

Stylometric Analysis of Church Slavonic Texts of Czech Origin

Miroslav Vepřek

Palacký University Olomouc

Radek Čech

Masaryk University in Brno



Faculty
of Arts

MUNI
ARTS

Research goals

- applicability of stylometric methods to (Old) Church Slavonic texts
- heuristic approach (without prior assumptions)
 - may these methods lead to the identification of new perspectives?
- stylometric analysis as a means of testing existing hypotheses
 - based on philological research

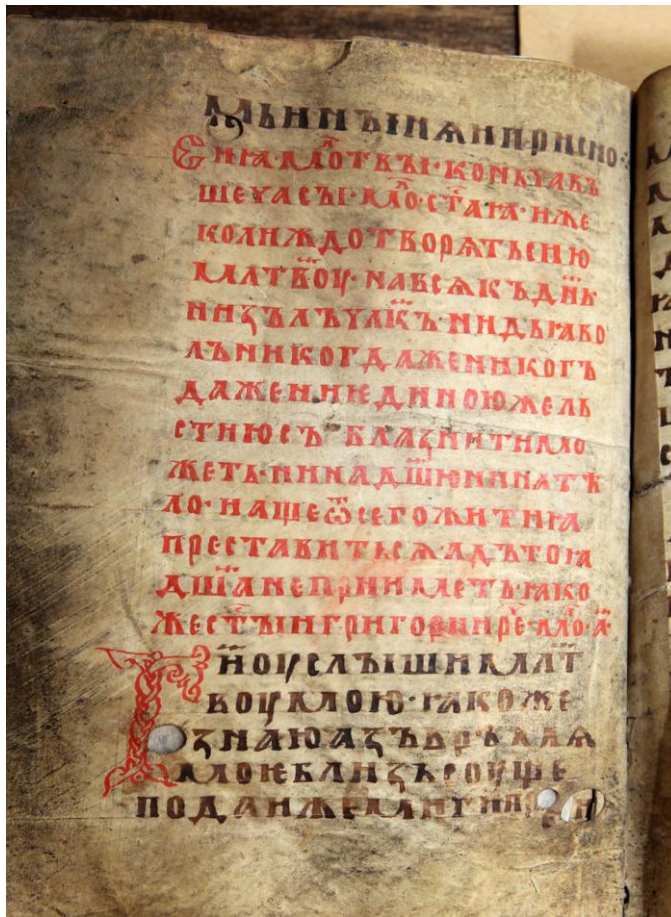
Material (Corpus) of Czech CS Texts

- Church Slavonic language
 - solely written language
 - used in different Slavic areas from the end of the 9th century in various variants (local „redactions“)
- 13 literary texts of various lengths
 - the *Forty Gospel Homilies of Gregory the Great* (93.358 tokens)
 - the *Prayer of Confession of Sins* (533 tokens)
- proven or at least hypothetical Czech origin (10th-11th centuries)
 - preserved in later copies in different redactions of CS (mostly East Slavic)

Material (Corpus) of Czech CS Texts

- original writings and translations from Latin
- various genres
 - legends, prayers, homilies, a pseudo-Gospel
- the *Codex Suprasliensis* and *Vita Methodii* added for a comparison

Digitalization



Fol. 54b: /14/ Мѣя септврѣа еѣ .кѣ. днѣ ѿбѣ/15/снѣ
сѣго бачеблая кнѣзѣ /16/ чвѣска. ѿ блѣн сѣчѣ.

/17/ Се нѣк свѣсѣ прѣческое словѣ /18/ еже глѣше гѣ нашѣ іѣхъ хѣ.
вѣдѣ /19/ во рече в послѣднѣа днѣ. іакѣ /20/ минимъ сѣча. вѣстанѣтъ
во /21/ вратѣ на врата сѣго ѿ сѣхъ /22/ на вѣхъ сѣхъ, ѿ врази домаш-
нинѣ. /23/ члѣвци во сѣхъ не минѣ вѣдоуѣ Fol. 55a: /1/ да вѣздастѣ ѿмѣ
вѣ по дѣлѣ /2/ ѿхѣ. вѣ же кнѣзѣ великѣ сѣа/3/воѣ в чѣхѣ живѣи ѿме-
немѣ /4/ вѣротиславѣ. ѿ жена сѣго доро/5/гомирѣ, родисѣа же сѣа пѣрѣ-
/6/вѣнца. ѿ іакѣ крѣпѣста ѿ наре/7/коша ѿмѣ сѣмѣ вѣчеславѣ. ѿ вѣхъ/8/зрасте
ѿтрокѣ іакѣ вѣ оуѣати /9/ сѣмѣ волѣсѣ. ѿ призѣа корѣти/10/славѣ кнѣзѣ
ѿппѣ сѣтера сѣ вѣсѣ /11/ клѣросѣ. ѿ пѣвѣшимъ литѣ/12/ргѣю вѣ цѣркѣ
сѣтѣмъ мѣрѣа. /13/ ѿ вѣземъ ѿтрока поставѣ на /14/ стѣпени прѣ вѣлатѣрѣ.
ѿ блѣн /15/ ѿ сѣ рекѣ гѣ іѣхъ хѣ блѣн ѿтроѣ /16/ сѣ вѣвѣнимѣ. ѿмѣ же
вѣвѣлѣ /17/ сѣнѣ всѣа правѣдники твоѣа. /18/ ѿ постѣригоша кнѣзи ѿни тѣ-
/19/мѣ минимѣ. іакѣ оуѣво блѣе/20/нимѣмъ ѿппѣ того. но мѣтѣа/21/ми
блѣговѣрѣнымѣ. нача ѿтро/22/кѣ рѣсти вѣлѣтѣю вѣжѣію хра/23/нимѣ. ѿ
вѣдѣа ѿ баба сѣвоа лю/24/дѣнинаѣ оуѣчити кнѣгамѣ словѣ- Fol. 55b: /1/ нѣ-
скимѣ. по сѣлѣдѣа попѣвѣ ѿ на/2/вычѣ разѣмѣ доверѣ. ѿсѣдѣи <и> во-
ро/3/тиславѣ в вѣдѣучѣа. ѿ нача ѿ/4/трокѣ оуѣчитѣсѣа кнѣгамѣ ла-
/5/тынскимѣ. ѿ наоуѣчисѣа до/6/верѣ. в тожѣ вѣремѣ оуѣмѣе вѣро/7/ти-
славѣ кнѣзѣ. ѿ поставѣишѣ /8/ кнѣзѣа вѣчеслава на столѣ /9/ дѣдѣни
ѿ ѿтоле болѣславѣ /10/ нача подѣ нимѣ ходитѣ. вѣа/11/шѣта во ѿба
мѣла. мѣти же /12/ сѣю дорогомѣрѣ оуѣтверѣди вѣ/13/мѣю ѿ люди сѣвоа
оуѣстрой. іакѣ /14/ вѣспѣитѣ сѣны сѣвоа. іакѣ на/15/ча вѣчеславѣ стрѣитѣ
люди /16/ сѣвоа ѿмѣашѣ же сѣстры .ѣ. ѿ /17/ вѣдѣста ѿ в рѣснаѣ кнѣже-
нѣа /18/ ѿ оуѣстройстѣа ѿ ѿ вѣзложѣ вѣ /19/ вѣлѣтѣа такѣ на вѣчеслава
кнѣ/20/зѣа. ѿ нача же оуѣмѣитѣ кнѣги /21/ лѣтынскимѣ. іакѣже доверѣ

- a) Г<о>сподъ и бл<аго>с<ло>ви ш<тъ>че: this is the usual initial formula of the Lives.
- b) Originally: врази члѣвкоу домашнѣи сѣго (Micah 7, 6 and Matthew 10, 36; cf. no. 6.12).
- c) Read: людмила.
- d) Read: воудѣчѣ, church, castle and village W of Prague, now Budeř.
- e) Read: розна (OCS разна) 'various'.
- f) = добрыи (the double grave accent = и).

Digitalization

- font Bukyvede
- unification of punctuation
- expansion of abbreviations
- special character for the numerical value of letters
- unification of variant graphemes: и (и, ѣ, ѣ, ѣ, ѣ), оу (оу, у, Ѹ), etc.

Digitalization

МѢ́ЦА́ то́го*. вѣ́. днѣ́. мѣ́нїе́ стѣ́го сте́фана. па́трїа́рха рѣ́мьскѣ́ и дру́жины́ его́.

МѢ́ЦА́ то́гоже · вѣ́ · днѣ́ · мѣ́нїе́ сва́тѣ́го сте́фана · па́трїа́рха рѣ́мьскѣ́го и
дру́жины́ его́ ·

мѣ́саца то́гоже · вѣ́ · днѣ́ · мѣ́нїе́ сва́тѣ́го сте́фана · па́трїа́рха
рѣ́мьскѣ́го и дру́жины́ его́ ·

Methods

- average token length
- moving average type-token ratio
- distances between texts
 - the most frequent words
 - Cosine delta distance

Average token length (ATL)

- length of the word measured in graphemes

$$ATL_{text} = \frac{\sum_{i=1}^N L_i}{N}$$

L ... length of the word

N ...number of words in text

Moving Average Type-Token Ratio (MATTR)

- vocabulary diversity (lexical richness)
- based on the type token ratio
- measures lexical diversity using moving windows

Peter loves Mary. John loves Mary too.

Moving Average Type-Token Ratio (MATTR)

- vocabulary diversity (Lexical richness)
- based on the type token ratio
- measures lexical diversity using moving windows

***Peter loves Mary. John loves** Mary too.*

$$TTR_1 = \frac{V}{N} = \frac{4}{5} = 0.8$$

Moving Average Type-Token Ratio (MATTR)

- vocabulary diversity (Lexical richness)
- based on the type token ratio
- measures lexical diversity using moving windows

*Peter **loves Mary. John loves Mary** too.*

$$TTR_2 = \frac{V}{N} = \frac{3}{5} = 0.6$$

Moving Average Type-Token Ratio (MATTR)

- vocabulary diversity (Lexical richness)
- based on the type token ratio
- measures lexical diversity using moving windows

*Peter loves **Mary. John loves Mary too.***

$$TTR_3 = \frac{V}{N} = \frac{4}{5} = 0.8$$

Moving Average Type-Token Ratio (MATTR)

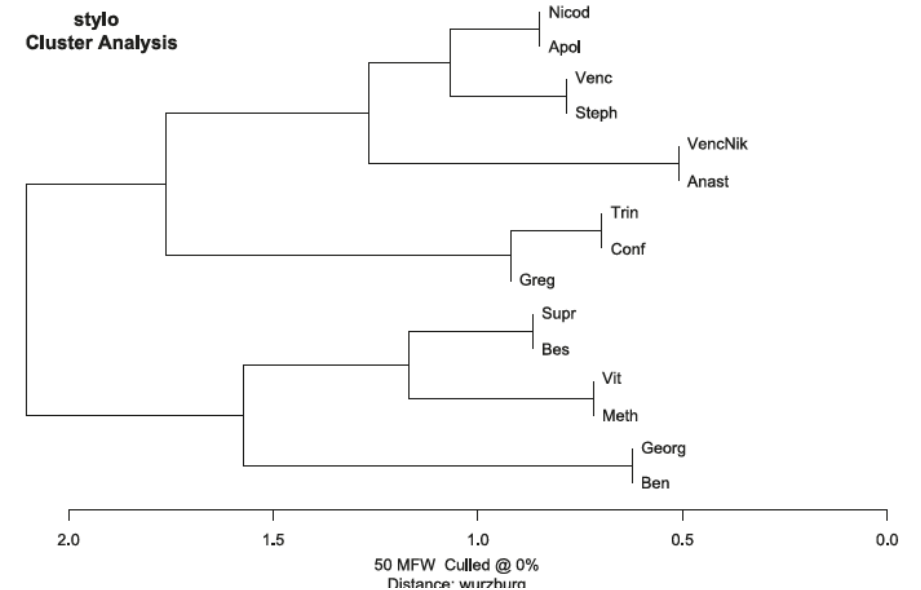
- vocabulary diversity (Lexical richness)
- based on the type token ratio
- measures lexical diversity using moving windows

Peter loves Mary. John loves Mary too.

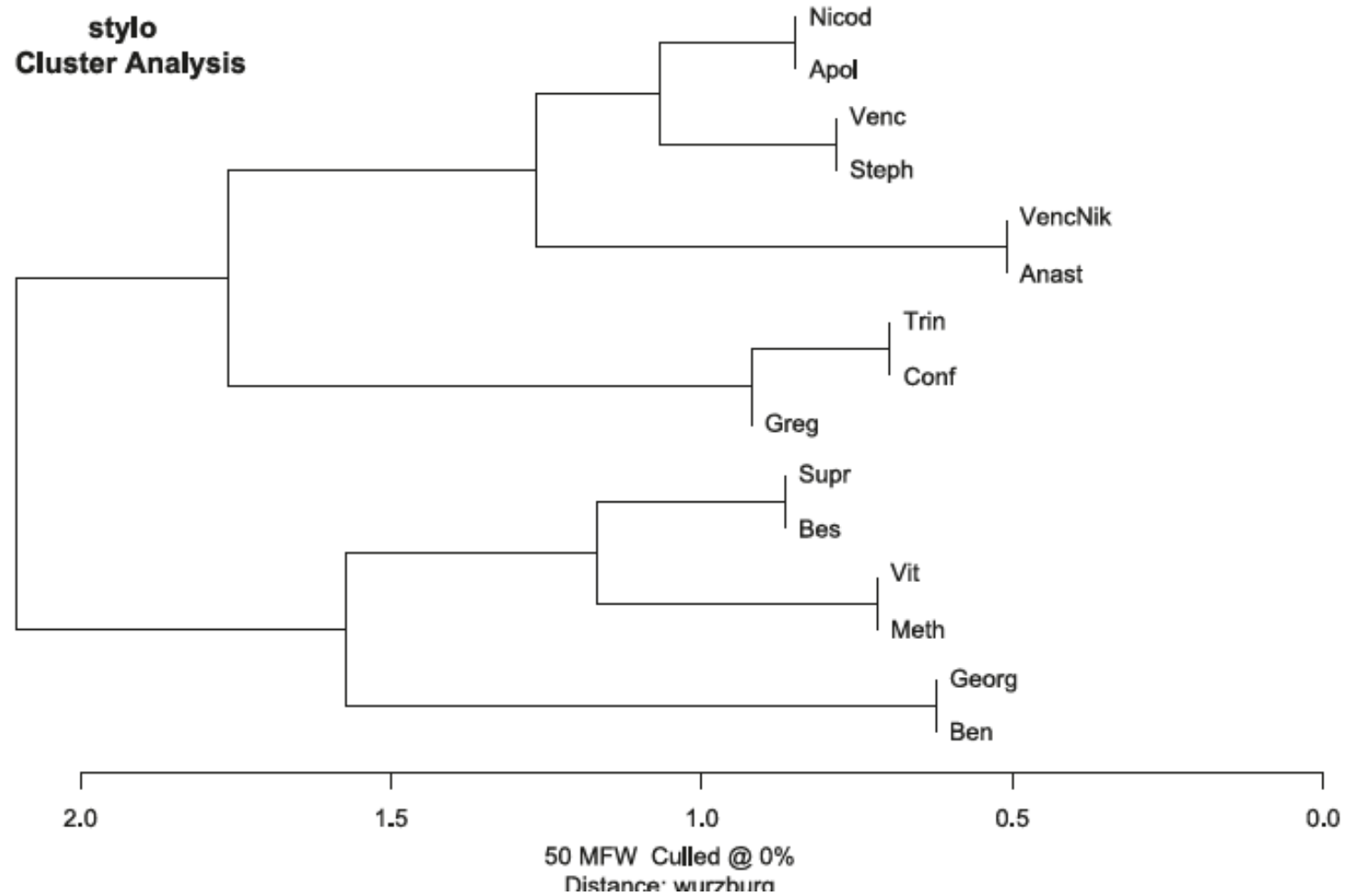
$$MATTR = \frac{0.8 + 0.6 + 0.8}{3} = 0.733$$

Distances between texts

- relative frequency of the most frequent words
 - in this analysis: 50 and 200
- Cosine Delta distance
- hierarchical cluster analysis
- Stylo software



Distances between texts

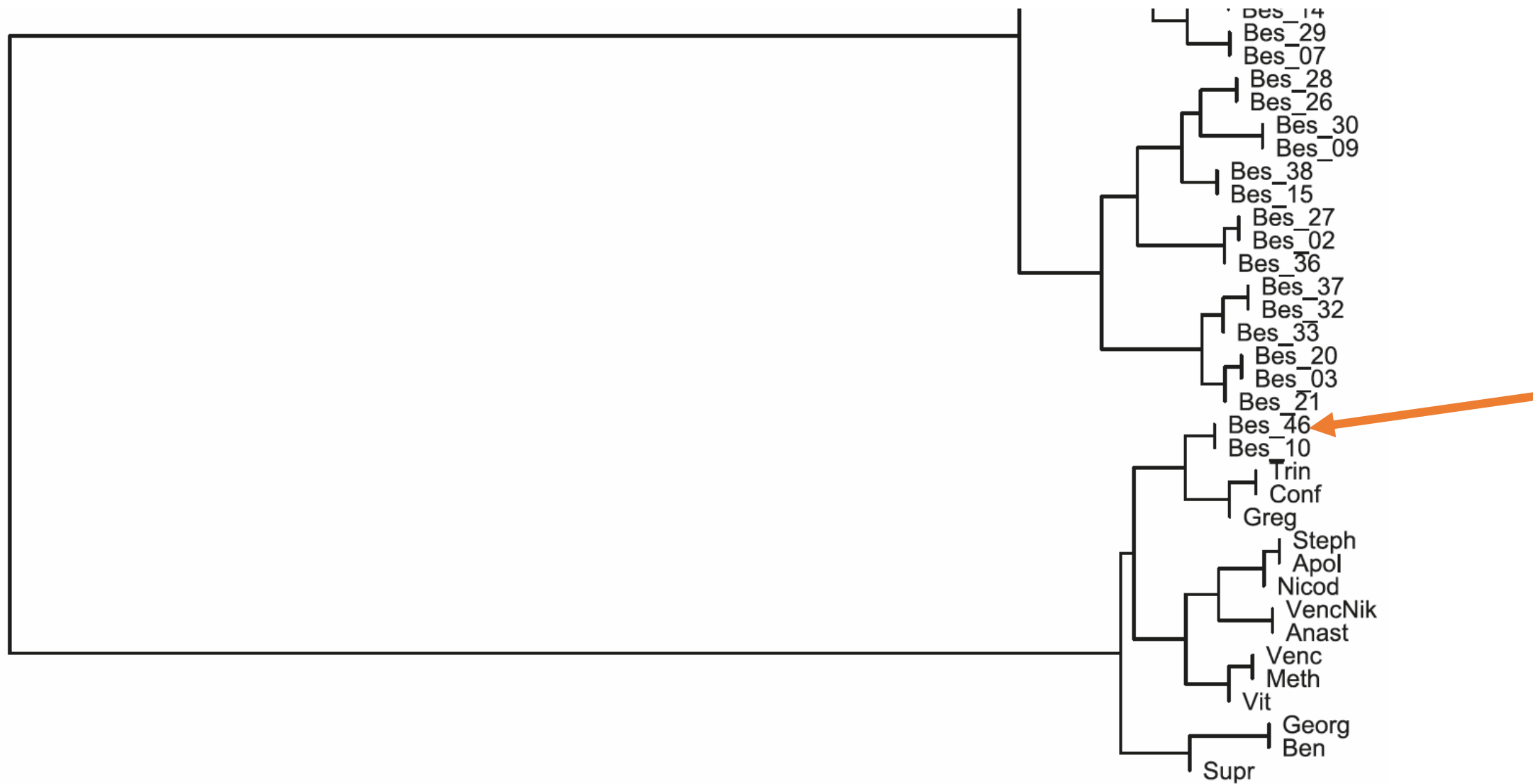


Results

- two sets of data:
 1. the *Forty Gospel Homilies of Gregory the Great (Bes)* as 1 file
 2. *Bes* divided into 46 parts (mainly according to chapters)
- reasons
 - to reduce the large differences in text length
 - to account for possible internal heterogeneity of *Bes*
- in general, **genre** characteristic is the main factor

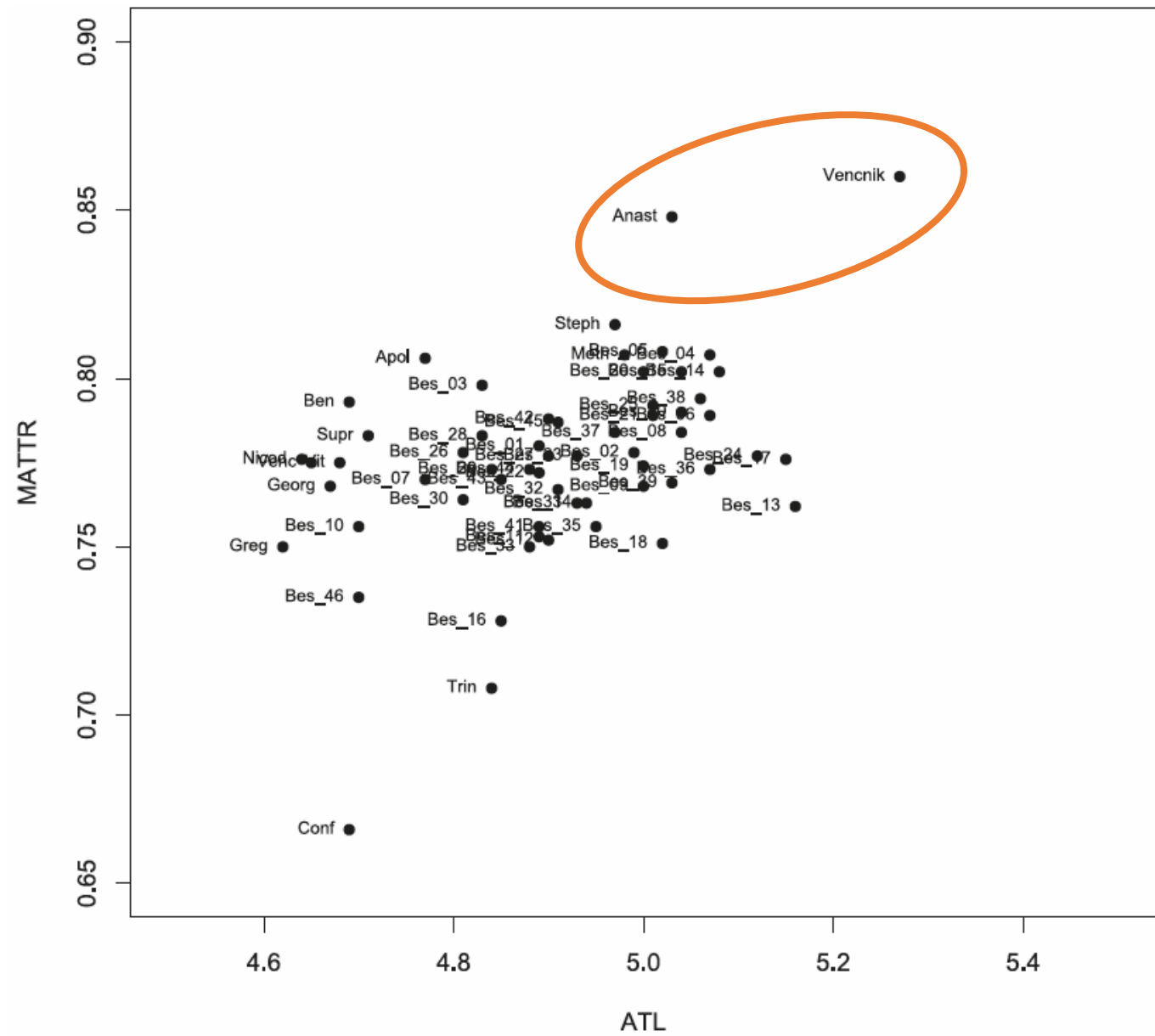
A) Prayers

- *Trin, Conf, Greg*
 - stable cluster at both MFW 50 and MFW 200
 - MATTR – relatively low lexical diversity
 - correlation between MATTR and ATL
- *Bes*, segment 46 (prayer added to the base text)
 - joining *Trin, Conf, Greg* cluster



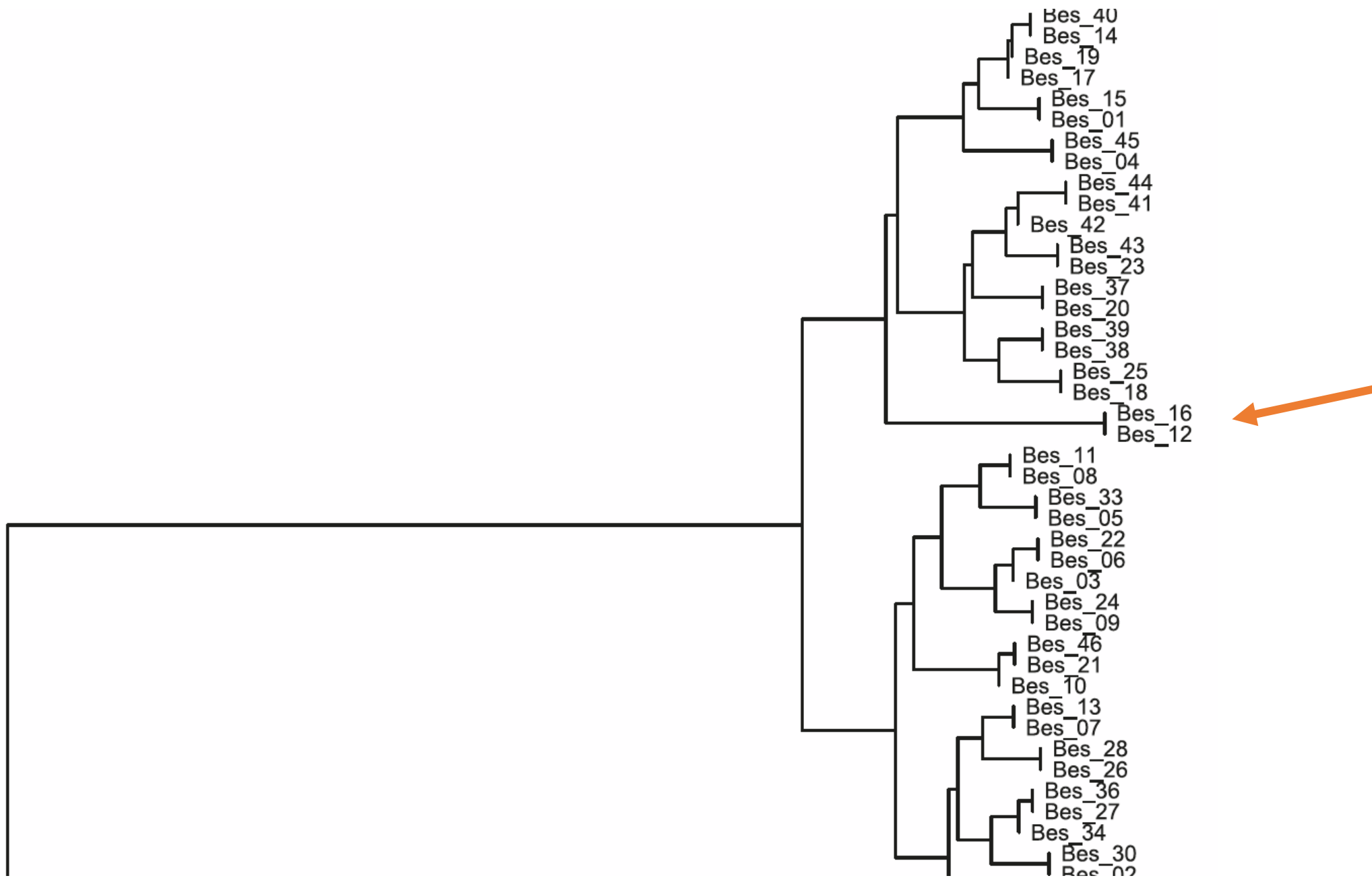
B) Legends

- two basic groups according to MFW
- *VencNik* and *Anast*
 - joint cluster in every analysis (MFW, MATTR, ATL)
 - supports the philological hypothesis of the same origin, perhaps even the same author/school
- *Ben* and *Georg*; *Meth* and *Vit*
 - possible influence of manuscript preservation?



C) „Segmented“ *Bes*

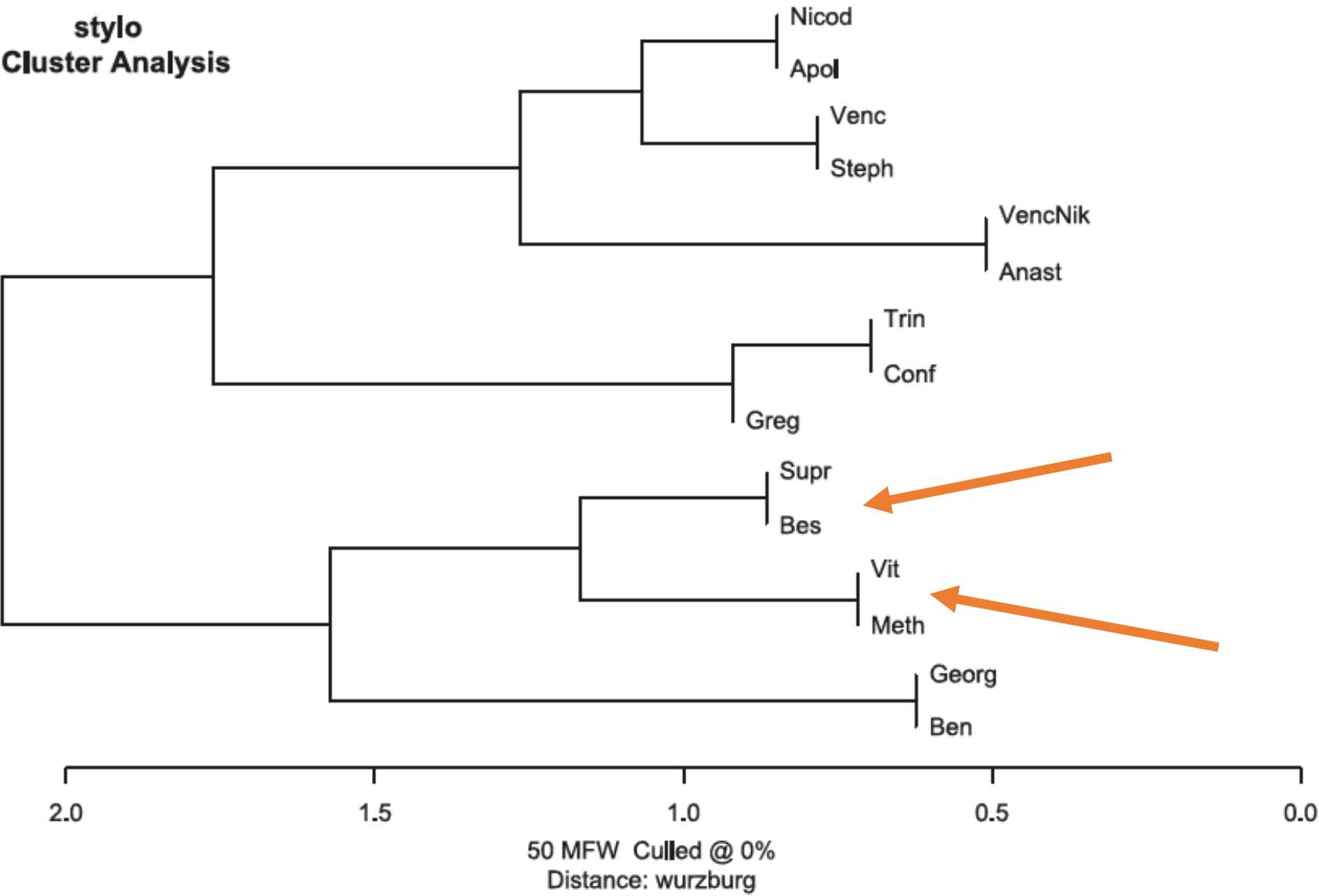
- some passages from a variant manuscript in the edition
 - they are not grouped in any parameter
 - no significant influence of manuscript characteristics!
- one homily (no. 9) twice in the CS translation (not an identical text!) – the same cluster



D) Texts of different origin

- *Supr* and *Meth*
 - methods cluster them mainly according to a genre
- *Supr* and *Bes*
 - both contain homilies
 - however, translations from Greek x Latin
- *Meth* and *Vit*
 - the same manuscript
 - the origin of *Vit* in the early 10th century
 - some scholars even suggested the Great Moravian origin

stylo
Cluster Analysis



Conclusions

- stylometric methods appear to be plausible and yield relevant results
- genre is the essential categorization factor
- similarities between texts may also reflect a possible common origin
- whether a text is translated or original is not the main criterion for relatedness
- this type of research can support existing hypotheses and highlight new connections between texts

Thank you!