

Discovering speech reductions across speaking styles and languages

Martine Adda-Decker, Lori Lamel

CNRS-LPP & *Spoken Language Processing Group* CNRS-LIMSI

November 30, 2017, Paris



- Temporal speech reduction
- Automatic speech processing as tools for linguistic studies
- Speech corpora & methodology
- Results
- Discussion

- Speech reduction: vowel reduction, consonant lenition, consonant cluster reduction, syllabic restructuring (e.g. Ernestus 2000, Duez 2003, Adda-Decker et al. 2005, Dilley & Pitt 2007, Van Son & Pols 2013)
- Traces of speech reduction in written language:
 - gonna be (going to be)
 - ça [sa] (cela [səla] 'that'), 'y a [ja] (il y a [ilija] 'there is')
 - ins [ins] (in das [ɪn das] 'in the')
- Temporal speech reduction (Adda-Decker & Snoeren 2010): includes any reduction process resulting in fewer segments in the produced speech

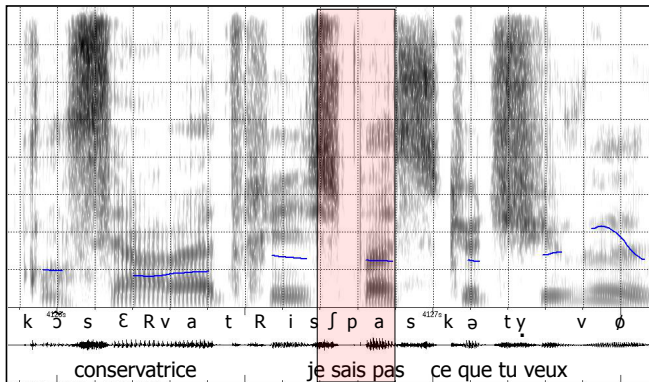
Speech reductions are known to frequently occur in spontaneous, interactive, less formal speech.

Speech reductions seem to first affect least informative speech portions (Jurafsky et al. 2001) and/or highly predictable from the context:

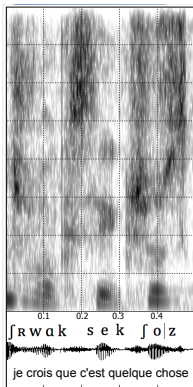
- function words
- discourse markers
- idioms
- dates...

Question: may speech reduction occur
in more formal (prepared) speech?
in more informative speech portions?

Examples - French - multiword reduction



Note : the sound [ʃ] stands for 2 words 'je' (I) and 'sais' (know)



'je crois que c'est quelque chose' (I believe it is something)

/ʒə kʁwa kə sɛ kɛlkə ʃoz/ [ʃɔwak sek ʃoz]

Temporal reduction phenomena raise issues:

- for automatic speech processing (both recognition, synthesis)
- for human processing in psycholinguistics
- for language learning/teaching...

Big data (large speech corpora) help to answer questions:

- Where do temporal reductions occur?
- Do they frequently occur?
- Are they conditioned by language, by speaking style?

Temporal reduction mainly appears in unstressed speech segments and is conditioned by speech style

1. careful speaking



clearly uttered

2. casual speech

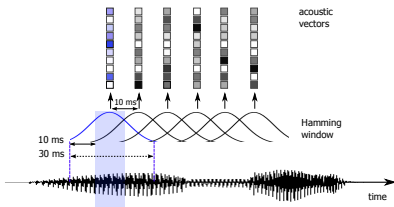


temporally reduced

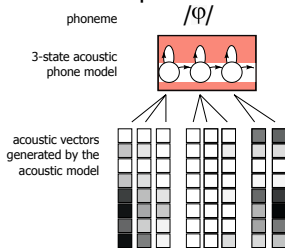
black boxes : stressed segments
grey boxes : unstressed segments

A bird-eye's view of acoustic modelling of speech in ASR (automatic speech recognition)

1. Parametrization of the signal

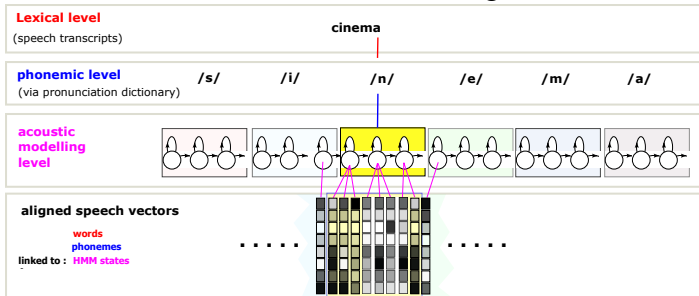


2. Acoustic phone modelling



A bird-eye's view of acoustic modelling of speech in ASR

3. Multi-level modelling



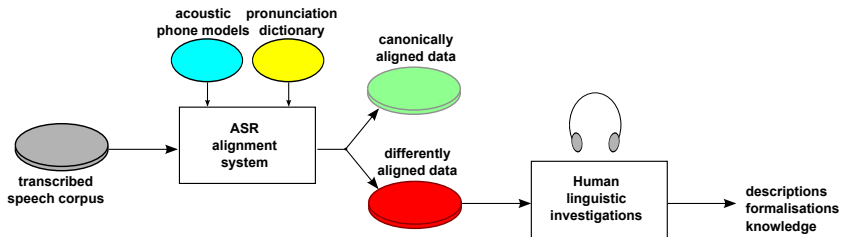
Pronunciations are given by a pronunciation lexicon
Variants may be introduced (e.g.: quatre [katʁə katʁ kat])

ASR: automatic speech recognition

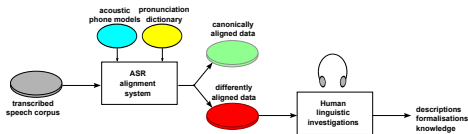
→ converts speech to text

alignment mode:

→ transcriptions constrain the matching process between speech and text



Define the meaning of canonically / differently aligned data!



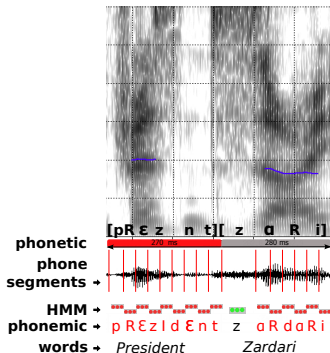
Define the meaning of canonically / differently aligned data!

- in case of (known/hypothesized) variants:
 - introduce relevant variants into pronunciation lexicon
 - canonically aligned = full form (reference) pronunciation
 - differently aligned = one of the variants
- for temporal reduction:
 - hypothesis: might occur anywhere (even though...)
 - introduce a segmental duration threshold on the aligned output
 - e.g. 30 or 40 ms segments

Speech sample : 'President Zardari' (strong temporal reduction)

Forced alignment using reference pronunciations:

→ sequence of short segments (30 ms duration) and false/misplaced labels



Phonetic: manual labelling

Phonemic: automatic labelling using the reference pron.

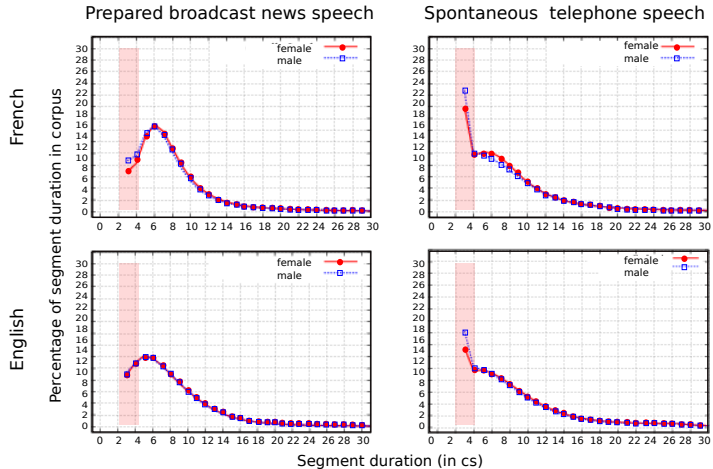
Different speech styles:

BN careful, prepared (news) and conversational (conv) speech

Casual telephone (tel) and face-to-face speech

| French | | | English | | |
|-------------------|---------------|----------|-------------------|---------------|----------|
| | # word tokens | duration | | # word tokens | duration |
| <i>BN-news</i> | 3600 k | 360 h | <i>BN-news</i> | 7200 k | 720 h |
| <i>BN-conv</i> | 600 k | 44 h | <i>BN-conv</i> | 1500 k | 124h |
| <i>Casual-tel</i> | 1000 k | 100 h | <i>Casual-Tel</i> | 25000 k | 2300 h |
| <i>Casual-f2f</i> | 350 k | 31 h | | | |

Phone segment duration distributions



Pink box : short segments (3-4 cs) involving potential temporal reduction
20% min. dur. segments for BN, 25% En / 30% Fr for casual speech

Case-study: English /t/

Position-dependent analysis in some typical frequent English words

word-initial position (stressed*)

word-medial position (stressed* / unstressed)

average phone durations are given in ms.

| | /t/ position | Broadcast Conversations | | | SWB/Fisher Conversations | | |
|----------|-----------------|----------------------------|---------------|----------------|-----------------------------|---------------|----------------|
| | | #token | avrg. dur. | % min. dur. | #token | avrg. dur. | % min. dur. |
| talking | w-init* | 814 | 95 | 5 | 4898 | 80 | 11 |
| trying | w-init* | 684 | 95 | 6 | 4464 | 85 | 11 |
| nineteen | w-med* | 560 | 80 | 7 | 706 | 89 | 8 |
| hotel | w-med* | 105 | 118 | 0 | 178 | 126 | 1 |
| little | w-med | 1041 | 59 | 41 | 9379 | 37 | 91 |
| getting | w-med | 803 | 59 | 52 | 5692 | 39 | 86 |
| exactly | w-med | 387 | 54 | 43 | 6328 | 39 | 85 |
| ninety | w-med | 323 | 70 | 20 | 210 | 43 | 84 |

Stress prevents from temporal reduction both in w-init. and w-medial pos., in both speech styles
In unstressed position, temporal reduction is high, especially in casual speech

Word-internal segment duration variation

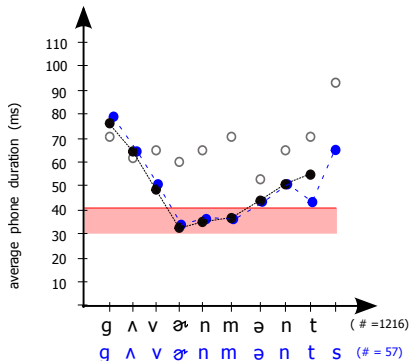


Case-study: English words 'government', 'governments'

word-initial word stress

empty circles: average phone duration (all segments pooled per phone)

coloured circles: word-position dependent average phone duration



word-internal segments in minimum duration region (30-40 ms)

Word-internal segment duration variation

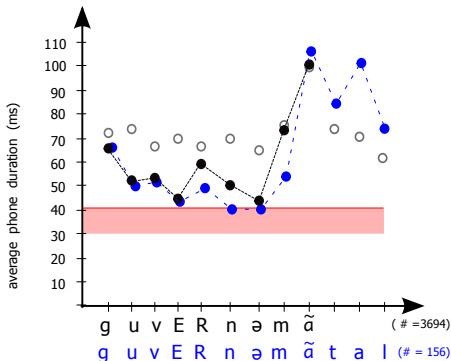


Case-study: French words 'gouvernement(s)', 'gouvernemental(e)'

stress on final syllable on these words

empty circles: average phone duration (all segments pooled per phone)

coloured circles: word-position dependent average phone duration



word-internal segments close to the minimum duration region (30-40 ms)

ASR multi-words with shortened pronunciation variants English Switchboard data

| <i>Multi-word</i> | #Total | <i>Full form + Variants</i> | #Align | %Align | Comments |
|-------------------|--------|---------------------------------|--------|--------|-------------------------|
| <i>did-not</i> | 2559 | did nɑt | 103 | 4.0 | full form |
| | | + didn̩t | 275 | 10.7 | n(ɑ→ ə) |
| | | + didn̩ | 1175 | 45.9 | + final-/t/ deletion |
| | | + dɪn | 1006 | 39.3 | + coda /d/ deletion |
| <i>we-have</i> | 3257 | wɪhæv | 1500 | 46.1 | full form |
| | | + wɪəv | 205 | 6.3 | onset /h/ del. + (æ→ ə) |
| | | + wɪv | 1552 | 47.7 | + V-deletion |

Multi-words with shortened pronunciation variants
English Switchboard data

| <i>Multi-word</i> | #Total | <i>Full form + Variants</i> | #Align | %Align | Comments |
|--------------------|--------|---------------------------------|--------|--------|-------------------------------|
| <i>going-to-be</i> | 750 | gɔŋg tubi | 73 | 9.7 | full form |
| | | + gɔŋəbi | 432 | 57.6 | complex: ɪŋg t → n |
| | | + gəbi | 245 | 32.7 | + complex: ɔŋə → ə |
| <i>wants-to</i> | 157 | wɔŋtstu | 15 | 9.6 | full form |
| | | + wɔŋstu | 78 | 49.7 | coda C-cluster simplification |
| | | + wɔŋtsə | 7 | 4.5 | onset /t/-deletion |
| | | + wɔŋsə | 57 | 36.3 | both /t/-deletions |

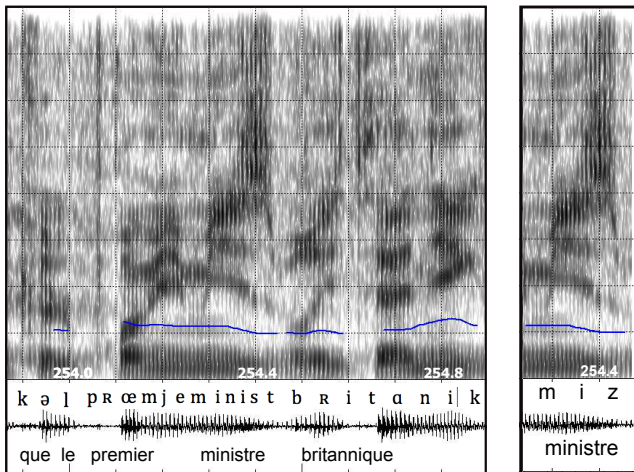
Words with shortened pronunciation variants
French casual speech corpus (NCCFr)

| Word | #Total | Full form + Variants | #Align | %Align | Comments |
|---------------------------------|--------|-------------------------|--------|--------|----------------------------|
| <i>parce (que)</i> 'because' | 2590 | pɑrsə | 4 | 0.2 | full form |
| | | + pɑrs | 45 | 1.7 | no final schwa |
| | | + pas | 1309 | 50.6 | + C-cluster simplification |
| | | + ps | 1232 | 47.6 | + vowel deletion |
| <i>quelques</i> 'some' | 56 | kɛlkə | 14 | 25 | full form |
| | | + kɛkə | 28 | 50 | + /l/-deletion |
| | | + kɛ(k—g) | 14 | 25 | + schwa deletion |

Example - French - reduction...



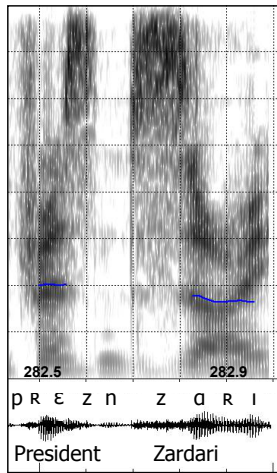
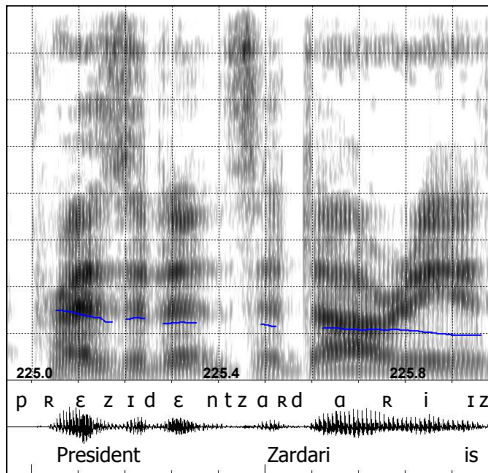
...in more formal, more informative speech portions



Example - English - reduction...



... in more formal, more informative speech portions



- Big Data: automatic speech processing generates ever increasing transcribed speech corpora
- Method: use of forced alignment to locate temporally reduced sequences in fluent speech → sequences of minimal duration segments reveal potentially reduced productions
- Temporal reduction observed in different speaking styles (broadcast conv., casual speech)
- Temporal reduction often involves unstressed stretches of speech
- Frequency and recency favor reduction

- Need to further foster interdisciplinarity between linguistics, phonetics, medicine, mathematics, signal processing, IT computer sciences...
- Foster team research using shared data and shared research questions
- Blend of expertise