

## The 2013 KIT Quaero Speech-to-Text System for French

*Joshua Winebarger*, Bao Nguyen, Jonas Gehring, Sebastian Stüker, and  
Alexander Waibel

KARLSRUHE INSTITUTE OF TECHNOLOGY (KIT)



## The Quaero Program

Decoding Strategy

Acoustic data & training techniques

Feature front-end developments

- Deep neural networks

- Tonal features

Pronunciation Modeling

- Pseudographemes

Language Modeling

Experiments

- Elision normalisation

- Acoustic Data Filtering

Results

Conclusion



## Overview

- French research program with German participation
- Centered on multilingual and multimedia processing
- Evaluation-driven, corpus-based
- Cooperation dynamic
  - Partners work separately on the same subjects
  - Regular exchange of findings and results

## Areas of research:

- Development of access
- Automatic extraction
- Analysis
- Classification
- Application of information

## Focus on:

- Innovation
- Applications

# The Quaero Speech-to-Text Task

- Automatic speech recognition (ASR) is a key Quaero technology

## Useful for:

- multimedia search
- content enrichment
- structuring multimedia and multilingual documents

## Domain:

- Broadcast news and conversation
- Challenging due to:
  - Disfluencies
  - Non-speech events

- Fall 2013 is the fifth and final Quaero ASR evaluation



The Quaero Program

**Decoding Strategy**

Acoustic data & training techniques

Feature front-end developments

- Deep neural networks

- Tonal features

Pronunciation Modeling

- Pseudographemes

Language Modeling

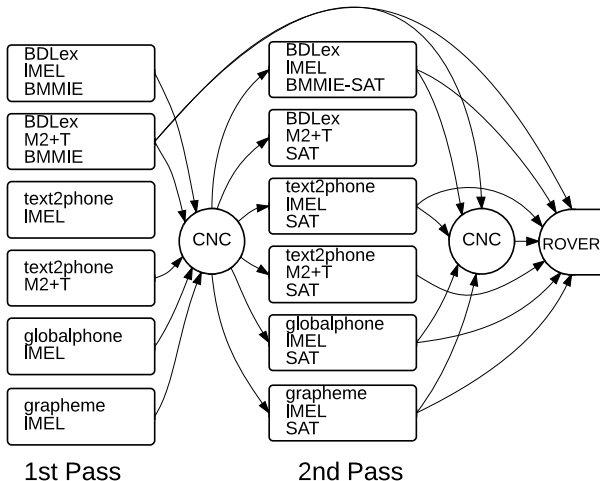
Experiments

- Elision normalisation

- Acoustic Data Filtering

Results

Conclusion



# Table of Contents

The Quaero Program

Decoding Strategy

Acoustic data & training techniques

Feature front-end developments

- Deep neural networks

- Tonal features

Pronunciation Modeling

- Pseudographemes

Language Modeling

Experiments

- Elision normalisation

- Acoustic Data Filtering

Results

Conclusion

The Quaero Program

Decoding Strategy

**Acoustic data & training techniques**

Feature front-end developments

- Deep neural networks

- Tonal features

Pronunciation Modeling

- Pseudographemes

Language Modeling

Experiments

- Elision normalisation

- Acoustic Data Filtering

Results

Conclusion



- **Strategy:** Multiple subsystems w. diverse pronunciation dictionaries (among them pseudographemes) and feature front-ends → system combination

## Innovations in front-ends

- tonal features
- deep bottleneck neural networks

## Front-end architecture

- log mel filter bank coefficients (**IMEL**)
- stacked combination of MFCC, MVDR and tonal features (**M2+T**)
- Acoustic modeling with HMMs
  - Generalized quinphones
  - Three states per phone without skip states
  - Train using classification and regression decision tree
- 8K acoustic models for phoneme systems
- 12K models for pseudographeme systems

### Training

- Incremental splitting of Gaussians (Merge and split)
  - Estimation of one global semi-tied covariance matrix after LDA
  - Two iterations of viterbi training
  - VTLN
  - Feature-space constrained MLLR SAT models for second-pass decoding
  - Discriminative training for some systems
- All subsystems:
    - 267 hours of speech
      - Quaero 2009-2011
      - Ester
  - Acoustic data filtered by decoding

Source	No filt.	Filt.
<i>Quaero</i>	194.1 hrs	187.7 hrs
<i>Ester</i>	107.8 hrs	80.3 hrs

### Training

- Incremental splitting of Gaussians (Merge and split)
  - Estimation of one global semi-tied covariance matrix after LDA
  - Two iterations of viterbi training
  - VTLN
  - Feature-space constrained MLLR SAT models for second-pass decoding
  - Discriminative training for some systems
- All subsystems: 267 hours of speech
    - Quaero 2009-2011
    - Ester
  - Acoustic data filtered by decoding

Source	No filt.	Filt.
<i>Quaero</i>	194.1 hrs	187.7 hrs
<i>Ester</i>	107.8 hrs	80.3 hrs

### Training

- Incremental splitting of Gaussians (Merge and split)
  - Estimation of one global semi-tied covariance matrix after LDA
  - Two iterations of viterbi training
  - VTLN
  - Feature-space constrained MLLR SAT models for second-pass decoding
  - Discriminative training for some systems
- All subsystems: 267 hours of speech
    - Quaero 2009-2011
    - Ester
  - Acoustic data filtered by decoding

Source	No filt.	Filt.
<i>Quaero</i>	194.1 hrs	187.7 hrs
<i>Ester</i>	107.8 hrs	80.3 hrs

The Quaero Program

Decoding Strategy

Acoustic data & training techniques

**Feature front-end developments**

- Deep neural networks

- Tonal features

Pronunciation Modeling

- Pseudographemes

Language Modeling

Experiments

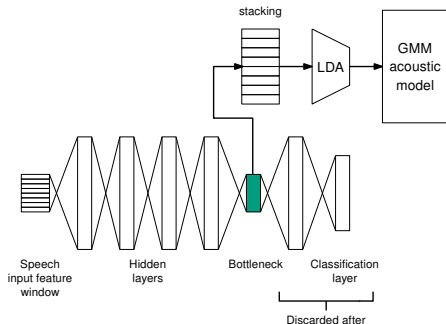
- Elision normalisation

- Acoustic Data Filtering

Results

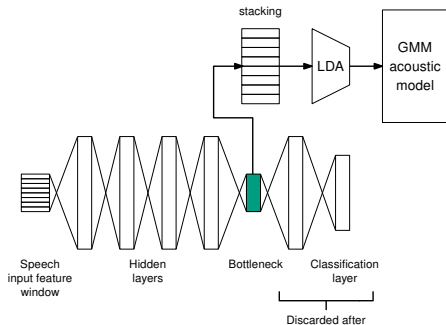
Conclusion

- Bottleneck features (BNFs) from multilayer perceptrons now common in ASR
  - Discriminative power, robustness
- We use architecture based Denoising Autoencoders (Gehring, et. al)
  - Many hidden layers
  - Significant reduction in WER ( $\approx 20\%$  relative)



## Architecture

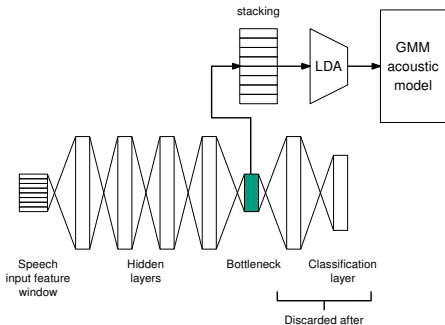
- Five fully-connected hidden layers
- 1200 units per layer
- Bottleneck layer of 42 units



- 1 Hidden layers pre-trained
  - Layerwise
  - Unsupervised
  - Stack of denoising autoencoders

- 2 Bottleneck + next hidden layer + classification layer init.'d w. random weights
- 3 These three layers are connected to the autoencoders



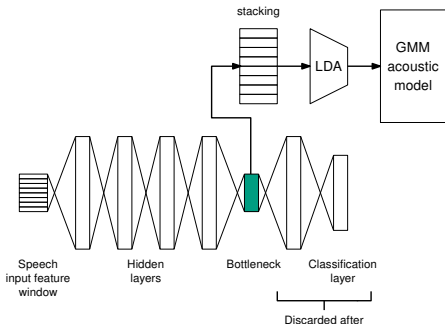


## 1 Hidden layers pre-trained

- Layerwise
- Unsupervised
- Stack of denoising autoencoders

## 2 Bottleneck + next hidden layer + classification layer init.'d w. random weights

3 These three layers are connected to the autoencoders

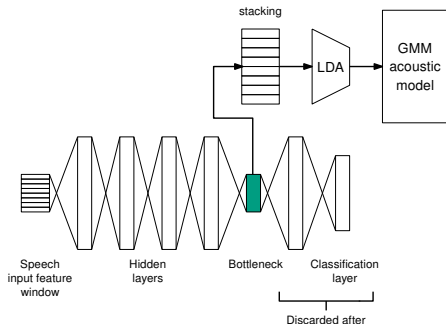


## 1 Hidden layers pre-trained

- Layerwise
- Unsupervised
- Stack of denoising autoencoders

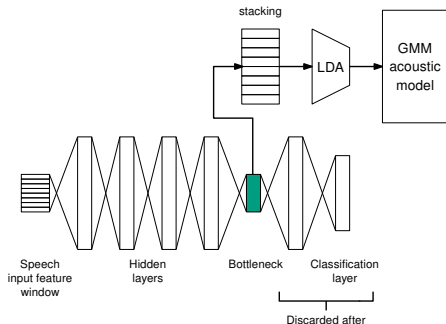
## 2 Bottleneck + next hidden layer + classification layer init.'d w. random weights

## 3 These three layers are connected to the autoencoders



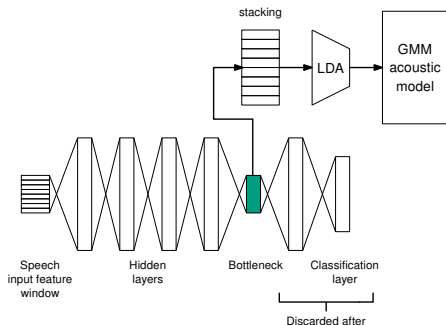
## 4 Entire network trained with supervision

- Estimation of context-dependent polyphone states
- 'newbob' learning rate schedule : about 14-18 epochs



- Compute DBNF features
- Feature stacking
- LDA

- Relative reductions in WER (vs. MFCC):
  - IMEL DBNF : 22%
  - M2+T DBNF : 24%



- Compute DBNF features
- Feature stacking
- LDA

- Relative reductions in WER (vs. MFCC):
  - **IMEL DBNF** : 22%
  - **M2+T DBNF** : 24%

- Tonal features traditionally thought not helpful for non-tonal languages
- Recent experiments at KIT show this to be untrue
- French M2+T-based DBNFs use two tonal features:
  - Pitch
  - Fundamental Frequency Variation (FFV)
- Tonal features bring relative improvements in WER:
  - 5% for Vietnamese (Tonal)
  - 3% for English (Non-tonal)
  - 3% for French (Non-tonal)

- Tonal features traditionally thought not helpful for non-tonal languages
- Recent experiments at KIT show this to be untrue
- French M2+T-based DBNFs use two tonal features:
  - Pitch
  - Fundamental Frequency Variation (FFV)
- Tonal features bring relative improvements in WER:
  - 5% for Vietnamese (Tonal)
  - 3% for English (Non-tonal)
  - 3% for French (Non-tonal)

- Tonal features traditionally thought not helpful for non-tonal languages
- Recent experiments at KIT show this to be untrue
- French M2+T-based DBNFs use two tonal features:
  - Pitch
  - Fundamental Frequency Variation (FFV)
- Tonal features bring relative improvements in WER:
  - 5% for Vietnamese (Tonal)
  - 3% for English (Non-tonal)
  - 3% for French (Non-tonal)



## Pitch

- Cepstrogram & Autocorrelation function
- Tracks across local maxima in Dynamic-Programming way to give highest cumulative path

## Fundamental Frequency Variation (FFV)

- No explicit segmentation into speech and silence necessary
- Computation of “vanishing point product” across range of  $\tau$
- Filterbank  $\rightarrow$  7 coefficients

## Pitch

- Cepstrogram & Autocorrelation function
- Tracks across local maxima in Dynamic-Programming way to give highest cumulative path

## Fundamental Frequency Variation (FFV)

- No explicit segmentation into speech and silence necessary
- Computation of “vanishing point product” across range of  $\tau$
- Filterbank  $\rightarrow$  7 coefficients

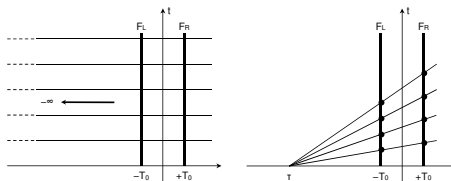


Figure: Visualization of the vanishing point product employed in FFV.

The Quaero Program

Decoding Strategy

Acoustic data & training techniques

Feature front-end developments

- Deep neural networks

- Tonal features

**Pronunciation Modeling**

- Pseudographemes

Language Modeling

Experiments

- Elision normalisation

- Acoustic Data Filtering

Results

Conclusion

# Pronunciation modeling

## Overview

- Subsystems based on four pronunciation dictionaries:
  - BDLex lexicon
  - Globalphone dictionary
  - Rule-based pronunciation generation (text2phone)
  - Pseudographeme approach
- Missing pronunciations : G2P

Dictionary	WER (CIns.)
BDLex	25.4
text2phone	25.6
globalphone	26.6
pseudographeme	27.0

- AMs for BDLex & Globalphone : bootstrap from German models
- AMs for text2phone : bootstrap from BDLex models
- AMs for pseudographemes. : flatstart

- Traditional grapheme-based pronunciation modeling:
  - symbols of written word used as sub-units of pronunciation
- Using graphemes instead of phonemes is:
  - Convenient
  - Effective in reducing WER in combination w. phoneme systems
  
- French orthography is relatively regular
  - However no one-to-one mapping between graphemes and phones for all graphemes
  - Clusters of graphemes can produce same sounds (esp. vowels)
  - Some graphemes can be silent
- This problem can be handled by creating simple rules:
  - Treat certain groups of graphemes as the same grapheme

- Traditional grapheme-based pronunciation modeling:
  - symbols of written word used as sub-units of pronunciation
- Using graphemes instead of phonemes is:
  - Convenient
  - Effective in reducing WER in combination w. phoneme systems
  
- French orthography is relatively regular
  - However no one-to-one mapping between graphemes and phones for all graphemes
  - Clusters of graphemes can produce same sounds (esp. vowels)
  - Some graphemes can be silent
- This problem can be handled by creating simple rules:
  - Treat certain groups of graphemes as the same grapheme

{ é }, { a i }, { é e }, { u é }	→	é
{ e n t }	→	ent
{ e }, { è }, { ê }	→	e
{ a u }, { e a u }	→	au
{ o n }	→	on
{ o i }	→	oin
{ g n }	→	gn

Table: *Some selected rules for merging graphemes*

- French orthography is relatively regular
  - However no one-to-one mapping between graphemes and phones for all graphemes
  - Clusters of graphemes can produce same sounds (esp. vowels)
- This problem can be handled by creating simple rules:
  - Treat certain groups of graphemes as the same grapheme

{ é }, { a i }, { é e }, { u é }	→	é
{ e n t }	→	ent
{ e }, { è }, { ê }	→	e
{ a u }, { e a u }	→	au
{ o n }	→	on
{ o i }	→	oin
{ g n }	→	gn

Table: *Some selected rules for merging graphemes*

délaissées ( <i>adj.</i> abandoned)	→	d é l é s é s
faisceaux ( <i>n.</i> bundles)	→	f é s c {au} x
fondièrement ( <i>adv.</i> fundamentally)	→	f {on} c i e r e m {ent}
joignent ( <i>v.</i> join)	→	j {oi} {gn} {ent}
pointée ( <i>adj.</i> pointed)	→	p {oin} t é

Table: *Selected entries from the grapheme dictionary with accompanying English translations*



The Quaero Program

Decoding Strategy

Acoustic data & training techniques

Feature front-end developments

- Deep neural networks

- Tonal features

Pronunciation Modeling

- Pseudographemes

**Language Modeling**

Experiments

- Elision normalisation

- Acoustic Data Filtering

Results

Conclusion

- 4-gram case-sensitive LM
- Modified Kneser-Ney smoothing
- Individual LMs built for each text source
- LMs interpolated using weights estimated on tuning set

## Sources

Quaero transcripts, other Quaero sources  
Gigaword  
Europarl transcripts  
Other newspapers  
Assorted small oral corpora  
CFPP 2000

The Quaero Program

Decoding Strategy

Acoustic data & training techniques

Feature front-end developments

- Deep neural networks

- Tonal features

Pronunciation Modeling

- Pseudographemes

Language Modeling

**Experiments**

- Elision normalisation

- Acoustic Data Filtering

Results

Conclusion

- Elision: Final unstressed vowel immediately before another vowel-beginning word *sometimes* dropped and both words joined

- Elision: Final unstressed vowel immediately before another vowel-beginning word *sometimes* dropped and both words joined
- Examples of elision:

“la	apparence”	→	l’apparence	(“the appearance”)
DEF-ART	NOUN			
“de	être”	→	d’être	(“to be”)
PREP	VERB			
“je	aimerais”	→	j’aimerais	(“I would like”)
SUBJ-PRONOUN	VERB			

- Elision: Final unstressed vowel immediately before another vowel-beginning word *sometimes* dropped and both words joined

- Examples of elision:

“la            apparence”    →    l’apparence    (“the appearance”)  
DEF-ART            NOUN

“de            être”            →            d’être            (“to be”)  
PREP            VERB

“je            aimerais”    →            j’aimerais       (“I would like”)  
SUBJ-PRONOUN    VERB

- However note:

“ça            arrive”        (“that happens”)  
DEM-PRONOUN    VERB

“déjà        estimé”       (“already estimated”)  
ADV            VERB

## Elision poses a challenge depending on how we consider elided words:

- “*join*”: Elisions are one token. → Large OOV, limited power of vocabulary
- “*sep*”: Elisions are two tokens. → Reduces OOV but decreases language and polyphone context
- “*sep. topN*”: Treat top  $N$  most common elisions as one token

## Quaero transcripts

	# Tokens	Total	Elided
Total words	23M	100%	-
Elided words	165K	7.2%	100%
• top50 words	100K	4.4%	61%
• non-top50 words	65K	2.8%	39%

## Elision poses a challenge depending on how we consider elided words:

- “*join*”: Elisions are one token. → Large OOV, limited power of vocabulary
- “*sep*”: Elisions are two tokens. → Reduces OOV but decreases language and polyphone context
- “*sep. topN*”: Treat top  $N$  most common elisions as one token

## Quaero transcripts

	# Tokens	Total	Elided
Total words	23M	100%	-
Elided words	165K	7.2%	100%
● top50 words	100K	4.4%	61%
● non-top50 words	65K	2.8%	39%



- Treat AM & LM training data to reflect these approaches
- Compute resulting OOV and WER

Elision treatment	OOV	WER(CI)
<i>join</i>	5.0%	29.6%
<i>sep.</i>	0.61%	26.7%
<i>sep. w/ top50</i>	0.68%	26.2%

Table: *OOV rate for 250K-word vocabularies on eval2011 data.*

# Filtering Background & Approach

- Initial system prototyping done around small (140 hrs) Quaero data
- Additional Quaero and Ester data did not initially improve performance
- It was decided to filter the additional acoustic training data

## Training steps

- 1 Divide training data into six parts
  - 1 Quaero-core (140hrs) (2010, 2011) (unfiltered)
  - 2 Quaero-fast (2009, 2010)
  - 3 Quaero-careful (2009, 2010)
  - 4 Ester1
  - 5 Ester2
  - 6 Ester2-dev
- 2 Decode on each part
- 3 Compute WER for each utterance
- 4 Reject utterances w.  $WER > X\%$
- 5 Retrain

# Filtering Background & Approach

- Initial system prototyping done around small (140 hrs) Quaero data
- Additional Quaero and Ester data did not initially improve performance
- It was decided to filter the additional acoustic training data

## Training steps

- 1 Divide training data into six parts
  - 1 Quaero-core (140hrs) (2010, 2011) (unfiltered)
  - 2 Quaero-fast (2009, 2010)
  - 3 Quaero-careful (2009, 2010)
  - 4 Ester1
  - 5 Ester2
  - 6 Ester2-dev
- 2 Decode on each part
- 3 Compute WER for each utterance
- 4 Reject utterances w.  $WER > X\%$
- 5 Retrain

- How do we choose  $X\%$ ?

## Quaero data

- Reject those top 10% of utterances with highest WER
  - → reject those with  $WER > 75\%$
- As alternative, reject those top 25%
  - → reject those with  $WER > 50\%$

## Ester data

- Ester data is different from Quaero data
  - BN vs. BC
  - Telephone and/or strongly accented speech
  - → Reject more utterances ( $WER > 37\%$ )

## Effect on WER

Training material	WER (CI)
baseline (Quaero core)	27.0
baseline + unfiltered add. Quaero & Ester	29.0
baseline + filtered add. Quaero (75%)	26.9
baseline + filtered add. Quaero (50%)	26.5
baseline + filt'd add. Q. (50%) + Ester 1, 2, & dev (37%)	26.3

## Effect on DB Size

Source	Utts. (K)			Hours		
	Unfilt.'d	Kept	%	Unfilt.'d	Kept	%
Quaero core	32.5	32.5	100	140	140	100
add. Quaero + Ester	31.7	18.1	56.9	162	128	78.9
Total	64.2	50.6	78.72	302	267	88.7

## Effect on WER

Training material	WER (CI)
baseline (Quaero core)	27.0
baseline + unfiltered add. Quaero & Ester	29.0
baseline + filtered add. Quaero (75%)	26.9
baseline + filtered add. Quaero (50%)	26.5
baseline + filt'd add. Q. (50%) + Ester 1, 2, & dev (37%)	26.3

## Effect on DB Size

Source	Utts. (K)			Hours		
	Unfilt.'d	Kept	%	Unfilt.'d	Kept	%
Quaero core	32.5	32.5	100	140	140	100
add. Quaero + Ester	31.7	18.1	56.9	162	128	78.9
Total	64.2	50.6	78.72	302	267	88.7

The Quaero Program

Decoding Strategy

Acoustic data & training techniques

Feature front-end developments

- Deep neural networks

- Tonal features

Pronunciation Modeling

- Pseudographemes

Language Modeling

Experiments

- Elision normalisation

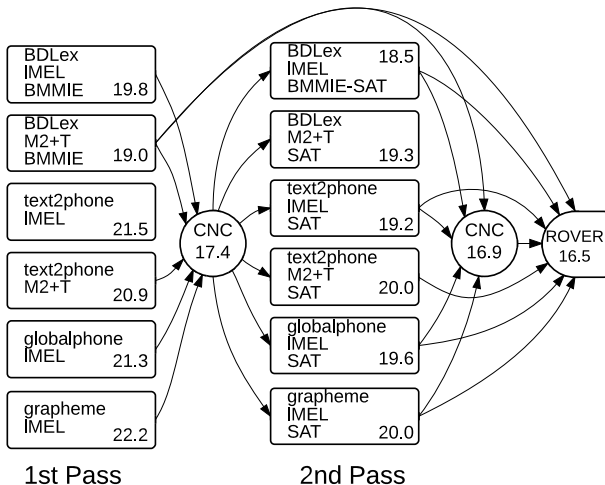
- Acoustic Data Filtering

**Results**

Conclusion

# System Results on Eval 2012

Results





The Quaero Program

Decoding Strategy

Acoustic data & training techniques

Feature front-end developments

- Deep neural networks

- Tonal features

Pronunciation Modeling

- Pseudographemes

Language Modeling

Experiments

- Elision normalisation

- Acoustic Data Filtering

Results

Conclusion

(Sub)System	eval2011	eval2012	eval2013
1st Single Best	18.6	19.0	18.6
2nd Pseudograph.	n.a.	20.0	18.9
ROVER	15.6	16.5	15.4

## Conclusions

- Created pseudographeme approach for French
- DBNFs & tonal features give significant gains
  - IMEL DBNFs better with SAT than M2+T
- Gains from filtering acoustic training data
- Proper elision normalisation is important

# Thank you!