# Representation of speech in CorpAfroAs

## Transcriptional strategies and prosodic units*

Shlomo Izre'el and Amina Mettouchi
Tel Aviv University and LLACAN, Paris

This paper surveys the transcriptional aspects of CorpAfroAs, a spoken corpus of Afroasiatic languages, with a focus on the representation of phonemes, morphemes, words, and longer units. We discuss the distinction between prosodic, phonological and morphosyntactic word, as well as that between intonation unit, paratone and period. Segmentation and transcription choices are analyzed and their outcome in terms of scientific breakthroughs is presented : the comparison between phonological and morphosyntactic word allows the systematic study of sandhi and other similar phenomena, and of the syntax/phonology interface. The segmentation into prosodic units allows the study of interfaces with syntax, information structure, and discourse.

## 1. Introduction

The spoken medium is acoustic, linear and temporally extended. Therefore, visual transmission is necessary in order to enable research, except, perhaps, for those focused on individual, small units. Even in this latter case, one needs to transmit sound into the visual medium in order to publish the results. The linguist must therefore use a transcript of the spoken text.

Transcribing a text is not a trivial undertaking, as has been noted time and again by those who have attempted an accurate transmission of speech into the written medium, i.e., its visualization. Transcribing a recording is a time consuming endeavor, and an hour of transcription may take — depending on the nature of

---

the speech — rate of speech, numbers of speakers, setting (naturally occurring or spontaneous), environment, genre, etc. — at least many dozens of hours of painstaking work, in some cases the amount of time invested will climb to hundreds of hours. While orthographic transcription is often used in languages with a written system and written tradition, transcription in the standard orthography has by its very nature a very limited range of uses for the analyst, and indeed seems to be most useful for lexical studies and discourse analysis, yet even there only with at least minimal prosodic notation. Other domains of linguistic analysis can hardly profit from using transcription in only the standard orthography of any speech, without having access to the sound stretch itself. This applies not only to phonetic or phonological analyses, but practically to all other domains, such as morphology, morphophonology, prosody, and even syntax. Notably, the standard orthography of a language is by definition related to only one (demographic or contextual) variety of linguistic forms used by speakers of that language. Moreover, the vast majority of the languages represented in CorpAfroAs have no orthographic standards or orthographies at all, which implies that other transcription systems should be used.

Therefore, CorpAfroAs is so constructed as to present to its end users both sound and transcription linked and aligned to the extent that each relevant unit of language can be easily retrieved and accessed together. The first (tx) tier presents a broad phonetic transcription, whereas the second tier (mot) brings forth a basically phonological representation of the same stretch. It is this second tier that forms the basis of all upper-level analyses, i.e., morphological, POS, and beyond. Both the phonetic and the phonological levels take cognizance of the prosodic structure of the language in terms of intonation units; the phonetic tier further exhibits lower level units, i.e., phonological words. One should note, however, that phonological words are not necessarily a lower level in the prosodic hierarchy, as will be clarified in section 2.1.

In what follows, we shall first discuss the representation of the segmental strings in different tiers. Then we shall discuss in some detail the theoretical basis of segmentation into prosodic units and its implications.

## 2.   Visualization of the spoken: Phones and segmental phonemes

The first (tx) tier presents a broad phonetic transcription of the speech stretch as actually perceived by the transcriber. In terms of sound, this tier conveys the sound segments at the surface level, i.e., after all phonological rules have been performed. Operations can be present in the creation of allophones, assimilation (total or partial), elision, or lengthening, shortening, etc. Analysis at the tx level is thus mostly phonetic, although it has much to do with the phonology of the language, as each represented segment actually stands for a class of phones which are related on both the phonetic and the phonological level (Wells 2006; Esling 2010: 680). Units at the tx level are phonological words (see below, §2.1).

A more abstract level of representation is presented at the mot tier. Each character at this level is thus ideally representing a phoneme. This transcription line does not represent any abstraction beneath the morphophonological level, i.e., it represents phonemic strings following the operation of morphophonemic rules. Analysis at this level is thus purely phonological, as allophonic variation, sandhi phenomena and their like are usually not shown. Ex. 1 from Gawwada will serve well to demonstrate the differences between the broad phonetic representation in the tx tier and the morphophonemic transcription represented in the mot tier:[1]

(1)   tx:    ħaːmos kujaʕte ogaːjb ano raːdonaba aqasi raːdonesi apaqasana
      mot:   ħaːmosí kujaʕte ʔokaːjpa ʔano raːtonepa ʔapaqasí raːtonesí ʔapaqasana
      mb:    ħaːmo=s-í kujaʕ-t-e ʔokaːj-i=pa ʔano raːton-e=pa ʔapaq-a=s-í raːton-
             e=s-í ʔapaq-a=s-a=n-a
      ge:    Haamo=DEICT-SPEC day-SING-F come-PFV.1SG=LINK IDP.1SG radio-
             F=LINK listen-IPFV.1SG=DEICT-SPEC radio-F=DEICT-SPEC listen-IPFV-
             1SG=DEICT-GEN=MOV-OUT
      ft:    'This Haamo came as I was listening to the radio. And as I was
             listening to the radio,' (GWD_MT_NARR_003_012)

In Ex. 2 from Moroccan Arabic, there are two occurrences of the definite article /əl/, one in each of the two words in this example. In the first occurrence, the vowel that usually precedes the consonant is now found following it: ləħbəq. It is thus duly represented in the tx tier, whereas the order has been reversed in the mot tier (əlħbəq), thus following the accepted representation of the Arabic definite article. In the second occurrence, the definite article is represented in both the tx and the mot tiers, as showing the morphophonemic change of /l/ to /s/ which occurs in adjacency to the following word, beginning with /s/. As this change is morphophonemic, it is similarly represented in the mot tier. The underlying phonemic string / əl / is represented in this case only in the mb tier.

(2)   tx:    ləħbəq wussuːsaːn
      mot:   əlħbəq wəssuːsaːn
      mb:    əl=ħbəq w=əl=suːsaːn
      ge:    DEF=basil and=DEF=lily
      ft:    'the basil and the lily' (ARY_AB_narr_1_004)

---

1.  In this section, we dispense with the prosodic notation of boundaries, which will be dealt with in §2.2.

Divergences from the principled system as characterized above can be discerned in some treatments of the languages represented in CorpAfroAs, notably with regard to vocalic epentheses, where theories may differ regarding their actual status. As Ex. 2 from Moroccan Arabic demonstrates, the theoretical premise that lies behind the representation /əl/ for the definite article in Moroccan Arabic is that the initial schwa is part of the phonemic string that forms this morpheme. In Hebrew, epenthesis usually takes the form [e]. However, scholars differ in their analysis and representation of various morphemes as regards the status of this vowel in the morphemic string, notably in the domain of prepositions. Note Ex. 3 from Hebrew, where the vowel [e] in the preposition [be] is interpreted as epenthetic:

(3)  tx:   beotobusim
     mot:  beotobusim
     mb:   b=otobus-im
     ge:   in=bus-M.PL
     ft:   'By buses.' (HEB_IM_NARR_7_SP1_0948)

While strict methodology would require the representation of /b/ as [b] rather than [be] in the **mot** tier, the reading of such a string will be misleading: [botobusim]. Therefore, it has been decided to copy the epenthetic [e] also to the **mot** tier.

Note also Ex. 4, exhibiting epenthetic vowels in Kabyle:

(4)  tx:   ikkrəd jufad jəssisulaʃiθənt
     mot:  ikkərdd jufadd jəssis ulaʃitənt
     mb:   i-kkər=dd j-ufa=dd jəssi-s ulaʃ=tənt
     ge:   SBJ3SG.M-stand_up\PFV=PROX SBJ3SG.M-find\PFV=PROX daughter\
           PL-KIN3SG NEGEXS=ABSV3SG.M
     ft:   "The father woke up and found that his daughters were no longer
           there" (KAB_AM_NARR_01_0902)

A strictly accurate representation of the phonemic string of the first word would yield /ikkrdd/. The long cluster of consonants would be hard to interpret. In this case, the final morpheme, /dd/, is represented as a cluster, still immediately following the final consonant of the verbal stem. It will be interpretable when compared to the **mb** tier. However, the **mb** tier cannot provide readability to the consonant cluster of the verbal stem, which has therefore been represented in all tiers along with an epenthetic vowel. The represented form, *kkər*, will serve as a basic allomorphic representation to the morpheme (=verbal stem) /kkər/, which also has the variants *əkkr* and *kkr*. The epenthesis can therefore appear in various places in the **tx** tier, but in the **mot** tier the stem is represented in a single form as above. From the technical point of view, only one record (=stem representation) will thus

be used in the ELAN lexicon, but the other forms will appear as variants of that record. In a similar vein, representation of the phonological structure of the absolutive clitic /tnt/ in Kabyle (at the end of this IU) will be unreadable, so that the epenthetic schwa, which is usually used in the pronunciation of this clitic [θənt] has been kept also in the **mot** tier: /tənt/.[2]

## 3.  Prosodic segmentation: Prosodic units and their representation

### 3.1  Phonological word

A *phonological word* is a unit consisting of one syllable or more which has at least one defining property chosen from the following areas: (1) Segmental features: internal syllabic and segmental structure; phonetic realization in terms of this; word boundary phenomena; pause phenomena. (2) Prosodic features: stress (or accent) and/or tone assignment; prosodic features such as nasalization, retroflexion, vowel harmony. (3) Phonological rules: some rules apply only within a phonological word; others (external sandhi rules) apply specifically across a phonological word boundary (Dixon & Aikhenvald 2002: 13).

There is no consensus over the definition of either a phonological word or a prosodic word. Definitions differ among linguistic schools, as well as within schools. For example, scholars of the generative school "differ in how function and content words are parsed into Prosodic Words, and also in how different types of morphemes are parsed into Prosodic Words" (Shattuck-Hufnagel & Turk 1996: 216–218). The issue of cliticization is often brought into account in determining the scope of the notion of prosodic word, without there being a consensus about its relevance to the definition of the notion of phonological word (*op. cit.*, §§3.1;3.2.4; Aikhenvald 2002; Vogel 2006: 532–3). Yet cliticization in itself is a complex feature in that the behavior of clitics should be regarded as language specific (Aikhenvald 2002; Schiering, Bickel & Hildebrandt 2010). Furthermore, cliticization is not invariable, and either content word or function words may have — under different conditions — both full and reduced versions (Zwicky 1977, 1995; Aikhenvald 2002: 72–75; Anderson 2005: §4). In any case, the relationship between prosody and morphosyntax plays a large role in the determination of prosodic words (Vogel 2006).

Units at the **tx** level are phonological ones. Units at the **mot** level are morphosyntactic words (or, as commonly called, "grammatical words"), preparing the

---

2. The fricative [θ] is a phonetic realization of the phoneme /t/, which is therefore used in the mot tier.

ground for morphological and morphophonological analyses which are operated while moving down to the **mb** tier, representing the morphemic structure of the language under scrutiny.

Whereas a phonological word may be defined on phonological or prosodic terms, a morphosyntactic word is defined on morphosyntactic terms as follows: it consists of a morpheme or several morphemes that (1) always occur together (rather than scattered through the clause); (2) occur in a fixed order; (3) have a conventionalized coherence and meaning (following Dixon & Aikhenvald 2002:19). As noted by Julien (2006:619), the rather commonly used term "grammatical word" to denote the nonphonological and nonlexical meaning of "word" is not strictly correct because phonology is, of course, also a part of grammar. Therefore, we will use the term "morphosyntactic word" instead (cf., e.g., Vogel 2006; Matthews 2007 s.v.; Crystal 2008: s.v.).

Ex. 5 from Hebrew is a clear illustration of the difference between morphosyntactic words (in the **mot** tier) and phonological (or prosodic) words (in the **tx** tier). Boundaries between either phonological words or morphosyntactic words are represented by spaces on the relevant tier. The vertical lines in Figure 1 show the boundaries between phonological words.

(5)  **tx:**   χalomʃlanu zʃtijelanu galeʁja
     **mot:**  χalom ʃelanu ze ʃetihje lanu galeʁja
     **mb:**   χalom ʃel=anu ze ʃe=t-ihje l=anu galeʁj-a
     **ge:**   dream of=POSS.1PL DEM.SG.M NMLZ=3SG.F-be\NFCT to=POSS.1PL
              gallery-F
     **ft:**   'Our dream is that we will have our (own) gallery.' (HEB_IM_
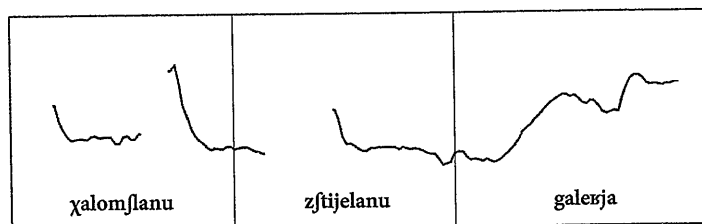              CONV_3_SP1_041)



Figure 1.  Phonological words in Hebrew (Ex. 5)

However, there are cases where morphosyntactic words and phonological words do not show morphosyntactic unity, and vice versa, a mismatch which is commonly attested in some languages (cf. Caink 2006:492). An interesting case is exhibited by Juba Arabic, an expanded pidgin of Southern Sudan. Given the lack of inflectional morphology in that language, phonological words often coincide with

grammatical words. If a prosodic word is defined by a stretch with only a single (main) stress, then reduplicated items can be seen as morphosyntactic words consisting of two prosodic (=phonological) words; e.g., bigídu~gídu 'pierce repeatedly' (JA_SM_CONV_2_SP2_299). On the other hand, a single prosodic word may consist of two morphosyntactic words; e.g. [jaʃán] /ja aʃán/ 'then because' (JA_SM_CONV_2_SP2_372).

Ex. 5 is, admittedly, an ideal representation of the tiers in the CorpAfroAs tier template. In practice, the **mot** tier exhibits a compromise between the morphosyntactic structure of the phonemic string, its actual pronunciation and its intermediary status between the transcription proper, on **tx**, and the morphemic analysis on **mb**.

Moreover, there are problems in determining and segmenting a text into phonological words. Such problems are not only the result of the diversity of languages represented on CorpAfroAs, or divergences in theoretical orientations of the respective schools involved and among individual scholars, but they are also inherent to the very issue of the definition of "phonological word", "prosodic word", and the relationship between those entities. Ex. 6 and Figure 2 exhibit cliticization is in Moroccan Arabic:

(6)   ħʃuːmiːja=u daːχla suːq ṛaːs=ha=u
      shyness=and come_in market head=her=and
      '(She was always in) modesty and minding her own business and ...'
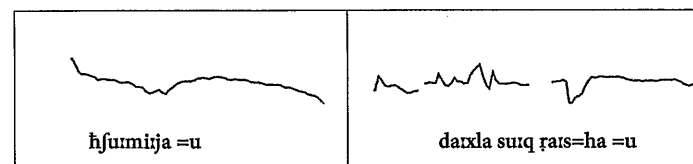      (ARY_AB_narr_1_029–030)



Figure 2.  Encliticization of the connective in Moroccan Arabic (Ex. 6)

The connective /u/ is usually regarded as having the tendency to cliticize to the following word (or unit).[3] However, in the two occurrences of the connective in this example, it is decisively cliticized to the preceding word. As suggested by the gloss, the analyst considers the connective to be a clitic not only on the phonological level, but also on a morphosyntactic level.

Another sort of problem can be illustrated by Ex. 7 (with Figure 3) from Ts'amakko:

---

3.  Thus, in a way, following the tradition of written Standard Arabic.

(7)   ˈqajto ˈχumɓiɣa ˈfugaɗɛ
      q'ajto χumɓi=ka pug-aɗ-aj
      time all=CONTR inflate-MID-IPFV.2SG
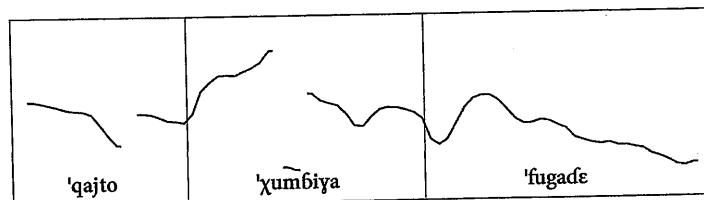      (What is that you eat and) 'you always get satiated?' (TSB_NARR_001_051)



**Figure 3.** Prosodic or phonological words (Ex. 7)

There are three content words in this intonation unit: q'ajto, χumɓi and pugaɗaj. From the prosodic point of view, they seem to form three prosodic words, where the contrastive focus marker is cliticized to the second content word, thus forming with it a single phonological word. This is indicated by the lack of stress on the clitic, as well as by the voicing and fricativization of its first consonant (k→ɣ\v_). One should, however, take into consideration the possibility that the stress on the first content word, namely, q'ajto, is a secondary stress, with the consequence that it be regarded as a single prosodic word with the following χumɓiɣa. The decision bears on the analysis of information structure of this string, i.e., whether the phonological compound as a whole is focused or only the second content word. From the perceptual point of view, the level of accent of the first word seems as prominent as in the third word, with only the second word showing more prominence. Therefore, the conclusion seems to be that the focal point of this intonation unit is on the second word, which conforms to the position of the segmental contrastive focus marker /ka/.

As we have seen, the initial consonant of the element /ka/ is fricativized in the process of cliticization. An interesting question then arises when one looks at the fricativization of the initial consonant of the last word, i.e., p→f. Should this change be interpreted as the result of cliticization or prosodic proximity between the second and the third word? In our opinion, this change can hardly suggest that we should regard the second and the third word as forming together a single phonological word, all the more so a single prosodic word. We should allow ourselves the liberty to interpret word-initial assimilation as this one as an external sandhi phenomenon.

Boundaries between prosodic word in particular and phonological words in general are not easy to detect (cf. Dixon & Aikhenvald 2002: 16; Fletcher 2010: §2; Basebøl 2000). As is clear from this example, sandhi phenomena may pose difficulties also in drawing morphosyntactic boundaries.

Giving attention to difficulties in boundary notations and the segmentation into prosodic words in particular and phonological words more generally, we propose that the units as represented in the tx tier may not be necessarily regarded as a lower level than Intonation Units in the prosodic hierarchy, although a rather widespread consensus may claim that they do, because prosodic and phonological units are usually not distinguished:

> The Phonological Word (or Prosodic Word) is located within the phonological hierarchy between the constituents defined in purely phonological terms (i.e., mora, syllable, foot) and those that involve a mapping from syntactic structure (i.e., clitic group, phonological phrase, intonational phrase, utterance). (Vogel 2006: 531)

The annotation of prosody in CorpAfroAs stops at the indication of boundaries. In their essence, the words contained in the tx tier are phonological and not strictly prosodic. The issue of determining prosodic or phonological words must be subject to further research.

### 3.2   Intonation unit

The units of the next level are *intonation units*. It has long been recognized that spoken language organizes itself in segments of speech that can be accounted for by their suprasegmental structure. The suprasegmental unit according to which segmentation of the spoken language can be made has been conceived to be dependent mainly on tone, or rather pitch, and has therefore been termed "tone group", "intonation group", "tone unit", "intonation(al) phrase", "intonation unit", or the like (e.g., Beckman & Pierrehumbert 1986; Halliday 1989; Selkirk 1984; Chafe 1994; Cruttenden 1997; Brazil 1997; Hirst & Di Cristo 1998; Fox 2000; Halliday 2004), where the identified prosodic stretch may be identical or different in some respects among the various approaches. Different paths have been used to explain the concept. Whatever approach is taken, it seems that there is a wide consensus that the intonation unit (henceforth: IU) encapsulates a functional, coherent segmental unit, be it syntactic, semantic, informational, or the like. IUs are therefore the first level of units where alignment of sound and transcription is made in CorpAfroAs.

It seems commonly accepted that an IU is a coherent intonation contour, and some would define the IU in these terms (Chafe 1994; Du Bois *et al.* 1992, 1993; Tao 1996; etc.). An example of a prototypical coherent intonation contour can be seen in the pitch curve in Figure 4, depicting the intonation contour of a single IU from Beja cited as Ex. 8:

(8)   'hoːɖaːbiˈjaːjiːha
      hoːj ɖaːbiˈjaj iːha
      hoːj ɖaːb-iˈja-j iː-ha
      3ABL run-PFV.3SG.M-SUFX.PROG AOR.3SG.M-be
      'He managed to run away from there.' (BEJ_MV_NARR_01_shelter_092)
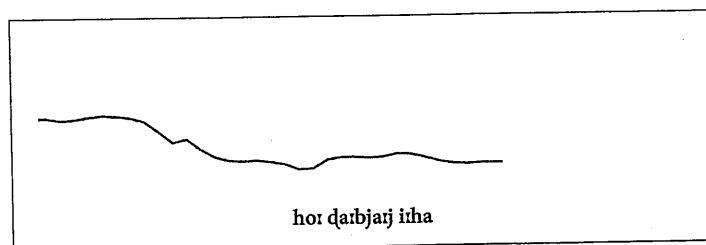


**Figure 4.** A coherent intonation contour (Ex. 8)

"A coherent intonation contour", while quite easily perceivable, is rather hard to define in itself by acoustic, formal terms, nor is it easy to define an IU by any other internal criteria (Cruttenden 1997). In practice, segmentation of a discourse flow into IUs is made by detecting their boundaries, whereas internal criteria are brought into consideration only secondarily (Cruttenden 1997). This practice has been used successfully in transcribing large corpora (Du Bois *et al.* 1992, 1993; Du Bois 2004; Cresti & Moneglia 2005; cf. also Cheng, Greaves & Warren 2005, following the methodology of Brazil 1997). Theory has also inclined towards the delimitation of the intonation unit — or "intonational phrase" — by reference to "boundary tones": "Each *intonational phrase* provides an opportunity for a new choice of tune, and ... some parts of the tune serve to mark the *phrase boundaries*" (Pierrehumbert & Hirschberg 1990, 272); "Rappelons que le rapport de dominance dépend uniquement des tons finals; il est insensible aux éléments intonatifs apparaissant ailleurs dans le groupe" (Blanche-Benveniste *et al.* 1990: 172). A useful account of the study of prosodic structures will be found in Fox 2000; see also Beckman and Venditti 2010.

Segmentation into IUs in CorpAfroAs was carried out applying both external and internal criteria, i.e., by detection boundaries of IUs and by looking at the internal structure of the pitch contour. Following previous research in various languages, we have decided to use four major perceptual and acoustic cues for boundary recognition as follows: (1) final lengthening; (2) initial rush; (3) pitch reset; (4) pause (cf. Cruttenden 1997; Du Bois *et al.* 1992; Hirst & Di Cristo 1998). It should be noted that the threshold over which we consider that the pause is significant have been set between 100 and 200 ms, depending on genre, language and rhythm of speech (see CorpAfroAs' manual). The internal criteria used

— apart from an impressionistic-perceptual conception of a contour — were: (1) declination (Cruttenden 1997: §§4.4.4.4, 5.5.1; Wichmann 2000: §5.1.1; Fox 2000: §5.5.5; also called "downdrift", Fox 2000: §4.2.2.3); (2) tonal parallelism, or isotony (Wichmann 2000: §4.3; Du Bois 2004). One may perhaps note at this juncture that the number of (morphosyntactic) words within an IU as exhibited in the CorpAfroAs texts is small, ranging between 1 and 7 (in extreme cases), with an average of ca. 2 to 4, depending on language and genre (for other languages see, *inter alia*, Chafe 1994: 64–65, 148).

None of the four cues for prosodic boundaries is in itself a necessary or sufficient cue for the existence of an IU boundary, and languages may differ in their most prominent cue for delimitation of IUs (Hirst & Di Cristo 1998). This is the case also with the Afro-Asiatic languages represented in CorpAfroAs. Previous research on Hebrew has shown that tempo, notably final lengthening, is the highest in hierarchy among acoustic features presented at an IU boundary, whereas pause occupies the last position in this hierarchy (Amir, Silber-Varod & Izre'el 2004; endorsed in the CorpAfroAs research). Pauses, however, have been shown to be a prominent cue in perception of IU boundaries in both Hebrew and Kabyle (Mettouchi *et al.* 2007), as is the case with some other language in the CorpAfroAs sample (e.g., Ts'amakko, Juba Arabic). Some CorpAfroAs researchers have noted different hierarchies among acoustic features for their languages while working on transcription and segmentation; e.g. in the Ts'amakko and Juba Arabic subcorpora, pitch reset is the most frequent cue, whereas pause is the most perceptually prominent; the Moroccan Arabic subcorpus seems to favor pause as its most frequent cue, whereas the most perceptually prominent cue is pitch reset). Minor (= non-terminal) boundaries and major (= terminal) boundaries may differ in this hierarchy. Furthermore, pause may be interpreted as indicating a major boundary, thus overpowering the final tone movement in some cases. Genre or style of speech, among other features, may also exhibit divergent hierarchies.

In Ex. 9 (and Figure 5) from Hebrew, the boundary between the first and the second IUs shows all four cues: lengthening of the last syllable of the first IU, fast-rate production of the first syllables of the following IU, pitch reset from the level of 240 HZ at the end of the first IU to 145 Hz at the beginning of the second IU, and a 210 ms pause between the two units. All first three cues are presented also at the boundary between the second and the third IUs, but in this case there is no pause present. As for the internal criteria, this stretch rather clearly exhibits declination of the F0 contour on the second and third IUs, as well as, with some complication, also on the first IU. The final tone being high for the first two IUs, declination naturally stops before the respective final rises. One should further note that declination affects not only any single IU, but a sequence of IUs, forming together — as in this case — a paratone (see below).

(9)  χaʃuv ʃehu javin / ʃemeoto ʁega ʃehu halaχ / hakvutsa niʁet tov joteʁ //
'It is important that he understand, that since the minute he left — the group looks better.' (OM[=Omer 4.2: 1350"-1354"; CoSIH text)
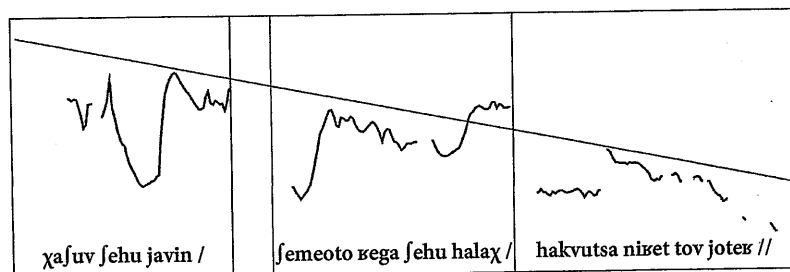


**Figure 5.** Intonation units: boundaries; declination (Ex. 9)

Isotony (Du Bois 2004), or tonic parallelism (Wichmann 2000), can be used to perceive an intonation contour, as it repeats itself in two or more adjacent IUs. This structure occurs notably in lists, but is found not infrequently also elsewhere, as in Ex. 10 (and Figure 6) and in Ex. 11 (and Figure 7), the first from Hebrew, the second from Ts'amakko:

(10)  haaχot baa / taktak nagaba / anijodaatma / tiplaba / vze /
'The nurse arrived, / just touched her / — whatever — / took care of her / and so on, / ...' (C514_309"-402"; CoSIH text)
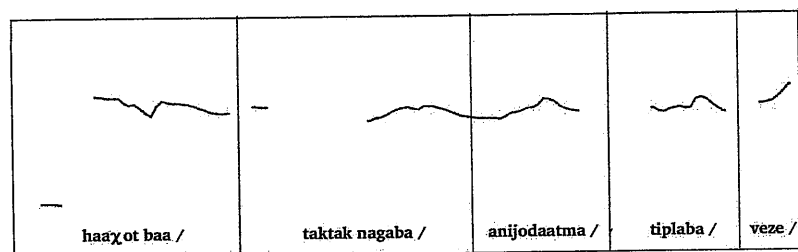


**Figure 6.** Isotony (tonic parallelism) in Hebrew (Ex. 10)

(11)  miːntea kaliatʃːo / kiniu ʃabbaje kiː //
miːnt-e=ʔa kaliatʃi-o kiniu ʃab~b-a=je kij-i
forehead-F=DEF anus-M M.POSS.3SG.F tie~PUNCT-JUSS.3SG.M=EMPH say-PFV.3SG.M
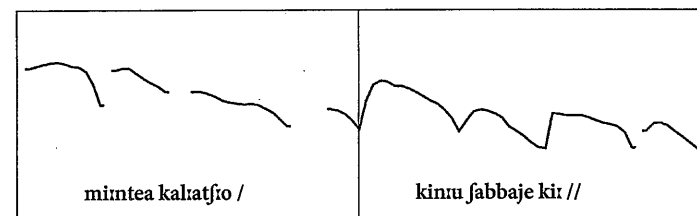'He said: "Let me tie your anus on the forehead."' (TSB_GS_NARR_001_128)



**Figure 7.** Isotony (tonic parallelism) in Ts'amakko (Ex. 11)

The final tone of an IU carries with it functional load in terms of discourse structure and information structure, with implications for syntax. For C-ORAL-ROM, the basic structural unit of spoken language is an "utterance", which is defined operatively as follows: "The operative definition of the utterance is such that every expression marked by a prosodic terminal break is an utterance" (Cresti & Moneglia 2005: 210). An utterance can include more than a single IU (referred to as an information unit), where the non-final IUs end with a "non-terminal" break. For the C-ORAL-ROM project,

> a prosodic break is considered terminal if a competent speaker assigns to it, according to his perception, the quality of concluding a sequence ... a prosodic break is considered non-terminal if a competent speaker assigns to it, according to his perception, the quality of being non-conclusive. (Cresti & Moneglia 2005: 17)

The reasoning behind this choice is the same as the one determined for CorpAfroAs:

> [T]he annotation of terminal and non-terminal breaks does not describe the prosodic movement that actually occurs in correspondence with a specific speech segment, but rather it selects the specific segment where, according to perception, a significant movement occurs. At the same time the annotation does not specify which proper speech act is performed by a sequence of word, but rather, specifies which sequence of words performs an act, for prosodic reasons. ... Once the relevant domain for prosodic movements and speech acts is determined, this will probably allow a better interpretation of both the relevant prosodic movements and the functional, dialogical value of the speech event. The same consideration can hold for syntactic features. Utterances cannot be identified and defined on the basis of syntactic properties as clauses can, for instance, but once an utterance is identified on the basis of a terminal break, any kind of morpho-syntactic and lexical evaluation can be driven on it. (Cresti & Moneglia 2005: 20)

It must be noted, however, that while Cresti and Moneglia have based their segmentation into prosodic units on speech act theory (*op. cit.*: 15 and note 17 on p. 67; 210), CorpAfroAs deliberately remains non-aprioristic in theoretical persuasion, left for its creators and end users for further research according to one's own individual stance.

CorpAfroAs concurs with the functional dichotomy between major and minor prosodic breaks, indicating terminal and continuing boundary tones by perception. Indicating boundary tones or breaks by perception has been proven reliable for C-ORAL-ROM (Cresti & Moneglia 2005: §1.2 and Appendix; Danieli *et al.* 2004). As it is not based solely on acoustic features but rather indicates functionality of the respective boundary tones as perceived by the annotator, the notation adopted for CorpAfroAs seems to be the best method for determining functional breaks, without any aprioristic ideas about the type of function involved. Still, for most subcorpora of CorpAfroAs, a concomitant acoustic check was carried out during the segmentation process and backed the perceptual indication of boundary breaks. In some cases the acoustic check served to refine prosodic notation; in other cases, it was an essential tool in the process, which was carried out using textgrids of Praat (see the CorpAfroAs manual, and Mettouchi & Chanard 2010). It should be noted that sometimes distinguishing between minor and major boundaries is not so easy, as there are cases where the final tone seems to be ambiguous. Major boundaries are usually better perceived than minor ones. On the other hand, syntax and discourse structure tend to influence this perception (cf. Mettouchi *et al.* 2007).

CorpAfroAs indicates minor boundaries by a single slash /, major boundary by a double slash //. Questions are indicated at the rx tier by the notation Q, irrespective of their segmental or prosodic structure. In Ex. 12 (Figure 8) from Hebrew, the first IU presents a minor boundary, both the second and the third major boundary, where the first of the two carries a rise indicating a yes/no question and the last one carries a falling tone:

(12)   ma /(Q) ataʁaχavta alsus //(Q) ʃloʃa jamim //
      sp2: "What? You rode on a horse?" — sp1: "Three days." (HEB_IM_
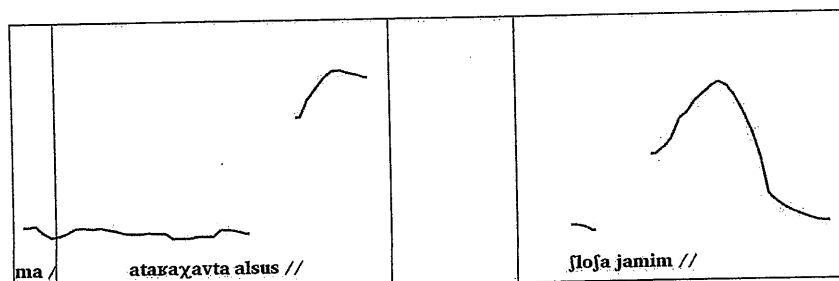      NARR_7_SP2_098 — SP1_0334)



**Figure 8.** Minor and major notations; notation of a question (Ex. 12)

Not all IUs fall into the two categories terminal vs non-terminal. Indeed there are incomplete IUs in all of the subcorpora of CorpAfroAs. IUs that have not come to completion can be of two types:

1.  An IU that has been truncated abruptly and can be perceived by prosodic cues like a shortening of a syllable or a part of a syllable, an additional glottal stop, along with a perceivable incompletion of a coherent intonation contour. Many an IU of this type will also end with a truncated word. This type of IU will be termed fragmentary (or truncated) and is marked in CorpAfroAs by a double crosshatch sign (##; a truncated word is indicated by a single crosshatch sign #). In Ex. 13 from Moroccan Arabic, a fragmentary IU which ends with a truncated word is indicated by both a single crosshatch sign and the double crosshatch sign, as explained in the CorpAfroAs manual:

(13)   ɣaːtəlqa fiːha tlaːta djə# ##
      ɣaː=t-lqa f=ha tlaːta djə# ##
      ɣaːtəlqa fiːha rəbʕa dəlkudjaːt //
      ɣaː=t-lqa f=ha rəbʕa d=əl=kudj-aːt //
      FUT=2-find\IPFV in =OBL.3SG.F three djə# ##
      FUT=2-find\IPFV in =OBL.3SG.F four of=DEF=hill-PL
      'You will find in it three of th- You will find in it four of the hills.' (ARY_AB_
      narr_1_406–407)

In Ex. 13, the speaker has corrected the number from '3' to '4', having noticed her mistake only after she had already started to utter the following word. The truncation of the (phonological) word *dəlkudjaːt* is accompanied and perceived by the palatalization of the dental stop [d] (/d/ 'of') and by a glottal stop following the schwa (not indicated in the transcription), which is the first segment of the definite article *əl*.

2.  An IU that seems to have been meant to continue and therefore shows a non-abrupt intonation contour, mostly (or always) carrying a continuing boundary tone. Still, the following IU seems not to be a continuation of this IU but starts a new stretch of speech. This new stretch of speech can be perceived as such by some prosodic cues (notably a long pause or hesitation phenomena; cf. Silber-Varod 2010; 2011; with previous references), by its syntactic structure, or by its semantic or pragmatic contents. In such instances, speakers may continue the stretch of speech, restart it or some part of it, or start a stretch of speech similar to the one already found in the suspended unit or any other unit before it, for example by rephrasing it. Alternatively, they can start a new sequence altogether. An IU of this type will not be regarded as truncated or fragmentary, but has been termed "suspended". We should notice that prosodic structures of the so-called suspended IUs seem not to differ from prosodic contours of minor IUs. In fact, speakers tend

at times to use suspension also as a discourse strategy, and therefore it would be a mistake to look at such IUs as representing cognitive failure. In Ex. 14, uttered by the same speaker who has contributed our Ex. 13, the first IU is truncated, the second is suspended.

> (14)   ɣiːɾ hiːja kaːɪt# ## (pause = 0.515 sec)
>        ɣiːɾ hiːja kaː=t-# ##
>        only 3SG.F REAL=3F-# ##
>        kaːɪnt diːɪːma fə fəːɪ /
>        kaːɪn-ət diːma fə fə /
>        be\PFV-3F always in in /
>        'She was always in in …' (ARY_AB_narr_1_026–027)

One should note that suspension is not a prosodic feature and is not recognized by prosodic cues. As mentioned above, the stretch of speech following a suspended unit cannot always be regarded as a direct continuation of the discourse presented at the suspended unit, either from the syntactic point of view or from the semantic point of view. In such cases, the discourse can resume (although it does not have to), either in close proximity to the suspended unit or at some distance from it, e.g., after a short or long parenthesis. Therefore, we have preferred the term "suspension" over other terms, such as "abandoned (unit)". The term "suspension" or "suspended (unit)" was chosen because the discourse can resume, after a false start, after a parenthesis (that can be long) or not resumed.

Segmentation into IUs was implemented in CorpAfroAs because text-sound indexation was a necessity. Several options were available: interpausal indexation, indexation of episodes, random chunking, periods, paratones etc. The choice of Intonation Units was favored in view of its possible correlation with morphosyntactic issues, one of the aims of the CorpAfroAs project being to analyze prosody and morphosyntax in several Afroasiatic languages. This type of segmentation proved very valuable for the study of grammatical relations: Mettouchi (to appear) shows that in Kabyle (Berber), the grammatical relations subject and object are only transparently coded within the intonation unit containing the verb. In the same paper, the author also shows that the interaction of IU boundaries, linear order and state are the building blocks of information structure constructions in Kabyle. Other papers using CorpAfroAs data (see for instance Malibert & Vanhove this volume, Caron, Lux, Manfredi & Pereira, this volume) use the IU as a unit for the study of information structure and syntactic dependency. This kind of segmentation is therefore relevant for syntactic and pragmatic studies.

## 3.3 Paratone

The next level is a *paratone*. The term "paratone", or "paratone group", coined on the analogy of the term "paragraph", has been used by some authors for the idea of a coherent formal sequence of intonation units (Crystal 2008 s.v.). Fox (1973 and subsequent studies), along lines suggested by Palmer (1922: section XI; 1924: 21–23), conceives a paratone (or a paratone group) as a larger prosodic unit than a tone group (in our terminology: intonation units), where "one or more major tone-groups are optionally preceded and/or followed by minor tone-groups" (Fox 2000: 318). Brown (1977: §5.2.1), who worked on read aloud news items, has defined a "paratone" on the basis of the organizational pattern of tone groups:

> If we go on to study the organization of a whole news item we shall find that the final tonic syllable in the complete item is marked by an even bigger pitch movement. So all the tonic syllables of what we might call the 'paratone', after the model of 'paragraph', are grouped together. The function of this patterning is to signal to the listener which tone groups are joined together in some larger structure and where the end of the larger structure comes. (Brown 1977: 86–7)

As analyzed and exemplified, Brown's notion of "paratone" suggests a sequence of IUs forming a sentence-like stretch (Brown 1977: §5.1; cf. Brown, Currie & Kenworthy 1980: §2.3). The "paratone" is further defined as a prosodic unit that encompasses a discourse where a new topic is being introduced (Brown, Currie & Kenworthy 1980: §2.3 and §3.6.ii; Brown & Yule 1983: §3.6.2). According to Brown (1990: 92), "[t]he most obvious phonetic cues [for the recognition of a paratone] are the high placing of the onset of a paratone, the brevity of the pauses within it, and the gradual drift down in overall pitch height towards a low ending". In these terms, Brown's "paratone" is closer to the notion of the oral "paragraph" as described by Wichmann (2000; cf. her discussion of Brown's "paratone" in §5.2.1) and the notion of "period" as suggested below, §2.4.[4] Noticing this ambiguity in Brown's definition and criteria, Yule (1980) has suggested the notion of "major paratone" for a single-topic related stretch, whereas the notion of "minor paratone" has been left somewhat ambiguous (see further Brown, Currie & Kenworthy,

---

4.  The term 'période' is employed in the French tradition for the notion of a unit that is larger than a clause or a comparable unit of the spoken language, but the definition of this unit has been different among scholars (Avanzi, Benzitoun and Glikman 2007). Work in computational linguistics has come up with a set of parameters to detect *périodes* automatically (Lacheret and Victorri 2002). It seems to us that this set of parameters may fit — mutatis mutandis — a prosodic unit which is located in hierarchy between an intonation unit and what we have defined below as 'period'. However, the relationship between a 'période' defined in these or similar terms and a 'paratone' as defined here is still to be sought.

71, who define the difference between major and minor paratone by the strength of their respective prosodic cues).

In CorpAfroAs, a "paratone" can be defined as one or more IUs ending in a major (terminal) final boundary, where any (optional) previous IU carries a minor (continuing) boundary tone. In this we follow the path of C-ORAL-ROM, for which a similar sequence has been defined as "utterance", as we have seen above. As the paratone frequently conveys a unified and coherent idea, and as translation may need to capture the whole idea conveyed by a paratone rather than by any individual IU internal to this paratone, the current **ft** tier has for some languages been supplemented with an **mft** tier aligned on paratones rather than intonation units. A comparative study of the respective merits of those two alignments still has to be done. Empirically, it appears that a paratone-aligned translation is highly desirable for V-final languages, but no theoretical investigation of the reasons behind this preference has been conducted yet.

At the time of writing this report, no significant theoretical studies on paratones have been conducted on any of the CorpAfroAs languages apart from Hebrew. In a preliminary study based on data from *The Corpus of Spoken Israeli Hebrew (CoSIH)*, Izre'el (forthcoming) suggests an interface between prosodic, discursive and syntactic units where the paratone, which encapsulates a discourse unit termed 'utterance', is the default domain of the clause rather than the IU, as is accepted by many authors (e.g., Chafe 1994: 65–6; Kibrik & Podlesskaya 2006). This suggestion is based on empirical research on spontaneous, everyday speech, mostly conversations, where a significant percentage of the attested IUs encapsulate only part of the components within a clause, whereas the paratone can neatly be regarded as the default domain of a clause, although it can consist of several clauses forming together a single discourse unit.

Table 1 summarizes the interface between prosodic, discursive and syntactic units in spoken Hebrew.

Table 1. Prosodic, discursive and syntactic units in Hebrew.

| Syntactic units | Discourse units | Prosodic units |
|---|---|---|
| Clause/Clause Cluster | Utterance | Paratone |
| Phrase/Clause (/Clause Cluster) | Speech Group (one of two or more in an utterance) | Prosodic Group (one of two or more in a paratone) |

A prototypical paratone can be seen in Ex.15 (Figure 9):

(15)  tɪpaˈʃut mɔˈʁaχtala / ɛʲˈtɔnɪ / ˈdɛvɛk //
at paˈʃut mɔˈʁaχat al=ha= /itɔn / ˈdɛvɛk //
you_SGF simply spread_SGF on=the=newspaper glue
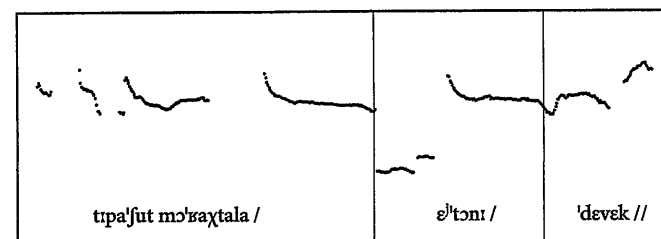'Do you just spread glue on the newspaper?' (C714_sp1_20–22; CoSIH text)



Figure 9. A prototypical paratone (Ex. 15)

Another typical paratone is shown in Ex. 16 (Figure 10). As is the case with the preceding example, Ex. 16 too exhibits a paratone composed of three IUs and consisting of a single clause:

(16)  vʷuaˈsa ɛtaisuˈmim / mbχinatsfiˈʁɔt / vɛaˈkɔl vɛaˈkɔl //
v=hu=aˈsa ɛt=ha=jisuˈmim / mi=bχinat=sfiˈʁɔt / v=ha=ˈkɔl v=ha=ˈkɔl //
and=he_did ACC=the=applications from=aspect_F_of counts and=the=all and=the=all
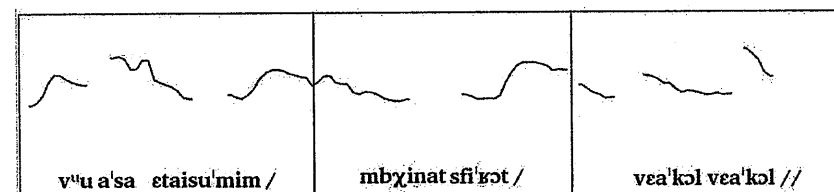'And he made the applications as regards counting and all?' (C612_4_sp2_115–7; CoSIH text)



Figure 10. A Typical paratone (Ex. 16)

Ex. 9, already cited above and cited again here as Ex. 17, consists of three IUs which include more than a single clause, forming together a discourse unit:

(17)  χaʃuv ʃehu javin / ʃemeoto ʁega ʃehu halaχ / hakvutsa niʁet tov joteʁ //
'It is important that he understand, that since the minute he left — the group looks better.' (OM[=Omer 4.2: 1350"-1354"; CoSIH text)
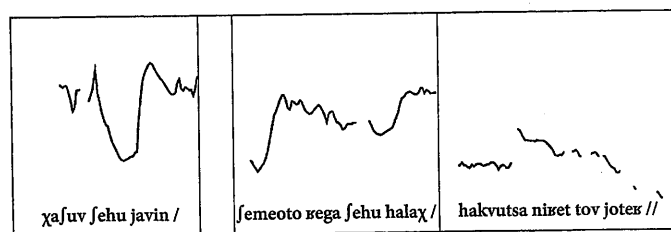
**Figure 11.** A paratone consisting of three intonation units including more than a single clause (Ex. 17)

Although much further research on the concept of paratone and its discourse parallel unit, the utterance, is still wanting, some further few comments may be in order now.

Although in general a paratone would be delineated by a perceivable major boundary, there are cases where a stretch of speech does not seem to carry a perceivable terminal tone, yet the continuing IU does not readily form part of one and the same paratone with that stretch. These are usually cases of fragmentary or suspended IUs. As explained above, a fragmentary IU is one that ends abruptly and has a perceivable prosodic cue(s) for truncation; a suspended IU is one that seems like a coherent minor IU, yet the following IU does not seem to be its direct continuation, either in prosodic terms or from the point of view of syntax, semantics or pragmatics. Therefore, paratones can also be perceived as either fragmentary or suspended.

While the end boundary of a paratone is easy to delineate, determining its beginning is somewhat more complex. As is obvious from the above, a new paratone may follow another paratone that — as defined — carries a major boundary tone, i.e., follows a major boundary. A paratone can further start after a fragmentary or suspended paratone (=IU), as is the case in the Ex. 18 from Hebrew:

(18)   veaz / kʃeχazaʁnu le / ulanbatoʁ / lakaχnu ta / tʁansibiʁit χazaʁa le / (pause)
       'And then / when we returned to / Olan Bator / we took the / Trans-Sibirian
       back to /' (suspension)
       lo jaʁadnu ##
       'We didn't get off' –– (truncation)
       lo imʃaχnu lebejdʒin / jaʁadnu po bebe# datong /
       'We did not continue to Beijin; we got off here, in Be- Datung,' (HEB_IM_
       NARR_7_SP1_0815–0823)

Ex. 18 exhibits two new starts. The suspended paratone and the new start following it are recognized by pause, rhythm change (length at the suspension point and rush in the following IU, in itself fragmented, with an immediate restart of yet a new paratone with a change in the lexicon. The truncated prosodic unit is

recognized as a separate unit by only a pitch reset at its right boundary, so its independent status is somewhat questionable.

In contrast to the above, a suspended unit can be shown to be an integral part of a single paratone, albeit not necessarily a coherent one. In Ex. 19, again from Hebrew, the speaker continues with a very similar topic as the one she was speaking about. Further, the speaker repeats the last word of the suspended unit and continues from there both syntactically and semantically, and in some way also prosodically. Moreover, the suspended unit ends with a level boundary tone which signals stronger continuation than a rising tone (Silber-Varod 2011).

(19)   ma jaani hem baim ## (pause) baim beeze ʃaloʃ babokeʁ e / aχaʁe
       miklaχat //
       what meaning they come ## (pause) come like three in_the_morning uh /
       after shower //
       'What? You mean, they come like three in the morning after shower?'
       (OCD6/1_41':20''-41':23''; CoSIH text)

Of course, the beginning of a discourse or a conversational turn will also start with a paratone. While this seems an obvious conclusion from the definition of a paratone, there are cases where a single paratone will be divided between interlocutors (Lerner 1996, 2004). Ex. 20 from Hebrew presents such a case.[5]

(20)   sp2:   az ze lo haja madʁiχ / ze haja paʃut em /
              'So, this was not a guide, it was just uh …'
       sp1:   miʃehu ʃe [ose et ze //
              someone who [does_this.'
       sp2:   [miʃehu / ʃe hovil otχem //
              '[someone who took you.'
       (HEB_IM_NARR7_sp2_166–167; sp1_0651)

This is an especially interesting case, as the speaker that started the paratone also continues it, but his interlocutor catches in the middle and continues the same paratone himself.

A significant prosodic cue for delineating paratones is the seemingly universal feature of declination. As declination is apparently a natural feature, it is discernible also in IUs (see above, §2.3). However, declination transcends IUs and is observable also in paratones, as well as in periods (see below, §2.4). In such cases, a pitch reset may occur between IUs comprising the paratone, but the overall curve will usually be lower in each IU than in the one preceding it. Ex. 9 above nicely

---

5. An opening bracket [ indicates the starting point of an overlap. (CorpAfroAs does not indicate overlaps as they are visually represented by the sound-transcription alignment.)

shows the feature of declination as it is observable in the paratone depicted there and in each of the three IUs that comprise this paratone.

Special cases are paratones with the insertion of parenthetical units. Some parentheses end with a major boundary, but they still show some prosodic cues like low pitch or reduced loudness that may enable us to regard the following units as continuing of an on-going paratone. Ex. 21 (and Figure 12) from Hebrew will illustrate the case:

(21)   [a] jeʃʃam paʁk /(pause) <creak> [b] *lo jodea ma* // [c] kama dunamim
tovim / (pause) [d] male male gumχot / im male / psalim ktanim / bealafim /
pesel eχad anak / tsiv# tsavua / lo tsavua / hakol budot //
ʻ[a] There is a park over there, <creak> — [b] *I'm not sure* — [c] (the size
of it is) a good number of acres; [d] (it has) many many alcoves, with
many small statues, by the thousand; (there was) one huge statue; (there
was another one) col- colored, (still one other) not colored; all (these are)
Budha(-statue)s.'
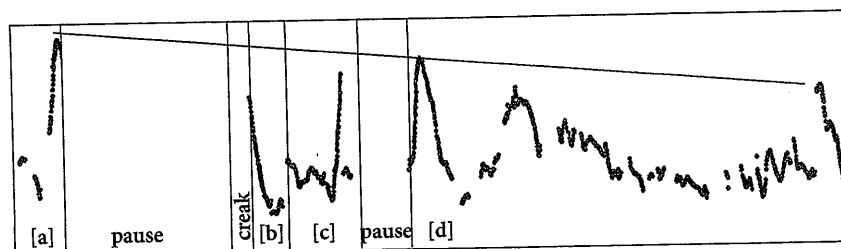(HEB_IM_NARR7_sp1_0837–0850)



Figure 12.  Paratone with a parenthesis inside (Ex. 21)

The second IU [b] (italicized), *lo jodea ma* // ʻI am not sure' ends in a major boundary, yet it is marked as a parenthesis by a low (and descending) pitch. The following unit [c], kama dunamim tovim / ʻa good number of acre' is still uttered in a low pitch, yet it rises at the end of the unit, indicating a return to the paratone stretch by a continuing (=minor) boundary. Indeed, the following IU [d] continues the paratone in both segmental and prosodic terms (note the declination throughout the paratone from the first IU [a]).

Parentheses in general, and the relationship between paratones and parenthetical units in particular, deserve special research (cf., *inter alia,* Barth-Weingarten, Dehé & Wichmann 2009; Debaisieux & Martin 2009).

Summing up, a paratone may be recognized by the following internal (1, 2) or external (3, 4) cues:

1.  If a paratone consists of either a single IU or of more than a single IU, it will show declination of the intonation curve throughout the entire stretch of the paratone. A change in the downdrift direction may occur if the last (or only) IU is an interrogative one ("yes/no" question) or other prosodically marked stretches such as exclamations or commands.
2.  If a paratone consists of more than a single IU, each of the non-final IUs composing this paratone will carry a minor boundary tone.
3.  A paratone begins following an IU ending in a major boundary tone; at the beginning of a discourse or at the beginning of a turn (unless shared by two interlocutors); following a fragmentary or a suspended IU (and therefore recognized mostly by non-prosodic features).
4.  A paratone ends in a major boundary tone. If fragmentary or suspended, the final boundary of a paratone can be discerned by prosodic cues (e.g., a long pause) or by noticing a new start in non-prosodic terms.
5.  A parenthesis ending in a major boundary may under certain conditions be inserted into a paratone.

The notion of paratone, as well as the prosodic and segmental criteria for defining paratone, still need much further research. It may perhaps be noted at this juncture that the number of IUs in a paratone as exhibited in the CorpAfroAs texts is usually small, depending on language and genre. In a significant number of cases, a paratone will consist of only a single IU; e.g., in the Hebrew part of CorpAfroAs, 37% of the paratones in the narrative texts and 49% of the conversational texts consist of only a single IU.

### 3.4  Period

A Period is the highest level in the prosodic hierarchy. A period will be defined as a speech stretch that shows declination along its paratones ("supradeclination" according to Wichmann 2000: §5.2.2), as well as by other prosodic means, e.g., isotony at specific defined stretches (cf. Martin 2009: §4.3). Contrary to the paratone, the period does not require that internal unit boundaries be continuing (minor) ones. A period encapsulates a "passage" in segmental terms (i.e., it shows some unity in syntactic, pragmatic or discursive structure, which is larger than an utterance). In a way, then, a spoken period can be compared to a written paragraph (Yule 1980; Brown & Yule 1983: §3.6.2; Wichmann 2000; cf. the discussion of "paratone" and "major paratone" in §2.3 above).

There is no reference to periods in the texts compiled and analyzed for CorpAfroAs, and the question remains a research topic for the future. Still

the following two examples, Ex. 22 (Figure 13) from Lybian Arabic and Ex. 23 (Figure 14) from Hebrew, will illustrate what can be referred to as a period.

(22)  haːda / ssaħləb / (pause) əːɪɪ jəʃərbuːh lamma fəʃʃte // (pause) ṣagaʃ //
      'This / salep — / they drink it during the winter. // Cold. //' (AYL_CP_narr_003_068–072)
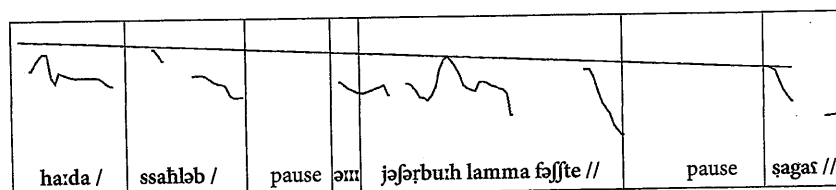


**Figure 13.**  Period in Lybian Arabic (Ex. 22)

(23)  ʃam amʁu ʃebenladen joʃev // ɪədaati biʃvilze hu nasa / χaʃvu ʃejitfesu oto //
      basof lo tafsu oto //
      'They said Ben Laden is residing there.// I think that this is why he went there. / They thought he will be caught. / At the end they did not catch him. //' (HEB_IM_NARR_7_SP2_021–024)
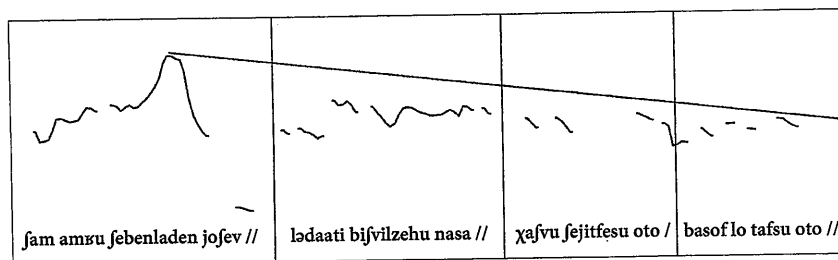


**Figure 14.**  Period in Hebrew (Ex. 23)

## 4.  Conclusions

This survey of the phonetic and transcriptional aspects of CorpAfroAs allows to sketch a portrait of the Corpus in terms of the choices that were implemented. First of all, the priority was given to the close relationship between the **tx** tier and the sound file, mirroring the structure of the software, in which **tx** is indexed to the sound file represented by the waveform window in ELAN. The transcription in **tx** was therefore meant to reproduce as faithfully as possible the spoken monologue or interaction, allowing the end-user to recognize the elements of the speech continuum. However, the length of the corpus does not allow detailed phonetic

representation, therefore, a degree of phonologization of the transcription was introduced, resulting in a broad phonetic transcription. In this tier, words are phonological (as opposed to morphosyntactic). The segmental string was segmented into prosodic units, defined by their boundaries and by their coherent internal contour. Intonation units were chosen over syntactic units (clauses or phrases) because they are the only organic units of speech. At a later stage, the corpus could be further segmented into other units if needed for further research on the correspondence between syntactic and prosodic units.

The **tx** tier was in turn further morpho-phonologized so that the **mot** tier should be composed of morphosyntactic words, morphemically transcribed. This level opens the way for a tokenization into morphemes in the **mb** tier. Those morphemes are then glossed in **ge** and **rx**. Finally, a free translation was given, which is currently aligned with respect to Intonation Units, but should ideally be aligned with respect to paratones rather than individual intonation units, because the latter often provide too small translation chunks which are difficult to organize together to form a coherent translation in the target language, English.

The process which led us to those decisions was based on some assumptions about the nature of speech, and on the research questions that interested us: the comparison between **tx** and **mot** for instance, allows the systematic study of sandhi and other similar phenomena, and of the syntax/prosody interface. The segmentation into prosodic units allows the study of various interfaces: syntax, information structure, discourse.

## References

Aikhenvald, Alexandra Y. 2002. Typological parameters for the study of clitics, with special reference to Tariana. In *Word: A Cross-linguistic Typology*, Robert M. W. Dixon & Alexandra Y. Aikhenvald (eds), 42–78. Cambridge: CUP.

Amir, Noam, Silber-Varod, Vered & Izre'el, Shlomo. 2004. Characteristics of intonation unit boundaries in spontaneous spoken Hebrew: Perception and acoustic correlates. In *Speech Prosody 2004, Nara, Japan, March 23-26, 2004: Proceedings*, Bernard Bel & Isabelle Marlien (eds), 677–680. <http://www.isca-speech.org/archive/sp2004/sp04_677.pdf>

Anderson, Stephen R. 2005. *Aspects of the Theory of Clitics* [Oxford Studies in Theoretical Linguistics 11]. Oxford: OUP. DOI: 10.1093/acprof:oso/9780199279906.001.0001

Avanzi, Mathieu, Benzitoun, Christophe & Glikman, Julie. 2007. Comment se comprendre sans se méprendre? L'exemple de trois termes problématiques: Période, parataxe et subordination inverse. In *Actes du 4ème Colloque Doctorants et Jeunes Chercheurs en Sciences du Langage (Coldoc'07) : Le vocabulaire scientifique et technique en Sciences du Langage, Nanterre, 20-21 juin 2007*. <http://www2.unine.ch/repository/default/content/sites/structuration_periodes/files/shared/articles_AM/AM_2007_benzitoun-glikman.pdf>

Barontini, Alexandrine. 2012. Moroccan Arabic Corpus. Corpus recorded, transcribed and annotated by Alexandrine Barontini. In Amina Mettouchi & Christian Chanard (eds). *The CorpAfroAs Corpus of Spoken Afro-Asiatic Languages.* DOI:http://dx.doi.org/10/1075/scl.68.website. Accessed on 10 January 2012.

Barth-Weingarten Dagmar, Dehé, Nicole & Wichmann, Anne (eds). 2009. *Where Prosody Meets Pragmatics* [Studies in Pragmatics 8]. Bingley: Emerald.

Basebøl, Hans. 2000. Word boundaries. In *Morphologie: ein internationales Handbuch zur Flexion und Wortbildung = Morphology: An International Handbook on Inflection and Word-formation* [Handbücher zur Sprach- und Kommunikationswissenschaft – Janbooks of Linguistics and communication Science 17(1)], Gert Booij, Christian Lehmann & Joachim Mugdan, in collaboration with Wolfgang Kesselheim & Stavros Skopeteas (eds), #40, 377–388. Berlin: Walter de Gruyter.

Beckman, Mary E. & Pierrehumbert, Janet B. 1986. Intonational structure in Japanese and English. *Phonology Yearbook* 3: 255–309. DOI: 10.1017/S095267570000066X

Beckman, Mary E. & Venditti, Jeniffer J. 2010. Tone and intonation. In *The Handbook of Phonetic Sciences*, 2nd edn [Blackwell Handbooks in Linguistics], William J. Hardcastle, John Laver & Fiona E. Gibbon (eds), 603–650. Chichester: Wiley-Blackwell. DOI: 10.1002/9781444317251.ch16

Blanche-Benveniste, Claire, Bilger, Mirelle, Rouget, Christine & Karel van den Eynde. 1990. *Le français parlé: Études grammaticales*, Participation de Piet Mertens [Sciences du Language]. Paris: CNRS Éditions.

Brazil, David. 1997. *The Communicative Value of Intonation in English.* Cambridge: CUP.

Brown, Gillian. 1977. *Listening to Spoken English* [Applied Linguistics and Language Study]. London: Longman.

Brown, Gillian. 1990. *Listening to Spoken English*, 2nd edn [Applied Linguistics and Language Study]. London: Longman.

Brown, Gillian, Currie, Karen L. & Kenworthy, Joanne. 1980. *Questions of Intonation.* London: Croom Helm.

Brown, Gillian & Yule, George. 1983. *Discourse Analysis* [Cambridge Textbooks in Linguistics]. Cambridge: CUP. DOI: 10.1017/CBO9780511805226

Caink, Andrew D. 2006. Clitics. In *Encyclopedia of Language and Linguistics*, 2nd edn, Keith Brown (ed.), 491–495. Oxford: Elsevier. DOI: 10.1016/B0-08-044854-2/00110-3

Chafe, Wallace. 1994. *Discourse, Consciousness, and Time: The Flow and Displacement of Conscious Experience in Speaking and Writing.* Chicago IL: The University of Chicago Press.

Cheng, Winnie, Chris Greaves & Martin Warren. 2005. *A Corpus-driven Study of Discourse Intonation: The Hong Kong Corpus of Spoken English* [Studies in Corpus Linguistics 32]. Amsterdam: John Benjamins. DOI: 10.1075/scl.32

*CoSIH: The Corpus of Spoken Israeli Hebrew (CoSIH):* <http://humanities.tau.ac.il/~cosih/english/index.html>

Cresti, Emanuela & Moneglia, Massimo (eds). 2005. *C-ORAL-ROM: Integrated Reference Corpora for Spoken Romance Languages* [Studies in Corpus Linguistics 15]. Amsterdam: John Bejnamins. DOI: 10.1075/scl.15

Cruttenden, Alan. 1997. *Intonation*, 2nd edn [Cambridge Textbooks in Linguistics]. Cambridge: CUP. DOI: 10.1017/CBO9781139166973

Crystal, David. 2008. *A Dictionary of Linguistics and Phonetics*, 6th edn. Oxford: Blackwell. DOI: 10.1002/9781444302776

Danieli, Morena, Garrido, Juan María, Moneglia, Massimo, Panizza, Andrea Quazza, Silvia & Swerts, Marc . 2004. Evaluation of consensus on the annotation of prosodic breaks in the Romance corpus of spontaneous speech 'C-ORAL-ROM'. In *Speech Corpus Production and Validation, LREC 2004: Fourth International Conference on Language Resources and Evaluation, 24th May, 2004, Lisbon*, Christoph Draxler, Henk van den Heuvel & Florian Schiel (eds). 1513–1516. < http://www.lrec-conf.org/proceedings/lrec2004/pdf/371.pdf >

Debaisieux, Jeanne-Marie & Martin, Philippe. 2010. Les parenthèses: Étude macrosyntaxique et prosodique sur corpus. In *La parataxe, Tome 1: Entre dépendance et intégration, Tome 2: Structures, marquages et exploitations discursives*, Marie-José Béguelin, Mathieu Avanzi & Gilles Corminboeuf (eds). Bern: Peter Lang.

Dixon, Robert M. W. & Aikhenvald, Alexandra Y. (eds). *Word: A Cross-linguistic Typology.* Cambridge: CUP.

Du Bois, John W., Cumming, Susanna, Schuetze-Coburn, Stephan & Paolino, Danae. 1992. *Discourse Transcription* [Santa Barbara Papers in Linguistics 4]. Santa Barbara CA: Department of Linguistics, University of California, Santa Barbara.

Du Bois, John W., Cumming, Susanna, Schuetze-Coburn, Stephan & Paolino, Danae. 1993. Outline of discourse transcription. In *Talking Data: Transcription and Coding in Discourse Research*, Jane A. Edwards, & Martin D. Lampert (eds), 45–89. Hillsdale NJ: Lawrence Erlbaum Associates.

Du Bois, John W. 2004. *Representing Discourse. Part 2: Appendices and Projects.* Santa Barbara CA: Linguistics Department, University of California. <http://www.linguistics.ucsb.edu/projects/transcription/representing>

Esling, John H. 2010. Phonetic Notation. In *The Handbook of Phonetic Sciences*, 2nd edn [Blackwell Handbooks in Linguistics], William J. Hardcastle, John Laver & Fiona E. Gibbon (eds), 678–702. Chichester: Wiley-Blackwell. DOI: 10.1002/9781444317251.ch18

Fletcher, Janet. 2010. The Prosody of Speech: Timing and Rhythm. In *The Handbook of Phonetic Sciences*, 2nd edn [Blackwell Handbooks in Linguistics], William J. Hardcastle, John Laver & Fiona E. Gibbon (eds), 523–602. Chichester: Wiley-Blackwell.

Fox, Anthony. 1973. Tone sequels in English. In *Archivum Linguisticum* 4 (new series): 17–26.

Fox, Anthony. 2000. *Prosodic Features and Prosodic Structure: The Phonology of Suprasegmentals.* Oxford: OUP.

Halliday, Michael A.K. 1989. *Spoken and Written Language*, 2nd edn. Oxford: OUP.

Halliday, Michael A. K. 2004. *An Introduction to Functional Grammar.* 3rd edn revised by Christian M. I. M. Matthiessen. London: Arnold.

Hirst, Daniel & Di Cristo, Albert (eds). 1998. *Intonation Systems: A Survey of Twenty Languages.* Cambridge: CUP.

Izre'el, Shlomo. Forthcoming. Basic units of language: Prosody, discourse and syntax. In *Researching Spoken Hebrew*, Einat Gonen (ed.). <http://www.academia.edu/2195304/> (In Hebrew; English version in preparation).

Julien, Marit. 2006. Word. In *Encyclopedia of Language and Linguistics*, 2nd edn, Keith Brown (ed.), 617–624. Oxford: Elsevier. DOI: 10.1016/B0-08-044854-2/00130-9

Kibrik, Andrej A. & Podlesskaya, Vera I. 2006. Problema segmentacii ustnogo diskursa i kognitivnaja sistema govorjashchego (Segmentation of spoken discourse and the speaker's cognitive system). In *Kognitivnye issledovanija*, Vol. 1, Valerij D. Solovyev (ed.), 138–158. Moscow: Institut psixologii RAN. <http://www.philol.msu.ru/%7Eotipl/new/main/people/kibrik-aa/files/Segmentation_discourse@Cognitive_studies_2006.pdf>; English summary: Discourse as a kind of cognitive activity: The principles of segmentation. In *The Second*

*Biennial Conference on Cognitive Science, June 9-13, 2006, St. Petersburg, Russia, Abstracts*, Vol. 2, 501–503.

Lacheret, Anne & Victorri, Bernard. 2002. La période intonative comme unité d'analyse pour l'étude du français. In *Verbum 24/1-2: Y a-t-il une syntaxe au-delà de la phrase?*, Michel Charolles, Pierre Le Goffic & Mary-Annick Morel (eds), 55–72.

Lerner, Gene H. 1996. On the 'semi-permeable' character of grammatical units in conversation: Conditional entry into the turn space of another speaker. In *Interaction and Grammar*, Elinor Ochs, Emanuel A. Schegloff & Sandra A. Thompson (eds), 238–276. Cambridge: CUP. DOI: 10.1017/CBO9780511620874.005

Lerner, Gene H. 2004. Collaborative Turn Sequences. In *Conversation Analysis: Studies from the First Generation* [Pragmatics & Beyond New Series 125], Gene H. Lerner (ed.), 225–256. Amsterdam: John Benjamins. DOI: 10.1075/pbns.125.12ler

Malibert-Yatziv, II -II. 2012. 'Hebrew Corpus'. Corpus recorded, transcribed and annotated by II-II Malibert-Yatziv. In Amina Mettouchi & Christian Chanard (eds). *The CorpAfroAs Corpus of Spoken AfroAsiatic Languages*. DOI: http://dx.doi.org/10.1075/scl.68.website. Accessed on 10 January 2012.

Manfredi, Stefano. 2012. 'Juba Arabic Corpus', Corpus recorded, transcribed and annotated by Stefano Manfredi. In Amina Mettouchi & Christian Chanard (eds). *The CorpAfroAs Corpus of Spoken AfroAsiatic Languages*, DOI: http://dx.doi.org/10.1075/scl.68.website. Accessed on 10 January 2012.

Martin, Philippe. 2009. *Intonation du français* [Collection U • Linguistique]. Paris: Armand Colin.

Matthews, Peter. H. 2007. *Oxford Concise Dictionary of Linguistics*, 2nd edn [Oxford Paperback Reference]. Oxford: OUP.

Mettouchi, Amina. 2012. 'Kabyle Corpus'. Corpus recorded, transcribed and annotated by Amina Mettouchi. In Amina Mettouchi & Christian Chanard (eds). *The CorpAfraAs Corpus of Spoken AfroAsiatic Languages*. DOI: http://dx.doi.org/10.1075/scI.68.website. Accessed on 10 January 2012.

Mettouchi, Amina. To appear. The Interaction of State, Prosody and Linear Order in Kabyle (Berber): Grammatical relations and information structure. In *Data and Perspectives in Afroasiatic*, Alessandro Mengozzi & Mauro Tosco (eds). Amsterdam: John Benjamins.

Mettouchi, Amina & Chanard, Christian. 2010. From fieldwork to annotated corpora: The CorpAfroAs project. *Faits de Langues – Les Cahiers* 2: 255–266.

Mettouchi, Amina, Lacheret-Dujour, Anne, Silber-Varod, Vered & Izre'el, Shlomo. 2007. Only prosody? Perception of speech segmentation. In *Nouveaux cahiers de linguistique française* 28: Interfaces discours – prosodie: Actes du 2ème Symposium international and Colloque Charles Bally, 207–218. <http://clf.unige.ch/display.php?numero=28&idFichier=109>; sound files and transcriptions: < http://clf.unige.ch/annexe.php?article=108>

Palmer, Harold E. 1922. *English Intonation; with Systematic Exercises*. Cambridge: Heffer & Sons.

Palmer, Harold E. 1924. *A Grammar of Spoken English: On a Strictly Phonetic Basis*. Cambridge: Heffer & Sons.

Pereira, Christophe. 2012. 'Tripolinian Arabic Corpus'. Corpus recorded, transcribed and annotated by Christophe Pereira. In Amina Mettouchi & Christian Chanard (eds). *The CorpAfroAs Corpus of Spoken AfroAsiatic Languages*. DOI: http://dx.doi.org/10.1075/scI.68.website. Accessed on 10 January 2012.

Pierrehumbert, Janet & Hirschberg, Julia. 1990. The meaning of intonational contours in the interpretation of discourse. In *Intentions in Communications* [Systems Development Foundation Benchmark Series], Philip R. Cohen, Jerry Morgan & Martha E. Pollak (eds), 271–311. Cambridge MA: The MIT Press.

Savà, Graziano. 2012. 'Ts'amakko Corpus'. Corpus recorded, transcribed and annotated by Graziano Savà. In Amina Mettouchi & Christian Chanard (eds). *The CorpAfroAs Corpus of Spoken AfroAsiatic Languages*. DOI: http://dx.doi.org/10.1075/scl.68.website. Accessed on 10 January 2012.

Schiering, René, Bickel, Balthasar & Hildebrandt, Kristine A. 2010. The prosodic word is not universal, but emergent. *Journal of Linguistics* 46: 657–709. DOI: 10.1017/S0022226710000216

Selkirk, Elisabeth. 1984. *Phonology and Syntax: The Relation between Sound and Structure*. Cambridge MA: The MIT Press.

Shattuck-Hufnagel, Stefanie & Turk, Alice E. 1996. A prosody tutorial for investigators of auditory sentence processing. *Journal of Psycholinguistic Research* 25: 193–247. DOI: 10.1007/BF01708572

Silber-Varod, Vered. 2010. Phonological aspects of hesitation disfluencies. *Proceedings of Speech Prosody 2010, Chicago*. < http://www.speechprosody2010.illinois.edu/papers/100020.pdf >

Silber-Varod, Vered. 2011. *The SpeeCHain Perspective: Prosodic-Syntactic Interface in Spontaneous Spoken Hebrew*. PhD dissertation, Tel-Aviv University. <http://www.openu.ac.il/Personal_sites/vered-silber-varod/download/Vered%20Silber-Varod%20Dissertation-7.pdf>

Tao, Hongyin. 1996. *Units in Mandarin Conversation: Prosody, Discourse, and Grammar* [Studies in Discourse and Grammar 5]. Amsterdam: John Benjamins. DOI: 10.1075/sidag.5

Tosco, Mauro. 2012. 'Gawwada Corpus'. Corpus recorded, transcribed and annotated by Mauro Tosco. In Amina Mettouchi & Christian Chanard (eds). *The CorpAfroAs Corpus of Spoken AfroAsiatic Languages*. DOI: http://dx.doi.org/10.1075/scI.68.website. Accessed on 10 January 2012.

Vanhove, Martine. 2012. 'Beja Corpus'. Corpus recorded, transcribed and annotated by Martine Vanhove. In Amina Mettouchi & Christian Chanard (eds). *The CorpAfroAs Corpus of Spoken AjraAsiatic Languages*. DOI: http://dx.doi.org/10.1075/scI.68.website. Accessed on 10 January 2012.

Vogel, Irene. 2006. Phonological words. In *Encyclopedia of Language and Linguistics*, 2nd edn, Keith Brown (ed.), 531–534. Oxford: Elsevier. DOI: 10.1016/B0-08-044854-2/00043-2

Wells, John C. 2006. Phonetic transcription and analysis. In *Encyclopedia of Language and Linguistics*, 2nd edn, Keith Brown (ed.), 386–396. Oxford: Elsevier. DOI: 10.1016/B0-08-044854-2/00014-6

Wichmann, Anne. 2000. *Intonation in Text and Discourse: Beginnings, Middles and Ends* [Studies in Language and Linguistics]. Harlow: Pearson Education.

Yule, George. 1980. Speakers' topics and major paratones. *Lingua* 52: 33–47. DOI: 10.1016/0024-3841(80)90016-9

Zwicky, Arnold M. 1977. *On Clitics*. Bloomington IN: Indiana University Linguistics Club.

Zwicky, Arnold M. 1995. What is a clitic? In *Clitics: A Comprehensive Bibliography, 1892-1991* [Library & Information Sources in Linguistics Series 22], Joel Ashmore Nevis, Brian D. Joseph, Dieter Wanner & Arnold M. Zwicky (eds), xii–xx. Amsterdam: John Benjamins. DOI: 10.1075/lisl.22