

# Mining GitHub for fun and profit

Georgios Gousios // @gousiosg  
TU Delft

## GitHub

- Project forge (Git, Issues)
- Project wiki (Git based)
- Social network (star repository, follow developer)
- Forks and Pull requests
- Comments on commits, issues and pull requests

## GitHub is cool because

- api.github.com
- process and product data
- interconnection

## Mining the GitHub API:



Only current state



Missing data schema



5k reqs/hr



**GHTorrent**: a scalable, queryable, offline mirror of GitHub data

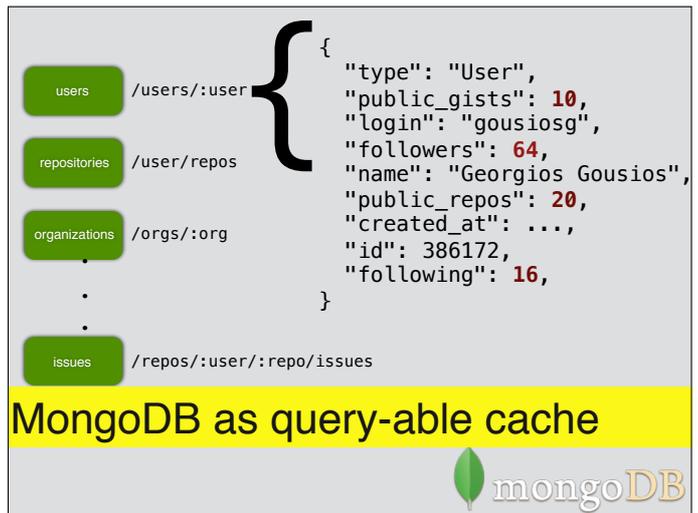
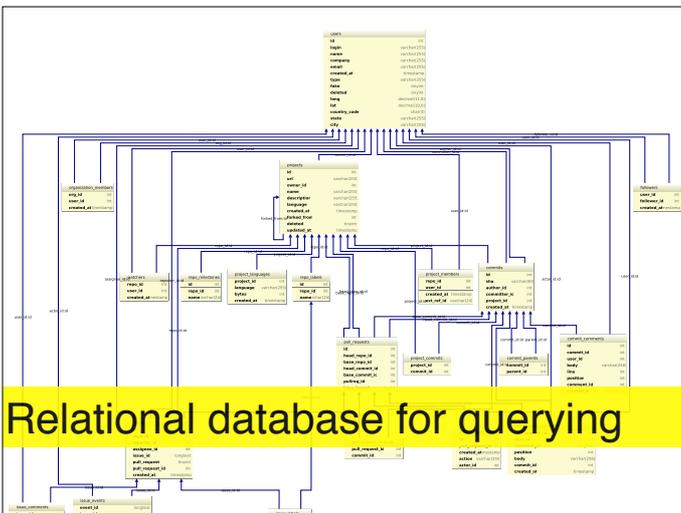
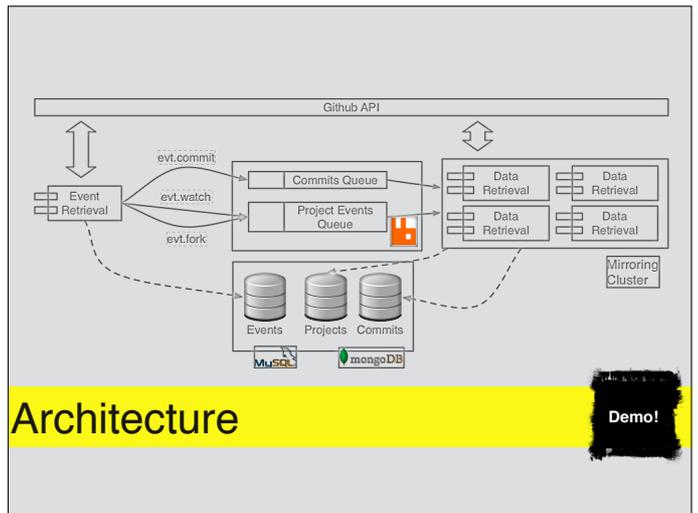
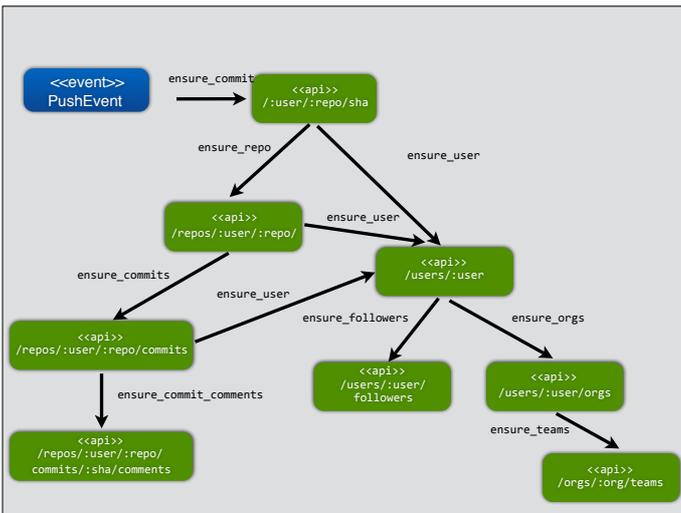
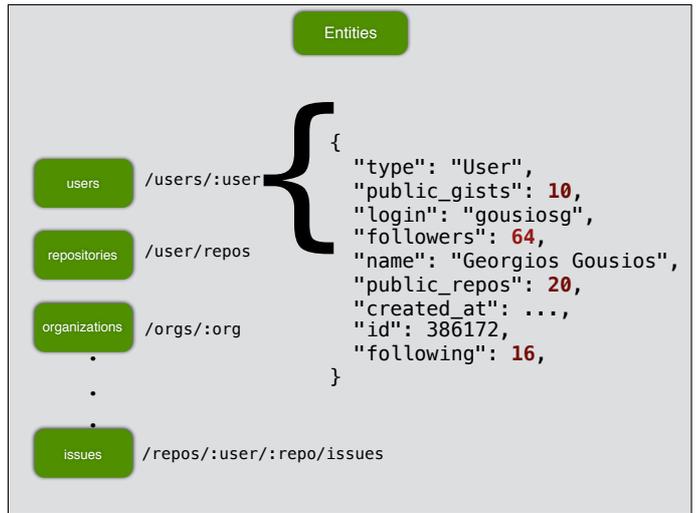
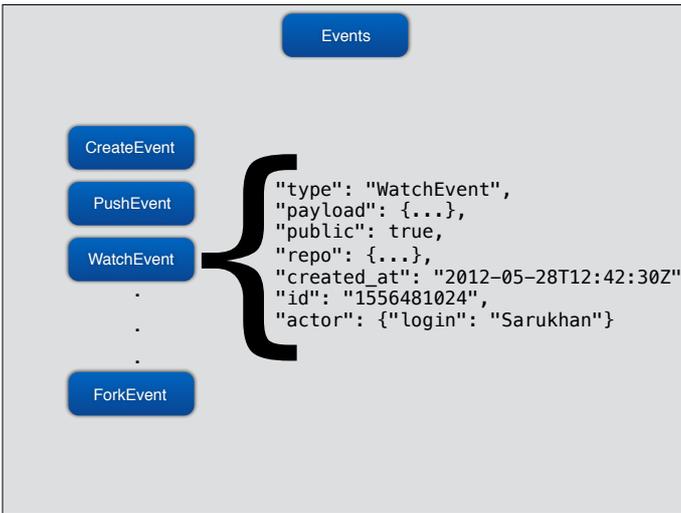
## api.github.com

Entities

- static view
- interlinked
- current state

Events

- dynamic view
- generated by user actions
- affect current entity state
- can be browsing roots



```

1 select country_code, count(*)
2   from users
3   where fake is false
4     and deleted is false
5     and country_code is not null
6   group by country_code
7   order by count(*) desc
8
9

```

Conn: rs0 DB: ghtorrent Table Structure History 1. F

Export Show SQL Refresh Auto Refresh

country_code	count(*)
us	283211
gb	57550
cn	50819
in	46204

**User geolocation**

de	45473
br	30329
fr	33003
ru	24441

```

SQL*
1 select u.login, p.name, p.language, pl.language as main_language,
2   pl.bytes, pl.created_at as date_measured
3   from project_languages pl, projects p, users u
4   where u.id = p.owner_id
5     and pl.project_id = p.id
6     and p.name = 'scala'
7     and u.login = 'scala';|
8
9

```

Connection DB: ghtorrent Table: project\_langu History 1. Result

Export Show SQL Refresh Auto Refresh

login	name	language	main_language	bytes	date_measured
scala	scala	Scala	scala	15934585	2015-10-29 10:46:57
scala	scala	Scala	java	1023702	2015-10-29 10:46:57
scala	scala	Scala	javascript	241109	2015-10-29 10:46:57
scala	scala	Scala	shell	57442	2015-10-29 10:46:57
scala	scala	Scala	python	44027	2015-10-29 10:46:57
scala	scala	Scala	c#	37216	2015-10-29 10:46:57
scala	scala	Scala	html	6396	2015-10-29 10:46:57
scala	scala	Scala	c	141	2015-10-29 10:46:57

**Full language retrieval**

## Operation modes

- **Normal operation:** Follow the event timeline, apply dependency-based retrieval
- **Bi-monthly updates:** Refresh the state of all repos/users (cater for deleted repos, changed user locations etc)
- **Full retrievals:** Get all info for a repo/user, in case stuff is missing



### MySQL database dumps

**New!** As of MySQL dump [mysql-2015-09-25](#), we are distributing CSV files (one file per table) instead of [mysql dump](#) based backups. The provided directory including a restore script and instructions on how to do the restore. See more information [here](#).

You can also [query MySQL](#). It is always loaded with the latest dump.

- 2016-02-01 (34969 MB)
- 2016-01-16 (34178 MB)
- 2016-01-08 (33942 MB)
- 2015-09-25 (32273 MB)
- 2015-08-07 (31537 MB)
- 2015-06-18 (33476 MB)
- 2015-04-01 (25075 MB)
- 2015-01-04 (16583 MB)
- 2014-11-10 (14418 MB)
- 2014-08-18 (11485 MB)
- 2014-04-02 (7013 MB)
- 2014-01-02 (5646 MB)
- 2013-10-12 (4312 MB)

### MongoDB database dumps

The MongoDB backups come in different formats:

- Bi-monthly and by collection, from August 2015 up to Nov 2015
- Daily from 2015-12-01, includes all collections

## Periodic dumps of DBs

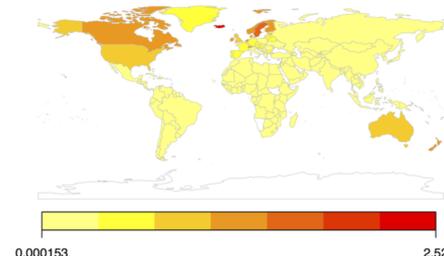
2016-02-08, 2016-02-07, 2016-02-06, 2016-02-05, 2016-02-04, 2016-02-02, 2016-02-01, 2016-01-31, 2016-01-30, 2016-01-29, 2016-01-28, 2016-01-27, 2016-01-26, 2016-01-25, 2016-01-22, 2016-01-21, 2016-01-20, 2016-01-19, 2016-01-18, 2016-01-17, 2016-01-16, 2016-01-15, 2016-01-14, 2016-01-13, 2016-01-12, 2016-01-11, 2016-01-10, 2016-01-09, 2016-01-08, 2016-01-07, 2016-01-06, 2016-01-05, 2016-01-04, 2016-01-03, 2016-01-02, 2016-01-01, 2015-12-31, 2015-12-30, 2015-12-29, 2015-12-28, 2015-12-25, 2015-12-24, 2015-12-23, 2015-12-22, 2015-12-21, 2015-12-20, 2015-12-19, 2015-12-18, 2015-12-17, 2015-12-16, 2015-12-15, 2015-12-14, 2015-12-13, 2015-12-12, 2015-12-11, 2015-12-10, 2015-12-09, 2015-12-08, 2015-12-07, 2015-12-06, 2015-12-05, 2015-12-04, 2015-12-03, 2015-12-02, 2015-12-01.

**Georgios Gousios**  
@gousiosg

Follow

World programmer density map, calculated by geocoding all @github users that declared a location (w/ @ghtorrent)

**World Programmer Density**  
(GitHub Users per 1000 inhabitants)



RETWEETS 381

LIKES 187





# Python developers in Estonia

```
select distinct(u.login), u.location
from users u, commits c,
     projects p, project_commits pc
where date(c.created_at) between date('2016-06-01')
                                and date('2016-07-01')
and pc.commit_id = c.id
and p.id = pc.project_id
and c.author_id = u.id
and u.country_code = 'ee'
and p.language = 'Python';
```

login	location
plaas	Estonia
si642	Estonia
riwol	Tallinn, Estonia
livenson	Estonia
mittya	Moon
rbozz	Tallinn
emdoon	Tallinn
argonodots	Estonia
Neobiz	Tallinn, Estonia
mittelek	Tallinn, Estonia
asollettav	Tallinn, Estonia
Jyrm042	Estonia, Tallinn
Huudlejev	Estonia
raidoz	Estonia
vruusmann	Estonia
ampr1	Estonia
Isaffre	Estonia
mika351	Tallinn
Sonmepah	Estonia
ropod7	Estonia
siias	Estonia
jaakerisalu	Tallinn, Estonia
kendax	Estonia
artizirk	Estonia
avaikarvamus	Tallinn, Estonia

27 rows in set (0.80 sec)

The screenshot shows the GHTorrent website with a navigation menu (Docs, Downloads, Query, Visualizations, Reports, Datasets, Hall of Fame). The main content area is titled 'Querying MongoDB' and includes a 'Query' input field. Below the input field, there is a 'Submit' button and a 'Sponsors' section listing Microsoft and TU Delft. A 'Obtaining access' section provides instructions on how to connect via SSH. A terminal window at the bottom shows the command `mongo -u ghtorrent -p ghtorrent github` and the output 'MongoDB shell version: 3.8.3'. A yellow banner at the bottom of the terminal window reads 'Query MongoDB programmatically' with a 'Demo!' button.

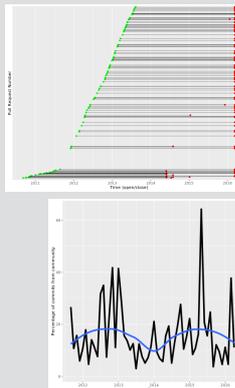
The screenshot shows the GHTorrent website with a navigation menu (Docs, Downloads, Services, Visualizations, Reports, Datasets, Hall of Fame). The main content area is titled 'Streaming updates from GHTorrent' and includes a 'Connection details' section with instructions on how to obtain access. A 'Declaring queues' section explains the AMQP protocol and provides a list of steps to connect to the server. A 'Connecting to the server' section provides instructions on how to connect to the server. A terminal window at the bottom shows the command `ch = conn.create_channel` and `exchange = ch.topic('ght-streams')`. A yellow banner at the bottom of the terminal window reads 'Streaming updates'.

The screenshot shows the GHTorrent dashboard with a navigation menu (GHTorrent stats, Zoom Out, Last 1 hour, Refresh every 1m). The dashboard displays several real-time analytics graphs: 'API requests per minute', 'API errors per minute', 'Remaining API Quota', 'API response time (median, mean, 95%)', and 'Total requests last hour'. A large red number '87542' is displayed in the center of the dashboard. A yellow banner at the bottom of the dashboard reads 'Real time analytics' with a 'Demo!' button.

# Openess reports

How open is your project to community contributions?

- 5k projects
- every 15 days



```
$ gem install sqlite3 bundler
$ git clone https://github.com/gousiosg/github-mirror
$ cd github-mirror
$ bundle install
$ mv config.yaml.standalone > config.yaml
$ ruby -Ilib bin/ght-retrieve-repo -t token rails rails
```

Roll your own dataset

Demo! <http://ghtorrent.org/pullreq-perf>



VAGRANT

```

$ apt-get install vagrant
$ git clone https://github.com/ghtorrent/ghtorrent-
vagrant.git
$ cd ghtorrent-vagrant
$ vagrant ssh
$ ruby -Ilib bin/ght-retrieve-repo -t token rails rails

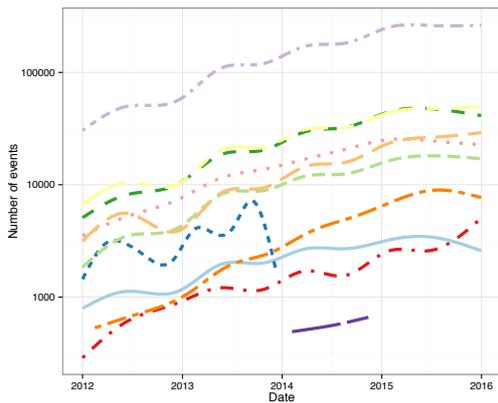
```

Even simpler with Vagrant

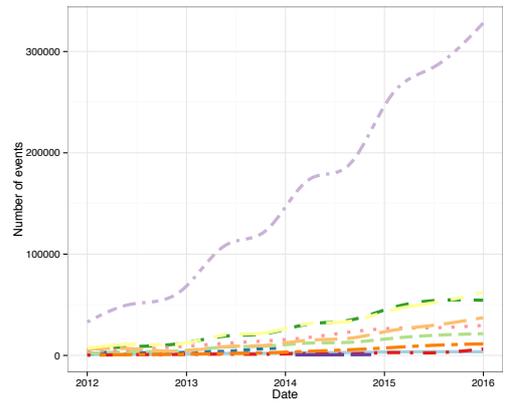
# Growth

	MongoDB	MySQL	Diff 2016/2013
Events	476		11.1x
Users	6,7	9,2	8.4x
Repos	28	25,5	21.8x
Commits	367	362	12.3x
Issues	24,1	25,3	10.3x
Pull requests	11,9	11,1	9.7x
Issue comments	42	43	14.6x
Watchers	51	37	6.6x

G. Gousios, The Evolution of GHTorrent: Growing an Open Access Dataset 10x



G. Gousios, The Evolution of GHTorrent: Growing an Open Access Dataset 10x



G. Gousios, The Evolution of GHTorrent: Growing an Open Access Dataset 10x

# Stats and Impact

- Since Feb 2012 **180+** users, **100+** institutions
- 12TB** in MongoDB
- 4.5B** rows in MySQL
- 2GB** per hour
- 170k** API reqs/hour
- 50** user donated API keys
- 100+** papers
- 3** data mining challenges
- 3** best paper awards

G. Gousios, "The GHTorrent dataset and tool suite," in MSR, 2013: 233-236  
G. Gousios and D. Spinellis, "GHTorrent: GitHub's Data from a Firehose," in MSR, 2012, 12-21

# Impact

MSR 2016: 35% of all papers on Github are done with GHTorrent



Valerio Cosentino, Javier Luis, and Jordi Cabot. 2016. Findings from GitHub: methods, datasets and limitations. MSR '16, pp 137-141

# MSR 2014 challenge

- Pull requests
- Project developer attraction
- Sentiment analysis of commit comments
- Biodiversity of project ecosystems
- Project openness
- Design discussions
- Co-evolution of documentation and popularity



looking for include code (and instructions to use it) in a GitHub repository, an academic write-up of your analysis, or an informal prose write-up. If you're not linking to a repository, you should submit a [Gist](#) with your documentation.

After the submission deadline on August 25th, GitHub employees will review and vote on all entries to pick the three top winners. We'll send out notifications to those top three by mid-September.

## Data Sources

GitHub activity data is available from several publicly-available sources. Here are a few links to get you started:

- Our very own [API](#).
- The [GitHub Archive](#), providing historical archives of our public timeline data.
- [Google BigQuery](#), where GitHub's public timeline is a featured public dataset; see the [GitHub Archive home page](#) for getting started instructions.
- [GHTorrent](#), which maintains a functional model of GitHub activity data and offers archives for download.

## Pro Tips

There are a few things we're looking for when we score your entry:

- **Innovation/Story:** Does your entry tell a good, data-driven story? Does it reveal interesting insights about GitHub activity? We love it when we're surprised by new insights hidden in our own data.
- **Accuracy:** Is your analysis accurate? Do accompanying visualizations clearly and unambiguously convey your conclusions?
- **Completeness:** Is your entry a code submission? If so, is your code well-organized and documented? Can others easily understand and reproduce your analysis from the materials you've submitted?

## The Prizes

**GitHub API rate limit and network latency.** The mining program was able to check about one repository ID per second, which was slowed by network latencies, or, once 5000 API calls had been made in one hour, was throttled by GitHub. The apparent solution was to create multiple accounts, each providing 5000 API calls per hour. After contacting GitHub to request help with this issue, they indicated that they do not want users to create multiple accounts for mining projects because it can put a strain on their servers, slowing the service for regular users. An alternative strategy was proposed by GitHub of using ghtorrent<sup>3</sup> to find Python projects without using the API. However, at this point in the project, enough data had been acquired to begin analysis and a determination was made to stop development of this mining program and focus on analysis. Future mining efforts are encouraged to obtain repository information from this database instead of crawling through all projects using the GitHub API, like `tour.de.source` did.

Chapman, Carl Allen, *Usage and refactoring studies of python regular expressions*, Graduate Theses and Dissertations, Paper 15139. Iowa State University, 2016

TravisTorrent, a [GHTorrent](#) partner project, provides free and easy-to-use [Travis CI](#) build analyses to the masses through its open database.

TravisTorrent provides access to a database of hundreds of thousands of analyzed travis builds in **under 10 seconds**.

**Access TravisTorrent now**

**What does TravisTorrent do?**

We access the [Travis CI](#) API and for each build, combine vanilla API data (such as build number and build result), with an analysis of the build log (such as how many tests were run, which test failed, ...) and repository and commit data from [GitHub](#) (such as latency between pushing and building), acquired through [GHTorrent](#).

**How to cite?**

Moritz Beller, Georgios Gousios and Andy Zaidman. *Oops, my tests broke the build: An analysis of Travis CI builds with GitHub*. No. e1984v1. PeerJ Preprints, 2016.

From: **Gianugo Rabellino** [Gianugo.Rabellino@microsoft.com](mailto:Gianugo.Rabellino@microsoft.com)  
Subject: [REDACTED]  
Date: 9 Sep 2015 18:49  
To: [Georgios Gousios](mailto:Georgios.Gousios@cs.ru.nl) [g.gousios@cs.ru.nl](mailto:g.gousios@cs.ru.nl)

Dear Mr. Gousios,

[REDACTED]

[REDACTED] At Microsoft we have been using your GHTorrent project and we can't thank you enough for the business critical insights your data have and will be able to provide.

[REDACTED]

**Data Lake**  
Batch, real-time, and interactive analytics made easy

- ✓ Store and analyze data of any kind and size
- ✓ Develop faster, debug and optimize smarter
- ✓ Interactively explore patterns in your data
- ✓ No learning curve—use U-SQL, Spark, Hive, HBase, and Storm
- ✓ Managed and supported with an enterprise-grade SLA
- ✓ Dynamically scales to match your business priorities
- ✓ Enterprise-grade security with Azure Active Directory
- ✓ Built on YARN, designed for the cloud

Try it now

**NEW!**

Microsoft / [ghinsights](#)

Code Issues 0 Pull requests 0 Wiki Pulse Graphs

GHInsights is a data processing pipeline using Azure Data Factory and Azure Data Lake. It processes GitHub data from the ghtorrent project. The resulting processed data is available in Azure Data Lake for users to query, generate reports, and analyze GitHub projects.

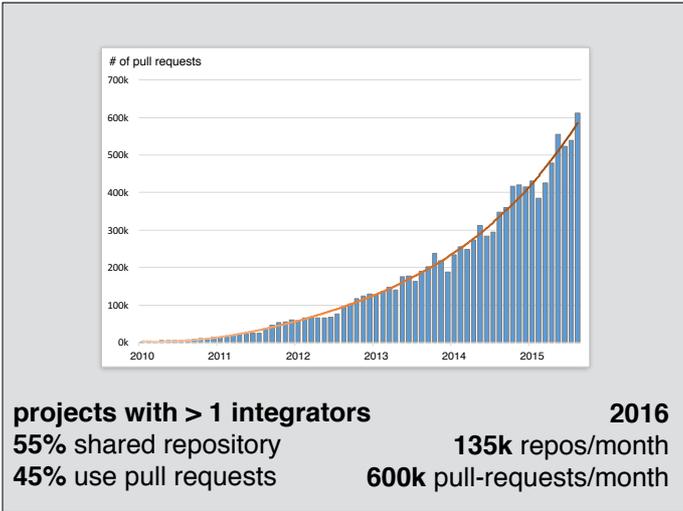
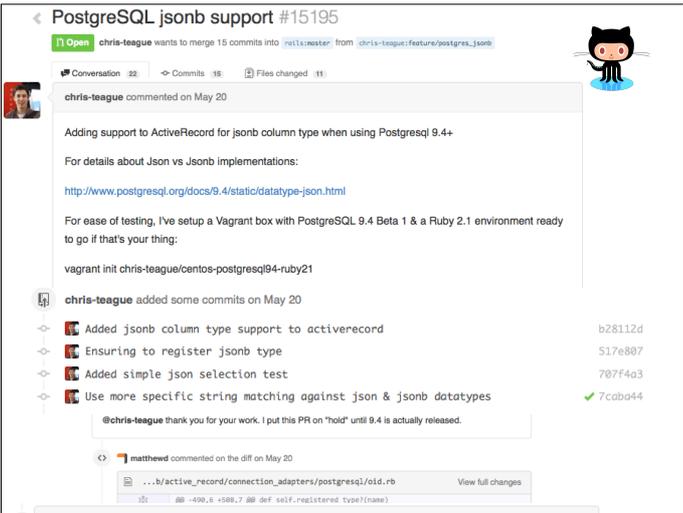
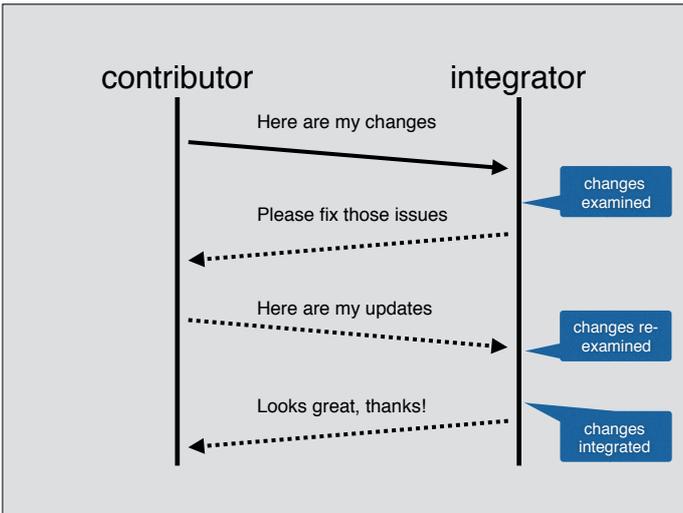
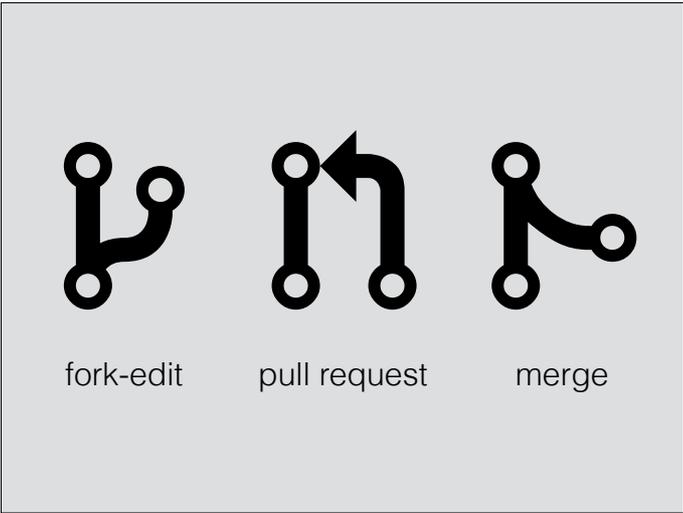
71 commits 1 branch 0 releases 4 contributors

Branch: master New pull request Create new file Upload files Find file Clone or download

- kelevis committed on GitHub Merge pull request #13 from riv/master Latest commit 5128996 12 days ago
- DataExport Removing this readme, content now merged into main file 3 months ago
- ghinsights Fixed the commit 12 days ago
- gitignore Adding a .gitignore to not include environment specific configurations. 2 months ago
- LICENSE Adding license and link in readme 3 months ago
- README.md add usql pointer 3 months ago
- README.md

**GHInsights**

GHInsights is a dataset and processing pipeline for GitHub event and entity data. It enables you to create your own insights on all or a portion of the activity and content on GitHub. Fundamentally GHInsights is based on [GHTorrent](#), an open, collaborative project for gathering exposing GitHub interactions. GHInsights takes that data and makes it available in [Azure Data Lake](#). This gives you an easily accessible dataset and scalable computes resources so you can create the insights you need without having to gather and manage the many terabytes of data involved.



## Large scale collaboration

Repository	Commits	Pull requests	Code reviews	Issues	Issue comments	Total
isaacs/npm	100	167	291	2568	3302	6428
torvalds/linux	5968	67	178	0	0	6213
symfony/symfony	1021	1261	1752	1844	2160	8036
jquery/jquery-mobile	212	431	384	2888	3008	6923
joyent/node	657	833	1127	2304	2805	7726
CocoaPods/Specs	2658	2584	1364	515	268	7389
gitter/gitter	605	871	1142			
angular/angular.js	875	1306	1751			
rails/rails	2699	2315	4090			
mxcl/homebrew	3426	3125	4492			

**23,500**

http://www.gousios.gr/blog/The-triumph-of-onlin

angular / angular.js  
 327 Open ✓ 5,360 Closed

rails / rails  
 468 Open ✓ 12,609 Closed

joyent / node  
 303 Open ✓ 3,347 Closed

Too successful?

"Lack of knowledge of git from contributors; most don't know how to resolve a merge conflict."

"Sifting through the GitHub information flood to find what, if any, I should address."

"Dealing with loud and trigger-happy developers."

Too successful?

## How does the PR process look like?

Across 5k repos/ ~1M pull requests

- 85% merged, 70% with merge button
- 80% < 150 lines, < 7 files, 3 commits
- 66% < 1 day to merge
- 80% 4 comments, 3 participants

Mostly rejected due to observability/awareness issues (not technical!)

**An Exploratory Study of the Pull-based Software Development Model**

**Abstract**  
 This study presents an exploratory study of the pull-based software development model. It aims to understand the characteristics of pull requests (PRs) and the factors that influence their acceptance and merging. The study is based on a large dataset of PRs from various open-source projects. The results show that most PRs are small, focused on specific tasks, and are merged quickly. However, a significant portion of PRs are rejected, often due to lack of observability or awareness of the project's state. The study also identifies several factors that affect PR acceptance, such as the quality of the PR, the involvement of the project maintainers, and the project's maturity.

G. Gousios, M. Pinzger, and A. van Deursen, "An Exploratory Study of the Pull-based Software Development Model," ICSE 2014, pp. 345-355

## Which factors affect PR acceptance?

?

## Which factors affect PR acceptance?

- Do we know the submitter?
- Can we handle the workload?
- What does the PR look like?
- How ready is our project for PRs?

G. Gousios et al., "Distributed software development with pull requests". Unpublished.

## Which factors affect the time to process PRs?

- Can we handle the workload?
- Do we know the submitter?
- How ready is our project for PRs?
- What does the PR look like?

G. Gousios et al., "Distributed software development with pull requests". Unpublished.

# What do *integrators* actually believe?

Survey of **650** integrators

Generally, few complains about the process

Points of pain are mostly social (workload, drive-by PRs, explaining rejection)

Needed tools

Quality analysis

Impact analysis

Work prioritization



G. Gousios, A. Zaidman, M.-A. Storey, and A. van Deursen, "Work Practices and Challenges in Pull-Based Development: The Integrator's Perspective." ICSE 2015, pp. 358-368.

# What do *contributors* actually believe?

Survey of **640** contributors.  
Similar issues, reversed

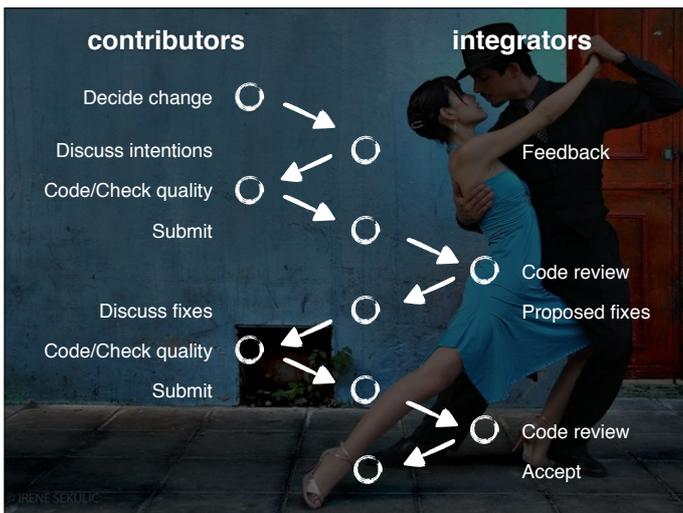
Awareness

Asynchrony

Responsiveness



G. Gousios, M.-A. Storey, and A. Bacchelli, "Work Practices and Challenges in Pull-Based Development: The Contributor's Perspective." ICSE, 2016.



# What are the integrators' biggest challenges?

Quality

Prioritization

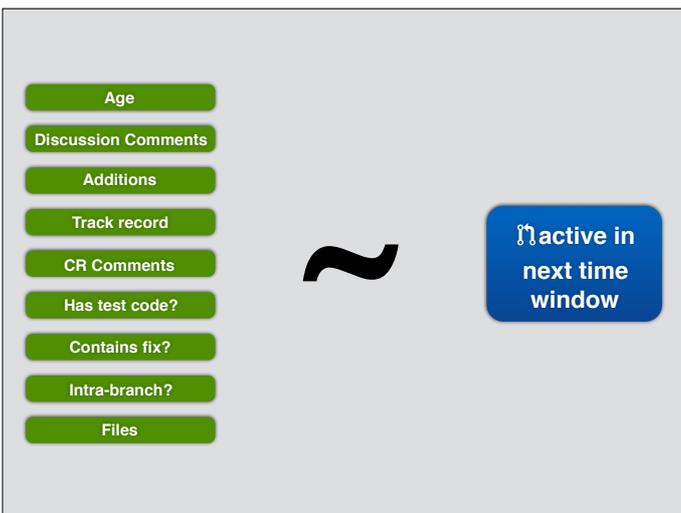
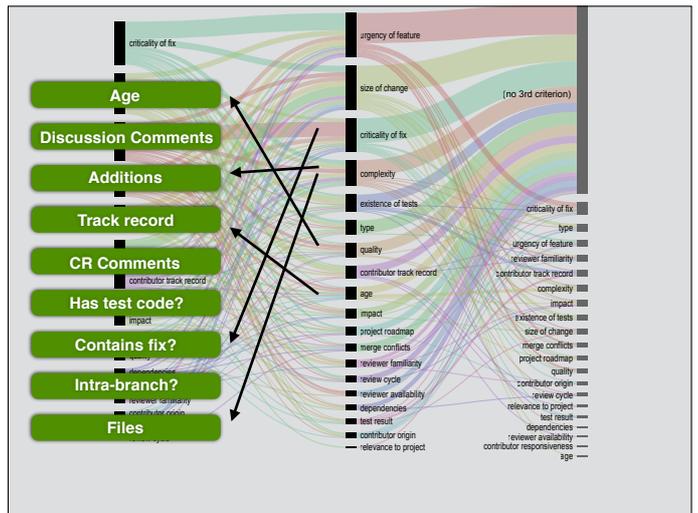
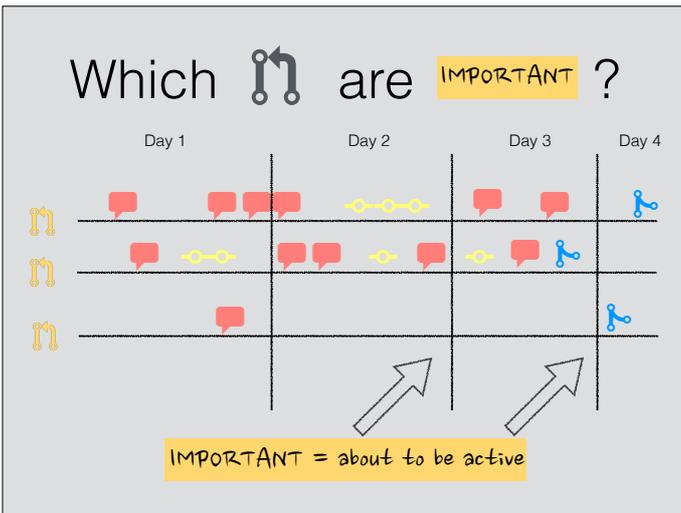
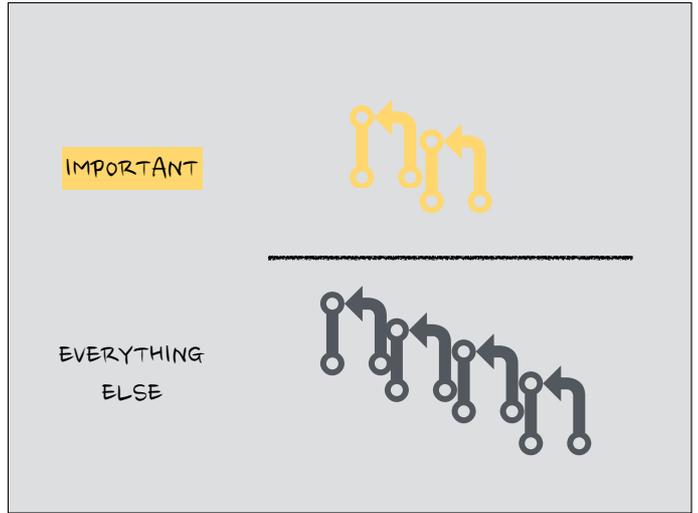
## Work Practices and Challenges in Pull-Based Development: The Integrator's Perspective

Georgios Gousios<sup>1</sup>, Andy Zaidman<sup>2</sup>, Margaret Anne Storey<sup>3</sup>, Arin van Deursen<sup>4</sup>  
<sup>1</sup> Radboud University Nijmegen, the Netherlands  
<sup>2</sup> Email: g.gousios@cs.ru.nl  
<sup>3</sup> Delft University of Technology, the Netherlands  
<sup>4</sup> Email: a.vandeursen@cs.ru.nl

Abstract—In the pull-based development model, the integrator has the central role of accepting and merging contributions. This work focuses on the role of the integrator and investigates how they manage the workflow of pull requests. We study the workflow and challenges of the integrator, and how they manage the workflow of pull requests. We study the workflow and challenges of the integrator, and how they manage the workflow of pull requests. We study the workflow and challenges of the integrator, and how they manage the workflow of pull requests.

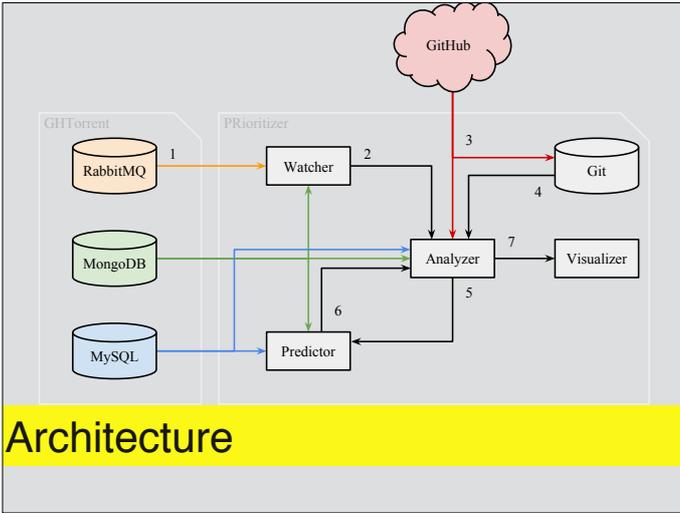
# Characteristics of PRs

- They change fast (66% are processed within a day)
- Attention is only needed after an event occurred
- They are integrated with other tools



	Precision	Recall	AUC	Accuracy
Random Forests	0.66	0.63	<b>0.89</b>	<b>0.86</b>
Naive Bayes	0.34	0.79	0.75	0.60
Logistic regression	0.36	0.84	0.81	0.62

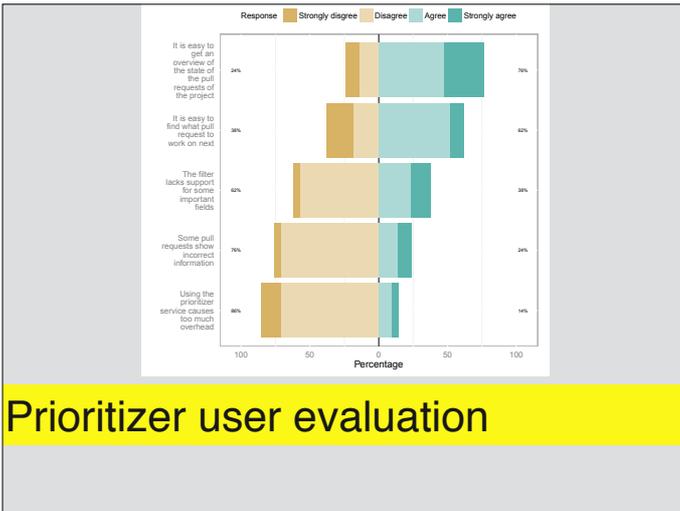
E. van der Veen, G. Gousios, and A. Zaidman, "Automatically Prioritizing Pull Requests," MSR, 2015, pp. 357-361.



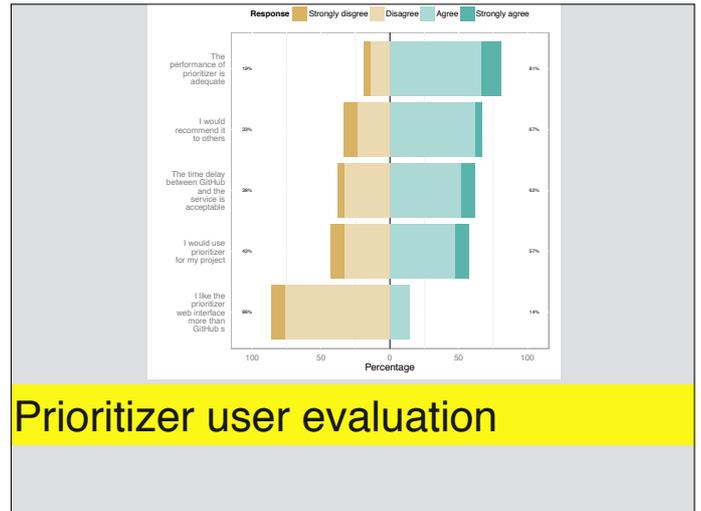
## Architecture

The screenshot shows a GitHub pull request list. At the top, there are filters for Branch, Mergeable, Author, Tests, Conflicts, and Sort. A dropdown menu is open, showing sorting options: Clear, Newest, Oldest, Largest, Smallest, Most commented, Least commented, Most conflicts, Least conflicts, Most contributed, Least contributed, Best accept rate, and Worst accept rate. The pull request list below shows items like 'spray:wip-15681-mathias-final', 'sirthias (core member) 2% contributed commit(s)', and '+str #15588 Java 6 Synchronous File Sink / Source'.

## User interface



## Prioritizer user evaluation



## Prioritizer user evaluation

The screenshot shows a GitHub pull request page for the repository Microsoft/ChakraCore. It displays a list of pull requests with details such as the author, title, status (Open, Closed), and mergeability. The first pull request is 'Force-win7 on Legacy Daily Builds because the test machinery doesn't automatically detect Windows Server 2008 R2 as Windows 7'. The page also shows repository statistics like 337 watchers, 4,626 stars, and 477 forks.

## Pourquoi

The image features the GitHub logo (Octocat) in the center. Below it, the text 'Based upon...' is displayed in a large, black, sans-serif font.

# Promises & Perils

How good is GitHub data for SE research?

- 5 promises
- 13 perils
- Ways to select projects for research
- Examined MSR 2014 papers, all susceptible to perils!



E. Kalliamvakou, G. Gousios, K. Blincoe, L. Singer, D. M. German, and D. Damian, "An in-depth study of the promises and perils of mining GitHub," *Empirical Software Engineering*, pp. 1–37, 2015

# Issue lifetime predictions

Work dual to PRs

Need to train in short term time windows

Difficult to generalise across projects

Contextual features more important than generic ones



Riivo Kikas, Marlon Dumas, and Dietmar Pfahli, *Using dynamic and contextual features to predict issue lifetime in GitHub projects*. MSR '16, pp 291–302

# Reviewer recommendation

Exploit @mention networks to propose top-3 reviewers for incoming pull requests.

Accuracy ~60% on top-3 recommendation



Yue Yu, Huaimin Wang, Gang Yin, Tao Wang, *Reviewer recommendation for pull-requests in GitHub: What can we learn from code review and bug assignment?* IST, 2016

# Automated code review

Examine how “natural” the PR code is WRT the project’s code base.

Accepted PRs are significantly similar to the project

More debated PRs are significantly less similar



Vincent J. Hellendoorn, Premkumar T. Devanbu, and Alberto Bacchelli, 2015. *Will they like this?: evaluating code contributions with language models*. MSR '15, pp157-167

# Not only OSS

“...GitHub’s transparency and popular workflow can promote open collaboration, allowing organizations to increase code reuse and promote knowledge sharing across their teams.”



Eirini Kalliamvakou, Daniela Damian, Kelly Blincoe, Leif Singer, and Daniel M. German, *Open source-style collaborative development practices in commercial projects using GitHub*. ICSE 2015, pp 574-585.

# Gender and Tenure diversity

“Our study suggests that, overall, when forming or recruiting a software team, increased gender and tenure diversity are associated with greater productivity.”



Vasilescu, Posnett, Ray, van den Brand, Serebrenik, Devanbu, and Filkov, *Gender and Tenure Diversity in GitHub Teams*. CHI, 2015

# Geographical diversity

Traces of bias on contributions from certain countries

- Contributors perceive it
- Integrator's do not



A. Rastogi, N. Nagappan, and G. Gousios, "All contributors are equal: some contributors are more equal than others.", TR, 2016

# Gender & Contributions

When gender is identifiable: women rejected more often

When gender is *not* identifiable: women accepted more often



Terrell J. Kofink A, Middleton J, Raineart C, Murphy-Hill E, Parmin C. (2016) Gender bias in open source: Pull request acceptance of women versus men. *PeerJ PrePrints* 4:e1733v1

# Productivity

Study how context-switching between repos affects productivity. If

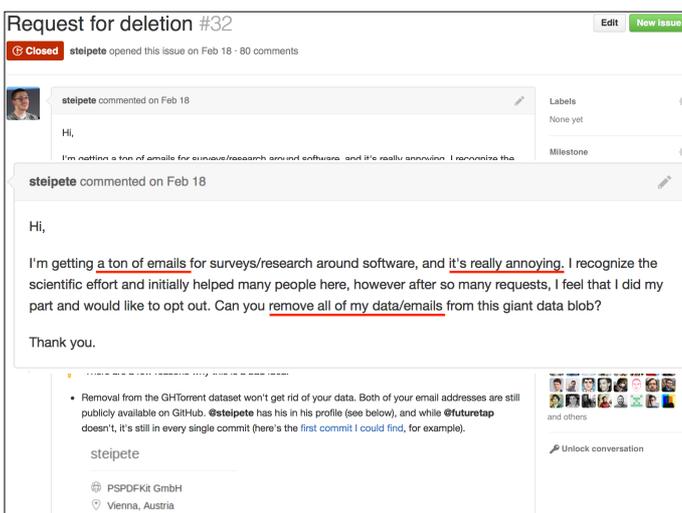
- rate of switching
- number of projects involved in are high, productivity goes down



Bogdan Vasilescu, Kelly Blincoe, Qi Xuan, Casey Casalnuovo, Daniela Damian, Premkumar Devanbu, and Vladimir Filkov. *The sky is not the limit: multitasking across GitHub projects*. ICSE 2016, pp 99—1005



## The #issue32 incident



The screenshot shows a GitHub issue page for 'Request for deletion #32'. The issue is marked as 'Closed' and was opened by 'steipete' on Feb 18. The main comment from 'steipete' reads: 'Hi, I'm getting a ton of emails for surveys/research around software, and it's really annoying. I recognize the scientific effort and initially helped many people here, however after so many requests, I feel that I did my part and would like to opt out. Can you remove all of my data/emails from this giant data blob? Thank you.' Below the comment, there is a list of users who commented, including 'PSPDFKit GmbH' from Vienna, Austria. The issue is currently locked.

# I am not a lawyer!

- Other commenters are no lawyers either
- The law is complicated and open to interpretation

## Two important issues

- *Copyright*: Who owns the data?
- *Privacy*: How does GHTorrent protect users from personal data misuse?

## Copyright — General terms

- For *original content*, the *publisher* maintains full copyright by default
  - Licenses *restrict* the effect of copyright
- Events (e.g. the fact that an issue comment was created) are *not* copyrightable, but their content may be

## Copyright — GitHub's POV

- GitHub: *We claim no intellectual property rights over the material you provide to the Service.* (TOS F.1)
- Structure of API responses is GitHub's IP
- Several fields in API responses may contain copyrighted material

## Copyright situation example

```
← {
  "id": "4141500869",
  "type": "IssueCommentEvent",
  "actor": {},
  "repo": {},
  "payload": {
    "action": "created",
    "issue": {
      "id": 58442053,
      "number": 128,
      "title": "Issue in CopyrightedProjectName",
      "user": {},
      "labels": [],
      "state": "closed",
      "body": "Added data holding classes and a
map manager. Will add a system soon"
    },
    "comment": {
      "created_at": "2016-06-14T05:51:16Z",
      "updated_at": "2016-06-14T05:51:16Z",
      "body": "continuing in #141 \r\n"
    }
  }
}
→
```

©GitHub

©Project Name

©Issue initiator

©Issue comment

**Privacy** is the ability of an individual or group to seclude themselves, or information about themselves, and thereby express themselves selectively.



## Privacy provisions — EU

- *Personal data* identify a person uniquely
  - *Facts* are not personal data
- GHTorrent processes personal data, therefore is a *controller*
- Controllers must
  - get consent for processing (except in the case of *legitimate interest*)
  - include mechanisms for opting out

## Privacy provisions — USA

- No single law/directive
- Consent only required for specific types of data storage (e.g. social security numbers)
- Offering an opting out mechanism

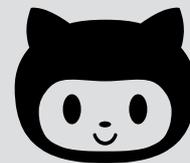
## What did GHTorrent do?

- Stopped distributing user names and emails in MySQL data dumps
- Researchers can “sign” a form to get access to private data
- Created an opt-out process
- In the process of creating Terms of Fair Use

## Research ethics

Can we, in the name of science,

- send emails to developers?
- create developer profiles?
- recommend work to developers?
- rank developers based on contributions?
- compare project characteristics?
- characterise community practices?



## Where is GHTorrent going?

## Challenges

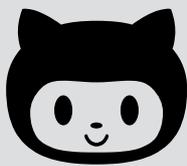
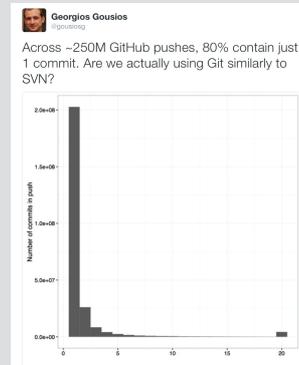
- 3Vs of big data: **V**olume, **V**elocity and **V**ariety
- Hardware not enough
- Software not optimal
- Services not always up an running

A screenshot of a GitHub repository report for Georgios Gousios. The report is titled "GHTorrent Report" and shows various statistics for the repository. The statistics are organized into sections: "Mongo", "Mongo", and "Ubuntu". Each section lists metrics such as commit\_comments, commits, events, followers, forks, issue\_comments, issue\_events, issues, org\_members, pull\_request\_comments, pull\_requests, repo\_collaborators, repo\_labels, repos, users, and watchers. The report is dated 21 Apr 2015 03:00. The screenshot also shows the GitHub logo and the repository name "Georgios Gousios" at the top.

# What can we do?

- Proper storage subsystem
  - MongoDB sharding? HDFS? Hyperdex?
- Proper querying infrastructure
  - Spark?
  - Datalake? ← Working with Microsoft
- Distributed operation and eventual consistency?

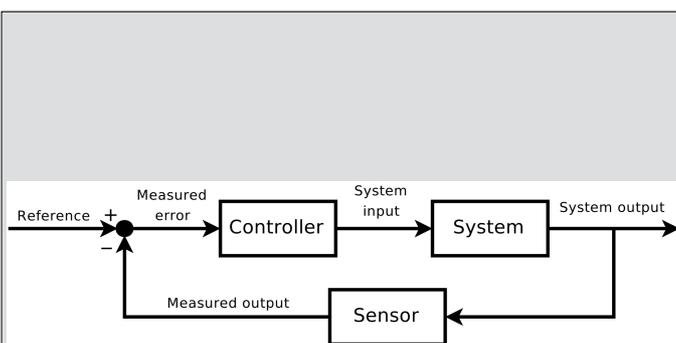
# DataLake: sneak preview



*A dataset accessible to anybody who wants to do research on extremely interesting big data*

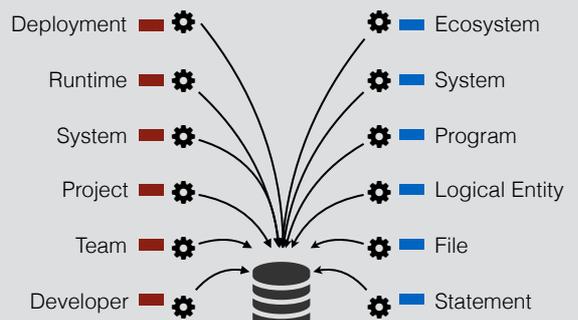


The future: Feedback Driven SE

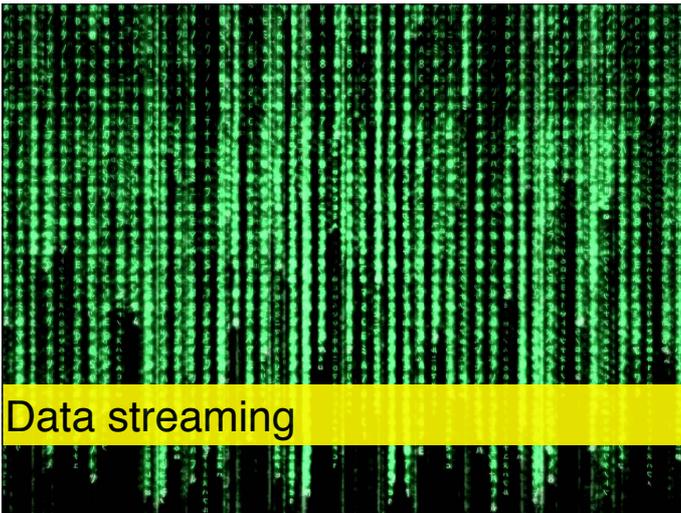
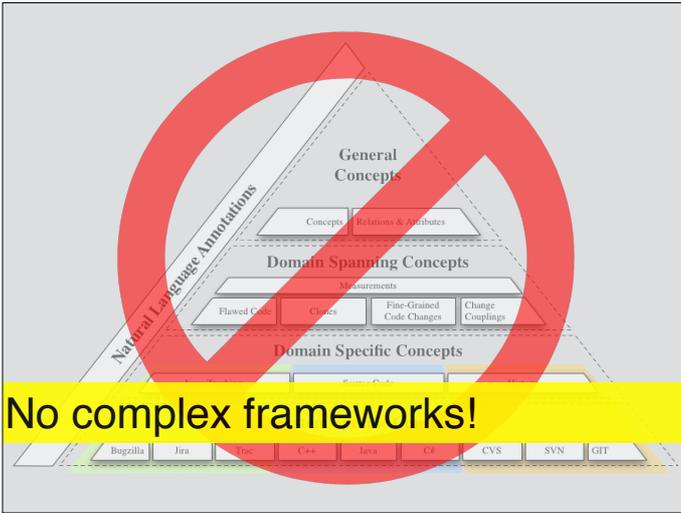


Feedback-driven development

## Processes

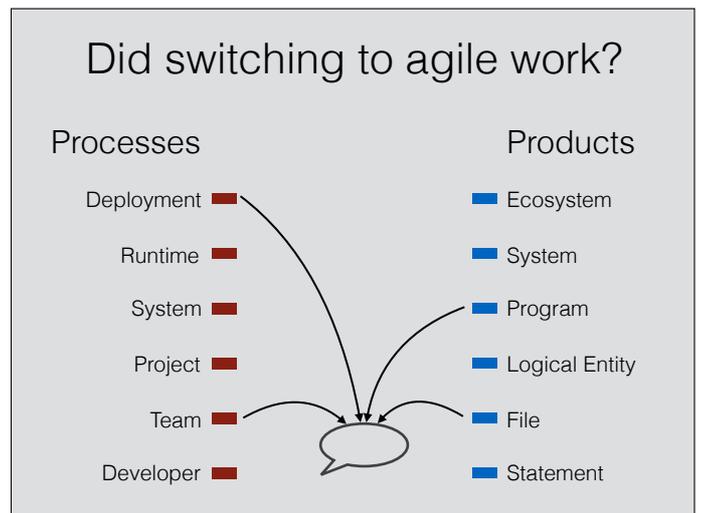
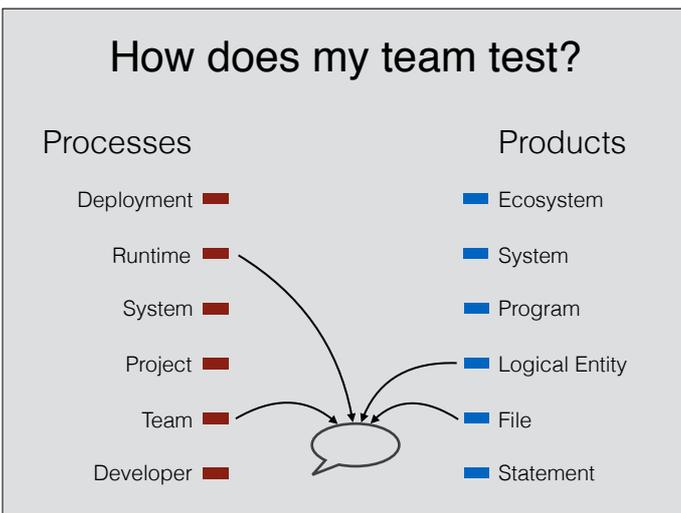


Data silos



```
Observable.from(stacktraces)
  .flatMap{s => findDev(s) match {
    case Some(s) => Observable.just(s)
    case None => Observable.error(new Exception("Not found"))
  }}
  .groupBy(x => x)
  .map { case (dev, xs) =>
    (dev, xs.scan(0)((count, _) => count + 1))
  }
  .flatMap { case (dev, xs) =>
    xs.map(c => (dev, c))
  }
```

**Queries on streams**



# How is my incident respon

## Processes

- Deployment
- Runtime
- System
- Project
- Team
- Developer

## Products

- Ecosystem
- System
- Program
- Logical Entity
- File
- Statement

