# Query languages and linguistics categories: annotating and searching spoken electronic corpora in indigenous languages

*Amina Mettouchi (EPHE, PSL & CNRS LLACAn)*
*Christian Chanard (CNRS LLACAn)*

Computer science and Linguistics have had regular interactions ever since the exponential development of information technology. Their relationship can be one of application (of IT methods to linguistic data), but more fundamentally, they are related together by the fact that at the core of programming, there is language. This study will explore the topic of the conference under that angle: how can query languages used in computer queries (Fig.1) help formulate, test and falsify linguistic claims about morphosyntactic constructions?

For such questions to be investigated, it is essential to work on annotated electronic corpora. Their contribution to linguistics is in the systematicity of data treatment: for instance, viewed from a computer programming perspective, a "word" is a string of characters surrounded with a blank space at each end, and therefore, for a linguist transcribing recordings into a computer programme such as ELAN, the potential for accurate statistics of word counts relies from the start on decisions concerning word boundaries, which is one of the first steps in the analysis of an undescribed language. And every step of the transcription and annotation process of spoken data is a potential locus for fertile interaction of programming and linguistic analysis.

This will be demonstrated through the investigation of electronic spoken corpora from several indigenous languages (CorpAfroAs http://corpafroas.huma-num.fr & CorTypo http://cortypo.huma-num.fr), whose annotation template (Fig.2) allows the formulation of complex queries to answer such questions as e.g. "How often do nominal subjects co-occur with nominal objects in the same intonation unit?".

What is crucial however in this perspective, is for annotation to be not only systematic, but also transparent, and as non-aprioristic as possible. Indeed, a query will only retrieve results on the basis of annotation, so that defining what a nominal subject/object and a prosodic boundary are language-internally, is absolutely crucial for the results of the query to be valid and non-trivial. The talk will explore the way the annotation/query process itself, as an instance of the scientific method (Fig.3), participates in the analysis, yielding explicit, testable and falsifiable definitions, such as the following for Kabyle (Berber) nominal direct object "a noun in the absolute state, immediately following the verb, or following <the verb followed by a noun in the annexed state> or following <the verb followed by an adverb> or following <the verb followed by a postverbal negator>, within the prosodic group of the verb", a definition which in turn needs transparent definition and consistent annotation of such annotations/categories as Verb, Noun, Absolute/Annexed state, Adverb, and Postverbal Negator.

This incremental procedure allows the testing and refining of hypotheses in view of the grammatical analysis of the language, while also questioning the process of categorization itself, and ultimately, the nature of linguistic data, and the specificity of linguistics as a science, namely the fact that its metalinguistic toolkit is linguistic in nature. Beyond the material discussed here, this reflexive nature of linguistics makes it relevant to many fields in the humanities, such as philosophy, human cognition, and the anthropology and sociology of science.

(Figures on page 2)

**Figure 1: Query language 'search N followed by PRO or CONJ in the same unit'**

[rx=\bN\b < tx=.] {rx=1 & tx=0} [rx=PRO|CONJ <tx=.] or
tx=0:[rx=\bN\b]{rx=1} [rx=PRO|CONJ]



**Figure 2: Annotation template in CorpAfroAs**



| ref | identifier for the annotation unit (time-associated) |
|---|---|
| tx | transcription in broad IPA into phonological words |
|    mot | intermediary tier with segmentation into morphosyntactic words |
|       mb | morphophonological transcription into morphemes |
|       ge | morpheme-by-morpheme gloss of mb |
|       rx | part-of-speech and other information relevant for retrieval purposes |
| ft | free translation into English |

**Figure 3: The Scientific Method in Corpus Annotation**