



MINISTÈRE DE L'ÉDUCATION
NATIONALE, DE L'ENSEIGNEMENT
SUPÉRIEUR ET DE LA RECHERCHE

PRESERVATION OF AFRICAN CULTURAL HERITAGE BY AUTOMATIC TRANSCRIPTION OF AFRICAN LANGUAGES

NIMAAN A.¹, NOCERA P.²

1 Institut des Sciences et des Nouvelles Technologies de Djibouti
2 Laboratoire Informatique d'Avignon

LABORATOIRE
D'INFORMATIQUE
CERI

339 Chemin des Meinajariès
BP 1228
84911 AVIGNON CEDEX 09
Tél. + 33 (0)4 90 84 35 09
Fax. + 33 (0)4 90 84 35 01
secretariat-lia@univ-avignon.fr
<http://www.lia.univ-avignon.fr>





Introduction

In Africa, the cultural, scientific and historical patrimonies are transmitted orally through generations.

This ancestral knowledge is disappearing

Most of the concerned countries, regional and international organizations (UNESCO) are elaborating policies to save this human richness.

LABORATOIRE
D'INFORMATIQUE
CERI

339 Chemin des Meinajariès
BP 1228
84911 AVIGNON CEDEX 09
Tél. + 33 (0)4 90 84 35 09
Fax. + 33 (0)4 90 84 35 01
secretariat-lia@univ-avignon.fr
<http://www.lia.univ-avignon.fr>





Introduction

MINISTÈRE DE L'ÉDUCATION
NATIONALE, DE L'ENSEIGNEMENT
SUPÉRIEUR ET DE LA RECHERCHE

Two main issues:

- Saving : recording and digitalizing the oral patrimony
- Using : access the databases

LABORATOIRE
D'INFORMATIQUE
CERI

339 Chemin des Meinajariès
BP 1228
84911 AVIGNON CEDEX 09
Tél. + 33 (0)4 90 84 35 09
Fax. + 33 (0)4 90 84 35 01
secretariat-lia@univ-avignon.fr
<http://www.lia.univ-avignon.fr>



Our objective

- Build a spoken document retrieval system for the Djibouti Republic.
 - Find the best speech representation for document retrieval

LABORATOIRE
D'INFORMATIQUE
CERI

339 Chemin des Meinajariès
BP 1228
84911 AVIGNON CEDEX 09
Tél. + 33 (0)4 90 84 35 09
Fax. + 33 (0)4 90 84 35 01
secretariat-lia@univ-avignon.fr
<http://www.lia.univ-avignon.fr>



Outline

- The Somali language
- Corpus
- Speech Recognition System
- Experiments and improvements
- Conclusion

LABORATOIRE
D'INFORMATIQUE
CERI

339 Chemin des Meinajariès
BP 1228
84911 AVIGNON CEDEX 09
Tél. + 33 (0)4 90 84 35 09
Fax. + 33 (0)4 90 84 35 01
secretariat-lia@univ-avignon.fr
<http://www.lia.univ-avignon.fr>



Djibouti languages



- 4 languages are spoken in Djibouti
- French and Arabic are official languages
- Somali and Afar are natives and widely spoken

Somali language



- 15 millions of speakers
- Djibouti, Somalia, Ethiopia and Kenya
- Cushitic language within the Afro-asiatic family

Somali language



- 22 consonants and 20 vowels
- Tone accent language
- Written system uses Roman letters (since 1972)
- Words are composed by the concatenation of sub-words/syllables (called “Roots” in our work).

Examples :

Birlab → Bir (CVC) lab (CVC)

Agoon → Ag (VC) oon (VVC)

Ugub → Ug (VC) ub (VC)

Corpus

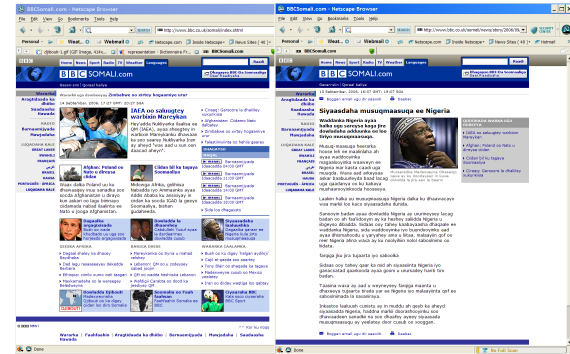
No corpus available for Somali language

- Collect text from internet (text)
- Record speech (audio)

LABORATOIRE
D'INFORMATIQUE
CERI

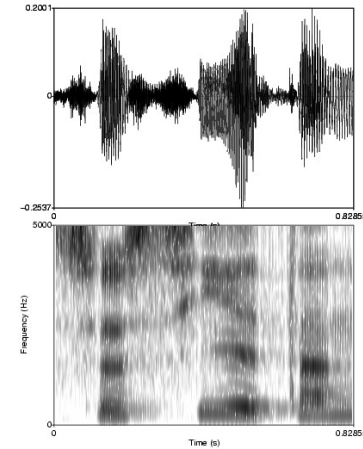
339 Chemin des Meinajariès
BP 1228
84911 AVIGNON CEDEX 09
Tél. + 33 (0)4 90 84 35 09
Fax. + 33 (0)4 90 84 35 01
secretariat-lia@univ-avignon.fr
<http://www.lia.univ-avignon.fr>

Text Corpus



- Only newspapers on internet are available
- WARGEYS corpus : from newspapers (2002-2004)
- 3 millions of words downloaded from internet (broadcast news)
- 121 k different words

Audio Corpus



■ ASAAS

- read speech
- somali news from 2002-2004
 - 10 hours
 - 10 speakers
- training **ASAAS-train** (8 hours and 57 minutes) ;
- test **ASAAS-test** (1 hour and 29 minutes) ;

■ RTD

- extract of the Djibouti oral tradition
- 1 hour of speech
- 5 speakers

Speech Recognition System

- Acoustic Models
 - Hidden Markov Models
 - Non contextual
 - 20 Vowels
 - 22 Consonants

LABORATOIRE
D'INFORMATIQUE
CERI

339 Chemin des Meinajariès
BP 1228
84911 AVIGNON CEDEX 09
Tél. + 33 (0)4 90 84 35 09
Fax. + 33 (0)4 90 84 35 01
secretariat-lia@univ-avignon.fr
<http://www.lia.univ-avignon.fr>

Speech Recognition System

■ Language Model

- Trigram Language Model
- 20k words lexicon
- OOV = 4,90% (Out Of Vocabulary)
- 723 k bigrams
- 1 750 k trigrams

LABORATOIRE
D'INFORMATIQUE
CERI

339 Chemin des Meinajariès
BP 1228
84911 AVIGNON CEDEX 09
Tél. + 33 (0)4 90 84 35 09
Fax. + 33 (0)4 90 84 35 01
secretariat-lia@univ-avignon.fr
<http://www.lia.univ-avignon.fr>

First Experiments

■ Speech Recognition System

LIA continuous speech recognition system : Speeral

Correct	Sub	Del	Ins	Error rate
76.4	20.8	2.8	4.7	28.2

Error Analyses

Ref : **GUDDO**MIYAHA gobolka oo uu **WERI**YAHAYAGU wax
Hyp : **GUDO**MIYAHA gobolka oo uu **WARI**YAHAYAGU wax

Ref : ka **WEYDII**YAY ARIMAHA AY ka wada hadleen WUXUU
Hyp : ka **WAYDII**YAY MASTAR ILAAHAY ka wada hadleen UGU

Ref : sheegay in waqti kale ay ** ***** U **BALLA**MEEN
Hyp : sheegay in waqti kale ay KU TIMID **BALA**MEEN

Ref : **DHAMMA**YSTIRKA HESHIISYO ***** hore U
Hyp : **DHAM**AYSTIRKA BISHII SIIYO hore UGU

Ref : dhexmaray oo * aanu **FAAH** **FAAHIN**
Hyp : dhexmaray oo U aanu **** **FAAHFAAHIN**

Standardization problem

- several transcriptions for the same word
- common problem for most of African languages
- Problem of standardization : spelling error
- 3 problems
 - double consonant
 - GUDDOOMIYAHA vs. GUDOOMIYAHA : director
 - compound-word or not
 - ISGAADSIIN vs. IS GAADSIIN : communication
 - same word written in two different ways
 - WEYDIIYAY vs. WAYDIIYAY : to ask

Output Standardization

- WER = 21,5%
- relative gain 24%

Standardized corpora	Correct	Sub	Del	Ins	Error rate
none	76.4	20.8	2.8	4.7	28.2
ASSAS_test	83.7	14.7	1.6	5.2	21.5

LM Standardization

- WER = 20,2%
- relative gain 28%

Standardized corpora	Correct	Sub	Del	Ins	Error rate
none	76.4	20.8	2.8	4.7	28.2
ASSAS_test	83.7	14.7	1.6	5.2	21.5
WARGEYS	85.1	13.4	1.5	5.3	20.2

Djibouti Broadcast News Recognition

- Djibouti Radio Archives
 - 1 hour ;
 - 7 803 words (2 378 different words)
 - spontaneous speech ;
 - thematic and temporal mismatch with the WARGEYS corpora ;
- OOV = 12,48%.

	Correct	Sub	Del	Ins	Error rate
Djibouti BN	46.6	46.4	7.0	8.7	62.1

Syllables decoding

- Lexicon (4 400 syllables)
- Trigram language model :
 - 189 000 bigrams of syllables,
 - 996 000 trigrams of syllables.
- OOV syllables rate (0,03%)

Syllables decoding

Syllables (Roots) decoding => (RER = Root Error Rate)

	Correct	Sub	Del	Ins	RER
Djibouti BN	57.2	32.3	10.5	4.2	47.0

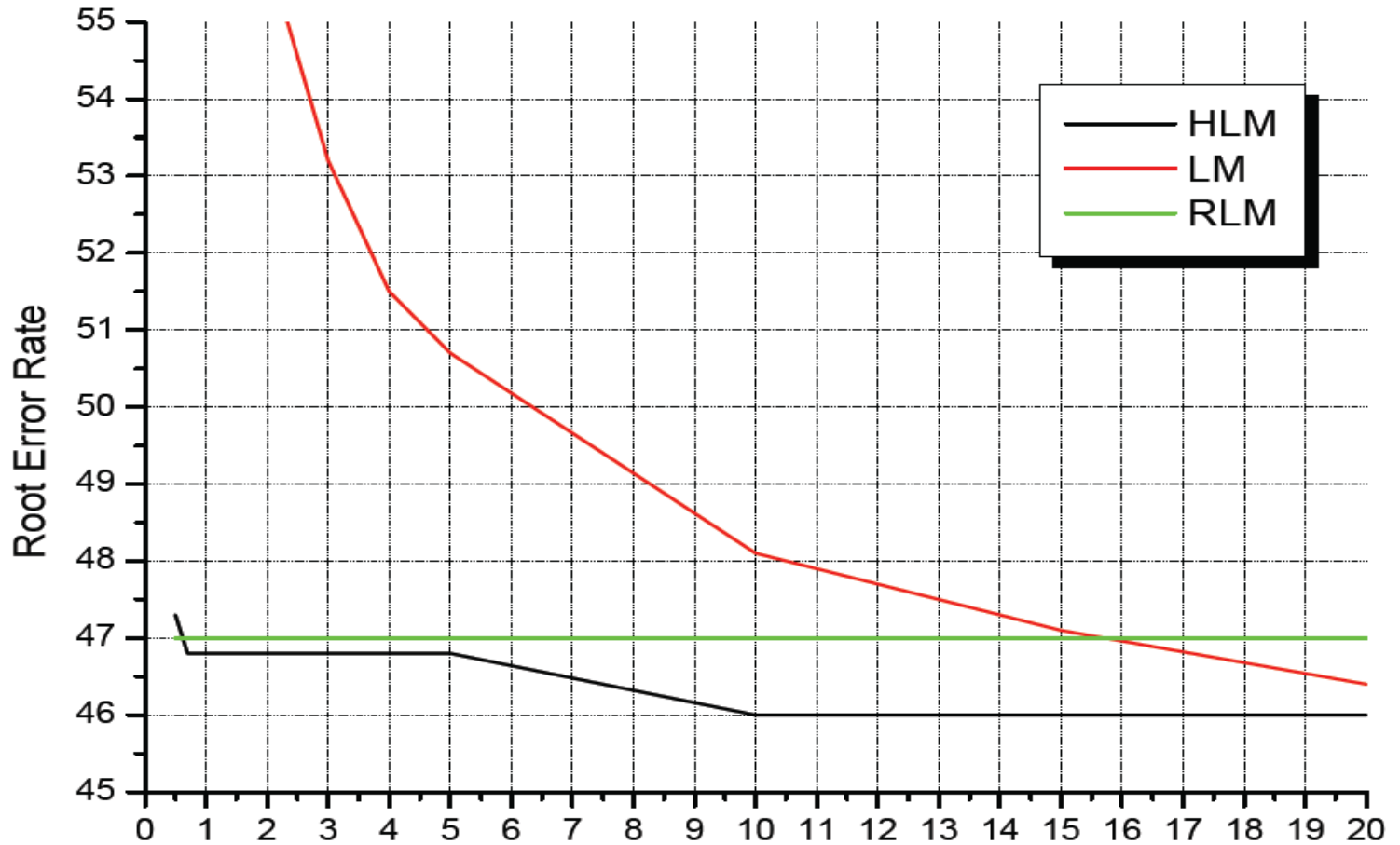
OOV Recognized by Syllabic Recognition

asnaamtaasi	as naam ta si
tafaraaruqa	taf ar aar uq a
faaqidi	aq ad i
shiinaha qudhooda	bish iib a qudh ood a
laba dakhare	lab ad a sar e

Hybrid Decoding

- Take into account n most frequent words
- Others words are transformed in syllables
- WARGEYS corpora is also transformed in words and syllables
- Hybrid language model (with lexicon n)
HLM n
- No difference between syllables and words
- Syllables Results is used

Hybrid Decoding



Conclusion

- Large vocabulary recognition system for Somali language
- problem of standardization
 - WER relative gain = 28%
- oral patrimony archives \neq training data
 - thematic and temporal mismatch
- syllables and hybrid decoding seems more appropriate than words