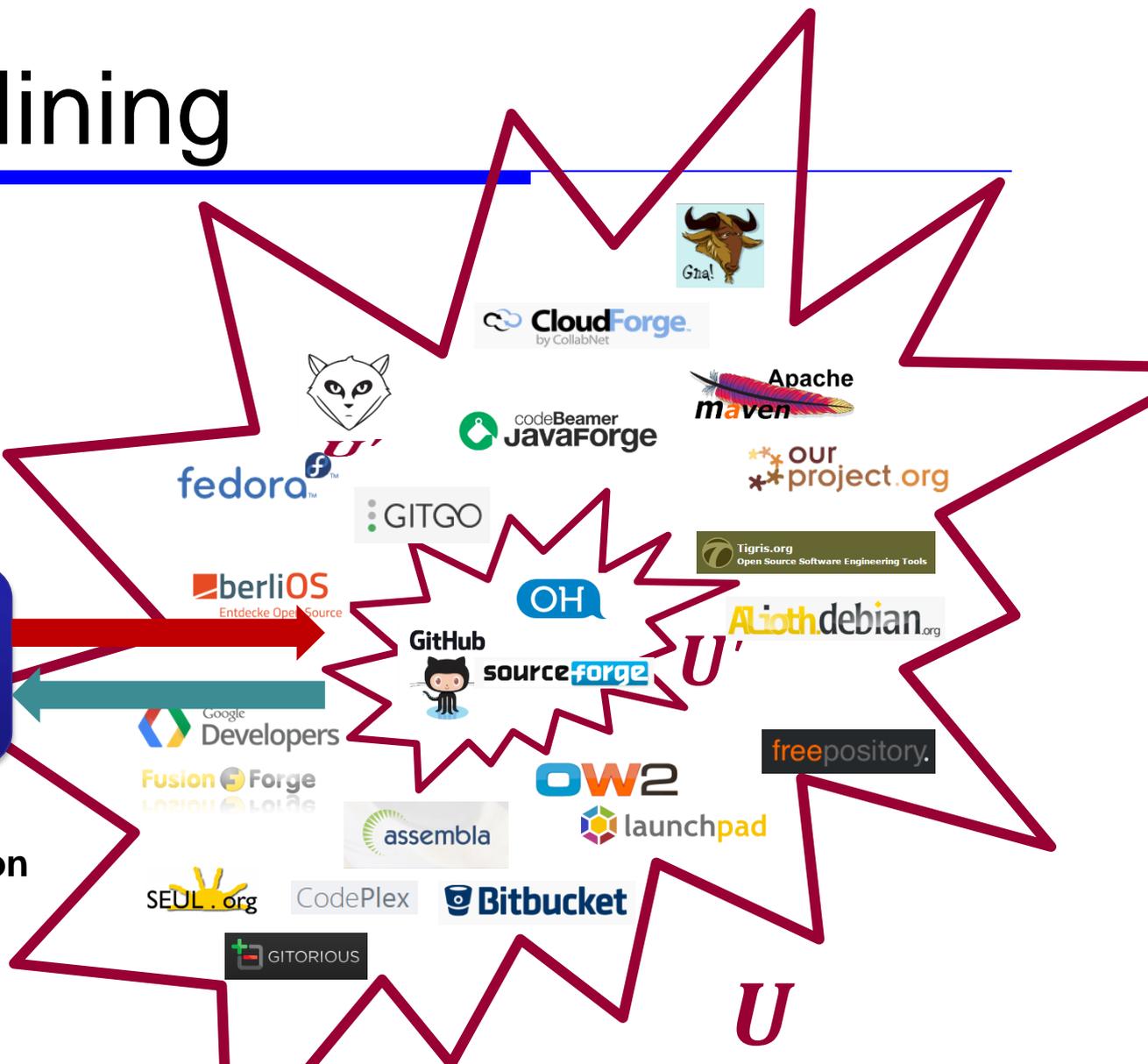
Indirect Mining

Want to know about U'



Mashup Engine

Query Decomposition and Result Composition



D) Why Do We Mine U ?

Objectives of mining U

- Reuse and knowledge transfer
 - We do not want to reinvent the *wheel*
- Historical Archive
 - Frontier's wisdom

...



Discussion!

- Is it interesting research topics?
- Can we get useful research results?
- Is it feasible research target?



Thank you



