# Issues in digital text representation, online dissemination, sharing and reuse for African tone languages

**Emmanuel Ngué Um**
University of Yaoundé I – Cameroon
ngueum@gmail.com

## Abstract

In tone languages of Africa, tones may encode meaning either as separate linguistic units or in association with segmental morphemes. In mainstream text representation models, however, the linguistic autonomy of tones is often overridden by the graphical layout of characters. In these models, accents which mark tones cannot be easily parsed for their linguistic information apart from the segments which bear them. This paper suggests a model or representation based on TEI-XML where both tones and segments can be represented as a unique string of characters, therefore making text information easily parsable.

## 1 Introduction

Language Documentation (LD) and description have generated textual resources for minority tone languages of Africa. With the event of computers, text resources are being created mostly in the form of digital-born texts. But visual layout of texts in these languages follows a variety of representation models, even within the same writing system. Some of the issues at stake are IPA vs Latin characters; omission vs surface representation of tones as accents; surface vs deep representation of tones, one tier vs multiple tier representation of linguistic analysis.

Even when widely shared standards exist, implementation of these standards from one project to another, or from one software to another is not consistent. This raises the issue of which technology is best suited for optimal graphical rendering of the linguistic information encoded through texts in tone languages of Africa. The question to know how to better reach a compromise between end-user-friendly orthographies of previously unwritten languages on the one hand, text-to-human transparency of linguistic information conveyed through texts, and computability of derived text corpora on the other hand is left unanswered.

The objectives for which text corpora are produced also dictate models of text representation. Lehman (2001) suggests that textual resources stemming from LD should be designed in such a way that they "represent the language for those who do not have access to the language itself". In the case of African tone languages, this entails representation of tone sequences as anchored in the melody pattern of speech production. The reason for representing tones in writing is that they form independent prosodic units of meaning which deserve appropriate analytical attention.

Taken from the point of view of a theory of text representation, this approach is undoubtedly more analytic than early missionary orthography models which are devoid of tone markers. However from the point of view of the text layout model and digital processability, tone representation as it has been implemented so far, tend to create fuzzy boundaries and associations between the text structure and the linguistic information. This is so because current models of tone representation in text production have built up from the Latin scripting framework where the linguistic information is encoded in a one-grid string of characters.

Sometimes tone markers represented as accents and associated with letter characters form pre-composed or ad hoc binary or ternary graphical unites. This actually adds an additional and somehow artificial layer of information.

The purpose of this paper is to bring out the limitations of tone representation and analysis as an upper layer of character strings, and to suggest an alternative model based on the TEI-XML markup language. Section 2 will review some theoretical and technical issues involving tones

in African languages, and the implication of a tone model of representation for the linguistic analysis. Section 3 will deal with the limitations of existing models of tone representation with regard to granularity in linguistic analysis, fuzziness in text and linguistic information mapping, and parsability of textual information. In section 4 I will propose a prototypical model for enhanced representation of tone in African tone languages texts. This model is meant to address the limitations pointed out in existing models.

## 2 Theoretical and technical issues in graphical representation of tones in African languages

### 2.1 Theoretical issues

When tones are represented in texts, they may either surface as accents on syllable nuclei (1), or as alphabetical labels standing on a separate tier and linking to syllable nuclei by means of association lines (2). In (2), labels 'H' and 'L' respectively stand for 'high' and 'low' tones.

(1) màlép má ńsóbì[1]
    "water is poured"

(2) màlép má ńsóbì
    | |   | | | |
    L H   H H H L

In (1), tones may either be interpreted as supra-segments, or as auto-segments. A supra-segmental approach to tone is one which poses tones and their Tone Bearing Units (TBU) as inherently bound linguistic units. In this approach, tones are not separated from or interpreted without their corresponding TBUs. An illustration of a supra-segmental analysis of tones is provided in (3).

(3) mà-lép má ń-sóbì
    CL6-water SM.CL6 PRES-PERF.pour
    "water is poured"

The nasal prefix [n] along with its associated high tone may be interpreted as a unique aspectual marker for the perfective. The same interpretive stance applies to other remaining building blocks, namely: class 6 noun préfix [ma] which

is associated with a low tone[2]; class 6 subject marker [ma] which is associated with a high tone; both roots [lep] and [sobi], where [lep] is assumed to inherently bear a high tone, whereas [sobi] is assumed to bear a sequence of two level tones H-L.

The supra-segmental approach to tone analysis, though still prevailing in some scholarly circles in Africa, is unable to account, for example, for why a word may bear a low tone in some situations, and a high tone in other situations. This is indeed the case with the root [sobi]. When standing alone, the verb *sobi* "to be poured" surfaces with a different tone melody, namely L-L[3], as in (4).

(4) /sòb-ì/ "to be poured"

Only if we consider tone to be an independent linguistic unit from its associated segment, can we consistently account for pitch variation in tone melody across word forms in a text. The strength of auto-segmental phonology (Goldsmith, 1990) lies in the fact that it views tone and it associated syllable as two distinct linguistics units, as in (2). Because at the abstract level of analysis tones exist independently from their TBUs, changes affecting one of the two associated units may not necessarily affect the other, as can be seen in the alternation between [sòbì] and [sóbì].

In spite of its theoretical robustness and validity, auto-segmental phonology has not significantly impacted traditional models of text representation and analysis in tone languages of Africa. Given that linear analysis of textual information relies on the text layout, it follows that much of the linguistic information encoded through texts in African tone languages is sometimes overridden by the constraints and limitations of a one-grid and linear representation of the linguistic information. This does not favor automatic processing of textual information; a situation which is further hampered by technical issues such as character encoding and typesetting technologies.

---

[1] Most examples provided in this paper are taken from Basaa, a Bantu language spoken in the Center and Coastal regions or Cameroon.

[2] In spite of the fact that noun class prefixes are generally represented as toneless morphemes in Bantu languages.

[3] Another possible reading of the tone melody associated with /sòb-ì/ is /L-ø/ where the final syllable nucleus is assumed to be toneless.

## 2.2 Technical issues in digital text representation of African tone languages

Within both the scholarly community and the speech groups of minority tone languages of Africa, there is currently much emphasis on developing typesetting technologies such as virtual keyboards and extension of Unicode character sets. The focus here is on facilitating end-user creation and editing of texts, and less so on ensuring digital exploration, automatic processing, software-independence, and universal sharing of these texts.

In the first place, surface representation of tones as accents is a thorny issue. This may be a convenient graphical method with regard to human cognitive ability to process visual information, but much less convenient from a computational perspective. This is so for many reasons.

(1) Graphical clustering of tone markers with letters (a, e, n, o, i, etc.) is perceived as a binary unit of writing by humans, but as two separate digital objects by the computer, given that Latin characters and accents bear distinct digital codes[4].

(2) When accented letters such as [é], [è] or [ê] are available as pre-composed characters and therefore assigned unique unicode code points, stability and consistency of character sets is achieved at the expense of parsability of textual information. For example, a search algorithm could be meant to determine the frequency of high tones in a text, as compared to low tones. If high tones are only counted as acute accents, the algorithm would only parse those characters with an acute accent. This would be misleading because in an auto-segmental approach to tone analysis, a contour high-low tone shape must be interpreted as a sequence of independent 'high' + independent 'low'. A workaround solution could be to count contour tones as representing tokens of both high and low tones. However this workaround solution would not hold if, in addition to looking for high tone frequency, the search algorithm also encompassed retrieving specific meaning associated with high tones.

(3) Deciphering of pitch associated with a given accent, whether acute, grave, circumflex or caron by humans does not require neither intrinsic knowledge of the linguistic information encoded by the tone marker or the digital code

point associated with it. In (3) for example, acute accent on the syllable [so] in [ńsóbì] signals both high tone and perfective aspect; however, this grammatical knowledge is not a prerequisite for proper pronunciation and semantic interpretation of the root morpheme [sob] by a human. For the computer, however the grammatical information encoded through the high tone on the syllable [so] in [ńsóbì] has to be labelled differently from the lexical information associated with the same syllable in the lexical form of the word in (4). If applied to text-to-machine modeling, it might be necessary to assign persistent digital codes for each tone level. This calls for explicit and unambiguous markup of the linguistic information associated with each single tone in a text, as opposed to merely representing tones as iconic pitch signals on syllables.

(4) Graphical representation of tones as accents does not facilitate linear glossing and markup of tone morphemes in a one-to-one relationship with their corresponding meanings, as is the case with 'ordinary' morphemes made up of 'conventional' Latin character strings.

Efforts to make it possible to parse tones automatically in text resources of African tone languages have been attempted with *TOOLBOX*, one of SIL's fieldwork software for text creation, editing, and analysis. This functionality is achieved thanks to a specific data input method designed by Buseman, as cited in McGill (2009). Buseman's method has been chiefly motivated by the need to overcome the software's understanding of tones as intrinsic parts of the string of characters. Overall, Buseman's method consists in separating tone marking and the corresponding segmental string of characters triggering TBUs, and then glossing them each on their own, as in (5)[5], which I have borrowed from McGill (2009: 244).

```
(5) \tx      dùkwá
    \mb      dukwa  -L       -H
    \ge      go     IMP
    \ft      'go'
```

Enhancement and enrichment of Buseman's method has been suggested by McGill (2009) in view of further development of *TOOLBOX* as well as development of new software. Among other suggestions, McGill advocates the implementation of SIL's *TonePars* program (Black 1997) into future software development for corpora creation for endangered languages. "The TonePars program […] allows for modeling an auto-segmental approach to tone"[6].

It is desirable that further refinement of *Toolbox* and/or its offspring *Fieldwork's Language Explorer (FLEx)* should include auto-segmental modeling of tone. This could be achieved through multiple-tier mapping where both character strings and tone markers could be separately analyzed for the linguistic information they encode. The resulting text file could then be transformed into parsable *XML* files. However, even if such software development could be implemented in *TOOLBOX* or similar linguistic analysis tools, many issues would still be left unresolved.

(1) Text representation and processing of tone languages with *TOOLBOX* is usually biased towards an exclusively scholarly stance, over other possible language usage frameworks, where surface representation of tones might be optional or even unecessary. As a matter of evidence, it should be noted that non-marking of tones in orthography in tone language communities is the norm rather than the exception. An edifying illustration is the case of *Google Translate*, where none of the five African tone languages which are available for online translation[7] (namely Hausa, Igbo, Shona, Xhosa and Yoruba) uses tone marking in their writing system. Plain Latin scripts devoid of tone markers are preferred over the IPA-based writing systems mostly advocated by linguists, which marks tones.

(2) Assuming *TonePars* algorithms are incorporated into *TOOLBOX*, the issue of character encoding is left unresolved. *TonePars* interprets the content value associated with a given tone, whether high or low, on the basis of its graphical shape, namely acute or grave. However, there may exist multiple character input possibilities for representing the same toned-segment graphically. For example <é> (the character <e> bear-

ing and acute accent), may have as possible typesetting inputs:

(a)  The unique pre-composed character <é> (*Unicode C1 Controls and Latin-1 Supplement*, code point 00E9);

(b)  <e> (*Unicode C0 Controls and Basic Latin*, code point 0065) and acute accent < ´ > (*Unicode C1 Controls and Latin-1 Supplement*, code point 00B4);

(c)  <e> (*Unicode C0 Controls and Basic Latin*, code point 0065) and acute accent < ´ > (*Unicode Spacing Modifier Letters*, code point 02CA);

(d)  <e> (*Unicode C0 Controls and Basic Latin*, code point 0065) and acute accent < ´ > (*Unicode Combining Diacritical Marks*, code point 0301);

(e)  etc.

In other words, the existence of multiple graphical representations of the same prosodic reality creates the possibility of arbitrary digital encoding of tones in a program such as TOOLBOX. This is definitely not good practice with regards to current standards in data sharing and dissemination as advocated by such data consortia as the Text Encoding Initiative (TEI), the Data Research Alliance (RDA), Component MetaData Infrastructure (CMDI), etc.

(3) As a consequence of the above two issues, text corpora created with tools such as TOOLBOX and FLEx cannot lend themselves appropriately to open data sharing and re-use, in a global text 'industry' where sustainability of data infrastructures rely heavily on interoperability. It should be noted, in addition, that mainstream linguistic analysis methods implemented in these software, namely standard morpheme-by-morpheme glossing of text information, does not provide a scheme for rich metadata input for linguistic information.

## 3   Limitations of existing models of tone analysis

Coming back to example (3) which is repeated in (6) for convenience, it appears that linguistic information is mapped in a one-to-one correspondence with the text building blocks.

(6)  mà-lép má ń-sóbì
     CL6-water SM.CL6 PRES-PERF.pour
     "water is poured"

---

This mode of text representation is common practice in linguistics scholarship, and there exists a set of standards for glossing linguistic information. The Leipzig Glossing Rules are one such standards. "They consist of ten rules for the 'syntax' and 'semantics' of interlinear glosses, and an appendix with a proposed 'lexicon' of abbreviated category labels" (Comrie, Haspelmath & Bickel 2015: 1). Rule 2 of the Leipzig Glossing Rules states that "there must be exactly the same number of hyphens in the example and in the gloss" (Comrie, Haspelmath & Bickel 2015: 2). Rule 4 further stipulates that "when a single object-language element is rendered by several metalanguage elements (words or abbreviations), these are separated by periods" (Comrie, Haspelmath & Bickel 2015: 3). This last rule justifies why 'SM' and 'CL6', then 'PERF' label and the lexical meaning 'pour' are clustered into two binary glosses in (3). This clustering indicates that the two labels making up a binary gloss are encoded by the same morpheme.

Single morphemes which encode more than one meaning are commonplace across the world's languages. What is problematic in clustering for example the 'PERF' label and the lexical meaning of the root 'pour' in (3) and (6) is that, these two meanings are indeed separately encoded by two linguistic units: 'PERF' is encoded by the high tone on the first syllable of the root, while 'pour' (the lexical meaning) is encoded by deep lexical structure of the verb root [sobi]. However, simply graphically mapping a piece of complex linguistic information such as 'PERF.pour' with textually 'unparsed' structural bundle such as [sòbì] without explicit and unambiguous assignment of this information is fuzzy.

Because morpheme-by-morpheme glossing tends to prioritize synchronization of analysis with on-line stretching of text for the sake of grammatical information tracking, only information needed for ad hoc understanding of the grammatical sequencation of the test is deemed relevant. It is questionable, for example, why tones in African tone languages are only provided with glosses (analytical labels such as 'PERF', 'PRES', etc.) when they appear to trigger variation in the tone shape of a given root as in the [sòbì] ~ [sóbì], or in grammatical morphemes such as the nasal prefix [ń]. To put it in simple terms, there is structural inconsistency in glossing some tones, and not others. After all, if the high tone signals grammatical meaning in situations such as [sóbì], there are good reasons

to believe that the low tone which occupies the same position in [sòbì] also signals some form of linguistic information. However, the primitive information associated with this low tone is over-ridden in standard linguistic glossing. To further demonstrate why glossing every single tone in a text is a pre-requisite to systematic and accurate parsing of the linguistic information associated with tones, lets take this other example (7).

(7) ndʒɔ̀ŋ     ì
 CL9.palm.oil.residue SM.CL9
 ǹ-sóbì
 PRES-PERF.pour
 "palm oil residue is poured"

In (7) the nasal prefix in [ǹsóbì] is now associated with a low tone, as opposed to a high tone in (3) and (6). This variation is linked with the noun class of the subject word [ndʒɔ̀ŋ], namely class 9. If we take a third (8) and fourth (9) example where the subject nouns both belong to class 10 while surfacing with different tone shapes, this becomes even more glaring.

(8) ɓàs    í
 CL10.salt  SM.CL9
 ń-sóbì
 PRES-PERF.pour
 "salt is poured"

(9) láŋ     í
 CL10.palmist oil  SM.CL9
 ń-sóbì
 PRES-PERF.pour
 "palmist oil is poured"

These examples clearly show that some critical linguistic information is overridden in the glossing of the nasal prefix which attaches to the verb root, whether it surfaces with a high tone as in (3), (6), (8) and (9), or with a low tone as in (7). Therefore a rigorous model for linguistic information analysis and retrieval would be one which systematically accounts for every bit of linguistic information encoded by every tone, in conjunction with the linguistic information encoded by segmental morphemes. The syllabic nasal prefix in the examples which precede should be glossed distinctively for both its segmental component which encodes tense, and its tonal component which encodes noun class. Likewise, tones surfacing on the subject noun roots, the subject concord markers, and the verb

roots should each be glossed for the linguistic information they encode.

If this assumption is valid - and there is robust evidence that it is -, then it becomes inescapable that existing models of tone analysis in African tone languages are flawed by inadequate analysis matrices.

The issue discussed here is just one of the many inconsistencies which may be observed in the analysis of tones in African languages using standard analysis tools and glossing models. Other issues are: proper deep representation of contour tones as in (10), (11), and (12)[8]; floating tones; tone spreading; tone shift; upstep; downstep, low tone rising, etc; all of which will not receive attention in this paper due to space constraint.

(10)  m-ùt            à
      CL1.person   SM.CL1
      bí-lɔ̂
      PST-come
      "somebody has come"

(11)  m-ùt            à
      CL1.person   SM.CL1        ǹ-lɔ̂
      PRES-PERF.come
      "somebody has come"

(12)  mὲ    ǹ-tɛ́hɛ́
      I       PREST-see
      m-ût
      CL1-person
      "I have seen somebody"

Examples (10) and (11) show how tone shape changes on the verb root [lɔ]. It should be reminded that this verb root is associated with a low tone lexically. Tone shape in (10) is therefore a 'normal' one, yet it deserves proper glossing. It is relevant for automatic analysis of tone-induced linguistic information, to identify the low tone on the root of the verb form [bílɔ̂] as a lexical tone, in order to make consistent parsing possible. As for the contour high-low shape in (11) namely [lɔ̂], it has quite a simple logical explanation. This is the result of the association of the high tone marking perfective aspect in (3), (6), (7), (8), and (9), with the underlying lexical low tone. If every tone were glossed consistently, then changes affecting tone shapes in and across words in a text could be modeled with more ac-

curacy, and anticipation of these changes would me made much easier for a computer program.

Another illustration of change in tone melody across words is seen in (10), (11) and (12), where the noun form [mùt] surfaces with a lexical low tone in (10) and (11), then with a contour falling tone in (12). The contour tone melody here is the result of a process known in Bantu languages as metatony (Nurse 2006, Hyman and Linnet 2011, Makasso 2012). This process triggers a high tone at post-verbal positions to signal prosodic conjunction between the verb and the object (Makasso 2012: 15). While metatony signals prosodic relationship across words in a verb phrase, this phenomenon also encodes syntactic information worth being accounted for in an optimal analysis scheme.

## 4    A prototype of the TEI-based Model for Enhanced Linguistic Annotation for African Tone Languages

TEI (the Text Encoding Initiative) provides comprehensive guidelines for the development of text encoding standards and schemes for virtually any text encoding project. Chapter 15 of the Guidelines (TEI consortium, 2016: 504-517) deals with annotation of language corpora. Section 17 (TEI consortium, 2016: 570-587) deals with 'Linguistic Segments Categories' and 'Linguistic Annotation', among other issues. Section 18 (TEI consortium, 2016: 588-620) is about annotation of 'Feature Structures' and 'Atomic Feature Values', among other issues. Specific structural phenomena such as tones, which are pervasive in Bantu languages are not yet addressed in the Guidelines.

TEI is a community-driven initiative and maintained by dynamic contributing members who share experiences and shape its further development. The Model for Enhanced Linguistic Annotation for African Tone Languages which I propose here is destined to be submitted to the Technical Council of the TEI consortium for review, enrichment and standardization. As in other TEI-related annotation schemes, the current model is not intended to address every specific aspect having to do with tone annotation in African languages. On the contrary, the model is a starting point, and therefore aims to stimulate discussions and contributions that could help build up a standard, digitally processable, interoperable and sustainable framework for the development of text corpora and text resources in African tone languages. Each individual lan-

---

[8] Glossing in these examples follow standard glossing models, namely the Leipzig Glossing Rules

guage encoding project for African tone languages could draw from this general scheme to tailor the model to their specific needs. The model is expected to be refined and adjusted as more stake-holders join in, either linguists, TEI experts, developers, projects managers, etc.

TEI being based on the XML encoding language, I will not deal here with modules such as Data Type Definition (DTD), XML schemas and name spaces. The model complies with the over-all TEI XML infrastructure, namely as concerns validation protocols and standards. For the sake of brevity, I will limit the presentation to tone encoding at the @word element. However, the model also encompasses issues such as optional non ASCII character encoding as well as phonological features, lexical, morphological, syntactic and discourse encoding.

I have adopted ad hoc vocabulary[9] for naming the following attributes:

a. @type describes a word or morpheme type; word types can be "nouns", "verbs", "adjectives", etc.; morpheme types can be "prefixes", "roots", "suffixes", "segmental", "tonal", etc.

b. @gloss describes the semantic value of a morpheme; this may apply to tense, aspects, plural, noun class markers, etc; gloss values are labels of linguistic information which conform to a standard glossing scheme such as the Leizig Glossing Rules.

c. @category describes a linguistic class to which a specific morpheme relates; a morpheme may fit into the "lexical", "grammatical", "syntactic", or "prosodic" classes.

d. @nounClass describes an integer for the noun class of a noun prefix, following standard Bantu grammatical reconstructions for noun classes (Meeussen 1967, etc.)

e. @segmental describes segmental morphemes, that is a word's building blocks devoid of their tone associates.

f. @tonal describes tone morphemes. The content of a tone element is represented by the unicode code point for the accent whose shape depicts the pitch level of the tone. Thus, the acute accent standing for high tone is represented in the model by the Unicode code point 02CA, while the grave accent is encoded by the Unicode code point 02CB. However, for the sake of economy and consistency, the content of tone elements is assumed to have earlier been described in the XML schema or DTD file as @high and @low entities with the value of each corresponding to its Unicode codepoint.

The model provides a comprehensive and linear analysis of every bit of linguistic unit which contributes to the meaning of the utterance. The analysis can be extended to more granular units such as consonant or vowel features for specific characters, depending on the needs of the encoder, without having to call multiple tier analysis into play (see Table 1).

---

[9] The vocabulary used in the model for attributes names, attributes values, and entity names is NOT yet standardized. It is expected that the TEI community along with linguistics will come together, perhaps within the framework of existing TEI Expert Groups, and work towards standardizing the model and its related vocabulary and other structural aspects.

**Table 1**       Implementation of the model for source utterance : [mùt à ǹlɔ̂]

```
<w type = "noun">
   <m type = "prefix" nounClass ="1">m<\m>
         <m type = "root">
                  <m type = "segmental">ut<\m>
                  <m type = "tone" category = "lexical">&low<\m>
   <\m>
<\w>
<w type = "nounParticle">
   <m type = "subjectMarker">
         <m type = "segmental" nounClass = "1">a<\m>
         <m type = "tone" category = "grammatical" nounClass = "1">&low<\m>
   <\m>
<\w>
<w type = "verb">
   <m type = "prefix">
         <m type = "segmental" gloss = "PRES">n<\m>
         <m type = "tone" nounClass = "1">&low<\m>
   <\m>
   <m> type= "root">
         <m type = "segmental">lɔ<\m>
         <m type = "tone" category = "grammatical" gloss = "PERF">&high<\m>
         <m type = "tone" category = "lexical">&low<\m>
   <\m>
<\w>
```

## 5   Conclusion

The above model of text representation solves the following problems otherwise not addressed in any existing annotation scheme:

(1) It brings out every unit of meaning distinctively, whether segmental or tonal. In the present model, linguistic information is analyzed unambiguously and consistently; whereas in mainstream linguistic glossing models the contour H-L tone in [lɔ̂] would have been represented as a complex gloss 'PERF.pour' without explicit specification as to which linguistic element encodes aspect and which other encodes lexical meaning;

(2) the XML encoding framework which the model builds on allows extensibility of the linguistic analysis; in this respect, a given encoding project may be reusable;

(3) it assigns persistent unicode code points[10] to each level tone, therefore making it easier for conversion, transliteration, compression and parsing of characters from one orthography scheme to another;

(4) the model equality reduces multiple representation of segments and tones into one linear string of text, making the text more conducive to interoperability across APIs[11];

(5) it forces granularity in linguistic analysis; in traditional grammatical analysis, the contour tone in the root 'lɔ' would not have been subject to binary glossing partly because of shortage of glossing space; in the present annotation model however, glossing of tones may apply to level and melody, toneless units, as well as in monosyllabic or polysyllabic words.

(6) Inasmuch as the model forces granularity upon units of meaning in a text, it triggers fine-grained description of tone phenomena, and therefore is likely to stir up further in-depth research in the prosody of African tone languages.

(7) Because the model makes the text more easily parsable, it may not only be implemented in text to to-text applications, but only in text-to-speech modeling.

---

10

http://www.unicode.org/versions/Unicode9.0.0/UnicodeStandard-9.0.pdf

[11] Application Programming Interface such as programming languages.

# References

Black, H. Andrew. 1997. TonePars: A computational tool for exploring autosegmental tonology. SIL Electronic Working Papers 1997-007. http://www.sil.org/silewp/1997/007/SILEWP1997-007.html.

Comrie, B., Haspelmath, M. and Bickel, B. 2015. The Leipzig Glossing Rules. Conventions for Interlinear morpheme-by-morpheme glosses. Weblink: http://www.eva.mpg.de/lingua/resources/glossing-rules.php

Goldsmith, John. 1990. *Autosegmental and metrical phonology*. Oxford: Blackwell.

Hyman, Larry, and Florian Lyonnet. 2011. Metatony in Abo (Bankon), A42. *UC Berkeley Phonology Lan Annual Report (2011),* 168-182.

Leben, William. 1973. *Suprasegmental phonology*. Ph.D. thesis, MIT, Cambridge, MA.

Lehmann, Christian. 2001. Language documentation: a program. In Walter Bisang (ed.) *Aspects of typology and universals*, 83-97. Berlin: Akademie Verlag.

Makasso, Emmanuel-Moselly. 2012. Metatony in Basaa. In *Selected Proceedings of the 42nd Annual Conference on African Linguistics*, ed. Michael R. Marlo et al., 15-22. Somerville, MA: Cascadilla Proceedings Project. www.lingref.com, document #2754.

Meeussen, A. E. 1967. Bantu grammatical reconstruction. *Africana Linguistica*, (3), 79–121. Nurse, Dereck. 2006. Focus in Bantu: verbal morphology and function. *ZAS Papers in Linguistics* 43, 189-207.

Stuart McGill (2009). Documenting grammatical tone using Toolbox: an evaluation of Buseman's interlinearisation technique. In Peter K. Austin (ed.) Language Documentation and Description, vol 6. London: SOAS. pp. 236 – 250.

TEI Guidelines. Weblink: http://www.tei-c.org/release/doc/tei-p5-doc/en/Guidelines.pdf.