

# Debsources

## Live and Historical Views on Macro-Level Software Evolution

Stefano Zacchioli

[zack@pps.univ-paris-diderot.fr](mailto:zack@pps.univ-paris-diderot.fr)

Laboratoire PPS, Université Paris Diderot

30 June 2014

European Open Symposium on Empirical Software Engineering  
Lille, France



 **Matthieu Caneill, Stefano Zacchiroli**

**Debsources: Live and Historical Views on Macro-Level Software Evolution**

*ESEM 2014: 8th International Symposium on Empirical Software Engineering and Measurement*

preprint available at:

<http://upsilon.cc/~zack/research/publications/debsources-esem-2014.pdf>

# Macro-level software evolution

*“Software evolution in the large”*

— *Gonzalez-Barahona et. al, 2009*

The study of **software evolution**, at the scale of **software collections**, at the granularity they allow (e.g., releases of individual **software components**).

## Debian

- popular FOSS distribution
- 20+ years of history
- one of the largest **curated software collections**
- good proxy of popular/  
relevant FOSS projects
- popular ESE/MSR subject
- *ad hoc* software ecosystem (customs, house tools, etc.)

Goal: ease macro-level software evolution studies on FOSS as a whole, using Debian as a proxy.

# Macro-level software evolution

*“Software evolution in the large”*

— *Gonzalez-Barahona et. al, 2009*

The study of **software evolution**, at the scale of **software collections**, at the granularity they allow (e.g., releases of individual **software components**).

## Debian

- popular FOSS distribution
- 20+ years of history
- one of the largest **curated software collections**
- good proxy of popular/**relevant FOSS projects**
- popular ESE/MSR subject
- *ad hoc* software ecosystem (customs, house tools, etc.)

Goal: ease macro-level software evolution studies on FOSS as a whole, using Debian as a proxy.

# Macro-level software evolution

*“Software evolution in the large”*

— *Gonzalez-Barahona et. al, 2009*

The study of **software evolution**, at the scale of **software collections**, at the granularity they allow (e.g., releases of individual **software components**).

## Debian

- popular FOSS distribution
- 20+ years of history
- one of the largest **curated software collections**
- good proxy of popular/**relevant FOSS projects**
- popular ESE/MSR subject
- *ad hoc* software ecosystem (customs, house tools, etc.)

Goal: ease macro-level software evolution studies on FOSS as a whole, using Debian as a proxy.

# What is Debsources?

- 1 an **infrastructure** to publish Debian source code on the Web
- 2 a **notable instance**,<sup>1</sup> indexing *all* Debian source code to date

For end users:

- browse/search source code
- syntax highlighting
- pinpoint code lines, add msgs

For data miners:

- Debian macro-level evolution
- 20+ years of history
- live/perennial monitoring

DEBSOURCES

package name Search package Code regex Search code

Home Search Documentation About

## Debian Sources

All Debian source are belong to us — Anonymous [1]

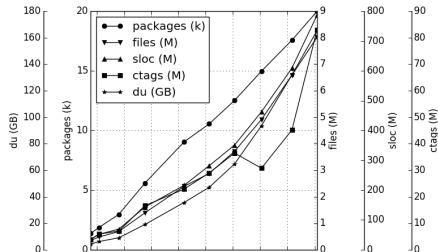
Browse through the source code of the Debian operating system. [Read more](#)

Browse by prefix

Search

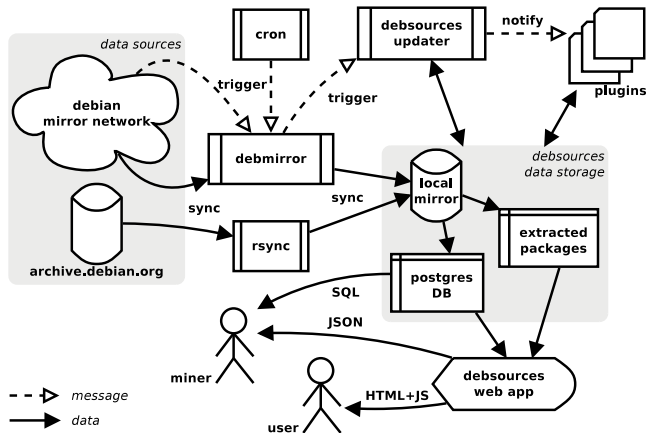
by package name: Search package

the source code (via [codessearch](#)): Search code



<sup>1</sup><http://sources.debian.net>

# Architecture



In essence, Debsources does the *heavy lifting* of maintaining a general purpose [always up to date] storage for Debian source code, enabling plugin authors to focus on *data extraction*.

## Plugin — example (sloccount)

```
def add_package(session, pkg, pkgdir, file_table): # plugin excerpt
    if 'hooks.fs' in conf['backends']:
        if not os.path.exists(slocfile): # run sloccount only if needed
            try:
                cmd = ['sloccount'] + SLOCCOUNT_FLAGS + [pkgdir]
                with open(slocfile_tmp, 'w') as out:
                    subprocess.check_call(cmd, stdout=out,
                                         stderr=subprocess.STDOUT)
            except subprocess.CalledProcessError:
                if not grep(['^SLOC total is zero,', slocfile_tmp]):
                    # rationale: sloccount fails when it can't find source c
                    raise
            finally:
                os.rename(slocfile_tmp, slocfile)
    if 'hooks.db' in conf['backends']:
        slocs = parse_sloccount(slocfile)
        db_package = dbutils.lookup_package(session, pkg['package'],
                                           pkg['version'])
        if not session.query(SlocCount).filter_by(package_id=db_package.id)\
            .first():
            # ASSUMPTION: if *a* loc count of this package has already been
            # added to the db in the past, then *all* of them have
            for (lang, locs) in slocs.iteritems():
                sloccount = SlocCount(db_package, lang, locs)
                session.add(sloccount)
```



# Dataset

data model (excerpt)

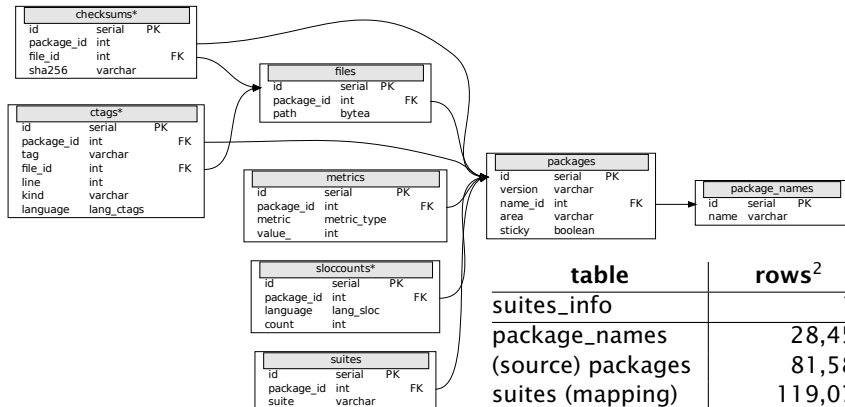


table	rows <sup>2</sup>
suites_info	16
package_names	28,454
(source) packages	81,582
suites (mapping)	119,078
metrics* (e.g., du)	81,582
sloccounts*	290,961
checksums*	33,495,057
ctags*	317,853,685

<sup>2</sup>snapshot, 9 March 2014

# Replication study



Gonzalez-Barahona, Robles, Michlmayr, Amor, and German

Macro-level software evolution: a case study of a large software compilation

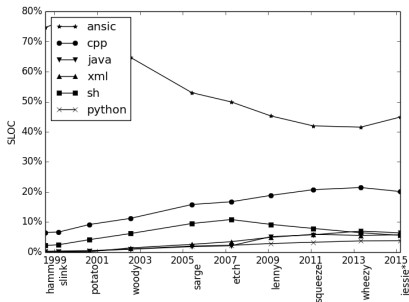
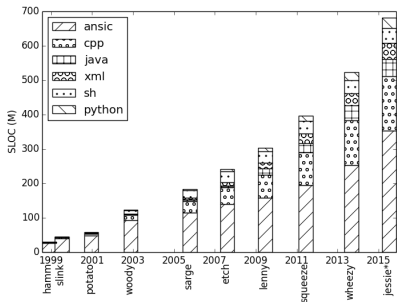
*Empirical Software Engineering*, 14(3):262–285, 2009

original	replica	(remarks)
total (release) size	✓	
package size	✓	
package maintenance	✓	more precise diff size evaluation
programming lang.	✓	more lang. (make, SQL, XML)
file sizes (per lang.)	✓	relying on file extension (bug)
dependencies	✗	no <i>binary</i> packages, yet

We have obtained *slightly different results* than our reference study, but *confirmed the general trends* and updated them in light of *7 extra years of evolution history*.

# Replica highlight #1: programming languages

top-5 most popular programming languages in Debian over time



Recent trends (post-*etch*, 2007):

- C still leads, steady (absolute) growth
- C stops losing (relative) ground to C++
- decrease of Perl & Shell popularity
- Python rises (more maintainable glue code?)
- Lisp halves its popularity
- Java no longer under-represented

## Replica highlight #2: package maintenance

Changes between Debian releases: 'c' for common, 'u' for unchanged (upstream), and 'm' for modified packages (common \ unchanged):

<i>from</i>	<i>to</i>									
	<i>slink</i>	<i>potato</i>	<i>woody</i>	<i>sarge</i>	<i>etch</i>	<i>lenny</i>	<i>squeeze</i>	<i>wheezy</i>	<i>jessie*</i>	<i>sid*</i>
<i>harm</i>	1324c 842u	1198c 463u	1079c 270u	958c 175u	864c 148u	782c 124u	719c 100u	670c 81u	648c 75u	663c 75u
<i>slink</i>		1657c 742u	1455c 384u	1281c 252u	1155c 210u	1037c 172u	941c 136u	881c 113u	852c 105u	872c 105u
<i>potato</i>			2456c 935u	2118c 551u	1881c 436u	1683c 352u	1497c 271u	1399c 220u	1359c 210u	1387c 211u
<i>woody</i>				4588c 1688u	3953c 1156u	3497c 908u	3021c 633u	2787c 520u	2680c 486u	2752c 494u
<i>sarge</i>					7671c 3832u	6828c 2597u	5903c 1717u	5353c 1369u	5102c 1240u	5259c 1272u
<i>etch</i>						9230c 4578u	8041c 2906u	7216c 2205u	6881c 1948u	7088c 2000u
<i>lenny</i>							10836c 5272u	9631c 3676u	9181c 3153u	9457c 3249u
<i>squeeze</i>								13117c 6812u	12464c 5425u	12902c 5622u
<i>wheezy</i>									16543c 10132u	17042c 10519u
<i>jessie*</i>										19795c 19593u
	<i>from previous suite to</i>									
	<i>slink</i>	<i>potato</i>	<i>woody</i>	<i>sarge</i>	<i>etch</i>	<i>lenny</i>	<i>squeeze</i>	<i>wheezy</i>		
<b>modified pkgs</b>	556m	1305m	3127m	4462m	2879m	3287m	4129m	4453m		
<b>changed files per pkg</b>	54.6%	64.4%	65.3%	67.5%	58.9%	59.8%	60.4%	57.2%		

Debsources is a flexible platform to monitor large FOSS collections over long periods of time. Its main instance and dataset are valuable resources for scholars interested in macro-level software evolution.

# Thanks! Questions?

Stefano Zacchioli

[zack@pps.univ-paris-diderot.fr](mailto:zack@pps.univ-paris-diderot.fr)

<http://upsilon.cc/zack>



Matthieu Caneill, Stefano Zacchioli

Debsources: Live and Historical Views on Macro-Level Software Evolution

*ESEM 2014: 8th International Symposium on Empirical Software Engineering and Measurement*