

Cross-linguistic investigations of the organizational scales in phonological systems

*Christophe Coupé, Egidio Marsico ,
Yoon Mi Oh & François Pellegrino*

Laboratoire Dynamique du Langage
Université de Lyon & CNRS, France



Issue

How do the phonological units of a given language carry the weight of differentiating the words of its lexicon?

- Segments
- Syllables
- Features

Approach (I)

The notion of **Functional Load (FL)**

- i. Relates to the role a phonological contrast plays in keeping words distinct in a given language (e.g. Trubeztkoy, 1939)
- ii. Has been considered with respect to language evolution (Martinet, 1977; King, 1967; Surendran & Niyogi, 2003), language acquisition (Van Severen et al., 2012) etc.
- iii. Has been seen as a useful, supplement to “standard” phonological descriptions.
- iv. Mostly centered on phonemes, but some investigations into the FL of features (Surendran & Niyogi, 2003)

Approach (II)

A quantitative approach to FL to shed light on the organization of phonological systems

A cross-linguistic, corpus-based approach

Vowels, consonants (but also tones/stress)

Material: digital lexicons

(with frequency of use of forms)

Language	ISO 639-3 Code	Source
Cantonese	YUE	A linguistic corpus of mid-20th century Hong Kong Cantonese (Research Centre on Linguistics and Language Information Sciences, 2013)
English	ENG	WebCelex (Max Planck Institute for Psycholinguistics, 2013, 2014)
Japanese	JPN	The corpus of spontaneous Japanese (NINJAL, 2011)
Korean	KOR	(Leipzig corpora collection)
Mandarin	CMN	Chinese Internet Corpus (Sharoff et al, 2006)
German	DEU	WebCelex (Max Planck Institute for Psycholinguistics, 2013, 2014)
Swahili	SWH	Helsinki corpus of Swahili (Gelas, Besacier, & Pellegrino, 2012)
Italian	ITA	PAISÀ Corpus (Lyding et al., 2014)
French	FRA	Lexique 3.80 (New et al., 2001)

Simple vs. complex syllabic structures (Maddieson et al., 2013), different morphological types

20,000 most frequent words (inflected forms) considered, except for Cantonese (5,000 forms) & Italian (15,788 forms)

Methodology (I)

Carter (1987)'s quantitative definition of FL

Language L considered as a source of sequences of independent words w_i taken from a finite set N_L

FL of a contrast x/y = quantification of the perturbation induced by merging x and y in terms of increase of homophony and of changes in the distribution of word frequencies

$FL(x,y)$: relative difference in entropy between the observed state L and a fictive state L_{xy}^* in which the contrast is neutralized

$$FL(x, y) = \frac{H(L) - H(L_{xy}^*)}{H(L)} \quad H(L) = - \sum_{i=1}^{N_L} p_{w_i} \cdot \log_2(p_{w_i})$$

Methodology (II)

Form	Frequency
pal	300
pil	200
bal	150
bil	150
pul	100
bul	100
TOTAL	1000

Contrast a-i

Form	Frequency
p <i>a</i> l	300
p <i>i</i> l	200
b <i>a</i> l	150
b <i>i</i> l	150
pul	100
bul	100
TOTAL	1000

Form	Frequency
p <i>a</i> l	500
b <i>a</i> l	300
pul	100
bul	100
TOTAL	1000

$$H(L^*_{ai}) = 1.69$$

$$FL(a-i) = (2.47 - 1.69) / 2.47 = \mathbf{31.8\%}$$

Inventory: /a i u p b l/

$$N_L = 6 \quad H(L) = 2.47$$

Phoneme /a/

From contrasts to units

$$FL(x) = \frac{1}{2} \sum_y FL(x, y)$$

$$FL(a) = \frac{1}{2} (FL(a-i) + FL(a-u)) = \frac{1}{2} (31.8 + 23.1) = \mathbf{27.45\%}$$

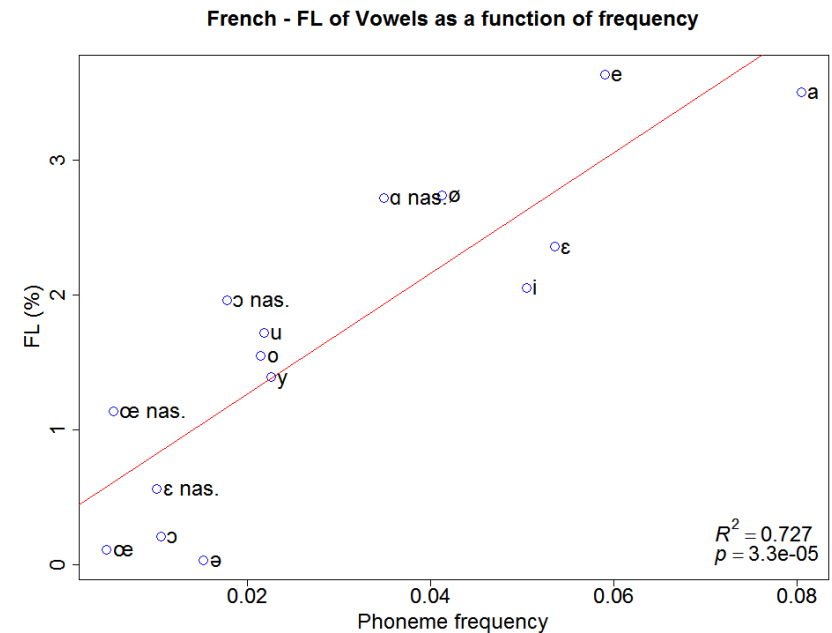
First scale of organization: segments

Language	ISO 639-3 Code	Phonological system	
Cantonese	YUE	V	13
		C	19
		T	6
English	ENG	V	22
		C	28
		S	2
Japanese	JPN	V	10
		C	16
Korean	KOR	V	8
		C	22
Mandarin	CMN	V	7
		C	25
		T	5
German	DEU	V	22
		C	24
		S	1
Swahili	SWH	V	5
		C	30
Italian	ITA	V	8
		C	25
		T	1
French	FRA	V	15
		C	21

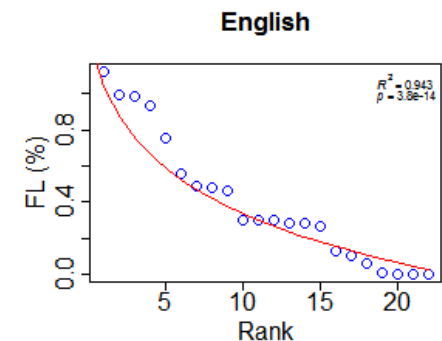
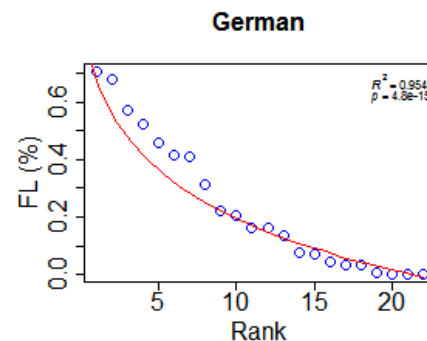
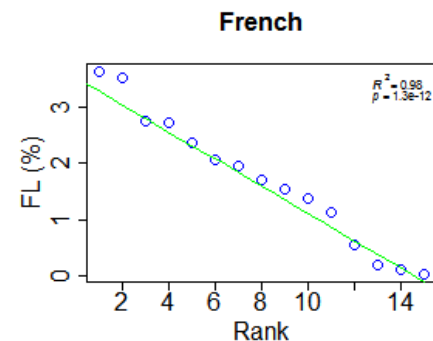
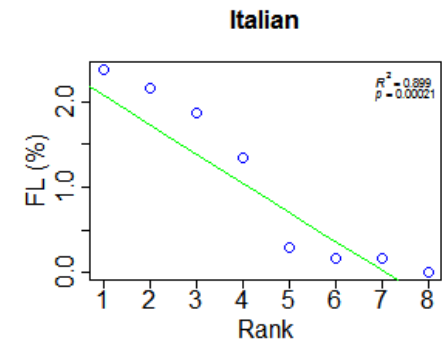
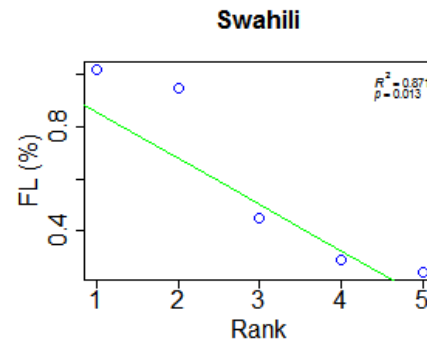
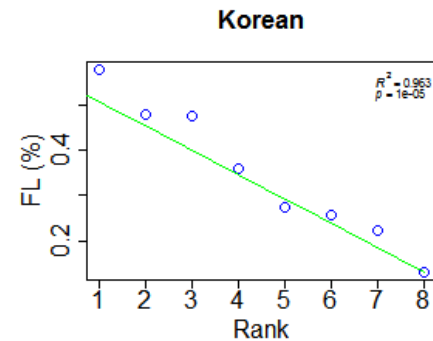
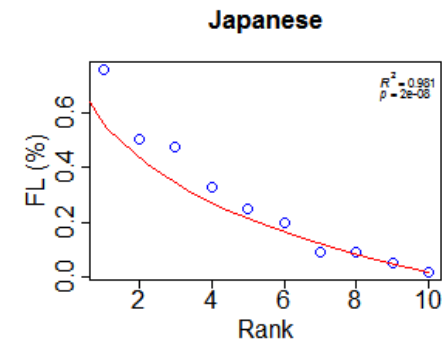
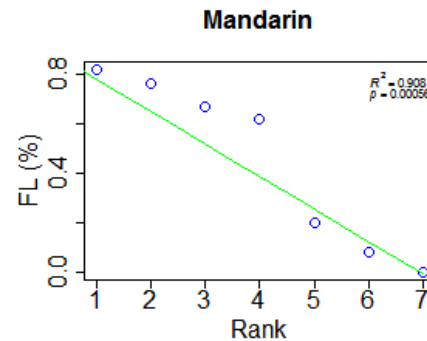
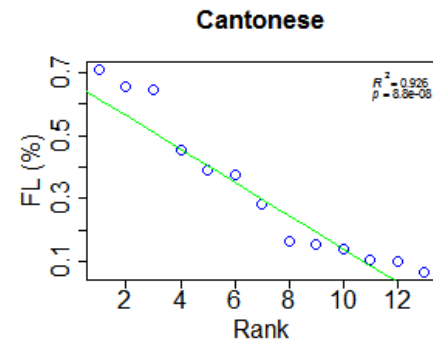
Vowels include diphthongs

Languages differ in terms of the phonemes they use.

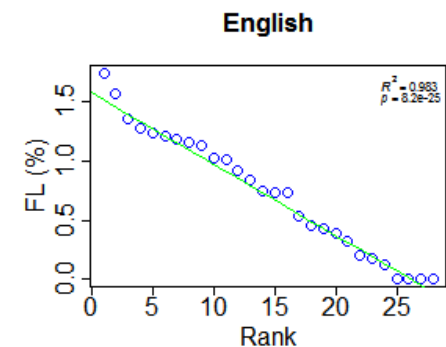
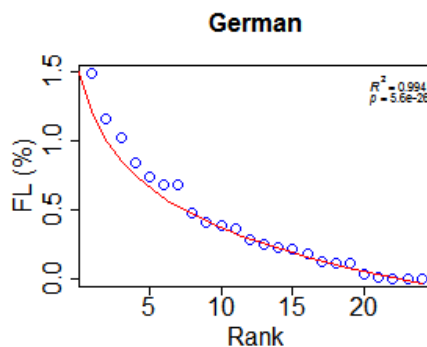
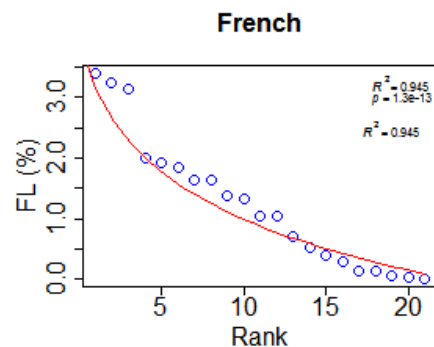
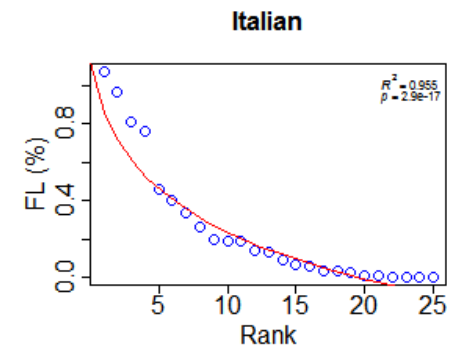
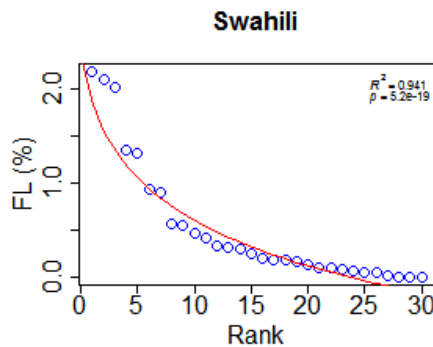
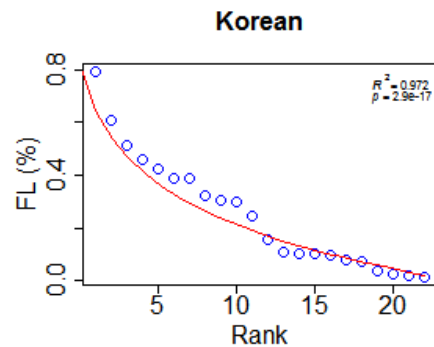
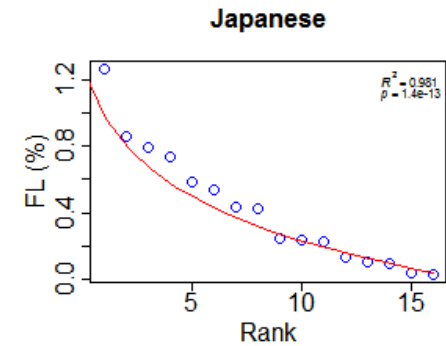
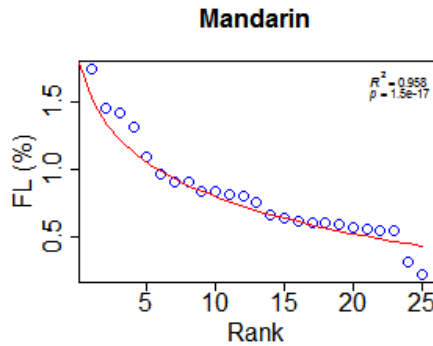
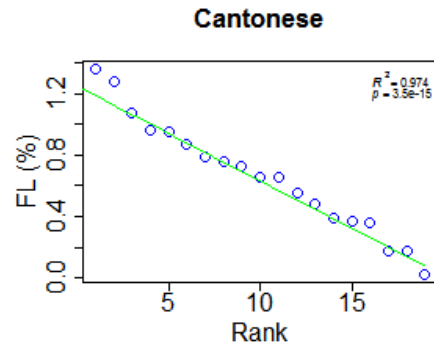
Some languages rely on much more vowels or consonants than others



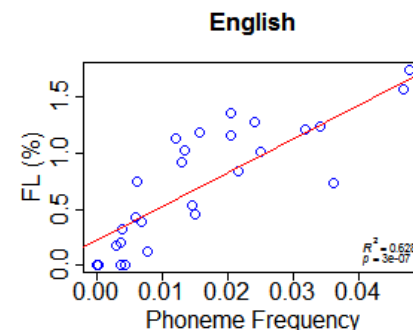
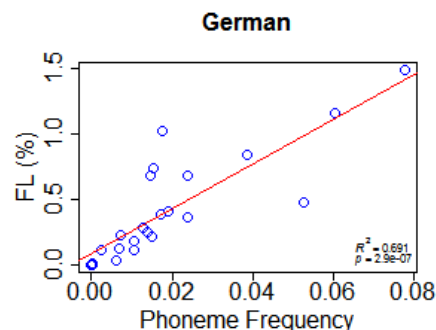
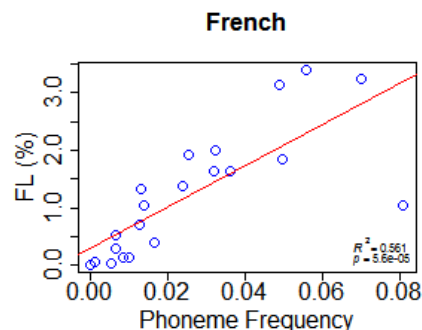
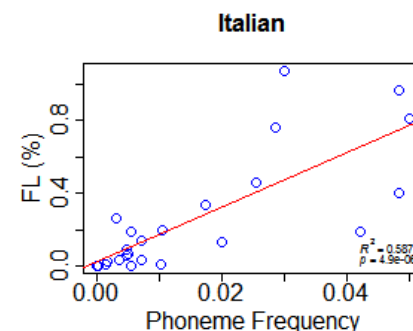
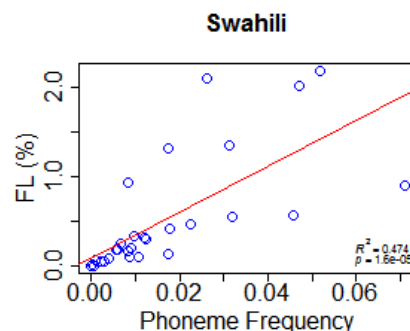
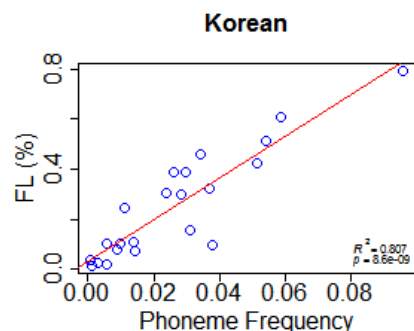
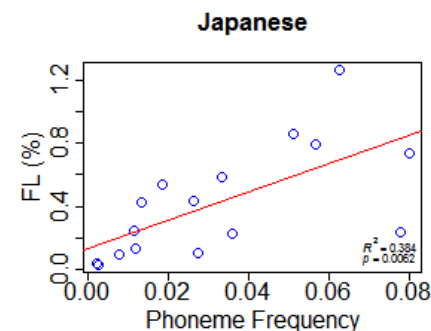
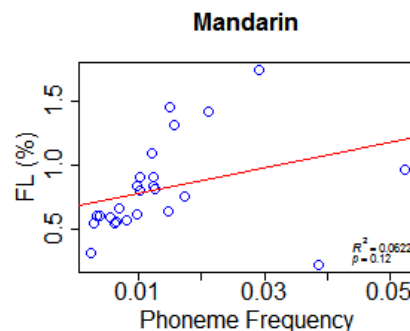
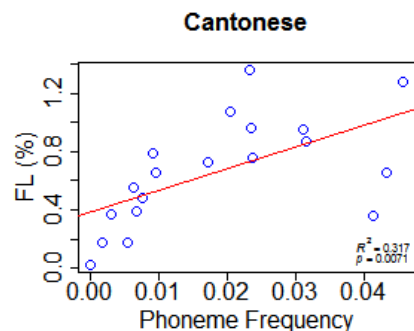
FL of vowels in 9 languages



FL of consonants in 9 languages



FL of consonants as a function of frequency



Discussion on segments

In all languages, some segments carry a heavier burden of differentiating the words than others

→ No principle of uniformity regarding functional load

Frequency is a partial predictor of FL, but other factors underlie the distinctive function within a system

Left for discussion time:

The segments with high FL differ widely from one language to another
There seem to be no *strong* cross-linguistic tendency.

Supra-segmental phonemes: the FL of tones is equivalent to the FL of vowels in tonal languages (yue & cmn)

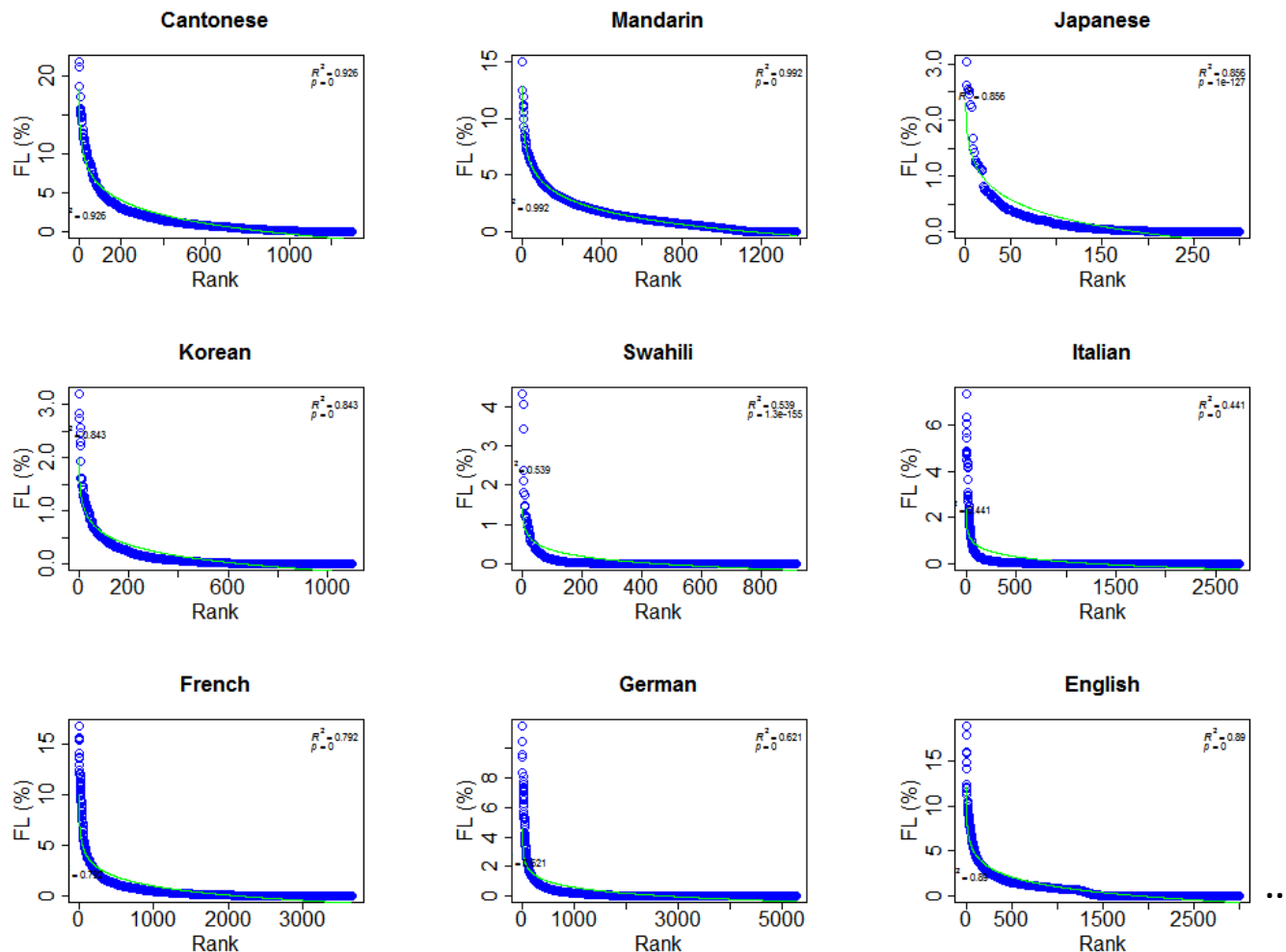
Second scale of organization: syllables

Language	ISO 639-3 Code	Phonological system		Number of different syllables
Cantonese	YUE	V	13	1298
		C	19	
		T	6	
English	ENG	V	22	8298
		C	28	
		S	2	
Japanese	JPN	V	10	300
		C	16	
Korean	KOR	V	8	1100
		C	22	
Mandarin	CMN	V	7	1283
		C	25	
		T	5	
German	DEU	V	22	5256
		C	24	
		S	1	
Swahili	SWH	V	5	914
		C	30	
Italian	ITA	V	8	2735
		C	25	
		T	1	
French	FRA	V	15	3666
		C	21	

Languages differ widely in terms of number of syllables

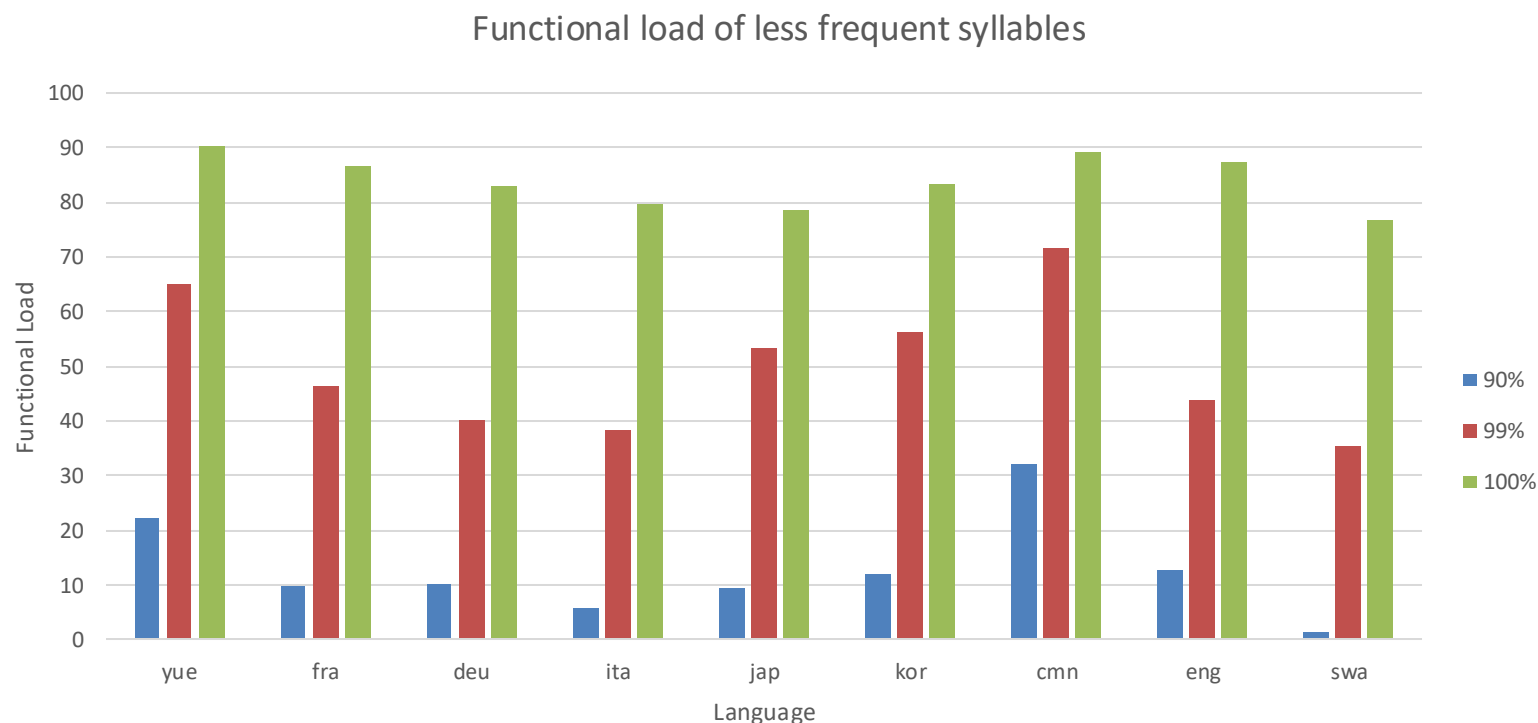
The FL of syllables can be computed in the same way they were for segments

FL of syllables according to rank



Relative to segments, the distribution of syllabic FL is more skewed (esp. in languages like **deu**, **ita** or **swa**): in each language, a large number of syllables have a very low FL.

FL of less frequent syllables



The green bars refer to a situation where words are only differentiated by their number of syllables, not by the nature of these syllables

Even when 90% of the (less frequent) syllables are merged into one, words can still be well differentiated thanks to their structure and to the most frequent syllables
Cantonese and Mandarin are more affected than other languages

Discussion on syllables

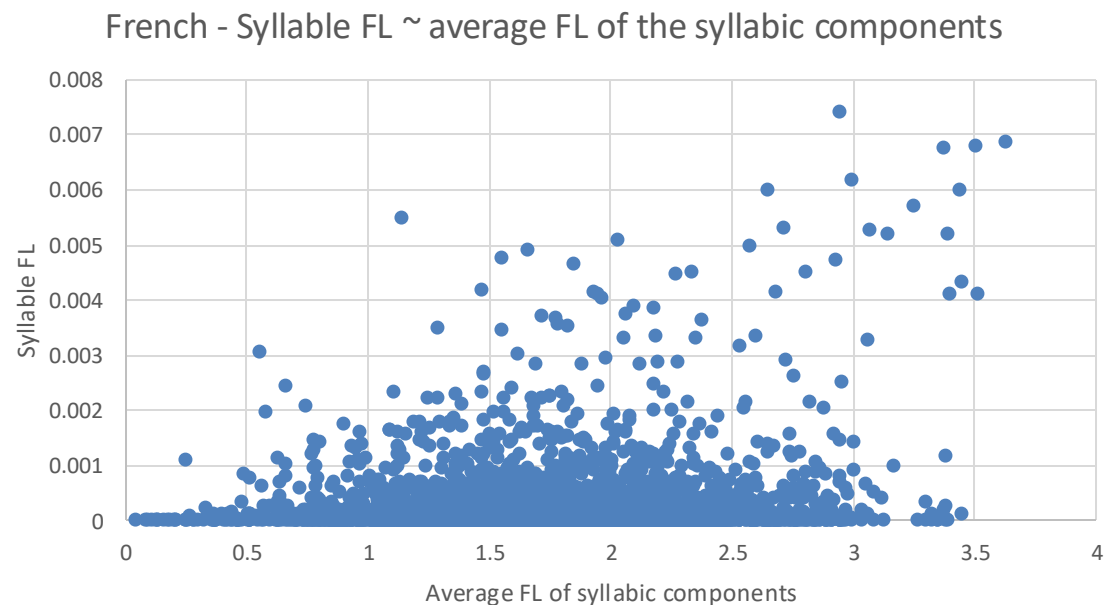
Non-linear and partial relationship between FL and frequency of use (sometimes strong, e.g. yue & cmn) (*not shown*)

Why do we have so many syllables, when most of them have very low FL?

Is this distribution, at the syllabic level, functional? or...

Does the FL of syllables derive from the FL of their components?

The FL of syllables does not correlate with either the average, the maximum value or the product of the FL of its components
Is there another relationship?



Third scale of organization: features

Language	ISO 639-3 Code	Phonological system		Number of features
Cantonese	YUE	V	13	12
		C	19	18
English	ENG	V	22	27
		C	28	19
Japanese	JPN	V	10	10
		C	16	15
Korean	KOR	V	8	10
		C	22	17
Mandarin	CMN	V	7	11
		C	25	19
German	DEU	V	22	21
		C	24	18
Swahili	SWH	V	5	9
		C	30	19
Italian	ITA	V	8	10
		C	25	18
French	FRA	V	15	12
		C	21	17

E.g.

/i/: high front unrounded

/p/: bilabial voiceless stop

(Mostly articulatory) description of segments in terms of features based on UPSID
(Maddieson & Precoda, 1990)

Obvious differences among languages. Some articulatory dimensions and specific features are always present, others are not.

How to compute the FL of features?

A diagram showing the relationships between the vowels of the Latin alphabet. The vowels are arranged in a circle: i, u, u, o, o, ɔ, ɔ, ɛ, ɛ, e, e, i. The relationships are indicated by lines connecting the vowels: i to u, u to o, o to ɔ, ɔ to ɛ, ɛ to e, e to i. There are also lines connecting i to e, u to ɔ, and o to ɛ.

Vocalic inventory of Korean

Aperture: 3 sets of merges

ε, e, i
o, u
Λ, Ψ

Anteriority: 2 set of merges

ε, \wedge
i, \cup

E.g. To compute the FL of **aperture**, the initial lexicon is contrasted with a lexicon where 3 different sets of merges create homophony and modify the distribution of word frequencies. Note: consistent combinations of features (e.g. front-unrounded) not considered here

FL of articulatory dimensions

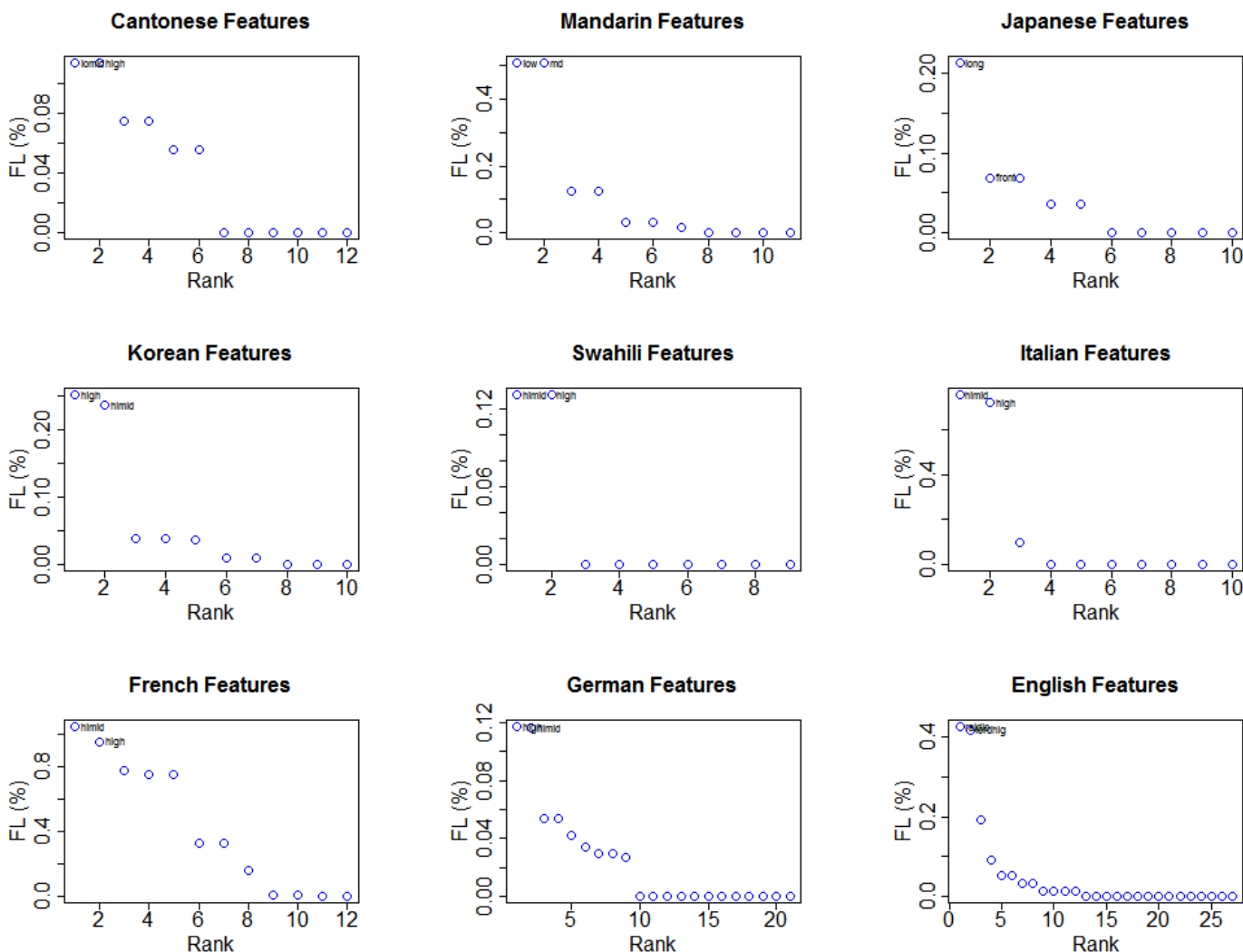
	V						C		
	<i>Aperture</i>	<i>Anteriority</i>	<i>Roundedness</i>	<i>Length</i>	<i>Aspiration</i>	<i>Nasalization</i>	<i>Place</i>	<i>Manner</i>	<i>Voicing</i>
yue	0,23	0,15	0,11		0,33		1,72	1,07	
cmn	1,02	0,06	0,25		0,67		2,07	1,04	0,18
jap	0,07	0,14		0,21			0,79	0,26	0,36
kor	0,53	0,08	0,02		0,01		0,67	1,01	0,05
swa	0,26						1,84	1,66	0,13
ita	1,62						0,22	2,03	0,12
fra	2,66	0,66	1,50			0,16	1,56	3,04	0,89
deu	0,31	0,11	0,06	0,03			0,91	2,72	0,11
eng	1,19	0,11					0,99	2,45	0,59

	Primary articulatory dimension with the highest FL
	Primary articulatory dimension with the 2nd highest FL
	Primary articulatory dimension with the 3rd highest FL

Regarding vowels and primary articulatory dimensions, **aperture** carries the heaviest load in 8 of the 9 languages. **Secondary features can have a high / the highest FL.**

Regarding consonants, languages seem to choose **either place or manner** as the primary way to differentiate between words. **Voicing** always comes after except in Japanese.

FL of vocalic features according to rank



The distribution is never uniform.

Most often, 1 or 2 features have a higher FL than others (but consistent combinations)

Discussion & perspectives

Cross-linguistic tendencies:

Whatever the organizational scale, phonetic units do not evenly carry the same FL. A few units carry a heavy load, while most others only carry a very light load.

Partial explanations can be provided at each level : morphosyntactic rules, ease of production versus sufficient perceptive contrast etc. (language-specific or not)

Cross-linguistic diversity:

Languages widely differ in the units they shoulder with the task of keeping words distinct

Are there ‘vertical’ integrative processes taking place in addition to ‘horizontal’ constraints?

The FL of syllables does not seem to derive straightly from the FL of their components

Perspectives:

Investigate whether the FL of segments derive from the FL of their features

Consider pairs of features in addition to single features.

Thank you for your attention



Yoonmi Oh



Egidio Marsico



François Pellegrino

Oh, Y., Pellegrino, F., Coupé, C. & Marsico, E., 2013. “ Cross-language Comparison of Functional Load for Vowels, Consonants, and Tones”. Proc. of Interspeech 2013, Lyon, France, 25-29 August.

Oh, Y., Coupé, C., Marsico, E. & Pellegrino, F. (accepted, to appear). “Bridging Phonological System and Lexicon: Insights from a Corpus Study of Functional Load”. *Journal of Phonetics*.

Funding: LABEX ASLAN (ANR-10-LABX-0081) of Université de Lyon - Program "Investissements d'Avenir" (ANR-11-IDEX-0007) of the French government