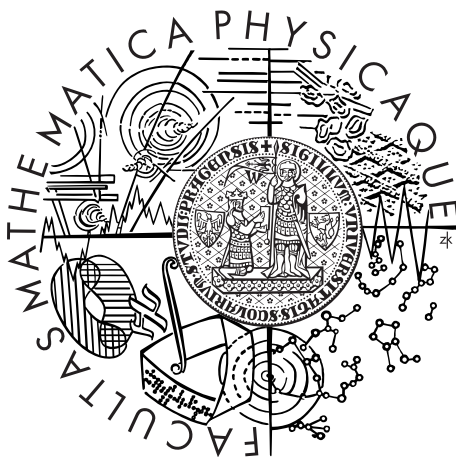


Univerzita Karlova v Praze  
Matematicko-fyzikální fakulta

## DIPLOMOVÁ PRÁCE



Bc. Ondřej Odcházal

## Automatické doporučování ilustračních snímků

Ústav formální a aplikované lingvistiky

Vedoucí diplomové práce: RNDr. Pavel Pecina, Ph.D.

Studijní program: Informatika

Studijní obor: Matematická lingvistika

Praha 2014

Děkuji svému vedoucímu, panu Pecinovi za pomoc a cenné připomínky v průběhu celého vývoje projektu. Také bych rád poděkoval své rodině a přítelkyni hlavně za velkou trpělivost.

Prohlašuji, že jsem tuto diplomovou práci vypracoval(a) samostatně a výhradně s použitím citovaných pramenů, literatury a dalších odborných zdrojů.

Beru na vědomí, že se na moji práci vztahují práva a povinnosti vyplývající ze zákona č. 121/2000 Sb., autorského zákona v platném znění, zejména skutečnost, že Univerzita Karlova v Praze má právo na uzavření licenční smlouvy o užití této práce jako školního díla podle §60 odst. 1 autorského zákona.

V ..... dne .....

Podpis autora

Název práce: Automatické doporučování ilustračních snímků

Autor: Bc. Ondřej Odcházal

Katedra: Ústav formální a aplikované lingvistiky

Vedoucí diplomové práce: RNDr. Pavel Pecina, Ph.D., Ústav formální a aplikované lingvistiky

Abstrakt:

Klíčová slova: vyhledávání obrazových informací

Title: Automatic suggestion of illustrative images

Author: Bc. Ondřej Odcházal

Department: Institute of Formal and Applied Linguistics

Supervisor: RNDr. Pavel Pecina, Ph.D., Institute of Formal and Applied Linguistics

Abstract:

Keywords: information retrieval, image retrieval

# Obsah

<b>1 Úvod</b>	<b>4</b>
Úvod	4
1.1 Práce s daty . . . . .	4
1.2 Extrakce klíčových slov . . . . .	4
1.3 Překlad do češtiny . . . . .	4
1.4 Detekce jazyka . . . . .	4
1.5 Webová aplikace . . . . .	5
1.6 Testování . . . . .	5
<b>2 Zadání</b>	<b>6</b>
<b>3 Poskytnutá data</b>	<b>7</b>
3.1 Profimedia dataset . . . . .	7
3.2 Profimedia dataset s detekovanými frázemi . . . . .	8
3.3 Vektor podobnosti obrázků . . . . .	8
<b>4 Vyhledávání relevantních obrázků</b>	<b>9</b>
4.1 Extrakce klíčových slov . . . . .	9
4.2 TF-IDF . . . . .	10
4.3 TF-IDF pro extrakci klíčových slov . . . . .	10
4.4 Vyhledávání obrázků . . . . .	11
<b>5 Překlad obrázkových popisků</b>	<b>12</b>
5.1 Strojový překlad . . . . .	12
5.2 Frázový statistický strojový překlad . . . . .	12
5.3 Charakteristika dat pro překlad . . . . .	13
5.4 Překlad vět . . . . .	13
5.5 Překlad slov . . . . .	14
5.6 Překlad frází . . . . .	14
5.7 Řešení . . . . .	14
5.8 Závěr . . . . .	15
<b>6 Backend</b>	<b>16</b>
6.1 Databáze . . . . .	16
6.2 Programovací jazyk . . . . .	16
6.2.1 NodeJS . . . . .	17
6.2.2 Go . . . . .	17
6.2.3 Ruby a Ruby on Rails . . . . .	18
6.3 Komunikace frontend-backend . . . . .	18
6.3.1 Formát dat . . . . .	18
6.3.2 REST API . . . . .	18
6.3.3 Websocket . . . . .	19
6.4 Shrnutí . . . . .	19

<b>7 Stemmer</b>	<b>20</b>
7.1 Český stemmer . . . . .	20
<b>8 Detekce jazyka</b>	<b>21</b>
8.1 Algoritmus na detekci . . . . .	21
<b>9 Podobné obrázky</b>	<b>22</b>
9.1 Bootstrap implementace . . . . .	22
9.2 Předgenerování výsledků . . . . .	23
9.3 Geohash . . . . .	24
<b>10 Frontend</b>	<b>26</b>
10.1 Možnosti programování frontendu . . . . .	26
10.2 JavaScriptové frameworky . . . . .	26
10.2.1 jQuery . . . . .	26
10.2.2 Google Closure . . . . .	27
10.3 Stylování uživatelského rozhraní a CSS . . . . .	27
<b>11 Instalace a zprovoznění</b>	<b>28</b>
11.1 Instalace . . . . .	28
11.2 Práce s metadaty k obrázkům . . . . .	28
11.3 Překlad metadat . . . . .	28
11.3.1 Export slov a frází . . . . .	30
11.4 Jazykový korpus . . . . .	30
11.5 Příprava dat pro detekci jazyků . . . . .	31
11.6 Import dat do databáze . . . . .	31
11.7 Import metadat obrázků . . . . .	31
<b>12 Anotační rozhraní</b>	<b>32</b>
12.1 Instalace rozhraní . . . . .	32
12.2 Přidání uživatelů . . . . .	32
12.3 Import anotačních dat . . . . .	32
12.4 Export anotačních dat . . . . .	33
12.5 Import obrázků a textů . . . . .	34
12.6 Anotační proces . . . . .	34
<b>13 Testování výsledků</b>	<b>36</b>
13.1 Vyloučení narušitele (o_test1) . . . . .	36
13.2 Detekce správného obrázku (o_test2) . . . . .	37
13.3 Shrnutí . . . . .	38
<b>14 Závěr</b>	<b>39</b>
<b>Závěr</b>	<b>39</b>
<b>Seznam použité literatury</b>	<b>40</b>
<b>Seznam tabulek</b>	<b>41</b>
<b>Seznam použitých zkratk</b>	<b>42</b>



# 1. Úvod

Cílem diplomové práce je implementovat kompletní webovou aplikaci pro doporučování a vyhledávání ilustračních obrázků v textu. Vytvořit takovou aplikaci přináší mnoho rozličných úkolů a problémů. Tato kapitola se bude snažit tyto problémy načrtnout. Další kapitola se bude jednotlivými problémy zabývat podrobně.

## 1.1 Práce s daty

Zadaná data obsahují 20 milionů anotací obrázků. Základním úkolem je být schopen takové množství dat vůbec nahrát do databáze a být schopný obsloužit mnoho požadavků za minutu. Bude zmíněn současný stav vývoje databázového software pro práci s velkými daty zejména s ohledem na snadnost hledání a škálovatelnost.

## 1.2 Extrakce klíčových slov

Extrakce klíčových slov je důležitý podobor NLP. V práci budou rozebrány algoritmy pro extrakci klíčových slov. Bude kladen zejména důraz na rychlost a nenáročnost na zdroje. Z uživatelských testování společnosti Google vychází, že rychlost načtení stránky je jedním z klíčových vlastností pro spokojenost uživatele. Klíčová slova budou mít v aplikaci dvě využití. Pokud uživatel zadá pouze text článku, extrahovaná klíčová slova se použijí na vyhledávání relevantních obrázků. Prvních několik klíčových slov bude navíc použita jako nápověda uživateli, ten pak může tato klíčová slova využít k exaktnímu omezení množiny klíčových slov.

## 1.3 Překlad do češtiny

Popisky klíčových slov jsou v angličtině. Tato práce řeší překlad množiny klíčových slov do češtiny. Kromě překladu je pro hledání také nutno implementovat algoritmus na stemming. Celá aplikace je navržena tak, aby případný další jazyk mohl být přidán co nejjednodušeji.

## 1.4 Detekce jazyka

Jednou z drobností, kterou ocení uživatel aplikace je detekce jazyka. Uživatel bude mít možnost zadat jazyk vstupního článku exaktně, ale aplikace bude také jazyk vstupního textu sama detekovat. Budou prozkoumány možnosti detekce jazyka. Opět se nejedná o nějakou klíčovou funkci aplikace. Výstup detekce bude moci být uživatelem měněn (podobně jako funguje Google Translate<sup>1</sup>), důraz bude tedy kladen na rychlost a jednoduchost.

---

<sup>1</sup><https://translate.google.com/>



## 1.5 Webová aplikace

Všechny předchozí komponenty se spojí v jedné webové aplikaci. Webový vývoj zažívá bouřlivý rozvoj. Na backendu jsou nové zejména způsoby práce s velkým množstvím dat v distribuovaném prostředí. Ve frontendové části probíhá rozvoj pomocí implementace nových technologií, známých pod hlavičkou HTML5, do moderních prohlížečů. Práce bude rozebírat všechny možnosti tvorby moderních webových aplikací.

## 1.6 Testování

Aplikace bude otestována na několika úrovních. Extrakce klíčových slov bude otestována pomocí korpusu článků a klíčových slov. Bude vytvořena komplexní webová aplikace pro testování doporučených obrázků. Tato aplikace bude vydělena ze samotné webové aplikace a bude používána i nezávisle.

## 2. Zadání

Většina zpravodajských serverů často opatřuje publikované články tzv. ilustračními snímky, jejichž úkolem je vizuálně dokreslovat obsah článku a upoutat na něj čtenářovu pozornost. Ilustrační snímky většinou pocházejí z rozsáhlých fotografických databází, jsou vybírány autory článku a s obsahem článku souvisejí jen relativně volně. Výběr ilustračních snímků probíhá nejčastěji na základě porovnávání klíčových slov specifikovaných autorem textu a popisků, kterými jsou obrázky v databázi opatřeny (typicky svými autory).

Proces výběru ilustračních snímků (dotazování ve fotografické databázi) je obtížný jednak pro samotný vyhledávací systém (hledání relevantních fotografií na základě uživatelských dotazů), jednak pro autory, kteří musí dotazy vytvářet. Konstrukce dotazů spočívá v několika krocích: uživatel nejdříve musí identifikovat ústřední téma (či témata) článku, které chce ilustrovat vhodnou fotografií, a ta potom popsat vhodnými klíčovými slovy, zvolit a zkombinovat je tak, aby vedla k nalezení vhodného obrázku. Tento proces by mohl být zjednodušen tím, že konstrukce dotazů pro vyhledávání bude prováděna automaticky pouze na základě textu článku.

Cílem diplomové práce je navržení a implementace komfortní webové aplikace pro automatické navrhování ilustračních snímků na základě textu článku, bez nutnosti explicitně konstruovat vyhledávací dotazy. Součástí práce bude i uživatelská evaluace celého systému. Pro experimenty bude použita kolekce ilustračních snímků od společnosti Profimedia.

## 3. Poskytnutá data

Tato kapitola popisuje data, která slouží jako vstup projektu.

### 3.1 Profimedia dataset

Společnost Profimedia poskytla pro výzkumné účely korpus více než dvaceti milionů ilustračních obrázků spolu s jejich textovým popisem.

Textové popisky byly dále očištěny[2] a poskytnuty pro tento projekt ve formě souboru `profi-text-cleaned.csv`. Soubor je ve formátu CSV a obsahuje 20 014 394 řádků. Každý řádek obsahuje 4 složky:

**locator**

Identifikátor obrázku v databázi Profimedia. Desetimístný řetězec číslic.

**title**

Název obrázku v anglickém jazyce.

**description**

Pro všechny řádky souboru prázdná položka.

**keywords**

Mezerami oddělená klíčová slova obrázku.

Ukázka 3.1 Příklad obsahuje příklad jednoho řádku souboru `profi-text-cleaned.csv`.

```
1 "0000000980","hradec kings holy ghost cathedral","",  
  "outdoors nobody urban scenes architecture houses  
  towers czech czech republic europe buildings build  
  history historical churches church fronts holy  
  ghost cathedral spirit ceska republika cathedrals  
  sv hradec kralove"~M
```

Ukázka 3.1: Řádek souboru `profi-text-cleaned.csv`

Na příkladu jsou vidět některé problémy, které data z datasetu Profimedia mají. Některá klíčová slova, jako například „czech republic“ k sobě patří a tvoří frázi. V souboru ovšem tyto víceslovné fráze nejsou vyznačené. Některým slovům chybí diakritika, například „ceska“. Některé fráze vznikly asi strojovým překladem z cizího jazyka do angličtiny. To je vidět na frázi „hradec kings“, která zřejmě původně byla názvem českého města „Hradec Králové“. Všechna slova v souboru obsahují pouze malá písmena, což například znesnadňuje detekci názvů.

Všechny popsané nedostatky dat negativně ovlivňují možnosti pro kvalitní vyhledávání v poskytnutých datech. Jedním z cílů práce je co nejvíce těchto nedostatků opravit.

## 3.2 Profimedia dataset s detekovanými frázemi

Některé problémy s datasetem Profimedia byly odstraněny v rámci bakalářské práce Bc. Jana Botorka[1]. Jedním z výsledků této práce je soubor `keyword-clean-phrase-export.csv`. Ukázka 3.2 obsahuje příklad jednoho řádku tohoto souboru.

```
1 "0000000980";"Hradec kings holy ghost  
cathedral";"outdoors,nobody,urban  
scenes,architecture,houses,towers,czech,czech  
republic,Česká republika,Europe,  
buildings,build,history,historical,  
churches,church,fronts,cathedrals sv,holy ghost  
cathedral,spirit,Hradec Králové"
```

Ukázka 3.2: Řádek souboru `keyword-clean-phrase-export.csv`

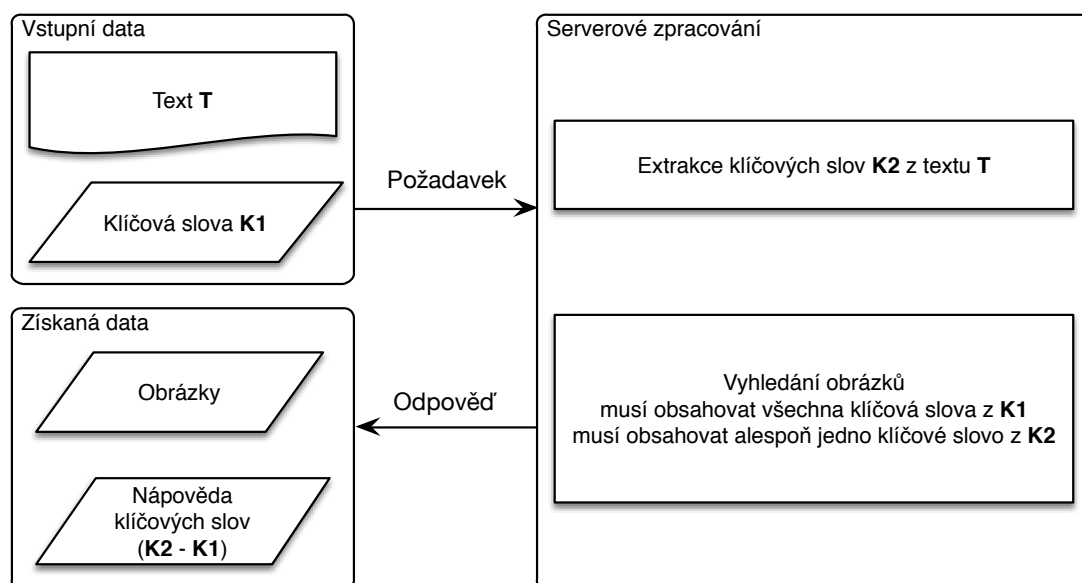
Soubor obsahuje lépe zpracovaná data z datasetu Profimedia. Nejdůležitější změnou pro tuto práci je detekce klíčových frází. Jednotlivé fráze jsou od sebe odděleny čárkou. Detekce probíhala pomocí databáze WordNet[9] a Wikipedie. Ne vždy se detekce povedla, takže v některých řádcích nejsou žádné fráze detekovány. Detekce frází je důležitá zejména pro překlad dat do jiných jazyků.

## 3.3 Vektor podobnosti obrázků

Kromě textových popisků jsou součástí datasetu Profimedia i samotné obrázky. Jedním z cílů práce je umožnit nad databází obrázků vyhledávání na základě vizuální podobnosti. K tomuto účelu byl poskytnut soubor `profi-neuralnet-20M.data.gz`, který je komprimovaný formátem ZIP a má velikost 129 GB.

Soubor obsahuje pro každý obrázek z datasetu Profimedia vektor 4096 reálných čísel. Pomocí vektorové vzdálenosti lze určit, jak jsou si podobné dva obrázky mezi sebou.

## 4. Vyhledávání relevantních obrázků



Obrázek 4.1: Základní tok dat mezi klientem a serverem.

Aplikace poskytuje uživateli dvě možnosti vyhledávání. Základní uživatelským scénářem je vložit do rozhraní text. Aplikace by v takovém případě měla poskytnout relevantní obrázky k danému textu. V dalším uživatelským scénářem je přímý požadavek na klíčová slova obrázku. Uživatel by měl mít možnost zadat přímo klíčová slova, které nalezené obrázky musí obsahovat. Oba scénáře by mělo navzájem být možné propojit pomocí nápovědy klíčových slov – uživatel zadá do rozhraní text, dostane výsledné obrázky a nápovědu klíčových slov, kterými může množinu nalezených obrázků více omezit. Klíčová slova ze vstupního textu se dále používají jako nápověda uživateli, který jimi může množinu hledaných obrázků dále omezit. Celý proces hledání vhodných obrázků popisuje diagram 4.1.

### 4.1 Extrakce klíčových slov

Extrakce klíčových slov je velmi důležitou složkou celého vyhledávání. Článek [8] shrnuje základní techniky extrakce klíčových slov z textu. Algoritmy na extrakci klíčových slov lze v zásadě rozdělit do dvou kategorií – „s korpusem“ a „bez korpusu“. Metody pracující bez korpusu jsou zajímavé a mohou dosahovat podobných výsledků jako metody s korpusem. My však máme k dispozici dataset Profimedia, takže o metody pracující bez korpusu se tato práce dále nezajímá.

## 4.2 TF-IDF

TF-IDF je jeden ze základních vyhledávacích algoritmů. Algoritmus využívá korpusu dokumentů  $D$  a dvou složek  $TF$  a  $IDF$ , lze ho vyjádřit jako rovnost:

$$TFIDF(t, d, n, N) = TF(t, d) \times IDF(n, N) \quad (4.1)$$

Složka  $TF$  znamená *TERM FREQUENCY* a pokud  $t$  je slovo a  $d \in D$  je dokument, je  $TF$

$$TF(t, d) = \sum_{slovo \in d} \begin{cases} 1 & \text{pokud } slovo = t \\ 0 & \text{jinak} \end{cases} \quad (4.2)$$

Jedná se tedy o frekvenci slova v dokumentu.

Složka  $IDF$ , tedy *INVERSE DOCUMENT FREQUENCY* vyjadřuje, jak moc daný termín popisuje dokument. Pokud je  $N$  počet všech dokumentů v  $D$ , tedy  $N = |D|$  a  $n$  je počet dokumentů, ve kterých se vyskytuje slovo  $t$ , je  $IDF$  tohoto slova

$$IDF(n, N) = \log \left( \frac{N}{n} \right) \quad (4.3)$$

$$IDF(n, N) = \log \left( \frac{N - n}{n} \right) \quad (4.4)$$

Čím je tedy slovo v korpusu častější, tím více se s logaritmem snižuje jeho informační hodnota. Slova, která jsou velmi běžná většinou klíčovými slovy nejsou.

Výsledný vzorec pak jde shrnout jako:

$$TFIDF(t, d, n, N) = \left( \sum_{slovo \in d} \begin{cases} 1 & \text{pokud } slovo = t \\ 0 & \text{jinak} \end{cases} \right) \times \log \left( \frac{N - n}{n} \right) \quad (4.5)$$

## 4.3 TF-IDF pro extrakci klíčových slov

Algoritmus TF-IDF můžeme použít pro extrakci klíčových slov z textu. Všechna znaky vstupního textu převedeme na malá písmena a text rozdělíme na slova. Dále nás nezajímá diakritika a různé speciální znaky, které může text obsahovat. Všechna slova převedeme na stemy. O převodu slov na stemmy se podrobněji píše v Kapitole 7.

Algoritmus TF-IDF pak použijeme na každé slovo vstupního textu a získáme tak skóre jeho významnosti v textu. Slova s nejvyšším skóre pak označíme za klíčová. Algoritmus musíme upravit pro naše účely. Zaprvé nás zajímají pouze taková slova, která existují v datasetu Profimédie. Slova která se v korpusu nenachází dostanou skóre 0. Část  $TF$  v algoritmu znamená četnost slova v textu.

Složitější je situace s  $IDF$ . Nabízí se použít frekvenci slov z datasetu Profimédie. Ukázalo se, že klíčová slova obrázků z datasetu nejsou pro tento účel vhodná. Klíčová slova totiž obsahují mnoho názvů a obecně méně běžných slov. Naopak obsahují velmi málo běžných slov. To pak způsobuje, že algoritmus nad

tímto datasetem dává vysoké skóre běžným slovům. Tato slova ale typicky nejsou vhodnými klíčovými slovy daného textu. Je tedy potřeba použít pro vzorec DF jiný korpus. V naší aplikaci jsme použili korpus Wikipedie, jejíž data jsou volně dostupná pro mnoho jazyků pod otevřenou licenci.

Pokud je  $w$  slovo vstupního textu,  $Freq(w)$  je četnost slova ve vstupním textu a  $Wiki(w)$  je četnost slova v korpusu Wikipedie, můžeme každému slovu vstupního slova přiřadit *Score*:

$$Score(w) = \begin{cases} Freq(w) \times (\frac{C}{Wiki(w)}) & w \text{ je v datasetu Profimédie} \\ 0 & w \text{ není v datasetu Profimédie} \end{cases} \quad (4.6)$$

$C$  je experimentálně zjištěná konstanta. V našem případě je její hodnota 10000000.

Nyní tedy máme skóre udávající význam slova pro každé slovo vstupního textu. Slova s nejvyšším skóre vrátíme uživateli jako náповědu pro explicitní požadavek obrázků s klíčovými slovy.

## 4.4 Vyhledávání obrázků

Samotné vyhledávání obrázků má dva druhy vstupních dat. Prvním typem je text článku, který uživatel vloží do uživatelského rozhraní. Máme tedy k dispozici řetězec s textem článku. Dále může uživatel zadat explicitní klíčová slova, které má hledaný obrázek obsahovat. Tato vstupní data získáváme jako pole řetězců. Úkolem algoritmu na vyhledávání obrázků je vrátit uživateli všechny relevantní obrázky v pořadí podle relevance.

Text si nejprve zpracujeme pomocí algoritmu uvedeném v sekci 4.3. Získáme skóre významnosti pro všechna slova ve vstupním textu. Pro vyhledávání použijeme pouze několik slov s nejvyšším skóre.

Nyní můžeme jako množinu relevantních obrázků označit obrázky, které ve svých klíčových slovech mají všechna uživatelem zadaná explicitní klíčová slova a alespoň jedno klíčové slovo získané extrakcí klíčových slov z textu.

Množina relevantních obrázků může být velká a obsahovat obrázky, které jsou relevantní jen velmi málo. Je proto důležité množinu relevantních obrázků správně seřadit. Máme množinu  $K_{text}$  klíčových slov získaných z textu. Každé  $w \in K_{text}$  má skóre  $Score(w)$ . Relevantní obrázek má množinu klíčových slov  $K_{img}$ . Relevanci obrázku vůči uživatelskému dotazu můžeme ohodnotit funkcí *Rank*:

$$Rank(K_{text}, K_{img}) = \sum_{w \in K_{text}} \begin{cases} Score(w) \times \frac{1}{|K_{img}|} & w \in K_{img} \\ 0 & w \notin K_{img} \end{cases} \quad (4.7)$$

Vzorec bere v úvahu počet klíčových slov, které obrázek obsahuje. Snižuje relevanci obrázků, které obsahují mnoho klíčových slov a zvýhodňuje tím ve vyhledávání obrázky, které mají menší počet přesnějších klíčových slov.

## 5. Překlad obrázkových popisků

Popisky obrázků v datasetu Profimédie jsou v angličtině. Jedním z cílů aplikace je poskytnout kromě angličtiny vyhledávání i v jiném jazyce. Zaměřujeme se na češtinu, ale podobné úvahy a postupy platí většinou i pro jiné jazyky. Jsou v podstatě dvě možnosti, jak implementovat vyhledávání v češtině. Můžeme buď překládat do angličtiny hledané texty a klíčová slova, nebo předpřeložit dataset Profimédie.

My jsme se rozhodli pro druhou možnost. Musíme tedy přeložit všechny texty v datasetu Profimédie. Vzhledem k tomu, že obrázků je více než 20 milionů, nepřipadá lidský překlad v úvahu z časových i finančních důvodů. Jedinou reálnou možností je použít překlad strojový.

### 5.1 Strojový překlad

Strojový překlad jako obor zažívá velký rozvoj. Potřeba překladu stále celosvětově prudce stoupá. To jednak zvyšuje poptávku po kvalitním strojovém překladu a druhak to strojovému překladu dává velké množství lidských překladů, které jsou pro kvalitní strojový překlad k dispozici. Lze brát v úvahu i pravidlové překladové systémy, které ke své práci tolik dat většinou nepotřebují. Obecnější a pro většinu jazykových párů nejlepší řešení však nabízí frázový statistický strojový překlad.

S tím se zvyšuje poptávka po lepším strojovém překladu a zároveň také vzniká více lidských překladů, které pak jde použít jako zdrojová data k překladům strojovým. Existují různé metody strojového překladu. Na jednu stranu existují pravidlové systémy, které podle pravidel zapsaných překladateli převádí text ze zdrojového do cílového jazyka. Na druhou stranu jedním z nejobecnějších řešení je frázový překlad.

### 5.2 Frázový statistický strojový překlad

Frázový statistický strojový překlad[6] ke své práci potřebuje databázi lidsky přeložených frází. Fráze jsou několikáslovné kusy přeloženého textu. Typicky se získávají z paralelního korpusu textů ve zdrojovém a cílovém jazyce. Soubor s frázovými daty obsahuje položku s textem fráze ve zdrojovém a cílovém jazyce spolu s pravděpodobností, že je takový překlad fráze správný.

Druhým důležitým vstupem strojového frázového systému je korpus cílového jazyka, ze kterého se vytvoří jazykový model. Slouží zejména k tomu, aby k sobě poskládané fráze v cílovém jazyce dobře „seděly“. Pokud máme překladový i jazykový model, můžeme vyjádřit pravděpodobnost překladu pomocí základního vzorce statistického strojového překladu. Překládáme řetězec  $f$  ve zdrojovém jazyce. K dispozici máme překladový model  $p(f|e)$ , který udává pravděpodobnost toho, že řetězec  $f$  ve zdrojovém jazyce je překladem řetězce  $e$  v cílovém jazyce. Jazykový model  $p(e)$  nám vrací pravděpodobnost řetězce  $e$  v cílovém jazyce. Chceme získat takový řetězec  $\tilde{e}$ , pro který je pravděpodobnost  $p(\tilde{e}|f)$  nevyšší. Pomocí Bayesova pravidla můžeme k hledání takové fráze použít překladový a jazykový model:



$$\tilde{e} = \arg \max_{e \in e^*} p(e|f) = \arg \max_{e \in e^*} p(f|e)p(e) \quad (5.1)$$

V praxi je potřeba u kvalitního strojového překladu vyřešit mnoho dalších problémů. K lepším výsledkům překladu potřebujeme reordering model, který umožňuje přesouvat pozice frází mezi zdrojovým a cílovým textem.

Dobrá frázový překlad potřebuje ke svému chodu miliony frází. Vyhledávání nejpravděpodobnějšího překladu v takovém množství dat je náročná úloha. Algoritmy, které umožňují vyhledávat ve velkém množství překladových dat jsou implementována v knihovně Moses[7], která je dostupná pod volnou licencí včetně základních překladových modelů pro některé jazykové páry.

### 5.3 Charakteristika dat pro překlad

Pro správné fungování vyhledávání v českém jazyce je potřeba přeložit klíčová slova v Profimedia datech do češtiny. Slova v korpusu jsou oddělena mezerou. Je ale zřejmé, že některá slova vedle sebe k sobě patří — jsou to fráze — zatímco některá nikoliv. Naskytují se tedy zhruba tři možnosti, jak takový text přeložit.

### 5.4 Překlad vět

<p><b>Zdrojový dokument</b></p> <p>"0000000102", "young woman cleaning teeth", "", "single faces people humans young youth hands indoors interiors woman women females blond fair young adult s girls close view beauty home home dental bathrooms person portrait adult years half length portrait open mouth hygiene teeth dental care years cleaning toothbrush underwear bras" ^M</p>
<p><b>Moses</b></p> <p>"0000000102", "young žena čištění teeth", "", "single čelí mladí lidé , lidé mládež rukou uvnitř interiors žena žen , žen , blondák spravedlivé mladé dívky zavřít dospělé s cílem krásy vnitřní vnitřní stomatologické koupelny portrét dlouhé roky polovina dospělé osoby portrét otevřené úst hygienické zuby zubní kartáček prádlo bras" ^M péče let čištění</p>
<p><b>Překladač Google</b></p> <p>"0000000102", "Mladá žena čištění zubů", "", "jednotlivé tváře lidí, lidé younge mládeží ruce interiéry ženě ženám ženy ženskému blond fair mladý dospělý s dívek close view krása domov domácí zubní koupelny osoba portrét dospělý let poloviční délka portrét otevřená ústa hygienické zubů zubní péče roky čistící kartáček na zuby spodní prádlo podprsenky " ^M</p>

Obrázek 5.1: Ukázka překladu metadat k obrázku pomocí Mosese a Překladače Google.

První možností je přistupovat k souboru klíčových slov u každého obrázku jako k větě a použít frázový strojový překlad — buď Moses, nebo Překladač Google[4] — k překladu z angličtiny do češtiny. Tento přístup má několik problémů. Frázový překlad se snaží aplikovat fráze z překladového modelu. V našem souboru klíčových slov ale mohou být vedle sebe slova, která tvoří frázi pouze zdánlivě. Například můžeme mít fotku dítěte, které stojí před automobilem značky Seat se dvěma klíčovými slovy vedle sebe — „child“ a „seat“. Frázový překlad

z angličtiny do češtiny pochopí tato dvě slova jako fráze, které do češtiny přeloží frází „dětské sedadlo“, která ovšem neodpovídá popisku obrázku.

Dalším problémem tohoto přístupu je časová a finanční náročnost takového řešení. Frázový překlad je dosti náročný algoritmus a překlad dvaceti milionů vět může být dosti obtížný. Překlad pomocí Mosesa nás omezuje výpočetní složitostí. Na jednom stroji by takový překlad trval řádově několik dní. Pokud bychom k překladu 20 milionů vět použili Překladač Google, jsme zase omezení cenou za přístup k překladovému API společnosti Google.

Pro práci s Mosesem jsme zvolili překladový a jazykový model pro překlad z angličtiny do češtiny, které jsou poskytnuty přímo s Mosesem ve verzi XX. Obrázek 5.1 porovnává překlad metadat k jednomu z obrázků datasetu Profimedia v Mosesovi s distribuovanými překladovými daty s překladem pomocí Překladače Google. Je vidět, že Překladač Google poskytuje kvalitnější překlad.

## 5.5 Překlad slov

Jednodušším přístupem k překladu klíčových slov je přístup slovníkový, teda překlad každého slova zvlášť. Nejprve je potřeba ze souboru klíčových slov u všech obrázků vyextrahovat slova. Ty pak lze přeložit přímo s použitím slovníku, nebo pomocí frázového strojového překlad (ten použije jednoslovné fráze) Mosesem, či Překladačem Google. Výhodou oproti předchozímu přístupu je menší množství dat a tedy i nižší finanční a časová náročnost takového překladu. Takový systém překladu ale nedokáže detekovat fráze a kvalita překladu slovo od slova je typicky horší než překlad delších frází. Vezměme si například anglická slova „weather“ a „vane“. Slovníkový překlad nám slova přeloží jako „počasí“ a „lopatka“. Pokud se ale tato slova nachází vedle sebe, je pravděpodobnějším překladem slovo „korouhvička“.

## 5.6 Překlad frází

Poslední možností je oba předchozí principy zkombinovat — nejprve detekovat v souboru klíčových slov fráze a ty pak přeložit statistickým strojovým překladem. Detekci frází z datasetu Profimedia provedl ve své bakalářské práci[1] Jan Botorek. Zkoušel detekovat N-gramy v databázi WordNet[9] a Wikipedii. Výsledky této detekce frází lze použít právě ke zlepšení překladu z angličtiny do cizích jazyků. Na slova která nejsou detekována ve frázi se použije slovníková metoda. Na překlad detekovaných frází lze použít přímo statistický strojový překlad.

## 5.7 Řešení

Řešením bylo nakonec použít poslední možnost. Bližší informace jsou v Kapitole [?]. Slova byly přeloženy pomocí Překladače Google, který překládá mnohem lépe než dostupný model pro překladový nástroj Moses.

Překlad frází by ovšem pomocí Překladače Google byl příliš drahý a pomocí Mosesa zase příliš pomalý. Použili jsme tedy jednoduchou metodu, která používá pouze překladový model. Fráze z datasetu Profimedia byla přeložena pouze

tehdy, když se její překlad nacházel v překladovém modelu. Fráze, které se v překladovém modelu nenacházely, byly přeloženy slovo od slova. Metodou přesné shody s překladovým modelem bylo přeloženo 18006 z celkového počtu 899244 detekovaných frází.

## 5.8 Závěr

Překlad klíčových slov z korpusu Profimedia není typickou překladovou úlohou — nepřekládají se celé věty. Přesto je dokonalý výsledek nemožný. I lidští překladatelé s citem pro jazyk by v této překladové úloze dávali rozdílné výsledky. Strojový překlad zdaleka není na takové úrovni, aby dokázal z širšího kontextu vybrat správný překlad. Navíc lidský překladatel může u překladu klíčových slov využít přímo obrázek, ke kterému se klíčová slova vztahují. Může tak z obrázku posoudit, jestli má slovo „single“ přeložit jako „jednolůžkový“, nebo ve významu „jeden“. Lepší překlad by mohl přinést překladový model natrénovaný na speciálnější množině dat bližší korpusu Profimedie. Vytvořit takový model by ovšem bylo nad rámec této práce.

Navrhovaný mechanismus překladu poskytuje uspokojivý, i když značně nedokonalý, překlad klíčových z angličtiny do češtiny. Jednoduše jde zevšeobecnit i pro překlad do dalších jazyků pro které máme slovník a překladový model.

## 6. Backend

Backend je jednou z klíčových částí aplikace. Kromě obsluhy uživatele statickými soubory (HTML, CSS, Javascript) je jeho hlavní úlohou poskytnout API pro vyhledávání.

### 6.1 Databáze

Úkolem databáze je uložit data a umožnit jejich prohledávání. Uživatelé této aplikace nemají možnost databázi modifikovat. Zápis do databáze provede administrátor pouze jednou před startem aplikace. Důležitým požadavkem je důraz na rychlost a snadnou škálovatelnost. V posledních několika letech vzniklo mnoho nových databází v kategorii vágně označené jako NoSQL. Tato kategorie databází se těžko popisuje, na každou popsanou vlastnost existuje NoSQL databáze, která tuto podmínku nesplňuje. Obecně ale lze říct že NoSQL databáze nepracují s prvky v tabulkovém uspořádání. Jejich výhodou oproti standardním relačním databázím může být vyšší výkon a snadná škálovatelnost.

Nevýhodou je většinou obtížnější práce s daty. Většina NoSQL databází například mapodporuje databázové transakce a vůbec celý ACID. Pro práci s databázovými daty se často používá model Map Reduce. Algoritmus Map Reduce vyvinula a publikovala společnost Google, která na něj má i patent. Google však oznámil, že algoritmus MapReduce postupně přestává používat<sup>1</sup>.

První verze aplikace byla postavena na databázi CouchBase. Výhodou CouchBase je dobrý výkon, velmi dobrá dokumentace a také podpora knihoven pro mnoho jazyků přímo od společnosti Couchbase. Brzy se však ukázalo, že implementace textového vyhledávání v Couchbase by byla velmi náročná.

Vhodnější pro daný účel se ukázala knihovna Elasticsearch. Nejedná se v pravém smyslu o databázi. Jejím hlavním cílem je poskytnout snadné vyhledávání nad daty. Je postavená nad knihovnou Apache Lucene a poskytuje snadnou škálovatelnost. Komunikace mezi knihovnou a klientem probíhá pomocí RESTového HTTP API. Hlavním podporovaným formátem dat je JSON. Elasticsearch poskytuje programátorovi velké možnosti v nastavení prohledávání. V textovém vyhledávání může uživatel databáze použít všechny tokenizery a stemmery z knihovny Lucene. Důležitou vlastností je, že u hledaných slov může uživatel databáze určit váhu jednotlivých slov. V průběhu implementace aplikace navíc vyšla verze knihovny 1.0.

### 6.2 Programovací jazyk

Se vzrůstající popularitou webů vzniká stále větší množství webových frameworků a dokonce programovacích jazyků zaměřených primárně na programování pro web.

---

<sup>1</sup><http://www.datacenterknowledge.com/archives/2014/06/25/google-dumps-mapreduce-favor-new->

### 6.2.1 NodeJS

Jedním z nových trendů je tvorba webového backendu v Javascriptu pomocí knihovny NodeJS. Frontendoví vývojáři jsou prakticky nuceni Javascript používat. Pokud se v Javascriptu tvoří i backendová část aplikace, může snadněji dojít ke sdílení kódu i pracovních pozic. Trend psaní všech aplikací v Javascriptu glosoval již před sedmi lety Jeff Atwood ve svém pravidle:

```
1 Any application that can be written in JavaScript ,  
  will eventually be written in JavaScript .  
2 Vše co může být napsáno v Javascriptu , bude v  
  Javascriptu napsáno .
```

Nevýhody použití Javascriptu jako backendového programovacího jazyka jsou poměrně zřejmé. Javascript byl navržen pro programování webového frontendu, jeho standardní knihovna je v porovnání s ostatními jazyky velmi chudá, podpora objektového programování je celkem nepřímá. Obsáhlý kód v Javascriptu může být poměrně nepřehledný a jazyk svádí k vytvoření takzvaného „callback hell“, který může mít strukturu jako na obrázku (source <http://strongloop.com/strongblog/node-js-callback-hell-promises-generators/>). Pokud chce navíc programátor sdílet kód z backendu i na frontendu, musí být kód kompatibilní s podporovanými prohlížeči. Funkce pro jednodušší práci s polem podporuje webový prohlížeč Internet Explorer až od verze 9 [<http://kangax.github.io/compat-table/es5/>]. V součtu nám převažily nevýhody NodeJS frameworku nad výhodami a myšlenku vývoje backendu v Javascriptu jsme opustily.

```
1 doAsync1(function () {  
2   doAsync2(function () {  
3     doAsync3(function () {  
4       doAsync4(function () {  
5       })  
6     })  
7   })  
8 })
```

### 6.2.2 Go

Go (známý také jako golang)[?] je programovací jazyk od společnosti Google. První stabilní verze byla zveřejněna v Květnu 2012. Go je kompilovaný, staticky typovaný jazyk s garbage collectorem. Syntaxe je inspirována jazykem C a přizpůsobena pro rychlou kompilaci. Velkou výhodou je snadná práce s vlákny pomocí „go rutin“. Programátoři v Javě, nebo C++ může překvapit poněkud netypická podpora práce s objekty.

První verze aplikace byla napsána právě v jazyce Go. V praxi se ukázaly všechny výše uvedené výhody. Hlavní nevýhodou se však ukázal nedostatek kvalitních knihoven. Přestože je jazyk velmi mladý, stal se velmi populárním, o čemž svědčí například počet repozitářů na GitHubu, nebo otázek na StackOverflow. Bude to ovšem trvat ještě nějaký čas, než se knihovny pro go a jejich vývoj stabilizují. Jeden z nyní nejpopulárnějších webových frameworků pro go — Martini — v

době začátku práce na aplikaci ani neexistoval. Právě nedostatek knihoven pro Go vedl k volbě jiného jazyka. Rozhodujícím pro opuštění backendu v jazyce Go byla neexistující knihovna pro práci s Elasticsearch. API je sice postaveno na protokolu HTTP, takže šlo s Elasticsearch komunikovat bez použití specializované knihovny, v praxi se to ovšem ukázalo být problematické. Vývoj aplikace si žádal testování různých nastavení a rychlé prototypování v kódu.

Z interních zdrojů máme informaci o tom, že se Elasticsearch chystá vytvořit knihovnu i pro jazyk Go. Pro naše účely se ukázalo výhodnější přepsat aplikaci do dynamického jazyka s lepší podporou pro Elasticsearch.

### 6.2.3 Ruby a Ruby on Rails

Ruby je dynamicky typovaný jazyk, silně inspirovaný Perlem s důslednou podporou objektového programování. Popularita Ruby vzrostla zejména kvůli webovému frameworku Ruby on Rails, který je v Ruby napsaný. RoR zpopularizovaly koncept Model-View-Controller při tvorbě webových aplikací. Backend aplikace byl nakonec kompletně přepsán jako aplikace v Ruby on Rails. Elasticsearch poskytuje pro práci s Ruby vlastní knihovnu. Další výhodou se ukázala podpora knihovny Rake, což je jakási obdoba Makefile skriptů v Ruby. Pomocí Rakefilu jdou napsat přehledné úlohy pro manipulaci s daty.

Standartní implementace jazyka Ruby byla v minulosti kritizována pro svou pomalost. V průběhu práce na této aplikaci vyšla verze Ruby 2.0, která Ruby dosti zrychluje. Tato rychlost byla pro aplikaci shledána dostatečnou.

## 6.3 Komunikace frontend-backend

Backend poskytuje služby frontendu pomocí API. Existuje několik přístupů a technologií, jak data mezi backendem a frontendem posílat.

### 6.3.1 Formát dat

V současnosti jsou pro posílání dat nejběžnější dva formáty – XML a JSON. XML je klasický formát se spoustou možností a nástrojů. Moderní aplikace však stále více přechází k formátu JSON. Formát JSON je velmi úsporný datový formát, který vychází z datových typů v JavaScriptu. Právě úspornost je jednou z jeho největších výhod oproti XML – stejné informace mají v JSON typicky kratší zápis než obdoba v XML. Většina moderních browserů umí formát JSON parsovat a práce v JavaScriptu je pak vzhledem ke kompatibilitě datových typů velmi pohodlná.

Zejména kvůli poslednímu důvodu používá tato práce pro komunikaci mezi frontendem a backendem právě formát JSON.

### 6.3.2 REST API

REST je zkratka pro Representational state transfer. Jedná se o obecnou architekturu rozhraní. V tomto kontextu nás ale zajímá hlavně její navázání na protokol HTTP. REST využívá metody protokolu HTTP pro změnu, nebo získání stavu datových objektů. Ukázkou REST API s použitím formátu JSON může

být například knihovna Elasticsearch. Mějme například frekvenční data pro český stem „lid“ a instanci Elasticsearch na adrese `http://localhost:9200/`. Ukázky HTTP requestů XX - XX ukazují, jak taková data uložit, získat, změnit, nebo smazat.

REST API v kombinaci s JSON formátem používá tato aplikace jak ke komunikaci mezi frontendem a backendem, tak mezi backendem a databází Elasticsearch.

### 6.3.3 Websocket

Alternativou k REST API a protokolu HTTP je protokol WebSocket. WebSocket je stejně jako protokol HTTP postaven nad protokolem TCP. Na rozdíl od HTTP ale poskytuje duplexní spojení. Klient je stále spojený se serverem a oba mohou posílat zprávy bez ohledu na druhou stranu. Spojení přes WebSocket má většinou menší latenci než použití protokolu HTTP<sup>2</sup>. Protokol je podporován všemi moderními verzemi webových prohlížečů a podpora existuje i v knihovnách pro backendové programovací jazyky.

První verze aplikace používali pro komunikaci mezi klientem a serverem právě protokol WebSocket. Nakonec však převážily nevýhody takového řešení nad výhodami. Jednou z nevýhod je nutnost udržovat spojení s klientem na serverové i klientské straně. Přináší to několik netriviálních problémů. Například v okamžik, kdy se toto spojení přeruší. Pokud se naproti tomu přeruší spojení server-klient při HTTP requestu, může klient zkusit vyslat stejný požadavek znovu.

Druhým problémem je emulace HTTP požadavků v protokolu WebSocket. Z klienta můžeme odeslat HTTP dotaz na server a dostaneme k němu přiřazenou odpověď. V protokolu WebSocket pošleme serveru zprávu a za nějaký čas můžeme dostat zprávu od serveru jako odpověď. Párování došlých zpráv z odeslanými zprávami-požadavky ale musíme implementovat vlastnoručně, například pomocí unikátních ID v těle zprávy. Navíc musíme umět řešit situaci, kdy žádná odpověď ze serveru nedorazí. Například nastavením timeoutu pro čekání na odpověď.

WebSocket využijí zejména aplikace, které potřebují, aby server mohl posílat klientovi kdykoliv zprávy, nebo co nejnižší latenci. Takovými aplikacemi mohou být například různé chatovací služby, nebo online hry. Pro jiné většinou asi převáží nevýhody WebSockets oproti HTTP protokolu.

## 6.4 Shrnutí

Architektura je shrnuta na diagramu XX. Aplikace má backend napsaný v jazyce Ruby a frameworku Ruby on Rails. Data jsou uložena v databázi Elasticsearch. Backend komunikuje s frontendem i databází pomocí REST API, data jsou přenášena ve formátu JSON.

---

<sup>2</sup><http://www.websocket.org/quantum.html>

## 7. Stemmer

Důležitou vlastností systémů na vyhledávání v textu je, aby dokázali najít i v jiných tvarech slova. Pokud uživatel hledá slovo „praha“, většinou očekává, že se mu zobrazí i výsledky obsahující slovo „praze“. Je tedy potřeba mít nějaký algoritmus, který k sobě slova jako „praha“ a „praze“.

První možností je použít lematizér. Úkolem lematizéru je ke každému slovu přiřadit jeho základní tvar. U podstatných jmen je to většinou první pád jednotného čísla („praha“), u sloves infinitiv.

Alternativou lematizátoru může být stemmer. Ten většinou používá jednoduché heuristiky k odstranění koncovek slov. Stemmer nemusí vrátit regulérní slovo jazyka, například pro slova „praha“ a „praze“ může stemmer vrátit slovo „prah“. Výhodou stemmeru oproti lematizéru je, že většinou používá pouze jednoduché heuristiky. Je tedy většinou rychlejší a méně náročný na zdroje než lematizér. Ve vyhledávacích aplikacích je stemmer dostačující.

Pro angličtinu je nejznámějším stemmerem pro angličtinu je Porterův stemmer popsáný[?] Martinem Porterem již v roce 1980. Kromě oficiální implementace existují porty do různých jazyků včetně Ruby. Tato aplikace využívá implementaci z Ruby Gemu „stemmify“<sup>1</sup> v licenci MIT.

### 7.1 Český stemmer

Pro jazyky jako je čeština, která má bohatší morfologii než angličtina, je tvorba stemmeru náročnější. Testovali jsme několik implementací českých stemmerů. Jako nejkvalitnější byla nakonec vybrána implementace českého stemmeru v knihovně Lucene. Tuto implementaci využívá i Elasticsearch.

V rámci této práce byl portován soubor `CzechStemmer.java` ze zdrojového kódu knihovny Lucene do jazyka Ruby. Výsledkem je Ruby Gem pod licencí MIT, který lze použít nezávisle na zbytku aplikace. Instaluje se příkazem

```
1 gem install czech-stemmer
```

Knihovna obsahuje pouze jednu třídu `CzechStemmer` s funkcí `stem`, která přijímá i vrací řetězec:

```
1 require 'czech-stemmer'  
2  
3 CzechStemmer.stem("praha") # => "prah"  
4 CzechStemmer.stem("praze") # => "prah"  
5 CzechStemmer.stem("předseda") # => "předsd"  
6 CzechStemmer.stem("mladými") # => "mlad"
```

---

<sup>1</sup><https://rubygems.org/gems/stemmify>



## 8. Detekce jazyka

Aplikace může přijímat vstup ve více jazycích. Jedním ze způsobů, jak uživateli zpříjemnit práci s aplikací, je použít automatický detektor jazyka. Při zadání vstupu aplikace se aplikace sama pokusí detekovat jazyk zadaného textu. Tato detekce nemusí být vždy stoprocentně správná, takže by uživatel měl mít možnost volbu jazyka manuálně změnit.

### 8.1 Algoritmus na detekci

Aplikace používá algoritmus založený na statistice nejčastějších N-gramů pro daný jazyk[?]. Algoritmus nejprve použije velký korpus textů pro každý detekovaný jazyk. Pro češtinu a angličtinu lze použít například korpus Wikipedie. Z tohoto korpusu získá uspořádaný seznam  $K$  nejčastějších N-gramů. Při samotné detekci jazyka textu, pak samotný algoritmus získá stejný uspořádaný seznam nejčastějších N-gramů pro vstupní text. Tento seznam pak porovnává se seznamy nejčastějších N-gramů pro každý jazyk.

Nechť  $A$  a  $B$  jsou seznamy N-gramů s  $K$  položkami,  $A[w]$  je pořadí N-gramu  $w$  v seznamu  $A$ . Potom lze vzdálenost mezi seznamy  $D(A, B)$  vyjádřit vztahem:

$$D(A, B) = \sum_{w \in A} \begin{matrix} |A[w] - B[w]| & \text{pokud } w \in B \\ K & \text{jinak} \end{matrix} \quad (8.1)$$

Nyní pokud máme množinu  $S$  všech seznamů N-gramů pro jednotlivé jazyky a  $X$  je seznam N-gramů pro vstupní text, vrátí algoritmus jako jazyk takový jazyk, pro který má seznam  $C \in S$  minimální hodnotu  $D(X, C)$ .

Naše implementace používá v seznamech N-gramů trigramy. Pro texty delší než několik slov funguje velmi spolehlivě. Pro několikaslovné texty se může stát, že v seznamu nejfrekventovanějších trigramů pro vstupní text není ani jeden trigram, který by se nacházel v seznamech pro jednotlivé jazyky. Pak algoritmus může vrátit špatně detekovaný jazyk.

Vylepšení by jistě přineslo spolu s použitím trigramů použít i bigramy a unigramy. Také konstanta  $K$  by šla zvětšit z používané hodnoty 50 výš. Implementace detekce jazyka ale probíhá na klientovi, takže všechna tato vylepšení algoritmu by zvýšila množství dat, která si klient musí stáhnout. Navíc uživatel má vždy možnost detekovaný jazyk manuálně změnit a typicky pracuje spíše s delšími texty. Implementace s použitím 50 nejčastějších trigramů se tedy zdá dostačující pro daný účel.

## 9. Podobné obrázky

Jednou ze služeb, které výsledná aplikace poskytuje, je vyhledávání podobných obrázků. Uživatel rozhraní najde pomocí textových dotazů nějaké ilustrační obrázky a má možnost u každého z nalezených obrázků získat obrázky vizuálně podobné. Tato kapitola pojednává o tvorbě backendové služby, která vyhledávání podobných obrázků umožňuje.

Vstupními daty je soubor s vektory pro každý obrázek datasetu Profimedia. Vektor má 4096 složek s reálnými nezápornými čísly. Vektory jsou vizuální deskriptory obrázků. Tyto deskriptory byly vygenerovány pomocí software Caffe a jsou odezvami předposlední vrstvy hluboké konvoluční neuronové sítě natrénované pro klasifikaci obrázků do 1000 kategorií.

Mějme obrázek  $I_1$  s deskriptorem  $D_1$  a obrázek  $I_2$  s deskriptorem  $D_2$ . Míru podobnosti obrázků *Similarity* pak můžeme definovat jako

$$Similarity(I_1, I_2) = \sum_{i=1}^{4096} |D_1[i] - D_2[i]| \quad (9.1)$$

V praxi se dají výsledky této míry klasifikovat zhruba do 3 kategorií. Tyto vypočítané kategorie popisuje tabulka 9. Ukázky jednotlivých kategorií podobných obrázků poskytuje obrázek 9.2. Rozdělení na tyto 3 kategorie je samozřejmě velmi volné a umělé. Nelze například zaručit, že obrázek, který by některý uživatel mohl označit za podobný, nebude mít míru *Similarity* vyšší než 1500.

Kategorie	<i>Similarity</i>
téměř shodné	0 – 500
podobné	500 – 1500
nepodobné	> 1500

Obrázek 9.1: Kategorie podobnosti obrázků podle *Similarity*.

Otázkou zůstává, které výsledky uživatel očekává jako výsledky vyhledávání podobných obrázků. Pravděpodobně nechce získat obrázky z kategorie „nepodobné“. Pak je otázkou, jestli uživatel chce jako výsledek získat obrázky z kategorie „téměř shodné“. V korpusu je spousta podmnožin obrázků, které byly vyfoceny v rámci jedné série. Často se liší jen malou změnou úhlu fotky, nebo pouze kompresí uloženého obrázku. Pokud bychom vraceli obrázky seřazené podle *Similarity*, uživatelé u takovýchto podmnožin nedostanou příliš rozmanité obrázky. Vracet obrázky seřazené podle *Similarity* tedy nemusí být vždy vhodné, náhodné pořadí obrázků z kategorií „téměř shodné“ a „podobné“ tedy může být lepším výsledkem. Tato úvaha může mít pozitivní dopad na efektivitu algoritmu na hledání podobných obrázků.

### 9.1 Bootstrap implementace

První implementace nahrála všechny deskriptory do databáze Elasticsearch. V Elasticsearch lze implementovat vyhledávání pomocí vlastní porovnávací funkce. Jako porovnávací funkce tedy byla zvolena funkce *Similarity*. Bohužel přes



Obrázek 9.2: Dvojice obrázků s různou kategorií podobnosti.

některé menší optimalizace algoritmu se vyhledávání nepodařilo implementovat příliš efektivně, takže fungovalo v přijatelném čase pouze pro pár tisíc deskriptorů. Naše aplikace ovšem potřebuje pracovat s více než dvaceti miliony deskriptorů. Šly by použít další optimalizace a sharding databáze na více strojů a by fungoval na větším množství dat. To by ovšem bylo příliš drahé a navíc se objevily jiné možnosti řešení.

## 9.2 Předgenerování výsledků

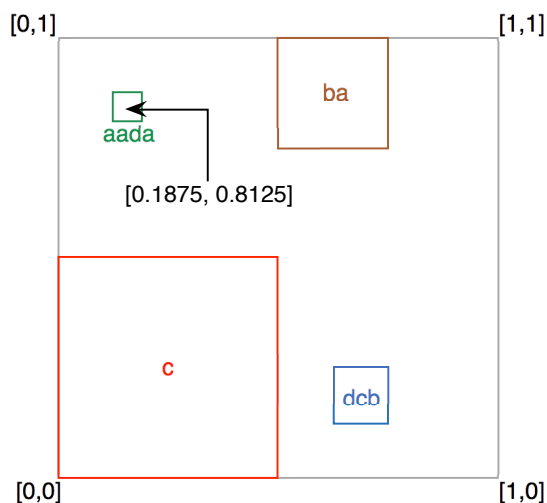
Další efektivnější variantou je ke každému obrázku vygenerovat podobné obrázky předem. Služba která by vracela podobné obrázky by pouze vrátila položku z databáze a neprováděla žádný výpočet. Tato služba by tak byla velmi rychlá a nenáročná na zdroje. Jako problém se ovšem ukázalo právě předgenerování obrázků. Přes snahu o co nejrychlejší implementaci v jazyce go s použitím více vláken a několika heuristik, by vygenerování podobných obrázků pro každý z 20 milionů obrázků trvalo na běžném počítači několik týdnů.

## 9.3 Geohash

Další pokus o implementaci vyhledávání podobných obrázků se inspiruje algoritmem Geohash[11]. Algoritmus Geohash byl vyvinut v rámci služby geohash.org a jedná se o způsob zakódování prostorových dat. Jeho hlavním využitím je efektivní vyhledávání bodů (určených zeměpisnými souřadnicemi) v oblasti (na mapě). Algoritmus Geohash využívá i Google, nebo databáze Elasticsearch.

Zjednodušení algoritmu Geohash popíšeme na jednotkové podmnožině (jednotkovém čtverci)  $\mathbb{R}^2$ . Algoritmus postupuje tak, že čtvercovou plochu rozdělí na 4 čtverce a pojmenuje je písmeny „a“, „b“, „c“ a „d“. Každý ze čtverců rekurzivně rozdělí na další čtyři čtverce, kterým dá jméno pomocí suffixů „a“, „b“, „c“ a „d“ k vlastnímu jména. Rekurzi provádíme do nějaké předem určené hloubky. Bodu v jednotkovém čtverci přiřadíme jméno podle čtverce s nejdelším jménem, který bod obsahuje. Každý bod v jednotkovém čtverci tedy bude mít jméno, které má stejnou délku jako hloubka rekurze. Jako oblast pak můžeme označit jakoukoliv množinu pojmenovaných čtverců. Bod leží v oblasti právě tehdy, když je jméno nějaké ho čtverce z množiny oblasti prefixem jména bodu. K ukládání názvů bodů pak lze použít prefixový strom.

Obrázek 9.3 ukazuje některé čtverce a jejich názvy. Pokud má algoritmus hloubku rekurze 4, má bod ležící na souřadnicích  $[0.1875, 0.8125]$  název *aada* a leží tedy ve čtvercích *a*, *aa*, *aad* a *aada*.

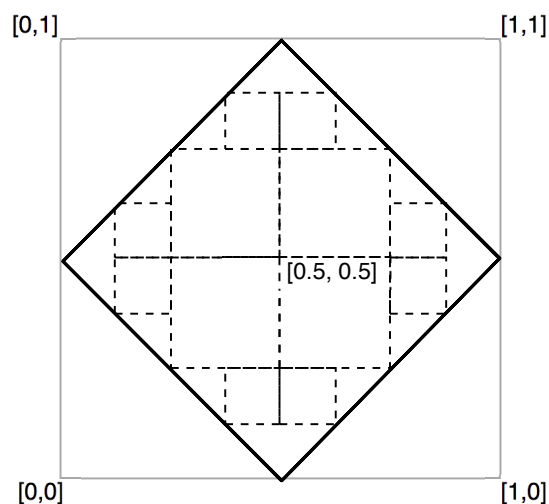


Obrázek 9.3: Ukázka čtverců algoritmu Geohash v jednotkovém čtverci.

Nyní bychom chtěli algoritmus Geohash použít pro hledání podobných obrázků v Profimedia datasetu. Každému deskriptoru bychom nejprve přiřadily název. Množina podobných obrázků by pak obsahovala obrázky, které mají *Similarity* s porovnávaným obrázkem nějak omezenou. Problémem je, že pomocí Geohash čtverců nedokážeme takovou oblast přesně definovat. Můžeme se ale pokusit co nejvíce oblast pomocí čtverců aproximovat.

Pokud bychom chtěli použít Geohash pro hledání podobných vektorů přímo, narazili bychom na několik problémů. Například oblast s ohraničenou mírou *Similarity* nejde popsat přesně pomocí čtverců. Lze libovolně přesně aproximovat, ale s každou přesnější aproximací vzrůstá počet potřebných čtverců a zhoršuje se

efektivita algoritmu. Obrázek 9.4 ukazuje jak bychom mohli použít Geohash algoritmus pro hledání podobných deskriptorů ve dvojrozměrné dimenzi. Na obrázku je popsána situace, kdy hledáme deskriptory, které mají od bodu  $[0.5, 0.5]$  vzdálenost 0.5. Tučně ohraničený čtverec vyznačuje hledanou oblast, 12 čárkovaných čtverců tvoří aproximovanou oblast.



Obrázek 9.4: Ukázka čtverců algoritmu Geohash v jednotkovém čtverci.

Ve dvou dimenzích by toto řešení fungovalo poměrně dobře. Deskriptory obrázků ovšem mají 4096 složek a počet potřebných čtverců pro stejnou míru aproximace roste s každou dimenzí exponenciálně.

# 10. Frontend

Možnosti tvorby webových aplikací se posledních několik let rapidně zvětšují. Prohlížeče implementují stále nové technologie, které rozšiřují možnosti.

## 10.1 Možnosti programování frontendu

Programátor webového frontendu si dnes může vybírat z několika paradigmat tvorby webové stránky. Standardem je dnes programovací jazyk Javascript. Spíše historicky bylo výhodné místo Javascriptu používat pro frontendový vývoj různé pluginy. Nejznámější je Adobe Flash, nebo Microsoft Silverlight. Programovat pro tyto pluginy mělo velké výhody. Například snadné přehrávání videa, zobrazení stránky v celoobrazovkovém režimu, nebo přístup k uživatelské webkamerě. V době, kdy byla Javascriptová API velmi chudá a rozdílně implementovaná mezi prohlížeči, nabízel Flash zejména tvůrcům webových her konzistenci mezi všemi platformami.

Velkou nevýhodou těchto pluginů byla jejich proprietálnost. Ostatní firmy se bály pustit cizí plugin do svých výrobků. Zlomový moment pro ústup Flashe ze slávy bylo uvedení telefonu iPhone, který podporu pro Flash nenabízel. Software-oví giganti Microsoft, Apple a Google začali místo proprietárních řešení tlačit otevřenou specifikaci, která se později nazvala HTML5. HTML5 je zastřešující termín pro spoustu technologií, které snaží webová API specifikovat a vytvářet nové.

## 10.2 JavaScriptové frameworky

Druhým hnacím motorem moderního Javascriptu jsou frameworky. Ještě před několika lety byla práce na interaktivních webech velmi náročná. Prohlížeče se zásadně lišily v implementaci práce s eventy a DOMem. Nové frameworky do jisté míry odstínily programátora od odlišných implementací Javascriptu v prohlížečích.

### 10.2.1 jQuery

Nejpopulárnějším frameworkem současnosti je jQuery. Ten nabízí jednoduché rozhraní a pro velkou většinu menších webových aplikací je zcela dostačující. Čím je ale aplikace větší, tím začíná být její vývoj s pomocí jQuery náročnější. Zkusme ukázat, jak by v jQuery vypadalo přidání CSS třídy `red` oknu s id `okno`:

```
1 $( "#okno" ).addClass( "red" );
```

Kód má několik problémů. Pokud neexistuje žádné okno s id `okno`, jQuery nevrátí žádnou chybu. Stačí malý překlep a chyba v kódu se hledá dost obtížně. jQuery nenabízí žádnou funkci typu `vratObjektPodleId`. Pokud by tyto funkce nabízel, velikost jeho kódu by se zvětšila. Pokud chce programátor použít knihovnu jQuery, musí si ji uživatel stránky celou stáhnout. Verze XX má XX bajtů.

## 10.2.2 Google Closure

Jiný přístup k vývoji frontendových aplikací přináší Google. Pro svou první webovou aplikaci Gmail vyvinul sadu nástrojů, kterou později vydal jako open source pod názvem Google Closure. Kromě Gmailu ji Google využívá v Google Vyhledávání, Google Mapách, nebo Google Dokumentech. Skládá se ze tří částí - Closure Compiler, Closure Library a Closure Templates. Closure Library je obsáhlá knihovna funkcí pro práci s DOMem, Eventy, matematickými výpočty a spoustou dalších věcí, které webový programátor může využít.

Closure Compiler je inteligentní minifikátor Javascriptového kódu. Odstraňuje funkce, které nejsou volány, přejmenovává všechny názvy funkcí a proměnných na co nejkratší řetězce a v ADVANCED módu se snaží i o pokročilejší optimalizace kódu (například kód funkcí, které jsou volány pouze jednou, je vlože inline). Nejlepších výsledků dosahuje s použitím speciálních anotací, které například vynucují typ proměnné a jsou schopny udělat z Javascriptu typovaný jazyk. Closure Compiler tyto anotace vyhodnocuje a při chybném přiřazení hodnoty vrátí chybu. To umožňuje tvořit více bezpečný javascriptový kód.

Třetí částí Google Closure jsou Closure Templates, šablonovací systém pro Javascript a Javu. Pomocí Closure Templates se snadno vytváří zanořené HTML šablony. Všechny uživatelské vstupy jsou escapované což zabraňuje sniffing útokům.

Části Templates a Compiler jdou použít odděleně v jakémkoliv Javascriptovém projektu. Používat Closure Library bez Compileru nedává příliš smysl, uživatel by při návštěvě webu musel stahovat ohromné množství zbytečných dat.

Tato práce na frontendu používá všechny části knihovny Google Closure. Díky tomu si uživatel při první návštěvě webu musí stáhnout pouze jediný soubor, který má pouze XX Kb.

## 10.3 Stylování uživatelského rozhraní a CSS

Specifikace HTML5 rozšiřuje i možnosti vizualizace pomocí CSS stylů. Nejviditelnějšími novinkami je podpora kulatých rohů, stínování, nebo barevných přechodů. Webový programátor nyní může ke stránce načíst i vlastní font. Všechny tyto možnosti velmi rozšířily možnosti webovým grafikům.

Dalším trendem, který v CSS světě probíhá, je používání CSS preprocesorů.

# 11. Instalace a zprovoznění

Celé anotační rozhraní je webová aplikace napsaná v jazyce Ruby a frameworku Ruby on Rails. Je k dispozici pod svobodnou licencí MIT. K jejímu spuštění potřebujete ruby verze alespoň 2.0 (nižší verze nejsou otestované), javu a ke stažení zdrojového kódu git. Program jde spustit na Linuxu a Macu.

## 11.1 Instalace

Zdrojový kód je volně dostupný na webu GitHubu<sup>1</sup>. Stáhnou tedy lze příkazem

```
1 git clone https://github.com/hypertornado/diplomka
```

Tento příkaz vytvoří adresář diplomka. Závislosti aplikace nainstalujete pomocí bundleru:

```
1 bundle install
```

Instalace může vyžadovat přístup administrátora. Dále je potřeba stáhnout knihovnu elasticsearch<sup>2</sup> do adresáře bin/elasticsearch. Stačí verze 1.0 a vyšší. Ve verzi 1.2.1 jsme objevili menší chybu<sup>3</sup>, která je způsobena chybou v Javě a jde obejít nastavením delšího hostname počítače.

Nyní je možné celý projekt spustit. Nejprve se spustí elasticsearch databáze pomocí příkazu `rake es:start`, poté je možné spustit samotnou aplikaci příkazem `rails server`. Po spuštění severu je uživatelské rozhraní dostupné ve webovém prohlížeči na adrese `http://localhost:3000`. Po načtení stránky se zobrazí uživatelské rozhraní, ale veškeré AJAXové dotazy skončí chybou. V databázi nejsou importována data.

## 11.2 Práce s metadaty k obrázkům

Metadata k obrázkům a obrázky samotné jsou poskytovány firmou Profimedia a nejsou volně dostupné. Ke zprovoznění aplikace je nutné vložit CSV soubor `keyword-cleaned-phrase-export.csv` do adresáře data.

## 11.3 Překlad metadat

Soubor obsahuje metadata k obrázkům v angličtině. Jedním z úkolů této práce je poskytnout doporučení obrázků i v jiných jazycích, primárně v českém jazyce. Bylo tedy nutné metadata přeložit. Pokoušeli jsme se použít volný nástroj na překlad Moses. Ve verzi 2.1<sup>4</sup> nabízí volně dostupné modely pro překlad z češtiny do angličtiny. I na SSD disku trvá několik hodin, než se překladový model načte

<sup>1</sup><https://github.com/hypertornado/diplomka>

<sup>2</sup><http://www.elasticsearch.org/downloads/1-0-3/>

<sup>3</sup><https://github.com/elasticsearch/elasticsearch/issues/6611>

<sup>4</sup><http://www.statmt.org/moses/RELEASE-2.1/models/en-cs/model/>



do paměti. Překlad jednoho segmentu s tímto modelem byl poměrně pomalý (překlad metadat k jednomu obrázku trval zhruba 3 sekundy) a také dosti nepřesný. Například řádek

```
1 "0000000003","little baby smiling","", "child children
    baby babies infants kids childhood single faces
    body naked naked facial expressions smile smiling
    viewing watching laying fun amusing amusement
    amused amuse dallying frolics playing wantoning
    open"^M
```

byl do češtiny přeložen takto:

```
1 "0000000003","little|UNK|UNK|UNK dítě
    smiling","", "child|UNK|UNK|UNK děti , dětské děti
    kojence děti dětství jednotného čelí orgán nahé
    naked|UNK|UNK|UNK pořídili vyjádření usmívat usmívá
    odůvodnění , která zábavné sledovat zábavné
    zábavných i pobavena tím amuse|UNK|UNK|UNK
    dallying|UNK|UNK|UNK frolics|UNK|UNK|UNK hrát
    wantoning|UNK|UNK|UNK open"^M|UNK|UNK|UNK
```

Je vidět poměrně velké množství nepřeložených slov (koncovky |UNK) a překlad je relativně nepřesný. Je pravděpodobné, že by s lepšími daty šel natrénovat lepší překladový a jazykový model. Strojový překlad není hlavním tématem této práce, takže bylo jednodušší komerční automatický překlad Google Translate, který přeloží ukázková metadata takto:

```
1 "0000000003", "malé dítě s úsměvem", "", "dítě děti
    dítě děti kojenci děti dětství jednotlivé plochy
    těla nahá naked výrazy obličeje, úsměvu, usmívavý
    sledování sledování kterým zábava zábavné zábavní
    pobavený pobavit laškoval frolics hrát wantoning
    otevřený" ^ M
```

I z této ukázky je zřejmé, že Google poskytuje kvalitnější překlad, než anglicko-český překladový model v releasu Moses. Google poskytuje překlad zdarma přes webové rozhraní. Pokud se ale do překladového formuláře nahraje nespecifikované větší množství dat, přestane překlad fungovat. Google pro překlad poskytuje placené API. Platí se XX amerických dolarů za přeložené slovo. Korpus Profimedia obsahuje zhruba (wc vystup 20119222 347129204 4811998848), takže by celý překlad stál XX dolarů. Jelikož se slova v textu opakují, lze použít překlad slovo po slovu. Ještě lepší by bylo překládat přímo klíčové fráze. Bohužel korpus Profimedia jednotlivé fráze neodděluje. Je ale možné použít učící algoritmus a klíčové fráze detekovat.

### 11.3.1 Export slov a frází

Nejprve příkazem `rake data:export_profimedia_words_for_translation` vy-exportujeme do souboru `data/word_list.txt` seznam všech slov použitých v metadatech k obrázkům. Tento příkaz běží několik hodin i na moderním počítači s SSD diskem. Z Profimedia dat získáme seznam 352862 slov. Tento soubor je nutné přeložit z angličtiny do dalších podporovaných jazyků, v našem případě češtiny.

## 11.4 Jazykový korpus

Tato práce potřebuje jazykové korpusy pro podporované jazyky ze dvou důvodů. První je potřeba pro určení relativní jazykové frekvence slov v algoritmu TF-IDF. Zadruhé je potřeba jazykový korpus pro rozpoznávání jazyků. Přirozeným jazykovým korpusem by mohla být samotná Profimedia data, ale pro oba účely jsou tato data nepoužitelná. Klíčová slova a názvy obrázků jsou odlišnými druhy textů, než je průměrný článek. Je tedy potřeba jazyková data získat jinde.

Wikipedia se jako korpus velmi hodí. Textový obsah je pod licencí Creative Commons a jeho struktura je velmi podobná běžnému publicistickému článku. Data se dají získat stažením přímo ze serverů wikipedie (pro češtinu zde), nebo pomocí sdílených torrentů. Práce wiki dumpy všech podporovaných jazyků v adresáři data pod názvem typu `wiki_dump_jazyk.xml`. Pro práci s daty z wikipedie je potřeba nejprve převést XML data do textového formátu. K tomu slouží python skript `lib/WikiExtractor.py` od Wikipedie. Pro převod anglických dat lze použít příkaz:

```
1 rake wiki:extract_words_from_wiki en
```

Stejný příkaz je potřeba spustit i pro ostatní podporované jazyky. Příkaz vytvoří v `data/wiki_en` adresářovou strukturu. Není potřeba převést všechna data, pouze tolik, abychom dostali reprezentativní korpus. Pro angličtinu stačí převést zhruba 10000 článků, pro podobné množství českých dat je potřeba převést zhruba 20000 článků (údaj o exportovaných článcích je průběžně vypisován na konzoli).

Nyní je možné vytvořit seznam slov v korpusu s frekvencemi. Ve skutečnosti nás zajímají pouze stemy slov, ne jednotlivé tvary. Příkaz

```
1 rake wiki:frequency_list_from_wiki
```

vytvoří seznam slov a jejich frekvencí s cestou `data/wiki_freq_list_count.en`. Ve skutečnosti nás nezajímají všechna slova z korpusu wikipedie, ale pouze ta, která se vyskytují mezi klíčovými slovy v Profimedia korpusu. Příkazem

```
1 rake data:create_tf_df_list
```

se z korpusu profimedia dat exportují stemy všech slov v profimedia korpusu spolu s hodnotami udávajícími frekvenci termínu v celém korpusu (TF) a frekvenci toho, v kolika popiskách obrázků se slovo objevilo (DF). Pro angličtinu jsou tato data uložena do souboru `data/tf_df_list_en.txt`. Nyní můžeme spárovat

informaci o stemech z wikipedie a nahrát je do databáze. Párování se provede příkazem:

```
1 rake data:pair_profimedia_and_wiki_data
```

který vytvoří soubor `data/paired_wiki_and_profimedia.txt` se statistikami všech stemů z Profimedia korpusu.

## 11.5 Příprava dat pro detekci jazyků

Automatická detekce jazyka probíhá ve frontendové části. Ke svému chodu ale potřebuje data z korpusu, konkrétně seznam nejčastějších trigramů pro každý podporovaný jazyk. Příkaz

```
1 rake trigrams:extract_most_frequent_trigrams
```

vytvoří pro každý jazyk seznam padesáti nejčastějších trigramů. Seznam pro angličtinu je uložen v souboru `data/most_frequent_trigrams_en.txt`. Příkaz

```
1 rake trigrams:trigrams_to_javascript_classes
```

pracuje se seznamy nejčastějších trigramů ze kterých vytvoří javascriptovou třídu `oo.diplomka.Languages.Data` v notaci Google Closure. Ta je uložena v souboru `public/js/js/diplomka/languages/data.js`.

## 11.6 Import dat do databáze

Veškerá data jsou importována do databáze elasticsearch. Aplikace očekává Elasticsearch připojený na portu 9200. Samotná data se importují příkazem:

```
1 rake es:import_image_metadata
```

Tento vytvoří v elasticsearch index `diplomka`. Poté nastaví explicitní mapování v celém indexu. Pokud se mapování explicitně nenastaví, elasticsearch se snaží data namapovat podle jednoduchých heuristik. Například text rozdělí na slova, která pak převede na malá písmena a upraví defaultním stemmerem. Aplikace se však o stemming a převedení na malá písmena stará sama, mapování tedy říká, aby nahraný text oddělil pouze pomocí mezer na slova a dále nezpracovával. Vyhledávací dotaz je pak rozdělen také pouze pomocí mezer. Nastavené mapování lze ověřit pomocí API `elasticsearch` na adrese `http://localhost:9200/diplomka/_mapping/`.

Po vytvoření indexu může začít samotný import dat. Data z každého řádku Profimedia dat je převeden do formátu JSON. Data obsahují položky `locator`, `title`

## 11.7 Import metadat obrázků

## 12. Anotační rozhraní

V rámci této práce bylo implementováno anotační rozhraní pro vyhodnocování algoritmů, které přiřazují vhodné obrázky k textům. Anotační rozhraní je velmi univerzální. Anotátor má ve webové aplikaci v levém sloupci novinový text a v pravém sloupci galerii obrázků. Jeho úkolem je označit obrázky, které se k danému textu hodí a obrázky které se k textu nehodí. Má také možnost nechat obrázek neoznačený, pokud by se nemohl rozhodnout ani pro jednu variantu.

### 12.1 Instalace rozhraní

Celé rozhraní je aplikace napsaná v jazyce Ruby a frameworku Ruby on Rails. Aplikace je volně šiřitelná pod licencí MIT. Pro zprovoznění anotační aplikace je potřeba UNIXový systém (Linux, Mac). Zdrojový kód aplikace je uložen na severu GitHub<sup>1</sup> a nejlépe jde stáhnout pomocí git. Pro běh serveru je potřeba verze ruby 2.0 a vyšší. Celá aplikace se zprovozní následujícím pořadím BASH příkazů:

```
1 git clone https://github.com/hypertornado/cemi_anotace
2 cd cemi_anotace
3 bundle install #nainstaluje vsechny ruby zavislosti
4 rake db:migrate #vytvori sqlite databazi s tabulkami
5 rails server #spusti anotacni server na portu :3000
```

### 12.2 Přidání uživatelů

Po spuštění serveru je možné přidat anotátory v administračním rozhraní. Přístup je zaheslován HTTP autentifikací. Defaultní uživatelské jméno je cfo a heslo cfo85. Administrátorské přístupové údaje lze změnit v souboru `ROOT_APLIKACE/app/controllers`. Uživatelé mají pouze dvě datové položky, uživatelské jméno (Name) a heslo (Password). Uživatele jde přidávat, mazat a upravovat. Nepředpokládá se, že by anotovaná data byla vysoce citlivá, heslo je proto v databázi uloženo v plaintextu.

### 12.3 Import anotačních dat

Data pro anotaci lze nahrát pomocí příkazu

```
1 rake data:import
```

Příkaz očeká existenci souboru `ROOT_APLIKACE/public/annotation_inputs.csv`. Ten musí mít speciální formát, kdy je každý řádek rozdělen mezerami na šest sloupců s následujícími položkami:

---

<sup>1</sup>[https://github.com/hypertornado/cemi\\_anotace](https://github.com/hypertornado/cemi_anotace)

## INDEX

unikátní číslo jedné anotace

## LABEL

interní popis pokusu

## PRIORITY

priorita, celé číslo  $\geq 0$ . Určuje prioritu s jakou se má anotace přiřadit. Čím vyšší číslo, tím vyšší priorita.

## PREFER\_USER

uživatelské jméno preferovaného anotátora. Pokud není žádný anotátor preferován, použije se pomlčka

## TEXT\_FILE

cesta k textovému souboru s referenčním článkem

## IMAGE\_FILES

seznam cest k obrázkům. Cesty nemohou obsahovat mezery a jsou oddělené středníkem.

Ukázka importovaných dat:

```
1 1 basics 0 - ./text/aha/aha-00263.txt.gz  
   img/1.jpg;img/2.jpg;img/3.jpg  
2 2 basics 1 - ./text/aha/aha-00006.txt.gz  
   img/1.jpg;img/2.jpg  
3 3 basics 0 - ./text/aha/aha-00009.txt.gz  img/2.jpg
```

## 12.4 Export anotačních dat

Hotové anotace lze exportovat příkazem

```
1 rake data:export
```

Tento příkaz vypíše na konzolu řádky, které mají tabulátorem oddělené položky:

### INDEX

ID anotace. Stejně jako u importovaných dat.

### USER

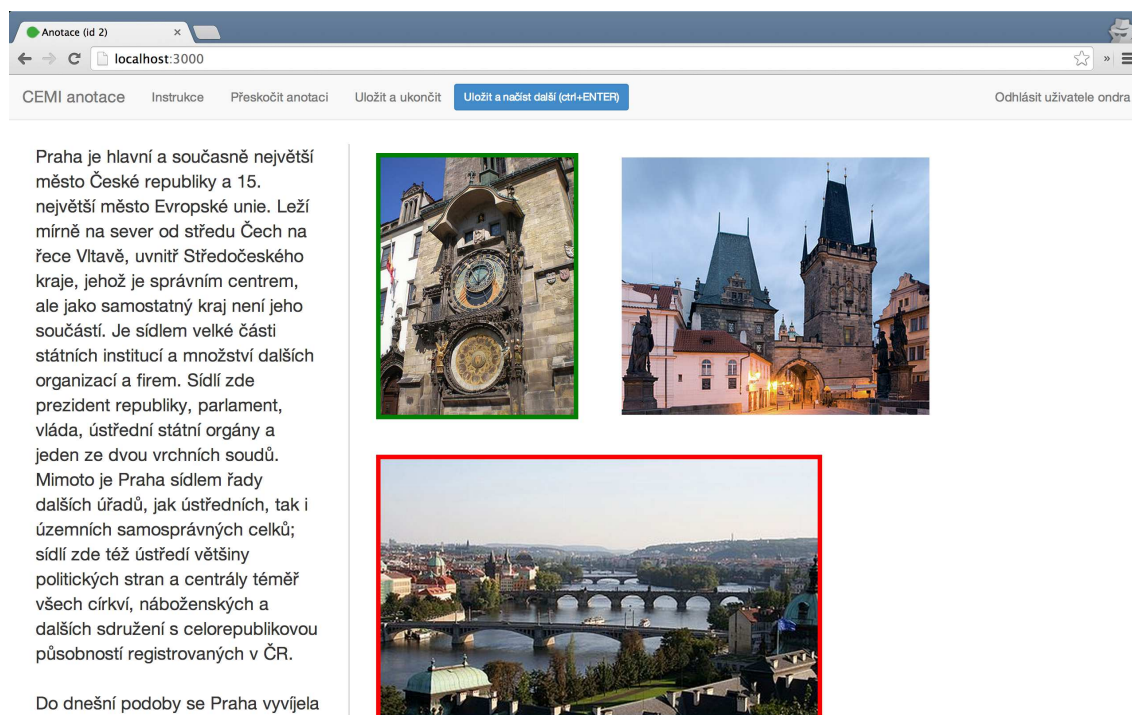
Jméno anotátora, který anotaci vytvořil.

### TIME

Čas uložení hotové anotace ve formátu UNIX timestamp.

### SKIPPED

Pokud uživatel anotaci přeskočil, je hodnota True, jinak False.



Obrázek 12.1: Anotační rozhraní. Vhodné obrázky jsou označené zeleným rámečkem, nevhodné červeným.

## APPROPRIATE

Seznam obrázků které anotátor označil jako vhodné k textu ve formátu relativních cest oddělených středníkem.

## NOT APPROPRIATE

Seznam obrázků které anotátor označil jako nevhodné k textu ve formátu relativních cest oddělených středníkem.

## 12.5 Import obrázků a textů

Anotační texty a obrázky musí být nahrány do adresáře `ROOT_APLIKACE/public` tak, aby jejich cesty odpovídali cestám v souboru `ROOT_APLIKACE/public/annotation_inputs.csv`. Pokud tedy importovaný soubor obsahuje cestu k obrázku `img/1.jpg`, musí být nahrán odpovídající soubor do `ROOT_APLIKACE/public/img/1.jpg`.

Obrázky musí být ve formátu, který podporují webové prohlížeče, tedy hlavně JPEG a PNG. Texty musí být uloženy v textových souborech s kódováním utf-8 a komprimovaný pomocí gzip<sup>2</sup>.

## 12.6 Anotační proces

Úkolem anotátora je přiřadit vhodné a nevhodné obrázky. Po přihlášení do anotačního rozhraní vidí v levé části test a v pravé obrázky. Levým tlačítkem myši může označit obrázky, které odpovídají textu, pravým tlačítkem myši označí

<sup>2</sup><http://www.gzip.org/>

obrázky, které textu neodpovídají. Pokud si anotátor není jistý, nechá obrázek neoznačený. Uživatel může použít klávesovou zkratku `ctrl+ENTER` k uložení a načtení další anotace. Může také anotaci přeskočit (pak se označí jako přeskočená a nepřihadí se jinému anotátorovi), nebo uložit a ukončit.

## 13. Testování výsledků

Výsledný program se skládá z několika komponent. Funkčnost každé komponenty jde testovat. Některé části jdou otestovat automaticky, například rozpoznávání jazyka. Na některé části je potřeba testování s uživateli.

Klíčovou částí celého projektu je algoritmus, který k textu přiřadí ilustrační obrázek. Úspěšnost tohoto algoritmu byla otestována na uživateli. Rozhraní může fungovat pro více jazyků. Testování však proběhlo pro vstupní texty v češtině. Jedním z důvodů je snadnější získání anotátorů v českém jazyce. Hlavním důvodem je ovšem pravděpodobná vyšší nepřesnost algoritmu v češtině. Dataset Profimedia má data v anglickém jazyce a pro české použití musel být přeložen. Pokud se tedy ukáže, že algoritmus funguje dobře pro češtinu, měl by pro angličtinu fungovat ještě lépe. Pro testování bylo využito testovací rozhraní popsané v Kapitole 12.

### 13.1 Vyloučení narušitele (o\_test1)

První metodou testování bylo „vyloučení nepřítele“ (anglicky „intruder detection“). Tato metoda se používá k evaluaci automaticky detekovaných shluků slov a je popsána například v [3]. V naší variantě se evaluují obrázky přiřazené k textu. Nejprve se k danému textu najdou pomocí testovaného algoritmu 4 nejvíce odpovídající obrázky. K nim se přidá jeden náhodně vybraný obrázek z datasetu a poté se náhodně zamíchá pořadím těchto obrázků. Anotátor vidí v anotačním rozhraní text a 5 obrázků. Jeho úkolem je označit obrázek, který danému textu, podle jeho názoru, odpovídá nejméně. Pokud algoritmus přiřazuje obrázky textu funguje správně, měl by být uživatel schopen označit obrázek, který byl do sady vybrán náhodně a rozlišit ho od obrázků, který vybral algoritmus přiřazující obrázky.

Texty pro testování algoritmu pochází z online článků českých webových serverů. Jedná se o texty článků, které obsahují ilustrační obrázky z datasetu Profimedia. Takto omezená množina novinových textů je pro náš účel velmi výhodná. Pro různé druhy článků dataset Profimedia neobsahuje žádné vhodné ilustrační obrázky. Jedná se například o politické zpravodajství, pro které v datasetu aktuální fotky událostí, nebo například archivní fotky osobností. Oproti tomu u článků, které již nějaký ilustrační obrázek z Profimedia obsahují, jsou možnosti vhodného obrázku daleko vyšší. Jedná se většinou o různé hobby a společenské články.

Pro naše testování jsme využili články ze serveru iDnes, který obsahoval nejvíce článků s ilustračními obrázky z Profimedia. Konkrétně jich máme k dispozici 4223.

Z těchto 4223 článků jsme náhodně vybrali pro testování 60 článků. Ke každému z článků byly získány algoritmem 4 doporučené ilustrační obrázky a byl přidán jeden obrázek náhodný. Každá tato testovací sada byla otestována dvěma anotátory. Bylo k dispozici 5 anotátorů. Dohromady to znamenalo pro každého anotátora 24 a dohromady 120 anotací. Anotace byly rozděleny tak, aby každá dvojice anotátorů měla právě 6 stejných testovacích sad.

Během anotace se zjistilo, že u dvou testovacích sad se jeden z obrázků nena-



	<b>o_test1</b>	<b>o_test2</b>	<b>celkem</b>
<b>anotátorů</b>	5	5	7
<b>textů</b>	58	120	178
<b>anotací</b>	116	240	356
<b>pozitivní shoda</b>	45/58 = 78%	76/120 = 63%	121/178 = 67%
<b>negativní shoda</b>	3/58 = 5%	15/120 = 13%	18/178 = 10%
<b>neshoda</b>	10/58 = 17%	29/120 = 24%	39/178 = 22%

Obrázek 13.1: Přehled výsledků uživatelského testování.

čítá. Tyto sady byly z testování vyřazeny a anotovaných testovacích sad je tedy pouze 58. Výsledky testování jsou shrnuty v tabulce 13.1 ve sloupci „o\_test1“. Po-  
brobné výsledky testování, včetně Cohenovy kappy pro všechny dvojice anotátorů je v příloze 14.

Výsledky ukazují, že pro 78 % testovacích sad měli anotátoři pozitivní shodu. Pozitivní shoda znamená, že se oba anotátoři shodli na stejném obrázku, který je pro vstupní text nejméně vhodný. Tento obrázek byl zároveň vybrán do testovací sady náhodně. Čím je toto procento vyšší, tím lépe náš algoritmus pracuje, protože uživatelé jsou schopni detekovat náhodný obrázek. Pro 5 % testovacích sad máme negativní shodu. Oba anotátoři vybrali obrázek, který nebyl přiřazen náhodně. Znamená to tedy, že tento obrázek nebyli schopni detekovat a algoritmus tedy pro vstupní text nepracuje dobře (pokud vyloučíme možnost, že náhodně přiřazený obrázek je k textu relevantní). U poslední skupiny testovacích dat – 17 % – byl pouze jeden z anotátorů schopen vybrat náhodně vybraný obrázek.

Získaná data ukazují, že algoritmus pracuje poměrně dobře. Pokud by algoritmus přiřazoval automaticky obrázky k textům i v praxi, čtenáři by mohli by poznali rozdíl od algoritmu, který přiřazuje obrázky nahodile.

Během testování se objevil jeden zásadní problém s testovací metodou. V datasetu Profimedia jsou i obrázky, které jsou si velmi vizuálně podobné, například fotky stejné osoby z různých úhlů. Tyto fotky mají často i stejné textové popisky. Mějme tedy 4 obrázky, které jsou si vizuálně velmi podobné a mají stejné textové popisky. Pokud algoritmus označí jako nejvhodnější obrázek k textu jeden z těchto obrázků, budou i na dalších třech doporučených pozicích vizuálně podobné obrázky (pokud tedy nemáme jinou množinu obrázků, která má stejné textové popisky, ale je vizuálně odlišná). Pokud se takový text objeví v naší testovací metodě, uvidí anotátor 4 velmi podobné obrázky a k nim jeden náhodně vybraný. Nejméně vhodný obrázek pak snadno označí, aniž by vůbec četl anotační text. Ukázalo se, že v našem testování k takovému problému opravdu došlo – anotátorovi se zobrazily obrázky, z nichž ani jeden nebyl vhodným obrázkem k danému textu, přesto anotátor snadno označil nejméně vhodný obrázek. Tento problém zkrusluje dosažené výsledky testování touto metodou.

## 13.2 Detekce správného obrázku (o\_test2)

Abychom předešli problémům s metodou popsaným v předchozí sekci, provedli jsme nové testování. Úkolem uživatelů bylo nyní vybrat obrázek, který se ke vstupnímu textu hodí nejvíce. Množinu obrázků nyní tvořil jeden výstup z algorit-

mu a 4 náhodně vybrané obrázky. Zvětšili jsme i testovací sadu. Bylo testováno 120 textů náhodně vybraných z iDnes datasetu. Každý z textů byl anotován dvěma anotátory. Na každém z pěti anotátorů tedy bylo 48 anotací. Výsledky testování jsou shrnuty v tabulce 13.1 ve sloupci „o\_test2“. Poborné výsledky testování, včetně Cohenovy kappy pro všechny dvojice anotátorů je v příloze 14.

Výsledky jsou oproti předchozí metodě o něco horší. Pro 63 % testovacích sad měli oba anotátoři pozitivní shodu, pro 13 % měli anotátoři negativní shodu a pro 24 % se anotátoři neshodli.

### 13.3 Shrnutí

Uživatelské testování ukazuje, že je algoritmus schopen v poměrně velkém procentu případů přiřadit automaticky „dostatečně vhodný“ ilustrační obrázek. To, že je ilustrační obrázek dostatečně vhodný ovšem neznamená, že je nejvhodnější. Výběr nejvhodnějšího ilustračního obrázku je ovšem velmi individuální a nedá se obecně měřit.

Testování ukázalo také na slabiny algoritmu. Možná největší slabinou je fáze překladu, díky níž popisky k některým obrázkům nesprávné. Tímto problémem netrpí algoritmus pro anglické vyhledávání, který pracuje s nepřeloženými popisky. Ten však nebyl uživatelsky testován.

Jedním z problémů testování je právě omezená doména testovaných textů. Můžeme říci, že algoritmus funguje dobře na nějaké doméně textů, ale pro obecné zpravodajské články může algoritmus fungovat velmi špatně. Toto ovšem není problém testovaného algoritmu, ale dat na kterých pracuje. Rozšíření domény obrázků v korpusu Profimedia by rozšířilo i doménu textů pro které algoritmus funguje dobře.

## 14. Závěr

# Seznam použité literatury

- [1] BOTOŘEK, Jan. *Tvorba nástroje pro zpracování textových popisů multimedialních dat* [online]. 2012 [cit. 2014-07-17]. Bakalářská práce. Masarykova univerzita, Fakulta informatiky. Vedoucí práce Petra Budíková. Dostupné z: [http://is.muni.cz/th/359815/fi\\_b/](http://is.muni.cz/th/359815/fi_b/).
- [2] BUDÍKOVÁ, Petra, Michal BATKO a Pavel ZEZULA. *Evaluation Platform for Content-based Image Retrieval Systems*. In *International Conference on Theory and Practice of Digital Libraries 2011, LNCS 6966*. Berlin: Springer, 2011. s. 130-142, 12 s. ISBN 978-3-642-24468-1.
- [3] CHANG, Jonathan, et al. *Reading tea leaves: How humans interpret topic models*. In: *Advances in neural information processing systems*. 2009. s. 288-296.
- [4] GOOGLE INC. *Překladač Google* [online]. 2014 [cit. 2014-07-25]. Dostupné z: <http://translate.google.cz>.
- [5] JIA, Yangqing. *Caffe: An open source convolutional architecture for fast feature embedding*. 2013 Dostupné z: <http://caffe.berkeleyvision.org>.
- [6] KOEHN, Philipp. *Statistical machine translation*. Cambridge University Press, 2009.
- [7] KOEHN, Philipp, et al. *Moses: Open source toolkit for statistical machine translation*. In: *Proceedings of the 45th Annual Meeting of the ACL on Interactive Poster and Demonstration Sessions*. Association for Computational Linguistics, 2007. s. 177-180.
- [8] LOTT, Brian. *Survey of Keyword Extraction Techniques*. UNM Education, 2012.
- [9] MILLER, George A. *WordNet: A Lexical Database for English*. Communications of the ACM, 1995, roč. 38, č. 11, s. 39-41.
- [10] PORTER, Martin F. *An algorithm for suffix stripping*. Program: electronic library and information systems. MCB UP Ltd, 1980, roč. 14, č. 3, s. 130-137.
- [11] WIKIPEDIA *Geohash* — *Wikipedia, The Free Encyclopedia* [online]. 2014 [cit. 2014-07-24] Dostupné z: <http://en.wikipedia.org/w/index.php?title=Geohash&oldid=609521802>.

# Seznam tabulek

# Seznam použitých zkratek

# Příloha 1

```
1 Výsledek testu o_test1: 103 / 116 = 88.79 %
2 Výsledek testu o_test2: 196 / 240 = 81.66 %
3 Celkové výsledky dohromady: 299 / 356 = 83.98 %
4
5
6 havel
7   o_test1: 20 / 22 = 90.90 %
8   o_test2: 36 / 48 = 75.0 %
9   celkem: 56 / 70 = 80.0 %
10  odchazel_o
11    o_test1
12      3    0
13      1    1
14      Cohen's kappa = 0.5454545454545455
15  odchazel_v
16    o_test2
17      10   2
18      0    0
19      Cohen's kappa = 0.0
20  paroubkova
21    o_test1
22      6    0
23      0    0
24      Cohen's kappa = NaN
25    o_test2
26      6    1
27      3    2
28      Cohen's kappa = 0.2727272727272727
29  pavlovic
30    o_test1
31      5    0
32      0    0
33      Cohen's kappa = NaN
34  rakosova
35    o_test1
36      6    0
37      0    0
38      Cohen's kappa = NaN
39    o_test2
40      9    0
41      1    2
42      Cohen's kappa = 0.7499999999999999
43  semerad
44    o_test2
```

```

45      6    2
46      0    4
47      Cohen's kappa = 0.6666666666666667
48
49 odchazel_o
50 o_test1: 21 / 23 = 91.30 %
51 celkem:  21 / 23 = 91.30 %
52   havel
53     o_test1
54       3    1
55       0    1
56       Cohen's kappa = 0.5454545454545455
57   paroubkova
58     o_test1
59       4    1
60       0    1
61       Cohen's kappa = 0.5714285714285715
62   pavlovic
63     o_test1
64       6    0
65       0    0
66       Cohen's kappa = NaN
67   rakosova
68     o_test1
69       6    0
70       0    0
71       Cohen's kappa = NaN
72
73 odchazel_v
74 o_test2: 42 / 48 = 87.5 %
75 celkem:  42 / 48 = 87.5 %
76   havel
77     o_test2
78       10   0
79       2    0
80       Cohen's kappa = 0.0
81   paroubkova
82     o_test2
83       10   1
84       0    1
85       Cohen's kappa = 0.625
86   rakosova
87     o_test2
88       9    1
89       0    2
90       Cohen's kappa = 0.7499999999999999
91   semerad
92     o_test2

```



```

93      11  0
94      0  1
95      Cohen's kappa = 1.0
96
97 paroubkova
98 o_test1: 20 / 24 = 83.33 %
99 o_test2: 43 / 48 = 89.58 %
100 celkem: 63 / 72 = 87.5 %
101 havel
102 o_test1
103 6 0
104 0 0
105 Cohen's kappa = NaN
106 o_test2
107 6 3
108 1 2
109 Cohen's kappa = 0.2727272727272727
110 odchazel_o
111 o_test1
112 4 0
113 1 1
114 Cohen's kappa = 0.5714285714285715
115 odchazel_v
116 o_test2
117 10 0
118 1 1
119 Cohen's kappa = 0.625
120 pavlovic
121 o_test1
122 4 0
123 2 0
124 Cohen's kappa = 0.0
125 rakosova
126 o_test1
127 5 1
128 0 0
129 Cohen's kappa = 0.0
130 o_test2
131 11 1
132 0 0
133 Cohen's kappa = 0.0
134 semerad
135 o_test2
136 10 2
137 0 0
138 Cohen's kappa = 0.0
139
140 pavlovic

```

```

141 o_test1: 21 / 23 = 91.30 %
142 celkem: 21 / 23 = 91.30 %
143 havel
144 o_test1
145 5 0
146 0 0
147 Cohen's kappa = NaN
148 odchazel_o
149 o_test1
150 6 0
151 0 0
152 Cohen's kappa = NaN
153 paroubkova
154 o_test1
155 4 2
156 0 0
157 Cohen's kappa = 0.0
158 rakosova
159 o_test1
160 3 1
161 1 1
162 Cohen's kappa = 0.24999999999999986
163
164 rakosova
165 o_test1: 21 / 24 = 87.5 %
166 o_test2: 39 / 48 = 81.25 %
167 celkem: 60 / 72 = 83.33 %
168 havel
169 o_test1
170 6 0
171 0 0
172 Cohen's kappa = NaN
173 o_test2
174 9 1
175 0 2
176 Cohen's kappa = 0.7499999999999999
177 odchazel_o
178 o_test1
179 6 0
180 0 0
181 Cohen's kappa = NaN
182 odchazel_v
183 o_test2
184 9 0
185 1 2
186 Cohen's kappa = 0.7499999999999999
187 paroubkova
188 o_test1

```

```

189      5    0
190      1    0
191      Cohen's kappa = 0.0
192      o_test2
193      11    0
194      1    0
195      Cohen's kappa = 0.0
196      pavlovic
197      o_test1
198      3     1
199      1     1
200      Cohen's kappa = 0.249999999999999986
201      semerad
202      o_test2
203      9     0
204      0     3
205      Cohen's kappa = 1.0
206
207      semerad
208      o_test2: 36 / 48 = 75.0  %
209      celkem:  36 / 48 = 75.0  %
210      havel
211      o_test2
212      6     0
213      2     4
214      Cohen's kappa = 0.66666666666666667
215      odchazel_v
216      o_test2
217      11    0
218      0     1
219      Cohen's kappa = 1.0
220      paroubkova
221      o_test2
222      10    0
223      2     0
224      Cohen's kappa = 0.0
225      rakosova
226      o_test2
227      9     0
228      0     3
229      Cohen's kappa = 1.0
230
231
232      o_test1
233      nesprávné označení (text, uživatel)
234      00207 o_paroubkova
235      01337 o_pavlovic
236      01595 o_paroubkova

```

237	01921	o_rakosova
238	03758	o_odchazel_ondrej
239	03758	o_paroubkova
240	03907	o_paroubkova
241	04107	o_rakosova
242	07604	o_havel
243	07604	o_odchazel_ondrej
244	09093	o_pavlovic
245	09093	o_rakosova
246	09551	o_havel
247	celkem 116 anotací 58 textů	
248	každý text oanotován dvěma anotátory	
249	oba správně : 45 / 58 = 77.58 %	
250	oba špatně : 3 / 58 = 5.172 %	
251	jeden špatně: 10 / 58 = 17.24 %	
252		
253	o_test2	
254	nesprávné označení (text, uživatel)	
255	00124	o_semerad
256	00405	o_havel
257	00405	o_paroubkova
258	00623	o_paroubkova
259	01907	o_rakosova
260	01911	o_odchazel_vojtech
261	01911	o_paroubkova
262	02505	o_odchazel_vojtech
263	02505	o_rakosova
264	03610	o_havel
265	03610	o_semerad
266	04211	o_odchazel_vojtech
267	04282	o_havel
268	04282	o_semerad
269	04368	o_odchazel_vojtech
270	04368	o_semerad
271	04480	o_semerad
272	04530	o_havel
273	05102	o_rakosova
274	05102	o_semerad
275	05692	o_semerad
276	06379	o_paroubkova
277	06777	o_rakosova
278	07326	o_rakosova
279	07326	o_semerad
280	07587	o_odchazel_vojtech
281	07587	o_rakosova
282	08015	o_havel
283	08015	o_semerad
284	08621	o_havel

285	12931	o_odchazel_vojtech
286	13674	o_havel
287	13674	o_semerad
288	13803	o_havel
289	13803	o_paroubkova
290	14257	o_havel
291	14257	o_rakosova
292	14329	o_rakosova
293	14329	o_semerad
294	14353	o_havel
295	14353	o_rakosova
296	14988	o_semerad
297	15726	o_havel
298	15770	o_havel
299	celkem 240 anotací 120 textů	
300	každý text oanotován dvěma anotátory	
301	oba správně : 76 / 120 = 63.33 %	
302	oba špatně : 15 / 120 = 12.5 %	
303	jeden špatně: 29 / 120 = 24.16 %	