A Phoneme-weighted Levenshtein Distance: A Measure of Cross-linguistic Phonological
Similarity

Serene Siow[1], Gonzalo García-Castro[2], Núria Sebastian-Galles[2], & Kim Plunkett[1]

[1] Department of Experimental Psychology, Oxford University

[2] Center for Brain and Cognition, Universitat Pompeu Fabra

Author Note

SS and GGC share first authorship.

Correspondence concerning this article should be addressed to Serene Siow, Anna
Watts Building, Radcliffe Observatory Quarter, Woodstock Rd, Oxford, Oxfordshire OX2
6GG, UK. E-mail: serene.siow@psy.ox.ac.edu

11                                        Abstract

12   Many studies in bilingualism have explored how form similarity between translation

13   equivalents (i.e. cognate- ness) impacts language processing. Cognates range from being

14   identical to sharing only a few key features. Although there exists some objective measures

15   of orthographic similarity between word-forms, there is still no representative equivalent for

16   phonological similarity. In infant bilingual research, a measure that reflects perceptual

17   similarity of spoken words as presented in infant-directed speech would be of great value. A

18   widely used metric for word similarity is Levenshtein distance (**???**). This measure

19   calculates word-form similarity by computing the smallest total number of substitutions,

20   insertions and deletions needed to change one word to the other. However, a challenge when

21   defining phonological similarity is that some phonemes (e.g. [k] and [g]) are perceptually

22   more similar than others (e.g. [k] and [o]). The standard Lev- enshtein distance calculation

23   treats each change equally. Thus it often overestimates phonological distance between

24   words. We are developing an adaptation of Levenshtein distance for phonological

25   representations, where edit operations are weighted according to the degree of phonemic

26   feature changes involved, type of change (e.g. vowel or consonant insertion) and position in

27   the word. To estimate the weights, we are collecting behavioural data in an auditory

28   translation elicitation task. Monolingual adult participants with no exposure to the target

29   language will be asked to guess and produce the translations of auditorily-presented words.

30   The only way for participants to guess the correct translation is by mapping the presented

31   unknown phonological form onto the phonology of its known translational equivalent.

32   Cross-linguistic phonological similarity will be operationalised as the probability that

33   participants correctly guess the corresponding trans- lation. The higher the probability, the

34   higher the perceptual similarity. Our weighted similarity metric will provide a useful tool

35   for identifying phonological similarity between words in infant-directed speech.

36        *Keywords:* keywords

37    Word count: X

A Phoneme-weighted Levenshtein Distance: A Measure of Cross-linguistic Phonological Similarity

# Methods

We report how we determined our sample size, all data exclusions (if any), all manipulations, and all measures in the study.

## Participants

## Material

## Procedure

## Data analysis

We used R (Version 3.6.3; R Core Team, 2020) for all our analyses.

# Results

# Discussion

# References

R Core Team. (2020). *R: A language and environment for statistical computing.* Vienna,

Austria: R Foundation for Statistical Computing. Retrieved from

https://www.R-project.org/