



1 – Introduction

We present the process of constructing a corpus of spoken and written material for Santome, a Portuguese-related creole language spoken on the island of São Tomé in the Gulf of Guinea (Africa). Santome is the result of language contact between Portuguese, the lexifier language, and several Benue-Congo languages, in particular Edo (Edoid) and Kikongo (Bantu) [2,3]. As the language lacks an official status, we faced the typical difficulties, such as language variation, lack of standard spelling, lack of basic language instruments, and only a limited data set.

For the corpus compilation we followed corpus linguistics standards and used UTF-8 character encoding and XML to encode meta information. We discuss how we normalized all material to one spelling, how we dealt with cases of language variation, and what type of meta data is used. We also present a POS-tag set developed for Santome.

2 – Corpus

The corpus contains written and oral sources and has about 184K words. The number of available sources is limited. A large part of the corpus can be placed in the domain of folklore (folktales, riddles, song texts). The written part contains a few newspaper articles and books, pamphlets, cultural magazines and a blog. The spoken corpus comprises transcriptions of recordings made in 1997 and 2001 with twenty native speakers. These spoken recordings have been freely transcribed in the sense that we have tried to match written text as much as possible. Since we do not yet have copyrights for all the materials used in the corpus, at a first stage we plan to make the corpus available for concordances in an online interface, CQPweb [4], which allows users to search for concordances.

3 – Language standardization

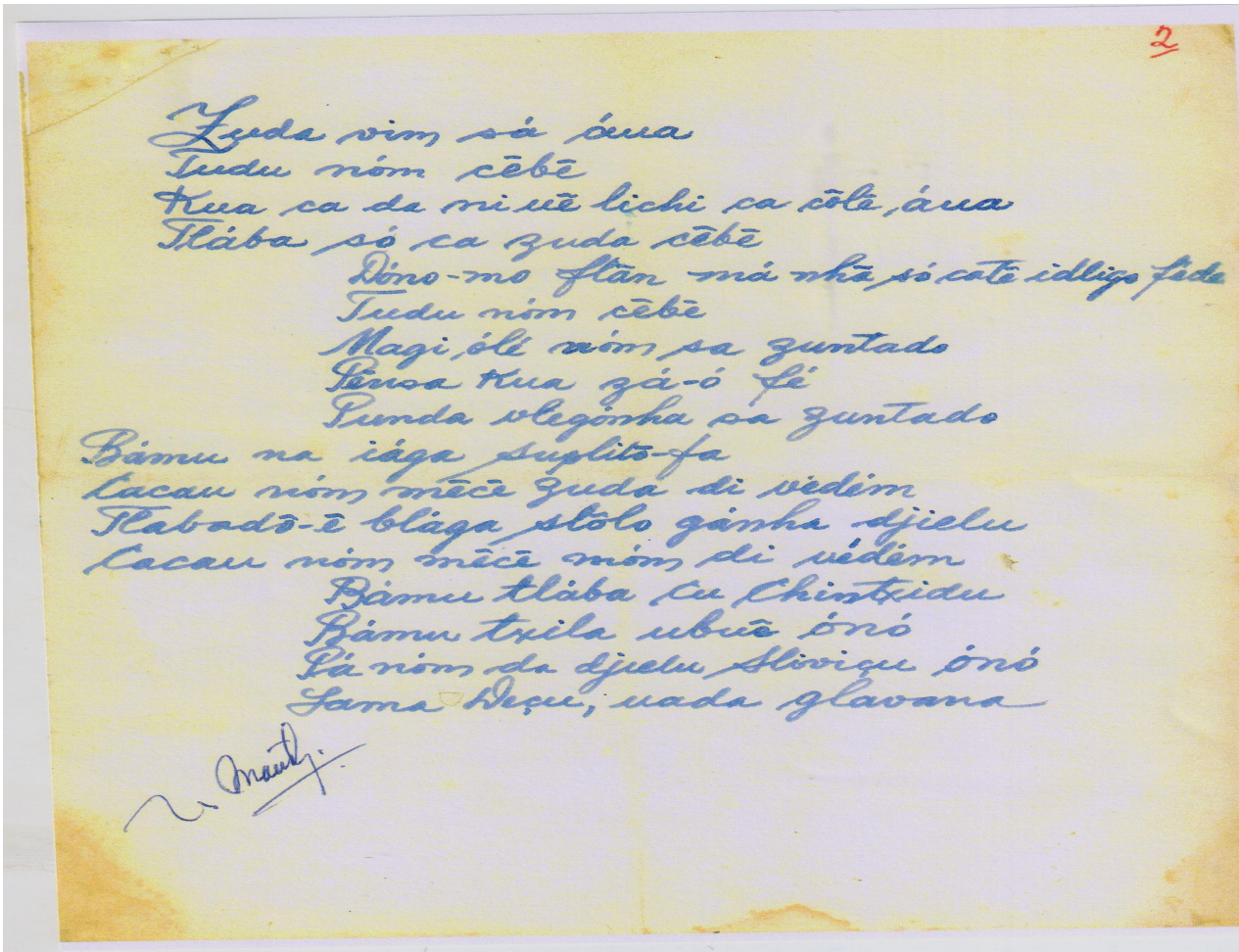
Despite the limited writing tradition in Santome, the orthographies have been highly variable. A word like [kwa] ‘thing’, for instance, has been written in the following ways: *cua*, *cuá*, *qua*, *quá*, *kua*, *kuá*, *kwa*, *kwá*. Therefore, we converted all orthography to ALUSTP, a phonology-oriented

(Original)	(Adapted)
<i>Inen piscadô nón</i>	<i>Inen pixkadô non</i>
<i>di tudu bôca plé</i>	<i>di tudu boka ple</i>
<i>di tēla cé non glavi ximentxi</i>	<i>di tela se non glavi ximentxi</i>
<i>cá chē ni ké d'inēn</i>	<i>ka xē ni ke dinen</i>
<i>cu amuelê cu buá vonté</i>	<i>ku amwêlê ku bwa vontê</i>
<i>chê bá nôtxi</i>	<i>xê ba nôtxi</i>
<i>chê bá Tláchia</i>	<i>xê ba Tlaxa</i>
<i>basta p'men bála blé d'omali</i>	<i>baxta pa inen ba ala blê d'omali</i>
<i>bá bucá vadô panhã cé</i>	<i>ba buka vadô panha se</i>

writing proposal developed in 2009 and ratified in 2010 by the Ministry of Education of Culture of S. Tomé and Príncipe [5]. The main principle of this proposal is a one-to-one phoneme-grapheme correspondence (see for instance the use of “k” [k] in the adapted version and its counterparts in the original orthography).

References

- [1] RGPH – 2001. (2003). Características educacionais da população – Instituto Nacional de Estatística. S. Tomé e Príncipe.
- [2] Ferraz, L. (1979). The creole of São Tomé. Johannesburg: Witwatersrand University Press.
- [3] Hagemeijer, T. (2011). The Gulf of Guinea creoles: genetic and typological relations». Journal of Pidgin and Creole Languages, 26:1, pp. 111-154.
- [4] Hardie, A (forthcoming) "CQPweb - combining power, flexibility and usability in a corpus analysis tool".
- [5] Pontífice, J. et al. (2009). Alfabeto Unificada para as Línguas Nativas de S. Tomé e Príncipe (ALUSTP). São Tomé.



4 – Meta data

We encode meta data about the corpus texts like author and date in a simple XML format that is compatible with the P5 guidelines of the Text Encoding Initiative (TEI). We use the following fields:

- author: The author of the text (if known)
- age: The age of the recorded speaker (spoken data)
- place of recording: geographical location of the recording
- date: The date of publication (if any), which can be exact or approximate. Unless we found evidence to the contrary, we assumed that publication dates are close to the date of writing.
- source: Book, newspaper article, (cultural) magazines, pamphlets, online, unknown.
- genre: See table
- Notes: Additional information, such as publisher, place of publication.

5 – POS Annotation

Some examples of POS annotation.

(1) Sentence from the corpus:

Ola nansê ka xka nda ku migu
When 2PL TAM TAM walk with friend
CJ PRS TAM TAM V PREP CN
sela nansê toma kwidadu ê
must 2PL take care PRT
MOD PRS V CN PRT
'When you are hanging out with friends, you must be careful. '

(2) Reduplication (fully or partially):

tlêxi-tlêxi 'in groups of three' RED:NUM
tlê-tlêxi 'all three' ' RED:NUM

(3) Ideophones

kabêsa wôlôwôlô 'foolish person' (lit. head+id.) CN ID
sola potopoto 'cry intensely' (lit.cry+id.) V ID
vlémê bababa 'intense red' (lit. red+id.) ADJ ID

6 - Final remarks

Corpora for the three other Gulf of Guinea creoles are under way and, together with the Santome corpus, will be used for comparative research leading to the reconstruction of properties of the proto-Gulf of Guinea creole. We further expect this corpus to be an important tool in language maintenance and revitalization namely

through the development of other language resources. Finally, a more wide-spread corpus-based approach to creole languages using large

amounts of data will endow comparative research on creoles with tools that allow to assess claims regarding typological similarity between these languages.

Acknowledgements

The Santome corpus is funded by the Portuguese Foundation of Science and Technology (FCT) as part of the project *The origins and development of creole societies in the Gulf of Guinea: An interdisciplinary study* (PTDC/CLE-LIN/111494/2009) and the FCT program Ciência 2007/2008 (Iris Hendrickx).

More information:

<http://www.gulfofguineaceoles.com/>

Tag	Category	Examples (Santome)
ADI	Adjective	glavi 'pretty', vlémê, 'red' bon 'good'
ADV	Adverb	oze 'today', za 'already', yôxi 'yes'
ART	Article	ũa, inen.
CJ	Conjunction	i 'and', ô 'or', maji 'but', ola 'when'
CN	Common Noun	mosu 'boy', ope 'foot', ke 'house'
COMP	Complementizer	kuma 'that', ku 'who, which, that', xi 'if'
CONN	Sentence connector	so 'and then', soku 'and then'
DEM	Demonstrative	se 'this', sala 'that', ise 'this one'
DGT	Digit	0, 1, 42, 12345, 67890, ...
EXC	Exclamative	Kê ... 'How ...'
FOC	Focus marker	so, soku.
FW	Foreign word	
ID	Ideophone	bababa, nwininwini, sũũũ, potopoto, ...
INDF	Indefinite	tudaxi 'everybody', nadaxi 'nothing'
INT	Interrogative	kuma 'how', andji, 'where'
ITJ	Interjection	kaki 'exclamation of surprise'
MOD	Modality Marker	sela 'must', milhon 'better...'
NEG	Negation marker	na, fa, fô, nantan, naxi.
ON	Onomatopoeia	plaplaplaplapa 'walking sound'
NUM	Numeral	ũa 'one', dôsu 'two', plumêlu 'first'
PM	Presentational marker	avia 'there was/were', ya 'here is/are'
PNM	Proper Name	Zon, Maya, Kolema, Txindadiji
PNT	Punctuation Mark	, . ? , , , , , , !
POSS	Possessive	mu(n) 'my', dinen 'their'
PP	Participle	fladu 'said', bixidu 'dressed'
PREP	Preposition	antê 'until', di 'of', djina 'since'
PRS	Personal	n 'I', ami 'I', bô 'you', non 'we'
PRT	Particle	an, en, ê, fa, fan, ô.
QNT	Quantifier	kada 'each', tudu 'all, every', yô 'many'
RED:xx	Reduplicated category	dôsu-dôsu 'in groups of two'
REFL	Reflexive	mu 'myself', dê 'himself/herself'
RV	Residual value	abbreviations, acronyms, etc. sun 'Mr.', san 'Misses', dôtêlô 'doctor'
STT	Social Title	
TAM	Tense-Mood-Aspect marker	ka, kâ, xka, ska, tava-ta.
V	Verb	mêsé 'want', pô 'may, can', sa 'be'

Contact

Av. Professor Gama Pinto, 2, 1649-003 Lisboa, Portugal
Tel.: +351 21 790 79 00 | URL: www.clul.ul.pt



CLUL

Centro de Linguística
da Universidade de Lisboa