

# A Taxonomy of Knowledge Gaps for Wikimedia Projects (Summary and Motivation)

If you are looking for a shorter (than the [full paper](#)) write-up to have an overview of our thoughts on how to approach measuring knowledge gaps, please read on.

## What is the problem?

Wikimedia projects aim to empower people from around the world to share in the sum of all knowledge. Knowledge gaps are a barrier to this goal. Today we do not have a shared definition and understanding of knowledge gaps on the Wikimedia projects. As a result, we do not have a way to measure the extent of knowledge gaps in the Wikimedia projects. We do not know how to measure certain types of gaps. We do not have established methods to understand the relationship between the different types of gaps.

## Why is it important?

Our decisions can suffer from the risk of being sub-optimal as we do not see the full spectrum of knowledge gaps. As a result, we tend to invest in more widely discussed or understood types of knowledge gaps. We also have very limited ways to measure progress towards addressing specific gaps across many years. This makes learning and course-correction harder.

## What is our solution?

We want to develop a Knowledge Gap Index that is a composite index which will measure the state of knowledge gaps across many Wikimedia projects.<sup>1</sup> The first step to do that is to build a taxonomy of knowledge gaps which gives us a shared framework for understanding the different types of knowledge gaps. Such a taxonomy also specifies what each gap type is (and is not). So far, we have built a taxonomy of knowledge gaps for the Wikimedia projects and we are seeking your input to improve it.

## What do we mean by knowledge gaps?

Before this research, the phrase knowledge gaps equated to content gaps for many. We want to change this by expanding the definition of knowledge gaps to include gaps in contributorship and readership on top of content. For those of you interested in a formal definition of knowledge gaps, here is what we propose:

*Knowledge gaps* are major differences in participation or coverage of a specific group of readers, contributors, or content.

## Why has developing the framework for understanding knowledge gaps been challenging?

Developing a framework for understanding and measuring knowledge gaps is hard because:

- We need to make judgement calls about what gaps to include and which ones to not include. We also need to decide how specific the gaps will need to be.
- The measurement of certain gap types can be complex due to technical reasons (how should we gather the data for it?) or privacy reasons. They can also be complex research topics that we know very little about from the methodological perspective.

---

<sup>1</sup>If this is the first time you are hearing about composite indices, don't be discouraged by the fancy name of it. A composite index captures a variety of indices that aim to measure specific dimensions of a problem and their relationship to the overall problem. For example, the [gender equality index](#) is a composite index that measures the state of gender equality considering work, money, power and other types of indices.

---

# 1 READERS

The readership dimension of knowledge gaps encompasses all those gaps related to readers’ access to Wikimedia sites. We define readers as all users who connect directly to the projects to access Wikimedia content excluding how content consumption happens outside of Wikimedia, e.g., voice assistants, search engines, or third-party apps.

<i>Facet</i>	<i>Gap</i>	<i>Description</i>
<p><b>Sociodemographics</b></p> <p><i>Objective:</i> readers with different social status, demographics, and cultural background can easily and safely accessing free knowledge</p>	<p><i>Gender</i></p> <p><i>Age</i></p> <p><i>Locale</i></p> <p><i>Language</i></p> <p><i>Income</i></p> <p><i>Education</i></p> <p><i>Background</i></p>	<p>Difference between readers of different gender identities in how and how much they access the sites.</p> <p>Difference between readers of different age in how and how much they access the sites.</p> <p>Differences in readership between rural areas, towns, and cities</p> <p>Differences in readership depending on readers’ ability to read one or more languages</p> <p>Difference on how readers with different income, wealth, or employment status access Wikimedia sites</p> <p>Differences in readership depending on readers’ educational background</p> <p>Differences in readership among people with different cultural, political and sexual preferences</p>
<p><b>Information Need</b></p> <p><i>Objective:</i> readers with different information needs can find and consume free knowledge</p>	<p><i>Motivation</i></p> <p><i>Information Depth</i></p> <p><i>Familiarity</i></p>	<p>Differences in readership depending on the reason behind readers’ visit to the site</p> <p>Differences in readership depending on the depth of information for which a reader is looking</p> <p>Differences in readership depending on one’s prior familiarity with a topic</p>
<p><b>Accessibility</b></p> <p><i>Objective:</i> readers with different technical setup and skills can easily access Wikimedia projects</p>	<p><i>Internet connectivity</i></p> <p><i>Device</i></p> <p><i>Tech Skills</i></p> <p><i>Disabilities</i></p>	<p>Contrasts among the ability of readers with different internet connections to access Wikimedia sites</p> <p>Difference in accessibility to the site depending on readers’ devices</p> <p>Differences in readership depending on readers’ general internet skill</p> <p>Disparities in ability to access the knowledge within Wikipedia depending on individual disabilities</p>

## 2 CONTRIBUTORS

The contributor dimension of knowledge gaps covers all gaps related to categories of people contributing to Wikimedia sites. We define contributors as all individuals who *edit* or otherwise *maintain* Wikimedia content. For the purpose of this taxonomy, this definition does not include technical contributors—i.e. the individuals who build the MediaWiki software on which Wikimedia sites run—though the software and choices made in its design certainly are highly impactful on what types of contributors feel supported and what content is created.

<i>Facet</i>	<i>Gap</i>	<i>Description</i>
<p><b>Sociodemographics</b></p> <p><i>Objective:</i> contributors with different social status, demographics, and cultural background can easily and safely access and contribute to free knowledge</p>	<p><i>Gender</i></p> <p><i>Age</i></p> <p><i>Locale</i></p> <p><i>Language</i></p> <p><i>Income</i></p> <p><i>Education</i></p>	<p>Differences between contributors of different gender identities in how and how much they contribute to the sites.</p> <p>Differences between contributors of different ages in how and how much they contribute to the sites.</p> <p>Differences between contributors of different locales (urban, rural) in how and how much they contribute to the sites.</p> <p>Differences between contributors of different reading abilities in a language in how and how much they contribute to the sites.</p> <p>Differences between contributors with different income, wealth, or employment status in how and how much they contribute to the sites.</p> <p>Differences between contributors of different educational backgrounds in how and how much they contribute to the sites.</p>
<p><b>Contextual</b></p> <p><i>Objective:</i> contributors with different motivations and roles can access and contribute to free knowledge</p>	<p><i>Motivation</i></p> <p><i>Role</i></p>	<p>Differences in contribution depending on one’s reason for contributing to the site.</p> <p>Differences in contribution depending on the type of editing that one chooses to do.</p>
<p><b>Accessibility</b></p> <p><i>Objective:</i> contributors with different technical resources and abilities can easily access and contribute to Wikimedia projects</p>	<p><i>Internet connectivity</i></p> <p><i>Device</i></p> <p><i>Tech Skills</i></p> <p><i>Disabilities</i></p>	<p>Disparities in ability to contribute to the knowledge within Wikipedia depending on one’s access to high-speed internet</p> <p>Disparities in ability to contribute to the knowledge within Wikipedia depending on one’s device.</p> <p>Disparities in ability to contribute to the knowledge within Wikipedia depending on one’s general internet skills</p> <p>Disparities in ability to contribute to the knowledge within Wikipedia depending on individual disabilities</p>

### 3 CONTENT

Wikipedia is incomplete by design. The opportunity to share new information with the world is a major motivating factor among both new and established Wikipedia contributors. However, when important information about a topic is absent, incomplete, biased, or otherwise inaccessible to readers, these content gaps can undermine Wikipedia’s ability to serve the needs of its global audience.

We characterize *gaps in content coverage* as follows.

<i>Facet</i>	<i>Gap</i>	<i>Description</i>
<p><b>Policy</b></p> <p><i>Objective:</i> content is consistent with core content policies</p>	<p><i>Verifiability</i></p> <p><i>Neutrality</i></p>	<p>Differences in the use of reliable sources in order to verify content.</p> <p>Biases in the content across Wikipedia articles .</p>
<p><b>Accessibility</b></p> <p><i>Objective:</i> content is accessible to different audiences</p>	<p><i>Multimedia</i></p> <p><i>Structured Data</i></p> <p><i>Readability</i></p>	<p>Differences in coverage with respect to the type of media used to share the content</p> <p>Differences in the use of information which is indexed and machine-readable</p> <p>Barriers for accessing or consuming information originating from content</p>
<p><b>Diversity</b></p> <p><i>Objective:</i> content covers knowledge that is underrepresented, marginalized, and locally relevant</p>	<p><i>Gender</i></p> <p><i>Geography</i></p> <p><i>Impactful topics</i></p> <p><i>Cultural context topics</i></p>	<p>Differences in content coverage depending on the gender identity of subjects</p> <p>Differences in coverage of topics related to geographic regions or population distribution</p> <p>Differences in coverage of topics that are of common interest</p> <p>Differences in coverage of topics related to the history, heritage, and characteristics of a current or former cultural group</p>