

CEPIS UPGRADE is the European Journal for the Informatics Professional, published bi-monthly at <<http://cepis.org/upgrade>>

Publisher

CEPIS UPGRADE is published by CEPIS (Council of European Professional Informatics Societies, <<http://www.cepis.org/>>), in cooperation with the Spanish CEPIS society ATI (*Asociación de Técnicos de Informática*, <<http://www.ati.es/>>) and its journal *Novática*

CEPIS UPGRADE monographs are published jointly with *Novática*, that publishes them in Spanish (full version printed; summary, abstracts and some articles online)

CEPIS UPGRADE was created in October 2000 by CEPIS and was first published by *Novática* and **INFORMATIK/INFORMATIQUE**, bimonthly journal of **SVI/FSI** (Swiss Federation of Professional Informatics Societies)

CEPIS UPGRADE is the anchor point for UPENET (UPGRADE European NETWORK), the network of CEPIS member societies' publications, that currently includes the following ones:

- **inforeview**, magazine from the Serbian CEPIS society JISA
- **Informatica**, journal from the Slovenian CEPIS society SDI
- **Informatik-Spektrum**, journal published by Springer Verlag on behalf of the CEPIS societies GI, Germany, and SI, Switzerland
- **ITNOW**, magazine published by Oxford University Press on behalf of the British CEPIS society BCS
- **Mondo Digitale**, digital journal from the Italian CEPIS society AICA
- **Novática**, journal from the Spanish CEPIS society ATI
- **OCG Journal**, journal from the Austrian CEPIS society OCG
- **Pliroforiki**, journal from the Cyprus CEPIS society CCS
- **Tölvumál**, journal from the Icelandic CEPIS society ISIP

Editorial Team

Chief Editor: Llorenç Pagés-Casas

Deputy Chief Editor: Rafael Fernández Calvo

Associate Editor: Fiona Fanning

Editorial Board

Prof. Vasile Baltac, CEPIS President

Prof. Wolfried Stucky, CEPIS Former President

Prof. Nello Scarabottolo, CEPIS President Elect

Luis Fernández-Sanz, ATI (Spain)

Llorenç Pagés-Casas, ATI (Spain)

François Louis Nicolet, SI (Switzerland)

Roberto Carniel, ALSI – Tecnoteca (Italy)

UPENET Advisory Board

Dubravka Dukic (inforeview, Serbia)

Matjaz Gams (Informatica, Slovenia)

Hermann Engesser (Informatik-Spektrum, Germany and Switzerland)

Brian Runciman (ITNOW, United Kingdom)

Franco Filippazzi (Mondo Digitale, Italy)

Llorenç Pagés-Casas (Novática, Spain)

Veith Risak (OCG Journal, Austria)

Panicos Masouras (Pliroforiki, Cyprus)

Thorvardur Kári Ólafsson (Tölvumál, Iceland)

Rafael Fernández Calvo (Coordination)

English Language Editors: Mike Andersson, David Cash, Arthur Cook, Tracey Darch, Laura Davies, Nick Dunn, Rodney Fennemore, Hilary Green, Roger Harris, Jim Holder, Pat Moody.

Cover page designed by Concha Arias-Pérez

"Upcoming Resolution" / © ATI 2011

Layout Design: François Louis Nicolet

Composition: Jorge Lácer-Gil de Rames

Editorial correspondence: Llorenç Pagés-Casas <pages@ati.es>

Advertising correspondence: <info@cepis.org>

Subscriptions

If you wish to subscribe to CEPIS UPGRADE please send an email to info@cepis.org with 'Subscribe to UPGRADE' as the subject of the email or follow the link 'Subscribe to UPGRADE' at <<http://www.cepis.org/upgrade>>

Copyright

© Novática 2011 (for the monograph)

© CEPIS 2011 (for the sections Editorial, UPENET and CEPIS News)

All rights reserved under otherwise stated. Abstracting is permitted with credit to the source. For copying, reprint, or republication permission, contact the Editorial Team

The opinions expressed by the authors are their exclusive responsibility

ISSN 1684-5285

Monograph of next issue (October 2011)

"Green ICT"

(The full schedule of CEPIS UPGRADE is available at our website)



The European Journal for the Informatics Professional

<http://cepis.org/upgrade>

Vol. XII, issue No. 3, July 2011

Monograph

Business Intelligence

(published jointly with *Novática**)

Guest Editors: *Jorge Fernández-González and Mouhib Alnoukari*

- 2 Presentation. Business Intelligence: Improving Decision-Making in Organizations — *Jorge Fernández-González and Mouhib Alnoukari*
- 4 Business Information Visualization — *Josep-Lluís Cano-Giner*
- 14 BI Usability: Evolution and Tendencies — *R. Dario Bernabeu and Mariano A. García-Mattío*
- 20 Towards Business Intelligence Maturity — *Paul Hawking*
- 29 Business Intelligence Solutions: Choosing the Best solution for your Organization — *Mahmoud Alnahlawi*
- 38 Strategic Business Intelligence for NGOs — *Diego Arenas-Contreras*
- 43 Data Governance, what? how? why? — *Óscar Alonso-Llombart*
- 49 Designing Data Integration: The ETL Pattern Approach — *Veit Köppen, Björn Brüggemann, and Bettina Berendt*
- 56 Business Intelligence and Agile Methodologies for Knowledge-Based Organizations: Cross-Disciplinary Applications — *Mouhib Alnoukari*
- 60 Social Networks for Business Intelligence — *Marie-Aude Aufaure and Etienne Cuvelier*

UPENET (UPGRADE European NETWORK)

67 From **Novática** (ATI, Spain)

Free Software

AVBOT: Detecting and fixing Vandalism in Wikipedia — *Emilio-José Rodríguez-Posada* — Winner of the 5th Edition of the *Novática* Award

71 From **Pliroforiki** (CCS, Cyprus)

Enterprise Information Systems

Critical Success Factors for the Implementation of an Enterprise Resource Planning System — *Kyriaki Georgiou and Kyriakos E. Georgiou*

CEPIS NEWS

77 Selected CEPIS News — *Fiona Fanning*

* This monograph will be also published in Spanish (full version printed; summary, abstracts, and some articles online) by *Novática*, journal of the Spanish CEPIS society ATI (*Asociación de Técnicos de Informática*) at <<http://www.ati.es/novatica/>>.

Free Software

AVBOT: Detecting and fixing Vandalism in Wikipedia¹

Emilio-José Rodríguez-Posada

© Novática, 2010

This paper was published, in Spanish, by **Novática**, issue no. 203, January-February 2010. **Novática**, <<http://www.ati.es/novatica>>, a founding member of **UPENET**, is a bimonthly journal published by the Spanish CEPIS society ATI (*Asociación de Técnicos de Informática* – Association of Computer Professionals).

Wikipedia is a project which aims to build a free encyclopaedia to spread the sum of all knowledge to every single human being. Today it can be said to be on the road to achieving that goal, having reached the 15 million articles milestone in 270 languages. Furthermore, if we include its sister projects (Wiktionary, Wikibooks, Wikisource, ...), it has received more than 1 billion edits in 10 years and now has more than 10 billion million page views every month. Compiling an encyclopaedia in a collaborative way has been possible thanks to MediaWiki software. It allows everybody to modify the content available on the site easily. But a problem emerges regarding this model: not all edits are made in good faith. AVBOT is a bot for protecting the Spanish Wikipedia against some undesired modifications known as vandalism. Although AVBOT was developed for Wikipedia, it can be used on any MediaWiki website. It is developed in Python and is free software. In the 2 years it has been in operation it has reverted more than 200,000 vandalism edits, while several clones have been executed, adding thousands of reverts to this count.

Keywords: AVBOT, Bot, Free Software, Mediawiki, Monitoring, Vandalism, Wikipedia, Wikis.

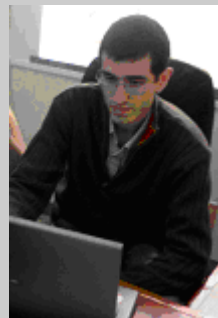
1 Introduction

Wikipedia [5] is the free encyclopaedia that anyone can edit. Any person using a web browser can modify its contents. This is made possible by MediaWiki [6], the software used by Wikipedia. The popularity of this free encyclopaedia has grown exponentially in its first 10 years, until it is now in the top 10 of most visited sites on the Internet, it appears in the first results of Google, and it has reached the milestone of 1 billion contributions [13]. Furthermore, when Microsoft closed Encarta in 2009, in a clear allusion to Wikipedia they said that "*the category of traditional encyclopaedias and reference material has changed*".

¹ This article was the winner of the 5th Edition of the Novática Award presented to the best article published in 2010 by Novática, journal of the Spanish society ATI. For details, in Spanish, about this edition of the Novática Award please visit <<http://www.ati.es/premio-novatica>>.

Author

Emilio-José Rodríguez-Posada holds a bachelor degree in Computing Engineering, is a scholarship holder in Libre Software and Open Knowledge Office at the *Universidad de Cádiz*, Spain, and is an assisting student in the Department of Computer Languages and Systems. As a free knowledge and wikis enthusiast, he is a veteran Wikipedia user. Most of his contributions to the free software community are related to wikis and his development of bots to perform maintenance tasks automatically. His best known project is AVBOT, a bot to revert vandalism edits in the Spanish Wikipedia,



which was awarded the "Best Community Project" in the *III Spanish National Free Software Contest*, co-sponsored by Novática. He has delivered several talks and workshops about bot development, Wikipedia and MediaWiki. He has also been a member of committees for some free software and knowledge conferences at the *Universidad de Cádiz*. He is currently the main developer of the StatMediaWiki project (a tool for the statistical analysis of wikis) and has developed WikiTeam (a set of tools for wiki preservation and a repository of wiki backups). <emiliojose.rodriguez@uca.es>

However, some persons add profanity or insults (vandalism) to articles, or remove pieces of text (blinking), thereby abusing this open editing model. All these undesired edits can easily be undone using the history of each article, but it is tedious a task for a human and, while the vandalism is being reverted, that Wikipedia page is compromised. As the visibility of Wikipedia has grown, the vandalism problem has worsened, revealing a need for new, automatic, non-intrusive anti-vandalism solutions. This will allow committed users to spend their time adding new content, improving available content, helping new community members, and performing other productive tasks, rather than reverting vandalism.

AVBOT [1][4], an acronym for "Anti-Vandalism BOT", offers an automatic solution to most of the above mentioned acts of vandalism. Generally speaking, computers can easily detect undesired contributions, due to the fact that modifications include profanities or nonsensical text. AVBOT analyses recent edits in the Spanish Wikipedia, searching for bad-faith contributions, and undoes them. AVBOT will never replace a person, due to the fact that human intelligence is better

at analysing texts, but it can revert a great deal of vandalism, saving a lot of community time. Since AVBOT was created, it has repaired more than 200,000 acts of vandalism [3] and this number is growing.

2 Development

AVBOT is developed in Python [7] and it requires the pywikipediabot [8] and python-irclib [9] packages. Its workflow can be described as follows: spotting recent changes in articles, analysing changes, and decision making.

In the first step it obtains a list of recent changes in articles. This log is published in real time in an IRC channel [10] to which AVBOT is connected 24/7. Next, it discards those edits by administrators, maintenance bots and veteran users, since those users are trusted. The rest of the edits are analysed to detect possible vandalism.

“Wikipedia is the free encyclopaedia that anyone can edit. Any person using a web browser can modify its contents, via the MediaWiki software”

The bot compares the text of the previous version and that of the amended one. Then, an analysis module checks for profanity using a configurable list of regular expressions. Each regular expression has a score attached so, if the total sum of the scores for an edit is over a certain threshold, the edit is classified as vandalism and the Wikipedia article is reverted to its previous version. Both regular expressions and scores have been improved using the experience gained over the years the tool has been tested on the Spanish Wikipedia.

Another type of bad faith edit is blanking, in which users remove article content totally or partially. This form of vandalism can be detected by comparing the percent of text removed in the edit and reverting if it is excessive. The exact percentage threshold depends on article size and has been estimated by analysing blanking vandalism over the years.

Since the Spanish Wikipedia has a great many contributions (about 30,000 edits/day), it was necessary to optimize AVBOT in several ways, such as by the use of threads for parallel processing and the pre-compilation of lists of regular expressions.

When AVBOT reverts an edit, it also sends a message to the author asking him or her to stop. If the user continues to vandalize, a report will be generated. This report is sent to the administrators and they will investigate the situation [11].

However, as the bot may fail (the error ratio is under 0.5%, 1 error every 200 undone edits), AVBOT watches to see whether its reverts are discarded by human users. In that case, a notice is sent to the bot operator for future checking. This enables the bot's reversion skills to be improved.

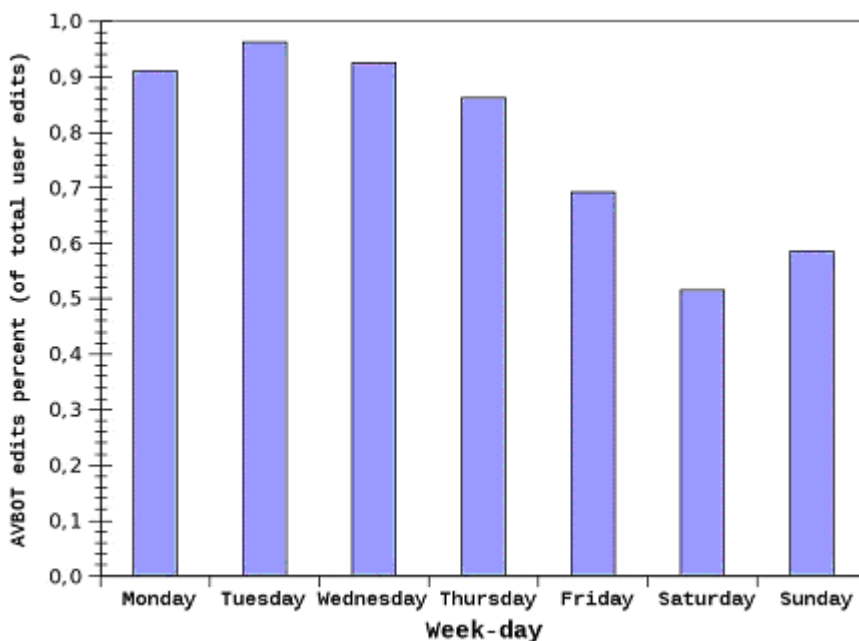


Figure 1: Percentage of AVBOT Edits (of Total User Edits) in the Course of a Week.

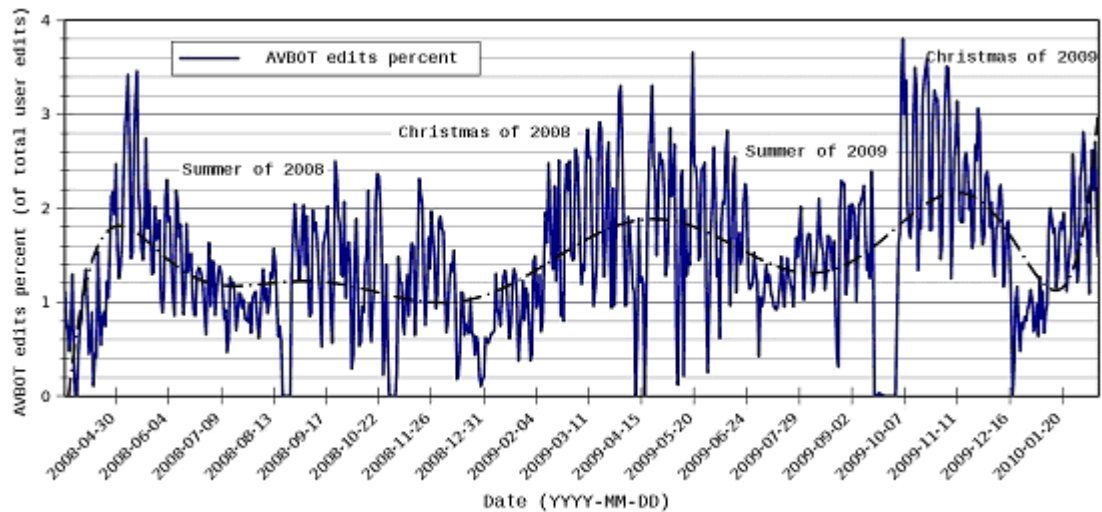


Figure 2: Percentage of AVBOT Edits (of Total User Edits) from April 2008 to January 2010. Low Vandalism Activity at Holiday Times and at Weekends is observed.

Speaking of AVBOT errors, most of them are related to polysemy. A good example is "The Ugly Duckling", which contains the word "ugly", a word used to attack biographies, although it is correct in this case. The error is fixed by adding a regular expression with a counterweight.

Finally, this software includes other features such as an exclusion list for pages which we do not want to monitor, for example, talk pages or Wikipedia internal maintenance pages, because slang is allowed in those places. AVBOT checks new pages (about 400 new articles are created in

Spanish Wikipedia every day), tagging for deletion those which contain test edits by anonymous users and other useless contributions.

3 Results Analysis

AVBOT's results have been satisfactory, with one error for every 200 correct reverts. Precompiled regular expressions used in vandalism analysis allow a very fast response (under a millisecond), so the only latency is due to network communication. Total time between when the article is vandalized and when it is restored is less than 5 seconds, making the effect of vandalism almost invisible.

Moreover, a detailed analysis of vandalism distribution throughout the week shows that it is concentrated from Monday to Thursday (see Figure 1), and during the European evening, because in that time frame both Spanish and Latin American users are active. On the other hand, there is low vandalism activity during holidays like summer, Christmas Day and New Year's Eve.

4 Community

Every bot which runs on the Spanish Wikipedia must pass a request for approval. AVBOT was approved with 19 votes in favour and 0 against [12].

Since the beginnings of AVBOT, the Spanish Wikipedia users have suggested fixes and new features. The

regular expressions list has been improved by community members adding new patterns.

Also, since AVBOT is free software and the source code is publicly available [2], some users have downloaded the code and run clones so, if the main copy of AVBOT does not work momentarily, Spanish Wikipedia is always protected.

The history of the project has been narrated in the official blog [1], which was used as a changelog for the *III Spanish National Free Software Contest*, where AVBOT was presented with the award for the "Best Community Project".

Software customization is possible by using a number of options which

“ As the visibility of Wikipedia has grown, the vandalism problem has worsened, revealing a need for new, automatic, non-intrusive anti-vandalism solutions ”

“ AVBOT, an acronym for "Anti-Vandalism BOT", is a program that offers an automatic solution to most of the acts of vandalism ”

“ AVBOT has been working in the Spanish Wikipedia since 2008, and has reverted more than 200,000 vandalism edits ”

enable AVBOT to run in any MediaWiki community with minor changes. Users from other Wikipedia languages have been interested in using the bot.

With regard to the future, the next tasks to be performed are code internationalization and adaptation to other languages, by creating new regular expressions lists for other Wikipedias and improving the documentation. The project is open to collaboration for everyone.

5 Conclusions and Future Work

The problem of vandalism has grown over the years as the popularity of Wikipedia has grown. AVBOT has been working in the Spanish Wikipedia since 2008, and has reverted more than 200,000 vandalism edits [3], saving the community from a huge workload. The wiki world evolves quickly, so it will be necessary to adapt the software to new changes, adding new features and improving existing ones. Adding machine learning to the bot would be a good step, training it with human reverted vandalism edits. Meanwhile, source code internationalization will help extend the community of users.

Another future line of work is the creation of tools for the analysis and assessment of wikis. A number of educational projects are currently being developed at the *Universidad de Cádiz*, where students are generating free knowledge in the classroom [15]. For example, with the support of the Libre Software and Open Knowledge Office, the Computer Science course "*Programación Funcional*" launched WikiHaskell, a wiki where students create free documentation about libraries for the Haskell programming language. Some projects related to wikis have also been published: WikiUNIX, WikiJuegos and IberOgre [16]. Most of these wikis do not allow anonymous edits during the academic year, but they

are open to edit for everyone later, so AVBOT would help to avoid vandalism. Moreover, StatMediaWiki [14], a tool for wiki analysis, has been developed by Emilio José Rodríguez to assess wiki contributions. Recently, he also founded WikiTeam [17], a set of tools for wiki preservation and a repository of wiki backups.

References

- [1] Official AVBOT website. Latest news, downloads and documentation. <<http://avbot.blogspot.com>>.
- [2] AVBOT source code repository. <<http://code.google.com/p/avbot>>.
- [3] Online demonstration of AVBOT reverts. <<http://es.wikipedia.org/wiki/Especial:Contribuciones/AVBOT>>.
- [4] Final project thesis for Bachelor in Computing Engineering of Emilio José Rodríguez-Posada. "Detector y corrector automático de ediciones maliciosas en Wikipedia". <<http://dx.doi.org/10498/8691>>.
- [5] Wikipedia, the free encyclopaedia. <<http://www.wikipedia.org>>.
- [6] MediaWiki. Wiki software for websites. <<http://www.mediawiki.org>>.
- [7] Official Python website. <<http://www.python.org>>.
- [8] Official pywikipediabot website. Framework for bots in MediaWiki. <<http://pywikipediabot.sourceforge.net>>.
- [9] Official Python-irclib website. Framework for IRC in Python. <<http://python-irclib.sourceforge.net>>.
- [10] Recent changes IRC channel for Spanish Wikipedia. <<irc://irc.wikimedia.org/es.wikipedia>>.
- [11] Wikipedia: Vandalismo en curso. Administrators' noticeboard. <http://es.wikipedia.org/wiki/Wikipedia:Vandalismo_en_curso>.
- [12] AVBOT bot request for approval. <http://es.wikipedia.org/wiki/Wikipedia:Bot/Autorizaciones/Archivo_2008#AVBOT>.
- [13] Wikimedia projects edits counter. <<http://toolsserver.org/~emijrp/wikimediacounter>>.
- [14] Official StatMediaWiki website. A tool for the statistical analysis of MediaWiki wikis. <<http://stat.mediawiki.forja.rediris.es>>.
- [15] Manuel Palomo, Inmaculada Medina, Emilio José Rodríguez and Noelia Sales. *Tecnologías wiki y conocimiento abierto en la Universidad*. Proceedings of the 5th Open Source World Conference, Cáceres 2009. pp 16-19. ISBN 978-84-692-8739-2.
- [16] Wikis @ OSLUCA. Some wikis managed by Libre Software and Open Knowledge Office at University of Cádiz. <<http://osl.uca.es/wikis>>.
- [17] Official WikiTeam website. Tools for wiki preservation and a repository of wiki backups. <<http://code.google.com/p/wikiteam>>.