



# DH@Madrid Summer School 2016

Quantitative research on versification:  
The Corpus of Czech Verse

Tecnologías digitales aplicadas estudio de la poesía

Petr Plecháč  
[plechac@ucl.cas.cz](mailto:plechac@ucl.cas.cz)



European  
Research  
Council

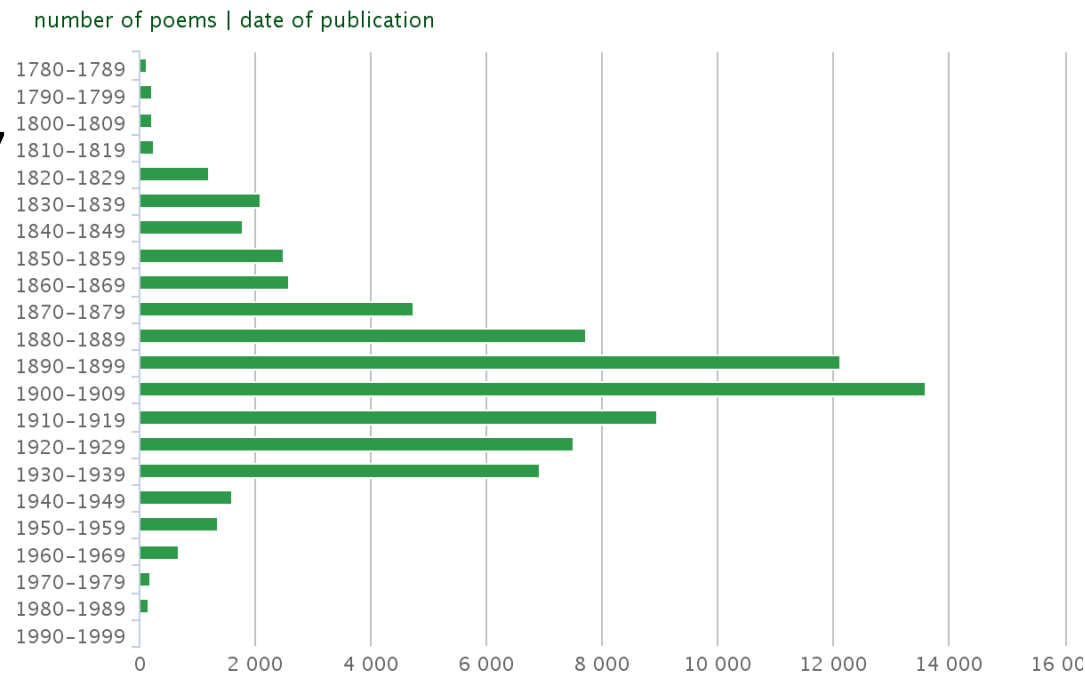
Madrid, 27 junio al 1 de julio de 2016



# Corpus of Czech Verse

- Czech poetry of 19th and the beginning of 20th century

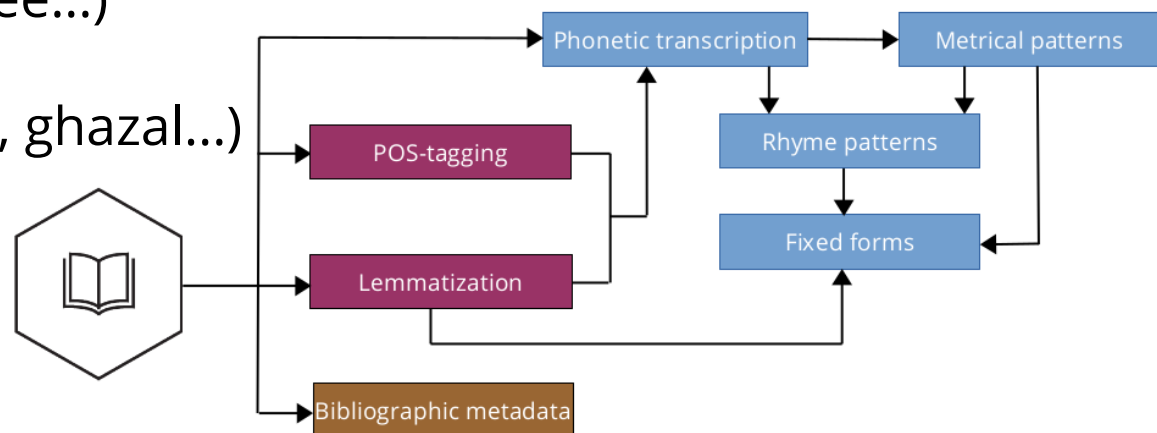
- Books: 1689
- Poems: 76 699
- Verse lines: 2 664 989
- Words: 14 592 037



# Corpus of Czech Verse

## Various levels of annotation

- Lemmatization
- Morphological tagging (POS & others)
- Phonetic transcription
- Metres (iamb, trochee...)
- Rhymes
- Fixed forms (sonnet, ghazal...)



# How to work with the data?



# On-line tools



## **Babel**

Console application – direct SQL queries to the corpus.



## **Database of Czech Metres**

Searching and statistics on metrical and stanzaic level of annotation



## **Euphonometer**

Measuring (statistically defined) euphony of text



## **Frequency lists**

of lemmata / words retrieved from particular books



## **Gunstick**

Searching and statistics on rhymes



## **Hex**

Searching and statistics on keywords in particular poems





Gunstick





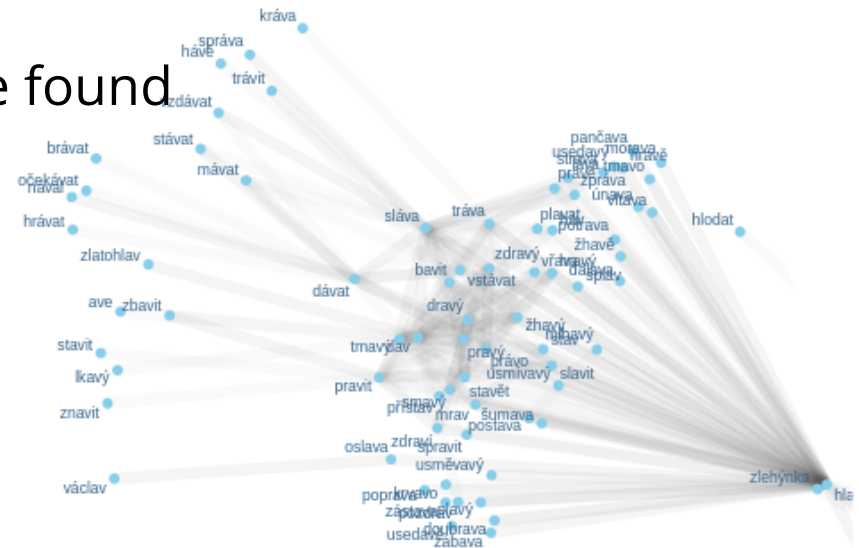
# Gunstick

[http://versologie.cz/gunstick/index\\_en.php](http://versologie.cz/gunstick/index_en.php)

Search for RHYMES on word specified by the user  
(over one million rhyme pairs)

Results:

- Frequency of rhymes on particular words
- Their distribution in time
- List of lines in which rhymes were found  
+ links to full text of poems





# Gunstick

[http://versologie.cz/gunstick/index\\_en.php](http://versologie.cz/gunstick/index_en.php)

## QUERY: láska (love)

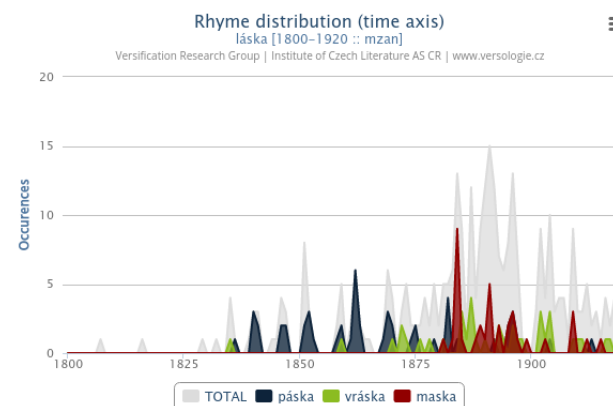
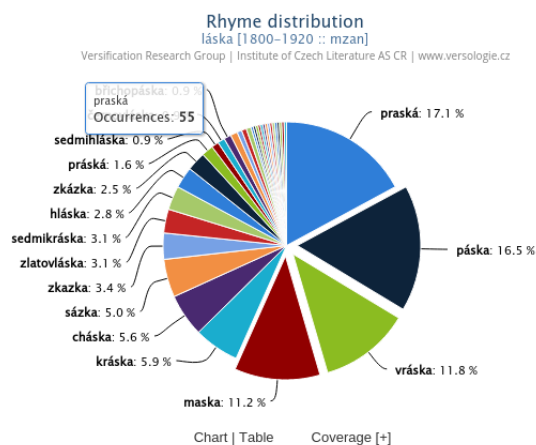
**LÁSKA** [1800-1920 :: mzan] OCCURRENCES: 321 (DIFFERENT WORDS: 35)

[NEW SEARCH](#) | [HELP](#)

[On-line tools](#) » [Gunstick](#)

[To se mi líbí](#) 0

[Tweet](#) [G+](#) 0



Rhyme	Line 1	Line 2	Author	Poem	Collection	Year	End of a line
páska	bud' velebena věčná Boha Láska,	nás pevná k nebi víry pojí páska	Bouška, Sigismund	Ó kniho Bohem psaná...	Pietas	1897	z
páska	čistá, pravá láska.	ta nebeská páska,	Chládek, František	Svatost a láska.	Bázně samouka Františka Chládky	1884	z







# Gunstick

[http://versologie.cz/gunstick/index\\_en.php](http://versologie.cz/gunstick/index_en.php)

## QUERY: láska (love)

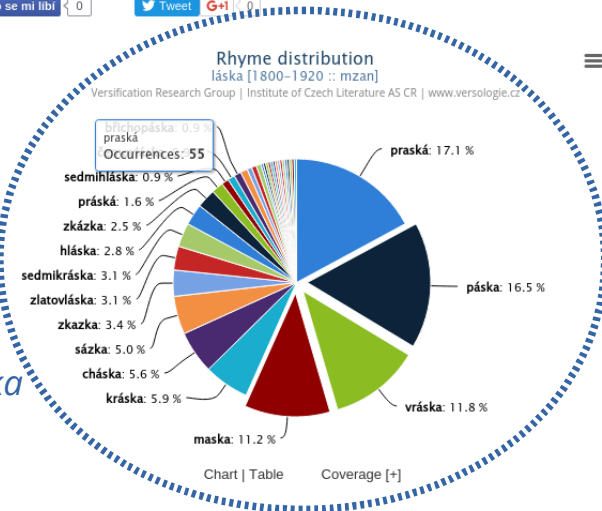
**LÁSKA** [1800-1920 :: MZAN] OCCURRENCES: 321 (DIFFERENT WORDS: 35)

[NEW SEARCH](#) | [HELP](#)

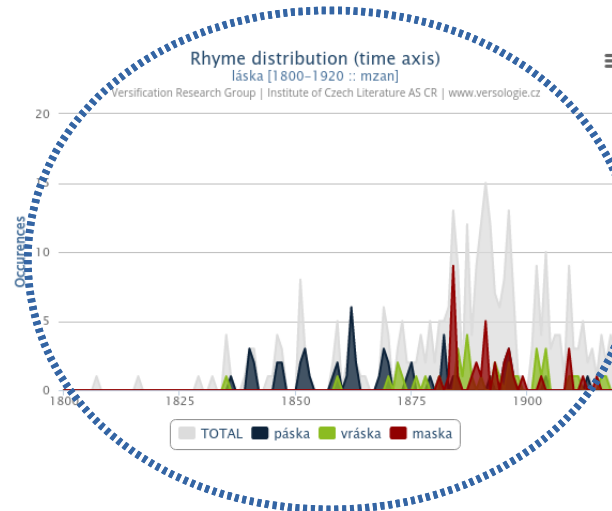
[On-line tools](#) » [Gunstick](#)

[To se mi líbí](#) 0

[Tweet](#) [G+](#) 0



Frequency of  
different words  
rhyming with *láska*



Distribution of  
selected rhymes  
in time (grey =  
all rhymes with  
*láska*)

List of lines +  
links to full  
texts

Rhyme	Line 1	Line 2	Author	Poem	Collection	Year	End of line
páska	bud' velebena věčná Boha Láska,	nás pevná k nebi víry pojí páska	Bouška, Sigismund	Ó kniho Bohem psaná...	Pietas	1897	z
páska	čistá, pravá láska.	ta nebeská páska,	Chládek, František	Svátost a láska.	Básně samouka Františka Chládky	1884	z

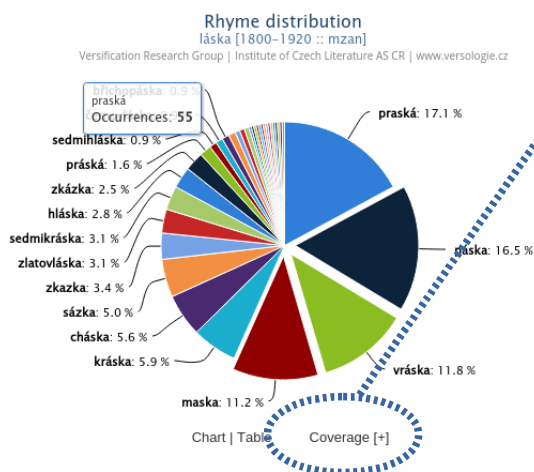




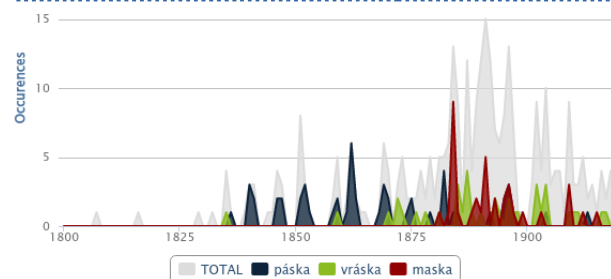
## QUERY: láska (love)

LÁSKA [1800-1920 :: MZAN] OCCURRENCES: 321 (DIFFERENT WORDS: 35)

To se mi líbí 0 Tweet G+ 0



Corpus is not balanced in each year!  
Increasing grey area in second chart does not necessarily mean rhymes with *láska* are becoming more and more frequented during the century. One needs to check COVERAGE.



Rhyme	Line 1	Line 2	Author	Poem	Collection	Year	End of a line
páská	bud' velebena věčná Boha Láska,	nás pevná k nebi víry pojí páská	Bouška, Sigismund	Ó kniho Bohem psaná...	Pietas	1897	z
páská	čistá, pravá láska.	ta nebeská páská,	Chládek, František	Svátost a láska.	Bázně samouka Františka Chládky	1884	z





# Gunstick

[http://versologie.cz/gunstick/index\\_en.php](http://versologie.cz/gunstick/index_en.php)

## QUERY: láska (love)

Chart shows  
WHAT PERCENTAGE  
of all rhymes found in each year  
contain the word *láska*.

=>

If there's some trend it is rather  
increase until the middle of the  
century followed by decrease

Rhyme distribution  
láska [1800-1920 :: mzan]  
Versification Research Group | Institute of Czech Literature AS CR | www.versologie.cz

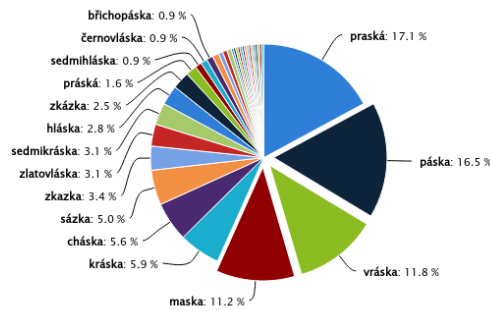
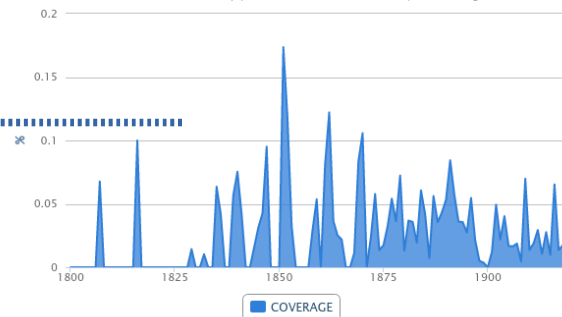


Chart | Table Coverage [-]

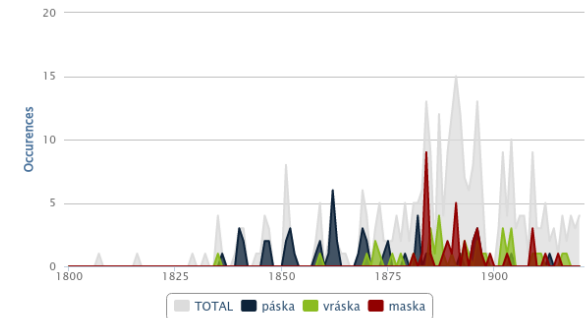
Coverage  
láska [1800-1920 :: mzan]

Versification Research Group | Institute of Czech Literature AS CR | www.versologie.cz



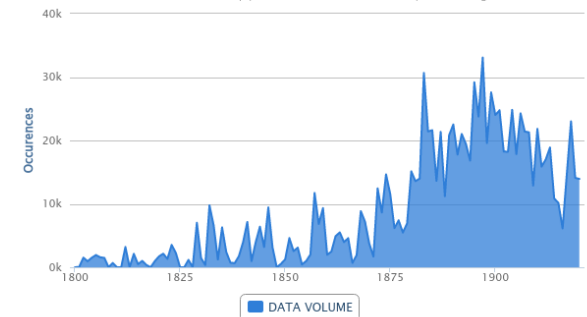
Rhyme distribution (time axis)  
láska [1800-1920 :: mzan]

Versification Research Group | Institute of Czech Literature AS CR | www.versologie.cz



Data volume

Versification Research Group | Institute of Czech Literature AS CR | www.versologie.cz





Hex





# Hex

<http://versologie.cz/hex/>

## Search for KEYWORDS

Keyword: Statistically defined „topic“ of the text (poem)

CONCEPTION → Keyword is a word (lemma) which occurs in a text (poem) more often than one would expect to occur by chance alone

DEFINITION → Keyword of a poem is a word (lemma) which occurs in it with statistically higher frequency than in the rest of the corpus (statistical test:  $\chi^2$  with Yates correction)





Hex

<http://versologie.cz/hex/>

## Search for KEYWORDS

- (1) Search for poems containing specified keyword
- (2) Search for keywords of specified poems





# Hex

<http://versologie.cz/hex/>

## Query (search for occurrence of keywords)

### HEX – KEYWORDS IN CZECH POETRY

#### HLEDAT KLÍČOVÉ SLOVO

slovo

autor

od

od

#### PROCHÁZET DATABÁZI

- ☒ podstatná jména
- ☒ přídavná jména
- ☐ zájmena
- ☐ číslovky
- ☒ slovesa
- ☐ příslovce
- ☐ předložky
- ☐ spojky
- ☐ částice
- ☐ citoslovce

Minimální četnost:

3 ▼

Hladina významnosti ( $\alpha$ ):

0.001 ▼

Hledat





# Hex

<http://versologie.cz/hex/>

## Query (search for occurrence of keywords)

### HEX – KEYWORDS IN CZECH POETRY

**HLEDAT KLÍČOVÉ SLOVO** **PROCHÁZET DATABÁZI** **Hledat**

Specify word .....   
(Limit to certain author) .....   
  
(Limit to a time span: from / to) .....   
.....

☒ podstatná jména  
☒ přídavná jména  
☐ zájmena  
☐ číslovky  
☒ slovesa  
☐ příslovce  
☐ předložky  
☐ spojky  
☐ částice  
☐ citoslovce

Minimální četnost:   
Hladina významnosti ( $\alpha$ ):

Set minimum number of occurrences of the lemma  
Set the alpha-level of statistical test

Limit query to only some parts-of-speech  
Due to rich homonymity in Czech, e.g.: *stát* = (1) state, country (noun), (2) stand (verb)







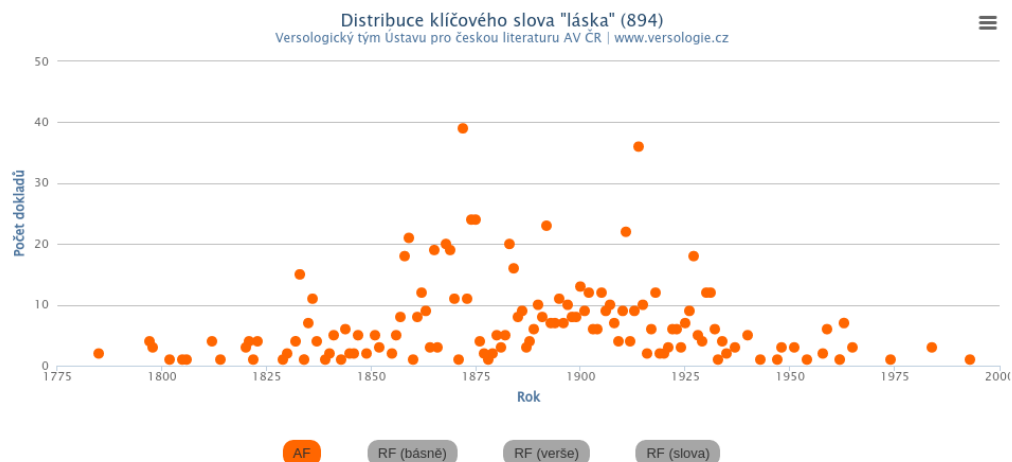
# Hex

<http://versologie.cz/hex/>

## Results – query: láska (love)

DISTRIBUCE KLÍČOVÉHO SLOVA "LÁSKA" (894) AUT: | SB: | - ( $A = 0.001$ ;  $N \geq 3$ ; POS: NAV)

NOVÉ HLEDÁNÍ | O APLIKACI



DOKLADY KLÍČOVÉHO SLOVA "LÁSKA" (894) AUT: | SB: | - ( $A = 0.001$ ;  $N \geq 3$ ; POS: NAV)

Autor	Báseň	Sbírka	Rok	AF	RF	$\phi$
Ambrož, Vilém	Tré hlasů. (III.) »	Pestré kvítí »	1883	5	61728	0.0025
Ambrož, Vilém	Ecce lignum crucis! »	Pestré kvítí »	1883	4	33898	0.0015





# Hex

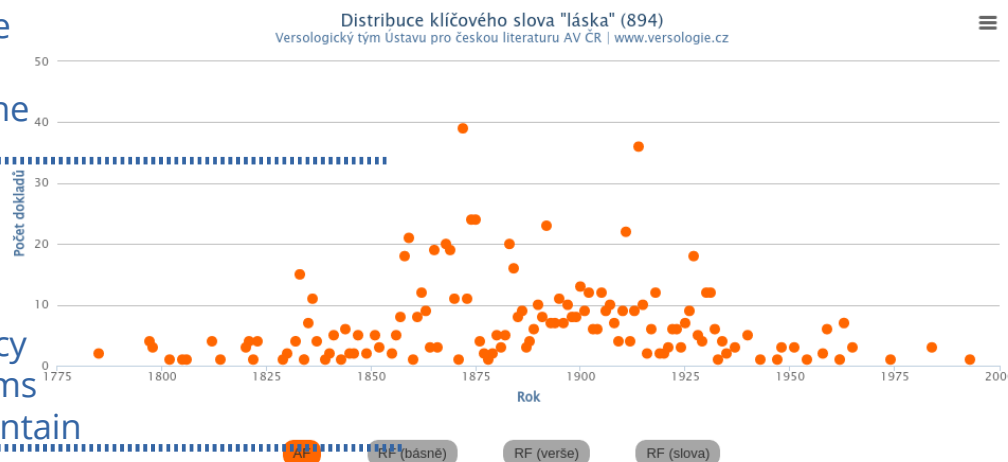
<http://versologie.cz/hex/>

## Results – query: láska (love)

DISTRIBUCE KLÍČOVÉHO SLOVA "LÁSKA" (894) AUT: | SB: | - (A = 0.001; N ≥ 3; POS: NAV)

NOVÉ HLEDÁNÍ | O APLIKACI

Each point represents the number of poems found in each year containing the keyword *láska*



Switch to relative frequency (what PORTION of all poems published in given year contain the word *láska*)

List of poems + links to full-texts + links to the lists of keywords found in them under same conditions

DOKLADY KLÍČOVÉHO SLOVA "LÁSKA" (894) AUT: | SB: | - (A = 0.001; N ≥ 3; POS: NAV)

Autor	Báseň	Sbírka	Rok	AF	RF	φ
Ambrož, Vilém	Tré hlasů. (III.) »	Pestré kvítí »	1883	5	61728	0.0025
Ambrož, Vilém	Ecce lignum crucist »	Pestré kvítí »	1883	4	33898	0.0015





Hex

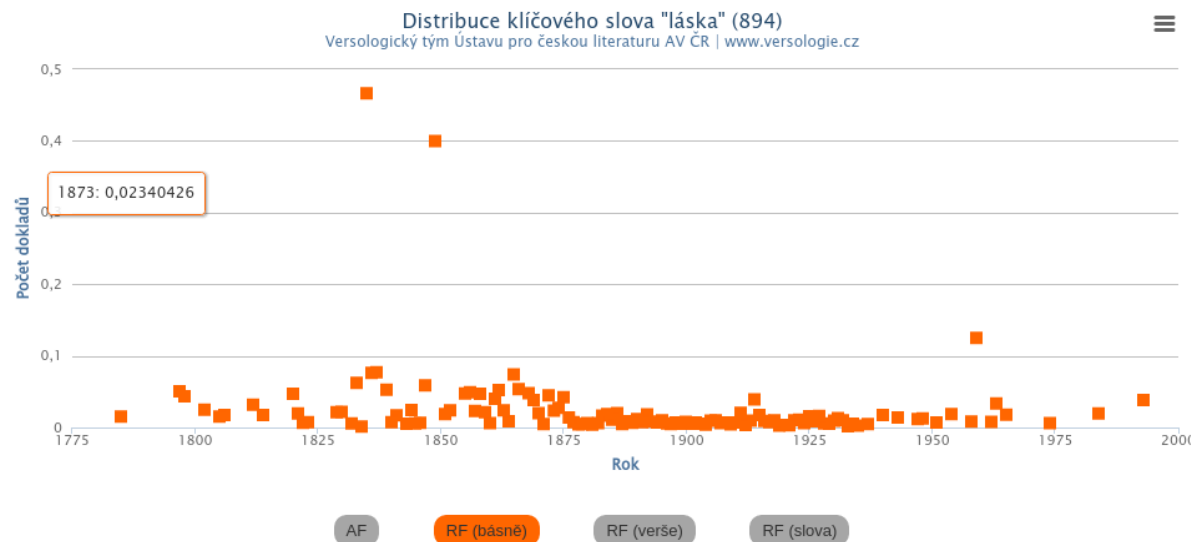
<http://versologie.cz/hex/>

## Results – query: láska (love) – RELATIVE FREQUENCIES

No clear trend at all, just some outliers...

DISTRIBUCE KLÍČOVÉHO SLOVA "LÁSKA" (894) AUT: | SB: | – (A = 0.001; N ≥ 3; POS: NAV)

NOVÉ HLEDÁNÍ | O APLIKACI



What such findings may be used for?

Among others: the attribution of texts  
with unknown or disputed authorship...

... let's go on-line:

<http://versologie.cz/talks/uned>

