

# Wikipedia Community Health Dashboards

Andrea Denina<sup>1</sup>, Paolo Aliprandi<sup>1</sup>, Marc Miquel-Ribé<sup>2</sup>, David Laniado<sup>3</sup>,  
Cristian Consonni<sup>4,\*</sup>

<sup>1</sup>DISI, University of Trento, Trento, Italy

<sup>2</sup>Universitat Pompeu Fabra, Mataró, Catalunya, Spain

<sup>3</sup>Big Data & Data Science Unit, Eurecat - Centre Tecnològic de Catalunya, Barcelona, Catalunya, Spain

<sup>4</sup>European Commission, Joint Research Centre (JRC), Ispra, Italy

andrea.denina@studenti.unitn.it, paoloalipro@gmail.com, mmiquelr@tecnocampus.cat, david.laniado@eurecat.org,  
cristian.consonni@acm.org

## Abstract

Wikipedia relies on the contributions of hundreds of thousands of volunteers across hundreds of language communities. The long-term health of Wikipedia’s editor base is essential for the platform’s stability and growth, and assessing it requires metrics that are reproducible, comparable over time, and regularly updated. Prior work introduced the *Community Vital Signs*, a set of six language-independent indicators capturing editor retention, stability, generational balance, role renewal, and global participation. In this paper, we present a live demonstration of the *Wikipedia Community Health Dashboards*, a set of interactive visualizations for exploring these indicators across all Wikipedia language editions. The dashboards are powered by a fully automated, open-source pipeline that updates the metrics monthly. We describe the system architecture and showcase how the dashboards enable researchers and community members to monitor and compare community health over time.

**Demo** — <https://vitalsigns.wmcloud.org/>

**Code** — <https://vitalsigns.wmcloud.org/code>

**Data** — <https://vitalsigns.wmcloud.org/data>

## 1 Introduction

Wikipedia exists thanks to a vast, multilingual ecosystem of volunteer communities whose work underpins one of the most widely used knowledge resources in the world. As these communities evolve, understanding their levels of participation, renewal, and organizational capacity becomes essential for the well-being of the community. Community health is a central theme in research, policy discussions, and community governance.

Despite the importance of these questions, systematic and comparable measures of community health across Wikipedia’s hundreds of language editions have remained limited and difficult to obtain. Existing tools typically focus on narrow aspects of editing activity—most notably re-

tion (Morgan and Halfaker 2018) or initiatives to stimulate editing and content contributions (Warncke-Wang et al. 2023)—and rarely provide longitudinal, cross-lingual, editor-level indicators that capture how communities evolve over time. Moreover, many analyses require substantial data-processing expertise, which limits their accessibility for Wikimedia community members and organizers, Wikimedia affiliates, and researchers interested in assessing the impact of community initiatives. This also makes it more difficult for community organizers to allocate resources for community health initiatives as there are no instruments to evaluate their long-term impact.

To address this gap, we introduce the *Wikipedia Community Health Dashboards*, a set of interactive visualizations that provide up-to-date, language-independent indicators of the health and sustainability of editor communities across all Wikipedia editions. The metrics presented are based on the *Community Vital Signs* (Miquel-Ribé, Consonni, and Laniado 2022), a framework of six indicators capturing fundamental aspects of community functioning: editor retention, population stability, balance between newcomers and experienced editors, renewal of administrators and special-function roles, and global participation across Wikimedia projects. All indicators are computed monthly from the MediaWiki History Dumps, enabling consistent, reproducible, and longitudinal comparisons across language editions. In this demo, we present the dashboards and the underlying pipeline that automates the computation of the *Community Vital Signs* for more than 358 Wikipedia language editions and we release the full implementation as open-source software. The dashboards allow users to interactively explore and compare community health trends over time and across language editions, supporting reflection and discussion around community sustainability. We demonstrate the system through use cases drawn from multiple Wikipedia communities that we also showcased in community-facing events, where it has been employed to support discussion, planning, and evaluation of community initiatives.<sup>1</sup>

The paper is organized as follows: in Section 2 we briefly

\**Disclaimer:* The view expressed in this paper is purely that of the authors and may not, under any circumstances, be regarded as an official position of the European Commission. Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

<sup>1</sup>Community Health Metrics project page on Meta-Wiki, [https://meta.wikimedia.org/wiki/Community\\_Health\\_Metrics](https://meta.wikimedia.org/wiki/Community_Health_Metrics), accessed on 2026-04-21.

summarize the *Community Vital Signs* framework which provides the metrics visualized in the dashboards; then, in Section 3 we describe the system architecture and data-processing pipeline. In Section 4, we illustrate the key functionalities of the dashboards and we focus on two selected use cases. Finally, Section 5 concludes the paper.

## 2 Overview of Community Vital Signs

The *Community Vital Signs* (Miquel-Ribé, Consonni, and Laniado 2022) are six language-independent metrics designed to measure several aspects of the health and sustainability of Wikipedia’s editor communities. Inspired by the medical analogy, where vital signs describe the basic functions that sustain life, these indicators assess the vital functions of a community: its capacity to grow, renew itself, and remain open and collaborative over time.

Three indicators describe the general population of active editors who create content: (i) *retention* refers to the ability to keep newcomers active after their first contributions; (ii) *stability* measures the consistency of participation within the active editor base over time; and (iii) *balance* captures the equilibrium between new and experienced editors, ensuring both renewal and continuity. The other three indicators focus on specific community roles and on the degree of coordination among communities: (iv) *administrators* measures the renewal and distribution of users with administrative rights; (v) *specialists* represent editors who perform technical or organizational roles, such as template editors or coordinators; and (vi) *global participation* reflects the engagement of local communities within broader Wikimedia spaces and projects. In addition, the dashboard displays general statistics about the (vii) *activity* of the community, i. e., the number of editors that have contributed in a given period. Appendix A reports the definitions of each *Vital Sign* indicator.

**Data sources.** The computation of the *Community Vital Signs* relies on the *MediaWiki History Dumps* (MWHHD),<sup>2</sup> a public dataset that records the full history of user, page, and revision events across all Wikimedia projects since 2001. Data dumps are updated monthly, partitioned by project (English Wikipedia (enwiki), French Wikipedia (frwiki), etc.). Additional information are available in Appendix B.

## 3 System Architecture

The data-processing pipeline—illustrated in Figure 1—integrates data extraction, processing, and visualization into a single workflow, which automates the entire computation process, from retrieving the latest *MediaWiki History Dumps*, to extracting and aggregating editor-level data, and finally computing the six *Vital Signs* indicators. In addition, we implemented a monitoring component to observe in real time the status of the computation.

**System Components.** The pipeline is built using Apache Airflow, an open-source workflow orchestration platform. In Airflow, a pipeline is defined as a directed acyclic graph

<sup>2</sup>The documentation for MWHHD is available at <https://w.wiki/HG7a>, while the data can be downloaded at <https://w.wiki/7kTP>, accessed on 2026-04-21.

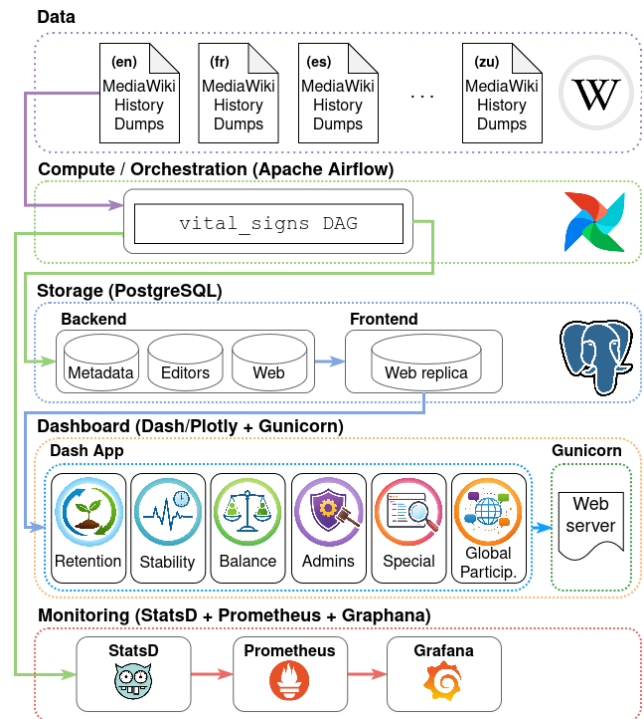


Figure 1: Architecture of the *Vital Signs* pipeline.

(DAG) composed of various tasks. Each task within the DAG is a modular step in the workflow, which is coordinated by Airflow’s scheduler according to its dependencies and periodicity. All components of the pipeline are containerized with Docker and orchestrated via Docker Compose. Each service in the pipeline runs within isolated but interconnected containers to ensure reproducibility and simplify both their development and deployment. For storage, we use PostgreSQL databases both as the metadata backend for Airflow and to store the results of metric computations. For monitoring, Airflow exports runtime metrics through StatsD, which forwards them to a StatsD-Exporter instance that translates them into the Prometheus export format. Prometheus periodically scrapes these metrics, storing them as time series and making them available for visualization and alerting. Grafana acts as the frontend monitoring interface, providing interactive dashboards that display task performance, execution times, and resource utilization in real time. Finally, the Dash web application provides the presentation layer, served through the Gunicorn WSGI server, allowing users to explore and compare the computed *Vital Signs* across different language communities.

**Deployment.** The full pipeline is deployed on a Wikimedia Cloud VPS instance with 32 GB of RAM, 16 cores and 512 GB of storage (20 GB for system storage and 492 GB for data storage). The Wikimedia Cloud infrastructure provides access to the MediaWiki History Dumps for all wikis (316 GB in TSV format, BZ2 compressed).

## 4 Results and Discussion

**Dashboards.** The dashboards are available at <https://vitalsigns.wmcloud.org>. The source code is available at <https://vitalsigns.wmcloud.org/code> (doi:10.5281/zenodo.18248819). The code is released under the MIT license.<sup>3</sup> All data can be downloaded from <https://vitalsigns.wmcloud.org/data>, which also contains the datasheet for the dataset. Data, charts, and other content is released under the Creative Commons CC0 dedication (public domain).<sup>4</sup>

**Controls.** Figure 2 shows the controls of the dashboards: users can select the metric, language editions, aggregation level, and visualization parameters. Results are displayed in the main view area. It is possible to restrict the time range of the visualization by clicking on the plot.

**Use case 1: *Balance* for Italian Wikipedia.**<sup>5</sup> Figure 3 (a) shows the *Balance* vital sign for Italian Wikipedia, i. e. the proportion of active editors over time by cohort of first edit. Editors are grouped by *lustrum* (5-year period) of their first recorded edit. The visualization highlights a clear generational shift in the composition of active editors. Early cohorts progressively decrease in relative weight, while newer cohorts increasingly contribute to the active population, reflecting a continuous renewal process. At the same time, a few older editor cohorts remain present throughout the timeline, indicating a stable core of long-term contributors, indicating a mature community that is capable of renewing.

**Use case 2: *Admin* for German Wikipedia.**<sup>6</sup> Figure 3 (b) shows the *Admin* vital sign for German Wikipedia, reporting the temporal distribution of sysop flags granted by editor cohort. The figure displays annual counts of newly appointed administrators (left) and cumulative total (right), stratified by the lustrum of editors' first edit. Admin appointments are concentrated in the early years of the project, with the vast majority of sysop flags (142 over 168, or 84.5%) granted to editors who joined German Wikipedia between 2001 and

<sup>3</sup><https://github.com/WikiCommunityHealth/vital-signs-pipeline/blob/main/LICENSE.md>

<sup>4</sup><https://creativecommons.org/publicdomain/zero/1.0/>

<sup>5</sup><https://flink.rtrace.io/koTS>

<sup>6</sup><https://flink.rtrace.io/Vw7b>

### Select the parameters

Vital Sign

Language

Editors  
 Active  Very Active

Time aggregation  
 Yearly  Monthly

Flag

Retention rate

Percentage or number (y-axis)  
 Percentage  Number

Figure 2: Dashboard parameter selection interface. Users can select a vital sign to display, Wikipedia language editions, editor group, time aggregation, and scale of the y-axis.

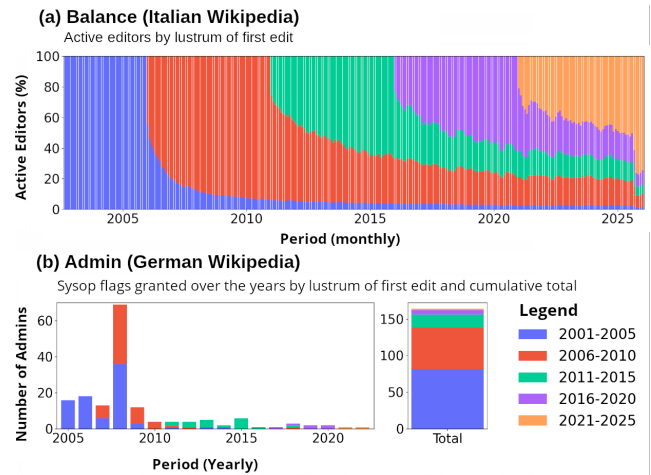


Figure 3: (a) The *Balance* vital sign for Italian Wikipedia. (b) The *Admin* vital sign for German Wikipedia. The legend refers to both plots.

2010. Afterwards, the number of new administrators declines sharply. As a result, the current administrator pool is dominated by early generations of editors with a slow renewal. This pattern suggests a contraction in turnover that may raise concerns about the long-term sustainability and governance capacity of the German Wikipedia community.

## 5 Conclusions

In this demo, we presented the *Community Health Dashboards*, a set of interactive visualizations that provide up-to-date, language-independent indicators of the health and sustainability of editor communities across all Wikipedia editions. We described the underlying data processing pipeline, for which we have provided the source code, that automates the computation of the *Community Vital Signs* metrics for more than 128 Wikipedia language editions. While many tools developed in academic contexts are not maintained beyond their initial publication, we have presented a self-updating system deployed on the Wikimedia Cloud community infrastructure. By design, the tool is intended to integrate into the broader Wikipedia ecosystem, with the potential to provide sustained value to both researchers and Wikimedia communities over time.

## References

- Miquel-Ribé, M.; Consonni, C.; and Laniado, D. 2022. Community Vital Signs: Measuring Wikipedia Communities' Sustainable Growth and Renewal. *Sustainability*, 14(8).
- Morgan, J. T.; and Halfaker, A. 2018. Evaluating the impact of the Wikipedia Teahouse on newcomer socialization and retention. In *Proc. of OpenSym*, OpenSym '18. New York, NY, USA: ACM. ISBN 9781450359368.
- Warncke-Wang, M.; Ho, R.; Miller, M.; and Johnson, I. 2023. Increasing Participation in Peer Production Communities with the Newcomer Homepage. *Proc. ACM Hum.-Comput. Interact.*, 7(CSCW2).

## Checklist

### 1. For most authors...

- (a) Would answering this research question advance science without violating social contracts, such as violating privacy norms, perpetuating unfair profiling, exacerbating the socio-economic divide, or implying disrespect to societies or cultures? **Yes**
- (b) Do your main claims in the abstract and introduction accurately reflect the paper's contributions and scope? **Yes**
- (c) Do you clarify how the proposed methodological approach is appropriate for the claims made? **Yes**
- (d) Do you clarify what are possible artifacts in the data used, given population-specific distributions? *N/A*, we do not sample any data. We compute our metrics using the full data available in the MediaWiki History Dump dataset.
- (e) Did you describe the limitations of your work? *N/A*, the dashboards display the *Community Vital Signs* which are indicators of community health rather than explanations. For a discussion of the limitations of the *Community Vital Signs* themselves refer to the work that introduced them (Miquel-Rib'e, Consonni, and Lanido 2022).
- (f) Did you discuss any potential negative societal impacts of your work? **No**, there is no specific discussion of potential negative societal impacts. However, such risks are low because data is aggregated and public. Furthermore, this work has been developed over more than five years with ongoing feedback from various Wikipedia communities. See the list of presentations given about the metrics on the "Community Health Metrics"<sup>1</sup> page on Meta-Wiki.
- (g) Did you discuss any potential misuse of your work? **No**, misuse is not discussed explicitly. However, metrics discussions are a common part of community debates on Wikipedia. Data are always presented aggregated to avoid potential misuses targeted at single Wikipedia contributors.
- (h) Did you describe steps taken to prevent or mitigate potential negative outcomes of the research, such as data and model documentation, data anonymization, responsible release, access control, and the reproducibility of findings? **Yes**, we processed data about users participating in Wikipedia, but all the raw data are public and we only present data in aggregate form. We adhere to the ethical guidelines about research on Wikipedia and privacy as outlined by previous work.<sup>2</sup>
- (i) Have you read the ethics review guidelines and ensured that your paper conforms to them? **Yes**

### 2. Additionally, if you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- (a) If your work uses existing assets, did you cite the creators? **Yes**
  - (b) Did you mention the license of the assets? **Yes**
  - (c) Did you include any new assets in the supplemental material or as a URL? **Yes**
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? **Yes**
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **Yes**. As mentioned in the paper, our data source is the MediaWiki History Dumps and the dashboards only present aggregated indicators at community level. Regarding offensive content, the work does not surface article text or edit content; it only processes metadata about editing activity, such as counts of editors, retention, cohorts, and roles. Therefore, potentially offensive content that may exist in Wikipedia articles is neither analyzed nor shown.
  - (f) If you are curating or releasing new datasets, did you discuss how you intend to make your datasets FAIR? **Yes**. Our project is developed to be part of the Wikimedia ecosystem, which has the principles of openness and transparency among its core inspiring values, and as such it strictly adheres to these principles. Although we had not space to explicitly frame the data release along the FAIR principles in the paper, we are following them:
    - Findability and Accessibility: all the code and data generated in the project are released under an open source, permissive license, and made available on public repositories. We archived the code on Zenodo to obtain a permanent DOI for the code release to improve findability.
    - Interoperability and Reusability: the dataset is released as a SQLite database. SQLite is a free and open source software and its file format is widely used due to its stability, cross-platform nature, and widespread adoption, even being recommended by the Library of Congress (LoC) for long-term data preservation.<sup>3</sup>
  - (g) If you are curating or releasing new datasets, did you create a Datasheet for the Dataset? **Yes**, the datasheet is available at <https://vitalsigns.wmcloud.org/data>.
- ### 3. Additionally, if you used crowdsourcing or conducted research with human subjects
- (a) Did you include the full text of instructions given to participants and screenshots? **NA**, the study does not involve human participants or crowdsourcing.
  - (b) Did you describe any potential participant risks, with mentions of Institutional Review Board (IRB) approvals? **NA**, the study does not involve human participants.

<sup>1</sup>[https://meta.wikimedia.org/wiki/Community\\_Health\\_Metrics](https://meta.wikimedia.org/wiki/Community_Health_Metrics), accessed on 2026-04-09.

<sup>2</sup><https://osf.io/preprints/osf/uyxnf.v2>, accessed on 2026-04-09.

<sup>3</sup><https://www.loc.gov/preservation/digital/formats/fdd/fdd000461.shtml>, accessed on 2026-04-09.

- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? NA, no participants were recruited.
- (d) Did you discuss how data is stored, shared, and de-identified? Yes. Data storage, aggregation, and sharing are described through the system architecture and the use of public MediaWiki History Dumps, with only aggregated outputs shared.

### **Generative AI Usage Statement**

We confirm that all text in this paper was written by the authors. In preparing this paper, we utilized ChatGPT 5.2 (Pro version) for grammar and spelling checks and to improve the clarity of the author-written text. No Generative AI tools were used to generate content and citations from scratch. The content and intellectual contributions remain entirely those of the human authors. We acknowledge the contributions of the mentioned tool in enhancing the writing process while maintaining full academic integrity and abiding by the AAAI policy.

## A Definition of *Vital Signs*

Table 2 reproduces a summary of the *Community Vital Signs* metrics originally introduced by Miquel-Ribé, Consonni, and Laniado (2022). The table provides an overview of the seven metrics presented in the dashboards used to assess Wikipedia community health—Retention, Stability, Balance, Special Functions, Administrators, and Global Participation—together with the corresponding indicators and their operational definitions. Each row specifies how a given indicator is computed (e.g., retention rate, number of very active editors by generation, administrators by lustrum), offering a concise reference framework for interpreting the metrics used throughout the analysis.

## B MediaWiki History Dump

The MediaWiki History Dump (MWHD) dataset is distributed as BZ2 compressed TSV files and forms the foundation for extracting and aggregating editor-level information used to compute the *Vital Signs* indicators. Table 1 reports some statistics about the MediaWiki History Dump for the 2025–11 release<sup>7</sup> files for the top-10 Wikipedia languages by number of users as of November 2025.<sup>8</sup>

| Project                          | Language   | Size (GB) | No. files |
|----------------------------------|------------|-----------|-----------|
| arwiki                           | Arabic     | 5.5       | 24        |
| dewiki                           | German     | 24.2      | 25        |
| enwiki                           | English    | 132.2     | 300       |
| eswiki                           | Spanish    | 14.8      | 25        |
| frwiki                           | French     | 20.6      | 25        |
| itwiki                           | Italian    | 12.3      | 25        |
| jawiki                           | Japan      | 9.5       | 24        |
| ptwiki                           | Portuguese | 6.2       | 25        |
| ruwiki                           | Russian    | 13.3      | 24        |
| zhwiki                           | Chinese    | 7.6       | 24        |
| <b>Total (top-10 languages)</b>  |            | 246.2     | 557       |
| <b>Total (all 358 languages)</b> |            | 383.3     | 1394      |

Table 1: Size and of 2025–11 release of the MediaWiki History Dump files for 10 Wikipedia language editions and totals for all available language. For each edition a file per year is available, with the exception that for arwiki, the file for 2001 is missing; for jawiki, ruwiki, and zhwiki the file for 2002 is missing; and for enwiki there are monthly files.

<sup>7</sup>[https://dumps.wikimedia.org/other/mediawiki\\_history/2025-11/](https://dumps.wikimedia.org/other/mediawiki_history/2025-11/), accessed 2026-04-21.

<sup>8</sup>[https://meta.wikimedia.org/wiki/List\\_of\\_Wikipedias](https://meta.wikimedia.org/wiki/List_of_Wikipedias), accessed 2026-04-21.

| <b>Metric</b>        | <b>Indicator</b>        | <b>Indicator Definition</b>  |
|----------------------|-------------------------|--|
| Retention            | Retention rate          | Percentage of new editors who edit at least once 60 days after their first edit.   |
| Stability            | Stability               | Number of active editors by the number of months they have been active in a row.   |
| Balance              | Balance                 | Number and percentage of very active editors by year and by generation (lustrum of the first edit).  |
| Special functions    | Technical editors       | Number of very active editors in technical namespaces (i.e., editors who performed more than 100 edits in one month in namespaces Mediawiki and Templates), broken down by year and by generation.       |
|                      | Coordinators            | Number of very active editors in coordination namespaces (i.e., editors who performed more than 100 edits in one month in namespaces Wikipedia and Help), broken down by year and by generation.         |
| Admins               | Admins by year          | Number of admins by year of flag granted and by generation.  |
|                      | Admins by lustrum       | Total number of active admins by generation at the current month.  |
|                      | Admins by lustrum       | Total number of active admins by generation at the current month.  |
| Global participation | Meta-wiki participation | Ratio between the number of active editors in Meta-wiki that have as primary a given language edition, divided by the number of active editors in that Wikipedia language edition during the same month. |
|                      | Primary language        | Distribution of the primary language edition of the editors contributing to a given language edition.  |
| Activity             | Active users            | The number of editors with 5 or more edits in a month  |

Table 2: *Vital Signs* definition