Please cite as: List, Johann-Mattis and Sims, Nathanial A. (2019): Towards a sustainable handling of inter-linear-glossed text in language documentation. [Preprint under review. Not peer-reviewed]

Towards a sustainable handling of inter-linear-glossed text in language documentation

JOHANN-MATTIS LIST*, Max Planck Institute for the Science of Human History, Germany

NATHANIEL A. SIMS*, The University of California, Santa Barbara, America

Efforts on language documentation have been increasing in the past. While the amount of digital data of the world's languages is increasing, only a small amount of the data is sustainable, since data reuse is often exacerbated by idiosyncratic formats and a negligence of standards that could help to increase the comparability of linguistic data. The sustainability problem is nicely reflected in the current practice of handling inter-linear-glossed text, one of the crucial resources produced in language documentation. Although large collections of glossed texts have been produced so far, the current practice of data handling greatly exacerbates the reuse of data. In order to address this problem, we propose a first framework for the computer-assisted, sustainable handling of inter-linear-glossed text resources. Building on recent standardization proposals for word lists and structural datasets, combined with state-of-the-art methods for automated sequence comparison in historical linguistics, we show how our workflow can be used to lift a collection of inter-linear-glossed Qiang texts (an endangered language spoken in Sichuan, China), and how the lifted data can assist linguists in their research.

CCS Concepts: • Applied computing \rightarrow Language translation.

Additional Key Words and Phrases: Sino-Tibetan, inter-linear-glossed text, computer-assisted language comparison, standardization, Qiang

ACM Reference Format:

Johann-Mattis List and Nathaniel A. Sims. 2019. Towards a sustainable handling of inter-linear-glossed text in language documentation. 1, 1 (November 2019), 14 pages. https://doi.org/+++

1 INTRODUCTION

With many of the world's spoken languages being threatened by extinction, efforts on language documentation have been increasing in the past, as reflected in a constantly growing amount of various resources, ranging from short grammatical sketches, via short wordlists, up to extensive dictionaries, detailed grammars, and corpora in various forms and formats. Depending on the original interests of the researchers, but also on the funding upon which scholars base their research, language documentation follows a range of rather different purposes, as reflected in *typological surveys*, surveys oriented towards *historical language comparison*, *language revitalization efforts*, efforts reflecting *political motives* (such as the dialect surveys conducted by Chinese scholars in the 1950s [23]), and efforts reflecting *missionary goals* (such as surveys conducted by religious organizations).

*Both authors contributed equally to this research.

Authors' addresses: Johann-Mattis List, list@shh.mpg.de, Max Planck Institute for the Science of Human History, Kahlaische Str. 10, Jena, Thüringen,
 07745, Germany; Nathaniel A. Sims, nsims@ucsb.edu, The University of California, Santa Barbara, 3432 University Drive, Santa Barbara, America.

- Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not
 made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components
 of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to
 redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
- ⁴⁹ © 2019 Association for Computing Machinery.

50 Manuscript submitted to ACM

While the amount of digitally available data on the worlds' languages is steadily increasing, with more and more languages being documented, only a very small proportion of the language resources that are produced account for *sustainability*. Sustainability – in the context of scientific research – is hereby understood as a resource that complies to the principles of FAIR data as outlined by Wilkinson et al. [28]: resources should be *findable*, *accessible*, *interoperable*, and *reusable*.

Due to the different objectives of scholars working in the field of language documentation, we face a situation where specifically the re-usability of language resources is largely exacerbated. This starts from the fact that some resources are still only produced in print, and even if they are produced digitally, they are rarely *machine-readable*, as they are shared in form of PDF documents, which cannot be converted to computer-friendly resource formats, such as spreadsheet tables or lightweight databases. Even if the data are shared in tabular, basically machine-readable form, they are often not *interoperable*, because they lack *standardization*, and in order to access one specific resource, huge efforts are needed in order to lift the data to a level where they could be easily reused in computer-based or computer-assisted frameworks oriented towards *cross-linguistic comparison*.

One might argue that it is not the primary purpose of language resources, such as, for example, dictionaries, to be parsed by a computer application, but rather by humans who want, for example, to teach an endangered language in school. But it is important to keep in mind that even humans tend to prefer digital dictionaries over resources written in prose and printed only on paper, and the easier a given resource can be searched, the more lasting will be its impact, specifically among younger generations. In addition, the current lack of sustainability of linguistic resources makes it very difficult, if not even impossible at times, to develop targeted applications in the field of *natural language processing* (NLP), specifically for endangered and poorly documented languages.

- 78 Most NLP applications are not only "blind" to language-specific aspects, since - specifically for poorly documented 79 languages - the resources are lacking, but additionally - since large language resources used for the study of big 80 languages (English, Chinese) are often of poor quality – ignore linguistic knowledge to a large degree. In order to 81 side-step the problem of lack of documentation, researchers in NLP now have started to try and impute missing data 82 83 from cross-linguistic typological databases, given that the data-hungry business of NLP can often not cope with datasets 84 small in size [25]. In fact, prediction (or retrodicton) of missing features can indeed be useful, not only in the typological 85 sphere but also for the lexicon, as scholars report in an ongoing experiment of word prediction of Kho-Bwa languages 86 (Tibeto-Burman) [2]. But in order to allow for a successful integration of linguistic resources that could help NLP 87 applications to improve it approaches, specifically also when dealing with smaller and endangered languages, it is 88 89 important to improve on the general sustainability in language documentation. 90
- While some steps in this direction have been already undertaken in the future, with new standards being proposed for the handling of word lists and structural data in historical linguistics and language typology [5], or initial frameworks having been developed for the handling of rhyme annotation [21], we want to draw the attention to *inter-linear-glossed text* as one of the crucial resources produced by language documentation efforts. Although large collections of interlinear-glossed text have been produced so far, and scholars use it across all subfields of linguistics, including opposing camps, the current usage practice largely lacks sustainability, being – despite its formal nature – mostly oriented towards manual digestion.

In the following, we want to propose a first framework for the computer-assisted, sustainable handling of inter linear-glossed text (IGT). After discussing our general strategy to increase the sustainability of linguistic resources,
 which follows closely the recommendations of the Cross-Linguistic Data Formats initiative (https://cldf.clld.org, [5],
 Section 2), we will present a detailed (but still rudimentary) proposal for the standardization of inter-linear-glossed
 Manuscript submitted to ACM

53

54

55

56

57 58

59

60

61

62 63

64

65

66

67 68

69

70

71

72 73

74

75

76

text (Section 2), and illustrate, how this framework can be successfully applied to lift the data of a small corpus of Qiang texts (Section 4), an endangered language, spoken in the northwest part of Sichuan Province in China [8, pp. 1-5]. We conclude by discussing further application possibilities for our framework and point to problems that need to be addressed in the nearer future (Section 5).

2 SUSTAINIBILITY OF LINGUISTIC RESOURCES

105 106

107

108

109 110

111 112

113

114

115 116

117

118

Given that linguists create linguistic resources with different purposes in mind, the resources – specifically those on endangered and low-resource languages – differ widely. While it is clear that there are generally different type of resources, and that not all linguists plan to create a dictionary of the languages they want to document, the problem does not lie in the broad categories (dictionary, grammar, text corpus, wordlist), but in the way in which the broad categories most scholars would agree upon are created and shared.

As an example, consider the seemingly simple problem of creating comparative wordlists for a couple of languages of 119 120 interest. While the basic format, according to the standard notion of the linguistic sign, would require a triple of language, 121 concept, and form, we find standardization issues in all three of these basic components. Language names, although 122 referring to the same language variety, may vary widely, both for historical reasons (e.g., because language names in 123 the past may have had a derogatory attitude), but also for reasons that are not always made explicit in published studies. 124 Concepts are usually denoted with help of *elicitation glosses* [19], but elicitation glosses that are intended to denote 125 126 the same concepts vary widely, even if the same language for elicitation has been used [14]. Word forms, finally, are 127 the least standardized of all items one encounters in wordlists, given that scholars usually do not provide phonetic 128 transcriptions, but rather turn to orthographies, where available, or make use of quasi-phonological transcriptions that 129 130 they consider more convenient for typing, but which are rarely explained with respect to the intended phonetic values. 131

While the problems may seem severe, initial standardization efforts have been done in the past years, and they have also shown that is possible to successfully enhance existing datasets, by applying a procedure that could be called *retro-standardization*. Instead of changing existing resources manually, semi-automatically, or automatically, retro-standardization adds several annotation layers to existing datasets that allow for an easy conversion of the original data into a format that is machine-readable and cross-linguistically comparable.

These efforts have been most prominently propagated by the Cross-Linguistic Data Formats initiative (CLDF, 138 https://cldf.clld.org, [5]). The basic idea of CLDF is to address comparability problems involving linguistic data by 139 introducing reference catalogs, i.e. meta-databases that offer information for those entities which are crucial for cross-140 141 linguistic comparison. As the most prominent example, the Glottolog catalog (https://glottolog.org) offers information 142 on language names, geographic locations, and basic genealogical classifications [6]. In order to make sure that it is clear 143 which languages a given resource documents, all that needs to be done is to list the Glottocodes, the identifiers provided by 144 Glottolog, for each language that occurs in the resource. Similarly, the Concepticon project (https://concepticon.clld.org, 145 146 [17]), offers standard identifiers for elicitation glosses and links existing concept lists to those identifiers in order to 147 illustrate the huge variation that can be encountered in concept elicitation. For word forms, the recent Cross-Linguistic 148 Transcription Systems initiative (CLTS, https://clts.clld.org, [16]) provides standard identifiers for speech sounds which 149 150 are themselves linked to different transcription systems and thus offer a convenient way to check if a given transcription 151 complies to the standard defined by a given system [1]. 152

CLDF reference catalogs do not stop with providing identifiers to which the original data could be linked. In addition, specific tools are provided that facilitate the process of linking. While identifying languages in Glottolog is already made easy by the web application, the Python API that comes along with it allows scholars proficient in Python programming Manuscript submitted to ACM

to use the data provided with Glottolog inside of Python scripts. Concepticon offers commandline tools that allow for 157 158 an automated mapping of elicitation glosses to the Concepticon identifiers in multiple languages, which can as well be 159 applied from within Python scripts. CLTS offers a range of strategies to normalize transcription data, specifically when 160 provided in the broad version of the IPA that is at the core of the reference catalog. Additionally, scholars can make use 161 162 of orthography profiles [24] that allow for a semi-automated conversion of transcriptions in a given resource into the 163 standards supported by CLTS. All in all, these tools, which are well-documented and also illustrated in several online 164 tutorials, greatly facilitate the process of retro-standardization [13]. 165

With respect to inter-linear-glossed text, the situation is still different. Although annotation tools exist, as, for 166 167 example provided by the Summer Institute of Linguistics' FieldWorks program (https://software.sil.org/fieldworks/), 168 their application is exacerbated by a lack of cross-platform support (with many tools working only on Windows 169 machines), but also by a large degree of freedom offered by the respective software. Since the majority of IGT is still 170 produced in research articles, and not in form of standardized databases, errors in the glossing procedure are still 171 172 rather common, as can be seen when checking a random resource provided by ODIN, the largest agglomeration of 173 inter-linear-glossed text examples taking from linguistic resources [10]. 174

Our strategy for working towards an increase of sustainability in language documentation, with a specific focus 175 on inter-linear-glossed text is two-fold, following the idea of retro-standardization, as it has been proposed by the 176 177 CLDF initiative. First, we want to increase scholar's awareness regarding available standards and the advantages of 178 using them. Second, we want to make it as easy as possible for scholars to produce their data in the way they know, 179 while encouraging them to open backdoors for quick retro-standardization of their data. The basic idea is to provide 180 initial standards that come close to the formats which scholars already use, but are strict enough to allow for a quick 181 182 processing by a machine. The advantage of such an approach is that data can be automatically checked for errors which 183 may be easily introduced in typing, while at the same time opening a door for quick retro-standardization with help of 184 computer tools which we will present in detail in the following sections. 185

186 187

193

3 PROPOSALS FOR STANDARDIZING INTER-LINEAR-GLOSSED TEXT

In the following, we will present our proposals for a flexible standardization framework of inter-linear-glossed text in detail. After briefly discussing the role that inter-linear-glossed text plays in language documentation, we will explain the basic ideas behind the CLDF initiative in more detail, and then present a workflow for the retro-standardization of resources that offer inter-linear-glossed text.

¹⁹⁴ 3.1 Inter-linear-glossed text ¹⁹⁵

Inter-linear-glossed text is a commonly used way of presenting the structure by which phrases in foreign languages 196 197 are built. The basic idea is to gloss each word of a phrase in a certain language by grammatical and lexical glosses in 198 order to elucidate how the respective language expresses a certain circumstance. Technically, IGT demands at least 199 two separators. First, words in the language that is being glossed need to be distinguished, which could be done by a 200 simple white-space character, wich is often represented by a tab-stop, in ordert to support a visual alignment of the 201 202 original text and the glosses. Second, all meaningful grammatical and lexical units, that is, the morphemes inside a 203 word need to be marked, which is usually done with the help of the dash character ("-"). Apart from this, there are 204 different rules to distinguish lexical from grammatical glosses. The most common way consists in writing grammatical 205 glosses in abbreviated form in capital letters, and providing a legend for the meaning of the abbreviations. Lexical 206 207 glosses are usually not standardized and simply follow the analysis of the researcher with respect to the utterance under 208 Manuscript submitted to ACM

Towards a sustainable handling of inter-linear-glossed text in language documentation

209 210

211

252

question. Table 1 provides an example of a piece of IGT in German along with the lexical and grammatical glosses and the translation.

212	Die Katze sitz-t auf den Matratz-en								
213	ARTIC.NM.SGL.F cat sit-3.SG on ARTIC.DT.PLR.F matress-PLR								
214	The cat sits on the matresses.								
215 216	Table 1. Simple example sentence of IGT in German.								
217									
218									
219									
220	Although there have been efforts to standardize IGT with respect to the usage of grammatical glosses, one can								
221	encounter a lot of variation with respect to the implementation of the principle. Scholars tend to provide their own								
222 223	abbreviations in the introduction or the appendix of the work, and they also tend to use their own transcription systems								
224	(if the language under question has no standardized orthography). Ideally, the information on the grammatical glosses								
225	and the transcription systems are exemplified in the studies providing IGT, but the fact that IGT is not following any								
226	strict principles – and is barely checked by computational methods for internal consistency – results in a large variation								
227	that makes it difficult to make actual use of large IGT collections such as the ones provided, for example, by the ODIN								
228 229	project [10].								
230	While it cannot be denied that there is a certain awareness of the problem of incomparability of IGT from a cross-								
231	linguistic perspective, with quite a few journals demanding IGT to follow the popular <i>Leipzig Glossing Rules</i> [3], the lack								
232	of a computer-assisted <i>testing</i> whether a given sample of IGT provided in an article or a database conforms to a given								
233	standard makes it extremely difficult to compare IGT corpora <i>across</i> the studies in which it was originally proposed.								
234 235									
235	Since most linguists digest IGT examples piece by piece, without expecting to use them for corpus studies or extended								
237	NLP applications. As a result, the majority of IGT corpora produced at the moment is largely incomparable and not								
238	amenable for quantitative comparison, at least not beyond the scope of the resource in which they were originally								
239	produced. This is extremely unfortunate, given the wealth of information that IGT could offer for cross-linguistic								
240	investigations. Although there are large resources of digitally available IGT, as it is provided, for example, by the								
241 242	PanGloss project (https://lacito.vjf.cnrs.fr/pangloss/), the Dictionaria project (https://dictionaria.clld.org), or the ODIN								
243	corpus [10], there is no way to unify the available resources in a common framework. This is a pity, since IGT offers –								
244	at least in theory – many possibilities for interesting analyses that could drastically increase the amount of resources								
245	that scholars who work on quantitative applications in NLP, historical linguistics, and linguistic typology have at their								
246	disposal. In cases where dictionaries are lacking, one could use larger IGT collections of the same language to construct								
247 248	wordlists for cross-linguistic comparison. Where grammatical surveys are lacking, IGT could help to extract structural								
240	features about a certain language. Finally, if the transcriptions in which IGT is shared were standardized, it could give								
250	hints not only to phoneme inventories but also to the potential usage frequency of the phonemes employed by a given								
251	language.								

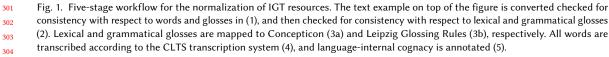
253 3.2 Workflow for retro-standardization of inter-linear-glossed text resources 254

255 Our workflow for the retro-standardization of inter-linear-glossed text is rather straightforward and seeks to standardize 256 those aspects of a given resource for which reference catalogs as propagated by the CLDF initiative are supported. A 257 minimal example of inter-linear-glossed text consists of two entities. First, there is a text that is divided into sentences, 258 which are themselves divided into phrases. Phrases again consist of a sequence of words which are themselves divided 259 260 Manuscript submitted to ACM

into *morphemes* (or *morphs*). Second, a sequence of glosses is aligned to the text, with each gloss providing lexical or
 grammatical semantic information for each morpheme.

While general rules for text glossing have long since been proposed[3], these rules only standardize the outer appearance of inter-linear morpheme glossing, while they do not provide any additional recommendations with respect to the way in which, for example, the text should be written, or which elicitation glosses should be used. Since, with the Concepticon project and the CLTS initiative, new reference catalogs are available by now, we think it is time to see to which degree these catalogs can be used to enrich the information that is provided in collections of inter-linear glossed text.

(1)	Word	Gloss	(2)	Morpheme	Lexical Gloss	Grammatical (
	Die	ARTIC.NM.SGL.F		Die		ART.NOM.SG
	Katze	cat		Katze	cat	
	sitz-t	sit-3.SGL		sitz	sit	
	auf	on		t		3.SG
	den	ARTIC.DT.PLR.F		auf	on	
	Matratze-n	matress-PLR		den		ART.DAT.PL.F
				Matratze	matress	
				n		PL
(3a)	Lex. Conce	pt Concepticon	(3b)	Gram. Concept	Leipzig Gloss	ing Rules
	cat	1208 CAT		ARTIC	ART	
	sit	1416 SIT		NDA	11014	
	511	1410 311		NM	NOM	
	on	1741 ABOVE		NM SGL	NOM SG	
					-	
	on	1741 ABOVE		SGL	SG	
(4)	on	1741 ABOVE	7 (5)	SGL	SG PL	
(4)	on matress	1741 ABOVE 105 MATRESS	(5)	SGL PLR 	SG PL 	
(4)	on matress Word	1741 ABOVE 105 MATRESS CLTS Transcription	(5)	SGL PLR Word	SG PL Cognacy	
(4)	on matress Word Die	1741 ABOVE 105 MATRESS CLTS Transcription d i:	(5)	SGL PLR Word d i:	SG PL Cognacy 1	
(4)	on matress Word Die Katze	1741 ABOVE 105 MATRESS CLTS Transcription d i: k a ts ə	(5)	SGL PLR Word d i: k a ts ə	SG PL Cognacy 1 2	
(4)	on matress Word Die Katze sitz-t	1741 ABOVE 105 MATRESS CLTS Transcription d i: k a ts ə s i ts + t	(5)	SGL PLR Word d i: k a ts ə s 1 ts + t	SG PL Cognacy 1 2 3 4	



Following the general idea of the CLDF initiative of linking resources to the major reference catalogs which have been proposed so far, our workflow towards a retro-standardization of IGT resources thus consists of the following five steps. In a first step, we *standardize* a given IGT resource by making sure that the basic principle of glossing is followed consistently. Starting from a digital IGT resource, we thus check that all *words* in a phrase have at least one *glossed complex* that explains them (1). In a second step, we make sure that each *morpheme* in a word is given a distinct *gloss* Manuscript submitted to ACM

Towards a sustainable handling of inter-linear-glossed text in language documentation

(be it grammatical or lexical) (2). In a third step, we try to extract *concept lists* for grammatical and lexical glosses, by creating a concordance of each pair of a morpheme and its corresponding gloss in the IGT resource. By automatically distinguishing lexical from grammatical elicitation glosses, this creates two concept lists, one grammatical concept list, and one lexical concept list (3). Having created the concept lists, we try to link the entries in the lexical concept list to the Concepticon resource, and the grammatical concept list to the abbreviations and additional instructions that are usually provided along with a given resource of IGT. In the future, we hope to be able to further link the grammatical glosses to reference catalogs similar to Concepticon, but devoted to abbreviations and elicitation glosses for grammatical concepts in linguistic resources (see, for example, the idea of creating a Grammaticon as a counterpart of the Concepticon by Haspelmath [7]). In a fourth step, we try to normalize the transcription system by linking each sound segment that occurs in a given IGT resource to the standard transcription systems (called B(road-coverage)IPA) proposed by the CLTS initiative (4). In a last step, we try to identify language-internal cognate words in the IGT resource by clustering all morphemes that show a certain degree of phonetic similarity and are glossed by the same elicitation gloss into the same word family (5).

Once having enriched a given IGT resource in this way, we can present the data in a combined form, in which each instance of the original IGT is accompanied by the additional information that we added during the retro-standardization process. To illustrate how this information can be successfully combined, we create a light-weight web-application in which scholars can *query* the resource for grammatical and lexical concepts, and word forms. Figure 1 illustrates this workflow in a schematic way.

4 APPLICATION EXAMPLE WITH DATA FROM QIANG (TIBETO-BURMAN)

In the following, we will illustrate how our workflow can be applied to a concrete IGT resource. The supplementary material provides all data and code needed to replicate the experiments we have carried out in this context, but since our work also includes steps of manual refinement, scholars may come to different results when following our example.

4.1 Materials: An inter-linear-glossed corpus of Qiang texts

Qiang 羌 (also called Rma) is a Tibeto-Burman language spoken by both ethnic Qiang and ethnic Tibetans in the mountainous area along the upper Min river 岷江 in the Rgnaba-Tibetan-Qiang Autonomous Prefecture of western Sichuan, China. Qiang is not a traditionally written language. It is an endangered language that is in many places being replaced by local varieties of Mandarin [4]. The present Qiang data come from a collection of texts from LaPolla and Huang's 2003 description of the Ronghong variety spoken in northwestern Mao County 茂 [9]. The grammar includes an appended six transcribed and annotated texts recorded by three different native speakers. The authors give a free translation into English and Chinese for the texts, but do not provide a line-by-line translation.

In order to make the data amenable for digital treatment, the texts were first digitized and stored in a simple text format which closely renders the format of the glossed text in the original PDF version of the resource, but uses tabstops as standard separators on the word level. In a second stage, these data were parsed into the basic input format currently required by our software package.

4.2 Methods: A Python package for IGT processing

The code needed to apply the workflow for the retro-standardization of IGT resources is provided in form of a small Python library (*pyigt*), available from the supplementary material accompanying this study. The code makes use of third-party libraries for a variety of tasks, specifically the LingPy Python library for quantitative tasks in historical Manuscript submitted to ACM linguistics (http://lingpy.org, [18]), which we use not only for data handling, but also for the automated detection of
 language-internal cognates [20, 22]. With respect to the design, our *pyigt* library resembles *PoePy*, a Python library for
 the quantitative handling of rhyme data (https://github.com/lingpy/poepy, [15, 21]). In the following, we will illustrate
 all steps of our workflow in detail.

4.2.1 Input formats. The input format required for our workflow is a plain text file in tab-separated form, with the first line providing the column headers and the following lines representing each one phrase of a give IGT resource. The first column of this tabular data schema is reserved for a numerical identifier (ID), while the order of the remaining columns is arbitrary, following the header. Assuming that a given IGT resource needs to provide at least two separation levels for the phrase, our tool expects white-space as a word separator, and the dash character - as a morphem-level separator, both in the word forms (PHRASE) and in the glosses (GLOSS). In order to group phrases to sentences, an identifier for sentences should be submitted (SENTENCE_ID), and texts can be distinguished by supplying a text identifier (TEXT). Figure 2 shows the first lines of the IGT resource on Qiang.

ID	TEXT	SENTENCE_ID	PHRASE_ID	PHRASE	GLOSS
1	Text 1	1	1	zəp-le: n̯i-ke: pe-ji	earth-DEF:CL WH- INDEF:CL become-CSM
2	Text 1	1	2	qe¹lotşu-ʁa, mutu-la mujuqu zguə-zi	in.the.past-LOC heaven- LOC sun nine-CL
3	Text 1	1	3	we-i, zəp-le: ə-tɛhəqha-z- əi. mə na ylu	exist-HS earth-DEF:CL DIR-burn-CAUS-HS older.brother COM younger.sister
4	Text 1	2	4	jə-tş-ŋuəni, zuamə-фu o- zgu-ta	two-CL-TOP cypress-tree one-CL-LOC
5	Text 1	2	5	i-pi-χua-ηi, ĥo-mu-xtεu- wei. steke-ta mi pe¹zə-s	DIR-hide-because-ADV DIR-NEG-burn-HS later- LOC people raise(child)- NOM

Fig. 2. Data representation in the standard input format employed in the workflow.

4.2.2 Consistency checks on IGT data (1). Once the data is prepared in the format as specified in the preceding section, it can be directly parsed by our library and checked for inconsistencies. This check, which is often only done by eyeballing glossed text resources before publication, turned out to be very useful, since it helped us to identify a couple of inconsistencies in the digital version of the data, which were introduced during the process of digitization.

4.2.3 Creation of lexical and grammatical concordances (2). Once the data has passed the first stage of consistency check, lexical and grammatical concordances can be prepared. In this stage, our workflow checks additionally, if the glosses match also at the morpheme-level with the words in the resource. In addition, given that grammatical functions often appear in complexes (such as case, number, and genus in many European inflecting languages), this stage introduces a third separator on the level of the glosses, which is used to separate multiple grammatical functions from each other. While the Leipzig Glossing Rules recommend to use a dot for this purpose, the Qiang resource consistently used a colon for this purpose.

The computation of the grammatical and lexical concordances yielded a total of 302 distinct grammatical forms linked to 53 grammatical concepts, and as many as 968 lexical forms linked to 591 lexical concepts. The most frequently occurring grammatical form was the interjection [fia], which we found as many as 355 times in the data, and the Manuscript submitted to ACM

most frequently expressed grammatical meaning is represented by numerous directional prefixes (708 examples). The 417 418 most frequently occurring lexical form was [jə] "say", with 139 occurrences, and the most frequently expressed lexical 419 meaning turned out to be "one" with 206 examples (representing different forms). All in all, this analysis did not yield 420 any surprises, but it helped us to further eliminate problems in the glosses, as we could identify erroneous glosses that 421 422 go back to the process of digitization as well as spelling errors in the original resource. An example for a problem in 423 the digitization is the wrong rendering of the word uncle's as unclefls, which is due to the internal rendering of the 424 apostrophe character in the PDF copy of the grammar. We did not identify many obvious errors (e.g. in spelling) going 425 back to the original source itself, which shows that the resource was thoroughly prepared. An example for a spelling 426 error is the elicitation gloss "daugher" which occurs two times in the original data and obviously refers to "daughter". 427 428

4.2.4 Mapping lexical and grammatical concepts to reference catalogs (3). Having extracted lexical and grammatical
 concept lists, we can *map* the lexical concepts to the Concepticon reference catalog. To ease the mapping procedure, the
 Concepticon Python API offers an automated mapping routine that checks a given elicitation gloss in a resource against
 those elicitation glosses that have been used in the 275 resources that have so far been linked to the Concepticon. As a
 result, the process of concept mapping is greatly enhanced, and it did not take us much time to manually refine the
 automated mappings.

Having linked the lexical concepts to Concepticon has the advantage of enabling us to check to which degree the 438 concepts in the resource could be used in other applications. Word lists, for example, are important for historical 439 440 language comparison, but aggregating word lists from different resources is extremely tedious. Once different resources 441 are linked to the Concepticon reference catalog, however, aggregation is simple, since we can automatically check to 442 which degree different resources overlap with respect to the concepts they employ. Thus, of the 591 concepts reflected 443 in the Qiang resource, we find an overlap of 112 concepts compared to the comparative word list collection established 444 445 by Sagart et al. for their phylogenetic study on Sino-Tibetan languages [26]. A comparison with the concept list of 446 100 basic vocabulary items proposed by Morris Swadesh [27] shows that the Qiang resource only covers 56 of these 447 concepts. This information is crucial, as it can help scholars who seek to create comparative wordlists from different 448 resources to check quickly if the coverage across different datasets is high enough. 449

450 In a similar way, the grammatical concepts offer extremely valuable information, as they can give immediate hints 451 with respect to the grammatical categories which are expressed in a given language. Since no reference catalog for 452 elicitation glosses pointing to grammatical concepts has been established so far, we compared the grammatical concepts 453 in the resource with the list of abbreviations listed in the original resource. In a second step, we added the standard 454 455 abbreviations suggested by the Leipzig Glossing Rules to the grammatical concept list. While the Qiang resource mostly 456 coincided with the Leipzig Glossing Rules, we find a few interesting cases of divergence. Thus, while the abbreviation 457 PRS is used by LaPolla and Huang in order to refer to a prospective aspect suffix, the abbreviation refers to the present 458 tense in the Leipzig Glossing Rules. On the other hand, Lapolla and Huang use INDEF to refer to an indefinite marker, 459 460 while the Leipzig Glossing Rules suggest to abbreviate this as INDF. While these comparisons may seem pedantic, 461 they greatly exacerbate an automated comparison across resources. Furthermore, the similarity of abbreviations used 462 in different IGT resources but referring to completely different things shows that a careful comparison of linguistic 463 resources can only be done when referring to the original list of abbreviations. In order to guarantee the future 464 465 comparability of linguistic resources, we need a reference catalog for grammatical elicitation glosses, as well as general 466 efforts to advocate these standards when producing IGT resources. 467

468

								I	Puln	non	ic (Cons	ona	ants										
Place →		Labial				Coronal										Dorsal					Laryngeal			
↓ Manner	Bila	bial	Lab den		Ling lab		De	ntal	Alve	eolar		ato- eolar	Retr	oflex	Alveolo- palatal	Pal	atal	Ve	elar	Uv	ular	Pharyn / Epiglo		Glotta
Nasal	ņ	m		ŋ	'n	ņ			ņ	n n ^j			η	η	ň, n.	'n	ŋ	ŋ	ŋ	Ņ	N			
Stop	$\mathbf{p} \ \mathbf{p}^{\mathrm{h}}$	b	p		ţ	ď			t th	d			t	d		С	Ĵ	k k ^h	g	$\mathbf{q} \ \mathbf{q}^{\mathrm{t}}$	G	2		?
Sibilant affricate									ts ts ^h	dz	t∫		tş tş ⁿ	dz	te dz									
Non-sibilant affricate	pф	bβ	pf	þv			tθ	dð	tΘ	dğ	tı	dı.				сç	đ	kx	as	qχ		2ħ	21	2h
Sibilant fricative									S	z	S	3	ş	z	6 Z									
Non-sibilant fricative	φ	β	f	V	ĕ		θ	ð	Θ	ğ	Ļ	Ţr				Ç	j	x	x	χ	R	ħ	Ŷ	h ĥ
Approximant	фт						- Û]	Ţ	L			Ĵ	J		j	j	ů	щ					
Flap or tap		-Y		\mathbf{V}_{-}		1 2			ţ	1			ť	T							Ğ		3	
Trill		В				Ĩ			ŗ	r			UF	Ţr				1		Ŗ	R	Н	5	
Lateral affricate									tł	dłz			tŀ			сĄг		kĭ	GĻ					
Lateral fricative									4	-13			-l-			- Au	À.	Ļ	Ļ					
Lateral approximant					<u> </u>					1				L.		- Ą	A		L		Ē	-		
Lateral flap																	À.		Ľ					

Fig. 3. Consonant chart produced by the EDICTOR tool from the standardized transcriptions.

4.2.5 Standardizing transcriptions (4). As discussed in detail by Anderson et al. [1], the current linguistic practice of phonetic transcription bears not only many pitfalls, but can barely seen as reflecting a coherent standard. In order to standardize the transcription system employed in a given resource, it is important to identify all distinct sound segments in the data, which can at times be represented by more than just one transcription symbol. While this may sound trivial at first sight, the procedure can turn out to be very tedious, specifically in those cases where a consistent description of the transcription system employed in a given resource is missing.

What has turned out to be extremely helpful in retro-standardizing transcription systems so far is the application of *orthography profiles*, an idea proposed by Moran and Cysouw [24], which consists of a simple table, in which all *graphemes* in a given resource are contrasted with their standardized counterpart. While the original preparation of orthography profiles is tedious, the LingPy software package offers a convenient algorithm for their first creation which also tries to link the transcription symbols to the standard proposed by the CLTS initiative, and which we implemented in our workflow. Once an initial, automated orthography profile has been produced, it can be easily manually corrected.

When adjusting the original transcriptions, it turned out that we did not have to correct many of the transcriptions in the original data. The most notable deviations from the standard transcription system proposed by the CLTS reference catalog was the usage of a normal [h] in order to mark aspiration (which should be represented by a superscript [h]). In addition, we found that the authors often used the letter [a] instead of the letter [a] in order to denote an unrounded open back vowel, although the former variant is not described in the phonology section of the grammer. We also found instances where orthographical spelling was used instead of the phonetic transcriptions, as in the case of zz, which reflects – at least according to the phonological description in the grammar – to a voiced alveolar affricate [dz].

Figure 3 shows a classical IPA chart of all the consonants in the Qiang resource, which was automatically created from the standardized transcriptions with help of the EDICTOR (https://digling.org/edictor/, a web-based tool for the creation of etymological dictionaries [12], which supports the standards proposed by the CLTS reference catalog. As can be seen from this chart, the data does not provide any surprises, but it helps to evaluate a given transcription system and to compare the one we extracted from the glossed texts with the one reported in the grammar.

520 Manuscript submitted to ACM

527

528

529

530 531

541 542 543

544

545

546

547 548

549

550

551

552 553

554

555

556

557 558

559

560

561 562

572

4.2.6 Identifying language-internal cognates (5). Once created and manually corrected, the orthography profile allows us to convert the original transcriptions into the standardized transcription system and segment the data into sound segments at the same time. This has the great advantage that the data in this form can be easily fed to algorithms for automated sequence comparison as they are provided by LingPy, and as they are needed for the final step of our retro-standardization workflow.

Since IGT resources taken alone never indicate whether two word forms that diverge slightly represent the same lexeme or not, the lexical and grammatical concordances which we created cannot replace a dictionary. What is needed, as a final step, is to make sure that all word forms which stem from the same lexeme, but which differ due to inflection or allomorphic variation, are assigned to the same lexeme entry.

ID	DOCULECT	CONCEPT	CONCEPT TYPE	FORM	TOKENS	OCCURRENCES	WORD FORMS	CROSSID
537	Qiang	market	lexicon	tşhaq	<mark>ts⁺ ªα</mark> q	2	tşhaq ta	606
538	Qiang	market	lexicon	tşhə	ts ^h ə	1	tşhə zeků ta	606
539	Qiang	market	lexicon	tşhaq	<mark>tş^h a</mark> q	2	tşhaq ta	606

Fig. 4. Three slightly diverging word forms denoting "market" in the IGT resource.

In order identify the lexemes in our data which are reflected by different word forms, we make use of methods for automated sequence comparison in order to produce an initial clustering of similar lexemes into language-internal cognate sets [11]. The result of this analysis is a Qiang wordlist that can be conveniently inspected in the aforementioned EDICTOR tool.

The benefits of this conversion become immediately evident when inspecting the data in detail. As can be seen from the example in Figure 4, we can find three different word forms in the column FORM which all denote the concept "market" in the corpus, which occur together as many as five times. While the two word forms, the first and the third, only differ by their vowel, the second form differs also in the lack of a final consonant. When comparing the differences with our standardized version of the transcription in the field TOKENS, one can see that the difference between [a] and [α] has been accounted for through our orthography profile, in which we already made the decision that [a] is meant to reflect [α]. The segmented form as rendered by the EDICTOR tool still lists this form with a super-script *a*, since we deliberately marked all cases of *a* being meant to represent [α] in our orthography profile.¹ For the form [ts ϑ], it is difficult to judge if this is a distinct word or a transcription problem. In any case, what we can clearly see from this example, is, that the procedure of retro-standardizing IGT resources can directly help to improve the resources by pointing to transcription problems.

4.2.7 Exporting the data. As a final step of our workflow, the Python library allows to export the retro-standardized
 resource to a web-based application that can be used to browse through the IGT examples, searching for lexical and
 grammatical glosses as well as specific word forms. Given that resources in book form are hard to inspect efficiently,
 this concordance browser offers a very convenient way for typologists and comparative linguists to dive deeper into
 a given resource. The concordance browser is available from the supplementary material accompanying this study.
 Figure 5 illustrates its basic usage.

 $[\]frac{1}{1}$ This is done by writing the original sound segment and the interpreted sound segment separated by a slash in the replacement column of an orthography profile, thus, underlyingly, the form reads [ts a/a q] and is rendered as superscript by the EDICTOR.

CONC	ORDAN	CE BRO	WSER		
hand	word form				
Found 1	1 matches				
	XI Text 6, SENI				
tsoqpi,	trile-upe	lo	tse-ze	japeq-tu	
				a p a q	
this:family	1pl- grandfather	also	this-CL	hand-LOC	
				hand	
TEM 2 (TE	XT Text 6, SENT	ENCE 19, PHR	ASE 46)		
da-329-u'	da-ûnsirj	fantşənşə	zmetşi	tsoqpi	japa-q-ta- quani
					j e p
DIR- set.out-2sg	1sg-TOP	†(anywaycis)	emperor	this family	hand-top- TOP

Fig. 5. Searching for occurrences of "hand" in the IGT resources of Qiang with help of the automatically generated *Concordance Browser.*

4.3 Examples

604 In order to illustrate how the concordance browser constructed from the retro-standardized dataset can be used to 605 shed light on actual linguistic questions, consider the annotation of the hearsay marker [(j)i]. When searching for the 606 grammatical concept "HS", referring to the hearsay marker in Ronghong Qiang, a search with help of the concordance 607 browser yields 24 results, of which the majority of examples has the form [i] (7 occurrences) or [ji] (6 occurrences), 608 as in [oqpi fio-pə-i], glossed as family DIR-become-HS, which can be translated as "became a family". However, in 609 610 several of these examples, the form corresponding to the hearsay marker appears as [wei], thus containing a bilabial 611 glide initial which is not present in any of the other examples. While it is difficult to confirm this for all 8 examples it 612 seems there that this form reflects an under-analyzed [-w] morpheme which LaPolla and Huang identify as being part 613 of the 'non-actor person marking suffixes' elsewhere in their grammar (see e.g., page 120, 143). We therefore think that 614 615 it is possible that this morpheme is incorrectly being marked as the HS marker, at least in some of the examples, as, for 616 example, in [fio-mu-xtcu-wei], glossed as DIR-NEG-burn-HS, which can be translated as '(they) weren't burned' (Text 617 1, Phrase 5), or in [de-l-wei], glossed as DIR-give-HS, '(god) gave it to them' (Text 2, Phrase 5). 618

The analysis of the hearsay marker in the Ronghong variety of Qiang is but one small example of how our retrostandardization can help to shed light on a given IGT resource. If more resources were retro-standardized in the way illustrated here, we think, the great service that inter-linear-glossed text provides for typologists and comparative linguistics, can further be increased.

624 Manuscript submitted to ACM

5

596

597

625 5 OUTLOOK

626

637

638 639

640

641

642

643 644 645

646

647 648

649

650 651

652 653

654

655

656

657

658

659

660

661

662

663

664

668

669

676

In this study we have proposed an initial framework for the consistent handling and the retro-standardization of 627 IGT resources in language documentation studies. By illustrating how a concrete resource of a highly endangered 628 629 Sino-Tibetan language can be successfully retro-standardized and presented in a way that facilitates not only the 630 linguistic but also the computational investigation of the language data, we have tried to show that retro-standardization 631 as well as a sustainable data handling is not per se impossible, as scholars often fear, but can even be carried out 632 much more quickly and efficiently than usually assumed. The workflow we propose integrates neatly into previous 633 634 standardization efforts in the field of computational historical linguistics and computational linguistic typology and 635 requires only a minimal amount of familiarity with the command line in order to be applied successfully. 636

In the future, we hope to expand our workflow further. First, we want to integrate it more closely with different formats currently used in larger IGT collections, such as PanGloss, ODIN, or the Dictionaria project. Second, we want to discuss with colleagues to which degree it might be possible to establish a reference catalog for grammatical elicitation glosses. Third, we want to integrate our workflow more closely with the CLDF initiative and ideally make a full-fledged proposal to integrate IGT resources along with concept lists, word lists, and list of grammatical elicitation glosses into the standard formats of linguistic data resources currently offered by CLDF.

SUPPLEMENTARY MATERIAL

The supplementary material contains the source code, the data, and additional instructions on how to use them in order to replicate the analyses discussed here. It can be downloaded from the Open Science Foundation at https://osf.io/n4vrk/?view_only=719c26b98c89443fbb6543234e702f19.

REFERENCES

- Cormac Anderson, Tiago Tresoldi, Thiago Costa Chacon, Anne-Maria Fehn, Mary Walworth, Robert Forkel, and Johann-Mattis List. 2018. A Cross-Linguistic Database of Phonetic Transcription Systems. Yearbook of the Poznań Linguistic Meeting 4, 1 (2018), 21–53.
- [2] Timotheus A. Bodt and Johann-Mattis List. 2019. Testing the predictive strength of the comparative method: An ongoing experiment on unattested words in Western Kho-Bwa languages. Papers in Historical Phonology 4, 1 (2019), 22–44.
- [3] Bernard Comrie, Martin Haspelmath, and Balthasar Bickel. 2015. Leipzig Glossing Rules. Conventions for interlinear morpheme-by-morpheme glosses. Max Planck Institute for Evolutionary Anthropology, Leizpig. https://www.eva.mpg.de/lingua/pdf/Glossing-Rules.pdf
- [4] Jonathan Evans and Jackson T. S. Sun. 2017. Contraction. In Encyclopedia of Chinese language and linguistics, Rint Sybesma (Ed.). Vol. 1. Brill, Leiden and Boston, 517–526.
- [5] Robert Forkel, Johann-Mattis List, Simon J. Greenhill, Christoph Rzymski, Sebastian Bank, Michael Cysouw, Harald Hammarström, Martin Haspelmath, Gereon A. Kaiping, and Russell D. Gray. 2018. Cross-Linguistic Data Formats, advancing data sharing and re-use in comparative linguistics. Scientific Data 5, 180205 (2018), 1–10.
- [6] Harald Hammarström, Robert Forkel, and Martin Haspelmath. 2019. Glottolog 4.0. Max Planck Institute for the Science of Human History, Jena. https://glottolog.org
- [7] Martin Haspelmath and Robert Forkel. 2017. Toward a standard list of grammatical comparative concepts: The Grammaticon. Talk held at the
 database workshop of the ALT Meeting 2017. http://dynamicsoflanguage.edu.au/storage/alt-2017-database-workshop-book-of-abstracts-forkel haspelmath-haynie-skirgard.pdf
 - [8] Randy J. LaPolla. 1996. A grammar of Qiang with annotated texts and glossary. City University of Hong Kong, Hong Kong.
 - [9] Randy J. LaPolla and Chenglong Huang. 2003. A grammar of Qiang with annotated texts and glossary. De Gruyter Mouton, Berlin and New York.
- - [11] Johann-Mattis List. 2014. Sequence comparison in historical linguistics. Düsseldorf University Press, Düsseldorf.
- [12] Johann-Mattis List. 2017. A web-based interactive tool for creating, inspecting, editing, and publishing etymological datasets. In *Proceedings of the* 15th Conference of the European Chapter of the Association for Computational Linguistics. System Demonstrations. Association for Computational
 Linguistics, Valencia, 9–12.
- 675 [13] Johann-Mattis List. 2017. Historical Language Comparison with LingPy and EDICTOR.

Manuscript submitted to ACM

677 678	[14]	Johann-Mattis List. 2018. Towards a history of concept list compilation in historical linguistics. <i>History and Philosophy of the Language Sciences</i> 5, 10 (2018), 1–14. http://hiphilangsci.net/2018/10/31/concept-list-compilation/
679	[15]	Johann-Mattis List. 2019. PoePy. A Python library for handling annotated rhymes. Max Planck Institute for the Science of Human History, Jena.
680	[]	https://github.com/lingpy/poepy/releases
681	[16]	Johann-Mattis List, Cormac Anderson, Tiago Tresoldi, Christoph Rzymski, Simon Greenhill, and Robert Forkel. 2019. Cross-Linguistic Transcription
		Systems. Max Planck Institute for the Science of Human History, Jena.
682	[17]	Johann Mattis List, Simon Greenhill, Christoph Rzymski, Nathanael Schweikhard, and Robert Forkel. 2019. Concepticon. A resource for the linking of
683		concept lists (Version 2.1.0). Max Planck Institute for the Science of Human History, Jena. https://doi.org/10.5281/zenodo.3351275
684	[18]	Johann-Mattis List, Simon Greenhill, Tiago Tresoldi, and Robert Forkel. 2019. LingPy. A Python library for quantitative tasks in historical linguistics.
685		Max Planck Institute for the Science of Human History, Jena. http://lingpy.org
686	[19]	Johann-Mattis List, Simon J. Greenhill, Cormac Anderson, Thomas Mayer, Tiago Tresoldi, and Robert Forkel. forthcoming. CLICS ² . An improved
687		database of cross-linguistic colexifications: Assembling lexical data with help of cross-linguistic data formats. Linguistic Typology 22, 2 (forthcoming).
688		http://clics.clld.org
689	[20]	Johann-Mattis List, Simon J. Greenhill, and Russell D. Gray. 2017. The potential of automatic word comparison for historical linguistics. PLOS ONE
690		12, 1 (2017), 1–18.
691 692	[21]	Johann-Mattis List, Nathan W. Hill, and Christopher J. Foster. 2019. Towards a standardized annotation of rhyme judgments in Chinese historical phonology (and beyond). <i>Journal of Language Relationship</i> 17, 1 (2019), 26–43.
693	[22]	Johann-Mattis List, Philippe Lopez, and Eric Bapteste. 2016. Using sequence similarity networks to identify partial cognates in multilingual wordlists.
		In Proceedings of the Association of Computational Linguistics 2016 (Volume 2: Short Papers). Association of Computational Linguistics, Berlin,
694		599-605.
695	[23]	Anatole Lyovin. 1969. Review of Hànyū făngyīn zihuì by Běijīng Dàxué. Language 45, 3 (1969), 687–697. http://www.jstor.org/stable/411456
696	[24]	Steven Moran and Michael Cysouw. 2018. The Unicode Cookbook for Linguists: Managing writing systems using orthography profiles. Language
697		Science Press, Berlin. http://langsci-press.org/catalog/book/176
698	[25]	Yugo Murawaki. 2019. Bayesian learning of latent representations of language structures. Journal of Computational Linguistics 45, 2 (2019), 199–228.
699	Fa (1	https://doi.org/10.1162/COLIa00346
700	[26]	Laurent Sagart, Guillaume Jacques, Yunfan Lai, Robin Ryder, Valentin Thouzeau, Simon J. Greenhill, and Johann-Mattis List. 2019. Dated language
701		phylogenies shed light on the ancestry of Sino-Tibetan. <i>Proceedings of the National Academy of Science of the United States of America</i> 116 (2019), 10317–10322. Issue 21.
702	[27]	Morris Swadesh. 1955. Towards greater accuracy in lexicostatistic dating. International Journal of American Linguistics 21, 2 (1955), 121–137.
703	[27]	arXiv:1263939
704	[28]	Mark D. Wilkinson, Michel Dumontier, Ilsbrand J. Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Boiten,
705		Luiz B. da Silva Santos, Philip E. Bourne, et al. 2016. The FAIR Guiding Principles for scientific data management and stewardship. Scientific Data 3
706		(2016), 1–8.
707		
708		
709		
710		
711		
712		
713		
714		
715		
716		
717		
718		
719		
720		
721		
722		
723		
724		
725		
726		
727		

728 Manuscript submitted to ACM