

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

PRODUCED WITH:  
L<sup>A</sup>T<sub>E</sub>X using the ClassicThesis package.

## 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

---

<b>I</b>	<b>THOUGHTS</b>	<b>1</b>
1	INTRODUCTION	2
1.1	Structure	3
2	FOUNDATION	5
2.1	Wikipedia	5
2.1.1	History	5
2.1.2	Wikimedia Foundation	6
2.1.3	Editing policy	6
2.2	Contributions	6
2.3	Georeferences	7
3	HYPOTHESES	9
3.1	Article creation	9
3.2	Participation	10
3.3	Text survival	11
3.4	Scope and limitations	12
3.4.1	Political events	12
3.4.2	Article location	12
3.4.3	Cross-language article growth	12
<b>II</b>	<b>METHODS</b>	<b>13</b>
4	APPARATUS	14
4.1	Wikipedia's Data Structures	14
4.1.1	Toolkits	15
4.1.2	Mediawiki API	15
4.2	Collective Authorship	16
4.2.1	Relevant Edits	16
4.3	Georeferences	16
4.3.1	IP Look-up	16
4.3.2	Information Extraction	17
4.3.3	Geographic Profiling	17
4.3.4	Consolidation	18
4.4	Visualization	18
4.4.1	Goals	19
4.4.2	Design	19
4.5	Data Model and System Overview	19
4.6	Analysis	19
5	EXPERIMENTS	20
5.1	Data Set	20
5.2	Application	20
<b>III</b>	<b>RESULTS</b>	<b>21</b>
6	RESULTS	22

7	CONCLUSION	23
7.1	Limitations	23
7.2	Further Research	23
IV	APPENDIX	24
	Bibliography	25

## 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.<sup>1</sup>*

— 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]

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, MediaWiki<sup>2</sup>, 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 knowledge 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

---

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

<sup>2</sup> <http://www.mediawiki.org> (visited on 10/31/2011)



medium could have a direct influence on the political decision-making process. As ordinary people become producers of journalism the need arises to analyze these contributions. To investigate, this thesis focuses on the geographic origins of contributions to Wikipedia articles.

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 article<sup>3</sup> as well as single articles covering the revolution in each of the affected countries, e.g. Egypt<sup>4</sup> and Libya<sup>5</sup>. The events in the latter country also exemplify how divided the political actors can be. While nearly all revolutionaries welcomed the airstrikes, one faction was concerned about foreign meddling and another one just opposed the deployment of ground troops.[8]

The collective authorship could be equally divided. Despite Wikipedia's core policy to oblige everyone to write from a *neutral point of view*<sup>6</sup> (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] These clashes of opinion create a potential for further investigation into the geopolitics of article contribution. 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 contributors. 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 chapter **FOUNDATION** provides background information about **Wikipedia**, article editing (**Contributions**) and the applica-

<sup>3</sup> [http://en.wikipedia.org/wiki/2010-2011\\_Middle\\_East\\_and\\_North\\_Africa\\_protests](http://en.wikipedia.org/wiki/2010-2011_Middle_East_and_North_Africa_protests) (visited on 10/31/2011)

<sup>4</sup> [http://en.wikipedia.org/wiki/Egyptian\\_Revolution\\_of\\_2011](http://en.wikipedia.org/wiki/Egyptian_Revolution_of_2011) (visited on 10/31/2011)

<sup>5</sup> [http://en.wikipedia.org/wiki/2011\\_Libyan\\_uprising](http://en.wikipedia.org/wiki/2011_Libyan_uprising) (visited on 10/31/2011)

<sup>6</sup> [http://en.wikipedia.org/wiki/Wikipedia:Neutral\\_point\\_of\\_view](http://en.wikipedia.org/wiki/Wikipedia:Neutral_point_of_view) (visited on 12/08/2011)

tion of geographic data ([Georeferences](#)). The first part ends with [HYPOTHESES](#) where I propose the research questions that this thesis hopes to answer.

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 to date. Nielsen [10] divides these roughly into four 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. Wikipedia citation in research, use in court, trend-spotting, natural language processing and automatic translation tools, 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.

**COMMUNICATION** Studies in this category look at Wikipedia as an online collaboration tool for education and research.

This thesis falls into the first category, content production, 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's perspective, I will introduce its model of collective authorship and present prior research of concerning location and geography.

### 2.1 WIKIPEDIA

Wikipedia is an online encyclopedia with editions in over 260 languages. Counting 3.6 million articles, the English version is by far the biggest. However, other language editions differ sharply in size and usage.[11] If articles covering the same topic exist in other language editions, these are connected by interwiki links.

#### 2.1.1 History

*when was it founded, idea behind it, first steps, road to success, anecdotes*

### 2.1.2 Wikimedia Foundation

*Foundation structure, independence of language wikis, funding*

### 2.1.3 Editing policy

*Image: graphic of article UI*

Anyone with a browser and internet access can edit Wikipedia's articles<sup>1</sup>. In collaboration, people all over the world contribute and improve the content. Each edit creates a new revision of the article and is stored in the revision history. Naturally, each article available today started from an empty page and is the result of a succession of edits.

In the revision history each entry consists of the text change, the date of submission, the user and an optional comment explaining the change. Each revision can not only be examined by other users and but also reverted. To minimize the potential for *edit wars*[9] Wikipedia urges its users to discuss controversial topics on the article's talk page.

*Image: revision history*

Contributions to an article can be done anonymously or as a registered user. A registered user gains privileges like the ability to create articles or the use of the social network features in Wikipedia. With the initial registration a *user page* is created where the user is allowed to publish a profile and interact with other registered users.[13]

## 2.2 CONTRIBUTIONS

Wikipedia's articles are continuously edited by its users. The nature of an edit can range from simple spelling or grammar correction, over improving the content of a sentence to writing or removing of paragraphs or even whole articles. This collective authorship makes it difficult to determine an individual author's contributions, in other words, it is not easy to tell who wrote what.

Research in this area is motivated by finding ways to identify individual authors with a good reputation in order to assign a trust score to them. This is based on the assumption that trusted authors consistently produce high quality contributions that outlive contributions of lower quality. Kramer, Gregorowicz, and Iyer [14] devised a method to assign trust scores to the authors of an article by examining the wealth of information contained in the article's revision history. They looked at an article as be-

<sup>1</sup> Some articles can be locked because of sustained vandalism or content disputes.[12]

ing a set of phrases. The author who first wrote a sentence gets the credit for that phrase and will gain trust if it survives future edits.<sup>2</sup>

A similar approach of calculating the longevity of text chunks was followed by Adler and De Alfaro [15]. They adapted standard text-diff algorithms to the peculiarities of the wiki revision system, e.g. keeping track of text chunks that were removed at one point and then reinserted in later revisions. Based on these algorithms a reputation system was implemented by Adler et al. [16] which offers an API<sup>3</sup> that can annotate a Wikipedia article. The annotated text is the result of splitting the original text into chunks and attributing them with their respective authors, the number of the revision where the chunk was added and a trust value for the author.<sup>4</sup>

## 2.3 GEOREFERENCES

*"Combined approaches (i.e., where quantitative spatial analysis models are calibrated with surveyed locations) may prove useful." [17, p. 85]*

In order to analyze the localness of contributions, it is necessary to geotag them, i.e. applying geospatial metadata like coordinates to each contribution, derived from the author's location. In his doctoral thesis Hardy [17] used the Wikipedia corpora to study the spatial behavior of article production. The dataset was limited to anonymous users and articles that were geotagged.

For each anonymous contribution an IP address, belonging to the point of Internet access, is stored in the revision that is created. Various methods to determine the geographic location from a given IP address have been studied by Muir and Oorschot [18]. Various visualizations<sup>5 6</sup> of edit distributions use geolocation databases like MaxMind<sup>7</sup> and Quova<sup>8</sup>.

For registered users, the IP address is not stored with the revision. Therefor IP geolocation services cannot be used. Lieberman and Lin [19] found an interesting approach by assuming users prefer to edit geographic articles in their proximity. The

<sup>2</sup> Kramer, Gregorowicz, and Iyer [14] define a sentence as an n-gram—a sequence of n words—and use a sliding window model to follow it across revisions to prevent simple rearrangements of text from counting as a new sentence.

<sup>3</sup> <http://www.wikitrust.net/vandalism-api> (visited on 10/31/2011)

<sup>4</sup> Adler et al. also released a Firefox add-on that highlights untrustworthy passages when viewing Wikipedia articles: <https://addons.mozilla.org/en-US/firefox/addon/wikitrust/>.

<sup>5</sup> <http://infodisiac.com/blog/2011/05/wikipedia-edits-visualized/> (visited on 10/31/2011)

<sup>6</sup> <http://sonetlab.fbk.eu/wikitrip/> (visited on 10/31/2011)

<sup>7</sup> <http://www.maxmind.com> (visited on 10/31/2011)

<sup>8</sup> <http://www.quova.com> (visited on 10/31/2011)

approximated user location was derived from the center of the convex hull around those articles.

Registered users are also given the opportunity to create a personal profile in their *user page*. The user can choose prose or structured boxes to reveal information like his general interests, spoken languages, but also his location. When entity names can be extracted from location information they can lead to coordinates as shown by Hecht and Gergle [20].

## HYPOTHESES

---

To gain insight on how Wikipedia is being used during and after political events, and ultimately whether articles covering the events are written by people that are most affected, I will propose a set of hypotheses aimed at different aspects of article production.

For this thesis I will define an event as an occurrence that has an effect on the public affairs of a country. Let then an occurrence be marked by a location, a start date and, if not ongoing, an end date. Further, let the beginning of an event be defined as the first quarter of time between its start and its end<sup>1</sup>. In conjunction with the spatial distribution this division into time intervals allows for a detailed look into [Article creation](#), the level of [Participation](#) in following intervals as well as [Text survival](#).

### 3.1 ARTICLE CREATION

The use of Wikipedia concerning a political event starts with the creation of the article describing the event. Due to the website's popularity I would expect the delay between the start date and the date of article creation to be rather short. Moreover, considering the rising year-on-year Wikipedia usage in numbers of pages viewed[6], this delay should become shorter and shorter suggesting an increased use of Wikipedia as a news channel. This leads to the following hypotheses:

HYPOTHESIS 1. Articles are created with only a short delay after the start date of the event.

HYPOTHESIS 2. The more recent an article, the shorter is the delay between the event start and article creation.

A user has the chance to create a new article in any of the 260-odd language editions of Wikipedia. Although the English version is by far the biggest and most used, it would be interesting to see whether it is the prime choice to create the first article for a new event. Shortly after the first article has been created, articles covering the same topic will be created across various editions of Wikipedia. These articles are then being linked, mostly manually, via inter-wiki links.<sup>2</sup> When studying knowledge diversity

<sup>1</sup> For ongoing events the end date will be the date of the analysis.

<sup>2</sup> Adar, Skinner, and Weld [21] found that between two languages the inter-linking is not symmetrical, i.e. the number of out-links does not match the

across language editions, Hecht and Gergle [22] found that the English edition is not the superset of concepts of all editions as was previously believed. This means authors retain knowledge they consider important only to their compatriots. For a citizen of a country where English is not the first language creating an article becomes a political decision: should the author make the information available to fellow citizens or to a world-wide readership. However, picking English, even though it is not an official language of the country where the event is happening, would further support Wikipedia's role as a news channel, thus:

HYPOTHESIS 3. Articles are being created first in the English Wikipedia.

Regarding the localness of contributions, Hardy [17, p. 57] has established that Wikipedians write about places in their proximity more often than distant ones. His sample included only articles that have a geotag. Naturally, articles about geographic places like towns and sights<sup>3</sup> will dominate this sample. Since my thesis is concerned with political events, I find this point to be worth revisiting<sup>4</sup> for the people most affected should be the ones creating the article, thus:

HYPOTHESIS 4. Articles about political events are created by people in the events' proximity<sup>5</sup>.

### 3.2 PARTICIPATION

A Wikipedia article usually has more than one author. Once it has been created, users from around the globe can edit an article collectively. Viégas, Wattenberg, and Dave [24] tried to find patterns in the revision history that would reveal certain aspects of collaboration or the lack thereof, e.g. discussions and vandalism.<sup>6</sup> In respect to authorship, the researchers found the proportions of anonymous contributions differed strongly from page to page while showing no preference to any topic. This

---

number of in-links. Links are either missing on one side or the respective topics are not congruous and the user intentionally left out one direction.

<sup>3</sup> According to Kittur, Chi, and Suh [23] articles about *geography and places* are third biggest group.

<sup>4</sup> In addition, Hardy [17, p. 61] considered only anonymous users. Since creating an article is only allowed for registered users, his method has to be extended.

<sup>5</sup> Hardy [17, p. 57] defined proximity not in absolute terms, rather he considered the likeliness of authors being located less far than the average distance between an article and all its contributors.

<sup>6</sup> Using their history flow visualization Viégas, Wattenberg, and Dave [24] first identified patterns in single articles and later tried to statistically confirm their prevalence by analyzing the complete English corpus.



inconclusive result and the age of the sample<sup>7</sup> merits further investigation.

In 2007, Kittur et al. [25] found that a core of registered users is still doing the bulk of all edits. However, anonymous users contribute considerable amounts of text. For accounts of political events, due to their dynamic nature, I expect a strong participation by unregistered users while the events are still unfolding:

HYPOTHESIS 5. In the beginning of the event anonymous users contribute more than registered users.

HYPOTHESIS 6. For the duration of the event there are more local contributions than distant ones.

When the political event is considered “over” its end date in the article changes from “present” to a calendar date. In this unbounded and final phase I would expect the flood of contributions to subside and the content to consolidate when editors tighten prose or remove text they believe to be irrelevant. Looking at the whole lifespan of an article I would also expect registered users to outnumber anonymous ones as suggested by Kittur et al. [25] and the spatial distribution of contributions to become less local. Thus the final hypotheses:

HYPOTHESIS 7. Articles of a political event that has ended will continuously shrink in size.

HYPOTHESIS 8. After an event has ended, there will be more contributions from registered users than from anonymous ones.

HYPOTHESIS 9. After an event has ended, the spatial distribution of the contributors will become less local.

### 3.3 TEXT SURVIVAL

In 3.2 the contributions are only treated in volume giving credit to each contributor. However, when multiple authors write the same article, they do not only add text but also modify or even delete parts. A user who reads an article will only see the text that has survived all edits after it was added. Viégas, Wattenberg, and Dave [24] found that early contributions have a high survival rate. Recognizing this *first-mover advantage*, I suspect that accounts of political events show a strong localness in the beginning. Thus the key hypotheses from 3.2 have to be extended to reflect the spatial distributions of the contributions that make up the article:

<sup>7</sup> Viégas, Wattenberg, and Dave [24] used a dataset from May 2003.

HYPOTHESIS 10. For the duration of the event the article text contains more local contributions than distant ones.

HYPOTHESIS 11. After an event has ended, the spatial distribution of the surviving contributions will become less local.

### 3.4 SCOPE AND LIMITATIONS

#### 3.4.1 *Political events*

The same hypotheses may be applicable to other types of articles than political ones. The key requirements are that they have a location attribute and a time interval. This is easily fulfilled by disaster articles, e.g. Fukushima Daiichi nuclear disaster<sup>8</sup>.

#### 3.4.2 *Article location*

Although *location* is central to more abstract concepts like *Culture*<sup>9</sup> these subjects clearly defy being attributed with *a* location. Nevertheless, an analysis of the spatial distribution of contributors could be interesting.

#### 3.4.3 *Cross-language article growth*

The growth rates of articles covering the same topic across various language editions could be analyzed to further investigate issues like language barrier—locals contributing only in their language—and information arbitrage as suggested by Adar, Skinner, and Weld [21].

<sup>8</sup> [http://en.wikipedia.org/wiki/Fukushima\\_Daiichi\\_nuclear\\_disaster](http://en.wikipedia.org/wiki/Fukushima_Daiichi_nuclear_disaster) (visited on 10/31/2011)

<sup>9</sup> <http://en.wikipedia.org/wiki/Culture> (visited on 10/31/2011)

Part II

METHODS

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 [17, 20], 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*<sup>1</sup> in Berlin is tagged with the coordinates 52°30′58.58″N 13°22′39.80″E. However, these articles are 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*<sup>2</sup>, that clearly happened in *Egypt*<sup>3</sup>.

*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.: WikiTrust<sup>??</sup> oder Wiki-Watcher<sup>4</sup>

Die Methoden zur Datenextraktion und Visualisierung werden anschließend in eine Software integriert. Die Gewinnung

<sup>1</sup> [http://en.wikipedia.org/wiki/Brandenburg\\_Gate](http://en.wikipedia.org/wiki/Brandenburg_Gate) (visited on 10/31/2011)

<sup>2</sup> [http://en.wikipedia.org/wiki/Egyptian\\_Revolution\\_of\\_2011](http://en.wikipedia.org/wiki/Egyptian_Revolution_of_2011) (visited on 10/31/2011)

<sup>3</sup> <http://en.wikipedia.org/wiki/Egypt> (visited on 10/31/2011)

<sup>4</sup> 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>

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 verfügbar<sup>5</sup>. Eine solche Kopie enthält alle Artikel inklusive Versionsgeschichte und ist damit jedoch sehr groß<sup>6</sup>.

**ARTIKELEXPORT** 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,<sup>7</sup> welche einen direkten Zugang zu einer replizierten, schreibgeschützten 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 *Toolkits*

*Tools and servers to access the articles.*

- *Revert analyzer and framework to map a processing function on dumps:* <https://bitbucket.org/halfak/wikimedia-utilities/>
- *Page view statistics for a given article:* [http://stats.grok.se/json/en/201105/2011\\_Egyptian\\_Revolution\\_using\\_data\\_from http://dumps.wikimedia.org/other/pagecounts-raw/](http://stats.grok.se/json/en/201105/2011_Egyptian_Revolution_using_data_from_http://dumps.wikimedia.org/other/pagecounts-raw/)

#### 4.1.2 *Mediawiki API*

- *check if a user is a bot*
- *all articles that are member of a category*
- *revisions, includes user names*
- 
- 

<sup>5</sup> <http://dumps.wikimedia.org>

<sup>6</sup> Eine Kopie der englischen Wikipedia-Datenbank umfasst derzeit 5,4 Terabyte.

<sup>7</sup> <http://toolserver.org>

## 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 [17] 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.<sup>8</sup> Die registrierten Nutzer können jedoch auf ihrer *user page* Informationen über ihre Person entweder als Freitext oder strukturiert in *user boxes* veröffentlichen. Letztere sind 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*

---

<sup>8</sup> Eine kleine, von der Wikipedia-Community gewählte Nutzerschaft mit der Berechtigung *checkuser* kann die Adressen demaskieren.

- *Accuracy*
- *Active prevention by proxies and anonymizers:*  
J.A. Muir and P.C. van 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 verfügbaren<sup>9</sup> Online-Diensten wie *Quova*<sup>10</sup> oder *geo-plugin*<sup>11</sup> 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

- *Automatic annotation of entities:* <http://wikipedia-miner.cms.waikato.ac.nz/demos/annotate/>, also has services for categories. Alternative: <http://tagme.di.unipi.it/>
- *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?*

#### 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

<sup>9</sup> Die vorgestellten Dienste haben ein tägliches Kontingent an Anfragen. Hilfstechiken wie Caching können diese Einschränkungen jedoch mindern.

<sup>10</sup> <http://developer.quova.com>

<sup>11</sup> <http://www.geoplugin.com/webservices>

- 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 Blasendiagramm<sup>12</sup> 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

<sup>12</sup> [http://en.wikipedia.org/wiki/Motion\\_chart](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*

### 4.6 ANALYSIS

- *article usage (page views) vs. article production (edits)*
- 
- 
- 
-

# 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*<sup>1</sup> 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.

---

<sup>1</sup> [http://en.wikipedia.org/wiki/Category:Revolutions\\_by\\_country](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](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> (visited on 08/10/2011).
- [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> (visited on 08/10/2011).
- [12] Wikipedia. *Protection policy*. URL: [http://en.wikipedia.org/wiki/Wikipedia:Protection\\_policy](http://en.wikipedia.org/wiki/Wikipedia:Protection_policy) (visited on 11/16/2011).
- [13] Wikipedia. *Why create an account?* URL: [http://en.wikipedia.org/wiki/Wikipedia:Why\\_create\\_an\\_account%3F](http://en.wikipedia.org/wiki/Wikipedia:Why_create_an_account%3F) (visited on 08/10/2011).
- [14] 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.

- [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] B.T. Adler et al. „Assigning trust to wikipedia content.“ In: *Proceedings of the 4th International Symposium on Wikis*. ACM. 2008, pp. 1–12.
- [17] D. Hardy. „Volunteered geographic information in Wikipedia.“ PhD thesis. UNIVERSITY OF CALIFORNIA, SANTA BARBARA, 2011.
- [18] J.A. Muir and P.C.V. Oorschot. „Internet geolocation: Evasion and counterevasion.“ In: *ACM Computing Surveys (CSUR)* 42.1 (2009), p. 4.
- [19] M.D. Lieberman and J. Lin. „You are where you edit: Locating Wikipedia users through edit histories.“ In: *ICWSM'09* (2009), pp. 106–113.
- [20] 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.
- [21] E. Adar, M. Skinner, and D.S. Weld. „Information arbitrage across multi-lingual Wikipedia.“ In: *Proceedings of the second ACM international conference on Web search and data mining*. ACM. 2009, pp. 94–103.
- [22] B. Hecht and D. Gergle. „The Tower of Babel meets Web 2.0: User-generated content and its applications in a multi-lingual context.“ In: *Proceedings of the 28th international conference on Human factors in computing systems*. ACM. 2010, pp. 291–300.
- [23] A. Kittur, E.H. Chi, and B. Suh. „What’s in Wikipedia?: mapping topics and conflict using socially annotated category structure.“ In: *Proceedings of the 27th international conference on Human factors in computing systems*. ACM. 2009, pp. 1509–1512.
- [24] Fernanda B. Viégas, Martin Wattenberg, and Kushal Dave. „Studying cooperation and conflict between authors with history flow visualizations.“ In: *Proceedings of the SIGCHI conference on Human factors in computing systems*. CHI '04. Vienna, Austria: ACM, 2004, pp. 575–582. ISBN: 1-58113-702-8. DOI: <http://doi.acm.org/10.1145/985692.985765>.
- [25] A. Kittur et al. „Power of the few vs. wisdom of the crowd: Wikipedia and the rise of the bourgeoisie.“ In: *World Wide Web* 1.2 (2007), p. 19.



- [26] J.A. Muir and PC van Oorschot. *Internet geolocation and evasion*. Tech. rep. Citeseer, 2006.
- [27] M. Duckham and L. Kulik. „A formal model of obfuscation and negotiation for location privacy.“ In: *Pervasive Computing* (2005), pp. 152–170.
- [28] L. Xiao et al. „Information extraction from the web: System and techniques.“ In: *Applied Intelligence* 21.2 (2004), pp. 195–224.
- [29] O. Etzioni et al. „Unsupervised named-entity extraction from the web: An experimental study.“ In: *Artificial Intelligence* 165.1 (2005), pp. 91–134.
- [30] B. Snook et al. „On the complexity and accuracy of geographic profiling strategies.“ In: *Journal of Quantitative Criminology* 21.1 (2005), pp. 1–26.
- [31] 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.