

A Large-Scale Characterization of How Readers Browse Wikipedia

TIZIANO PICCARDI, EPFL, Switzerland
 MARTIN GERLACH, Wikimedia Foundation, USA
 AKHIL ARORA, EPFL, Switzerland
 ROBERT WEST*, EPFL, Switzerland

Despite the importance and pervasiveness of Wikipedia as one of the largest platforms for open knowledge, surprisingly little is known about how people navigate its content when seeking information. To bridge this gap, we present the first systematic large-scale analysis of how readers browse Wikipedia. Using billions of page requests from Wikipedia’s server logs, we measure how readers reach articles, how they transition between articles, and how these patterns combine into more complex navigation paths. We find that navigation behavior is characterized by highly diverse structures. Although most navigation paths are shallow, comprising a single pageload, there is much variety, and the depth and shape of paths vary systematically with topic, device type, and time of day. We show that Wikipedia navigation paths commonly mesh with external pages as part of a larger online ecosystem, and we describe how naturally occurring navigation paths are distinct from targeted navigation in lab-based settings. Our results further suggest that navigation is abandoned when readers reach low-quality pages. Taken together, these insights contribute to a more systematic understanding of readers’ information needs and allow for improving their experience on Wikipedia and the Web in general.

CCS Concepts: • **Information systems** → **World Wide Web; Information retrieval**; • **Applied computing** → **Digital libraries and archives**.

Additional Key Words and Phrases: wikipedia, web navigation, server logs, information needs

ACM Reference Format:

Tiziano Piccardi, Martin Gerlach, Akhil Arora, and Robert West. 2023. A Large-Scale Characterization of How Readers Browse Wikipedia. *ACM Trans. Web* 1, 1, Article 1 (January 2023), 22 pages. <https://doi.org/10.1145/3580318>

1 INTRODUCTION

Evolution has optimized humans for information seeking, and humans have in turn optimized the world around them to facilitate access to information. Many of the most consequential evolutionary, cultural, and technological advances in humans—from the development of language and writing systems to modern telecommunication—have enhanced their ability to find, ingest, process, and transfer information. Given the central importance of information seeking to human nature—epitomized by the view of humans as *informavores* [43]—, understanding the dynamics of how humans seek information and engage with knowledge is of key significance across disciplines, both

*Robert West is a Wikimedia Foundation Research Fellow.

Authors’ addresses: Tiziano Piccardi, EPFL, Lausanne, Switzerland, tiziano.piccardi@epfl.ch; Martin Gerlach, Wikimedia Foundation, USA, mgerlach@wikimedia.org; Akhil Arora, EPFL, Lausanne, Switzerland, akhil.arora@epfl.ch; Robert West, EPFL, Lausanne, Switzerland, robert.west@epfl.ch.

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. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

1559-1131/2023/1-ART1 \$15.00
<https://doi.org/10.1145/3580318>

in the basic and applied sciences. In the basic sciences, biologists, psychologists, anthropologists, among others, stand to gain fundamental insights into how humans function, whereas in the applied sciences, such insights can enable the design of more effective tools and information environments, such that humans can more readily find relevant knowledge in an ever-surging flood of information.

However, closely observing humans as they seek information is challenging, since it requires measuring predominantly cognitive behaviors at a great level of detail. As a consequence, although much work has been dedicated to shedding light on human information seeking behavior (see Sec. 2), it has faced important limitations: surveys [87] and thinking-out-loud studies [48] are prone to cognitive biases, as humans generally perform poorly at introspection [50]. Lab-based experiments [41] typically involve small samples consisting of biased populations (e.g., university students) and are thus frequently not representative and might lack statistical power. Studies based on surrogate tasks (e.g., navigation games [79]), although measuring navigation-related skills, do not capture real-world, self-motivated information seeking and may thus lack external validity [51]. Finally, studies based on aggregated versions of real-world information seeking traces (specifying page-to-page transition counts instead of full traces [14, 19]), although capturing local, page-level choices accurately, may lack relevant trace-level information (e.g., relating the start of a trace to its end).

In this work, we provide a complementary perspective in the context of encyclopedic information seeking—an important special case of human information seeking—by leveraging a large-scale dataset of digital traces compiled from one month’s worth of English Wikipedia’s complete server logs, which offer unprecedented opportunities for observing humans interacting with knowledge in great detail.

Wikipedia is a primary source of encyclopedic knowledge and plays a unique role in the global knowledge ecosystem, fulfilling a wide range of information needs [40, 66]. It is the largest encyclopedia ever built, with almost 60M articles in more than 300 languages. It is freely accessible across the globe and attracts more than 1.5B unique devices generating billions of pageloads every month, and it is the most popular website (except for search engines, Facebook, and YouTube) in 43 countries (more than any other website) [24]. Wikipedia thus reaches an audience whose representativeness far surpasses that of lab-based studies. Since Wikipedia’s server logs contain a record of all pageloads, the logs are uniquely suited for providing a geographically and temporally complete mirror of real-world, self-motivated encyclopedic information seeking.

In contrast to prior work, which has leveraged Wikipedia’s server logs to shed light on specific aspects of reader behavior (including reasons for visiting Wikipedia [40, 66], engagement with citations and external links [44, 54, 55], studying variation in dwell time [71], and measuring geo-localized collective behavior [72]), this paper is the first to employ the logs in a principled, broad analysis with the goal of systematically elucidating the nature and structure of encyclopedic information seeking pathways. By analyzing billions of navigation traces extracted from the logs (Sec. 3) at various levels of aggregation, we consider three research questions:

RQ1 How do readers reach Wikipedia articles? (Sec. 4)

RQ2 How do readers transition from one article to the next? (Sec. 5)

RQ3 What are the properties of entire reading sessions? (Sec. 6)

We find that Wikipedia navigation traces expose a wide variety of structures. While shallow sessions consisting of single pageloads dominate, we observe a long tail of long, complex traces, whose depth and shape vary systematically with topic, device type, and time of day. Although it is known that search engines play a key role in driving readers to Wikipedia, we further highlight their importance for navigation between pages, showing that browsing Wikipedia does not happen in isolation, but is embedded in sessions where users transition fluidly to and from the external

Web, frequently via search engines. We describe the interaction between article content and reader navigation, finding strong evidence that users stop navigating when reaching articles of low quality or the periphery of the network. Finally, we show important differences between in-the-wild Wikipedia usage on the one hand and targeted navigation behavior captured by lab-based studies on the other hand.

These findings are complemented with a description of best practices for analyzing readers' navigation using Wikipedia's server logs. We examine different ways of aggregating user sessions as well as their impact on the conclusions drawn.

Our results have important implications for Wikipedia and beyond. Understanding how readers explore content on Wikipedia is critical for framing its role in fulfilling information needs and for making design decisions regarding its structure, format, accessibility, and supportive tools such as recommender systems. Going beyond Wikipedia, these findings may help deepen our understanding of how humans navigate information when seeking knowledge.

2 RELATED WORK

Information-seeking behavior. Over time, information-seeking behavior has received attention from sociologists, cognitive psychologists, and, more recently, computer scientists, thanks to the availability of digital trace data. The study of information seeking investigates the strategies used by humans to find a piece of information to satisfy an information need [83]. Whereas the definition of "need" is unclear and relatively hard to formalize, seeking behavior is observable and easier to model, especially in information systems [84, 85]. A complementary hypothesis from cognitive psychology argues that humans are *informavores* and seek information with the same dynamics used by animals in searching for food [43]. This idea inspired the formulation of information foraging theory [57], which describes humans as behaving akin to predators in the information space, relying on "information scent" [9] to find the paths that maximize the chances of leading them to the desired piece of information [45, 62, 69]. Similarly, complementary models applied to the Web describe users' navigation as driven by the relatedness of a link or image with the desired goal [31]. Finally, additional cognitive models include "berrypicking" [5], which describes the search for information as a dynamic process where users collect small portions of information bit by bit, and "exploratory search" [59], which describes the information-seeking behavior of users unfamiliar with the topic of their search or with an unclear goal.

Navigation on the Web and log analysis. Characterizing user navigation on the Web is a challenging task because of the limited availability of data. Previous work focused on modeling navigation patterns based on server logs of large websites or by using modified browser versions. A common finding is that people frequently revisit the same content multiple times [1, 70]. This repeat consumption behavior, which is abandoned when the person becomes bored of the content [6], makes human mobility on the Web predictable [34]. Although researchers found that Web users are not strictly Markovian (the page visited next does not depend exclusively on the current page) [10], many prediction models approximate the navigation of users on a network with Markov chains [12, 42, 58] and hybrid models [4, 28, 30, 49].

A significant effort in investigating Web navigation has focused on search engines and how people find content from relevant keywords [27]. Log-based analyses of the navigation following a Web search show that people's behavior exhibits a high level of variability [81] and that different search queries and origins are associated with different navigation patterns [7, 26]. Beyond characterizing users' information needs, digital trails can be exploited to improve search engine results [15, 17, 68, 82], e.g., by using the collective interest of a destination page as a metric of relevance [80]. Similarly, navigation traces have proven useful as a tool to improve website navigability by identifying

missing links [35, 52, 78] and other usability issues that normally require the work of domain experts [18]. Finally, navigation logs can be used to compute the semantic relatedness of pages by studying what content is typically accessed together [11, 67].

Reader behavior on Wikipedia. Researchers have also studied how readers behave when reading Wikipedia. Recent work focuses on the interaction with external links [55] and references [44, 54], and on the reading time of articles [71]. Researchers have concluded that Wikipedia users have reading patterns that fall in different categories, such as exploration, focus, trending, and passing [39], and that readers prefer links that lead to the periphery of the network, about semantically similar content and located at the top of the article [14, 36]. Other studies have investigated the inter-event time in the navigation logs of Wikipedia and found strong regularities in the temporal rhythms, which suggest a reasonable rule of thumb for segmenting sessions after inactivity periods of one hour [22].

These studies are complemented by investigations of the motivations for visiting Wikipedia [40, 66], which describe a variety of factors such as current events, media coverage of a topic, personal curiosity, work or school assignments, or boredom.

Closest in spirit to the present work, multiple approaches have been used to study human navigation on Wikipedia. The public clickstream [86] contains transition counts for pairs of articles. Although the clickstream constitutes an aggregated and filtered version of the server logs, it has been shown that it can serve as a useful approximation in many practical applications [3]. It has been used to study how different topics relay more traffic than others [13, 19], and how readers' navigation paths tend to start general and become incrementally more focused at every step [61].

Other approaches to understanding readers' navigation have identified different types of curiosity during Wikipedia exploration by relying on data shared by volunteers [41], while yet others have characterized human navigation as manifested in digital traces obtained via Wikipedia navigation games such as Wikispeedia [79], where players start from a random article and are tasked to reach a target page in as few clicks as possible by following links only. These trajectories, denoted as *targeted navigation* here, show how efficient people are at finding short paths [23, 76, 77]. In contrast to natural navigation, targeted navigation posits an unambiguous definition of success (i.e., reaching the target article), which allows researchers to study how users drift away from the best path and when they abandon their search [32, 64]. Targeted navigation behavior as observed in navigation games may, however, differ from natural navigation behavior, which limits the utility of such traces for studying the real-world usage of Wikipedia.

3 MATERIALS AND METHODS

The data sources exploited in this study include user traces mined from Wikipedia's server logs and features extracted from articles.

3.1 Pageloads

To study how readers navigate Wikipedia, we analyze the server logs of the English language edition collected for four weeks between 1 and 28 March 2021. This data contains an entry for each time a Wikipedia page is loaded. It is continuously and automatically collected for analytic purposes on Wikimedia's infrastructure and deleted after 90 days.

We limit our analysis to the pageload requests for articles (MediaWiki namespace 0), filtering out requests from bots. To protect readers' privacy, we remove sensitive information in several steps: discarding pageloads from readers who edited or were logged in during the time of data collection; discarding all requests from countries with at least one day with fewer than 300 pageloads; generating (pseudo) user identifiers by hashing IPs and user agent strings, as done in previous work

Origin	Desktop	Mobile	Total
Search engines	45.97%	48.77%	47.71%
Wikipedia			
Articles	35.64%	35.75%	35.72%
Main page	1.65%	0.70%	1.06%
Lang. switching	1.62%	0.50%	0.92%
Categories	0.59%	0.25%	0.39%
Search page	0.38%	0.22%	0.29%
Special pages	0.07%	0.01%	0.03%
Portals	0.03%	0.01%	0.02%
Others	0.07%	0.01%	0.03%
Unspecified origin	12.64%	13.03%	12.88%
External websites	1.36%	0.70%	0.95%

Table 1. Statistics of referrers of single pageloads.

[52]; and dropping IP, user agent, and fine-grained geo information. In total, these anonymization steps lead to the removal of around 3% of the data. In addition, we perform the following filtering steps. First, we drop pageloads of the *Main_Page* article, as it does not represent any specific entity. These requests may, e.g., come from users who set Wikipedia as the browser’s default page. Second, we remove traffic from massively common IPs, which would make it hard to study individual users’ activities, by dropping all user identifiers with more than 2,800 pageloads, or on average 100 per day, thus removing 28k (0.0019%) user identifiers. After the above steps, each request entry includes the anonymous user identifier, the page title, the timezone-corrected timestamp, the access method (mobile or desktop), and the referrer URL. The final dataset contains 6.52B pageloads associated with 1.47B user identifiers.

3.2 Article features

To characterize the content viewed by readers, we collect a set of article features. To ensure alignment between the server logs and the articles’ content, we compute the features for the revisions of the public snapshot released at the end of March 2021.

We obtain article features such as the number of outgoing links, the PageRank, article quality score, and topic. We assign the quality of the articles using the *articlequality* model of ORES¹ [21], Wikipedia’s official scoring platform. This model offers a way to obtain a score [20] that summarizes the structural properties of the article, such as the number of sections, references, and the presence of infoboxes. To represent articles semantically, we use two approaches: (1) the probabilities for 64 manually curated topics obtained from the ORES [21] *articletopic* model let us assign topical labels to articles; (2) the crosslingual WikiPDA [56] topic model lets us place articles in a 300-dimensional topic space.

4 RQ1: HOW DO READERS REACH WIKIPEDIA ARTICLES?

In this work, we use the term “*n*-gram” to designate a sequence of *n* subsequent Wikipedia pageloads from the same user, where the “vocabulary” consists of all articles available on Wikipedia. We start

¹<https://www.mediawiki.org/wiki/ORES>

Device	AB	AA	ABC	ABA	ABB	AAB	AAA
Desktop	0.900	0.099	0.749	0.121	0.047	0.049	0.031
Mobile	0.880	0.119	0.719	0.143	0.055	0.053	0.027
Total	0.888	0.111	0.732	0.134	0.052	0.052	0.029

Table 2. Frequencies of bigram and trigram patterns.

our analysis with unigrams ($n = 1$) to investigate individual pageloads and enumerate how readers can reach Wikipedia articles. We classify Web traffic according to HTTP referrers and quantify the frequency of each referrer type (Table 1). In total, 4B (61.5%) pageloads have external or empty referrers and are thus entry points to Wikipedia.

Search engines. The most common way to reach the content of Wikipedia is through external search engines, at 3.1B pageloads (45.9% of all recorded traffic, or 77.5% of external traffic). This volume reflects the significant value offered by Wikipedia in fulfilling the information needs of search engine users [2, 73].

Wikipedia. Clicks from other articles account for 35.7% of all traffic. Interestingly, as observed in previous work [47], 6.6% of these pageloads happen through links that do not exist in the link network itself, but likely through other interactions such as Wikipedia’s search drop down menu. Content can also be reached from other pages on the Wikipedia platform: (1) the main page, (2) category pages, (3) Wikipedia’s internal search, (4) portals, or (5) other Wikipedia pages, including talk pages or pages in other languages (language switching).

Unspecified origin. In 12.9% of all traffic, we observe an empty referrer field (20.9% of external traffic). Multiple reasons can produce a request without an explicit origin, including direct access via the browser history, redirects from apps, bookmarks, search toolbars, or when the link source has explicitly turned on the *noreferrer* property.

External websites. In total, only 0.95% of the requests originated from external websites that are not search engines nor Wikipedia domains (1.55% of the external traffic). Among those, the most common sources are Facebook (15.6%), Reddit (9.6%), YouTube (8.0%), and Twitter (4.3%).

Others. Other external visits (0.015% of external traffic) come from Android Web views and custom embedded visualizations, with the most common being the Telegram and Reddit sync apps, and Facebook on Android devices.

5 RQ2: HOW DO READERS TRANSITION FROM ONE ARTICLE TO THE NEXT?

Next, we move from unigrams ($n = 1$) to bigrams ($n = 2$), in order to understand how readers transition between Wikipedia articles. We study events aggregated by user identifier and sorted by time to investigate the properties of consecutive pageloads and their inter-event time. We consider two subsequent pageloads from the same user identifier as a bigram if they are separated by less than one hour [22].

Since here we are not interested in the exact article visited, we instead represent each pageload in a bigram with an alias indicating if the reader loaded the same page or different pages. The pattern “AA” means that the user revisited sequentially the same article, whereas “AB” indicates a load of two different pages. Here it is important to note that the Wikipedia server instructs the browser to disable the cache, such that the server logs contain essentially all pageload events, including cases when the readers reloaded an article, e.g., by using the back button.

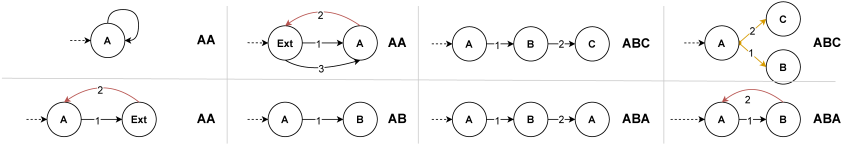


Fig. 1. Examples of patterns in the logs and the multitude of client-side behaviors that can generate these patterns. Black arrows represent forward link clicks, red arrows represent back-button clicks, yellow arrows represent clicks that open multiple tabs from the same source page. “Ext” represents external (non-Wikipedia) pages. Numbers represent the order of clicks.

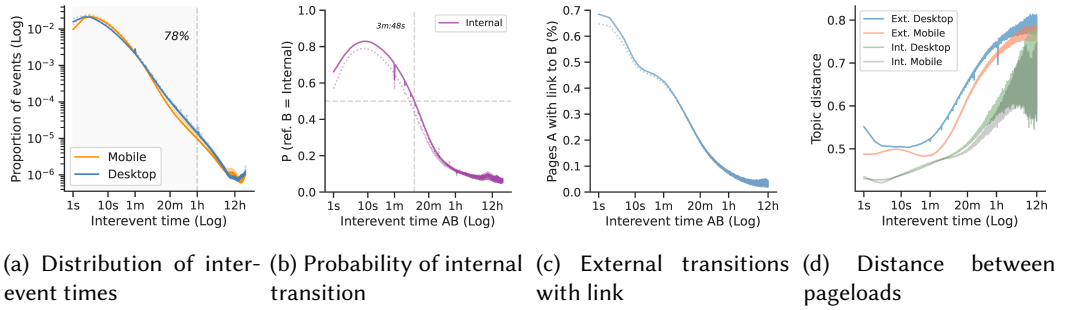


Fig. 2. Statistics of bigrams as a function of the inter-event time between two pageloads. Dotted curves represent the distributions with AA patterns included.

Bigrams. The logs contain 3.95B instances of bigrams. The emerging patterns, described next, are summarized in Table 2. The most frequent bigram pattern (“AB” in Table 2) corresponds to transitions between two different articles. It can happen both through internal and external navigation (cf. Fig. 1). This pattern represents around 89% of all bigrams. The other possible bigram pattern (“AA” in Table 2), corresponds to the consecutive reload of the same article. Representing 11% of all bigrams, it is rather common (84% share the same referrer). This pattern appears at least once in 37% of the navigation histories of readers with at least two pageloads in the month of data collection. The pattern can be generated by different client behaviors (cf. Fig. 1), including repeated consumption as described in previous work [6, 70], user activities involving external navigation, or artificial reloads by the browser when a tab unloaded from memory is restored.

Trigrams. Finally, we also briefly consider the 2.98B trigrams present in the logs. The most common trigram pattern (73%, “ABC” in Table 2) represents transitions between three different articles. A variety of behaviors can generate this pattern, including sequential clicks or multitab behavior (cf. Fig. 1). The second most common trigram pattern (13%, “ABA” in Table 2) can be generated by intentionally revisiting the same page in a forward manner or by clicking the back button (cf. Fig. 1). In 89% of ABA instances, the first and last event also share the same referrer. The remaining trigram patterns (ABB, AAB, AAA) are combinations of the bigrams described above.

Dynamics of transitions. In order to understand the dynamics of these transitions, we investigate the inter-event time between the two pageloads in each bigram. The interval between two consecutive pageloads peaks at very short times, with a median of 74 seconds (63 and 93 seconds for mobile and desktop devices, respectively). However, as Fig. 2a shows, the distribution is long-tailed, with 22% of pairs separated by more than one hour.

Investigating the referrer of the second page of the bigrams reveals that readers frequently do not use internal links to transition between two articles, but external pages by leaving and re-entering Wikipedia. These external transitions are not rare: in 35.2% (or 40.1% when including AA patterns) of the bigrams with less than one hour between the two events, the second page was reached through external navigation. This observation is corroborated by Fig. 2b, which shows that for pairs with an inter-event time greater than 3 minutes and 48 seconds, transitions via internal links are even less common than transitions via external navigation. External transitions tend to be semantically coherent: considering all 1.4B AB-type bigrams where the second page is reached via search, in 18% of the cases, the first page explicitly contained the link. This proportion increases to 30% [56%] when considering pairs with an inter-event time of less than one hour [less than 10 seconds] (Fig. 2c). The topical coherence of these transitions is also visible in Fig. 2d, which plots the average topical distance (measured by the cosine of WikiPDA vectors, cf. Sec. 3.2) as a function of inter-event time, showing that external navigation recorded within a few minutes from the previous pageload shows topical distance comparable to internal navigation.

6 RQ3: WHAT ARE THE PROPERTIES OF ENTIRE READING SESSIONS?

Using our insights about navigation at the unigram, bigram, and trigram levels, we can now characterize entire navigation sessions. We start by introducing two different approaches to conceptualizing navigation sessions (Sec. 6.1) and discuss how each captures different aspects of reader navigation. We then describe the properties of reader navigation by focusing on three aspects of the resulting sessions: contextual features defining when and how sessions start (Sec. 6.2); static properties, such as the structural features of sessions (Sec. 6.3); and finally, the dynamic properties of the sessions, such as the evolution in the content consumed over the course of navigation (Sec. 6.4).

6.1 Conceptualizing reader sessions

Grouping all pageloads of the same user, there is no unique way to operationalize the notion of a reading session. Based on different previously employed approaches, we identify two distinct notions of a session: (1) *navigation trees* connect pageloads hierarchically based on referrer information, whereas (2) *reading sequences* order pageloads linearly based on temporal information. These capture different aspects of how readers navigate, and which approach is better suited depends on the context and the phenomenon one aims to observe. From the original 6.52B pageloads, we obtain 3.7B navigation trees and 2.51B reading sequences.

Navigation trees [52] describe how readers traverse Wikipedia by following internal links. We generate a tree by connecting pageloads via the referrer contained in HTTP headers. Pages reached through internal transitions (i.e., using internal links) are added as children of the most recent load of the article in the referrer, while pageloads with external or *Main_Page* referrers generate a new tree. If a page is loaded multiple times from the same referrer, the parent node retains only the first instance as a child. This method has the advantage of representing coherent sessions created through clicks on internal links—regardless of the time spent on one article—and of capturing multitab behavior [25]. The downside is the difficulty of capturing content consumption over time for subsequent pages not reached through internal clicks, even if close in time (a common pattern, cf. Sec. 5). Since this aggregation method does not model temporally linear consumption, loading articles by opening multiple tabs or backtracking to select a different path leads to the same navigation tree.

Reading sequences describe how readers consume content in temporal order. They are defined as linear sequences of all pageloads by the same user ordered by time. Sequences are split if the inter-event time between two consecutive pageloads separated by external navigation exceeds a

threshold value of one hour, following recommendations from previous studies [22] and common practice [40, 66]. Within such sessions, we keep only the first pageload of each article, in order to only capture the first exposure of the respective content. This method generates topically less coherent sessions, capturing the temporal and linear sequence of pageloads of a reader within a defined period of time, both via internal and external transitions (e.g., multiple external searches). This method has the disadvantage of being a simplification of how readers explore the link network, and a fixed threshold of one hour may not be ideal in every context.

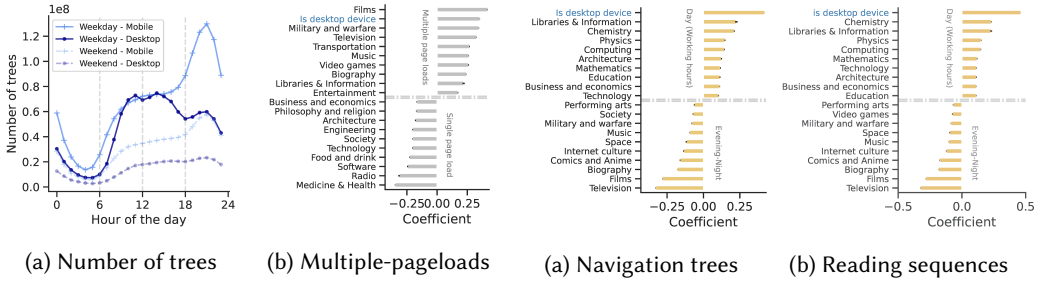


Fig. 3. The total number of trees started at different times of day (Fig. 3a) and feature contributions to the logistic model predicting if the reading sequence is composed of more than one pageload (Fig. 7d).

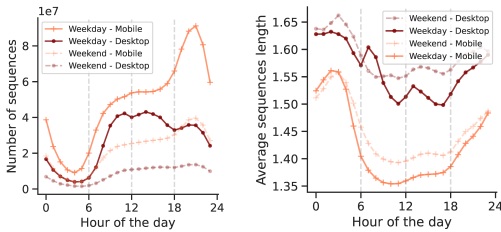


Fig. 5. Total count (Fig. 5a) and average length (Fig. 5b) of reading sequences started at different times of day.

6.2 Contextual properties: time and device

We study the context of a session by focusing on the time of the first pageload and the device used to access Wikipedia. This section focuses on navigation trees, but reading sequences give qualitatively similar results (cf. Fig. 4b, Fig. 6a).

Time. To remove confounding via different timezones, we use geolocation information to normalize the time of all pageloads to local time. The distribution of session starting times follows a regular circadian rhythm (Fig. 3a and Fig. 5a). Both access methods (desktop and mobile) show a similar pattern during the day, with a substantial increase of mobile sessions in the evening. Wikipedia has fewer sessions during weekends, but with similar temporal distributions as working days. The

Fig. 4. Feature contributions to the logistic model predicting if the reading session started during daytime, for navigation trees (Fig. 4a) and reading sequences (Fig. 4b).

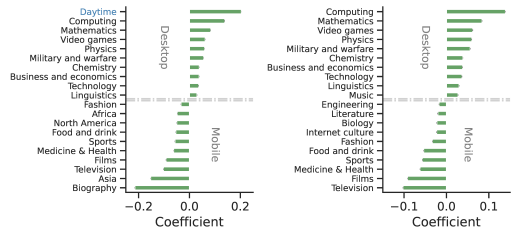


Fig. 6. Feature contributions to a logistic model predicting if the session is started from a mobile or desktop device.

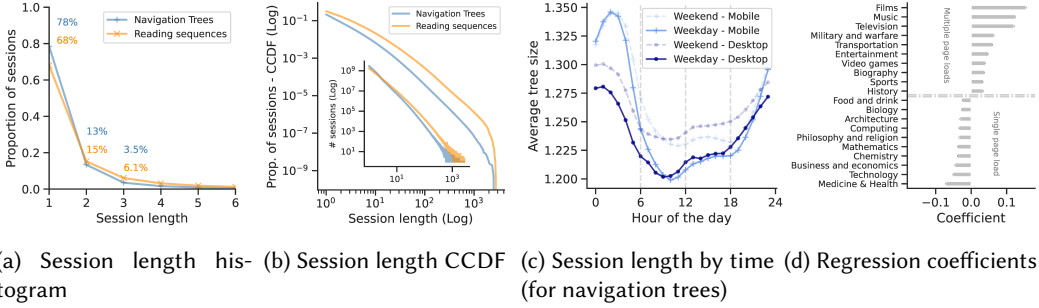


Fig. 7. Session-length statistics.

desktop distribution shows dents at 12:00 and 18:00, mirroring work rhythms with a lunch break around noon and the end of work in the evening (and possibly commuting).

In order to understand which features are associated with requests at different times of day, we fitted a logistic regression model to predict if a pageload was observed during the day or evening/night. We represent each pageload by its topic probabilities (obtained from ORES, cf. Sec. 3.2) and the type of device (desktop or mobile). Binarizing the target variable by representing daytime (9:00–18:00) as the positive class, we obtain an AUC/ROC of 0.586 on a held-out test set. Inspecting the fitted feature weights (Fig. 4a) shows that desktop devices and articles associated with STEM and education are associated with sessions starting during the day, whereas topics about entertainment are predictors of sessions starting during the evening or night.

Device. Fig. 3a indicates that people prefer different devices at different times of day. Next, we study whether specific topics are associated with device types by representing each pageload with the vector of topic probabilities (obtained from ORES) and a feature indicating if the page was loaded during the daytime. We again fit a logistic regression to predict the device used, with an AUC of 0.639. Inspecting feature importance shows that people tend to access STEM and business content from desktop devices, and biographies, entertainment, and medicine from mobile devices (Fig. 6).

6.3 Static properties: structure of sessions

Session length. We measure session length as the number of pageloads in the navigation tree or the reading sequence, respectively. Most sessions consist of a single pageload (Fig. 7a), but the length distribution also exposes a long tail (Fig. 7b). Therefore, we summarize session lengths via the geometric mean (arithmetic mean in parentheses). By construction, reading sequences tend to be longer because, unlike navigation trees, they merge both external and internal transitions.

In the case of reading sequences, the average session length shows differences with respect to the access method, with an average length of 1.41 (1.99) for mobile, and 1.54 (2.40) for desktop. This difference is less pronounced for navigation trees, where mobile sessions contain on average 1.23 (1.5) articles, vs. 1.24 (1.5) for desktop. The average session length varies during the day, with readers engaging in longer sessions during the evening and night, for both navigation trees and reading sequences (Fig. 7c and Fig. 5b).

To understand what properties are associated with short sessions consisting of a single pageload, we fitted a logistic regression to predict if the reader will continue after loading the first page in a navigation tree (results are qualitatively identical for reading sequences), representing each first pageload with its topic probabilities (obtained from ORES), device type, and time of day, and

Tree size			
Top 10 (larger trees)	Bottom 10 (smaller trees)		
1.377	Films	1.152	Earth and environment
1.373	Entertainment	1.148	Food and drink
1.340	Television	1.145	Biology
1.327	Military and warfare	1.138	Technology
1.324	Music	1.128	Physics
1.295	Comics and Anime	1.122	Software
1.284	History	1.114	Medicine & Health
1.272	Biography	1.112	Computing
1.269	Sports	1.104	Mathematics
1.264	Transportation	1.100	Chemistry

Table 3. Top and bottom 10 topics with respect to (geometric) average tree size (geographical topics excluded).

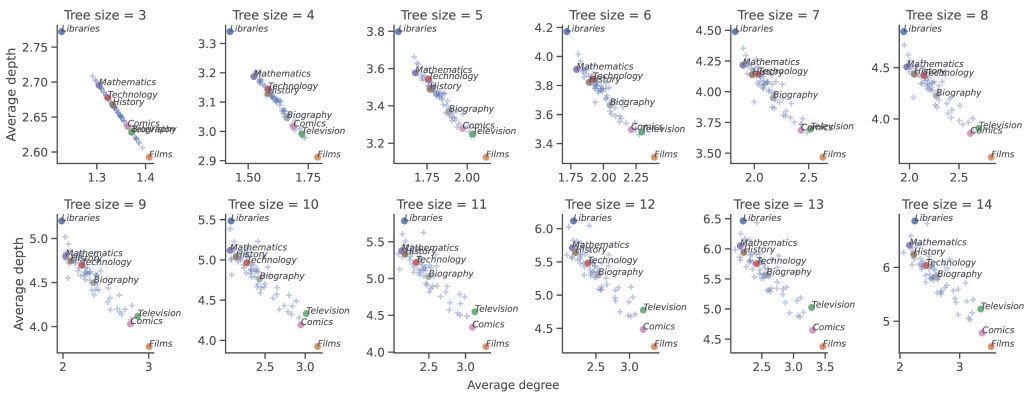


Fig. 8. Relation between the average depth and average degree for navigation trees of different sizes.

obtaining a model with an AUC/ROC of 0.606 on a held-out test set. Inspecting the coefficients of the regression (Fig. 7d), we find that longer [shorter] sessions are associated with topical content around entertainment [STEM and medicine]. This observation is corroborated by the substantial difference in average navigation tree size across topics (Table 3).

Shape of navigation trees. In order to better understand how readers navigate the link network, we analyze the shape of navigation trees (in contrast, the shape of reading sequences is, by construction, always a linear chain). The three most common patterns (Fig. 9, left) are described as follows, in order of decreasing frequency: (1) a linear chain of pageloads; (2) fanning out from one page to several different pages, e.g., by opening multiple tabs or rolling back and selecting a different path; (3) a combination of the two (one-step chain followed by fanning out). These three patterns remain the most frequent for all tree sizes (Fig. 9, right).

We further characterize the different strategies associated with navigation trees in terms of tree depth (i.e., average length of paths from the root to the leaves) and breadth (i.e., average out-degree of non-leaves in the tree) for trees starting with different topics. Noting that the two metrics are almost perfectly anti-correlated and that the relative ordering of topics is stable across all tree sizes (Fig. 8), we define an aggregate tree-breadth ranking for each topic based on the average rank across tree sizes (Table 4). This shows that entertainment topics are associated with wider trees

Top 10 (wider trees)			Bottom 10 (deeper trees)		
Rank (mean)	SD	Root topic	Rank (mean)	SD	Root topic
1.00	0.00	Films	27.42	2.72	Linguistics
2.50	0.87	Television	29.42	0.95	Earth and environment
3.58	0.76	Entertainment	29.50	1.19	Space
4.50	1.85	Comics and Anime	30.08	2.78	History
4.67	1.31	Education	31.92	1.11	Computing
6.58	1.98	Video games	32.92	1.55	Software
7.92	2.43	Literature	34.67	1.75	Chemistry
8.50	2.36	Fashion	34.75	1.30	Physics
8.83	1.07	Performing arts	35.50	1.26	Mathematics
10.42	2.29	Internet culture	35.67	1.65	Libraries & Information

Table 4. Rank with respect to average degree of navigation trees, by topic (geographical topics excluded). A separate rank was computed per tree size (3–15), and arithmetic means over tree sizes are reported, alongside standard deviations (SD).

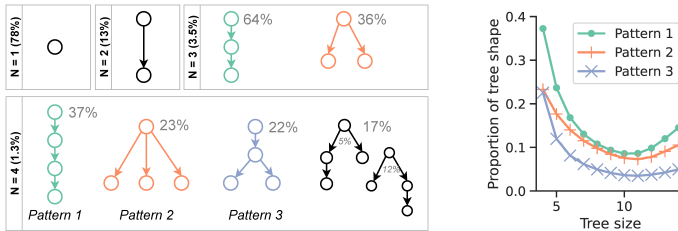


Fig. 9. Shape of navigation trees. Frequency of patterns for trees size $N \leq 4$ (left panel). Dominance of top three patterns (see main text) for larger trees (right panel).

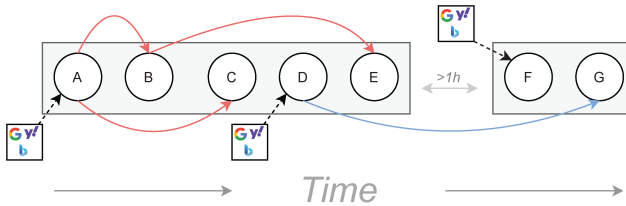


Fig. 10. This set of log events yields three navigation trees, represented by arrows and composed of ABCE, DG, and F. The reading sequences method creates two sessions represented as gray boxes: ABCDE and FG. Square boxes are clicks from external origins.

with higher branching, and STEM topics are characterized by deeper trees with a more chain-like structure.

6.4 Dynamic properties: within-session article-property evolution

To shed light on navigation dynamics, we track the evolution of different article properties within sessions. Our evolution analysis revolves around three domains: topic space (distance from the first and previous articles), quality, and network centrality (out-degree and PageRank). Here, reading sequences are represented as defined above, whereas a navigation tree is represented by the linear

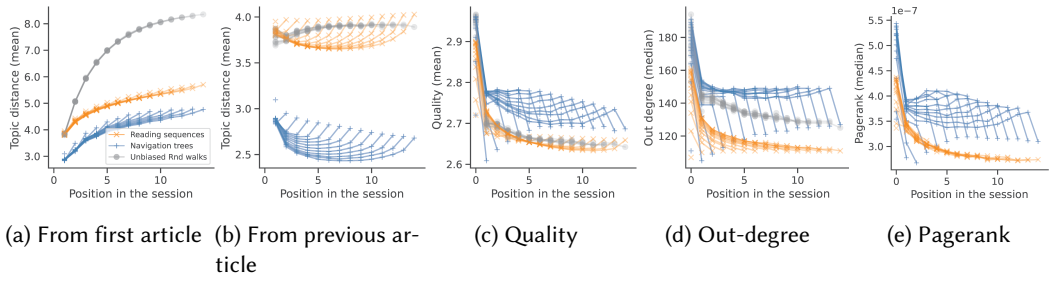


Fig. 11. Within-session evolution of five article properties. Each curve represents sessions of different lengths.

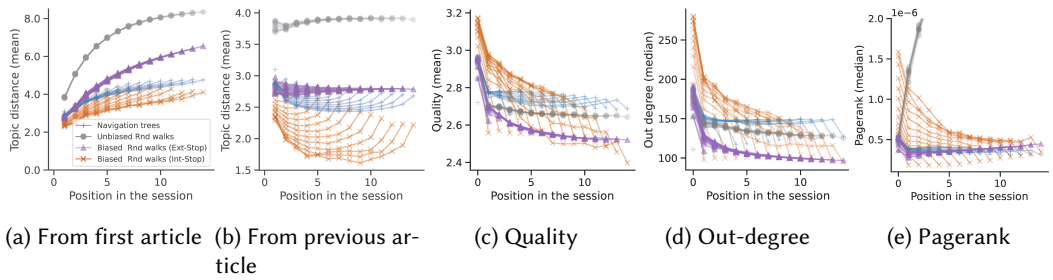


Fig. 12. Property evolution of the trajectories generated by the three random walk models, compared with natural navigation as captured by navigation trees. Each curve represents sessions of different lengths.

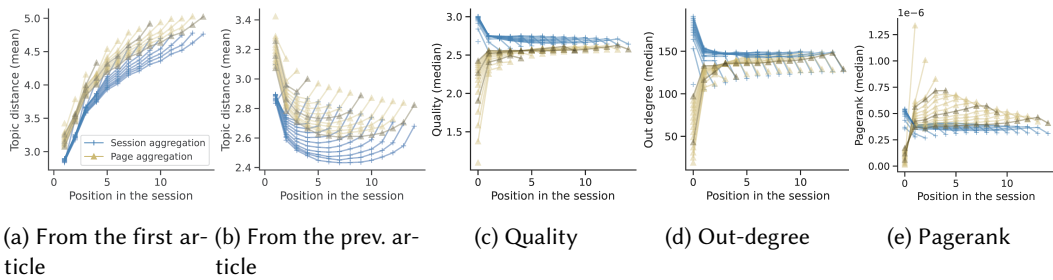


Fig. 13. Comparison of the evolution of five different properties when aggregating navigation trees by session (micro-average, blue) and by starting page (macro-average, gold). Gray trajectories highlighted for readability. Each curve represents sessions of different lengths.

path from the root to the temporally last leaf, from where the reader ceased to click further via internal links.

It is important to note that these two approaches can produce different sequences of pageloads: e.g., a pageload in position 1 of a navigation tree could be in position 4 of a reading sequence (as in Fig. 10). Also, the last pageload of each sequence can have different interpretations: for navigation trees, the reader stopped link-based navigation on that page, whereas for reading sequences, the reader did not load a Wikipedia page for at least one hour.

In order to better interpret our observations, we compare them with three null models corresponding to different random walkers. The null models serve as a comparison to assess to which degree the observed properties of the navigation dynamics are due to chance. We randomly sample

120M paths from the navigation trees, and run (from the tree’s starting article) (1) an *unbiased random walker* that selects the next step with uniform probability from the available links and generates a sequence of the same length as the original path; (2) an *extrinsic-stop biased random walker* that selects the next step based on the pairwise transition probabilities obtained from the public clickstream and generates a sequence of the same length as the original path; (3) an *intrinsic-stop biased random walker* that selects the next step—or stops—based on the pairwise transition probabilities from the public clickstream [61]. We consider sessions up to length 15, stratifying by session length.

Topic space. We measure the topical distance between articles via the Kullback–Leibler (KL) divergence of their respective WikiPDA topic distribution vectors (Sec. 3.2). For robustness, we tried different topic models (WikiPDA and ORES) and different distance metrics (KL divergence, Euclidean, cosine, and Wasserstein), obtaining qualitatively similar results. First, we study how readers diffuse in topic space starting from the first article, which plays a special role, as it represents the entry point to Wikipedia. On average, readers diffuse in topic space, moving further from the first article with every step (Fig. 11a). Reading sequences and navigation trees exhibit the same trend, with a shift due to the tendency of reading sequences to ignore external navigation. All the random walkers show similar increasing trajectories (Fig. 12a), diffusing faster than natural navigation when the random walker is unbiased, or biased but extrinsically stopped.

Second, we measure the semantic step size in topic space by tracking how the topical distance to the previous article evolves. Both navigation trees and reading sequences exhibit a U-shape, suggesting that readers tend to first reduce their semantic step size, before diverging and finally abandoning (Fig. 11b). The discrepancy between navigation trees and reading sequences is consistent with the previous observation on diffusion from the first article. Interestingly, this U-shape is similar to the trajectories generated by the intrinsic-stop biased random walker (Fig. 12b), as also reported in previous work [61]. In contrast, the other two random walk models show that by selecting a random link or stopping at predefined lengths, the average distance from the previous article tends to stabilize to an equilibrium value.

Quality. The evolution of article quality shows a sharp drop at the beginning, for both reading sequences and navigation trees (Fig. 11c). This behavior can be interpreted as a form of regression to the mean, since many sessions start from popular pages with high quality, which thus contribute more to the distribution. By moving one step in the link network, readers naturally reach a page that is, on average, of lower quality. The intuition is confirmed by the behavior of the unbiased random walker, which shows the same drop with the first step (Fig. 12c).

In contrast to reading sequences, navigation trees show a sharp drop in quality with the *last* pageload. This indicates that readers have a higher chance to stop Wikipedia-internal navigation when reaching a low-quality page, and as a result, continue navigating in a different branch of the tree or via an external transition.

Compared to the random walkers (Fig. 12c), readers tend to navigate across pages with less variance in quality. The random walkers’ traces support the hypothesis that there are articles with a higher chance of terminating the navigation: while the unbiased and extrinsic-stop biased walkers show no termination pattern, the intrinsic-stop biased walker shows a final drop as in human navigation. The organic stopping of this random walker, mirroring readers’ behavior more closely, increases the chances to abandon the navigation on pages of low quality that, according to the clickstream data, relay less traffic.

Network centrality. Finally, we are interested in how reader sessions evolve in the network with respect to different centrality measures. We start with out-degree (the number of outgoing links

in article bodies). Similar to article quality, the out-degree shows a sharp drop with the first step (Fig. 11d) for navigation trees and reading sequences, likely caused by the presence of many sessions starting from pages with a particularly high out-degree. We also find a sharp drop for the last pageload in the sequence of the navigation trees, suggesting that readers have a higher chance of stopping Wikipedia-internal navigation upon reaching a page with low out-degree.

In the case of the random walkers, we draw similar conclusions as for article quality. Whereas unbiased random walks and extrinsic-stop biased random walks show a decrease and stabilization of out-degree, the intrinsic-stop random walker, as humans, terminates on pages of lower degree (Fig. 12d). Compared to random walkers, human navigation is more stable: after the initial drop, they have a higher chance to stay on pages with around 150 links.

Finally, we characterize how the PageRank of visited articles changes during sessions. We observe that the PageRank mirrors the evolution of quality and out-degree with regard to the initial drop (Fig. 11e). Readers tend to enter more frequently on popular pages with high centrality and naturally move to a less central node in one step. Also for this case, a drop is visible in the last step of the navigation trees, indicating that, when the readers reach an article leading to the network periphery, they have higher chances to stop the Wikipedia-internal navigation. The random walkers (Fig. 12e) show that unbiased walks naturally converge in a few step to the most central pages with very high PageRank. The extrinsic-stop biased walker, on the contrary, after an initial drop, tends to move to central nodes at a much lower speed. Finally, the intrinsic-stop biased walker, again, shows a final drop from a stable value before abandoning the navigation, similar to human readers.

Aggregation by page. The quantities in Fig. 11 correspond to a micro-average over all sessions, where the average behavior can be dominated by sessions starting from the most popular pages since the overall distribution of pageviews is highly skewed. Therefore, we also calculate a macro-average by aggregating on a starting-page level to make each first article contribute equally. The diffusion in topic space is qualitatively similar in both aggregation methods (Fig. 13a and Fig. 13b). In contrast, for quality, out-degree, and PageRank, the overall trend is inverted, i.e., instead of a sharp drop, we observe a sharp increase in these metrics after the first step (Fig. 13c, Fig. 13d, and Fig. 13e). This discrepancy could be caused by the presence of many low-quality [53] and low-degree articles, such that readers at the first step tend to move to better articles in search of information (a sort of regression to the mean). Interestingly, the drop towards the last pageload in a session appears across both aggregation methods.

7 DISCUSSION

7.1 Summary of findings

We have provided a systematic characterization of the navigation pathways of Wikipedia readers through a large-scale study of the site's server logs. Starting from the raw logs, we aggregated the data in navigation trails to quantify how readers reach, and transition between, pages. First, the most common way to reach a page is through an external search engine, followed in frequency by internal navigation from other Wikipedia articles; other sources, such as external websites (mostly social media sites) and other Wikipedia content (such as categories or special pages), are much less frequent, but still substantial in absolute numbers. Second, readers frequently transition between pages via external search engines instead of using direct Wikipedia links. These external transitions are characterized by larger topical jumps and larger inter-event times between pageloads; they must, however, still be considered semantically meaningful, for, in many cases, a link for internal navigation—even if not taken—would still be available. Third, by analyzing sequential patterns, we find that consecutive reloads and revisits of previously visited articles are common (10% or more each).

We continued by characterizing how readers combine the above patterns into extended navigation sequences. First, we introduced two approaches to capture paths of readers: *navigation trees* based only on internal navigation, and *reading sequences* based on the time-ordered pageloads including internal and external transitions. Second, we described how sessions are affected by their context in terms of device type and time of day. We find that topics related to STEM [entertainment] are more associated with working [evening and night] hours. Third, we measured the size and structure of sessions. While most sessions consist of a single pageload (68–78% depending on the aggregation method), the size distribution shows a long tail with tens of millions of sessions consisting of 10 or more pageloads. The topic not only affects the size but also the shape of trees: while sessions starting from articles on entertainment generally consist of more pageloads, such trees are also broader (higher branching factor) than sessions starting, e.g., from STEM topics, which are smaller and deeper. Fourth, we investigated the within-session evolution of article properties. In topic space, longer sessions diffuse ever further away from the origin, with semantic step size following a characteristic U-shape pattern suggesting that readers reduce their semantic step size first, before diverging in ever larger steps and finally abandoning the session. The first and last pageload of a session show special behavior regarding the evolution of article quality and network centrality. More popular (and thus higher-quality and higher-centrality) pages are naturally more common as first articles, thus engendering a form of regression to the mean with the second step. An inverted effect appears when sessions are aggregated at the starting-page level, such that every starting article is represented equally. Either way, articles at the end of the navigation are typically lower-quality pages, suggesting that readers stop following the internal navigation when they reach these pages, which thus act as network sinks.

7.2 Implications

Complexity of navigation behavior. Our results show that the navigation paths extracted from Wikipedia’s server logs constitute a non-trivial dataset requiring extreme care in order to avoid drawing spurious conclusions. First, in contrast to existing pre-processing pipelines for sequence analysis (e.g., tokenization, stopword removal, stemming, etc., in NLP), we still lack an understanding of universal best practices for navigation paths, and as a result we had to investigate and compare alternative strategies for conceptualizing sessions—namely, reading sequences vs. navigation trees. Second, operationalizing navigation paths makes strong assumptions: while navigation trees from pure internal navigation are more topically coherent with more complex structure, reading sequences from temporally ordering all of the user’s pageloads are less coherent but provide a linear sequence that is not broken by external searching (which is common). The latter typically introduces an additional cutoff for sessions if consecutive pageloads are separated by more than one hour [22]; however, our analysis suggests other potential data-informed choices, such as the time separation of internal and external transitions at approximately four minutes (Fig. 2b). Naturally, the suitable choice depends on the question of interest. Third, our analysis shows that the data can exhibit Simpson’s paradoxes; e.g., the inversion of the within-session evolution of page properties such as PageRank (Fig. 13) depends on the aggregation level. Fourth, the prevalence of trivial patterns (e.g., reload or revisit) points to potential caveats when applying prediction models to session-based recommendation [75].

Diversity. There is extraordinary diversity in the ways readers browse Wikipedia, modulated by topic, device, time of day, etc. This reflects the diversity found in previous studies on the different motivations and information needs of readers across the globe [29, 40, 66]. This heterogeneity indicates caution against simplistic models aiming to capture a single average behavior.

Online ecosystem. The usage of Wikipedia is embedded in a larger online ecosystem. Multiple studies have shown the importance of Wikipedia to search engines [46, 73], as a gateway to the Web [54, 55], and as a main educational resource for online learning more generally [33]. Our results show that this interplay between external and internal (with respect to Wikipedia) also plays a crucial role on an intra-session level when navigating encyclopedic information.

Navigation in the wild. The navigation of readers on Wikipedia differs from targeted navigation in lab-based settings [23, 76, 77]. We do not observe typical strategies characterized by, e.g., navigation via hubs (an initial increase, followed by a drop, in out-degree) or gradually decreasing the step size in semantic space towards a target. Instead, we find a range of other patterns, such as a U-shape for the step-size in semantic space and an immediate sharp drop followed by largely constant centrality measures (out-degree, PageRank). This highlights conceptual limitations of targeted-navigation experiments with respect to generalizing their results to how humans seek knowledge more generally.

Furthermore, our results provide a more nuanced picture on the conclusions derived from publicly available data, most notably the Wikipedia clickstream [86], which provides aggregate data on the number of times a link was clicked. For example, we can observe that the overall tendency to navigate towards peripheral nodes [14] is mainly driven by the first step after reaching Wikipedia, with subsequent steps showing much smaller differences in centrality measures (with the exception of the last step, see below). One possible interpretation is a regression-to-the-mean effect as popular pages (the starting points of navigation) are generally skewed towards higher centrality and quality.

Content and navigation. Our results contribute to describing the relation between content and navigation, expanding the prior understanding of how readership and popularity are influenced by visual position [14] or quality [88]. Our results go beyond the population level, suggesting that upon encountering low-quality pages readers tend to stop navigating along a specific branch in the navigation tree (and continuing along a different branch or stopping altogether). This is specifically important in the context of knowledge gaps in Wikipedia [60], in order to address the uneven representation of, e.g., articles on women, where a better understanding of the interaction between content, readers, and editors [16, 65] is crucial to allow for more informed decision-making in designing interventions.

More generally, this finding is aligned with the definition of information scent used in information foraging theory [9]. The theory states that, in analogy to animals following the scent of food, when seeking information, we rely on our intuition—or “built-in” foraging strategies—to pick the path that maximizes information intake while minimizing the investment of time and energy. In this view, readers foraging for information follow the scent with higher chances of leading to the desired content; when scent loses intensity, they move to more promising information sources. Additionally, our work can have implications for developing theoretical frameworks to describe navigation patterns. Understanding how readers follow specific trails can inform researchers about the distinct properties of the information scent that guide our search for information online. These findings can be instrumental in developing novel theories on how humans move in information networks.

Best practices for Wikipedia log analysis. One of the challenges in conducting analyses like those presented here is the lack of standard pipelines to preprocess and aggregate the logs. Unlike fields such as NLP or computer vision, where preprocessing steps are de facto standardized, modeling behavior from access logs does not have a unique standard procedure yet, and researchers are forced to make many modeling choices. The purpose of the study and limitations of the data, such as privacy concerns, may influence how sessions are defined and, consequently, the results obtained.

Our work fills this gap in the case of Wikipedia by providing best practices for processing server logs to study reader sessions. We describe two complementary approaches based on trees and temporal sequences and demonstrate their relative advantages and disadvantages. This operationalization is crucial for developing systematic approaches in future studies to understand reader navigation better, capture their information needs, and improve their experience on Wikipedia and the Web more generally.

7.3 Limitations and future work

Limitations. In terms of limitations, we capture navigation paths only via events in the server logs. Moving forward, how people engage with content could be more accurately observed via client-side instrumentation. The aggregation based on IP addresses and user agent information also has limitations; e.g., we had to discard the sessions of large organizations with shared IP addresses.

The navigation logs suggest that Wikipedia fulfills various information needs and readers exhibit diverse navigation patterns. Using large-scale digital traces offers important advantages over other methods when we are interested in the quantitative measurement of behavioral phenomena [63]. However, purely log-based analysis also has limitations, and it should be considered a complementary, and not a substitutive, approach. Previous work indicates that big-data analyses are not immune to biases introduced by algorithmic dynamics [38, 74], data collection problems, preprocessing errors, and measurement errors [37, 74].

Finally, we only focused on a single language, English. While this already revealed a rich spectrum of phenomena, additional variation can be expected from a comparison across languages [40].

Future work. To overcome these limitations, future work should capture the variation in navigation across Wikipedia's over 300 languages. Moreover, in order to better serve the different information needs of readers, a better understanding is needed regarding how patterns in navigation correspond to underlying motivations [66] and other traits such as curiosity [41]. By enriching the behavioral patterns with qualitative feedback, we can better understand user objectives and design ways to facilitate more efficient access to the desired information.

In line with previous studies [7, 26, 26], future work should also investigate the relationship between search queries (from external and internal origin) and subsequent navigation behavior on encyclopedic platforms such as Wikipedia. Finally, in order to capture encyclopedic information seeking more generally, researchers should capture navigation beyond individual platforms to take into account the interdependence of Wikipedia with the rest of the Web.

7.4 Conclusion

Seventy-seven years ago, in 1945, Vannevar Bush sketched his vision of an information management device—the “memex”—that would allow users to not only retrieve documents quickly, but to also easily interlink documents [8]. With the advent of the Web, the hyperlink structure envisioned by Bush has since become a reality—but Bush's vision went further: he saw the trails taken by users as first-class citizens of the hypertext environment, as important as the text content itself: “Wholly new forms of encyclopedias will appear, ready made with a mesh of associative trails running through them, ready to be dropped into the memex and there amplified” [8]. In this regard, our technological reality has not caught up with Bush's vision yet, and the present work should be seen as a small step toward achieving it: we have started by describing the “associative trails running through” Wikipedia, and we hope that its future versions will build on these insights to incorporate tools and features that will allow readers to continually benefit from each other's encyclopedic trail blazing.

AVAILABILITY OF DATA AND CODE

The underlying data from Wikipedia’s server logs are not publicly available due to privacy reasons. Code is available at https://github.com/epfl-dlab/how_readers_browse_wikipedia.

ACKNOWLEDGMENTS

We thank Leila Zia for insightful discussions. West’s lab is partly supported by grants from Swiss National Science Foundation (200021_185043), Swiss Data Science Center (P22_08), H2020 (952215), Microsoft Swiss Joint Research Center, and Google, and by generous gifts from Facebook, Google, and Microsoft.

REFERENCES

- [1] Ashton Anderson, Ravi Kumar, Andrew Tomkins, and Sergei Vassilvitskii. 2014. The Dynamics of Repeat Consumption. In *Proc. International World Wide Web Conference (WWW)*.
- [2] Dan Andreescu, Kinneret Gordon, Isaac Johnson, and Nicholas Perry. 2021. Searching for Wikipedia. <https://techblog.wikimedia.org/2021/06/07/searching-for-wikipedia/>. accessed: 13 October 2021.
- [3] Akhil Arora, Martin Gerlach, Tiziano Piccardi, Alberto García-Durán, and Robert West. 2022. Wikipedia Reader Navigation: When Synthetic Data Is Enough. In *WSDM ’22: Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. Association for Computing Machinery, New York, NY, USA, 16–26.
- [4] Mamoun A Awad and Latifur R Khan. 2007. Web navigation prediction using multiple evidence combination and domain knowledge. *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans* 37, 6 (2007), 1054–1062.
- [5] Marcia J Bates. 1989. The design of browsing and berrypicking techniques for the online search interface. *Online review* (1989).
- [6] Austin R. Benson, Ravi Kumar, and Andrew Tomkins. 2016. Modeling User Consumption Sequences. In *Proc. International World Wide Web Conference (WWW)*.
- [7] Mikhail Bilenko and Ryen W White. 2008. Mining the search trails of surfing crowds: identifying relevant websites from user activity. In *Proceedings of the 17th international conference on World Wide Web*. 51–60.
- [8] Vannevar Bush et al. 1945. As we may think. *The atlantic monthly* 176, 1 (1945), 101–108.
- [9] Ed H Chi, Peter Pirolli, Kim Chen, and James Pitkow. 2001. Using information scent to model user information needs and actions and the Web. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 490–497.
- [10] Flavio Chierichetti, Ravi Kumar, Prabhakar Raghavan, and Tamas Sarlos. 2012. Are Web Users Really Markovian?. In *Proc. International World Wide Web Conference (WWW)*.
- [11] Alexander Dallmann, Thomas Niebler, Florian Lemmerich, and Andreas Hotho. 2016. Extracting semantics from random walks on wikipedia: Comparing learning and counting methods. In *Proc. Conference on Web and Social Media (ICWSM)*.
- [12] Mukund Deshpande and George Karypis. 2004. Selective markov models for predicting web page accesses. *ACM transactions on internet technology (TOIT)* 4, 2 (2004), 163–184.
- [13] Dimitar Dimitrov, Florian Lemmerich, Fabian Flöck, and Markus Strohmaier. 2018. Query for Architecture, Click through Military: Comparing the Roles of Search and Navigation on Wikipedia. In *Proc. Conference on Web Science (WebSci)*.
- [14] Dimitar Dimitrov, Philipp Singer, Florian Lemmerich, and Markus Strohmaier. 2017. What Makes a Link Successful on Wikipedia?. In *Proc. International World Wide Web Conference (WWW)*.
- [15] Carsten Eickhoff, Jaime Teevan, Ryen White, and Susan Dumais. 2014. Lessons from the journey: a query log analysis of within-session learning. In *Proceedings of the 7th ACM international conference on Web search and data mining*. 223–232.
- [16] Wikimedia Foundation. 2019. Medium-term plan 2019: The model for engagement. https://meta.wikimedia.org/wiki/Wikimedia_Foundation_Medium-term_plan_2019#The_model_for_engagement. accessed: 13 October 2021.
- [17] Steve Fox, Kuldeep Karnawat, Mark Mydland, Susan Dumais, and Thomas White. 2005. Evaluating implicit measures to improve web search. *ACM Transactions on Information Systems (TOIS)* 23, 2 (2005), 147–168.
- [18] Ruili Geng and Jeff Tian. 2015. Improving Web Navigation Usability by Comparing Actual and Anticipated Usage. *IEEE Transactions on Human-Machine Systems* 45, 1 (2015), 84–94.
- [19] Patrick Gildersleve and Taha Yasseri. 2018. Inspiration, Captivation, and Misdirection: Emergent Properties in Networks of Online Navigation. *Complex Networks IX* (2018), 271–282.
- [20] Aaron Halfaker. 2017. Interpolating quality dynamics in Wikipedia and demonstrating the Keilana effect. In *Proc. International Symposium on Open Collaboration (OpenSym)*.

- [21] Aaron Halfaker and R. Stuart Geiger. 2019. ORES: Lowering Barriers with Participatory Machine Learning in Wikipedia. *Proc. Human-Computer Interaction (HCI)*.
- [22] Aaron Halfaker, Os Keyes, Daniel Kluver, Jacob Thebault-Spieker, Tien Nguyen, Kenneth Shores, Anuradha Uduwage, and Morten Warncke-Wang. 2015. User Session Identification Based on Strong Regularities in Inter-Activity Time. In *Proc. International World Wide Web Conference (WWW)*.
- [23] Denis Helic. 2012. Analyzing user click paths in a Wikipedia navigation game. In *Proc. International Convention MIPRO*.
- [24] Hostinger Tutorials. 2022. The most visited website in every country (that isn't a search engine). <https://www.hostinger.com/tutorials/the-most-visited-website-in-every-country>
- [25] Jeff Huang and Ryen W White. 2010. Parallel browsing behavior on the web. In *Proceedings of the 21st ACM conference on Hypertext and Hypermedia*. 13–18.
- [26] Luis-Daniel Ibáñez and Elena Simperl. 2022. A comparison of dataset search behaviour of internal versus search engine referred sessions. In *ACM SIGIR Conference on Human Information Interaction and Retrieval*. 158–168.
- [27] Daxin Jiang, Jian Pei, and Hang Li. 2013. Mining search and browse logs for web search: A survey. *ACM Transactions on Intelligent Systems and Technology (TIST)* 4, 4 (2013), 1–37.
- [28] Honey Jindal, Neetu Sardana, and Raghav Mehta. 2020. Efficient web navigation prediction using hybrid models based on multiple evidence combinations. *International Journal of Computers and Applications* 42, 7 (2020), 715–728.
- [29] Isaac Johnson, Florian Lemmerich, Diego Sáez-Trumper, Robert West, Markus Strohmaier, and Leila Zia. 2020. Global gender differences in Wikipedia readership. In *Proc. Conference on Web and Social Media (ICWSM)*.
- [30] Faten Khalil, Jiuyong Li, and Hua Wang. 2009. An integrated model for next page access prediction. *International Journal of Knowledge and Web Intelligence* 1, 1-2 (2009), 48–80.
- [31] Muneo Kitajima, Marilyn H Blackmon, and Peter G Polson. 2000. A comprehension-based model of web navigation and its application to web usability analysis. In *People and computers XIV—Usability or else!* Springer, 357–373.
- [32] Tobias Koopmann, Alexander Dallmann, Lena Hettinger, Thomas Niebler, and Andreas Hotho. 2019. On the Right Track! Analysing and Predicting Navigation Success in Wikipedia. In *Proc. Conference on Hypertext and Social Media (HT)*.
- [33] Sean Kross, Eszter Hargittai, and Elissa M Redmiles. 2021. Characterizing the Online Learning Landscape: What and How People Learn Online. *ACM Hum.-Comput. Interact.* 5, CSCW1 (Feb. 2021), 19.
- [34] Juhi Kulshrestha, Marcos Oliveira, Orkut Karacalik, Denis Bonnay, and Claudia Wagner. 2020. Web Routineness and Limits of Predictability: Investigating Demographic and Behavioral Differences Using Web Tracking Data. *arXiv preprint arXiv:2012.15112* (2020).
- [35] Daniel Lamprecht, Dimitar Dimitrov, Denis Helic, and Markus Strohmaier. 2016. Evaluating and Improving Navigability of Wikipedia: A Comparative Study of Eight Language Editions. In *Proc. International Symposium on Open Collaboration (OpenSym)*.
- [36] Daniel Lamprecht, Kristina Lerman, Denis Helic, and Markus Strohmaier. 2017. How the structure of Wikipedia articles influences user navigation. *New Review of Hypermedia and Multimedia* 23, 1 (2017), 29–50.
- [37] David Lazer, Eszter Hargittai, Deen Freelon, Sandra Gonzalez-Bailon, Kevin Munger, Katherine Ognyanova, and Jason Radford. 2021. Meaningful measures of human society in the twenty-first century. *Nature* 595, 7866 (2021), 189–196.
- [38] David Lazer, Ryan Kennedy, Gary King, and Alessandro Vespignani. 2014. The parable of Google Flu: traps in big data analysis. *science* 343, 6176 (2014), 1203–1205.
- [39] Janette Lehmann, Claudia Müller-Birn, David Laniado, Mounia Lalmas, and Andreas Kaltenbrunner. 2014. Reader Preferences and Behavior on Wikipedia. In *Proc. Conference on Hypertext and Social Media (HT)*.
- [40] Florian Lemmerich, Diego Sáez-Trumper, Robert West, and Leila Zia. 2019. Why the world reads Wikipedia: Beyond English speakers. In *Proc. International Conference on Web Search and Data Mining (WSDM)*.
- [41] David M Lydon-Staley, Dale Zhou, Ann Sizemore Blevins, Perry Zurn, and Danielle S Bassett. 2021. Hunters, busybodies and the knowledge network building associated with deprivation curiosity. *Nature Human Behaviour* 5, 3 (2021), 327–336.
- [42] Nizar R Mabroukeh and Christie I Ezeife. 2009. Semantic-rich markov models for web prefetching. In *Proc. International Conference on Data Mining Workshops (ICDMW)*. IEEE, 465–470.
- [43] Fritz Machlup. 1983. The study of information: Interdisciplinary messages. (1983).
- [44] Lauren A Maggio, Ryan M Steinberg, Tiziano Piccardi, and John M Willinsky. 2020. Meta-Research: Reader engagement with medical content on Wikipedia. *Elife* 9 (2020), e52426.
- [45] M Mangel, WH Satterthwaite, P Pirolli, B Suh, and Y Zhang. 2013. Invasion biology and the success of social collaboration networks, with application to Wikipedia. *Israel Journal of Ecology and Evolution* 59, 1 (2013), 17–26.
- [46] Connor McMahon, Isaac Johnson, and Brent Hecht. 2017. The Substantial Interdependence of Wikipedia and Google: A Case Study on the Relationship Between Peer Production Communities and Information Technologies. *Proc. Conference on Web and Social Media (ICWSM)*.

- [47] Blagoj Mitrevski, Tiziano Piccardi, and Robert West. 2020. WikiHist. html: English Wikipedia's Full Revision History in HTML Format. In *Proc. Conference on Web and Social Media (ICWSM)*.
- [48] Jack Muramatsu and Wanda Pratt. 2001. Transparent Queries: investigation users' mental models of search engines. In *Conference on Research & Development in Information Retrieval (SIGIR)*.
- [49] Meera Narvekar and Shaikh Sakina Banu. 2015. Predicting user's web navigation behavior using hybrid approach. *Procedia Computer Science* 45 (2015), 3–12.
- [50] Richard E Nisbett and Timothy D Wilson. 1977. Telling more than we can know: Verbal reports on mental processes. *Psychological review* 84, 3 (1977), 231.
- [51] Alexandra Olteanu, Carlos Castillo, Fernando Diaz, and Emre Kiciman. 2019. Social data: Biases, methodological pitfalls, and ethical boundaries. *Frontiers in Big Data* 2 (2019), 13.
- [52] Ashwin Paranjape, Robert West, Leila Zia, and Jure Leskovec. 2016. Improving Website Hyperlink Structure Using Server Logs. In *Proc. International Conference on Web Search and Data Mining (WSDM)*.
- [53] Tiziano Piccardi, Michele Catasta, Leila Zia, and Robert West. 2018. Structuring Wikipedia articles with section recommendations. In *Conference on Research & Development in Information Retrieval (SIGIR)*.
- [54] Tiziano Piccardi, Miriam Redi, Giovanni Colavizza, and Robert West. 2020. Quantifying engagement with citations on Wikipedia. In *Proc. International World Wide Web Conference (WWW)*.
- [55] Tiziano Piccardi, Miriam Redi, Giovanni Colavizza, and Robert West. 2021. On the Value of Wikipedia as a Gateway to the Web. In *Proc. International World Wide Web Conference (WWW)*.
- [56] Tiziano Piccardi and Robert West. 2021. Crosslingual Topic Modeling with WikiPDA. *Proc. International World Wide Web Conference (WWW)*.
- [57] Peter Pirolli and Stuart Card. 1999. Information foraging. *Psychological review* 106, 4 (1999), 643.
- [58] Peter L T Pirolli and James E Pitkow. 1999. Distributions of surfers' paths through the World Wide Web: Empirical characterizations. *World Wide Web* 2, 1 (1999), 29–45.
- [59] Yan Qu and George W Furnas. 2008. Model-driven formative evaluation of exploratory search: A study under a sensemaking framework. *Information Processing & Management* 44, 2 (2008), 534–555.
- [60] Miriam Redi, Martin Gerlach, Isaac Johnson, Jonathan Morgan, and Leila Zia. 2020. A Taxonomy of Knowledge Gaps for Wikimedia Projects (Second Draft). (Aug. 2020). [arXiv:2008.12314 \[cs.CY\]](https://arxiv.org/abs/2008.12314)
- [61] Giovanna Chiara Rodi, Vittorio Loreto, and Francesca Tria. 2017. Search strategies of Wikipedia readers. *PLOS ONE* 12, 2 (02 2017), 1–15.
- [62] Dana Rotman, Sarah Vieweg, Sarita Yardi, Ed Chi, Jenny Preece, Ben Shneiderman, Peter Pirolli, and Tom Glaisyer. 2011. From slacktivism to activism: participatory culture in the age of social media. In *CHI'11 Extended Abstracts on Human Factors in Computing Systems*.
- [63] Matthew J Salganik. 2019. *Bit by bit: Social research in the digital age*. Princeton University Press.
- [64] Aju Thalappillil Scaria, Rose Marie Philip, Robert West, and Jure Leskovec. 2014. The Last Click: Why Users Give up Information Network Navigation. In *Proc. International Conference on Web Search and Data Mining (WSDM)*.
- [65] Aaron Shaw and Eszter Hargittai. 2018. The Pipeline of Online Participation Inequalities: The Case of Wikipedia Editing. *The Journal of communication* 68, 1 (Feb. 2018), 143–168.
- [66] Philipp Singer, Florian Lemmerich, Robert West, Leila Zia, Ellery Wulczyn, Markus Strohmaier, and Jure Leskovec. 2017. Why we read Wikipedia. In *Proc. International World Wide Web Conference (WWW)*.
- [67] Philipp Singer, Thomas Niebler, Markus Strohmaier, and Andreas Hotho. 2013. Computing Semantic Relatedness from Human Navigational Paths: A Case Study on Wikipedia. *Int. J. Semant. Web Inf. Syst.* 9, 4 (Oct. 2013), 41–70.
- [68] Adish Singla, Ryen White, and Jeff Huang. 2010. Studying trailfinding algorithms for enhanced web search. In *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*. 443–450.
- [69] Bongwon Suh, Lichan Hong, Peter Pirolli, and Ed H Chi. 2010. Want to be retweeted? large scale analytics on factors impacting retweet in twitter network. In *2010 IEEE second international conference on social computing*. IEEE, 177–184.
- [70] Linda Tauscher and Saul Greenberg. 1997. Revisitation Patterns in World Wide Web Navigation. In *Proc. Conference on Human Factors in Computing Systems (CHI)*.
- [71] Nathan TeBlunthuis, Tilman Bayer, and Olga Vasileva. 2019. Dwelling on Wikipedia: Investigating Time Spent by Global Encyclopedia Readers. In *Proc. International Symposium on Open Collaboration (OpenSym)*.
- [72] Michele Tizzoni, André Panisson, Daniela Paolotti, and Ciro Cattuto. 2020. The impact of news exposure on collective attention in the United States during the 2016 Zika epidemic. *PLoS computational biology* 16, 3 (March 2020), e1007633.
- [73] Nicholas Vincent and Brent Hecht. 2021. A Deeper Investigation of the Importance of Wikipedia Links to Search Engine Results. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW1 (April 2021), 1–15.
- [74] Claudia Wagner, Markus Strohmaier, Alexandra Olteanu, Emre Kiciman, Noshir Contractor, and Tina Eliassi-Rad. 2021. Measuring algorithmically infused societies. *Nature* 595, 7866 (2021), 197–204.

- [75] Shoujin Wang, Longbing Cao, Yan Wang, Quan Z Sheng, Mehmet A Orgun, and Defu Lian. 2021. A Survey on Session-based Recommender Systems. *ACM Comput. Surv.* 54, 7 (July 2021), 1–38.
- [76] Robert West and Jure Leskovec. 2012. Automatic Versus Human Navigation in Information Networks. *Proc. Conference on Web and Social Media (ICWSM)* (2012).
- [77] Robert West and Jure Leskovec. 2012. Human Wayfinding in Information Networks. In *Proc. International World Wide Web Conference (WWW)*.
- [78] Robert West, Ashwin Paranjape, and Jure Leskovec. 2015. Mining Missing Hyperlinks from Human Navigation Traces: A Case Study of Wikipedia. *Proc. International World Wide Web Conference (WWW)*.
- [79] Robert West, Joelle Pineau, and Doina Precup. 2009. Wikispeedia: An Online Game for Inferring Semantic Distances between Concepts. In *Proc. International Joint Conference on Artificial Intelligence (IJCAI)*.
- [80] Ryen W White, Mikhail Bilenko, and Silviu Cucerzan. 2007. Studying the use of popular destinations to enhance web search interaction. In *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*. 159–166.
- [81] Ryen W White and Steven M Drucker. 2007. Investigating behavioral variability in web search. In *Proc. International World Wide Web Conference (WWW)*. 21–30.
- [82] Ryen W White and Jeff Huang. 2010. Assessing the scenic route: measuring the value of search trails in web logs. In *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*. 587–594.
- [83] Tom D Wilson. 1981. On user studies and information needs. *Journal of documentation* (1981).
- [84] Tom D Wilson. 1997. Information behaviour: an interdisciplinary perspective. *Information processing & management* 33, 4 (1997), 551–572.
- [85] Tom D Wilson. 1999. Models in information behaviour research. *Journal of documentation* (1999).
- [86] Ellery Wulczyn and Dario Taraborelli. 2015. Wikipedia clickstream. https://meta.wikimedia.org/wiki/Research:Wikipedia_clickstream.
- [87] Paula Younger. 2010. Internet-based information-seeking behaviour amongst doctors and nurses: a short review of the literature. *Health Information & Libraries Journal* 27, 1 (2010), 2–10.
- [88] Kai Zhu, Dylan Walker, and Lev Muchnik. 2020. Content Growth and Attention Contagion in Information Networks: Addressing Information Poverty on Wikipedia. *Information Systems Research* 31, 2 (June 2020), 491–509.