

# Lexémy na Wikidatech a slovník pro kontrolu pravopisu

Wikikonference 2019

**Stanislav Horáček**

česká skupina kolem LibreOffice  
The Document Foundation

Wikipedista:Strepon

23. listopadu 2019

# Úvod

český slovník pro kontrolu pravopisu

- Hunspell
- licence GNU GPL

*Toto je český slovník pro kontrolu pravopisu založený na českém slovníku pro ispell, verze z 29. 10. 2006, který vytvořil Petr Kolar spolu s desítkami dalších přispěvatelů.*

# Úvod

## Český tvarotvorný slovník

slovní zásoba

- Masarykova univerzita
- únor 2019, licence public domain (= CC0)
- analýza jazykového korpusu
- podstatná jména, přídavná jména, slovesa
- ~60 000 základních tvarů
- [github.com/plin/slovník](https://github.com/plin/slovník)

# Úvod

## Slovníková data na Wikidatech

rozhraní

- oddělený prostor
- databáze slovní zásoby
- základní jednotka lexém
- ~3500 základních tvarů v češtině
- během roku 2018, licence CC0
- [www.wikidata.org/wiki/Wikidata:Lexicographical\\_data/cs](http://www.wikidata.org/wiki/Wikidata:Lexicographical_data/cs)

# Slovníková data na Wikidatech

projekt Wikimedia Foundation

# Slovníková data na Wikidatech

projekt Wikimedia Foundation

možné využití

- slovník pro kontrolu pravopisu
- výkladový slovník
- slovník dělení slov
- slovník synonym, antonym, ...
- cizojazyčné slovníky a překladače
- nástroje pro kontrolu gramatiky
- ...

# Lemma, významy, tvary

(L45253)

## duchaplně

 edit

CS

Language [Czech](#)

Lexical category [adverb](#)

### Statements

[+ add statement](#)

### Senses

L45253-S1

CS

duchaplným způsobem

 edit

Statements about L45253-S1

[+ add statement](#)

[+ add Sense](#)

### Forms

L45253-F1

duchaplně

CS

 edit

Grammatical features [positive](#)

Statements about L45253-F1

[+ add statement](#)

L45253-F2

duchaplněji

CS

 edit

Grammatical features [comparative](#)

Statements about L45253-F2

# Charakteristika, etymologie

(L10536)

# prasátko

 edit

CS

Language [Czech](#)

Lexical category [noun](#)

## Statements

grammatical gender



neuter

 edit

▼ 0 references

+ add reference

+ add value

derived from



prase

 edit

▼ 0 references

+ add reference

+ add value

+ add statement

## Senses

L10536-S1

CS

malé prase

 edit

Statements about L10536-S1

language style



diminutive

 edit

▼ 0 references

+ add reference



# Příznaky

(L58354)

# měďák

 edit

CS

Language [Czech](#)

Lexical category [noun](#)

## Statements

grammatical gender



masculine

 edit

▾ 0 references

+ add reference

## Senses

L58354-S1

Czech

mince z mědi

 edit

### Statements about L58354-S1

language style



Common Czech

 edit







▾ 0 references

+ add reference

+ add value

+ add statement

# Propojení s položkami Wikidat

(L46367)	<b>lopuch</b>	 edit
	CS	
Language <a href="#">Czech</a> Lexical category <a href="#">noun</a>		
<b>Statements</b>		
<a href="#">grammatical gender</a>	 masculine	 edit
	<a href="#">▼ 0 references</a>	<a href="#">+ add reference</a>
<b>Senses</b>		
L46367-S1	Czech	rod rostlin  edit
Statements about L46367-S1		
<a href="#">item for this sense</a>	 <a href="#">Arctium</a>	 edit
	<a href="#">▼ 0 references</a>	<a href="#">+ add reference</a>
		<a href="#">+ add value</a>
		<a href="#">+ add statement</a>





# Vztahy mezi lexémy

(L43081) **sosna** 

CS

item for this sense	 Pinus ▼ 0 references	 + add reference + add value
synonym	 L43080-S1 ▼ 0 references	 + add reference + add value

# Varianty

(L58010)	chromovaný	CS	 edit
	chrómovaný	CS-X-Q28861	
Language <a href="#">Czech</a>			
Lexical category <a href="#">adjective</a>			
<b>Statements</b>			
derived from	 chromovat/chrómovat		 edit
	▼ 0 references		
L58010-F3	chromovanému CS	chrómovanému CS-X-Q28861	 edit
Grammatical features <a href="#">singular, dative case, animate masculine</a>			
Statements about L58010-F3			
<a href="#">+ add statement</a>			

# Výslovnost

(L442)

být

CS



L442-F14

být

CS



Grammatical features [infinitive](#)

Statements about L442-F14

[pronunciation audio](#)



[Cs-být.ogg](#)

1.1 s; 14 KB

language of work or name

Czech

[0 references](#)

[+ add reference](#)

[+ add value](#)

[IPA transcription](#)



bi:t

language of work or name

Czech

[0 references](#)



[+ add reference](#)

# Příklady užití

(L205210)

## Pavel

CS

 edit



S Pavlem si dobře rozumíme. (Czech)

 edit

demonstrates form

Pavlem

demonstrates sense

L205210-S1

▾ 0 references

+ add reference



My o Pavlu, a Pavel za dveřmi. (Czech)

 edit

demonstrates form

Pavel

Pavlu

demonstrates sense

L205210-S1

▾ 0 references

+ add reference



Ve třídě máme dva Pavly. (Czech)

 edit

demonstrates form

Pavly

demonstrates sense

L205210-S1

▾ 0 references

+ add reference

# Překlady?

(L2080)

# voda

CS

 edit

## Senses

L2080-S1

Czech  
Russian

tekutina  
жидкость

 edit

### Statements about L2080-S1

item for this sense



liquid water

 edit

▼ 0 references

+ add reference

+ add value

translation



L189-S1

 edit

▼ 0 references

+ add reference

+ add value

# Citování?

(L10856)

ječmen

CS

 edit

L10856-F3

ječmena

CS

 edit

Grammatical features [singular, genitive case](#)

Statements about L10856-F3

described by source



Internet Language Reference Book

 edit

URL

<http://prirucka.ujc.cas.cz/?slovo=je%C4%8Dmen>

▼ 0 references


+ add reference

+ add value



# Vytvoření nového lexému

[www.wikidata.org/wiki/Special:NewLexeme](http://www.wikidata.org/wiki/Special:NewLexeme)

Special page  

## Create a new Lexeme

By clicking "Create", you agree to the [terms of use](#), and you irrevocably agree to release your contribution under the [Creative Commons CC0 License](#).

**Create a new Lexeme**

Lemma

 \*

Language of Lexeme

 \*

Lexical category

 \*

# Šablony pro různé slovní druhy

[tools.wmflabs.org/lexeme-forms/](https://tools.wmflabs.org/lexeme-forms/)

[Wikidata Lexeme Forms](#) [Documentation](#) [Wikimedia Toolforge](#) [Source code](#)

## české podstatné jméno (rod střední)

### 1. pád, jednotné číslo

To je mé .

### 2. pád, jednotné číslo

Má strach z mého .

### 3. pád, jednotné číslo

Dej to mému .

### 4. pád, jednotné číslo

Vidím jedno .

### 5. pád, jednotné číslo

Kam kráčíš, ?

### 6. pád, jednotné číslo

Pověz mi něco o tvém .

### 7. pád, jednotné číslo

Seznámil jsem se s tvým .

### 1. pád, množné číslo

To jsou má .

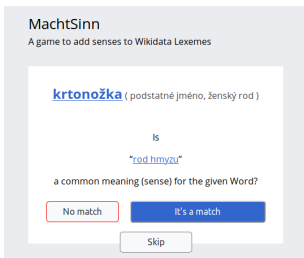
### 2. pád, množné číslo

Má strach z mých .



# Nástroje

MachtSinn – propojení významů s položkami Wikidat  
[tools.wmflabs.org/machtsinn/](https://tools.wmflabs.org/machtsinn/)



## Senses

L54480-S1

Czech

rod hmyzu



## Statements about L54480-S1

item for this sense



Gryllotalpa



< 0 references

+ add reference

+ add value

## Přístup přes API

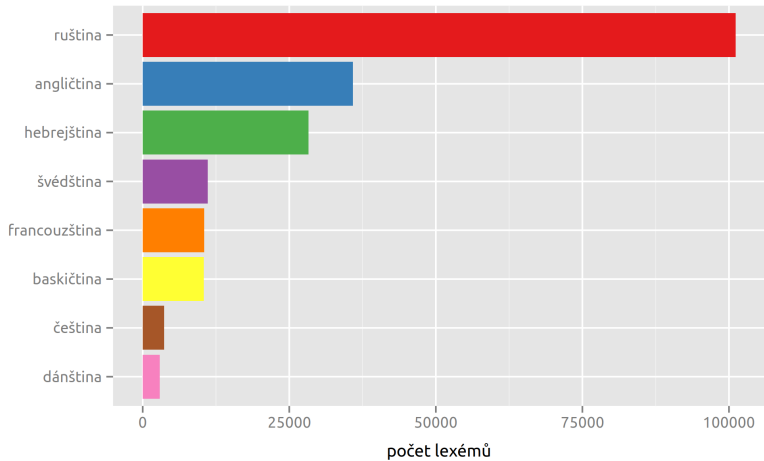
dotazovací jazyk SPARQL pro Wikidata

- [query.wikidata.org](http://query.wikidata.org)

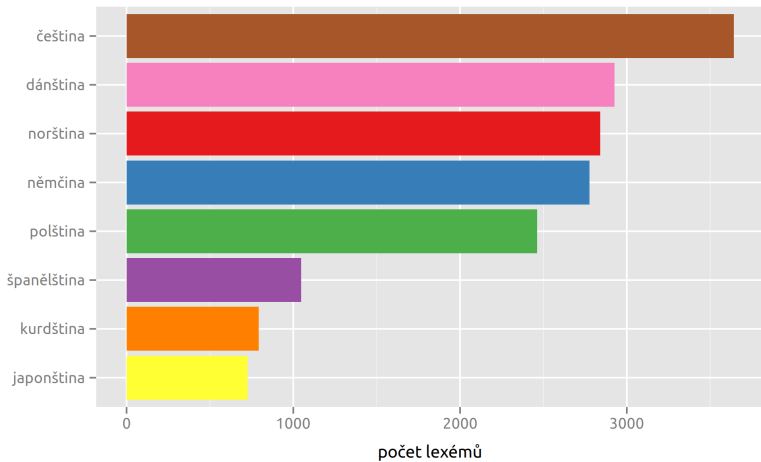
Pywikibot

- [www.mediawiki.org/wiki/Manual:Pywikibot](http://www.mediawiki.org/wiki/Manual:Pywikibot)
- neúplná podpora pro lexémy na Wikidatech
- [phabricator.wikimedia.org/T189321](https://phabricator.wikimedia.org/T189321)
- nefunkční generátory – WikidataSPARQLPageGenerator

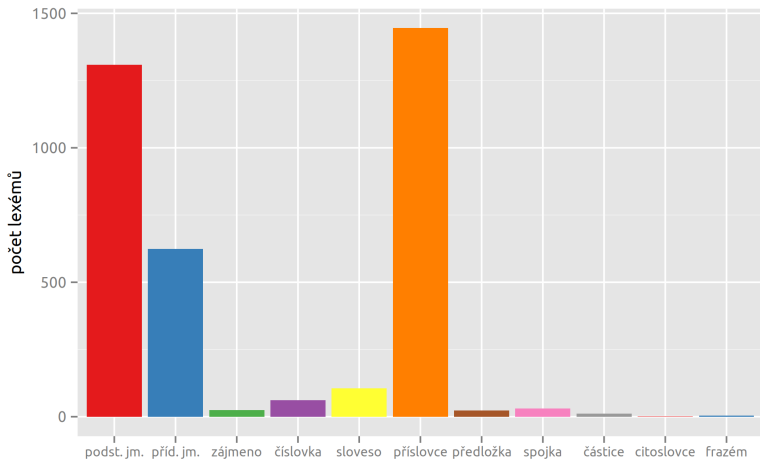
# Lexémy podle jazyků



# Lexémy podle jazyků



# Slovní druhy pro češtinu






Český tvarotvorný slovník  
+ slovníková data z Wikidat  
= české CC0 slovníky

**experimentální!**

# České CC0 slovníky

## rozšíření pro LibreOffice

[extensions.libreoffice.org/extensions/czech-cc0-dictionaries-ceske-cc0-slovniky](https://extensions.libreoffice.org/extensions/czech-cc0-dictionaries-ceske-cc0-slovniky)

Log in Register Search Site Search

[Home](#) [Extensions](#) [Templates](#) [Events](#)

[Home](#) [Extensions](#) [Czech CC0 dictionaries / České CC0 slovníky](#)

## Czech CC0 dictionaries / České CC0 slovníky

Czech spell check dictionary licensed under the Creative Commons CC0 License / Slovník kontroly pravopisu pro češtinu zveřejněný pod licencí Creative Commons CC0 👍 13 👤 1

### Project Description

The dictionary uses data from the [Czech morphological dictionary](#) created by the Masaryk University in Brno and from [Czech Wikidata lexemes](#).

Repository of this dictionary can be found at [GitLab](#).

The dictionary is considered as experimental (many words missing), with a lot of space for improvements.

For more details and instructions how to contribute, see the webpage [ceskeslovniky.cz](#).

---

Slovník využívá data [Českého tvaroslovného slovníku](#) vytvořeného na Masarykově univerzitě v Brně a [česká slovníková data z Wikidat](#).

Repozitář s tímto českým slovníkem je k dispozici na [GitLabu](#).

Slovník je třeba považovat za experimentální (mnoho slov chybí) a je zde značný prostor pro jeho vylepšování.

Podrobnosti a návod jak se zapojit naleznete na stránkách [ceskeslovniky.cz](#).

### Project Resources

- [External Project Page](#)

#### Screenshot



### Install Instructions

To install an extension, follow these steps:

- Download an extension and save it anywhere on your computer.
- In LibreOffice, select Tools --> Extension Manager from the menu bar.
- In the Extension Manager dialog click Add.

# České CC0 slovníky

srovnání úspěšnosti

– procentuální podíl slov označených jako chybná

	červen 2019	listopad 2019	GNU GPL
Dobrodružství Sherlocka Holmese	9,02	6,78	2,63
Evangelium podle Jana	7,46	4,61	0,67
LibreOffice Writer: Praktický průvodce	6,00	4,56	3,33
R.U.R.	16,82	11,87	8,37
Ústava České republiky	7,89	5,57	0,90
Wikikonference 2019	14,68	12,61	8,51

# Jak vylepšit

- doplňování slov na Wikidata
- import z Tvarotvorného slovníku
- import z Wikislovníku
- import z Wikidat
- nové šablony
- nové nástroje
- kampaň pro určité texty

# Shrnutí

## slovníková data na Wikidatech

- databáze slovní zásoby pod licencí CC0
- [wikidata.org/wiki/Wikidata:Lexicographical\\_data/cs](https://wikidata.org/wiki/Wikidata:Lexicographical_data/cs)
- [wikidata.org/wiki/Wikidata\\_talk:Lexicographical\\_data](https://wikidata.org/wiki/Wikidata_talk:Lexicographical_data)
- [wikidata.org/wiki/Wikidata:Mezi\\_bajty](https://wikidata.org/wiki/Wikidata:Mezi_bajty)

## nový český slovník kontroly pravopisu

- experimentální, pod licencí CC0
- Tvarotvorný slovník a lexémy z Wikidat
- budoucnost?
- [ceskeslovniky.cz](https://ceskeslovniky.cz)