# Reader Preferences and Behavior on Wikipedia

Janette Lehmann
Universitat Pompeu Fabra
Barcelona, Spain
lehmannj@acm.org

Claudia Müller-Birn
Freie Universität Berlin
Berlin, Germany
clmb@inf.fu-berlin.de

David Laniado
Barcelona Media
Barcelona, Spain
david.laniado@gmail.com

Mounia Lalmas
Yahoo Labs
London, UK
mounia@acm.org

Andreas Kaltenbrunner
Barcelona Media
Barcelona, Spain
kaltenbrunner@gmail.com

## ABSTRACT

Wikipedia is a collaboratively-edited online encyclopaedia that relies on thousands of editors to both contribute articles and maintain their quality. Over the last years, research has extensively investigated this group of users while another group of Wikipedia users, the readers, their preferences and their behavior have not been much studied. This paper makes this group and its activities *visible* and *valuable* to Wikipedia's editor community. We carried out a study on two datasets covering a 13-months period to obtain insights on users preferences and reading behavior in Wikipedia. We show that the most read articles do not necessarily correspond to those frequently edited, suggesting some degree of non-alignment between user reading preferences and author editing preferences. We also identified that popular and often edited articles are read according to four main patterns, and that how an article is read may change over time. We illustrate how this information can provide valuable insights to Wikipedia's editor community.

## Keywords

Wikipedia; reader; reading behavior; reading interest; article quality; editor; engagement

## Categories and Subject Descriptors

H.1.2 [**User/Machine Systems**]: Human information processing; H.5.3 [**Group and Organization Interfaces**]: Computer-supported cooperative work, Organizational design

## General Terms

Human Factors; Measurement;

## 1. INTRODUCTION

Peer-production communities have transformed the way people use and experience the Web. The collective action of these communities usually evolves around a digital artifact, such as an online encyclopedia or a piece of software. Wikipedia is a famous example of a peer-production community, and is the focus of our study.

Wikipedia is a multilingual, web-based, free encyclopedia, written collaboratively by a large number of volunteers. Since its creation in 2001, Wikipedia has grown into one of the most visited websites, attracting 530 million unique visitors monthly (October 2013).[1] As of November 2013, Wikipedia was available in 287 languages and comprised about 30 million articles. The English Wikipedia, the largest language version, had more than 30,653 active contributors[2] working on over 4.5 million articles. It was ranked as the 8th most popular website on the Internet in the US,[3] where popularity is measured by the number of page views.

Scholars have attributed the success of Wikipedia to its *production side*, that is the quality of its articles and authors' participation [12, 22, 32]. Thus, Wikipedia's production side has been the focus of numerous studies. A literature review by Okoli et al. [15] covering 477 research studies on Wikipedia showed that 42% of the studies mostly centered on issues related to participation, i.e. how editors create and edit articles, resolve disputes, or organize their community. Only 20% of the studies related to readers in Wikipedia, the *usage side* of Wikipedia, such as examining the popularity of articles or topics in Wikipedia. Less than 1% of the reviewed studies looked at users' reading preferences and only one study investigated reading behavior [15].

One reason for the limited focus on Wikipedia readers might be how scholars consider the role of passive users, i.e. the readers, in online communities. Readers are often considered to not provide any visible contribution to the community, and have been referred to as "lurkers" or "free-riders" who are "more resource-taking than value-adding" [13]. When scholars showed interests in this user group, it was mostly because reading is often seen as the prerequisite for becoming a contributor [14, 16, 17]. For example, Halfaker et al. [8] carried out several experiments to encourage Wikipedia readers to become contributors.

An exception is the work by Antin et al. [2], who claim that reading can be seen as a form of participation and is therefore valuable: the fact that a user is reading an article and not editing could be interpreted as an indication of an article's quality, such as its reliability [1]. Thus, reading activity – the usage side – can provide valuable insights to editors – the production side.

---

[1] http://reportcard.wmflabs.org/
[2] Registered (and signed in) users who made 5 or more edits in a month.
[3] http://www.comscore.com/Insights (for desktop access). The ranking is 9th when accounting for mobile access.

Other peer-production communities, such as open source software development projects, have included the usage side into their definition of success. They use measures that typically revolve around quantifications of volume related to the number of accesses to a particular project's product or outcome [5, 10].

Inspired by these perspectives, we conjecture that the same paradigm can be used in the context of Wikipedia. Instead of looking exclusively at the production side (the editors), we analyze the usage side (the readers) and discuss how our analysis can inform Wikipedia's production side. We explore *users' reading preference and behavior*, and Wikipedians' *editing preference*, which enables such a connection. Through examples, we demonstrate how readers can provide valuable insights to Wikipedia's editor community and that they are *not* resource-taking *but* value-adding.

But first, we review existing literature on reading preference and reading behavior and show that current knowledge is limited and rather exploratory.

## 2. RELATED WORK

Few studies about reading preference of users on Wikipedia exist. Spoerri et al. [20] examined readers' interests with respect to the topics they read about. The analysis, based on view count, showed that the most accessed articles were in the areas of entertainment (music, films, TV series), politics/history (politicians such as George W. Bush, historical events such as World War II), and geography (places such as Paris or countries such as USA). This aligns with the study reported by Waller et al. [26], who investigated search queries from Australians to Wikipedia. In general, people are more interested in "lighter" topics such as entertainment than in more "serious" or advanced topics. In this paper, we also show that readers in English Wikipedia have similar interests. However, a survey carried out on university students regarding the specific websites they have in mind when searching for information, reveals that 34% of the students would use Wikipedia for factual information and only 6% indicated thinking about Wikipedia when searching for entertainment related information [24].

Preference for both information searching and entertainment is a shared characteristic of Wikipedia readers and editors according to West et al. [27]. The authors leveraged data from a browser toolbar to investigate differences in Web usage between Wikipedia editors, readers and Internet users who did not access Wikipedia. They found that editors are "information-hungry" and "entertainment-loving", as they spend more time on news and search, but also on YouTube and other entertainment sites; Wikipedia readers' preferences are in a middle ground between those of editors and users not accessing Wikiepdia.

A comparison of reading behavior and editing activity in Wikipedia was performed by Reinoso et al. [19]. The authors compared for different language editions of Wikipedia the number of page views and the number of edits performed on them. For languages such as English, German and Spanish, the number of views and edits were highly correlated. This was not the case for Japanese and Dutch.

Reading behavior has been studied by Ratkiewicz et al. [18], who explored the dynamics of the popularity of Wikipedia topics. Popularity was defined as the number of hyperlinks linking to an article and the number of clicks to it. The authors found that almost all articles experience a burst just after their creation and the majority of articles receive little attention thereafter. Only few articles show intermittent bursts later in their lifetime. Ten Thij et al. [25] built a model to explain bursts in reading behavior caused by featuring an article on Wikipedia's main page.

Finally, two studies looked at how readers navigate within Wikipedia. Helic [9] analyzed users' click paths on Wikigame, where users must find the way (clicking links) from one randomly selected Wikipedia article to another. The author showed that users are very efficient at navigating; indeed users easily found short paths between the randomly selected articles. Gyllstrom et al. [7] investigated different browsing patterns on Wikipedia. They found out that users' browsing behavior depends more on the page topic than on the linking structure. They suggested that understanding different browsing strategies can help editors to better present or organize their content.

These studies demonstrate the still limited knowledge on user reading behavior on Wikipedia, in particular in relation to Wikipedia's production side. In the following sections, we carry out an analysis to gain insights on reading preference and behavior in Wikipedia, and discuss how these insights can add value to Wikipedia's peer-production side. We start by describing the datasets used in our work.

## 3. DATASETS

Our study is based on data collected over a period of 13 months (September 2011 to September 2012) from various sources for the English Wikipedia. In the first part of our analysis (Section 4.1), we use all Wikipedia articles to determine and study the most popular topics. To work on a more homogeneous dataset and avoid the effect of structural differences between different types of articles, we then focus for the rest of our analyses on a specific sub-set of articles – namely biography articles which contain descriptions of persons, such as actors, singers and historical figures. Biography articles form the most popular topic in Wikipedia. This approach was already followed in previous research [6]. To detect biographies, we considered all articles belonging to the Wikipedia category "Living people", as well as to the categories "Births by year" and "Deaths by year" and recursively to their subcategories. We then removed categories that did not contain biographies, and articles that were lists of biographies.

**Page view data.** As a measure of page popularity we use the *page view* data provided by the Wikimedia Foundation.[4] The dataset contains for each page in any Wikimedia project the number of requests per hour. We used this dataset for our study on reading preferences. For the 13-month period under consideration, we aggregated the hourly views for each month, to have monthly views for each article. The resulting dataset comprises a total of 4.3 million articles. The most visited page is the Main page, with 600 million page requests. Within this dataset we identified 1.02 million biography articles having 460 million page views in total.

**Browsing data.** Page popularity is only one criterion that can be considered when studying readers in Wikipedia. For example, accounting for the time spent on a page and the pages accessed during a visit on Wikipedia provides additional insights on reading activity. This information can be obtained from activity log data containing the entire navigation trace of users.

Since activity log data are not provided by the Wikimedia Foundation, we collected anonymized activity log data (tuples of browser cookie, URL, referral URL and timestamp) for a sample of users who gave their consent to provide browsing data through the Yahoo toolbar.[5] We identified in these browsing data users who have accessed the English Wikipedia by requesting for the following two types of URLs:

---

[4]http://dumps.wikimedia.org/other/
pagecounts-raw/
[5]A toolbar is a browser extension that provides additional functionalities and direct access to selected websites.

```
http://en.wikipedia.org/wiki/PAGE
http://en.wikipedia.org/w/index.php?title=PAGE
```

where $PAGE$ refers to the title of the page that was viewed. We identified these page titles in Wikipedia and resolved redirects to avoid duplicate entries.[6] We detected 288K biography articles, accessed by 387K users, and a total of 4.5M million clicks for our 13-month sample.

**Article characteristics.** To characterize Wikipedia articles from the editors' side, we computed their length and edit count. We retrieved these data through the Wikimedia Tool Labs [28]. Depending on the time window of our analysis (we used several), we computed for each article its text length (the size in bytes of the last revision of the article for the given time window) and number of edits (the number of revisions of the article during that time window).

To identify articles that have been considered of high quality by the community through its internal quality assessment system, we checked for each article whether it was included in the Wikipedia lists of Featured [29] or Good articles [30], or assigned as an A-class article [31] at the end of our 13-month period. These articles have been assessed by Wikipedia's editors using a set of pre-defined criteria developed over the course of the Wikipedia project, such as being well-written, comprehensive, and neutral. We found that 0.37% of the 1.02 million biography articles were assessed of high quality. 3% of these articles are A-class articles, 74% are good articles, and 23% are featured. In the rest of this paper, we refer to these articles as high quality articles ($HQA$).

## 4. READING PREFERENCE

In the first part of our study we look at the *reading preferences* of users on Wikipedia. First, we identify what are the most read topics on Wikipedia and show that articles belonging to the most popular topics do not necessarily correspond to those frequently edited by Wikipedia editors. In other words, reading preferences do not always align with *editor preferences*. We then characterize the difference between reader and editor preferences using a *preference* matrix. All studies in this section are based on the page view data provided by the Wikimedia Foundation.

### 4.1 Popular topics

In the first part of our analyses, we study the popularity of topics in Wikipedia. We select the 500 most read articles, measured by the number of article views over our data period.[7] We manually assigned a topic to each article using a three-round process.[8] In the first round, we collaboratively coded the articles (about 50) by using Wikipedia categories as reference point until we obtained an almost stable set of topics for these articles. In the second round, we separately coded the remaining articles. In the third round, we checked the assigned topics and discussed all ambiguous cases. To ensure a shared understanding of the existing topics, the second and third rounds were iterative. Newly introduced topics were cross-validated over the entire dataset. This process resulted into 12 distinct topics listed in the first column of Table 1. A description of the topics is provided in the fourth column.

---

[6] http://en.wikipedia.org/wiki/Wikipedia: Redirect

[7] We selected the 500 most read articles only, as we did not observe significant changes in our results by considering more articles.

[8] The hierarchical and overlapping structure of Wikipedia's category system prevent us to automatically determine the main category (the topic) of an article in a straightforward manner.

**Table 1: Article topics, percentages of articles in each topic, and percentage of high quality articles ($\%HQA$) in each topic for the 500 most popular articles (measured using page views).**

| Topics | %Articles | %HQA | Description |
|---|---|---|---|
| Biography | 44.2% | 31.2% | *Biographies of persons* |
|   Media personality | 18.8% | 24.5% | |
|   Musician | 11.6% | 37.9% | |
|   Sportsperson | 6.8% | 35.3% | |
|   Historical figure | 4.2% | 33.3% | |
|   Politic./businessp. | 1.8% | 33.3% | |
|   Criminal/victim | 0.4% | 0.0% | |
|   Misc | 0.4% | 50.0% | |
|   Publisher/writer | 0.2% | 100.0% | |
| Entertainment | 17.4% | 32.2% | *Cinema and TV* |
|   Series | 10.8% | 22.2% | |
|   Movie | 5.4% | 55.6% | |
|   Misc | 1.2% | 16.7% | |
| List | 7.6% | 0.0% | *"List of" articles* |
| Tech | 5.0% | 12.0% | *Web, software, electronics, etc.* |
| History | 4.4% | 22.7% | *Wars, monuments, incidents, etc.* |
| Misc | 3.8% | 15.8% | *Further articles* |
| Health | 3.4% | 23.5% | *Diseases, medicine, etc.* |
| Leisure | 3.2% | 18.8% | *Games, novels, etc.* |
| Sport | 3.0% | 66.7% | *Sports, sport events, etc.* |
| Places | 2.8% | 21.4% | *Regions, buildings, etc.* |
| Adult | 2.6% | 7.7% | *Articles about adult content* |
| Culture/Belief | 2.6% | 7.7% | *Religions, festivals, etc.* |

From Table 1, we see that a large percentage of users access Wikipedia to read about entertainment-related topics such as TV series, movies, and biographies of actors and singers. Articles related to history, health and tech content (such as web services and software) are also frequently accessed. This is in accordance with previous studies [21, 26].

The third column of Table 1 shows the percentage of high quality articles per topic. The lower the topic popularity, the smaller the number of high quality articles belonging to that topic. Indeed, we observe a Spearman's rank correlation coefficient of $\rho = 0.72$ ($p$-value $< 0.01$), suggesting a high correlation between topic popularity and the percentage of high quality articles. However, there are some exceptions. For instance, for the topics "Health" and "Sport", although the percentage of articles belonging to these topics is relatively low, many articles are of good or high quality (23.5%, and 66.7%, respectively). On the other hand, the percentage of high quality articles in the "Tech" area is low (12.0%), albeit this being the fourth most popular topic in our dataset.

These observations suggest some degree of non-alignment between *users'* reading preferences and *authors'*[9] editing preferences. To examine this further, we define several measures to characterize these two preferences next.

### 4.2 Reading and editing preferences

Table 1 provides a first indication of some non-alignment between reading and editing preferences in Wikipedia. In this section, we define various measures to study this.

**Measuring reading preferences.** In this paper, we define reading preference as the popularity of articles, measured by the number of page views. Previous studies suggest that popularity is a dynamic phenomenon that can partly be characterized by bursty behavior of page views [18, 25]. Our goal is to determine a value that best represents the popularity of an article by filtering out such bursty behavior. Thus, we calculate the monthly article popularity measured by the number of page views in each month from September 2011 to September 2012. Then, we measure the median rank of ar-

---

[9] In this paper, we use author and editor interchangeably.

**Table 2: Spearman's rank correlation coefficient $\rho$ between measures capturing reader and editor preferences.**

| | $ArticleLength_a$ | $\#Edits_a$ |
|---|---|---|
| $Popularity_a$ | 0.22 | 0.16 |

ticle popularity ($Popularity_a$) by their monthly popularity, which is less sensitive to outliers.

**Measuring editing preferences.** To determine editors' preferences, i.e. the articles they are mostly working on, we use three measures, each indicating a particular angle regarding editors' preferences. First, we employ the number of edits ($\#Edits_a$), a common measure of editing activity. For each article, we calculate the number of revisions over the whole period range. This measure, however, does not provide information about the effect of an edit, such as its informativeness and quality. We therefore propose to use article length ($ArticleLength_a$) as a measure for the informativeness of an article. The fact that an article is long suggests that a number of editors spent time and effort writing about the topic of the article, to make it more informative. We calculate the length of an article for a given time period using the latest version of the article in that period. Finally, editing may lead to the article being identified by the community as good, featured, or A-class (the pinnacle of the editing process). This would happen when the article is considered to provide comprehensive information on a topic.[10] We use the available data provided by Wikipedia – whether an article is a good or featured article, or belongs to the A-class articles ($HQA_a$) at the end of our data period – as a measure for article quality.

We compare reader and editor preferences by measuring the correlation between the reader preference measure $Popularity_a$ and the editor preference measures $ArticleLength_a$ and $\#Edits_a$. As discussed in Section 3, we focus on biography articles, which form the most popular article topic in Table 1. Table 2 reports the Spearman's rank correlation coefficient $\rho$ for the three metrics. We observe low correlations: 0.22 for $ArticleLength_a$, and 0.16 for $\#Edits_a$. These values suggest some non-alignment between reader and editor preferences. To further investigate this, we built a linear regression model using $ArticleLength_a$, $\#Edits_a$, and $HQA_a$ as features to predict the number of page views of an article. Our model predicted the number of page views with a coefficient of determination of $R^2 = 0.24$ ($R^2 = 1.0$ would represent a perfect fitting model), further indicating that readers and editors preferences diverge in many cases.

Next, we introduce a *preference* matrix, which allows visualizing the differences in reading and editing preferences using the above defined measures.

## 4.3 Preference matrix

For each article, we calculate its popularity (our reading preference measure) and its length (an editing preference measure).

The distributions of popularity values and article length values indicate whether articles are popular or not, and whether articles are long or short. We determine the upper and lower quartiles of both distributions since we want to identify articles with extreme values. We remove all articles that fall into the interquartile range of the article length or popularity distribution (the middle $25 - 75\%$ of both distributions). This means that we only consider articles that differ significantly from those having an average length or popularity.

[10]See criteria for featured articles http://en.wikipedia.org/wiki/Wikipedia:Featured_article_criteria



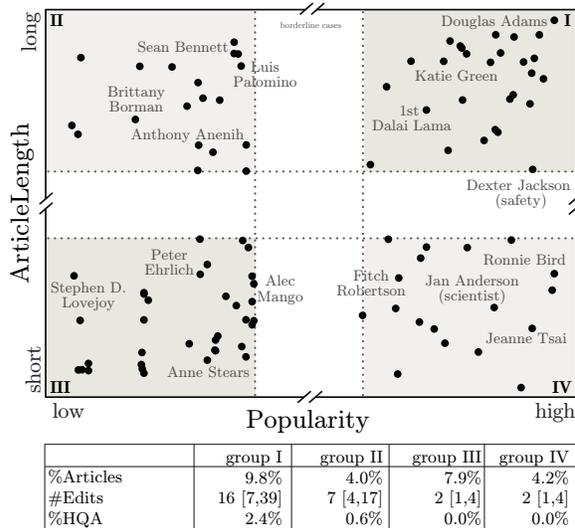| | group I | group II | group III | group IV |
|---|---|---|---|---|
| %Articles | 9.8% | 4.0% | 7.9% | 4.2% |
| #Edits | 16 [7,39] | 7 [4,17] | 2 [1,4] | 2 [1,4] |
| %HQA | 2.4% | 0.6% | 0.0% | 0.0% |

**Figure 1: (Top) Preference matrix defined by article popularity and length. (Bottom) Percentage of articles belonging to a group ($\%Articles$), median and interquartile range of the number of edits per group over the whole data range ($\#Edits$), and percentage of high quality articles in each group ($\%HQA$).**

This results in the four groups of articles shown in Figure 1. The horizontal axis represents article popularity (the reading preference) and the vertical axis represents article length (one of the measures characterizing editing preference). The values of both measures are transformed into an ordinal scale to overcome scaling issues, i.e. we ranked all values for article popularity and article length. Each dot in the matrix represents an article and the position corresponds to its popularity and length. We only show a random sample of 100 articles in Figure 1 to improve legibility.

Under the preference matrix, we report the percentage of articles belonging to each group, and the other two editing preference measures, namely, the percentage of high quality articles, and the median and interquartile range of the number of edits. We see that featured articles tend to be long, confirming previous work [33] and suggesting a relationship between article length and article quality.

Many articles belong to group $I$ (9.8%) and group $III$ (7.9%). Whereas group $I$ contains very long and often read articles, articles in group $III$ are short and seldom read. In both groups, we have articles for which editing and reading preferences align.

A divergence between reader and editor preferences can be observed for articles belonging to groups $II$ and $IV$. Group $II$ articles (4% of all articles) tend to be not read very often, even though they are very long (probably very informative). This group also contains a low number of high quality articles. For instance, it contains the biographies of the Nigerian politician "Anthony Anenih" and the American football player "Sean Bennett". We speculate that not many users read these articles because the person in question is not popular nowadays (e.g., former American football player) or is of interest only to a specific user community (e.g., users interested in Nigerian politics). This is further accentuated by the lower edit activity in this group (median of 7 edits) compared to group $I$ (median of 16 edits). In fact, many of the articles forming group $II$ are on topics that were popular in the past and heavily edited during that time.
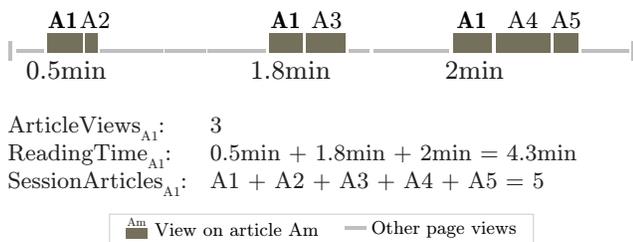
ArticleViews$_{A1}$:    3
ReadingTime$_{A1}$:    0.5min + 1.8min + 2min = 4.3min
SessionArticles$_{A1}$:    A1 + A2 + A3 + A4 + A5 = 5

**Figure 2: Example of a reading session on Wikipedia showing the reading behavior measures for article $A1$.**

Finally, $4.2\%$ of all articles belong to group $IV$. For these articles as well, reader and editor preferences do not align. Even though articles are regularly accessed by readers, they are short (and have seldom been edited) and none of them is of high quality. Taking the examples of "Jan Anderson (scientist)" and "Ronnie Bird", we see that these articles are often viewed, but are short and have hardly been edited during the last 13 months (median of 3 edits per article). Additionally, none of these articles is considered to be of high quality, even though readers access them very often.

To summarize, we observe differences between what readers access and what editors work on. The most edited articles tend to be long (groups $I$ and $II$) and the number of high quality articles in these groups is higher compared to the other two groups. However, only articles in group $I$ are very popular, suggesting that article quality does not drive popularity.

The opposite can be observed for articles in groups $III$ and $IV$. These groups contain shorter articles, and fewer high quality articles. Moreover, articles in these groups tend to be edited less. This indicates that editors rarely added content to them in the past, reflecting low interests in these articles. Whereas articles in group $III$ neither meet authors nor readers interests, we can see that readers are interested in articles of group $IV$ despite the scarce attention these receive from editors.

Next, we analyze how users read articles during their browsing sessions, and how this matches with the editing activity.

## 5. READING BEHAVIOR

Here we study how users read Wikipedia articles. That is, we look at the biographies that users read when *visiting* Wikipedia. We introduce three measures to characterize patterns of reading behavior: $ArticleViews_a$, $ReadingTime_a$ and $SessionArticles_a$.

**Reading session.** We separated the reading activity of users by sessions, where a session is a sequence of pages visited by a user until he or she goes offline. Following [3], a user is said to have gone offline – meaning that the session has ended – if more than 30 minutes have elapsed between two successive activities of that user. All reading activities of a user on Wikipedia during this session form what we refer to as a *reading session*.

During a reading session, a user spends time reading an article $a$. We use $ReadingTime_a$ to refer to the time spent on article $a$. This user may return to article $a$ several times during a session, and visit several other articles. To capture these, we define two additional measures: $ArticleViews_a$ is the number of times the article $a$ was viewed and $SessionArticles_a$ is the number of articles viewed during the reading session where article $a$ was read.

In Figure 2 we show an example of a reading session of a user. The figure depicts how the user visits articles on Wikipedia as well as several other webpages. The starting point on Wikipedia is arti-

cle $A1$. After reading this article, the user follows a link to article $A2$. Then, he or she probably clicks on an external link to navigate to other websites. After a while, the user returns to article $A1$ and then uses a link to reach article $A3$. This pattern occurs another time during this session. At the bottom of Figure 2 the values for the aforementioned three measures are given.

**Data processing.** We use the browsing data as it enables us to access the readers' entire navigation traces (Wikipedia articles and other webpages) during their browsing sessions. To have a more homogeneous and robust dataset, we discarded articles with lower values of length or popularity, and focused our analysis on articles belonging to group $I$ of the preference matrix (see Figure 1), which contains the large majority of articles in our browsing data ($83.47\%$). These articles allow for a reliable interpretation of any observed difference between reading interests and editing preferences since their length and popularity are high enough.

We characterize the reading behavior of an article $a$ by calculating per month the average of $ArticleViews_a$, $ReadingTime_a$ and $SessionArticles_a$. We also calculate $Popularity_a$, the popularity measure defined in the previous section. Therefore, for each article $a$ we obtain 13 vectors, one for each month of the 13-month period. We refer to each vector ($ArticleViews_a$, $ReadingTime_a$, $SessionArticles_a$, $Popularity_a$) as a *behavior vector*.

We generate behavior vectors of an article for the months where it was visited in at least 10 reading sessions. This enables us to derive stable values for the three measures calculated based on reading sessions. This results into 9,726 articles and 49,921 behavior vectors. To ensure that the two datasets (page view and browsing data) are comparable, i.e., no strong bias in the browsing data is influencing our results, we ranked the articles according to their overall popularity in both datasets, and found that their rankings correlate (Spearman's rank correlation coefficient was of 0.64).

## 5.1 Reading patterns

We use the k-means algorithm to cluster the behavior vectors. Since our dataset does not follow a normal distribution, and thus to avoid the extensive influence of heavy outliers, we do not use the value of each measure, but the corresponding article rank. The number of clusters is determined by a minimal cluster size such that each cluster contains at least $20\%$ of the 49,921 behavior vectors. Since the clustering is performed with the behavior vectors of the articles, an article can occur in multiple clusters. This allows us to analyze changes in the reading pattern of an article across the 13-month period; we return to this in Section 5.2.

We obtain four clusters, shown in Figure 3, each corresponding to a pattern of reading behavior. The first row displays the name given to each pattern. The second row contains the cluster centers normalized by the z-score. Each bar corresponds to one measure. The vertical axis shows how many standard deviations a rank value is above or below the mean rank, on average. This means that bars above zero indicate higher ranks for the respective measure whereas bars below zero indicate lower ranks. The third row contains the number of articles and behavior vectors within each cluster. Since the sizes of the clusters are similar, there is no dominant reading pattern. As some subtopics are more predominant (e.g., there are more articles about media personalities and musicians), the fourth row shows the percentage of articles of the various biography subtopics (e.g., musicians, sportspersons) that belong to a given cluster.[11] We only show the subtopics with the largest percentage. The last three rows of Figure 3 report the values

[11]For each cluster, we sampled at random a subset of 500 articles, and determined the sub-categories of these articles by using the three-round process described in Section 4.1. We manually catego-

| Focus | Trending | Exploration | Passing |
|---|---|---|---|
| | | | |



|  | Focus | Trending | Exploration | Passing |
|---|---|---|---|---|
| Reading behavior | | | | |
| | 5,278 articles<br>10,605 behavior vectors | 3,876 articles<br>14,267 behavior vectors | 4,826 articles<br>11,579 behavior vectors | 5,684 articles<br>13,470 behavior vectors |

Article topic

| | CA | | CA | | CA | | CA |
|---|---|---|---|---|---|---|---|
| artist/writer | 43% | historical fig. | 42% | sportsperson | 28% | media pers. | 27% |
| historical fig. | 41% | criminal/victim | 38% | musician | 26% | sportsperson | 27% |
| polit./businessp. | 37% | musican | 38% | media pers. | 23% | musician | 19% |

| | Focus | Trending | Exploration | Passing |
|---|---|---|---|---|
| ArtLen | 28K [16K,51K] | 38K [21K,69K] | 26K [15K,45K] | 16K [10K,27K] |
| #Edits | 11 [5,23] | 20 [9,41] | 10 [5,21] | 8 [3,18] |
| %HQA | 7.7% | 16.9% | 10.5% | 5.1% |

■ $Popularity_a$   ■ $ArticleViews_a$   ■ $ReadingTime_a$   ☐ $SessionArticles_a$

CA - Percentage in topic    ArtLen - Article length    #Edits - Number of edits    %HQA - Percentage of high quality articles
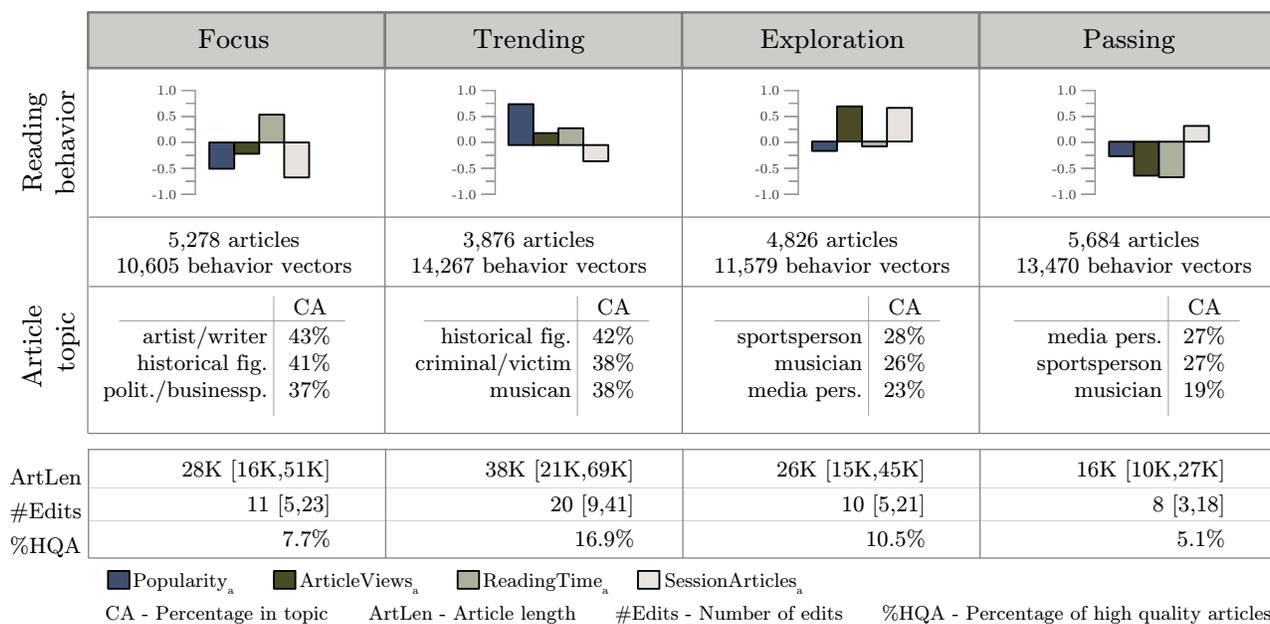
**Figure 3: (1st row) Article clusters and reading characteristics. (2nd row) Number of articles and behavior vectors per cluster. (3rd row) Most dominant article topics per cluster. (4th row) Median and interquartile range of the article length ($ArtLen$), number of edits per cluster (#$Edits$), and percentage of high quality articles in each group (%$HQA$).**

of the three editing preference measures. For each behavior vector, we calculated the length of the corresponding article, using the latest revision of the article for the given month, and the number of edits made during the month (we report median and interquartile range). We also determined the percentage of high quality articles.

We discuss now each of the identified reading patterns and relate them to the editing preferences. The patterns "Focus", "Trending", and "Exploration" are what content portals aspire to: users spending time reading their articles and/or reading many articles.

**Focus.** Articles following this pattern are characterized by an expected encyclopedic reading behavior: people spend a lot of time reading the article (high $ReadingTime_a$), but access very few other articles (low value of $SessionArticles_a$) within the session. Users have a specific information need (e.g., they want to learn something about "Jacques Cousteau"). Articles in this cluster have a lower than average popularity, and are about artists/writers, historical figures, and politicians/businesspersons.

The high reading time indicates a strong interest in the content of the article. Hence, we would expect many of these articles to be marked as good, featured, or A-class, as the quality of these articles seems important. However, the percentage of high quality articles in "Focus" (%$HQA = 7.7\%$) is lower than for the "Trending" and "Exploration" clusters. Moreover, although we observe an appropriate article length ($ArtLen = 28K$), the number of edits (#$Edits = 11$) suggests that editors are not interested in improving these articles. Indeed, the article about "Jacques Cousteau" is long (a median of $30K$ characters), but it is neither featured nor good nor A-class, and the number of edits is low (a median of $5.5$ edits per month).

**Trending.** Many biographies about historical figures, musicians and criminals/victims follow this pattern: articles are visited very often (high $Popularity_a$). Users read only a few other articles

rized all articles based on the subtopics of the category biography as shown in Table 1.

(low $SessionArticles_a$), similarly to the "Focus" reading pattern, but they spend less time reading the articles. This suggests that users are probably "quickly looking up" for information about something that is currently trending or has recently happened. For example, users read about the politician "Ron Paul" when he was a candidate for the presidency of the United States, but only to catch up on any recent news about him.

"Trending" articles exhibit the highest edit activity and the highest percentage of high quality articles compared to the other two clusters (#$Edits = 20$ and %$HQA = 16.9\%$). These articles not only attract users to read them but also authors to edit them, which is in accordance with a previous study by Reinoso [19], and also aligns with the work from Keegan et al. [11] about breaking news and current events in Wikipedia. The high percentage of high quality articles suggests that editors do not only work on the articles to increase the quality, but also to "update" information caused by recent or continuous events related to the article topic. Indeed, we saw in our dataset that featured articles are also edited frequently (a median of 19 edits per month). Featured articles are usually only changed in case new information becomes available.

Returning to our previous example, the politician "Ron Paul", we observed a median of 81 edits per months during the time the article was trending (December 2011 until May 2012). In the other months (when "Ron Paul" was not competing for the presidential primaries), the article belonged to the "Focus" cluster and had only 20 edits per month. We return to this later in this section.

**Exploration.** This pattern primarily contains biographies describing sportspersons, musicians, and media personalities that have an average popularity. The number of articles viewed in a session ($SessionArticles_a$) is the highest compared to the other clusters, indicating that users explore many other articles in a reading session. Looking into the articles that were visited, we saw that articles requested during the same session belong mostly to
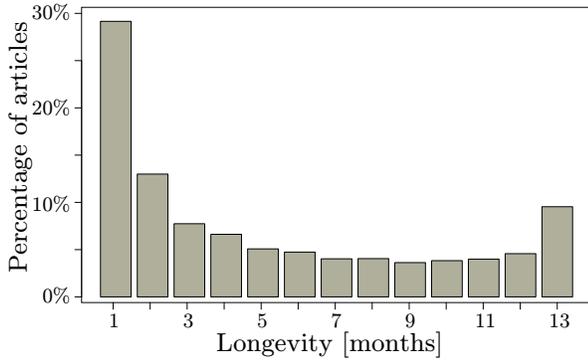
**Figure 4: Stability of articles:** *longevity* **is the number of months in which an article was visited in at least 10 reading sessions. For each longevity value, we plot the percentage of articles with that value.**



**Figure 5: Stability of articles:** *membership* **is the average fraction of the months of its longevity an article remains in its "home" (i.e. most frequent) cluster. For each longevity value, we plot the mean and standard deviation of the membership values of the articles.**

the same topic[12] (e.g., users who read the article about the actor "Al Pacino" also read articles about his movies). The high value of $ArticleViews_a$ indicates that users return regularly to the article under consideration, suggesting that they use it as a basis to navigate to other articles on the same topic. This hypothesis is supported by the low reading time of the focal article.

The editing preferences are comparable to the "Focus" pattern, in terms of number of edits (moderate values of 10 edits per article) and article length. The difference between the "Focus" and the "Exploration" reading patterns may be explained by external factors that influence the consumption of online content by users, such as the death of a famous artist [18].

**Passing.** Many biography articles about sportspersons, musicians, and media personalities belong to this cluster. The number of articles viewed in a session ($SessionArticles_a$) is above average, suggesting that users read different articles. Users browse many articles in the same session, but in contrast to "Exploration" they seem to only pass through the focal article (low $ReadingTime_a$), and do not return to it (low $ArticleViews_a$).

An example is the article about "Jackie Jackson", member of "The Jackson 5". When users are reading about "The Jacksons", they also view this article, but then quickly move to other related articles. The question is whether users do not spend much time on the article, because they are not interested in reading more about "Jackie Jackson", or because there is not much information provided about her (her article has a median text length of 9K).

Indeed, compared to the other clusters, the "Passing" cluster has a lower percentage of high quality articles ($\%HQA = 5.1\%$), and has shorter articles ($ArtLen = 16K$) and the lowest number of edits per article (a median of 8).

To summarize, we observe that articles exhibit different reading patterns. These seem to be mainly driven by the topics of the articles and therefore the interests of users, and less by their quality. Thereby, users show their interest in an article in different ways, e.g., by exploring also related articles ("Exploration" cluster) or by spending time reading the article ("Focus" cluster). Sometimes, the interest in an article is driven by external factors, as shown with

---

[12]We extracted all wikilinks between the articles in each reading session and found that on average over 76% of the articles visited in a session are connected to one another. This applies even for long reading sessions containing more than 10 articles (the average becomes 70%).
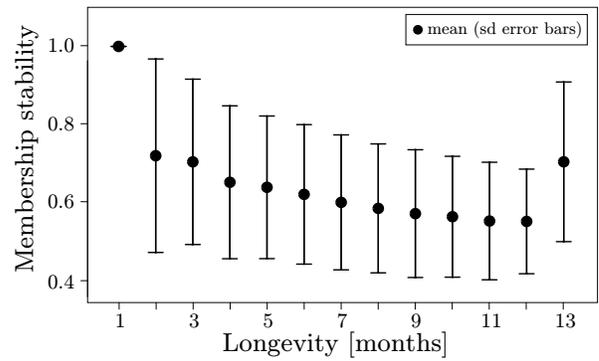
articles belonging to the "Trending" cluster (e.g., users read biographies about currently trending persons). On the other hand, for articles belonging to the "Passing" cluster, the question is whether the reading behavior is partly caused by a lower quality of the articles. In future work, we will investigate this. Overall, our results show that popularity and reading time are not the only factors that should be taken into account when measuring user reading preferences – how often users return to an article and how many other related articles they read provide further information about their reading preferences.

Three out of the four clusters constitute reading patterns where users are interested in the articles they are reading. However, editors seem to focus on articles in mainly one cluster, the "Trending cluster". The editing activity, the article length and the percentage of high quality articles is higher in that cluster, than in the "Focus" and "Exploration" cluster. This shows again a non-alignment between reader and editor preferences.

Finally, as shown in our example of "Ron Paul", an article can be in several clusters, depending on the month under consideration. That is, articles can transition between patterns across the 13-month period. We study this next.

## 5.2 Changes in reading patterns

The analysis conducted in the previous section used measures calculated on a monthly basis (the behavior vectors) to identify reading patterns. As a result, articles can belong to more than one cluster. In this section, we use this fact to study how articles might move (if they do) between reading patterns, and discuss possible reasons for these transitions. First, we determine how stable articles are in terms of their popularity and the way they are read across the 13-month period. We then look at typical transitions between reading patterns.

**Stability.** We calculate the number of months in which an article was visited in at least 10 reading sessions. We refer to this as the article *longevity*, denoted $Longevity_a$ for article $a$. In Figure 4, we plot on the x-axis the longevity values and on the y-axis the percentage of articles for a given longevity value. Almost 30% of the articles $(2,836)$ have a longevity value of 1, meaning that these articles have been accessed in at least 10 reading sessions only in a single one-month period. Another 13% of the articles $(1,264)$ have been accessed in at least 10 reading sessions in two different months. This percentage decreases continuously for larger numbers

of months, but increases again for 11 and more months. About 10% (928) of articles are read at least 10 times a month over the whole 13-month period. This suggests that there are articles that are frequently accessed over a long time period.

We examine the stability of an article in terms of which clusters it belongs to (i.e., the reading patterns it exhibits). We calculate the number of months an article $a$ remains in its "home" cluster, which is the predominant reading pattern exhibited by the article. Then, we normalize this value by dividing it with the corresponding $Longevity_a$ value. We refer to this as the article *membership* stability, denoted $Membership_a$ for article $a$. In Figure 5, the y-axis shows the average and standard deviation of the *membership* values for all articles for a given longevity. For example, an article with a longevity of 3 (the article was visited in at least 10 reading sessions during three, not necessarily consecutive, months) has a membership value of 0.7 on average. This means that on average the article was read 70% of its lifetime according to its most frequent reading pattern.

Figure 5 suggests that the higher an article's longevity, the lower its membership stability. This means that the longer - in terms of months - the article is accessed frequently, the higher the probability that its reading pattern changes. However, the average membership stability values are always above 0.5, indicating that many articles remain in their "home" cluster for at least 50% of their lifetime. It is interesting to note that the membership stability increases again for articles with a longevity value of 13. This means that high longevity implies high membership stability.

Next, we look at changes of reading patterns of articles and explore possible reasons behind these changes.

**Transitions.** We study the most frequent changes, i.e., transitions, between reading patterns (clusters). A $Transition_a$ exists for an article $a$ if one behavior vector of $a$ belongs to cluster $C$ at month $m$ and another behavior vector belongs to cluster $D$ at month $m + 1$, where the clusters represent two distinct reading patterns. We selected two cases to explore transitions. We consider all articles and then only articles with a $Longevity_a$ value of 13 – the set of highly stable articles in terms of their monthly access rate.

In Figure 6, we visualize the transitions between the four clusters by two networks, one for each case. Each vertex represents one cluster (i.e., reading pattern) and the size of a vertex corresponds to the number of articles in that cluster. The undirected edges in the network depict the transitions between the clusters. We use an undirected network since we observed a similar number of transitions in both directions. The largest difference we observed is smaller than 2.0%, which can be explained by the fact that an article usually belongs to one cluster (e.g., "Exploration"), moves to another cluster for a short time (e.g., to "Trending" because something happened with the person under consideration), and then moves back to the original cluster.

Each edge has a weight, which is the percentual amount of transitions between two clusters; for example, an edge weight of 23% means that 23% of all transitions in the network take place between these two clusters.

The complete network (left side of Figure 6) show how external factors, such as recent or continuous events related to a person, drive changes in reading patterns. This is the case for example for the biography article of the Facebook co-founder "Chris Hughes". Before March 2012, users tended to "pass by" this article (when reading about Facebook). However, this changed in March 2012 when Chris Hughes became the owner of the "The New Republic" magazine, attracting some media attention. Users started reading this article in a more "explorative" manner, using it as a starting point to access other articles related to the person.



**All articles**
9,726 articles
49,921 behavior vectors.

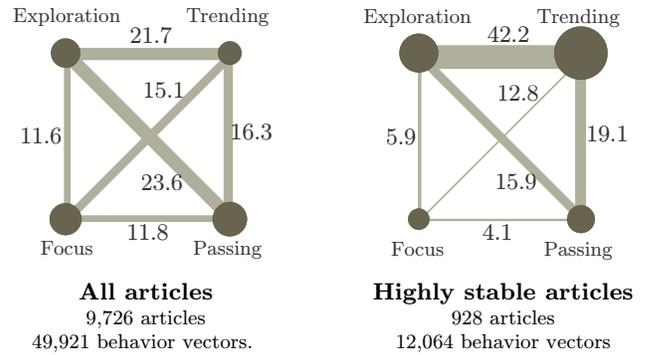**Highly stable articles**
928 articles
12,064 behavior vectors

**Figure 6: Transitions between reading patterns considering all articles (left) and only articles with a high stability (right). The vertex size represents the number of articles that belong to the cluster, and the edge weight represents the percentage of transitions between two clusters.**

The edge weights differ a lot in the network. We see a strong connection between the "Passing", "Exploration", and "Trending" clusters, indicating that many articles adopt all three reading patterns and sway between clusters. A transition can be even long-lasting, as in the case of the article "Jacqueline Kennedy Onassis". Until April 2012, the article was in the "Trending" cluster, but then lost its popularity and moved to the "Exploration" cluster. We assume that the article started to trend when her audio tapes, recorded after her husband's assassination, were released.

We observe also that articles belonging to the "Focus" cluster are isolated – the likelihood that an article is moving from or to the "Focus" cluster is low. Articles that are read in this way can be considered as the most stable ones, as their reading behavior hardly varies. An example of such an article is "Franklin D. Roosevelt".

Looking at stable articles only (right side of Figure 6) we see a different pattern. Compared to the network comprising all articles, we observe that "Focus" becomes even more isolated, showing again its special characteristic – a constant reading pattern. The transitions between "Exploration" and "Trending" become stronger, whereas the transitions between "Exploration" and "Passing" become weaker. Returning to Figure 3 we see that the "Exploration" and "Trending" clusters contain the same type of biographies: media personalities, musicians, and sportspersons. We conclude that these two clusters indeed describe reading patterns for the same type of articles.

## 6. DISCUSSION

We analyzed readers' preferences and reading behavior. We did so by connecting them to editors' preferences, allowing us to relate the usage side of Wikipedia to its production side. Our goal was to provide insights about how the reading experience and the editing process on Wikipedia could be enhanced. We discuss now our main results, position them in light of our goal, and present examples of potential applications.

Using the page view data provided by the Wikimedia Foundation, we studied reading preferences of users on Wikipedia. Our results confirm other works showing the dominance of *entertainment-related topics* among the most read topics on Wikipedia [21, 26]. The encyclopedic character of Wikipedia does not exempt it from following the known prominence of consuming and interacting with entertainement-related content observed on the Web.

We then introduced a preference matrix, which enabled us to differentiate four groups of reading versus editing preferences. These groups provide valuable insights into Wikipedia's quality system, in particular groups $II$ and $IV$, where the preferences do not align.

Group $II$ articles are often edited but not often read, whereas group $IV$ contains articles that are popular, but hardly looked at by editors. Being aware of these divergences can help Wikipedia editors making an *informed* decision about which articles to focus next. As opposed to tools such as WikiDashboard [23], which allows readers to evaluate article quality on the basis of an author history, the preference matrix can provide editors with a visualization of user reading preferences. This might draw their attention to articles or topics they have not edited before. Moreover, task recommendation services, such as the SuggestBot [4], could use the preference matrix as input to recommending tasks.

In the second part of our work, we studied the reading behavior of users and identified four main reading patterns: "Focus", "Trending", "Exploration" and "Passing". Information about the reading behavior of users can be useful in many ways, such as for the selection of articles for the main page, or the Article Feedback Tool (AFT).[13] Knowing which articles follow, for instance, the "Focus" pattern might help making the Article Feedback Tool more efficient since using this tool over the entire Wikipedia corpus failed. Editors complained about the low quality of the feedback made on articles. The fact that an article is often read and users spent time on it may indicate that users are interested in the article. As such, their feedback (if any) is likely to be more constructive and valuable.

In conformity with the work from Gyllstrom et al. [7], we showed that the reading behavior depends less on the article quality, but more on the article topic and therefore the interests of the reader. The quality of articles does not greatly influence what users choose to read. In general, the editing activity and the quality of articles reflect mostly the authors' interests and not the readers' interests. The exception to this are articles belonging to the "Trending" pattern, which are both accessed by many users and edited by many authors, compared to the three other reading patterns.

Understanding reader preferences and behavior can support editors in their work in several ways. The identified browsing patterns provide information about which articles the readers are interested in and how these and related content are read. If readers are interested in an article topic, they tend to look up information ("Trending"), spend a lot of time in reading the article ("Focus"), or consume the article content, but also related information ("Exploration").

This information can be used to improve the structure and presentation of the article content. For instance, the "Exploration" pattern corresponds to a navigation "way" to consume Wikipedia content. One article is focal, but also acts as a source to explore other articles. Knowing that these articles are consumed in this way, Wikipedia editors may add more links, keeping the users engaged by providing additional and relevant content. From an interface design perspective, navigation tools could be provided to guide users with the aim to enhance their reading experience.

Additionally, reading pattern can help editors to decide which articles to edit next. For instance, for the "Focus" pattern, we observed the highest reading time per article compared to all other patterns, but the proportion of high quality articles is lower than in the "Trending" and "Exploration" cluster. With respect to Wikipedia's production side, articles following the "Focus" pattern may greatly benefit from improvement in their quality, as users are very interested in them.

Also articles belonging to the "Passing" cluster may benefit from improvements. We assumed that users are not interested in the article, and therefore only pass through the focal article during their reading session. Another explanation is that these articles are not very informative (they are often short), and have rarely (likely as a consequence) been marked as good or chosen to be featured. In future work, we will investigate this hypothesis.

Finally, we looked at the stability of the reading patterns. We found that many articles are stable (remain in the same cluster), and that changes of the reading patterns are of temporal nature (e.g., in case of an event) or due to the time passing (e.g., interest in the person is decreasing). Studying the transitions between the reading patterns revealed two main findings. First, we observed a strong connection between the "Exploration" and "Trending" clusters, indicating that many articles adopt both reading patterns. Second, we observed that the "Focus" cluster represents a reading pattern that is isolated from the others. Articles in this cluster usually do not change their reading pattern. It indicates that this pattern represents articles with a high stability.

The above observations can inform the Wikipedia editor community in two ways. The stability of articles allows them to make long-lasting decisions for their editorial work. For instance, when adapting an article for explorative reading, these adaptions, such as adding links, are useful for the consumption of that article later on. On the other hand, transitions between reading patterns inform editors about recent trends (e.g., when an article is moving from "Passing" to "Exploration", indicating an increased interest from the reader side). Such articles can be candidates to be placed on the front page to raise awareness.

In future research, we will look more closely into the structural aspects of the article and how these may affect user reading behavior. It will also be important to look at different topics (beyond biographies) and how a specific reading behavior might depend on the topic.

# 7. CONCLUDING REMARKS

This paper provides new insights about how users consume content on Wikipedia: their reading preferences and behaviour. This paper also attempts to connect Wikipedia's readers (usage side) and Wikipedia's editors (production side). Using several measures to characterize reading preferences and behavior, we learn how users consume Wikipedia content, and illustrate how this information could inform Wikipedia editors about their editing tasks, for instance which articles to prioritize and why.

Identifying how an article is read can be used to determine which articles are more "engaging" than others, for instance, as measured by the average time spent on the article or the number of articles accessed from it. Articles that are more engaging are likely to promote a successful reading experience and even encourage users to return to them or to other articles. Readers that regularly return to Wikipedia are more likely to recognize the effort of Wikipedia's community and might even develop a sense of belonging to that community [17]. This in itself may further engage Wikipedia editors as they feel that their work is recognized and appreciated.

# 8. ACKNOWLEDGMENTS

---

[13] http://en.wikipedia.org/wiki/Wikipedia:
Article_Feedback_Tool

# 9. REFERENCES

[1] B. T. Adler, K. Chatterjee, L. de Alfaro, M. Faella, I. Pye, and V. Raman. Assigning trust to Wikipedia content. In *Proc. WikiSym*, 2008.

[2] J. Antin and C. Cheshire. Readers are not free-riders: reading as a form of participation on Wikipedia. In *Proc. CSCW*, 2010.

[3] L. D. Catledge and J. E. Pitkow. Characterizing browsing strategies in the World-Wide Web. In *Proc. WWW*, 1995.

[4] D. Cosley, D. Frankowski, L. Terveen, and J. Riedl. Suggestbot: using intelligent task routing to help people find work in wikipedia. In *Proc. IUI*, 2007.

[5] K. Crowston, H. Annabi, J. Howison, and C. Masango. Towards a portfolio of floss project success measures. In *Proc. ICSE*, 2004.

[6] L. Flekova, O. Ferschke, and I. Gurevych. What Makes a Good Biography? Multidimensional Quality Analysis Based on Wikipedia Article Feedback Data. In *Proc. WWW*, 2014.

[7] K. Gyllstrom and M.-F. Moens. Surfin' Wikipedia: An Analysis of the Wikipedia (Non-random) Surfer's Behavior from Aggregate Access Data. In *Proc. IIIX*, 2012.

[8] A. Halfaker, O. Keyes, and D. Taraborelli. Making Peripheral Participation Legitimate: Reader engagement experiments in Wikipedia. In *Proc. CSCW*, 2013.

[9] D. Helic. Analyzing user click paths in a wikipedia navigation game. In *MIPRO*, pages 374–379, 2012.

[10] A. Iriberri and G. Leroy. A life-cycle perspective on online community success. *ACM Comput. Surv.*, 41(2), 2009.

[11] B. Keegan, D. Gergle, and N. Contractor. Hot Off the Wiki Structures and Dynamics of Wikipedia's Coverage of Breaking News Events. *American Behavioral Scientist*, 57(5):595–622, 2013.

[12] A. Kittur and R. E. Kraut. Harnessing the wisdom of crowds in wikipedia: quality through coordination. In *Proc. CSCW*, 2008.

[13] P. Kollock. The economies of online cooperation: Gifts and public goods in cyberspace. In *Communities in Cyberspace*, pages 220–239. Routledge, 1990.

[14] J. Lave and E. Wenger. *Situated Learning: Legitimate Peripheral Participation (Learning in Doing: Social, Cognitive & Computational Perspectives)*. Cambridge University Press, September 1991.

[15] C. Okoli, M. Mehdi, M. Mesgari, F. Å. Nielsen, and A. Lanamäki. The People's Encyclopedia Under the Gaze of the Sages: A Systematic Review of Scholarly Research on Wikipedia. http://ssrn.com/abstract=2021326, 2012.

[16] J. Preece, B. Nonnecke, and D. Andrews. The top five reasons for lurking: improving community experiences for everyone. *Comp. in Human Behavior*, 20(2), 2004.

[17] J. Preece and B. Shneiderman. The reader-to-leader framework: Motivating technology-mediated social participation. *TOCHI*, 1, 2009.

[18] J. Ratkiewicz, S. Fortunato, A. Flammini, F. Menczer, and A. Vespignani. Characterizing and modeling the dynamics of online popularity. *Phys. Rev. Lett.*, 105:158701, Oct 2010.

[19] A. J. Reinoso. *Temporal and behavioral patterns in the use of Wikipedia*. PhD thesis, Universidad Rey Juan Carlos, 2011.

[20] A. Spoerri. Visualizing the overlap between the 100 most visited pages on Wikipedia for September 2006 to January 2007. *First Monday*, 12(4), 2007.

[21] A. Spoerri. What is popular on wikipedia and why? *First Monday*, 12(4), 2007.

[22] B. Stvilia, M. B. Twidale, L. C. Smith, and L. Gasser. Information Quality Work Organization in Wikipedia. *J. Am. Soc. Inf. Sci. Technol.*, 59(6):983–1001, Apr. 2008.

[23] B. Suh, E. H. Chi, A. Kittur, and B. A. Pendleton. Lifting the veil: improving accountability and social transparency in Wikipedia with wikidashboard. In *Proc. CHI*, 2008.

[24] C. Tann and M. Sanderson. Are web-based informational queries changing? *J. Am. Soc. Inf. Sci. Technol.*, 60(6):1290–1293, 2009.

[25] M. ten Thij, Y. Volkovich, D. Laniado, and A. Kaltenbrunner. Modeling and predicting page-view dynamics on Wikipedia. *CoRR*, abs/1212.5943, 2012.

[26] V. Waller. The search queries that took Australian Internet users to Wikipedia. *Information Research*, 16(2):476, 2011.

[27] R. West, I. Weber, and C. Castillo. Drawing a data-driven portrait of Wikipedia editors. In *Proc. WikiSym*, 2012.

[28] Wikimedia. Tool Labs. https://wikitech.wikimedia.org/wiki/Nova_Resource:Tools/Help.

[29] Wikipedia. List of featured articles. http://en.wikipedia.org/wiki/Wikipedia:Featured_articles.

[30] Wikipedia. List of good articles. http://en.wikipedia.org/wiki/Wikipedia:Good_articles.

[31] Wikipedia. Version 1.0 Editorial Team/Assessment. http://en.wikipedia.org/wiki/Wikipedia:Version_1.0_Editorial_Team/Assessment.

[32] D. Wilkinson and B. Huberman. Assessing the value of cooperation in wikipedia. *First Monday*, 12, 2007.

[33] T. Wöhner and R. Peters. Assessing the quality of Wikipedia articles with lifecycle based metrics. In *Proc. WikiSym*, 2009.