**Description**

**CSB1.0:** Continuous Speech Berber is a set of sentences representing Berber language. We present six sentences of duration 3s and 4 seconds.

*Emphasis in Berber*

A typical feature of Afrasian languages is the series of emphatic consonants, in contrast with voiced and voiceless series. Emphasis is glottalic in Ethiopic Semitic, pharyngeal in Arabic or uvular in Tuareg. In Niger Tuareg, but not in Tahaggart, the presence of an emphatic phoneme in a root or an affix triggers a complete spread of emphasis on the whole consonantal skeleton

**Speech dataset**

*Tamazight* (the Berber word for language) covers a vast geographical area: all of North Africa, the Sahara, and a part of the West African Sahel. But the countries principally concerned are, by order of demographical importance: Morocco (35 to 40% of the total population), Algeria (25% of the population), Niger and Mali (Tuaregs) [7].

We chose six speakers for the Berber language (kabyle) which are: Nacera, Zohra,Gassi, Kaci, Ghania and Smail, all native to the region of Kabylia. These speakers repeated with an average speed and an average energy six sentences in continuous speech. The sound card we used is called MobilePre USB. MobilePre USB is a preamplifier mobile integrating an audio interface perfect for the records on computer (laptop or mobile).

The dataset includes speech signals from six (6) different subjects. The speech signals are acquired during 2 s or 3s with different sessions to consider all variations at a sampling rate 16 KHz. After, a manual segmentation with "wavesurfer software" is done to consider only the part of produced voice. During the recording, each repetition has been analyzed to ensure that the entire sequences have been properly stored and avoiding any external interference. The chosen sequences are as follows:

"assagui yelha lhal";
 "azeka anrouh arthemourth";
"anzoum Remthane g nevdhou";
"yetchak izeme yellozen";
"yeghine wagroud aghour yemasse" and
"thelha thazalith gakhem rebi".

The first four sentences have duration of two seconds while the last sequences have equal duration of 3 seconds. The segmentation has been performed manually by using the recording software Wavesurfer in order to take only the essential part that has been produced.