



The systematic collection of speech corpora for all eleven official South African languages

Marissa van Rooyen

Centre for Text Technology (CTexT™)

Research Unit: Languages in South African Context (School of Languages)
North West University, Potchefstroom Campus (PUK)
South Africa

E-mail: 13017527@student.nwu.ac.za



NORTH-WEST UNIVERSITY
YUNIBESITI YA BOKONE-BOPHIRIMA
NOORDWES-UNIVERSITEIT
POTCHEFSTROOM CAMPUS

Introduction

- South African Situation
 - 11 official languages that need to be treated equally
 - Completed resources
 - Spell checkers and CALL-products mostly
 - No workable solution for ASR
- Technology can ease lives and bridge gap
- Goal of this paper: to show a way in which data collection can be organised and managed effectively and quickly

Outline

- Scope of project
- Basic method (5 steps)
- Management of assistants
- Data management
- Suggestions
- Conclusion

Scope of this project

- Developing a speech-driven, telephone-based information system
- All languages in 4 phases
- Phases include recording, transcription and quality control for each language
- 200 speakers per language (male:female, cell:landline, 18-35:36-65)

Basic method

- 5 Basic steps to collecting data
- 2. Recruit speakers and make appointment
- 3. Send prompt sheet
- 4. Record speaker
- 5. Transcribe recording
- 6. Final quality control

1. Recruit speakers

- Where to find them:
 - Family and friends
 - Competition
 - Language Boards etc.
- Get demographics and verify first language
- Appointment

2. Send prompt sheet

- Automatically generated by PromptSheetGen
- Includes questions, dates, times, numbers and sentences
- Each one is unique – reference number

PromptSheet
Generated by PromptSheetGen

Ref :: Ses otho_1.1.1.1

Section [1]

- Thusa o eme ho fihlela moqophiso a o kapa ho araba dipoto tse ka tliso tse lishomo le matso o mana (14).
- O se ka wa bale dipoto haholo, di araba feela.

Section [2]

1. Phae ya hao ya lapang ke efe?
2. Dilemo tsa hao di kang?
3. Na o matho o matona kapa o matahahadi?
4. Na o bua ka sefeloana kapa mabala wa fatiso?
5. Thusa o bale nomoro o lehakarang le letshahadi hodimo hukung ya leqephe lena.
6. Le nang letsoali la hao la tsewala?
7. Na o na le malikani?
8. O tseletsewa kang?
9. O dula tarapang / matsiang ote peale?

Section [3]

Thusa o bale nako e latelang ::

1. Mesele o lishomo le matso o mathano ka mona ho hana ya basupa.

Section [4]

Thusa o pelete lentse le latelang ::

1. seletsang

Section [5]

Thusa o bale letsatsi le latelang ::

1. 23 Hlakola

Section [6]

Thusa o bale nomoro e latelang ya founu: ::

- 025 057 858 5

Section [7]

Thusa o bale ho latelang: ::

1. Letaba la Matlato a Naha Pophapoko ya Afrika Boroa

Section [8]

- Thusa o bale dipoto tse latelang jwalo ka ha di hlophatsewa mona ka tliso. Ema nako e ka bang matsoewana e mmedi pakeng tsa potolo ka mgisa.

Section [9]

Dipolelo :: ::

1. Sefate se ommeng se thibile tse a e patisang.
2. Naga mgisa le ba le nakano ho kgathana.
3. Ke ka hana e buang ka letsoeli.
4. Ho kgathana ha Mopresidente.
5. Wkgwang senole e na le monamo o tlaming ka hahlehlile.
6. Waseho o na boadiba ba utlawa ka masepa.
7. Ha ke utlawa e doko seo ka tseba hana ho senyehile.
8. Kageho o fumano manyalla wa ho basetsa molamu sefatsang.
9. Hlakomelang kgafuta o e pakeng tsa mgisa e mmedi.
10. Pheaphatheha, qadika, hahlehlile, shahana le nyanyanya.
11. Letaba la letsoeli lena le tataelwa matsuwi kotsahlang.
12. Mohlomong ho nang baka tse ding tseo re di selang.
13. Bashepanya ba sena ba nku lebalang.
14. Ka ya moa ba ditsang ho hana o na o le haufi le tsaka.
15. Phepang ke efe maharang a ho basetsa le ho tshela?
16. Tshapo o tsepeka ka batho ba bangata.

Basic method

- 5 Basic steps to collecting data
- Recruit speakers and make appointment
- Send prompt sheet
- **Record speaker**
- Transcribe recording
- Final quality control

3. Record speaker

- One-A-LOG
- Mute button
- First unofficial QC
 - Listen to answers
 - Listen to quality of voice and reading ability
 - Listen for noise or interruptions that would reduce usability

4. Transcription

- Praat
- According to predefined protocol
 - Other voices in surroundings
 - Non-speech sounds
 - Filled pauses
- Cut out all noisy parts if possible
- Second unofficial QC

5. Final QC

- Listen to recording and read transcript
- Adhere to transcription and recording protocols
- Must be assistant's first language

Assistants to do the job

- Advantages of assistants:
 - Personal touch necessary in rural areas
 - No overshooting
 - QC throughout stages
- Skilled in language and use of computer
- University students

Management of assistants

- Proved difficult
 - Skill
 - Time
 - Motivation
- Improvement from Phase 1
 - More assistants (from 2 to 7)
 - More languages (trilingual)
 - More pay (hourly vs piecework)
- Only paid for quality work
- From 8 months to 3

Data management

- Database with criteria as fields

The screenshot shows a software window titled "Refgen" with a data entry form. The form is organized into several sections:

- Reference:** A dropdown menu.
- Name:** Edwin
- Surname:** Ramavhoya
- Age:** 18 - 35
- Gender:** Male
- Call_Type:** Cellphone
- Landline:** (empty)
- Cell:** 0783004985
- Fax:** 0866086675
- Email:** (empty)
- Recorded:** Yes (dropdown)
- Transcribed:** No (dropdown)
- Recorded By:** Malebo (dropdown)
- Appointment (yyyy/mm/dd):** 2008/03/26
- Record Date (yyyy/mm/dd):** 2008/03/27
- Prompt sheet:** W:\Phase 4 Promptsheet\T sivenda\T sivenda_1

Buttons: "Open Prompt" and "New Search".

Record: 1 of 200

Data management (2)

- Common storage location
 1. Keep recordings in batches (date it was recorded)
 2. Move to assistant's folder on common drive for transcription
 3. Move to folder for QC
 4. Bring everything back to one folder per language for delivery
- Connected to server – automatic backup
- Assistant can work from any station
- Spreadsheet for easy tracking

Suggestions for further improvement

- Relational database
 - User rights on common storage location
 - Faster
- Dedicated recruiters

Conclusion

- Overview of method and practices we used with success
- Time halved from 1st to 2nd phase – might be further reduced
- Quality remains no. 1 priority

Acknowledgements

- CText-staff
- The Meraka Institute (CSIR)
- Every assistant, speaker and recruiter

THANK YOU!