

The Maaloula Aramaic Speech Corpus (MASC): From Printed Material to a Lemmatized and Time-Aligned Corpus

Ghattas Eid, Esther Seyffarth, Ingo Plag

Heinrich-Heine-Universität Düsseldorf

ghattas.eid@hhu.de, esther.seyffarth@hhu.de, ingo.plag@uni-duesseldorf.de



Motivation and corpus content

Western Neo-Aramaic is an endangered Neo-Aramaic (Semitic) variety spoken in the Syrian villages of Maaloula, Jubbaadin, and Bakhaa. ‘Maaloula Aramaic’ refers to the dialect of Maaloula.

Neither an electronic text corpus nor a speech corpus with time-aligned transcriptions was available. Therefore, we designed the Maaloula Aramaic Speech Corpus (MASC, Eid et al., 2022), the first electronic speech corpus of this variety.

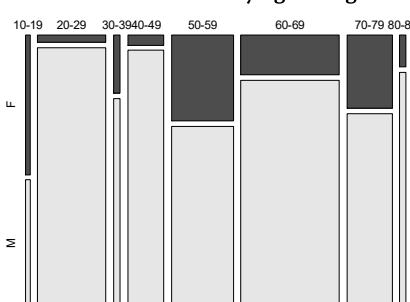
MASC is a multi-purpose corpus available to the scientific community at <https://doi.org/10.5281/zenodo.6496714>

Data included in the corpus:

1. the transcriptions of tape-recorded narratives that Werner Arnold collected in the 1980s (Arnold 1991a, 1991b):

- 173 monologues (different text types)
- 64,845 tokens, 12,220 types
- tokens per narrative (M = 375, Min = 19, Max = 4,340)
- 45 speakers (32 males, 13 females) aged 13 - 89 years
- tokens per speaker (M = 1,441, Min = 42, Max = 10,688)

Distribution of tokens by age and gender



2. the audio files of these narratives, available at the Semitisch Tonarchiv website of Heidelberg University (Arnold, 2003): 173 mp3 files (10 hours)

- Permission to use the data was granted by Harrassowitz Verlag and Werner Arnold.

1. Transcriptions

These are the digitized transcriptions that contain no annotation (except for a few informative tags). The digitization process involved using an optical character recognition (OCR) program and collaborating with a native speaker consultant to produce an error-free text.

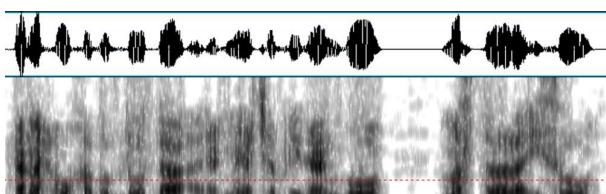
Extract from the MASC dataframe

Root	Lemma	LemmaFreq	Word_Form	Word_formFreq
?mr	amar ȳmar	1925	amella	185
?mr	amar ȳmar	1925	amelle	393
?mr	amar ȳmar	1925	amellen	2
...

3. Time-aligned transcriptions

The original and denoised audio files are included in our corpus. They can be opened in Praat (Boersma & Weenink, 2021) together with their TextGrid files to conduct acoustic analyses. For the time alignment process, we used the WebMAUS tool (Schiel, 1999, 2015) by BAS Web Services (Kisler et al., 2017).

Screenshot from Praat displaying the TextGrid tiers



anah hōxa b̄-blōta nmī?cabrill ?inbō maštra ra?isō l̄-blōta.

anah hōx b̄ blōta nmī?cabrill ?inbō maštra ra?isō
a n h o: b̄ b l o: t a n m i: ? s ? n b o: m a s̄ r a ? i s o:
a o: t a n m i: t ? n b o: <p> t r i: o: /

2. Lemmatized transcriptions

Each word is followed by the citation form of its lemma as it appears in Arnold's (2019) dictionary.

A lemmatized sentence:

anah<anah> hōxa<hōxa> b̄<b->-blōta<blōta>
nmī?cabrill<?cabr yī?čbar> ?inbō<?enapṭa>
maštra<maštra> ra?isō<ra?isa> l̄<l>-blōta<blōta>.

These files are the result of a lemmatization process that was completed with the help of our language consultant.

4. SQLite database

