# Vector Semantics: Lecture 3

András Kornai
SZTAKI Computer Science Research Institute

20 March 2024

# Main takeaways from studying frequency

- Zipf and Herdan laws work well Kornai, 1999a; Kornai, 2002
- It's all about information Brown et al., 1992
- Entropy measures maximum average compressibility over the wire (Jelinek, 1997) ML Ch 7.1
- Kolmogorov complexity can achive more compression, but only semi-computable (Li and Vitányi, 1997; Vitanyi and Li, 2000) ML Ch 7.3
- Minimum Description Length philosophy used in linguistics (Borbély and Kornai, 2019)
- With a twist: amortizing the universal component (Kornai, Zséder, and Recski, 2013)
- Engineering takeaway (1): do the frequent things first
- Engineering takeaway (2): OOV is a persistent problem

# WORDS

- Minimum free forms (Bloomfield, 1926)
- Phonological words: units between pauses
- (Orthographical words)
- Lexemes (also called lexical entries) can be MWEs like *as is*
- Subentries (Kornai, 2023)
- Subsubentries csinál

# THE STRUCTURE OF THE LEXEME

- Pronunciation (phonology database key)
- Part of speech (syntax db key)
- Definition (semantics db key)
- Bunch of ancillary info: etymology, variants, style, topic, frequency, hyphenation ...
- Headword usually derived via orthography
- Easily extended to bilingual/multilingual
- But what to do with technical vocabulary? Millions of "words" for chemical compounds, animal species, ...

# COVERAGE

- Ideally, we'd want the dictionary frequency-ordered
- But high coverage remains elusive, OOV is a big problem
- Common vocabulary often used in L2 instruction (Kornai, 2021)
- It is less trivial to define than 'most frequent' (Thorndike, 1921), corrected frequency
- Our interest is more with basic vocabulary (Ogden, 1944), Simple Wikipedia (Yasseri, Kornai, and Kertész, 2012)
- Everybody tries to build a basic list: https://concepticon.clld.org has 450+ sources

# LEXICON OR ENCYCLOPEDIA

- In many topics, technical vocabulary is key
- Proper names and named entities
- PER, LOC, ORG – hundreds of millions of entries in each category
- *hutch for sale, as is*

# HUTCH, AS IS

# General principles

- Universality
- Reductivity
- No encyclopedic knowledge
- Read VS Ch. 1.2

📄 Bloomfield, Leonard (1926). "A set of postulates for the science of language". In: *Language* 2, pp. 153–164.

📄 Borbély, Gábor and András Kornai (June 2019). "Sentence Length". In: *Proceedings of the 16th Meeting on the Mathematics of Language*. Toronto, Canada: Association for Computational Linguistics, pp. 114–125. URL: https://www.aclweb.org/anthology/W19-5710.

📄 Brown, P. et al. (1992). "An Estimate of an Upper Bound for the Entropy of English". In: *Computational Linguistics* **18**/1, pp. 31–40.

📄 Jelinek, Frederick (1997). *Statistical Methods for Speech Recognition*. MIT Press.

📄 Kornai, András (2002). "How many words are there?" In: *Glottometrics* 2.4, pp. 61–86.

📄 — (2021). "Vocabulary: Common or Basic?" In: *Frontiers in Psychology*. DOI: 10.3389/fpsyg.2021.730112. URL: https://www.frontiersin.org/articles/10.3389/fpsyg.2021.730112/full.

📄 — (2023). "Poliszémia politópokkal". In: *Általános Nyelvészeti Tanulmányok* 35. Ed. by Beáta Gyuris, pp. 311–326. ISSN: HU 0569-1338.

📄 — (1999a). "Zipf's law outside the middle range". In: *Proceedings of the Sixth Meeting on Mathematics of Language*. Ed. by J. Rogers. University of Central Florida, pp. 347–356.

📄 Kornai, András, Attila Zséder, and Gábor Recski (2013). "Structure Learning in Weighted Languages". In: *Proceedings of the 13th Meeting on the Mathematics of Language (MoL 13)*. Sofia, Bulgaria: Association for Computational Linguistics, pp. 72–82. URL: http://www.aclweb.org/anthology/W13-3008.

📄 Li, Ming and Paul Vitányi (1997). *An Introduction to Kolmogorov Complexity and Its Applications*. Springer.

📄 Ogden, C.K. (1944). *Basic English: a general introduction with rules and grammar*. K. Paul, Trench, Trubner.

📄 Thorndike, Edward L. (1921). *The teacher's word book*. New York Teachers College, Columbia University.

📄 Vitanyi, Paul M. B. and Ming Li (2000). "Minimum description length induction, Bayesianism, and Kolmogorov complexity". In: *IEEE Transactions on Information Theory* 46.2, pp. 446–464.

📄 Yasseri, Taha, András Kornai, and János Kertész (2012). "A practical approach to language complexity: a Wikipedia case study". In: *PLoS ONE* 7.11. DOI: e48386.doi:10.1371/journal.pone.0048386.