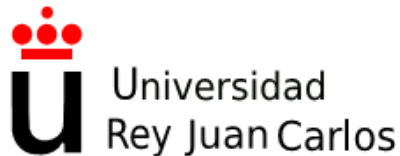# Doctoral thesis

## Temporal and behavioral patterns in the use of Wikipedia

**Author:**

Antonio J. Reinoso Peinado

Ingeniero en Informática

**Director:**

Jesús M. González Barahona

Doctor Ingeniero de Telecomunicación

GSyC

Universidad Rey Juan Carlos

# License

**http://gsyc.es/~ajreinoso/phdthesis**
**http://gsyc.es/~ajreinoso/phdthesis/slides.pdf**

**Antonio J. Reinoso**
**ajreinoso@libresoft.es**

**Temporal and Behavioral patterns in the use of Wikipedia**
**Doctoral thesis Móstoles (Madrid) September 2011**

# Summary

Introduction

Research objectives, Motivation and Contributions

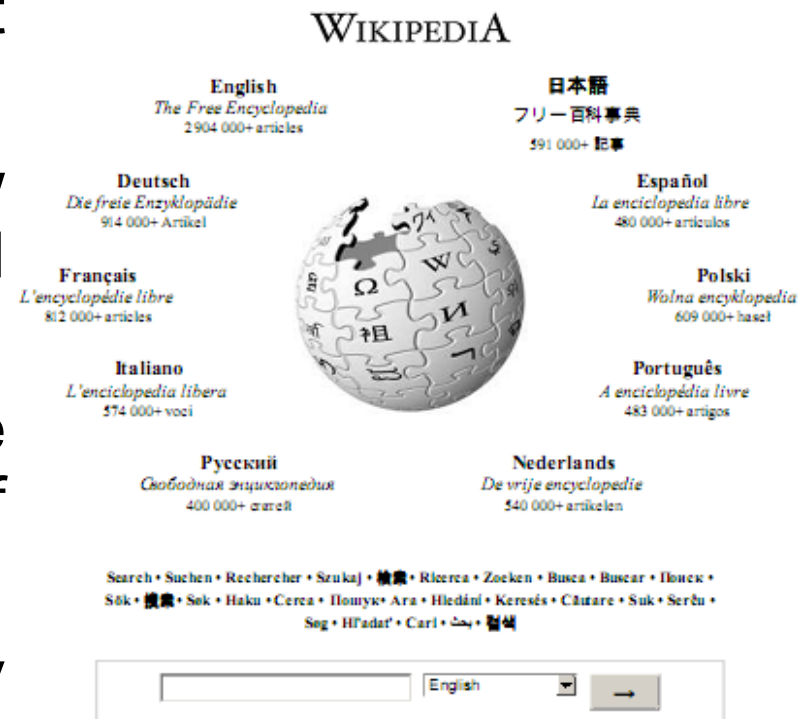State of the art

Methodology

Results

Further work

Questions

# Introduction. The Wikipedia project
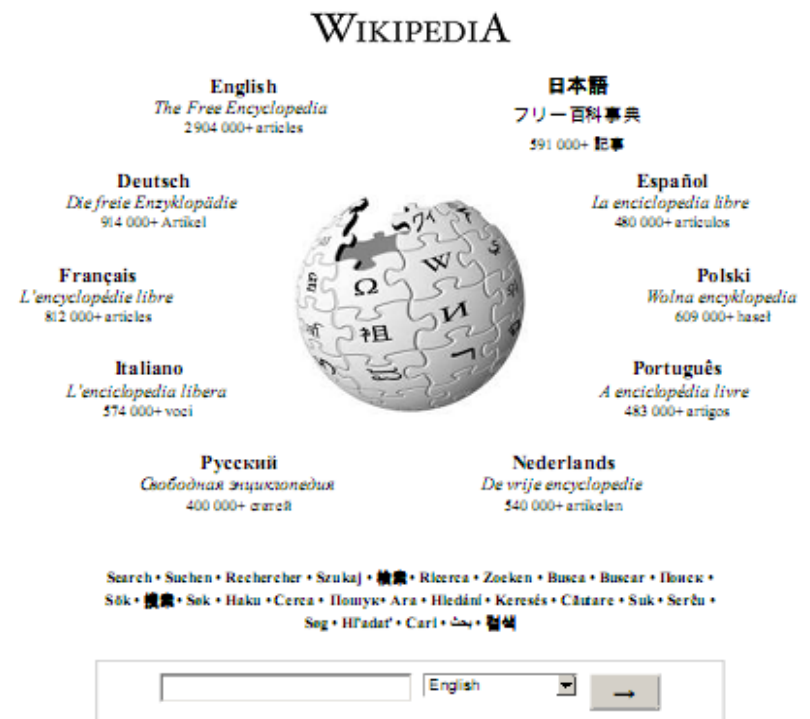
Some relevant features:

- Wikipedia is probably the most successful wiki-based platform.

- Wikipedia represents a new form of free management and distribution of knowledge.

- It is being built with the collaborative efforts of thousands of volunteers.

- Wikipedia is not supported by any well-known authority.

# Introduction. The Wikipedia project
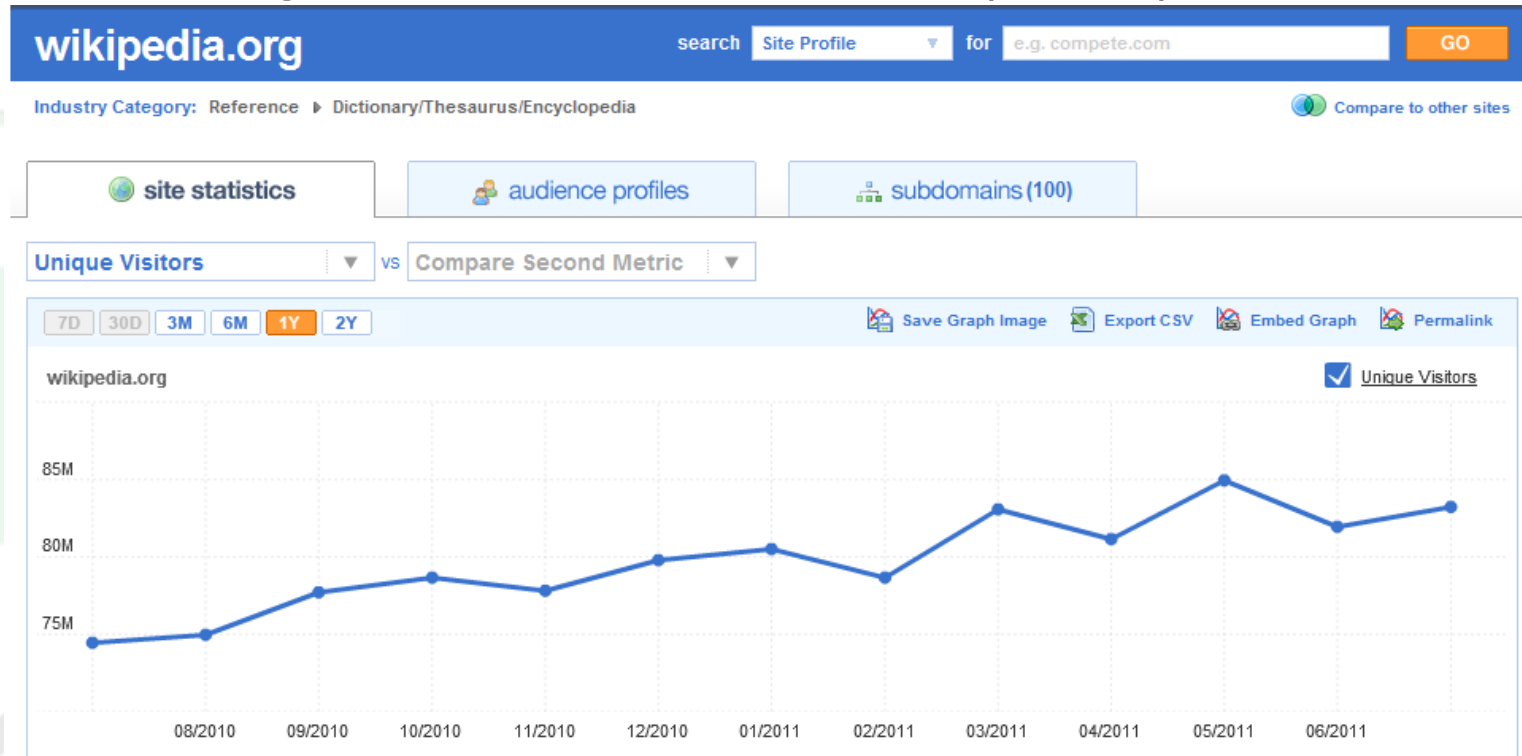
Wikipedia in figures:

- 282 editions (languages)

- By 19 million and a half articles

- By 83,000 active editors

- More than 30 million registered users

**Antonio J. Reinoso**
**ajreinoso@libresoft.es**

**Temporal and Behavioral patterns in the use of Wikipedia**
**Doctoral thesis Móstoles (Madrid) September 2011**

# Introduction. The Wikipedia project

The Wikipedia's audience places its web site within the top-seven most visited pages all over the Internet (Alexa)[1]



- Near 453 million pageviews a day (about 13,500 million a month)[2]
- By 85 million unique visitors in USA during May, 2011 (Compete)[3]

[1]http://www.alexa.com/siteinfo/wikipedia.org
[2]http://stats.wikimedia.org/EN/Sitemap.htm
[3]http://siteanalytics.compete.com/wikipedia.org/?metric=uv

Antonio J. Reinoso
ajreinoso@libresoft.es

Temporal and Behavioral patterns in the use of Wikipedia
Doctoral thesis Móstoles (Madrid) September 2011

# Introduction. The Wikipedia project.

This relevance has made Wikipedia to became a subject of research regarding...

- Its dimension as a mass phenomenon.

- The concern with the quality and reliability of its contents.

- The evolution and growth tendency of such kind of project.

## Research objectives

We are interested in the way in which users interact, communicate and make use of Wikipedia. So...

How do users utilize Wikipedia ?

What do users use Wikipedia for?

# Research objectives

In detail...

## How do people use Wikipedia ?

Traffic characterization: Types and frequencies of users' interactions with Wikipedia

Temporal patterns.

Behavioral patterns.

## What do people use Wikipedia for?

Contents that attract users' attention in visits and edits

Topics involved in search operations

# Motivation

Why ?

- Few studies addressing the use of Wikipedia by its users.

- Available data feed (Openness principles).

- Use of an innovative approach based on the analysis of the traffic (previous ones were based on dumps files).

- Validation: Possibility of comparing some of our results with similar ones from other sources.

- Sociological aspects derived from temporal and behavioral patterns.

- Easy reproducibility of the analysis.

# What for ?

So, what for ?

- Precise traffic characterization may lead to improvements in the systems in charge of managing users' requests.

- Time models may permit to forecast the evolution of users' requests.

- To quantify the degree of collaboration exhibited by some communities.

- To determine the origin of contributions to Wikipedia: elite of authors vs. general visitors.

# Research questions (I)

Focusing on raw traffic...

- Validation of log analysis: Can we trust the results obtained from the analysis of traffic containing users' requests?

- Detailed characterization of user's requests: Can we determine the exact composition of the traffic to Wikipedia?

- Traffic-size relationship: Is there any relationship between the traffic to each Wikipedia edition and its size?

# Research questions (II)

Focusing on users' requests...

- Temporal patterns: Do the different kinds of requests present cyclical evolutions over time (periodicity)?

- Behavioral aspects: Are the users' requests the reflect of different kind of behaviors when browsing Wikipedia?

- Degree of users' participation and reluctance: Can we establish the degree of participation and collaboration of the different communities of users?

# Research questions (III)

Focusing on contents...

- Audience of featured contents. Does the promotion of high quality articles to a featured status have an impact on the traffic they attract?

- Popular contents: What kind of contents are the most visited and contributed in each Wikipedia edition?

- Search operations : What are the topics more frequently involved in search operations? How do search operations influence visits related to the same contents?

# State of the art. Wikipedia research

From farther to closer approaches:

- Communities and generation of knowledge.

  - Surowiecki's "The wisdom of the crowds" [Sur04]

  - Stalder and Hirsh "Open intelligence" [SH02]

- The wikis and Wikipedia as research topics.

  - Wikis to involve users in the process of generation of knowledge. Ebersbach [EG04] [EG05]

  - Quality: Analysis of **credibility** by Korfiatis [KNP+06] and Chesney [Che06]. **Comparison** approach by Giles [Gil05] and Luyt [LKSY07]

  - Author reputation: Adler y Alfaro [Ada07] analyzed longevity of editions.

  - Evolution: Buriol's Wikigraph [BCD+06]

  - Featured articles: Viegas [VWM07]

  - Consensus and vandalism: Priedhorsky [PCS+07]

# State of the art. Wikipedia research

- Analysis based on logged information.

  - Web servers. Arlitt & Williamson [AW96] [AW97]

  - Squids: Khunkitti  [KI01]

- Use of wikis and Wikipedia

  **Academic research**

  - Wiki-based networks: SNA and DNA by Müller [MMB08]

  - Surveys on academic environments: Head and Heisenberg [HE10], Schweitzer [Sch08]

  - Contributions: long-tail dist. Kittur [KPSM07] Chi [Chi07]

  - Most popular topics: Spoerry [Spo07]

  - WMF servers' workload: Urdaneta [UPvS07b]

  - Authoring, coordination and survival of contributions: Ortega [OGB07] [Ort09]

## Non-Academic research

- Quantitative information: WMF itself, Mituzas's pageviews and Zachte's portal.

- Several visualizations of Mituza's logs but most of them unmaintained and not-updated.

- External sources: Alexa or ComScore.
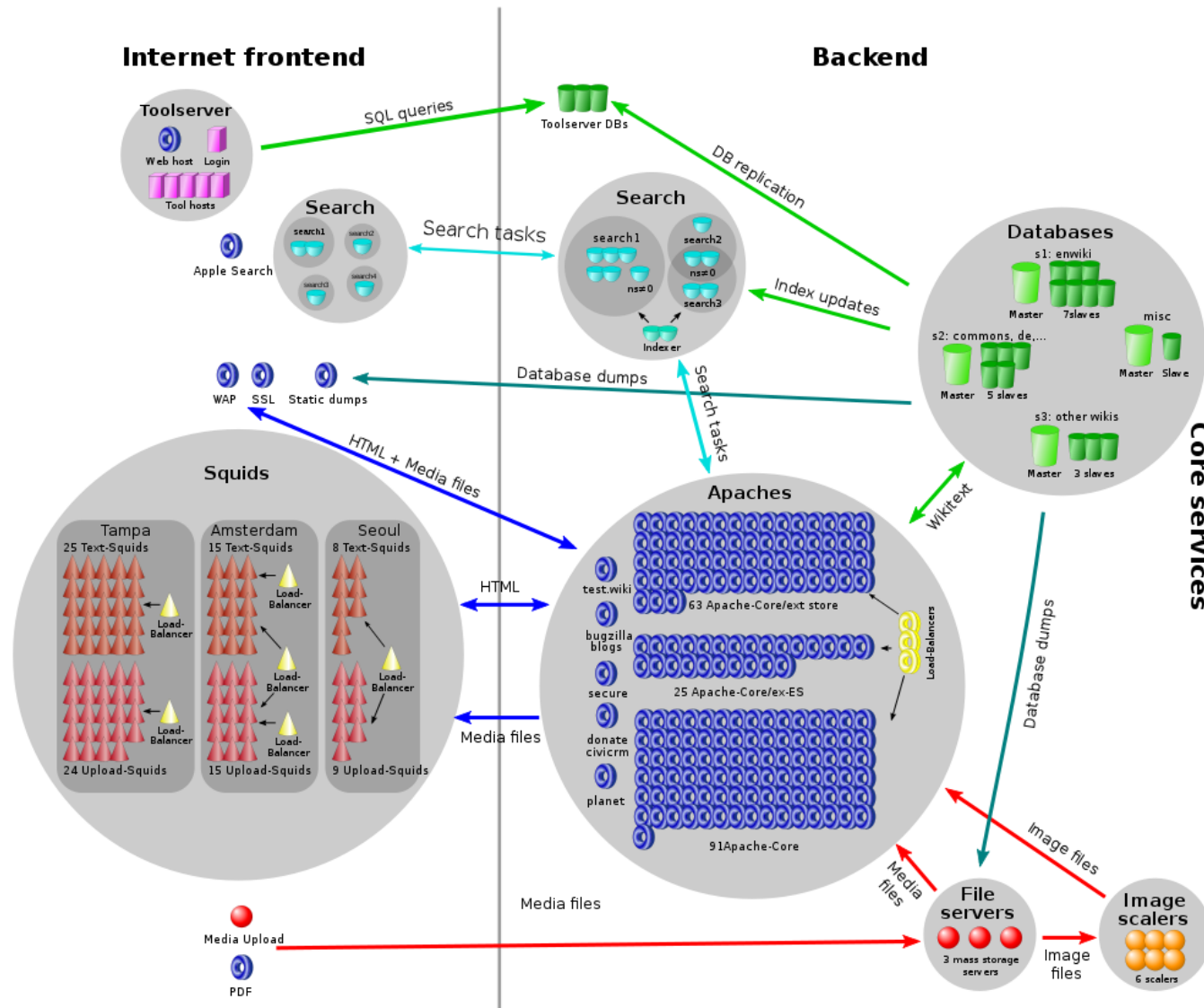
# Methodology. Introduction

Basic idea: To analyze the users' requests submitted to Wikipedia in order to:

- Perform a characterization of the whole traffic.

- Filter and store the information elements of the ones considered of interest for the analysis.

In the following:

- The Wikipedia hardware architecture.

- Description of the data feed.

- Processing users' requests:  The ad-hoc developed *WikiSquilter* application.

**Antonio J. Reinoso**
**ajreinoso@libresoft.es**

# Methodology. The Wikimedia Foundation server architecture

# Methodology. The Wikimedia Foundation server architecture

Squid servers:

- Usually work as **proxy servers** .

- May also work as **reverse proxies** caching contents previously requested to reply to new demands.

- Working in this way, Squids avoid the operation of both database and web servers placed behind them.

- WMF Squids deal with all the traffic to all its projects.

- After having sent a response, Squids register the corresponding user's URL demanding it.

- Every Squid send its lines to a central aggregator from where we receive them: **heterogeneity**.

Antonio J. Reinoso
ajreinoso@libresoft.es

# Methodology. The WMF Squid log format

Squids register different information about users' requests according to their log format:

| Field | Received |
|---|---|
| Squid Hostname | √ |
| Sequence number | |
| GMT Time | √ |
| Request service time (ms.) | √ |
| Reply size including HTTP headers | |
| Request method | √ |
| URL | √ |
| Squid hierarchy status | |
| MIME content type | |
| Referrer header | |
| User-Agent header | |

Sequence number: allows to detect packet losses.

Save field: Indicates whether a request caused a write operation.

Users' privacy is always preserved: logs are anonymized

Antonio J. Reinoso
ajreinoso@libresoft.es

# Methodology. The WMF Squid log format

```
May 6 13:46:04 208.80.152.138 22260437
```
**2010-05-06T13:42:43.827**
**http://en.wikipedia.org/wiki/April –** 2 GET


Most important fields

- Squid datetime

- URL specifying a user request.

- Field indicating a save operation (save) or a read one (-)


Received Squid lines are stored in a file which is daily rotated. In average, each of these files contain 40M lines and are about 900 MB.

# Methodology. The data feed in detail

Our data feed consists in Squid log lines corresponding to users' requests:

- 1/100 sample of the whole traffic directed to all the WFM projects during 2009 (15,000 million log lines)

Our analysis focuses on:

- The ten largest (articles and traffic volume) Wikipedia editions: DE, EN, ES, FR, IT, JA, NL, PL, PT, RU.

- Some specific **namespaces**: Main, Talk, User, User_Talk and Special.

- Most commonly requested **actions**: read, edit, submit, history, save and search.

# Methodology. The Wikipedia interface

Articles visited for reading are in the **Main** namespace:



http://en.wikipedia.org/wiki/Squid

Antonio J. Reinoso
ajreinoso@libresoft.es

Temporal and Behavioral patterns in the use of Wikipedia
Doctoral thesis Móstoles (Madrid) September 2011

# Methodology. The Wikipedia interface

The **Discussion** namespace gathers information devoted to improve the quality of the article or to broaden its contents.



Namespaces are translated into editions' languages:
http://**en**.wikipedia.org/wiki/**Talk:Squid**
http://**es**.wikipedia.org/wiki/**Discusi%C3%B3n**:**Teuthida**
http://**it**.wikipedia.org/wiki/**Discussione:Teuthida**
http://**ja**.wikipedia.org/wiki/**%E3%83%8E%E3%83%BC%E3%83%88:%E3%82%A4%E3%82%AB**

# Methodology. The Wikipedia interface

Edits:



http://en.wikipedia.org/w/index.php?title=**Squid**&action=**edit**
http://en.wikipedia.org/w/index.php?title=**Talk:Squid**&action=**edit**

**Antonio J. Reinoso**
**ajreinoso@libresoft.es**

# Methodology. The Wikipedia interface

Preview, changes and save:

```
| pages = 138-145
| doi = 10.1666/0022-3360(2006)080[0138:TCJFTU]2.0.CO;2
| postscript = <!--None-->
}}</ref>}}
| image = Sepioteuthis lessoniana (Bigfin reef squid).jpg
| image_caption = [[Bigfin Reef Squid]], ''Sepioteuthis lessoniana''
| authority = [[Adolf Naef|A. Naef]], 1916
| subdivision_ranks = [[Suborders]]
| subdivision =
†[[Plesioteuthididae]] <small>(''[[incertae sedis]]'')</small><br>
[[Myopsina]]<br>
[[Oegopsina]]
```

Content that violates any copyrights will be deleted. Encyclopedic content must be **verifiable**.

By clicking the "Save Page" button, you agree to the Terms of Use, and you irrevocably agree to release your contribution under the CC-BY-SA 3.0 License and the GFDL. You agree that a hyperlink or URL is sufficient attribution under the Creative Commons license.

Edit summary (Briefly describe the changes you have made)

[ Save page ]  [ Show preview ]  [ Show changes ]  Cancel | Editing help (opens in new window)

If you do not want your writing to be edited, used, and redistributed at will, then do not submit it here. All text that you did not write yourself, except brief excerpts, must be available under terms consistent with Wikipedia's Terms of Use before you submit it.

http://en.wikipedia.org/w/index.php?title=Squid&action=submit
http://en.wikipedia.org/w/index.php?title=Squid&action=submit **save**

# Methodology. The Wikipedia interface

## History review:



http://en.wikipedia.org/w/index.php?title=Squid&action=**history**

Antonio J. Reinoso
ajreinoso@libresoft.es

# Methodology. The Wikipedia interface

Search operations:



http://en.wikipedia.org/wiki/Special:Search?search=Linux
http://**en**.wikipedia.org/w/index.php?title=**Special**%3ASearch&search=Linux
http://**es**.wikipedia.org/w/index.php?title=**Especial**%3ASearch&search=Linux

Antonio J. Reinoso
ajreinoso@libresoft.es

# Methodology. The WikiSquilter application

The WikiSquilter project is the software tool developed to parse and filter the information from the Squid log lines.

- Tailored Java written application.
    - Multithreaded capabilities.
    - Good performance database drivers.
- Three basic functionalities:
    - **Parsing**: The application parses the information elements from the log lines.
    - **Filtering**: Determination of the information elements considered of interest.
    - **Storage**: Filtered elements are normalized and stored in a MySQL relational database.

# Methodology. The WikiSquilter application

- Strong adherence to SE principles:

  - **Robustness:** 15,000 million lines successfully processed.

  - **Extensibility**: Highly modular and coupling reduced to the minimum.

  - **Efficiency**: Multithreaded approach, filter mechanism based in a hash structure, and

  - **Flexibility**: Parameters of the analysis fully configurable. The structure acting as the filter is built upon the specifications of an XML file.

# Methodology. The WikiSquilter application

```xml
<filter_cfg>
    <WikiMediaProject dbCode="0" name="WIKIPEDIA">
  <NNSS_INDEXES>
    <NSINDEX>ARTICLE</NSINDEX>
    <NSINDEX>INDEX</NSINDEX>
    <NSINDEX>ARTICLE_TALK</NSINDEX>
    <NSINDEX>USER</NSINDEX>
    <NSINDEX>USER_TALK</NSINDEX>
    <NSINDEX>SPECIAL</NSINDEX>
  </NNSS_INDEXES>
  <Language dbCode="EN" name="ENGLISH">
    <NameSpaces><NS>Talk</NS><NS>User</NS><NS>User_Talk</NS><NS>Special</NS></NameSpace
  </Language>
  <Language dbCode="DE" name="GERMAN">
    <NameSpaces>
      <NS>Diskussion</NS><NS>Benutzer</NS><NS>Benutzer_Diskussion</NS><NS>Spezial</NS>
    </NameSpaces>
  </Language>
  <Language dbCode="JA" name="JAPANESE">
    <NameSpaces>
      <NS>%E3%83%8E%E3%83%BC%E3%83%88</NS>   <NS>%E5%88%A9%E7%94%A8%E8%80%85</NS>
      <NS>%E5%88%A9%E7%94%A8%E8%80%85%E2%80%90%E4%BC%9A%E8%A9%B1</NS>
      <NS>%E7%89%B9%E5%88%A5</NS>
    </NameSpaces>
  </Language>
  <Actions> <Action>edit</Action>  <Action>save</Action> </Actions>
  <Methods>  <Method>GET</Method>   <Method>POST</Method> </Methods>
</filter_cfg>
```

**Antonio J. Reinoso**
**ajreinoso@libresoft.es**

## Parsing

```
May 6 13:46:04 208.80.152.138 22260437 2010-
05-06T13:42:43.827
http://en.wikipedia.org/wiki/April - 2 GET
```

The application parser an analyzes each **log line** to extract:

- The Squid date and time.

- The URL as a block.

- Whether the URL caused a save operation.

- The response time.

- The  request method.

Antonio J. Reinoso
ajreinoso@libresoft.es

# Methodology. The WikiSquilter application

**Parsing**

`http://en.wikipedia.org/wiki/April`

`http://en.wikipedia.org/wiki/Talk:April`

The application tokenizes the **URL** to determine:

- The Wikimedia Foundation Project.
- The language edition.
- The targeted namespace.
- The article's title.

# Methodology. The WikiSquilter application

## Parsing

Requests for actions are a bit difficult to scan

```
http://en.wikipedia.org/w/index.php?
title=London\&action=history
```

```
http://de.wikipedia.org/w/index.php?
title=Diskussion:Berlin\&action=edit
```

```
http://it.wikipedia.org/w/index.php?
title=Utente\%3AAjreinoso\&action=history
```

URLs specifying actions are first assigned to a fictitious *Index* namespace.

If language, project and action are of interest, the URL is re-parsed to determine the article's title and namespace.

## Overall performance

After all our efforts:

- The traffic corresponding to a whole month is processed in 1 day and six hours in a **quad-core system** running under **Ubuntu server** and equipped with **8 GB of RAM** memory.

- Such traffic consists in about **1,300 million** log lines.

- **80%** of the total time corresponds to the indization of the tables.

- In average, the application analyzes **60,000 log lines per second**.

# Methodology. Answering the research questions

**Once data are available in the database, we undertake...**

- The **validation of part of our results** using several sources (WMF info, Ortega's WikiXRay and Alexa).

- The finding of **temporal patterns** and periodicity (autocorrelation and cross-correlation).

- The study of the **behavioral habits** (participation, reluctance...) exhibited through the requests.

- The attention attracted by **featured contents**. Two perspectives: promotion and inclusion in main pages (Comparative tests such as Wilcoxon rank-sum test).

- The **most visited, contributed and searched topics**. (Grouping by md5 hash of articles' titles and searched strings followed by manual classification based on Spoerry's one)

- **Traffic characterization** has been already performed!!

# Results. Validation.

As some sources are not sampled, our results should maintain a ratio with them similar to the sampling factor (1%).

**Visits**

| Lang. | Jan. | Feb. | Mar. | Apr. | May. | Jun. |
|-------|------|------|------|------|------|------|
| DE (Reinoso) | 10,821,625 | 6,833,171 | 8,034,636 | 6,945,878 | 7,612,949 | 7,249,244 |
| DE (Mituzas) | 1,271 M | 982 M | 978 M | 817 M | 875 M | 909 M |
| Ratio | 0.009 | 0.007 | 0.008 | 0.009 | 0.009 | 0.008 |
| EN (Reinoso) | 47,369,841 | 43,136,627 | 51,845,199 | 48,242,580 | 48,085,156 | 43,950,168 |
| EN (Mituzas) | 5,615 M | 5,944 M | 6,092 M | 5,989 M | 6,066 M | 5,819 M |
| Ratio | 0.0084 | 0.0073 | 0.0085 | 0.0081 | 0.0079 | 0.0076 |

**Edits**

| Lang. | Jan. | Feb. | Mar. | Apr. | May. | Jun. |
|-------|------|------|------|------|------|------|
| DE (Reinoso) | 11,041 | 9,457 | 10,341 | 8,361 | 8,052 | 7,754 |
| DE (Zachte) | 876 K | 752 K | 802 K | 655 K | 684 K | 701 K |
| DE (Ratio) | 0.0126 | 0.0126 | 0.0129 | 0.0128 | 0.0118 | 0.0111 |
| EN (Reinoso) | 53,121 | 46,778 | 54,564 | 47,921 | 47,692 | 42,282 |
| EN (Zachte) | 4,300 K | 4,200 K | 4,400 K | 4,000 K | 4,300 K | 4,000 K |
| EN (Ratio) | 0.0124 | 0.0111 | 0.0124 | 0.0120 | 0.0111 | 0.0106 |

# Results. Validation.

Results also match at a much finer grain (level of articles)

# Results. Representativeness.

Do the considered elements correspond to a significant part of the traffic to Wikipedia?
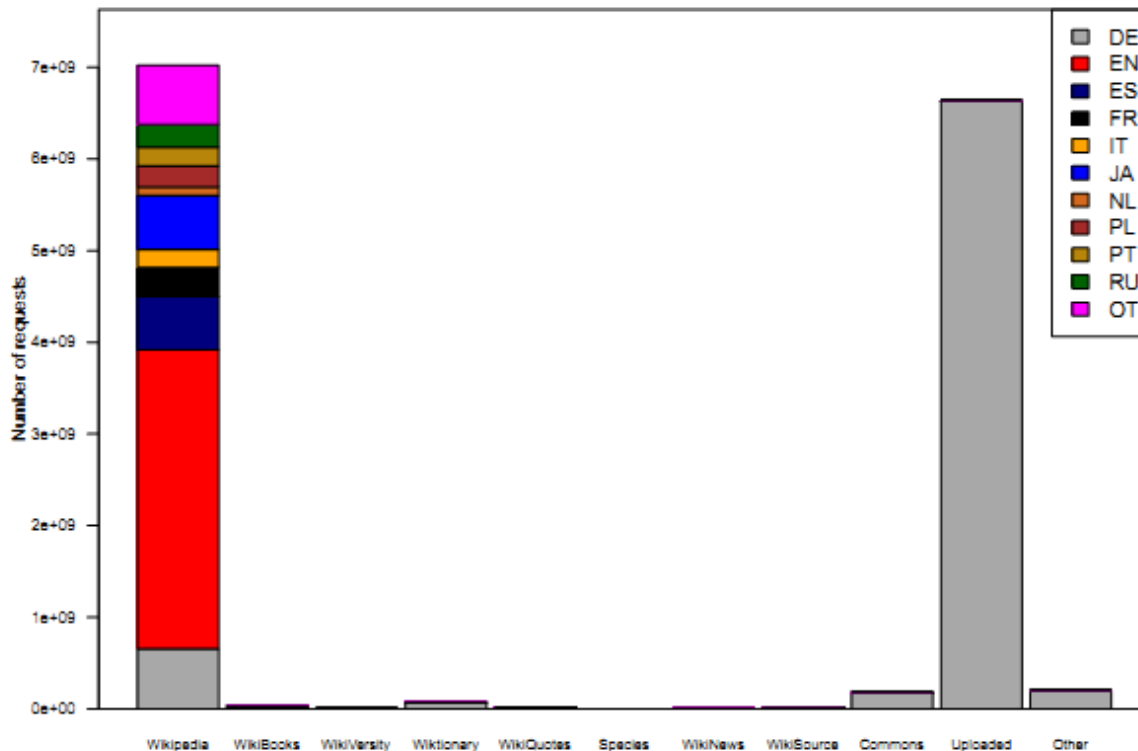
- **90%:** Traffic attracted by our considered Wikipedias.

- **85%:** Requests directed to our namespaces.

- **94%:** Filtered edit requests from the total of save op.

- **99%:** Search operations which have always to be filtered.

## Traffic characterization.

Traffic is computed in terms of number of requests disregarding amount of information or transfer rates.

- Requests to Wikipedia and to previous uploaded images and media constitute the 96% of the traffic.



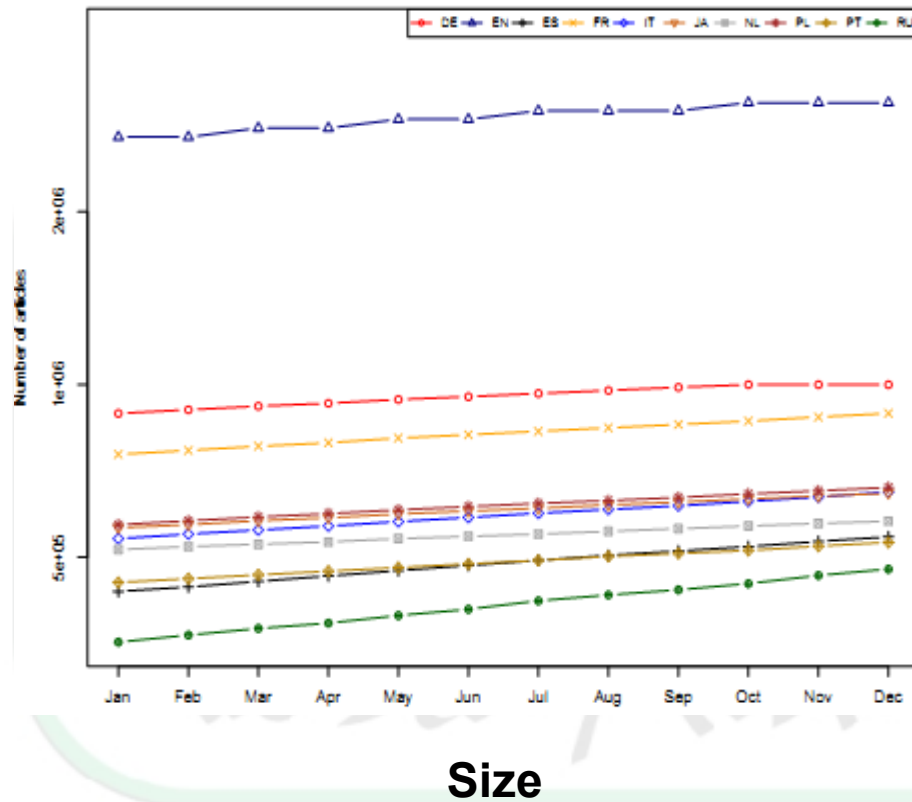| Wikipedia edition | Percentage |
|---|---|
| DE | 9.40% |
| EN | 46.45% |
| ES | 8.25% |
| FR | 4.54% |
| IT | 2.79% |
| JA | 8.38% |
| NL | 1.34% |
| PL | 3.30% |
| PT | 2.87% |
| RU | 3.51% |
| REST | 9.17% |

# Results. Traffic characterization

Percentage of the traffic corresponding to each type of request.

| Ed. | Visits to articles | Actions (except search) | Edit op. | Search op. | Api calls | Skins /css | icons | mw ext. | Undet. |
|---|---|---|---|---|---|---|---|---|---|
| EN | 21.51% | 22.52% | 0.27% | 4.75% | 6.53% | 34.62% | 4.38% | 3.47% | 6.95% |
| DE | 16.54% | 20.87% | 0.23% | 4.09% | 7.69% | 30.74% | 3.46% | 14.72% | 5.98% |
| ES | 13.58% | 33.90% | 0.31% | 4.12% | 6.02% | 32.13% | 3.68% | 3.89% | 6.80% |
| FR | 18.24% | 23.15% | 0.33% | 4.00% | 6.05% | 36.87% | 4.42% | 4.23% | 7.04% |
| IT | 19.80% | 21.81% | 0.43% | 4.44% | 5.77% | 37.57% | 4.49% | 3.07% | 9.69% |
| JA | 20.69% | 25.15% | 0.37% | 4.22% | 3.95% | 36.01% | 4.19% | 2.81% | 9.22% |

Antonio J. Reinoso
ajreinoso@libresoft.es

# Results. Traffic characterization

Relationship between traffic and editions' sizes.



**Size**



**Traffic**

The Spanish and Russian Wikipedias are the smallest regarding their number of articles but attract much more traffic than other larger editions.

Antonio J. Reinoso
ajreinoso@libresoft.es

Comparison among the traffic composed by our filtered requests, the traffic to all the WMF projects, and to Wikipedia.
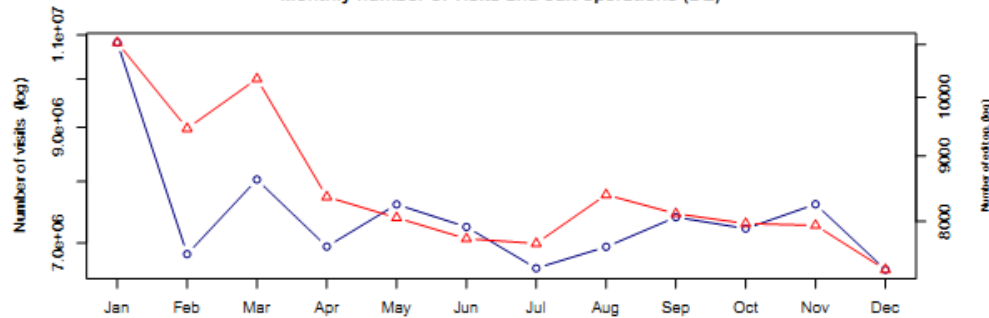


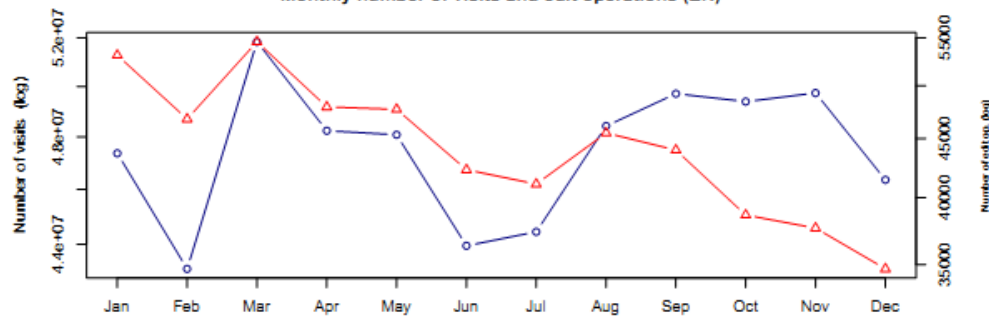Wikipedia traffic is positively correlated to the traffic to all the Wikimedia Foundation projects.

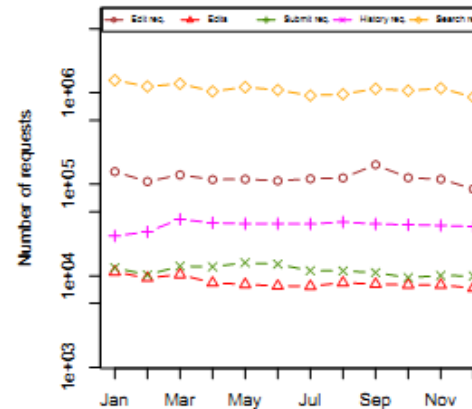# Results. Temporal patterns

## Monthly evolution.

# Results. Temporal patterns

## Weekly evolution.



Number of daily requests of each type during every whole week of 2009 (ES)

- Periodicity in visits, searches and requests to edit.
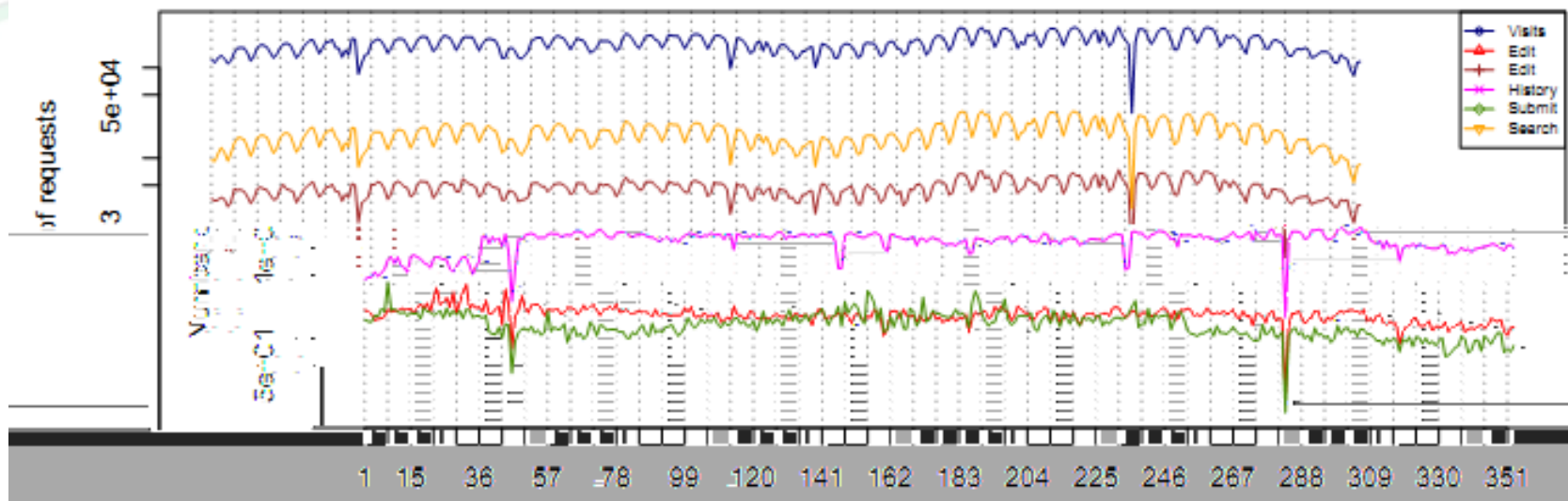- It is very difficult tu pronounce about the rest of actions: Atypical character and too small number of requests.

Antonio J. Reinoso
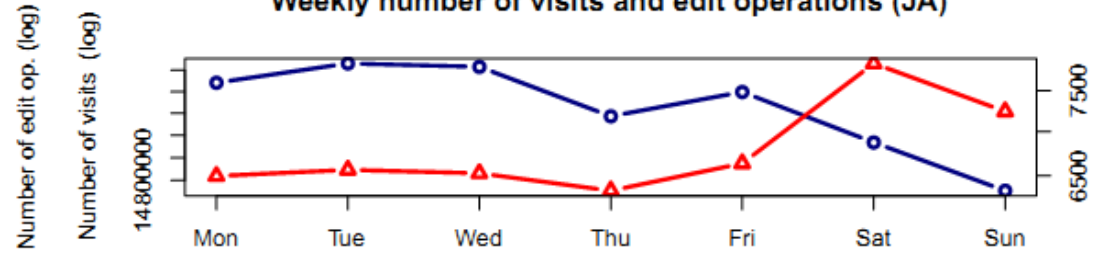ajreinoso@libresoft.es

# Results. Temporal patterns

Evolution of visits and edits throughout the days of the week.



- Close evolutions in the DE, EN, ES, IT and RU editions.
- Edits raise in weekends in FR, JA, NL and PL. Elite of authors?

# Results. Behavioral patterns

Relationships among the different types of requests:

- Positive correlation of **visits** and **edits** throughout the days of the week in DE, EN, ES, IT and RU editions. Negative in the case of Dutch and Japanese editions.

- Positive correlation of **edits** and **req. to edit** also in DE, EN, ES, IT and RU editions: Massive collaboration.



Corr. of visits and edits for the days of the week (DE)

Pearson cor.= 0.900001476433106  p=  0.005751

# Results. Behavioral patterns

Analysis of the **edits/visits ratio** for the different Wikipedias:

- A good indicator of proactivity and participation.



- Three groups: G1(**NL**, PL, **IT**, FR, RU), G2(ES, PT, EN, DE) and G3(JA).

- Both types of Wikipedias (having or not an elite of users) present high ratios of edits to visits.

# Results. Behavioral patterns

Analysis of the **performed edits/requested edits ratio:** Users' reluctance:

| Edition | Edits | Edit requests | Percentage of finished edits |
|---|---|---|---|
| IT | 57447 | 632295 | 9.09% |
| FR | 76377 | 941017 | 8.12% |
| NL | 29799 | 379450 | 7.85% |
| PL | 31199 | 419411 | 7.44% |
| RU | 60516 | 814103 | 7.43% |
| DE | 102442 | 1426027 | 7.18% |
| EN | 533879 | 8026886 | 6.65% |
| PT | 28469 | 584498 | 4.87% |
| ES | 66547 | 1666890 | 3.99% |
| JA | 47546 | 2079305 | 2.29% |

**Wikipedias having the highest ratios of edits to visits are also the ones having highest percentages of finished edits.**

Significant increase in the number of visits to the "Today's featured articles" during the month of their presentation for all the considered editions **except the Spanish one**.

Boxplots picturing the visits to the featured articles just promoted in the same periods.



Different patterns of visits as a result of the different dynamics in the promotion processes.
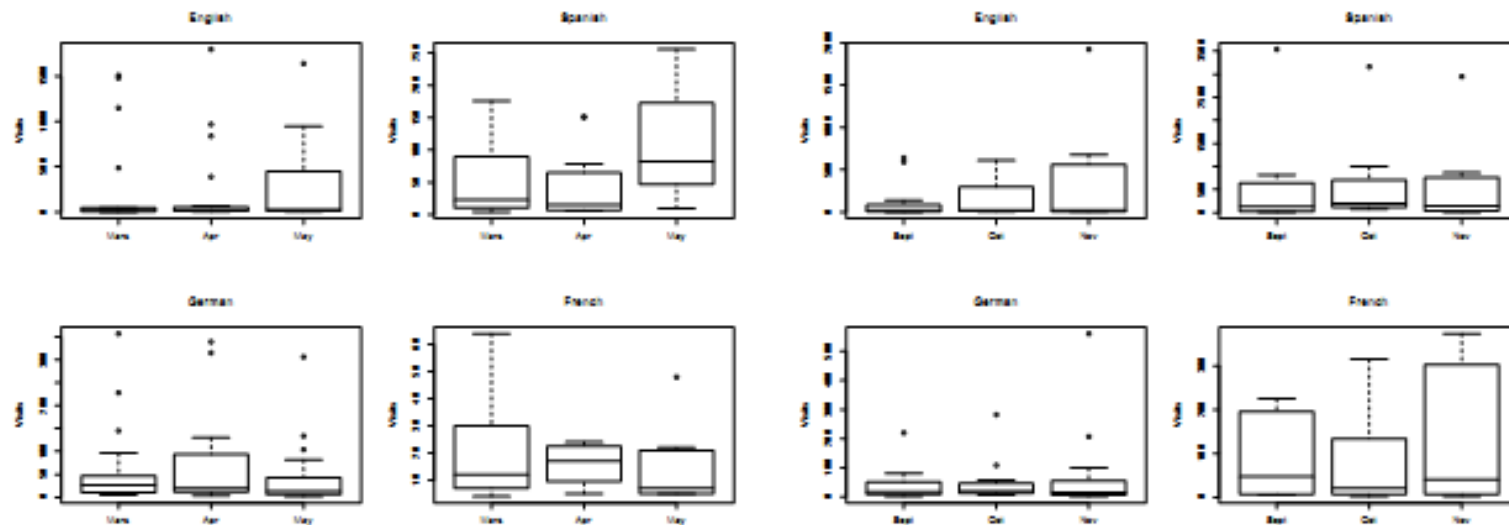
# Results. Most popular contents

Classification of contents most visited and edited in several Wikipedias:

| Category | DE (Visited) | DE (Edited) | EN (Visited) | EN (Edited) | ES (Visited) | ES (Edited) | FR (Visited) | FR (Edited) |
|---|---|---|---|---|---|---|---|---|
| MAIN | 47.28% | 0.00% | 74.05% | 0.00% | 7.41% | 0.00% | 57.77% | 0.00% |
| CUR | 5.53% | 20.27% | 6.18% | 28.30% | 7.76% | 5.94% | 8.18% | 11.58% |
| GEO | 11.60% | 14.40% | 1.55% | 11.16% | 11.66% | 18.47% | 9.51% | 24.73% |
| ICT | 5.97% | 7.64% | 2.26% | 2.27% | 10.66% | 1.17% | 2.79% | 0.58% |
| ENT | 16.64% | 16.17% | 10.92% | 31.63% | 14.48% | 50.53% | 9.00% | 23.74% |
| POL | 5.25% | 13.18% | 2.36% | 10.37% | 4.31% | 4.88% | 2.15% | 6.29% |
| SCI | 2.97% | 6.42% | 0.95% | 1.36% | 22.72% | 6.16% | 1.72% | 3.72% |
| ART | 2.25% | 17.17% | 0.16% | 12.33% | 17.70% | 12.10% | 4.63% | 28.21% |
| SEX | 2.50% | 0.22% | 1.47% | 0.00% | 0.18% | 0.00% | 0.61% | 0.25% |
| UNDETERMINED | 0.00% | 4.54% | 0.09% | 2.57% | 3.13% | 0.74% | 3.63% | 0.91% |

# Results. Future work.

The work developed in this thesis can be extended in serveral ways.

1. Visits and edit distributions.

2. Study of time series.

3. Geolocation.

4. Consensus process.

5. Access through different interfaces and devices.

6. Automatic categorization.

# Results. Related publications.

- **Temporal characterization of the requests to Wikipedia.**
  5th International Workshop on new Challenges in Distributed Information Filtering and Retrieval (DART'11)

- **A quantitative examination of the impact of featured articles in Wikipedia.**

  International Conference on Software and Data Technologies (ICSOFT'11)

- **A statistical approach to the impact of featured articles in Wikipedia.**

  International Conference on Knowledge Engineering and Ontology Development (KEOD'10)

- **A quantitative approach to the use of  Wikipedia.**

  IEEE Symposium on Computers and Communications (ISCC'09)

- **Quantitative analysis and characterization of Wikipedia requests.**

  ACM WikiSym 2008: 4th International Symposium on Wikis (WikiSym'08)

- **Workshop on interdisciplinary research on Wikipedia and Wiki communities.**

  ACM WikiSym 2008: 4th International Symposium on Wikis (WikiSym'08)

# Results. Related publications.



## Wikipedians' weekends in international comparison

A paper titled "Temporal characterization of the requests to Wikipedia" examined how search requests, read accesses and edits on Wikipedia change over time, and relate to those at the entirety of Wikimedia sites (based on squid logs for the whole year of 2009, provided by the Wikimedia Foundation). Among findings are differences between language versions of Wikipedia, such as that the "the number of edits tends to raise in weekends" for the French, Japanese, Dutch and Polish Wikipedia, but not for other languages. Another paper, titled "Circadian patterns of Wikipedia editorial activity: A demographic analysis"[9], similarly analyzed "34 Wikipedias in different languages [trying] to characterize and find the universalities and differences in temporal activity patterns of editors", with the underlying data provided by the German Wikimedia chapter from the toolserver. They found that "in contrast to diurnal [daily] pattern, which is universal to a great extent, weekly activity patterns of WPs show remarkable differences. We could, however, identify two main categories, namely 'weekends' and 'working days' active WPs."[10]

Any questions...?