# The Pauses & Lexical Stress Processing Pipeline (PLSPP)

Source code: https://gricad-gitlab.univ-grenoble-alpes.fr/lidilem/plspp

Sylvain Coulange

Univ. Grenoble Alpes, Laboratory of Linguistics and Didactics of Foreign and Mother Tongues (LIDILEM), Grenoble, France
Univ. Grenoble Alpes, CNRS, Institute of Engineering, Grenoble Computer Science Laboratory (LIG), Grenoble, France
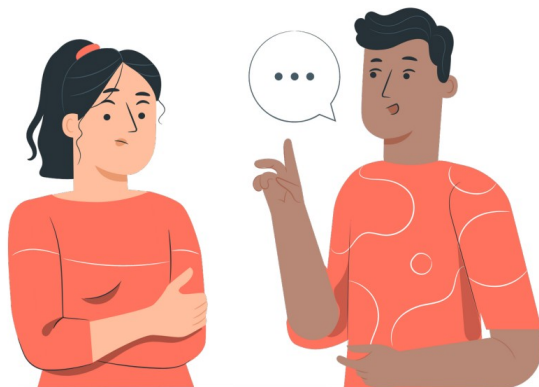Doshisha University, Spoken Language Processing Laboratory (SLPL), 610-0394 Kyoto, Japan

sylvain.coulange@univ-grenoble-alpes.fr

LIDILEM
Université
Grenoble Alpes

LIG

SLPL
Spoken Language Processing Laboratory

UGA
Université
Grenoble Alpes

同志社大学
Doshisha University

# Assessing L2 pronunciation: From nativelikeness to intelligibility

**Native speaker
as a target**

➡️

**Be (easily) understood**

**"Intelligibility"**

**"Comprehensibility"**

Isaacs, T., Trofimovich, P., and Foote, J. A. (2018) Developing a user-oriented L2 comprehensibility scale for english-medium universities. Language Testing 35(2), 193–216.
Jenkins, J., Baker, W., & Dewey, M. (Eds.). (2017) The Routledge Handbook of English as a Lingua Franca (1st ed.). Routledge.
Frost, D., O'Donnell, J. (2018) Evaluating the essentials, the place of prosody in oral production. In J. Volín (ed.). Pronunciation of EFL.
Council of Europe (2020) Common European framework of reference for languages. Strasbourg, France.
Walker, R., Low, E., & Setter, J. (2021) English pronunciation for a global world. Oxford: Oxford University Press

# Assessing L2 pronunciation: From nativelikeness to intelligibility

Rhythm

Beats

Speech flow

**Parameters related to L2 English comprehensibility:**

● Hesitation markers position (pauses, false starts, repetitions…)
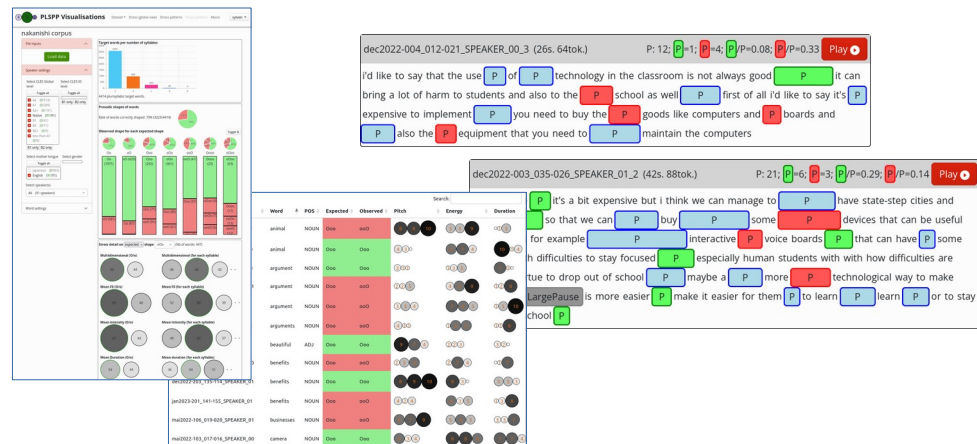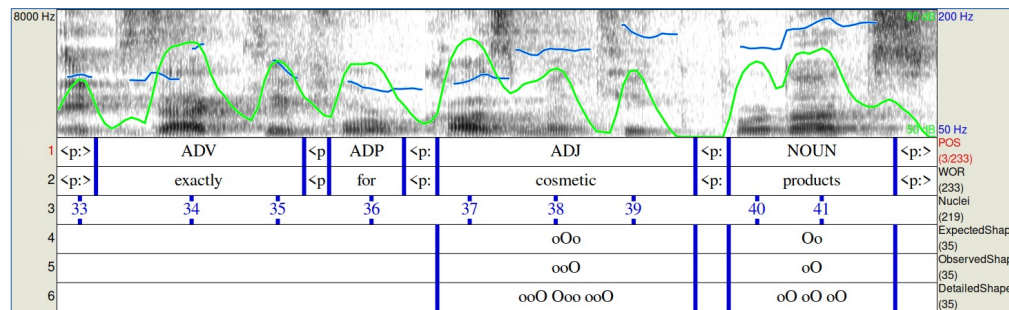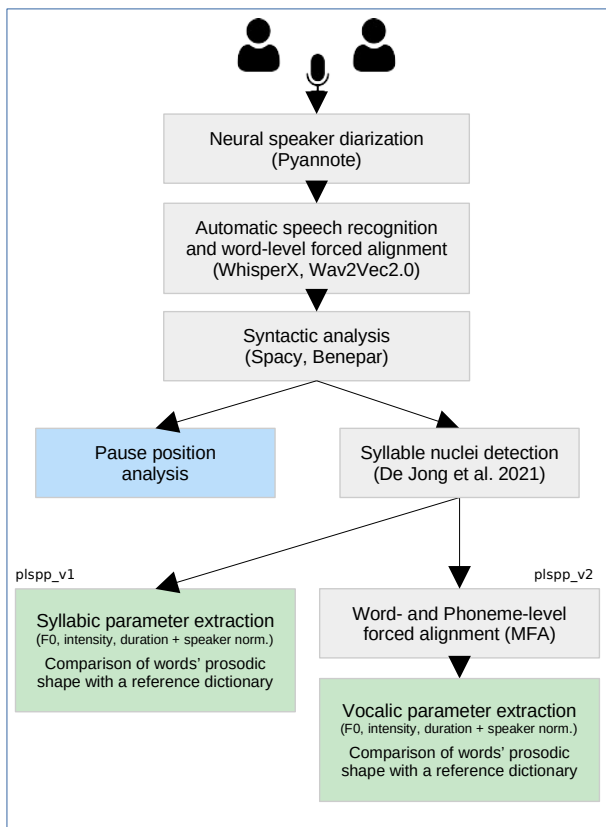
● Lexical stress (presence, position, quality)

● Speech rate (not too fast, not too slow)

● Pitch variation (make the speech sound lively and engaging)

● Phonemes quality (depending on their functional load)

Isaacs, T., Trofimovich, P., and Foote, J. A. (2018) Developing a user-oriented L2 comprehensibility scale for english-medium universities. Language Testing 35(2), 193–216.
Jenkins, J., Baker, W., & Dewey, M. (Eds.). (2017) The Routledge Handbook of English as a Lingua Franca (1st ed.). Routledge.
Frost, D., O'Donnell, J. (2018) Evaluating the essentials, the place of prosody in oral production. In J. Volín (ed.). Pronunciation of EFL.
Council of Europe (2020) Common European framework of reference for languages. Strasbourg, France.
Walker, R., Low, E., & Setter, J. (2021) English pronunciation for a global world. Oxford: Oxford University Press

# Assessing L2 pronunciation: From nativelikeness to intelligibility

**Rhythm**

**Beats**

**Speech flow**

**Parameters related to L2 English comprehensibility:**

- Hesitation markers position (pauses, false starts, repetitions…)
- Lexical stress (presence, position, quality)

**PhD**

Université Grenoble Alpes (France) - 3rd year

Doshisha University (Japan)

**Semi-automatic diagnosis of spontaneous English as a foreign language: the role of rhythm in speaker comprehensibility**

# The Pauses & Lexical Stress Processing Pipeline (PLSPP)

# Pauses processing



Customisable fixed duration
threshold (here 180ms-2s)

file: dec2022-003_039-040_SPEAKER_01_5
Speaker total speech duration: 6'33''

# Pauses processing



**3 categories:**
- Pauses between clauses
- Pauses between phrases
- Pauses within phrases

Customisable fixed duration threshold (here 180ms-2s)

# Pauses processing



**3 categories:**
- Pauses between clauses
- Pauses between phrases
- Pauses within phrases

Customisable fixed duration threshold (here 180ms-2s)

file: dec2022-003_039-040_SPEAKER_01_5
Speaker total speech duration: 6'33''

# Pauses processing

# Stress processing (PLSPP v1)

LIDILEM
Université
Grenoble Alpes

LIG

SLPL
Spoken Language Processing Laboratory

UGA
Université
Grenoble Alpes

同志社大学
Doshisha University

# Stress processing (PLSPP v1)



*(F0, Intensity, Duration)*

Percentiles
of the speaker's
distribution

F0
normalization

Intensity
normalization

Duration
normalization

|      | **37**    | **38**    | **39**   |
|------|-----------|-----------|----------|
| **F0:**  | (19),  | (80),  | (90)  |
| **dB:**  | (95),  | (28),  | (28)  |
| **dur:** | (55),  | (61),  | (78)  |

**Observed pattern**

(56),   (56),   (65)

○   ○   **O**

**Reference pattern**

○   **O**   ○

# Stress processing (PLSPP v1)



*(F0, Intensity, Duration)*

F0
normalization

Intensity
normalization

Duration
normalization

Percentiles
of the speaker's
distribution

|       | 37      | 38      | 39     |
|-------|---------|---------|--------|
| F0:   | (19),   | (80),   | (90)   |
| dB:   | (95),   | (28),   | (28)   |
| dur:  | (55),   | (61),   | (78)   |

Mean F0 (for each syllable)

19   80   90

Mean intensity (for each syllable)

95   28   28

Mean duration (for each syllable)

55   61   78

# Stress processing (PLSPP v1)



*(F0, Intensity, Duration)*

Percentiles of the speaker's distribution

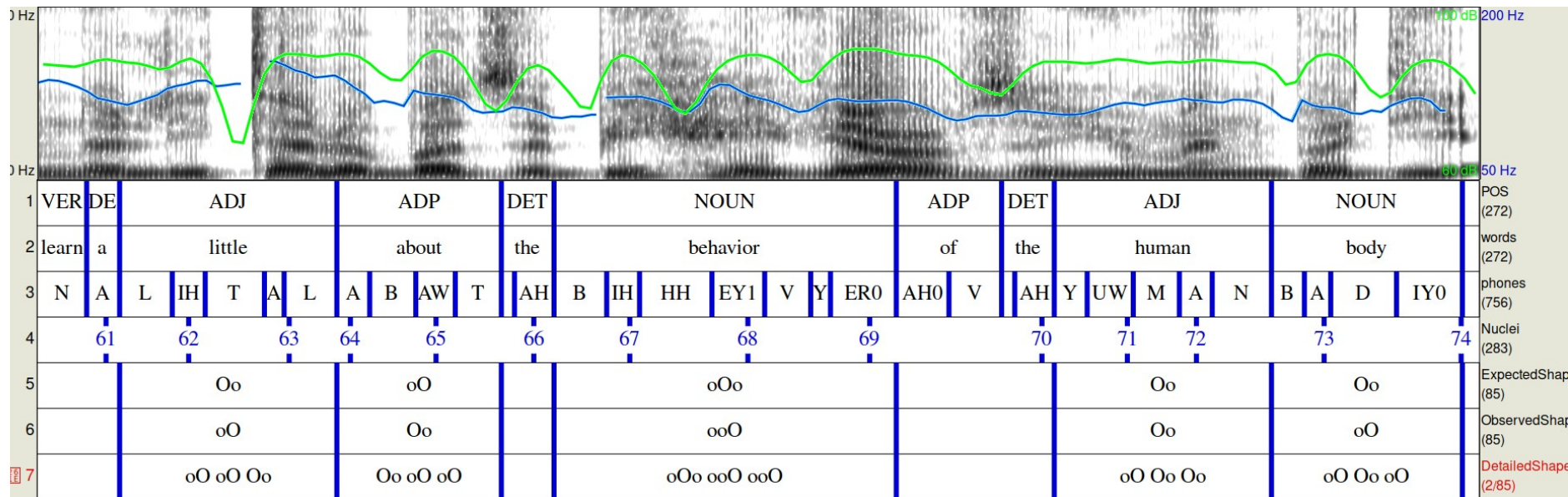Multidimensional (for each syllable)

Mean F0 (for each syllable)

Mean intensity (for each syllable)
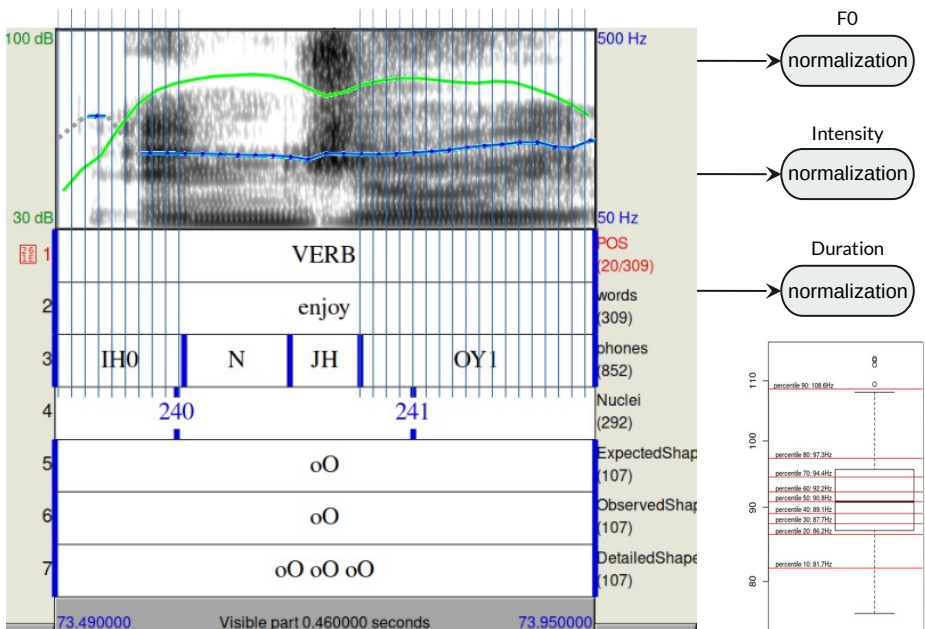
Mean duration (for each syllable)

# Stress processing (PLSPP v2)



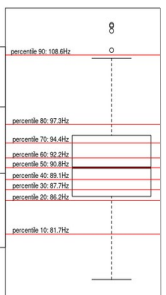| | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| VER | DE | | ADJ | | | ADP | | | DET | | | NOUN | | | | | ADP | | DET | ADJ | | NOUN | |

POS (272)

| learn | a | little | | about | | the | behavior | | | | of | | the | human | | body | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

words (272)

| N | A | L | IH | T | A | L | A | B | AW | T | AH | B | IH | HH | EY1 | V | Y | ER0 | AH0 | V | AH | Y | UW | M | A | N | B | A | D | IY0 |

phones (756)

| 61 | 62 | 63 | 64 | 65 | 66 | 67 | 68 | 69 | 70 | 71 | 72 | 73 | 74 |

Nuclei (283)

| Oo | oO | oOo | Oo | Oo |
ExpectedShap (85)

| oO | Oo | ooO | Oo | oO |
ObservedShap (85)

| oO oO Oo | Oo oO oO | oOo ooO ooO | oO Oo Oo | oO Oo oO |
DetailedShape (2/85)

14

# Stress processing (PLSPP v2)

(F0, Intensity, *Duration*)

Percentiles
of the speaker's
distribution

time_step = 10ms
*(customizable)*

## F0
- mean(F0s)
- (Min, max, sd)

  *pitch linear interpolation*

## Intensity
- max(dBs)

## Duration
- Length of vowel interval

# Visualizations

https://plspp.univ-grenoble-alpes.fr/

https://gricad-gitlab.univ-grenoble-alpes.fr/lidilem/plsppviz

# Studies using PLSPP

## PLSPP v1

### CLES Spontaneous speech

UGA Université Grenoble Alpes | CLES

Multispeaker spontaenous speech
University students (B1~B2)
L1: **French**

- Coulange S, Kato T, Rossato S, Masperi M. (2024). Enhancing Language Learners' Comprehensibility through Automated Analysis of Pause Positions and Syllable Prominence. Languages 9(3):78
- Coulange, S., Kato, T., Rossato, R., Masperi, M. (2023). Automatic Measurement of Lexical Stress in Spontaneous L2 English Speech of French Learners. Phonetic Society of Japan, Sep 2023, Sapporo, Japan. pp. 126-131
- Coulange, S., Kato, T. (2023). Pause position analysis in spontaneous speech for L2 English fluency assessment. Acoustic Society of Japan, Sep. 2023, Nagoya, Japan. pp. 991-994

**Corpus:**
- Coulange, S., Fries, M.-H., Masperi, M., Rossato, R. (2024). A corpus of spontaneous L2 English speech for real-situation speaking assessment. LREC-COLING 2024, 20-25 May, Torino, Italy.

### CLES-jp Spontaneous speech

同志社大学 Doshisha University | WASEDA University

Multispeaker spontaenous speech
University students (A2~C1)
L1: **Japanese**, **English**

**Corpus:**
- Coulange, S., Konishi, T., Kato, T., Sugahara, M., Rossato, R., Masperi, M. (2024). A corpus of spontaneous dialogues in L2 English by French and Japanese L1 speakers for automated assessment of fluency. 6th International Symposium on Learner Corpus Studies in Asia and the World (LCSAW6), Feb. 2024, Kobe, Japan.
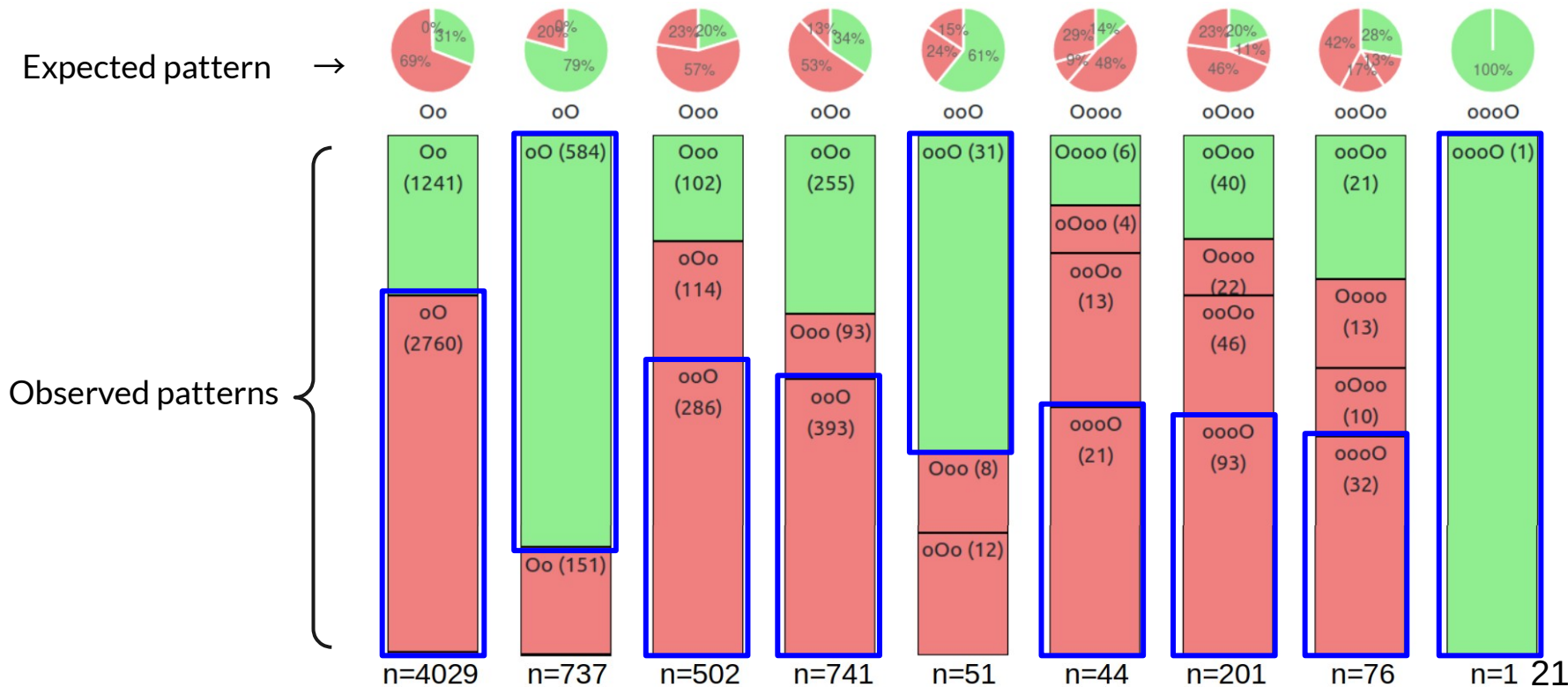
## PLSPP v2

### Fluency evaluation

UGA Université Grenoble Alpes

Read-aloud
University students (B1-B2)
L1: **French**

**PIC** Prosodie, Intelligibilité, Communication (PIC) (Frost, D.)

- Paper coming soon :)

### Fluency evaluation

神戸学院大学 KOBE GAKUIN UNIVERSITY

Read-aloud
University students (A1-B2)
L1: **Japanese**, **English**

### PLSPP v3...

- Nakanishi, M., Coulange, S. (2024). Measuring speech rhythm through automated analysis of syllabic prominences. Prosodic features of language learners' fluency (Speech Prosody WS), July 1, Leiden.

### Stress awareness vs. stress production

同志社大学 Doshisha University

Carrier phrases
L1: **Japanese**, **Korean**, **English**

- Sugahara, M., Coulange, S., Kato, T. (2024). English Lexical Stress in Awareness and Production: Native and Non-native Speakers. The 19th Conference on Laboratory Phonology, June 27-29, Seoul.
- Sugahara, M., Coulange, S., Kato, T. (2023). Stress awareness vs. stress production: Comparison of primary stress assignment to English words between Japanese and Korean university students. 347th regular meeting of the Phonetic Society of Japan, Nov 25, online.

### Automatic vs. native speakers' evaluation of lexical stress

同志社大学 Doshisha University

Text recitation
Elementary school children (A2-B1)
L1: **Japanese**

- Kimura, T., Coulange, S., Kato, T. (2024). Automatic estimation and native speakers' evaluation of lexical stress positions in English recitation speech produced by Japanese elementary school children. Spring Meeting of the Acoustic Society of Japan, Mar 6-8, Tokyo.

# Current PhD experiment: Corpus

**Corpus:**



- ✔ L2 English spontaneous speech from 176 French learners recorded during CLES certification speaking session.

- ✔ Situation: 2 or 3 candidates discussing a polemical topic (role play) during 10min.

- ➢ Total 11 hours of continuous speech (per speaker: mean 3'44", min 32", max 6'51)

- ➢ Speaking B1 level: 34%, B2 level: 66%

- ➢ Speech duration: B1≈B2, Nb tokens: B1<B2, Nb pauses: B1<B2, Silence proportion: B1≈B2

**Hypothesis:**

- **Pauses:**
  - More random pauses with B1
  - More structurant pauses with B2

- **Stress:**
  - Stress position accuracy B2>B1
  - Lower contrast stressed/unstressed
  - Stress shift to last syllable

18

# Current PhD experiment: Corpus

**Corpus:**

✔ L2 English spontaneous speech from 176 French learners recorded during CLES certification speaking session.

✔ Situation: 2 or 3 candidates discussing a polemical topic (role play) during 10min.

➢ Total 11 hours of continuous speech (per speaker: mean 3'44", min 32", max 6'51)

➢ Speaking B1 level: 34%, B2 level: 66%

➢ Speech duration: B1≈B2, Nb tokens: B1<B2, Nb pauses: B1<B2, Silence proportion: B1≈B2

**Hypothesis:**

- **Pauses:**
  - More <u>random pauses</u> with B1
  - More <u>structurant pauses</u> with B2

  *(intra-phrase)*

  *(inter-clause)*

- **Stress:**
  - Stress position accuracy B2>B1
  - Lower contrast stressed/unstressed
  - Stress shift to last syllable

19

# Stress position



Expected pattern →

Observed patterns

# Stress position

# Stress position



B1 speakers
spk=59
words=1873

B2 speakers
spk=117
words=4551

# Current PhD experiment: Stress position analysis



Proportion of target words with correct stress position per speaker (n=176)

➢ *Mean stress position accuracy: 35.4 %*

➢ *Stress accuracy per speaker: 0 % ～ 68.4 %*

➢ *Stress accuracy per CEFR level: B1 = 29.6 %    B2 = 36 % (p<.001)*

# Stress quality: dimension

全ての話者（176 人）

# Stress quality: dimension

## Expected Ooo

**Multidimensional (for each syllable)**
⑧ | 25 | 29

**Mean F0 (for each syllable)**
⑪ | 58 | 30

**Mean intensity (for each syllable)**
⑦ | 19

**Mean duration (for each syllable)**
⑤ ⑬ | 38

→

**Multidimensional (O/o)**
⑧ | 27

**Mean F0 (O/o)**
⑪ | 44

**Mean intensity (O/o)**
⑦ ⑫

**Mean Duration (O/o)**
⑤ | 26

→

**Multidimensional (O/o)**
45 | 54

**Mean F0 (O/o)**
40 | 59

**Mean intensity (O/o)**
50 | 51

**Mean Duration (O/o)**
45 | 52

25

# Current PhD experiment: Stress quality analysis

# Current PhD experiment: Stress quality analysis

Stress position accuracy:

**65%**

**58%**

**60%**

Stress position accuracy:

**21%**

**16%**

**19%**

# Current PhD experiment: Main observations

● First prototype of the Pauses and Lexical Stress Processing Pipeline

● Analysis of B1 and B2 speaking level French-L1 university students
   11 hours of speech     6350 target words     21 831 pauses

➢ Pause position:
   ○ Great variation of number of pauses within phrases among speakers, less with pauses between clauses
   ○ B2 speakers make less pauses within phrases than B1 speakers (p<0.01)
   ○ Difference between B1 and B2 is small
   ○ High intra-speaker variability

➢ Lexical stress position:
   ○ *Mean stress position accuracy: 35.4 %*
   ○ *Stress accuracy per speaker: 0 % ～ 68.4 %*
   ○ *Stress accuracy per CEFR level:*
     *B1 = 29.6 %     B2 = 36 %   (p<0.001)*
   ○ *Frequent stress shift to the last syllable*
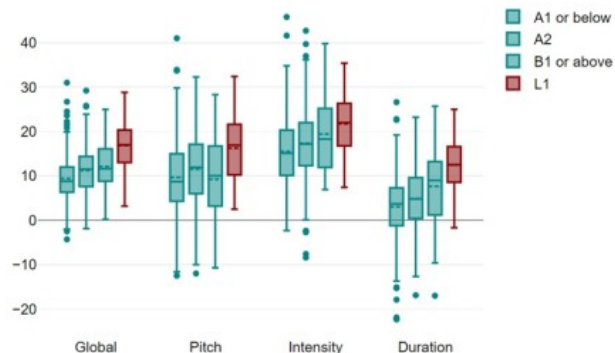
➢ Lexical stress quality:
   ○ *Low accuracy speakers: lengthening of the last syllable*
     *tendency to make it higher*
     *No change in intensity*
   ○ *High accuracy speakers: the expected syllable is higher in F0 and intensity*
     *No change in duration*

# Nakanishi & Coulange (2024)

- 34 hours read-aloud speech

- 877 Japanese-L1 samples (42 speakers, <A1-B2)

- 91 Native English samples (7 professional narrators)

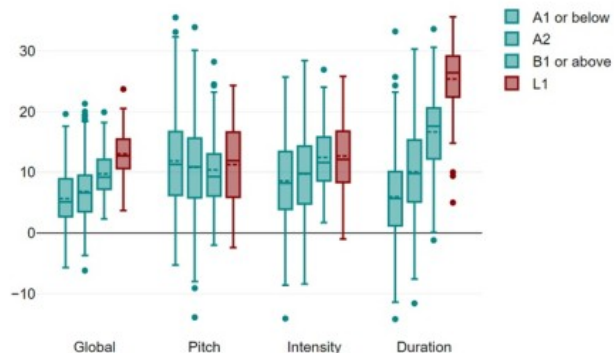- PLSPP extension to monosyllabic words

> analysis of contrast between **content and function words**

Figure 1. Syllabic Contrast Scores within Monosyllabic Words by CEFR Level.

Global scores between groups ($p < .001$)
A1 <*** A2 *n.s.* B1 < *** L1

Figure 2. Lexical Contrast Scores between Content and Function Words by CEFR Level.

Global scores between groups ($p < .001$)
A1 <*** A2 <*** B1 < *** L1
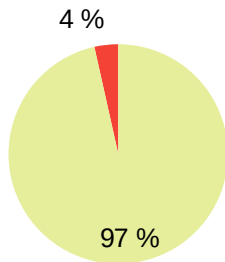
# Pipeline Evaluation & Limitations:

- As the pipeline combines several modules, errors can occur at different levels, often leading to incorrect annotations.

- ▲ Syllable detection and word alignment often mismatches, leading to a limited nb. of target words (only **41%** of polysyllabic words in the study below were **target words**).

- ▲ Manual evaluation of random 100 target words showed that **17%** were miss-recognized or miss-aligned, potentially leading to wrong judgments that can be problematic in a real assessment context.

- ▲ **Intrinsinc vowel length** and **word ending lengthening** need to be considered in order to improve stress estimation.

- ▲ Some cases of **vowel devoicing** also impacted F0 measures (tackled with linear interpolation for now)
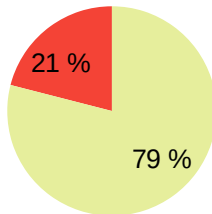
# Word alignment precision

Number of target words with <u>totally wrong alignment</u>,
<mark>among the first 200 plain target words</mark> in the visualization interface:
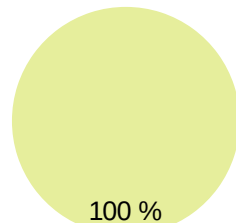
plspp : 7 target words

4 %

97 %

plspp_mfa : 42 target words

21 %

79 %

Corpus PIC (Frost, D.)
(280 speakers Read speech ~1min20s/spk)

Plspp: 0 words

100 %

plspp_mfa: 7 words

4 %

97 %