

Most Frequent Words as a Tool for Authorship Recognition

Radek Čech

MUNI
ARTS



Ján Mačutek



Acknowledgment

- supported by the project GA ČR GA17-02545S, 2017–2019

Pavel Kosek (MU), Petra Mutlová (MU), Miroslav Kubát (OU)

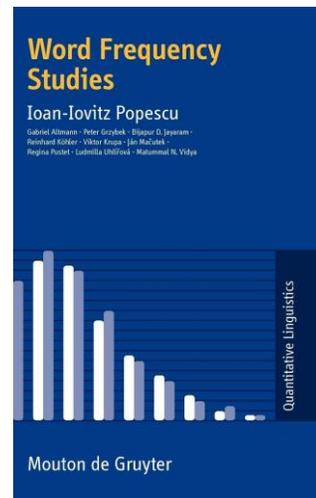
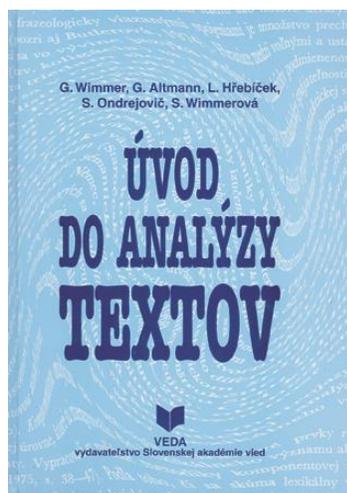
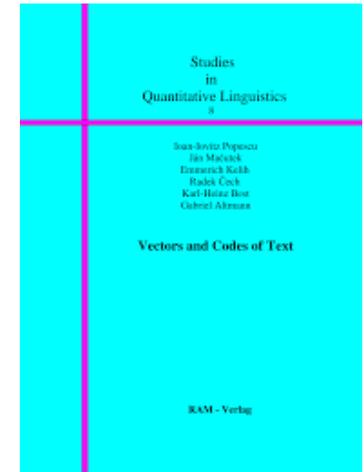
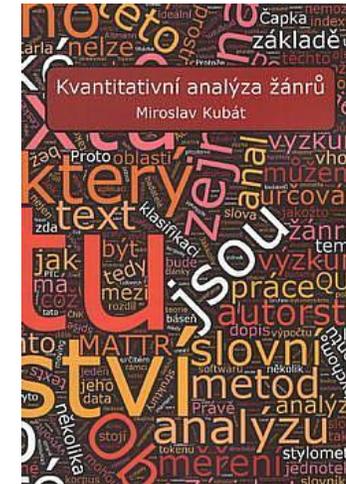
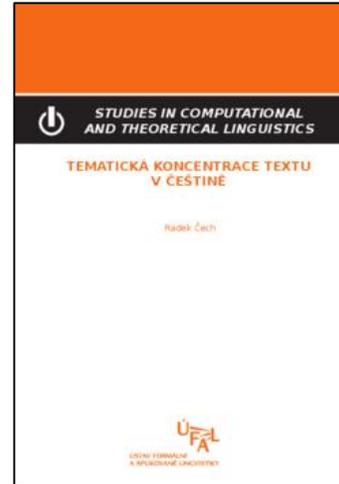


Background

- theory of text
 - quantitative properties
 - regularities
 - mechanisms
 - laws

Background

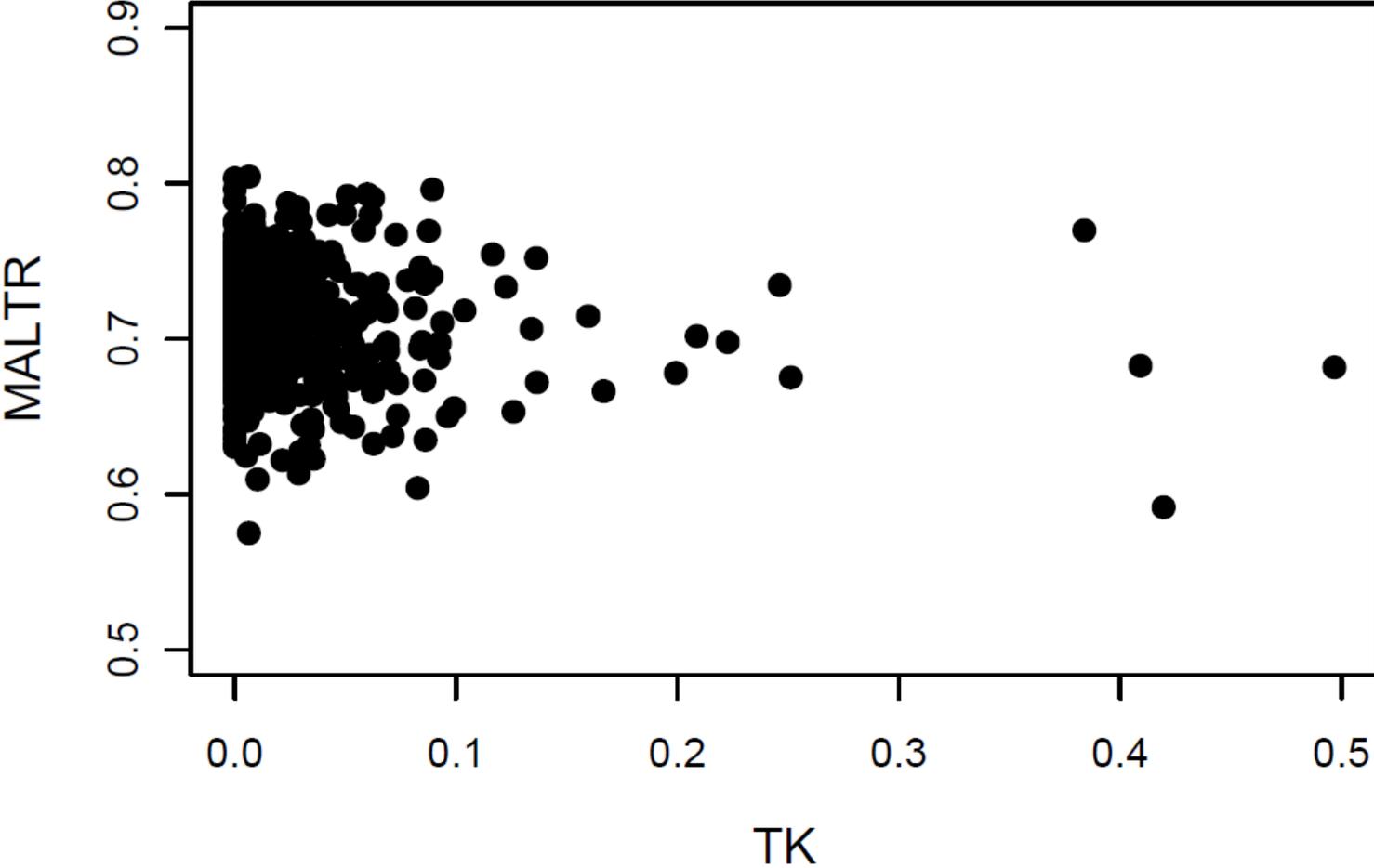
- theory of text
 - quantitative properties
 - regularities
 - mechanisms
 - laws



Theory of text (and language)

- relationships among particular properties (stochastic)
 - lexical diversity – text length
 - lexical diversity – thematic concentration
 - entropy – verb distances
 - distances between hapax legomena – lexical diversity
 - distribution of clause/sentence lengths – distances between hapax legomena
 - autosemantic text structure – thematic concentration

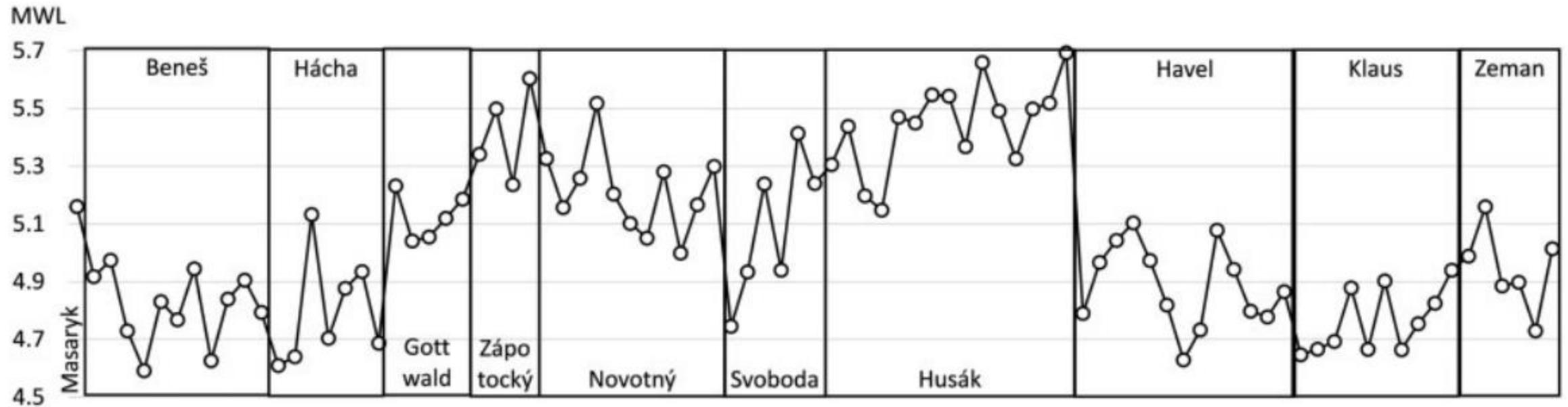
Theory of text (and language)



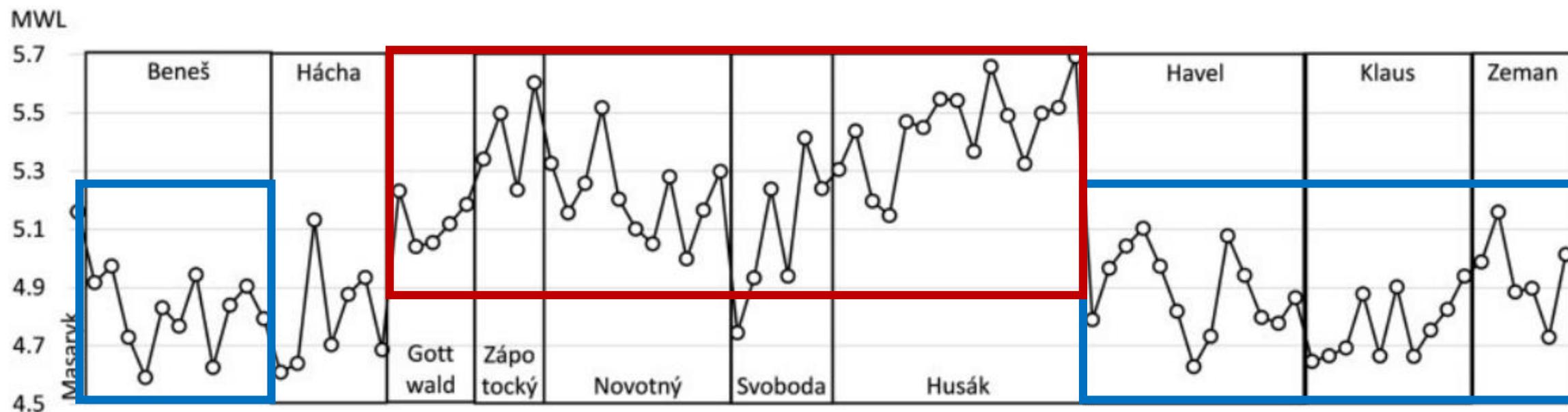
Theory of text (and language)

- relationships between particular properties and pragmatics
- stylometry
 - text type
 - genre
 - period
 - gender
 - political background / ideology
 - authorship

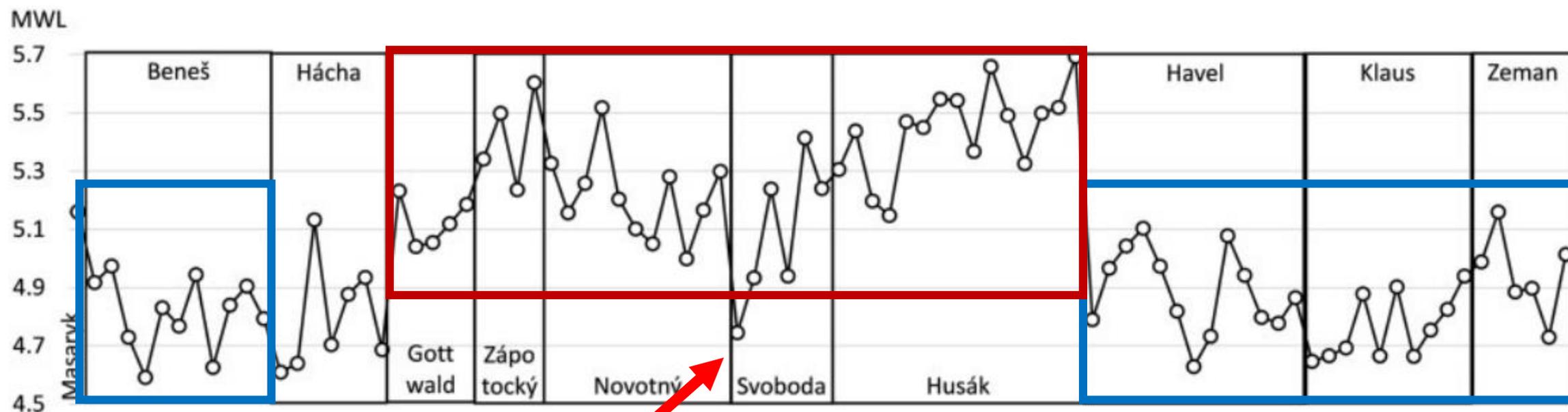
Mean word length



Mean word length



Mean word length

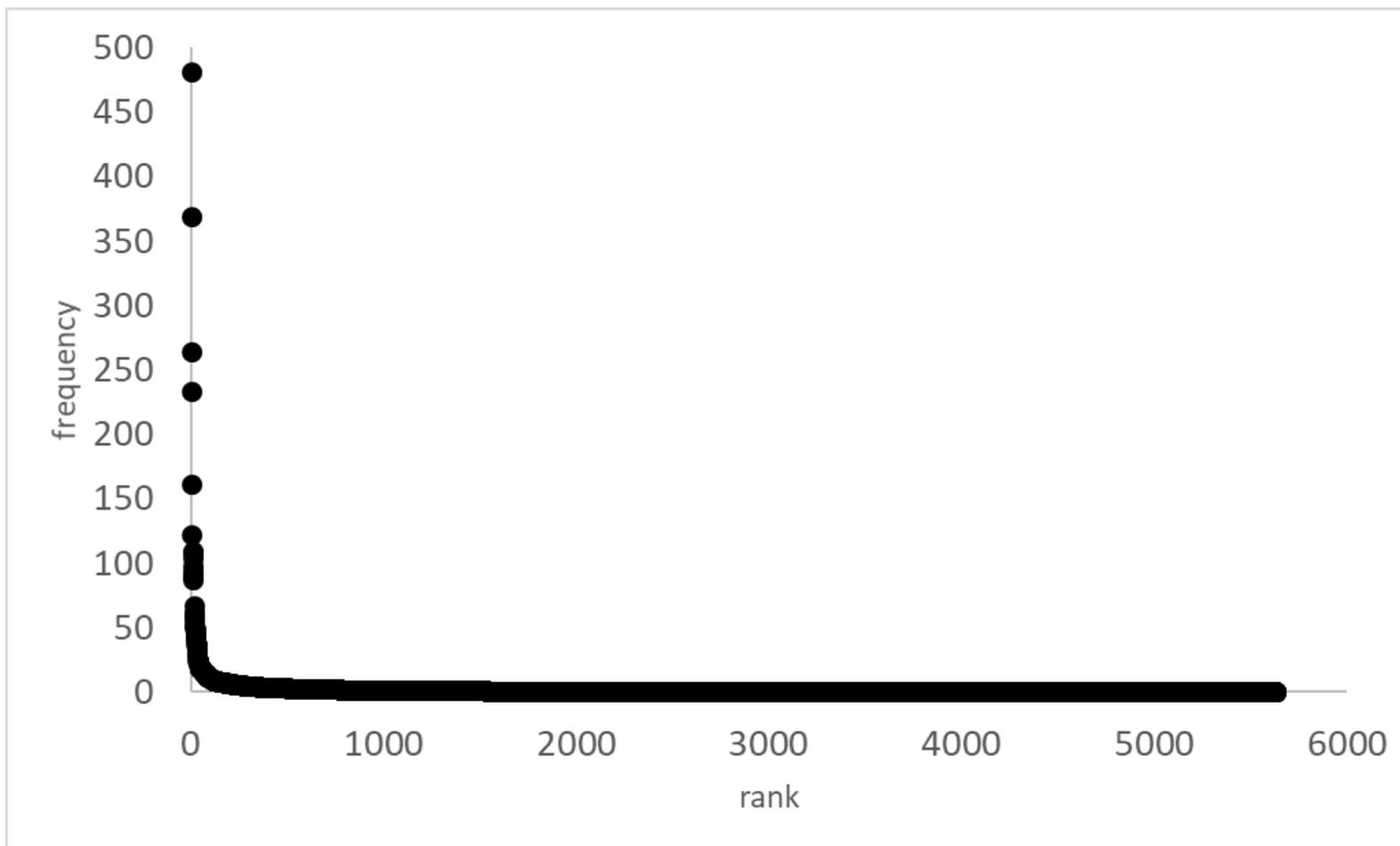


From theory to the application

Word frequencies and authorship

- lexical diversity / vocabulary richness
- proportion of hapax legomenon
- ...
- proportion of the most frequent words

Most frequent words



J. Škvorecký: Eva byla nahá

Most frequent words

Škvorecký: Eva byla nahá		
pořadí	slovo	f
1	a	481
2	se	369
3	na	264
4	v	234
5	jsem	161
6	s	122
7	z	110
8	američan	108
9	to	104
10	řekl	98
11	ale	94
12	do	91
13	že	89
14	řekla	87
15	dívka	67

Most frequent words

Škvorecký: Eva byla nahá		
pořadí	slovo	f
1	a	481
2	se	369
3	na	264
4	v	234
5	jsem	161
6	s	122
7	z	110
8	američan	108
9	to	104
10	řekl	98
11	ale	94
12	do	91
13	že	89
14	řekla	87
15	dívka	67

Most frequent words

Škvorecký: Eva byla nahá		
pořadí	slovo	f
1	a	481
2	se	369
3	na	264
4	v	234
5	jsem	161
6	s	122
7	z	110
8	američan	108
9	to	104
10	řekl	98
11	ale	94
12	do	91
13	že	89
14	řekla	87
15	dívka	67

Hrabal: Perlička na dně		
pořadí	slovo	f
1	a	2239
2	se	1203
3	to	1037
4	na	879
5	ale	514
6	tak	504
7	do	467
8	si	459
9	jsem	456
10	v	446
11	že	440
12	je	432
13	já	363
14	když	296
15	jak	283

Hašek: Osudy...I.		
pořadí	slovo	f
1	a	7045
2	se	6061
3	na	3927
4	že	3469
5	to	3075
6	v	2585
7	je	1801
8	do	1749
9	s	1667
10	si	1534
11	když	1387
12	z	1375
13	tak	1308
14	jsem	1286
15	švejk	1188

Most frequent words

Škvorecký: Eva byla nahá		
pořadí	slovo	f_rel
1	a	0.037
2	se	0.028
3	na	0.020
4	v	0.018
5	jsem	0.012
6	s	0.009
7	z	0.008
8	američan	0.008
9	to	0.008
10	řekl	0.007
11	ale	0.007
12	do	0.007
13	že	0.007
14	řekla	0.007
15	dívka	0.005

Hrabal: Perlička na dně		
pořadí	slovo	f_rel
1	a	0.054
2	se	0.029
3	to	0.025
4	na	0.021
5	ale	0.012
6	tak	0.012
7	do	0.011
8	si	0.011
9	jsem	0.011
10	v	0.011
11	že	0.011
12	je	0.010
13	já	0.009
14	když	0.007
15	jak	0.007

Hašek: Osudy...I.		
pořadí	slovo	f_rel
1	a	0.035
2	se	0.030
3	na	0.020
4	že	0.017
5	to	0.015
6	v	0.013
7	je	0.009
8	do	0.009
9	s	0.008
10	si	0.008
11	když	0.007
12	z	0.007
13	tak	0.007
14	jsem	0.006
15	švejk	0.006

Most frequent words

Škvorecký: Eva byla nahá		
pořadí	slovo	f_rel
1	a	0.037
2	se	0.028
3	na	0.020
4	v	0.018
5	jsem	0.012
6	s	0.009
7	z	0.008
8	američan	0.008
9	to	0.008
10	řekl	0.007
11	ale	0.007
12	do	0.007
13	že	0.007
14	řekla	0.007
15	dívka	0.005

Hrabal: Perlička na dně		
pořadí	slovo	f_rel
1	a	0.054
2	se	0.029
3	to	0.025
4	na	0.021
5	ale	0.012
6	tak	0.012
7	do	0.011
8	si	0.011
9	jsem	0.011
10	v	0.011
11	že	0.011
12	je	0.010
13	já	0.009
14	když	0.007
15	jak	0.007

Hašek: Osudy...I.		
pořadí	slovo	f_rel
1	a	0.035
2	se	0.030
3	na	0.020
4	že	0.017
5	to	0.015
6	v	0.013
7	je	0.009
8	do	0.009
9	s	0.008
10	si	0.008
11	když	0.007
12	z	0.007
13	tak	0.007
14	jsem	0.006
15	švejk	0.006

Distances between words / texts

$$\Delta_{(AB)} = \frac{1}{n} \sum_{i=1}^n \left| \frac{A_i - \mu_i}{\sigma_i} - \frac{B_i - \mu_i}{\sigma_i} \right|$$

n ... the number of MFW

A, B ... texts for the comparison

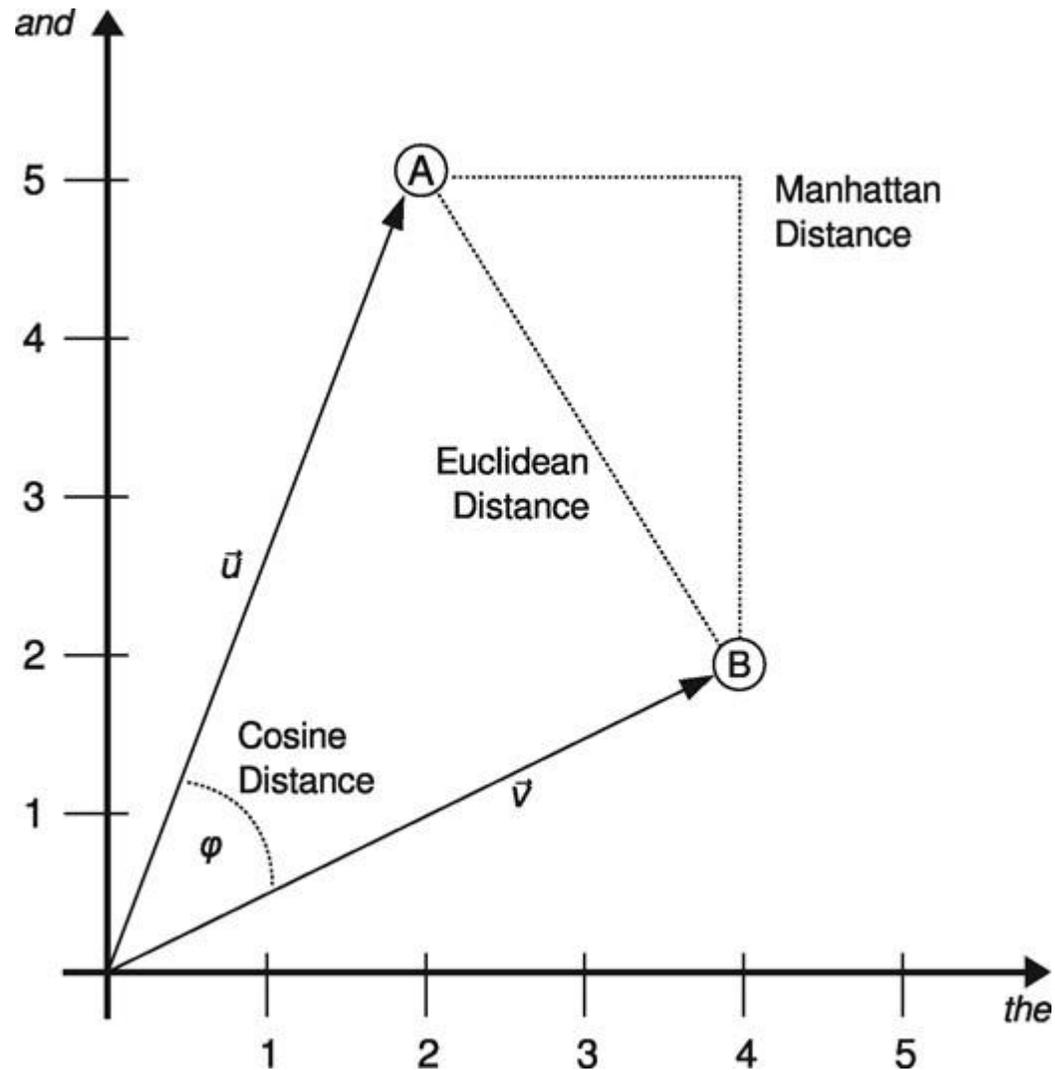
A_i ... the frequency of a given word in the text A

B_i ... the frequency of a given word in the text B

μ_i ... the average frequency of a given word in corpus

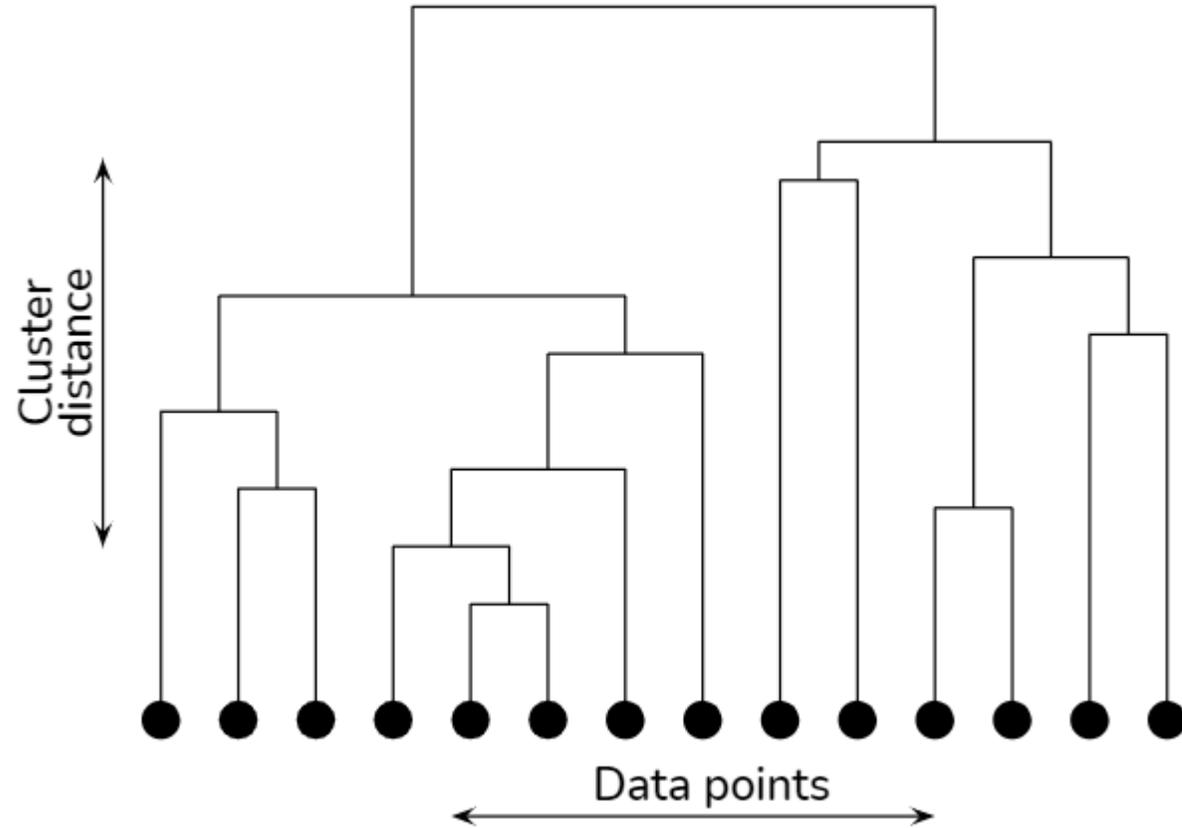
σ_i ... the standard deviation of the frequency of a given word

Distances between words / texts

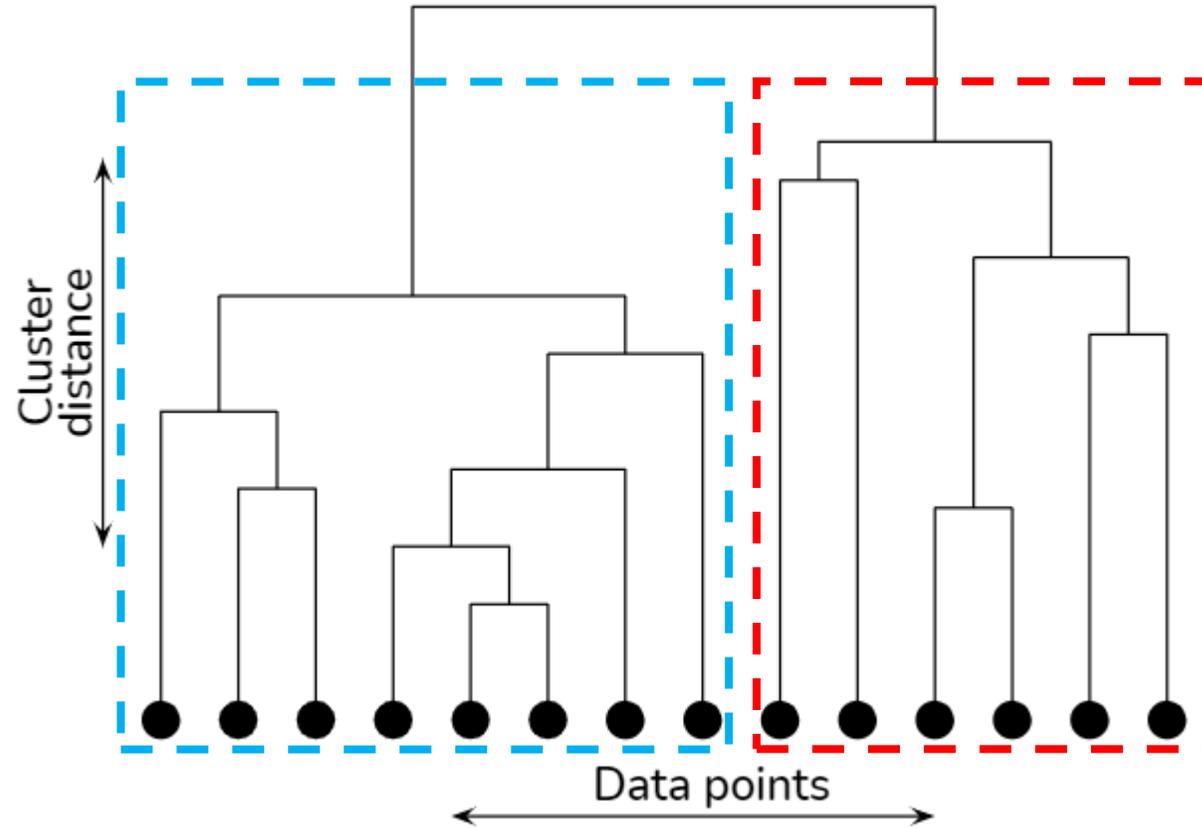


Evert et al. (2017)

Clustering



Clustering



Stylo

- Eder, M., Rybicki, J., & Kestemont, M. (2016). Stylometry with R: a package for computational text analysis. *The R Journal*, 8(1).

Stylometry with R | stylo | set parameters

INPUT & LANGUAGE	FEATURES	STATISTICS	SAMPLING	OUTPUT	
INPUT:	plain text	xml	xml (plays)	xml (no titles)	html
	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
LANGUAGE:	English	English (contr.)	English (ALL)	Latin	Latin (u/v > u)
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	Polish	Hungarian	French	Italian	Spanish
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
	Dutch	German	CJK	Other	Native encoding
	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input checked="" type="checkbox"/>
<input type="button" value="OK"/>					

Stylo

Stylometry with R | stylo | set parameters

INPUT & LANGUAGE	FEATURES	STATISTICS	SAMPLING	OUTPUT
FEATURES:	words <input checked="" type="radio"/>	chars <input type="radio"/>	ngram size <input type="text" value="1"/>	preserve case <input type="checkbox"/>
MFV SETTINGS:	Minimum <input type="text" value="100"/>	Maximum <input type="text" value="100"/>	Increment <input type="text" value="100"/>	Start at freq. rank <input type="text" value="1"/>
CULLING:	Minimum <input type="text" value="0"/>	Maximum <input type="text" value="0"/>	Increment <input type="text" value="20"/>	List Cutoff <input type="text" value="5000"/>
				Delete pronouns <input type="checkbox"/>
VARIOUS:	Existing frequencies <input type="checkbox"/>	Existing wordlist <input type="checkbox"/>	Select files manually <input type="checkbox"/>	List of files <input type="checkbox"/>

OK

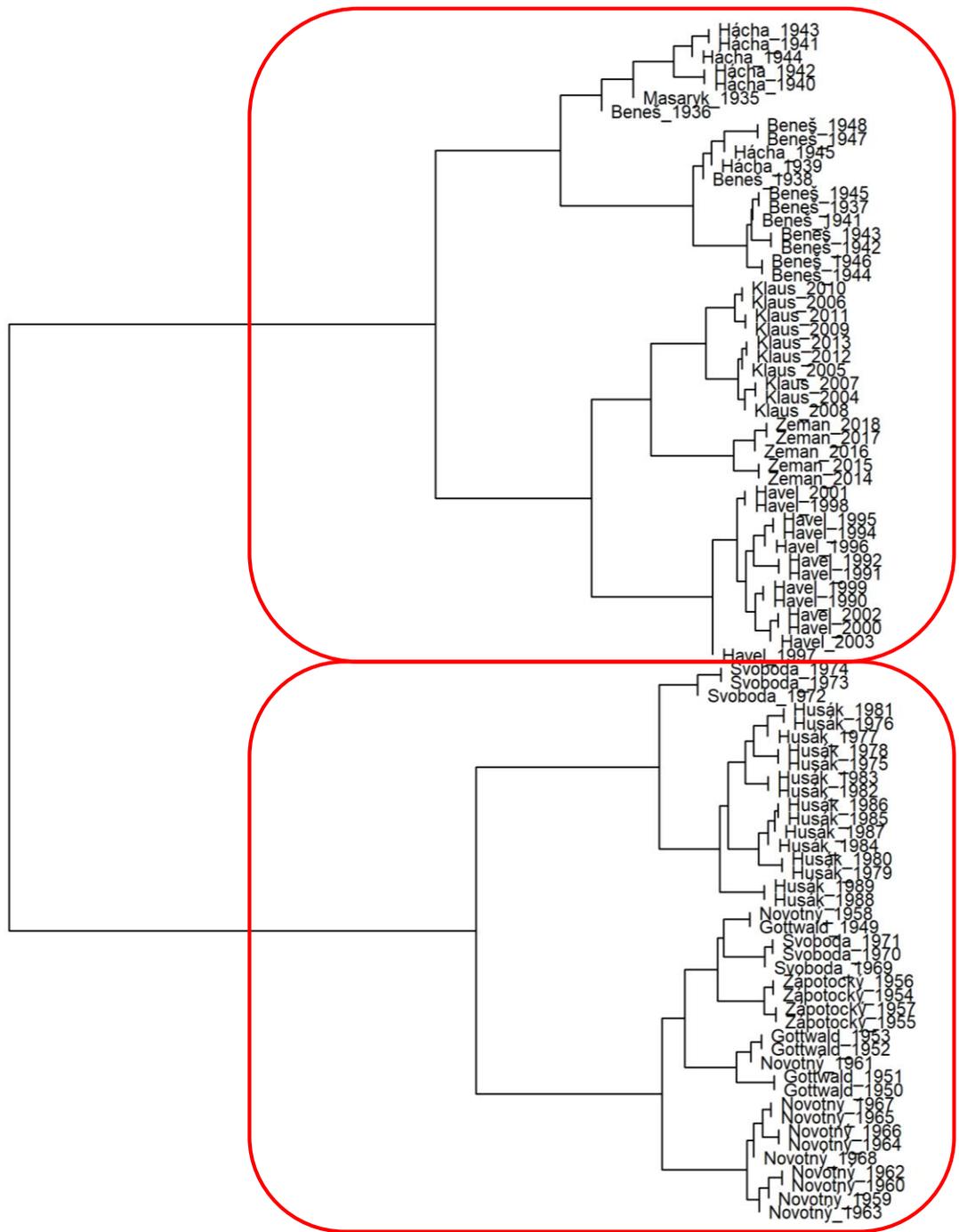
Stylometry with R | stylo | set parameters

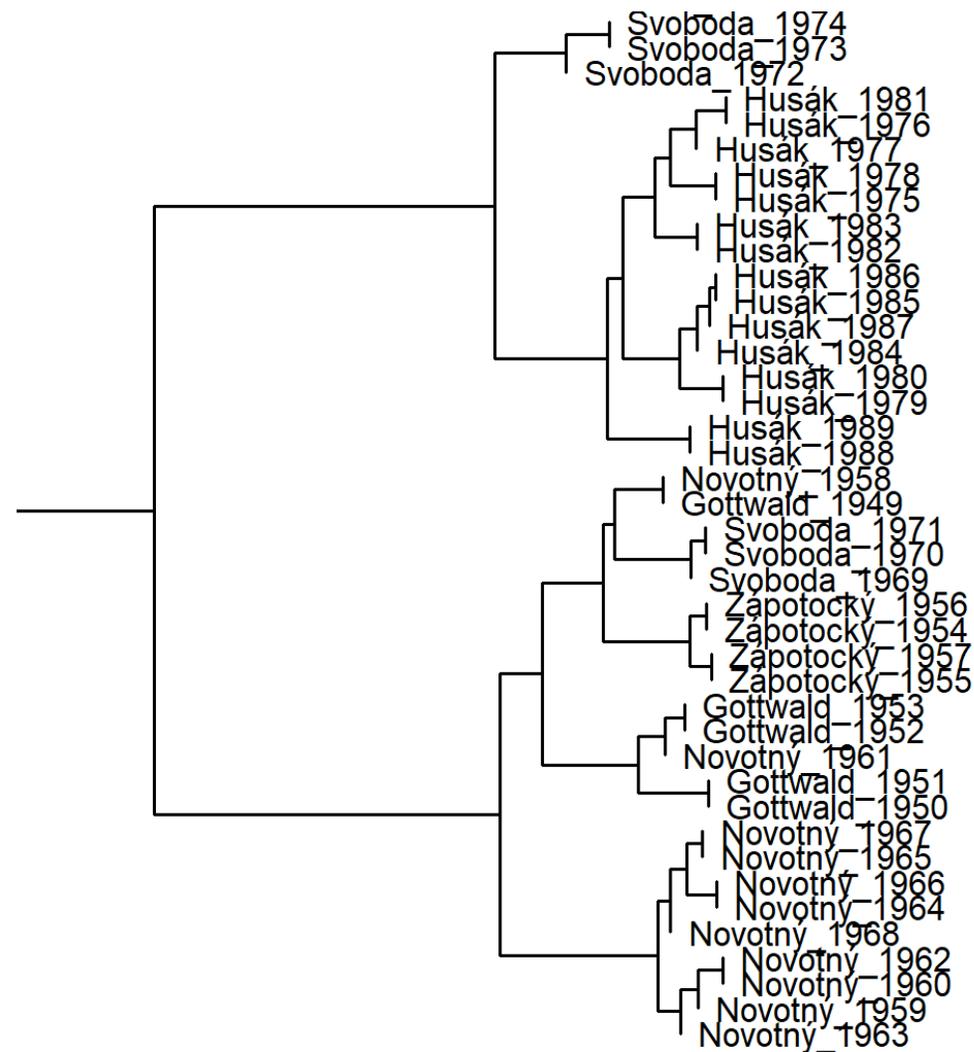
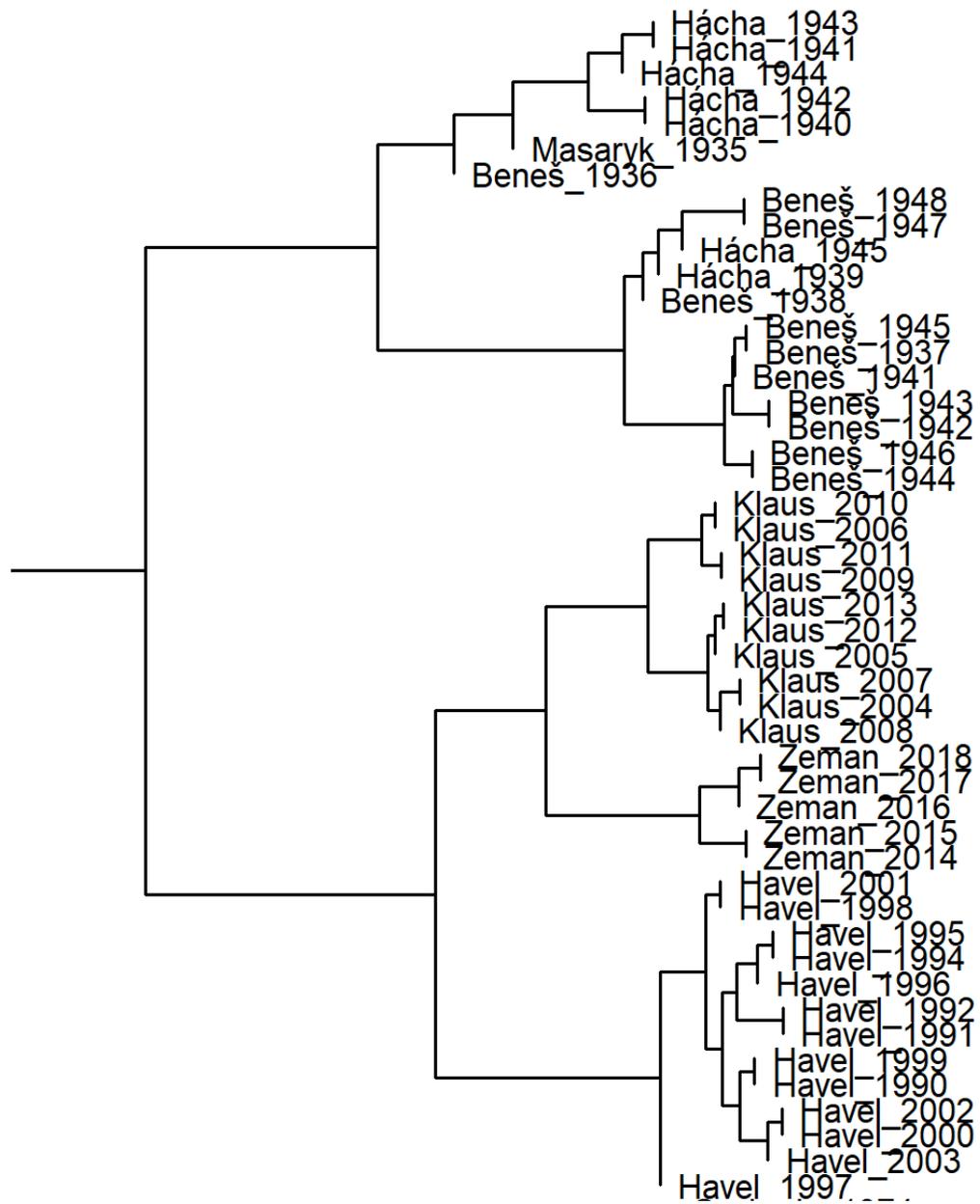
INPUT & LANGUAGE	FEATURES	STATISTICS	SAMPLING	OUTPUT	
STATISTICS:	Cluster Analysis <input type="radio"/>	MDS <input type="radio"/>	PCA (cov.) <input checked="" type="radio"/>	PCA (corr.) <input type="radio"/>	tSNE <input type="radio"/>
	Consensus Tree <input type="radio"/>	Consensus strength <input type="text" value="0.5"/>			
DELTA DISTANCE:	Classic Delta <input checked="" type="radio"/>	Cosine Delta <input type="radio"/>	Eder's Delta <input type="radio"/>	Eder's Simple <input type="radio"/>	Entropy <input type="radio"/>
	Manhattan <input type="radio"/>	Canberra <input type="radio"/>	Euclidean <input type="radio"/>	Cosine <input type="radio"/>	Min-Max <input type="radio"/>

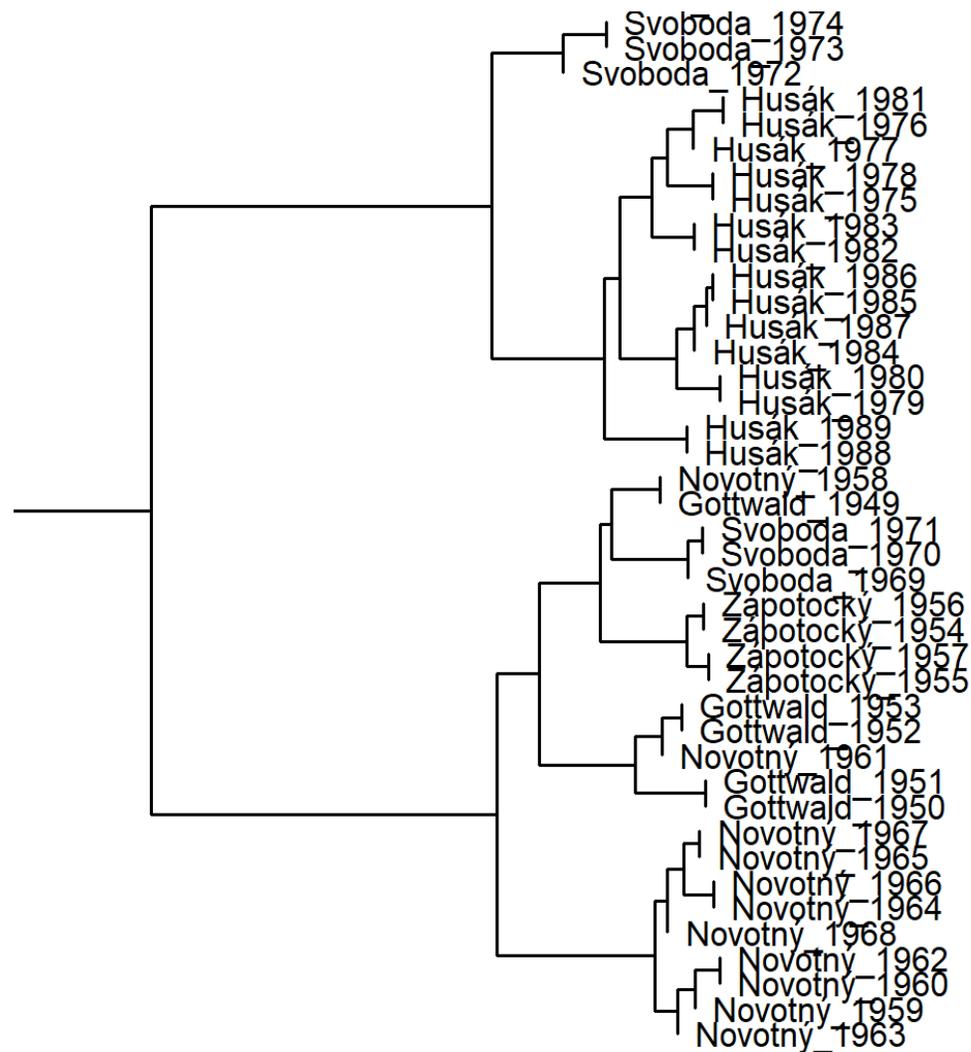
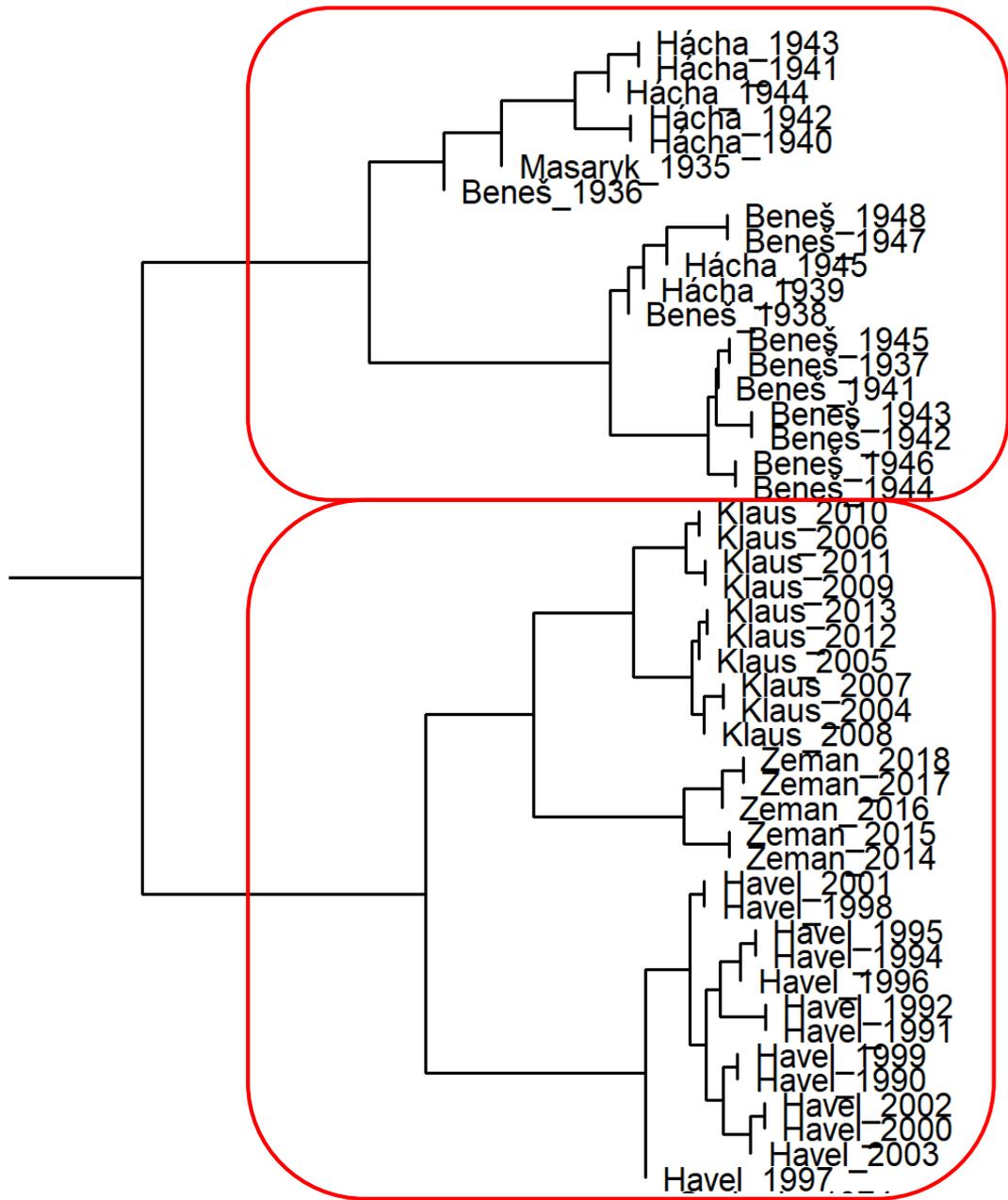
OK

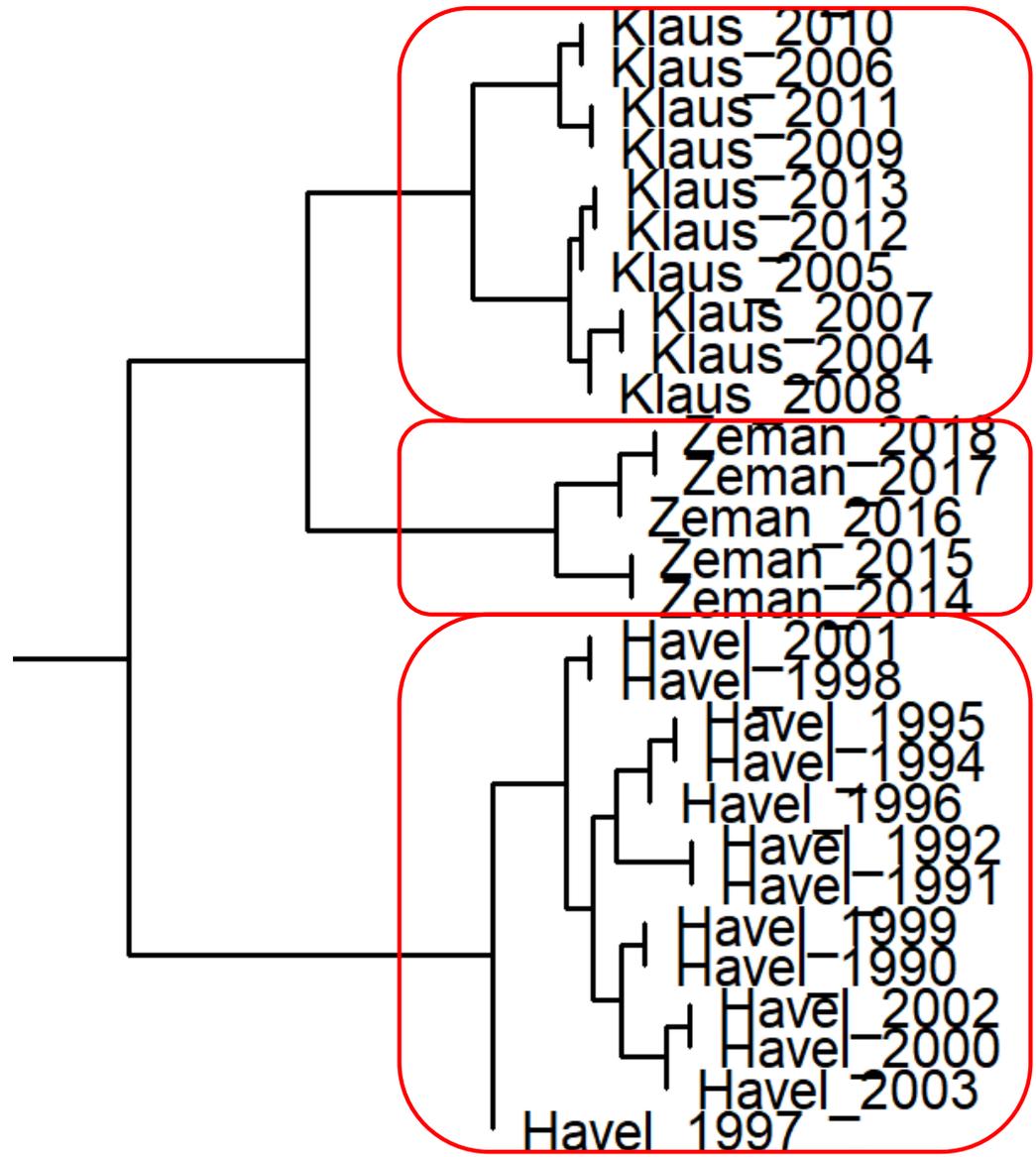
Prezidentské projevy

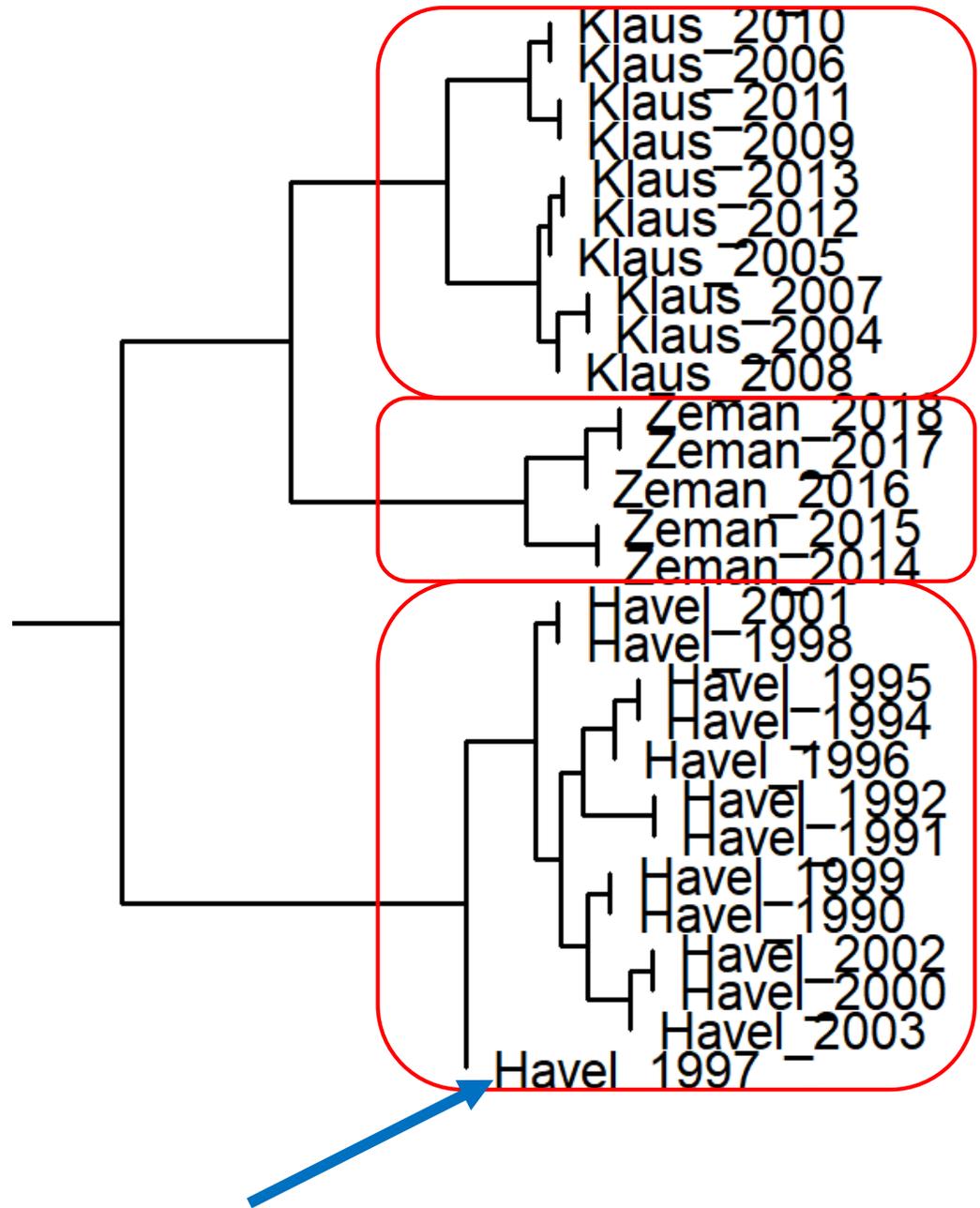
- Kubát, M., Mačutek, J., Čech, R. (2021). Communists spoke differently. An analysis of Czechoslovak and Czech annual presidential speeches. *Digital Scholarship in the Humanities*, 36, 138-152.
- presidential speeches: 1935–2018
- 100 MFW, culling = 60 %

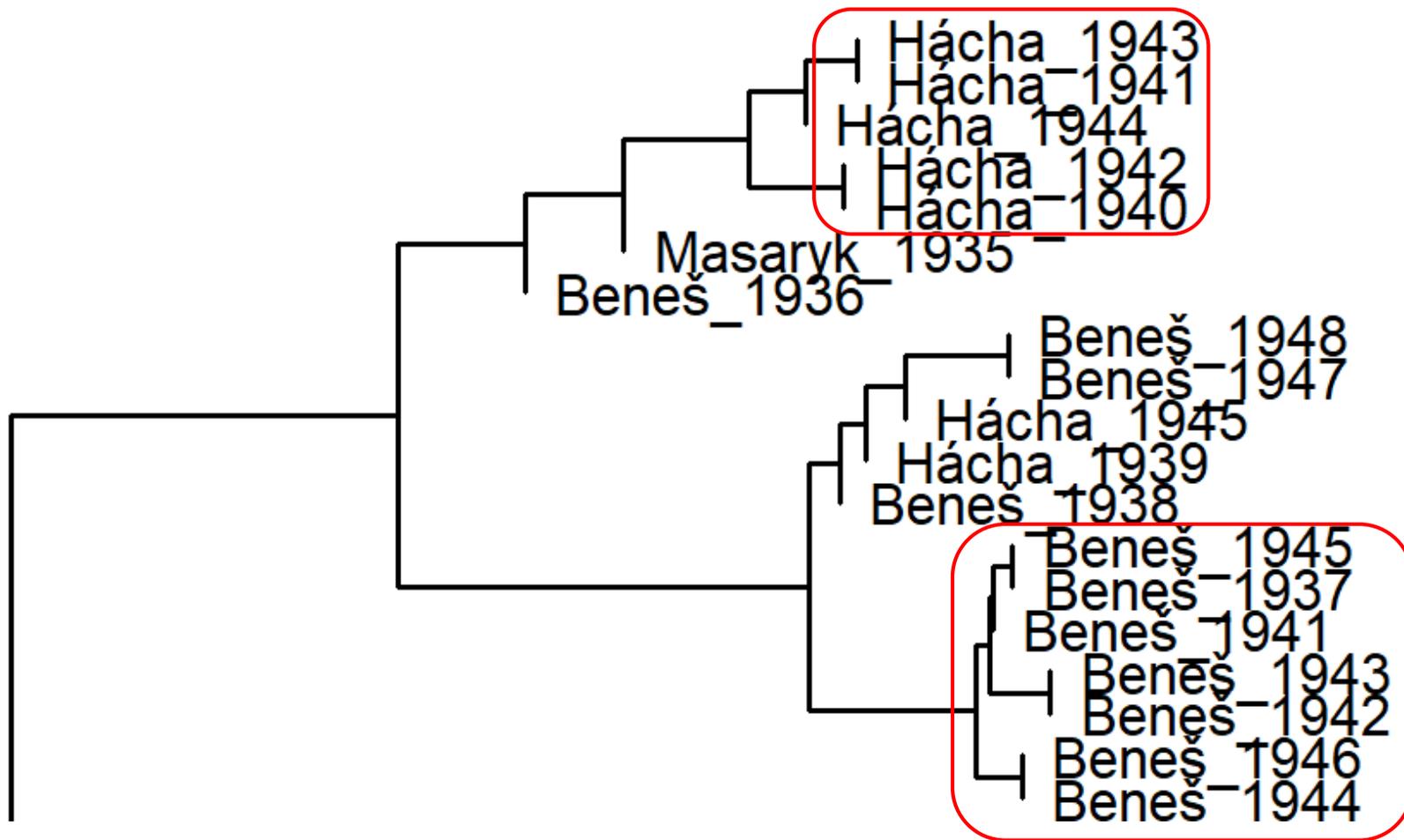


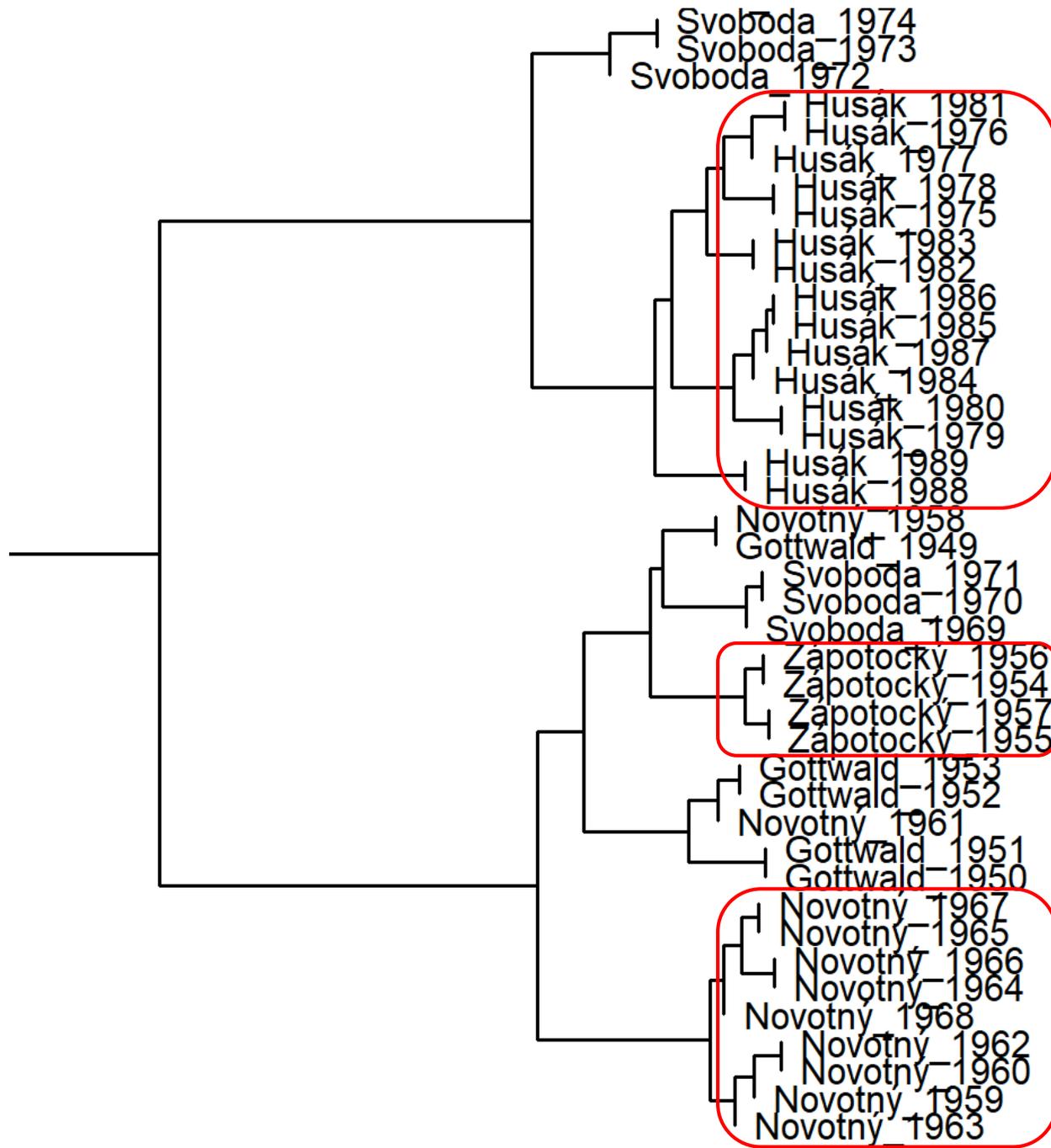


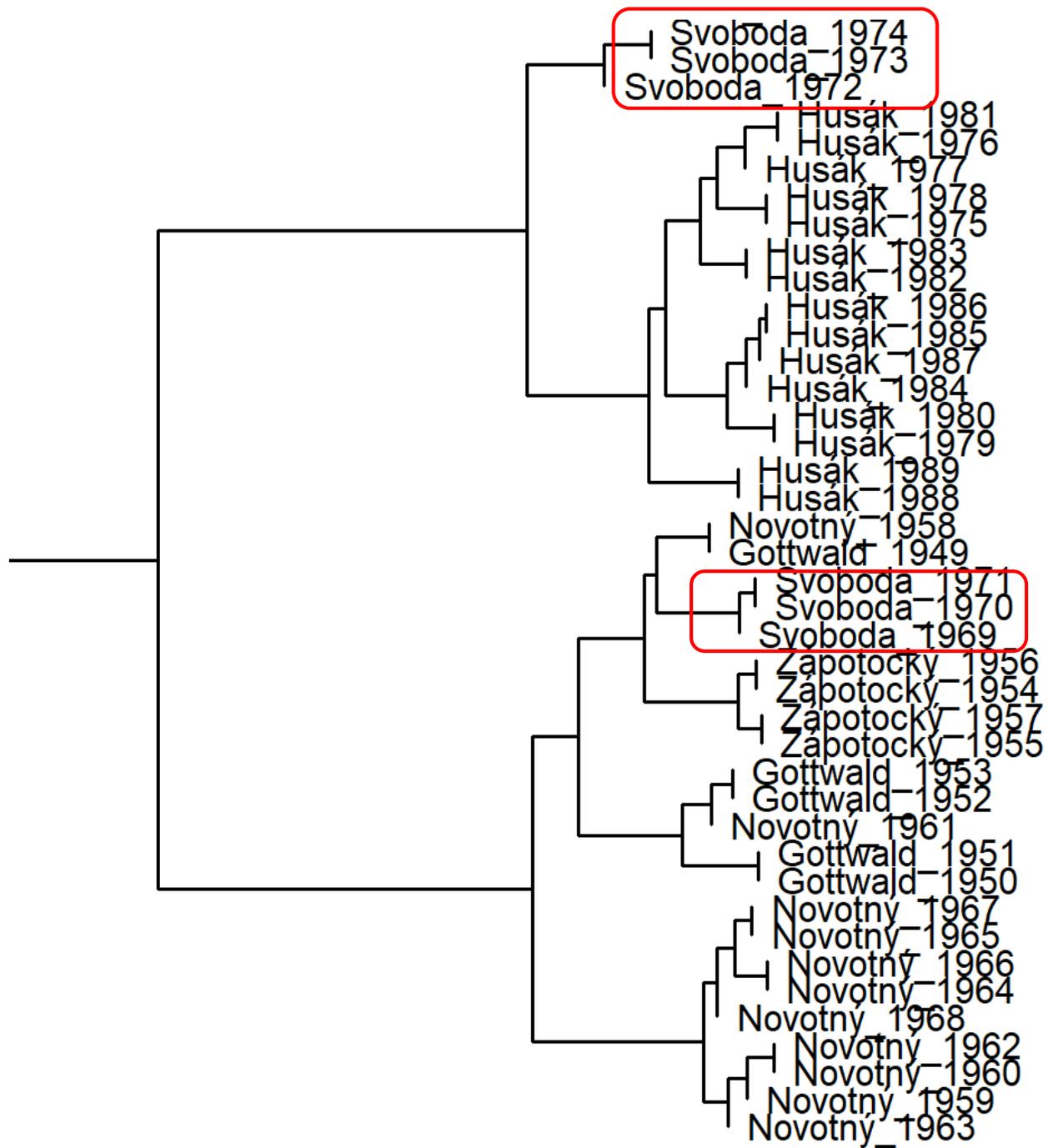












Bible svatováclavská & comments

- Kosek, P., Čech, R. (2018). Stylové aspekty Bible svatováclavské – stylometrická analýza. In Zand, G., Newerkla, S.M. (eds.). Jezuitská kultura v českých zemích / Jesuitische Kultur in den böhmischen Ländern. Host, 195-209.

Bible svatováclavská & comments

- New Testament
 - Konstanc → Šteyer
- Old Testament
 - Šteyer → Barner
 - Job – divide
- „[...] (*domníváme se, že by se stylistickým rozdílem textu dalo zjistit, odkud překládal už jen Šteyer*) [...]“ Vintr (1997)

Bible svatováclavská & comments

- comments
 - Šteyer: New Testament + Genesis
 - Barner: Old Testament

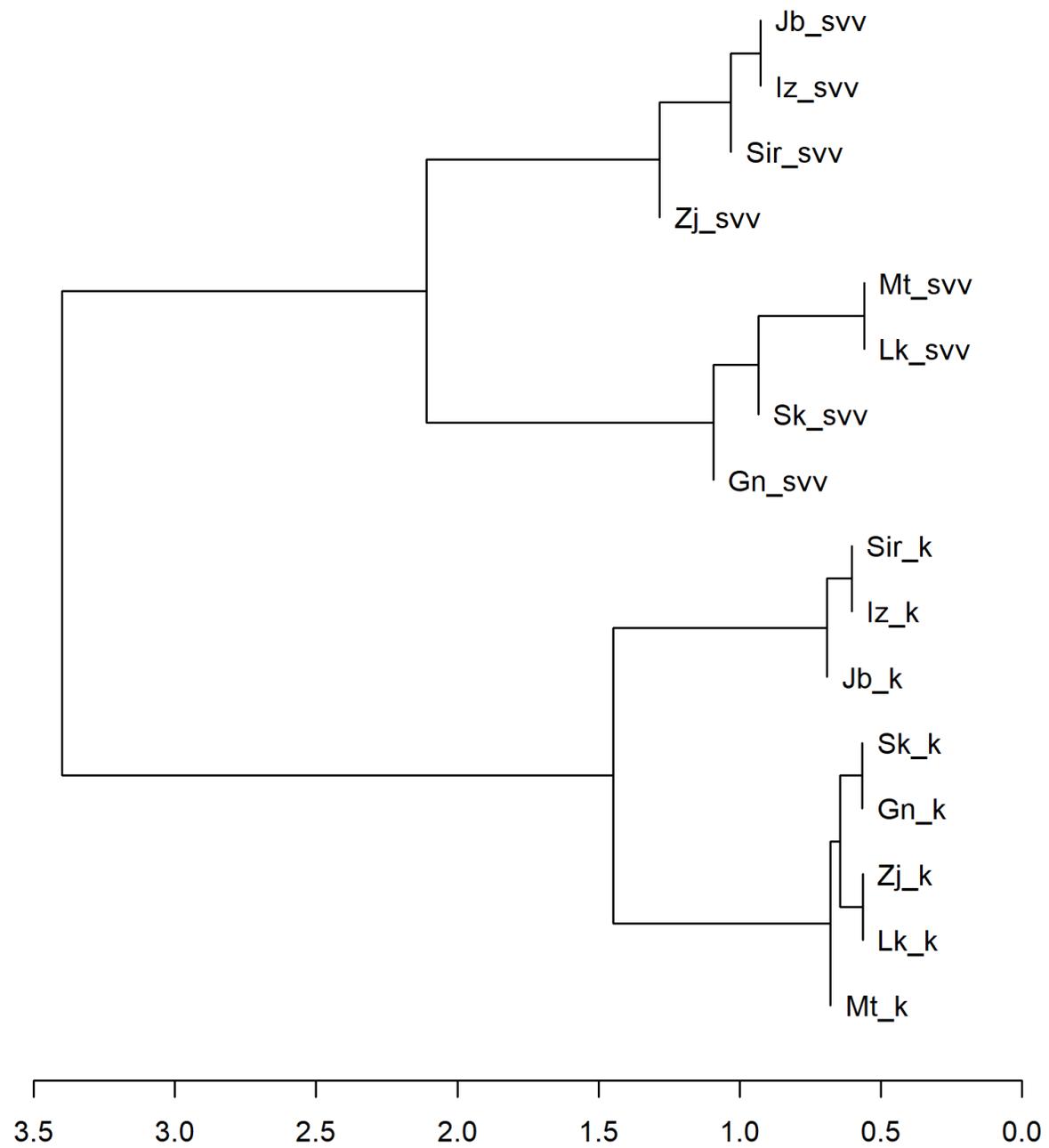
Bible svatováclavská & comments

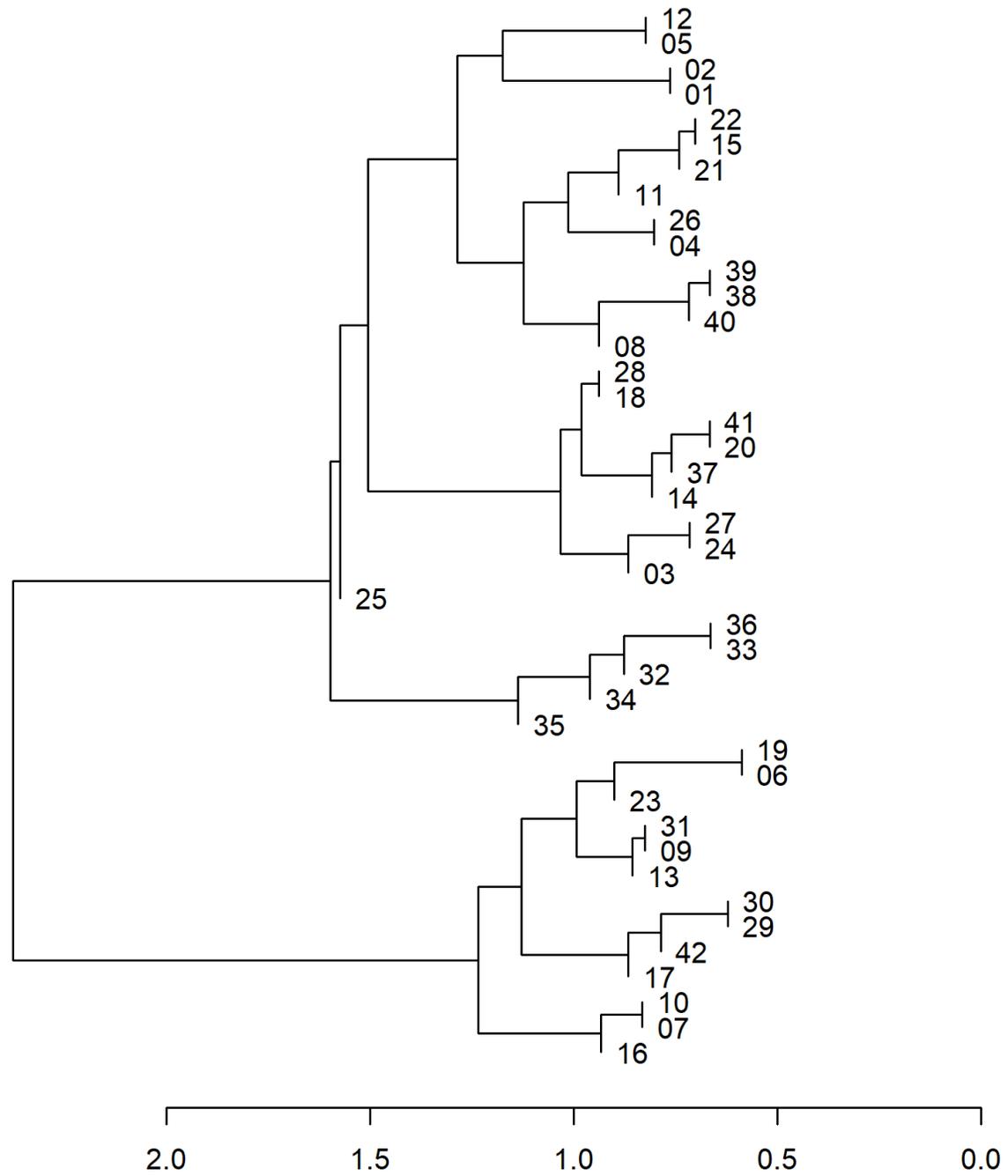
- translation
- sacred text
- editing

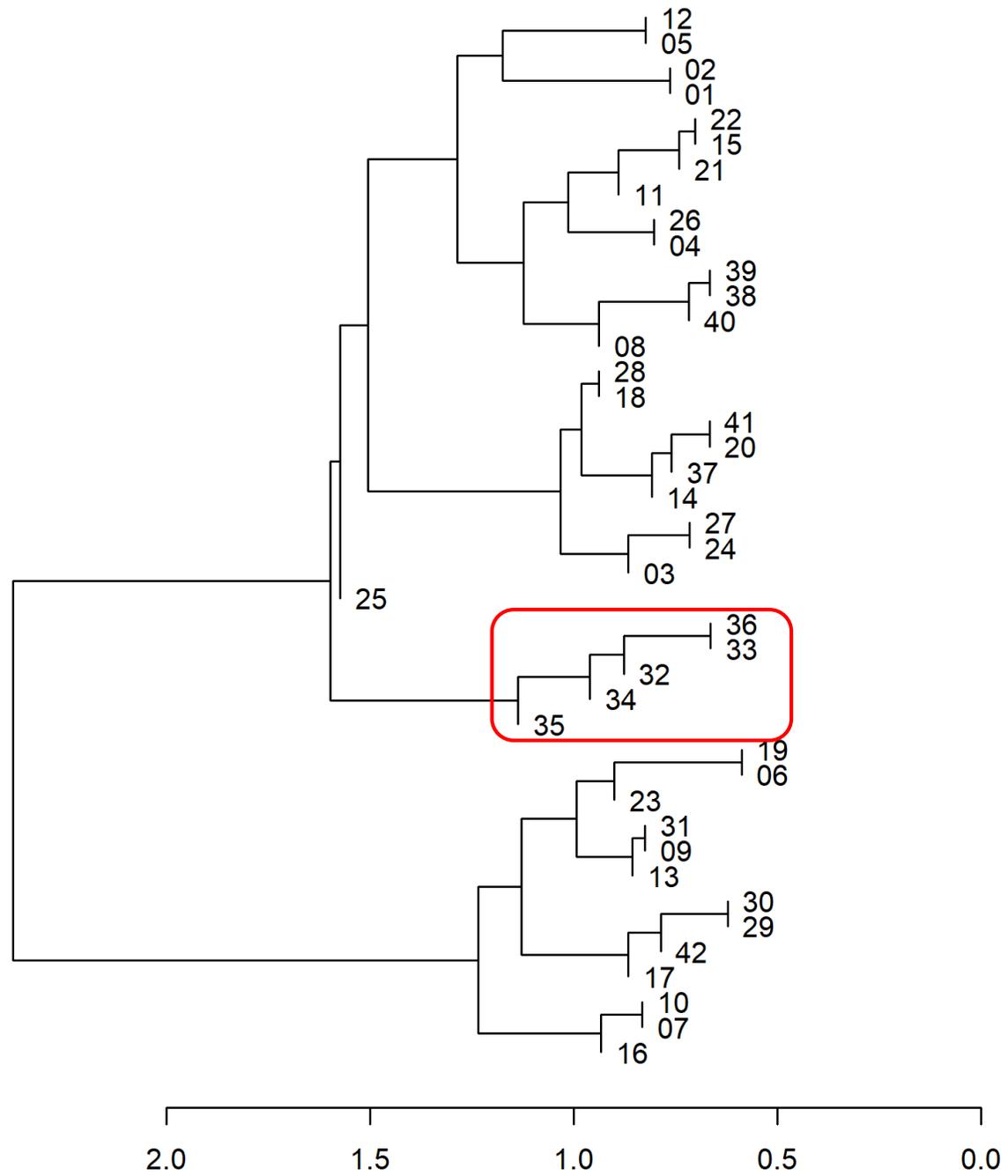
Bible svatováclavská & comments

- New Testament
 - Mt, Lk, Sk, Zj
- Old Testament
 - Gn, Jb, Iz, Sir

- 100 MFW
- culling = 0

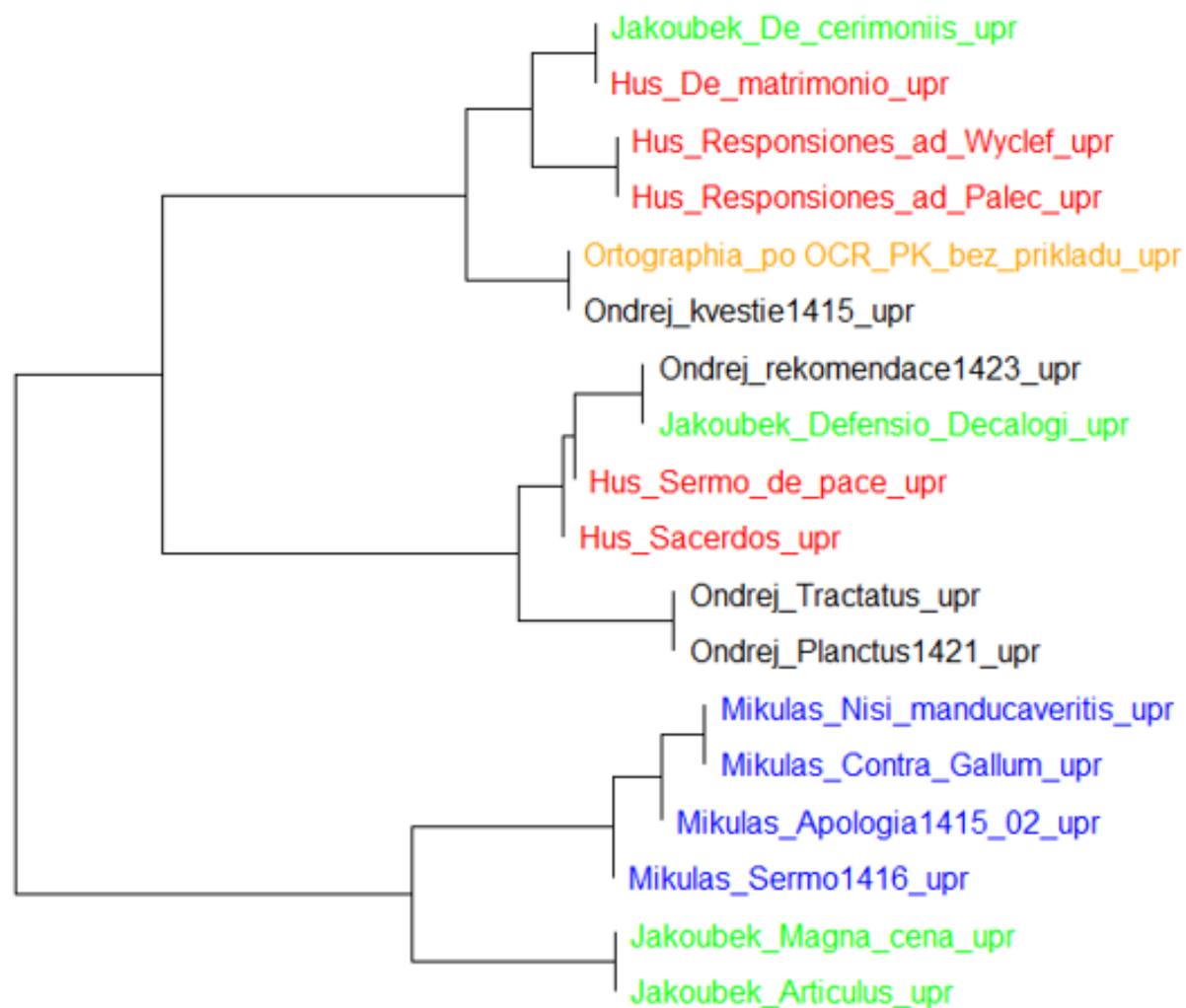






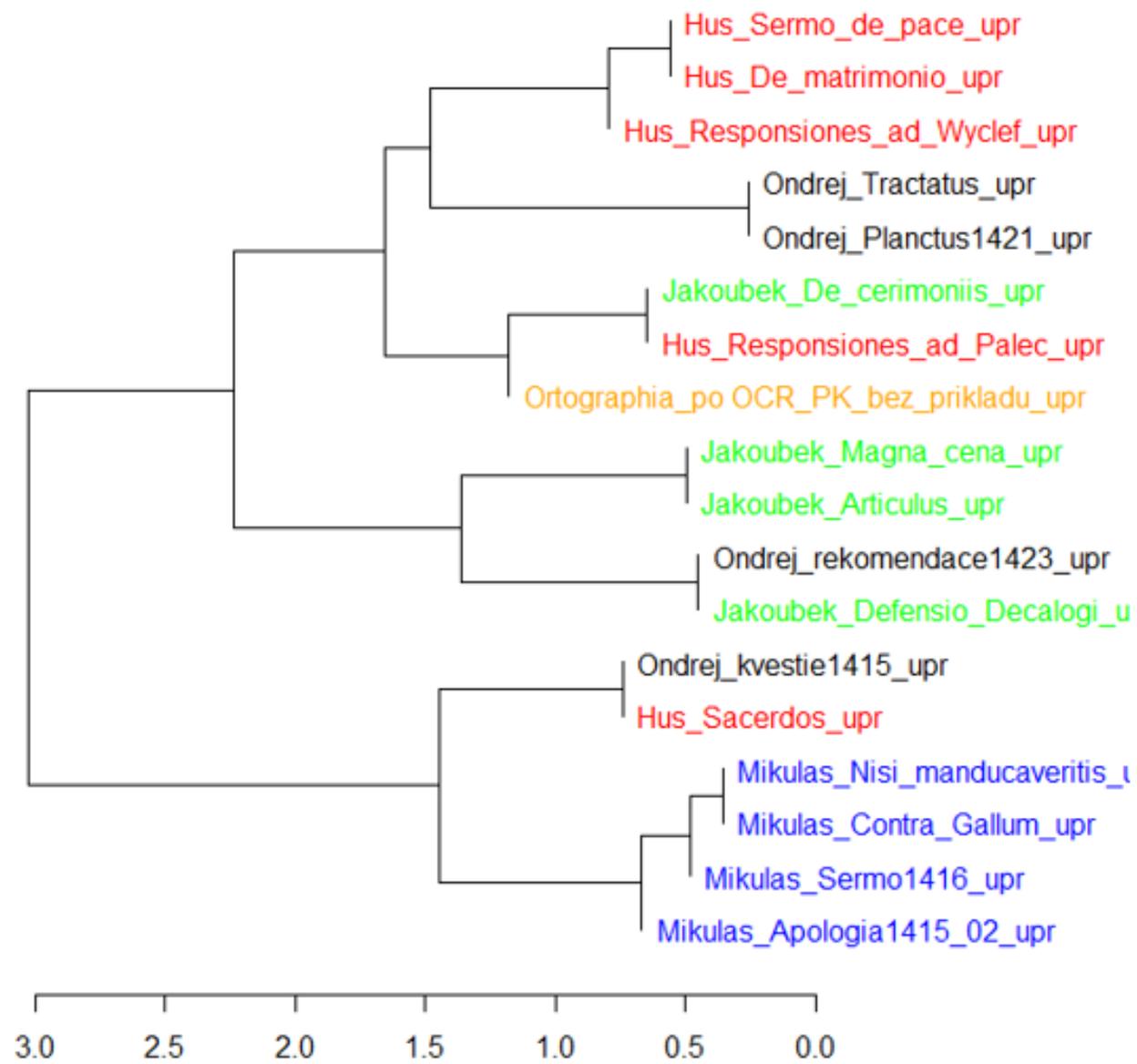
Ortographia Bohemica

- Jan Hus
 - De matrimonio, Responsiones ad Palec, Sacerdos, Sermo de pace
- Jakoubek ze Stříbra
 - Articulus, De cerimoniais, Defensio Decalogi, Magna cena
- Mikuláš z Drážd'an
 - Apologia 1415, Contra Gallum, Nisi manducaveritis, Sermo 1416,
- Ondřej z Brodu
 - Kvestie 1415, Planctus 1421, Recomendace 1423, Tractatus,
- Petr z Pulky
 - Confutatio, Epistola, Sermo ite, De congruitate



2.5 2.0 1.5 1.0 0.5 0.0

100 MFW Culled @ 30%
Distance: wurzburg



20 MFW Culled @ 30%
Distance: wurzburg

Ortographia Bohemica

- paraphrases, non-attributed quotations → analysis of 'authorial' parts
- Latin

Syntactic functions

- Čech, R., Mačutek, J., Kubát, M., Koščová, M. (2022). Does an author leave a syntactic footprint? In: Misuraca, M., Scepi, G., Spano, M. Proceedings of the 16th International Conference on statistical analysis of textual data. Naples: Vadistat Press, 221-228.
- presidential speeches: 1935–2018

Syntactic functions

- syntactic annotation (UD pipe)
- syntactic relations
 - UD → 37 universal syntactic relations
 - subject, object, indirect object, nominal modifier, adverbial modifier etc.

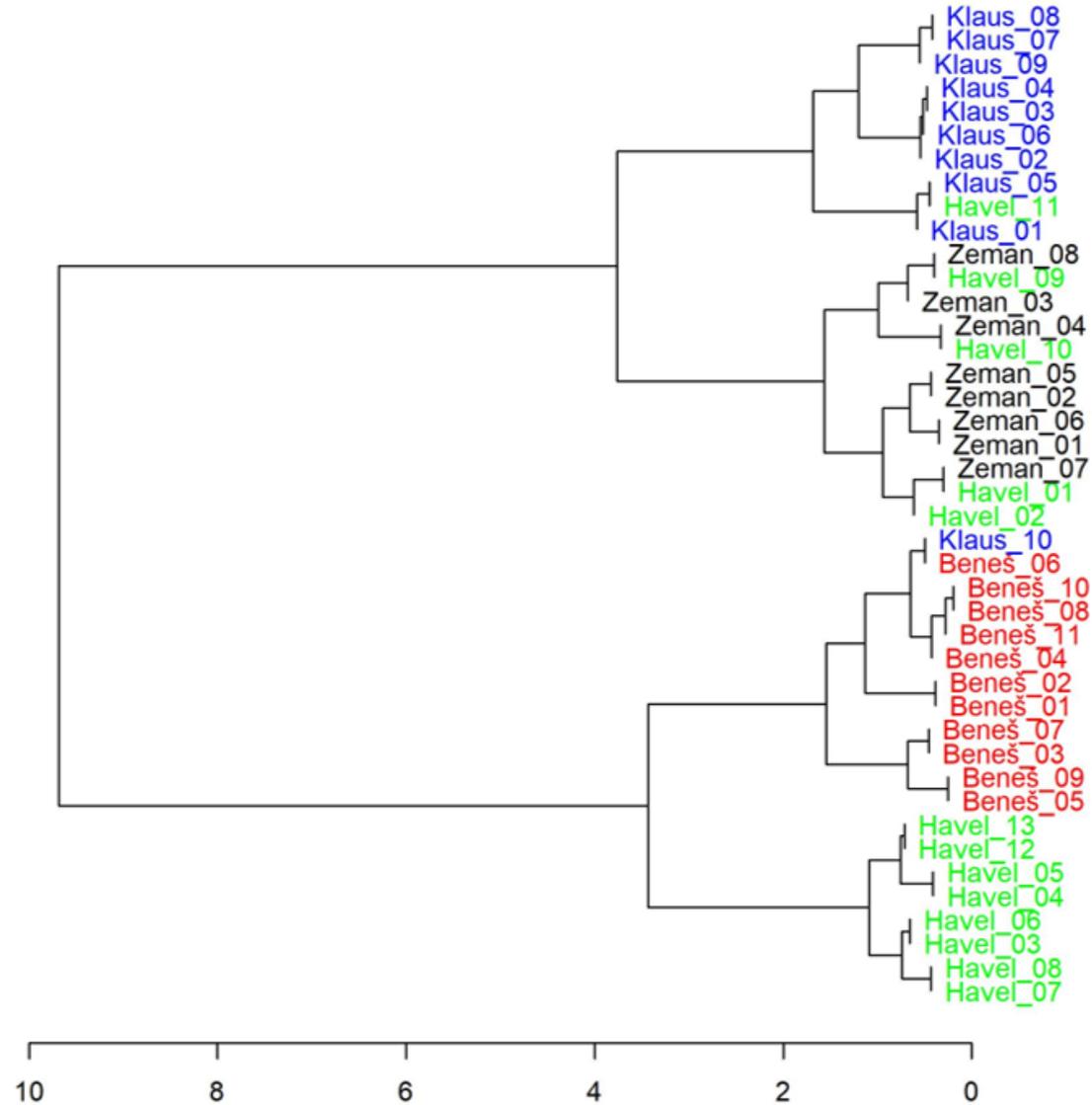
Syntactic functions

- Milí spoluobčané, když jsem vás před pěti lety v tento den poprvé oslovil, řekl jsem, že do našich srdcí se opět vrátila naděje
- amod nsubj mark aux obj case nummod obl case det obl obl advcl
root aux mark case det obl expl:pv advmod ccomp obj

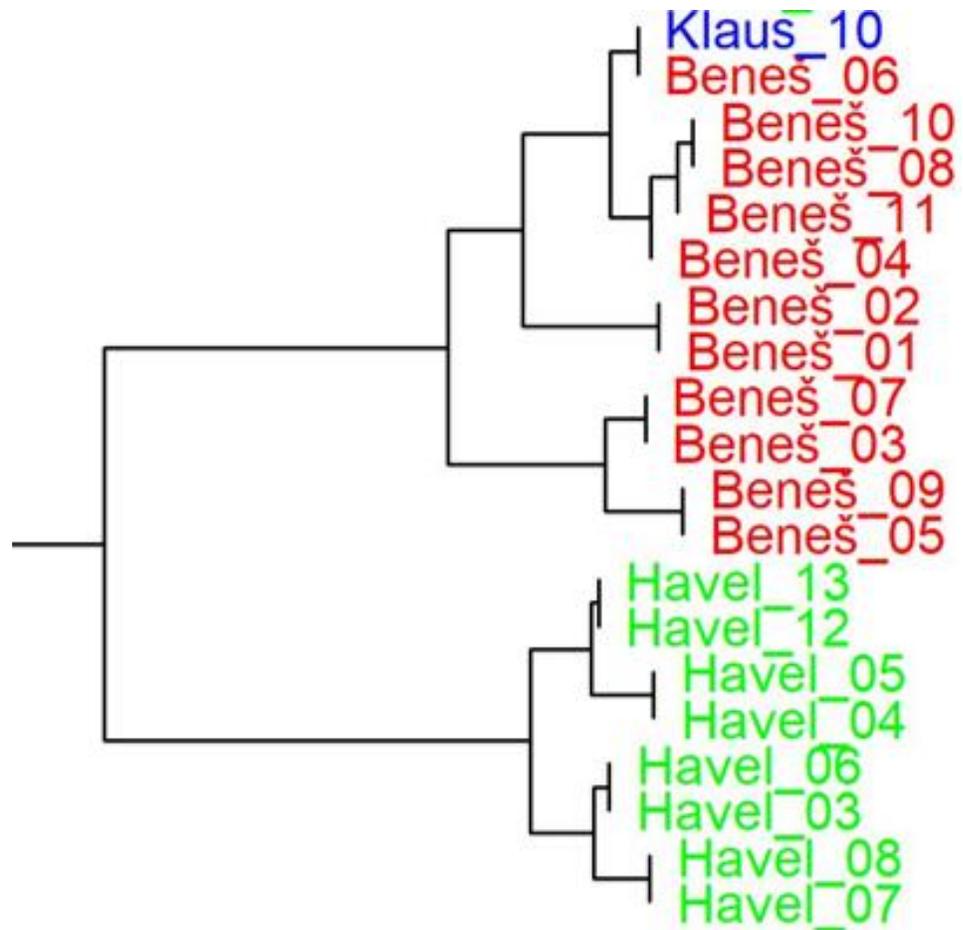
Syntactic functions

- democratic addresses
 - 22 most frequent syntactic relations
- communist addresses
 - 18 most frequent syntactic relations
- configuration
 - culling = 0
 - n-gram size = 1

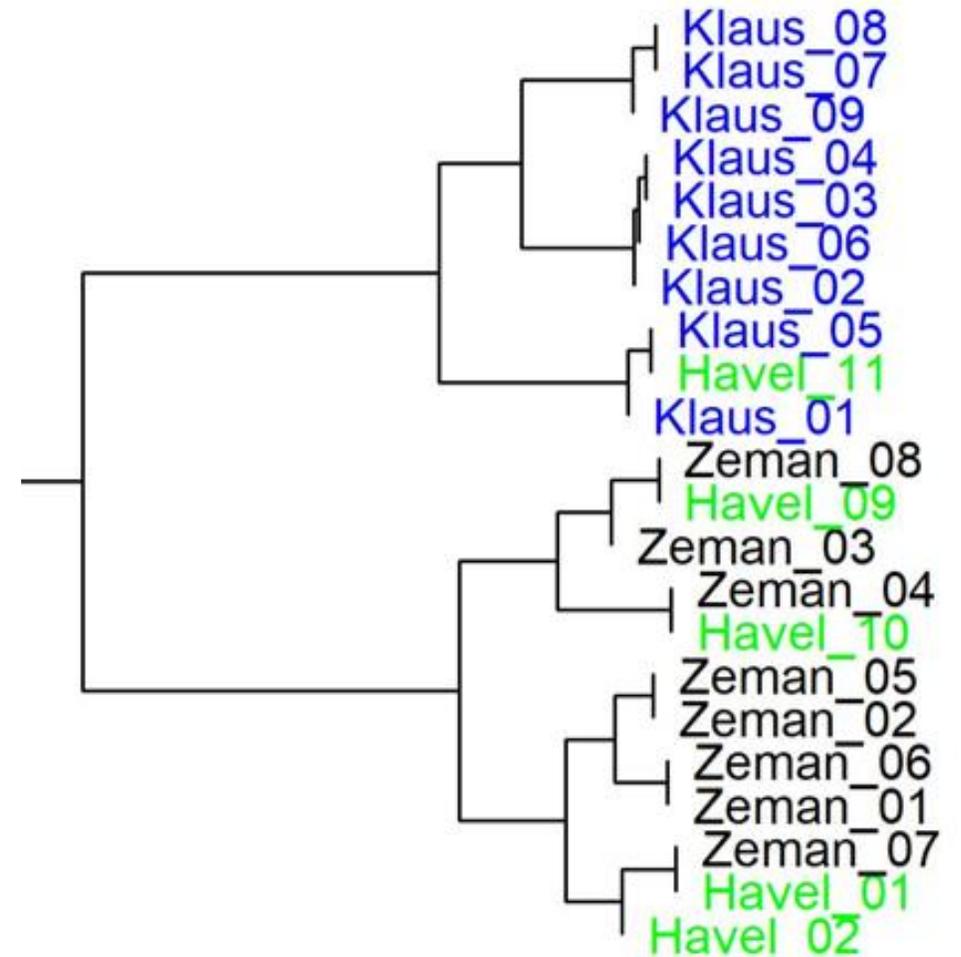
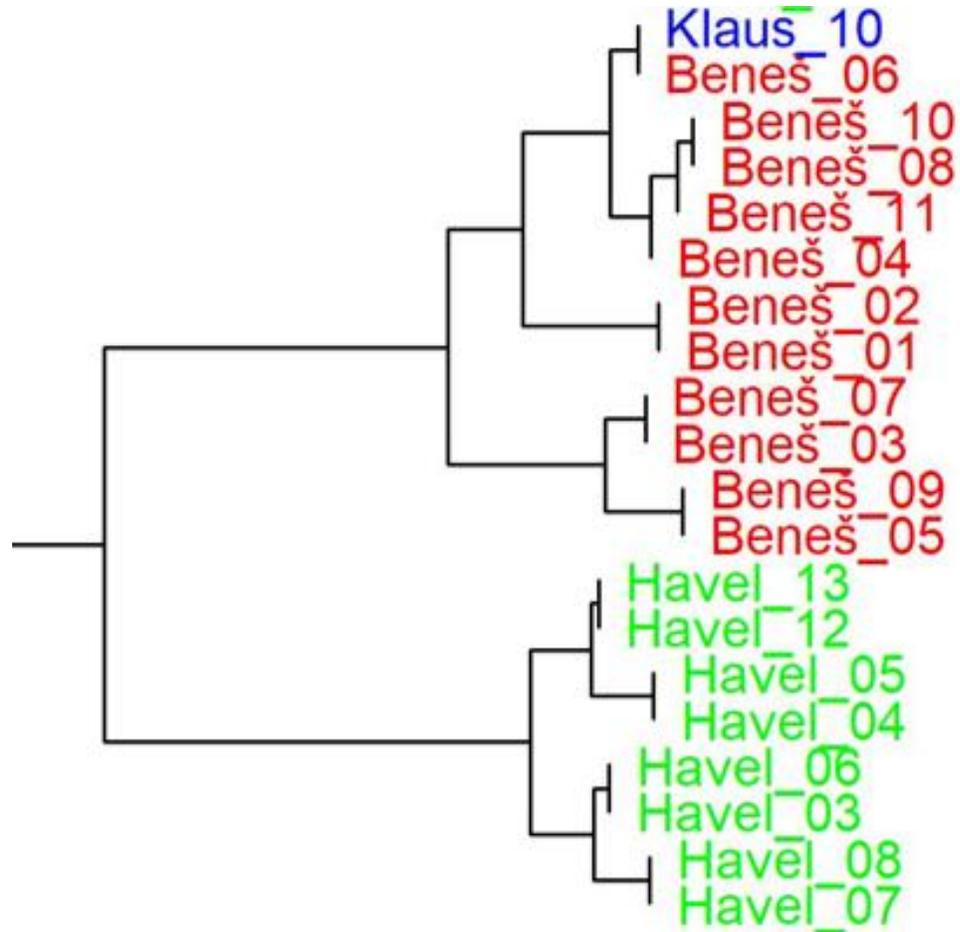
Results – democratic presidents



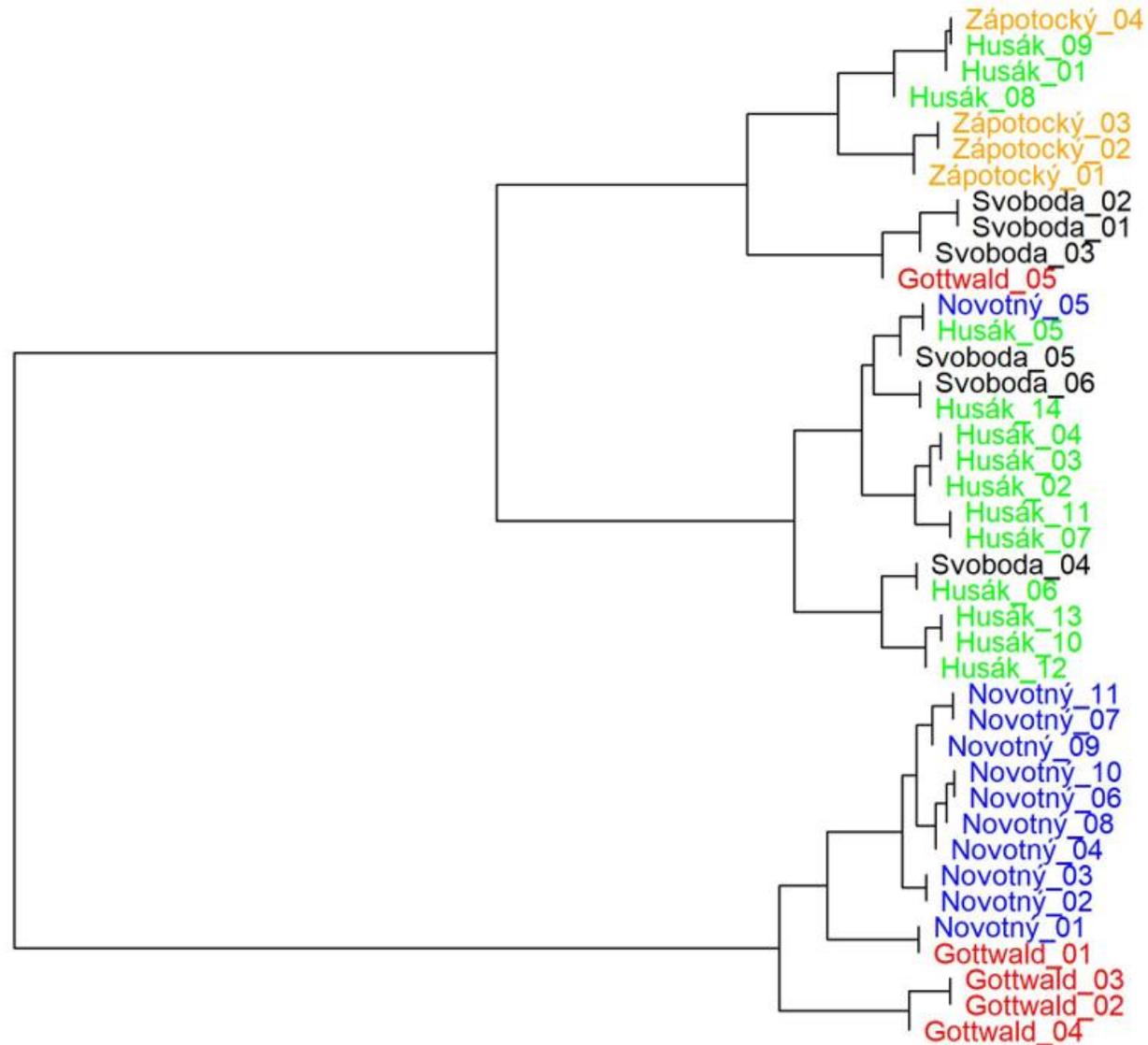
Results – democratic presidents



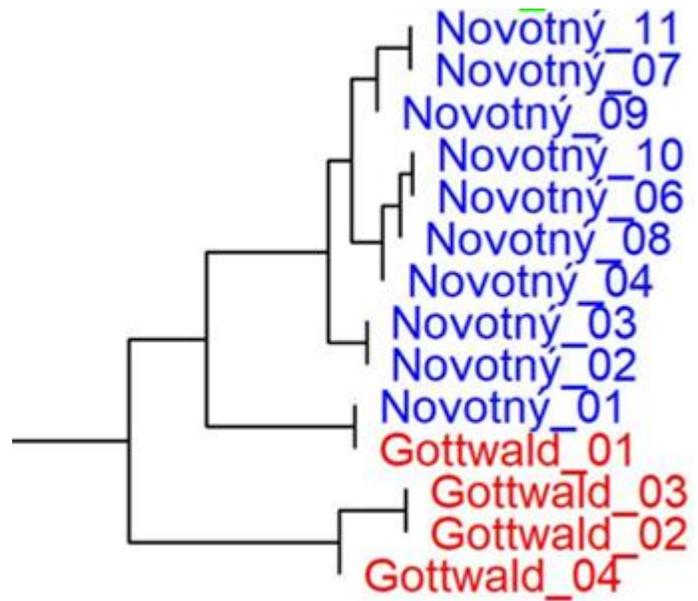
Results – democratic presidents



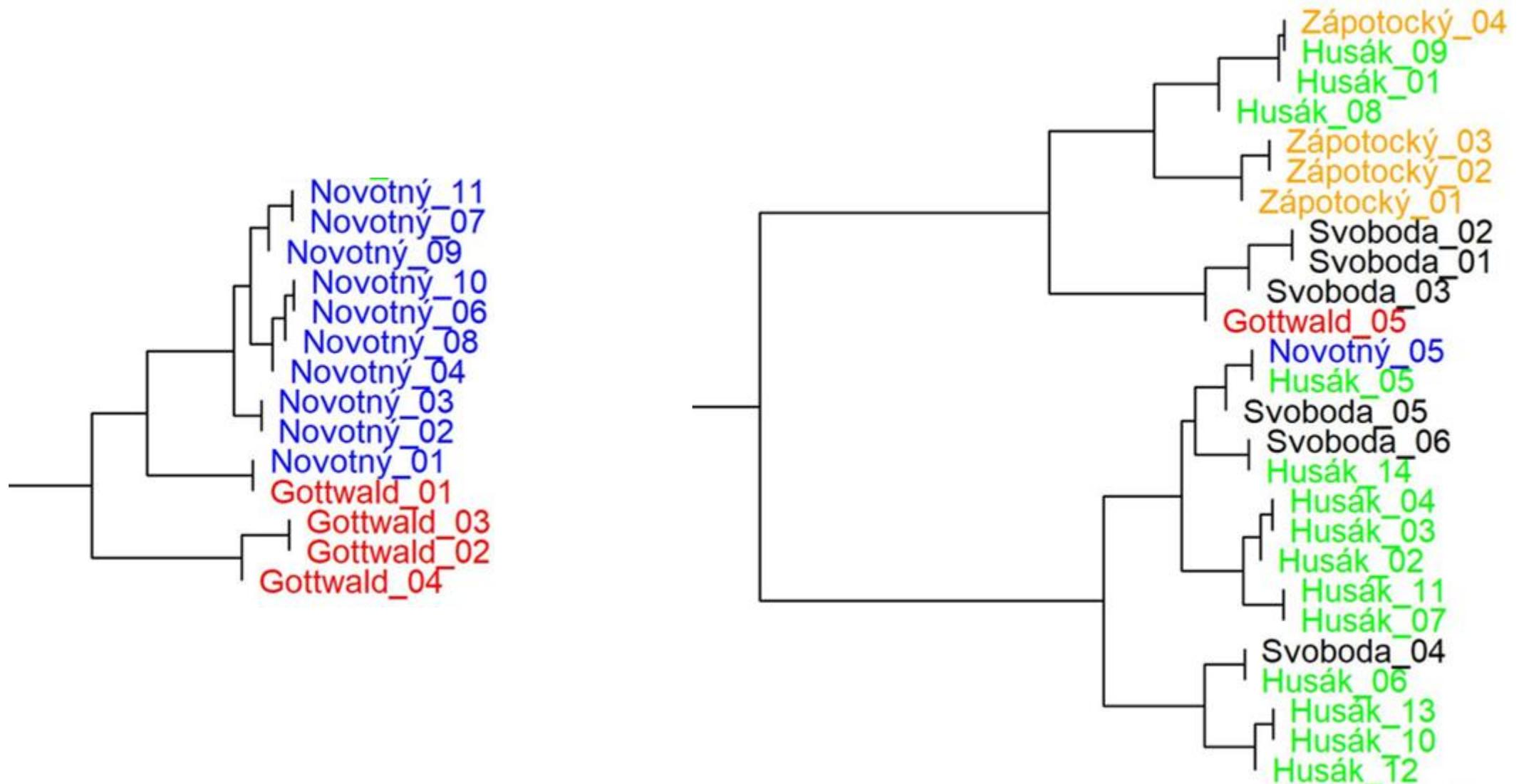
Results – communist presidents



Results – communist presidents



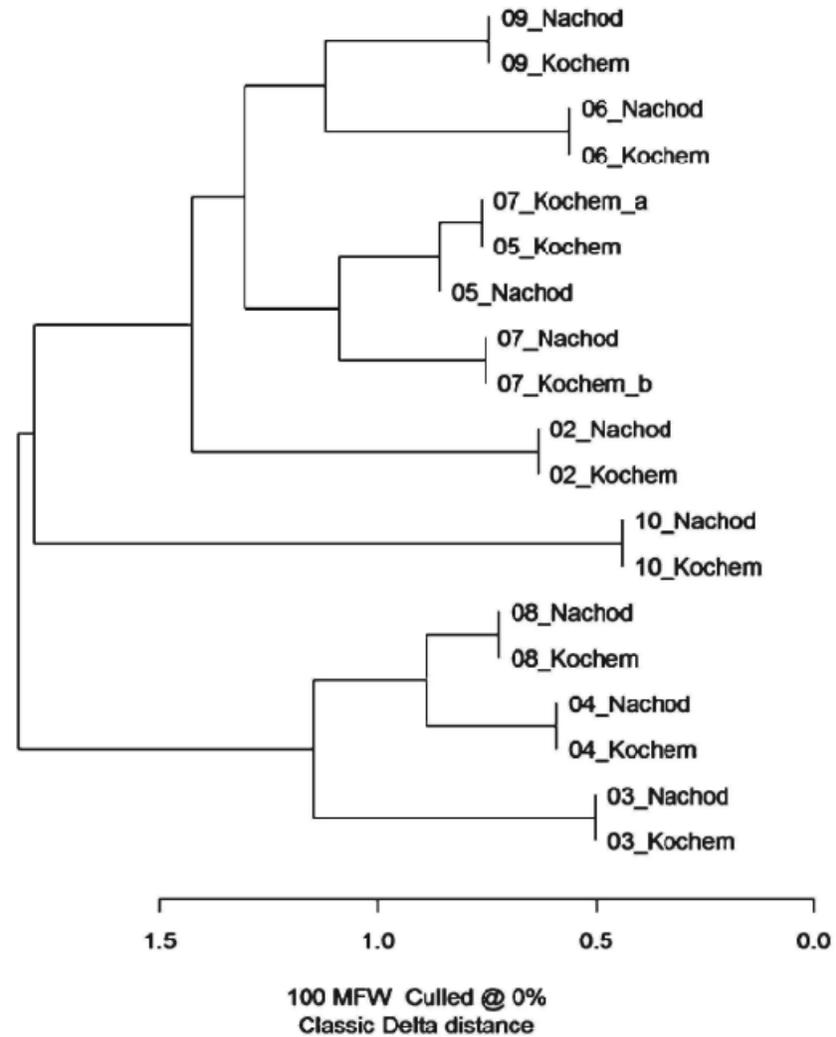
Results – communist presidents



Baroque Prayers

- Kubát, M., Netolická, Š., Čech, R., Mačutek, J. (2021). Martin of Cochem's Golden Key of Heaven and its Czech Relatives: Quantitative Analysis of Baroque Prayers. *Bohemistika*, 21, 283-294.

Baroque Prayers



Thank you
for your attention!