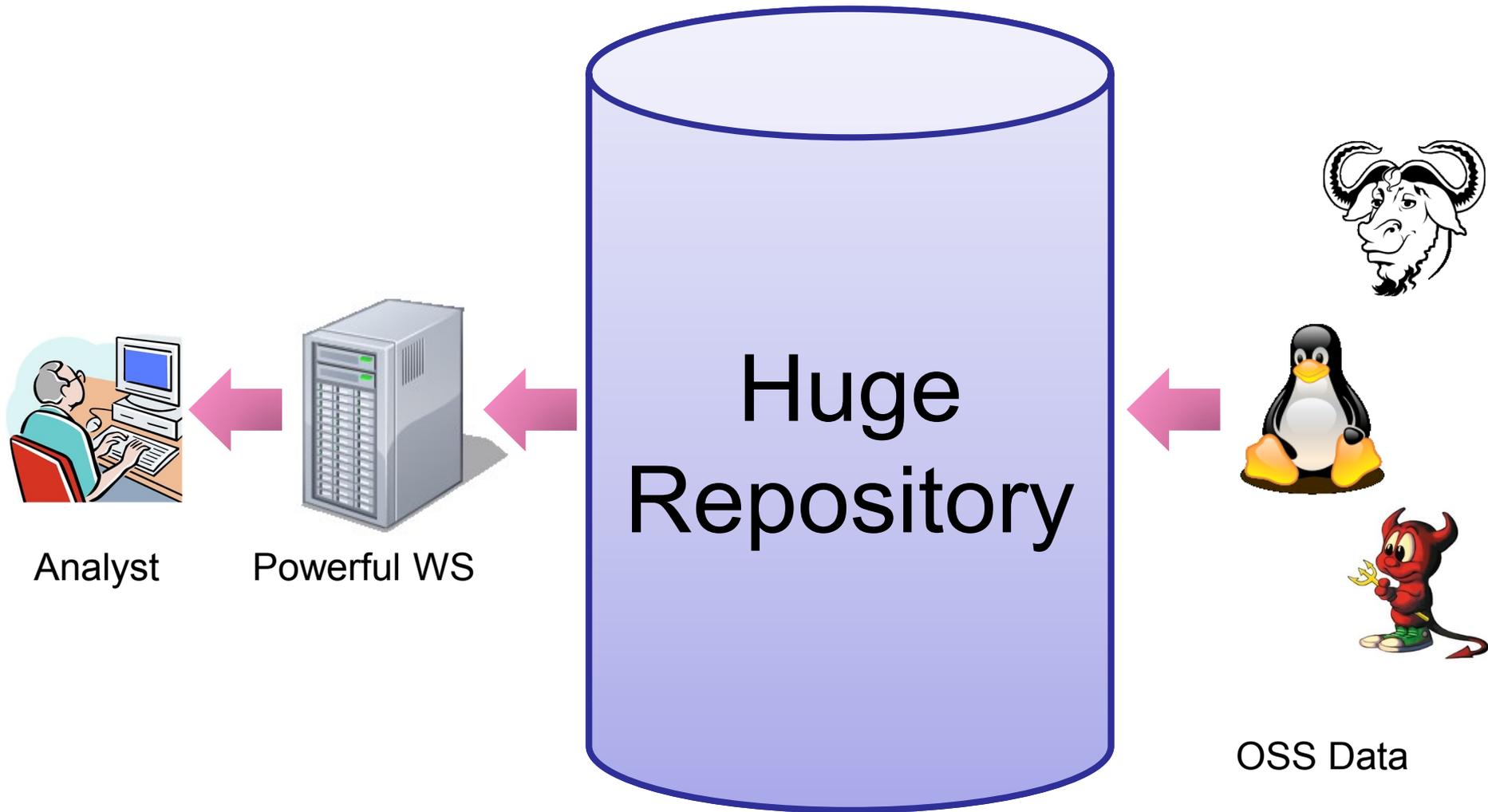


Mining Resources in the Internet Space

Katsuro Inoue
Osaka Univeristy



Mining OSS Data

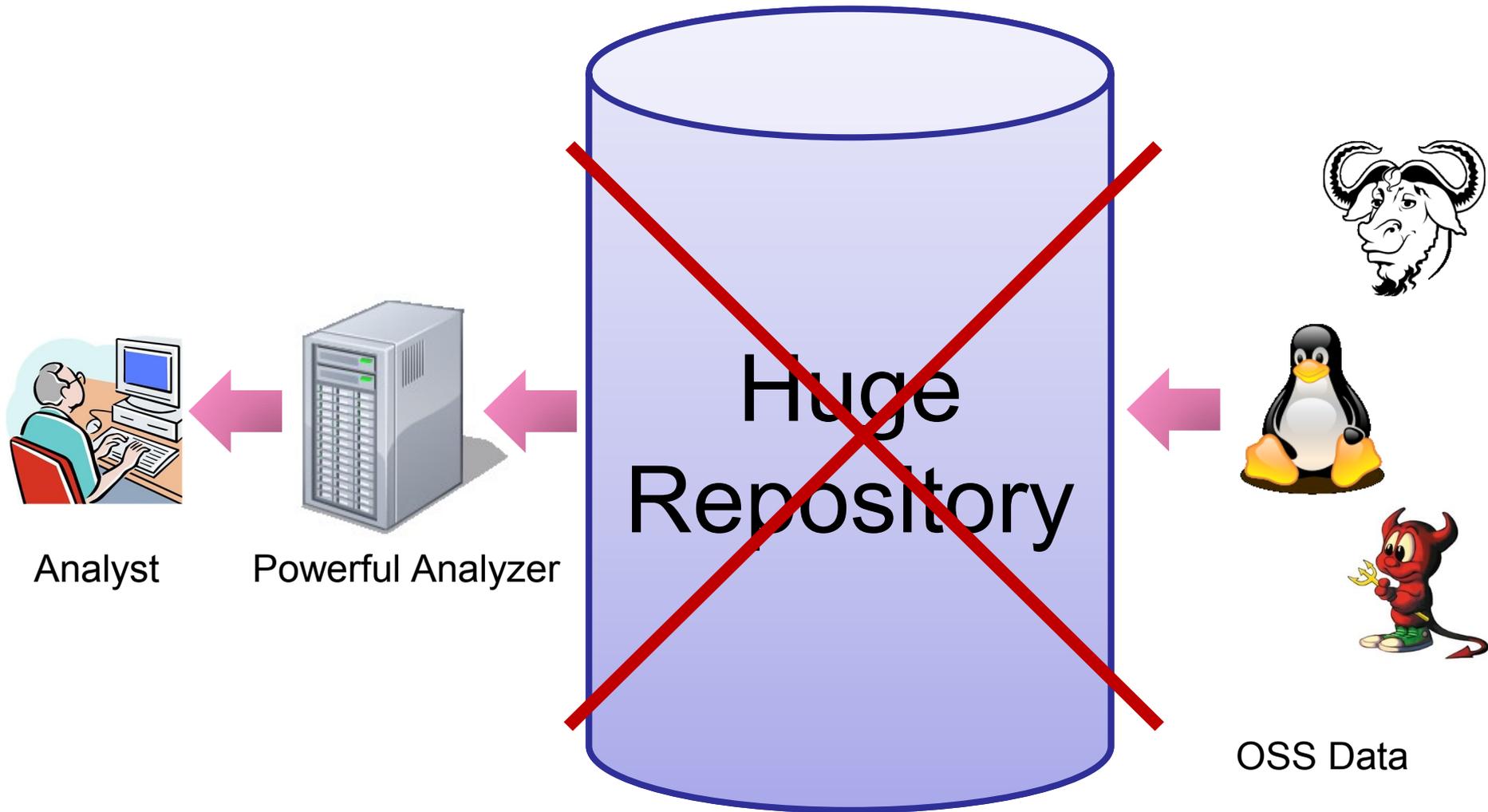


Maintaining Huge Repository

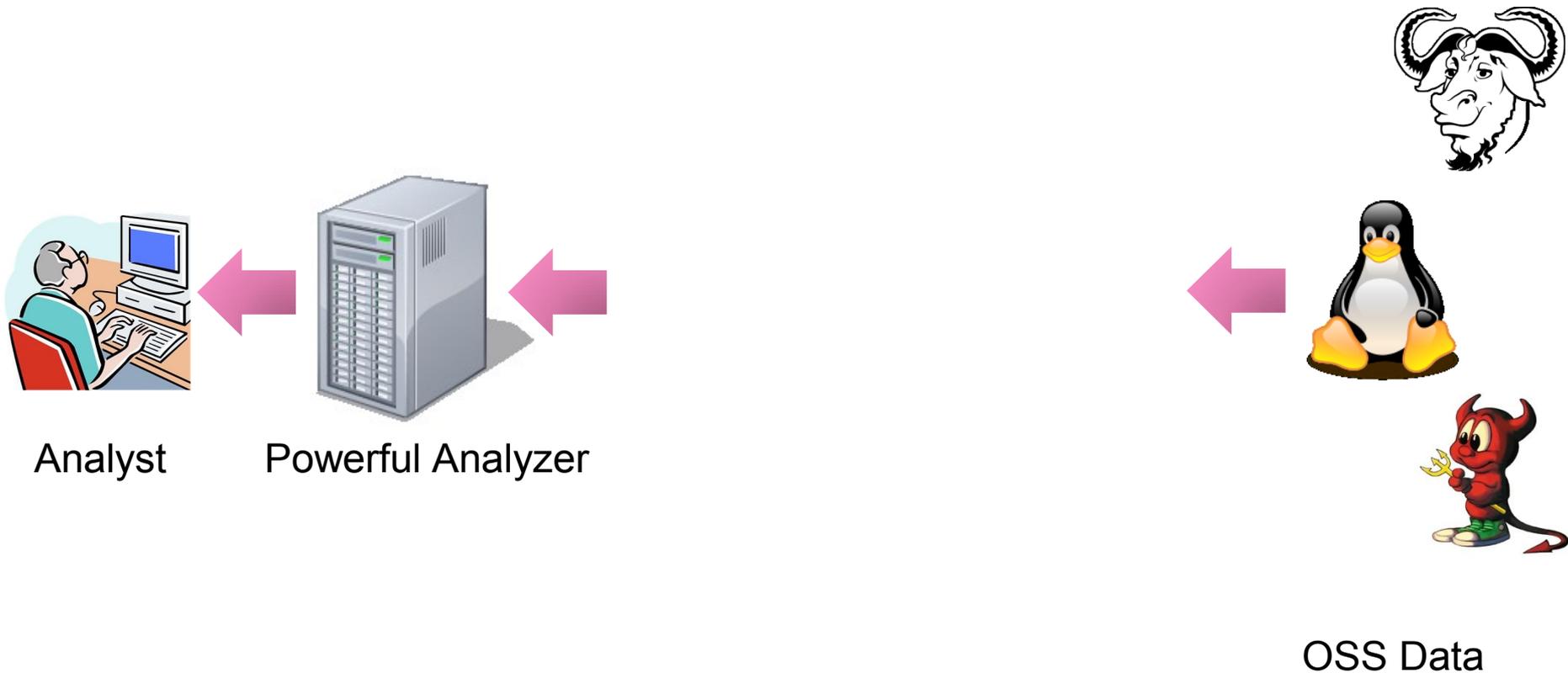
- Keeping the repository updated is important
 - Time consuming effort
 - Automatic crawling is not easy



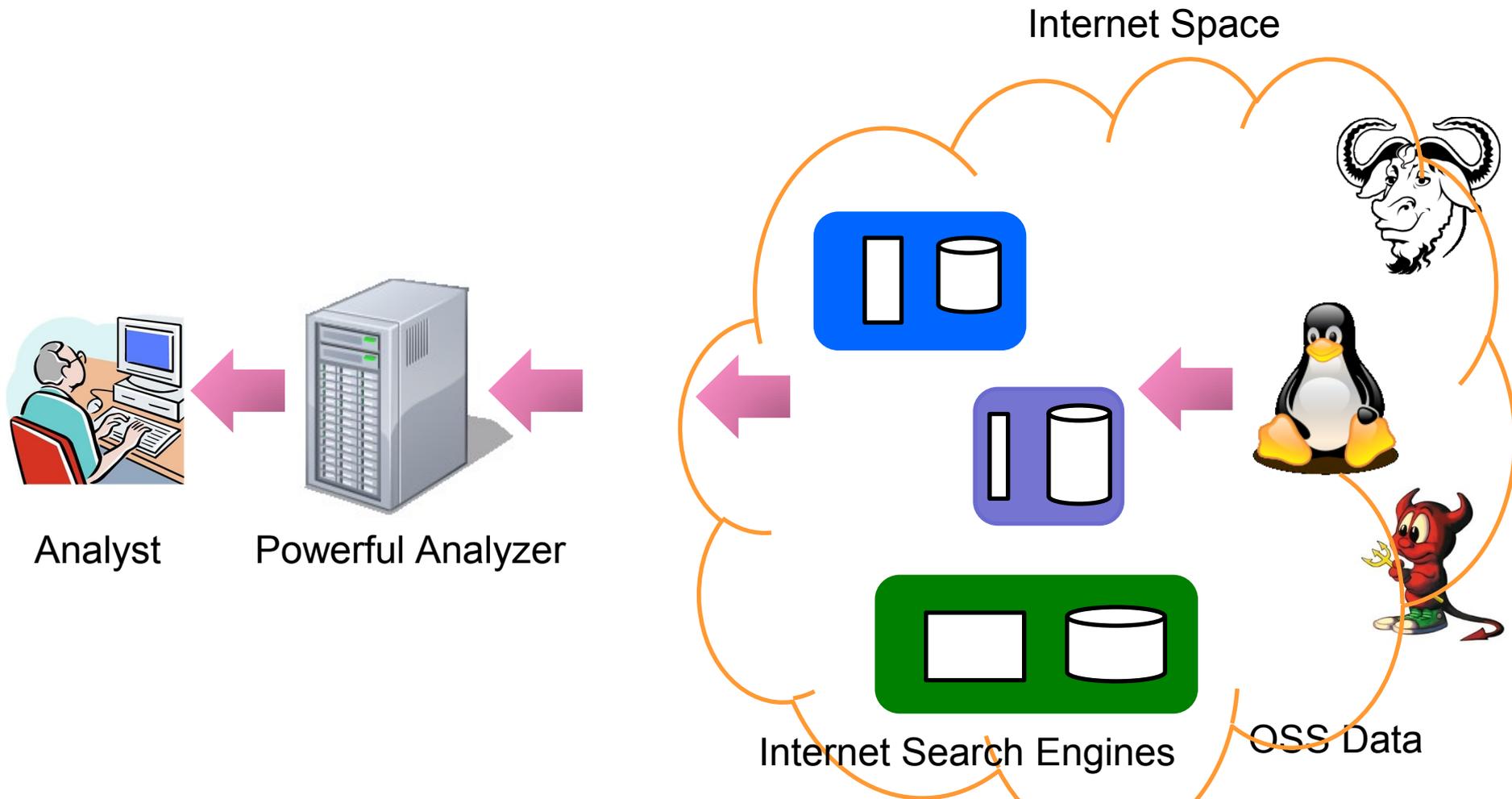
Mining OSS Data



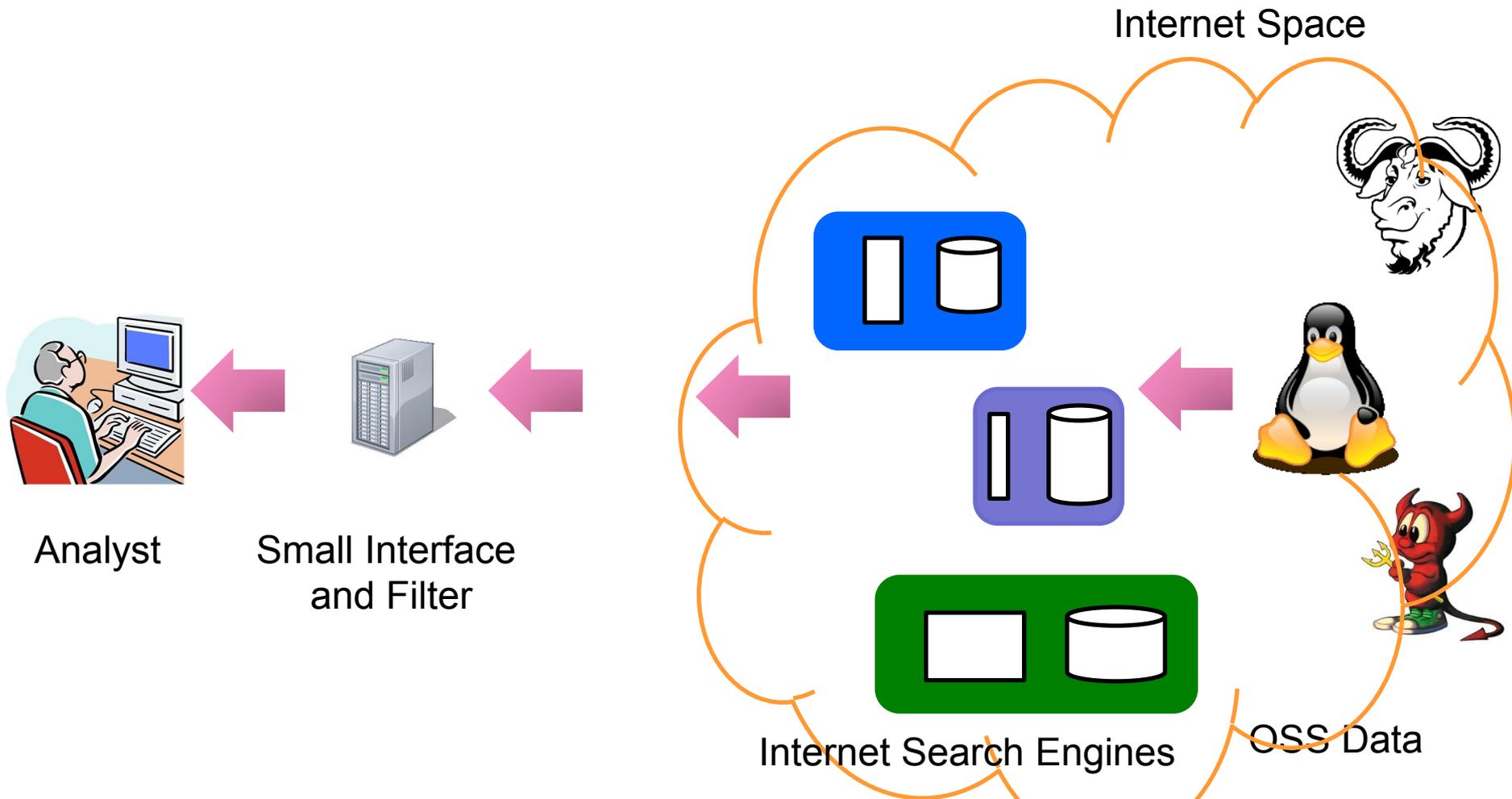
Mining OSS Data



Mining OSS Data



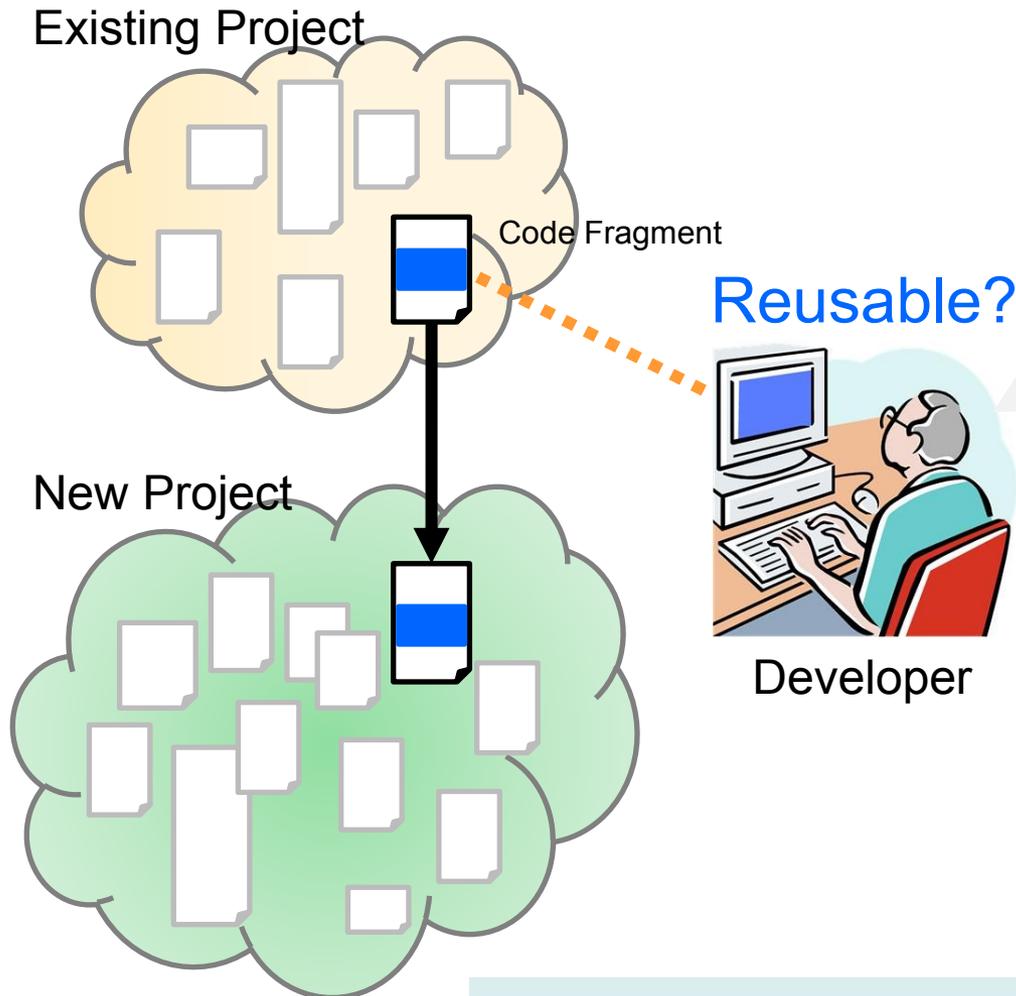
Mining OSS Data



Code History Tracking Model and Ichi Tracker



Developer's Concerns



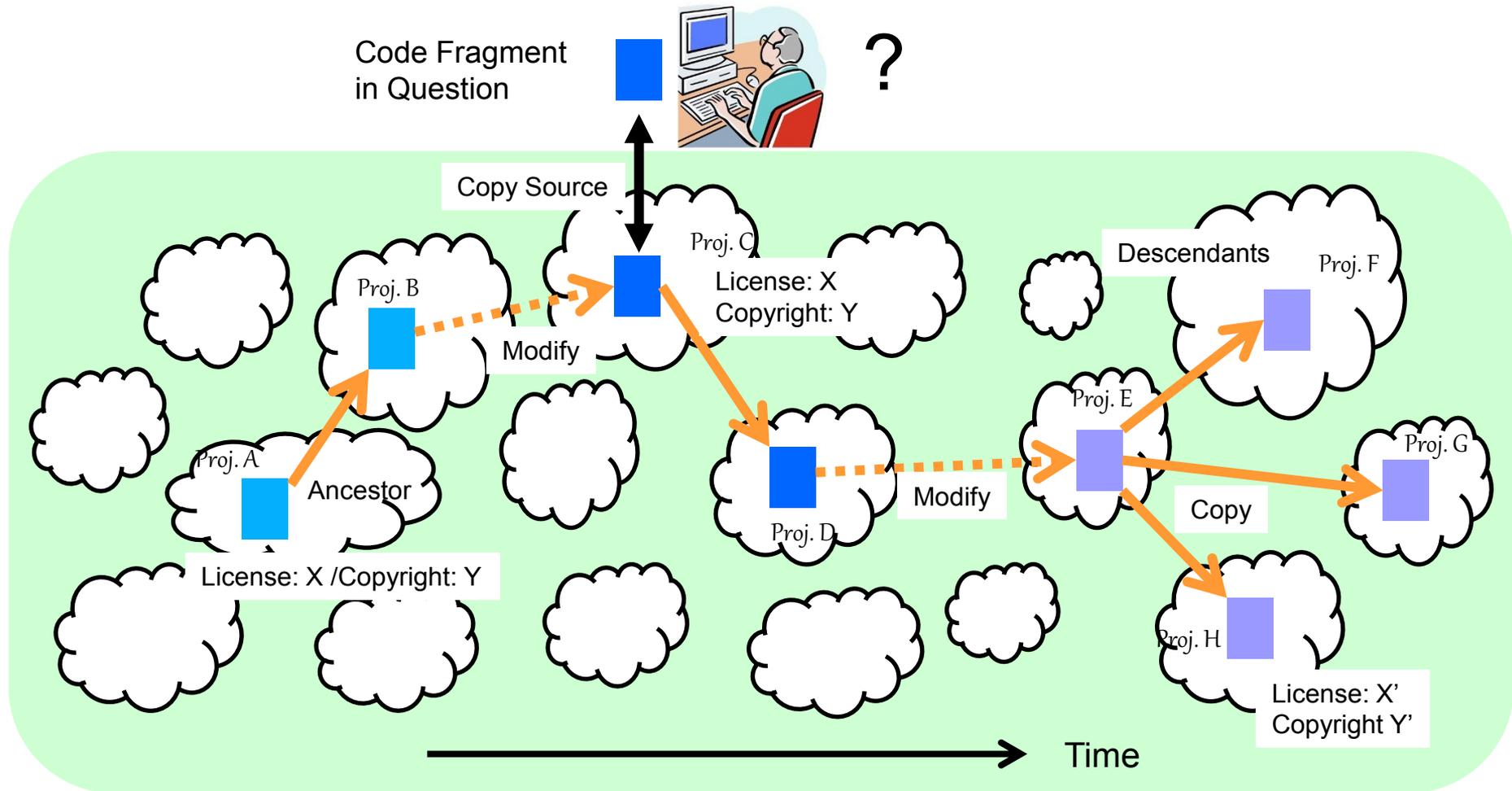
- Origin
 - Who?
 - When?
 - License?
 - Copyright?
- Evolution
 - Maintenance?
 - Popularity?
 - Newer version?
 - ...

Concerns

To ease concerns, a support system is needed



Code History Tracking System



OSS Repositories



Design Policy of Ichi Tracker

OSS repository

- Target many OSS projects from old to new ones
- No crawling, no maintenance
 - Do not have local repository, but use external code search engines



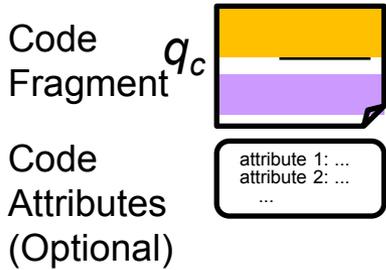
Output quality

- Find not only exactly same code fragments, but also similar ones
- Lower false positive results
- No real-time response
 - Use code clone filtering to improve the output quality

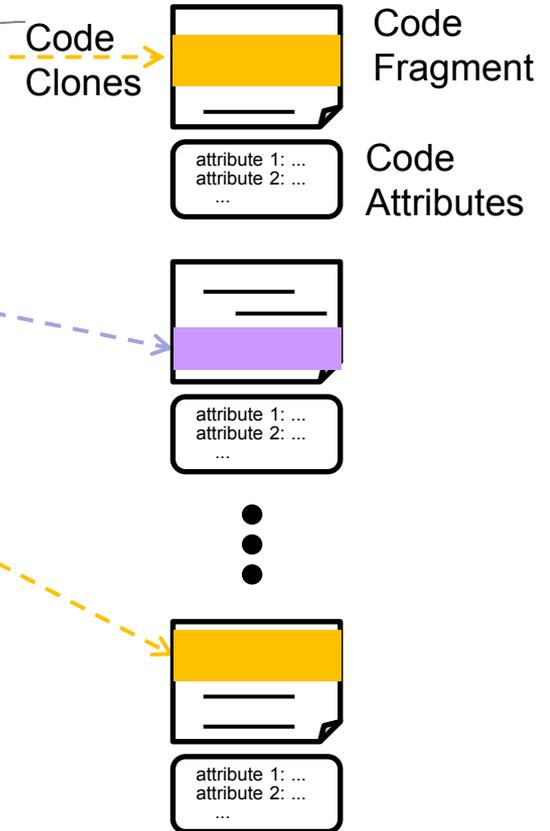


Code History Tracking Model

Input Query Q



Output Results R



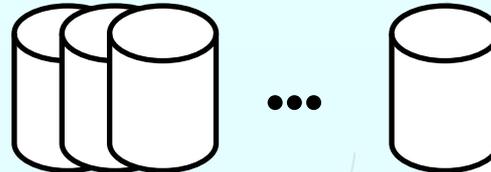
Search Query SQ

Search Results SR

Code Search Engines



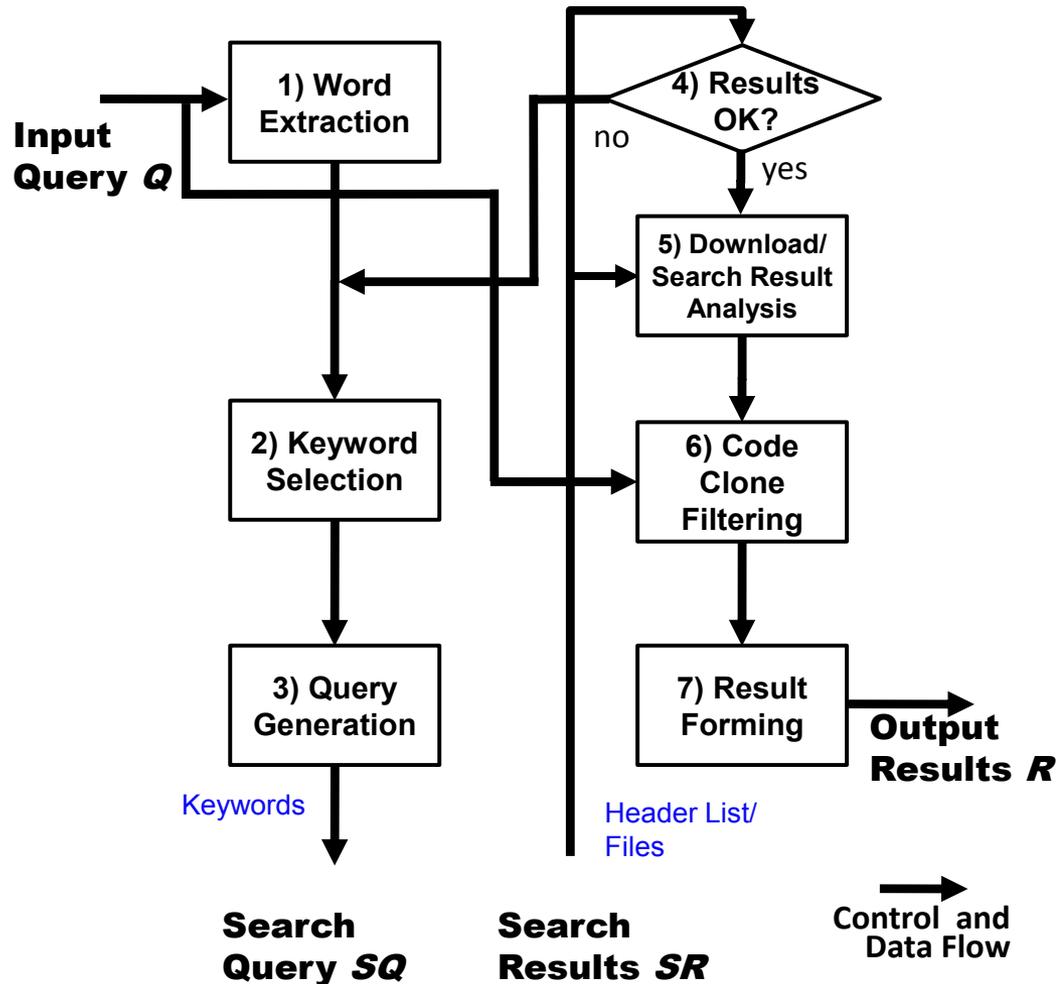
Internet



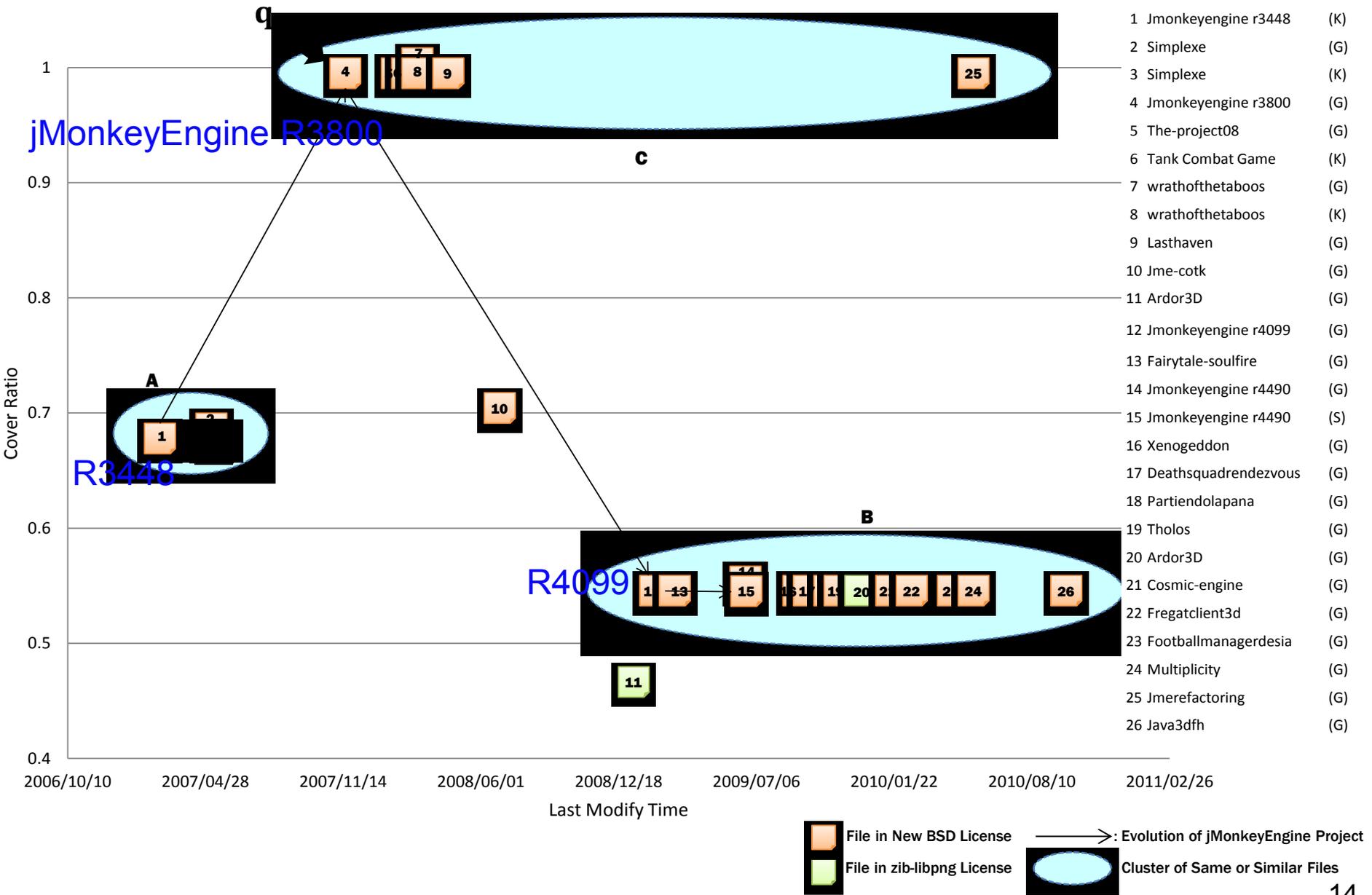
Open Source Repositories



Process of Ichi Tracker



Evolution Pattern of Texture.java



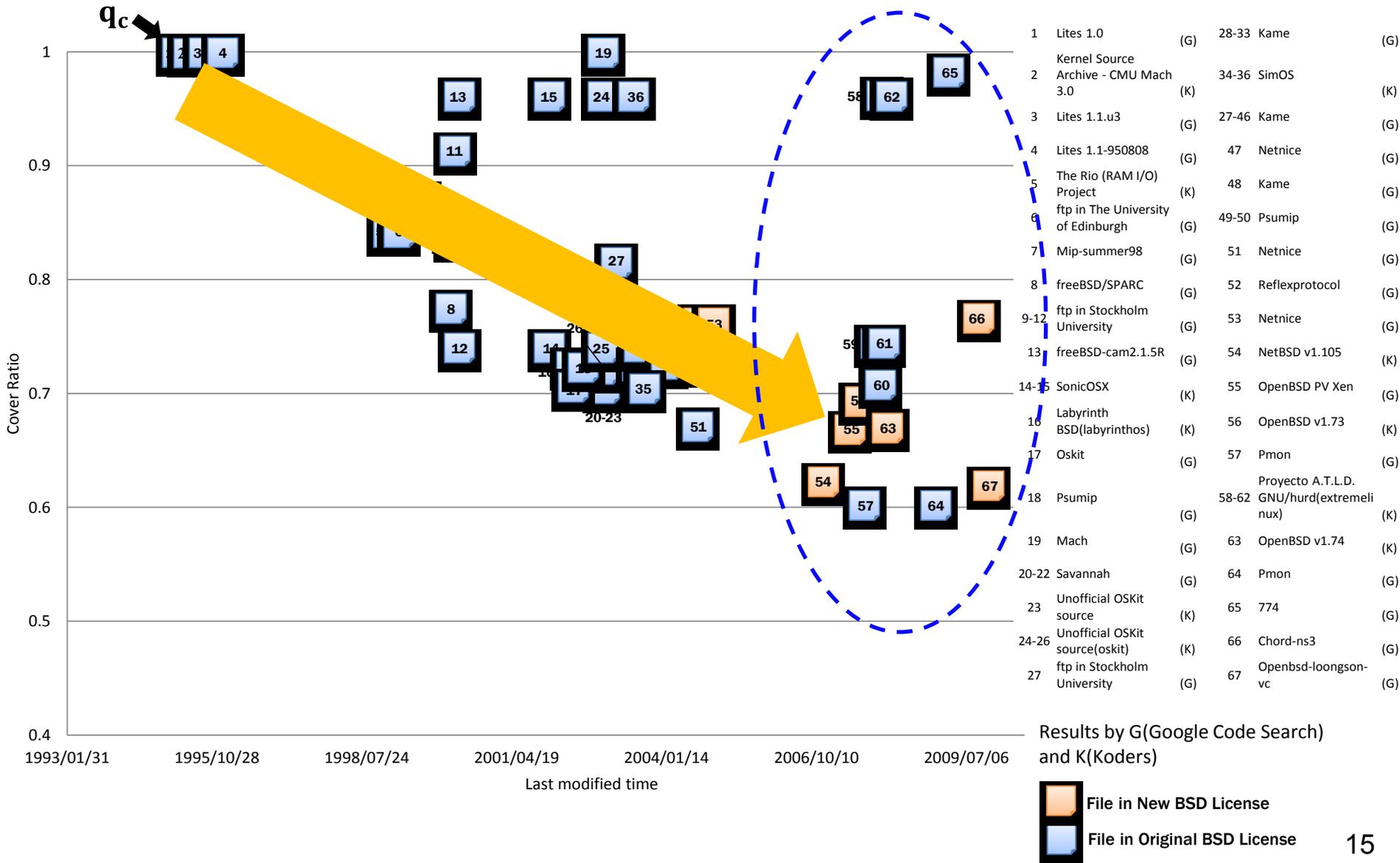
- 1 Jmonkeyengine r3448 (K)
- 2 Simplexe (G)
- 3 Simplexe (K)
- 4 Jmonkeyengine r3800 (G)
- 5 The-project08 (G)
- 6 Tank Combat Game (K)
- 7 wrathofthetaboos (G)
- 8 wrathofthetaboos (K)
- 9 Lasthaven (G)
- 10 Jme-cotk (G)
- 11 Ardor3D (G)
- 12 Jmonkeyengine r4099 (G)
- 13 Fairytale-soulfire (G)
- 14 Jmonkeyengine r4490 (G)
- 15 Jmonkeyengine r4490 (S)
- 16 Xenogeddon (G)
- 17 Deathsquadrendevous (G)
- 18 Partindolapana (G)
- 19 Tholos (G)
- 20 Ardor3D (G)
- 21 Cosmic-engine (G)
- 22 Fregatclient3d (G)
- 23 Footballmanagerdesia (G)
- 24 Multiplicity (G)
- 25 Jmerefactoring (G)
- 26 Java3dfh (G)

jMonkeyEngine R3800

R3448

R4099

Evolution Pattern of kern_malloc



Oldest Files Found for Each Query File

Query File in SSHTools	Project Name of Found File	Cover Ratio of Found File	License	Copyright	Last Modified Time
SocketProxySocket.java	CVS client interface in Java	0.88	GPL 2	1998-99 Mindbright Technology AB	2001/11/12
Sftp.java	Apache Ant	0.62	Apache 1.1	2000-2002 Apache Software Foundation	2002/4/15
StringScanner.java	Programmer's Friend 4.1	0.81	CPL 1.0	2000-2003 Manfred Duchrow	2002/9/29
GeneralUtil.java	Gruntspud CVS Client	0.73	GPL 2	2002 Brett Smith	2003/11/12
Base64.java	Base64 notation	0.62	Public Domain	No copyright	2004/1/6
CharBuffer.java	Java Telnet daemon	0.81	GPL 2	2000 Dieter Wimberge	2004/1/16

SSHTools
(all 339 files)

License:	GPL 2
Copyright:	2002-2003 Lee David Painter and Contributors
Last Modified Time:	2007/6/23

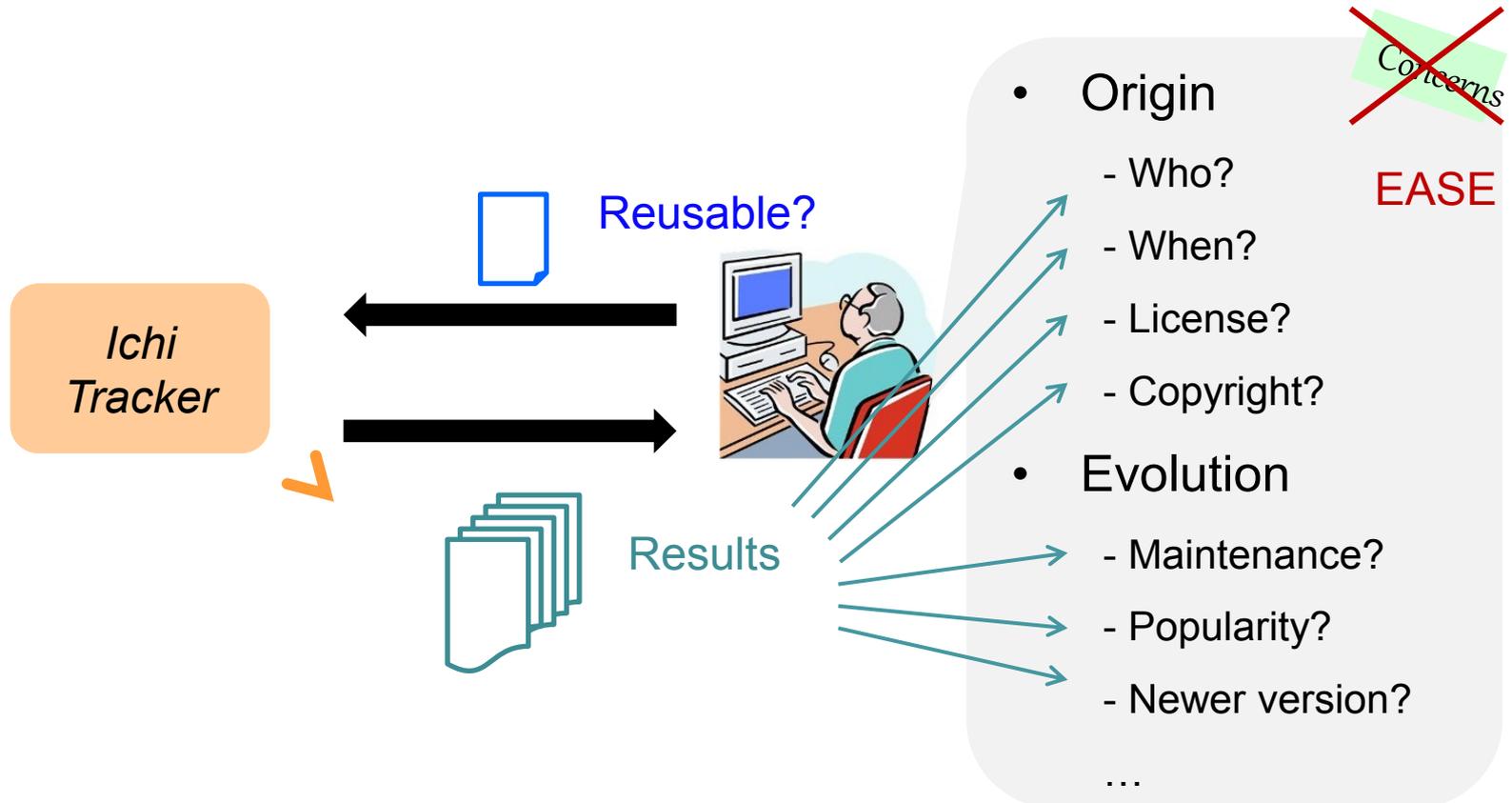
We have found many cases of different projects, different licenses and different copyrights

-> We need to check very carefully when we reuse SSHTools code



Usefulness

With simple check of the output of Ichi Tracker, we can get useful information for the history and evolution of code



Approach and Process

- Choosing good code search engines is a key to get high quality results

GCS >> Koders > SPARS/R

(GCS, Koders, SPARS/R) >> (Google, Bing)

Need good code search engine available!

- Keyword selection strategy:

Incremental strategy: try 1, 2, ... keywords until the header list becomes less than 50

- Decrement strategy, random, less frequently-used keywords, comment keywords, short keywords, ...

Less effective



Other Issues

- Performance
 - Case Studies A and B: 1 to 4 min.
 - Heavily depend on the code search engines and network performance

Acceptable as non-interactive support system
- Quality of search result
 - Non-removed rate at the code clone filtering is an indicator of effectiveness of keyword search
e.g., 0.46 (Case Study A(1) default setting)
 - The final output contains no false positives results



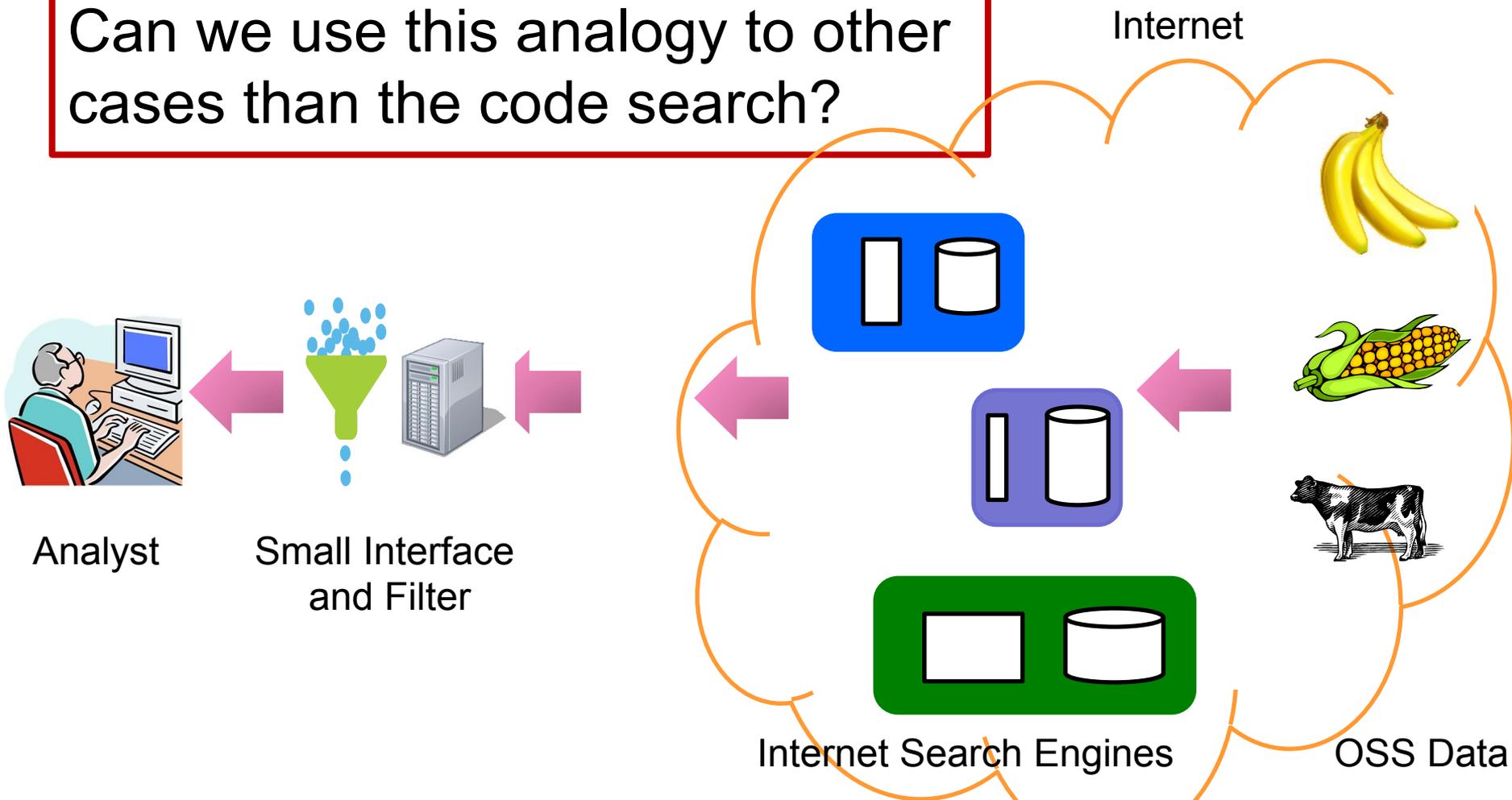
Pros and Cons

- + No huge repository
- + No maintenance cost and inconsistency
- + Easy to get the initial results
- Hard to control search results (engines) in details
- Fragile results depending on search engines



Mining Resources in the Internet Space

Can we use this analogy to other cases than the code search?



Thank you!

