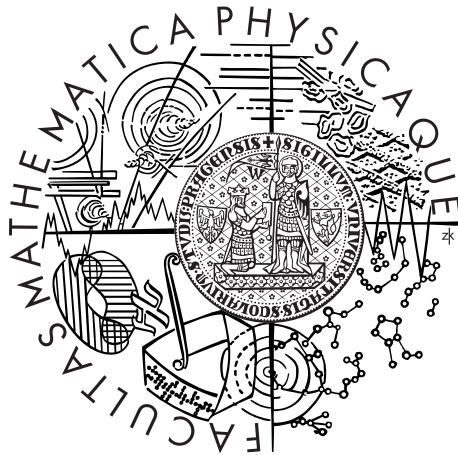


Univerzita Karlova v Praze
Matematicko-fyzikální fakulta

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, RNDr. Pavlovi Pecinovi, Ph.D. za pomoc a cenné připomínky v průběhu celého vývoje projektu. Také bych rád poděkoval své rodině, přítelkyni a kamarádům za pomoc a 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: Cílem této práce je implementace webové aplikace určené k doporučování ilustračních obrázků. Aplikace dostane na vstupu novinový článek v češtině nebo angličtině a na základě textu sama doporučí vhodné ilustrační obrázky. Implementovaná aplikace také vyhledává obrázky dle vizuální podobnosti. Práce se věnuje teoretickým aspektům extrakce klíčových slov, nebo detekce jazyka textu. Dále jsou rozebrány některé možnosti efektivního hledání podobných vektorů, které jsou využity v komponentě pro vyhledávání vizuálně podobných obrázků. Práce dále popisuje moderní možnosti vývoje webového frontendu i backendu. Algoritmus doporučování obrázků je v rámci práce otestován na uživateli.

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: The objective of this thesis is to implement a web application designed for recommendation of stock photos. The application gets the input from newspaper articles in Czech or English and, based on the text itself, suggests appropriate stock photos. The implemented application also searches images according to visual similarity. The thesis deals with theoretical aspects of keywords extraction and language of text detection. Further it analyzes possibilities of efficient search for similar vectors that are used in the search component for visually similar images. It also describes the possibilities in development of modern web frontend and backend. The quality of algorithm for recommending stock photos is tested on users.

Keywords: information retrieval, image retrieval

Obsah

| | | |
|----------|---|-----------|
| 1 | Úvod | 4 |
| 1.1 | Práce s daty | 4 |
| 1.2 | Vyhledávací algoritmy | 4 |
| 1.3 | Překlad | 4 |
| 1.4 | Webová aplikace | 4 |
| 1.5 | Vyhledávání podobných obrázků | 4 |
| 1.6 | Testování | 5 |
| 1.7 | Ostatní problémy | 5 |
| 2 | Poskytnutá data | 6 |
| 2.1 | Profimedia dataset | 6 |
| 2.2 | Profimedia dataset s detekovanými frázemi | 7 |
| 2.3 | Vektor podobnosti obrázků | 7 |
| 3 | Vyhledávání relevantních obrázků | 8 |
| 3.1 | Extrakce klíčových slov | 8 |
| 3.2 | TF-IDF | 8 |
| 3.3 | TF-IDF pro extrakci klíčových slov | 9 |
| 3.4 | Vyhledávání obrázků | 10 |
| 4 | Překlad obrázkových popisků | 11 |
| 4.1 | Strojový překlad | 11 |
| 4.2 | Frázový statistický strojový překlad | 11 |
| 4.3 | Charakteristika dat pro překlad | 12 |
| 4.4 | Překlad vět | 12 |
| 4.5 | Překlad slov | 13 |
| 4.6 | Překlad frází | 13 |
| 4.7 | Řešení | 13 |
| 4.8 | Závěr | 14 |
| 5 | Stemmer | 15 |
| 5.1 | Český stemmer | 15 |
| 6 | Detekce jazyka | 16 |
| 7 | Podobné obrázky | 18 |
| 7.1 | Bootstrap implementace | 19 |
| 7.2 | Předgenerování výsledků | 19 |
| 7.3 | Geohash | 20 |
| 7.4 | Řešení | 21 |
| 7.5 | Implementace | 22 |
| 8 | Backend | 23 |
| 8.1 | Databáze | 23 |
| 8.2 | Programovací jazyk | 23 |
| 8.2.1 | NodeJS | 24 |

| | | |
|-----------|--|-----------|
| 8.2.2 | Go | 24 |
| 8.2.3 | Ruby a Ruby on Rails | 25 |
| 8.3 | Komunikace mezi backendem a frontendem | 25 |
| 8.3.1 | Formát dat | 25 |
| 8.3.2 | REST API | 25 |
| 8.3.3 | WebSocket | 26 |
| 8.4 | Shrnutí | 27 |
| 9 | Frontend | 28 |
| 9.1 | HTML5 | 28 |
| 9.2 | JavaScriptové frameworky | 28 |
| 9.2.1 | jQuery | 28 |
| 9.2.2 | Google Closure | 29 |
| 9.3 | CSS | 29 |
| 9.4 | Uživatelské rozhraní | 30 |
| 9.4.1 | Rozhraní pro vyhledávání obrázků | 30 |
| 9.4.2 | Rozhraní pro detail obrázku | 31 |
| 9.5 | Závěr | 31 |
| 10 | Instalace a zprovoznění | 33 |
| 10.1 | Instalace aplikace | 33 |
| 10.2 | Vložení vstupních dat | 33 |
| 10.3 | Překlad | 33 |
| 10.4 | Jazykový korpus | 34 |
| 10.5 | Import dat do databáze | 35 |
| 10.6 | Data pro detekci jazyků | 36 |
| 10.7 | Vyhledávání podobných obrázků | 36 |
| 10.8 | Podpora dalších jazyků | 37 |
| 10.9 | Demo aplikace ve virtuálním stroji | 37 |
| 11 | Anotační rozhraní | 38 |
| 11.1 | Instalace rozhraní | 38 |
| 11.2 | Přidání uživatelů | 38 |
| 11.3 | Import anotačních dat | 38 |
| 11.4 | Export anotačních dat | 39 |
| 11.5 | Import obrázků a textů | 40 |
| 11.6 | Anotační proces | 40 |
| 12 | Testování výsledků | 42 |
| 12.1 | Vyloučení narušitele (o_test1) | 42 |
| 12.2 | Detekce správného obrázku (o_test2) | 44 |
| 12.3 | Shrnutí | 46 |
| 13 | Závěr | 47 |
| | Závěr | 47 |
| | Seznam použité literatury | 49 |
| | Příloha 1 | 53 |

| | |
|------------------------------|-----------|
| Příloha 2 | 56 |
| 13.1 Zdrojové kódy | 56 |
| 13.2 Demo | 56 |

1. Úvod

Cílem diplomové práce je implementovat webovou aplikaci pro doporučování ilustračních obrázků k danému textu. Název aplikace je „Stock Photo Finder“. V této práci ji budeme dále označovat slovem „aplikace“. Vytvořit takovou aplikaci přináší mnoho rozličných úkolů a problémů.

1.1 Práce s daty

Vstupní data z datasetu Profimedia obsahují více než 20 milionů anotací obrázků. Základním úkolem je být schopen takové množství dat nahrát do databáze a v tak velkých datech vyhledávat. Kapitola 8 rozebírá některé moderní programovací jazyky a databáze a jejich vhodnost pro práci s velkými daty.

1.2 Vyhledávací algoritmy

Klíčovou součástí aplikace je algoritmus pro extrakci klíčových slov a algoritmus pro vyhledávání v datasetu Profimedia. Tyto algoritmy jsou popsány v Kapitole 3. Praktická implementace algoritmů a proces získávání potřebných dat je popsán v Kapitole 10.

1.3 Překlad

Popisky klíčových slov jsou v datasetu Profimedia v anglickém jazyce. Cílem aplikace je umožnit vyhledávání doporučených ilustračních obrázků k textu i pro texty v jiných jazycích, než je angličtina. Kapitola 4 rozebírá překladové problémy z teoretického i praktického hlediska. Sekce 10.3 ukazuje praktickou při překladu do češtiny. Sekce 10.8 ukazuje postup pro zprovoznění aplikace i v dalších jazycích.

1.4 Webová aplikace

Výsledná aplikace je webovou aplikací. Kapitola 8 popisuje problémy při vývoji backendové části. Kapitola 9 popisuje moderní technologie pro vývoj frontendové části webových aplikací.

1.5 Vyhledávání podobných obrázků

Aplikace umožňuje kromě textového vyhledávání hledat ilustrační obrázky i podle vizuální podobnosti. Kapitola 7 popisuje algoritmy, které toto vyhledávání umožňují.

1.6 Testování

Výsledky detekce algoritmu byly testovány na uživateli. pro testování byla vyvinuta speciální webová aplikace popsána v Kapitole 11. Samotné testování algoritmu a jeho výsledky jsou popsány v Kapitole 12.

1.7 Ostatní problémy

Při práci na aplikaci bylo řešeno množství drobnějších problémů. Kapitola 5 popisuje stemming a vlastní implementaci stemmeru pro český jazyk. Kapitola 6 popisuje algoritmus pro automatickou detekci jazyka textu. Kapitola 11 popisuje anotační rozhraní, které bylo vytvořeno pro testování výsledků této aplikace.

2. Poskytnutá data

Tato kapitola popisuje data, která slouží jako vstup projektu.

2.1 Profimedia dataset

Společnost Profimedia¹ poskytla pro výzkumné účely korpus více než dvaceti milionů ilustračních obrázků spolu s jejich textovým popisem.

Textové popisky byly očištěny[10] 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.

Následující ukázka obsahuje příklad jednoho řádku souboru `profi-text-cleaned.csv`

```
"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
```

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 nijak explicitně vyznačené. Některým slovům chybí diakritika, například „ceska“. Některé fráze vznikly nejspíše strojovým překladem z cizího jazyka do angličtiny. To je vidět na frázi „hradec kings“, která původně zřejmě byla název 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.

K samotným obrázkům aplikace přistupuje přes API aplikace MUFIN².

¹<http://www.profimedia.cz/>

²<http://mufin.fi.muni.cz/profimedia>

2.2 Profimedia dataset s detekovanými frázemi

Některé problémy s datasetem Profimedia byly odstraněny v rámci bakalářské práce Bc. Jana Botorka[8]. Jedním z výsledků této práce je soubor `keyword-clean-phrase-export.csv`. Ukázka jednoho řádku tohoto souboru:

```
"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"
```

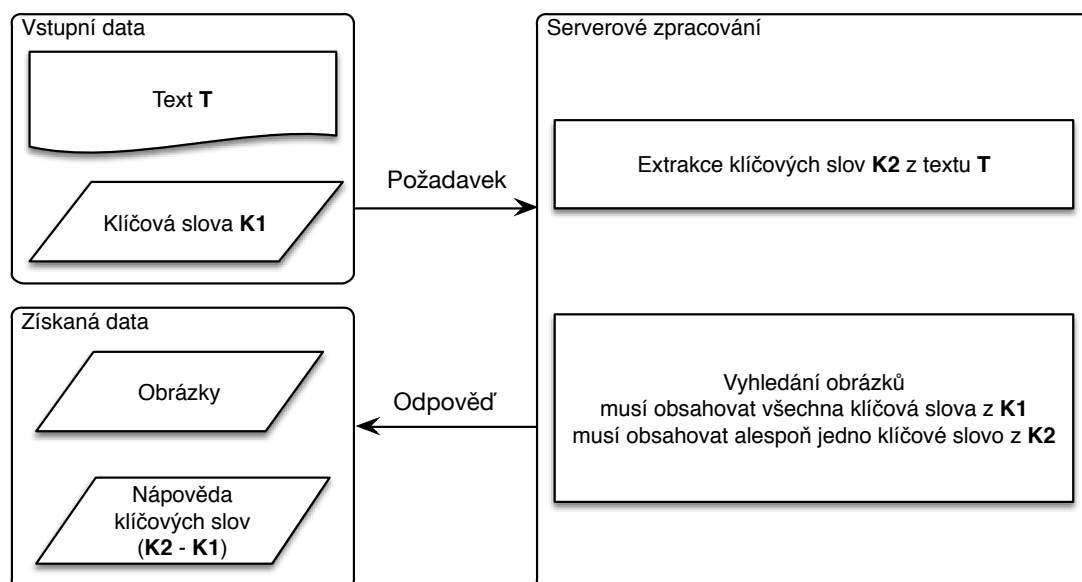
Soubor obsahuje lépe zpracovaná data z datasetu Profimedia. Nejdůležitější změnou pro tuto práci je detekce klíčových frází. Jednotlivé fráze jsou od sebe odděleny čárkou. Detekce probíhala pomocí databáze WordNet[35] a Wikipedie[41]. Ne vždy se detekce povedla, takže v některých řádcích nejsou žádné fráze detekovány. Detekce frází je důležitá zejména pro překlad dat do jiných jazyků.

2.3 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 GZ a má velikost 129 GB.

Soubor obsahuje pro každý obrázek z datasetu Profimedia vektor 4 096 reálných čísel. Pomocí vektorové vzdálenosti lze určit, jak jsou si podobné dva obrázky mezi sebou.

3. Vyhledávání relevantních obrázků



Obrázek 3.1: Základní tok dat mezi klientem a serverem.

Aplikace poskytuje uživateli dvě možnosti vyhledávání. Základním uživatelským scénářem je vložit do rozhraní text. Aplikace by v takovém případě měla poskytnout relevantní obrázky k danému textu. v dalším uživatelském scénáři je přímý požadavek na klíčová slova obrázku. Uživatel by měl mít možnost zadat přímo klíčová slova, které metadata k nalezeným obrázkům musí obsahovat. Oba scénáře by mělo navíc 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. Celý proces hledání vhodných obrázků popisuje diagram 3.1.

3.1 Extrakce klíčových slov

Extrakce klíčových slov je velmi důležitou složkou celého vyhledávání. Článek [31] 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.

3.2 TF-IDF

TF-IDF[26] je jeden ze základních vyhledávacích algoritmů[32]. Algoritmus využívá korpusu dokumentů D a dvou složek TF a IDF , lze ho vyjádřit jako rovnost

$$TFIDF(t, d, n, N) = TF(t, d) \times IDF(n, N). \quad (3.1)$$

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

$$TF(t, d) = \sum_{slovo \in d} \begin{cases} 1 & \text{pokud } slovo = t \\ 0 & \text{jinak} \end{cases}. \quad (3.2)$$

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

$$IDF(n, N) = \log\left(\frac{N}{n}\right). \quad (3.3)$$

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

$$TFIDF(t, d, n, N) = \left(\sum_{slovo \in d} \begin{cases} 1 & \text{pokud } slovo = t \\ 0 & \text{jinak} \end{cases} \right) \times \log\left(\frac{N - n}{n}\right). \quad (3.4)$$

3.3 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. Odstraníme slova ze stop-slovníku, který obsahuje nejběžnější slova pro každý z podporovaných jazyků¹. Dále nás nezajímá diakritika a různé speciální znaky, které může text obsahovat. Slova převedeme na stemy. O převodu slov na stemy se podrobněji píše v Kapitole 5.

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 Profimedie. 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 Profimedie. 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 na tomto datasetu přiřazuje 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 Wikipedie², jejíž data jsou volně dostupná pro mnoho jazyků pod otevřenou licencí.

¹Seznam slov ve stopslovnících je převzat z knihovny Lucene[2].

²<http://dumps.wikimedia.org/>

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 $Score$:

$$Score(w) = \begin{cases} Freq(w) \times \left(\frac{C}{Wiki(w)}\right) & w \text{ je v datasetu Profimédie} \\ 0 & w \text{ není v datasetu Profimédie} \end{cases} \quad (3.5)$$

C je experimentálně zjištěná konstanta. v našem případě je její hodnota 10 000 000.

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

3.4 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 3.3. 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 $Rank$:

$$Rank(K_{text}, K_{img}) = \sum_{w \in K_{text}} \begin{cases} Score(w) \times \frac{1}{|K_{img}|} & w \in K_{img} \\ 0 & w \notin K_{img} \end{cases} \quad (3.6)$$

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.

4. Překlad obrázkových popisků

Popisky obrázků v datasetu Profimedie 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 Profimedie.

My jsme se rozhodli pro druhou možnost. Musíme tedy přeložit všechny texty v datasetu Profimedie. 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.

4.1 Strojový překlad

Strojový překlad jako obor zažívá velký rozvoj. Potřeba překladu stále celosvětově prudce stoupá[11]. To jednak zvyšuje poptávku po kvalitním strojovém překladu a dále to strojovému překladu dává k dispozici velké množství lidských překladů, které jsou pro kvalitní strojový překlad nezbytné. Lze brát v úvahu i pravidlové překladové systémy, které ke své práci tolik dat většinou nepotřebují[3]. Obecnější a pro většinu jazykových párů nejlepší řešení však nabízí frázový statistický strojový překlad[28].

4.2 Frázový statistický strojový překlad

Frázový statistický strojový překlad ke své práci potřebuje databázi lidsky přeložených frází. Fráze jsou několikaslavné 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:

$$\tilde{e} = \arg \max_{e \in e^*} p(e|f) = \arg \max_{e \in e^*} p(f|e)p(e) \quad (4.1)$$

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.

Dobrý 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[29], která je dostupná pod volnou licencí včetně základních překladových modelů pro některé jazykové páry.

4.3 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.

4.4 Překlad vět

| |
|---|
| <p>Zdrojový dokument</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>Moses</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>Překladač Google</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 4.1: 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[23] — 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 frázi „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 Mosese 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 2.1¹. Obrázek 4.1 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.

4.5 Překlad slov

Jednodušším přístupem k překladu klíčových slov je přístup slovníkový, tedy 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“.

4.6 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[8] Jan Botorek. Zkoušel detekovat N-gramy v databázi WordNet[35] 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.

4.7 Řešení

Výsledná aplikace nakonec používá poslední možnost. Bližší informace jsou v Kapitole 10. Slova byly přeloženy pomocí Překladače Google, který překládá s lepšími výsledky než dostupný model pro překladový nástroj Moses.

Překlad frází by ovšem pomocí Překladače Google byl příliš drahý a pomocí Mosesem 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 18 006 z celkového počtu 899 244 detekovaných frází v Profimedia datasetu.

¹<http://www.statmt.org/moses/RELEASE-2.1/models/en-cs/model/>

4.8 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, zda má slovo „single“ přeložit jako „jednolůžkový“, nebo ve významu „jeden“. Lepší překlad by mohl přinést překladový model natrénovaný na speciálnější množině dat bližší korpusu Profimedia. 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.

5. Stemmer

Důležitou vlastností systémů na vyhledávání v textu je prohledávání ve více tvarech hledaného slova. Pokud uživatel hledá slovo „praha“, ve většině případů očekává, že se mu zobrazí i výsledky obsahující slovo „praze“. Je tedy potřeba mít nějaký algoritmus, který k sobě slova „praha“ a „praze“ přiřadí.

První možností je použít lemmatizér[32]. Úkolem lemmatizéru je ke každému slovu přiřadit jeho základní tvar. U podstatných jmen je to většinou první pád jednotného čísla („praha“), u sloves infinitiv.

Alternativou lemmatizátoru může být stemmer. Stemmer nemusí, na rozdíl od lemmatizéru, nemusí vrátit regulérní slovo jazyka. Například pro slova „praha“ a „praze“ může stemmer vrátit slovo „prah“. Výhodou stemmeru oproti lemmatizéru je, že většinou používá pouze jednoduché heuristiky. Je tedy často rychlejší a méně náročný na zdroje než lemmatizér. ve vyhledávacích aplikacích je stemmer dostačující.

Pro angličtinu je nejznámějším stemmerem pro angličtinu je Porterův stemmer[36] popsaný Martinem Porterem již v roce 1980. Kromě oficiální implementace existují porty do různých jazyků včetně Ruby[33]. Aplikace využívá implementaci z Ruby Gemu „stemmify“¹ v licenci MIT².

5.1 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[2]. Tuto implementaci využívá i Elasticsearch[6].

V rámci této práce byl portován soubor `CzechStemmer.java` ze zdrojového kódu knihovny Lucene do jazyka Ruby. Výsledkem je Ruby Gem³ pod licencí MIT, který lze použít nezávisle na zbytku aplikace. Instaluje se příkazem

```
gem install czech-stemmer
```

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

```
require 'czech-stemmer'

CzechStemmer.stem("praha") # => "prah"
CzechStemmer.stem("praze") # => "prah"
CzechStemmer.stem("předseda") # => "předsd"
CzechStemmer.stem("mladými") # => "mlad"
```

¹<https://rubygems.org/gems/stemmify>

²<http://opensource.org/licenses/MIT>

³<https://rubygems.org/gems/czech-stemmer>

6. 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. Podobnou možnost nabízí uživatelům také Překladač Google[23].

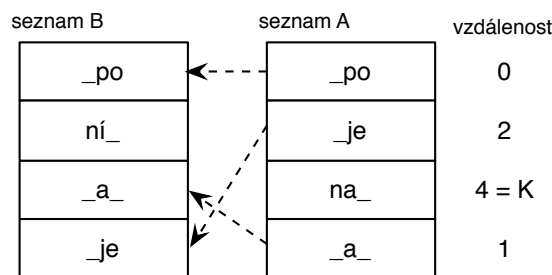
Aplikace používá algoritmus založený na statistice nejčastějších N-gramů pro daný jazyk[12]. Algoritmus nejprve použije velký korpus textů pro každý detekovaný jazyk. pro češtinu a angličtinu lze použít například korpus Wikipedie. z tohoto korpusu získá uspořádaný seznam K nejčastějších N-gramů. Při 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:

$$D(A, B) = \sum_{w \in A} \begin{cases} |A[w] - B[w]| & \text{pokud } w \in B \\ k & \text{jinak} \end{cases} \quad (6.1)$$

Na Obrázku 6.1 je ukázka porovnání dvou seznamů pomocí funkce $D(A, B)$. Trigram $na_$ se v seznamu B nenachází, proto dostane hodnotu $K = 4$.

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 $D(X, C)$.



Obrázek 6.1: Ukázka porovnání dvou trigramových seznamů. $D(A, B) = 7$

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.

7. Podobné obrázky

Jednou ze služeb, které výsledná aplikace poskytuje, je vyhledávání podobných obrázků. Uživatel rozhraní najde pomocí textových dotazů 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 Profimedie. Vektor má 4096 složek s reálnými nezápornými čísly. Vektory jsou vizuální deskriptory obrázků.

Vektory byly získány pomocí software Caffe[25]. Tento software slouží ke klasifikaci obrázků do vizuálních kategorií. Například můžeme natrénovat algoritmus na detekci kategorií „dítě“, nebo „ovoce“. pro klasifikaci obrázků z Profimedie bylo použito 1 000 takových kategorií. Informace o těchto kategoriích pak byly převedeny do vektorů, jejichž vektorová vzdálenost nám určuje míru podobnosti obrázku. Vektory jsou odezvami předposlední vrstvy hluboké konvoluční neuronové sítě natrénované pro klasifikaci obrázků z datasetu Profimedie do 1000 kategorií [7].

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

$$Similarity(I_1, I_2) = \sum_{i=1}^{4096} |D_1[i] - D_2[i]|. \quad (7.1)$$

V praxi se dají výsledky této míry klasifikovat zhruba do 3 kategorií. Tyto vypočítané kategorie popisuje tabulka 7. Ukázky jednotlivých kategorií podobných obrázků poskytuje obrázek 7.1. Rozdělení na 3 kategorie podobnosti podle míry *Similarity* je pouze přibližné. 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 |

Tabulka 7.1: 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 jen 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“ může uživatel chápat jako lepší výsledek. Tato úvaha umožňuje použít algoritmy, které nevrací obrázky seřazené podle *Similarity*, ale jsou výrazně rychlejší.



Obrázek 7.1: Dvojice obrázků s různou kategorií podobnosti.

7.1 Bootstrap implementace

První implementace nahrála všechny deskriptory do databáze Elasticsearch[6]. v Elasticsearch lze implementovat vyhledávání pomocí vlastní porovnávací funkce. Jako porovnávací funkci tedy zvolíme funkci *Similarity*. Bohužel i přes některé menší optimalizace algoritmu se vyhledávání nepodařilo implementovat efektivně, takže fungovalo v přijatelném čase pouze pro pár tisíc deskriptorů. Naše aplikace ovšem potřebuje pracovat s více než dvaceti miliony deskriptorů. Bylo by možné použít další optimalizace a sharding databáze na více strojů, aby algoritmus fungoval efektivněji i na větším množství dat. To by ale bylo příliš drahé a navíc se objevily jiné možnosti řešení.

7.2 Předgenerování výsledků

Efektivnější variantou je ke každému obrázku vygenerovat nějaké množství podobných obrázků 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

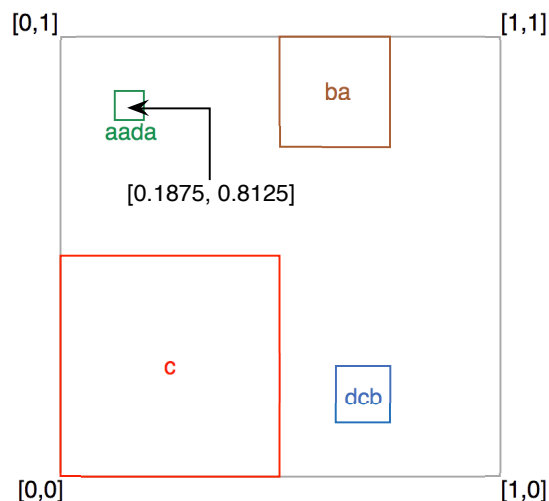
s použitím více vláken a optimalizačních heuristik, by vygenerování podobných obrázků pro každý z 20 milionů obrázků v datasetu Profimedia trvalo na běžném počítači několik týdnů.

7.3 Geohash

Další pokus o implementaci vyhledávání podobných obrázků se inspiruje algoritmem Geohash[39]. 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 nějakého čtverce z množiny oblasti prefixem jména bodu. k ukládání názvů bodů pak lze použít prefixový strom.

Obrázek 7.2 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*.

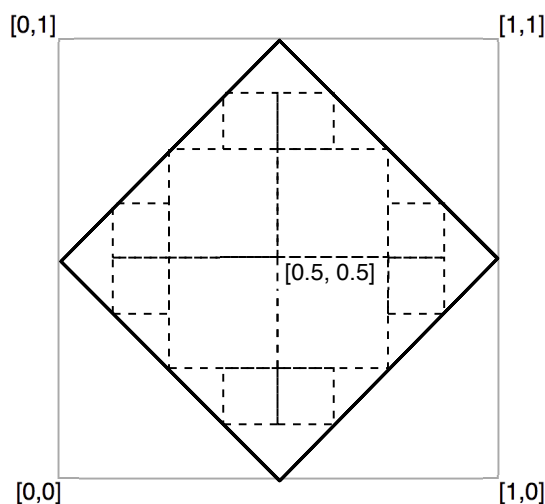


Obrázek 7.2: Ukázka čtverců algoritmu Geohash v jednotkovém čtverci.

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 7.3 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 7.3: 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.

7.4 Řešení

Naše řešení využívá princip algoritmu Geohash. Jde nám o to, převést hledání podobných vektorů na fulltextové vyhledávání. Vektory reálných čísel převedeme na množinu slov oddělenou mezerami — řetězec. Slova pochází z množiny $\{0, 1, \dots, 4095\}$. Každé složce vektoru deskriptoru tedy odpovídá právě jedno slovo z množiny slov. Nyní přiřadíme každému deskriptoru nějakou nějakou podmnožinu slov. Způsob přiřazení je klíčovým faktorem algoritmu, který ovlivňuje jeho efektivitu. Základním metodou přiřazení je přiřadit vektoru slova, která odpovídají nenulovým vektorům. Na deskriptorech obrázků z datasetu Profimedia to znamená, že každý vektor bude mít přiřazeno zhruba 1500 slov. To je stále příliš mnoho slov pro efektivní fulltextové vyhledávání. Elasticsearch v základní konfiguraci ani neumožňuje vyhledávat s dotazy delšími než tisíc slov. Musíme tedy množinu slov dále zmenšit. Můžeme například místo nenulovosti zvýšit požadovanou mez velikosti složky vektoru.

V našem řešení používáme jinou metodu. Seřadíme složky vektorů podle velikosti a slova přiřadíme pouze K složkám s nejvyšší hodnotou. pro $K = 4$ pak každému obrázku přiřadíme řetězec, například "2013 432 1065 3433". Při vyhledávání podobného obrázku nejprve získáme jeho deskriptor a jeho přiřazený řetězec. Pomocí řetězce najdeme fulltextovým vyhledáváním L obrázků. Výsledky seřadíme podle podobnosti s hledaným obrázkem mírou *Similarity* a vrátíme obrázky, které jsou v kategorii „podobné“, nebo „téměř shodné“. pro správné fungování je nutné dobře nastavit hodnoty K a L . Jejich zvýšení vede k větší přesnosti na úkor rychlosti algoritmu.

7.5 Implementace

Celá služba je implementována nezávisle na zbytku projektu. Jedná se o program `similar_img_finder` napsaný v jazyce Go. Jako databáze je použita Elasticsearch. Technické detaily služby vyhledávání podobných obrázků jsou uvedeny v Sekci 10.7.

8. Backend

Moderní webové aplikace lze rozdělit na back a frontend. Backend je část aplikace, která běží na serveru. Pomocí svých služeb poskytuje přístup k databázi, k souborům uloženým na serveru a zpracovává uživatelské operace. Frontend těchto služeb využívá. Tato kapitola popisuje proces výběru backendových technologií pro tuto práci. Některé technologie se osvědčily, jiné se ukázaly pro daný účel nevhodné.

8.1 Databáze

Úkolem databáze je uložit data a umožnit jejich prohledávání. Důležitou vlastností naší aplikace je, že uživatelé 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 letech vzniklo několik nových databází v kategorii vágně označené jako NoSQL[38]. Tato kategorie databází se těžko popisuje, na každou popsanou vlastnost existuje NoSQL databáze, která danou vlastnost nespĺňuje. Obecně ale lze říct že, NoSQL databáze nepracují s prvky v tabulkovém uspořádání a oproti standardním SQL databázím nekladou tolik omezujících požadavků na data. Jejich výhodou oproti standardním relačním databázím může být vyšší výkon a snadná škálovatelnost.

Nevýhodou je většinou obtížnější práce s daty. Většina NoSQL databází například nepodporuje databázové transakce a vůbec celý ACID. pro práci s daty v NoSQL databázi nelze použít klasické SQL dotazy. Místo nich se používá například model MapReduce[16] vyvinutý ve formě Google.

První verze aplikace byla postavena na databázi CouchBase[14], což je právě jedna z NoSQL databází podporující MapReduce mechanismus. Výhodou CouchBase je vysoký výkon, snadná škálovatelnost, velmi dobrá dokumentace a také existence oficiálních knihoven pro nejrozšířenější programovací jazyky. Ukázalo se však, že pro implementaci vyhledávání pro naši aplikaci je model MapReduce nedostatečný a implementace vyhledávání v CouchBase by byla prakticky nemožná.

Vhodnější pro daný účel se ukázala knihovna Elasticsearch[6]. Nejedná se v pravém smyslu o databázi. Jejím hlavním cílem je poskytnout vyhledávání nad daty. Je postavená nad knihovnou Apache Lucene ke které přidává snadnou horizontální i vertikální škálovatelnost a komunikaci pomocí REST HTTP JSON API. Díky tomu, že je Elasticsearch postaven na knihovně Lucene[2], může programátor využít velkou množinu možností, které Lucene poskytuje. Například v textovém vyhledávání může využít všechny tokenizery a stemmery z knihovny Lucene. v průběhu implementace aplikace navíc vyšla stabilní verze 1.0.

8.2 Programovací jazyk

Volba programovacího jazyka je při tvorbě backendové části klíčová. Tato sekce rozebírá některé výhody a nevýhody třech programovacích jazyků.

8.2.1 NodeJS

JavaScript je jazyk především pro programování frontendu a dlouho byl brán mezi programátory jako nutné zlo — bez použití problematických rozšíření nejde ve webovém prohlížeči programovat jiným jazykem. Každý větší webový projekt tak byl nucen používat nejméně 2 technologie — jednu pro frontend a jednu pro backend. Knihovna NodeJS toto paradigma obrací, umožňuje programovat v JavaScriptu i na backendu. Výhodou je sdílení stejného kódu mezi backendem a frontendem. To se hodí například při validacích dat, které potřebujeme kontrolovat na frontendu i backendu.

Původní záměr byl využít pro vývoj naší aplikace právě NodeJS. Nakonec však převážily nevýhody takového řešení. Javascript byl navržen pro programování webového frontendu, jeho standartní knihovna je v porovnání s ostatními jazyky velmi chudá, podpora objektového programování není přímočará. Obsáhlý kód v Javascriptu může být poměrně nepřehledný a jazyk svádí k vytvoření takzvaného „callback hell“. Ukázka takového problémového kódu je v článku [24]

```
doAsync1(function () {
  doAsync2(function () {
    doAsync3(function () {
      doAsync4(function () {
    })
  })
})
})
```

Pokud chce navíc programátor sdílet kód z backendu i na frontendu, musí být kód kompatibilní s podporovanými prohlížeči. Pokud aplikace potřebuje podporovat starší verze prohlížečů, nemůže sdílený kód obsahovat novinky z páté verze ECMAScript[44].

8.2.2 Go

Go (známý také jako golang)[22] je staticky typovaný programovací jazyk od společnosti Google. Syntaxe je inspirována jazykem C s důrazem rychlou kompilací. Velkou výhodou je snadná práce s vlákny pomocí „goroutines“. Programátoři v Javě, nebo C++ může překvapit poněkud netypická podpora práce s objekty.

První verze aplikace byla napsána právě v jazyce Go. Později byla přepsána většina v jazyce Ruby. Výhodou řešení v jazyce Go je velká výkonost a stabilita aplikace. Zkompilovaná aplikace je pouze jeden binární soubor, takže oproti řešení v nekompilovaných jazycích se lépe distribuuje a odpadají problémy se závislostmi.

Go je stále poměrně mladý jazyk, první stabilní byla vydána teprve v roce 2012[20]. Největším problémem takto mladých jazyků je nedostatek kvalitních knihoven. Pokud už na daný problém existuje knihovna, těžko se odhaduje, jestli její autor projekt neopustí a jestli bude knihovna podporována i další rok. pro naši aplikaci se ukázal jako problém neexistence kvalitní knihovny pro komunikaci s databází Elasticsearch. S tou jde komunikovat pomocí HTTP requestů, takže jsme pro naši potřebu vytvořili jakousi mikroknihovnu. Experimentální povaha aplikace ovšem vyžadovala rychlé prototypování a úprava kódu na komunikaci

s Elasticsearch začala být příliš omezující. pro naše účely se ukázalo výhodnější přepsat hlavní část aplikace do dynamického jazyka s lepší podporou pro Elasticsearch.

V jazyce Go je napsána nezávislá komponenta vyhledávání podobných obrázků, která je popsána v Kapitole 7.

8.2.3 Ruby a Ruby on Rails

Ruby[33] je dynamicky typovaný jazyk, inspirovaný jazyky jako Perl, nebo Smalltalk, s důslednou podporou objektového programování. Popularita Ruby vzrostla zejména díky webovému frameworku Ruby on Rails (Rails)[37], který je v Ruby napsaný. Rails zpopularizovaly architekturu Model-View-Controller[40] pro tvorbu webových aplikací. Backend aplikace s výjimkou služby pro vyhledávání podobných obrázků, byl nakonec přepsán právě do Rails. Elasticsearch poskytuje pro práci s Ruby oficiálně podporovanou knihovnu[5]. Další výhodou je knihovna Rake, což je obdoba příkazu Make pro skripty v Ruby. Pomocí Rakefilu jdou napsat přehledné úlohy pro manipulaci s daty.

8.3 Komunikace mezi backendem a frontendem

Tato sekce popisuje možnosti komunikace mezi frontendem a backendem.

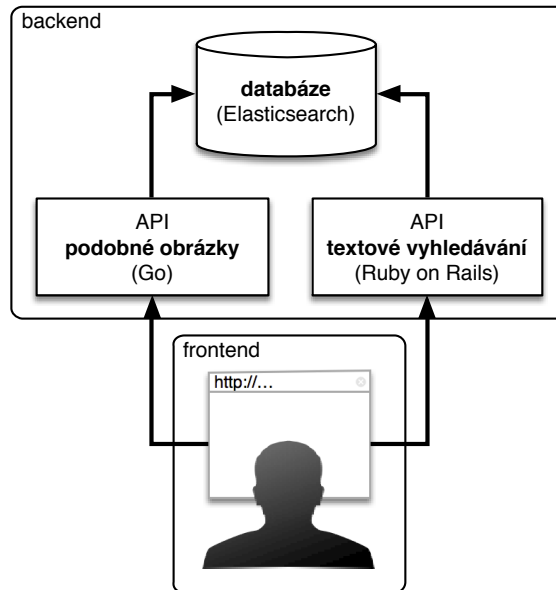
8.3.1 Formát dat

Nejprve je potřeba zvolit formát dat. v současnosti existují dva nejrozšířenější textové komunikační formáty — XML[9] a JSON[15]. XML umožňuje popis velkého množství vztahů mezi entitami. JSON je oproti XML úspornější. XML se hodí v případech, kdy potřebujeme zaznamenat nějaké složité datové struktury. pro většinu webových aplikací je výhodnější použít JSON. Oproti XML je zápis dat v JSON většinou kratší, což je výhodné zejména při přenosu dat po síti. JSON je založen na datových typech JavaScriptu a moderní prohlížeče umí tento formát zpracovávat nativně. Na frontendu je tedy práce s JSON oproti práci s XML jednodušší. z těchto důvodů poskytuje backend naší aplikace služby frontendu ve formátu JSON.

8.3.2 REST API

Dále je potřeba specifikovat způsob výměny dat. Jednou z nejrozšířenějších metod je REST, zkratka pro anglické „Representational state transfer“. Jedná se o architekturu rozhraní mezi klientem a serverem. pro specifikaci rozhraní využívá protokol HTTP. REST využívá k definici práce s daty metody protokolu HTTP „GET“, „PUT“, „POST“ a „DELETE“.

Rozhraní REST[19] s použitím formátu JSON je v aplikaci využíváno pro komunikaci s databází Elasticsearch i pro komunikaci mezi backendem a frontendem.



Obrázek 8.1: Architektura aplikace.

8.3.3 WebSocket

Alternativou k REST API a protokolu HTTP je protokol WebSocket[18]. WebSocket je stejně jako protokol HTTP postaven nad protokolem TCP. Na rozdíl od HTTP ale poskytuje duplexní spojení. Klient je tedy 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[30]. Protokol je podporován všemi moderními verzemi webových prohlížečů a podpora existuje i v knihovnách pro backendové programovací jazyky.

První verze aplikace používaly 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 technických problémů. Například v okamžiku, kdy se toto spojení přeruší. Pokud se naproti tomu přeruší spojení server-klient při HTTP požadavku, 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 s odeslanými zprávami — požadavky — ale musíme implementovat sami, například pomocí unikátních identifikátorů v těle zprávy. Naše implementace navíc musí vyřešit situaci, kdy žádná odpověď ze serveru nedorazí. Například nastavením timeoutu pro čekání na odpověď.

Řešení pomocí protokolu WebSocket využijí zejména aplikace, které potřebují komunikovat obousměrně mezi klientem a serverem, nebo je pro ně důležitá nízká latence získání odpovědi. Takovými aplikacemi mohou být chatovací služby, nebo online hry. pro většinu ostatních aplikací budou zatím spíše převažovat nevýhody obtížnějšího technického řešení nad přínosy technologie WebSockets.

8.4 Shrnutí

Architektura je zobrazena na Obrázku 8.1. Část aplikace, která se stará extrakci klíčových slov a textové vyhledávání obrázků, je napsána v jazyce Ruby a frameworku Ruby on Rails. Frontend dále komunikuje se službou pro vyhledávání podobných obrázků, která je napsána v jazyce Go. Data jsou uložena v databázi Elasticsearch. Komunikace mezi všemi komponentami probíhá pomocí HTTP REST API, data jsou přenášena ve formátu JSON.

9. Frontend

Frontend je část webové aplikace, která komunikuje přímo s uživatelem, posílá dotazy na služby backendu a zobrazuje uživateli výsledky.

9.1 HTML5

Historicky měli frontendoví weboví developeři na výběr mezi několika technologiemi. Proprietární řešení jako Adobe Flash[1] nebo Microsoft Silverlight[34] byly dříve jednou z mála možností, jak vyvíjet složitější frontendové aplikace. Jejich místo dnes nahrazuje skupina otevřených standardů pod souhrnným označením HTML5.

HTML5 je označení jednak pro nejnovější verzi standardu značovacího jazyka HTML[42], jednak je to termín zastřešující další moderní webové technologie, jako je například CSS3[43]. HTML5 poskytuje frontendovým vývojářům široké spektrum nových API. Bez proprietárních pluginů je dnes možné v moderních prohlížečích například přehrávat video, nebo vykreslovat 2D grafiku.

Další zlepšení, které frontendovým vývojářům zlehčuje práci, jsou nové verze standardu ECMAScript[17], což je standard jazyka JavaScript. JavaScript je nyní daleko produktivnější jazyk, než byl ještě před pár lety a může využívat více API ze standardu HTML5.

9.2 JavaScriptové frameworky

Jednotlivé webové prohlížeče se dlouhá léta lišily v implementaci JavaScriptu, a to do takové míry, že pro podporu více prohlížečů bylo použito nějakého frameworku, který sjednocuje API, téměř nezbytné. Dnes už do značné míry nejsou při vývoji front-endu frameworky nezbytné, ale stále jsou velmi užitečné.

9.2.1 jQuery

Nejpopulárnějším frameworkem současnosti je jQuery[27]. Ten nabízí jednoduché rozhraní pro práci s DOM a javascriptovými událostmi a pro velké množství menších webových aplikací je dostačující. ve větších aplikacích se začne projevit nižší výkon aplikace používající jQuery a také obtížné odstraňování chyb v aplikaci — debugování. Například můžeme pomocí jQuery chtít přidat CSS třídu „red“ HTML elementu s id „okno“

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

Tento kód proběhne bez nahlášené chyby i v případě, že žádný element s id „okno“ neexistuje, nebo když uděláme překlep v zápisu řetězce. Pokud chce web knihovnu jQuery, musí ji celou nejdříve stáhnout. Aktuální verze 1.11.1 má velikost téměř 100 KB. Knihovna jQuery je použita v implementaci anotačního rozhraní popsaném v Kapitole 11.

9.2.2 Google Closure

Jiný přístup k vývoji frontendových aplikací přináší Google. pro svou první webovou aplikaci Gmail¹ vyvinul sadu nástrojů, kterou později vydal jako sadu nástrojů Google Closure[21]. Kromě služby Gmail používá Google sadu nástrojů i pro vývoj Google Vyhledávání², Google Map³, nebo Google Dokumentů⁴. Skládá se ze tří částí - Closure Compiler, Closure Library a Closure Templates.

Closure Library je javascriptová knihovna funkcí pro práci s DOMem, javascriptovými událostmi, matematickými výpočty a spoustou dalších věcí, které webový programátor může využít.

Closure Compiler je minifikátor javascriptového kódu. Minifikace javascriptového kódu je důležitá, protože snižuje množství dat, které musí webová stránka před spuštěním stáhnout. Compiler odstraňuje nedůležité mezery v kódu, odstraňuje funkce, které nejsou volány, přejmenovává názvy funkcí a proměnných na co nejkratší řetězce a v ADVANCED módu se snaží i o pokročilejší optimalizace kódu (například inlining funkcí). Nejlepších výsledků dosahuje s použitím typových anotací. Anotace jsou silně inspirovány jazykem Java a díky nim lze z JavaScriptu vytvořit typovaný objektový jazyk s privátními metodami nebo dědičností tříd. Closure Compiler tyto anotace vyhodnocuje a vrací chyby a varování, podobně jako jiné kompilované jazyky. To umožňuje snadněji psát spolehlivý javascriptový kód.

Třetí částí sady nástrojů Google Closure jsou Closure Templates, šablonovací systém pro JavaScript a Javu. Pomocí Closure Templates se snadno vytváří zanořené HTML šablony se zadanými parametry. Closure Templates automaticky ošetřují vstupní data, což zabraňuje bezpečnostním útokům na stránku.

Částí Templates a Compiler lze 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.

Frontendová část aplikace je napsána s použitím všech částí Google Closure. Díky tomu je frontendová část serveru tvořena jediným javascriptovým souborem velikosti 65 KB.

9.3 CSS

CSS ve verzi 3 přidává mnoho nových vlastností, které je možné HTML elementům přiřadit. Nejviditelnějšími jsou kulaté okraje nebo stínování. pro programátory je zajímavá vlastnost „calc“ pomocí níž lze zadat hodnoty v CSS jako výpočet. Pokud například chceme, aby se element roztáhl na polovinu obrazovky a měl vedle sebe ještě 10 pixelů místo, nastavíme mu vlastnost

```
width: calc("50% - 10px");
```

Bez této vlastnosti bylo možné dříve podobné chování simulovat pouze pomocí JavaScriptu. Vlastnosti jako „calc“, umožňují programátorům snadněji psát

¹<http://www.gmail.com/>

²<http://www.google.com/>

³<http://maps.google.com/>

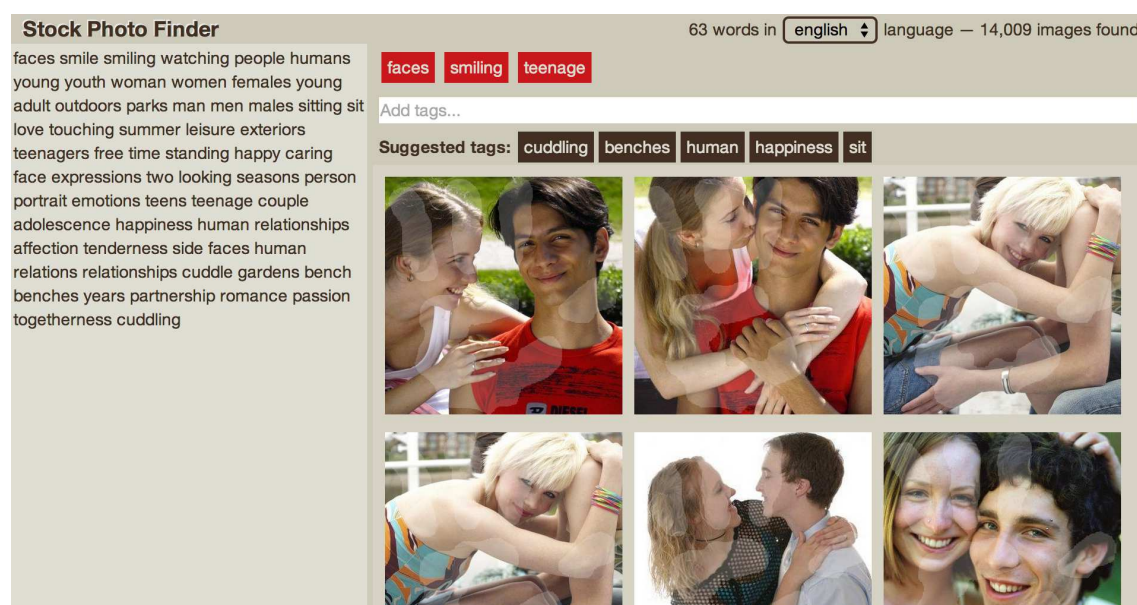
⁴<http://docs.google.com/>

responzivní aplikace. To jsou takové aplikace, které se dobře přizpůsobují velikosti okna prohlížeče. Responzivnost webových aplikací začala být důležitá zejména s rozvojem mobilních technologií. Stále větší počet uživatelů přistupuje na weby ze svých mobilních telefonů nebo tabletů a tito uživatelé chtějí mít při používání aplikace podobný komfort jako na při používání na počítači. Při vývoji front-endové části aplikace byl kladen důraz právě na snadnost používání aplikace v různých rozlišeních monitoru. pro mobilní telefony ovšem frontend aplikace optimalizován není.

9.4 Uživatelské rozhraní

Uživatelské rozhraní aplikace tvoří jediná webová stránka. Uživatel pracuje ve dvou základních režimech — vyhledávání obrázků a detail obrázku.

9.4.1 Rozhraní pro vyhledávání obrázků



Obrázek 9.1: Uživatelské rozhraní pro vyhledávání ilustračních obrázků.

Rozhraní pro vyhledávání obrázků slouží uživateli k vyhledání obrázku podle jím zadaných textových parametrů. Ukázka rozhraní je na Obrázku 9.1. Rozhraní lze rozdělit na 3 části — záhlaví, levý panel a pravý panel. Přetahováním myši za oblast rozdělující levý a pravý panel je možné měnit jejich velikost.

V levém panelu se nachází oblast pro vložení textu. Zde může uživatel vložit text článku, ke kterému hledá ilustrační obrázky. Uživatelské rozhraní samo detekuje změny v textu a zobrazuje ilustrační obrázky, aniž by uživatel musel stisknout tlačítko.

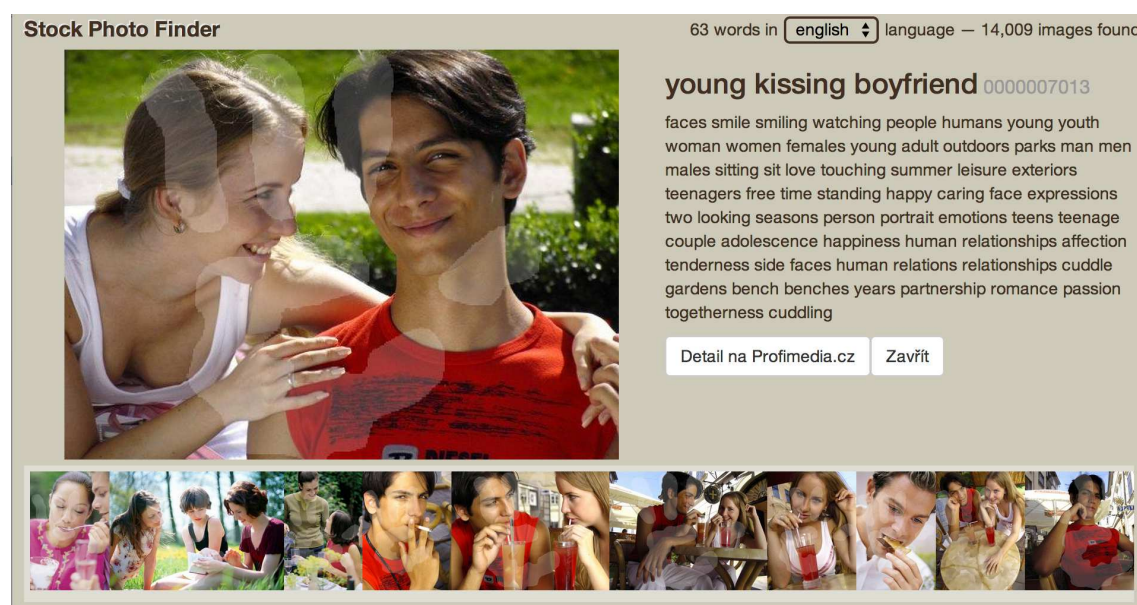
Záhlaví obsahuje statistiky o počtu nalezených obrázků a počtu slov v textu. v rozbalovacím menu je vybrán automaticky detekovaný jazyk vstupního textu. Uživatel má možnost tento jazyk manuálně změnit.

Pravý panel je rozdělen na horní a spodní část. ve spodní části jsou zobrazeny náhledy nalezených obrázků. Po najetí myši nad obrázek se zobrazí u obrázku i

jeho název a klíčová slova z datasetu Profimedia. Celý levý panel funguje na principu „nekonečného“ skrolování — po naskrolování na spodek panelu se načtou další obrázky z výsledků vyhledávání (pokud ještě nějaké nezobrazené výsledky vyhledávání existují). Po kliknutí na obrázek myši se zobrazí rozhraní s detailem obrázku.

Horní část pravého panelu obsahuje komponenty pro explicitní zadání klíčových slov, které uživatel požaduje u vyhledaných obrázků. Klíčová slova může uživatel přidat buď přímo zápisem do textového vstupu „Add tags...“ a stisknutím klávesy ENTER, nebo může kliknout na některé z detekovaných klíčových slov pod textovým vstupem. Nad textovým vstupem se na červeném pozadí nachází explicitně požadovaná klíčová slova. Kliknutím na explicitně vybrané slovo se slovo odstraní z požadované množiny klíčových slov.

9.4.2 Rozhraní pro detail obrázku



Obrázek 9.2: Uživatelské rozhraní s detailem obrázku.

Rozhraní s detailem obrázku vyvolá uživatel kliknutím na obrázek ve vyhledávacím rozhraní. Ukázka tohoto rozhraní je na Obrázku 9.2.

V levé horní části vidí uživatel celý obrázek. Pravá horní část rozhraní obsahuje textové popisky obrázku v datasetu Profimedia — název, id a klíčová slova. Pod nimi jsou dvě tlačítka. První je odkazem na detail obrázku na webu Profimedia, druhým tlačítkem uživatel zavře detail obrázku a vrátí se zpět do vyhledávacího rozhraní.

Ve spodní části rozhraní je seznam podobných obrázků k obrázku v detailu. Po kliknutí na některý z podobných obrázků se zobrazí detail příslušného obrázku.

9.5 Závěr

Frontend aplikace zpřístupňuje uživateli prohledávání datasetu Profimedia. Při vývoji byl kladen důraz na uživatelskou přívětivost a použití nejmodernějších fron-

tendových technologií. v celém rozhraní jsou tak například pouze dvě tlačítka. Rozhraní se samo přizpůsobuje velikosti obrazovky uživatele. pro uživatele by měla být komfortní i automatická detekce jazyku vloženého textu nebo možnost měnit velikost panelů ve vyhledávacím rozhraní.

10. Instalace a zprovoznění

Backend aplikace je webová aplikace napsaná v jazyce Ruby a frameworku Ruby on Rails. Je k dispozici pod svobodnou licencí MIT. k jejímu spuštění je nutná instalace jazyka Ruby verze alespoň 2.0 (nižší verze nejsou otestované) a Javy 7. Aplikaci lze spustit na platformách Linux a Mac.

10.1 Instalace aplikace

Zdrojový kód je volně dostupný na webu GitHub¹, nebo na přiloženém DVD. Stáhnou aplikaci lze příkazem

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

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

```
bundle install
```

Instalace může vyžadovat přístup administrátora. Dále je nutné stáhnout knihovnu Elasticsearch² do adresáře bin/elasticsearch. Stačí verze 1.0 a vyšší.

Nyní je možné celý projekt spustit pomocí Bash příkazů. 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 AJAXové dotazy skončí chybou. v databázi nejsou importována data.

10.2 Vložení vstupních dat

Metadata k obrázkům a obrázky samotné jsou poskytovány firmou Profimedia a nejsou volně dostupné. Ke zprovoznění aplikace je nutné vložit CSV soubor `profi-text-cleaned.csv` obsahující dataset Profimedia do adresáře `data`. Dále potřebujeme do adresáře `data` vložit soubor `keyword-cleaned-phrase-export.csv` z bakalářské práce Jana Botorka, který obsahuje dataset Profimedia s detekovanými frázemi.

10.3 Překlad

Pro překlad použijeme metodu popsanou v Kapitole 4. k překladu budeme potřebovat frázovou tabulku pro anglicko-český překlad. Použijeme tu, která je vydána s překladovým nástrojem Moses verze 2.1³.

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

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

³<http://www.statmt.org/moses/RELEASE-2.1/models/en-cs/model/>

Nyní můžeme použít příkaz

```
rake data:export_profimedia_words_for_translation
```

který vytvoří soubor `word_list.txt` s 352 862 slovy. Nyní je potřeba vytvořit soubor s překladem všech slov. pro naši aplikaci jsme využili překladový nástroj Překladač Google. Při tak velkém objemu dat jsme byly nuceni použít přístup přes placené API, které společnost Google poskytuje. Překladové API jsme použili skrz překladový nástroj firmy Memsources. Tímto strojovým překladem jsme získali soubor `word_list_translated_cs.txt` se všemi slovy z datasetu Profimedia přeloženými do češtiny. Příkazem

```
rake data:create_word_dictionary
```

můžeme seznam anglických a českých slov spárovat do jednoho souboru `word_dictionary_en_cs.txt`.

Kvůli lepší kvalitě překladu jsme se snažili přeložit i detekované fráze z datasetu Profimedia. Ze souboru `keyword-cleaned-phrase-export.csv` se fráze exportují příkazem

```
rake data:export_profimedia_phrases_for_translation
```

do souboru `phrases_list.txt`. Příkaz

```
rake data:translate_phrases
```

pak tyto exportované fráze porovná s frázemi přítomnými v překladovém modelu Moses. Fráze, které se v překladovém modelu nachází, jsou pak i s překladem uloženy v souboru `phrase_table_en_cs.txt`.

Získali jsme tedy překladový slovník pro slova (`word_dictionary_en_cs.txt`) i fráze (`phrase_table_en_cs.txt`). Algoritmy, které tato data používají jsou implementovány v souboru `language_tool.rb`.

10.4 Jazykový korpus

Aplikace potřebuje jazykový korpus pro podporované jazyky ze dvou důvodů. Zaprvé je korpus potřeba pro správné fungování algoritmu TF-IDF popsaném v Kapitole 3. Zadruhé je potřeba jazykový korpus pro získání nejfrekventovanějších trigramů jazyka. Tato data se pak využívají v algoritmu pro detekci jazyka, který je popsán v Kapitole 6.

Samotný dataset Profimedia se jako korpus pro naše účely nehodí. Struktura textů v datasetu neodpovídá běžnému textu a navíc je dataset pouze v angličtině. Jako korpus jsme použili data z Wikipedie. Ta je možné stáhnout jako databázový dump ve formátu XML a pod licencí Creative Commons. Stáhli jsme tedy soubory `wiki_dump_cs.xml` a `wiki_dump_en.xml`. Wikipedie nabízí pro převod z XML do obyčejného textu vlastní pythonovský skript `lib/WikiExtractor.py`[4]. pro extrakci dat z anglického XML souboru lze použít příkaz

```
rake wiki:extract_words_from_wiki en
```

který vytvoří adresářovou strukturu s textovými soubory. Není potřeba převést všechna data, pouze tolik abychom dostali reprezentativní korpus. pro angličtinu stačí převést zhruba 10 000 článků, pro podobné množství českých dat je potřeba převést zhruba 20 000 článků (údaj o exportovaných článcích je průběžně vypisován na konzoli).

Nyní je možné vytvořit frekvenční seznam slov z korpusu Wikipedie. ve skutečnosti potřebujeme frekvenční seznam stemů slov. Ten získáme příkazem

```
rake wiki:frequency_list_from_wiki
```

Nyní příkazem

```
rake data:create_tf_df_list
```

získáme pro každý stem v Profimedia datasetu jejich celkovou četnost (TF) a četnost metadat obrázků, kde se stem nachází.

Nakonec spárujeme frekvenční informace o stemech z Profimedia datasetu a korpusu Wikipedie příkazem

```
rake data:pair_profimedia_and_wiki_data
```

Získali jsme dva soubory — `paired_wiki_and_profimedia_cs.txt` pro češtinu a `paired_wiki_and_profimedia_en.txt` pro angličtinu. Soubory obsahují tabulátorem odsazené statistiky četnosti pro všechny stemy v datasetu Profimedia.

10.5 Import dat do databáze

Nyní máme připravená všechna data pro nahrání do databáze. Jako databáze nám slouží Elasticsearch, který musí být připojen na portu 9200. pro práci s Elasticsearch se velmi hodí plugin `elasticsearch-head`, který umožňuje prohlížení databáze ve webovém prohlížeči.

Nejprve importujeme frekvenční data pro všechny stemy příkazem

```
rake es:import_word_data
```

Pro každý z podporovaných jazyků takový import trvá zhruba jednu hodinu. Dále příkazem

```
rake es:import_image_metadata
```

importujeme metadata pro všechny obrázky z datasetu Profimedia. Data jsou importována pro všechny podporované jazyky. pro neanglické jazyky jsou použity

pro překlad slovníky slov a frází vytvořené v Sekci 10.3. Nahrání těchto dat do databáze může trvat několik hodin. Nahrávací skript zobrazuje odhad zbývajících času.

10.6 Data pro detekci jazyků

Algoritmus pro automatickou detekci jazyka, popsáný v Kapitole 6 potřebuje ke své práci seznam nejfrekventovanějších trigramů pro každý jazyk. Tento seznam 300 nefrekventovanějších slov pro každý jazyk získáme z korpusu Wikipedie příkazem

```
rake trigrams:extract_most_frequent_trigrams
```

Detekce jazyka probíhá ve frontendové části aplikace (viz Kapitola 9). Příkaz

```
rake trigrams:trigrams_to_javascript_classes
```

vytvoří ze seznamů nejfrekventovanějších seznamů javascriptovou třídu `oo.diplomka.Languages.Data` v notaci Google Closure. Ta je uložena v souboru `public/js/js/diplomka/languages/data.js`. Tuto třídu pak využívá frontend pro automatickou detekci jazyka vloženého textu.

10.7 Vyhledávání podobných obrázků

Služba vyhledávání obrázků je nezávislá na zbytku aplikace. Princip jejího fungování je popsán v Kapitole 7. Zdrojový kód služby je v jazyce go a je dostupný na GitHub⁴. pro kompilaci služby je potřeba go alespoň verze 1.2. Příkaz `go build` v adresáři se zdrojovými kódy služby zkompiluje program do souboru `similar_img_finder`.

Pro zprovoznění služby je nejprve nutné naimportovat data do databáze příkazem

```
./similar_img_finder -a import -n 1000000 -e 9200 -k 500 -f data.gz
```

Parametr `-n` určuje, kolik dat se má naimportovat, parametr `-e` určuje na kterém portu běží databáze Elasticsearch (můžeme použít stejnou instanci databáze na které běží hlavní aplikace), parametr `-k` odpovídá hodnotě K z Kapitoly 7 a parametr `-f` je cesta ke komprimovanému souboru s importovanými deskriptory.

Nyní můžeme spustit službu hledání podobných obrázků na portu 8585 příkazem

```
./similar_img_finder -a server -p 8585 -l 100 -e 9200
```

⁴https://github.com/hypertornado/similar_img_finder

Parametr `-l` odpovídá hodnotě L z Kapitoly 7.

Pokud nyní chceme získat id podobných obrázků k obrázku s id „0000000003“, můžeme využít JSON HTTP API a získat výsledky na adrese `localhost:8585/similar?id=0000000003`. Služba vrátí pole obsahující id podobných obrázků a míru *Similarity* pro každý vrácený obrázek.

10.8 Podpora dalších jazyků

Popsaný postup vytvoří aplikaci, která funguje pro anglický a český vstupní text. Aplikace je ale navržena tak, aby rozšíření do dalších jazyků nebylo obtížné. Popíšeme postup pro přidání dalšího jazyka — francouzštiny. Podle standardu ISO 639-1 je kód francouzštiny „fr“. Upravíme tedy konstantní pole aplikace v souboru `application.rb` — do pole `SUPPORTED_LANGUAGES` přidáme řetězec „fr“ a jako název jazyka přidáme do pole `LANGUAGE_NAMES` řetězec „french“.

Pro každý přidáný jazyk potřebujeme implementovat stemmer. pro podporu francouzského stemmeru budeme muset rozšířit metodu `stem_word` v souboru `language_tool.rb`. pro francouzštinu je v Ruby volně dostupných několik knihoven pro převod slov na stemy. Dále ve stejné třídě můžeme pro nově přidáný jazyk rozšířit datovou strukturu `@stopwords` o další stop-slova.

Dále pro francouzštinu stáhneme dump Wikipedie jako soubor `wiki_dump_fr.xml` a zopakujeme celý postup ze Sekcí 10.3, 10.4, 10.5 a 10.6 této kapitoly.

Podobnými kroky můžeme implementovat poměrně přímočaře i podporu pro další jazyky vstupních článků.

10.9 Demo aplikace ve virtuálním stroji

Na přiloženém DVD se nachází virtuální stroj s demoverzí aplikace. Demoverze z důvodů licence datasetu Profimedia a také omezené kapacitě DVD pracuje pouze s prvními 1 000 obrázky v datasetu. Podrobnosti o datech na přiloženém DVD obsahuje příloha

11. 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.

11.1 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 GitHub¹ a nejlépe lze stáhnout pomocí gitu. pro běh serveru je potřeba verze ruby 2.0 a vyšší. Celá aplikace se zprovozní následujícím pořadím BASH příkazů:

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

11.2 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/admin_controller.rb`. Uživatelé mají pouze dvě datové položky — uživatelské jméno (Name) a heslo (Password). Uživatele lze přidávat, mazat a upravovat. Nepředpokládá se, že by anotovaná data byla vysoce citlivá, heslo je proto v databázi uloženo v plaintextu.

11.3 Import anotačních dat

Data pro anotaci lze nahrát pomocí příkazu

```
rake data:import
```

Příkaz očekává 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:

¹https://github.com/hypertornado/cemi_annotace

INDEX

unikátní číslo jedné anotace.

LABEL

interní popis pokusu.

PRIORITY

priorita, celé číslo ≥ 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 basics 0 - ./text/aha/aha-00263.txt.gz
  img/1.jpg;img/2.jpg;img/3.jpg
2 basics 1 - ./text/aha/aha-00006.txt.gz img/1.jpg;img/2.jpg
3 basics 0 - ./text/aha/aha-00009.txt.gz img/2.jpg
```

11.4 Export anotačních dat

Hotové anotace lze exportovat příkazem

```
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.

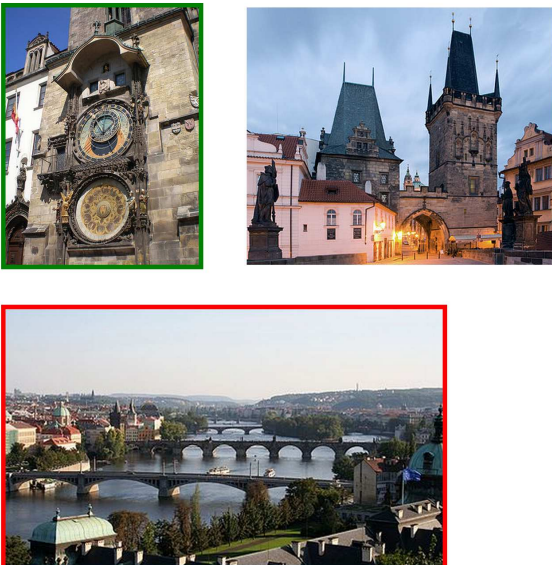
Anotace (id 2) x

localhost:3000

CEMI anotace Instrukce Přeskočit anotaci Uložit a ukončit Uložit a načíst další (ctrl+ENTER) Odlásit uživatele ondra

Praha je hlavní a současně největší město České republiky a 15. největší město Evropské unie. Leží mírně na sever od středu Čech na řece Vltavě, uvnitř Středočeského kraje, jehož je správním centrem, ale jako samostatný kraj není jeho součástí. Je sídlem velké části státních institucí a množství dalších organizací a firem. Sídlí zde prezident republiky, parlament, vláda, ústřední státní orgány a jeden ze dvou vrchních soudů. Mimoto je Praha sídlem řady dalších úřadů, jak ústředních, tak i územních samosprávných celků; sídlí zde též ústředí většiny politických stran a centrály téměř všech církví, náboženských a dalších sdružení s celorepublikovou působností registrovaných v ČR.

Do dnešní podoby se Praha vyvíjela



Obrázek 11.1: Anotační rozhraní. Vhodné obrázky jsou označeny 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.

11.5 Import obrázků a textů

Anotační texty a obrázky musí být nahrány do adresáře `ROOT_APLIKACE/public` tak, aby jejich cesty odpovídaly 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 gzip².

11.6 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

²<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.

12. Testování výsledků

Výsledný program se skládá z několika komponent. Funkčnost každé komponenty lze testovat. Některé části je možné 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. Výhodou 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 11.

12.1 Vyloučení narušitele (o_test1)

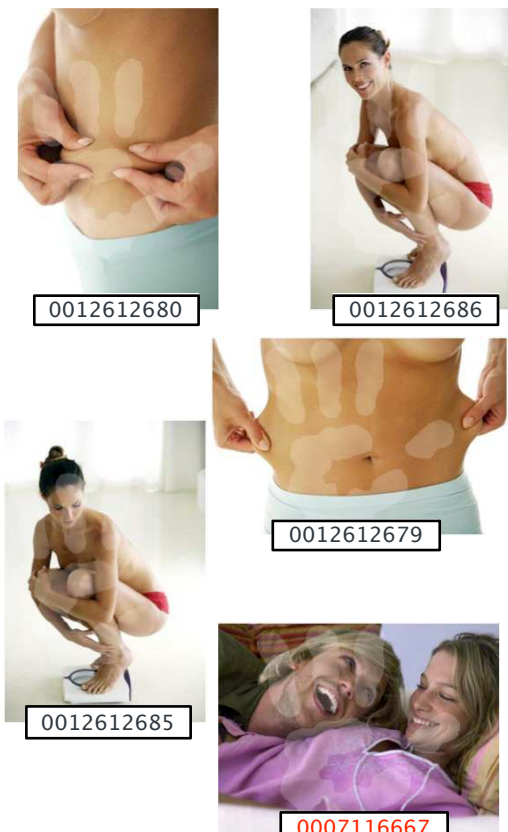
idn-07792.txt.gz

Ideální žena podle mužů: krásná violoncellistka s láskou k fotbalu

Sníme o ideálních mužích, ale víme, jakou mají oni představu o ideální ženě? Možná by si přáli po svém boku štíhlou vlnadnou blondýnu s dlouhými nohama schopnou zapařit, fandit fotbalu a ujet za odpoledne sto kilometrů na kole... Zeptaly jsme se jich.

Žena v červených lodičkách kope do fotbalového míče. I

Šimon, 27 let, student, zadaný
Ideální žena je plná kompromisů. Pro mě je to určitě slečna, která je atraktivní, protože jinak by si jí samozřejmě člověk nikdy nevyšiml, ale ne na úkor její inteligence a intelektu. Což si myslím, že dneska je bohužel většina případů... Ale myslím si, že když má slečna vnitřní kouzlo, spojené se špetkou atraktivity a inteligence, tak je ideální...



Obrázek 12.1: Ukázka testu vyloučení narušitele pro článek idn-07792.

První metodou testování bylo „vyloučení nepřítele“ (anglicky „intruder detection“). Tato metoda se používá k evaluaci automaticky detekovaných shluků slov a je popsána například v [13]. 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 přiřazující 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 korpusu online článků stažených z českých webových serverů. Jedná se o články, 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 hobby a společenské články.

Pro naše testování jsme využili články ze serveru iDnes¹, 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 nenačítá. Tyto sady byly z testování vyřazeny a anotovaných testovacích sad je tedy pouze 58. Výsledky testování jsou shrnuty v tabulce 12.1 ve sloupci „o_test1“. Pochopitelné výsledky testování včetně Cohenovy kappy pro všechny dvojice anotátorů jsou v příloze 13.

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 správně. Pokud by algoritmus přiřazoval automaticky obrázky k textům i v praxi, čtenáři by byli schopni výstupy tohoto algoritmu odlišit od algoritmu, který přiřazuje obrázky k textům náhodně.

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

¹<http://www.idnes.cz>

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

Tabulka 12.1: Přehled výsledků uživatelského testování.

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 může zkreslovat výsledky testování touto metodou.

Dobře je tento problém vidět na ukázce anotace v Obrázku 12.1. Text článku je anketa mezi muži o ideální ženě. Obrázek ze sady fotek s tématem hubnutí by asi nebyl úplně ideální ilustrační obrázek k danému článku. Možná že náhodně přidaný obrázek s id „0007116667“ by článek ilustroval lépe. Oba anotátoři ho však označili za nejméně vhodný, protože ostatní obrázky jsou ze stejné série.

Zajímavé je analyzovat, proč obrázky s tematikou hubnutí přiřadil algoritmus jako vhodné k článku o ideální ženě. Algoritmus v článku označil jako klíčová postupně slova „ideální“, „žena“, „let“, „svobodný“ a „ráda“. První nalezený obrázek má id „0012612680“. Slovo „ideální“ se v českém překladu nachází dvakrát (z překladu frází „ideal figure“ a „ideal body measurements“). Dvakrát se v klíčových slovech obrázku nachází „žena“. Stejně tak slovo „let“ z anglického „years“, které není pro daný obrázek příliš relevantní. Slovo „svobodný“ se do českých metadat obrázku dostalo nesprávným překladem. Anglická metadata obrázku obsahují frázi „skin folds upper body freely“. Slovo „freely“ by v takovém kontextu mělo být přeloženo spíše jako „volný“ než jako „svobodný“. Obrázek s id „0012612680“ ukazuje velkou část problémů, které algoritmus vyhledávání ilustračních obrázků k textu má, zejména s přeloženými texty.

12.2 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 algoritmu 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 12.1 ve sloupci „o_test2“. Poborné výsledky

testování, včetně Cohenovy kappy pro všechny dvojice anotátorů je v příloze 13.

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.

Zajímavý je rozbor anotací, u kterých správný obrázek neuhodl ani jeden z obrázků. Některé nesprávně označené výsledky jsou způsobeny špatnou volbou náhodných obrázků. v anotaci textu „idn-01911“ ukázané na Obrázku 12.2 o novém softwarovém centru společnosti Microsoft se náhodně objevili 3 obrázky na kterých je počítač.



Obrázek 12.2: Obrázky v evaluaci textu idn-01911. Algoritmem doporučený je 2. obrázek v pořadí.

Podobný problém má anotace na Obrázku 12.2, kde se u článku o úbytku plochu polí a zvětšování plochy lesů objevily náhodně obrázky pole a lesa. Tyto problémy při testování algoritmu by šly odstranit lepším výběrem náhodných obrázků. Například by obrázky mohly být náhodně vybrány z množiny, která neobsahuje žádné z doporučených klíčových slov textu.



Obrázek 12.3: Obrázky v evaluaci textu idn-02505. Algoritmem doporučený je 4. obrázek v pořadí.

Někdy jsou špatné výsledky způsobeny chybným nastavením důležitosti jednotlivých klíčových slov. Obrázek 12.4 ukazuje obrázky u anotace článku o povýšení a kariéře žen. Algoritmus pro detekci klíčových slov správně určil jako první klíčové slovo „ženy“. Na dalších pozicích jsou „manažerských“, „muže“, „řídících“ a „zaměstnání“. Doporučený obrázek kromě klíčového slova „ženy“ obsahuje všechna ostatní detekovaná klíčová slova. Anotátoři však měli ve výsledcích obrázek se ženou, kterému dávali přednost.



Obrázek 12.4: Obrázky v evaluaci textu idn-04282. Algoritmem doporučený je 3. obrázek v pořadí.

Problémem se ukázaly být také příliš krátké popisky obrázků. Obrázek 12.5 ukazuje anotační obrázky pro text vracení záloh za elektřinu. Třetí obrázek má v klíčových slovech jediné slovo „backup“, které je přeloženo jako „záloha“. Vyhledávací algoritmus preferuje ve výsledcích obrázky s kratšími popisky klíčových slov, takže se tento obrázek s jednoslovným popisem stal doporučeným k danému textu. Zdá se z více případů, že tyto obrázky s velmi krátkými popisky by mohly být pro lepší výsledky z vyhledávání odstraněny.



Obrázek 12.5: Obrázky v evaluaci textu idn-04368. Algoritmem doporučený je 3. obrázek v pořadí.

12.3 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. Při testování se ukázalo, že některé náhodně vybrané obrázky, byly anotátory označeny jako lepší ilustrační obrázky pro daný text, než obrázky automaticky přiřazené algoritmem. To ovšem neznamená, že by tyto obrázky byly vždy pro daný test nevhodné. Ukazuje se, že některé druhy textů mohou obsahovat poměrně širokou škálu přijatelných ilustračních obrázků.

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. Další množina problémů je způsobena nesprávnými popisky u obrázků. Některé obrázky jsou popsány desítkami klíčových slov, jiné zase pouze jedním.

Jedním z problémů našeho testování je 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 Profimedie by rozšířilo i doménu textů, pro které algoritmus funguje dobře.

13. Závěr

Podářilo se nám implementovat webovou aplikaci pro vyhledávání ilustračních obrázků z textu. Obrázky pochází z datasetu Profimedia. Každý obrázek má přiřazena klíčová slova, která aplikace používá k vyhledávání. Podporuje český a anglický jazyk, ale je snadno rozšiřitelná na podporu jiných jazyků. Uživatel může nejen zadat text článku, ale může si i vynutit přímo klíčová slova u hledaných obrázků.

Klíčovým algoritmem pro vyhledávání je algoritmus extrakce klíčových slov. Ten je založen na modifikované metodě TF-IDF. Extrahovaná klíčová slova se používají k vyhledávání doporučených obrázků. Druhé využití vyextrahovaných klíčových slov z textu je jako nápověda uživateli.

Zdrojová data v datasetu Profimedia jsou v angličtině. Aby aplikace byla schopná pracovat s českými texty, bylo nutné klíčová slova obrázků přeložit. Použili jsme hybridní metodu, která kromě slovníkového překladu využívá i překlad detekovaných víceslovných frází v textu. Menší kvalita takového překladu způsobuje občasné problémy algoritmu při doporučování obrázků k českým textům. Při tak velkém množství dat by ovšem kvalitnější překlad byl velmi náročný na zdroje, nebo příliš drahý.

Součástí práce je i rozbor některých moderních možností tvorby frontendu a backendu webových aplikací. Práce vysvětluje některé novinky v HTML5 a ukazuje využití javascriptových frameworků pro různé účely. Na backendové straně aplikace ukazuje některé nové komunikační a datové protokoly. Rozebrány jsou některé možnosti uložení velkého množství prohledavatelných dat.

Kromě textového vyhledávání umožňuje aplikace i vyhledávání podle vizuální podobnosti obrázků. Uživatel aplikace vidí v detailu každého z nalezených obrázků množinu podobných obrázků. v rámci aplikace byla vytvořena nezávislá služba, která k obrázkům v datasetu Profimedia hledá obrázky. Podobnost obrázků je určena vzdáleností vektorů se 4 096 složkami. Jelikož obrázků v datasetu Profimedia je více než 20 milionů, bylo největším úkolem implementovat službu doporučování obrázků tak, aby vracela kvalitní výsledky co nejrychleji.

Algoritmus doporučování ilustračních obrázků byl testován s lidskými anotátory. Anotátoři dostávali české texty z online médií a množinu obrázků. Některé obrázky byly získané algoritmem na doporučení obrázků k textu. Jednu anotaci dostali nezávisle na sobě vždy dva anotátoři. v prvním úkolu dostali anotátoři jeden náhodný a čtyři doporučené obrázky. z 58 anotací označili anotátoři správně náhodně vybraný výsledek ve 45 případech. Ukázalo se ovšem, že metoda testování má některé problémy, které umožnily anotátorům vybrat náhodný obrázek aniž by ostatní obrázky byly k textu relevantní. Druhá metoda testování otočila poměr a anotace obsahovaly čtyři náhodné a jeden doporučený obrázek. Úkolem anotátorů bylo označit právě jeden nenáhodný obrázek. Každou anotaci opět dostali nezávisle na sobě vždy dva anotátoři. Ze 120 anotací se oba anotátoři shodli na správném obrázku v 76 případech. Rozbor anotací u kterých se oba anotátoři neshodli ukazuje různé příčiny. Některá klíčová slova k obrázkům jsou chybně přeložena. Některé obrázky obsahují špatná klíčová slova. Klíčová slova některých obrázků jsou příliš krátká. v několika anotacích by některé z náhodně vybraných obrázků mohly být ilustrační obrázky k textu také, což negativně

ovlivnilo výsledky testování v neprospěch algoritmu. pro anglické články testování neproběhlo, ale je pravděpodobné, že dosažené výsledky by byly lepší — u anglických klíčových slov nebylo potřeba použít automatický překlad.

Testování ukázalo, že aplikace ve vysokém procentu nabízí uživatelům přijatelné ilustrační obrázky. Pokud jsou výsledky nevyhovující, může uživatel využít nápovědu doporučených klíčových slov a hledaný výsledek snadno přesněji specifikovat. ve výsledku by tedy aplikace měla být schopna usnadnit proces vyhledávání ilustračních obrázků ke zpravodajským článkům, což bylo hlavním cílem práce.

Seznam použité literatury

- [1] ADOBE. *Adobe Flash Player* [software]. 2014. Dostupné z: <http://get.adobe.com/cz/flashplayer/>.
- [2] APACHE SOFTWARE FOUNDATION. *Apache Lucene 4.9.0* [software]. 2014. Dostupné z: <http://lucene.apache.org/>.
- [3] APERTIUM. *Apertium: A free/open-source machine translation platform* [software]. 2014. Dostupné z: <http://www.apertium.org/>.
- [4] ATTARDI, Giuseppe. *Wikipedia Extractor* [software]. 2013. Dostupné z: http://medialab.di.unipi.it/wiki/Wikipedia_Extractor.
- [5] ELASTICSEARCH BV. *elasticsearch-ruby* [software]. 2014. Dostupné z: <https://github.com/elasticsearch/elasticsearch-ruby>.
- [6] ELASTICSEARCH BV. *Elasticsearch 1.3.0* [software]. 2014. Dostupné z: <http://www.elasticsearch.org/>.
- [7] BATKO, Michal, et al. *Mufin: A multi-feature indexing network*. In: Proceedings of the 2009 Second International Workshop on Similarity Search and Applications. IEEE Computer Society, 2009, 158-159.
- [8] BOTOŘEK, Jan. *Tvorba nástroje pro zpracování textových popisů multimediálních 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/.
- [9] BRAY, Tim, et al. *Extensible markup language (XML)*. World Wide Web Consortium Recommendation REC-xml-19980210. <http://www.w3.org/TR/1998/REC-xml-19980210>, 1998.
- [10] 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, 130-142, ISBN 978-3-642-24468-1.
- [11] CASTELLANOS, Melissa. *Interpreters and Translators Among Most In-Demand and High-Paying Jobs of 2014* [online]. 2014 [cit. 2014-07-17]. Dostupné z: <http://bit.ly/1s6nvF9>.
- [12] CAVNAR, William B., et al. *N-gram-based text categorization*. In Proceedings of Third Annual Symposium on Document Analysis and Information Retrieval (SDAIR-94), 1994, 161-175.
- [13] CHANG, Jonathan, et al. *Reading tea leaves: How humans interpret topic models*. In: *Advances in neural information processing systems*. 2009, 288-296.
- [14] COUCHBASE. *Couchbase 2.5.1* [software]. 2014. Dostupné z: <http://www.couchbase.com/>.

- [15] CROCKFORD, Douglas. *The application/json media type for javascript object notation (json)*. 2006. Dostupné z: <http://tools.ietf.org/html/rfc4627.txt>.
- [16] DEAN, Jeffrey; GHEMAWAT, Sanjay. *MapReduce: simplified data processing on large clusters*. Communications of the ACM, 2008, 51(1), 107-113.
- [17] ECMAScript. *ECMAScript Programming Language*. 2014. Dostupné z: <http://www.ecmascript.org/>.
- [18] FETTE, Ian; MELNIKOV, Alexey. *The WebSocket Protocol*. 2011. Dostupné z: <http://tools.ietf.org/html/rfc6455>.
- [19] FIELDING, Roy. *Representational state transfer*. Architectural Styles and the Design of Network-based Software Architecture, 2000, 76-85.
- [20] GERRARD, Andrew. *Go version 1 is released* [online]. 2012 [cit. 2014-07-24]. Dostupné z: <http://blog.golang.org/go-version-1-is-released>.
- [21] GOOGLE INC. *Closure Tools* [software]. 2014. Dostupné z: <https://developers.google.com/closure/>.
- [22] GOOGLE INC. *The Go Programming Language* [software]. 2014. Dostupné z: <http://golang.org/>.
- [23] GOOGLE INC. *Překladač Google* [online]. 2014 [cit. 2014-07-25]. Dostupné z: <http://translate.google.cz>.
- [24] HARTER, Mark. *Managing Node.js Callback Hell with Promises, Generators and Other Approaches*. 2014. Dostupné z: <http://bit.ly/1uK01Ki>.
- [25] JIA, Yangqing. *Caffe: An open source convolutional architecture for fast feature embedding*. 2013. Dostupné z: <http://caffe.berkeleyvision.org>.
- [26] JONES, Karen Sparck. *A statistical interpretation of term specificity and its application in retrieval*. Journal of documentation, 1972, 28(1), 11-21.
- [27] THE JQUERY FOUNDATION. *jQuery* [software]. 2014. Dostupné z: <http://jquery.org/>.
- [28] KOEHN, Philipp. *Statistical machine translation*. Cambridge University Press, 2009.
- [29] 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, 177-180.
- [30] LUBBERS, Peter; GRECO, Frank. *HTML5 Web Sockets: A Quantum Leap in Scalability for the Web* [online]. 2013 [cit. 2014-07-25]. Dostupné z: <http://www.websocket.org/quantum.html>.
- [31] LOTT, Brian. *Survey of Keyword Extraction Techniques*. UNM Education, 2012.

- [32] MANNING, Christopher D.; RAGHAVAN, Prabhakar; SCHÜTZE, Hinrich. *Introduction to information retrieval*. Cambridge: Cambridge university press, 2008.
- [33] MATSUMOTO, Yukihiro, et al. *Ruby Programming Language* [software]. 2014. Dostupné z: <https://www.ruby-lang.org/>.
- [34] MICROSOFT. *Microsoft Silverlight* [software]. 2014. Dostupné z: <http://www.microsoft.com/silverlight/>.
- [35] MILLER, George A. *WordNet: A Lexical Database for English*. Communications of the ACM, 1995, 38(11), 39-41.
- [36] PORTER, Martin F. *An algorithm for suffix stripping*. Program: electronic library and information systems. MCB UP Ltd, 1980, 14(3), 130-137.
- [37] RAILS CORE TEAM. *Ruby on Rails* [software]. 2014 [cit. 2014-07-24]. Dostupné z: <http://rubyonrails.org/>.
- [38] STRAUCH, Christof; SITES, Ultra-Large Scale; KRIHA, Walter. *NoSQL databases*. Lecture Notes, Stuttgart Media University, 2011.
- [39] WIKIPEDIA *Geohash* — *Wikipedia, The Free Encyclopedia* [online]. 2014 [cit. 2014-07-24]. Dostupné z: <http://en.wikipedia.org/w/index.php?title=Geohash>.
- [40] WIKIPEDIA *Model-view-controller* — *Wikipedia, The Free Encyclopedia* [online]. 2014 [cit. 2014-07-28]. Dostupné z: <https://en.wikipedia.org/w/index.php?title=MVC>.
- [41] WIKIPEDIA *Wikipedia, The Free Encyclopedia* [online]. 2014 [cit. 2014-07-28]. Dostupné z: <https://www.wikipedia.org>.
- [42] THE WORLD WIDE WEB CONSORTIUM (W3C) *HTML5* [online]. 2014 [cit. 2014-07-28]. Dostupné z: <http://www.w3.org/TR/html5/>.
- [43] THE WORLD WIDE WEB CONSORTIUM (W3C) *CSS3* [online]. 2014 [cit. 2014-07-28]. Dostupné z: <http://www.w3.org/Style/CSS/current-work>.
- [44] ZAYTSEV, Juriy. *ECMAScript 5 compatibility table* [online]. 2014 [cit. 2014-07-28] Dostupné z: <http://kangax.github.io/compat-table/es5/>.

Seznam tabulek

| | | |
|------|--|----|
| 7.1 | Kategorie podobnosti obrázků podle <i>Similarity</i> | 18 |
| 12.1 | Přehled výsledků uživatelského testování. | 44 |

Příloha 1

Tato příloha obsahuje podrobná data z uživatelského testování. Obrázek 13.1 obsahuje statistiky správně oanotovaných obrázků pro oba testy. Navíc obsahuje seznam všech textů ve kterých anotátoři správný obrázek neuhádli spolu se jménem dotyčného anotátora.

Obrázek 13.2 obsahuje statistiky pro každého z anotátorů. Statistiky jsou jak celkové tak jednotlivé pro všechny dvojice anotátorů, kteří spolu anotovali alespoň jeden stejný text. pro každou dvojici anotátorů je uvedena matice 2×2 se statistikou shody a neshody anotátorů. z této matice je pak spočítáno Cohenovo kappa pro každou z dvojic anotátorů.

| o_test1 | o_test2 |
|---|---|
| nesprávné označení (id textu, uživatel) | nesprávné označení (id textu, uživatel) |
| 00207 o_paroubkova | 00124 o_semerad |
| 01337 o_pavlovic | 00405 o_havel |
| 01595 o_paroubkova | 00405 o_paroubkova |
| 01921 o_rakosova | 00623 o_paroubkova |
| 03758 o_odchazel_ondrej | 01907 o_rakosova |
| 03758 o_paroubkova | 01911 o_odchazel_vojtech |
| 03907 o_paroubkova | 01911 o_paroubkova |
| 04107 o_rakosova | 02505 o_odchazel_vojtech |
| 07604 o_havel | 02505 o_rakosova |
| 07604 o_odchazel_ondrej | 03610 o_havel |
| 09093 o_pavlovic | 03610 o_semerad |
| 09093 o_rakosova | 04211 o_odchazel_vojtech |
| 09551 o_havel | 04282 o_havel |
| celkem 116 anotací 58 textů | 04282 o_semerad |
| každý text oannotován dvěma anotátory | 04368 o_odchazel_vojtech |
| oba správně : 45 / 58 = 77.58 % | 04368 o_semerad |
| oba špatně : 3 / 58 = 5.172 % | 04480 o_semerad |
| jeden špatně: 10 / 58 = 17.24 % | 04530 o_havel |
| | 05102 o_rakosova |
| | 05102 o_semerad |
| | 05692 o_semerad |
| | 06379 o_paroubkova |
| | 06777 o_rakosova |
| | 07326 o_rakosova |
| | 07326 o_semerad |
| | 07587 o_odchazel_vojtech |
| | 07587 o_rakosova |
| | 08015 o_havel |
| | 08015 o_semerad |
| | 08621 o_havel |
| | 12931 o_odchazel_vojtech |
| | 13674 o_havel |
| | 13674 o_semerad |
| | 13803 o_havel |
| | 13803 o_paroubkova |
| | 14257 o_havel |
| | 14257 o_rakosova |
| | 14329 o_rakosova |
| | 14329 o_semerad |
| | 14353 o_havel |
| | 14353 o_rakosova |
| | 14988 o_semerad |
| | 15726 o_havel |
| | 15770 o_havel |
| | celkem 240 anotací 120 textů |
| | každý text oannotován dvěma anotátory |
| | oba správně : 76 / 120 = 63.33 % |
| | oba špatně : 15 / 120 = 12.5 % |
| | jeden špatně: 29 / 120 = 24.16 % |

Obrázek 13.1: Výsledky textů a id textů se špatně detekovanými obrázky.

```

havel
o_test1: 20 / 22 = 90.90 %
o_test2: 36 / 48 = 75.0 %
celkem: 56 / 70 = 80.0 %
odchazel_o
o_test1
3 0
1 1
Cohen's kappa = 0.5454545454545455
odchazel_v
o_test2
10 2
0 0
Cohen's kappa = 0.0
paroubkova
o_test1
6 0
0 0
Cohen's kappa = NaN
o_test2
6 1
3 2
Cohen's kappa = 0.2727272727272727
pavlovic
o_test1
5 0
0 0
Cohen's kappa = NaN
rakosova
o_test1
6 0
0 0
Cohen's kappa = NaN
o_test2
9 0
1 2
Cohen's kappa = 0.7499999999999999
semerad
o_test2
6 2
0 4
Cohen's kappa = 0.6666666666666667

odchazel_o
o_test1: 21 / 23 = 91.30 %
celkem: 21 / 23 = 91.30 %
havel
o_test1
3 1
0 1
Cohen's kappa = 0.5454545454545455
paroubkova
o_test1
4 1
0 1
Cohen's kappa = 0.5714285714285715
pavlovic
o_test1
6 0
0 0
Cohen's kappa = NaN
rakosova
o_test1
6 0
0 0
Cohen's kappa = NaN

odchazel_v
o_test2: 42 / 48 = 87.5 %
celkem: 42 / 48 = 87.5 %
havel
o_test2
10 0
2 0
Cohen's kappa = 0.0
paroubkova
o_test2
10 1
0 1
Cohen's kappa = 0.625
rakosova
o_test2
9 1
0 2
Cohen's kappa = 0.7499999999999999
semerad
o_test2
11 0
0 1
Cohen's kappa = 1.0

paroubkova
o_test1: 20 / 24 = 83.33 %
o_test2: 43 / 48 = 89.58 %
celkem: 63 / 72 = 87.5 %
havel
o_test1
6 0
0 0
Cohen's kappa = NaN
odchazel_o
o_test1
4 0
1 1
Cohen's kappa = 0.5714285714285715
odchazel_v
o_test2
10 0
1 1
Cohen's kappa = 0.625
pavlovic
o_test1
4 0
2 0
Cohen's kappa = 0.0
rakosova
o_test1
5 1
0 0
Cohen's kappa = 0.0
o_test2
11 1
0 0
Cohen's kappa = 0.0
semerad
o_test2
10 2
0 0
Cohen's kappa = 0.0

pavlovic
o_test1: 21 / 23 = 91.30 %
celkem: 21 / 23 = 91.30 %
havel
o_test1
5 0
0 0
Cohen's kappa = NaN
odchazel_o
o_test1
6 0
0 0
Cohen's kappa = NaN
paroubkova
o_test1
4 2
0 0
Cohen's kappa = 0.0
rakosova
o_test1
3 1
1 1
Cohen's kappa = 0.24999999999999986

rakosova
o_test1: 21 / 24 = 87.5 %
o_test2: 39 / 48 = 81.25 %
celkem: 60 / 72 = 83.33 %
havel
o_test1
6 0
0 0
Cohen's kappa = NaN
odchazel_v
o_test2
9 1
0 2
Cohen's kappa = 0.7499999999999999
odchazel_o
o_test1
6 0
0 0
Cohen's kappa = NaN
odchazel_v
o_test2
9 0
1 2
Cohen's kappa = 0.7499999999999999
paroubkova
o_test1
5 0
1 0
Cohen's kappa = 0.0
o_test2
11 0
1 0
Cohen's kappa = 0.0
pavlovic
o_test1
3 1
1 1
Cohen's kappa = 0.24999999999999986
semerad
o_test2
9 0
0 3
Cohen's kappa = 1.0

semerad
o_test2: 36 / 48 = 75.0 %
celkem: 36 / 48 = 75.0 %
havel
o_test2
6 0
2 4
Cohen's kappa = 0.6666666666666667
odchazel_v
o_test2
11 0
0 1
Cohen's kappa = 1.0
paroubkova
o_test2
10 0
2 0
Cohen's kappa = 0.0
rakosova
o_test2
9 0
0 3
Cohen's kappa = 1.0

```

Obrázek 13.2: Výsledky textů a id textů se špatně detekovanými obrázky.

Příloha 2

Tato příloha obsahuje informace o datech na přiloženém DVD. Přiložené DVD obsahuje dva soubory komprimované formátem ZIP — `zdrojove_kody.zip` a `demo.vdi.zip`.

13.1 Zdrojové kódy

Po rozbalení souboru `zdrojove_kody.zip` získáme čtyři adresáře obsahující zdrojové soubory. Adresář `czech-stemmer` je knihovna jazyka Ruby (Gem) pro stemmování českých slov popsané v Kapitole 5. Adresář `similar_img_finder` obsahuje zdrojové kódy služby pro vyhledávání podobných obrázků. Adresář `stock_photo_finder` obsahuje aplikaci se všemi pomocnými skripty. Adresář `cemi_annotace` obsahuje zdrojové kódy k webové aplikaci anotačního rozhraní popsaném v Kapitole 11.

13.2 Demo

Po rozbalení souboru `demo.vdi.zip` získáme soubor `diplomka.vdi`, což je virtuální stroj s operačním systémem Ubuntu 14.04 LTS¹. Tento stroj lze spustit například programem VirtualBox².

Na stroj se lze přihlásit s uživatelským jménem „odchazel“ a heslem „odchazel“. pro nastartování dema aplikace je nutné mít najednou spuštěné Bash skripty:

```
bash start_elasticsearch.bash #spusti databazi elasticsearch
```

```
bash start_app.bash #spusti webovou aplikaci
```

```
bash start_similar.bash #spusti sluzbu pro doporucovani podobnych  
obrazku
```

Nyni lze ve webovém prohlížeči otevřít aplikaci na adrese `http://localhost:3000`. Demo obsahuje pouze prvních 1000 obrázků z datasetu Profimedia. pro získání alespoň minimálních výsledků je dobré použít texty s klíčovými slovy „smile“, nebo „children“.

¹<http://www.ubuntu.com/>

²<https://www.virtualbox.org/wiki/Downloads>