Highly Parallel Texts Enriched with Highly Useful Metadata? A Wikipedia Case-Study Combining Machine Learning and Social Technology

Ahmad Aghaebrahimian, Andy Stauder and Michael Ustaszewski Department of Translation Studies, University of Innsbruck, Innsbruck, Austria

This is an un-refereed preprint version of an article accepted for publication in *Digital Scholarship in the Humanities*. The refereed, revised and copyedited version of record "Ahmad Aghaebrahimian, Andy Stauder, Michael Ustaszewski: Automatically extracted parallel corpora enriched with highly useful metadata? A Wikipedia case study combining machine learning and social technology. In: *Digtal Scholarship in the Humanities*, DOI: fqaa002" is available online at: https://doi.org/10.1093/llc/fqaa002

Full correspondence details: Dr. Michael Ustaszewski Universität Innsbruck Institut für Translationswissenschaft Herzog-Siegmund-Ufer 15 A-6020 Innsbruck | Austria Telefon +43 512 507 42482 E-Mail michael.ustaszewski@uibk.ac.at Web www.uibk.ac.at/translation

1 Highly Parallel Texts Enriched with Highly Useful Metadata? A Wikipedia

2 Case-Study Combining Machine Learning and Social Technology

3

4 Abstract

The extraction of large amounts of multilingual parallel text from web resources is a widely used technique in natural language processing. However, automatically collected parallel texts usually lack precise metadata, which are crucial to accurate data analysis and interpretation. The combination of automated extraction procedures and manual metadata enrichment may help address this issue. Wikjipedia is a promising candidate for the exploration of the potential of said combination of methods, because it is a rich source of translations in a large number of language pairs and because its open and collaborative nature makes it possible to identify and contact the users who produce translations.

This article tests to what extent translated texts automatically extracted from Wikipedia by means of neural networks can be enriched with pertinent metadata through a self-submission-based user survey. Special emphasis is placed on data usefulness, defined in terms of a catalogue of previously established assessment criteria, most prominently metadata quality. The results suggest that from a quantitative perspective the proposed methodology is capable of capturing metadata otherwise not available. At the same time, the crowd-based collection of data and metadata may face important technical and social limitations.

19

20 **1** Introduction

Parallel corpora are an indispensable resource for a wide range of language-related research and engineering problems. The availability of parallel corpora is limited, however, especially for less-common language pairs or certain subject domains. The extension of existing and compilation of new parallel corpora is therefore a prerequisite to fully reaping the benefits of the increasingly sophisticated and powerful data-driven approaches, be it in the humanities or in natural language processing (NLP). 26 Spurred by a high demand for multilingual training data, most notably in the field of machine transla-27 tion, the automatic extraction of parallel text from existing multilingual websites has received consid-28 erable attention in NLP as a more time and cost-efficient alternative to the manual compilation of 29 parallel corpora from the ground up. The extracted data are mainly used for domain adaptation of NLP 30 tools and systems and to alleviate data bottlenecks faced by less-resourced languages. While auto-31 matic harvesting procedures usually yield large amounts of parallel data, they lack metadata that pre-32 cisely describe the extracted parallel texts or text fragments, mainly because the mined websites do 33 not make such metadata available. Metadata are a key component of knowledge discovery, prediction, 34 and decision making based on (big) data analytics (Greenberg, 2017). Metadata are vital for accurate 35 data interpretation by both humans and machines (ibid.), which was also confirmed from the perspec-36 tive of digital humanities in an analysis of the importance and, conversely, the impact of a lack of 37 metadata, highlighting that corpus size cannot make up for insufficient or inexistent metadata 38 (Koplenig, 2017).

39 The present paper addresses the lack of metadata in automatically harvested data. Based on the 40 use case of Wikipedia, we test to what extent translated texts automatically extracted using machine 41 learning can be enriched with metadata of high usefulness, collected with the help of social technol-42 ogy. The proposed methodology leverages the wealth of openly available administrative metadata in 43 a collaboratively built knowledge resource in combination with human intervention to obtain value-44 added data: highly parallel text data with pertinent highly useful metadata. The extracted data are 45 analyzed as to whether they are compliant with the principles of what Greenberg (2017) calls smart 46 metadata (see Section 3). The choice of using Wikipedia as a data source is motivated by its size, mul-47 tilinguality, diverse topics, openness, dynamic growth, and transparency. It contains large amounts of 48 CC-licensed material, unproblematic in terms of rights and data protection issues, and it enables inter-49 action with data producers. Yet, Wikipedia translations have largely been neglected by research (Jones, 50 2018; Shuttleworth, 2018), at least in fields other than NLP (see Section 2). This may be due to the non-51 trivial assessment of the extent and exact location of translated material in Wikipedia, as well as to the

52 closely related difficulty of pinpointing fully parallel text pairs among interlanguage-linked articles that 53 may evolve independently from each other over time (Shuttleworth, 2018). The present study ad-54 dresses these difficulties in order to tap into Wikipedia as a multilingual translation resource. Unlike 55 many automatic harvesting approaches documented in the literature, this study aims to extract fully 56 parallel source-target text pairs at the document level rather than isolated sentence pairs. The meth-57 odology was tested within the scope of a collaborative, open-access, open-ended, open-domain corpus-building initiative in the area of translation that aims to make high-quality data freely available for 58 59 reuse (Ustaszewski and Stauder, 2017), but the proposed combination of automatic extraction and 60 enrichment by humans might be of relevance to other mono- and multilingual data collection projects 61 as well.

62 2 Related Work

63 Multilingual websites hold great potential for parallel data harvesting, but it must be borne in mind 64 that multilingual content can exhibit varying degrees of parallelity across languages. The following par-65 allelity levels can be distinguished (Babych *et al.*, 2019; Sharoff *et al.*, 2013a):

parallel: translations, i.e. source-target text pairs that can be aligned at sentence and phrase level
 strongly comparable: texts that can be aligned at the document level, such as heavily edited trans lations and their source texts, or independently produced but closely related texts about the same
 subject

weakly comparable: texts from the same domain or genre but about different sub-topics; collec tions of such texts can usually only be aligned at sub-corpus rather than document level

• non-comparable, or unrelated: texts that cannot be aligned across languages

A slightly different typology was suggested by Fung and Cheung (2004) who distinguish parallel, noisy parallel, comparable, and very-non-parallel corpora. The comparable corpus type widely used in contrastive and translation studies (Zanettin, 2012) mostly comprises *weakly comparable* materials according to the above typology. The object of this study, Wikipedia, makes multilingual articles on the 77 same topic easily accessible through so called *interlanguage links*¹. The degree of parallelity of inter-78 language-linked articles may vary widely (Babych et al., 2019), and the exact amount of parallel and 79 comparable data in Wikipedia is unknown (O'Hagan, 2016). In a manual evaluation of 200 randomly 80 sampled French-English article pairs from 2009, the proportions of parallel, noisy parallel, comparable, 81 and very non-parallel texts were 14%, 11%, 29%, and 46%, respectively (Patry and Langlais, 2011), but 82 the authors pointed out the difficulty of drawing clear-cut boundaries between the four degrees of 83 comparability. In a more recent human evaluation study with eight language pairs, including seven 84 less-resourced ones, 52.5% of interlanguage-linked document pairs were judged to be highly similar 85 and 18.8% to be different on a five-point Likert scale, confirming that Wikipedia articles in different 86 languages on the same topic are not necessarily (very) similar (Babych et al., 2019). The study also 87 elicited factors influencing the perceived parallelity of document pairs, indicating that pairs judged to 88 be similar have similar structure, overlapping named entities, alignable fragments and contain trans-89 lation equivalents. Altogether, these quantitative figures seem to confirm that while there are coordi-90 nated Wikipedia translation initiatives, translation is mainly self-motivated and only one of several 91 mechanisms for the multilingual expansion of Wikipedia (Shuttleworth, 2018). Moreover, Wikipedia 92 content is dynamic, which means that articles between which a translation relation obtained at a cer-93 tain point in time may evolve independently from each other. Nevertheless, the sheer size of Wikipedia 94 makes it a promising and valuable source of parallel text data, as will be discussed below.

Since comparable multilingual data is much more abundant on the web than strictly parallel data, the compilation and exploitation of comparable corpora as a means to compensate for the lack of parallel corpora has become a prolific line of research (Sharoff *et al.*, 2013b; Skadiņa *et al.*, 2019). Along this line, the identification and extraction of parallel sentences from comparable corpora has played an important role, most notably from Wikipedia and with the purpose of improving the performance of statistical machine translation systems (e.g. Barrón-Cedeño *et al.*, 2015; Labaka *et al.*, 2016; Smith

¹ <u>https://en.wikipedia.org/w/index.php?title=Help:Interlanguage_links&oldid=885377785</u>

101 et al., 2010). While most approaches are based on comparable corpora, Fung and Cheung (2004) ex-102 ploited non-comparable corpora for this task. Ture and Lin (2012), on the other hand, aim to maximize 103 the recall of extracted sentences; their system determines the degree of comparability between inter-104 language-linked Wikipedia article pairs using a more exhaustive and computationally more costly ap-105 proach instead of relying on heuristics only, thus yielding 5.8 million English-German sentence pairs. 106 More recently, neural networks have been used for parallel sentence extraction (e.g. Chu et al., 2016; 107 España-Bonet et al., 2017; Grégoire and Langlais, 2018) – an approach that has been adopted in the 108 present study, too (see Section 4). Web resources other than Wikipedia are also being used for parallel 109 sentence extraction, for instance in a system that crawls the web to detect entry points to multilingual 110 websites and subsequently uses intra-site crawlers and alignment procedures (Barbosa et al., 2012).

111 However, there is also a good deal of research that focuses on identifying text data that exhibit high 112 degrees of parallelity at the document level, as opposed to extracting isolated sentence pairs. Such 113 approaches consist in finding interlingually corresponding text pairs from bi- or multilingual collections 114 of candidate documents, for example Wikipedia (e.g. Enright and Kondrak, 2007; Etchegoyhen and Azpeitia, 2016; Mohammadi, 2016; Morin et al., 2015; Patry and Langlais, 2011). They offer algorithmic 115 116 alternatives to the manual identification of translated Wikipedia articles based on the information con-117 tained in Wikipedia's administrative pages and the structural mark-up assigned by editors to internally 118 organize the content of the encyclopaedia (Shuttleworth, 2018). As mentioned in Section 1, the pre-119 sent study is in the vein of these works that aim at extracting fully parallel source-target text pairs, or 120 bitexts, at the document level. However, it also aims to enrich the extracted bitexts with metadata 121 provided by the producers of the respective texts.

122 **3 Data Usefulness**

The goal of the present study is to obtain parallel text data from Wikipedia that fulfil two criteria of usefulness. First, the extracted bitexts need to have the highest degree of interlingual document-level parallelity according to the typologies reviewed in Section 2. This means that text pairs are to share the same content, which, for the sake of practicality, is operationalized in terms of alignability at the 127 sentence or phrase levels (Babych et al., 2019). For the purpose of this study, sentence alignment does 128 not need to be in a bijective (i.e. one-to-one), monotonic (i.e. non-crossing) relation, which is the case when both sides of a bitext are structured strictly in the same order such that the n^{th} segment of the 129 source side corresponds to the *n*th segment on the target side (Tiedemann, 2011). Thus, cross-align-130 131 ments (n^{th} source segment corresponds to m^{th} target segment, $n \neq m$) as well as one-to-many, many-132 to-one and many-to-many alignments (n source segments correspond to m target segments, $n \neq m$) 133 are permitted. The decision to loosen the bijectivity and monotonicity constraints mirrors the fact that 134 (conscious) alterations of text segmentation (e.g. splitting one long source sentence into several 135 shorter target sentences) or structure (e.g. changing the order of semantic elements) are common 136 types of so-called shifts (Gambier, 2010) in real-world translations. Similarly, as a result of translation 137 shifts a given piece of information may have no counterpart in either the target text (omission of 138 source text information) or source text (addition of explanatory information in the target text). There-139 fore, the presence of *n*-to-zero (omission) or zero-to-*n* (addition) relations is permitted in the extracted 140 bitexts, provided the proportion of sentence pairs exhibiting such relations does not exceed an arbi-141 trarily chosen threshold of 5%.

142 The second criterion is concerned with the usefulness of the metadata with which the extracted 143 bitexts are enriched. Greenberg (2017) defines five principles of smart metadata that help provide 144 context and meaning for data: good quality, accessibility, trust, actionability, and preservation. The 145 collection of metadata within the present study aims to comply with the first three principles, whereas 146 the remaining two principles are of more concern to the storage and (re-)use of the collected data and 147 metadata. The following paragraphs describe how these three principles are reflected in this study's 148 metadata collection, which complements the automatic extraction of parallel texts from Wikipedia by 149 means of a dedicated self-submission web interface that invites the producers of their respective texts 150 to interactively contribute pertinent metadata (see Section 4.2).

151 *Principle 1: Good quality*

Based on Bruce and Hillmann (2004), Greenberg (2017) lists five indicators of good quality metadata:
accuracy, completeness, conformance to expectations, logical consistency and coherence, and timeliness. Each of these is discussed in the following.

To ensure accuracy of the collected metadata, metadata contributors are guided through the submission process by means of a lean, easy-to-use web interface. Instructions and tooltips in plain English are available throughout the interface. Instead of free text fields, dropdowns – each one with few, clear-cut choices – are used in the interface for the sake of consistency and manageability. In addition, trained staff are to verify and curate the submitted metadata.

160 Completeness is to be ensured with the help of a set of metadata labels that capture the relation 161 between translated texts and their originals, including the circumstances under which they were pro-162 duced. The label set strikes a middle ground between maximum detail and economic feasibility. The 163 label set used in the web interface is limited to text-extrinsic information known to data producers 164 only (e.g. translator's age at time of writing) and kept to a manageable size to minimize submission 165 effort and maximize completion. Text-intrinsic metadata (e.g. subject matter) are to be complemented 166 by trained project staff at a later stage in the project, of which this study is a part. To mitigate the 167 legally motivated lack of mandatory fields, the interface at least encourages contributors to make com-168 plete metadata submissions with visual cues.

To meet the criterion of conformance to (user) expectations, the set of translation-specific metadata labels is, in part, a combination of labels used by existing translation corpora. Translation corpora are, as a rule, mostly narrow in focus, which makes necessary the aforementioned combination in order to cater to a wide variety of translation audiences and stakeholders (Ustaszewski and Stauder, 2017). So, the label set aims to make findable what a large community would reasonably *expect* to find in a translation corpus (c.f. the FAIR Data Principles, Wilkinson *et al.*, 2016).

The criterion of logical consistency and coherence has two facets: an internal one, which in our case boils down to the use of controlled vocabulary for the respective metadata fields, and an external one, which results from the fact that the developed metadata label set combines features from a number of well-established repositories in the field of translation. Those postulating the latter facet (Bruce and Hillmann, 2004) seem to perceive communities (such as the translation community) as large systems which should use common standards, making their mode of operation coherent and consistent. The internal facet – keeping the labels used by *one* corpus internally consistent and coherent – can reasonably be assumed to be more easy to achieve due to the smaller number of people involved. Whether the latter facet can truly be regarded as coherence and consistency, or rather interdependence of several systems, depends on what one views as a whole.

Lastly, the criterion of timeliness consists, on the one hand, in having up-to-date metadata and, on the other, not publishing data that haven't been labelled with meaningful metadata yet. Due to the nature of the present study, which collects data of objects that cannot change (i.e. one-time snapshots of Wikipedia articles), the currency aspect cannot be applied here. As for the second criterion: publishing only data that already have metadata, this is to be satisfied by not making available the data as long as the pertinent metadata have not been collected, verified and complemented by trained staff.

191 Principle 2: Accessibility

On the technical side, the collected metadata are to be stored alongside the extracted parallel texts and made available through a still-under-construction faceted search interface which enables users to compile and download corpora on-demand, tailored to their specific needs. On the legal side, both the extracted texts and collected metadata are to be made freely available under either the CC-BY-SA license or a permissive license specially drafted for the purpose of this study by experts on legal aspects in the field of language data. The required permissions are obtained from the (meta-)data contributors through the metadata collection interface, where they can select licenses as part of the specially drafted contributor agreement² and specify the desired degree of anonymity, which ensures conform ity with the repository's privacy policy³. Hence, legal clearance and transparency as a prerequisite to
 accessibility is of paramount importance to metadata usefulness, especially in the context of open access corpora.

203 Principle 3: Trust

Good quality metadata are trusted metadata and produced by reliable sources (Greenberg, 2017). In the present study, this means that the text producers themselves, i.e. translators of Wikipedia articles, contribute metadata through a self-submission interface. The underlying assumption is that the volunteer community of Wikipedians is open to the idea of sharing (meta-)data and contributing to a collaborative open-science initiative. No less important, collected metadata are to undergo revision by trained project staff prior to integration into the open-source repository.

210 **4 Method**

211 4.1 Automated Parallel Data Extraction

The extraction of bitexts from Wikipedia requires locating candidate text pairs, which was accomplished with the help of Wikipedia interlanguage links, which link corresponding articles on the same topic across languages, e.g. the English article on 'water' with the German article on 'Wasser'. When someone creates a version of an article in a different language, it is common practice to create an interlanguage link between the original article and the new text in the other language, by embedding a so-called *Interlanguage link template*⁴ in the latter (Wikipedia 2019). All of the interlanguage links

² <u>https://transbank.info/contributor-agreement</u> [anonymized for review]

³ <u>https://transbank.info/privacy</u> [anonymized for review]

⁴ <u>https://en.wikipedia.org/wiki/Template:Interlanguage_link</u>

are stored in a dedicated SQL database⁵, exported by Wikipedia at regular intervals. On Wikipedia, the interlanguage links appear in a sidebar on the left. We downloaded the SQL database from March 2018 to extract all interlanguage links. The links available in the database are unique identifiers, each pointing to two associated article revisions in different languages. As an example and to provide a better idea of the numbers involved, **Table 1** lists the respective numbers of articles in the top 13 languages that have been linked to respective English versions.

224 Table 1 Number of interlanguage-linked Wikipedia article pairs for English (March 2018)

Language Pair	# interlanguage-linked comparable article pairs
English-French	1,491,578
English-German	1,247,102
English-Italian	1,123,058
English-Spanish	1,096,328
English-Russian	975,983
English-Swedish	918,314
English-Dutch	906,950
English-Polish	906,105
English-Portuguese	900,508
English-Farsi	866,408
English-Chinese	703,217
English-Arabic	643,360
English-Japanese	631,855

225

226 The article revisions (=versions) to which the links in the aforementioned SQL database point are 227 contained in database dumps for each language version of Wikipedia that are made available twice a 228 year. These database dumps are the data source in which we searched for text pairs that constitute 229 respective originals and translations. This is done by identifying pairs of article revisions with at least 230 95% of parallel sentences, as was described in Section 3. We used Wikipedia dump files (June 2018) of 231 languages contained in the total set of language pairs from the SQL database. The text pairs could also be retrieved through the MediaWiki Action API⁶, but for reasons of performance and ease of use we 232 233 chose to download Wikipedia dumps and process them directly. Retrieving text pairs from data dumps

⁵ <u>https://dumps.wikimedia.org/wikidatawiki/latest/wikidatawiki-latest-langlinks.sql.gz</u>

⁶ <u>https://www.mediawiki.org/wiki/API</u>

makes it easier to browse different revisions of the same article, which is an essential step in filtering parallel data from comparable data. The reason for this is that due to the dynamic and collaborative nature of Wikipedia, individual articles may undergo editing by a number of different users. Therefore, although the interlanguage links provide information as to which pages are comparable, the degree of comparability of these pages may change over time due to changes made independently from other language versions.

240 The next step was retrieving parallel texts from the set of texts that had been identified as compa-241 rable with the help of interlanguage links. For this purpose we developed a deep neural network ar-242 chitecture to identify fully parallel texts among the candidate pairs. This approach offers a language-243 independent, robust and highly scalable state-of-the-art solution to parallel sentence extraction 244 (Aghaebrahimian, 2018). As has been mentioned in Section 3, the criterion for identifying parallel texts 245 was that they contained at least 95% of sentences that had been identified as parallel. The threshold 246 used by the neural network methodology for deciding whether two sentences were parallel or not was 247 determined with the help of human evaluation: human raters were presented sentence pairs from the 248 Europarl corpus (Koehn 2005) and had to decide whether the two sentences of each pair were parallel 249 or not (Aghaebrahimian, 2018).

250 4.2 Metadata Enrichment

251 4.2.1 Identification of Translators

For each extracted Wikipedia translation, we identified the user who had written it with the help of the pointers from the SQL database article's revision history. In this way we could retrieve the name of each user who has used an interlanguage link template⁷ for linking a revision written by him/her to the corresponding Wikipedia page in another language. Material of users that had the email feature of Wikipedia turned off, and thus could not be contacted, was left out. Also, "translations" that were

⁷ <u>https://en.wikipedia.org/wiki/Template:Interlanguage_link</u>

257 obviously unusable, e.g. because they consisted only of a title and had no text body⁸, were filtered out. 258 Each one of the users has a unique identifier through which they can be contacted via Wikipedia's 259 internal mailing system, as long as they have activated this feature for their respective Wikipedia ac-260 count. When users create a translation of a Wikipedia article, and if they follow Wikipedia's recom-261 mendation, they add an interlanguage link template containing their own username to the newly cre-262 ated revision. If someone else decides to modify this translation or add another translation for the 263 same source page, another revision is created which contains another interlanguage link template. So 264 in this way extracting the name of the translators is deterministically feasible.

265 Users identified as translators of the extracted Wikipedia articles were contacted manually via Wikipe-266 dia's built-in email facility in order to invite them to provide metadata about themselves and their 267 translations using an on-line interface designed by us for the purpose (see Section 4.2.2). We compiled 268 a list of translators' user IDs, taken from the interlanguage-link SQL database mentioned in Section 4.1. 269 The list was compiled as a one-to-many list, meaning that it maps each user ID to the titles of the 270 revisions constituting their translations. To make sure that the users could only access their own texts, 271 we generated user-specific log-in names and passwords using which they could log into their person-272 alized page in our interface.

273 Contacting the users via email is a time-consuming process since the Wikipedia email system does 274 not allow sending more than ten emails per 24 hours per IP address to prevent spamming through the 275 system. Moreover, as has been mentioned, sending emails is only possible to users who activated this 276 for their account. Due to all these limitations, the invitation emails containing the credentials for the 277 metadata submission interface were sent manually from five Wikipedia accounts. 934 Wikipedia users 278 were contacted in this way. A more detailed account of the process follows in the next subsection.

⁸ Some users apparently create interlanguage-linked pages and only translate the respective title of an article, maybe to get the translation process going. We did not research this phenomenon. It should only be noted that we filtered out such empty articles, because they could easily be identified as non-translations outright.

279 4.2.2 Metadata Submission Interface

280 The contacted users were invited to contribute metadata about themselves and the extracted texts 281 through the metadata submission interface, for which they were given personalized user accounts, 282 each of which was associated with all the extracted texts produced by the respective user. The inter-283 face consisted of four components, through which contributing users were taken sequentially: (1) a 284 login page for Wikipedia translators; (2) an instructions page, including opt-in checkboxes for reading 285 and accepting the project's terms of service and privacy policy; (3) a form for entering basic personal 286 information, including choices of the preferred degree of anonymity; (4) a form for entering text-re-287 lated metadata, displayed on separate pages for each text produced by a user. While Component 3 288 collected immutable personal data (native language, gender), Component 4 collected personal infor-289 mation relative to a given text (e.g. age or education level at the time of text production) and infor-290 mation on the text production circumstances (e.g. type of translation tools used), as shown in Fig. 1. 291 The contribution of metadata was not remunerated; instead, the contacted Wikipedians were en-292 couraged to contribute to a collaborative open science initiative. The contributors were given control 293 as to their choice of license, and full transparency regarding the project's legal framework was a major 294 concern in the design of the interface.

295

(a) Person-related

Welcome, user!			Transbank
Basic Profile			
Identification: 🕦	User-Specific Number O	Real Name () O Nick Name ()	Anonymous 1
Native Language(s) 🕦		Gen	der
Basque × Select Your Native Langua	age(s)	+ Fe	emale
b) text-related			Next
Wikipodia Translators (Contributing to TransBank	<	TransBank
Virkipedia Translators of Qualification Profile 1/3 Please provide us with the following deta "Mandibular_yan_kesici_diş" Highest Level of Education	ills about your qualifications <u>AT THE TIME OF TRANSLA</u>	TING the following article into Turkish.	n (Type) A
Virkpeura Translators (Qualification Profile 1/3 Please provide us with the following deta "Mandibular_yan_kesici_diş" Highest Level of Education (Master's degree	Completed Education (Type)	TING the following article Into Turkish. Ongoing Education translation trai	n (Type) ()
Virkpeura Translators (Qualification Profile 1/3 Please provide us with the following deta "Mandibular_yan_kesici_diş" Highest Level of Education (Master's degree	ils about your qualifications <u>AT THE TIME OF TRANSLA</u> Completed Education (Type)	TING the following article into Turkish. Ongoing Education translation trai	n (Type) ()
With pedia Translators (Qualification Profile 1/3 Please provide us with the following deta "Mandibular_yan_kesici_diş" Highest Level of Education (Master's degree Tools (Select one	ills about your qualifications <u>AT THE TIME OF TRANSLA</u> Completed Education (Type) Less than 3 years of academic t Other tertiary education (Medici Remuneration No No	TING the following article into Turkish. Ongoing Education translation trai	n (Type) () () s a free service.
Virkipeura Translators (Qualification Profile 1/3 Please provide us with the following deta "Mandibular_yan_kesici_diş" Highest Level of Education ① Master's degree Tools ① Select one	Ills about your qualifications <u>AT THE TIME OF TRANSLA</u> Completed Education (Type) Completed Education (Type) Less than 3 years of academic 1 Other tertiary education (Medici Remuneration No Pe(s) Workload Into Your Native Lange	TING the following article into Turkish.	n (Type) ① S a free service.
Virkipeura Translators (Qualification Profile 1/3 Please provide us with the following deta "Mandibular_yan_kesici_diş" Highest Level of Education (Master's degree Tools (Select one Workload From Your Native Languag 0 - 5 hours per week	ils about your qualifications <u>AT THE TIME OF TRANSLA</u> Completed Education (Type) Completed Education (Type) Less than 3 years of academic to Other tertiary education (Medici Other tertiary education (Medici Remuneration No Pe(s) V O - 5 hours per week	TING the following article into Turkish.	n (Type))
Virkipeura Translators (Qualification Profile 1/3 Please provide us with the following deta "Mandibular_yan_kesici_diş" Highest Level of Education (Master's degree Tools (Select one Workload From Your Native Language 0 - 5 hours per week	ills about your qualifications <u>AT THE TIME OF TRANSLA</u> Completed Education (Type) ① Less than 3 years of academic 1 Other tertiary education (Medici Remuneration ① No Je(s) ① Workload Into Your Native Lang ① - 5 hours per week	TING the following article into Turkish.	Type)
Virkipeura Translators (Qualification Profile 1/3 Please provide us with the following deta "Mandibular_yan_kesici_diş" Highest Level of Education () Master's degree Tools () Select one Workload From Your Native Languag () - 5 hours per week Association Membership ()	ils about your qualifications <u>AT THE TIME OF TRANSLA</u> Completed Education (Type) ① Less than 3 years of academic t Other tertiary education (Medici Remuneration ① No ye(s) ① Workload Into Your Native Lang 0 - 5 hours per week Sworn/Certified Translator ①	TING the following article into Turkish.	n (Type) () s a free service. anslator () years
Virkipeura Translators (Qualification Profile 1/3 Please provide us with the following deta "Mandibular_yan_kesici_diş" Highest Level of Education ① Master's degree Tools ① Select one Workload From Your Native Languag 0 - 5 hours per week Association Membership ① Select One	ils about your qualifications <u>AT THE TIME OF TRANSLA</u> Completed Education (Type) Less than 3 years of academic i Other tertiary education (Medici Remuneration No No ye(s) Workload Into Your Native Lang O - 5 hours per week Sworn/Certified Translator Select One	TING the following article into Turkish.	Type) Type) Type) Type) Type) Type) Type Type Type Type Type Type Type Type
Virkipeura Translators (Qualification Profile 1/3 Please provide us with the following deta "Mandibular_yan_kesici_diş" Highest Level of Education ① Master's degree Tools ① Select one Workload From Your Native Languag 0 - 5 hours per week Association Membership ① Select One I have read and agree to the Cont	Ils about your qualifications <u>AT THE TIME OF TRANSLA</u> Completed Education (Type) Less than 3 years of academic 1 Other tertiary education (Medici Remuneration Remuneration No Y No Y Y No Y Y Y Y Y Y Y Y Y Y Y Y	TING the following article into Turkish.	(Type)

296

297 Fig. 1 Screenshot of person-related (a) and text-related metadata (b) submission form for Wikipedia translators

298 **5** Results and discussion

We extracted an arbitrary quantity of bitexts containing about 50,000,000 tokens from Wikipedia, based on the information on which text was connected to which via an interlanguage link. From this material, about 4,800 parallel bitexts texts could be extracted using the described neural-network based translation identification methodology. The translators of 3,104 of these bitexts had the email functionality of their Wikipedia accounts enabled and could thus be contacted, in principle. In total, invitations to contribute metadata were sent to 934 translators – some translators had produced more
than one text, therefore the numbers of bitexts and translators are not the same. Due to limitations
we already mentioned about the number of emails that can be sent per 24 hours per IP address, we
used five Wikipedia accounts to send the emails containing the log-in information to translators. The
procedure described in Section 4.2 took almost a month but made the email distribution possible.

Of the 934 translators that were contacted, a total of 216 logged into our interface and provided us with the metadata for each of their articles. As they had been identified as translators, every one of them had translated at least one article, the most productive user, however, had written 136 article translations, with more than 38,000 contributions to Wikipedia in general. On average, users that had provided metadata through the interface had translated 4.1 articles each (Mdn = 1.0, SD = 10.5).

314 We received metadata for 881 different bitexts and 44 different language pairs, with Russian-315 Ukrainian being the most numerous one (191 bitexts), 15 language-pairs yielding only one article pair 316 each (see Table 2). It must be stated at this point that we do not know if there are language pairs with 317 even more texts as this is not entirely clear from the SQL database, because, firstly, the database con-318 tains this information in an unstructured way and we stopped parsing it after we had reached our 319 arbitrarily chosen number of texts, and, secondly it is not clear how many users follow the recommen-320 dation of including an interlanguage-link template when creating a translation. As far as our data are 321 concerned, on average, 20.0 bitexts (SD = 39.4, Mdn = 3.5) were observed for each language pair. On 322 average there were 5.8 different translators per language pair (Mdn = 2.0, SD = 8.8). The average num-323 ber of texts per translator varied greatly across language pairs: the most productive translators were 324 observed for Spanish-Asturian (35.5 texts per translator), French-Portuguese (8.0) and Russian-Ukrain-325 ian (6.4), while the mean number of translations per translator averaged across all language pairs was 326 2.8 (Mdn = 1.4, SD = 5.3).

327

328

329 Table 2 Quantitative summary of extracted and metadata-enriched bitexts by language pair

Language Pair	# bitexts	# trans- lators	Word count source lan- guage	Word count target language	Language Pair	# bitexts	# trans- lators	Word count source language	Word count tar- get language
Russian –					English – Rus-				
Ukrainian	191	30	182682	181521	sian	3	1	1968	1721
Spanish –	1/12	1	214420	200875	French – Roma-	2	2	1572	1667
Fnglish -	142	4	214420	200873	Spanish — Ital-	5	5	1373	1007
Spanish	121	44	171063	187488	ian	2	1	4061	3791
English – French	64	24	98232	108640	Italian – French	2	1	1552	1618
Spanish – Catalan	41	18	53420	53307	French – Cata- Ian	2	2	2984	3185
English - Por- tuguese	41	21	68489	72400	Belarusian — Ukrainian	2	2	1269	1265
Spanish – Galician	36	9	49751	46975	German – Por- tuguese	1	2	1308	1405
English – Ro- manian	35	6	67275	70974	Portuguese – Galician	1	1	3509	3391
English – Greek	26	5	27771	28758	English – Nor- wegian	1	1	473	415
Ukrainian – Russian	24	9	13804	13817	English – Astu- rian	1	1	2701	2696
French – Portuguese	24	3	20839	20252	French – Italian	1	1	470	405
English – Catalan	23	12	49968	54294	Portuguese – French	1	1	2137	2558
English – Italian	15	6	17128	17658	Spanish – French	1	1	3769	3379
Russian — Belarusian	15	4	10864	10911	Italian – English	1	1	681	685
Catalan – Spanish	11	6	9853	9854	German – Ro- manian	1	1	1678	1832
French — Spanish	11	8	13346	13476	Bulgarian – Cat- alan	1	1	1779	2109
English – Ukrainian	7	5	8921	7283	German – French	1	1	2823	3265
Spanish – Portuguese	6	4	5841	5535	English – Turk- ish	1	1	122	110
French – English	4	3	5673	5247	Romanian – Catalan	1	1	3599	4112
German – English	4	2	3460	3472	French – Dutch	1	1	404	459
Russian – English	4	1	3985	4760	Belarusian — Russian	1	1	1357	1347
Portuguese – Spanish	4	3	8464	8892	TOTAL	881	-	1,147,394	1,169,826
Swedish – Norwegian	3	1	1928	2022	MEAN/PAIR	20.0	5.8	26,077.1	26,587.0

330

Through the metadata submission interface described in Section 4.2 we collected two items of person-related and eleven items of text-related metadata. The former were considered immutable and therefore users were prompted only once to submit them via Component 3 of the interface, whereas Component 4 for the collection of text-related metadata was displayed for every single text produced by a given user. Since in Section 3 completeness and accuracy were identified as two fundamental criteria of metadata quality and hence of data usefulness, the viability of the proposed data collection
methodology needs to be assessed in terms of these two criteria.

338 Metadata completeness

339 The mean submission rates for the eleven text-related and two person-related metadata fields of the 340 submission interface are reported in **Table 3**. Due to technical issues, the data for one of the fields (the 341 tools field, where participants could state whether they used electronic translation tools such as trans-342 lation memory systems) were not collected correctly and were therefore excluded from the analysis. 343 The two items of immutable person-related metadata, native language and gender, which can arguably 344 be considered more prone to contributors' data privacy concerns, were provided by 100% and 88.4% 345 of users, respectively. For 287 out of 881 texts (35.6%), users provided all of the requested metadata. 346 Summing up, 216 out of 934 contacted Wikipedia translators (23.1%) submitted metadata, thus provid-347 ing metadata-enrichment for 881 of the 3104 (28.4%) extracted parallel texts for which we were able to request metadata from their translators. 287 out of these 3104 texts (9.2%) were fully annotated 348 349 by the translators with both person-related and text-related metadata. At the text level, the overall 350 completeness rate for the 881 metadata-enriched texts is 92.2%.

Table 3 Degree of completeness of filling in metadata fields: e.g., "83.4%" means a field has been filled in 83.4% of the time.

Category	#	Ν	Min		Mdn	Mean	SD	Max
	fields		(least	filled-in				(most filled-in
			field)					field)
text-related	11	881 texts	83.4%		88.1%	91.1%	7.3	100% (4 fields)
person-re-	2	216 us-	88.4%		94.2%	94.2%	8.2	100% (1 field)
lated		ers						

352

The reported completeness statistics have to be seen in the light of the voluntary nature of user submissions and the legally motivated absence of compulsory metadata fields in the submission interface. For the users who contributed less than 10 texts, the submission rate was 91.4%. Among the 19 users who contributed ten texts or more – their contributions account for 461 (52.3%) of all texts – the submission rate for text-related metadata was 90.9%, which can be seen as an indicator of commitment. However, since we cannot entirely rule out the possibility that parts of our invitations were not correctly received by the addressees due to the anti-spam filters in the Wikipedia messaging system, we cannot be sure whether the response rates and completeness statistics adequately mirror Wikipedians' willingness to contribute to open science initiatives.

362 Metadata accuracy

363 In the context of the present study metadata accuracy is closely related to the quality criterion of trust, 364 because the metadata collection was limited to text-extrinsic information known to data producers 365 only. We used information publicly available on Wikipedia user pages⁹ as a proxy to assess the reliabil-366 ity of submitted metadata. To this end, we randomly sampled 32 from the 216 contributing users 367 (14.8%) and manually compared the metadata submitted by them with their user pages, focusing on 368 gender and native language. Users who left either field empty (25 out of 216 users = 11.6% left the 369 gender field empty; the native language field, on the other hand, was filled in by 100% of users) were 370 not considered for the random sample, because, firstly, not providing certain metadata is a conscious 371 decision taken by users and in line with the study's legal framework, and, secondly, missing data is a 372 concern of data completeness rather than accuracy and reliability. Only explicit information from the 373 page text, highly likely to have been written by the users themselves (e.g. gender-specific occupations, 374 such as Spanish escritor 'male writer' vs. escritora 'female writer') or from the so-called userboxes¹⁰ 375 was used to verify users' gender and native language. Less reliable information was ignored, e.g. the 376 grammatical gender of usernames or gender-specific forms of the word user in the page header (e.g. 377 Spanish usuario vs. usuaria or German Benutzer vs. Benutzerin) – it is not clear if the grammatical gen-378 der actually reflects the gender of the respective user themselves in the case of usernames; also, the

⁹ https://en.wikipedia.org/wiki/Wikipedia:User_pages_

¹⁰ <u>https://en.wikipedia.org/wiki/Wikipedia:Userboxes</u>

male form of the word meaning *user* is often used generically in many languages. A major limitation to this analysis is that Wikipedia user pages are not standardized and users are entirely free to choose whether, what and how to publish personal data; however, they are the *only* source for cross-checking submitted metadata.

383 Table 4 shows the confusion matrix for user-submitted vs. manually cross-checked gender 384 metadata. In 18 cases, no gender information was found on the user pages, which means that in most 385 cases user submissions are the only way to collect pertinent metadata. However, in the remaining 14 386 cases for which gender data was available, the accuracy of submitted vs. cross-checked gender infor-387 mation was merely 0.36 (Kappa agreement = 0.06, p > 0.05), suggesting that user-submitted metadata 388 was less reliable than initially assumed. Further evidence for the lacking reliability of submitted 389 metadata is that women account for approximately 34% of all users in both the sample and full dataset, 390 while estimates suggest that less than 10% of Wikipedia editors are women (Ford and Wajcman, 2017). 391 Contradictory gender information was also observed in spot checks for several other users not part of 392 the random sample, including the most productive of all users. It is surprising that nine of the sampled 393 users submitted gender information that contradicted their user pages, since users were free to leave 394 the gender field empty in the submission interface. We can only speculate about the reasons for these 395 contradictions: a lack of intrinsic motivation to contribute to an open science project, or data privacy 396 concerns despite a transparent data privacy policy that contains the option to leave fields empty and 397 to freely choose from various degrees on anonymity. A further source of contradiction may be ascribed 398 to the availability of three rather than two gender options in the metadata submission field (male, 399 female, other), which was to increase the social inclusiveness of the metadata survey. The option 400 'other', selected by four of the sampled users and by 12.5% of all users, may have been misinterpreted 401 as 'not specified', especially among non-native speakers of English. The reliability of these cases can 402 hardly be assessed, because it is unclear whether users who do not identify themselves with the tradi-403 tional binary distinction would openly specify their gender identity on their user pages. Note, however, 404 that there are numerous Wikipedia userbox templates¹¹ that allow specifying gender identities other

than male and female.

406 **Table 4** Confusion matrix of user-submitted vs. cross-checked gender metadata (F = female, M = male, O = other)

	Cross-checked								
Submitted		F	Μ	0	n.a.	Total			
	F	2	4	0	5	11			
	Μ	1	3	0	13	17			
	0	0	4	0	0	4			
	Total	3	11	0	18	32			

407

408 By contrast, the submitted metadata about users' native languages from the random sample was 409 in perfect agreement (accuracy = 1.0, Kappa = 1.0, $p \le 0.01$) with the information from user pages. Ten 410 of the sampled user pages did not contain native language information, whereas all users indicated 411 their native language in the submission interface. This, again, shows that the survey was the only way 412 to retrieve the metadata. Given the perfect response and agreement rates, it seems that information 413 about the native language is perceived to be less sensitive by users. Altogether, the discrepancies in 414 completeness and accuracy of the two metadata items under investigation suggest that the willingness 415 to provide metadata is dependent on the type of the requested information.

Spot checks revealed that there might be accuracy issues with regard to other metadata fields as well, for instance age. However, these instances were not investigated systematically, because the presented analyses suffice to show that user-submitted data may be inaccurate to a certain extent.

419 6 Conclusion and Outlook

Large amounts of parallel text data covering a wide range of different languages, including less common language pairs, can be extracted in a fairly straightforward manner from Wikipedia using stateof-the-art neural network approaches. Our methodology could have yielded many more parallel texts, but since the aim of the present study was to evaluate in how far the automatically extracted data can

¹¹ <u>https://en.wikipedia.org/wiki/Category:Gender_user_templates</u>

424 be manually enriched, the corpus had to be limited to a manageable size. Also, the availability of disk 425 space and memory to process and analyse the Wikipedia dumps might constitute limiting factors, de-426 pending on the available infrastructure. Noteworthy, our study was not restricted to particular lan-427 guage pairs, because we aimed to explore the multilingual potential of parallel text extraction from 428 Wikipedia and to capture the "dark matter" (Shuttleworth, 2017) of Wikipedia translations, i.e. texts 429 whose translational status is not explicitly declared. That said, our approach depends on certain administrative and structural information to select a set of candidate text pairs - information contained 430 431 in the discussed interlanguage-link database. This also means that the language pairs that ended up in 432 our sample were not arbitrarily chosen by us, but are due to the structure of this database and might 433 therefore not be representative of the characteristics of Wikipedia when it comes to translation. The-434 oretically, the neural network architecture could mine the entirety of Wikipedia for parallel texts, but 435 such an approach would be too time consuming, computationally too costly and thus impractical. Even 436 with those restrictions in mind, Wikipedia is a rich source of parallel texts in constant growth.

Wikipedia's structural and administrative information not only makes parallel text extraction feasible, it also allows identifying translators. The identification of source text authors was deliberately not tackled in the present study, because the collaborative nature of text production in Wikipedia challenges the traditional understanding of authorship. Attempts to pinpoint individual users as authors of particular texts are thus inevitably mere approximations. Similarly, metadata labelsets traditionally used in corpus linguistics are of limited value to capture collaborative text production in all its complexity.

As far as the metadata enrichment procedure is concerned, important technical and social conclusions can be drawn. On the technical side, Wikipedia's limitations on the messaging system and its anti-spam filter have repercussions on the scalability of metadata enrichment via self-submission, because users can only be contacted manually, which is very time-consuming in view of the daily messaging limits. Alternative ways to engage with the community of Wikipedia translators might circumvent this technical limitation, for instance through a closer collaboration with one of the numerous 450 translation initiatives and task forces in Wikipedia. On the social side, the usefulness of the collected 451 data did not fully meet the established quality criteria. Although the assessment of the collected 452 metadata was based on a small sample and focused only on a small subset of the metadata fields of 453 interest, it highlighted the need for metadata verification and curation. No less important, solutions to 454 improve metadata completeness and accuracy are crucial to crowd-based data collection initiatives. 455 Social technology has a great potential to foster user motivation and involvement, for instance through 456 the already mentioned closer collaboration with the user community in question, through gamification 457 or user rating systems.

458 Despite the identified quality issues, the study has shown that considerable quantities of metadata 459 otherwise not available can be collected with the help of social technology. The presented approach 460 tries to bridge the gap between well-established NLP techniques and the needs of research that re-461 quires high-quality metadata, including, but not limited to, the digital humanities. In the future, auto-462 matically extracted and user-enriched translation data may contribute to work on numerous research 463 questions, such as the improvement of multilingual NLP applications, the study of the multilingual ex-464 pansion of Wikipedia and the dynamics of cross-lingual and cross-cultural knowledge production, or 465 shed light on peripheral and thus underresearched translation practices, most prominently the novel 466 phenomenon of massively open translation (O'Hagan, 2016).

The findings of our study suggest that the combination of automated and manual procedures is a viable approach to collecting humanities data, part of which are translation data. Crowd-based selfsubmission makes it possible to expand data repositories and thus to extend the scope of research to previously untapped resources. However, such approaches are not without risk and require the careful consideration of questions such as how to ensure user involvement, quality control, and data protection.

- 473
- 474

475 Funding

476 This research is part of the project "TransBank: A Meta-Corpus for Translation Research", funded by

477 the *goldigital 2.0* programme of the Austrian Academy of Sciences (grant number GD 2016/56).

478 Acknowledgments

We thank the Wikipedia who voluntarily participated in the metadata collection. The computational
results presented have been achieved (in part) using the HPC infrastructure LEO of the University of
Innsbruck.

482

483 References

- 484 **Aghaebrahimian, A.** (2018). Deep Neural Networks at the Service of Multilingual Parallel Sentence
- 485 Extraction. In E. M. Bender, L. Derczynski, & P. Isabelle (Eds.), *Proceedings of the 27th International*
- 486 *Conference on Computational Linguistics* (pp. 1372–1383). Santa Fe, New Mexico: Association for

487 Computational Linguistics. Retrieved from <u>https://www.aclweb.org/anthology/C18-1116</u>

488 Babych, B., Su, F., Hartley, A., Aker, A., Paramita, M. L., Clough, P., & Gaizauskas, R. (2019). Cross-

489 Language Comparability and Its Applications for MT. In I. Skadiņa, R. Gaizauskas, B. Babych, N.

490 Ljubešić, D. Tufiş, & A. Vasiljevs (Eds.), *Theory and Applications of Natural Language Processing*.

491 Using Comparable Corpora for Under-Resourced Areas of Machine Translation (Vol. 1866, pp. 13–

492 53). Cham: Springer International Publishing. <u>https://doi.org/10.1007/978-3-319-99004-0_2</u>

493 Barbosa, L., Rangarajan Sridhar, V. K., Yarmohammadi, M., & Bangalore, S. (2012). Harvesting Paral-

- 494 lel Text in Multiple Languages with Limited Supervision. In M. Kay & C. Boitet (Eds.), *Proceedings*
- 495 of COLING 2012: Technical Papers (pp. 201–214). Mumbai: Association for Computational Linguis-
- 496 tics. Retrieved from https://www.aclweb.org/anthology/C12-1013
- 497 Barrón-Cedeño, A., España-Bonet, C., Boldoba, J., & Màrquez, L. (2015). A Factory of Comparable
- 498 Corpora from Wikipedia. In P. Zweigenbaum, S. Sharoff, & R. Rapp (Eds.), Proceedings of the
- 499 *Eighth Workshop on Building and Using Comparable Corpora* (pp. 3–13). Beijing: Association for
- 500 Computational Linguistics. <u>https://doi.org/10.18653/v1/W15-3402</u>

501	Bruce, T. R., & Hillmann, D. I. (2004). The Continuum of Metadata Quality: Defining, Expressing, Ex-
502	ploiting. In D. I. Hillmann & E. L. Westbrooks (Eds.), Metadata in practice (pp. 238–256). Chicago:
503	American Library Association.

- 504 Chu, C., Dabre, R., & Kurohashi, S. (2016). Parallel Sentence Extraction from Comparable Corpora
- 505 with Neural Network Features. In N. Calzolari, K. Choukri, T. Declerck, S. Goggi, M. Grobelnik, B.
- 506 Maegaard, ... S. Piperidis (Eds.), Proceedings of the Tenth International Conference on Language
- *Resources and Evaluation (LREC 2016)* (pp. 2931–2935). Paris: European Language Resources As sociation.
- 509 Enright, J., & Kondrak, G. (2007). A Fast Method for Parallel Document Identification. In C. Sidner, T.
- 510 Schultz, M. Stone, & C. Zhai (Eds.), Human Language Technologies 2007: The Conference of the
- 511 North American Chapter of the Association for Computational Linguistics; Companion Volume,
- 512 Short Papers (pp. 29–32). Stroudsburg, PA: Association for Computational Linguistics. Retrieved
- 513 from https://www.aclweb.org/anthology/N07-2008
- 514 España-Bonet, C., Varga, Á. C., Barrón-Cedeño, A., & van Genabith, J. (2017). An Empirical Analysis
- of NMT-Derived Interlingual Embeddings and Their Use in Parallel Sentence Identification. *IEEE*
- 516 Journal of Selected Topics in Signal Processing, 11(8), 1340–1350.
- 517 https://doi.org/10.1109/JSTSP.2017.2764273
- 518 Etchegoyhen, T., & Azpeitia, A. (2016). A Portable Method for Parallel and Comparable Document
- 519 Alignment. Baltic Journal of Modern Computing, 4(2), 243–255. Retrieved from
- 520 https://www.bjmc.lu.lv/fileadmin/user_upload/lu_portal/projekti/bjmc/Contents/4_2_15_Etche-
- 521 goyhen.pdf
- 522 Ford, H., & Wajcman, J. (2017). 'Anyone can edit', not everyone does: Wikipedia's infrastructure and
- 523 the gender gap. *Social Studies of Science*, 47(4), 511–527.
- 524 <u>https://doi.org/10.1177/0306312717692172</u>
- 525 Fung, P., & Cheung, P. (2004). Mining Very-Non-Parallel Corpora: Parallel Sentence and Lexicon Ex-
- 526 traction via Bootstrapping and EM. In *Proceedings of the 2004 Conference on Empirical Methods*

- 527 *in Natural Language Processing: 25-26 July2004, Barcelona* (pp. 57–63). Barcelona: Association for
 528 Computational Linguistics.
- 529 Gambier, Y. (2010). Translation strategies and tactics. In Y. Gambier & L. van Doorslaer (Eds.), Hand-
- 530 book of Translation Studies. Handbook of Translation Studies: Volume 1 (Vol. 1, pp. 412–418). Am-
- 531 sterdam: John Benjamins Publishing Company. <u>https://doi.org/10.1075/hts.1.tra7</u>
- 532 Greenberg, J. (2017). Big Metadata, Smart Metadata, and Metadata Capital: Toward Greater Synergy
- 533 Between Data Science and Metadata. *Journal of Data and Information Science*, 2(3), 19–36.
- 534 https://doi.org/10.1515/jdis-2017-0012
- 535 Grégoire, F., & Langlais, P. (2018). Extracting Parallel Sentences with Bidirectional Recurrent Neural
- 536 Networks to Improve Machine Translation. In E. M. Bender, L. Derczynski, & P. Isabelle (Eds.), Pro-
- 537 *ceedings of the 27th International Conference on Computational Linguistics* (pp. 1442–1453). Santa
- 538 Fe, New Mexico: Association for Computational Linguistics. Retrieved from
- 539 https://www.aclweb.org/anthology/C18-1122
- 540 Jones, H. (2018). Wikipedia, Translation, and the Collaborative Production of Spatial Knowledge. Alif:
- 541 *Journal of Comparative Poetics*. (38), 264–297.
- 542 Koplenig, A. (2017). The Impact of Lacking Metadata for the Measurement of Cultural and Linguistic
- 543 Change Using the Google Ngram Data Sets—Reconstructing the Composition of the German Cor-
- pus in Times of WWII. *Digital Scholarship in the Humanities*, *32*(1), 169-188.
- 545 <u>https://doi.org/10.1093/llc/fqv037</u>
- 546 Labaka, G., Alegria, I., & Sarasola, K. (2016). Domain Adaptation in MT Using Titles in Wikipedia as a
- 547 Parallel Corpus: Resources and Evaluation. In N. Calzolari, K. Choukri, T. Declerck, S. Goggi, M.
- 548 Grobelnik, B. Maegaard, ... S. Piperidis (Eds.), *Proceedings of the Tenth International Conference*
- 549 on Language Resources and Evaluation (LREC 2016) (pp. 2209–2213). Paris: European Language
- 550 Resources Association. Retrieved from <u>http://www.lrec-conf.org/proceed-</u>
- 551 ings/lrec2016/pdf/632 Paper.pdf

- 552 Mohammadi, M. (2016). Parallel Document Identification using Zipf's Law. In R. Rapp, P. Zweigen-
- baum, & S. Sharoff (Eds.), *Proceedings of the 9th Workshop on Building and Using Comparable Cor-*
- 554 *pora* (pp. 21–25). Portorož: Evaluations and Language resources Distribution Agency.
- 555 Morin, E., Hazem, A., Boudin, F., & Loginova-Clouet, E. (2015). LINA: Identifying Comparable Docu-
- 556 ments from Wikipedia. In P. Zweigenbaum, S. Sharoff, & R. Rapp (Eds.), *Proceedings of the Eighth*
- 557 *Workshop on Building and Using Comparable Corpora* (pp. 88–91). Beijing: Association for Compu-
- tational Linguistics. <u>https://doi.org/10.18653/v1/W15-3413</u>
- 559 O'Hagan, M. (2016). Massively Open Translation: Unpacking the Relationship Between Technology
- and Translation in the 21st Century. *International Journal of Communication*, *10*, 929–946.
- 561 Patry, A., & Langlais, P. (2011). Identifying Parallel Documents from a Large Bilingual Collection of
- 562 Texts: Application to Parallel Article Extraction in Wikipedia. In P. Zweigenbaum, R. Rapp, & S.
- 563 Sharoff (Eds.), Proceedings of the 4th Workshop on Building and Using Comparable Corpora: Com-
- 564 *parable Corpora and the Web* (pp. 87–95). Portland: Association for Computational Linguistics.
- 565 Retrieved from <u>https://www.aclweb.org/anthology/W11-1212</u>
- 566 Sharoff, S., Rapp, R., & Zweigenbaum, P. (2013). Overviewing Important Aspects of the Last Twenty
- 567 Years of Research in Comparable Corpora. In S. Sharoff, R. Rapp, P. Zweigenbaum, & P. Fung
- 568 (Eds.), Building and Using Comparable Corpora (pp. 1–17). Berlin, Heidelberg: Springer Berlin Hei-
- 569 delberg. <u>https://doi.org/10.1007/978-3-642-20128-8_1</u>
- 570 Sharoff, S., Rapp, R., Zweigenbaum, P., & Fung, P. (Eds.). (2013). Building and Using Comparable
- 571 *Corpora*. Berlin, Heidelberg: Springer Berlin Heidelberg. <u>https://doi.org/10.1007/978-3-642-</u>
- 572 <u>20128-8</u>
- 573 Shuttleworth, M. (2017). Locating foci of translation on Wikipedia. Translation Spaces, 6(2), 310-
- 574 332. <u>https://doi.org/10.1075/ts.6.2.07shu</u>
- 575 Shuttleworth, M. (2018). Translation and the Production of Knowledge in Wikipedia: Chronicling the
- 576 Assassination of Boris Nemtsov. (38), 231–263.

- 577 Skadiņa, I., Gaizauskas, R., Vasiļjevs, A., & Paramita, M. L. (2019). Introduction. In I. Skadiņa, R. Gai-
- 578 zauskas, B. Babych, N. Ljubešić, D. Tufiş, & A. Vasiljevs (Eds.), *Theory and Applications of Natural*
- 579 Language Processing. Using Comparable Corpora for Under-Resourced Areas of Machine Transla-
- 580 *tion* (Vol. 1866, pp. 1–11). Cham: Springer International Publishing. <u>https://doi.org/10.1007/978-</u>
- 581 <u>3-319-99004-0 1</u>
- 582 Smith, J. R., Quirk, C., & Toutanova, K. (2010). Extracting Parallel Sentences from Comparable Cor-
- 583 pora using Document Level Alignment. In R. Kaplan, J. Burstein, M. Harper, & G. Penn (Eds.), Hu-
- 584 man Language Technologies: The 2010 Annual Conference of the North American Chapter of the
- 585 Association for Computational Linguistics (pp. 403–411). Los Angeles: Association for Computa-
- 586 tional Linguistics. Retrieved from <u>https://www.aclweb.org/anthology/N10-1063</u>
- 587 **Tiedemann, J.** (2011). *Bitext Alignment*. <u>https://doi.org/10.2200/S00367ED1V01Y201106HLT014</u>.
- 588 Morgan & Claypool.
- 589 **Ture, F., & Lin, J.** (2012). Why Not Grab a Free Lunch? Mining Large Corpora for Parallel Sentences to
- 590 Improve Translation Modeling. In E. Fosler-Lussier, E. Riloff, & S. Bangalore (Eds.), *Proceedings of*
- 591 the 2012 Conference of the North American Chapter of the Association for Computational Linguis-
- 592 *tics: Human Language Technologies* (pp. 626–630). Stroudsburg, PA: Association for Computa-
- 593 tional Linguistics. Retrieved from <u>https://www.aclweb.org/anthology/N12-1079</u>
- 594 Ustaszewski, M., & Stauder, A. (2017). TransBank: Metadata as the Missing Link between NLP and
- 595 Traditional Translation Studies. In I. Temnikova, C. Orasan, G. Corpas Pastor, & S. Vogel (Eds.), Pro-
- 596 ceedings of the First Workshop on Human-Informed Translation and Interpreting Technology
- 597 (pp. 29–35).
- Wikipedia. (2019). Wikipedia:Translation. Retrieved from https://en.wikipedia.org/w/index.php?ti-tle=Wikipedia:Translation&oldid=888728057. Accessed 24 June 2019
- 600 Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J. J., Appleton, G., Axton, M., Baak, A., . . . Mons,
- 601 **B.** (2016). The FAIR Guiding Principles for scientific data management and stewardship. *Scientific*
- 602 *Data*, *3*, 160018. <u>https://doi.org/10.1038/sdata.2016.18</u>

- 603 Zanettin, F. (2012). Translation-driven corpora corpus resources for descriptive and applied transla-
- 604 *tion studies. Translation practices explained*. Manchester, UK, Kinderhook, NY: St. Jerome Pub.