

Using Edit Sessions to Measure Participation in Wikipedia

R. Stuart Geiger

School of Information
University of California, Berkeley
102 South Hall, Berkeley CA 94703
stuart@stuartgeiger.com

Aaron Halfaker

Grouplens Research
University of Minnesota
200 Union St. S.E., Minneapolis, MN 55455
halfak@cs.umn.edu

ABSTRACT

Many quantitative, log-based studies of participation and contribution in CSCW and CMC systems measure the activity of users in terms of output, based on metrics like posts to forums, edits to Wikipedia articles, or commits to code repositories. In this paper, we instead seek to estimate the amount of time users have spent contributing. Through an analysis of Wikipedia log data, we identify a pattern of punctuated bursts in editors' activity that we refer to as *edit sessions*. Based on these edit sessions, we build a metric that approximates the labor hours of editors in the encyclopedia. Using this metric, we first compare labor-based analyses with output-based analyses, finding that the activity of many editors can appear quite differently based on the kind of metric used. Second, we use edit session data to examine phenomena that cannot be adequately studied with purely output-based metrics, such as the total number of labor hours for the entire project.

ACM Classification Keywords

H.5.3 [Information Systems]: Group and Organization Interfaces—*computer-supported collaborative work*

Keywords

Labor; activity; work; work practices; Wikipedia; peer production; labor-hours; sessions; quantitative methods

INTRODUCTION

Measuring Wikipedia

In less than a decade, Wikipedia has grown from a frequently ridiculed experiment to one of the world's most popular websites. The online encyclopedia has reached near-ubiquity among Internet users and is often invoked as a synecdoche for user-generated content communities, crowdsourcing, peer production, and Web 2.0. As such, it is hardly surprising that a number of high-impact statistics demonstrating the project's unexpected success are frequently mentioned in the public sphere. As of April 2012, there have been 528 million edits made to the English-language version and a total of 1.29 billion edits

across all language versions [23]. Other commentators describe the project in terms of its article content, not the amount of work put into those articles, and such figures are equally daunting: 19 million encyclopedia articles contain 8 billion words in 270 languages, and the English-language Wikipedia alone has 3.9 million articles containing 2.5 billion words. [30]

While most of these and other statistics are backed up by a substantial amount of empirical research, estimations of the total number of labor-hours contributed to Wikipedia are one notable exception. However, this has not stopped champions of the project from stating with more and less certainty that Wikipedia is one of the largest projects in human history. Yet in his 2010 book *Cognitive Surplus*, [24] Clay Shirky makes the opposite argument: he first estimates that 100 million labor-hours have contributed to Wikipedia, but then compares this amount of time with the absolutely staggering statistic that Americans spend 200 billion hours watching television each year. Shirky's argument is that we spend a substantial amount of time on activities like television, which effectively waste our collective brainpower on acts of consumption as opposed to projects like Wikipedia, which foster creativity and collaboration.

While the total number of labor hours that have been contributed to Wikipedia is an interesting bit of trivia, measuring contributions in terms of labor hours is a novel approach to not only studying Wikipedia, but other CSCW and CMC platforms and communities. Casting contributions to a peer production platform like Wikipedia in terms of labor hours, as opposed to a metric based on the number of contributions or posts, radically reframes how we conceptualize users. If we are interested in measuring users in terms of how prolific or active they are, then previous quantitative methods are rather well-developed and deployed. This is especially the case considering the vast and often public records that log and document activity in wikis, open source software projects, message boards, forums, citizen science projects. However, if we are interested in understanding the volunteers of peer production projects in terms of their level of investment or engagement, we have been traditionally limited to surveys, interviews, time diaries, and other approaches which do not scale very well.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

CSCW '13, February 23–27, 2013, San Antonio, Texas, USA.
Copyright 2013 ACM 978-1-4503-1331-5/13/02...\$15.00.

In this paper, we introduce a novel method for estimating the amount of time that users contribute to projects like Wikipedia, using the concept of edit sessions – which are punctuated bursts of editing activity captured by log data – as a way of identifying continuous periods of work on the wiki. This approach is not merely useful for arriving at a far more accurate estimation of the total number of labor-hours contributed to Wikipedia, but more importantly, it can be used to measure the activity of editors at a variety of scales. This approach provides an alternative to how most large-scale quantitative studies of not only Wikipedia but other CSCW platforms have operationalized activity using output-based metrics. Most notably, our analysis of labor-hours contributed by different cohorts of editors sheds new light on the project's oft-discussed growth and decline. [8,15,27] Specifically, we find that editors who joined the project in 2006 have not only contributed more labor-hours to the project than any other annual cohort, but continue to out-perform all other cohorts in 2012.

Measuring labor and work in CSCW

The most common way in which work in Wikipedia is measured is through edit counts, where one edit to a wiki document is considered one fungible unit of work. In most of the accounts of Wikipedia's power law inequalities – that is, how 1% of the editors contribute 55% of the content – edit counts are used. [27] Total bytes or words contributed in edits have also been used in order to arrive at a more nuanced figure. In addition, output-based metrics that examine how long edits persist are becoming quite popular when studying Wikipedia editors. [1,9,21,22]

However, many Wikipedia researchers have been moving away from raw edit count metrics in recent years. The main reason behind this stems from the realization that all edits are not equal. For example, the kinds of tasks and activities that predict whether editors become administrators has been modeled, with the results indicating that in many cases, the kind of contributions made matters more than the raw number. [4] One trend is towards using structured traces and articulations of work, such as barnstars and warning templates, to qualitatively and quantitatively measure the kinds of work that editors are rewarded and punished for performing. [6,7,12,19] Another trend is to measure editors by the number of edits or words that persist in articles, so that a spammer who makes thousands of edits which are always instantly removed ranks lower than an less active editor whose few contributions form the basis of the project's most edited and viewed articles. [9,10,21,22]

This tendency to use output-based metrics is not unique to studies of Wikipedia, as many large-scale quantitative studies of discussion forums [13], learning environments [16], recommender systems [17], open source software development, and other platforms often reduce interaction to one or occasionally two fungible units of activity. These are usually based on whatever kind of contribution is natively supported in the software platform. Wilkinson's study of peer production communities [29] is an excellent

example of this, as he compares the power law distributions of activity in Wikipedia, Digg, Bugzilla, and Essembly. In Wikipedia, he examines articles created and edits to articles; in Bugzilla, bugs submitted and comments made; in Digg, stories submitted and 'diggs' to stories; and in Essembly, resolves proposed and votes cast.

These output-based metrics are quite useful in measuring work practices and the distribution of labor across content creation communities, which, as Wilkinson argues, often follows a power law distribution. However, we take from Barabasi's insight that that human activity is often not randomly or normally distributed, but instead occurs in bursts. [2] Reconceptualizing work and contributions in terms of time as opposed to content may seem counter-intuitive given that communities like Wikipedia are organized around producing content. Yet as we show in the later sections, labor-based metrics give us quite a different view of who Wikipedia's top contributors are, for example.

METHODOLOGY

Beyond 'editcountitis': the story of a mixed-methodological collaboration

Before detailing our quantitative methodology, we wish to note that this study was the result of a methodological collaboration between the authors: one of us is a qualitative ethnographer and the other is quantitative computer scientist. Furthermore, both of us have been long-term editors of Wikipedia and members of the Wikipedian community, in addition to having studied Wikipedia for some time. Our decision to measure the activity of Wikipedians in terms of labor-hours was inspired by a number of qualitative and ethnographic observations we made about the ways in which Wikipedians quantify and aggregate their own labor practices. We believe that our quantitative methodology independently verifies the veracity of the edit session metric, but we wish to detail our motivation to provide context as well as to inspire future lines of mixed-method research.

Wikipedians have developed a number of ways to measure the relative contributions made by each editor. The simplest and easiest of these is the edit count, which as previously mentioned, is also a metric widely used by Wikipedia researchers. However, we have found that many Wikipedians, particularly veteran editors and administrators, know quite well that a Wikipedian's edit counts do not necessarily reflect the amount of time, energy, and effort they have contributed to the project. A widely-circulated essay on "Editcountitis" succinctly summarizes this view:

Editcountitis is used humorously to suggest a belief that a Wikipedian's overall contribution level can be measured solely by their edit count. This is a phenomenon which some think may be harmful to processes such as requests for adminship, as well as to the Wikipedia community in itself. The problems with using edit counts to measure relative level of experience

are that it does not take into account that users could have an extensive edit history prior to registering an account (posting anonymously), and that major and minor edits are counted equally, regardless of whether the edit is a typo fix, or the creation of a full article. Furthermore, edit counts do not judge the quality of the edits, as insightful comments on talk pages and acts of vandalism are counted equally. Hence, it is not always a reliable way of telling how experienced or worthy a user truly is. [28]

Our edit session metric is a direct response to the claims of editcountitis by Wikipedians; we see our metric as one of many tools that researchers and editors can use to measure the labor of Wikipedians. A qualitative, ethnographic study of the ways in which Wikipedians measure and value their own activities would be quite revealing and could further ground this line of research. However, such a study is outside of the scope of this paper, which is to establish a new metric for quantifying the labor practices in peer production communities. We should note that we have begun conducting preliminary qualitative interviews with Wikipedians about these issues in order to inform this quantitative research, and in future research, we aim to introduce the edit sessions metric and study how it affects the relationships between Wikipedians and their work practices.

From edit sessions to labor hours

In this section, we describe and justify the metric we use to estimate the labor hours spent by editors working on Wikipedia: the *edit session*. Our intent is to estimate, in a consistent manner, the total amount of time a user has spent contributing to a site like Wikipedia. We use the concept of an activity session, a technique that is commonly used to track and categorize website visitor activity. [20] While sessions are usually tracked via page view data, we track editor activity based on revision histories and logging data. This metric only includes work that is done by editors on <http://en.wikipedia.org>, and we can only identify editor activity based on their editing history. Well-founded privacy concerns in the Wikipedian community prohibit us from using page view data to track individual users as they perform actions which do not result in an edit – such as reading a long discussion without making a comment – and we note this as a limitation later in the paper. However, for researchers who have access to these data, our methodology can certainly take advantage of access logs to provide an even more sophisticated analysis.

Defining edit sessions

Within the log data of Wikipedia, a user’s edits appear as a stream of events with associated timestamps. In order to divide the stream edit activity into contiguous sessions of edits, the time between edit events (inter-edit time) can be examined and a method for identifying boundaries in the stream must be employed. We define an edit session as a *sequence of edits made by an editor where the difference between the time at which any two sequential edits are*

saved is less than one hour. In other words, a set of edits S is an edit session if:

$$\forall e_1, e_2 \in S: I(e_1) = I(e_2) - 1 \rightarrow T(e_2) - T(e_1) \leq \alpha$$

where:

$I(e)$ = the index of edit e in a sequence of edits

$T(e)$ = the time at which edit e occurred in seconds

α = the maximum time between edits (one hour)

To gather edit sessions for the English Wikipedia, we used a backup copy of the project’s MySQL database to process all page revisions sorted by the time in which they occurred and processed them sequentially. While stepping forward through revisions, we identified the start and end of edit sessions by tracking the “user_text” (username for registered editors, IP address for anonymous editors) and last edit timestamps. When the last edit timestamp for an individual became stale (> 1 hour old), we conclude that an session ended at the time of their last edit. This processing approach allows us to compute the edit sessions for a user in the same amount of time and complexity as the commonly used edit counters ($O(n)$).

Figure 1 illustrates an example edit session produced by applying this method. Toby Bartels’ first edit in this session occurred at 00:11, with edits at 00:20, 00:29, and 00:47. While this editor made edits before and after this session, but they took place more than 60 minutes before the first

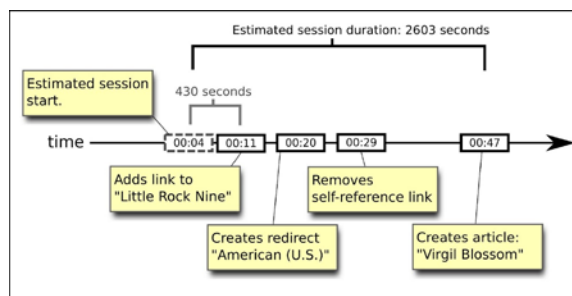


Figure 1. Estimated session length for Toby Bartels. Edits made by “Toby Bartels” are plotted and annotated over time for a session he completed on September 4th, 2010. The estimated session start time is plotted at 430 seconds before the user’s first edit.

edit at 00:11 and more than 60 minutes after the last edit at 00:47. As this example illustrates, there are a variety of tasks which the user performs during this time on a number of different pages.

Vetting the one hour cutoff

Preliminary work that looked at edit sessions in a more limited context simply placed a cutoff on inter-edit time at 1 hour under the assumption that the largest edits will take about an hour to complete, so any more time between edits meant that the editor left of site. We tested the validity of this cutoff time by analyzing the time between edits for a random sample of 1 million revisions, which were then filtered to exclude edits from unregistered and bot users as

well as the first edit by a registered user. With the remaining 821,749 revisions, we retrieved the time between the sampled revision and the previous revision from that user. This produced a long tail distribution that we suspected to be log-normal, so we bucketed inter-edit time logarithmically to produce the empirical histogram in Figure 2.

We suspected that the histogram was a result of summing at least two log-normal distributions representing within session time and between session time, so we fit the summed distributions using an expectation maximization (EM) algorithm. It turns out that we achieved a much better fit with three distributions than with two. We suspect that the smallest inter-edit times, overlaid in red, correspond to within-session edits and is on the order of minutes, making up a bulk of the revisions sampled. The second, overlaid in blue, corresponds to time between edit sessions and is on the order of days, while the third, overlaid in green, corresponds to extended breaks from the project (or "wikibreaks") and is on the order of months. As the dashed line in Figure 2 shows, the hour cutoff is just under the intersection of the inter-session and between-session distributions. To ensure that the 1 hour cutoff behaves reasonably over time, we graphed the fitted means and standard deviations for the history of Wikipedia (Figure 3).

The consistency of these fits indicate that the distribution of within-session and between-session inter-edit time stays largely consistent since the project gained critical mass and began its exponential growth phase. [27] These analyses suggest that the 1 hour cutoff appropriately divides within-session edits from between session edits and is useful throughout the history of the Wikipedia project. This consistency is especially striking considering that the pattern of inter-edit times has been generally unaffected by the many changes which have occurred in the community, such as the growing and shrinking size of the editor base, the development of policies and bureaucracies, the evolution of tasks in and around article writing, and various technical developments and new features. In addition, the distribution of the time between edits suggests that editors in Wikipedia follow a "Barabasi queueing" [2] model in which edits are largely concentrated in contiguous sessions as opposed to being normally or randomly distributed, which we suspect will also hold in other peer production projects.

Estimating edit session duration

From an edit session, we derive the session duration as a proxy to the amount of time an editor actually spent working on Wikipedia. We assume that, in between the edits they make, editors are performing legitimate wiki-work, and therefore, we can estimate their labor hours by measuring the time taken to complete their session. A naive way to derive session duration from an edit session is to simply find the difference in time between the first and last edits in the session. However, this approach does not account for the amount of time that the first edit in a

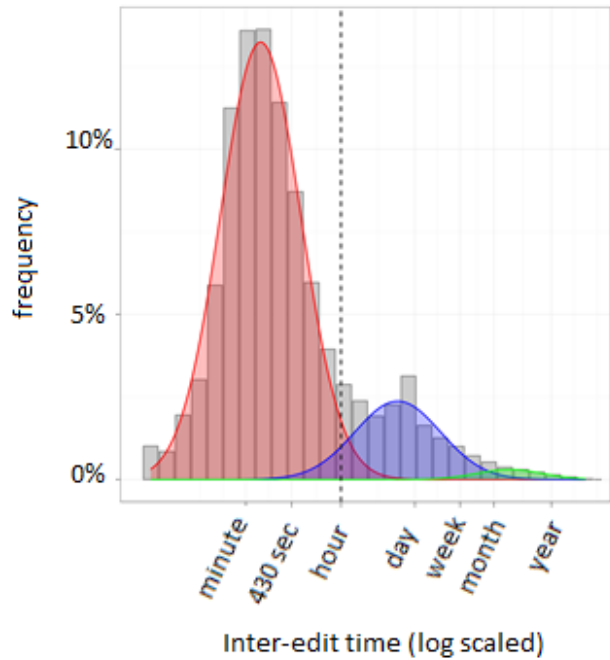


Figure 2. Histogram of time between users edits with an EM fit of three log-normal distributions corresponding to within-session, between-session, and extended session breaks.

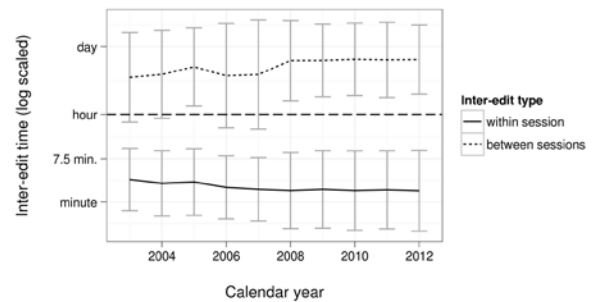


Figure 3. EM fitted means and standard deviations of within and between session inter-edit times over time (wikibreaks excluded).

session required to make, and therefore, sessions that contain only one edit would appear to have required zero labor-hours. To account for the time that the temporal bounds of edit sessions do not capture, we calculated the average time between edits across all sessions that contained more than one edit, which was 430 seconds. We combine the difference in time between the first and last session edits with 430 seconds (see "Estimated session start" in Figure 1) to produce an estimated session duration.

RESULTS

Analyses derived from edit session data

Out of the 528 million edits in the English-language version of Wikipedia from January 2001 to April 2012, we identified 423 million edits that were not made by automated bots. Iterating through these edits, we identified 114 million distinct edit sessions by 33.6 million distinct registered and anonymous editors. Of these, 60.6 million edits (14.3% of all edits) were made outside of a

continuous edit session and were assigned the duration of 430 seconds, the average time between edits in a multi-edit session, as explained in the methodology section. The median session length was 10 minutes, and as Figure 2 illustrates, the distribution of sessions is highly skewed towards short sessions, although some notable outliers exist. For example, our dataset includes a 1251 minute session by a Wikipedian who spent nearly 21 continuous hours contributing to articles, participating in discussions and sending messages to other users during a marathon of editing in December of 2006.

Most sessions are much shorter, as 83.4% of all edit sessions were less than 30 minutes in length. As Figure 4 illustrates, there is a slight increase in the mean session duration from 27.3 to 33.6 minutes between sessions performed in editor's first month and sessions performed in their 2nd year. The mean session time for edits performed in a editor's second year stays relatively consistent from that point forward. This indicates that as Wikipedians editors mature, they edit in slightly longer sessions, but do not substantially change their session habits. Due a long tail of session length, the median session length is substantially lower than the mean across editors of all tenure.

Edit session data can be used to reveal interesting aspects of the work practices of contributors to peer production projects. Figure 5 plots the average session length per day of the week, finding that there is a small but noticeable pattern. Edit sessions on weekends tend to be longer than those during the middle of the week, suggesting that for at least some Wikipedians, weekends are spent on longer and more complex tasks. Figure 6 plots the average session length per month, showing that the longest edit sessions are in the summer and winter months. Taking into account Western work and education cycles, these analyses lend support to a hypothesis that Wikipedians edit for longer periods of time during periods of leisure and vacation.

Edit counts versus labor hours

The rate at which an editor saves revisions to pages can vary substantially based on their wiki-work habits and the kind of activity they are engaged in. Even when working on a single article, some editors save changes every minute while others will not save their changes until they have finished the task at hand. Furthermore, 3rd party tools like AutoWikiBrowser and Huggle pre-script similar tasks and allow editors to make several revisions per minute, or even every few seconds in some cases. It is assumed that more edits typically corresponds to more labor, but only if editors are doing the same type of work with the same editing habits. To demonstrate this, we compared edit counts versus edit session metrics, first determining the top 20 editors in March 2012 by edit count (Figure 7). These editors performed a total of 259,516 edits, or 7.87% of all the edits from registered editors that month. Applying our metric to the top 20 editors by edit count, we found that they account for 4,276.5 labor hours in March 2012, or

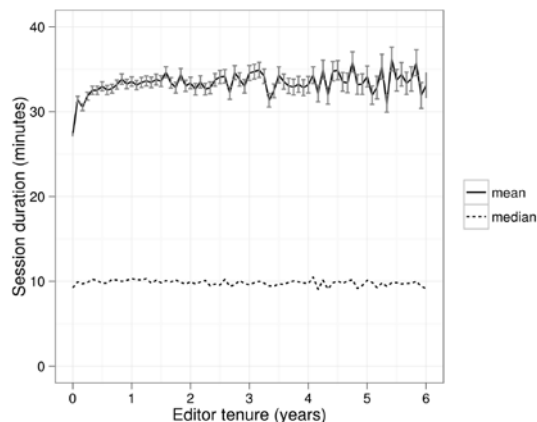


Figure 4. Mean and median session duration by number of years the editor has been registered.

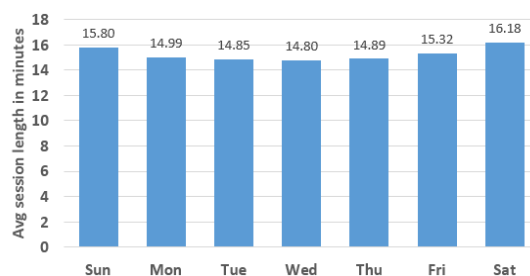


Figure 5. Average edit session length by day of week, 2001-12

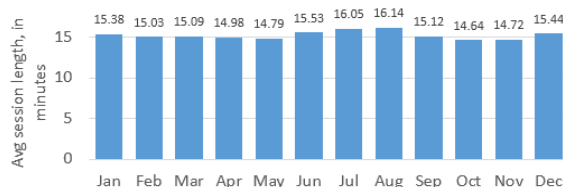


Figure 6. Average edit session length by month, 2001-12

Top editors by edits			Top editors by edit sessions		
	Username	edits	Username	edits	hours
1	Koavf	43997	Materialscientist	7472	453.2
2	Waacstats	33402	Jtmorgan	1231	402.9
3	Hmains	17176	Kwamikagami	9088	356.3
4	Rich Farmbrough	17169	TonyTheTiger	6152	344.0
5	Bgwhite	14531	ACP2011	2218	337.2
6	Courcelles	13832	Pinethicket	3894	317.7
7	Fortdj33	12919	Armbrust	6288	311.0
8	VasuVR	12095	P.T. Aufrette	6257	306.4
9	BD2412	9801	Koavf	43997	302.2
10	Cloudz679	9779	Derek R Bullamore	4228	290.0
11	Kwamikagami	9088	MathewTownsend	1807	280.8
12	Muboshgu	8098	Crisco 1492	2747	278.5
13	Tassedethe	7976	Alarbus	1669	277.5
14	Materialscientist	7472	Rich Farmbrough	17169	274.8
15	John of Reading	7415	Alan Liefting	5970	274.3
16	DBigXray	7405	BD2412	9801	273.2
17	Ssriram mt	7100	Sitush	4421	270.7
18	Woohookitty	7099	DBigXray	7405	270.2
19	Allens	6757	Allens	6757	270.1
20	Fram	6405	Cloudz679	9779	249.9

Figure 7. Top 20 editors in March 2012 by edits and sessions

1.5% of the total labor hours from registered editors that month.

As an alternative ranking, we retrieved all the editors who, according to our edit session metric, contributed a total of at least 8 hours of labor a day per day in March 2012. In all, 20 registered editors met this 248 hour threshold, and the lists (Figure 7) are quite different. There is a small amount of overlap (Jaccard index = 0.29) that suggests these metrics are making a similar but certainly not identical measurement. For sample, the highest editors in both rankings appear in both lists. However, the same is not true for either of the second-highest ranked editors. These differences suggest that edit count and edit sessions are measuring editor labor differently, but why choose one over the other? We contend that measuring labor with edit sessions benefits over edit count in two important ways: (1) labor hours should be comparable between editors performing a wide range of different tasks in Wikipedia and (2) measuring work in hours is more intuitive. To demonstrate this intuitive nature, we pose a question. Which tells us more about an editor such as MatthewTownsend and his investment and motivations: that he made 1,807 edits in March 2012, or that in that time period, he edited Wikipedia for an average of 9 hours per day, *every day*?

For researchers and community members who are interested in identifying the most dedicated and invested contributors, both the raw number of labor hours and the number of sessions which last longer than a standard work day provide an alternative metric. This is one of the many kinds of analyses that edit session data can be used to generate. For example, Figure 9 plots the raw number and proportion of edit sessions over 8 hours in length since 2004. This shows a different view of the project's much-discussed decline (Figure 8, see [15,27]), indicating that there is only a slight decline in the number of times Wikipedians engage in quite lengthy, dedicated editing sessions. This is useful because it indicates that core editors are continuing to invest substantial amounts of time in the project, suggesting that the decline may result from a loss of the more peripheral and casual contributors. In fact, the steady rise in the proportion of edit sessions lasting longer than 8 hours is in line with recent work demonstrating that the project's decline stems from a failure to recruit and retain newcomers as opposed to a mass exodus from the project's most longstanding and dedicated contributors. [8]

Comparing labor across user cohorts

An interesting use of our labor-hour metric is to compare the total labor hours between groups of Wikipedians depending on how long they have been editing Wikipedia. We bucket registered users into annual cohorts based on when each user made their first edit and tracked their aggregate labor hours over time. Figure 10 plots the stacked total labor hours for each cohort up to April, 2012. Editors who started in 2006 spend the most time editing Wikipedia. Even in 2012, editors who joined the project in 2006

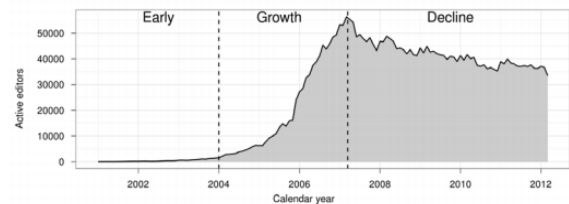


Figure 8. Wikipedia's growth and decline. The number of active editors, defined as over 5 edits per month, data from [8].

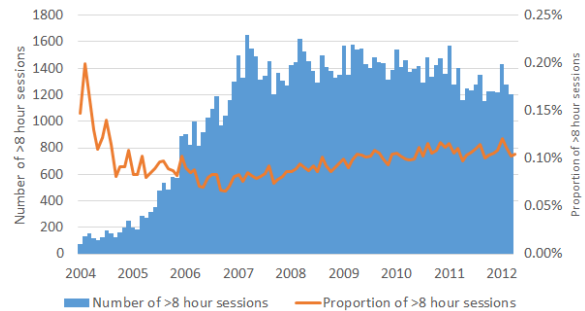


Figure 9. Number and proportion of extended edit sessions (>8 hours) over the history of Wikipedia since 2004.

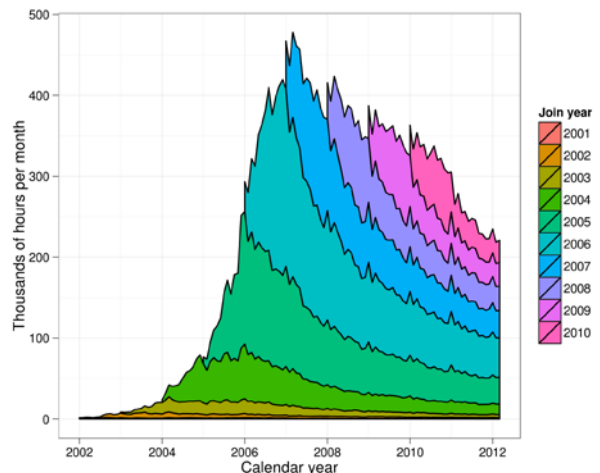


Figure 10. Labor-hours contributed to English Wikipedia per month, by registered user cohorts from 2001-2010.

collectively spend more time editing Wikipedia than any other cohort. The total labor hours contributed by the cohort of 2006 in March 2012 is 41,583 while the next highest (the 2007 cohort) is 30,598. In contrast, the 2011 cohort is quite low, at 16,122 hours.

Computing total labor hours contributed

Our approach also enables us to replicate Shirky and Wattenberg's back of the envelope estimations about the total number of labor hours contributed to Wikipedia. Labor-hour (or man-hour) calculations for large-scale projects are typically found in back-of-the-envelope calculations, not rigorous analysis of actual work performed – such analysis is usually impossible. Labor-hour calculations are typically derived by multiplying the

number of employees who work on a project in a given week by the average length of their work week, and then multiplying that figure by the number of weeks spent on the project. While this can be slightly complicated when some workers are employed for different amounts of time (part vs. full time) or when the number of workers changes in a new phase of the project, most modern megaprojects are administered in such a way that these labor patterns are well-documented.

For example, a well-documented and often-repeated labor-hour estimation is that of the Empire State Building, which took 3,000 laborers a total of 7 million labor-hours to construct. [14] Figures for the construction of the Channel Tunnel report a total 170 million labor-hours, [5] while estimations of the Great Pyramid at Giza range from 880 million [25] to 3.5 billion labor-hours. [26] The first edition of the *Encyclopedia Britannica* was written and published by 3 employees authoring 24 pages a week for 100 weeks, [11] which is around 12,000 labor-hours assuming 40 hour work week. Alternatively, labor hours have been used as the basis of studies of software development and project management, such as in Brooks' influential *The Mythical Man-Month* [3], where he argues that adding more labor to a project does not necessarily speed up the project – in fact, it can often slow a project down.

Summing the duration of all continuous editing sessions and single edit sessions, we identified 41,018,804 total labor-hours expended in the English-language version of Wikipedia. As Figure 11 shows, the number of labor-hours per month experiences a comparable exponential growth and decline as in editing that has been discussed extensively [8,15,27]. At the peak of the project's growth, approximately 675,000 labor-hours were contributed each month, but this has fallen to approximately 425,000 labor-hours in 2012. This graph also illustrates a similar distribution between the number of labor-hours contributed by registered and anonymous editors: 27.43% of all labor-hours were from anonymous editors, compared to 25.83% of all edits. Extrapolating to all language version of Wikipedia based on the total number of edits made to each project, we estimate that 61,706,883 total labor-hours have been contributed in edit sessions for non-English language Wikipedias, for a total of 102,673,683 total labor-hours to all Wikipedia versions.

OBJECTIONS AND LIMITATIONS

Although this approach simplifies all facets of the wiki-work into interactions that change the content of the wiki's pages, we argue that this approach to measuring work hours is robust to the most common types of wiki-work. Yet the most immediate objection to such a metric is that there is little ancillary data to support the assumption that a user is active during the time between edits. This is mitigated by the fact that in Wikipedia, nearly all actions within the MediaWiki software platform are represented by an edit to a page. [6,7] This is because the platform is notoriously lacking in built-in features. While the MediaWiki platform

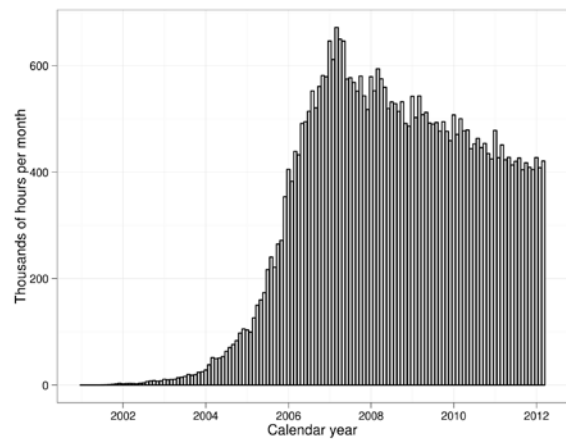


Figure 11. Total monthly labor hours over the history of Wikipedia for all registered, human editors.

supports discussion spaces, meta-discussion spaces, user-to-user messaging, user profiles, quality control procedures, task discovery and assignment mechanisms, administrative queues, newsletters, and many other tasks beyond simple encyclopedia editing, all of these features are represented as pages that can be edited.

For example, when an administrator blocks a user, it is customary to edit the user's talk page (where anyone can leave them a message) and leave a templated message indicating that the user has been blocked. [6] When a user requests that an administrator temporarily protect an article from editing, they do so by editing a specialized page, and it is customary for the responding administrator to edit this page as well, marking the request as approved or rejected and removing it from the queue. Almost all of these activities are represented as edits to pages and are, therefore, included in our analysis of edit sessions. As such, an edit session should capture article editing activity as well as communication and coordination activities across the system. However, this prolific logging of social and organizational activity is quite specific to Wikipedia and grounded through ethnographic fieldwork. Because of this, researchers seeking to use our edit session metric should first understand what kind of activity is logged.

Labor which occurs but is not measured

We must stress that Wikipedians perform a number of activities which are critically important to the encyclopedia project but are completely invisible in this calculation. The most revealing aspect of this can be seen in our algorithm's estimation for the number of hours spent in edit sessions by Wikipedia's co-founder Jimmy Wales: 869 hours – or just under 22 full-time work weeks – since January 2001. First, this estimation completely neglects the amount of work Jimmy Wales has put into Wikipedia behind the scenes. However, it is lower than expected given that he spends much of his time resolving disputes and building consensus, tasks which involve reading a substantial amount of existing dialogue before stating one's own opinion. A better figure to showcase Wales' dedication would be that out of the ~4,100 days since Wikipedia was

founded, he has made an edit almost every other day – 1,993 distinct days in total.

There are many different ways in which Wikipedians contribute labor to the encyclopedia project that never result in even a single edit. These kinds of activities are entirely unaccounted for in our analysis, and include: carefully reviewing articles for errors and finding none; looking up a source to see if it is accurate, and finding it is; reading various policies, guidelines, and the manual of style; and keeping up with various project-wide debates and controversies without weighing in. Furthermore, many editors do not edit in continuous sessions, spending dozens or even hundreds of hours in activities that result in just one edit. These include activities like: writing a full-length article in a text editor and submitting it in one edit; searching for historical documents in archives; reading an entire book to verify a source; reviewing hundreds of comments in a controversy before weighing in; traveling to a remote location to take a photograph. Another potential issue is that anonymous editors are treated as distinct editors, such that one hundred editors can all simultaneously be in a continuous edit session at the same library or institution, but they are treated as one editor. There are also a variety of activities that take place outside of <http://wikipedia.org>, such as the project's hundreds of mailing lists and IRC channels, editor-to-editor communication over personal channels like e-mail or IM, in person meetups, and the annual Wikimania conference. However, we must note that most of these activities are also not directly captured by current methodologies that rely on output-based metrics.

Possible over-estimations of labor

First, our inter-edit cutoff for sessions is one hour, and while we argue the validity of this value in section 2.2, it could be argued that this is too long of a time period. Hypothetically, an editor could make an edit, head to lunch, and then 59 minutes later, return home and respond to a message sent to them. If this occurs in less than 60 minutes, we assume they have been working the whole time. However, this raises a more fundamental question as to what labor is in a post-industrial society: is time spent 'on break' time spent working? Traditionally, labor-hours are the time that all laborers spend directly working on or supporting a project, which rarely includes the time each worker spends while on a scheduled break, paid vacation, commuting, etc.

Multitasking and rapid task switching are now ubiquitous [18], and many people edit Wikipedia while performing other tasks, such as watching television or even talking with friends and family on the telephone. This complicates our understanding of the edit session as a metric of discrete, continuous labor contributed to the encyclopedia project. For example, an editor who spends an hour editing Wikipedia and watching television may actually just be editing during the commercial breaks, spending the other 45 minutes focused on the TV. Or in an extreme case, an

editor could be gainfully employed at some unrelated workplace, and spend thirty seconds every half hour editing an article. If they did this all work day, it would appear in our metric as eight continuous hours of labor contributed to Wikipedia when only four minutes were expended. However, we can complicate this example even more: what if this user is patrolling a high-profile article for vandalism and has put themselves 'on call' for eight hours, using a notification tool to help them review every edit made to this article within one minute?

FUTURE WORK

This study has exclusively studied the English-language Wikipedia, and future research is necessary to further validate the use of edit sessions as a way of analyzing activity in both Wikipedia and other collaboration platforms. Diary studies or surveillance-based techniques, in which users are recorded or record their own behavior, could provide another form of validation to the edit session metric. Furthermore, it would be interesting to see what kinds of activity are not included in edit sessions but are recorded by Wikipedians, as discussed in section 4.2.

Future research can use the session approach to explore the differences and similarities between different classes and types of users, as well as the different kinds of activity which are performed. We noted that one of the longest edit sessions in our dataset was from a Wikipedian who was engaged more in communication and dispute resolution than article editing itself. Qualitative coding of a sample of edit sessions and both extremes could reveal substantial differences between how editors engage with a peer production project. Future research could also use statistical modeling to classify sets of similar users based on their edit session behaviors, asking, for example, if editors who are engaged in dispute resolution, counter-vandalism, or article construction tend to have many short sessions or a few long sessions.

The concept of the activity session is not unique to Wikipedia, and a study of contributions to fast-paced crowdsourcing platforms like Galaxy Zoo and NASA clickworkers would likely result in a drastically reduced value for the average time between sessions, as well as the three different distributions of breaks between edits. Yet we expect that this threefold distribution of breaks within sessions, breaks between sessions, and extended breaks would appear in any peer production community. Based on this framework, it would be quite revealing to compare the kind, number, and duration of between session and extended session breaks.

CONCLUSION

This paper has introduced and explained the concept of the edit session as a way of estimating the number of labor-hours that Wikipedians spend continuously contributing to the encyclopedia project. Our metric is a more robust and revealing way of operationalizing editors' contributions, activities, and levels of investment than pure edit counts and other output-based metrics. Edit sessions are also more

intuitive than edit counts, and labor-based metrics provide for a more robust comparison between editors, independent of the kinds of activities that editors perform. Session data also provides for interesting studies of interaction and activity in CSCW sessions, such as daily and seasonal periodicity as well as queuing behaviors. Examining other peer production communities using a labor-hours approach may also prove fruitful, and comparing other communities to Wikipedia might reveal interesting similarities and differences. Further research is also necessary to validate this method to real-world activity, possibly using diary studies or other modes of observational research.

ACKNOWLEDGMENTS

This work would not have been possible without the support of our research groups, the Wikimedia Foundation, NSF grants IIS 09-68483 and IIS 11-11201, as well as the constructive comments from the CSCW reviewers.

REFERENCES

- Adler, B.T., Alfaro, L. de, Pye, I., and Raman, V. Measuring Author Contributions to the Wikipedia. *Proc WikiSym 2008*, ACM Press (2008).
- Barabási, A.-L. The origin of bursts and heavy tails in human dynamics. *Nature* 435, 7039 (2005), 207–211.
- Brooks, F.P. *The Mythical Man-Month*. Addison-Wesley, 1995.
- Burke, M. and Kraut, R. Taking up the mop: identifying future wikipedia administrators. *Proc CHI 2008*, ACM (2008), 3441–3446.
- Cutler, D. Facts on the Channel Tunnel. *Reuters*, 2010. <http://www.reuters.com/article/2010/10/07/siemens-eurostar-tunnel-idUSLDE69623720101007>.
- Geiger, R. and Ribes, D. The work of sustaining order in wikipedia: the banning of a vandal. *Proc CSCW 2010*, ACM (2010).
- Geiger, R.S. and Ribes, D. Trace Ethnography: Following Coordination Through Documentary Practices. *Proc HICSS 2011.*, IEEE (2011).
- Halfaker, A., Geiger, R.S., Morgan, J.T., and Riedl, J. The Rise and Decline of an Open Collaboration System. *American Behavioral Scientist*. (In press). Accessed online 27 Aug 2012 at <http://halfaker.info/archive/halfaker12rise.pdf>
- Halfaker, A., Kittur, A., Kraut, R., and Riedl, J. A Jury of Your Peers: Quality, Experience and Ownership in Wikipedia. *Proc WikiSym 2009*, ACM Press (2009).
- Halfaker, A., Kittur, A., and Riedl, J. Don't bite the newbies: how reverts affect the quantity and quality of Wikipedia work. *Proc WikiSym 2011*, ACM (2011).
- Kogan, H. *The Great EB: The Story of the Encyclopædia Britannica*. University of Chicago Press, Chicago, 1958.
- Kriplean, T., Beschastnikh, I., and McDonald, D. Articulations of wikiwork: uncovering valued work in wikipedia through barnstars. *Proc CSCW 2008*, ACM (2008), 47–56.
- Lampe, C. and Resnick, P. Slash(dot) and burn: distributed moderation in a large online conversation space. *Proc CHI 2004*, ACM (2004), 543–550.
- Langmead, D. and Garnaut, C. *Encyclopedia of Architectural and Engineering Feats*. ABC-CLIO, 2001.
- van Liere, D. and Fung, H. *Editor Trends Study*. 2011. Accessed 13 April 2012 from http://strategy.wikimedia.org/wiki/Editor_Trends_Study
- Lipponen, L., Rahikainen, M., Lallimo, J., and Hakkarainen, K. Patterns of participation and discourse in elementary students' computer-supported collaborative learning. *Learning and Instruction* 13, 5 (2003), 487–509.
- Ludford, P.J., Cosley, D., Frankowski, D., and Terveen, L. Think different: increasing online community participation using uniqueness and group dissimilarity. *Proc CHI 2004*, ACM Press (2004), 631–638.
- Mark, G., Gonzalez, V.M., and Harris, J. No task left behind?: examining the nature of fragmented work. *Proc CSCW 2005*, ACM Press (2005).
- McDonald, D.W., Javanmardi, S., and Zachry, M. Finding patterns in behavioral observations by automatically labeling forms of wikiwork in Barnstars. *Proc WikiSym 2011*, ACM Press (2011).
- Mongomery, Alan, L., Li, S., Srinivasan, K., and Leichty, J.C. Modeling Online Browsing and Path Analysis Using Clickstream Data. *Marketing Science* 23, 4 (2004), 579–595.
- Panciera, K., Halfaker, A., and Terveen, L. Wikipedians are born, not made: a study of power editors on Wikipedia. *Proc GROUP 2009*, ACM (2009), 51–60.
- Priedhorsky, R., Chen, J., Lam, S.T.K., Panciera, K., Terveen, L., and Riedl, J. Creating, destroying, and restoring value in wikipedia. *Proc GROUP 2007*, ACM Press (2007), 259–268.
- S23. WikiStats: List of Wikipedias. 2012. http://s23.org/wikistats/wikipedias_html.php.
- Shirky, C. *Cognitive Surplus: Creativity and Generosity in a Connected Age*. Penguin, New York, 2010.
- Smith, C.B. Program Management B.C. *Civil Engineering* 69, 6 (1999).
- Smith, N. Classic Projects: Great Pyramid at Giza. *Engineering and Technology Magazine* 6, 1 (2011).
- Suh, B., Convertino, G., Chi, E., and Pirolli, P. The Singularity is Not Near: Slowing Growth of Wikipedia. *Proc WikiSym 2009*, ACM (2009).
- Wikipedia contributors. "Wikipedia:Editcountitis." *Wikipedia*. <http://enwp.org/W:ITIS>, accessed 7 Dec 2012.
- Wilkinson, D.M. Strong regularities in online peer production. *Proc EC 2008*, ACM Press (2008), 302.
- Zachte, E. *Wikipedia Statistics*. 2012. Accessed 13 Apr 2012, <http://stats.wikimedia.org/EN/Tables/WikipediaZZ.htm>

