

Wikipedia Culture Gap: Quantifying Content Imbalances Across 40 Language Editions

Marc Miquel-Ribé¹, David Laniado^{2*}

¹Universitat Pompeu Fabra – Department de Comunicació, Catalonia

²Eurecat – Centre Tecnològic de Catalunya

* Correspondence:

Dr. Marc Miquel-Ribé

marcmiquel@gmail.com

Keywords: content imbalance₁, cross-cultural studies₂, cultural diversity₃, online communities₄, wikipedias₅, digital humanities₆, data mining₇, big data₈.

Abstract

The online encyclopedia Wikipedia is the largest general information repository created through collaborative efforts from all over the globe. Despite the project's goal being to achieve the sum of human knowledge, there are strong content imbalances across the language editions. In order to quantify and investigate these imbalances, we study the impact of cultural context in 40 language editions. To this purpose, we developed a computational method to identify articles that can be related to the editors' cultural context associated to each Wikipedia language edition. We employed a combination of strategies taking into account geolocated articles, specific keywords and categories, as well as links between articles. We verified the method's quality with manual assessment and found an average precision of 0.92 and an average recall of 0.95. The results show that about a quarter of each Wikipedia language edition is dedicated to represent the corresponding cultural context. Although a considerable part of this content was created during the first years of the project, its creation is sustained over time. An analysis of cross-language coverage of this content shows that most of it is unique in its original language, and reveals special links between cultural contexts; at the same time, it highlights gaps where the encyclopaedia could extend its content. The approach and findings presented in this study can help to foster participation and inter-cultural enrichment of Wikipedias. The datasets produced are made available for further research.

1. Introduction

Wikipedia's most striking characteristic is the fact that it is a collaborative project: everybody can become a volunteer contributor and join the community. At present, there are 288 Wikipedia language editions, English being the largest with more than 5 million articles (and a total of 40 million articles counting all the languages). Wikipedia's goal is to provide the "sum of human knowledge", available to everyone for free, and at the moment it is already one of the most successful collaborative efforts in the Internet. Even though there is no central authority dictating the content to be created, the system is

Imbalances Across 40 Language Editions

35 based on the following content rules. Probably the most important rule is the 'Neutral Point of View'
 36 (NPOV), which roughly means "representing fairly, proportionately, and, as far as possible, without
 37 editorial bias, all of the significant views that have been published by reliable sources on a topic"¹.
 38 'Notability'², another core content rule, defines the criteria through which editors judge whether a
 39 specific topic deserves an article.

40 Although the above-mentioned norms exist in all Wikipedia language editions, their application and
 41 interpretation are constantly negotiated by the editors from each community. The fact that policies
 42 neither encourage or discourage languages' cultural differences and idiosyncrasies being reflected into
 43 content, results in a spontaneous creation of content. Moreover, each Wikipedia language edition is
 44 created in a decentralized way; as a result, editors themselves may not always be aware of the global
 45 product. In fact, each language edition has proven to be diverse in terms of both article content and
 46 absolute number of articles, up to the point that diversity has been often called "Systemic bias", which
 47 is referred to as "an imbalanced coverage of subjects and perspectives on the encyclopedia". This
 48 imbalance is often attributed to the lack of editors or resources in a particular language background.
 49 Among the reasons that explain why some languages do not have a Wikipedia language edition or have
 50 it underdeveloped, Van Dijk (2009) and Ensslin (2011) mention, among others, the reduced number of
 51 speakers, the digital divide, and the low online reputation of their language.

52 **Cultural contextualization.** In general, differences in the content of language editions are attributed
 53 by the current literature to contextual factors or to a process named by Hecht (2013, p. 47) as cultural
 54 contextualization, which "is the cause of some of the content diversity in multilingual Wikipedia".
 55 Cultural contextualization is also present in other user-generated projects such as OpenStreet Maps,
 56 Twitter or Flickr (Hecht 2013). The explanation of how it influences the final characteristics of content
 57 is rooted in the fields of Linguistics, and Cultural and Social Psychology. For instance, according to
 58 Clark (1996), the members of a cultural community usually share "facts, beliefs, procedures, norms,
 59 and assumptions". Hence, it is likely that the editors of each language community (and
 60 subcommunities, especially considering those languages with large geographical extension) may
 61 reflect in their articles the meanings they implicitly agree on, resulting in a great deal of diversity in
 62 such a worldwide project. Cultural contextualization occurs when there is a certain degree of freedom
 63 in content-based projects.

64 In Wikipedia, there is extensive literature on how cultural contextualization has shaped each language
 65 edition. Depending on whether the emphasis is put on the articles' text or on the Wikipedia's overall
 66 structure, effects can be classified into two main groups: Discourse and Structure.

67 *Discourse effects* are based on the idea that since each language edition constitutes a community (and
 68 perhaps few subcommunities), their editors tend to hold a shared cultural background and this
 69 ultimately limits the points of view adopted in the articles within one and the same language edition.
 70 (In the literature, the editor's point of view is referred to as: 'linguistic point of view', 'national point
 71 of view', or 'cultural bias'). In different language editions, the differences in the editors' point of view

¹ https://en.wikipedia.org/wiki/Wikipedia:Neutral_point_of_view

² <https://en.wikipedia.org/wiki/Wikipedia:Notability>

72 become more prominent, especially when it comes to controversial topics, where history and politics
73 are seen from opposite positions (Massa and Scrinzi 2011; Apic, Betts, and Russell 2011). For instance,
74 Rogers and Sendijarevic (2012) compared an article dedicated to ‘The Srebrenica Massacre’
75 throughout different Wikipedia language editions, including English and Balkan languages. The study
76 shows how the *same* article in different language editions adopts a different point of view to illustrate
77 facts; such points of view are sometimes unified, other times in total disagreement when it comes to
78 the terminology employed and its political connotations.

79 Likewise, in order to explore how contextualized Wikipedia language editions are, Bao et al. (2012)
80 developed a website which allows to explore similarities and differences in points of view of an article
81 whose concept exists across languages. Pentzold et al. (2017) showed that topics related to cultural
82 heritage such as ‘Bullfighting’ are framed differently in Catalan, Spanish and English language
83 editions, and have different focuses of controversies. Other studies point out that editors’ geographical
84 closeness to the subject of their articles impacts on the level of article exhaustiveness. Callahan and
85 Herring (2011) explored in the English and Polish Wikipedia how the biographical articles of well-
86 known people are more complete (in terms of features such as the number of pictures, education,
87 political ideology, controversies mentioned or family members names) in the language editions
88 associated to the territories where the person is from.

89 *Structural effects* are based on the idea that context and culture are relevant factors that affect editor
90 interests and consequently content coverage. Ronen et al. (2014) explored the relationships between
91 Wikipedia language editions by creating a network with all languages (global language network)
92 articles’ edits and assessed their centrality with eigenvector centrality. They found that English acts as
93 an influential central hub, followed by other well-spread languages such as French, Spanish, German,
94 among others. However, besides attributing it to visibility, they do not explain the factors which
95 influence each other. In this sense, Saimolenko et al. (2016), in order to explore to understand cultural
96 similarity understood as the significant interest of communities in contributing to articles about similar
97 topics, analyzed both edits in articles existing in various language editions and several cultural factors.
98 They found that cultural similarity is due to various factors affecting topic choices such as shared
99 language family, number of bilinguals, geographical proximity, among others.

100 In another study on common editing interests, Karimi et al. (2015) gathered all the editors’ edits from
101 English Wikipedia and analyzed their relationships in order to determine how close their affinities
102 were. Results showed that editors from close locations tend to have a higher coincidence in the articles
103 they edit than editors from distant geographical locations. The geographical factor was also used to
104 explain that Wikipedia language editions whose language-related territories are far from each other
105 tend to have less articles in common (i.e. their articles have no equivalence) than those whose territories
106 that are geographically close (Warncke-Wang et al. 2012).

107 Other studies show that editors tend to focus on their territories, either because geolocated articles are
108 edited by nearby editors or because they give them a higher visibility in the overall Wikipedia network
109 of articles. For instance, Hecht and Gergle (2010a) computed the location of each anonymous edit in
110 geolocated articles and discovered that many of the contributions were made from close distances.
111 Another effect detected by Hecht and Gergle (2009), called ‘Self-focus bias’, explains that the articles

Imbalances Across 40 Language Editions

112 located in the countries local to each language edition are linked to many more articles (i.e. they have
 113 more inlinks) than the articles located in the other countries.

114 All in all, this second group of effects shows that context has a key impact on Wikipedia *content*
 115 *coverage* and shows the relevance of geographical context to editors' activity.

116 However, in this research stream, one key perspective is missing. We argue that in order to estimate
 117 the impact of cultural context on content coverage, it would be necessary to know which articles relate
 118 to the cultural context of each language edition besides geolocated articles, including topics such as
 119 language, people, traditions, among others. This association would permit a more elucidated
 120 cartography on content coverage which would allow, first, to show whether the cultural context
 121 occupies a considerable part of each Wikipedia language edition, and second, to verify whether cultural
 122 context is at the base of the imbalances between Wikipedia language editions.

123 To do so, we advance the following three research questions about cultural context content:

- 124 • **RQ1.** *What is the extent of cultural context content in each Wikipedia language edition?*
- 125 • **RQ2.** *How have cultural context content articles been created over time?*
- 126 • **RQ3.** *What is their availability across different language editions?*

127

128 Therefore, in this work we aim to go one step further in the study of cultural contextualization, focusing
 129 on its structural effects and content coverage. We propose obtaining, for every Wikipedia language
 130 edition, a group of articles related to the editors' cultural context(s). In this way we are able to
 131 understand the relationship between the content imbalances and the representation of editors' cultural
 132 context in every Wikipedia language edition. We called culture gap the imbalances across language
 133 editions in content representing cultural context.

134 To the best of our knowledge, this is the first study that performs a perimetric analysis of the cultural
 135 context content. In particular, a valuable corpus is obtained to examine Wikipedia's cultural
 136 contextualization effects on content coverage more in depth than it has been done in previous studies.
 137 Moreover, the corpus also represents a valuable tool to understand the editors' culture and may be
 138 useful to both researchers and Wikipedia editors who want to increase cultural diversity.

139 In summary, our main contributions are the following:

- 140 • We provide a computational method to identify articles related to the cultural context of a given
 141 language community.
- 142 • We construct a dataset for 40 Wikipedia language editions comprising the articles representing
 143 their cultural contexts and make it publicly available for future research.
- 144 • We analyze the availability of the articles representing the cultural context of each language
 145 community across Wikipedia editions.

146 In this work we extend our previous study (Miquel-Ribé and Laniado, 2016) adding a rigorous manual
147 assessment of the accuracy of the method, an analysis of the creation of cultural context content over
148 time, and a deeper analysis of cross-language coverage.

149

150 **2. Methods**

151 ***Dataset Construction: Cultural Context Content (CCC)***

152 In this section, we describe the method employed to map Wikipedia articles to the cultural context(s)
153 in every language edition with the aim of constructing a dataset. First, we report the selection of the
154 list of languages to be included in the study. Second, we explain the criteria by which we include an
155 article into the dataset. Third and finally, we propose a mechanism to manually assess the performance
156 of the method.

157 **List of Languages.** For the study of cultural context, we consider that having a rich and diverse list of
158 languages increases its value. The selection of languages includes the 30 largest Wikipedia language
159 editions in terms of number of articles (as of July 2015³): Arabic, Catalan, Cebuano, Chinese, Czech,
160 Danish, Dutch, English, Finnish, French, German, Hungarian, Indonesian, Italian, Japanese, Korean,
161 Malay, Norwegian, Persian, Polish, Portuguese, Romanian, Russian, Serbian, Spanish, Swedish,
162 Turkish, Ukrainian, Vietnamese and Waray. To cover diversity, we take into account different
163 sociolinguistic factors and we decided to add 10 language editions to the initial list of 30; at least one
164 language edition per continent, having various linguistic roots, various speaking community sizes, and
165 various editing community sizes. The 10 added languages are: Afrikaans, Basque, Estonian, Guarani,
166 Greek, Hebrew, Icelandic, Macedonian, Nepali and Swahili. For the analysis of CCC article creation
167 over time, we select a reduced subset of 15 language editions.

168 **Cultural Context Content.** Once languages are selected, it is necessary to map the content of each
169 Wikipedia language edition to their cultural context concepts. The aim is to elaborate a method to
170 collect a comprehensive set of Cultural Context Content articles (from now on referred to as CCC) for
171 every Wikipedia language edition. The CCC encompasses a wide variety of topics to represent the
172 shared concepts linked to the corresponding territories. We formalize that CCC articles deal with
173 concepts that have been: a) originated in the context, or b) located in that context and have had a
174 considerable influence there. In addition, in some contexts where two languages are spoken, their
175 speakers may even share some of its concepts (they refer to the same objects or places), and at the same
176 time, geographically widespread languages may be spoken in geographically distant contexts because
177 of historical reasons. With this method, we created individual CCC datasets for every language,
178 including all the cultural contexts of their speakers. This implies that languages that are official in
179 several countries will conform a single dataset encompassing the diverse concepts of these contexts.

180 **Language-Territory Mapping.** After having taken all this into consideration, before being able to
181 elaborate the method, we still need a first ground-truth with some reliable and central concepts for each

³ https://meta.wikimedia.org/wiki/List_of_Wikipedias

Imbalances Across 40 Language Editions

language related cultural contexts. In this sense, we identify for each language: the language name, geographical entities (top political territories such as country and region names) where it is spoken, and its demonyms. To do so, it is necessary to use the ISO 639-2 and 639-3 codes already employed by Wikimedia Foundation to classify Wikipedia language editions (e.g., ‘ru’ for the Russian language Wikipedia: ru.wikipedia.org), as well as the ISO 3166 and 3166-2 codes to identify each country and its subdivisions at a regional level. These codes are widely used on the Internet in geolocation services. In this way it is possible to pair each of the selected language editions with its native words to specify the territories where it is official or indigenous, their inhabitants’ demonyms and the language names (e.g., eswiki españa mexico … español castellano) (see *Wikipedia Cultural Diversity Observatory* for the complete list⁴). This word list is generated by automatically crossing language ISO codes and the Ethnologue⁵ databases, which contain the territories where a language is spoken and their names in their corresponding language. This is especially relevant for those languages which are only spoken or official in a specific region of a country. The generated list is subsequently manually revised and extended (using information from the specific articles in the correspondent Wikipedia language edition) with second names for the same language and demonyms, which are introduced primarily in singular masculine, feminine and plural when available, and with information from the Wikidata database property ‘demonym’⁶.

Article Selection and Retrieval. Once the language-territory mapping keywords list is obtained, a computational implementation of the method is developed applying and integrating the **three strategies** described below. The method uses the databases of each Wikipedia language edition, which are updated in real time (we were granted access to them by the Wikimedia Foundation⁷). The first two strategies gather the articles considered totally reliable, while the third collects some undesired ones that need to be automatically filtered at a later stage.

The first strategy (i - geocoordinates) consists in examining the article location tags `{{#coordinates}}` and the information located in the geotags table, such as the geocoordinates and the ISO code, in order to obtain articles clearly located within the specified territories for each language edition. Articles satisfying this first criterion are directly retrieved from the databases of each Wikipedia language edition. Nonetheless, the implementation of geocoordinates is unequal throughout the different language editions and may contain errors. Therefore, articles with coordinates are verified using a *reverse geocoder* tool in Python⁸. Such tool returns an ISO code that needs to be verified in the ISO codes database to see whether the article is located in a territory associated to the language or not. As a last step, it is possible to add articles that are not tagged with coordinates and do not have a territory ISO code, but that can be matched to the corresponding articles in other language editions, where they are properly geolocated (e.g., an article about a city in Nepal which is not geolocated in the Nepali Wikipedia, but it is in the English Wikipedia).

⁴ <https://github.com/marcmiquel/WCDO>

⁵ <https://www.ethnologue.com>

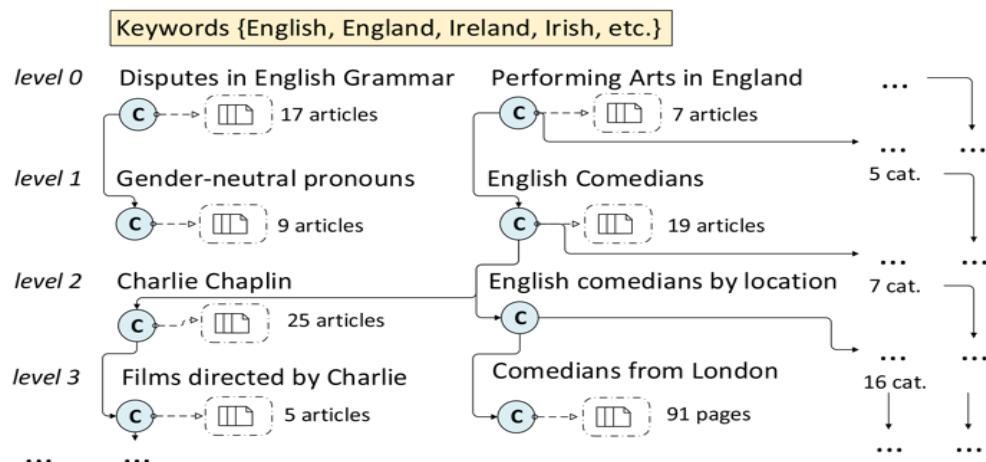
⁶ <https://www.wikidata.org/wiki/Property:P1549>

⁷ <http://wikitech.wikimedia.org>

⁸ <https://pypi.python.org/pypi/pygeocoder>

217 **The second strategy (ii - keywords)** implies examining the articles that contain in their title keywords
 218 related to the language or to the corresponding territories (e.g., "England National football team",
 219 "English law", etc.). These two criteria ensure a high reliability, but unfortunately, they cannot
 220 guarantee that all the articles which should belong to CCC are actually included.

221 **The third strategy (iii - categories)** aims at retrieving the articles more generally related to the
 222 identified keywords. Wikipedia articles are classified into categories that are named according to the
 223 topics developed in the articles. These categories are organized in a hierarchical tree structure. As the
 224 category hierarchy is manually curated and maintained by each community, it tends to be very rich,
 225 although it may contain some noise. The hierarchical structure can be leveraged to identify subareas of
 226 the encyclopedia related to a particular topic: starting from some category, it is possible to crawl down
 227 the classification structure, and gather all the articles belonging to the category and recursively to its
 228 subcategories. In a similar way to the second strategy, we start from the list of keywords associated to
 229 a language and the corresponding territories, and retrieve all the categories that include such keywords
 230 in their titles; for example: "Performing Arts in England" or "Disputes in English Grammar" in the
 231 English Wikipedia. These categories contain articles and other categories which contain in turn more
 232 specific articles (see Figure 1), until at a certain level the process of crawling and gathering articles
 233 finishes. The precise point where the process ends depends on how the category structures have been
 234 constructed (smaller Wikipedia language editions in number of articles also tend to have a less
 235 developed category graph).



236
 237 *Figure 1. Crawling down the category graph with keywords (strategy iii).*

238 The main advantage of this strategy is that it allows to obtain articles related to some keywords.
 239 However, the distance to the top also matters: while the category "Films directed by Charlie Chaplin,"
 240 is part of the category "Performing Arts in England", its content will be considerably more specific.
 241 The further from the top category containing the keyword, the more specific and less related to the
 242 original top category topic the articles will be. The drawback of this category crawling is that
 243 sometimes the categorization includes circular references or does not follow a specialization path (e.g.,
 244 occasionally a more general category appears under a more specific one, other times a category appears
 245 to be related to the immediately preceding one, but totally unrelated to the preceding ones). Such
 246 phenomena may produce interferences in the collection (e.g. the category "Wars involving the United

Imbalances Across 40 Language Editions

247 States" includes the category "World War II", which in turn leads to articles about the German army
 248 and makes them appear as related to the English Wikipedia cultural contexts). Because of this
 249 interference issue, when we use this method on the English Wikipedia we set a limit of five levels of
 250 iteration, i.e. when moving down towards more and more specific categories, we stop at the fifth level.
 251 As the category trees are simpler and less entangled in the other Wikipedias, in the rest of language
 252 editions we complete the iterations until the down category graph goes extinct.

253 **Filtering.** Considering that most, but not all of the articles collected using this third strategy can be
 254 considered CCC, we tackle possible interferences with a filter. In order to be effective, the filter has to
 255 discriminate whether the article is related to the editors' cultural contexts; as a proxy for assessing this
 256 thematic coherence, we look at the extent to which the links contained in the text of an article point to
 257 other CCC articles. As a starting point, we rely on the articles identified with the first two strategies:
 258 the geolocated articles and those including the keywords in their title, which we take as an initial
 259 reliable set of CCC articles. We then iteratively add to this set the articles from the bulk category
 260 crawling selection that have at least 15% of their links pointing out to these articles. While the
 261 algorithm usually converges and stops adding more articles after the third iteration, in large Wikipedia
 262 language editions such as the English it is necessary to limit the algorithm before too many articles,
 263 including false positives, start to be included; we decided to stop the algorithm after the fifth iteration.
 264 Using this procedure, we obtain a final CCC slightly smaller than we would obtain taking all the articles
 265 from the category crawling selection, and we are able to avoid most of the false positives.

266 **Method assessment.** In order to validate the method, there has to be agreement over the nature of CCC
 267 articles (i.e. it is valid) and whether this is a stable construct that the method can identify in a consistent
 268 way (i.e. it is reliable). To examine the method's validity and reliability, we select German and
 269 Japanese, and propose an inter-rater reliability test between 3 raters and the algorithm, calculating the
 270 Cohen's Kappa coefficient (Cohen 2016). In this way we can assess the agreement between human
 271 raters and test the accuracy of the automatic method as compared to expert human judgement.

272 We randomly selected 100 articles classified as CCC, and 100 non CCC articles from the German
 273 language edition. The same process was applied to the Japanese edition. We relied on Google translator
 274 to translate the text of each article into English, for the raters to understand the article content.
 275 Subsequently, the three raters manually classified the articles as positive or negative, i.e. as belonging
 276 or not to CCC.

277 The results of the inter-rater assessment are shown in Table 1, which reports, for the two language
 278 editions, the agreement between the algorithm and the raters, as well as the interrater agreement.
 279 Overall, the degree of agreement between human raters is beyond the 95% in all cases, and with a
 280 Kappa coefficient over 0.9, which confirms that there is agreement over what makes an article belong
 281 to the CCC category. The agreement between the human raters and the algorithm is slightly lower than
 282 between humans, but still satisfactory (nearly 90% agreement and a Kappa coefficient of 0,76 in
 283 average), confirming the reliability of the automatic method.

284 **Table 1. Inter-rater reliability tests for the Japanese and German Wikipedia language editions.** For
 285 each Wikipedia we crossed the ratings (CCC) from three raters. Legend: coincidence is the degree of
 286 coincidence in %, and K is the Cohen's Kappa coefficient.

Inter-rater reliability	Japanese		German	
	coincidence	K	coincidence	K
algorithm-rater1	0.86	0.71	0.90	0.8
algorithm-rater2	0.89	0.77	0.91	0.82
algorithm-rater3	0.86	0.72	0.89	0.77
rater1-rater2	0.97	0.94	0.96	0.93
rater1-rater3	0.97	0.93	0.95	0.9
rater2-rater3	0.96	0.91	0.98	0.95

287

288 Non-CCC articles given as positive by the algorithm are mostly articles about specific topics from
 289 adjacent countries, or articles related through incidental relationships, for instance a basketball player
 290 who competed for one of the countries associated to the language. The cases of disagreement between
 291 raters concerned articles partially related to a particular territory or language, which lend themselves
 292 to different interpretations. For instance, there was disagreement about the article “Bronvaux”, a French
 293 municipality in the region of Lorraine, close to the German border and historically disputed between
 294 the two countries. The article in the German Wikipedia is categorized as “Historical Territory
 295 (Germany)”, and this is why also the algorithm considered it as part of CCC for the German Wikipedia.
 296 One rater however considered that being located in France, the article should be part of CCC only for
 297 the French Wikipedia. In another case, there was disagreement on the article “ECN-T002” which refers
 298 to a mobile phone model released in 2009 by Toshiba, a Japanese company. While two raters
 299 considered this kind of creation to be part of the cultural context, the other one argued that technological
 300 products on the global market should not be associated with a cultural context. The algorithm did not
 301 assign the article to CCC. As the borders of cultural contexts are fuzzy, this kind of disagreements may
 302 be inevitable. Instead of confronting imported and original concepts, we argue that the selection of
 303 CCC articles should be seen as a *continuum* going from central to peripheral relevant concepts.

304 At this point, in order to evaluate the overall accuracy of the method we repeated the manual assessment
 305 procedure with one human rater for the rest of the Wikipedia language editions and computed the F1
 306 Score. The results are presented in Table 2, which details the percentage of false positives (FP) and
 307 false negatives (FN) with the resulting F1 score for each language edition. We observe that false
 308 positives are on average the 8.1%, and false negatives the 5.9%. The average value of F1 is 0.92. The
 309 selections with more false positives are Korean and Serbian (19% and 23%, respectively). The results
 310 for these two language editions appear to be affected by categorization issues. In the former case, the
 311 Korean Wikipedia has a category tree where subcategory relationships do not always reflect a strict
 312 hierarchical structure, so that the algorithm gets to include unrelated concepts; at the same time, the
 313 15% threshold on the outlinks is not always sufficient to filter out noise in this case due to the high
 314 presence of very short articles that only include very few links. In the latter case, the Serbian Wikipedia

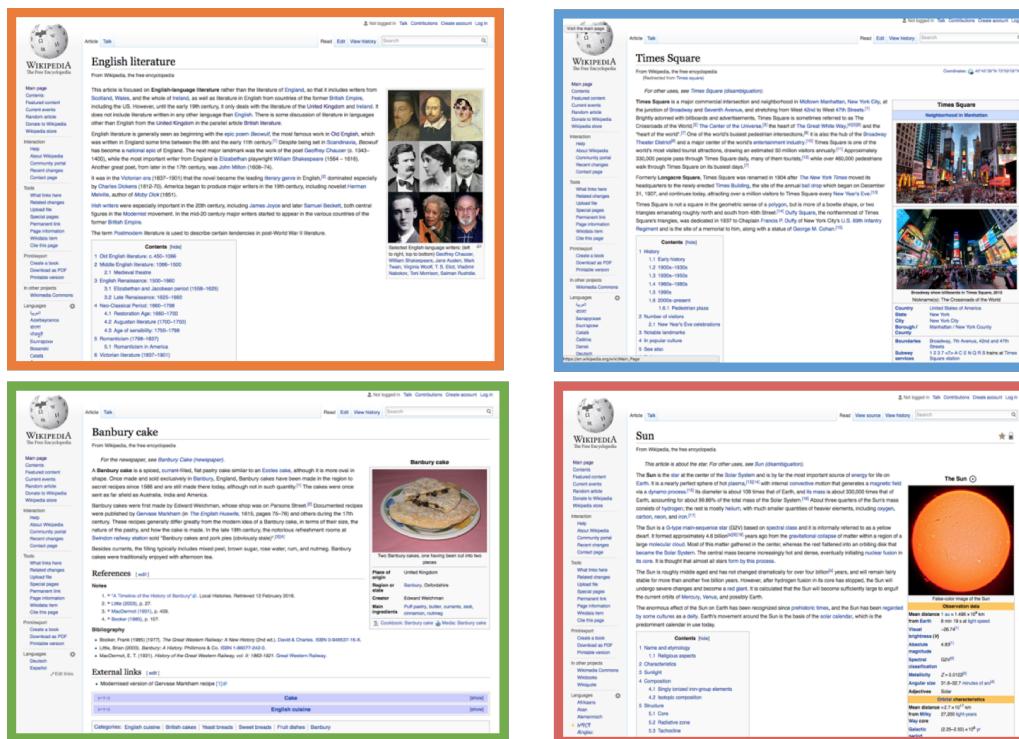
Imbalances Across 40 Language Editions

still employs ‘Yugoslavia’ as a label for its categories and tends to encompass also non-Serbian territories, therefore the false positives detected through the method assessment actually reflect an inconsistency in the data due to a geopolitical conflict.

3. Results

3.1 RQ1. Extent of the Cultural Context Content Articles

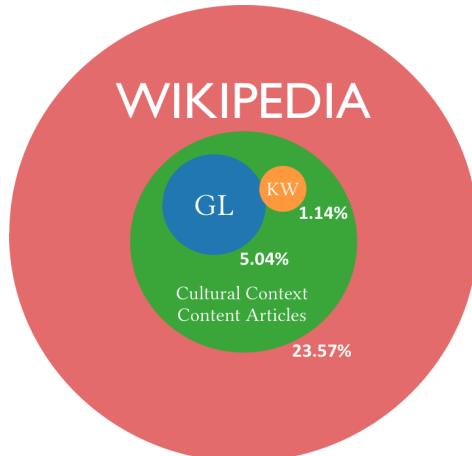
Before answering the first research question, it is worth introducing four prototypical articles from English Wikipedia that *represent the various types of content selected*: (1) CCC keywords, (2) CCC geolocated, (3) the rest of CCC as part of the constructed datasets, (4) the rest of Wikipedia (Figure 2). For instance, for the English language edition, a good example of CCC Keyword is ‘English Literature’, because it perfectly explains the content of the article. The articles from the category CCC keywords are often a synthesis of a topic aggregated by the demonym or the territory name (e.g. English writers’ biographies and works). An example of the CCC geolocated articles, the ‘Times Square’ article, contains the name of a geographical territory associated to the English Wikipedia. Even though this is a very iconic place, in CCC geolocated there are articles with all levels of notability – from small towns to nationally renowned companies and famous monuments. A good example of the rest of CCC articles is the ‘Banbury Cake’ article. After the CCC geolocated and CCC keywords articles, the rest of CCC articles dedicated to specific themes of local scope represent the majority in CCC. An example of an article from the rest of Wikipedia could be for instance the article ‘Sun’, an article containing universal knowledge not related to any cultural context in particular.



334

Figure 2. Examples of articles from English Wikipedia. CCC keywords (English Literature), CCC geolocated (Times Square), the rest of CCC (Banbury cake) and the rest of Wikipedia (Sun).

337 **Results.** The Venn diagram shown in Figure 3 presents the average proportion of CCC articles in the
 338 40 considered language editions, and the breakdown into articles identified via the first (geolocation)
 339 and second (keywords in the title) strategies. We observe that about 1 out of 5 CCC articles were
 340 identified via geo-coordinates, and only about one out of 20 via keywords in the title. The intersection
 341 between the two subgroups is rather small. The proportion of articles identified through the third
 342 strategy (category structure) are omitted, as they represent almost the totality of CCC (29.5% on
 343 average).



344

345 **Figure 3. Average proportion of CCC, and of CCC detected through geolocation and keywords.**
 346 *Sizes are in scale according to their proportion.*

347 As shown in Table 2, almost a quarter of each Wikipedia language edition (mean 23.2%, median
 348 24.2%, standard deviation 11.1%) belongs to CCC articles (**RQ1**). These results indicate that a non-
 349 negligible percentage of each Wikipedia is dedicated to concepts representing the cultural contexts
 350 associated with it. Table 2 shows the total number of articles and the percentage of articles classified
 351 as CCC for each of the 40 language editions considered. Furthermore, the table shows the breakdown
 352 according to the different strategies through which CCC articles have been identified, i.e. through
 353 Strategy 1 (through geolocation tags) or Strategy 2 (keywords in the title). As above, the percentage
 354 for Strategy 3 (category crawling) is not reported, as for most language editions it is very near or almost
 355 equal to the final percentage of articles included in the CCC set.

356 The comparison of CCC percentages across languages shows that there is no obvious pattern. While
 357 for its role as an international reference one could expect the English Wikipedia to have a lower
 358 proportion of articles associated with its specific cultural context, the results actually show that it has
 359 the second highest percentage of CCC articles (46.8%), exceeded only by the Japanese version
 360 (49.2%), while the proportion is below 40% for the rest of language editions.

361 The extremely low percentage (0.1%) for Cebuano and Waray-Waray reflects that these language
 362 editions have a high number of articles but are mostly made by bots that automatically translate articles
 363 from other languages versions. This observation points out that the creation of CCC articles is
 364 inherently connected to the presence of an active and engaged community.

Imbalances Across 40 Language Editions

Table 2. Percentage of CCC articles in Wikipedia language editions and CCC cross-language coverage. For each of the 40 editions considered, columns report: total number of articles (WP art); percentage of CCC articles over the entire Wikipedia (CCC %), and percentage of articles identified through the first strategy - GeoLocated tags (GL %) and through the second strategy - KeyWords in their titles (KW %); percentage of False Positives (FP %) and False Negatives (FN %), with resulting F1-score (F1) after manual evaluation; average number of Interlanguage links per article in the language edition (ILL WP) and in CCC (ILL CCC), percentage of CCC articles having no ILLs, and percentage of CCC articles having no ILLs with respect to WP articles having no ILLs (CCC NO ILL / WP NO ILL).

374

ISO cod.	Language	WP Art.	CCC %	GL %	KW %	FP %	FN %	F1	Avg. ILL WP	Avg. ILL CCC	CCC NO ILL %	CCC NO ILL / WP NO ILL
af	Afrikaans	35966	19.2	5.9	0.9	1	2	0.99	40.1	4.5	34.3	76.2
ar	Arabic	375282	26.9	3.2	2.4	2	18	0.91	12.9	3.6	59.5	54.5
eu	Basque	208630	10.1	1.7	0.4	4	1	0.97	14.4	1.3	50.6	73.1
ca	Catalan	467486	16.2	7.9	0.8	2	3	0.98	21.5	3.6	68.7	62.7
ceb	Cebuano	1211531	0.1	0.0	0.1	12	1	0.93	15	1.6	0.6	0.1
zh	Chinese	851670	32.9	6.3	1.2	10	11	0.90	6.3	11	58.2	63.4
cs	Czech	326187	25.9	9	1.2	2	3	0.98	4.8	8.9	60.3	71
da	Danish	205764	31.7	6.1	1.0	10	2	0.94	10.0	2.6	52.3	73.4
nl	Dutch	1828148	7.8	1.6	0.3	1	3	0.98	13	1.8	64.4	22.4
en	English	4917741	46.8	9.8	2.8	10	16	0.87	6.8	1.5	55	63.1
et	Estonian	136362	31.1	6.1	1.7	2	2	0.98	20.2	1.8	64.4	69.8
fi	Finnish	375347	21.9	2.3	1	1	4	0.98	6	2.9	70.2	70
fr	French	1642276	29.0	6.9	1.7	10	6	0.92	23.2	4.7	46.2	59.3
de	German	1834147	36.8	8.8	1.9	10	10	0.90	15	2.5	60.1	62.5
el	Greek	108090	33.5	6.4	0.6	9	8	0.91	17.9	4.2	46.1	71.9
gn	Guarani	3031	23.6	14	3.3	2	6	0.96	82.1	24.2	6.9	57.6
he	Hebrew	174667	31.7	2.1	1.6	14	2	0.91	20	4.8	50.3	79.5
hu	Hungarian	326146	18.5	1.9	1.5	10	4	0.93	16	2.9	54.8	61
is	Icelandic	39554	30.7	2.2	1.5	2	10	0.94	12	1.7	66.1	74.2
id	Indonesian	363529	27	1	0.6	7	4	0.94	33.7	2.4	36.7	73.8
it	Italian	1210801	19.2	3.6	0.7	10	5	0.92	9.3	3.5	54.5	48.4
ja	Japanese	973955	49.2	3.4	1	4	10	0.93	7.1	1.2	75.9	77.5
ko	Korean	320742	32.6	2.4	0.8	19	12	0.84	14.1	7.8	69.9	58.6
mk	Macedonian	82743	15.9	2.5	1.3	15	4	0.90	25.3	3.4	41	51.6
ms	Malay	275031	19.5	1.4	0.8	15	2	0.91	15.5	1.8	55.4	61.5
ne	Nepali	29114	29.7	11.8	2.2	2	19	0.90	22	3.3	41.4	29.5
no	Norwegian	415015	26.8	5.5	0.8	12	6	0.91	12.4	2.3	54.3	72.5
fa	Persian	460523	11	10.3	0.7	3	19	0.90	7.6	4.8	8.2	4.7
pl	Polish	1122218	23.2	9.4	1.1	9	3	0.94	9.4	1.3	58.1	54.2
pt	Portuguese	880529	19.1	2	1	6	1	0.96	11.2	2.4	64.5	58.4
ro	Romanian	329925	20.7	7.2	1.1	13	2	0.92	16.9	3.5	32.6	72.7
ru	Russian	1237127	31.2	11	1.1	14	5	0.90	8.3	2.2	44.6	56.7
sr	Serbian	321912	12.1	3.2	0.1	23	1	0.87	16	4.7	25	50.3
es	Spanish	1147742	27.7	5	2	13	6	0.90	9.3	3.4	44.3	64.9
sw	Swahili	29168	18.3	3.6	1	1	6	0.97	40	3.7	46.8	73.8
sv	Swedish	1970808	11.4	4.3	0.4	8	4	0.94	6	1.4	72.5	66
tr	Turkish	249061	33.9	4.4	2.1	6	4	0.95	16.2	3.4	36.4	70
uk	Ukrainian	581735	24.8	6.8	1	14	2	0.91	13	2.4	43.1	57
vi	Vietnamese	1137180	2.5	0.9	0.2	5	0	0.97	7.4	1.5	72.8	17
war	Waray	1259278	0.1	0.0	0.0	23	0	0.87	6.3	10.9	12.6	1.9
Avg.	Average	736654	23.3	5.1	1.1	8.1	5.9	0.92	16.6	4	49	57.1

375

376 **3.2 RQ2. Cultural Context Content Articles Over Time**

377 The considerable extent of CCC shows that editors engage in contributing with content related articles
378 to their context. One could think that topics about one's very near context may be finite or would stop
379 being notable, especially if compared to the amount of universal content which deserves being included
380 into an encyclopaedia. However, we hypothesize that editing CCC could be an activity sustained over
381 time, as editors may feel attached to their cultural context and keep enriching it in their language
382 edition. To verify this, we propose an analysis of how CCC has been created over time. Such analysis
383 may explain the most productive period and predict future scenarios. To investigate whether the
384 creation of cultural context content is consistent over time, we count the number of CCC articles created
385 every year in a Wikipedia language edition, since its creation until January 2016, and compare it to the
386 overall number of articles created every year.

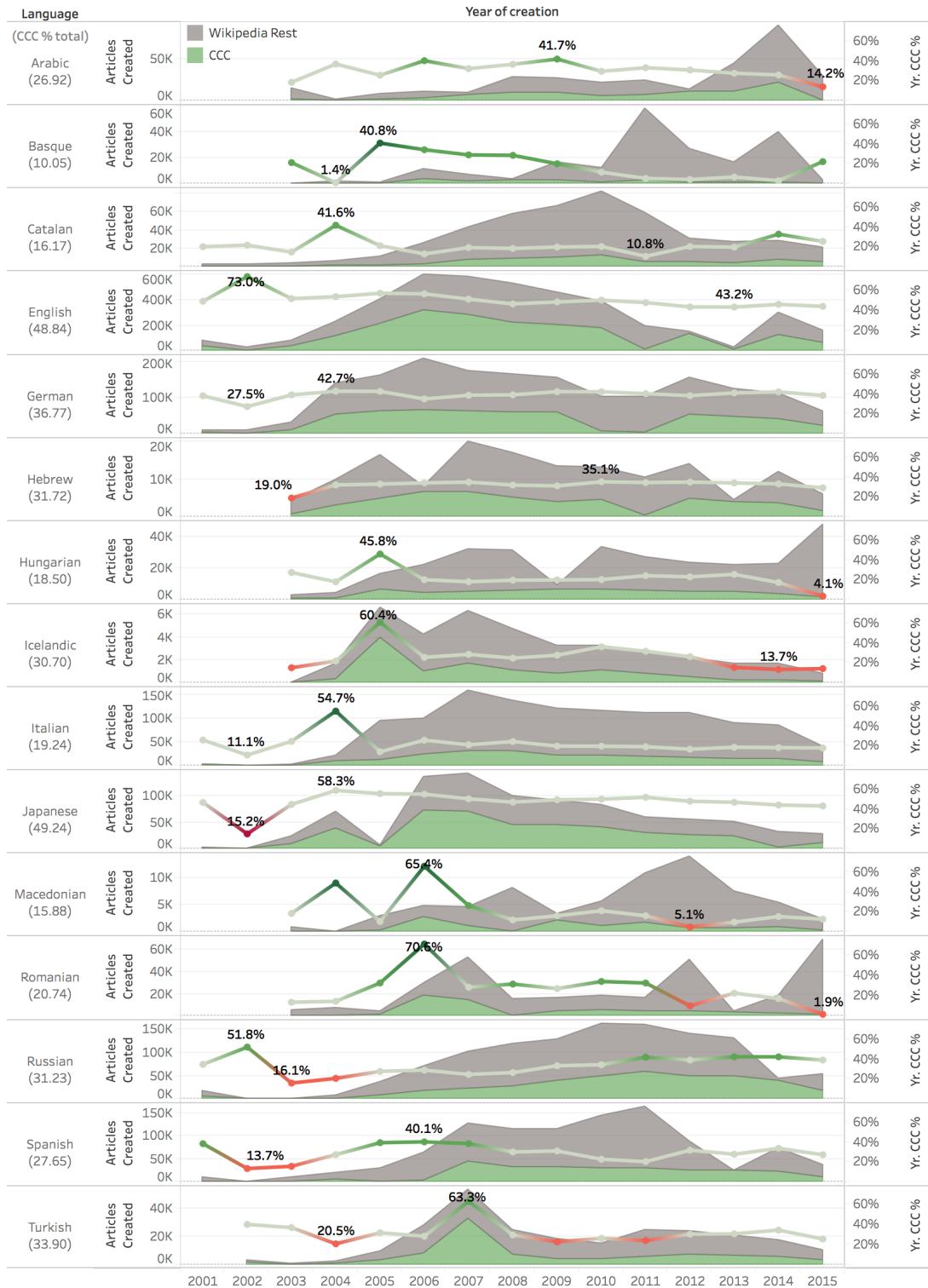
387 **Results.** Figure 4 shows the growth of each Wikipedia language edition in terms of number of articles.
388 CCC articles are depicted in green, while the rest of the articles in grey. Figure 4 also represents the
389 percentage of CCC created every year (green and red indicate respectively values above and below the
390 overall percentage). In general, CCC creation tends to remain as a stable part of the activity over the
391 years, although some general patterns can be noticed (**RQ2**). The most prolific period tends to be
392 located between 2005-2010, when Wikipedia language editions experienced their most important
393 growth. It is the same period when the highest percentages of CCC for most languages occurred, which
394 suggests that the most important bursts in content creation have been dominated by the local cultural
395 context.

396 Usually CCC has grown parallel to Wikipedia, but in those years, it also grew more proportionally,
397 occupying an important percentage of the entire amount of Wikipedia articles. After the years of
398 "content boom", the proportion of CCC tends to get stabilized for most of the languages and does not
399 decrease. Generally, large Wikipedia language editions with strong communities, such as the English
400 and the German ones, exhibit a more balanced growth, less affected by spikes in the creation of content,
401 as it happens for instance in the Icelandic or the Macedonian Wikipedia.

402

403

Imbalances Across 40 Language Editions



404

405 **Figure 4. CCC creation over the 15 years of Wikipedia.** For each language edition, the green area
 406 represents the absolute number of CCC articles created over years, and the grey area the rest of the
 407 articles created. The line shows the percentage of CCC over the total number of articles created
 408 during each year; it is depicted in grey when it is in line (less than 10% variation) with the final
 409 overall percentage of CCC in the encyclopaedia, in green or red when it is higher or lower.

4B03 RQ3. Cross-language Coverage of Cultural Context Content

411 To address our third research question concerning cross-language coverage of CCC, we look at the
 412 Interlanguage links (ILL), i.e. links that connect the same article in two different languages.
 413 Interlanguage links may be created either by editors or by automatic bots and allow one to map content
 414 coverage between language editions. Previous work has shown that many articles tend to be created
 415 first in large language editions, and then translated and re-adapted into smaller language editions
 416 (Warncke-Wang et al. 2012). In our case however we expect to find different patterns; as we focus on
 417 content that is specific to each cultural context, the presence or absence of interlanguage links is an
 418 indicator of the degree of uniqueness, while the interwiki links towards specific language versions
 419 show the coverage that each cultural context receives from other language communities.

420 3.3.1 Interlanguage links analysis.

421 **Results.** As seen in Table 2, the average number of ILLs per article is variable across languages, both
 422 in CCC articles and in the entire Wikipedia. The average for CCC articles is 4.15 times lower than the
 423 overall average (**RQ3**). Therefore, CCC is less shared across languages, and part of the language gap
 424 is due to the fact that the content representing the cultural context is not shared across languages.
 425 Namely, we can affirm that in the language gap there is a **culture gap** (where by culture gap is intended
 426 the CCC articles not shared across languages). Even though in most cases, the average number of ILLs
 427 in CCC is lower than in the entire Wikipedia, the ratio (avg. ILLs CCC / avg. ILLs WP) is also variable
 428 across languages. In fact, minor language editions like Icelandic, Afrikaans, Estonian and Swahili have
 429 between 7 and 11 times less ILLs in CCC than in their entire language edition. On the contrary,
 430 languages like English, French, Korean, German and Italian, which represent the largest Wikipedia
 431 language editions, show smaller differences between ILLs in CCC and overall. This suggests that both
 432 language status and development degree of a Wikipedia language edition may strongly influence
 433 whether its CCC articles are created into other languages.

434 In order to further investigate the culture gap in each language edition, we measure the percentage of
 435 articles with no ILLs in CCC and the entire WP. This allows us to observe the degree to which CCC
 436 articles are responsible for the differences in content imbalance between Wikipedia language editions.
 437 Results show that languages with a high percentage of articles with no Interlanguage Links (WP NO
 438 ILLs) also tend to have a high percentage of CCC articles. In fact, CCC articles with no ILLs account
 439 for the majority of Wikipedia content with no ILLs in most languages (mean 62.83%, median 63.25%
 440 and standard deviation 12.31%, without taking into account the results for languages such as
 441 Vietnamese, Waray-Waray and Cebuano, where the automatic program bot had a major contribution
 442 in the creation and translation of articles from other language editions). This confirms again that to a
 443 great extent the culture gap is responsible for the language gap between Wikipedia language editions
 444 described by (Warncke-Wang et al. 2012).

445 3.3.2 CCC cross-language availability.

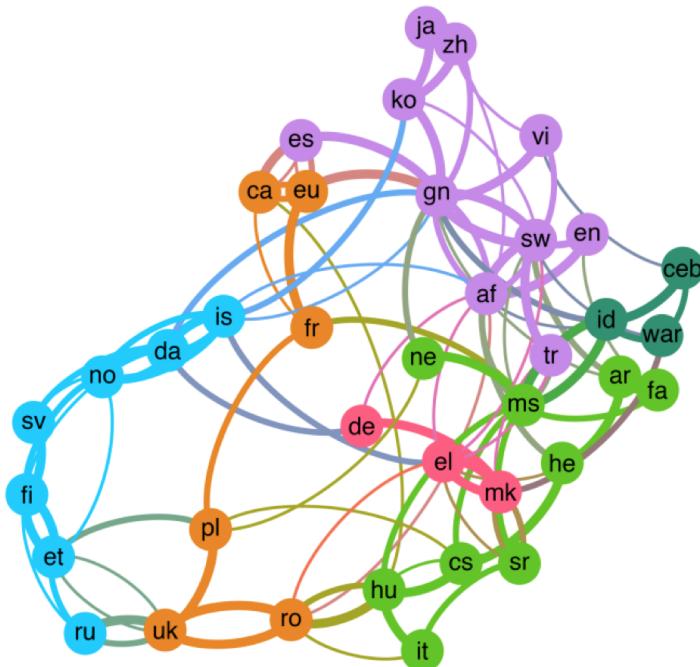
446 **Results.** Taking a closer look at CCC's Interlanguage links, it is possible to obtain a better
 447 understanding of the proximity between language communities, as indicated by the availability or
 448 expansion of CCC across Wikipedia language editions. To study such proximity, we compute the
 449 proportion of CCC articles from a particular language edition that can be found in other language

Imbalances Across 40 Language Editions

450 editions (e.g. the proportion of the Italian CCC articles in the Catalan Wikipedia). In Figure 5 we depict
 451 a network of languages and show which ones have a higher proportion of CCC articles represented in
 452 other language editions. In order to create this graph, for each CCC we select the three languages where
 453 it is represented in the highest proportion and draw the corresponding edges. The network is therefore
 454 directed, as a link from language A to language B implies that B is one of the three language editions
 455 with better coverage of A's CCC articles, and this relationship is not necessarily reciprocal as A could
 456 have a poor coverage of B's CCC. Following a standard convention in graph representation, edges are
 457 curved and drawn in clockwise direction. Colors are assigned according to the clusters identified by
 458 the Louvain community detection algorithm (Blondel et al. 2008) to highlight groups of language
 459 editions that are closer to each other.

460 Nordic languages form a cluster which includes Russian, while languages of the Iberian Peninsula are
 461 tightly connected to each other, as well as languages of Asia and Middle East. These results point out
 462 the relevance of geographic proximity and seem to confirm Tobler Law's idea according to which
 463 things near tend to be similar. These finding is in line with the comparison of biographical articles'
 464 availability in different languages (Aragón et al. 2012; Eom et al. 2015). However, some less expected
 465 relationships also become apparent, such as the relevance of Italian CCC articles in the Hungarian
 466 Wikipedia.

467



468

469 **Figure 5. Network graph of language proximity in terms of shared CCC articles.** Each node
 470 represents a Wikipedia language edition and has three outgoing links to the three language editions in
 471 which its CCC represents the highest percentage of the articles. Links are represented in clockwise
 472 direction. Colors represent clusters of language editions identified through the Louvain algorithm for
 473 community detection.
 474

475 3.3.3 Mapping the culture gap

476 **Results.** To see how well each Wikipedia language edition covers the CCC articles of other languages,
477 we created Figure 6. The entire table allows to see the culture gap of each language edition, and how
478 this also depends on linguistic and geographical proximity. However, it seems the factor of scale is
479 more important, since wide language editions (in number of articles and created by large communities
480 such as English, German, French, etc.) cover a higher percentage of the CCC articles of language
481 editions that are significantly smaller.

482 Generally, the culture gap highlights a common difficulty in achieving a representation of cultural
483 diversity, indicating that editors are often not able to cover concepts from other cultures. Few languages
484 cover a good percentage of the CCC of other languages. English is one of them, but still it only covers
485 on average a 33.71% of the CCC articles of other languages (median 28.27%, standard deviation
486 19.36%). Conversely the CCC of large language editions such as English or German are poorly
487 represented in the other language editions (barely 5% of the English and German CCC articles are
488 found in other language editions). Considering the dimensions of the English and German editions as
489 well as the difficulty of translating a large percentage of relevant articles for their culture into other
490 languages, such a gap is not surprising.

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

Lang	af	ar	ca	ceb	cs	da	de	el	en	es	et	eu	fa	fi	fr	gn	he	hu	id	is	it	ja	ko	mk	ms	ne	nl	no	pl	pt	ro	ru	sr	sv	sw	tr	uk	vi	war	zh	
af	1.05	0.11	0.63	0.21	1.57	0.67	1.34	0.62	0.41	0.22	0.75	0.11	0.17	0.53	13.99	1.49	0.24	0.21	0.53	0.41	0.22	0.55	0.52	0.46	0.52	0.98	0.29	0.14	0.40	0.46	0.36	0.43	0.23	1.37	0.71	0.40	0.37	1.71	0.62		
ar	4.13		2.43	6.96	1.60	5.57	3.46	7.99	4.21	5.17	0.90	4.94	30.45	1.47	4.34	18.74	12.15	2.57	1.31	1.37	4.05	1.61	2.90	4.64	2.17	1.89	1.90	1.21	0.49	2.56	3.48	1.91	3.20	1.33	3.32	6.57	1.88	1.60	4.71	3.29	
ca	5.89	7.57		16.46	2.33	7.00	4.88	13.46	4.26	18.10	1.47	19.08	1.04	1.25	15.28	31.61	8.65	3.53	0.94	2.96	7.06	1.31	2.40	3.30	1.33	2.13	3.14	1.97	0.88	3.96	3.96	2.42	4.06	1.75	5.32	5.33	2.30	1.33	8.78	3.43	
ceb	0.43	0.80	0.39		0.09	0.28	0.59	0.19	1.86	2.14	2.97	1.65	0.62	0.16	8.81	9.65	1.40	0.38	1.19	0.12	0.22	0.19	0.34	2.77	0.74	0.82	0.26	0.76	0.09	3.71	0.33	0.30	3.06	0.36	0.81	0.46	0.60	0.81	0.48	78.59	4.45
cs	3.77	5.14	1.08	3.80		6.84	5.22	6.75	3.13	2.85	2.72	3.29	0.73	2.04	5.46	17.34	14.57	6.62	0.60	3.95	3.75	1.18	2.60	2.91	1.18	1.75	1.97	2.11	2.45	2.05	4.27	3.23	6.39	2.09	2.94	3.70	3.39	1.51	1.50	2.97	
da	4.58	3.20	0.93	0.63	1.53		3.90	3.90	2.22	2.17	1.69	2.62	0.51	1.66	2.69	16.92	5.51	2.41	0.46	5.34	2.05	0.80	1.99	1.33	0.94	2.36	1.58	6.53	0.68	1.64	2.18	1.29	2.08	4.43	4.24	2.36	1.40	1.02	0.86	2.05	
de	28.20	13.13	5.60	58.65	21.21	22.29		22.10	14.79	14.17	12.10	9.91	2.47	7.62	22.87	45.32	21.73	14.11	2.63	17.09	18.96	4.43	6.11	8.18	4.00	7.75	11.32	9.46	8.81	12.26	13.43	9.13	14.39	10.04	13.88	16.29	8.28	4.01	56.10	9.43	
el	1.46	2.88	0.57	0.84	0.82	2.93	1.64		1.05	1.33	0.49	1.73	0.34	0.45	1.53	9.51	3.96	1.18	0.23	1.31	1.78	0.31	1.01	8.89	0.54	1.13	0.52	0.27	0.97	2.28	0.83	3.92	0.58	2.77	3.68	0.89	0.48	1.50	0.98		
en	62.36	33.63	12.66	67.30	24.34	29.75	28.27	42.22		38.12	22.12	22.34	91.36	18.45	40.23	86.29	41.85	24.94	10.00	22.95	30.58	16.86	17.75	25.69	15.40	44.74	22.30	28.27	31.79	25.26	39.53	17.58	49.62	18.87	50.63	31.62	14.40	20.31	68.52	25.65	
es	14.27	11.30	26.04	28.27	4.50	11.12	9.56	18.50	12.62		2.94	39.36	1.66	3.26	19.69	91.05	15.37	6.98	2.38	10.77	14.38	4.28	5.58	6.04	3.43	4.06	5.39	4.47	2.38	16.12	7.45	5.34	10.02	4.08	9.35	7.64	5.11	3.79	20.56	7.30	
et	2.09	1.90	0.30	1.06	1.06	4.52	1.95	3.54	1.08	0.99		1.78	0.45	3.02	1.43	16.92	3.13	1.53	0.35	1.88	1.67	0.42	1.39	1.06	0.74	0.80	0.68	1.86	0.62	0.88	1.65	2.49	1.61	1.74	2.42	1.84	2.11	0.66	2.14	1.48	
eu	2.95	2.96	4.58	16.88	0.74	3.34	1.67	4.10	1.63	8.47	0.68		0.46	0.84	10.81	17.20	3.89	1.34	0.83	2.38	1.75	0.56	1.19	1.07	0.75	2.02	0.93	0.74	0.24	1.60	1.84	0.70	1.30	0.67	2.17	1.85	0.97	0.83	4.50	2.46	
fa	4.87	10.32	1.79	11.39	1.71	5.76	4.12	6.29	6.28	4.37	2.49	3.01		0.97	4.63	26.29	9.04	3.43	1.13	1.57	3.72	1.95	3.15	2.09	2.00	2.82	1.48	1.59	1.02	3.62	3.04	2.38	3.81	1.57	6.28	6.93	2.06	2.18	8.57	3.88	
fi	7.71	5.42	1.52	3.38	3.14	8.25	4.58	12.19	4.58	3.84	5.50	4.21	0.83		4.98	27.55	8.49	4.50	1.01	5.35	4.23	1.67	2.77	3.02	1.88	3.53	2.18	4.61	1.13	3.19	3.47	4.39	4.11	8.31	6.54	4.05	2.76	1.89	3.85	3.62	
fr	24.92	20.66	10.19	41.35	13.10	15.04	17.07	22.79	16.26	20.65	6.83	24.01	2.66	8.49		56.64	21.30	18.16	3.30	12.67	22.53	7.78	7.42	21.85	4.14	15.92	14.65	7.90	12.42	14.88	12.61	8.43	41.58	7.46	14.92	12.05	6.52	6.10	28.69	10.84	
gn	0.12	0.09	0.02	0.21	0.01	0.16	0.04	0.03	0.04	0.24	0.02	0.19	0.00	0.01	0.03		0.10	0.02	0.02	0.04	0.02	0.01	0.09	0.04	0.03	0.12	0.01	0.02	0.01	0.27	0.00	0.02	0.04	0.01	0.26	0.06	0.03	0.05	0.00	0.00	0.07
he	3.84	7.21	0.80	2.11	1.60	5.04	3.09	5.67	2.43	2.14	0.86	2.23	0.63	0.69	2.43	16.50		2.74	0.43	1.26	2.27	0.81	2.04	2.23	0.96	1.34	1.06	1.01	0.97	1.59	3.44	1.88	2.88	1.04	3.00	3.41	2.10	0.93	0.86	2.03	
hu	3.35	4.10	3.72	4.01	11.02	5.75	5.58	6.44	2.69	6.15	1.53	6.69	0.51	1.53	11.85	15.39	5.89		0.62	4.12	8.28	1.05	2.32	4.02	1.00	1.46	2.01	1.87	1.25	2.51	17.03	2.40	9.69	1.72	2.68	3.84	2.88	0.91	1.71	2.57	
id	4.24	5.61	0.79	53.38	1.16	4.45	2.60	4.74	2.51	2.29	0.58	2.43	0.87	0.69	4.83	19.16	6.48	2.02		1.65	3.47	2.04	3.80	1.61	35.96	2.08	2.24	0.92	0.45	2.13	2.41	1.13	2.37	0.95	3.37	3.22	1.04	2.60	53.96	4.10	
is	0.65	0.76	0.09	0.00	0.24	2.71	0.63	1.54	0.42	0.31	0.30	0.86	0.07	0.24	0.44	8.81	1.39	0.71	0.14		0.42	0.16	0.60	0.25	0.33	0.28	0.24	0.94	0.11	0.37	0.49	0.20	0.43	0.63	0.90	0.58	0.34	0.31	0.86	0.48	
it	17.03	11.98	7.26	58.86	14.27	12.23	13.66	23.17	11.86	16.83	4.27	12.68	3.78	4.86	22.95	66.57	15.13	14.49	1.52	9.56		5.24	5.14	7.19	2.65	5.18	6.54	6.52	4.02	11.96	15.00	6.97	12.78	5.35	10.01	9.73	4.87	3.44	54.60	6.35	
ja	7.26	6.66	1.73	6.33	3.59	7.71	6.14	8.71	6.80	5.21	1.59	6.05	1.05	2.10	6.35	24.62	9.99	4.76	1.64	4.08	8.35		12.68	2.99	3.63	4.20	2.56	2.29	1.23	3.93	4.16	3.77	5.06	2.37	8.39	4.89	3.10	4.76	9.85	16.28	
ko	4.66	4.73	0.71	15.61	1.51	5.36	2.93	6.30	3.15	2.86	0.90	2.41	0.74	0.95	2.92	26.01	6.01	2.47	0.98	2.55	3.12	7.10		2.70	2.57	1.96	1.20	1.14	0.53	2.33	2.63	1.82	2.71	1.21	6.46	3.49	1.51	3.29	8.78	9.45	
mk	1.12	1.48	0.17	2.95	0.43	2.23	1.67	5.02	0.56	0.93	0.32	1.26	0.29	0.69	15.94	2.04	0.86	0.17	0.96	0.97	0.17	0.59		0.50	0.54	0.29	0.35	0.18	0.45	1.70	0.74	4.79	0.36	1.31	2.15	0.86	0.34	5.78	0.71		
ms	1.78	4.14	1.92	2.11	7.89	2.17	2.95	1.02	1.56	3.77	0.34	4.37	18.92	0.30	9.14	11.05	2.81	6.70	26.07	0.54	4.12	1.08	1.59	0.88		30.32	0.75	2.18	0.64	1.40	5.57	0.50	12.87	0.41	1.54	9.47	0.56	1.41	5.14	2.06	
ne	0.46	0.25	0.02	0.00	0.03	0.43	0.17	0.09	0.27	0.09	0.09	0.12	0.09	0.03	0.10	2.24	0.38	0.09	0.12	0.06	0.05	0.06	0.28	0.08	0.15		0.04	0.07	0.01	0.11	0.08	0.05	0.04	0.06	0.32	0.17	0.10	0.14	0.00	0.40	
nl	27.27	8.18	5.55	90.72	11.98	17.62	12.28	12.48	9.01	11.68	5.97	10.18	1.35	3.34	21.07	52.17	11.91	11.45	33.86	8.28	11.37	2.55	3.91	4.77	10.96	5.00		6.13	6.79	12.34	10.81	3.91	12.84	6.35	10.72	7.83	3.68	7.82	80.39	6.66	
no	8.95	5.47	1.43	18.14	2.73	18.66	5.27	7.47	4.54	4.48	3.26	3.69	0.89	4.52	4.58	41.96	8.30	3.83	0.84	11.84	4.49	1.35	3.37	2.63	1.85	3.61	2.79		1.23	5.92	3.08	2.88	4.10	8.20	9.07	4.06	2.28	2.18	11.78	5.66	
pl	17.57	11.02	5.96	9.49	16.73	12.46	13.03	16.04	9.76	12.10	14.27	10.38	2.36	5.07	17.																										

517 4. Conclusions

518 In this paper we have presented a study of the content imbalances in Wikipedia language editions as a
519 result of the impact of their cultural contexts. To this aim, we have proposed a method to analyze
520 Wikipedia content and select articles that specifically relate to the cultural context of each Wikipedia
521 language edition. We named such articles Cultural Context Content (CCC), whether they are about
522 geography, people, language, traditions, among other topics. The method relies on a combination of
523 different strategies in order to retrieve articles, leveraging characteristics such as geolocation, specific
524 keywords in the titles, associated categories or links to other articles. We applied it to 40 language
525 editions selected according to a diversity criterion. The method accuracy has been assessed manually
526 resulting in an average of 8.1% of false positives, 5.9% of false negatives, and an accuracy of $F1 = 0.92$.
527

528 **Limitations.** Our work is not exempt of limitations, some of which are intrinsic to the same nature of
529 cultural context, while other are more related to the constraints set by the Wikipedia data structure and
530 the method proposed. Although we established a language-territory mapping for each Wikipedia
531 language edition, the method aggregates all the articles from the different territories into a single
532 generated CCC dataset, despite in some cases they may be geographically distant and share few
533 elements in common but the language. We could consider this a limitation of the current dataset, since
534 it would be much better to have a more fine grained collection which would allow further investigation
535 and applications. For instance, the cross-language analysis proposed could be developed into more
536 depth to bring new insights on particular cultural contexts within a language (e.g., British or US with
537 respect to the English language edition) or even across different ones in the same territory (e.g.,
538 assessing differences and similarities between Ireland CCC in the English and in the Gaelic Wikipedia).

539 In regards to the method, even though the results from the manual assessment can be considered
540 satisfactory, we want to acknowledge several observations. First, we need to be cautious that the
541 generated list of keywords may not be as extensive as it would be desirable. Even though articles
542 usually employ the territory names and demonyms, there may exist specific cases which employ other
543 forms which our keywords do not capture. Even though the lack of these words may not imply missing
544 a significant number of articles, it would be necessary to obtain a systematic source for them, either a
545 database or a collaboratively created space in the same Wikipedia, especially when considering to
546 extend the method from 40 to all the 288 available language editions. Second, we are aware that using
547 the category crawling strategy in order to retrieve articles may not be as reliable when the
548 categorization is not exhaustive or precise. This is likely to be the main issue behind the lower accuracy
549 obtained for some languages such as Korean and Waray. In future developments of the method we plan
550 to introduce strategies that take into account Wikidata⁹, a rich complementary structured database
551 which contains a variety of properties and relationships between items. Other strategies to diminish
552 interference include using articles solidly included as CCC for another language as a negative ground-
553 truth. Machine learning approaches could also be used to improve accuracy.

⁹ wikidata.org

Imbalances Across 40 Language Editions

554 As a final general consideration, in this study we have pursued a quantitative approach, in order to be
 555 able to delimit and quantify the content specific to the context associated to each linguistic version of
 556 Wikipedia; as we deal with cultural contexts, boundaries are of course not always straightforward, and
 557 the manual assessment demonstrated that, although in few cases, even human experts may disagree on
 558 the definition of this task. While a comprehensive qualitative inspection of the results would help to
 559 more deeply grasp and interpret the findings of this work, such kind of analysis is out of the scope of
 560 this work and would be unviable at the scale of the whole dataset, which includes millions of articles
 561 in 40 different languages. However, we believe that making our method and the resulting datasets
 562 available has the potential to open up to more focused studies including qualitative methods and
 563 delving into specific aspects of the complex phenomenon of which we have here offered a first
 564 quantitative overview.

565 **Main Findings.** Our analyses offer new insights into how the cultural context impacts Wikipedia
 566 content coverage and imbalances across languages. In first place, the analysis of cultural context
 567 content in 40 language editions shows that its extent ranges from 7% to 49% of the total number of
 568 articles, with an average value of 23.53% (**RQ1**). This is a considerable extent, especially considering
 569 that CCC articles have generally been produced with no specific policy or guideline recommending it,
 570 but as an effect of editors' preferences. In second place, an analysis of the creation of CCC over time
 571 shows that this content grew constantly along with Wikipedia (**RQ2**). Even though certain relevant
 572 geographical places for the editors (cities, towns, rivers, etc.) can be finite, the degree of specificity
 573 that CCC can reach through very different topics implies that new content can continually appear. In
 574 third place the analysis based on ILLs in the 40 languages unveils and quantifies the culture gap: CCC
 575 articles are 4.15 times less shared between languages than the average content of each language edition
 576 (**RQ3**), and almost the half of CCC articles do not exist in any other language. This shows that the lack
 577 of correspondence of content between languages found by previous research (Warncke-Wang et al.
 578 2012; Hecht and Gergle 2010b) is due largely to CCC articles. The graphs provided to illustrate this
 579 culture gap can be useful to show editors from every language edition which cultural context content
 580 of other languages should be priorly imported or extended.

581 **Theoretical Implications.** Our study makes a contribution to the online communities and cultural
 582 contextualization literature. The imbalance of content across languages has been seen as a negative
 583 issue, since it hinders the goal of achieving "the sum of human knowledge", and is often explained by
 584 several demographic and territory factors. Some authors have demonstrated how editors tend to edit
 585 about geographically close territories (Hecht and Gergle 2010a), or how the overall article link
 586 structure in each Wikipedia language edition revolves around the countries where the language is
 587 spoken (Hecht and Gergle 2009). Proving also the centrality of geographical context, other authors
 588 have shown that editors' editing interests are similar when languages are geographically near (Karimi
 589 et al. 2015; Warncke-Wang et al. 2012; Saimolenko et al. 2016).

590 Differently from these studies, we wondered whether obtaining all the content associated with a
 591 language cultural context could explain in a more thorough way the impact of cultural context on the
 592 content coverage and imbalances across languages. In fact, the results from the interlanguage analysis
 593 of Cultural Context Content are in line with the results from previous literature that language

Imbalances Across 40 Language Editions

594 communities share common interests (Saimolenko et al. 2016). The fact that large part of the language
 595 gap is due to the CCC articles confirms that cultural and geographical context influences communities'
 596 common interests. At the same time, the same constitution of the CCC dataset and the filtering process
 597 showed us that the collection is a continuum, in which there are articles that are prominently at the core
 598 of the dataset, while others are less related to the collection central meanings and could even belong to
 599 other collections related to neighbour contexts. The CCC datasets allow for further investigations to
 600 analyze the content similarities between contexts and their relationships.

601 As a final consideration, it is important to note that these imbalances in content should not be only
 602 considered as a bias to be corrected, but also as a natural expression of cultural diversity, which
 603 represents a richness of the Wikipedia project. In this sense, our work can be seen as a first effort to
 604 quantify and describe such diversity, facing the delicate challenge of tracing the boundaries between
 605 cultural contexts.

606 **Implications for Practice.** As an important contribution of this paper, we make available both the
 607 code we used to process the Wikipedia language editions, the language-territories mapping with the
 608 keywords and ISO codes, as well as the generated datasets, in a website (*Wikipedia Cultural Diversity*
 609 *Observatory*¹⁰) aimed to both the academia and the Wikipedia communities. We believe this can
 610 encourage and motivate new research on cultural context content at the same time as it helps the
 611 Wikipedia language communities to bridge the culture gap. At the same time, we believe that
 612 increasing the coverage of each other's' cultural context content may be an important goal for fostering
 613 inter-cultural dialogue and enrichment. Therefore, we suggest that the translator and the article
 614 recommendation tool developed by the Wikimedia Foundation could include CCC from each language,
 615 or subparts of it (e.g. articles including keywords in their title) as preferential content to be translated
 616 and exported to other languages. Making editors more aware of the culture gap and offering them tools
 617 to bridge it will likely encourage them to enrich their Wikipedia language editions with inter-cultural
 618 content, enlightening their readers and helping to build a world more open to diversity.

619 **Future Lines of Research.** The datasets and methods proposed in this paper suggest several lines for
 620 future research. On the one hand, it would be interesting to investigate the overlap between CCC and
 621 other relevant groups of articles for a particular reason or metric, e.g. the ones receiving higher attention
 622 (in terms of edits and page views), or related to breaking news or to controversial topics. The CCC
 623 dataset can be also used as a basis for cross-cultural studies in the field of Digital Humanities. On the
 624 other hand, analogous strategies to the ones proposed here to collect the CCC datasets could be applied
 625 to identify other kinds of content, and therefore are relevant to studies aimed at providing a topical
 626 coverage of Wikipedia or other knowledge repositories.

¹⁰ [https://meta.wikimedia.org/wiki/Grants:Project/Wikipedia_Cultural_Diversity_Observatory_\(WCDO\)](https://meta.wikimedia.org/wiki/Grants:Project/Wikipedia_Cultural_Diversity_Observatory_(WCDO))

Imbalances Across 40 Language Editions**627 5. Conflict of Interest**

628 *The authors declare that the research was conducted in the absence of any commercial or financial*
 629 *relationships that could be construed as a potential conflict of interest.*

630 6. Funding

631 This work has been partially funded by a Project Grant from Wikimedia Foundation.

632 7. Acknowledgments

633 We want to thank Andreas Kaltenbrunner for his wise critical comments that helped to strengthen the
 634 methodology, and Laura Vincze for her precious help to improve style and clarity of the manuscript.
 635 We would also like to acknowledge all those who supported the study in a direct or indirect way,
 636 including the Wikimedia Foundation and all the Wikimedians who endorsed the Wikipedia Cultural
 637 Diversity Observatory project. We hope the methods and datasets presented here as part of this open
 638 source project may be of help for all who are working hard to improve Wikipedia's cultural diversity
 639 coverage.

640 8. References

- 641 Apic, G., Betts, M. J., & Russell, R. B. (2011). Content disputes in Wikipedia reflect geopolitical
 642 instability. *PloS One*, 6(6). doi:10.1371/journal.pone.0020902.g001
- 643 Aragón, P., Laniado, D., Kaltenbrunner, A., & Volkovich, Y. (2012). Biographical social networks on
 644 Wikipedia: a cross-cultural study of links that made history. (p. 19). Presented at the WikiSym '09:
 645 Proceedings of the 5th International Symposium on Wikis and Open Collaboration, New York,
 646 New York, USA: ACM Press. doi:10.1145/2462932.2462958
- 647 Bao, P., Hecht, B., Carton, S., Quaderi, M., Horn, M., & Gergle, D. (2012, May). Omnipedia: bridging
 648 the wikipedia language gap. In Proceedings of the SIGCHI Conference on Human Factors in
 649 Computing Systems (pp. 1075-1084). ACM.
- 650 Blondel, V. D., Guillaume, J. L., Lambiotte, R., & Lefebvre, E. (2008). Fast unfolding of communities
 651 in large networks. *Journal of statistical mechanics: theory and experiment*, 10, P10008.
- 652 Callahan, E. S., & Herring, S. C. (2011). Cultural Bias in Wikipedia Content on Famous Persons.
 653 *Journal of the American Society for Information Science and Technology*, 62:10, 1899–1915.
 654 doi:10.1002/asi.21577
- 655 Clark, H. (1996). Using Language. Cambridge University Press.
- 656 Cohen, J. (2016). A Coefficient of Agreement for Nominal Scales. *Educational and Psychological
 657 Measurement*, 20:1, 37–46. doi:10.1177/001316446002000104
- 658 Ensslin, A. (2011). “What an un-wiki way of doing things”: Wikipedia’s multilingual policy and
 659 metalinguistic practice. *Journal of Language and Politics*, 10:4, 535–561.
 660 doi:10.1075/jlp.10.4.04ens
- 661 Eom, Y.-H., Aragón, P., Laniado, D., Kaltenbrunner, A., Vigna, S., & Shepelyansky, D. L. (2015).
 662 Interactions of cultures and top people of Wikipedia from ranking of 24 language editions. *Plos
 663 One*, 10:3, e0114825. doi:10.1371/journal.pone.0114825
- 664 Hecht, B. J. (2013). *The Mining and Application of Diverse Cultural Perspectives in User-Generated
 665 Content*. Dissertation. Northwestern University. United States.

Imbalances Across 40 Language Editions

- 666 Hecht, B. J., & Gergle, D. (2010a). On the localness of user-generated content (pp. 229–232). Presented
 667 at the CSCW '10: Proceedings of the 2010 ACM conference on Computer supported cooperative
 668 work. doi:10.1145/1718918.1718962
- 669 Hecht, B., & Gergle, D. (2009). Measuring self-focus bias in community-maintained knowledge
 670 repositories (pp. 11–20). Presented at the C&T '09: Proceedings of the fourth international
 671 conference on Communities and technologies. doi:10.1145/1556460.1556463
- 672 Hecht, B., & Gergle, D. (2010b). The tower of Babel meets web 2.0: user-generated content and its
 673 applications in a multilingual context (pp. 291–300). Presented at the CHI '10: Proceedings of the
 674 SIGCHI Conference on Human Factors in Computing Systems, New York, New York,
 675 USA: ACM Request Permissions. doi:10.1145/1753326.1753370
- 676 Karimi, F., Bohlin, L., Samoilenco, A., Rosvall, M., & Lancichinetti, A. (2015). Quantifying national
 677 information interests using the activity of Wikipedia editors. *ArXiv Abs/1312.0976*, 1503, 5522.
- 678 Massa, P., & Scrinzi, F. (2011). Exploring linguistic points of view of Wikipedia. (pp. 213–214).
 679 Presented at the WikiSym '09: Proceedings of the 5th International Symposium on Wikis and Open
 680 Collaboration, New York, New York, USA: ACM Press. doi:10.1145/2038558.2038599
- 681 Pentzold, C., Weltevrede, E., Mauri, M., Laniado, D., Kaltenbrunner, A., & Borra, E. (2017). Digging
 682 Wikipedia. *Journal on Computing and Cultural Heritage*, 10(1), 1–19. doi:10.1145/3012285
- 683 Rogers, R., & Sendijarevic, E. (2012). Neutral or National Point of View? A Comparison of Srebrenica
 684 articles across Wikipedia's language versions. Proceedings of the *Wikipedia Academy Conference*
 685 2012, Berlin.
- 686 Ronen, S., Gonçalves, B., Hu, K. Z., Vespignani, A., Pinker, S., & Hidalgo, C. A. (2014). Links that
 687 speak: The global language network and its association with global fame. Proceedings of the
 688 National Academy of Sciences, 111(52), E5616-E5622.
- 689 Samoilenco, A., Karimi, F., Edler, D., Kunegis, J., & Strohmaier, M. (2016). Linguistic
 690 neighbourhoods: explaining cultural borders on Wikipedia through multilingual co-editing
 691 activity. EPJ data science, 5(1), 9.
- 692 Van Dijk, Z. (2009). Wikipedia and lesser-resourced languages. *Language Problems & Language*
 693 *Planning*, (3), 33. doi:10.1075/lplp.33.3.03van
- 694 Warncke-Wang, M., Uduwage, A., Dong, Z., & Riedl, J. (2012). In search of the ur-Wikipedia:
 695 universality, similarity, and translation in the Wikipedia inter-language link network. *OpenSym*
 696 '12: *Proceedings of the Eighth Annual International Symposium on Wikis and Open Collaboration*,
 697 20. doi:10.1145/2462932.2462959
- 698

699 9. Data Availability Statement

700 The datasets generated/analyzed for this study can be found in the ‘Wikipedia Cultural Diversity
 701 Observatory Repository’: <https://github.com/marcmiquel/WCDO>

702 10. Tables and Figures

703
 704
 705
 706
 707
 708
 709
 710
 711
 712

Imbalances Across 40 Language Editions

713
714
715
716
717
718719 **Table 1. Inter-rater reliability tests for the Japanese and German Wikipedia language editions.** For each
720 Wikipedia we crossed the ratings (CCC) from three raters. Coincidence is the degree of coincidence in %, and K is
721 the Cohen's Kappa coefficient.

Inter-rater reliability	Japanese		German	
	coincidence	K	coincidence	K
algorithm-rater1	0.86	0.71	0.90	0.8
algorithm-rater2	0.89	0.77	0.91	0.82
algorithm-rater3	0.86	0.72	0.89	0.77
rater1-rater2	0.97	0.94	0.96	0.93
rater1-rater3	0.97	0.93	0.95	0.9
rater2-rater3	0.96	0.91	0.98	0.95

722

Table 2. Percentage of CCC articles in Wikipedia language editions and CCC cross-language coverage. For each of the 40 editions considered, columns report: total number of articles (WP art); percentage of CCC articles over the entire Wikipedia (CCC %), and percentage of articles identified through the first strategy - GeoLocated tags (GL %) and through the second strategy - KeyWords in their titles (KW %); percentage of False Positives (FP %) and False Negatives (FN %), with resulting F1-score (F1) after manual evaluation; average number of Interlanguage links per article in the language edition (ILL WP) and in CCC (ILL CCC), percentage of CCC articles having no ILLs, and percentage of CCC articles having no ILLs with respect to WP articles having no ILLs (CCC NO ILL / WP NO ILL).

ISO cod.	Language	WP Art.	CCC %	GL %	KW %	FP %	FN %	F1	Avg. ILL WP	Avg. ILL CCC	CCC NO ILL %	CCC NO ILL / WP NO ILL
af	Afrikaans	35966	19.2	5.9	0.9	1	2	0.99	40.1	4.5	34.3	76.2
ar	Arabic	375282	26.9	3.2	2.4	2	18	0.91	12.9	3.6	59.5	54.5
eu	Basque	208630	10.1	1.7	0.4	4	1	0.97	14.4	1.3	50.6	73.1
ca	Catalan	467486	16.2	7.9	0.8	2	3	0.98	21.5	3.6	68.7	62.7
ceb	Cebuano	1211531	0.1	0.0	0.1	12	1	0.93	15	1.6	0.6	0.1
zh	Chinese	851670	32.9	6.3	1.2	10	11	0.90	6.3	11	58.2	63.4
cs	Czech	326187	25.9	9	1.2	2	3	0.98	4.8	8.9	60.3	71
da	Danish	205764	31.7	6.1	1.0	10	2	0.94	10.0	2.6	52.3	73.4
nl	Dutch	1828148	7.8	1.6	0.3	1	3	0.98	13	1.8	64.4	22.4
en	English	4917741	46.8	9.8	2.8	10	16	0.87	6.8	1.5	55	63.1
et	Estonian	136362	31.1	6.1	1.7	2	2	0.98	20.2	1.8	64.4	69.8
fi	Finnish	375347	21.9	2.3	1	1	4	0.98	6	2.9	70.2	70
fr	French	1642276	29.0	6.9	1.7	10	6	0.92	23.2	4.7	46.2	59.3
de	German	1834147	36.8	8.8	1.9	10	10	0.90	15	2.5	60.1	62.5
el	Greek	108090	33.5	6.4	0.6	9	8	0.91	17.9	4.2	46.1	71.9
gn	Guarani	3031	23.6	14	3.3	2	6	0.96	82.1	24.2	6.9	57.6
he	Hebrew	174667	31.7	2.1	1.6	14	2	0.91	20	4.8	50.3	79.5
hu	Hungarian	326146	18.5	1.9	1.5	10	4	0.93	16	2.9	54.8	61
is	Icelandic	39554	30.7	2.2	1.5	2	10	0.94	12	1.7	66.1	74.2
id	Indonesian	363529	27	1	0.6	7	4	0.94	33.7	2.4	36.7	73.8
it	Italian	1210801	19.2	3.6	0.7	10	5	0.92	9.3	3.5	54.5	48.4
ja	Japanese	973955	49.2	3.4	1	4	10	0.93	7.1	1.2	75.9	77.5
ko	Korean	320742	32.6	2.4	0.8	19	12	0.84	14.1	7.8	69.9	58.6
mk	Macedonian	82743	15.9	2.5	1.3	15	4	0.90	25.3	3.4	41	51.6
ms	Malay	275031	19.5	1.4	0.8	15	2	0.91	15.5	1.8	55.4	61.5
ne	Nepali	29114	29.7	11.8	2.2	2	19	0.90	22	3.3	41.4	29.5
no	Norwegian	415015	26.8	5.5	0.8	12	6	0.91	12.4	2.3	54.3	72.5
fa	Persian	460523	11	10.3	0.7	3	19	0.90	7.6	4.8	8.2	4.7
pl	Polish	1122218	23.2	9.4	1.1	9	3	0.94	9.4	1.3	58.1	54.2
pt	Portuguese	880529	19.1	2	1	6	1	0.96	11.2	2.4	64.5	58.4
ro	Romanian	329925	20.7	7.2	1.1	13	2	0.92	16.9	3.5	32.6	72.7
ru	Russian	1237127	31.2	11	1.1	14	5	0.90	8.3	2.2	44.6	56.7
sr	Serbian	321912	12.1	3.2	0.1	23	1	0.87	16	4.7	25	50.3
es	Spanish	1147742	27.7	5	2	13	6	0.90	9.3	3.4	44.3	64.9
sw	Swahili	29168	18.3	3.6	1	1	6	0.97	40	3.7	46.8	73.8
sv	Swedish	1970808	11.4	4.3	0.4	8	4	0.94	6	1.4	72.5	66
tr	Turkish	249061	33.9	4.4	2.1	6	4	0.95	16.2	3.4	36.4	70
uk	Ukrainian	581735	24.8	6.8	1	14	2	0.91	13	2.4	43.1	57
vi	Vietnamese	1137180	2.5	0.9	0.2	5	0	0.97	7.4	1.5	72.8	17
war	Waray	1259278	0.1	0.0	0.0	23	0	0.87	6.3	10.9	12.6	1.9
Avg.	Average	736654	23.3	5.1	1.1	8.1	5.9	0.92	16.6	4	49	57.1

Imbalances Across 40 Language Editions

3 **Figure 2.** Examples of articles from English Wikipedia. CCC keywords (English Literature), CCC
 4 geolocated (Times Square), the rest of CCC (Banbury cake) and the rest of Wikipedia (Sun).

5 **Figure 3.** Average proportion of CCC, and of CCC detected through geolocation and keywords.
 6 Sizes are in scale according to their proportion.

7 **Figure 4.** CCC creation over the 15 years of Wikipedia. For each language edition, the green area
 8 represents the absolute number of CCC articles created over years, and the grey area the rest of the
 9 articles created. The line shows the percentage of CCC over the total number of articles created
 10 during each year; it is depicted in grey when it is in line (less than 10% variation) with the final
 11 overall percentage of CCC in the encyclopaedia, in green or red when it is higher or lower.

12 **Figure 5.** Network graph of language proximity in terms of shared CCC articles. Each node
 13 represents a Wikipedia language edition, and has three outgoing links to the three language editions
 14 in which its CCC represents the highest percentage of the articles. Links are represented in clockwise
 15 direction. Colors represent clusters of language editions identified through the Louvain algorithm for
 16 community detection.

17
 18 **Figure 6.** Culture gap: 40 Wikipedia language editions coverage (% articles) of 40 Wikipedia
 19 language editions CCC. Each row shows the coverage of other language editions' CCC. The
 20 coverage is calculated as the percentage of a Wikipedia language edition CCC (column) which exists
 21 in another Wikipedia language edition (row). For an easy identification of values, cells are colored
 22 in a scale of colors from red (to indicate a percentage lower than 1%), to green in a continuum until
 23 93.67% (the highest value).