

# VECTOR SEMANTICS: LECTURE 1

András Kornai  
SZTAKI Computer Science Research Institute

13 February 2024

# MAJOR THEMES

- What is semantics?
- Vectors, matrixes
- Ordinary language philosophy V2.0: aims at ordinary (non-technical) language, but heavy use of technical tools
- From linguistics (morphology, lexicography)
- From math (algebra, information theory)
- Logic content: default logic, probabilities, negation, naive implication
- Discussion of big issues: can LLMs have consciousness? are they intelligent? are they safe?

# METAINFO

- **Time/Place** Tuesday 12:15-1:45 pm CET/ELTE Logic Dept
- **Course webpage, slack**  
<https://nessie.ilab.sztaki.hu/~kornai/2024/VectorSemantics/>  
<https://vectorsemantics.slack.com>
- **Prerequisites** None
- **Requirements** regular attendance (max. 3 classes can be skipped), good knowledge and understanding of the main text, and active participation in discussions. Participants need to write 2 short essays (3-3 pages) responsive to the assigned tasks. By the end of the semester they need to submit another essay (minimum 4 pages) which relies substantively on at least one chapter of the book and one of the assigned readings. The readings will be assigned in consultation with the professor, and will be in English. Over the semester a total of three essays need to be submitted. Grades are based on classroom activity (20%) and the results of these essays (20%+20%+40%).



spell 927250 spelling 666868 spells 375175 spelled 237181 spellings  
51680 spelt 36573 spellbound 17346 spellbinding 14765 **spelen 6823**  
speller 6687 spellchecker 6539 spellcheck 6059 **spel 5062** spellers  
4439 **spelunking 4089** spellcasting 4058 **spelung 3722** spellbook 3550  
spellcaster 3209 spellbinder 3125 spell's 3030 spellcasters 2970  
**speleothems 1871 speleology 1455 spelunkers 1345** spellchecking  
1313 spellcraft 1126 **speleological 1122** spelter 1043 spellcheckers  
990 **spell&quot 951 spelunker 930** spellwork 766 **speleothem 754**  
spelljamming 683 spellchecked 652 **spellen 643 speleologists 641**  
spellcast 601 **spells&quot 598 speleo 558 spellin 550 spelar 548** spell'  
486 **spela 475 spelvin 432 spelspiel 378 speler 373** spellbind 359  
**spelende 355 spelta 329 spelling&quot 327 spell&gt 325** spellmasters  
322 **spelunk 315** spellman 309 spelthorne 291 **spelletjes 278 spellyou**  
**264 spellex 252** spelljammer 249 **speleologist 248** spellserver 237  
spells' 225 spellchk 219 spellworking 217 spellbindingly 213 spelare  
209 speltoides 203 **spellin' 198** spelling's 195 spelling' 195 spellout  
188 speld 185 spello 183 spellbinders 182 spellmaker 180 spellchips  
176 **spelade 175** spellpoints 172 **speleogenesis 172 spelling 169 spelld**

# SPELEOLOGY

- *speleum* 'cave' + *ology* 'science of' = *speleology* 'science of caves'
- Yes, but what is the '+' and what is the '=' here?
- We will not look at the etymology because the ordinary speaker does not have access to it
- But we will look at the frequencies
- We will do a fair bit of linguistic analysis
- Starting with a crash course in morphology

# WHAT KIND OF SCIENCE IS SPELEOLOGY?

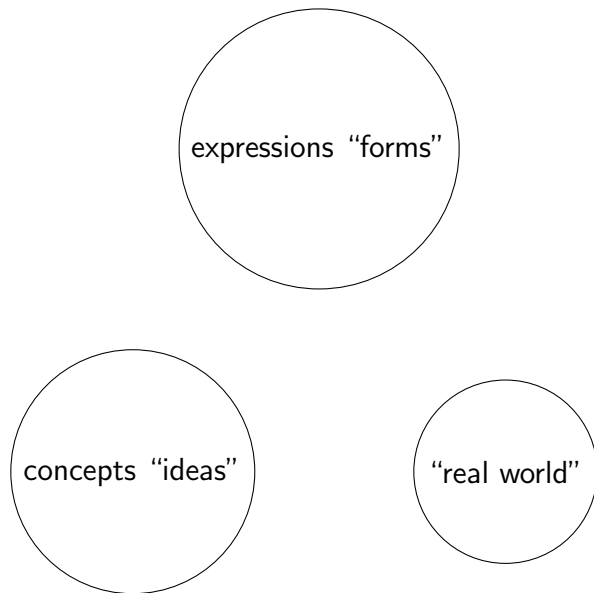
- Obviously, there are caves, and we deeply care about them
- But their formation is a matter of geology
- Their flora/fauna (very interesting!) is a matter of biology
- Their population is a matter of archeology
- So we don't have a unified science of speleology, all we have are theories/principles from other, more coherent theories that we try to apply/extend to caves
- Semantics is not any different

# WHAT KIND OF SCIENCE IS SEMANTICS?

- Obviously people talk to each other, and can understand each other well enough to cooperate
- Or go to war when the communication breaks down. The stakes are high!
- We will throw everything at the problem: logic, statistics, math, computer science, linguistics, semiotics, cognitive science, philosophy . . .
- And see what sticks – whatever works, works, the rest goes on the back burner
- The approach taken here is *irenic* and *syncretic*
- It will also be *bottom up* rather than *top down*



# THE OVERALL MODEL

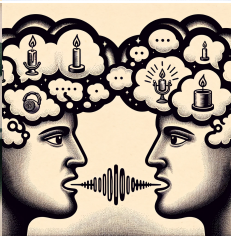
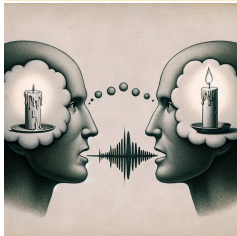
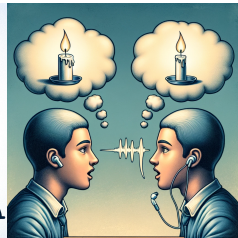
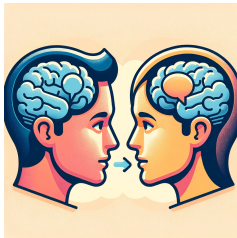


# COMMUNICATION



# THE PROMPT

Please create a sequence of three images: the first one should show a woman thinking about a candle. The second should show the woman saying the word "CANDLE" to the man by means of enclosing the text CANDLE inside a text bubble coming from her mouth. The third image should show the man thinking about a different candle.



irenic or ei·renic \(')i:renik, -rēn-\ also ireni·cal \-nəkəl\  
adj [Gk *eirēnikos*, fr. *eirēnē* peace (prob. of non-IE origin)  
+ *-ikos* -ic, -ical] : conducive to or operating toward peace,  
moderation, harmony, and conciliation and away from con-  
tention and partisanship esp. among disputants (<~ measures>  
<~ without being namby-pamby —*Chicago Theol. Seminary  
Register*) <the viewpoint is ~ and the author seeks to show  
the best features of each religion and church in turn —N.K.  
Burger> **syn** see PACIFIC

ireni·cal·ly \-nək(ə)lē\  
**adv** : in an irenic manner : in a way

**syn·cret·ic** \(')si|n|'kred·ik, sə|n|'k-, |ŋ|\  
-ic] **1** : characterized or brought about by syncretism : aiming  
at or making for syncretism : SYNCRETISTIC (<~ religious sect>  
**2** : having absorbed the functions of one or more other gram-  
matical cases <the Latin ablative is a ~ case>  
**syn·cre·tion** \sən'krēshən, sən'k-\  
instance of syncretism : act of syncretizing  
**syn·cre·tism** \'siŋkrə,tizəm, 'sɪŋk-\  
Gk *synkrētismos* federation of Cretan cities, fr. *synkrētizein* to  
unite against a common enemy] **1** : the reconciliation or union  
of conflicting (as religious) beliefs or an effort intending  
such; *specif* : a movement of a Lutheran party in the 17th

irenic	ety	C19: from Greek eir*\_enikos, from eir*\_
eirenic	alt	
	head	irenic
	or	eirenical
	syl	ei:ren+ic
	pron	<I1rEnik>, <-1ren->
	pos	adj.
irenic	0.	
	or	irenical
	syl	i:ren+ic
	pron	<I1rEnik>, <-1ren->
	pos	adj.
	qual	Chiefly theol.
	def	tending to conciliate or promote peace.
irenically	sub	
	head	irenic
	or	eirenically
	syl	i:1ren:i+cal+lv

spell 927250 spelling 666868 spells 375175 spelled 237181 spellings  
51680 spelt 36573 spellbound 17346 spellbinding 14765 **spelen 6823**  
speller 6687 spellchecker 6539 spellcheck 6059 **spel 5062** spellers  
4439 **spelunking 4089** spellcasting 4058 **spelung 3722** spellbook 3550  
spellcaster 3209 spellbinder 3125 spell's 3030 spellcasters 2970  
**speleothems 1871 speleology 1455 spelunkers 1345** spellchecking  
1313 spellcraft 1126 **speleological 1122** spelter 1043 spellcheckers  
990 **spell** 951 **spelunker 930** spellwork 766 **speleothem 754**  
spelljamming 683 spellchecked 652 **spellen 643 speleologists 641**  
spellcast 601 **spells** 598 **speleo 558 spellin 550 spelar 548** spell'  
486 **spela 475 spelvin 432 spelspiel 378 speler 373** spellbind 359  
**spelende 355 spelta 329 spelling** 327 **spell** 325 spellmasters  
322 **spelunk 315** spellman 309 spelthorne 291 **spelletjes 278 spellyou**  
**264 spellex 252** spelljammer 249 **speleologist 248** spellsserver 237  
spells' 225 spellchk 219 spellworking 217 spellbindingly 213 spelare  
209 speltoides 203 **spellin' 198** spelling's 195 spelling' 195 spellout  
188 speld 185 spello 183 spellbinders 182 spellmaker 180 spellchips  
176 **spelade 175** spellpoints 172 **speleogenesis 172 spelling 169 spelld**

# MORPHOLOGY

- Morphemes: minimal signs
- Words: minimal free forms
- Roots, stems, fully formed words
- Lexemes
- Clitics, phonological v written words
- Word meaning



# INFORMATION

- Measured in **bits**
- Can be computed by Shannon's formula  $H = -\sum_i p_i \log_2(p_i)$
- Property of distributions not individual items
- Counts the average number of the best Twenty Questions-style questions it takes to identify a particular item
- If something contains 21 bits of information, there is *no* clever girl who can get to it in 20 questions

# INFORMATION CONTENT OF SENTENCES

Number of possible binary trees over  $n$  words is

$$C_n \sim 4^n / n^{1.5} \sqrt{\pi}$$

or **< 2 bits per word** (see also Aronoff, 1985)

The word entropy  $H$  of a language with Zipf constant  $B$  is given by

$$H \approx H_k + \frac{1 - P_k}{\log(2)} (B/(B - 1) - \log(B - 1) + \log(k) - \log(1 - P_k))$$

which yields **12.67 bits** for English, **15.41 bits** for Hungarian.

Other definitions of logical structure than by parse tree are possible but they do not alter the picture significantly: logical structure accounts for **at most 15% of the information**.

# THE LINGUISTIC SIGN

- **The classic model** (Ferdinand de Saussure 1916)  
A two-part structure composed of *form* and *meaning*
- **The modern model** (Kracht, 2003)  
A three-part structure composed of *form*, *category* and *meaning*
- 4lang

```
mark jel nota znak shirushi 印 biao1j14 标记 1182 u N sign, visible %
mark_ jel01 nota znak shirusu 標寸 biao1zh14 标志 3331 u V =agt[sign], =pat[meaning], represent %
```

# COMING ATTRACTIONS/ESSAY TOPICS

- Word frequency distributions, Zipf's Law
- $n$ -dimensional space
- Hypernode graphs
- The syntax of  $\lambda$ -calculus formulas
- Basic units of meaning: are there any?
- Compositionality, semicompositionality
- Homonymy/polysemy



Aronoff, Mark (1985). “Orthography and Linguistic Theory: The Syntactic Basis of Masoretic Hebrew Punctuation”. In: *Language* 61.1, pp. 28–72.



Kracht, Marcus (2003). *The Mathematics of Language*. Berlin: Mouton de Gruyter.