

Are Audio or textual training data more important  
for ASR in less-represented languages?

Thomas Pellegrini, Lori Lamel

LIMSI-CNRS, UNIVERSITÉ PARIS SUD  
Spoken Language Processing Group

May 5th 2008

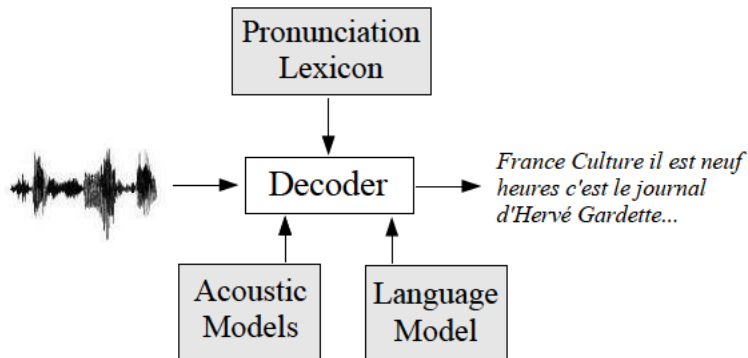
- ASR: Automatic Speech Recognition
- More than 30 years of research
- Primary applications: dictation, spoken language dialogue, and transcription for information archival and retrieval systems
- State-of-the-art systems for about 10 languages

→ Focus: ASR for less-represented languages

# Outline

- 1 Automatic Speech Recognition Overview
- 2 Less-represented Languages
- 3 Experiments
- 4 Conclusions

## Bird's eye view of Automatic Speech Recognition



# State-of-the-art overview

- NIST Evaluations (1992-)

<i>Task</i>	<i>Word Error Rate (WER in %)</i>
Read speech (1992-1995)	7
Broadcast news (1996-)	[10-20]
Conversational speech (1997-)	[20-70]

For which languages?

English, Mandarin, Arabic, Spanish

# Outline

- 1 Automatic Speech Recognition Overview
- 2 Less-represented Languages**
- 3 Experiments
- 4 Conclusions

# Quantitative definition

## Well-represented languages

Example : American English

Transcribed Audio: hundreds (thousand) hours

Texts: hundreds of million words (3000 M)

## Less-represented languages

Example: Amharic

Transcribed Audio: a few hours (37h)

Texts: a few million words (5M)

# Questions

- How much data is needed to train the models?
- What performance can be expected with a given amount of resources?
- Are Audio or textual training data more important for ASR in less-represented languages (languages with texts available)?

## Impact on the WER

- 1 Speech material (acoustic modeling)
  - 2 Texts (transcriptions)
  - 3 Texts (newspapers, newswires available on the Web)
- Case study: Amharic, Broadcast news transcription



# Outline

- 1 Automatic Speech Recognition Overview
- 2 Less-represented Languages
- 3 Experiments**
- 4 Conclusions

# Corpus sizes

- Audio

<i>Source</i>	<i>Train</i>	<i>Devtest</i>
<i>Deutsche welle</i>	24h06	1h20
<i>Radio Medhin</i>	11h08	0h37
<i># Speakers</i>	100	15
<i># Words</i>	233k	14.1k

- Texts

<i>Type</i>	<i># Words</i>	<i>Distinct Words</i>
<i>Web</i>	4.6M	200k
<i>Transcriptions</i>	247.1k	50k

# Configurations

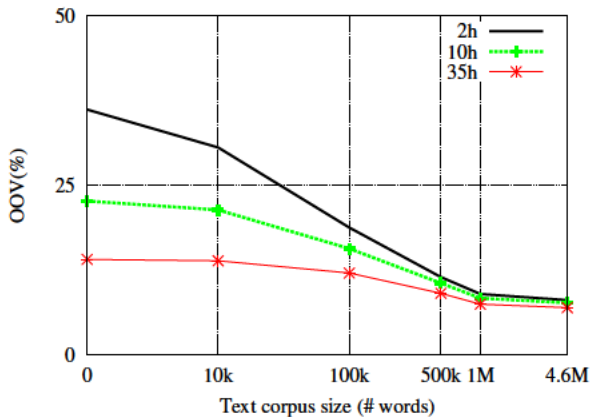
- Acoustic models: 2h, 5h, 10h, 35h
- Language models:
  - Transcriptions: 17k, 35k, 70k, 240k words +
  - Web texts: 10k, 100k, 500k, 1M, 4.6M words

## Lexicons

- Lexicon sizes

<i>Transcriptions</i>		<i>Words (texts)</i>				
<i>Hours</i>	<i>Words/Types</i>	<i>10k</i>	<i>100k</i>	<i>500k</i>	<i>1M</i>	<i>4,6M</i>
<i>2h</i>	<i>17k / 7k</i>	<b>11k</b>	36k	96k	142k	114k
<i>5h</i>	<i>35k / 12k</i>	16k	39k	98k	144k	116k
<i>10h</i>	<i>70k / 21k</i>	24k	45k	<b>103k</b>	148k	119k
<i>35h</i>	<i>240k / 50k</i>	52k	69k	120k	163k	<b>133k</b>

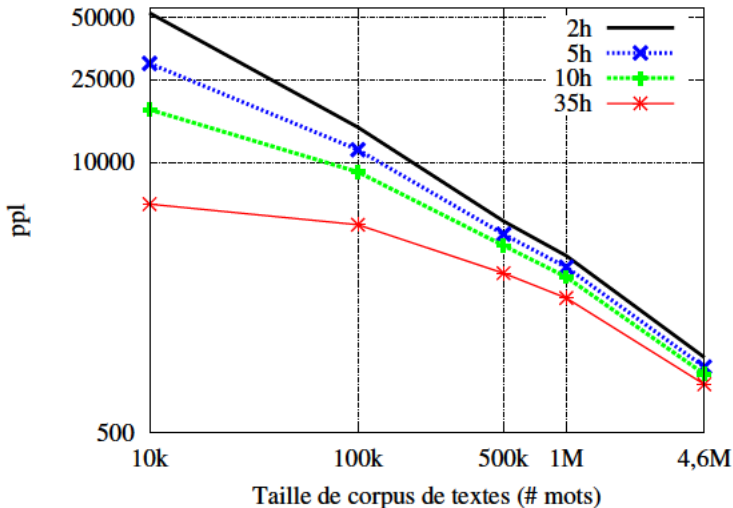
## OOV rates



- Out-Of-Vocabulary (OOV) rates

# Perplexities

- Normalized perplexities (500k word virtual lexicon)

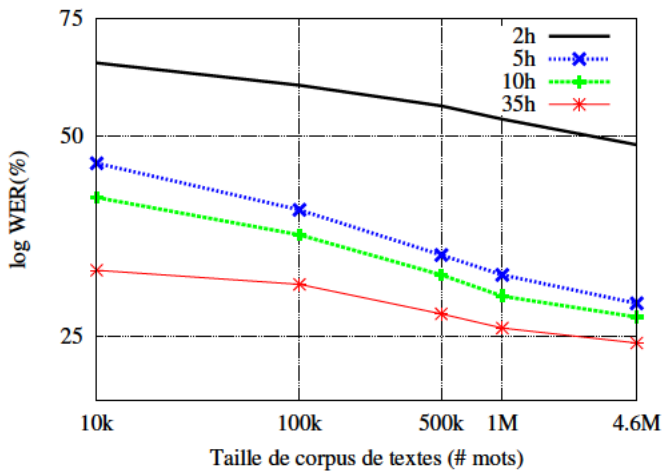


# Acoustic models

- Tied-state intra-word and cross-words HMMs, with 3 states per model and 32 Gaussians/state
- Specific acoustic models built for each audio training corpus subset

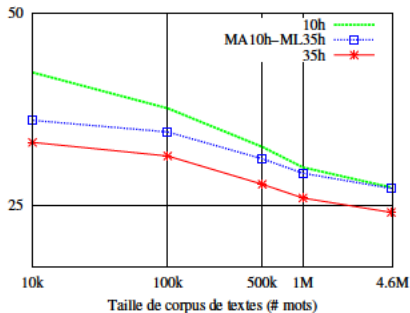
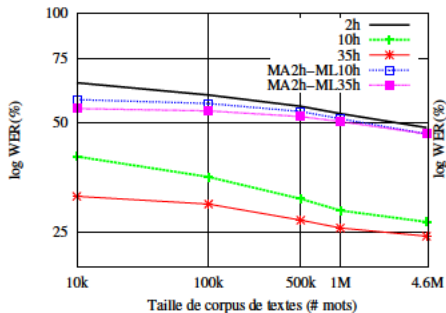
<i>Transcriptions</i>	<i>2h</i>	<i>5h</i>	<i>10h</i>	<i>35h</i>
<i># Contexts</i>	3027	4557	6323	10726
<i># Tied states</i>	1187	2286	3861	8554

## Performances





## Crossed experiments



- Influence of the transcription LM part

# Conclusions

- 2h of audio training data not sufficient
- Good operating point: 10h audio, 1M word texts
- Collecting texts may be more important than transcribing audio at this operating point, particularly if combined with unsupervised acoustic model training

# Perspectives

- Amharic: simple grapheme-to-phoneme conversion
- Similar experiments with other languages (e.g. syllable-based languages)
- Attention: these studies used broadcast news data, behaviour may be different for conversational speech or meetings
- Semi- and unsupervised training for acoustic models

Thank you for your attention!

# Interpolation coefficients

<i>Transcriptions</i>	<i># Words (texts)</i>				
	<i>10k</i>	<i>100k</i>	<i>500k</i>	<i>1M</i>	<i>4.6M</i>
<i>2h</i>	0.71	0.40	0.27	0.23	0.22
<i>5h</i>	0.84	0.53	0.38	0.33	0.28
<i>10h</i>	0.90	0.63	0.46	0.40	0.33
<i>35h</i>	0.95	0.81	0.61	0.52	0.43

- Interpolation coefficients

# Unsupervised experiments

- Lamel et al, 2002

