

WHERE IS HISTORY BEING WRITTEN?
GEOREFERENCING CONTRIBUTIONS
TO WIKIPEDIA

DAVID KALTSCHMIDT

Diplomarbeit

Dr. Claudia Müller-Birn
Prof. Dr. Robert Tolksdorf
Institut für Informatik
Freie Universität Berlin

David Kaltschmidt: *Where is history being written? Georeferencing
contributions to Wikipedia*
Diplomarbeit, © 2011

SUPERVISORS:
Dr. Claudia Müller-Birn
Prof. Dr. Robert Tolksdorf

LOCATION:
Berlin, Germany

YEAR:
2011

ABSTRACT

Wikipedia is more than an online encyclopedia. It is also a news channel as well as a self-updating history book. A global readership can follow political events as they unfold, written about by local people and later edited by other volunteers. This thesis describes a method to answer the question to what extent local volunteers write about events in their own country. First, the geographic origin of each individual article contribution is determined. In a second step, a given article is annotated with georeferences on a word level. The properties of these annotations then allow for a statistical geographic analysis of a single article or a category of articles.

ZUSAMMENFASSUNG

Als Online-Enzyklopädie ist Wikipedia nicht nur Nachschlagewerk sondern auch ein sich stetig wandelndes Geschichtsbuch. Eine global verteilte Nutzerschaft liest und schreibt über lokale Ereignisse noch während sie passieren. Diese Arbeit beschreibt eine Methode zur Bestimmung des Anteils an Beiträgen, die vom betreffenden Land ausgehen. In einem ersten Schritt werden die geographischen Ursprünge aller Beiträge eines Artikels ermittelt. Mit den daraus erhaltenen Georeferenzen wird der Artikel Wort für Wort annotiert. Basierend auf diesen Annotationen kann dann der lokale Autoren-Anteil bestimmt werden.

CONTENTS

I	THOUGHTS	1
1	INTRODUCTION	2
1.1	Structure	3
2	FOUNDATION	5
2.1	Wikipedia	5
2.2	Contributions	6
2.3	Georeferences	6
2.4	Visualization	7
3	HYPOTHESES	8
3.1	Research questions	8
3.2	Scope	9
3.2.1	Historic articles	9
3.2.2	Article location	9
II	METHODS	10
4	APPARATUS	11
4.1	Wikipedia's Data Structures	11
4.1.1	Zugriff	12
4.2	Collective Authorship	12
4.2.1	Relevant Edits	12
4.3	Georeferences	12
4.3.1	IP Look-up	13
4.3.2	Information Extraction	13
4.3.3	Geographic Profiling	14
4.3.4	Consolidation	14
4.4	Visualization	14
4.4.1	Goals	15
4.4.2	Design	15
4.5	Data Model and System Overview	15
5	EXPERIMENTS	16
5.1	Data Set	16
5.2	Application	16
III	RESULTS	17
6	RESULTS	18
7	CONCLUSION	19
7.1	Limitations	19
7.2	Further Research	19
IV	APPENDIX	20
	Bibliography	21

LIST OF FIGURES

LIST OF TABLES

LISTINGS

ACRONYMS

Part I

THOUGHTS

INTRODUCTION

If you are open to contributions from others, you generally end up with richer, better, more diverse and expert content than if you try to do it alone.¹

— Alan Rusbridger, editor of THE GUARDIAN

At the end of January 2011, when a wave of public protest spilled from Tunisia into Egypt, a small group of opposition parties and political activists called for a “Day of Rage” via Facebook, a social networking website. By January 25th their Facebook group had more than 80,000 supporters who drew attention to and helped organize the country-wide protests that followed. As people rallied the streets day after day, the Egyptian government first limited access to Twitter, a micro-blogging service, before cutting Egypt off the internet completely on January 28th.[2, 3]

In what came to be known as the Arab Spring, the use of online networks directly influenced the political development. While Facebook played a part in organizing the protests, Twitter acted as an information channel during the demonstrations. As the events unravelled, they were reflected by articles created on Wikipedia, an online encyclopedia. Updated by the minute, the articles covering the protests formed a well of news reports.[4] As ordinary people become producers of journalism the need arises to analyze these contributions. Specifically, this thesis focuses on the geographic origins of contributions to Wikipedia articles.

Wikipedia’s free access and open editing policy as well as a quality level—putting it “head to head”[5] with Encyclopedia Britannica—turned it into a hugely popular website[6]. The server software used for the website, MediaWiki², ensures that the effort to change an article is minimal. Given an Internet connection and a web browser, anyone can add or edit an account of current events in a related article and publish it in a matter of seconds.

This form of news production turns the encyclopedia into a news channel that is constantly updated and corrected by an army of volunteers. The result is a self-governed news source that lends itself the aura of authority and credibility of a knowl-

¹ The Economist. *The people formerly known as the audience*. URL: <http://www.economist.com/node/18904124> (visited on 08/10/2011)

² <http://www.mediawiki.org>

edge reference. At the same time a technophile public, that uses the Internet as an efficient means of news acquisition, can check facts on Wikipedia and act upon the consumed information.[7, p. 424–427] Therefore the collective authorship of such a news medium could have a direct influence on the political decision process.

Political events are often limited to a country or region. This is reflected by the Wikipedia articles covering the Arab Spring: there is an overarching parent article³ as well as single articles covering the revolution in each of the affected countries, e.g. Egypt⁴ and Libya⁵. The latter also exemplifies how divided the political actors can be. While nearly all revolutionaries welcomed the airstrikes, one faction was concerned foreign meddling and another one just opposed the deployment of ground troops.[8]

The collective authorship could be equally divided⁶ while at the same time creating a potential for further analysis. Where do the first reports of an event originate? As later iterations of revisions turn these reports into historical accounts, are these editors from the same country? And more generally, to what extent is a collection of these articles written by volunteers located at the respective location of the event.

In this thesis I will propose a method to answer these questions. By trying to determine the geographic origin of each edit to an article I will be able to calculate the geographic distribution of contributions. This distribution will then be used to answer the questions above for either a single article or a collection.

Complete summary and key findings.

1.1 STRUCTURE

complete over time, name the basic chapters and their function, one part = one paragraph

The next chapter **FOUNDATION** provides background information about **Wikipedia**, article editing (**Contributions**), the application of geographic data (**Georeferences**) and its presentation (**Visualization**). The first part ends with **HYPOTHESES** where I propose the research questions that this thesis hopes to answer.

³ http://en.wikipedia.org/wiki/2010-2011_Middle_East_and_North_Africa_protests

⁴ http://en.wikipedia.org/wiki/Egyptian_Revolution_of_2011

⁵ http://en.wikipedia.org/wiki/2011_Libyan_uprising

⁶ Despite Wikipedia's core policy to oblige everyone to write from a neutral point of view (NPOV), people regularly express opinions. The collision of opinions in a collectively written article can result in a prolonged series of an edit and its subsequent reversal by another person. The resulting edit pattern is known as an *edit war*.^[9]

In [APPARATUS](#) I will describe the tools available to be used in the method that will be applied to a host of articles in [EXPERIMENTS](#). The findings will be presented in [RESULTS](#). Followed by a discussion of their feasibility in [CONCLUSION](#).

FOUNDATION

Wikipedia is a phenomenon that has attracted researchers across all fields, notably computer science and sociology, who have written over 1,000 reports on the subject.[10] These roughly fall into three categories:

CONTENT PRODUCTION Covering all aspects of voluntary production such as motivation, collaboration, coverage and bias, quality and vandalism, actuality and geography.

INFORMATION USE Treating how the resulting corpus is being used, e.g. citing in research, use in court, trendspotting, thesaurus construction or categorization.

IMPROVEMENT These are studies concerned with the improvement of both the software used by Wikipedia and the content, e.g. automatic linking, bots, improved editors as well as quality and trust indicators.

This thesis falls into the first category as it examines the geography of article contributions that will become part of the Wikipedia's corpus. After a short overview of Wikipedia from a user perspective, I will introduce its model of collective authorship and present prior research of concerning location and geography.

2.1 WIKIPEDIA

this section should cover the basics to understand components of wikipedia

Die Online-Enzyklopädie Wikipedia gibt es in über 260 Sprachvarianten, von denen die englische mit derzeit 3,6 Millionen Artikeln mit Abstand die größte ist. Die Anzahl der Artikel in den anderen Sprachen sowie die Nutzung der jeweiligen Sprachvariante unterscheiden sich jedoch erheblich.[11] Wenn ein Artikel zum selben Thema in Wikipedias unterschiedlicher Sprachen vorhanden ist, sind diese Varianten in der Regel über sogenannte Interwiki-Links untereinander verlinkt.

language editions, chart

Die Artikel dieser Lexika werden von Freiwilligen auf der ganzen Welt geschrieben, gemeinschaftlich korrigiert und aktualisiert. Jede Änderung eines Artikels erzeugt eine neue Version, die der Versionsgeschichte des Artikels hinzugefügt wird und danach für alle Benutzer einsehbar ist.

article, draw nice graphic of UI

Jeder Eintrag in der Versionsgeschichte besteht dabei aus der Textänderung, dem Datum der Version, dem Benutzer sowie einem optionalen Kommentar über den Grund der Änderung. Jede Änderung kann mit Hilfe dieser Historie ausführlich begutachtet und bei Missfallen wieder revidiert werden. Dies kann mitunter sogenannte *edit wars* hervorrufen, in denen neue Beiträge von Nutzern mit entgegengesetzten Standpunkten sofort wieder revidiert werden.[9]

revision history, changesets

Die Mitarbeit an den Artikeln kann mit oder ohne vorherige Registrierung erfolgen. Autoren, die sich registrieren, erlangen sowohl bestimmte Privilegien wie zum Beispiel das Recht, neue Einträge zu erstellen, als auch den Zugang zu Wikipedias sozialem Netzwerk: Jeder Benutzer erhält nach der Registrierung eine *user page* auf der er Informationen über sich veröffentlichen und über die er mit anderen Nutzern Kontakt aufnehmen kann.[12]

2.2 CONTRIBUTIONS

introduce collective authorship and name some important concepts. prior research in:

- *text-longevity*
- *attribution*

Eine Analyse der Autorschaft bis auf Satzebene innerhalb eines Artikels wird von Kramer in [13] erforscht. Durch Auswertung der Versionsgeschichte lässt sich zu jedem Satz der Autor bestimmen, der dessen Hauptteil geschrieben hat. Eine automatische Auswertung eines Artikels bis auf Wortebene wird von Adler in [14] vorgestellt. Sie basiert auf dem von Adler selbst entwickelten Reputationssystem [15], das Textstellen eine hohe Vertrauenswürdigkeit zuweist, die von einem vertrauenswürdigen Autor geschrieben oder mindestens einmal bearbeitet worden sind.¹

Für eine Analyse der Artikel bis auf Satzebene werden Algorithmen wie in [13] auf ihre Anwendbarkeit untersucht.

2.3 GEOREFERENCES

explain this intermediate step to assign a location to a contribution

¹ Basierend auf diesen beiden Arbeiten wurde die Software WikiTrust implementiert, welches die Vertrauenswürdigkeit als weiß-orange *Heatmap* darstellt: zweifelhafte Textstellen werden orange hinterlegt und damit leicht erkennbar. Über ein API ist eine mit Vertrauenspunkten annotierte Version eines Artikels abrufbar: <http://www.wikitrust.net/vandalism-api>

- *pick up where he left off: D. Hardy. „Volunteered geographic information in Wikipedia.“ PhD thesis. UNIVERSITY OF CALIFORNIA, SANTA BARBARA, 2011*
He's more about the how, I'm about the where. Key findings:
- *"I find that as a group, anonymous contributors write about fewer places than registered contributors, despite outnumbering them five-to-one." [16]*
- *"I find that anonymous contributors are more likely to write about nearby places, and that the geographic effects fit an exponential distance decay function." [16]*
- *"Combined approaches (i.e., where quantitative spatial analysis models are calibrated with surveyed locations) may prove useful." [16, p. 85]*
- *WikiScanner*
- *Erik Zachte's: Wikipedia edits visualized²*
- *Indirect approach M.D. Lieberman and J. Lin. „You are where you edit: Locating Wikipedia users through edit histories.“ In: ICWSM'09 (2009), pp. 106–113*

Zur Bestimmung der Herkunft eines Autors bietet Wikipedia zwei direkte Ansätze: Für jeden Beitrag eines nicht registrierten Benutzers wird die IP-Adresse gespeichert, über die er Zugang zum Internet erlangt hat. Die registrierten Nutzer können jedoch auf ihrer *user page* Informationen über ihre Person entweder als Freitext oder strukturiert in *user boxes* veröffentlichen.

Ein zusätzlicher, indirekter Ansatz für die Bestimmung der Herkunft eines Nutzers wird von Lieberman in *You are where you edit: Locating Wikipedia users through edit histories*[17] beschrieben. Er basiert auf der Annahme, dass ein Nutzer mit Vorliebe an Artikeln über Orte in seiner geographischen Nähe mitarbeitet. Diese Artikel sind in der Regel mit geographischen Koordinaten versehen und erlauben so eine sehr grobe Bestimmung des Aufenthaltsortes und dessen Visualisierung auf einer Landkarte.

2.4 VISUALIZATION

Write about prior works of visualizing the aspects of attribution and georeference

- *Erik Zachte's: Wikipedia edits visualized³*
- *Wikitrust 1*

² <http://infodisiac.com/blog/2011/05/wikipedia-edits-visualized/>

³ <http://infodisiac.com/blog/2011/05/wikipedia-edits-visualized/>

3

HYPOTHESES

My effort to answer this thesis' main question, where the contributions to historic articles are coming from, will follow the path laid out by a set of questions:

3.1 RESEARCH QUESTIONS

A central part of this thesis is the determination of the geographic origins of individual contributions. In prior research regarding the spatial distribution of contributors, most notably „Volunteered geographic information in Wikipedia,“ Hardy [16, see 3.4.2] took only anonymous users into account. Since anonymous contributions reveal the IP address their location can easily be approximated by using IP geolocation services. He has firmly established that anonymous contributors are more likely to write about nearby places. However, the above approach did not consider contributions done by registered users who, though being outnumbered by anonymous users, provide the majority of edits. Based on alternative approaches like Lieberman and Lin [17] and Engelhardt, Freytag, and Schulz [18], I will investigate whether registered users can be included for a spatial analysis. The search for the location of such a user will be satisfied when it yields a country which will be used as the origin of this user's contributions. A key question is thus:

QUESTION 1. Can a country of origin can be assigned to registered users?

Hardy [16] also stopped short of how the spatial distribution changes over time. The knowledge of time and location of the contributions suggest the following question:

QUESTION 2. Does the spatial distribution of contributors change over time?

Once the contributions have been georeferenced, the question of localness of the contributions can only be answered if the article itself can be geographically designated. A handful of studies exist, notably [16, 19], that perform a spatial analysis on Wikipedia articles that are geo-tagged. These articles contain geographic coordinates as part of the content, e.g. the article about the *Brandenburg Gate*¹ in Berlin is tagged with the coordinates 52°30'58.58"N 13°22'39.80"E. However, these articles are

¹ http://en.wikipedia.org/wiki/Brandenburg_Gate

geographic in nature treating cities, rivers and places of interest. This thesis hopes to expand the spatial analysis of contributions to articles that have a location property only by association, e.g. the *Egyptian Revolution of 2001*², that clearly happened in *Egypt*³. To allow for an analysis of historic events it is necessary to investigate the articles in a similar manner as the registered users:

QUESTION 3. Can a location be assigned to an historic article?

Finally, an article with a location can be related to the spatial distribution of its contributors. In the context of historic articles this suggests the final question:

QUESTION 4. Are historic articles written by their respective citizens?

3.2 SCOPE

3.2.1 *Historic articles*

The same questions may be applicable to other types of articles than historic ones. However, the latter satisfy the requirement of having a location attribute and were therefor referred to throughout this thesis.

3.2.2 *Article location*

Although *location* is central to more abstract concepts like *Culture*⁴ these subjects clearly defy being attributed with a location. Nevertheless, an analysis of the spatial distribution of contributors could be interesting.

² http://en.wikipedia.org/wiki/Egyptian_Revolution_of_2011

³ <http://en.wikipedia.org/wiki/Egypt>

⁴ <http://en.wikipedia.org/wiki/Culture>

Part II

METHODS

APPARATUS

What tools do I have and how can they be extended.

4.1 WIKIPEDIA'S DATA STRUCTURES

ARTIKEL Ein Artikel hat mindestens einen Autor und ist gegebenenfalls in mehreren Sprachen vorhanden.

VERSIONSGESCHICHTE Diese Historie liefert Informationen wie Benutzername oder IP-Adresse, Datum der Version sowie die inkrementelle Textänderung.

USER PAGES & USER BOXES Auf den *user pages* kann ein registrierter Benutzer Informationen über sich veröffentlichen, die Aufschluss über seine Herkunft geben könnten.

EXTERNE QUELLEN Im Internet existieren zahlreiche Dienste, die Schnittstellen anbieten, um Informationen über Nutzer und deren Beiträge zu erhalten, z.B.: WikiTrust¹ oder Wiki-Watcher¹

Die Methoden zur Datenextraktion und Visualisierung werden anschließend in eine Software integriert. Die Gewinnung der von dieser Anwendung zu verarbeitenden Daten kann aus einer der folgenden Quellen erfolgen:

DB-KOPIE Monatlich angefertigte Moment-Aufnahmen der gesamten Wikipedia-Datenbank sind öffentlich verfügbar². Eine solche Kopie enthält alle Artikel inklusive Versionsgeschichte und ist damit jedoch sehr groß³.

ARTIKELEXPORTE Jeder einzelne oder mehrere Artikel der Wikipedia kann auch separat exportiert werden. Diese Daten umfassen ebenfalls die Versionsgeschichte und sind im Umfang bedeutend kleiner.

TOOLSERVER Die Wikimedia Deutschland e.V. stellt Server bereit,⁴ welche einen direkten Zugang zu einer replizierten, schreibgeschützten

¹ Das WikiWatcher-Teilprojekt *Poor Man's Check User* erlaubt eine Auflösung des Benutzernamens in eine IP-Adresse, wenn dieser Nutzer in der Vergangenheit beim Ändern eines Artikels das Session-Limit überschritten hatte. Inzwischen wurde diese Sicherheitslücke in der WikiMedia-Software jedoch behoben.

<http://wikiwatcher.virgil.gr/pmcu>

² <http://dumps.wikimedia.org>

³ Eine Kopie der englischen Wikipedia-Datenbank umfasst derzeit 5,4 Terabyte.

⁴ <http://toolserver.org>

Wikipedia-Datenbank ermöglichen. Die Nutzung eines solchen Servers vermeidet es zwar, eine eigene komplette Kopie der gesamten Wikipedia-Datenbank halten zu müssen, bedarf jedoch einer Anmeldung.

4.1.1 Zugriff

Tools and servers to access the articles.

4.2 COLLECTIVE AUTHORSHIP

Introduce types of authors (roles) as well as methods to determine contribution/attribution

- *Autoren*
- *Bots*
- *Wer überlebt?*
- *Algorithmen, welche Unterschiede?*

4.2.1 Relevant Edits

Are all edits relevant? Edit wars? Bots?

4.3 GEOREFERENCES

- *registered vs. unregistered vs. bots vs. admins*
- *incorporate key findings of [16] as laid out in chapter 2.3*
- *IPs of unregistered users: Geo lookup*
- *Autoren-Profile: Information Extraction*
- *Geographische Zuordnung vom user profile*

Zur Bestimmung der Herkunft eines Autors bietet Wikipedia zwei direkte Ansätze: Für jeden Beitrag eines nicht registrierten Benutzers wird die IP-Adresse gespeichert, über die er Zugang zum Internet erlangt hat. Der zweite Ansatz betrifft die registrierten Benutzer. Ihre IP-Adressen sind maskiert und nicht öffentlich zugänglich.⁵ Die registrierten Nutzer können jedoch auf ihrer *user page* Informationen über ihre Person entweder als Freitext oder strukturiert in *user boxes* veröffentlichen. Letztere sind

⁵ Eine kleine, von der Wikipedia-Community gewählte Nutzerschaft mit der Berechtigung *checkuser* kann die Adressen demaskieren.

definierte Einheiten mit denen der Nutzer persönliche Eigenschaften wie Herkunftsland, gesprochene Sprachen oder wissenschaftliche Interessen kodifizieren kann. Zusammen decken beide Ansätze jedoch nur einen Teil der Beiträge schreibenden Nutzerschaft ab.

4.3.1 IP Look-up

- *Services*
- *Accuracy*
- *Active prevention by proxies and anonymizers:*
J.A. Muir and P.C.V. Oorschot. Internet geolocation and evasion. Tech. rep. Citeseer, 2006
J.A. Muir and P.C.V. Oorschot. „Internet geolocation: Evasion and counterevasion.“ In: ACM Computing Surveys (CSUR) 42.1 (2009), p. 4
M. Duckham and L. Kulik. „A formal model of obfuscation and negotiation for location privacy.“ In: Pervasive Computing (2005), pp. 152–170

Mit frei verfügbaren⁶ Online-Diensten wie *Quova*⁷ oder *geoplugin*⁸ lässt sich für einen Großteil der IPs daraufhin das Herkunftsland bestimmen.

Im Bezug auf die Herkunft sind sowohl das Land als auch die Geo-Koordinaten interessant. Basierend auf der Versionsgeschichte würde für nicht registrierte Benutzer eine Gewinnung von Daten dann beispielsweise folgende Schritte durchlaufen:

IP \Rightarrow Geolocation-Dienst \Rightarrow Koordinaten und Land

4.3.2 Information Extraction

- *IE approach with Machine Learning* L. Xiao et al. „Information extraction from the web: System and techniques.“ In: Applied Intelligence 21.2 (2004), pp. 195–224
- *unsupervised IE:* O. Etzioni et al. „Unsupervised named-entity extraction from the web: An experimental study.“ In: Artificial Intelligence 165.1 (2005), pp. 91–134
- *if city is mentioned, determine country (needs disambiguation, e.g. Berlin)*
- *coordinates are optional?*

⁶ Die vorgestellten Dienste haben ein tägliches Kontingent an Anfragen. Hilfstechiken wie Caching können diese Einschränkungen jedoch mindern.

⁷ <http://developer.quova.com>

⁸ <http://www.geoplugin.com/webservices>

4.3.3 Geographic Profiling

- M.D. Lieberman and J. Lin. „You are where you edit: Locating Wikipedia users through edit histories.“ In: ICWSM’09 (2009), pp. 106–113
- B.J. Hecht and D. Gergle. „On the localness of user-generated content.“ In: Proceedings of the 2010 ACM conference on Computer supported cooperative work. ACM. 2010, pp. 229–232
- from other fields such as criminal research:
B. Snook et al. „On the complexity and accuracy of geographic profiling strategies.“ In: Journal of Quantitative Criminology 21.1 (2005), pp. 1–26
- feasibility, maybe just as enhancer

4.3.4 Consolidation

- settle for a resolution
- some examples on accuracy for different countries
- clustering of origins: areas of influence

4.4 VISUALIZATION

- Darstellung der geographischen Analyse
- per Wort, Satz, Artikel, Wort

Auf Basis der strukturierten Daten in Form von Artikeln, Sätzen, Ländern, Koordinaten und Sprachen sollen nun Visualisierungen gefunden werden, welche die Fülle an Informationen zugänglich machen. Mögliche Visualisierungen wären etwa:

V1 Revisionshistogramm à la Google Finance

V2 *Heatmap* einer Landkarte mit Ursprüngen der Revisionen

V3 Netzwerkgrafik, die Metriken desselben Artikels in verschiedenen Sprachvarianten anzeigt

V4 Dynamisches Blasendiagramm⁹ über die Entwicklung unterschiedlicher Sprachvarianten

V5 *Heatmap* des Artikels mit Stellen höchster Aktivität

V6 Landeskürzel für eine gegebene Textstelle

V7 Edit wars on map, linking two or more places

⁹ http://en.wikipedia.org/wiki/Motion_chart

4.4.1 Goals

- A. Kjellin et al. „Evaluating 2D and 3D visualizations of spatiotemporal information.“ In: ACM Transactions on Applied Perception (TAP) 7.3 (2010), pp. 1–23
- *Identify.* Characteristics of an object.
- *Locate.* Absolute or relative position.
- *Distinguish.* Recognize as the same or different.
- *Categorize.* Classify according to some property (e.g., color, position, or shape).
- *Cluster.* Group same or related objects together.
- *Distribution.* Describe the overall pattern.
- *Rank.* Order objects of like types.
- *Compare.* Evaluate different objects with each other.
- *Associate.* Join in a relationship.
- *Correlate.* A direct connection.

4.4.2 Design

4.5 DATA MODEL AND SYSTEM OVERVIEW

- *fetch article*
- *get revision history*
- *determine contributions*
- *transform to word attribution*
- *attach georeference*

5

EXPERIMENTS

5.1 DATA SET

- *Wahl einer Kategorie/Artikel*
- *Wieso repräsentativ für politische Ereignisse?*

Mithilfe der Export-Funktion von Artikeln lässt sich ein kleiner Datensatz generieren, an dem die Anwendung getestet werden kann. Über dieselbe Export-Funktion kann auch eine Kategorie wie zum Beispiel *Revolutions by country*¹ angegeben werden. Als Ergebnis erhält man eine Sammlung von Artikeln über politische Ereignisse.

5.2 APPLICATION

- *Beispielhafte Durchführung*
- *Sammlung der Ergebnisse*

Dabei könnte zum Beispiel sichtbar werden, dass sich ein bestimmter Artikel in verschiedenen Sprachvarianten unterschiedlich entwickelt. Falls ein Land mehrere offizielle Sprachen hat, könnte man diese entweder gruppiert oder einzeln im direkten Vergleich betrachten. Ebenso könnten sich in Anlehnung an die *edit wars* Streitpunkte anhand von Textstellen herauskristallisieren, die besonders umkämpft sind.

¹ http://en.wikipedia.org/wiki/Category:Revolutions_by_country

Part III

RESULTS

6

RESULTS

- *Statistische Auswertung*

Anhand eines ausgewählten Datensatzes von politischen Ereignissen wie *Revolutions by country* soll eine statistische Auswertung erfolgen, um die Frage zu beantworten, wer die Geschichte eines Landes schreibt.

CONCLUSION

- *Interpretation der Ergebnisse*
- *Vermutungen bestätigt*

7.1 LIMITATIONS

- *Mobile contributions, smartphones*
- *Privacy*

7.2 FURTHER RESEARCH

Part IV

APPENDIX

BIBLIOGRAPHY

- [1] The Economist. *The people formerly known as the audience*. URL: <http://www.economist.com/node/18904124> (visited on 08/10/2011).
- [2] The Economist. *Protest in Egypt: Another Arab regime under threat*. 2011. URL: <http://www.economist.com/node/18013760>.
- [3] sueddeutsche.de. *Krise in Ägypten - Die Kinder des 6. April und der Tag der Entscheidung - Politik*. 2011. URL: <http://www.sueddeutsche.de/politik/krise-in-aegypten-die-kinder-des-april-rufen-zum-protest-1.1053426>.
- [4] Wikipedia. *Revision history of 2011 Egyptian revolution*. URL: http://en.wikipedia.org/w/index.php?title=2011_Egyptian_revolution&dir=prev&action=history.
- [5] J. Giles. „Internet encyclopaedias go head to head.“ In: *Nature* 438.7070 (2005), pp. 900–901. ISSN: 0028-0836.
- [6] Wikipedia. *Wikipedia Page Views*. URL: <http://stats.wikimedia.org/EN/TablesPageViewsMonthly.htm>.
- [7] A. Chadwick. *Routledge handbook of Internet politics*. Taylor & Francis, 2009. ISBN: 0203962540.
- [8] The Economist. *Libya: A civil war beckons*. 2011. URL: <http://www.economist.com/node/18290470>.
- [9] B. Suh et al. „Us vs. them: Understanding social dynamics in Wikipedia with revert graph visualizations.“ In: *Visual Analytics Science and Technology, 2007. VAST 2007. IEEE Symposium on*. IEEE. 2007, pp. 163–170.
- [10] F.Å. Nielsen. „Wikipedia research and tools: Review and comments.“ In: (2011).
- [11] Wikipedia. *Statistics*. URL: <http://stats.wikimedia.org/EN/Sitemap.htm>.
- [12] Wikipedia. *Why create an account?* URL: http://en.wikipedia.org/wiki/Wikipedia:Why_create_an_account%3F.
- [13] M. Kramer, A. Gregorowicz, and B. Iyer. „Wiki trust metrics based on phrasal analysis.“ In: *Proceedings of the 4th International Symposium on Wikis*. ACM. 2008, pp. 1–10.
- [14] B.T. Adler et al. „Assigning trust to wikipedia content.“ In: *Proceedings of the 4th International Symposium on Wikis*. ACM. 2008, pp. 1–12.

- [15] B.T. Adler and L. De Alfaro. „A content-driven reputation system for the Wikipedia.“ In: *Proceedings of the 16th international conference on World Wide Web*. ACM. 2007, pp. 261–270.
- [16] D. Hardy. „Volunteered geographic information in Wikipedia.“ PhD thesis. UNIVERSITY OF CALIFORNIA, SANTA BARBARA, 2011.
- [17] M.D. Lieberman and J. Lin. „You are where you edit: Locating Wikipedia users through edit histories.“ In: *ICWSM’09 (2009)*, pp. 106–113.
- [18] S. Engelhardt, A. Freytag, and C. Schulz. „On the Geographic Allocation of Open Source Software Activities.“ In: *Jena Economic Research Papers* (2010), p. 009.
- [19] B.J. Hecht and D. Gergle. „On the localness of user-generated content.“ In: *Proceedings of the 2010 ACM conference on Computer supported cooperative work*. ACM. 2010, pp. 229–232.
- [20] J.A. Muir and PC van Oorschot. *Internet geolocation and evasion*. Tech. rep. Citeseer, 2006.
- [21] J.A. Muir and P.C.V. Oorschot. „Internet geolocation: Evasion and counterevasion.“ In: *ACM Computing Surveys (CSUR)* 42.1 (2009), p. 4.
- [22] M. Duckham and L. Kulik. „A formal model of obfuscation and negotiation for location privacy.“ In: *Pervasive Computing* (2005), pp. 152–170.
- [23] L. Xiao et al. „Information extraction from the web: System and techniques.“ In: *Applied Intelligence* 21.2 (2004), pp. 195–224.
- [24] O. Etzioni et al. „Unsupervised named-entity extraction from the web: An experimental study.“ In: *Artificial Intelligence* 165.1 (2005), pp. 91–134.
- [25] B. Snook et al. „On the complexity and accuracy of geographic profiling strategies.“ In: *Journal of Quantitative Criminology* 21.1 (2005), pp. 1–26.
- [26] A. Kjellin et al. „Evaluating 2D and 3D visualizations of spatiotemporal information.“ In: *ACM Transactions on Applied Perception (TAP)* 7.3 (2010), pp. 1–23.