



## Fulltextový vyhledávač

Štěpán Škrob

[<stepan.skrob@firma.seznam.cz>](mailto:stepan.skrob@firma.seznam.cz)

# O čem bude přednáška?

1. Úvod
2. Přehled architektury vyhledávače
3. Důležité datové struktury
4. Použité algoritmy
  - a) Crawlování
  - b) Indexace
  - c) Vyhledávání
5. Výkonové statistiky

# Úvod

- Vyhledávače jsou si prakticky velmi podobné, liší se pouze v implementačních detailech :-)
- Je to jako s auty...



# Přehled architektury vyhledávače

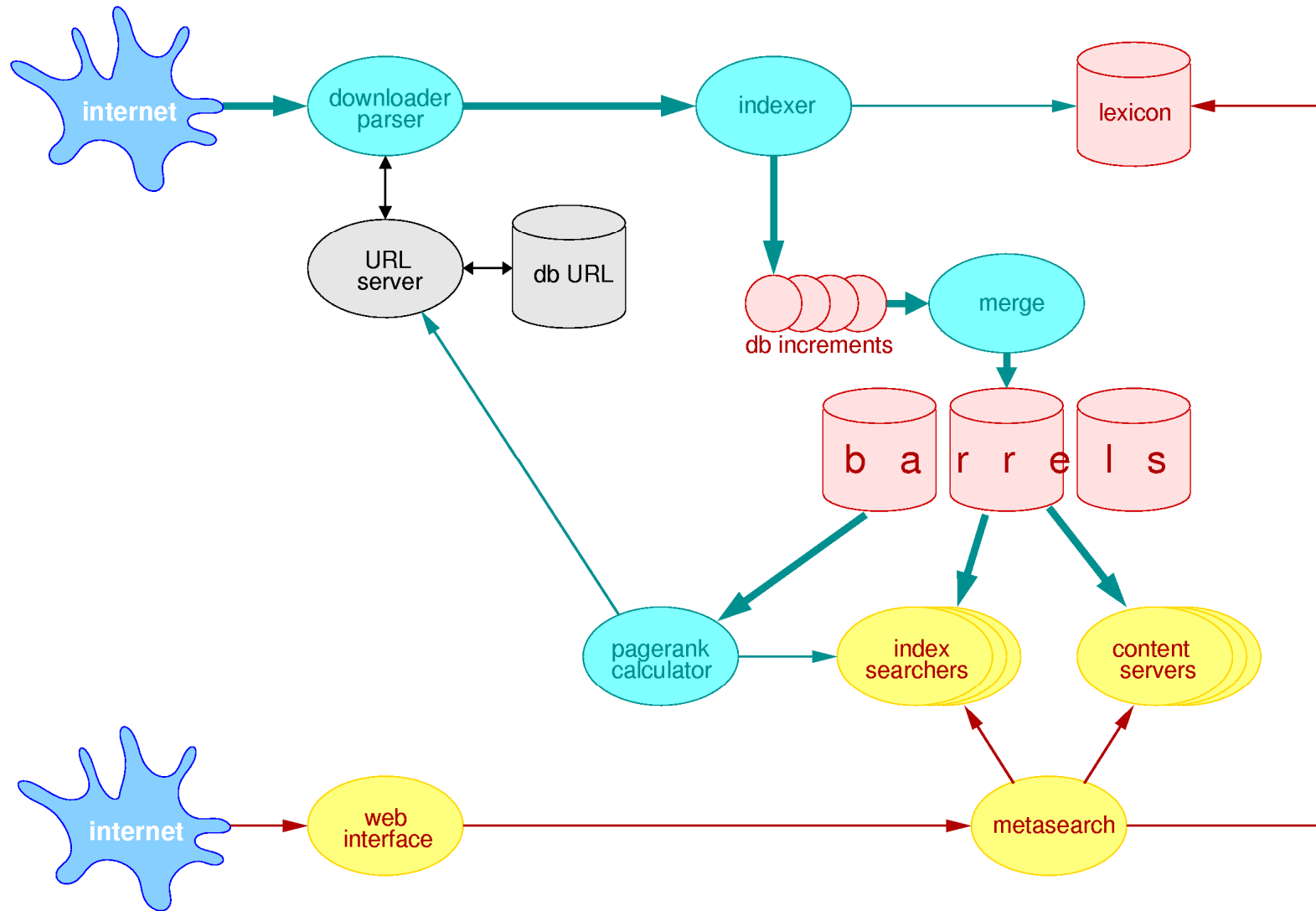
# Přehled architektury (1)

- Původní škálovatelnost do 100 M dokumentů
- Nyní až 1 G dokumentů
  - Rozdělení databáze do svazků
  - Zkrácení odezvy vyhledávače
- Současná databáze
  - Počet dokumentů 30 M
  - Datová velikost 100 GB

## Přehled architektury (2)

- Výkon robota (crawlera)
  - 2 M dokumentů / den
- Dvě úzká hrdla
  - Výběr dokumentu k reindexaci
  - Přeložení URL na url\_id
  - Obojí je náhodný přístup k datům...

# Přehled architektury (3)



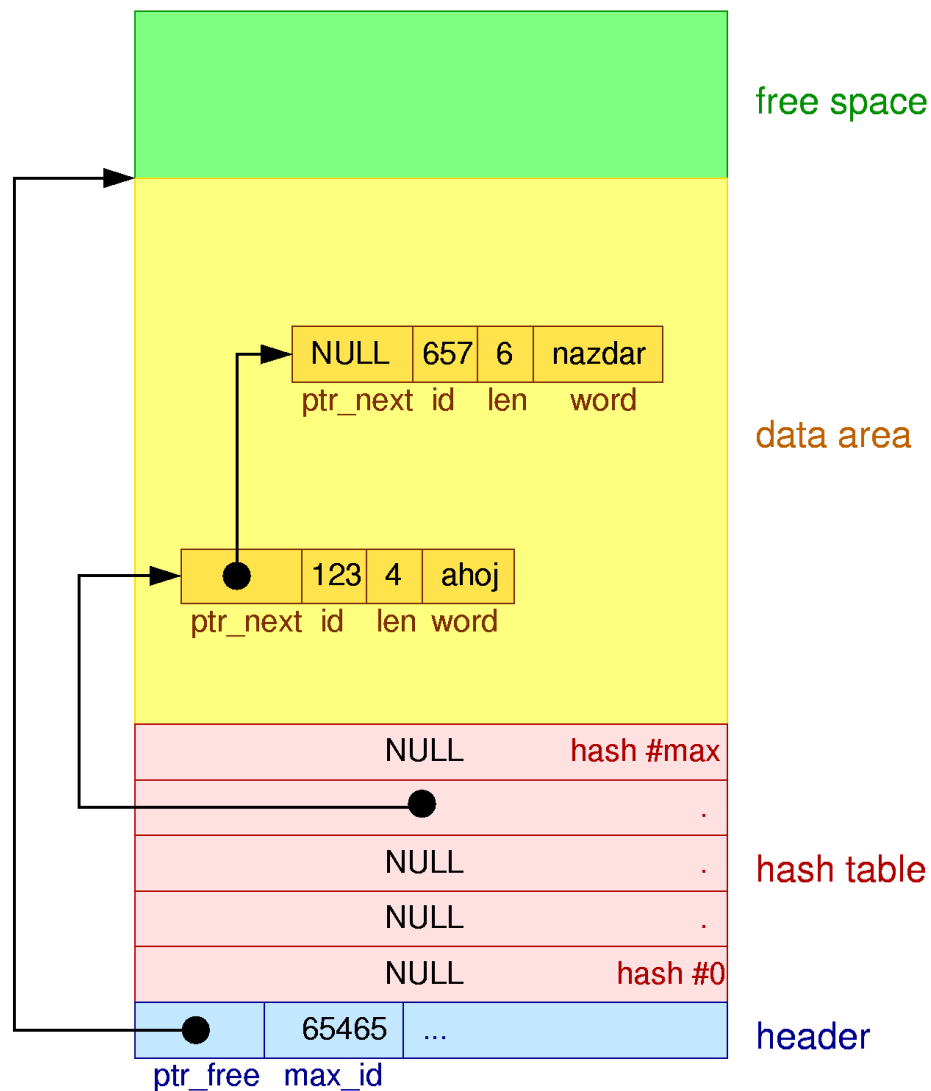
# Důležité datové struktury



# Důležité datové struktury

1. Lexikon
2. Invertovaný index
3. Texty dokumentů

# Lexikon (1)



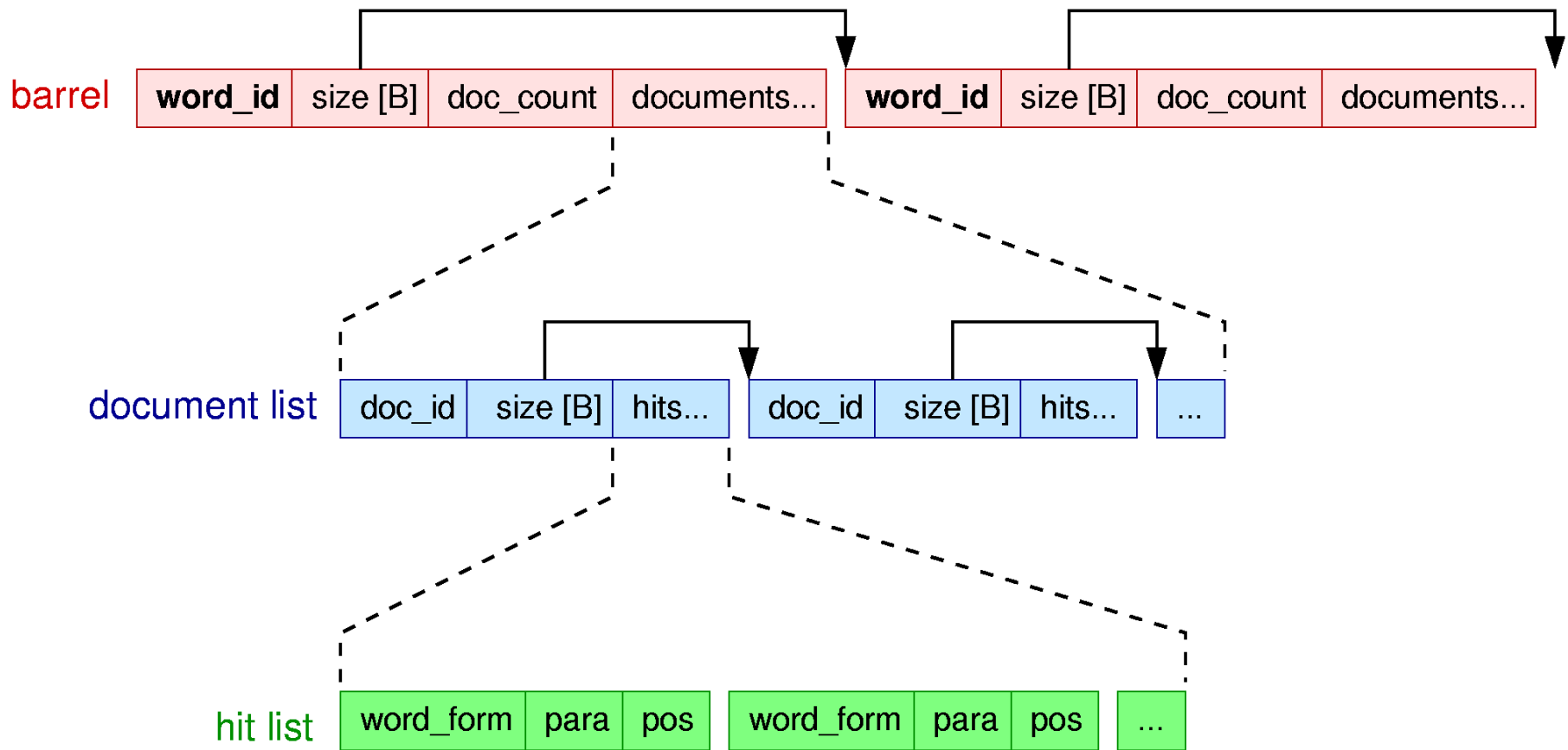
## Lexikon (2)

- Lemmatizátor je spojený s lexikonem do jedné komponenty
- Obsah lexikonu
  - Lemmatizovaná slova => lemmata
  - Nelemmatizovaná slova => termy

## Lexikon (3)

- Statistické údaje
  - 450 MB celková velikost
  - 22 miliónů slov
  - 16 miliónů hashů (využitých 75%)
  - Největší kolize 13 slov / hash
- Výkon
  - 100.000 slov / sec

# Invertovaný index (1)



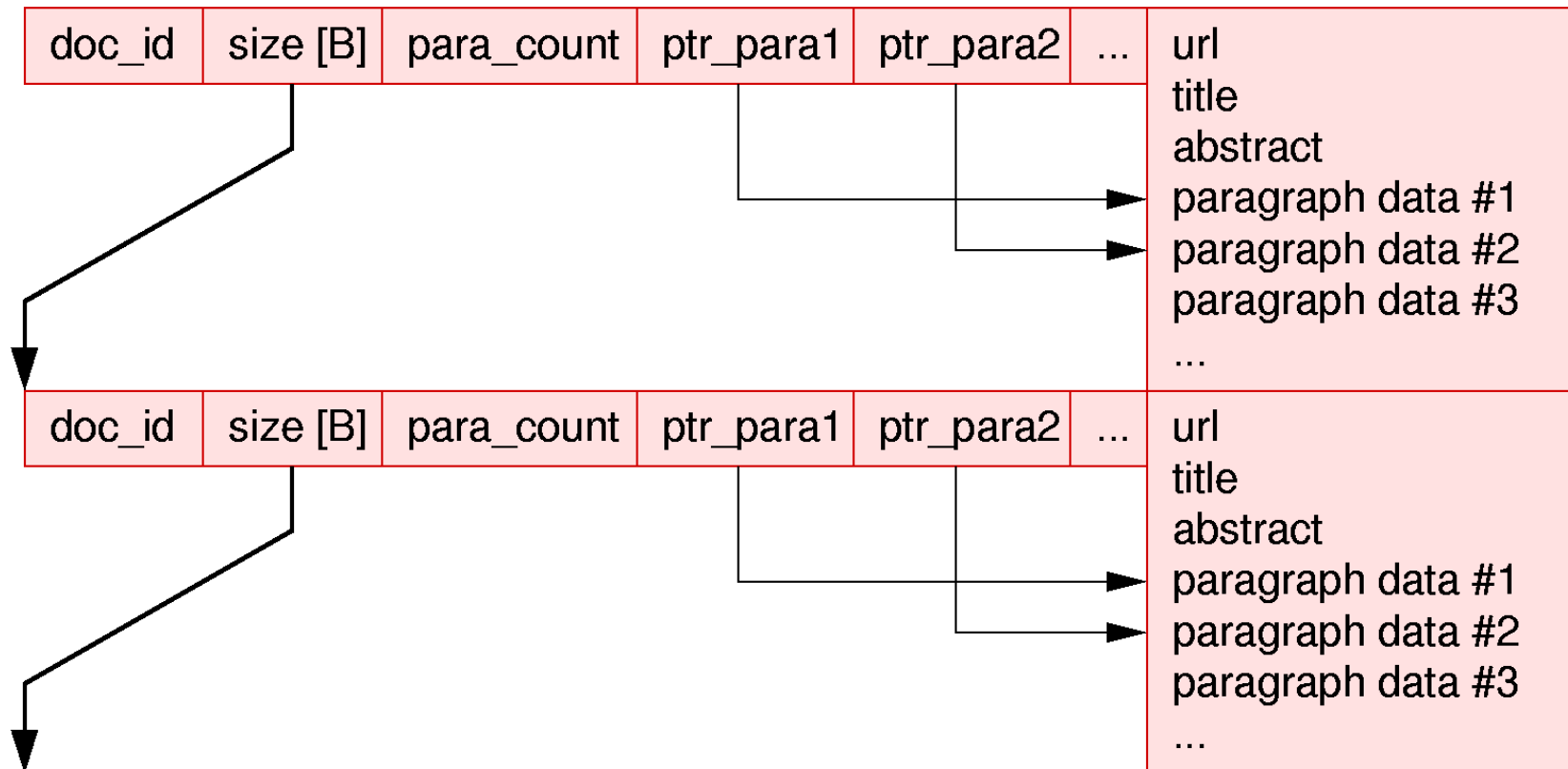
## Invertovaný index (2)

- Rozdílová + prostorová komprimace ideček
- Prostorová komprimace ostatních *int* čísel
- Uchovávané informace o slově:
  - Tvar slova (5 bitů)
  - Číslo odstavce
  - Pozice v odstavci

## Invertovaný index (3)

- Pro nalezení záznamu slova v barrelu (invertovaného indexu) slouží další index
  - Trvale přítomný v paměti (no I/O)
  - Obsahuje informace o souboru, pozici, délce
  - Seřazený podle word\_id => binární půlení
- Statistické údaje (30 M dokumentů)
  - Invertovaný index 50 GB
  - Trvale načtené struktury ~250 MB

# Texty dokumentů (1)





## Texty dokumentů (2)

- Texty všech zaindexovaných stránek včetně balastu
  - Zachované informace o odstavcích
  - Označované začátky a konce slov
  - Neobsahuje HTML značky
- Použití pro dynamický popis dokumentu
- Automatický abstrakt dokumentu

# Texty dokumentů (3)

- Statistické informace
  - Celková velikost 50 GB na 30 M dokumentů
  - Průměr 200 slov na dokument
  - Průměrná velikost slova 7 znaků

# Použité algoritmy

# Použité algoritmy

## 1. Crawlování

- a) Filtrace dokumentů
- b) Detekce českého jazyka

## 2. Indexace

- a) Lemmatizace
- b) Identifikace podobných dokumentů

## 3. Vyhledávání

- a) Doplnování diakritiky
- b) Výběr vhodného lemmatu
- c) Hodnocení relevance

# Filtrace dokumentů (1)

- Základní filtrace textových dokumentů podle Content-Type
  - Uhádnutého z URL
  - Zjištěného z HTTP hlaviček serveru
  - Pozor na CGI scripty a URL zakončená lomítkem (<http://novinky.cz/sport/>)

# Filtrace dokumentů (2)

- Vyřazení binárních souborů
  - Citlivost na řídicí znaky ASCII < 32
  - Hledání tokenů podobných slovům (mezery a délka slova)
- Vyřazení zcela duplicitních dokumentů
  - MD5 hash z textového derivátu

# Detekce českého jazyka

- Kombinace slovníkové a statistické metody
  - Slovník ze 100.000 článků novinky.cz
  - Frekvence slov
  - Charakterističnost slov pro češtinu (žížala...)
- Vysoká citlivost na přítomnost českého jazyka
- Nutná spolehlivá funkce i pro krátké texty
- Ruční nastavení vah, neumí se učit

# Lemmatizace textu (1)

- Lemmatizátor Lingea
  - Zaměření na spisovnou češtinu
  - Více správných možností => dvě lemmata (citron a citrón)
  - Názvy, jména, příjmení...
- Lemmatizace pouze podstatných a přídavných jmen



# Lemmatizace textu (2)

- Neznámá slova
  - Indexace termu
- Nejednoznačná slova
  - Zaindexování ve prospěch všech lemmat
  - Problémy: „Jana Králová“
- Text bez diakritiky
  - Doplnuje se pokud je celý odstavec bez diakritiky
  - Výskyt hlavně v diskuzích

# Určení podobnosti dokumentů (1)

- Zkoumá se podobnost každého dokumentu se všemi ostatními dokumenty
- Pro každý dokument se spočítá  $N$  hashů
  - $N$  hashovacích funkcí
  - Hashe se počítají pro plovoucí okno  $m$ -tic slov
  - Výstupem je  $\max(\text{hash})$  z každé funkce

## Určení podobnosti dokumentů (2)

- Podobné dokumenty jsou takové, kde  $X$  z  $N$  hashů je stejných
- Výstupem výpočtu je  $X$  souborů s identifikovanými podobnostmi
- Bohužel neplatí transitivita
  - Pokud existuje  $A \sim B$ ,  $B \sim C$
  - Pak neplatí, že  $A \sim C$

## Určení podobnosti dokumentů (3)

- Ze souborů vytvoříme třídy podobných dokumentů jako kdyby podobnost byla transitivní

# Doplňování diakritiky (1)

- Účel
  - Stále nezanedbatelné procento dotazů bez diakritiky
  - Nesprávně psané dlouhé samohlásky (á, é, í, ó)  
(Slávia, ale myslí Slavia)
  - Překlepy  
(Škoda místo Škoda)

# Doplňování diakritiky (2)

- Postup
  - Úplné odstranění diakritiky z dotazu
  - Doplnění diakritiky do očištěného dotazu
  - Nalezení možných lemmat pro doplněná slova
    - => Zvětšení nejednoznačnosti
    - => Nutný výběr

# Výběr vhodného lemmatu (1)

- Pro každé slovo dotazu existuje  $N$  možných lemmat
- Musíme vybrat jedno lemma pro každé slovo
  - Výkon: nelze hledat všechna slova
  - Smysl dotazu: musíme určit slovo které uživatel myslel

## Výběr vhodného lemmatu (2)

- Preference lemmat ze znalosti situace
- Typická podoba dotazů
  - Drtivá většina podstatná a přídavná jména (dovolená chorvatsko)
  - Slovesa vyjíměčně (běží liška k táboru)
  - 1. pád jedn. i mn. čísla (krby)



## Výběr vhodného lemmatu (3)

- Použití frekvenčního slovníku dvojic lemmat
  - Kombinace všech možností lemmatizace dotazu
  - Ohodnocení každé kombinace podle frekvence výskytu
- Výsledek – rozdílná lemmata pro:
  - Německých tancích
  - Německých lidových tancích

# Výběr vhodného lemmatu (4)

slovo (doplněná diakritika) #ID	lemma	slovní druh	word-barel		
			svazek	typ	velikost [KB]
německých (německých) #1176	německý	přídavné jméno	1	velký	1 522,4
			2	velký	1 513,7
			3	velký	1 541,1
			4	velký	1 163,0
tancích (tancích) #8793	tank	podstatné jméno	1	velký	356,2
			2	velký	327,3
			3	velký	251,1
			4	velký	196,6

- Výsledek z administračního rozhraní pro „Německých tancích“

# Výběr vhodného lemmatu (5)

slovo (doplněná diakritika) #ID	lemma	slovní druh	word-barel		
			svazek	typ	velikost [KB]
nemeckych (německých) #1176	německý	přídavné jméno	1	velký	1 522,4
			3	velký	1 541,1
			4	velký	1 163,0
lidovych (lidových) #1464	lidový	přídavné jméno	1	velký	1 706,4
			3	velký	354,1
			4	velký	449,3
tancich (tancích) #22810	tanec	podstatné jméno	1	velký	686,4
			3	velký	599,4
			4	velký	542,6

- Výsledek z administračního rozhraní pro „Německých lidových tancích“

# Hodnocení relevance (1)

- Vyhodnocuje se pozice a vzdálenost slov
- Hledá se uspořádání slov co nejpodobnější dotazu
- Příklady (dotaz A B):
  - A B      vzdálenost 0
  - A x B    vzdálenost 1
  - A x x B  vzdálenost 2
  - B A      vzdálenost 1

## Hodnocení relevance (2)

- V textu musejí být všechna slova
- Hodnocení slov v:
  - Titulku
  - URL
  - Textu
- Hledání složeného slova v URL  
([www.jakpsatweb.cz](http://www.jakpsatweb.cz))

# Hodnocení relevance (3)

- Celková relevance dokumentu
  - Nelineární kombinace částečných relevancí všech nalezených výskytů
  - Připočítání ranku dokumentu
- Výsledek seřazen podle celkové relevance
- Ze skupiny podobných dokumentů je vybrán ten nejrelevantnější

[www.seznam.cz](http://www.seznam.cz)

Seznam.cz, a.s. | Radlická 2 | 150 00 Praha 5 | Tel.: +420 234 694 111 | Fax: +420 234 694 115

# Výkonové statistiky

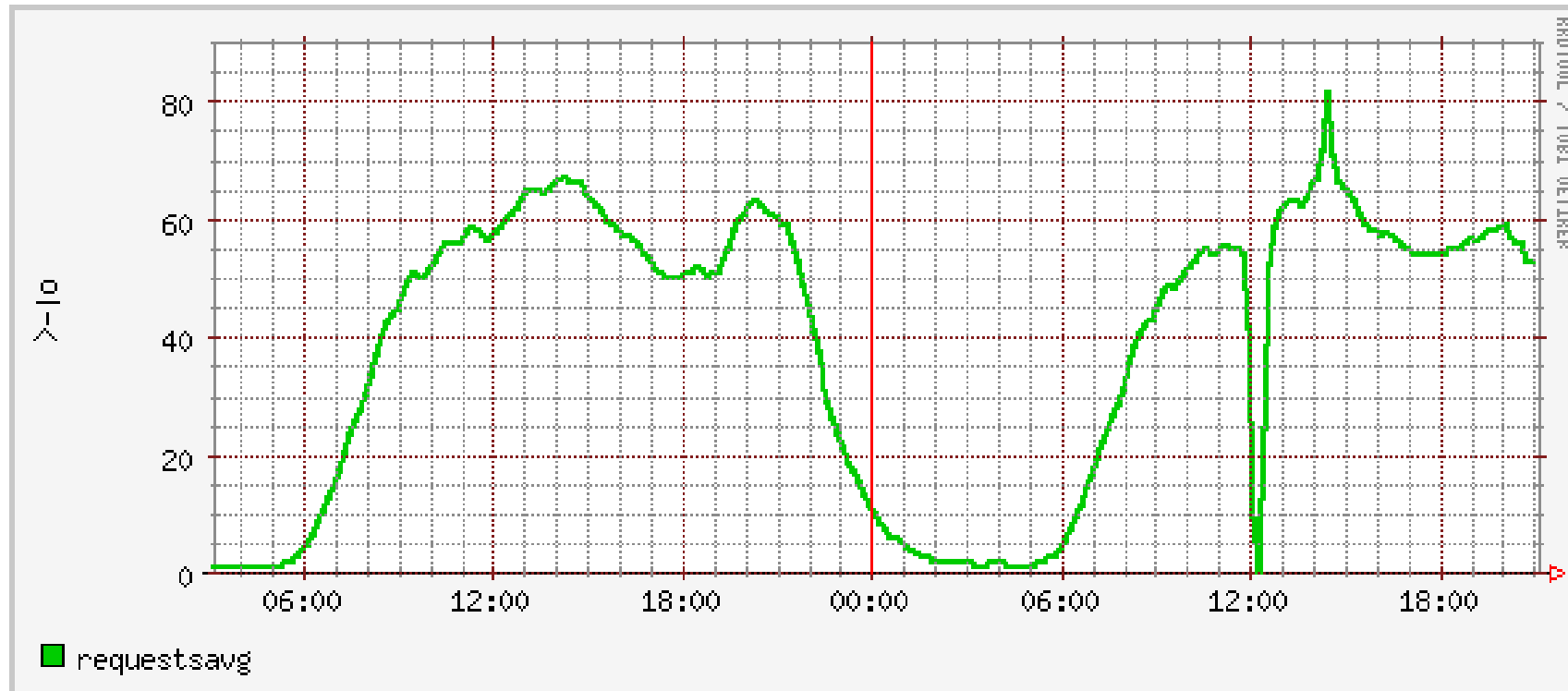
 SEZNAM

# Výkonové statistiky

1. Dotazy na web interface
2. Úspěšnost query cache
3. Odezva index searcheru
4. Odezva content serveru

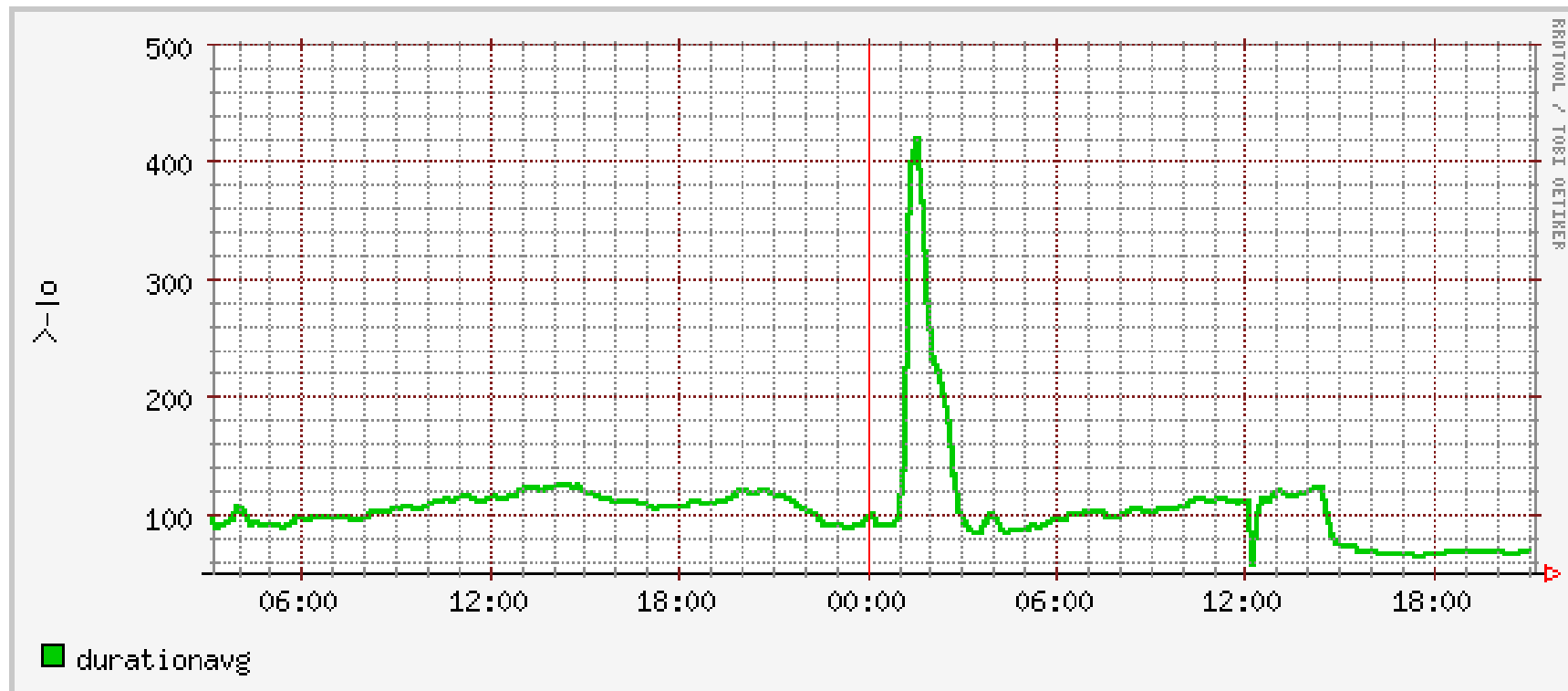


# Dotazy na web interface (1)



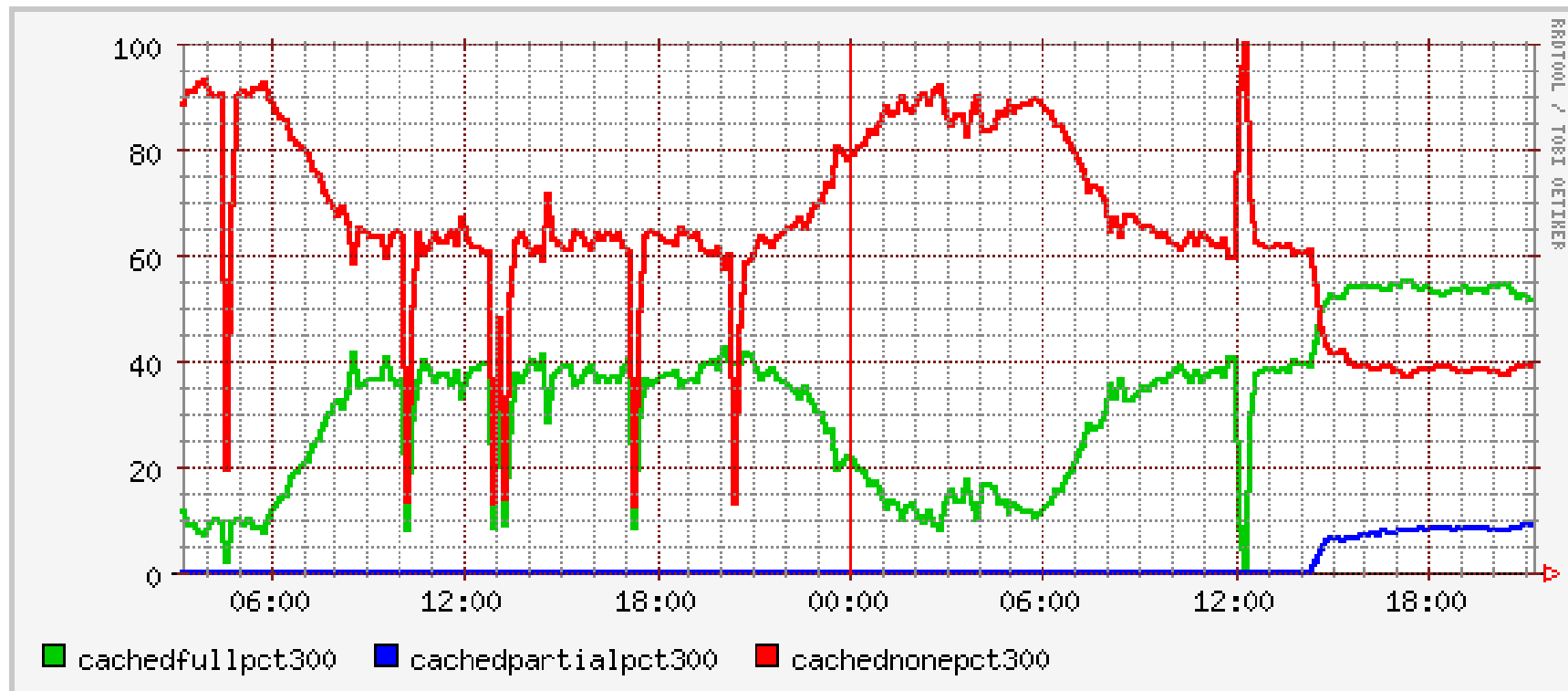
- Počet dotazů / sec
- 1 ze 3 metasearchů

# Dotazy na web interface (2)



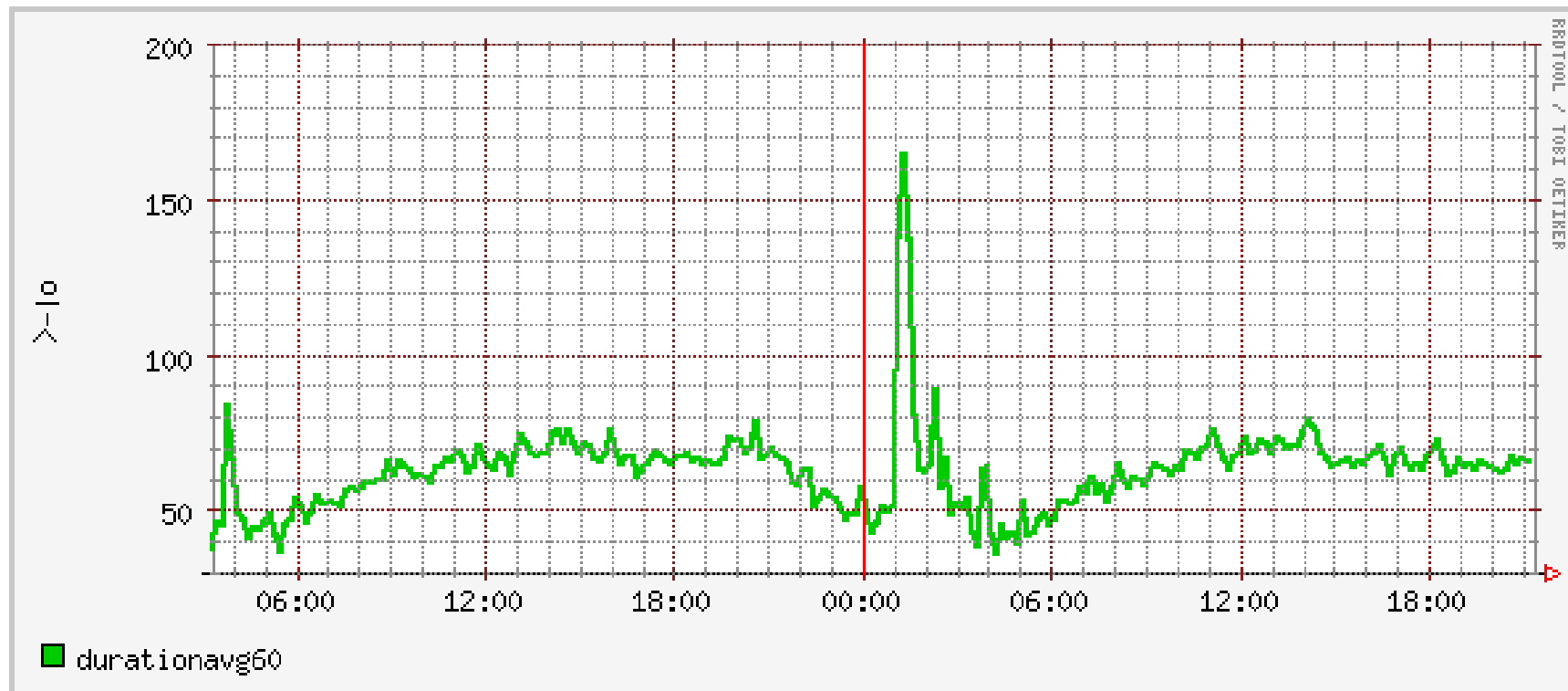
- Průměrná latence odpovědí v msec
- Špička = update db, pokles = zapnutí cache

# Úspěšnost query cache



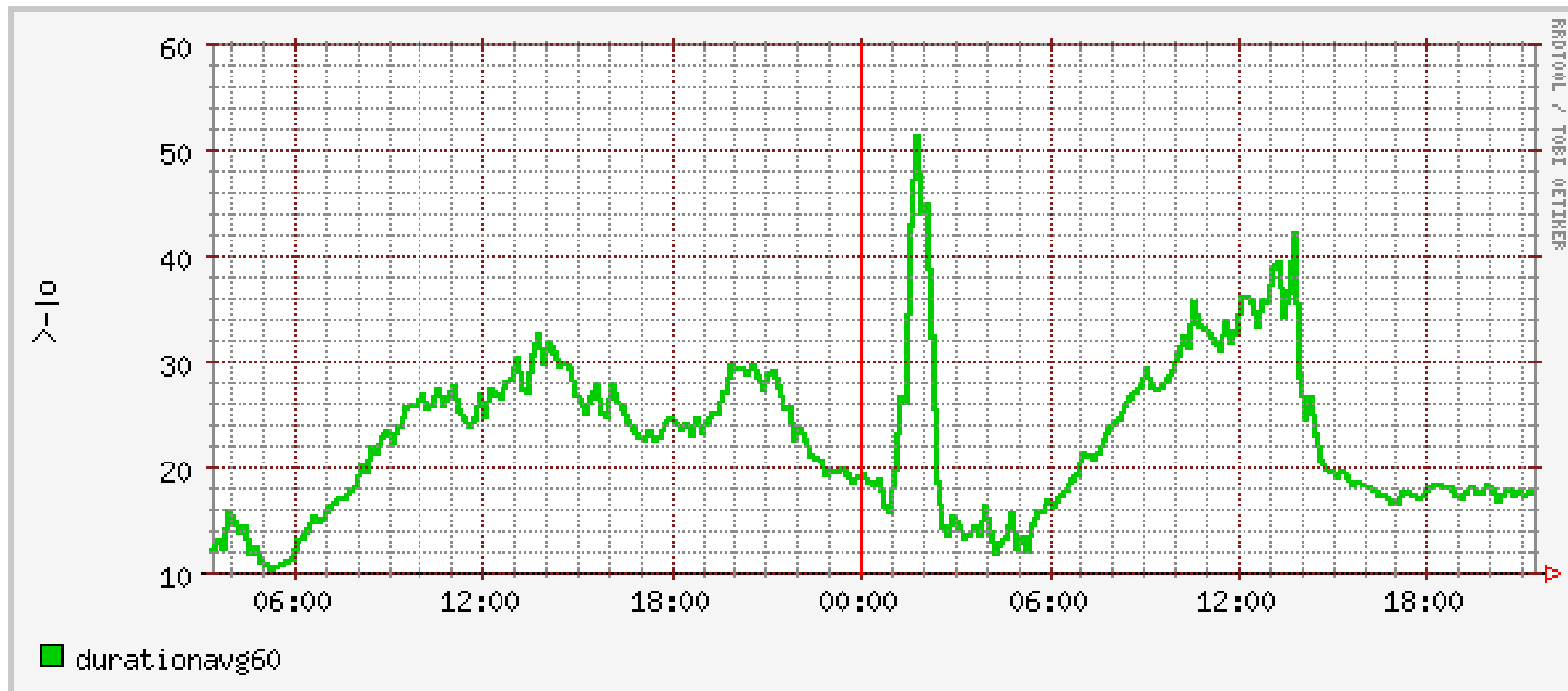
- Úspěšnost cache v %
- Skok na konci = zapnutí částečného cachování

# Odezva index searcheru



- Průměrná latence odpovědi v msec
- Špička = update db

# Odezva content serveru



- Průměrná latence odpovědi v msec
- Pokles na konci = zapnutí cachování

## Další rozvoj

- Zapojení pokročilých lingvistických metod
- Zvyšování počtu dokumentů v databázi
- Zlepšování výpočtu relevance
  
- ... Zde je šance pro Vás!

# Konec

Děkuji za pozornost.