# Cross-linguistic Variation in Phonemic Decomposition

**TSUNG-YING CHEN**
DEPARTMENT OF FOREIGN LANGUAGES AND LITERATURE
NATIONAL TSING HUA UNIVERSITY, TAIWAN

**JAMES MYERS**
GRADUATE INSTITUTE OF LINGUISTICS
NATIONAL CHUNG CHENG UNIVERSITY, TAIWAN

# Our thanks to...

- The Ministry of Science and Technology, Taiwan: MOST-103-2410-H-194-119-MY3 MOST-106-2410-H-194-055-MY3

- Todd Bailey, James Kirby, and Anna Veres-Székely for sharing their experimental data

- Assistants: You-Chu Chang, Kuei-Yeh Chen, Pei-Shan Chen, Yi-Hsin Lin, Mei-Jun Liu, Hsiao-Yin Pan, Si-Qi Su

- Our many participants

# Overview

- **Typological variation in syllable complexity and phonemic decomposition**

- **Cross-linguistic test (I): Wordlikeness judgments in English, Mandarin, and Cantonese**

- **Cross-linguistic test (II): Picture naming latencies in seven languages**

- **Implications for cross-linguistic psycholinguistics**

# Overview

- **Typological variation in syllable complexity and phonemic decomposition**

- Cross-linguistic test (I): Wordlikeness judgments in English, Mandarin, and Cantonese

- Cross-linguistic test (II): Picture naming latencies in seven languages

- Implications for cross-linguistic psycholinguistics

# Syllable complexity & Cross-linguistic variation

- **Languages vary in possible syllable structures** (Haspelmath et al., 2005)

  Simple = max CV (e.g., Hawaiian; *Mele Kalikimaka!*)

  Moderately complex = max CCVC (e.g., Mandarin; [ljaŋ])

  Complex = beyond CCVC (e.g., English; [stɹɛŋθs])

- **Languages thus also vary in the number of lexical syllable types**

  English: 12,000 (e.g., Levelt et al., 1999)
  Mandarin: 1,300 (including tones; e.g., Myers, 2015)

# Syllable complexity & Phonemic decomposition

- **Hypothesis:**

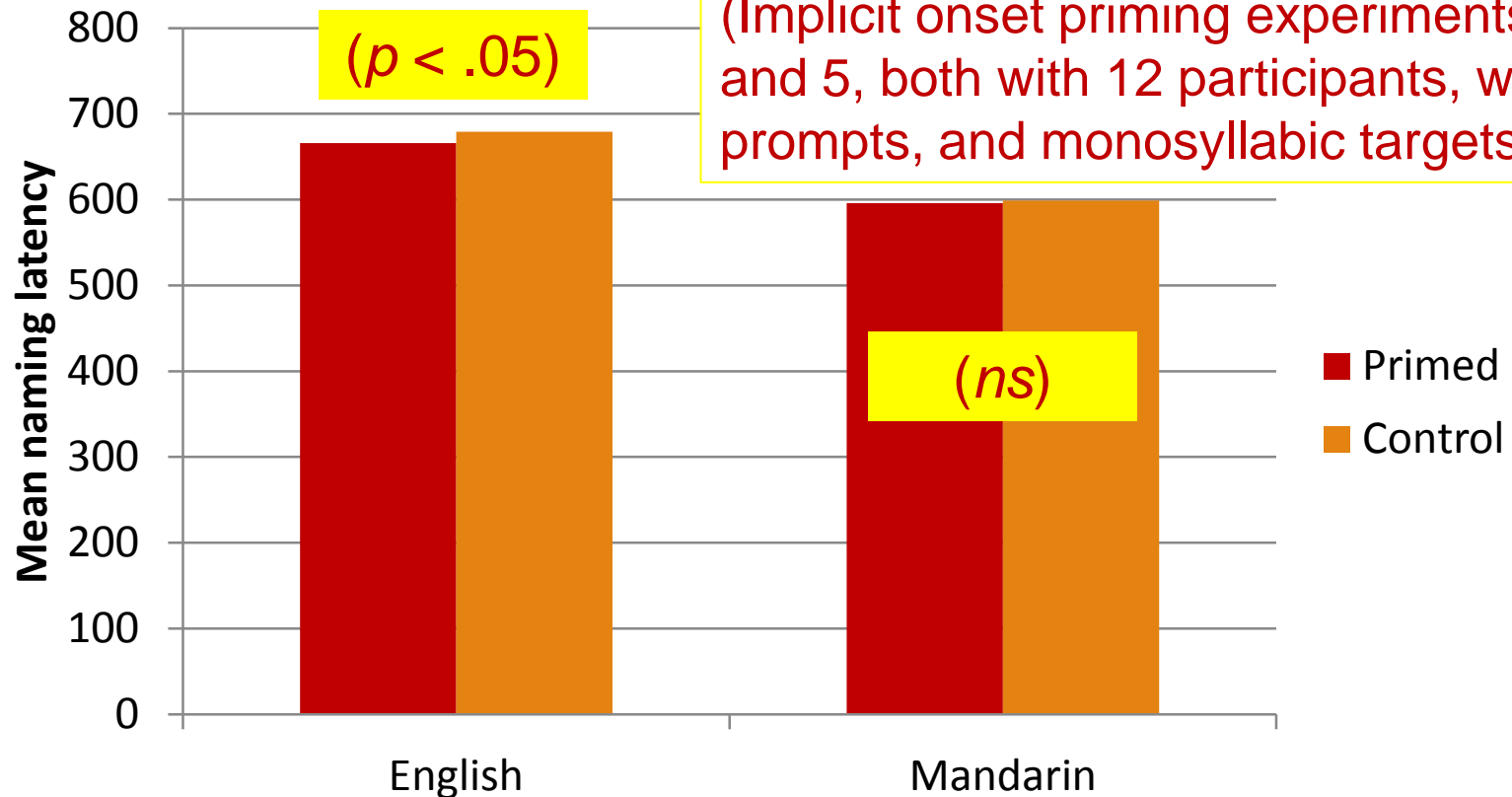  **Simpler/fewer syllables = Less phonemic decomposition**

- **Some suggestive evidence:**

**English** – Phoneme priming in production (O'Seaghdha et al., 2010) and phoneme > syllable advantage in perception (Norris & Cutler, 1998)

**Mandarin** – No phoneme priming in production (O'Seaghdha et al., 2010) and lexical syllable superiority effect in phoneme perception (Tseng et al. 1996)

# Phonemic decomposition in English vs. Mandarin

**O'Seaghdha et al. (2010)**



(Implicit onset priming experiments 3 and 5, both with 12 participants, written prompts, and monosyllabic targets)

# Phonemic decomposition: A simple diagnostic

- **Two lexical influences** (Luce & Large 2001)

**Phonotactic probability (PP)** – Probability of subsyllabic phoneme sequences, *depends on phonemic decomposition*

**Neighborhood density (ND)** – Overall similarity to lexical words, *does not depend on phonemic decomposition*

- **Predictions:**

  - **Effect sizes with strong phonemic decomposition: PP » ND** (e.g., English)

  - **Effect sizes with weak phonemic decomposition: ND » PP** (e.g., Mandarin)

# Overview

- Typological variation in syllable complexity and phonemic decomposition

- **Cross-linguistic test (I): Wordlikeness judgments in English, Mandarin, and Cantonese**

- Cross-linguistic test (II): Picture naming latencies in seven languages

- Implications for cross-linguistic psycholinguistics

# Wordlikeness judgments: Reanalyzing three studies

- **Nonword acceptability:** e.g., *blick* vs. *bnick*

  - **Higher PP = Higher acceptability**

  - **Higher ND = Higher acceptability**

  (Can be deconfounded via regression; Bailey & Hahn, 2001)

- **Test languages**

  **English:** Complex syllables

  **Mandarin:** Moderately complex

  **Cantonese:** Moderately complex

- **Predictions**
English (**PP » ND**), Mandarin and Cantonese (**ND » PP**)

# Wordlikeness judgments: Study procedures

- **English** (Bailey & Hahn, 2001, Exp 2)

  **24 participants**, **259 spoken monosyllabic nonwords**

  Nine-point Likert scale (1 = very atypical, 9 = very typical)

- **Mandarin** (Myers, 2015)

  **110 participants**, **3274 monosyllabic nonwords** written in Zhuyin Fuhao (Taiwan's onset/rime-based "pinyin")

  Binary scale (0 = 'unlike Mandarin', 1 = 'like Mandarin')

- **Cantonese** (Kirby & Yu, 2007)

  **10 participants**, **270 spoken monosyllabic nonwords**

  Seven-point Likert scale (1 = very poor, 7 = very good)

# Wordlikeness judgments: Quantification & analysis

- **Definition of predictors**
  **PP** – Transition probability in bigrams
  **ND** – Number of lexical monosyllables differing in just one element (tone ignored in Myers, 2015, to simplify bigrams)

- **Making judgment scales uniform**
  By-item mean judgments already in 0-1 range (Mandarin acceptance rates) or after rescaling (English, Cantonese), and transformed via arcsine square root.
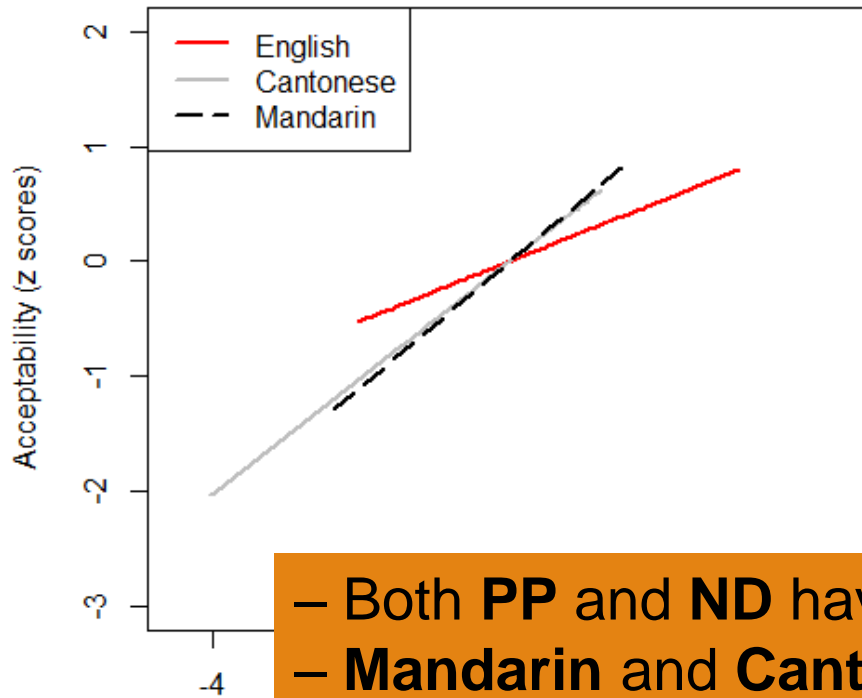
- **Standardizing**
  By-item ND, PP, judgments $z$-scored within each language

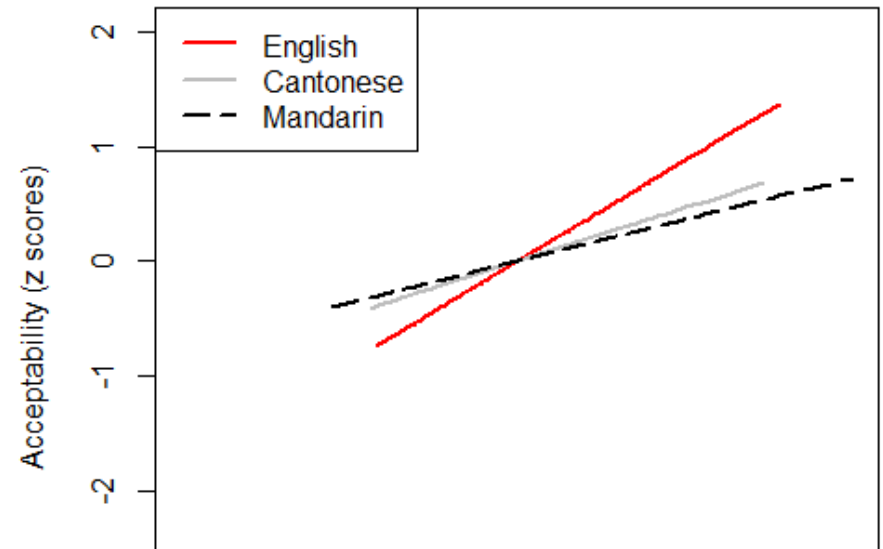- **Linear regression on by-item values**
  **Response ~ Language × (PP + ND)**

# Wordlikeness judgments: Results and discussion

**Neighborhood effects on acceptability**

**Phonotactic effects on acceptability**



- – Both **PP** and **ND** have **overall positive effects**
- – **Mandarin** and **Cantonese** behave **the same: ND » PP**
- – **English** has **weakest ND** and **strongest PP** effects

# Overview

- Typological variation in syllable complexity and phonemic decomposition

- Cross-linguistic test (I): Wordlikeness judgments in English, Mandarin, and Cantonese

- **Cross-linguistic test (II): Picture naming latencies in seven languages**

- Implications for cross-linguistic psycholinguistics

# Picture naming latencies:
## Seven test languages

▪ **Picture naming in seven languages** (Bates et al., 2003)

|  | **Syllable** | **OrthUnit** | **OrthDepth** |
|---|---|---|---|
| **Bulgarian** | Complex | Phoneme | Shallow |
| **English** | Complex | Phoneme | Mid |
| **German** | Complex | Phoneme | Mid |
| **Hungarian** | Complex | Phoneme | Shallow |
| **Italian** | ModComplex | Phoneme | Shallow |
| **Mandarin** | ModComplex | Syllable | Deep |
| **Spanish** | ModComplex | Phoneme | Shallow |

▪ **520 pictures, 30 participants for German, 50 participants for each of the other six languages.**

# Picture naming latencies: Quantifying variables

- **ND and PP were recalculated from free electronic dictionaries**
English (Lenzo, 2014), Mandarin (Denisowski et al., 2016), Spanish (Cuetos et al., 2011), the rest (Deri & Knight, 2016)

- **PP = Mean transition probability in bigrams**
(tone ignored in Mandarin)

- **(Inverse) ND (neighborhood sparsity) = PLD20**
(Yarkoni et al., 2008) Mean phonological Levenshtein (edit) distance from the twenty nearest lexical neighbors
(more effective measure for polysyllabic words)

# Picture naming latencies: Expected patterns

- **Different effects of phonotactics and neighbors on picture naming, depending on syllable types**


- **Higher PP = Stronger prelexical preparation**
  → Faster responses
  (Bulgarian, English, German, Hungarian) »
  <div align="right">(Italian, Mandarin, Spanish)</div>

- **Higher PLD20 (inverse ND) = Weaker postlexical activation**
  → Slower responses
  (Italian, Mandarin, Spanish) »
  <div align="right">(Bulgarian, English, German, Hungarian)</div>
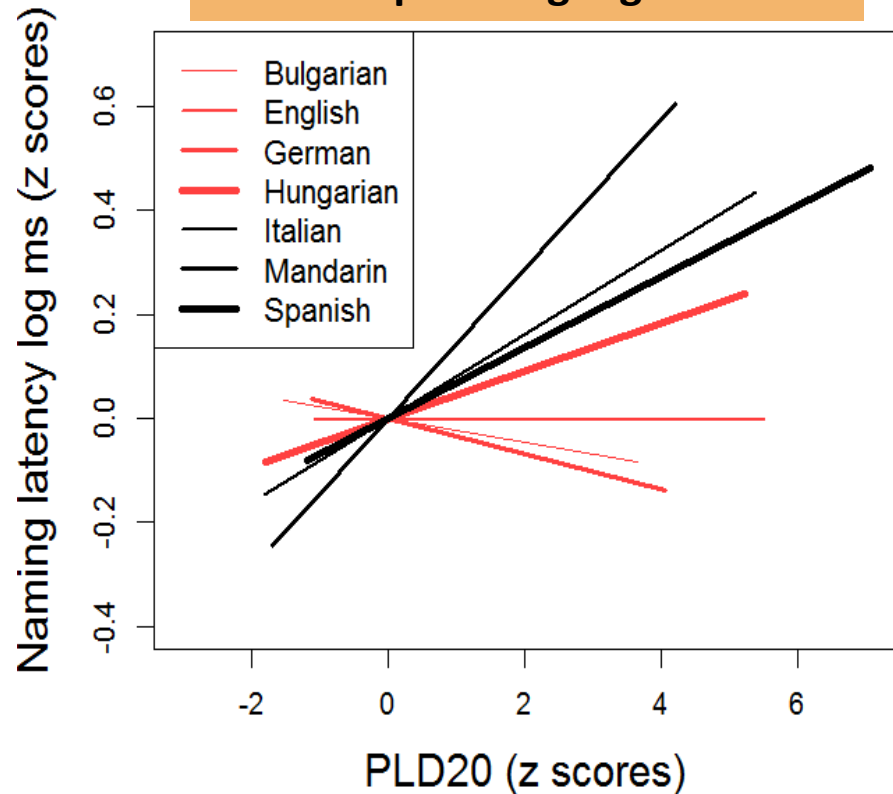
# Picture naming latencies: Statistical analysis

- **Linear mixed-effects regression**

  – **Dependent variable** – Reaction time (log-transformed)

  – **Independent variables** – Inverse ND (PLD20), PP, eight nuisance variables (e.g., lexical frequency), and their interaction with syllable complexity

  – **Random intercepts for pictures and languages**

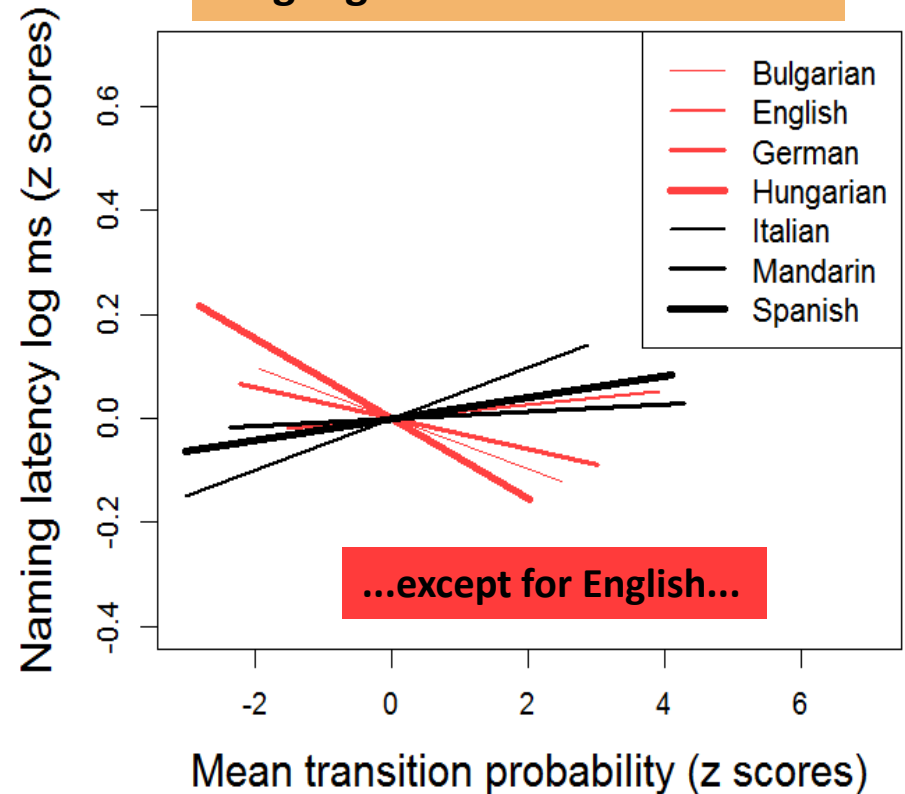  – **All variables were *z*-scored within language**

**Response ~ SylComplex x (Nuisances + PP + InvND)**

# Picture naming latencies:
## Results and discussion

**Stronger (inverse) ND effect for ModComplex languages**

**Stronger PP effect for Complex languages**



...except for English...

# Overview

- Typological variation in syllable complexity and phonemic decomposition

- Cross-linguistic test (I): Wordlikeness judgments in English, Mandarin, and Cantonese

- Cross-linguistic test (II): Picture naming latencies in seven languages

- **Implications for cross-linguistic psycholinguistics**

# Cross-linguistic psycholinguistics: Dealing with confounds

▪**Our databases are still too small:**
  – **Syllable complexity vs. inventory vs. orthography**
    Mandarin differs from Spanish and Italian in many ways
  – **Microvariation?**
    Are Mandarin and Cantonese really processed the same?

▪**Expanding the typological survey**
  – **Existing databases to exploit**
    Lexical decision latencies in English, Dutch, French, Malay…
  – **Collect our own wordlikeness judgments**
    Hakka and Southern Min (no orthographic influence?)
    Japanese (moderately complex, but different orthography)
    ... and as many other languages as we can manage...

# Cross-linguistic psycholinguistics: Making it feasible

- **Avoiding task-related confounds**
  – Different scales may be OK: binary vs. Likert scale
  – But task matters: wordlikeness vs. picture naming

- **Methodological consistency is thus crucial**

- **Yet no single team can test a sufficient number and variety of languages for a proper regression**

**Let the internet help:**
Web-based experimentation + Web-based data sharing

# Worldlikeness:
A Web application for typological psycholinguistics

- **https://Worldlikeness.org** (Chen & Myers 2017; Myers 2016)

# Worldlikeness: Overall architecture

**Experimenters** —— Create online experiments ——→ **Worldlikeness server** & typological database

Distribute

Web experiment ads
**Facebook, Twitter, etc.**

Invite

**Participants**

Participate in lab or online

Rewarded with a result report

Download data sets shared by experimenters and participants

**Researchers**

# Worldlikeness:
## Special features

- **Limited parameters to increase consistency**
  – Focused on wordlikeness

- **Privacy protections to encourage participation**
  – Fully anonymous
  – Full control of data authorization

- **Yet also facilitates and encourages data sharing**
  – Share more, do more
  – Most-open authorization option selected by default

- **Rapid data collection via Web crowdsourcing**
  – 16,000 judgments from 160 participants collected via Facebook in less than two weeks (Chen & Myers, in prep.)

# Thank you!

**TSUNG-YING CHEN**
chen.ty@mx.nthu.edu.tw

**JAMES MYERS**
Lngmyers@ccu.edu.tw

# References (1/4)

Bailey, T. M., & Hahn, U. (2001). Determinants of wordlikeness: Phonotactics or lexical neighborhoods? *Journal of Memory & Language, 44*, 569-591.

Bates, E., et al. (2003). Timed picture naming in seven languages. *Psychonomic Bulletin & Review, 10*(2), 344-380.

Chen, T.-Y., & Myers, J. (2017). Worldlikeness: A Web-based tool for typological psycholinguistic research. *University of Pennsylvania Working Papers in Linguistics, 23*(1), Article 4, 20-30.

Chen, T.-Y., & Myers, J. (In prep.). Worldlikeness: A Web application for typological psycholinguistic. Ms., National Tsing Hua University and National Chung Cheng University.

# References (2/4)

Cuetos, F., Glez-Nosti, M., Barbón, A., & Brysbaert, M. (2011). SUBTLEX-ESP: Spanish word frequencies based on film subtitles. *Psicológica, 32*, 133-143.

Denisowski, P. A. et al. (2016). CC-CEDICT (version cedict_ts.u8). https://www.mdbg.net/chindict/chindict.php?page=cedict. Retrieved November 11, 2016.

Deri, A., & Knight, K. (2016). Grapheme-to-phoneme models for (almost) any language. *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, 399-408.

Haspelmath, M., Dryer, M.S., Gil, D., & Comrie, B. (Eds.) (2005). *The world atlas of language structure*. Oxford University Press.

Kirby, J. P., & Yu, A. C. L. (2007). Lexical and phonotactic effects on wordlikeness judgments in Cantonese. *Proceedings of the International Congress of Phonetic Sciences, 16*, 1389-1392.

# References (3/4)

Lenzo, K. et al. (2014). Carnegie Mellon University Pronouncing Dictionary (version 0.7b). http://www.speech.cs.cmu.edu/cgi-bin/cmudict. Retrieved November 14, 2016.

Levelt, W. J., Roelofs, A., & Meyer, A. S. (1999). A theory of lexical access in speech production. *Behavioral and Brain Sciences, 22*(01), 1-38.

Luce, P. A., & Large, N. R. (2001). Phonotactics, density, and entropy in spoken word recognition. *Language & Cognitive Processes, 16*(5/6), 565-581.

Myers, J. (2015). Markedness and lexical typicality in Mandarin acceptability judgments. *Language & Linguistics, 16*(6), 791-818.

Myers, J. (2016). Meta-megastudies. *The Mental Lexicon, 11*(3), 329-349.

# References (4/4)

Norris, D., & Cutler, A. (1988). The relative accessibility of phonemes and syllables. *Perception & Psychophysics, 43*(6), 541-550.

O'Seaghdha, P. G., Chen, J.-Y., & Chen, T.-M. (2010). Proximate units in word production: Phonological encoding begins with syllables in Mandarin Chinese but with segments in English. *Cognition, 115*, 282-302.

Tseng, C.-H., Huang, K.-Y., & Jeng, J.-Y. (1996). The role of the syllable in perceiving spoken Chinese. *Proceedings of the National Science Council, Part C: Humanities and Social Sciences, 6*(1), 71-86.

Yarkoni, T., Balota, D., & Yap, M. (2008). Moving beyond Coltheart's N: A new measure of orthographic similarity. *Psychonomic Bulletin & Review, 15*(5), 971-979.

# Appendix:
## Bates et al. (2003) nuisance variables

- Lexical frequency

- Picture quality (via pretest judgments)

- Fricative onset

- Word length in phonemes

- Number of alternative names

- Number of names shared across pictures

- Naming consistency across participants

- Naming consistency within each participant