

CLD²: Language Documentation Meets Natural Language Processing for Revitalising Endangered Languages

Anonymous ACL submission

Abstract

Language revitalisation should not be understood as a direct outcome of language documentation, which is mainly focused on the creation of language repositories. Natural language processing (NLP) offers the potential to complement and exploit these repositories through the development of language technologies that may directly impact the vitality status of endangered languages. In this paper, we discuss the current state of the interaction between language documentation and computational linguistics, present a diagnosis of how the outputs of recent documentation projects for endangered languages are underutilised for the NLP community, and discuss how the situation could change from both the documentary linguistics and NLP perspectives. All this is introduced as a bridging paradigm called Computational Language Documentation and Development (CLD²). CLD² calls for (1) the inclusion of NLP-friendly annotated data as a deliverable of future language documentation projects; and (2) the exploitation of language documentation databases by the NLP community to promote the computerization of endangered languages at a global scale.

1 Introduction

There are around 6,500 mutually unintelligible languages in the world (Hammarström et al., 2018). However, several minority languages are in danger of being lost forever without leaving systematic records. In response to this, in the last decades *Documentary Linguistics* has become a major and vibrant field in Linguistics, which attempts to produce permanent records of the linguistic and cultural practices of the most threatened speech communities ((Himmelman, 2012; Austin, 2010; Woodbury, 2011), among many others).

The outcomes of documenting a language in the frame of contemporary *Documentary Linguistics* often comprise a large amount of transcribed, translated and parsed data, supported in audio and video.

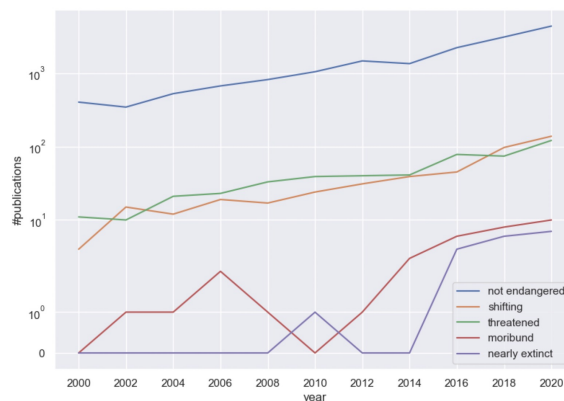


Figure 1: Number of publications in the ACL Anthology where languages are explicitly named in the title or abstract, and they are classified by their vitality from the Agglomerated Endangerment Status (Seifart et al., 2018). Vertical axis is in log-scale.

These data are often allocated in international language archives, from which scholars and members of speech communities could access them. Although field linguists often incorporate revitalisation components in their documentation projects, language *documentation* and language *revitalisation* are not equivalent in terms of their frames, methods and outcomes. Language revitalisation may surely take advantage of the records produced in language documentation projects, but creating a language repository alone cannot revert language endangerment or decay.

The concern about language endangerment has also reached contemporary approaches to Computational Linguistics, and in the last years, the “computerisation” of minority languages has become a growing field in NLP research (Berment, 2002). NLP developments’ potential for revitalising endangered languages is high, but there is still moderate interaction between *Documentary Linguistics* and NLP research for language revitalisation.

In this paper, we reflect on the necessary interactions between *Documentary Linguistics* and NLP

066	by introducing a paradigm that assumes language documentation and NLP developments as integral parts of a global response to language endangerment: Computational language Documentation and Development (CLD ²).	
067		
068		
069		
070		
071	2 Language documentation and language revitalisation	
072		
073	Language documentation (or documentary linguistics) emerged at the end of the last century as a research program whose primary motivation lies in the concern about the accelerating loss of language diversity in the world. As a response, language documentation aims to create permanent records of the linguistic and cultural practices of the most threatened speech communities (Himmelmann, 1998; Austin, 2010; Woodbury, 2011). These records are framed as databases, ideally including several hours of audio and video recordings of monologue and dialogue texts belonging to various genres and topics (e.g. traditional tales and myths, verbal art, jokes, historical facts, life stories, cultural knowledge, among others). A good portion of these recordings is transcribed, translated and parsed. Each transcribed sentence is expected to be time-aligned and to include an orthographic or IPA representation, a morphemic parse, glossing, information about parts of speech and a free translation.	
074		
075		
076		
077		
078		
079		
080		
081		
082		
083		
084		
085		
086		
087		
088		
089		
090		
091		
092		
093	Producing such linguistic databases is a long-term and time-consuming task that may take several years and requires considerable funding. The expectation is that these linguistic databases, conceptualised as multipurpose repositories allocated in international archives, will be preserved for posterity and thus will support community-based revitalisation projects in the future. Although it is true that language documentation projects very often incorporate revitalisation components, they are inevitably marginal since the documentation itself is the main focus of documentary linguistics. Therefore, the contribution of language documentation to language revitalisation is potentially significant but mainly indirect: the linguistic repositories produced in the frame of language documentation projects can indeed contribute to future revitalisation projects, but crafting and archiving a repository will not necessarily have a positive impact on the vitality status of an endangered language.	
094		
095		
096		
097		
098		
099		
100		
101		
102		
103		
104		
105		
106		
107		
108		
109		
110		
111		
112		
	3 Language documentation and computational linguistics	113
		114
	Most interactions between computational linguistics and documentary linguistics relate to the release of software tools for language documentation, processing and archiving (van Esch et al., 2019; Anastasopoulos et al., 2020). Computational linguists and computer scientists have developed advanced software tools to assist field linguists in the various processes of contemporary language documentation, making them less time-consuming, more efficient and more systematic. These tools have been crucial for the exponential growth of language documentation on a global scale.	115
		116
		117
		118
		119
		120
		121
		122
		123
		124
		125
		126
	Contemporary language documentation implies a large amount of technical sophistication for managing, annotating, processing and archiving lasting and large repositories (Himmelmann, 2006; Austin, 2006; Woodbury, 2003, among many others). This could not be achieved without the contribution of computer scientists (particularly software developers). In the last decades, we have witnessed the release of specialised software tools nowadays customary for language documentation, speech analysis and linguistic fieldwork. Field linguist’s Toolbox (before “Shoebox”) (Summer Institute of Linguistics, 2021a) and more recently Fieldworks (FLex) (Summer Institute of Linguistics, 2021b) are data management and analysis tools for field linguists developed by the Summer Institute of Linguistics, which are used in language documentation and taught in linguistics schools worldwide. Toolbox and Flex allow to create dictionaries, which can be used for morphosyntactic parsing and annotation of transcribed texts. Transcription is often conducted in a different and nowadays very popular software called ELAN (Max Planck Institute for Psycholinguistics, 2021), developed by the Max Planck Institute for Psycholinguistics. ELAN allows to visualise and play audio and video files in order to create time-aligned transcriptions and translations. ELAN can also be used for morphological parsing, but most linguists prefer to conduct such tasks in Toolbox or FLex since ELAN transcriptions can be easily exported into these programs. In Toolbox or Flex, each sentence in an ELAN file (containing a transcription and a free translation) can receive morphemic parsing, morpheme-by-morpheme glossing and parts of speech tags, among any other relevant information in the frame of a specific project. The resulting	127
		128
		129
		130
		131
		132
		133
		134
		135
		136
		137
		138
		139
		140
		141
		142
		143
		144
		145
		146
		147
		148
		149
		150
		151
		152
		153
		154
		155
		156
		157
		158
		159
		160
		161
		162
		163

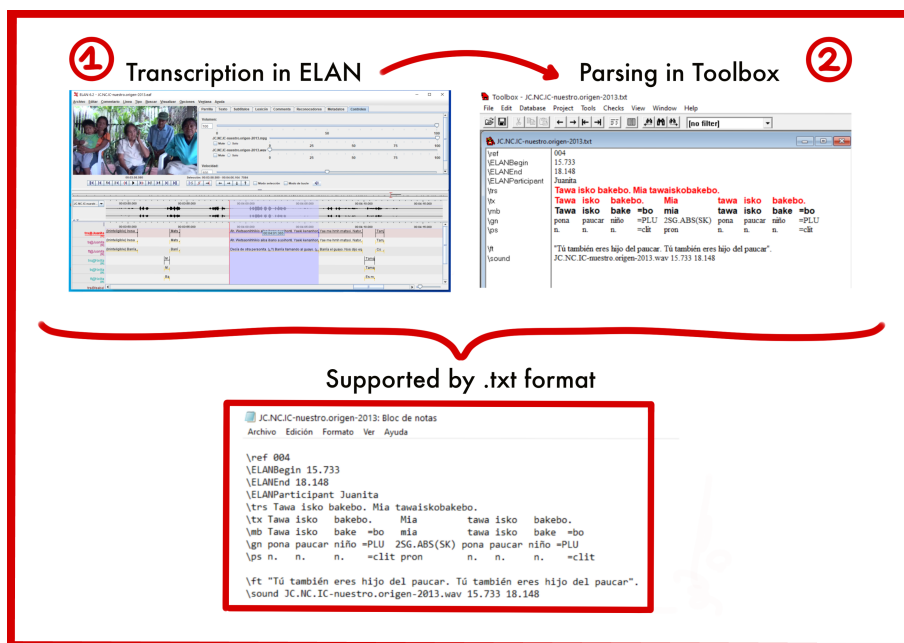


Figure 2: Graphic representation of the standard computational frame of language documentation: transcription is conducted in ELAN; ELAN files are imported into Toolbox or FLeX where they are fully parsed and glossed. Crucially, we are dealing with .txt files throughout the process, which enormously facilitates their manipulation in any programming language

Toolbox/Flex files are text files that can be opened back in ELAN, in PRAAT (a phonetics analyser) (Boersma and Weenink, 2001), or to be processed in Python or any other programming language as plain texts. This is shown in Figure 2.

In sum, there have been several attempts from the computational side trying to create or incorporate intelligent components in language documentation tools and procedures (Good et al., 2014; Arppe et al., 2017, 2019; van Esch et al., 2019; Anastopoulos et al., 2020). We find a one-direction application (computation into language documentation), but there are still few developments in the other direction (language documentation into computation). One of our takes in this paper is that language documentation can significantly contribute to computational linguistics by providing data and insights to develop NLP tools for endangered languages.

4 NLP has not really met endangered language documentation

As mentioned before, NLP has mainly focused on aiding the language documentation pipeline. However, has NLP taken advantage of the outputs of the documentation projects, especially for endangered languages?

4.1 Data

To address that question, we looked into the central repository of NLP publications: the ACL Anthology¹, the language inventory of massive multilingual datasets in NLP research (UniMorph (McCarthy et al., 2020), Universal Dependencies (Nivre et al., 2020), Tatoeba (Tiedemann, 2020))², and the central database of language documentation projects for endangered languages: The Endangered Languages Archive, or ELAR³, which is supported by the Endangered Languages Documentation Programme or ELDP⁴.

Besides, we work with the list of languages from Glottolog 4.4 (Hammarström et al., 2021), which is an extended inventory of living and extinct languages, including metadata such as geographical location and other properties. Moreover, we use the Agglomerated Endangerment Status (AES) classification proposed by Seifart et al. (2018) to distinguish the vitality status of the language inventory. The classes are, from more to less vital: not endangered, shifting, threatened, moribund, nearly

¹<https://aclanthology.org/>

²We chose these datasets as they are the most diverse collections according to their language inventory.

³<https://www.elararchive.org/>

⁴<https://www.eldp.net/>

212 extinct and extinct⁵.

213 4.2 Processing

214 With the language inventory and their vitality sta-
215 tus, we first identified all the publications in the
216 ACL Anthology (both conference and workshop
217 proceedings) whose title or abstract explicitly in-
218 cludes the name of a language⁶. We manually clean
219 false positives, such as concise language names
220 (less than five characters) that can be confused with
221 English words or acronyms.

222 A similar procedure is done with the ELAR
223 database: all the projects are extracted, the lan-
224 guage names are matched with the Glottolog in-
225 ventory, and we manually curated potential false
226 positives. From all the 570 projects published in
227 the ELAR database, we identified 307 language
228 names matching with the Glottolog database. With
229 this, we obtained geographical information for 286
230 languages.

231 The procedure is similar for the massively multi-
232 lingual (MM) datasets (Unimorph, Universal De-
233 pendencies and Tatoeba), and the language iden-
234 tifiers (ISO code or name) are matched with the
235 Glottolog inventory. Details of the considered lan-
236 guages are shown in Table 1.

237 4.3 Results

238 First, we look into how the NLP literature has con-
239 sidered endangered languages across time. Figure 1
240 shows that, in the current century, there is a consid-
241 erable growth of publications related to languages
242 with shifting or threatened status (from ten to a hun-
243 dred papers annually), but a very shy increase of the
244 moribund or nearly extinct languages (from zero to
245 ten annually), which are the most endangered ones.
246 Besides, this is highly contrasted by the continuous
247 increment of NLP publications for not endangered
248 languages (from hundreds to thousands annually).

249 Then, we observe the overlap of the language
250 coverage between the ELAR database, the ACL
251 Anthology and the language inventory of massive
252 multilingual datasets above-mentioned. Figure 3
253 shows the cross-over in a map. The very low over-
254 lapping was expected: from the ELAR inventory
255 (286), there are only 22 languages with at least

⁵We do not consider the extinct languages in our analysis

⁶We are aware that this was not an extended practice previ-
ously, but the Bender’s Rule (Bender, 2011) has remarked it
recently. Moreover, if a work does not specify which language
is working on, we can expect the target to be English or very
well-known established multilingual datasets.

256 one entry in the ACL Anthology (7.7%), and 12
257 languages included in at least one massive multi-
258 lingual NLP dataset (4.2%). Moreover, the geo-
259 localisation allows us to observe the potential of
260 these under-utilised resources in terms of represen-
261 tation for NLP research. Geographical areas such
262 as the Americas, Africa, South-East Asia or Aus-
263 tralia are better covered by language documentation
264 projects than NLP resources and studies. Regional
265 initiatives, such as Masakhane for Africa (Nekoto
266 et al., 2020), or AmericasNLP (Mager et al., 2021),
267 must look towards these still unexplored resources
268 for extending their language coverage.

269 4.4 Discussion

270 The NLP community is recently more aware of the
271 importance of language diversity in their research
272 (Bender, 2009, 2011). Typologically-diverse lan-
273 guage data allows to discuss results more broadly
274 and to identify potential flaws of the proposed meth-
275 ods in languages with typologically uncommon
276 grammatical properties and categories (O’Horan
277 et al., 2016; Ponti et al., 2019). Furthermore, it has
278 been pointed out that minority languages are in-
279 deed expected to exhibit unusual typological trends
280 and non-prototypical degrees of complexity (Trudg-
281 ill, 2011, 2010). Therefore, accessing and pro-
282 cessing databases of a wide sample of endangered
283 languages data would be beneficial for the NLP
284 agenda.

285 However, as we observed, this has not been a
286 priority. Why? We argue that this is mainly because
287 of the visibility, accessibility, and readability of the
288 data (from the NLP perspective):

289 **Visibility** Language documentation archives are
290 mostly known in the linguistic community. The
291 NLP community should look for data beyond the
292 usual repositories. Besides ELAR, other famous
293 repositories are the Archive of the Indigenous
294 Languages of Latin America (AILLA)⁷ from the
295 University of Texas, and The Language Archive
296 (TLA)⁸ from the Max Planck Institute for Psy-
297 cholinguistics.

298 **Accessibility** Most of the language documenta-
299 tion databases are open-source, but one often needs
300 to become a registered user in order to access the
301 materials allocated in the language archives. Fur-
302 thermore, some linguists block fully public access

⁷<https://ailla.utexas.org/>

⁸<https://archive.mpi.nl/tla/>

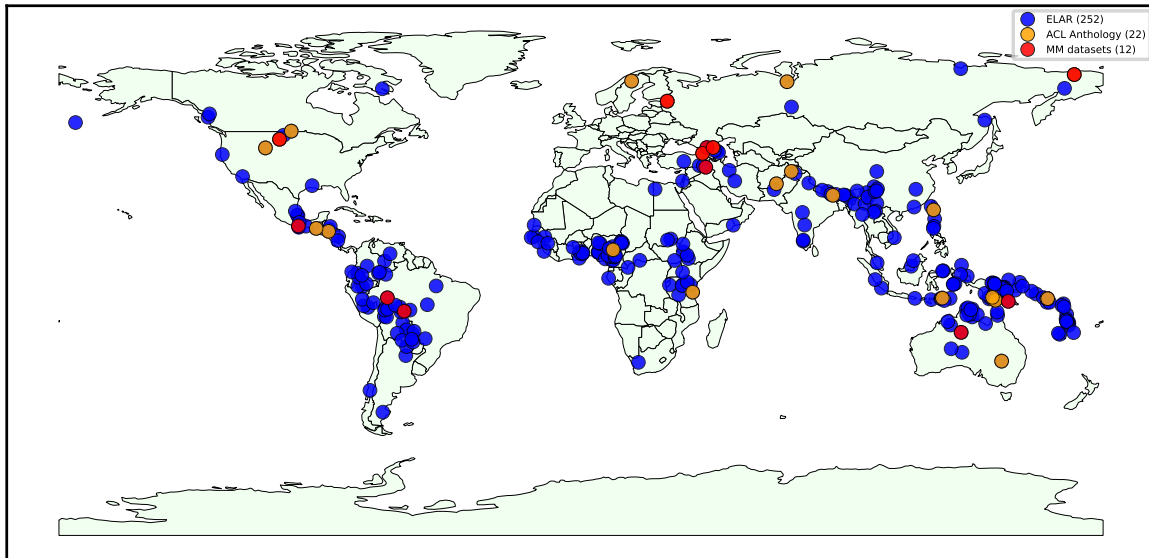


Figure 3: World map with languages in ELAR database and ACL Anthology. For the the present study, we only consider the languages of the ELAR database (570), whose names appear in *Glottolog* (version 4.4). This selection consists in 286 languages with geographical information. With this, 252 languages only belong to ELAR database (in blue); 22 languages belong to both ELAR database and ACL Anthology (in orange); and 12 languages belong to both ELAR database and massively multilingual (MM) datasets (Unimorph, Universal Dependencies and Tatoeba) (in red).

to their records as a way to protect speech community’s rights.

Readability Although most language documentation outputs video, audio and text files (plain texts or interlineal glossed texts, known as IGT), they are not labelled or processed for immediate use for NLP developments. If we observe the example in Figure 2, we can quickly identify potential resources for morphological segmentation and analysis, part-of-speech tagging, and machine translation. However, IGT is partially standardised, as not all the annotations follow the same label schema.

In sum, NLP is not taking advantage of all the resources potentially available for different applications. Moreover, from the three previously explained factors, readability is the hardest to overcome. One of our takes in this paper is to push the NLP community to focus more on the parsing and processing of the already published data, which is unlikely to be modified, unfortunately⁹. For instance, there should be paid more attention to IGT parsing research (Lewis and Xia, 2010; Round et al., 2020) or to the establishment of a more

⁹Most of the language documentation projects that are published might do not have extra funding allocated for any update, or new funding will be required for the job.

universally-readable IGT schema (Palmer and Erk, 2007). All this is complementary to the last point of Section 3, as we expect that, ideally, future deliverables of documentation projects could consider the annotation schema and resources that are more easily readable for NLP research.

5 CLD²: Computational Language Documentation and Development

Computational linguistics and language documentation share not only the assumption that technology plays an important role in the design and development of language-related projects, but also a crucial concern about language endangerment and loss. This concern is obvious from the perspective of language documentation, in the sense that it assumes itself as a response to language endangerment Himmelmann (2006, 5). A similar shift towards minority languages can be found in contemporary approaches to computational linguistics. Berment (2002) regrets that less than 1% of the world’s languages have been correctly “computerised”. That is, for Berment (2002), the fact that 99% of the world’s languages lack computational tools (NLP tools as spell-checking or machine translation) requires immediate attention. Since the seminal article by Krauss (1992), lan-

352 guage endangerment and language dormancy is a
353 major concern for both current language documenta-
354 tion and computational linguistics.

355 This paper takes the shared interest in linguistic
356 diversity found in language documentation and
357 computational linguistics further by proposing a
358 paradigm that assumes an intense and multifaceted
359 interaction between the two: Computational Lan-
360 guage Documentation and Development (CLD²).
361 CLD² assumes, following (Berment, 2002), that
362 “computerisation” should be understood as one
363 main task in language documentation and, at the
364 same time, proposes a basic protocol to carry out
365 this task. This basic protocol is based on a straight-
366 forward idea according to which any documenta-
367 tion project, in addition to its customary outcomes
368 (audio and video recordings, transcriptions, mor-
369 phological parsing and glossing, and free transla-
370 tions), should include NLP-friendly annotated data
371 as its deliverables:

- 372 1. A digital monolingual and parallel corpora¹⁰
373 ideally taken from a specific domain or dis-
374 course that is relevant for the language speaker
375 community;
- 376 2. A public representative set of sentences anno-
377 tated in universal frameworks for morphology
378 and syntax, such as Universal Morphology
379 (McCarthy et al., 2020) and Universal Depen-
380 dencies (Nivre et al., 2020)¹¹, which are well-
381 known in the NLP field; and
- 382 3. A communication describing the main char-
383 acteristics of the released Universal Depen-
384 dencies (Nivre et al., 2020) treebank and Uni-
385 versal Morphology (McCarthy et al., 2020)
386 dataset, so that NLPers can understand the
387 particularities and challenges of the data.

388 We attempt then to draw documentary and com-
389 putational linguists’ attention towards the potential-
390 ities of a more integral and systematic collabora-
391 tion between them. On the one hand, field linguists
392 may get involved in creating relevant products from
393 the NLP perspective (e.g. preparing representative

¹⁰Translations paired with English or another relevant lan-
guage spoken in the specific region, such as Spanish in Latin
America.

¹¹The identification of syntax dependencies and their an-
notation is not common in language documentation projects.
However Croft et al. (2017) have argued that the UD scheme
shares crucial principles with typological research. Indeed,
research on linguistic typology may benefit from the develop-
ment of an annotation scheme like UD and vice-versa.

treebanks taking as a starting point their own data).
On the other hand, NLPers can get involved in the
development of processes and protocols that may
contribute to the transformation of linguistic data
of the traditional sort into formats that may support
NLP developments.

According to Forcada (2006, 1), one feature for
a language to be considered as a minor is the few
to zero availability of machine-readable resources.
There are features such as the number of speakers
or literacy speakers that may support the definition
of a minor language in a general overview, but we
want to emphasise the computational perspective in
Forcada’s statement. Dictionaries, translated text
or annotated corpora, that are currently part of a
standard language documentation process, are in-
stances of machine-readable data. We consider that
linguistic corpora are insufficient to disentangle the
relationship between a language and its character-
isation as a minor. We claim the need to develop
more multiple resources to support a consistent re-
vitalisation of the language. However, we do not
mean that all language documentation processes
should include a massive technology development
by itself. The magnitude of such a project would
be cost-prohibitive. Nevertheless, we have iden-
tified some elements that might be included in a
documentation process that could drive a “comput-
erisation” effect in the studied language.

We want to emphasise the development of mul-
tipurpose linguistic databases, specifically aiming
at language technologies, whose implementation
will not radically increment the amount of expected
work for the linguist. Language technologies are
purpose-specific programmes that try to address
language-related tasks from spell- or grammar-
checking to automatic machine translation. Based
on such databases, NLPers and field linguists may
work together to develop NLP toolkits for minority
languages. An NLP Toolkit is a set of different
tools made to computerise a language fully. We
then take inspiration from the Basic Language Re-
source Kit (Krauwer, 2003) and also consider the
current state-of-the-art methods in NLP, such as
transfer learning. With transfer learning protocols,
especially multilingual pretraining (Lauscher et al.,
2020; Ebrahimi and Kann, 2021), CLD² projects
might automatise learning tasks by taking advan-
tage of larger amounts of multilingual data and
tools. A learning task in this context may refer to a
specific NLP or functionality, such as a dependency

445 parser, which has been trained to learn how to parse
446 the syntax in a textual sentence. Finally, we list the
447 main tools that such basic toolkits could have:

- 448 1. Morphological tools: such as morphological
449 analysis, to determine the base form or lemma
450 of an inflected word and its morphological fea-
451 tures; morphological segmentation, to identify
452 the canonical or surface morphemes (Mager
453 et al., 2020); and morphological reinflection
454 (Pimentel et al., 2021), which exploits Uni-
455 Morph data. Morphological knowledge is
456 usually crafted in language documentation
457 projects (see Figure 2), so these deliverables
458 could be the most manageable.
- 459 2. Spell-checker: to detect and automatic cor-
460 rect of spelling errors. Dictionary-based spell-
461 checkers can be easily retrieved from a docu-
462 mentation project with a lexicon as an output,
463 whereas rule-based ones can be adapted from
464 a finite-state morphological analyser. Data-
465 driven spell-checking is also possible to de-
466 velop from monolingual data only.
- 467 3. Syntactic parser: to analyse the relationships
468 between the words and phrases that compose
469 a text. A dependency syntax parser can be
470 developed using UD annotated data, and is
471 also benefited for transfer learning and pre-
472 training approaches (Lauscher et al., 2020).
473 Current language documentation projects do
474 not usually focus on this kind of annotation,
475 but we emphasise that it might be relevant for
476 research not only on NLP but also in linguistic
477 typology (Croft et al., 2017).
- 478 4. Part-of-Speech tagger and Named Entity
479 Recognition: both tasks are sequence taggers,
480 and are two of the tasks that have been bene-
481 fitted the most from multilingual pretraining,
482 and few- or zero-shot learning (Lauscher et al.,
483 2020; Ebrahimi and Kann, 2021). POS tag-
484 ging could be easily adapted from the cur-
485 rent glossing annotation, whereas NER anno-
486 tation can be quickly extended or marked in
487 the glosses.

488 Besides these tools, further developments that
489 can be achieved for endangered languages, such as
490 machine translation, are very appealing. However,
491 we also need to point out that, despite the progress
492 of the pretraining approaches and the use of few
493 labelled examples, a translation system (or other

kinds of NLP tools) should not be deployed with
494 low-quality outputs, as it can mislead the user. Lim-
495 itations of their usage should be assessed according
496 to the annotated data used and the purpose of the
497 systems.
498

6 Conclusion 499

CLD² calls for an enrichment of language documen-
500 tation projects by means of incorporating compo-
501 nents, outcomes and methods from NLP research,
502 as a strategy to promote the computarisation and
503 revitalisation of minority languages. This paper
504 shows that most of the interactions between com-
505 putational linguistics and language documentation
506 are framed as software developments that facilitate
507 the various processes involved in documenting a
508 language. The potential contributions of language
509 documentation and language repositories to NLP
510 research are under-exploited and deserve urgent at-
511 tention from the NLP community. At the same time
512 field linguists may also incorporate into the out-
513 comes of their projects, data crafted into paradigms
514 that can be automatically used for NLP develop-
515 ments (Universal Dependencies and/or Universal
516 Morphology, for instance).
517

This will benefit not only language documenta-
518 tion and computational linguistics scholars but also
519 typologists and speech communities, as research
520 in NLP has recently paid some attention to linguis-
521 tic typology as a substantial source of linguistics
522 knowledge to improve performance in different al-
523 gorithms and technologies (O’Horan et al., 2016;
524 Ponti et al., 2019). Indigenous communities, in
525 turn, are highly enthusiastic about the computer-
526 isation of their languages as a political strategy
527 that vindicates their languages and demonstrates
528 that they are as valuable as major European lan-
529 guages. CLD² can significantly contribute to this
530 aim by promoting productive exchanges among
531 field linguists, NLP researchers and members of in-
532 digenous communities as part of multi-component
533 projects that put language revitalisation at their
534 core.
535

References 536

Antonios Anastasopoulos, Christopher Cox, Graham
537 Neubig, and Hilaria Cruz. 2020. *Endangered lan-
538 guages meet Modern NLP*. In *Proceedings of
539 the 28th International Conference on Computa-
540 tional Linguistics: Tutorial Abstracts*, pages 39–45.
541

542	Barcelona, Spain (Online). International Committee for Computational Linguistics.	597
543		598
544	Antti Arppe, Jeff Good, Mans Hulden, Jordan Lachler, Alexis Palmer, and Lane Schwartz, editors. 2017. <i>Proceedings of the 2nd Workshop on the Use of Computational Methods in the Study of Endangered Languages</i> . Association for Computational Linguistics, Baltimore, Maryland, USA.	599
545		600
546		601
547		602
548		603
549		604
550	Antti Arppe, Jeff Good, Mans Hulden, Jordan Lachler, Alexis Palmer, Lane Schwartz, and Miikka Silverberg, editors. 2019. <i>Proceedings of the 3rd Workshop on the Use of Computational Methods in the Study of Endangered Languages Volume 1 (Papers)</i> . Association for Computational Linguistics, Honolulu.	605
551		606
552		607
553		608
554		609
555		610
556		611
557	Peter K Austin. 2006. Data and language documentation. <i>Essentials of language documentation</i> , 178:87.	612
558		613
559	Peter K. Austin. 2010. Communities, ethics and rights in language documentation. In Peter K. Austin, editor, <i>Language documentation and description</i> , volume 7, pages 34–54. London: School of Oriental and African Studies.	614
560		615
561		616
562		617
563		618
564	Emily M. Bender. 2009. Linguistically naïve != language independent: Why NLP needs linguistic typology. In <i>Proceedings of the EACL 2009 Workshop on the Interaction between Linguistics and Computational Linguistics: Virtuous, Vicious or Vacuous?</i> , pages 26–32, Athens, Greece. Association for Computational Linguistics.	619
565		620
566		621
567		622
568		623
569		624
570		625
571	Emily M Bender. 2011. On achieving and evaluating language-independence in nlp. <i>Linguistic Issues in Language Technology</i> , 6(3):1–26.	626
572		627
573		628
574	Vincent Berment. 2002. Several directions for minority languages computerization. In <i>COLING 2002: The 17th International Conference on Computational Linguistics: Project Notes</i> .	629
575		630
576		631
577		632
578	Paul Boersma and David Weenink. 2001. Praat, a system for doing phonetics by computer. <i>Glott International</i> , 5(9/10):341–345.	633
579		634
580		635
581	William Croft, Dawn Nordquist, Katherine Looney, and Michael Regan. 2017. Linguistic typology meets universal dependencies. In <i>TLT</i> , pages 63–75.	636
582		637
583		638
584	Abteen Ebrahimi and Katharina Kann. 2021. How to adapt your pretrained multilingual model to 1600 languages. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 4555–4567, Online. Association for Computational Linguistics.	639
585		640
586		641
587		642
588		643
589		644
590		645
591		646
592	Mikel Forcada. 2006. Open source machine translation: an opportunity for minor languages. In <i>Proceedings of the Workshop “Strategies for developing machine translation for minority languages”</i> , LREC, volume 6, pages 1–6.	647
593		648
594		649
595		650
596		651
	Jeff Good, Julia Hirschberg, and Owen Rambow, editors. 2014. <i>Proceedings of the 2014 Workshop on the Use of Computational Methods in the Study of Endangered Languages</i> . Association for Computational Linguistics, Baltimore, Maryland, USA.	652
		653
		654
		655
		656
		657
		658
		659
		660
		661
		662
		663
		664
		665
		666
		667
		668
		669
		670
		671
		672
		673
		674
		675
		676
		677
		678
		679
		680
		681
		682
		683
		684
		685
		686
		687
		688
		689
		690
		691
		692
		693
		694
		695
		696
		697
		698
		699
		700

652	Findings of the AmericasNLP 2021 shared task on open machine translation for indigenous languages of the Americas. In <i>Proceedings of the First Workshop on Natural Language Processing for Indigenous Languages of the Americas</i> , pages 202–217, Online. Association for Computational Linguistics.	710
653		711
654		712
655		
656		
657		
658	Max Planck Institute for Psycholinguistics. 2021. <i>ELAN (Version 6.2)</i> . The Language Archive, Nijmegen. https://archive.mpi.nl/tla/elan .	
659		
660		
661	Arya D. McCarthy, Christo Kirov, Matteo Grella, Amrit Nidhi, Patrick Xia, Kyle Gorman, Ekaterina Vylomova, Sabrina J. Mielke, Garrett Nicolai, Miikka Silfverberg, Timofey Arkhangelskiy, Natalya Krizhanovsky, Andrew Krizhanovsky, Elena Klyachko, Alexey Sorokin, John Mansfield, Valts Ernštreits, Yuval Pinter, Cassandra L. Jacobs, Ryan Cotterell, Mans Hulden, and David Yarowsky. 2020. <i>UniMorph 3.0: Universal Morphology</i> . In <i>Proceedings of the 12th Language Resources and Evaluation Conference</i> , pages 3922–3931, Marseille, France. European Language Resources Association.	
662		
663		
664		
665		
666		
667		
668		
669		
670		
671		
672		
673	Wilhelmina Nekoto, Vukosi Marivate, Tshinondiwa Matsila, Timi Fasubaa, Taiwo Fagbohunge, Solomon Oluwole Akinola, Shamsuddeen Muhammad, Salomon Kabongo Kabenamualu, Salomey Osei, Freshia Sackey, Rubungo Andre Niyongabo, Ricky Macharm, Perez Ogayo, Orevaoghene Ahia, Musie Meressa Berhe, Mofetoluwa Adeyemi, Masabata Mokgesi-Selinga, Lawrence Okegbemi, Laura Martinus, Kolawole Tajudeen, Kevin Degila, Kelechi Ogueji, Kathleen Siminyu, Julia Kreutzer, Jason Webster, Jamiil Toure Ali, Jade Abbott, Iroro Orife, Ignatius Ezeani, Idris Abdulkadir Dangana, Herman Kamper, Hady Elshahar, Goodness Duru, Ghollah Kioko, Murhabazi Espoir, Elan van Biljon, Daniel Whitenack, Christopher Onyefuluchi, Chris Chinenye Emezue, Bonaventure F. P. Dossou, Blessing Sibanda, Blessing Bassey, Ayodele Olabiyi, Arshath Ramkilowan, Alp Öktem, Adewale Akinfaderin, and Abdallah Bashir. 2020. <i>Participatory research for low-resourced machine translation: A case study in African languages</i> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 2144–2160, Online. Association for Computational Linguistics.	
674		
675		
676		
677		
678		
679		
680		
681		
682		
683		
684		
685		
686		
687		
688		
689		
690		
691		
692		
693		
694		
695		
696		
697	Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Jan Hajič, Christopher D. Manning, Sampo Pyysalo, Sebastian Schuster, Francis Tyers, and Daniel Zeman. 2020. <i>Universal Dependencies v2: An evergrowing multilingual treebank collection</i> . In <i>Proceedings of the 12th Language Resources and Evaluation Conference</i> , pages 4034–4043, Marseille, France. European Language Resources Association.	
698		
699		
700		
701		
702		
703		
704		
705		
706	Helen O’Horan, Yevgeni Berzak, Ivan Vulić, Roi Reichart, and Anna Korhonen. 2016. <i>Survey on the use of typological information in natural language processing</i> . In <i>Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers</i> , pages 1297–1308, Osaka, Japan. The COLING 2016 Organizing Committee.	710
707		711
708		712
709		
	Alexis Palmer and Katrin Erk. 2007. <i>IGT-XML: An XML format for interlinearized glossed text</i> . In <i>Proceedings of the Linguistic Annotation Workshop</i> , pages 176–183, Prague, Czech Republic. Association for Computational Linguistics.	713
		714
		715
		716
		717
	Tiago Pimentel, Maria Ryskina, Sabrina J. Mielke, Shijie Wu, Eleanor Chodroff, Brian Leonard, Garrett Nicolai, Yustinus Ghanggo Ate, Salam Khalifa, Nizar Habash, Charbel El-Khaissi, Omer Goldman, Michael Gasser, William Lane, Matt Coler, Arturo Oncevay, Jaime Rafael Montoya Samame, Gema Celeste Silva Villegas, Adam Ek, Jean-Philippe Bernardy, Andrey Shcherbakov, Aziyana Bayyr-ool, Karina Sheifer, Sofya Ganieva, Matvey Plugaryov, Elena Klyachko, Ali Salehi, Andrew Krizhanovsky, Natalia Krizhanovsky, Clara Vania, Sardana Ivanova, Aelita Salchak, Christopher Straughn, Zoey Liu, Jonathan North Washington, Duygu Ataman, Witold Kieraś, Marcin Woliński, Totok Suhardijanto, Niklas Stoehr, Zahroh Nuriah, Shyam Ratan, Francis M. Tyers, Edoardo M. Ponti, Grant Aiton, Richard J. Hatcher, Emily Prud’hommeaux, Ritesh Kumar, Mans Hulden, Botond Barta, Dorina Lakatos, Gábor Szolnok, Judit Ács, Mohit Raj, David Yarowsky, Ryan Cotterell, Ben Ambridge, and Ekaterina Vylomova. 2021. <i>Sign-morphon 2021 shared task on morphological reflection: Generalization across languages</i> . In <i>Proceedings of the 18th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology</i> , pages 229–259, Online. Association for Computational Linguistics.	718
		719
		720
		721
		722
		723
		724
		725
		726
		727
		728
		729
		730
		731
		732
		733
		734
		735
		736
		737
		738
		739
		740
		741
		742
		743
		744
	Edoardo Maria Ponti, Helen O’Horan, Yevgeni Berzak, Ivan Vulić, Roi Reichart, Thierry Poibeau, Ekaterina Shutova, and Anna Korhonen. 2019. <i>Modeling language variation and universals: A survey on typological linguistics for natural language processing</i> . <i>Computational Linguistics</i> , 45(3):559–601.	745
		746
		747
		748
		749
		750
	Erich Round, Mark Ellison, Jayden Macklin-Cordes, and Sacha Beniamine. 2020. <i>Automated parsing of interlinear glossed text from page images of grammatical descriptions</i> . In <i>Proceedings of the 12th Language Resources and Evaluation Conference</i> , pages 2878–2883, Marseille, France. European Language Resources Association.	751
		752
		753
		754
		755
		756
		757
	Frank Seifart, Nicholas Evans, Harald Hammarström, and Stephen C. Levinson. 2018. Language documentation twenty-five years on. <i>Language</i> , 94(4e):324–345.	758
		759
		760
		761
	Summer Institute of Linguistics. 2021a. <i>Field linguist’s Toolbox (Version 1.6.4)</i> . http://www.fieldlinguiststoolbox.org/?i=1 .	762
		763
		764
	Summer Institute of Linguistics. 2021b. <i>Fieldworks (Version 9.0)</i> . https://software.sil.org/fieldworks/ .	765
		766

- 767 Jörg Tiedemann. 2020. [The tatoeba translation chal-](#)
768 [lenge – realistic data sets for low resource and multi-](#)
769 [lingual MT](#). In *Proceedings of the Fifth Conference*
770 *on Machine Translation*, pages 1174–1182, Online.
771 Association for Computational Linguistics.
- 772 Peter Trudgill. 2010. Contact and sociolinguistic typol-
773 ogy. In Raymond Hickey, editor, *The Handbook of*
774 *Language Contact*, pages 299–319. Oxford: Wiley-
775 Blackwell.
- 776 Peter Trudgill. 2011. *Sociolinguistic Typology: social*
777 *determinants of linguistic complexity*. Oxford: Ox-
778 ford University Press.
- 779 Daan van Esch, Ben Foley, and Nay San. 2019. [Fu-](#)
780 [ture directions in technological support for language](#)
781 [documentation](#). In *Proceedings of the 3rd Work-*
782 *shop on the Use of Computational Methods in the*
783 *Study of Endangered Languages Volume 1 (Papers)*,
784 pages 14–22, Honolulu. Association for Computa-
785 tional Linguistics.
- 786 Anthony Woodbury. 2011. Language documentation.
787 In Peter Austin and Julia Sallabank, editors, *Hand-*
788 *book of Endangered Languages*, The Cambridge
789 Handbook of Endangered Languages, pages 159–
790 186. Cambridge: Cambridge University Press.
- 791 Anthony C Woodbury. 2003. Defining documentary
792 linguistics. *Language documentation and descrip-*
793 *tion*, 1(1):35–51.

794 **A AES status for massively multilingual** 795 **datasets**

AES status	Tatoeba	Unimorph	UD
not endangered	164	60	52
threatened	71	25	16
shifting	44	17	16
moribund	11	4	2
nearly extinct	7	4	1
extinct	24	17	11

Table 1: Agglomerated Endangerment Status (AES) (Seifart et al., 2018) statistics for MM databases (Tatoeba, Unimorph and Universal Dependencies).