

# **Tematická koncentrace textu**

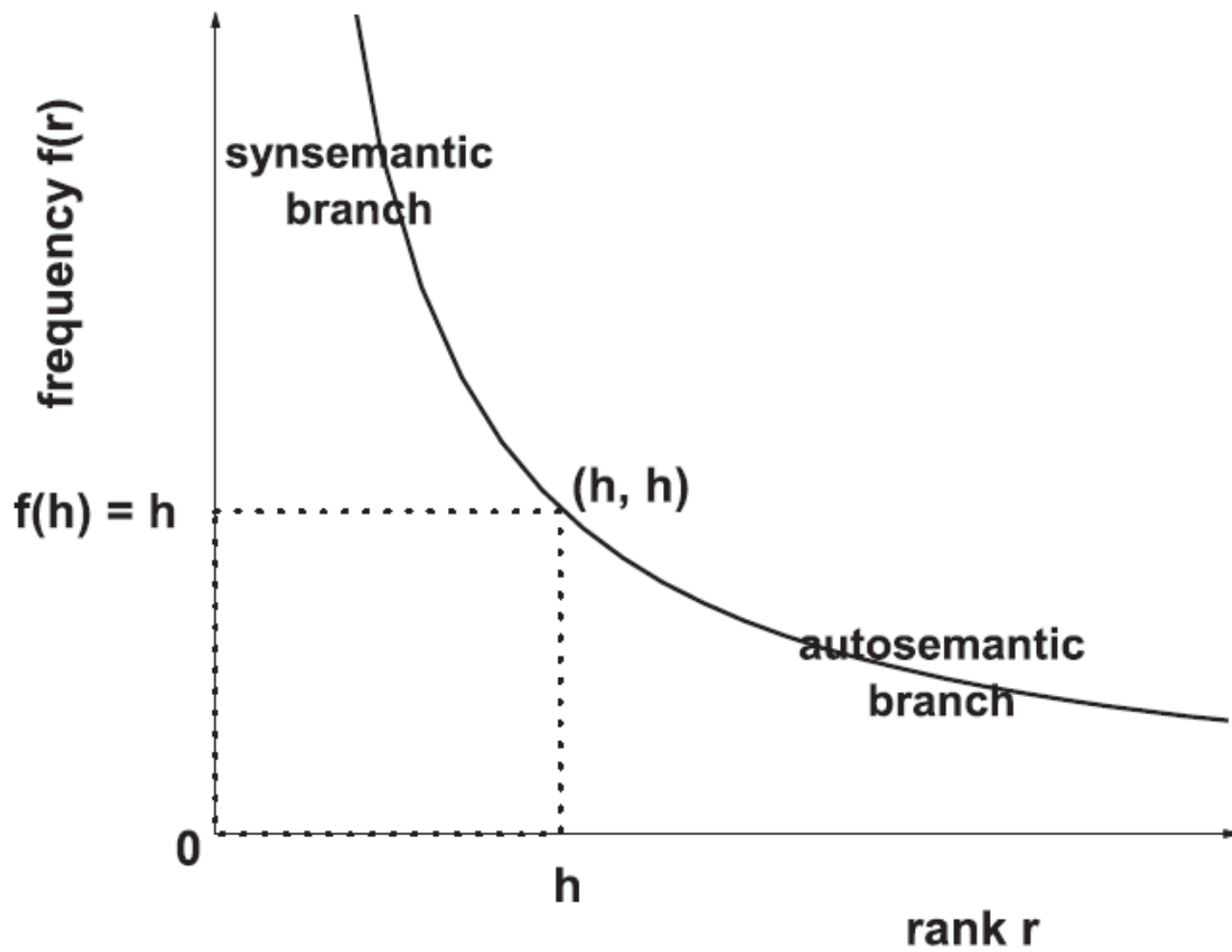
Radek Čech

# Cíle

- pomocí kvantitativní textové analýzy
  - odhalit ústřední téma textu (pokud existuje)
  - zjistit tematickou „sílu“ jednotlivých slov reprezentujících témata textu
  - zjistit tematickou koncentraci celého textu

# Metodologie

- Popescu et al. (2009)
  - tematická slova → substantiva, verba, adjektiva
  - východiskem frekvenční charakteristika textu
    - h-bod → určuje fuzzy hranici mezi synsémantickými a autosémantickými slovy



# Metodologie

- Popescu et al. (2009)
  - tematická slova → substantiva, verba, adjektiva
  - východiskem frekvenční charakteristika textu
    - h-bod → určuje fuzzy hranici mezi synsémantickými a autosémantickými slovy
  - tematická slova na h-bodem vyjadřují hlavní téma textu (lepší výsledky u lemmatizovaných textů)

# Postup

- text → J. Hus *Dcerka*
  - (text poskytl Oddělení vývoje jazyka UJČ)
- vytvoření frekvenční distribuce slovních tvarů
  - AntConc
    - [http://www.antlab.sci.waseda.ac.jp/antconc\\_index.html](http://www.antlab.sci.waseda.ac.jp/antconc_index.html)

- Corpus Files
- HusDcerkaH.txt

Concordance Concordance Plot File View Clusters Collocates Word List Keyword List

Hits Total No. of Word Types: 2416 Total No. of Word Tokens: 8574

Rank	Freq	Word	Lemma Word Form(s)
1	503	a	
2	200	jest	
3	176	v	
4	174	sě	
5	174	že	
6	112	neb	
7	98	i	
8	98	tak	
9	88	aby	
10	78	k	
11	71	člověk	
12	70	to	
13	58	když	
14	57	má	
15	53	bóh	
16	52	protož	
17	51	jenž	
18	49	od	

Search Term  Words  Case  Regex    Display Options  Treat all data as lowercase

Total No.     Sort by

Files Processed    Invert Order

# Frekvenční distribuce

pořadí	f	slovo	pořadí	f	slovo	pořadí	f	slovo
1	503	a	12	70	to	<b>23</b>	<b>41</b>	<b>die</b>
2	200	jest	13	58	když	24	41	ho
3	176	v	14	57	má	25	41	na
4	174	sě	<b>15</b>	<b>53</b>	<b>bóh</b>	26	39	ale
5	174	že	16	52	protož	27	38	jeho
6	112	neb	17	51	jenž	28	37	z
7	98	i	18	49	od	<b>29</b>	<b>36</b>	<b>boha</b>
8	98	tak	19	48	své	30	35	jako
9	88	aby	20	47	by	31	33	jakož
10	78	k	21	47	s	32	33	ty
<b>11</b>	<b>71</b>	<b>člověk</b>	22	43	li	<b>33</b>	<b>30</b>	<b>d'áble</b>



# H-bod

$$h = f(h)$$

$$h = \frac{f(i)j - f(j)i}{j - i + f(i) - f(j)}$$

$$h_{Dcerka} = \frac{33 \cdot 33 - 30 \cdot 32}{33 - 32 + 33 - 30} = 32.25$$

# Index tematické koncentrace

- $r'$  .... pořadí slova před h-bodem
- tematická váha:  $h-r'$ 
  - slova s nejnižším rankem mají největší váhu
- relativizace
  - tematická váha daného slova vynásobená jeho vlastní frekvencí je podělena součtem všech vah a nejvyšší možnou frekvencí → index tematické koncentrace slova

$$TC = 2 \sum_{r'=1}^T \frac{(h-r')f(r')}{h(h-1)f(1)}$$

# Postup

- výpočet normalizační konstanty  $C$

$$\frac{2}{h(h-1) f(1)} = \frac{2}{32.25(31.25) 503} = 0.00000394532 = C$$

- člověk ( $r' = 11, f = 71$ )

$$(h - r') f(r') C = (32.25 - 11) 71 \cdot C = 0.005952502$$

# Tematická koncentrace textu

slovo	$r$	$f$	$(h-r')f(r')C$
člověk	11	71	0.00592
bóh	15	53	0.00361
die	23	41	0.00150
boha	29	36	0.00046
$TC$			0.01152

$$tcu = 1000(TC)$$

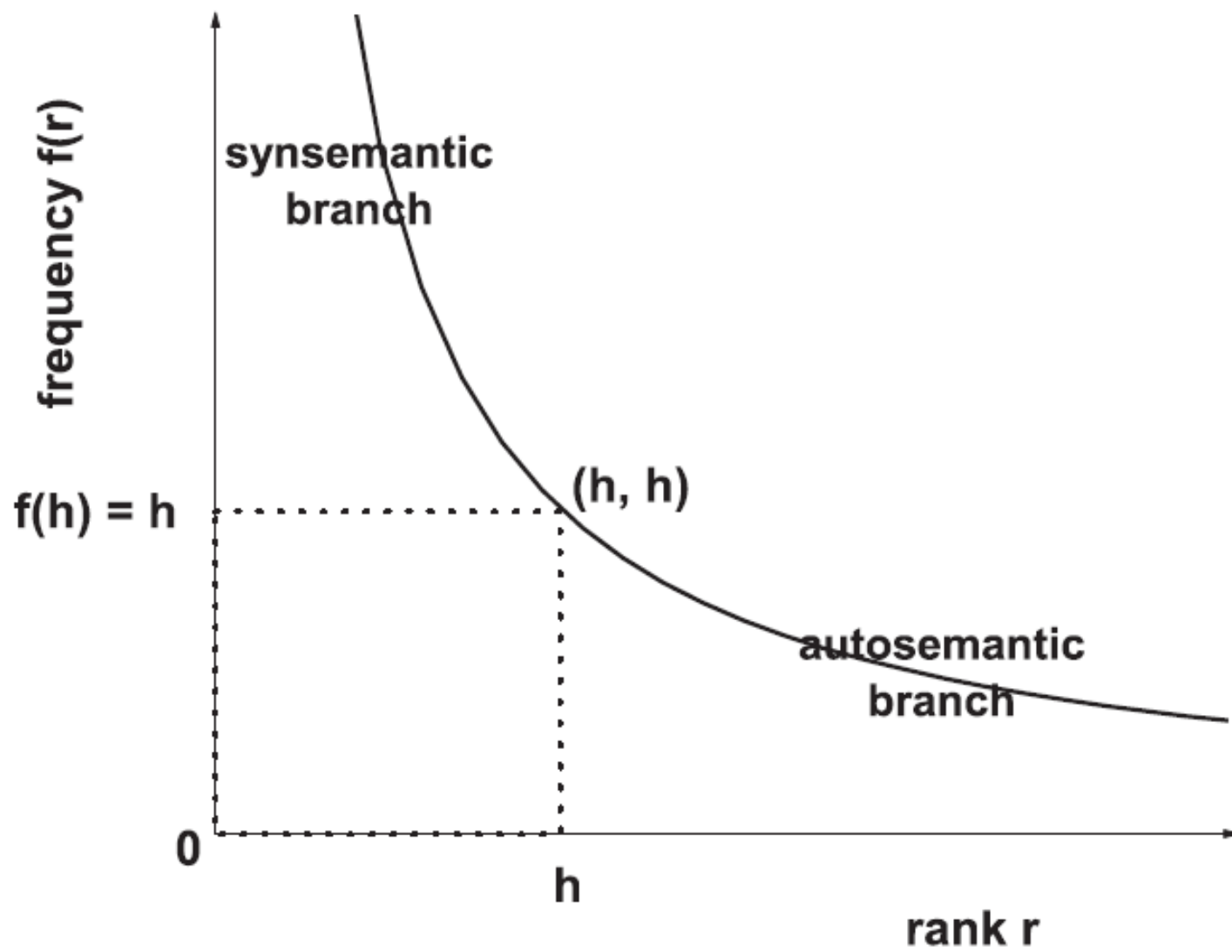
$$tcu_{Dcerka} = 1000(0.01152) = 11.52$$

# Tematická koncentrace textu

- k čemu je to dobré?
  - jednoduchá metoda umožňující zjistit, o čem daný text je a „jak moc“ se týká centrálního tématu (případně témat)
  - kvantifikace umožňuje porovnávat texty (za pomoci statistických testů)
- nedostatek
  - pokud žádné tematické slovo nad h-bodem,  $TC=0$   
→ všechny texty s  $TC=0$  jsou si rovny

# Index negativní TC

- východisko
  - h-bod fuzzy hranice
  - pokud text tematicky „roztříštěn“ → tematická slova mají relativně malou frekvenci
  - čím je text tematicky „roztříštěnější“, tím větší vzdálenost mezi h-bodem a prvním tematickým slovem (pořadí dáno frekvencí)



# Index negativní TC

- postup
  - u textů  $TC=0$  hledáno první tematické slovo za h-bodem

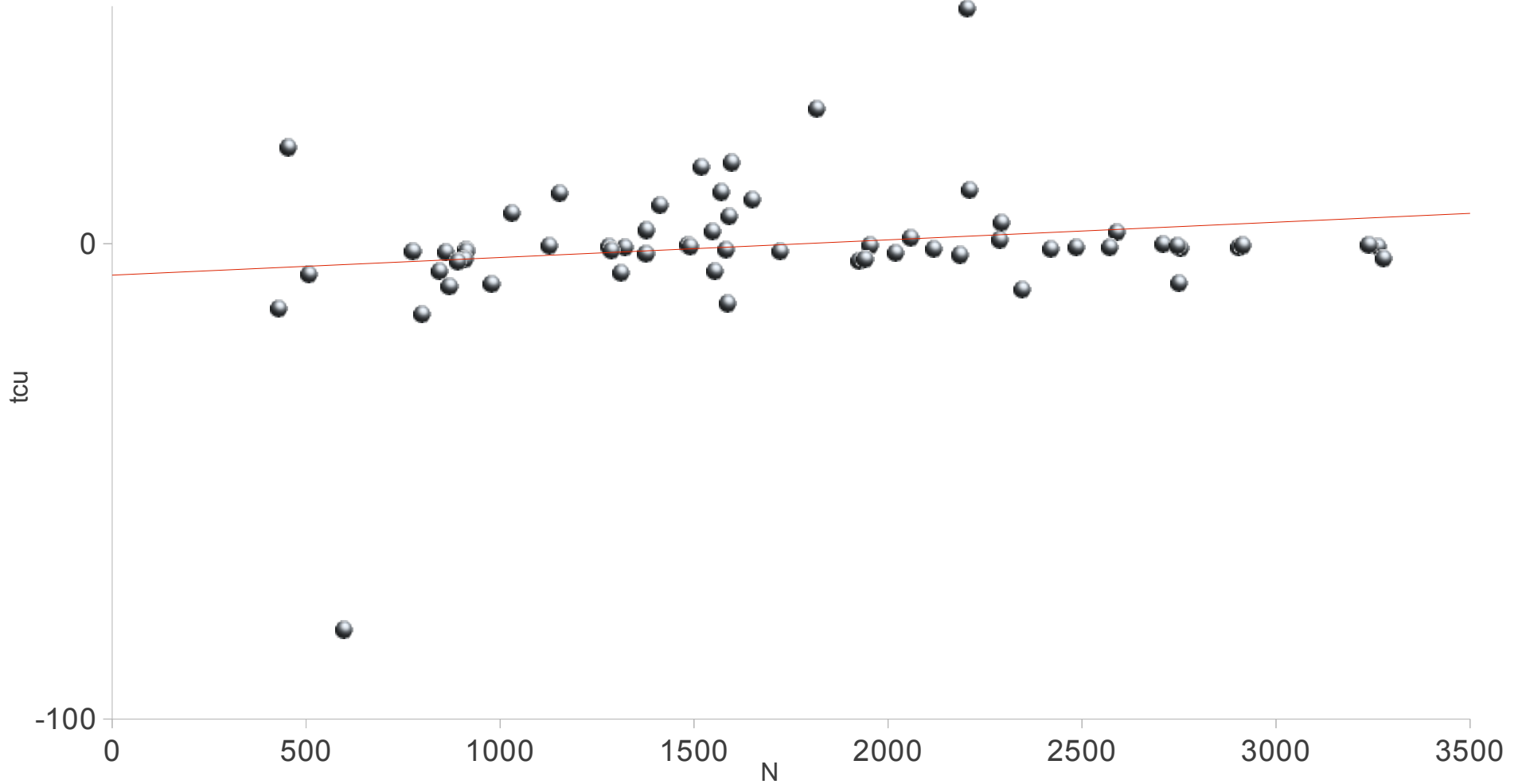
$$TC_{\text{neg}} = (h - r'') \frac{h}{f(r'')} C$$

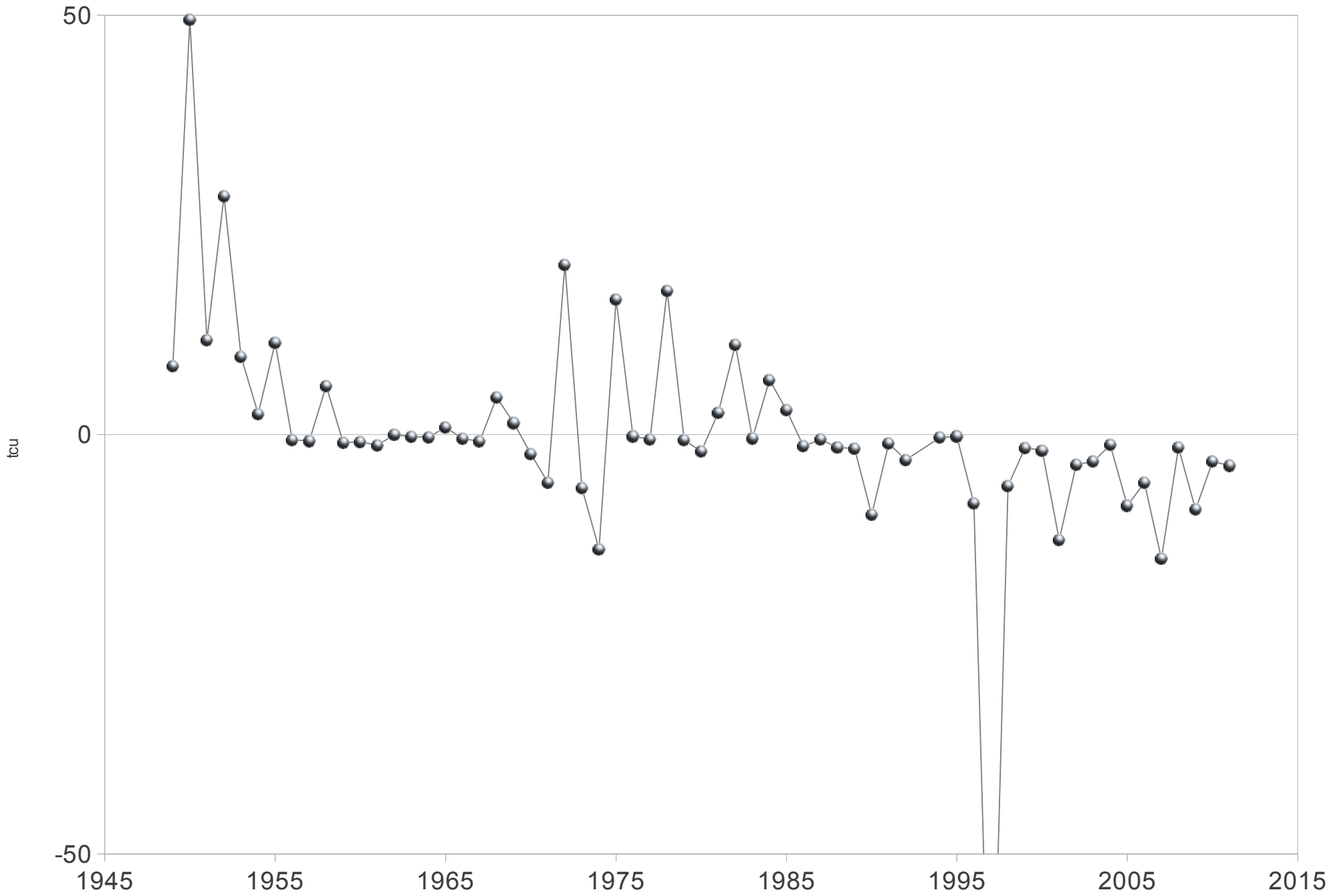


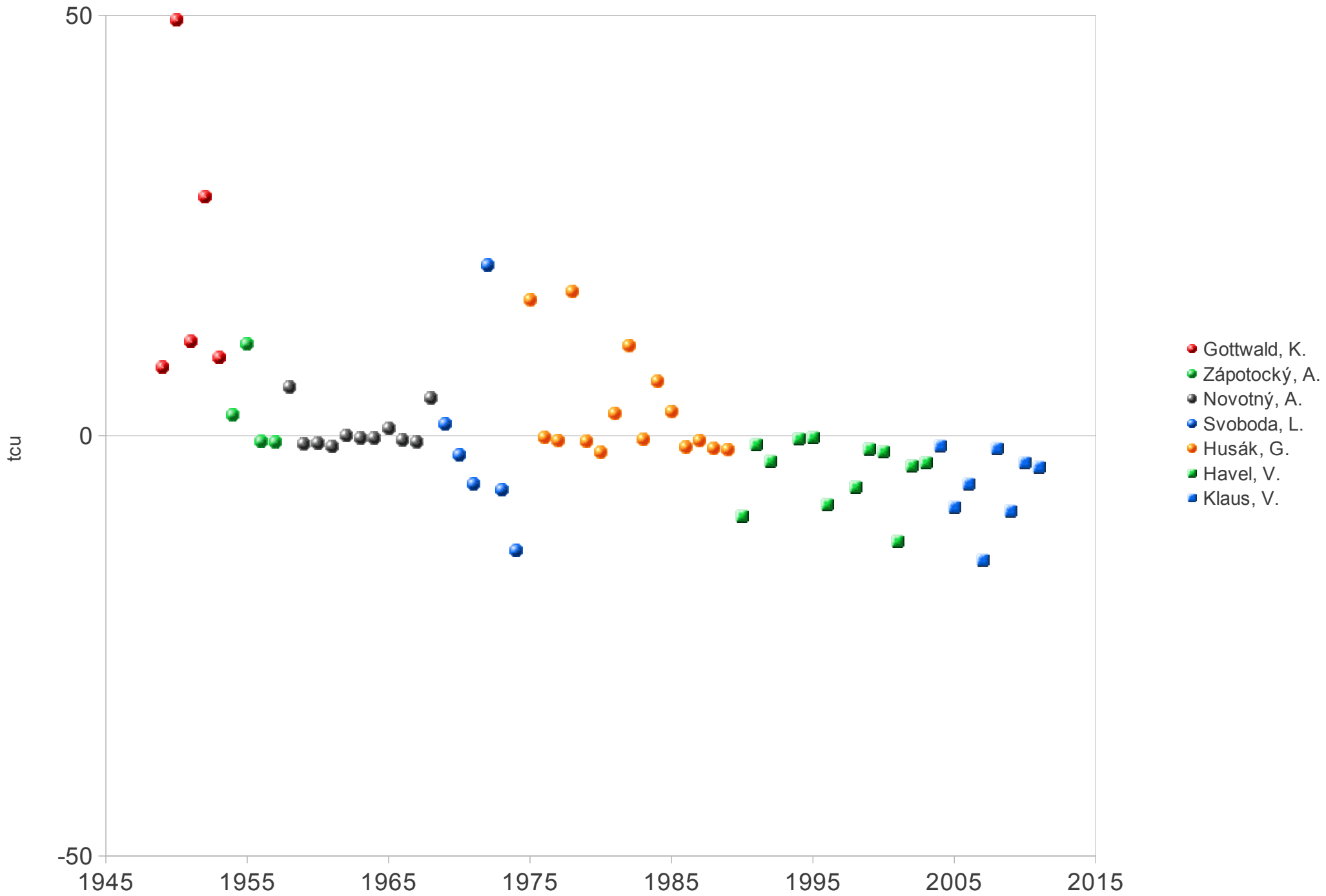
# Novoroční projevy československých a českých prezidentů

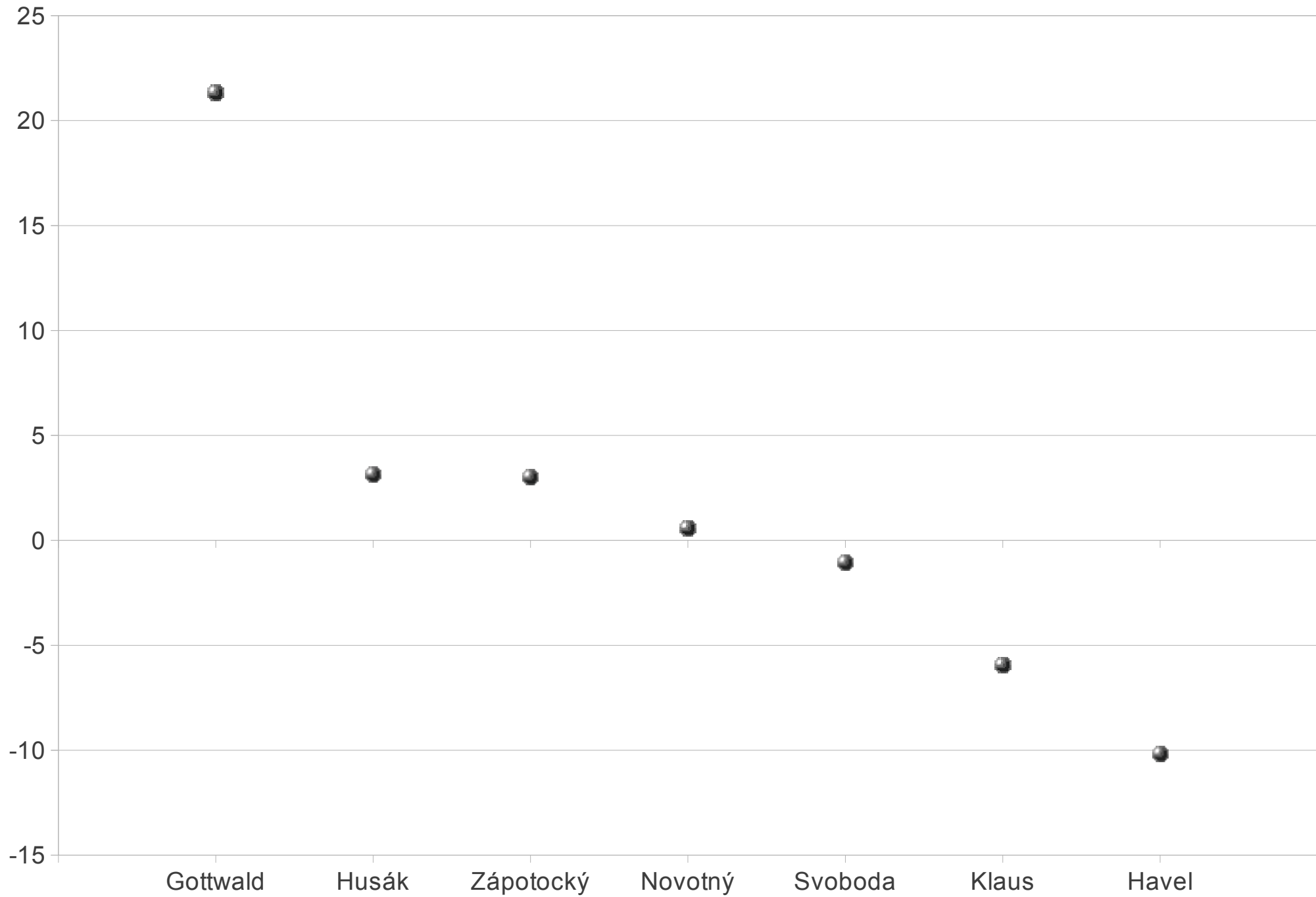
- 1949 – 2011 (kromě r. 1993)
- relativně stabilní žánr

# Délka textu versus TC

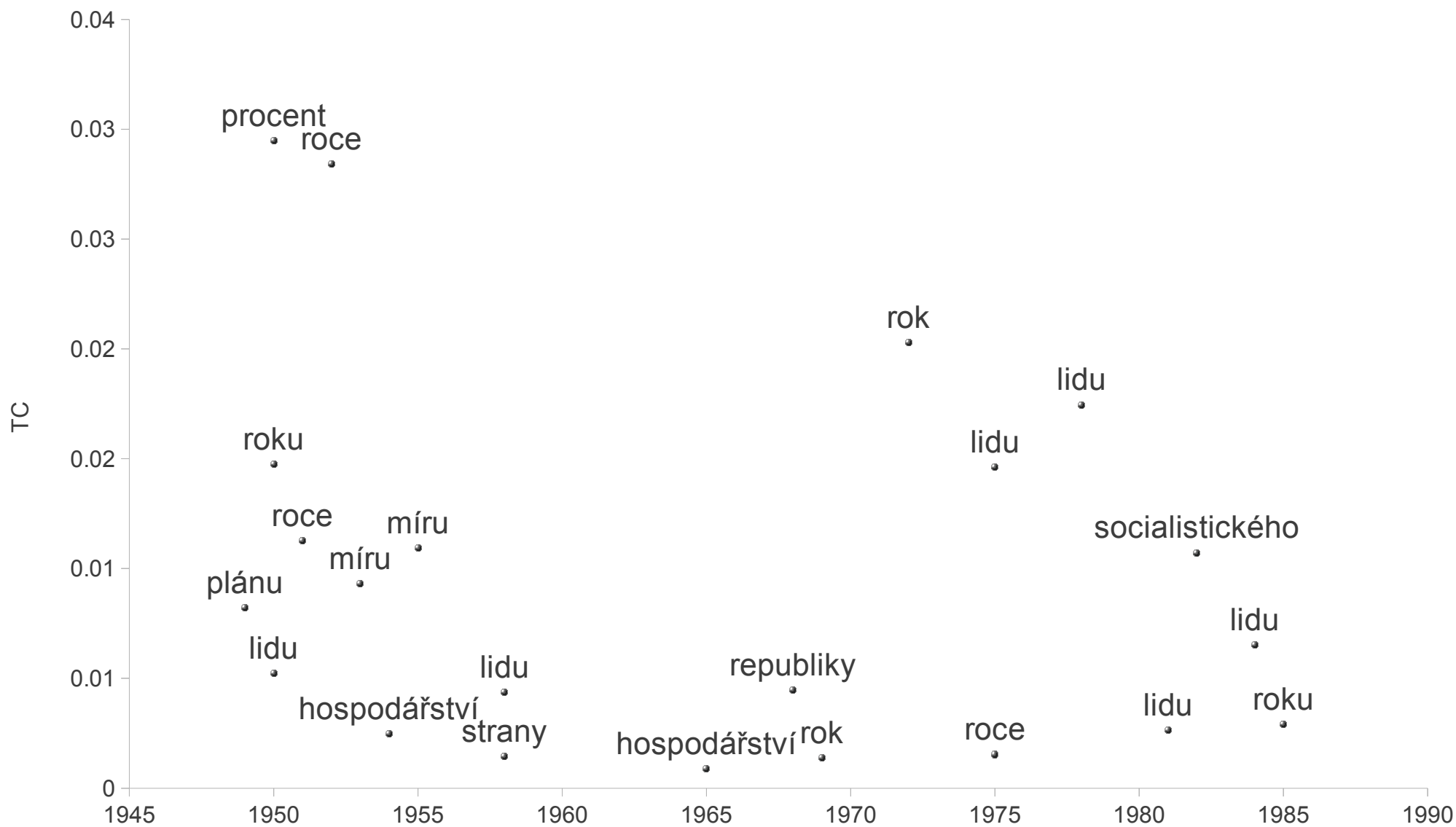


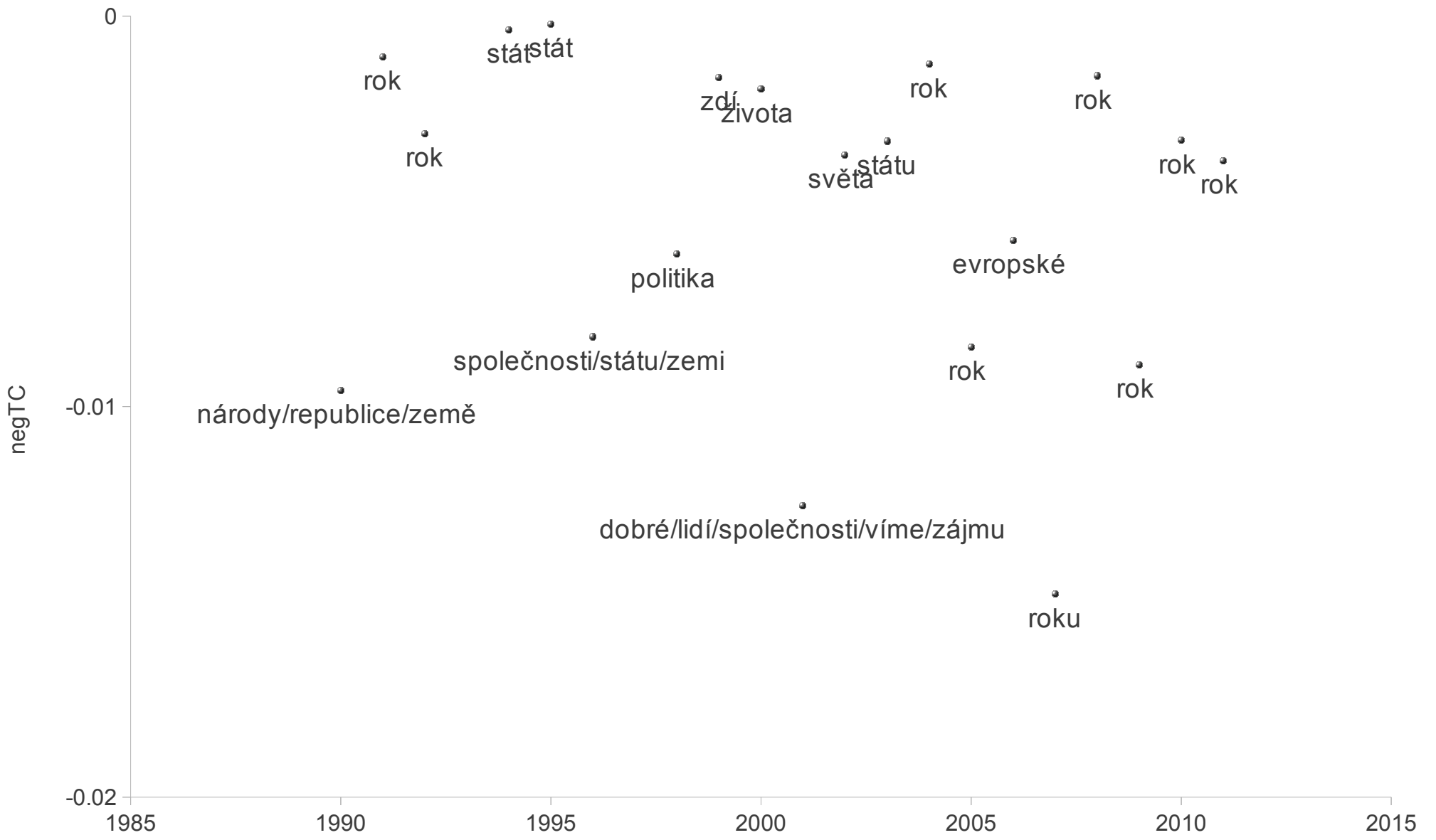






# Tematická slova





# Tematická slova

- AntConc
  - možnost detailní analýzy textu

The screenshot displays the AntConc 3.2.2.1u (Linux OS) 2011 interface. The main window shows a concordance search for the term "plánů" in the file "Gottwald\_1949.txt". The search results are displayed in a table with columns for Hit, KWIC, and File. The search term is "plánů", and there are 12 concordance hits. The search window size is set to 50. The interface also shows the search term "plánů" in the search box, the search options (Words, Case, Regex), and the search window size. The search results are displayed in a table with columns for Hit, KWIC, and File. The search term is "plánů", and there are 12 concordance hits. The search window size is set to 50. The interface also shows the search term "plánů" in the search box, the search options (Words, Case, Regex), and the search window size.

Hit	KWIC	File
1	yly splněny úkoly dvouletého <b>plánů</b> v průmyslu, vyjímaje průmysl výž	Gottwald_1949.txt
2	celkovém splnění dvouletého <b>plánů</b> nezměněn. V posledním čtvrtletí	Gottwald_1949.txt
3	eří se o splnění a překonání <b>plánů</b> zasloužili. Neméně zdatně sp	Gottwald_1949.txt
4	Jedním z nejslabších článků <b>plánů</b> bylo stavebnictví, kde až do úno	Gottwald_1949.txt
5	mí kapitalisté, kteří plnění <b>plánů</b> sabotovali. Převzali jsme od nic	Gottwald_1949.txt
6	y zlepšit v roce 1948 plnění <b>plánů</b> , i když s ním nemůžeme být ještě	Gottwald_1949.txt
7	ktelé splněním a překročením <b>plánů</b> umožnily nákup potřebných potrav	Gottwald_1949.txt
8	půda obdělána a oseta podle <b>plánů</b> až na brambory, kde nedostatek s	Gottwald_1949.txt
9	tek sadby způsobil nesplnění <b>plánů</b> . Potěšitelné je zvyšování počtu	Gottwald_1949.txt
10	anizace byl během dvouletého <b>plánů</b> plněn dobře. Dalším pokračováním	Gottwald_1949.txt
11	ás všech. Úspěšné splnění <b>plánů</b> v průmyslu a dopravě umožnilo ud	Gottwald_1949.txt
12	v boji o splnění pětiletého <b>plánů</b> . Nový rok přinese další upevnění	Gottwald_1949.txt



# Výhledy

- navrhnout statistické testy pro testování rozdílu mezi
  - jednotlivými texty
  - obdobími (např. totalita x demokracie)
  - autory

Děkuji za pozornost!

[radek.cech@osu.cz](mailto:radek.cech@osu.cz)  
[www.cechradek.cz](http://www.cechradek.cz)