

SLTU-2014, St. Petersburg, Russia, 14-16 May 2014

# The NCHLT Speech Corpus of the South African languages

Etienne Barnard<sup>1</sup>, Marelie H. Davel<sup>1</sup>, Charl van Heerden<sup>1</sup>, Febe de Wet<sup>2</sup> and Jaco Badenhorst<sup>2</sup>

<sup>1</sup>Multilingual Speech Technologies, North-West University, Vanderbijlpark, South Africa

<sup>2</sup>Human Language Technologies Research Group, Meraka Institute, CSIR, Pretoria, South Africa

{etienne.barnard, marelie.davel, cvheerden}@gmail.com, {fdwet, [jbadenhorst](mailto:jbadenhorst@csir.co.za)}@csir.co.za

## Abstract

The NCHLT speech corpus contains wide-band speech from approximately 200 speakers per language, in each of the eleven official languages of South Africa. We describe the design and development processes that were undertaken in order to develop the corpus, and report on associated materials such as orthographic transcriptions and pronunciation dictionaries that were released as part of the corpus. In order to benchmark speech recognition performance on the corpus, we have also developed both phone-recognition and word-recognition systems for all eleven languages; we find that high accuracies can be achieved for these speaker-independent but vocabulary-dependent recognition tasks in all languages.