# The Measurement of Mutual Intelligibility between West-Slavic Languages

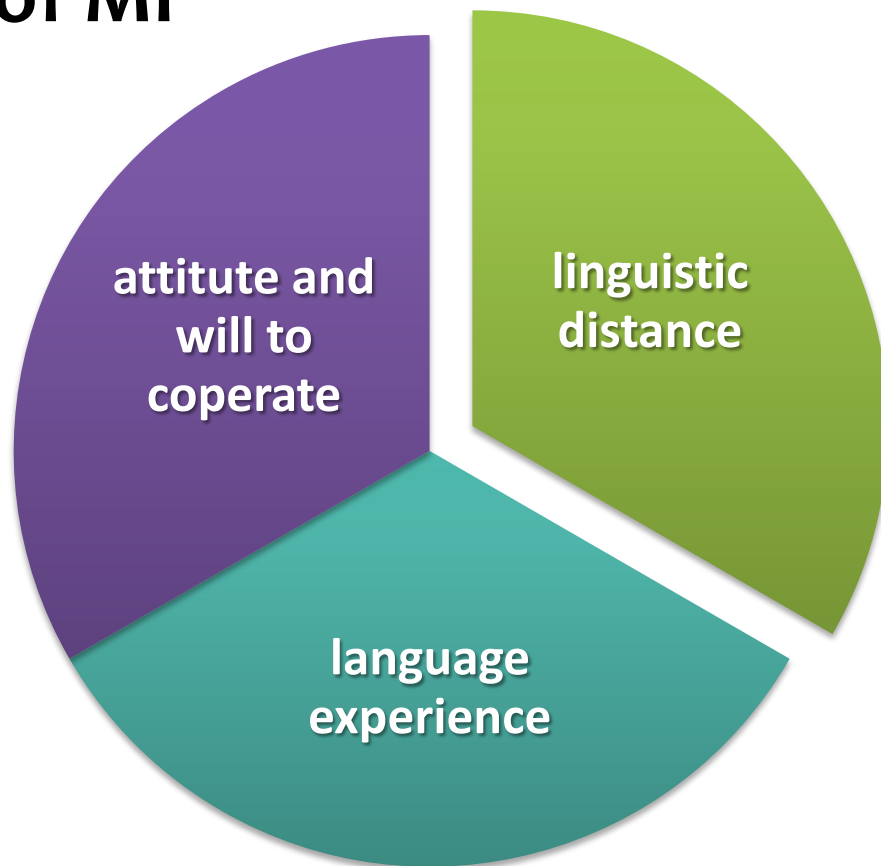Perlová Voda, September 2018

# Mutual Intelligibility (MI) → Semicommunication



HAUGEN, E. (1966). Semicommunication: The Language Gap in Scandinavia.

# MI languages & factors of MI

- Danish – Norwegian – Swedish

- Afrikaans – Frisian – Dutch

- Faroese – Icelandic

- Croatian – Serbian – Slovenian

- Belarusian – Russian – Ukrainian

- Italian – Spanish – Portuguese

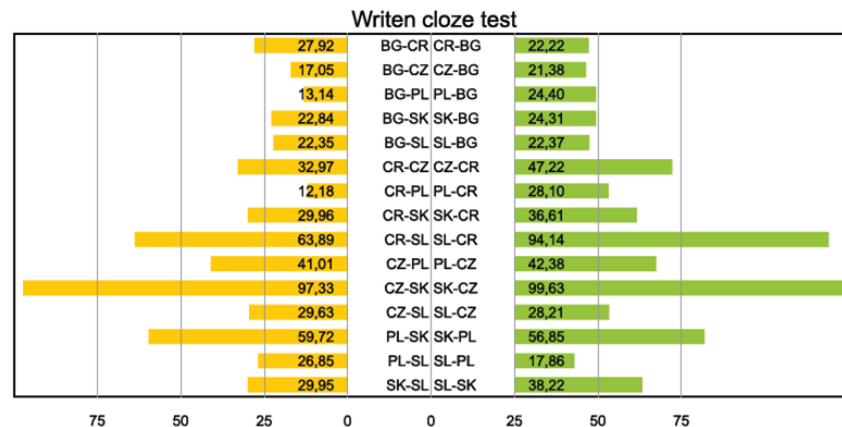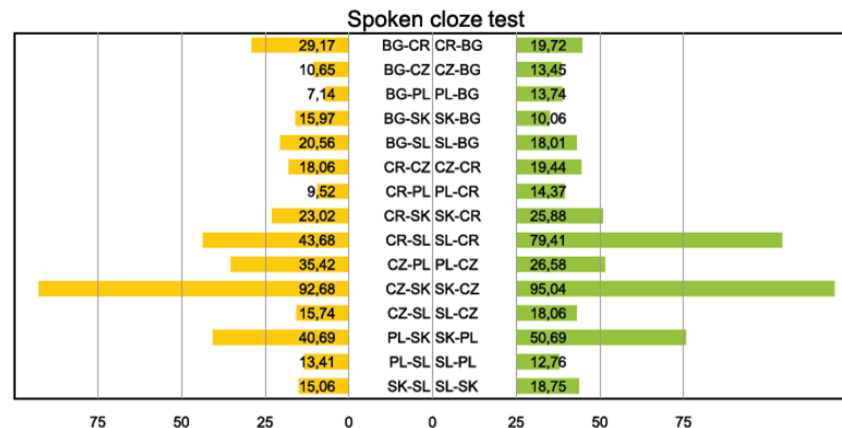- Turkish – Azerbaijani

- …

# Research objectives

- Overall mutual intelligibility between West-Slavic languages
- Asymmetry of mutual intelligibility between West-Slavic languages
- Mutual intelligibility of content and function words
- Mutual intelligibility of various styles of material (stylistics)

- Differences between spoken and written forms of West-Slavic languages in all above mentioned areas
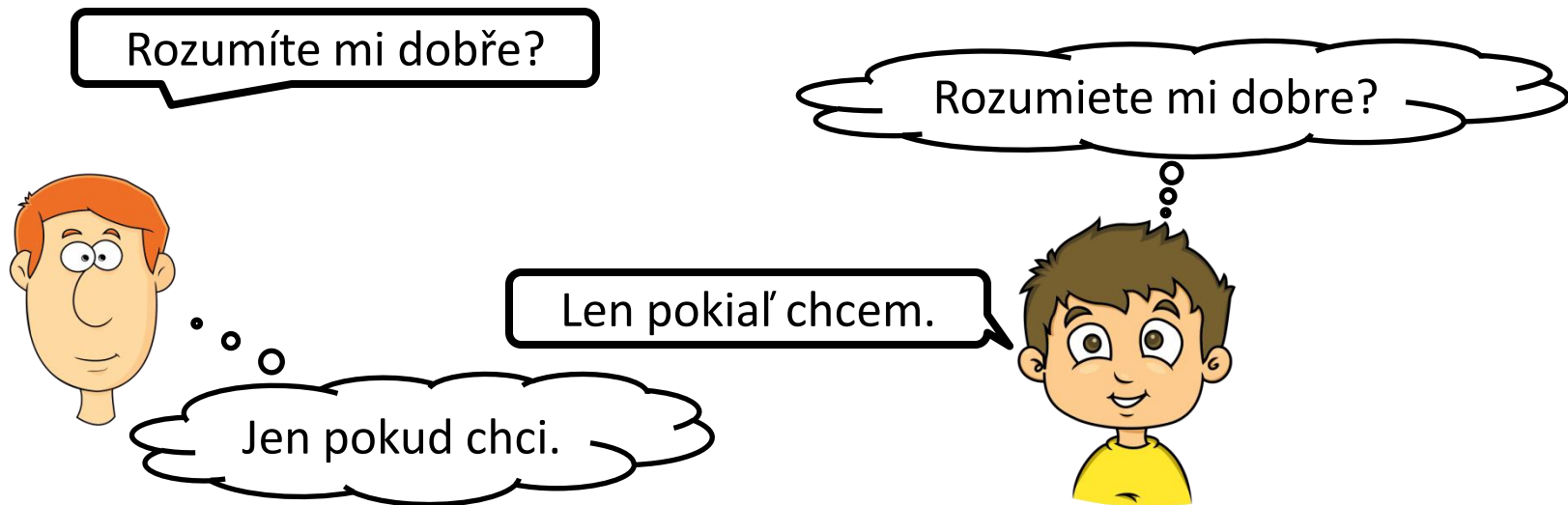
# Related works

- Dialectometry:
  - (2007) MOBERG J., GOOSKENS Ch., NERBONNE J., VAILLETTE N.

- Sociolinguistics research:
  - (2016) GOLUBOVIĆ, J.
  - (2012), (2009), (2000), (1987)



Spoken cloze test

| | | |
|---|---|---|
| 29,17 | BG-CR CR-BG | 19,72 |
| 10,65 | BG-CZ CZ-BG | 13,45 |
| 7,14 | BG-PL PL-BG | 13,74 |
| 15,97 | BG-SK SK-BG | 10,06 |
| 20,56 | BG-SL SL-BG | 18,01 |
| 18,06 | CR-CZ CZ-CR | 19,44 |
| 9,52 | CR-PL PL-CR | 14,37 |
| 23,02 | CR-SK SK-CR | 25,88 |
| 43,68 | CR-SL SL-CR | 79,41 |
| 35,42 | CZ-PL PL-CZ | 26,58 |
| 92,68 | CZ-SK SK-CZ | 95,04 |
| 15,74 | CZ-SL SL-CZ | 18,06 |
| 40,69 | PL-SK SK-PL | 50,69 |
| 13,41 | PL-SL SL-PL | 12,76 |
| 15,06 | SK-SL SL-SK | 18,75 |

Writen cloze test

| | | |
|---|---|---|
| 27,92 | BG-CR CR-BG | 22,22 |
| 17,05 | BG-CZ CZ-BG | 21,38 |
| 13,14 | BG-PL PL-BG | 24,40 |
| 22,84 | BG-SK SK-BG | 24,31 |
| 22,35 | BG-SL SL-BG | 22,37 |
| 32,97 | CR-CZ CZ-CR | 47,22 |
| 12,18 | CR-PL PL-CR | 28,10 |
| 29,96 | CR-SK SK-CR | 36,61 |
| 63,89 | CR-SL SL-CR | 94,14 |
| 41,01 | CZ-PL PL-CZ | 42,38 |
| 97,33 | CZ-SK SK-CZ | 99,63 |
| 29,63 | CZ-SL SL-CZ | 28,21 |
| 59,72 | PL-SK SK-PL | 56,85 |
| 26,85 | PL-SL SL-PL | 17,86 |
| 29,95 | SK-SL SL-SK | 38,22 |

# Method

- Levenshtein distance & Conditional entropy
- Inspired by psycholinguistics idea about process of semicommunication

# Conditional entropy (CE)

- Quantifies the amount of information needed to get the X when Y is given
- Lower entropy = better mutual intelligibility (smaller linguistic distance)
- Allows asymmetrical results (from the definition of CE)

$$H(X \mid Y) = - \sum_{x \in X, y \in Y} p(x, y) \log_2 (p(x \mid y))$$

- X … native language,    x … native phoneme/grapheme
  Y … foreign language,   y … foreign phoneme/grapheme

# CE - example

$$H(X \mid Y) = - \sum_{x \in X, y \in Y} p(x, y) \log_2 (p(x \mid y))$$

| CS | r | ɔ | z | ʊ | m | iː | t | ɛ | | m | ɪ | | d | ɔ | b | ɽ | ɛ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SK | r | ɔ | z | ʊ | m | i̯ɛ | t | ɛ | | m | i | | d | ɔ | b | r | ɛ |
| p(CS\|SK) | .50 | 1 | 1 | 1 | .67 | 1 | 1 | .75 | | .67 | 1 | | 1 | 1 | 1 | .50 | .75 |
| p(SK\|CS) | 1 | 1 | 1 | .50 | 1 | 1 | .50 | 1 | | 1 | .50 | | 1 | 1 | 1 | 1 | 1 |

Asymmetries: r{ r / ɽ , ɛ{ ɛ / ɪ , ʊ{ i̯a / ʊ , ...

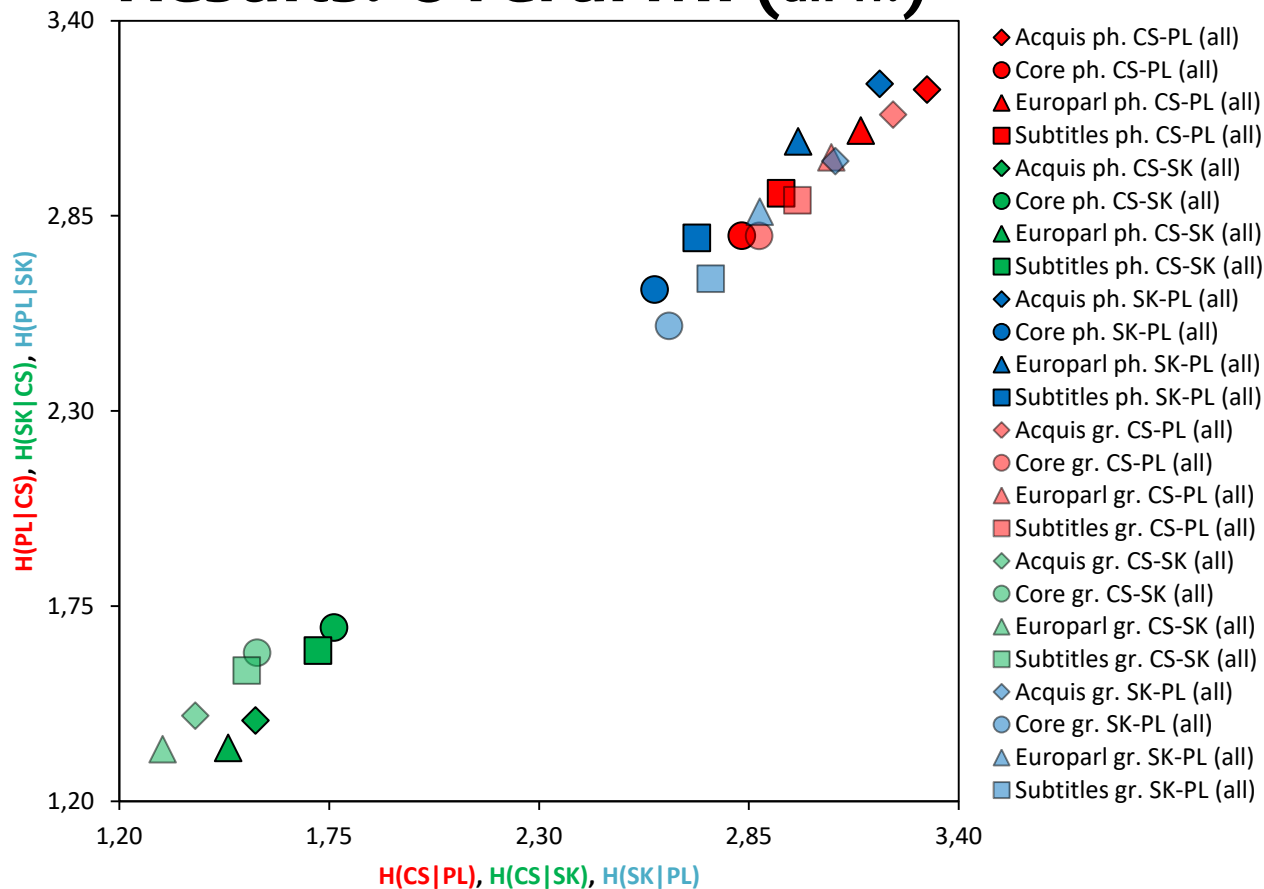| CS | j | ɛ | n | | p | ɔ | k | ʊ | t | | x | ts | ɪ | # |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SK | ɭ | ɛ | n | | p | ɔ | k | i̯a | ʎ | | x | ts | ɛ | m |
| p(CS\|SK) | 1 | .75 | 1 | | 1 | 1 | 1 | 1 | 1 | | 1 | 1 | .25 | .33 |
| p(SK\|CS) | 1 | 1 | 1 | | 1 | 1 | 1 | .50 | .50 | | 1 | 1 | .50 | 1 |

# Material

- corpora: **InterCorp v9 2016 (ČNK)**
- subcorpora: **Acquis, Europarl, Core, Subtitles**

- loaded from: **KonText v0.9.3**
- translations: **Treq v1.1**

- sample size: **2 000 most frequently used words**
- transcription: **IPA (semi-automatic)**

# Results: Overal MI (all w.)
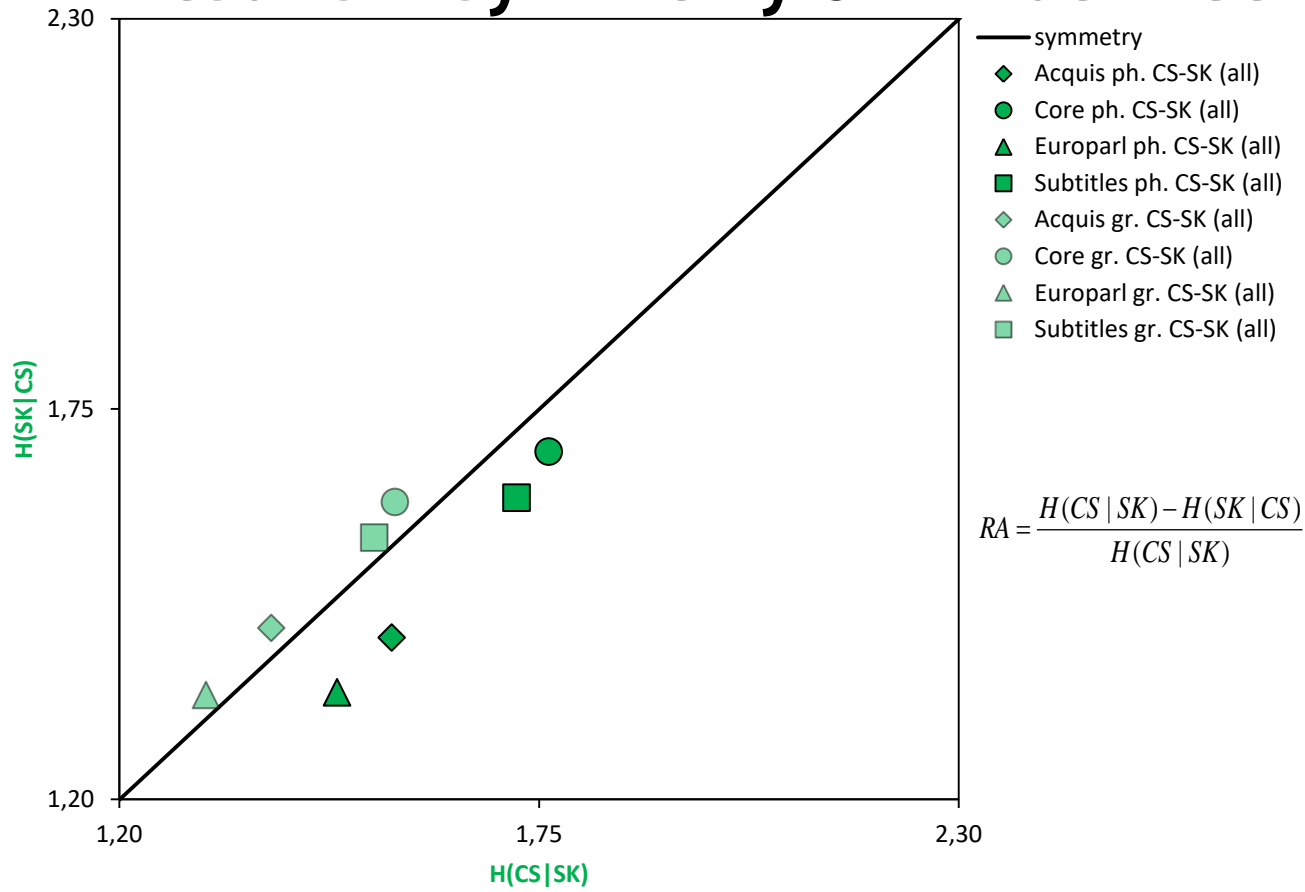


MI on phonetic layer ≈
MI on graphemic layer

***CS-SK* < SK-PL < CS-PL**

↳ Agree with socioling.
research

The most MI for:
*CS-SK* = Europarl, Acquis;
*CS-PL* = Core, Subtitles;
*SK-PL* = Core, Subtitles.

Subtitles ≈ middle of
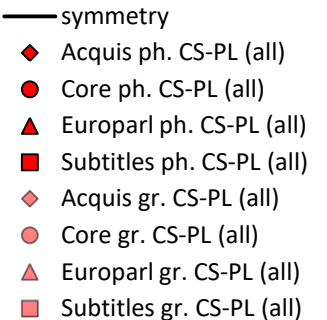groups

# Results: Asymmetry of MI between CS-SK (all w.)



symmetry
Acquis ph. CS-SK (all)
Core ph. CS-SK (all)
Europarl ph. CS-SK (all)
Subtitles ph. CS-SK (all)
Acquis gr. CS-SK (all)
Core gr. CS-SK (all)
Europarl gr. CS-SK (all)
Subtitles gr. CS-SK (all)

Phonetic layer:

**SK > CS** (RA = 0,068)

Graphemic layer:

**CS > SK** (RA = 0,029)

$$RA = \frac{H(CS \mid SK) - H(SK \mid CS)}{H(CS \mid SK)}$$

↳ Agree with socioling. research, except graph.

Same side for all subcorpora across layers

# Results: Asymmetry of MI between CS-PL (all w.)



Phonetic layer:

**PL > CS** (RA = 0,017)

Graphemic layer:

**PL > CS** (RA = 0,026)

↳ Agree with socioling. research, except phon.

Same side for all subcorpora across layers

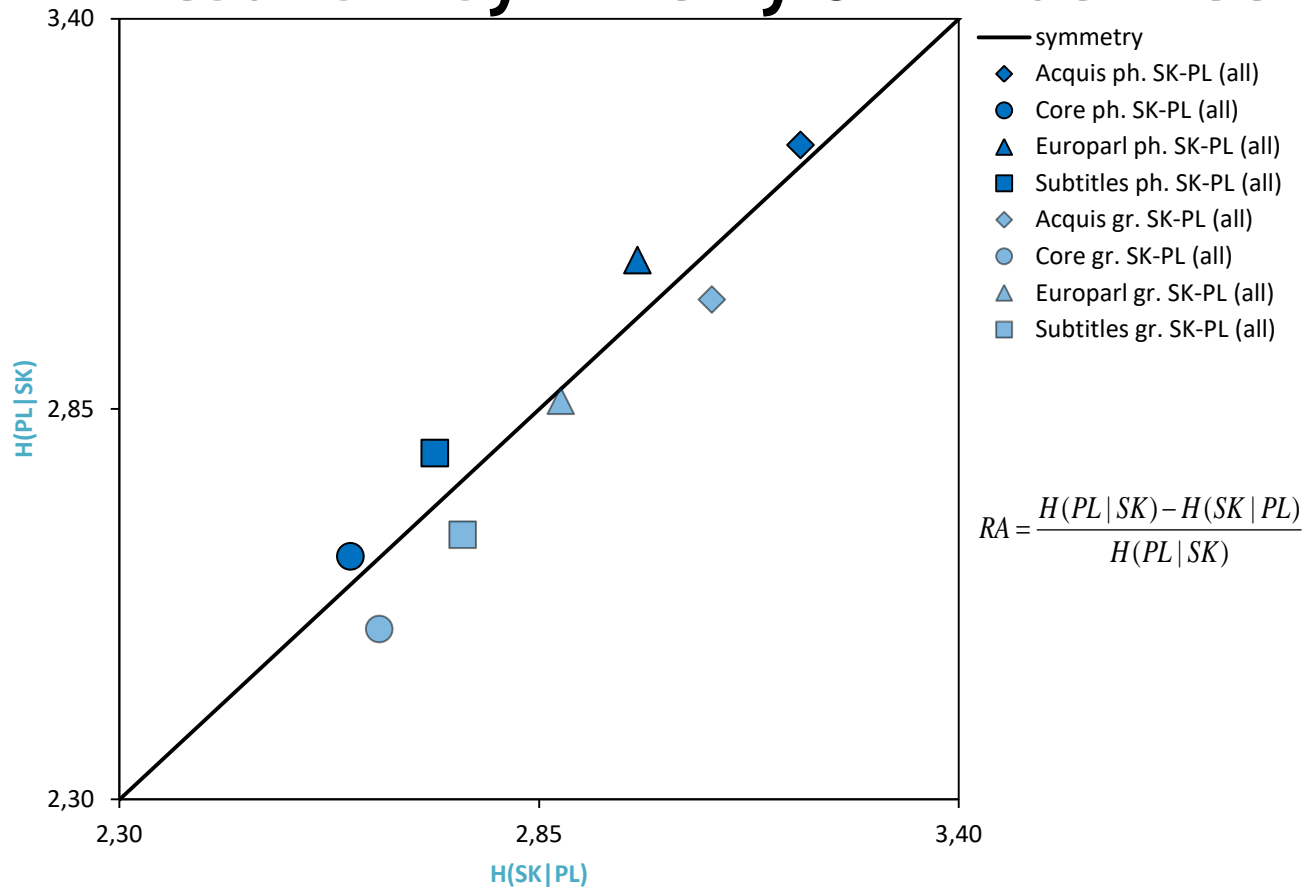$$RA = \frac{H(CS \mid PL) - H(PL \mid CS)}{H(CS \mid PL)}$$

Legend:
— symmetry
◆ Acquis ph. CS-PL (all)
● Core ph. CS-PL (all)
▲ Europarl ph. CS-PL (all)
■ Subtitles ph. CS-PL (all)
◆ Acquis gr. CS-PL (all)
● Core gr. CS-PL (all)
▲ Europarl gr. CS-PL (all)
■ Subtitles gr. CS-PL (all)

Axes:
H(PL|CS) (vertical), H(CS|PL) (horizontal), range 2,30 to 3,40

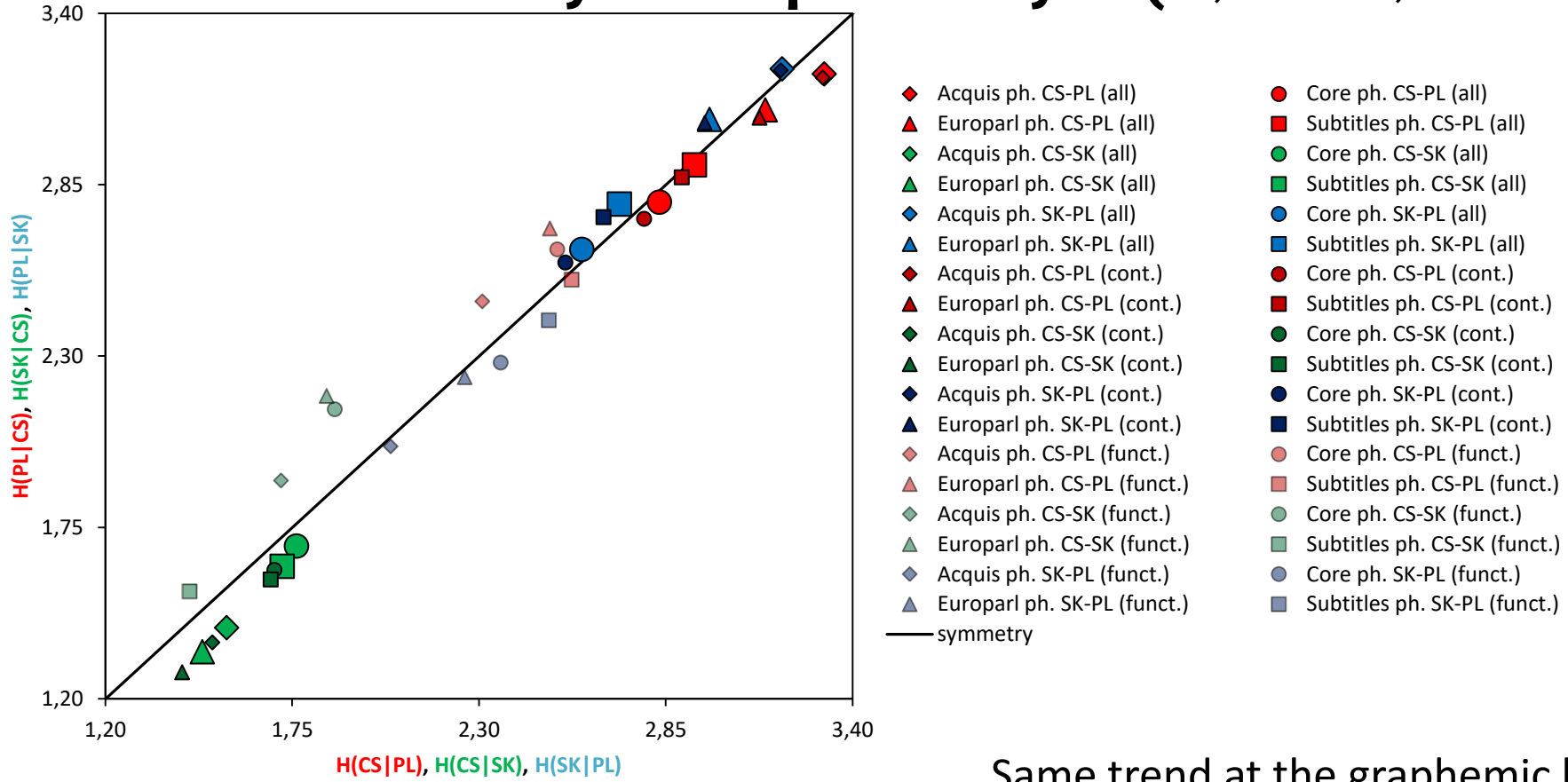# Results: Asymmetry of MI between SK-PL (all w.)



Phonetic layer:

**SK > PL** (RA = 0,019)

Graphemic layer:

**PL > SK** (RA = 0,025)

↳ Agree with socioling. research

Same side for all subcorpora across layers

$$RA = \frac{H(PL \mid SK) - H(SK \mid PL)}{H(PL \mid SK)}$$

Legend:
— symmetry
◆ Acquis ph. SK-PL (all)
● Core ph. SK-PL (all)
▲ Europarl ph. SK-PL (all)
■ Subtitles ph. SK-PL (all)
◆ Acquis gr. SK-PL (all)
● Core gr. SK-PL (all)
▲ Europarl gr. SK-PL (all)
■ Subtitles gr. SK-PL (all)

Axes: H(PL|SK) vertical, H(SK|PL) horizontal, range 2,30 to 3,40

# Results: MI & asym. on phon. layer (all, content, function w.)



Same trend at the graphemic layer...

# Future: What could be improved?

- Data
  - Usable parallel corpora aligned word-by-word

- Levenshtein method
  - CE without aligning by Lev. distance were not so different
  - Need to add constraints or additional rules for aligning

    example:    CS: x aː p ʊ #          CS: x aː p # ʊ

                    SK: x aː p ɛ m          SK: x aː p ɛ m

- Conditional entropy
  - Statistical validation of this method (realized only for Scandinavian languages)

Thank you.

# References

- HAUGEN, Einar. (1966). Semicommunication: The Language Gap in Scandinavia. *Sociological Inquiry*, **36**: 280-297.
- GOLUBOVIĆ, Jelena. (2016). Mutual intelligibility in the Slavic language area [Groningen]: University of Groningen
- BUDOVIČOVÁ, Viera. (1987). Semikomunikácia jako lingvistický problém. *Studia Academica Slovaca*, **16**: 49-66.
- MUSILOVÁ, Květoslava. (2000). *Město a jeho jazyk*. Bratislava: Vydavateľ´stvo slovenskej akadémie vied.
- NEKVAPIL, Jiří, SLOBODA, Marián & WAGNER, Petr. (2009). *Mnohojazyčnost v České republice*. Praha: Nakladatelství Lidové noviny.
- DOBROTOVÁ, Ivana, MUSILOVÁ, Květoslava. (2012). Z dějin kontaktu dvou blízkých slovanských jazyků. *Bohemistyka*, 2012(**1**): 35-60.
- Heeringa, Wilbert Jan. (2004). Measuring Dialect Pronunciation Differences using Levenshtein Distance. [Groningen]: University of Groningen
- MOBERG Jens, GOOSKENS Charlotte, NERBONNE John, VAILLETTE Nathan. (2007). Condition Entropy Measures Intelligibility among Related Languages. In: Peter Dirix, Ineke Schuurman, Vincent Vandeghinste and Frank Van Eynde (eds.). *Computational Linguistics in the Netherlands 2006: Selected papers from the 17th CLIN Meeting*. Utrecht: LOT, 51-66.