

# **The South African Human Language Technology Audit**

Language Resources and Evaluation, Special Issue: African Language Technology

Aditi Sharma Grover<sup>1,2</sup>, Gerhard B. van Huyssteen<sup>3</sup>, Marthinus W. Pretorius<sup>2</sup>

*Human Language Technology Research Group, CSIR<sup>1</sup>,*

*Graduate School of Technology Management, University of Pretoria<sup>2</sup>,*

*Centre for Text Technology (CTexT), North-West University<sup>3</sup>*

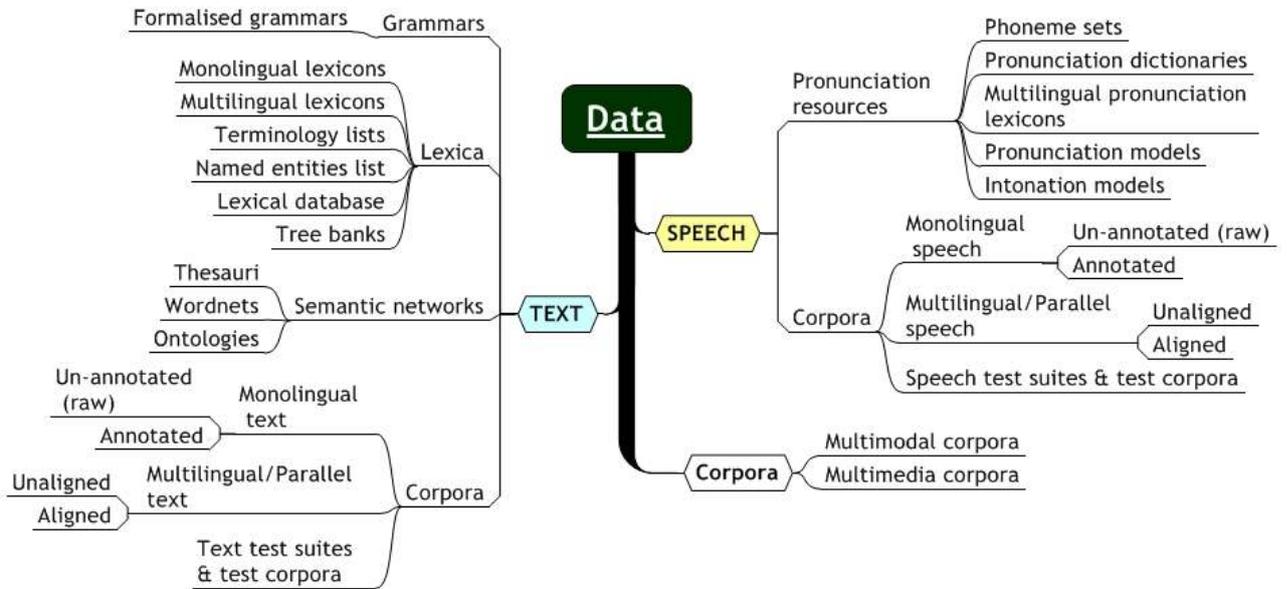
Corresponding author email: [asharma1@csir.co.za](mailto:asharma1@csir.co.za)

## **ONLINE RESOURCE**

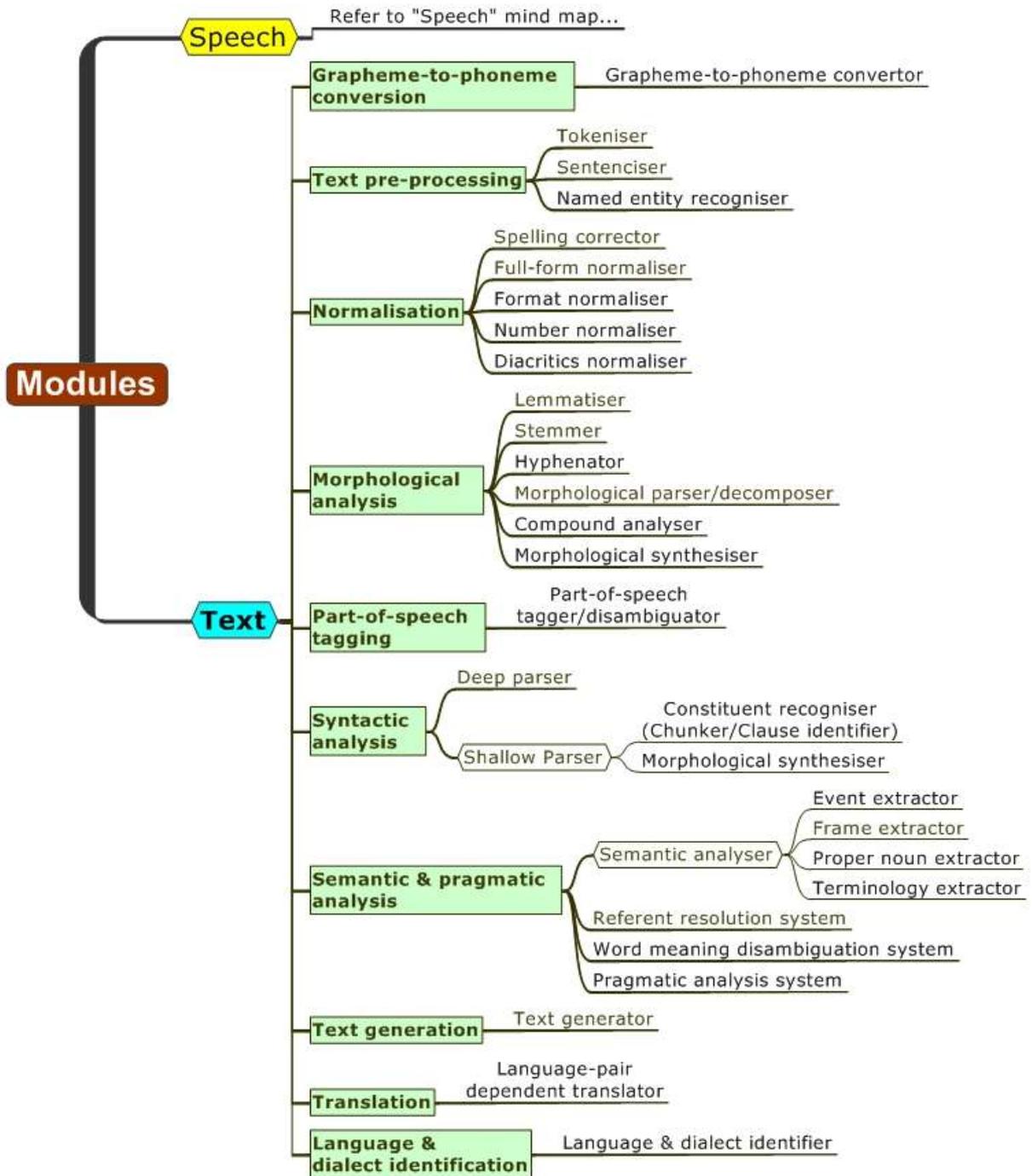
This file is the online resource cited in the above-mentioned article. It contains a series of figures and tables that further provide further information for certain sections of the article (as detailed below).

## SECTION 3.1

Figures A, B, C and D illustrate the detailed ontology of HLT components i.e. data, modules and applications developed during the SAHLTA.



**Fig.A** HLT components: Data ontology



**Fig.B** HLT components: Text modules ontology

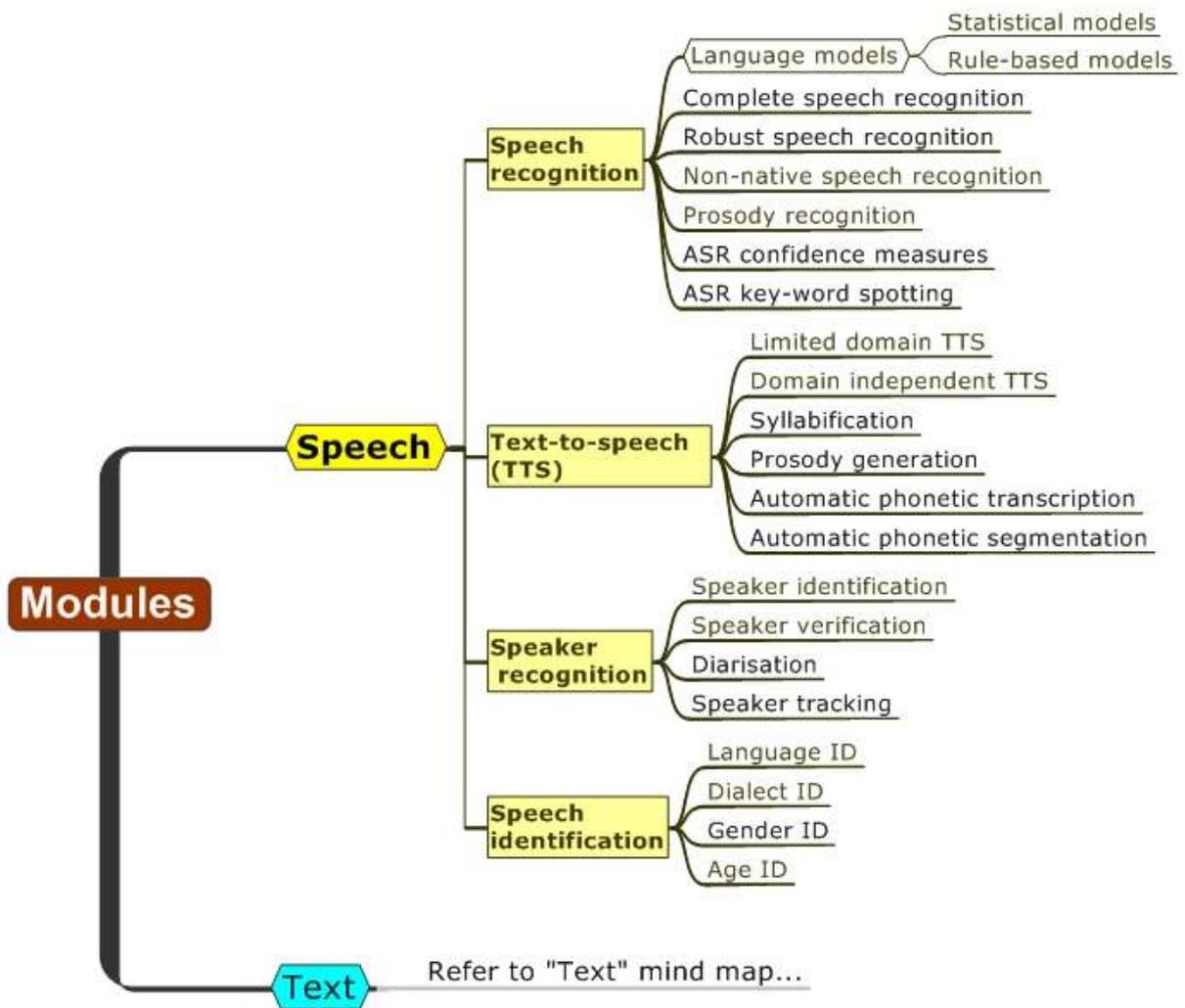
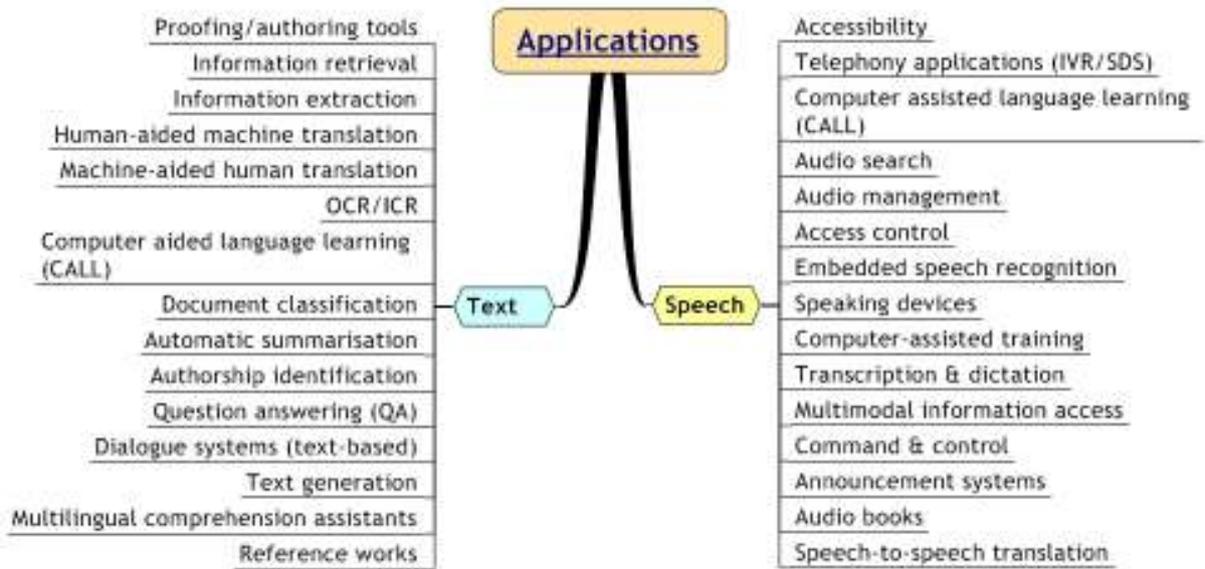


Fig.C HLT components: Speech modules ontology



**Fig.D** HLT components: Applications ontology

## SECTION 4.1

Tables A and B respectively show the prioritised HLT modules and data for South Africa determined in the SAHLTA workshop.

**Table A.** Prioritisation of HLT modules for South Africa

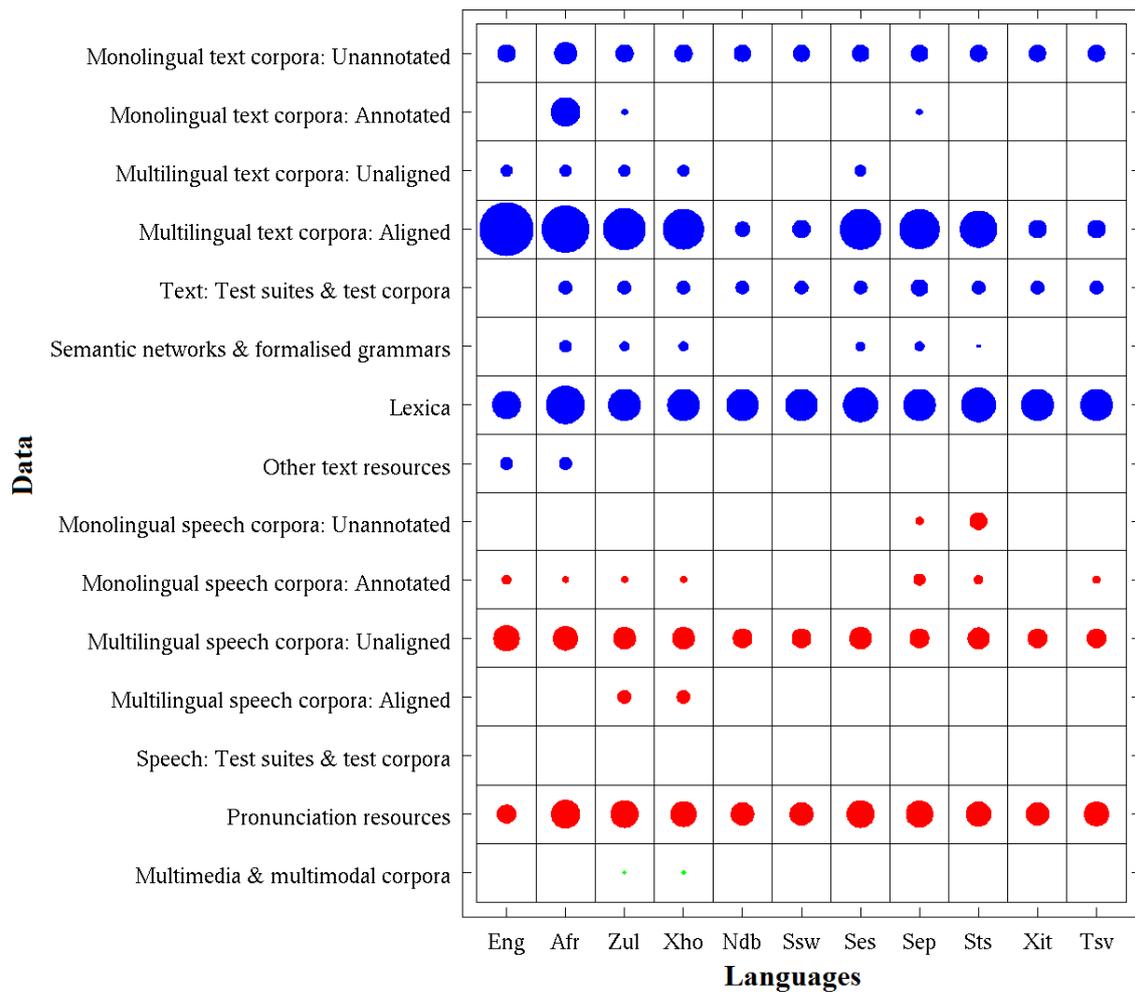
Priority	Text Modules	Speech Modules
1	Grapheme-to-phoneme (G-2-P) convertor	Complete speech recognition (domain-independent and limited) <ul style="list-style-type: none"> <li>Acoustic models</li> <li>Application based language models. Both statistical (implies large corpus) &amp; rule-based.</li> </ul> (Orthographically transcribed representative corpora & related pronunciation resources)
	Text pre-processing <ul style="list-style-type: none"> <li>Tokeniser</li> <li>Sentenciser</li> <li>Named entity recogniser</li> </ul>	Non-native speech recognition
	Normalisation <ul style="list-style-type: none"> <li>Spelling corrector</li> <li>Full-form normaliser</li> <li>Format normaliser</li> <li>Number normaliser</li> <li>Diacritics normaliser</li> </ul>	Complete TTS (limited domain and domain-independent) <ul style="list-style-type: none"> <li>G-2-P convertor</li> <li>Pre-processing NLP, POS tagger, Chunker</li> <li>Normalisation</li> <li>Prosody generation</li> <li>Automatic phonetic segmentation</li> <li>Syllabification</li> </ul> (Related pronunciation resources)
	Morphological analysis <ul style="list-style-type: none"> <li>Lemmatiser</li> <li>Morphological parser/decomposer</li> <li>Compound analyser</li> <li>Morphological synthesiser</li> </ul>	Confidence measures (ASR)
	Part-of-speech (POS) tagger	Speaker identification
	Syntactic analysis <ul style="list-style-type: none"> <li>Shallow parser- Constituent recogniser (chunker/clause identifier)</li> </ul>	Diarization
	Semantic and pragmatic analysis <ul style="list-style-type: none"> <li>Word meaning disambiguation system</li> </ul>	Language identification
	Language and dialect identifier	
2	Morphological analysis <ul style="list-style-type: none"> <li>Stemmer</li> <li>Hyphenator</li> </ul>	Speaker verification
	Syntactic analysis <ul style="list-style-type: none"> <li>Shallow parser (Relation finder)</li> </ul>	Dialect identification
	Semantic and pragmatic analysis <ul style="list-style-type: none"> <li>Referent resolution system</li> </ul>	Automatic phonetic transcription
	Language-pair dependent translator	
3	Syntactic analysis <ul style="list-style-type: none"> <li>Deep parser</li> </ul>	Prosody recognition
	Semantic & pragmatic analysis <ul style="list-style-type: none"> <li>Semantic analyser (Event, Frame, Proper noun &amp; Terminology extractors)</li> <li>Pragmatic analysis system</li> </ul>	Robust speech recognition
	Text generator	Speaker tracking
		Keyword-spotting

**Table B.** Prioritisation of HLT data for South Africa

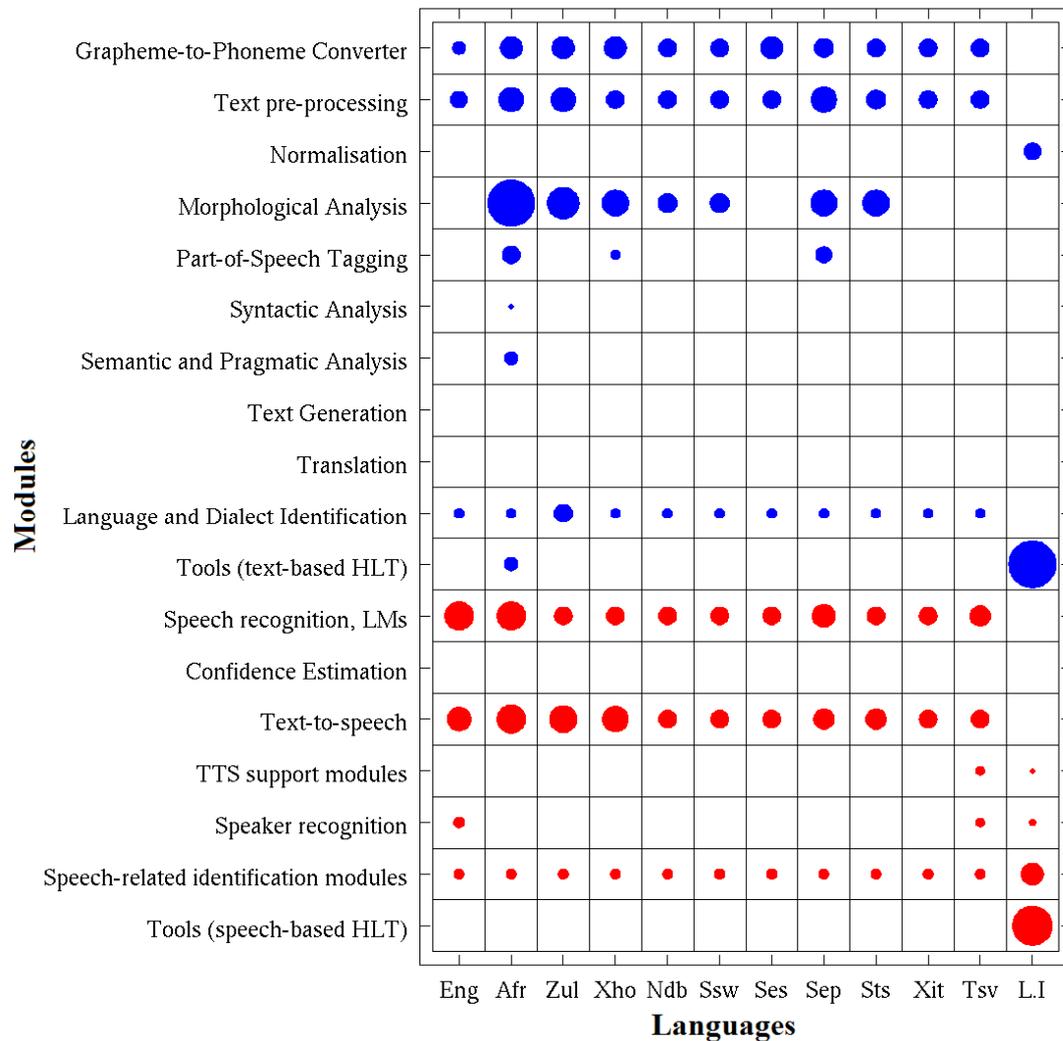
Priority	Text Data	Speech Data
<b>1</b>	Monolingual text corpora <ul style="list-style-type: none"> <li>Un-annotated</li> <li>Annotated (words+POS, sentence+constituent structure)</li> </ul>	Monolingual speech corpora <ul style="list-style-type: none"> <li>Annotated (with minimum orthographic transcription)</li> </ul>
	Multilingual text corpora <ul style="list-style-type: none"> <li>Unaligned</li> <li>Aligned</li> </ul>	Corpora of related domains (e.g. for audio search corpus of related searches)
	Text test suites & corpora	Speech test suites & corpora
	Lexica: <ul style="list-style-type: none"> <li>Monolingual lexicons</li> <li>Multilingual lexicons</li> <li>Terminology lists</li> <li>Named entity lists</li> </ul>	Pronunciation resources <ul style="list-style-type: none"> <li>Phoneme sets</li> <li>Pronunciation lexicons</li> <li>Pronunciation models</li> <li>Intonation models</li> </ul>
	Domain/application related text corpora	
<b>2</b>	Formal grammars	Pronunciation resources <ul style="list-style-type: none"> <li>Multilingual pronunciation lexicons</li> </ul>
	Semantic networks <ul style="list-style-type: none"> <li>Wordnets</li> </ul>	Annotated speech corpora with prosodic and phonetic annotations Multi-lingual speech corpora <ul style="list-style-type: none"> <li>Unaligned</li> <li>Aligned (parallel)</li> </ul>
<b>3</b>	Mono-lingual text corpora annotated at sentence+semantic level	Unannotated speech corpora
	Semantic networks <ul style="list-style-type: none"> <li>Thesauri</li> <li>Ontologies</li> </ul>	Annotated speech corpora with various types of annotations (age, emotions, speakers, etc.)
	Lexica <ul style="list-style-type: none"> <li>Lexical databases</li> </ul>	
		Multimodal corpora
		Multimedia corpora

### SECTION 4.2.3

Figures E and F show the HLT Component Indexes for data and modules respectively.



**Fig.E** HLT Component Index for data



**Fig.F** HLT Component Index for modules

Figures G, H, I, J illustrate a gap analysis, which identifies the gaps between the current status and the prioritised South African HLT components (from the SAHLTA workshop). This information could be highly informative for future road-mapping exercises, as well as to immediately identify areas or languages that should receive particular attention.

The gap analysis for text LRs and applications (Figures G and H) revealed that although all of the data resources identified as priority 1 are somewhat existent, many of them have restricted use due to insufficient maturity or accessibility, and are not always available for all languages. The majority of core HLT modules (priority 1) are somewhat existent, however they are only available in one or a few languages and in most cases are of an uncertain quality. Approximately one-third of the HLT modules that were identified as priority 1 (require definite attention) are currently non-existent. Note, in Figures G and H, components are classified as either ‘somewhat existent’ if in the detailed inventory analysis they were of fairly adequate quality as , or of somewhat uncertain quality as , and if items are ‘non-existent’ as .

Priority	TEXT		Languages											
1	Status	DATA	Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I
	Somewhat existent	Mono-lingual text corpora	Un-annotated	●	●	●	●	●	●	●	●	●	●	●
Annotated (words+POS, sentence+constituent structure)			○	●	◐	○	○	○	○	◐	○	○	○	-
Multi-lingual text corpora		Unaligned	◐	◐	◐	◐	○	○	◐	○	○	○	○	-
		Aligned	●	●	●	●	◐	◐	●	●	●	●	◐	-
		Test suites & corpora	○	◐	◐	◐	◐	◐	◐	◐	◐	◐	◐	-
Lexica		Mono-lingual lexicons	●	●	●	●	●	●	●	●	●	●	●	-
		Multi-lingual lexicons	◐	◐	○	○	○	○	◐	○	◐	◐	◐	-
		Terminology lists	●	●	●	●	●	●	●	●	●	●	●	-
		Named entity lists	○	◐	○	○	○	○	○	○	○	○	○	-
		Domain/application related text corpora	○											
Status	MODULES	Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I	
Somewhat existent	Grapheme to phoneme conversion	Grapheme-to-phoneme convertors	○	●	●	●	●	●	●	●	●	●	●	-
	Text pre-processing	Tokeniser	○	●	●	●	●	●	●	●	●	●	●	-
		Sentenciser	◐	◐	◐	○	○	○	○	◐	○	○	○	-
	Normalisation	Format normaliser	-	-	-	-	-	-	-	-	-	-	-	◐
		Diacritics normaliser	-	-	-	-	-	-	-	◐	-	-	◐	-
	Morphological analysis	Lemmatiser	○	●	○	○	○	○	○	○	○	○	○	-
		Morphological parser/decomposer	○	○	●	●	●	●	○	●	●	○	○	-
		Compound analyser	○	●	○	○	○	○	○	○	○	○	○	-
		Morphological synthesiser	○	○	●	●	●	●	○	●	●	○	○	-
	POS tagging	Part-of-speech (POS) tagger	○	●	○	◐	○	○	○	◐	○	○	○	-
Syntactic analysis	Shallow parser: Constituent recogniser (Chunker/Clause)	○	◐	○	○	○	○	○	○	○	○	○	-	
Semantic & pragmatic analysis	Word meaning disambiguation system	○	◐	○	○	○	○	○	○	○	○	○	-	
Language & dialect ID	Language & dialect identifier	◐	◐	◐	◐	◐	◐	◐	◐	◐	◐	◐	-	
Non-existent	Text pre-processing	Named entity recogniser	○											
	Normalisation	Spelling corrector	○											
		Diacritics normaliser	○											
		Full-form Normaliser	○											
		Number Normaliser	○											
Status	APPLICATIONS	Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I	
Somewhat existent	Proofing/authoring tools	●	●	●	●	●	●	●	●	●	●	●	●	-
	Human-aided machine translation	◐	◐	○	◐	○	○	◐	○	○	○	○	◐	-
	Machine-aided human translation	◐	◐	◐	○	○	○	○	◐	○	○	○	●	-
Non-existent	Information retrieval	○												
	Information extraction	○												

Fig.G Text-based HLTgap analysis (priority 1)

Priority	TEXT		Languages												
2	Status	DATA		Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I
	Somewhat existent	Semantic networks	Wordnets	○	◐	◐	◐	○	○	○	◐	◐	○	○	·
	Non- existent	Grammars	Formal grammars	○											
	Status	MODULES		Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I
	Somewhat Existent	Morphological analysis	Hyphenator	○	●	●	●	○	○	○	●	●	○	○	·
	Non- existent	Morphological analysis	Stemmer	○											
		Syntactic analysis	Shallow parser: Relation finder												
		Semantic & pragmatic analysis	Referent resolution system												
		Translation	Language-pair dependent translator												
	Status	APPLICATIONS		Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I
Somewhat Existent	Multilingual comprehension assistants		●	●	○	○	○	○	●	○	●	○	○	·	
	Computer assisted language learning (CALL)		●	●	●	●	●	●	●	●	●	●	●	·	
Non- existent	Optical/Intelligent character recognition (OCR/ICR) Authorship Identification		○												
3	Status	DATA		Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I
	Somewhat Existent	Lexica	Lexical databases	◐	◐	○	○	○	○	○	○	○	○	○	·
	Non- existent	Corpora	Mono-lingual text corpora annotated at sentence+semantic level	○											
		Semantic networks	Thesauri Ontologies												
	Status	MODULES		Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I
	Non- existent	Syntactic analysis	Deep parser	○											
		Semantic & pragmatic analysis	Semantic analyser												
			Event extractor												
			Frame extractor												
			Proper noun extractor												
Terminology extractor															
Pragmatic analysis system															
Text generation	Text generator														
Status	APPLICATIONS		Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I	
Somewhat Existent	Reference works		●	●	◐	◐	○	○	○	○	○	○	○	·	
Non- existent	Text generation Document classification Automatic summarisation Question Answering (QA) Dialogue systems (text-based)		○												

Fig. H Text-based HLT gap analysis (priority 2 and 3)

A similar gap analysis was performed for speech-based HLT components across all eleven languages, and is illustrated in Figures I and J. In slight contrast to the text domain, the majority of priority 1 data resources and speech modules (such as pronunciation dictionaries, speech recognition and text-to-speech (TTS) technologies) exist, though mostly in a very basic state. This is as a result of a large, government sponsored project in this domain, which recently made these basic speech domain modules and data available. However, a significant number of core data and modules – such as the more advanced speech recognition and TTS related modules – are still non-existent across the eleven languages. Telephony-based

services (IVR/SDS) is the only priority 1 speech application that has some significant activity (though not equally so in all languages); other applications such as accessibility and audio search are in their very early stages, whilst priority 1 applications like audio management and computer-assisted language learning (CALL) are non-existent.

Priority	SPEECH		Languages												
1	Status	DATA	Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I	
	Somewhat Existing	Pronunciation resources	Monolingual speech corpora (minimum orthographic transcription)	●	●	●	●	●	●	●	●	●	●	●	●
Phoneme sets			●	●	●	●	●	●	●	●	●	●	●	●	-
Pronunciation dictionaries			◐	●	●	●	●	●	●	●	●	●	●	●	-
Pronunciation models			○	●	●	●	●	●	●	●	●	●	●	●	-
Intonation models			-	-	◐	○	○	○	◐	◐	◐	○	○	○	-
Non-existent	Test suites & corpora	○													
Somewhat Existing	Speech recognition	Rule-based language models	●	●	○	○	○	○	○	◐	○	○	○	-	
		Complete ASR	●	●	●	●	●	●	●	●	●	●	●	●	-
	Text-to-speech	Grapheme-to-phoneme convertor	○	●	●	●	●	●	●	●	●	●	●	●	-
		Complete TTS - limited domain	◐	◐	○	◐	○	○	○	○	○	○	○	○	-
		Complete TTS- domain independent	●	●	●	●	●	●	●	●	●	●	●	●	-
		Normalisation	◐	◐	◐	◐	◐	◐	◐	◐	◐	◐	◐	◐	-
		Automatic phonetic segmentation	○	○	○	○	○	○	○	○	○	○	○	○	◐
	Speech related ID	POS tagger (see text LRs)	○												
		Chunker (see text LRs)	○												
	Non-existent	Speech recognition	Statistical language models	○											
Non-native speech recognition			○												
Confidence measures (ASR)			○												
Text-to-speech		Pre-processing NLP	○												
		Syllabification	○												
Speaker recognition	Prosody generation	○													
Somewhat Existing	APPLICATIONS	Speaker identification	○												
		Diarization (lang-independent)	○												
		Accessibility	◐	◐	◐	○	○	○	○	○	○	○	○	○	-
		Telephony applications (IVR/SDS)	●	●	●	●	○	○	●	●	●	○	○	○	●
		Audio Search	◐	○	○	○	○	○	○	○	○	○	○	○	-
Non-existent	Audio Management	○													
	Computer assisted language learning (CALL)	○													

Fig.I Speech-based HLT gap analysis (priority 1)

Priority	SPEECH		Languages													
2	Status	DATA		Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I	
	Somewhat Existent	Annotated speech corpora	Prosodic and phonetic annotations	<input type="radio"/>	-											
		Multi-lingual speech corpora	Unaligned	<input checked="" type="radio"/>	-											
		Multi-lingual speech corpora	Aligned/parallel	<input type="radio"/>	-											
	Non-existent	Pronunciation resources	Multilingual pronunciation lexicons	<input type="radio"/>												
	Status	MODULES		Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I	
	Somewhat Existent	Speaker recognition	Speaker verification	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>
		Non-existent	Text to speech	Automatic phonetic transcription	<input type="radio"/>											
	3	Status	APPLICATIONS		Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I
		Somewhat Existent	Access control		<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Embedded speech recognition			<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	-	
Speaking devices			<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	-	
Computer-assisted training			<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	-	
Status	DATA		Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I		
Somewhat Existent	Corpora	Multimedia corpora	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	-	
		Un-annotated speech corpora	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	-	
Non-existent	Corpora	Multimodal corpora Annotated speech corpora with other types of annotations (age, emotions, speakers, etc.)	<input type="radio"/>													
Status	MODULES		Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I		
Non-existent	Speech recognition	Prosody recognition	<input type="radio"/>													
		Robust speech recognition Keyword-spotting (ASR)	<input type="radio"/>													
Non-existent	Speaker recognition	Speaker tracking	<input type="radio"/>													
Status	APPLICATIONS		Eng	Afr	Zul	Xho	Ndb	Ssw	Ses	Sep	Sts	Xit	Tsv	L.I		
Somewhat Existent	Multimodal information access		<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	-	
	Speech-to-speech translation		<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	-	
Non-existent	Transcription and dictation		<input type="radio"/>													
	Command & Control		<input type="radio"/>													
	Announcement systems		<input type="radio"/>													
	Audio books		<input type="radio"/>													

Fig.J Speech-based HLT gap analysis (priority 2 and 3)