



## ZADÁNÍ DIPLOMOVÉ PRÁCE

<b>Název:</b>	Metoda pro sumarizaci a hodnocení významnosti informací na Webu dat
<b>Student:</b>	Bc. Marek Filteš
<b>Vedoucí:</b>	Ing. Milan Doj inovski
<b>Studijní program:</b>	Informatika
<b>Studijní obor:</b>	Webové a softwarové inženýrství
<b>Katedra:</b>	Katedra softwarového inženýrství
<b>Platnost zadání:</b>	Do konce zimního semestru 2018/19

### Pokyny pro vypracování

Web dat (angl. Web of Data) obsahuje množství informací, které popisují charakteristiky entit z r zných domén. Je však obtížné posoudit, které z t chto informací jsou d ležit jší než ostatní. Hlavním cílem této práce je navrhnout a naimplementovat metodu pro ur ení d ležitosti informací, která bude založena na datech z Wikipedie a DBpedie. Dalším cílem je navrhnout a naimplementovat webovou aplikaci “prohlíže ” pro sumarizaci a prohlížení entit.

Pokyny:

- Seznamte se s existujícími metodami pro sumarizaci entit.
- Navrh te metodu pro sumarizaci informací o entitách založenou na datech z Wikipedie a DBpedie. Jako hlavní zdroj dat použijte dataset DBpedia abstracts a ve Wikipedia lánkách analyzujte výskyt d ležitých informací.
- Ve vybraném prost edí naimplementujte webovou aplikaci pro prohlížení a sumarizaci entit. Aplikace má zobrazit top-N nejd ležit jších informací pro danou entitu.
- Ov te funk nost aplikace na vybrané podmnožin entit z DBpedie a prove te evaluaci metody.

### Seznam odborné literatury

Dodá vedoucí práce.

Ing. Michal Valenta, Ph.D.  
vedoucí katedry

prof. Ing. Pavel Tvrdík, CSc.  
d kan

V Praze dne 7. dubna 2017



ČESKÉ VYSOKÉ UČENÍ TECHNICKÉ V PRAZE  
FAKULTA INFORMAČNÍCH TECHNOLOGIÍ  
KATEDRA SOFTWAREVÉHO INŽENÝRSTVÍ



Diplomová práce

## **Metóda pre sumarizáciu a hodnotenie významnosti informácií na Webe dát**

*Bc. Marek Filteš*

Vedúci práce: Ing. Milan Dojčinovski

8. mája 2017



---

## Podakovanie

Chcem sa podakovať vedúcemu tejto diplomovej práce Ing. Milanovi Dojčínovskému za jeho cenné poznatky, rady a pripomienky, ktorými mi bol nápomocný pri jej tvorbe.



---

## Prehlásenie

Prehlasujem, že som predloženú prácu vypracoval(a) samostatne a že som uviedol(uviedla) všetky informačné zdroje v súlade s Metodickým pokynom o etickej príprave vysokoškolských záverečných prác.

Beriem na vedomie, že sa na moju prácu vzťahujú práva a povinnosti vyplývajúce zo zákona č. 121/2000 Sb., autorského zákona, v znení neskorších predpisov, a skutočnosť, že České vysoké učení technické v Praze má právo na uzavrenie licenčnej zmluvy o použití tejto práce ako školského diela podľa § 60 odst. 1 autorského zákona.

V Prahe 8. mája 2017

.....

České vysoké učení technické v Praze

Fakulta informačních technologií

© 2017 Marek Filteš. Všechny práva vyhrazené.

*Táto práca vznikla ako školské dielo na FIT ČVUT v Prahe. Práca je chránená medzinárodnými predpismi a zmluvami o autorskom práve a právach súvisiacich s autorským právom. Na jej využitie, s výnimkou bezplatných zákonných licencií, je nutný súhlas autora.*

### **Odkaz na túto prácu**

Filteš, Marek. *Metóda pre sumarizáciu a hodnotenie významnosti informácií na Webe dát*. Diplomová práca. Praha: České vysoké učení technické v Praze, Fakulta informačních technologií, 2017.



---

## Abstrakt

Táto práca sa zaoberá sumarizáciu entít sémantického webu. Najprv sa rieši otázka informácií, hodnotenia miery dôležitosti informácií ako aj všeobecne sumarizácie entít. Prechádza sa k sumarizácii entít sémantického webu. V praktickej časti sa zaoberá návrhom modelu a implementáciou nástroja pre sumarizáciu entít na základe datasetu DBpedia abstracts. Vygenerovaná báza znalosti je integrovaná v rámci implementácie webového prehliadača.

**Kľúčová slova** sumarizácia entít, dôležitosť informácií, sémantický web, DBpedia

---

## Abstract

This work deals with entity summarization of semantic web. Firstly, the issue of information, evaluation the importance of information as well as the general summation of entities. It goes to entity summarization of semantic web entities. The practical part deals with design of the model and implementation of the entity summary tool based on the DBpedia abstracts dataset. The generated knowledge base is integrated within the implementation of the web browser.

**Keywords** entity summarization, importance of information, semantic web, DBpedia

---

# Obsah

Úvod	1
<b>1 Teoretická časť</b>	<b>3</b>
1.1 Informácia . . . . .	3
1.2 Relevatnosť informácií . . . . .	4
1.3 Sumarizácia entít . . . . .	7
1.4 Základné technológie sémantického webu . . . . .	10
1.5 Existujúce riešenia sumarizácie entít sémantického webu . . . . .	14
<b>2 Metóda pre sumarizáciu a hodnotenie významnosti informácií</b>	<b>19</b>
2.1 Zdroje dát pre sumarizáciu entít . . . . .	19
2.2 Analýza sumarizácie entít . . . . .	22
2.3 Návrh modelu sumarizácie entít . . . . .	24
2.4 Implementácia nástroja . . . . .	29
2.5 Výstup nástroja . . . . .	35
2.6 Použitie nástroja . . . . .	37
<b>3 Webová aplikácia</b>	<b>39</b>
3.1 Integrácia výsledkov sumarizácie entít . . . . .	39
3.2 Návrh prehliadača . . . . .	39
3.3 Implementácia prehliadača . . . . .	41
3.4 Architektúra prehliadača . . . . .	46
<b>4 Experimenty a vyhodnotenie</b>	<b>49</b>
Záver	57
Literatúra	59
A Obsah priloženého CD	63

<b>B</b>	<b>Doplňujúce informácie testovania</b>	<b>65</b>
<b>C</b>	<b>Ukážky webovej aplikácie</b>	<b>67</b>
<b>D</b>	<b>Ukážky bázy znalostí</b>	<b>69</b>

---

## Zoznam obrázkov

1.1	Vzťah medzi indikatívnou, tématickou a kritickou sumarizáciou [1]	8
1.2	Definícia PageRank [2]	10
1.3	RDF trojica	11
1.4	Názorná ukážka štruktúry SPARQL dotazu.	11
1.5	Príklad jednoduchej ontológie	13
1.6	Ukážka výsledku summaClient[3]	17
2.1	Ukážka obsahu nif-context datasetu.[4]	21
2.2	Ukážka obsahu nif-page-structure datasetu.[4]	21
2.3	Ukážka obsahu nif-text-links datasetu.[4]	21
2.4	Model sumarizácie entít s využitím NIF DBpedia abstraktu	26
2.5	Definícia čiastkového skóre predikátu $p$ triedy $T$	27
2.6	Definícia globálneho skóre predikátu $p$ triedy $T$	27
2.7	Definícia miery identickosti dvojice predikátov $a, b$	28
2.8	Definícia celkového čiastkového skóre predikátu $p$ triedy $T$	29
2.9	Definícia výsledného skóre predikátu $p$ triedy $T$	29
2.10	Základný diagram tried nástroja	31
2.11	SPARQL dotaz pre získanie zdrojov na základe DBpedia PageRank	34
2.12	SPARQL dotaz pre získanie predikátov zdroja na základe extrahovaného prepojenia	34
2.13	SPARQL dotaz pre získanie všetkých predikátov na základe extrahovaného prepojenia	34
2.14	SPARQL dotaz pre získanie počtu použítí predikátu triedy.	35
2.15	SPARQL dotaz pre získanie predikátov globálnej štatistiky triedy	35
2.16	Štruktúra súboru s globálnou štatistikou	36
2.17	Štruktúra súboru s identickými predikátmi	36
2.18	Štruktúra súboru s čiastkovým výsledkom zdroja	37
2.19	Štruktúra súboru bázy znalostí	37
2.20	Spôsob inštalácie nástroja	37

3.1	Návrh dizajnu webovej aplikácie . . . . .	41
3.2	Základný diagram tried servera . . . . .	42
3.3	Základná štruktúra JSON požiadavku na server . . . . .	42
3.4	SPARQL dotaz pre získanie tried zdroja . . . . .	43
3.5	SPARQL dotaz pre získanie vlastnosti zdroja . . . . .	43
3.6	Základná štruktúra komponent klienta . . . . .	45
3.7	Základna architektúra webového prehliadača . . . . .	46
3.8	Základna komunikácia klienta so serverom . . . . .	47
4.1	Definícia miery identickosti pozície predikátu $p$ . . . . .	49
4.2	Graf porovnania pozícií BWD voči GKB pre triedu Film . . . . .	50
4.3	Graf porovnania miery dôležitosti BWD voči GKB pre triedu Film	51
4.4	Graf porovnania pozícií BWD voči UBES pre triedu Film . . . . .	51
4.5	Graf porovnania miery dôležitosti BWD voči UBES pre triedu Film	52
4.6	Graf porovnania pozícií BWD voči SUMMA pre triedu Film . . . . .	53
4.7	Graf porovnania miery dôležitosti BWD voči SUMMA pre triedu Film . . . . .	54
4.8	Graf porovnania pozícií BWD voči SUMMA pre triedu Person . . . . .	54
4.9	Graf porovnania miery dôležitosti BWD voči SUMMA pre triedu Person . . . . .	55
4.10	Graf porovnania pozícií BWD voči SUMMA pre triedu President . . . . .	55
4.11	Graf porovnania miery dôležitosti BWD voči SUMMA pre triedu President . . . . .	56
C.1	Ukážka zobrazenia zdroja triedy Person . . . . .	67
C.2	Ukážka zobrazenia zdroja triedy President . . . . .	68
C.3	Ukážka zobrazenia bázy znalosti . . . . .	68

---

## Zoznam tabuliek

1.1	Aspekty modelu PSP/IQ [5] . . . . .	5
1.2	Dimenzie kvality informácie [5] . . . . .	6
1.3	Tabuľka výsledkov dôležitosti vlastnosti filmov [6] . . . . .	16
B.1	Tabuľka mapovania Freebase predikatov na DBpedia predikáty . .	66
D.1	Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Plant . .	70
D.2	Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Animal .	71
D.3	Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Person .	72
D.4	Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:President	73
D.5	Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Monarch	74
D.6	Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Artist .	75
D.7	Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Film . .	76
D.8	Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Country	77
D.9	Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Company	78
D.10	Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Sport . .	79
D.11	Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Work . .	80
D.12	Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Eukaryote	81





---

# Úvod

Web dát (angl. Web of Data) obsahuje veľké množstvo informácií, popisujúce charakteristiky entít z rôznych domén. Je zložité vyhodnotiť, ktoré informácie sú dôležitejšie než ostatné. Problém dôležitosti informácií, do určitej miery, limituje použitie zdrojov prepojených dát (angl. Linked Data) uložených v databázach ako sú DBpedia, Freebase, Wikidata a ďalšie. DBpedia ako jedná z najdlhšie existujúcich, aktuálne vo verzií 2016-04, obsahuje informácie pre približne 6 miliónov zdrojov (angl. resource) a 1,3 miliárd základných faktov (RDF tripletov)[7]. Každý zdroj preto obsahuje v priemere približne 216 RDF trojíc [7]. Hľadanie informácií človekom v takýchto zdrojoch, si vyžaduje veľké množstvo času. Tento stav nie je vhodný. Vo všeobecnosti má užívateľ pri hľadaní niekoľko možných cieľov. Napríklad chce rýchlo nájsť zdroj, ktorý obsahuje informáciu, ktorú potrebuje. Zároveň potrebuje preskočiť zdroje, ktoré pre jeho hľadanie nie sú relevantné. Medzi ďalší možný cieľ patrí zoskupenie zdrojov na základe ich kľúčových informácií. V rámci týchto a iných možných cieľov hľadania, je bez stanovenia dôležitosti informáciám, prehľadávanie zdrojov nefektívne. Preto sa v tejto práci budeme venovať návrhu metódy pre satňovanie dôležitosti informáciám a ich prehľadaniu.

Hlavným cieľom tejto práce je poskytnúť používateľovi  $n$  najdôležitejších informácií jednotlivých zdrojov DBpedia na základe príslušnosti do určitej skupiny entít. Tieto informácie chceme získať z neštrukturovanej časti textu Wikipédie, tzv. abstraktu. Zo získaných vzťahov chceme vytvoriť bázu znalosti. Naše výsledky budú pre používateľa prehľadné pomocou webovej aplikácie, ktorá bude používať vytvorenú bázu znalosti.

Preto sa budeme v tejto práci venovať riešeniu problému, nazývaného sumarizácia entít (angl. entity summarization). V rámci sumarizácie entít budeme identifikovať mieru dôležitosti jednotlivých informácií entít. Pokúsime sa na základe existujúcich informácií vybudovať bázu znalosti, pomocou ktorej budeme môcť inetpretovať dôležité informácie. V prevej, teoretickej, časti práce sa zoznámime s existujúcimi riešeniami sumarizácie entít. Popíšeme základne pojmy, ktoré potrebujeme pochopiť. Identifikujeme hlavné technológie,

potrebné pre sumarizáciu entít, ktoré budeme používať. V časti návrhu si popíšeme proces návrhu metódy. Na základe metódy vytvoríme spustiteľný nástroj. Pomocou nástroja si vytvoríme bázu znalosti, ktorú použijeme pri tvorbe webového prehliadača. V časti porovnania a vyhodnotenia, získane výsledky miery dôležitosti porovnáme voči výsledkom iných riešení. V závere zhodnotíme dosiahnuté výsledky práce, možné použitie tohto riešenia, nedostatky a prípadne vylepšenia.

Výsledkom práce sa stane nástroj pre sumarizáciu entít a identifikáciu dôležitých informácií, báza znalosti ako aj webová aplikácia pre prehliadanie dôležitých informácií zdrojov DBpedia.

---

# Teoretická časť

Predtým ako sa dostaneme k hlavnej časti práce si musíme rozobrať potrebné predpoklady pre našu sumarizáciu entít. Rozoberme teoretické vedomosti. Ako stanoviť mieru dôležitosti informácií. Pozrime sa na to čo znamená sumarizácia entít. Čo v rámci jej riešenia musíme pochopiť. S akými technológiami budeme musieť pracovať. Zároveň si prezrime aktuálne dostupné riešenia sumarizácie entít a stanovenia miery dôležitosti informáciám.

## 1.1 Informácia

Informácie sú hnacou silou spoločnosti. Informácie sa už od dávnych dôb predávali medzi ľuďmi v rámci vzájomnej komunikácií. Vďaka výmene informácií sa ľudstvo vyvinulo k technologicky vyspelej civilizácii. Pozrime sa na pojem informácia v širšom a taktiež v užšom zmysle, potrebnom pre našu prácu.

Pojem informácia je veľmi široký a ťažko definovateľný. Predná definícia pojmu informácia neexistuje. Slovo informácia má pôvod v latinčine, kde znamená „uvádzať do tvaru“. Uvedme si aspoň dve základne prístupy chápania informácie. Existuje samozrejme viacero možností chápania informácií.

- Informácia ako oznámenie. Ide o komunikačný prostriedok, ktorý má význam z hľadiska príjemcu.
- Informácia ako správa. Súčasť technického prenosu.

Potrebu skúmať informácie aj z hľadiska technického, prinieslo až využívanie elektrickej energie. Z počiatku sa však vedci venovali viac problému prenosu, než otázke čo znamená informácia. Až publikácia "*Přenos informací*" od R. Hartleyho z roku 1928, priniesla jednu zo základných myšlienok technického chápania informácie. V základe popisuje prenos správy ako postupnosti symbolov. Väčšina neskorších prác a prístupov technického chápania informácií je založená na tomto princípe.

Keďže informácia, v technickom chápaní, vznikla abstrakciou, nemá hmotnú podobu. Vždy je ale spojená, s určitým fyzikálnym pochodom. Napríklad signálom. Vyjadrením informácie, ako postupnosti znakov, je správa. Správa môže byť textom, číslom a pod. Všetky správy majú svoju skladbu (syntaxi), obsah (sémantiku) a dôležitosť (pragmatickosť):

- Sémantický obsah vyjadruje význam správy. Správy s rovnakým množstvom informácií je možné zapísať v rôznych jazykoch.
- Pragmatický obsah vyjadruje dôležitosť (významnosť, užitočnosť) a prioritu pre príjemcu. Je zrejmé, že správa "Peter je doma." má inú dôležitosť pre jeho otca a inú pre osobu, ktorá ho nepozná.
- Syntaktický obsah je vzájomné usporiadanie znakov. Syntaktický obsah nie je závislý na existencii sémantického, či pragmatického obsahu. Sémantický obsah správa získava smerom k objektu a pragmatický smerom k príjemcovi.

Syntaktické usporiadanie znakov patrí pod kvantitatívnu vlastnosť informácie. Sémantický a pragmatický obsah spolu vytvárajú kvalitatívnu vlastnosť informácie. [8]

Je podstatné rozlíšiť tri pojmy: *dáta*, *informácia* a *znalosť*. Dáta predstavujú rozpoznateľné, vzájomne oddelené znaky. Tie zvyčajne bývajú uložené v pamäti. Systém, ktorý ich dokáže akceptovať a spracovať, je im schopný priradiť určitý význam. V takom prípade je ich možné, v rámci daného systému, považovať za informácie. Pokiaľ majú dané informácie pre systém organizačný význam, je možné ich nazvať znalosťami. [9]

### 1.2 Relevatnosť informácií

Najskôr sa pozrime ako popisujú problém relevantnosti autori prác o medziludskej komunikácii. Podľa Jany Hoffmannovej sa človek snaží predávať adresátovi len relevantné informácie. Príjemcovia venujú svoju pozornosť, len takým informáciám, ktoré pokladajú za podstatné a v daný moment použiteľné. V rámci vnímania stále očakávame prospešné informácie. Ide o očakávanie, ktoré slúži ako kritérium, podľa ktorého vyberáme z množiny alternatívnych interpretácií, tú ktorá je najpodstatnejšia a najzávažnejšia. [10] V tomto tvrdení sa síce jedná o medziludskú komunikáciu, avšak podobné princípy hodnotenia platia všeobecne pri vyhľadávaní informácií, pretože informácie majú hlavný význam pre človeka. Človek má svoje očakávania, na základe ktorých hodnotí relevantnosť získaných informácií.

Informácie môžu mať určitý význam jedine v definovanom kontexte. Podľa Roberta Stalnakera, je kontext jednoducho definovateľný ako určitý celok dostupných informácií. Tieto informácie sú dostupné všetkým členom konverzácie. Kontextová množina predstavuje možné situácie, ktoré sú kompatibilné

	Splnenie špecifikácií	Naplnenie očakávaní používateľa
Kvalita produktu	Spolahlivé informácie: charakteristiky dodaných informácií spĺňajú štandard informačnej kvality.	Užitočné informácie: charakteristiky dodaných informácií spĺňajú používateľove potreby.
Kvalita služby	Dôveryhodné informácie: transformácia dát na informácie spĺňa štandardy.	Použiteľné informácie: transformácia dát na informácie spĺňa používateľove potreby.

Tabuľka 1.1: Aspekty modelu PSP/IQ [5]

s dostupným celkom informácií. Na základe toho môžeme reprezentovať kontext ako množinu určitých svetov. Tieto svety predstavujú otvorené možnosti vzhľadom k tomu, čo je pevne dané. [11]

Mieru relevancie novej informácie pre jednotlivca musíme hodnotiť pomocou, určitého zlepšenia, ktoré je výsledkom získania novej informácie v danom kontexte. Ľudia vymedzujú mieru relevancie, spôsobom, kde je myšlienkový predpoklad relevantný v kontexte, len ak má kontextuálny účinok, na tento kontext. Psychické procesy, zahrňajúce úsilie pre dosiahnutie relevancie predstavuje negatívny efekt, vzhľadom ku efektu dosiahnutej relevancie. Myšlienkový predpoklad je pre človeka relevantný, v danej dobe, ak je relevantný aspoň v jednom kontexte. [12]

### 1.2.1 Relevatnosť informácií sémantického webu

Problém dôležitosti informácií sémantického webu (viac o sémantickom webe nájdete v kapitole 1.4.1) je zložitejší a je možné ho riešiť z rôznych hľadísk. Pokúsime sa postupne prepracovať k čo možno najvhodnejšiemu riešeniu, ako chápať relevatnosť informácií v prostredí sémantického webu.

Beverly K. Kahn a kolektív v práci “Information quality benchmarks: product and service performance” popisujú model hodnotenia informačnej kvality. Kvalitu informácií popisujú ako vhodnosť (angl. fitness) pre používateľa. Existujú však aj ďalšie konkrétnejšie pohľady. Patrí sem napríklad: dokonalosť, splnenie špecifikácie alebo naplnenie očakávaní používateľa. Vo svojom modeli PSP/IQ (product and service performance model for information quality) pristupujú ako k splneniu špecifikácie a naplnenie očakávaní používateľa. [5] Základné dimenzie modelu sú zobrazené v tabuľke 1.1.

Pomocou modelu PSP/IQ poskytujú popis dimenzií kvality informácií. Tento popis je nápomocný k lepšiemu pochopeniu požiadaviek pre poskytovanie vysoko kvalitných informácií. Tento popis je v tabuľke 1.2. Tieto dimenzie poukazujú na vysokú mieru komplexnosti problému relevancie informácií.

Relevancia informácií v prostredí sémantického webu môže byť ďalej určená pomocou hodnotenia kvality dát. Kvalita dát so zameraním na kvalitu datového produktu. Takáto kvalita môže byť posudzovaná v rôznej úrovni granularity [13]:

## 1. TEORETICKÁ ČASŤ

Dimenzia	Popisuje
Dostupnosť	Do akej miery sú informácie k dispozícii alebo ľahko a rýchlo vyhľadateľné.
Primerané množstvo informácií	Do akej miery je objem informácií vhodný pre danú úlohu.
Dôveryhodnosť	Do akej miery sú informácie považované za pravdivé a dôveryhodné.
Úplnosť	Do akej miery informácie nechýbajú a majú dostatočnú šírku a hĺbku pre danú úlohu.
Stručné zastúpenie	Do akej miery sú informácie zjednodušené reprezentované.
Dôsledné zastúpenie	Do akej miery sú informácie prezentované v rovnakom formáte.
Jednoduchá manipulácia	Do akej miery sa s informáciami ľahko manipuluje a uplatňujú na rôzne úlohy.
Bezchybovosť	Do akej miery sú informácie správne a spoľahlivé.
Interpretovateľnosť	Do akej miery sú definície jasné a informácie v príslušných jazykoch, symboloch a jednotkách.
Objektívnosť	Do akej miery sú informácie nezaujaté, bezprecedentné a nestranné.
Relevancia	Do akej miery sú informácie použiteľné a užitočné pre danú úlohu.
Reputácia	Do akej miery je informácia odporúčaná, pokiaľ ide o jeho zdroje alebo obsah.
Bezpečnosť	Do akej miery je prístup k informáciám primerane obmedzený, aby sa zachovala jeho bezpečnosť.
Včasnosť	Do akej miery sú informácie dostatočne aktuálne pre danú úlohu.
Pochopiteľnosť	Do akej miery sú informácie ľahko pochopiteľné.
Pridaná hodnota	Do akej miery sú informácie prospešné a poskytujú výhody z ich využívania.

Tabuľka 1.2: Dimenzie kvality informácie [5]

- Úroveň objektu (inštancie). Hodnotením relevancie inštancií tried v RDF databáze.
- Úroveň dokumentu. Pomocou analýzy obsahu a štruktúry prepojení.
- Úroveň podgrafov. Pomocou analýzy obsahu založenej na ontológií v kontexte výsledkov dotazu a následne výpočtom dôveryhodnosti na základe kontextu.

Pojem relevancia je často chápaný s trocha odlišným významom. Viacerí autori ho preto rozoberajú vo svojich prácach a snažia sa ho zadefinovať. Ide o komplexný a zároveň multidimenzionálny pojem. V oblasti informatiky sa objavuje určitý súlad v jeho chápaní. Kagolovsky Y. uvádza, že relevancia je vo všeobecnosti rozdelená do dvoch hlavných kategórií:

- Tématická relevancia. Je objektívna a týka sa terminológie.
- Používateľský orientovaná relevancia. Je subjektívna, závisí na potrebách používateľa.

Z hľadiska úrovne, medzi ktorými informácie prechádzajú pri ich získavaní, sa relevancia dá riešiť na viacerých z nich. Nižšie úrovne riešia interakcie s informačným systémom. Naproti tomu vyššie úrovne zabezpečujú interakcie používateľa. V rámci vyšších úrovni existujú ďalšie aspekty. Patrí sem:

- Kognitívny aspekt
- Situačný aspekt
- Afektívny aspekt
- Kontextový aspekt

[14][15] Kognitívny aspekt relevancie hovorí o zložitosti a množstve biochemických procesov potrebných k dosiahnutiu pochopenia informácií používateľom. V prípade, že rozoznáva málo relevantné informácie je táto kognitívna záťaž vyššia. Vyššia kognitívna záťaž predstavuje negatívny stav pri vyhľadávaní. [16] Afektívny aspekt relevancie je vzťahom medzi zámerom, cieľom a motiváciou používateľa a zdroja. Situačný aspekt relevancie predstavuje vzťah medzi situáciou, úlohou alebo problémom a zdrojom. [14] Kontextový aspekt bol spomenutý vyššie. Hovorí o množine svetov v rámci dostupného celku informácií.

R. Lachica s kolektívom ďalej, za účelom jasnejšieho chápania problému relevance, popisujú definície postulátov pre pojmy:

**Kvalita** Odzrkadľuje vnútornú hodnotu informácie. V prípade, že sú informácie nespoľahlivé alebo nepochopiteľné sú pokladané za nekvalitné. Nadalej môžu byť označené za relevantné alebo dôležité.

**Relevancia** Odráža hodnotu informácie z potrieb používateľa. Ide o vnímanie používateľom.

**Dôležitosť** Odzrkadľuje hodnotu informácie zo širšieho hľadiska. Z pohľadu ľudí, ktorý majú poznatky o tom čo môže používateľ potrebovať. Tieto poznatky nemusí mať samotný používateľ.

[14]

### 1.3 Sumarizácia entít

Sumarizácia entít predstavuje aktuálne rozsiahlu oblasť výskumu. Ide o všeobecný problém analýzy obsahu entít. Keďže väčšina informácií je reprezentovaná pomocou textov, je najčastejšou sumarizáciou, sumarizácia textov. Dôvodom, prečo sa tomuto problému začala venovať pozornosť, je nárast informovanosti. Ide predovšetkým o veľké množstvo informácií zdieľaných pomocou webových technológií. Preto je sumarizácia entít naobširnejšie spracovaná z pohľadu sumarizácie dokumentov na webe. Pozrime sa na základy sumarizácie entít na webe ako aj na konkrétnejší problém sumarizácie entít sémantického webu.

Sumarizácia dokumentov webu môže byť chápaná, ako text ktorý vznikol z jedného alebo viacerých textov. Poskytuje podstatné informácie z pôvodných textov. Nemal by však byť dlhší než je polovica celkového textu. Hlavným cieľom sumarizácie je prezentácia hlavných ideí v dokumente vo výrazne menšom rozsahu. Ak by mali všetky elementárne časti textu rovnakou mieru dôležitosti,

## 1. TEORETICKÁ ČASŤ

---

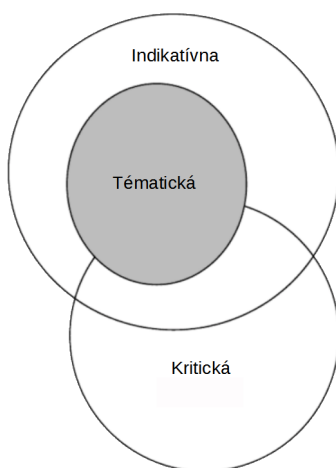
nebola by sumarizácia efektívna. V praxi je však väčšina textov rozdelená do určitých logických častí, pričom je vždy možné posúdiť ich významnosť. Výsledok sumarizácie je znovupoužitím časti textu (vety, slovné spojenia, slová a ďalšie), prípadne ich modifikovanou podobou. Existuje niekoľko typov sumarizácií, ktoré je možné definovať z rôznych nezávislých hľadísk:

Z hľadiska formy reprezentácie výsledku existuje extrakčná a abstrakčná sumarizácia.

- Extrakčná sumarizácia. Vo výsledku sumarizácie sú priamo vety z pôvodného textu.
- Abstrakčná sumarizácia. Výsledok nemusí obsahovať vety z pôvodného textu. Môže ich parafrázovať, zovšeobecňovať prípadne konkretizovať. Vyžaduje pokročilejšie metódy spracovania.

Z hľadiska účelu existuje indikatívna, tématická a kritická sumarizácia entít. Ich vzajomný vzťah je znázornený na obrázku 1.1.

- Indikatívna sumarizácia (angl. indicative summary) identifikuje o čom text informuje. Poskytuje skrátenú verziu obsahu.
- Tématická sumarizácia (angl. topic-oriented summary) sa zamierava na požadované témy čitateľa. Takáto sumarizácia odzrkadľuje pohľad čitateľa.
- Kritická sumarizácia poskytuje hlavné zistenia systematického prieskumu. Výsledkom je posudok o kvalite prieskumu a platnosti výkladu poznatkov.



Obr. 1.1: Vzťah medzi indikatívnou, tématickou a kritickou sumarizáciou [1]



Samotný proces sumarizácie môže obsahovať rôzne čiastkové procesy. Vo všeobecnosti je možné tieto procesy zaradiť do 4 základných skupín:

1. Extrakcia - identifikácia podstných častí textu.
2. Abstrakcia - preformulovanie časti na nové pojmy
3. Fúzia - kombinovanie extrahovaných častí
4. Kompresia - eliminovanie nepotrebných častí

[17] [1]

Od kedy bola publikovaná prvá metóda automatickej sumarizácie textu, vzniklo množstvo ďalších. Sumarizačné metódy môžeme rozdeliť do troch skupín:

1. Klasické metódy:
  - Heuristické metódy. Predstavujú najstaršie sumarizačné metódy. Jednou z nich bola metóda frekvencie termov vytvorená z roku 1958. Jej hlavnou myšlienkou je výber najčastejšie sa vyskytujúcich slov textu.
  - Štatistické metódy. Medzi typickú štatistickú sumarizáciu patrí Bayesova klasifikácia kde sa klasifikátor snaží o každej vete dokumentu rozhodnúť či patrí do výsledku, alebo nepatrí.
2. Metódy so spracovaním prirodzeného jazyka:
  - Metódy využívajúce súvislosti v texte. Patrí sem napr. metóda lexikálnych reťazcov.
  - Metódy modifikujúce pôvodný text.
3. Ďalšie matematické metódy:
  - Grafové metódy. Tieto metódy sú založené na algoritme PageRank, prípadne jeho modifikáciách TextRank a LexRank.
  - Algebrické metódy. Patrí sem napr. metóda latentnej sémantickej analýzy (LSA), ktorej princípom je rozklad matice na singulárne hodnoty.

[18]

Medzi najčastejšie metódy hodnotenia poradia dokumentov na webe patrí PageRank. PageRank vivinul Larry Page na Stanfordovej univerzite. PageRank patril ku kľúčovým hodnoteniam výsledov vo vyhľadávači Google. Základnou myšlienkou hodnotenia dokumentu je počet iných dokumentov, ktoré na neho odkazujú. Týmto smerom mu pridávajú dôležitosť. Na druhej strane

je jeho dôležitosť zníženia počtom vlastných odkazov na iné dokumenty. Definícia PageRanku: nech  $u$  je webová stránka a nech  $F_u$  je množina stránok  $u$ . Ďalej nech  $B_u$  je množina stránok odkazujúcich na  $u$  a  $N_u = |F_u|$  je množina odkazov z stránky  $u$  [2]. Potom:

$$R(u) = c \sum_{v \in B(u)} \frac{R(v)}{N(v)}$$

Obr. 1.2: Definícia PageRank [2]

Jedným z modelov, ktorý implementuje PageRank je model náhodného surfera. Surfista sa pohybuje medzi jednotlivými uzlami a s rovnakou pravdepodobnosťou vyberá, ktorú hranu na prechod použije. Zároveň s malou pravdepodobnosťou, skočí do náhodného uzla. Poradie uzlov sa získava pomocou stacionárneho rozloženia Markovského reťazca. Uzol s vyššou pravdepodobnosťou dosahu surfera je hodnotený vyšším skóre. Týmto spôsobom sa predpokladá, že najvyššie zaradené uzly zachytia hlavné témy pôvodných dát. [2]

### 1.4 Základné technológie sémantického webu

V tejto časti sa zoznámime s podstatnými technológiami pre realizáciu našej metódy sumarizácie entít. Pohybujeme sa vo webových technológiách. Presnejšie technológiách spojených s oblasťou sémantického webu.

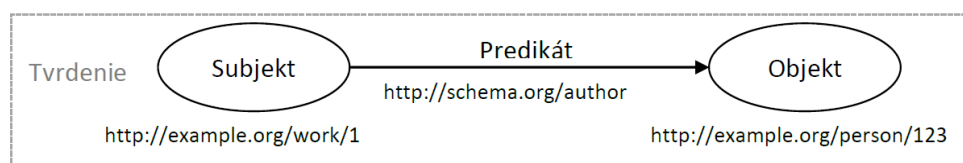
#### 1.4.1 Sémantický web

Sémantika je náuka o význame jednotlivých slov, morfému a iných znakov, prípadne tiež ich vzťahu ku skutočnosti, ktorú označujú. Slovo sémantika má pôvod v gréckom slove "séma", ktoré znamená význam alebo nosič významu. [19] V roku 2001 vyslovil Tim Berns-Lee myšlienku sémantického webu. Upozornil na skutočnosť, že súčasný internet je len chaotická zmes webových stránok a domén, ktorá začína rapídne narastať. Zároveň je stále zložitejšie nájsť správne relevantné informácie. Základom sémantického webu je popis informácií pomocou vzťahov medzi entitami. Sémantický web je založený na technológií Resource Description Framework (RDF). [20]

#### 1.4.2 Resource Description Framework

Resource Description Framework (RDF) je obecným frameworkom určeným pre popis, výmenu a znovu použitie metadat. Poskytuje jednoduchý model zápisu a zároveň nie je závislý na konkrétnej implementácii. Informáciu o objekte poskytuje tzv. tvrdenie (angl. statement). Základom tvrdenia je syntaktické vytvorenie trojice (angl. triples) medzi subjektom a objektom pomocou pedikátu (angl. predicate). Subjekt predstavuje zdroj, ktorému chceme priradiť

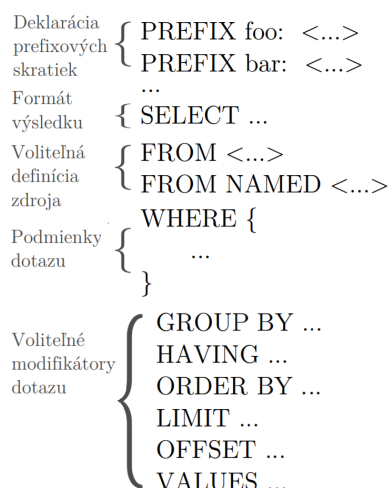
určitú vlastnosť. Predikát je priradením vlastnosti subjektu. Vlastnosť nadobúda hodnotu definovanú objektom. Týmto spôsobom je možné vyjadriť väčšinu potrebných informácií. Priradená hodnota vlastnosti môže predstavovať ďalší zdroj. Ak sa nejedná o ďalší zdroj, takýto objekt sa nazýva literál (angl. literal). Najlepšie ako si to je možné predstaviť, je znázornenie pomocou grafovej štruktúry. Subjekt spolu s objektom tvoria uzli grafu. Predikát je orientovanou úsečkou spájajúcou subjekt s objektom. Príklad je znázornený na obrázku 1.3. Dané tvrdenie udáva informáciu, že dielo identifikované pomocou `http://example.org/work/1` má autora (podľa slovníka `http://schema.org`) identifikovaného `http://example.org/person/123`. [21]



Obr. 1.3: RDF trojica

### 1.4.3 SPARQL

SPARQL predstavuje štruktúrovaný dotazovací jazyk, určený pre prácu s RDF uložiškami (tzv. triplestore). Ide o jednoducho použiteľný dotazovací jazyk. Existujú celkovo 4 základne konštrukcie dotazov, a to pomocou príkazov SELECT, ASK, CONSTRUCT a DESCRIBE. Názorná štruktúra dotazu je na obrázku 1.4



Obr. 1.4: Názorná ukážka štruktúry SPARQL dotazu.

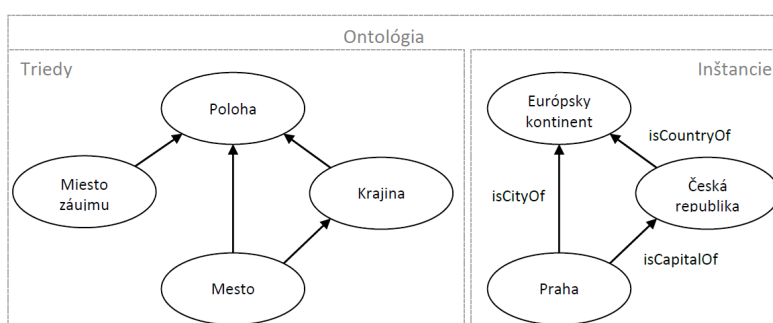
Na začiatku SPARQL dotazu je možné deklarovať skrátené prefixy, ktoré zastupujú v zvyšku dotazu uvedenú URI. Po nich, nasleduje definícia formátu výsledku pomocou kľúčových slov, ktoré už boli spomenuté vyššie. Kľúčové slová FROM poskytujú možnosť definovať obmedzenia na zdrojové grafy, z ktorých budú dáta pri tvorbe odpovede čerpané. Časť WHERE obsahuje podmienkové telo dotazu. V podmienkach sa používajú premenné, ktoré musia začínať symbolom "?". Pri tvorbe odpovede sa v nadom RDF zdroji dát na tieto premenné mapujú dostupné hodnoty, pričom musia byť splnené všetky podmienky z tela WHERE. Takéto hodnoty sa potom použijú pre tvorbu odpovede na základe požadovaného formátu výsledku. V rámci tvorby podmienok platí, že konjunkcia podmienok je vytvorená bodkov ".". OPTIONAL operátor slúži pre pripojenie časti (pravej) podmienky k predošlej (ľavej) podmienke, tak aby to nepoškodilo vyhodnotenie splnenia predošej časti. UNION operátor slúži k zlučeniu časti podmienok. Operátor MINUS odstráni výsledky získané podmienkov na pravej časti. Špeciálny operátor FILTER umožňuje redukovať výsledok len na tie ktoré spĺňajú danú podmienku. Tento operátor podporuje viaceré funkcie, ako sú logické, textové, číselné a pod. Okrem iného taktiež umožňuje filtrovanie trojích, ktorých objekty sú literálmi. [22]

### 1.4.4 Ontológia

Ontológia bola prebratá z filozofických vied, predstavovala náuku o bytí. Vo vzťahu k sémantickému webu sú ontológie chápané ako definície pojmov a vzťahov medzi nimi. Účelom je popis domény určitého ľudského záujmu. Takáto oblasť obsahuje jednotlivé triedy, ktoré sú prepojené vzťahmi. Objekty obsiahnuté v tejto doméne a ich prepojenie ontológia popisuje pomocou štyroch prvkov: entita, trieda, atribút a väzba. U niektorých definícií sa uvádza taktiež piaty prvok udalosť. Najznámejšou definíciou pojmu ontológia je definícia podľa T. Grubera: „*ontológia je explicitná špecifikácia konceptualizácie*“. Cieľom ontológie je definovať spoločné chápanie určitých tried a pojmov. Mala by podporovať porozumenie medzi viacerými ľuďmi, z rôznych profesií, komunikáciu počítačových systémov a taktiež zjednodušiť návrh znalostne orientovaných aplikácií. [23] Príklad jednoduchej ontológie je na obrázku 1.5

### 1.4.5 Linked Data

Pojem linked data uviedol v roku 2006 Tim Berners-Lee. Linked data definuje spôsob publikovania štruktúrovaných dát. Tieto dáta sú ďalej medzi sebou prepojené, čím navzájom získavajú pridanú hodnotu. Medzi základne technológie, ktoré sú pri linked data použité patria štandardné webové technológie: HTML, RDF a URI. Vďaka tomu sa rozširuje využitie publikovaných dát, strojovo čitateľných a ďalej zobraziteľných v rôznych podobách. Hlavným účelom linked data je sprístupnenie dát v jednoducho spracovateľnej forme. [24] Linked data používajú súbor overených postupov, pre publikovanie a prepoje-



Obr. 1.5: Príklad jednoduchej ontológie

nie štruktúrovaných dát na webe. Postupy boli prijaté v priebehu niekoľkých rokov, narastajúcim počtom poskytovateľov dát. To viedlo k vytvoreniu celosvetového dátového priestoru, pozostávajúceho z viac ako miliardy dát. Vďaka dodržiavaniu štandardov a overených postupov je umožnená interoperabilita dát a tzv. znovu použitie dát na webe. [25]

#### 1.4.6 Wikipedia

Cieľom Wikipédie je vytvorenie voľne dostupnej a spoľahlivej encyklopédie. Wikipédiu založil Jimmy Wales a charakterizuje ju ako *"snahu o vytvorenie a distribuovanie voľne dostupnej encyklopédie každej osobe na planéte v jej rodnom jazyku, a to v čo možno najlepšej kvalite."* Je založená a základe softvéru wiki. Wikipedia je systematicky tvorená kolaboratívnou činnosťou svojich užívateľov. Predstavuje jeden z najväčších znalostných zdrojov. Anglická verzia obsahuje viac než 5 miliónov, článkov. Články Wikipédie obsahujú predovšetkým voľný text. Obsahuje taktiež niekoľko rôznych typov štruktúrovaných informácií:

- informačné boxy
- kategorizačné informácie
- multimédia
- geologické informácie
- odkazy na interné aj externé zdroje

[26]

#### 1.4.7 DBpedia

DBpediu založili vedci z Free University of Berlin a The University of Leipzig v spolupráci s OpenLink Software. Prvá verzia bola vydaná v roku 2007. DBpediu je možné definovať viacerými spôsobmi. Je možné sa na tento projekt dívať ako na komunitu ľudí, venujúcich sa získavaniu informácií a znalosti

z Wikipédie. Alebo ako na službu, pomocou ktorej používateľa vedia vyhľadávať informácie z Wikipédie oveľa sofistikovanejším spôsobom. DBpedia obsahuje sadu nástrojov určených pre extrakciu štruktúrovaných informácií z Wikipédie. Tieto informácie tvoria bázu prepojených a strojovo čitateľných dát. DBpedia umožňuje jednoduchým spôsobom získať odpovede na otázky, ktoré nie je jednoduché získať pomocou textového hľadania. Napríklad získať zoznam všetkých účastníkov zimných olympijských hier, ktorý sa narodili po roku 1980.[27, 28]

### 1.4.7.1 Extrakcia informácií z Wikipédie

Extrahované sú predovšetkým informácie z informačných boxov. Informačné boxy Wikipédie obsahujú základne fakty, zvyčajne v tabuľkovom formáte. Tabuľka priradzuje atribútom ich hodnoty. Extrahovanie informácií z informačných boxov prebieha pomocou dvoch hlavných paralelných metód:

1. Generické extrahovanie informačných boxov. Cieľom generickej extrakcie je veľké pokrytie. URI zdroja DBpedia je tvorená z URL článku Wikipédie. Táto URI je postupne použitá ako subjekt pre všetky extrahované triplety. Napríklad pre článok <https://en.wikipedia.org/wiki/Mazda> je jeho URL použitá pre vytvorenie zdroja <http://dbpedia.org/resource/Mazda>. URI predikátov sú vytvorené zretazením menného priestoru <http://dbpedia.org/property/> a názvu atribútu informačného boxu. Napríklad atribút Revenue má potom vo výsledku URI predikátu <http://dbpedia.org/property/revenue>. Objekty sú tvorené z priradených hodnôt atribútov. Hodnoty sú predspracované a transformované na podporované reprezentácie hodnôt v RDF. Ide zvyčajne o objekty typu literál.
2. Mapovacie extrahovanie informačných boxov. Cieľom mapovania je lepšia kvalita dát. Rieši sa problém synonymných názvov atribútov a viacerých šablón, ktoré sú použité pre rovnaké typy. Šablóny Wikipédie sú mapované na ontológiu DBpedia. Ontológia je tvorená manuálne, usporiadaním najviac použitých šablón informačných boxov do hierarchie menších premís. Z šablón sa mapujú atribúty na vlastnosti ontológie. RDF trojice tejto metódy obsahujú predikát so základom <http://dbpedia.org/ontology/>.

## 1.5 Existujúce riešenia sumarizácie entít sémantického webu

Medzi podstatnú časť analýzy patrí oboznámenie sa s existujúcimi prácami. Relevantnými pre nás sú tie, ktoré riešia problém sumarizácie entít, stanovenia

dôležitosti informáciám, alebo aspoň niektorú jeho časť. Preto sa v nasledujúcej časti tejto práce pozrieme na niektoré aktuálne dostupné riešenia sumarizácie a identifikácie dôležitostí informácií entít v prostredí sémantického webu.

### 1.5.1 Dôležitosť získaných informácií pomocou fuzzy sémantických sietí

V práci sa R. Lachica s kolektívom venoval hodnoteniu zdrojov použitím sémantickej vzdialenosti. Vzdialenosť určili pomocou prechodov medzi jednotlivými asociáciami množiny tém. Celková vzdialenosť je určená od témy používateľa. Asociácie pokladajú za dôležité. Súlužia ako vstupné body do sémantického webu. Konečný výsledok hodnotenia je založený na sémantickej vzdialenosti a kvalite hodnotenia zdroja. Asociácie majú rôzne váhy. Pri ceste do podkategórie triedy sa váha znižuje výraznejšie.

V prípade konextu špecifickej témy sa poradie počíta podľa všetkých výstupných ciest z počiatočnej témy. Váha relevancie je znižovaná každým prechodom. Všetky témy nad určitou prahovou hodnotou sú ohodnotené ako relevantné. Tento princíp vyplýva z techniky CSA (The Constrained Spreading Activation). Pomocou druhého prechodu je zlepšene umiestnenie zdrojov, ktoré sú spojené s témou používateľa. U ostatných kontextových tém sa hodnotenia automaticky prispôbia. Získaným riešeniam je pomocou sumarizácie priradené skóre relevantnosti a zároveň priradené skóre kvality zdroja riešení.

V prípade vyhľadávania pomocou kľúčových slov, sa každé kľúčové slovo porovnáva s témami z množiny tém. Syntaktické zúženie termínov sa používa na vyhľadávanie zodpovedajúcich tém. Zväčšenie sémantického množstva termínov sa vykoná taktiež použitím CSA. Vďaka tomu sa získa zoznam subjektov, ktoré sa zhodujú s kľúčovými slovami z dotazu. Proces pokračuje získaním zdrojov. Poradie zdrojov sa stanoví pomocou najkratšej cesty. Cesty sú vypočítané pomocou algoritmu pre hľadanie najkratšej cesty v neorientovanom váženom grafe. V práci odkazujú na použitie Bellman-Fordovho algoritmu. Algoritmus slúži pre získanie najkratších ciest z počiatočného uzla do všetkých ostatných uzlov grafu. Ak neobsahuje cyklus zápornej dĺžky. Ak by ho graf obsahoval, algoritmus ho dokáže detekovať. [14]

### 1.5.2 Sumarizácia entít pomocou vedomostnej hry

Vo svojej práci sa v roku 2012 A. Thalhammer a ďalší venovali sumarizácií entít na základe kvízovej hry. Kvíz bol zameraný na jednu konkrétnu doménu. Doménu filmov. Kvíz podľa predpokladu mal poskytnúť kvalitnejšiu sumarizáciu než boli automaticky generované sumarizácie. Kvíz obsahoval filmy, ktorých RDF grafy boli získane z Freebase.<sup>1</sup> Dataset obsahoval veľke množstvo

---

<sup>1</sup>Freebase API bola súčasťou služieb poknýkaných spoločnosťou Google. Aktuálne je Freebase API ukončená. Presmerovali podporu znalostí z Freebase na službu Wikidata.

Pozícia	Freebase predikát	Pozícia	Freebase predikát
1	prequel	14	production company
2	film series	15	runtime
3	sequel	16	music
4	parodied	17	award
5	adapted original	18	actor
6	subject	19	story writer
7	genre	20	editor
8	initial release date	21	event
9	director	22	cinematographer
10	rating	23	budget
11	writer	24	film festival
12	featured song	25	film casting director
13	featured filming location		

Tabuľka 1.3: Tabuľka výsledkov dôležitosti vlastnosti filmov [6]

otvorených dát, čo bolo pre ich výber datsetu kľúčové. Medzi ďalšie zvažované alternatívy patrili DBpedia a Linked Movie Database (LinkedMDB).

Pre kvíz si zvolili 60 najlepšie hodnotených filmov na základe rebríčka IMDb. Princípom kvízu bolo zodpovedať otázku, výberom správnej odpovede. Implementovali 2 druhy otázok. S jedinou správnou odpoveďou a s  $n$  správnymi odpoveďmi. Jednotlivé otázky boli zostavené na základe existujúcich tripletov v datasete. Príkladom pre vytvorenie otázky s jedinou správnou odpoveďou je použitie tripletu:

*fb:en.pulp\_fiction test:hasActor fb:en.john\_travolta .*

Nesprávne možnosti boli generované na základe nájdenia subjektov, ktoré majú priradený rovnaký predikát ako triplet otázky (*test:hasActor*). Zároveň objekt tripletu otázky nie je v žiadnom triplete nesprávneho subjektu. Otázka znela: "John Travolta je hercom vo filme:". Otázka skúma mieru dôležitosti daného predikátu. Každá správna odpoveď pridáva skóre hráčovi. Výsledky hry boli spracované v zmysle percentuálnej úspešnosti zodpovedania otázky každého predikátu. Prehľadný výsledok ich metódy je v tabuľke 1.3. [6]

### 1.5.3 API rozhranie pre sumarizáciu prepojených dát

V roku 2015, A. Thalhammer a S. Stadtmüller riešili návrh a implementáciu Web API pre sumarizáciu entít. Účelom práce bolo poskytnutie nástroja, ktorý dokáže spracovať sumarizáciu entít, dodaných klientom z rôznych zdrojov. Sumarizácia je založená na hodnotení prichádzajúcich odkazov z Wikipédie. Výsledkom ich implementácie sa stala webová aplikácia summaServer a klientska knižnica summaClient.



V rámci implementácie riešili niekoľko problémov. Ako popisujú, znalostné bázy zvyčajne obsahujú podporu viacerých jazykov pre označenie zdrojov. Pridaním podpory jazykov je možné znížiť zbytočné požiadavky na iné zdroje alebo znalostné bázy. Pre potreby pokrytia širšieho aspektu vlastnosti, využili princíp  $n$ -árnych vzťahov. Takto umožnili definovať maximálny počet prechodov k získaniu cieľového objektu. Na základe zadanej množiny predikátov je možné získať prepojenia, ktoré sú zamerané len na určitú oblasť. Vo výsledku je hodnotenie priradené skupine tvrdení (angl. statement). Uvádzajú, že rozdiel medzi RDF trojicou, ktorá obsahuje prepojenie na ďalší zdroj a literálom je pri tvorbe sumarizácie v účelnosti využitia. Literály poskytujú užitočné informácie o špecifickej entite naproti tomu prechod na iný zdroj umožňuje preskúmanie vzťahov danej entity. Pre každé RDF prepojenie typu  $a \text{ link } b$  je možné vytvoriť obrátené prepojenie  $b \text{ linkBy } a$ . Uvádzanie týchto obrátených prepojení vo výsledku sumarizácie má zmysel, pretože zahŕňa informácie o danom zdroji. Výsledná sumarizácia obsahuje skóre založené na slovníku vRank. [3]

### Pulp Fiction

film director	<a href="#">Quentin Tarantino</a>
distributor	<a href="#">Miramax</a>
producer	<a href="#">Lawrence Bender</a>
editing	<a href="#">Sally Menke</a>
cinematography	<a href="#">Andrzej Sekula</a>
series of	<a href="#">Mia Wallace</a>
Wikipage disambiguates of	<a href="#">Vincent</a>
Wikipage disambiguates of	<a href="#">Pulp fiction</a>
Wikipage disambiguates of	<a href="#">Jules</a>
Wikipage disambiguates of	<a href="#">Coolidge</a>

Obr. 1.6: Ukážka výsledku summaClient[3]

### 1.5.4 Centralizovaná sumarizácia založená na príbuznosti a informovanosti

V práci z roku 2011 od G. Cheng, T. Tran a Y. Qu je riešený návrh a implementácia varianty metódy náhodného surfera. Využíva vzťahy popisov prvkov pre hodnotenie. Model je grafovo orientovaným datovým modelom, ktorý spĺňa RDF štandard.

Problém sumarizácie chápu ako výber  $n$  najlepšie hodnotených vlastností zdroja. Centralizovane orientované hodnotenie pokladajú za úspešne v rámci sumarizácie textov a ontológií. Preto tento spôsob hodnotenia aplikovali pri návrhu ich modelu. Centralizovane orientované hodnotenie si vyžaduje konštrukciu gafu, kde uzli zodpovedajú prvkom, ktoré majú byť hodnotené. Každý pár súvisiacich uzlov je spojený pomocou neorientovaných alebo orientovaných hrán. Vo výsledku sú uzly zoradené podľa ich centrality v grafe. Najčastejšie

vypočítané pomocou PageRank algoritmu. Pri použití náhodného surfera identifikovali dva problémy. Pojem centrálnosti môže byť príliš všeobecný. Identifikácia hlavných tém pôvodného popisu entity nebol jediným cieľom. Hľadala sa sumarizácia, ktorá najlepšie charakterizuje entitu a pomáha rozlišovať entitu od ostatných. To znamená, že pri meraní centrality by sa malo zväžiť, koľko informácií prináša určitá vlastnosť, ktorá môže prispieť k identifikácii entity. Druhý problém predstavuje hodnotenie susedných uzlov uzla za rovnako dôležité pri výbere cesty. Tento výber ma pre každý susedný uzol rovnakú pravdepodobnosť. Model neposkytuje stupeň príbuznosti na jemnozrnnejšej úrovni. Táto nepresnosť môže viesť k neoptimálnym výsledkom.

Vytvorili model podobný ako PageRank, ktorý simuluje správanie náhodného surfera pomocou dvoch druhov akcií. Surfer je cieľovo orientovaný. Prechádza množinou vlastností v snahe popísať požadovanú entitu. Surfer vykoná buď relačný pohyb s vyššou mierou pravdepodobnosti do uzla, ktorý obsahuje súvisiace informácie o téme, ktorá sa sumarizuje. Druhou akciou je informačný skok, pri ktorom s vyššou mierou pravdepodobnosti preskočí do uzla, ktorý poskytuje väčšie množstvo nových informácií. Tieto voľby sú reprezentované dvomi nerovnomerne rozdelenými pravdepodobnosťami, pričom jedna je daná súvislosťou medzi uzlami a druhá informatívnosťou uzlov. [29]

### 1.5.5 Zhodnotenie

V skratke sa pokúsime zhodnotiť získané informácie na základe prieskumu existujúcich riešení. Spôsob hodnotenia získaných informácií na základe fuzzy sématických sietí je viac zameraný na hodnotenie výsledkov, než na samotný proces ich získavania. Poskytuje návod pre stanovenie dôležitosti výsledku pomocou grafovej štruktúry. Pre našu prácu je viac informačný. Podáva však informácie o dôležitosti sumarizácie skupín entít. Sumarizácia entít pomocou kvízu, predstavuje zaujímavý nápad, ktorý zároveň dokáže poskytuje kvalitné informácie. Tie sú však závislé priamo na interakcií používateľa. Jej ďalšou nevýhodou je šírka pokrytia domén. Pre potreby vytvorenia bázy znalosti všetkých domén by si metóda vyžadovala veľké množstvo respondentov. Prístup k priradeniu miery dôležitosti predikátu pomocou jeho úspešnosti je možné v pozmenenej forme použiť aj v rámci nášho riešenia. Sumarizácia pomocou verejného API pre ľubovoľný RDF zdroj poskytuje vysokú mieru pokrytia. Pomocou neho je možné automatické spracovanie ľubovoľných zdrojov. Vo výsledkoch sa však objavujú duplicitné predikáty čo znižuje prehľadnosť dôležitých informácií. Táto skutočnosť poukazuje na to, že je vhodné pokúsiť sa tieto duplicity redukovať. Spôsob sumarizácie existujúcich zdrojov vieme ďalej použiť, ako určitú podpornú časť procesu našej sumarizácie. Práca zameraná na centralizovanú sumarizáciu na základe príbuznosti a informovanosti poskytuje potrebné informácie pre analýzu vzťahov pri prechode grafom. Pokladá vhodný návod ako odlíšiť akcie pohybu grafom v rámci sumarizácie.

# Metóda pre sumarizáciu a hodnotenie významnosti informácií

V úvode sa pozrieme na všetky technológie s ktorými budeme pracovať. Ďalej si predstavme základnú myšlienku a návrh nášho modelu. V implementačnej časti sa pozrieme na proces vytvorenia nástroja, ktorý je založený na našom modeli. V poslednej časti tejto kapitoly si popíšeme spôsob práce s nástrojom a uvedieme príklady ako sme vygenerovali našu bázu znalostí.

## 2.1 Zdroje dát pre sumarizáciu entít

V našom modeli budeme používať ako hlavný zdroj informácií pre sumarizáciu dataset NIF DBpedia abstracts. Budeme uvažovať dataset vo formáte NLP Interchange Format (NIF). Pre vyhľadanie informácií o zdrojoch na DBpedii použijeme DBpedia SPARQL endpoint. Preto sa v úvode návrhu modelu zoznámime s vyššie uvedenými technológiami.

### 2.1.1 DBpedia abstracts

Wikipedia je najdôležitejším a najkomplexnejším zdrojom otvorených, encyklopedických poznatkov. Obsahuje články, ktoré predstavujú obrovský zdroj otvoreného textu v prirodzenom jazyku. Ako sme už uviedli vyššie, extrakcia informácií do DBpedie spočíva v mapovaní informačných tabuliek, šablón a ďalších ľahko identifikovaných štruktúrovaných údajov. Textový popis a údaje, ktoré môže obsahovať, nie sú súčasťou extakcie, hoci predstavujú najväčšiu časť dokumentu. Z celého popisu je extrahovaná len prvá kapitola. Táto kapitola je celá extrahovaná a predaná do DBpedie pod označením *abstract*. Exisuje taktiež skrátená verzia tohto abstraktu. Prvý odstavec abstraktu je mapovaný ako hodnota predikátu *rdfs:comment* [30] Pomocou sumarizácie

## 2. METÓDA PRE SUMARIZÁCIU A HODNOTENIE VÝZNAMNOSTI INFORMÁCIÍ

---

založenej na použití abstraktov preto dokážeme pokryť takmer všetky zdroje jednotlivých entít. V aktuálnej verzii DBpedia 2016-04 obsahuje približne 77% zdrojov s predikátom <http://dbpedia.org/ontology/abstract>. [7]

### 2.1.2 NLP Interchange Format

Formát NLP Interchange Format (NIF) je založený na RDF/OWL formáte. Jeho účelom je dosiahnutie interoperability medzi nástrojmi na spracovanie prírodného jazyka (NLP), jazykovými zdrojmi a anotáciami. Hlavnou časťou NIF je slovník, ktorý umožňuje reprezentovať referencie ako zdroje RDF. Pomocou špeciálneho dizajnu URI je umožnené presné určenie anotácií časti dokumentu. Tieto identifikátory URI je možné použiť na pripojenie ľubovoľných anotácií k príslušnej sekvencii znakov. Použitím týchto identifikátorov URI môžu byť anotácie zverejnené na webe ako prepojené údaje. NIF sa skladá z štrukturálnej, koncepcnej a prístupovej interoperability. Štrukturálna interoperability znamená, že URI sa používajú na pripojenie anotácií v dokumentoch pomocou fragmentu identifikátorov. Koncepcná predstavuje vetne štruktúrovanú ontológiu (SSO), ktorá bola špeciálne vyvinutá pre napojenie existujúcich textových ontológií. Pomocou toho je možné pripojiť bežné anotácie k URI. Prístupová predstavuje popis REST rozhrania. URI NIF-1.0 je konštruovaný z URI zdroja. Jeho URI je použitá ako prefix. Vo všeobecnosti je to jednoduchý a spoľahlivý prístup ako anotovať webové zdroje existujúcich zdrojov. Pre NIF URI je doporučené použitie ukončenia prefixu pomocou lomky "/", mriežky "#" prípadne dotazovou komponentou (napr. ?Nif-id=123). [31] [32]

### 2.1.3 NIF DBpedia Abstracts

DBpedia sa v súčasnosti primárne zameriava na zastupovanie faktických poznatkov, ktoré sú obsiahnuté v informačných boxoch Wikipédie. Za účelom možnosti získania štruktúrovaných údajov z neštruktúrovaných textov, DBpedia poskytuje všetky informácie<sup>2</sup> extrahované zo zdrojového kódu HTML. Tieto informácie sú rozdelené do 3 datsetov [4]:

- nif-context - obsahuje text stránky ako kontext (vrátane indexov začiatku a konca). Príklad obsahuje je na obrázku 2.1
- nif-page-structure - obsahuje štruktúru stránky v sekciách a odsekoch (názvy, kapitoly a iné). Príklad obsahuje je na obrázku 2.2
- nif-text-links - obsahuje všetky odkazy na iné zdroje DBpedia, ako aj externé odkazy. Príklad obsahuje je na obrázku 2.3

Príklady obsahu jednotlivých datsetov, pre prehľadnosť sú vynechané typové označenia literálov:

---

<sup>2</sup>V rámci prvej publikácie (verzia 2016-04, z 12.10.2016) boli zverejnené len texty abstraktov. Až od verzie 2016-10 budú zverejnené všetky texty vo formáte NIF. [4]

## 2.1. Zdroje dát pre sumarizáciu entít

```
dbr:Anthropology?dbpv=2016-04&nif=context a nif:#Context ;
nif:isString "Anthropology is the study of humanity...";
nif:beginIndex "0" ;
nif:endIndex "634" ;
nif:sourceUrl <http://en.wikipedia.org/wiki/Anthropology> ;
nif:predLang <http://lexvo.org/id/iso639-3/eng> .
```

Obr. 2.1: Ukážka obsahu nif-context datasetu.[4]

```
dbr:Anthropology?dbpv=2016-04&nif=context nif:hasSection dbr:Anthropology?
dbpv=2016-04&nif=section_0_634 .

dbr:Anthropology?dbpv=2016-04&nif=section_0_634 a nif:Section ;
nif:beginIndex "0" ;
nif:endIndex "634" ;
nif:referenceContext dbr:Anthropology?dbpv=2016-04&nif=context ;
nif:hasParagraph dbr:Anthropology?dbpv=2016-04&nif=paragraph_0_330 ;
nif:hasParagraph dbr:Anthropology?dbpv=2016-04&nif=paragraph_331_634 ;
nif:firstParagraph dbr:Anthropology?dbpv=2016-04&nif=paragraph_0_330 ;
nif:lastParagraph dbr:Anthropology?dbpv=2016-04&nif=paragraph_331_63 .

dbr:Anthropology?dbpv=2016-04&nif=paragraph_0_330 a nif:Paragraph ;
nif:beginIndex "0" ;
nif:endIndex "330" ;
nif:referenceContext dbr:Anthropology?dbpv=2016-04&nif=context ;
nif:superString dbr:Anthropology?dbpv=2016-04&nif=section_0_634 .

dbr:Anthropology?dbpv=2016-04&nif=paragraph_331_634 a nif:Paragraph ;
nif:beginIndex "331" ;
nif:endIndex "634" ;
nif:referenceContext dbr:Anthropology?dbpv=2016-04&nif=context ;
nif:superString dbr:Anthropology?dbpv=2016-04&nif=section_0_634 .
```

Obr. 2.2: Ukážka obsahu nif-page-structure datasetu.[4]

```
dbr:Anthropology?dbpv=2016-04&nif=word_29_37 a nif:Word ;
nif:referenceContext dbr:Anthropology?dbpv=2016-04&nif=context ;
nif:beginIndex "29" ;
nif:endIndex "37" ;
nif:superString dbr:Anthropology?dbpv=2016-04&nif=paragraph_0_634 ;
<http://www.w3.org/2005/11/its/rdf#taIdentRef> dbr:Human ;
nif:anchorOf "humanity" .

dbr:Anthropology?dbpv=2016-04&nif=phrase_65_84 a nif:Phrase ;
nif:referenceContext dbr:Anthropology?dbpv=2016-04&nif=context ;
nif:beginIndex "65" ;
nif:endIndex "84" ;
nif:superString dbr:Anthropology?dbpv=2016-04&nif=paragraph_0_634 ;
<http://www.w3.org/2005/11/its/rdf#taIdentRef> dbr:Social_anthropology ;
nif:anchorOf "social anthropology" .
```

Obr. 2.3: Ukážka obsahu nif-text-links datasetu.[4]

### 2.1.4 DBpedia SPARQL endpoint

Okrem webovej reprezentácie jednotlivých zdrojov existuje ďalší verejne dostupný spôsob pre získavanie informácií z DBpedia. Predstavuje ho SPARQL endpoint, ktorý je dostupný na adrese <http://dbpedia.org/sparql>. SPA-

RQL endpoint je poskytovaný pomocou služby OpenLink Virtuoso. [33] Tento endpoint je vhodný, pre použitie v prípade, že používateľ, alebo klientska aplikácia vie vopred, aké informácie potrebuje. Okrem štandardného protokolu SPARQL endpoint podporuje niekoľko rozšírení. Vyhľadávanie v celkom texte na vybraných RDF predikátoch. Agregáčnych funkcií, ako napríklad *COUNT*. V rámci zabezpečenia má implementované obmedzenia pre určité dotazy. Napríklad dotaz, ktorý požaduje získanie celého obsahu úložiska. Ten je zamietnutý ako príliš náročný. Výsledky dotazov SELECT sú redukované na maximálne 1000 záznamov. [34]

### 2.2 Analýza sumarizácie entít

V našom návrhu sa pokúsime navrhnúť prístup pre sumarizáciu entít sémantického webu, ktorý bude možné použiť ako model pre ďalšie použitie. Vychádzajúc z poznatkov o údajoch, ktoré sú obsiahnuté na Wikipédii, sa zameriame na neštruktúrované textové informácie. V tomto prístupe nám bude nápomocný dataset DBpedia abstrakt.

Hlavnou myšlienkou prečo využiť informácie obsiahnuté v textovom popise je tá, že obsahujú pravdepodobne značne množstvo užitočných informácií o daných zdrojoch. Vyplyva to okrem iného, z množstva času, ktorý bol k vytvoreniu popisov venovaný ľuďmi. Ďalším dôležitým faktom je usporiadanie samotného textu. Každý článok na Wikipédii je štruktúrovaný do určitej miery podobne ako nejaká výstupná práca. Týka sa to v hlavnej miere úvodu, ten predstavuje abstrakt, ktorý poskytuje skrátený popis postatných informácií daného článku. Napríklad ak sa pozrieme na článok popisujúci mesto Praha, bude v prvej časti textu určite uvedená informácia, že Praha je hlavným mestom Českej republiky. Rovnako to bude pri iných hlavných mestách. Medzi ďalšie skutočnosti patrí to, že sú jednotlivé články medzi sebou prepojené odkazmi. To znamená, že Praha v abstrakte obsahuje odkaz na Českú republiku. Ak si to zovšeobecníme na všetky články, je možné, že pomocou analýzy vzťahov, ktoré sú v abstraktoch odhalíme ďalšie podstatné informácie, ktoré doposiaľ nie sú extrahované. Ak sa pozrieme na typ vzťahov, ktoré môžu existovať v rámci odkazovania sa medzi článkami, môžeme identifikovať dve typy. Prvým je výstupne prepojenie zo zdrojového článku. Príkladom je už spomenutý vzťah článku Praha k článku Česká republika. Druhým typom je vzťah opačný. Sumarizácia zdrojov má zmysel na úrovni typovej príslušnosti k určitej triede, definovanej pomocou ontológie. Výsledné hodnotenie bude v takom prípade priradené daným triedam ontológie.

Na základe identifikovaných tried ontológie bude potrebné zvoliť pre každú triedu dostatočný počet zdrojov. Ešte pred voľbou počtu zástupcov tried, je potrebné zvoliť spôsob výberu vhodných zdrojov. Je nutné definovať vhodnosť kandidátov. V tomto štádiu existuje viacero možností ako k výberu pristupovať. Uvedme si apoň niekoľko uvažovaných metód:

- Výber na základe najdlhšieho abstraktu.
- Výber na základe DBpedia PageRank
- Výber na základe najvyššieho počtu unikátnych predikátov

Varianty výberu na základe dĺžky abstraktu a počtu predikátov sú závislé na kompletnom prehľadaní zdrojov. Ide o optimalizačné vyhľadávanie. Navyše pri použití počtu predikátov je možné, že sa vhodnejšie zdroje dostanú na nižšie miesta. Zdroje totiž môžu obsahovať viacero "šumových" predikátov. Zároveň automatická identifikácia unikátnosti, predstavuje smotný veľký problém. Existuje viacero slovníkov, ktoré popisujú identické hodnoty pomocou odlišných predikátov. Výber na základe dĺžky abstraktu vyzerá síce vhodným kandidátom, ale ako sme už uviedli vyžaduje si kompletné prehľadanie zdrojov. Na druhej strane existuje vypracované hodnotenie DBpedia zdrojov pomocou PageRanku. Tvorbe PageRanku sa venovali A. Thalhammer a A. Rettinger. Vytvorili vlastný model PageRanku pre články Wikipédie a taktiež zdroje DBpedia, ktoré sú dostupné pomocou SPARQL dotazu. [35]

Ďalším problémom sú typické literálové vlastnosti. Vlastnosti, ktoré nemajú prepojenia. Tie sú obsiahnuté ako čisté časti textu. Takéto informácie nebude možné získať pomocou analýzov vzťahov článkov. Za účelom získania takýchto informácií budeme potrebovať rozšíriť proces hľadania prepojení. Znova pripadá do úvahy niekoľko prístupov ako to je možné realizovať. Pokúsme si načnúť niekoľko možnosti ako doplniť literálové vlastnosti k vzťahovým vlastnostiam:

- Identifikáciou slov, ktoré sú okrem abstraktu obsiahnuté vo zvyšku textu.
- Aplikáciou metód dolovania dát z textov.
- Sumarizáciou tried obsiahnutých v DBpedií

Zaujímavým prístupom je identifikácia slov, ktoré sú použité aj vo zvyšnej časti textu. Existuje pritom niekoľko problémov. Prvým z nich je problém mapovania hodnôt na existujúce slovníky. Menším problémom je skutočnosť, že budeme pracovať s datestom NIF DBpedia abstrakt, pričom DBpedia aktuálne ešte nepublikovala podobný NIF dataset pre celé články z Wikipédie <sup>3</sup>. Spôsob sumarizácie na základe tried zdrojov aktuálne obsiahnutých, predstavuje jednoduchší prístup, ktorý na druhej strane zaručí spoľahlivosť získaných informácií.

Za zmienku ešte stojí otázka počtu zástupcov jednotlivých tried pre sumarizáciu. Počet rozhodne nesmie byť príliš nízky a zároveň nemá zmysel zvoliť vysoký počet. Jednou z možností je použiť počet percentuálneho zastúpenia v danej triede. Tu však nastáva problém s nerovnomerným rozložením zdrojov

---

<sup>3</sup>Až od verzie DBpedia 2016-10 budú zverejnené všetky texty vo formáte NIF.

medzi triedami. Aktuálne DBpédia obsahuje približne 3,7 milóna zdrojov triedy *dbo:Agent* nadruhej strane len približne 22 tisíc pre triedu *dbo:Disease*[36]. Nepredpokladá sa že by sa tento rozdiel v rozložení mohol niekedy zmeniť. Stanovením konštantného počtu sa zníži náročnosť, ale u tried s výrazne vyšším zastúpením je možná nižšia kvalita získaných výsledkov.

Možným problémom je, že vo výsledku sa objavia duplicitné hodnoty, ktoré sú mapované na iné predikáty slovníkov. Bude vhodné, ak sa pokúsime takéto duplikácie redukovať. Jednou z duplicit, ktoré sa najpravdepodobnejšie objavia je duplicita spôsobená dvojitým spôsobom extrakcie informácií z Wikipédie do DBpedie. Aj napriek tomu, že hodnoty sú mapované na predikát s prefixom <http://dbpedia.org/ontology/>, sú často priradené aj pomocou predikátu <http://dbpedia.org/property/>. Ako bolo spomenuté vyššie v kapitole 14, hodnoty predikátátov typu *ontology* sú kvalitnejšie. Je možné, že existujú všetky 3 kombinácie výskytu týchto predikátov. Uvažujeme len kombinácie, kedy zdroj obsahuje aspoň jeden z týchto predikátov, alebo obsahuje obidve. Príkladom je zdroj <http://dbpedia.org/page/Mazda>, ktorý obsahuje predikáty <http://dbpedia.org/ontology/owner> a zároveň <http://dbpedia.org/property/owner>. V tomto prípade ide o nežiadúcu duplicitu hodnôt.

Báza znalosti by mala obsahovať zoradenú množinu predikátov pre každú definovanú triedu v rámci zvolenej ontológie. Radenie je závislé na pridelenom hodnotení predikátov. Musíme si rozobrať možné spôsoby určenia významnosti predikátov:

- Pomocou percentuálneho zastúpenia.
- Pomocou relatívnej pozícií v textoch.

Stanovenie dôležitosti pomocou percentuálneho zastúpenia je jednoduchým spôsobom, ktorý je vhodnejší na situácie, kedy nedokážeme nájsť sofistikovanejší spôsob stanovenia významnosti. Identifikácia pozície v texte je priamo závislá na možnosti spätného zistenia umiestnenia odkazu v texte abstraktu.

### 2.3 Návrh modelu sumarizácie entít

V rámci definovaných vstupov, čiastkových úloh a problémov uvedených v analýze si v tejto časti predstavíme konceptuálny model pre našu sumarizáciu entít.

Vďaka tomu, že existuje verejne dostupný dataset DBpedia abstrakt vo formate NIF, budeme môcť získať vzťahy zdrojov priamo z neho. Model bude založený sa využití NIF DBpedia datasetu. Výber ontológie pre výber zdrojov daných tried je spojený so závislosťou na DBpedií. Budeme preto používať triedy ontológie DBpedie. Výber kontrétných zdrojov, pre sumarizáciu, budeme realizovať pomocou získania existujúceho hodnotenia DBpedia Page-Rank. Počet zvolených zdrojov stanovíme pomocou konštanty. Zlepšíme tým



rýchlosť spracovania, hoci vznikne riziko zníženia kvality výsledkov u nadpočetných tried. Nepredpokládame však, že by takéto zníženie bolo výrazne. Spôsob, ktorým doplníme literálové vlastnosti si pre náš model zvolíme sumarizáciu tried obsiahnutých v DBpedii. Bude sa jednať o ďalšiu doplnkovú sumarizáciu. Ďalšou doplňujúcou úlohou v rámci modelu bude redukcia duplicitných predikátov. Pokúsime sa ich na základe podobnosti zoskupiť pod jeden záznam vo výsledku sumarizácie. Dokážeme tak poskytnúť kvalitnejšiu bázu znalosti. Zároveň bude lepšie splnený cieľ pre zlepšenie prehľadnosti informácií o zdrojoch.

Model navrhnutý na základe analýzy obsahuje 4 hlavné komponenty a 3 pomocné. Model je znázornený na obrázku 2.4. Hlavnými komponentami sú:

**Extraktor prepojení** Má za úlohu na základe získaných zdrojov určených pre sumarizáciu, vyhľadať ich záznamy v NIF DBpedia abstrakt data-sete. Jeho výstupom je množina prepojení, pre ktoré cez *Vyhľadávač predikátov* doplní predikáty, ktoré odpovedajú najdeným prepojeniam. Všetky získane informácie následne pomocou *Generátora čiastkových výsledkov* uloží.

**Generátor globálnej štatistiky** Má za úlohu paralelne s generovaním štatistiky zdrojov, vygenerovať jednu globálnu štatistiku literálov pre všetky triedy ontológie. Generuje ich na základe vstupných zdrojov a pomocou *Vyhľadávača predikátov* získa potrebné predikáty. Výsledok uloží v jednom súbore.

**Identifikátor identických predikátov** Má za úlohu pre všetky vstupné predikáty identifikovať tie, ktoré sú identické na základe stanovenej miery podobnosti. Výsledky priebežne ukladá do súboru.

**Generátor bázy znalosti** Má za úlohu spojiť všetky čiastkové štatistiky jednotlivých zdrojov, pripojiť globálnu štatistiku a redukovať identické predikáty.

Pomocné komponenty:

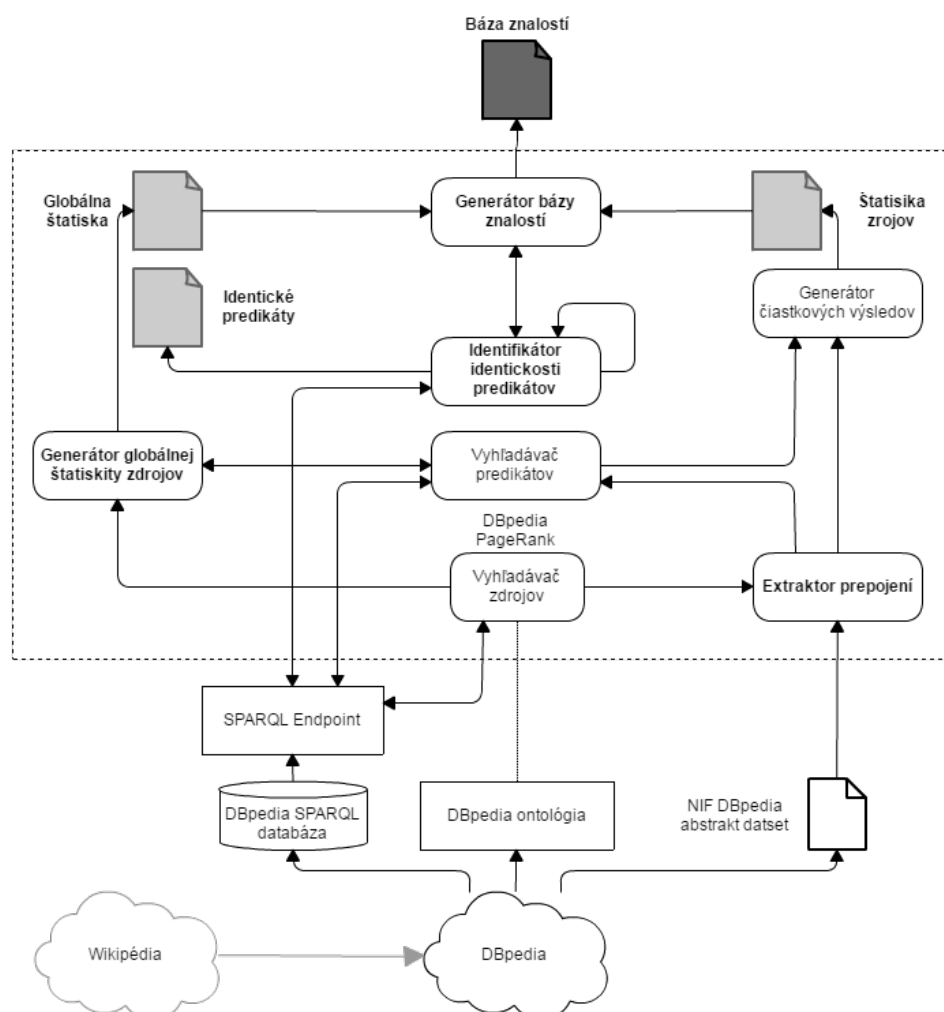
**Vyhľadávač zdrojov** Ma za ulohu na základe DBpedia PageRanku a triedy ontológie vrátiť zoradený požadovaný počet zdrojov.

**Vyhľadávač predikátov** Má za úlohu na základe zadaných kritérií vrátiť zoznam predikátov.

**Generátor čiastkových výsledkov** Má za úlohu uchovať všetky potrebné informácie daného zdroja pre následne spracovanie.

Ďalej si podrobnejšie popíšeme hlavné komponenty. Predovšetkým aké vnútorne procesy riešia a zároveň ako vypočítavajú jednotlivé hodnotenia.

## 2. METÓDA PRE SUMARIZÁCIU A HODNOTENIE VÝZNAMNOSTI INFORMÁCIÍ



Obr. 2.4: Model sumarizácie entít s využitím NIF DBpedia abstraktu

### 2.3.1 Extraktor prepojení

Komponenta bude pracovať s datasetom formátu NIF. Pre potrebu zabezpečenia zvládnutia veľkých datasetov s ním bude pracovať pomocou streamu. Rešpektuje dávkové spracovanie a pre zadanú množinu zdrojov vyhľadá potrebné údaje pomocou jediného prechodu datasetom. Potrebné údaje pre výstup sú definované na základe potrieb sumarizácie. Pre každý záznam zdroja vyhľadáva násobne (v zápise je použitý prefix nif:<sup>4</sup>):

<sup>4</sup>nif: zastúpuje <http://persistence.uni-leipzig.org/nlp2rdf/ontologies/nif-core#>

- Prepojenie. Hodnota priradená v NIF triplete pomocou predikátu `http://www.w3.org/2005/11/its/rdf#taIdentRef`
- Indexy pozície v texte, pre text prepojenia. Hodnoty priradené v NIF triplete pomocou predikátu `nif:beginIndex` a `nif:endIndex`
- Paragraf výskytu prepojenia. Hodnota priradená v NIF triplete pomocou predikátu `nif:hasSection`

Výsledky úkladá ako čiatočné štatistiky každého zdroja. Rozlišuje 2 druhy získaných predikátov. Predikáty striktné viazané na daný zdroj a predikáty viazané na celú triedu ontológie. V rámci ukladania realizuje výpočet čiastkového skóre. To je vypočítané na základe relatívnej pozície odkazovaného textu v celom abstrakte. Zároveň je znižované pozíciou daného paragrafu v abstrakte. Nech  $A$  je množina znakov abstraktu a  $D$  množina paragrafov abstraktu  $A$  potom počet znakov abstraktu  $A$  je rovný  $c_A = |A|$  a počet paragrafov  $c_D = |D|$ . Ďalej nech je  $L$  množina degradácie dôležitosti vplyvom pozície paragrafu potom degradácia pre paragraf  $D_i$  je stanovená ako  $L_{D_i} = 1 - \frac{i}{c_D} \times 0.1$  pričom  $c_D > 0$ . Nech pre odkaz  $o$  v paragrafe  $D_i$  je  $f$  index prvého znaku odkazovaného textu v abstrakte potom je pozičná váha odkazu  $o$  stanovená ako  $W_o = 1 - \frac{f}{c_A}$  pričom  $c_A > 0$ . Výsledne čiastkové skóre pre predikát  $p$  odkazu  $o$  obsiahnutého v paragrafe  $D_i$  je definovaný ako súčin pozičnej váhy a degradácie pozíciou paragrafu. Vzťah je znázornený na obrázku 2.5

$$S_{T_p} = W_o \times L_{D_i}$$

Obr. 2.5: Definícia čiastkového skóre predikátu  $p$  triedy  $T$ 

### 2.3.2 Generátor globálnej štatistiky

Komponenta je zodpovedná za vytvorenie doplnkovej štatistiky sumarizáciou zdrojov DBpedia na základe príslušnosti k určitej triede. Pre každú triedu ontológie získa potrebný počet jej zdrojov. Pre každý zdroj získava všetky jeho predikáty. Tie uchováva a navyšuje počet výskytu jednotlivých predikátov v rámci danej triedy. Mieru významnosti priradzuje na základe percentuálneho výskytu. Nech je  $C$  množina všetkých počtov výskytu predikátov potom  $C_p$  je počet výskytu unikátného predikátu  $p$  a  $S_p$  definujeme ako výsledné skóre predikátu  $p$  pre globálnu štatistiku danej triedy  $T$ . Vzťah je znázornený na obrázku 2.6.

$$G_{T_p} = \frac{C_p}{\max(C,1)}$$

Obr. 2.6: Definícia globálneho skóre predikátu  $p$  triedy  $T$

### 2.3.3 Identifikátor identických predikátov

Účelom komponenty je redukcia duplicitných predikátov na základe priradených hodôt v podobe objektov. Vstupom pre komponentu je množina predikátov zoradených zostupne na základe aktuálne priradeného skóre. Overí pre všetky dvojice predikátov, počet identicky hodnôt priradených v rámci jedného zdroja. Počet vyšetrovaných predikátov je redukovaný na maximalný povolený počet. Nutné opatrenie pred neúmerne vysokým časom potrebným na vyšetrenie úplne všetkých kombinácií. Z redukovanej množiny predikátov sa vygeneruju všetky kombinácie dvojíc. Pre výpočet kritéria, ktoré rozhodne o príslušnosti dvojice k vzájomne identickým predikátom, je použité percentuálne zastúpenie. Pre každý predikát sa získa počet všetkých použití v rámci zdrojov. Nech je  $c_a$  počet všetkých použití predikátu  $a$  a  $c_b$  počet všetkých použití predikátu  $b$ . Nech  $o_a$  je objekt mapovaný pomocou predikátu  $a$ , objekt  $o_b$  je mapovaný pomocou predikátu  $b$  a  $c_{ab}$  počet spoločného použitia predikátov  $a$  a  $b$  v jednotlivých tripletoch, tak že objekty  $o_a = o_b$  potom mieru identickosti predikátov definujeme ako podiel  $c_{ab}$  a maxima z celkového počtu použitia predikátov  $a$  a  $b$ . Vzťah je znázornený na obrázku 2.7.

$$I_{ab} = \frac{c_{ab}}{\max(c_a, c_b, 1)}$$

Obr. 2.7: Definícia miery identickosti dvojice predikátov  $a, b$

Na základe minimálnej miery identickosti sú dvojice ohodnotené, buď to ako identické alebo rôzne. Tieto výsledky sú priebežne ukladané. V prípade, že už pre danu dvojicu existuje zaradenie do jednej z týchto skupín a nie je explicitne zadaný požiadavok na ich opätovnú identifikáciu je dvojica vynechaná z procesu identifikácie.

### 2.3.4 Generátor bázy znalosti

Komponenta ma za úlohu spojenie čiastkových výsledkov jednotlivých zdrojov a globálnej štatistiky pre definovanú vstupnú triedu ontológie. Zároveň pri generovaní bázy znalosti aplikuje získané vedomosti o identickosti predikátov. Preto jednotlivé záznamy osahujú množinu predikátov, ktorá v prípade identických obsahuje viacero predikátov. Tejto množine predikátov je stanovené jednotné výsledne skóre dôležitosti. Skóre dôležitosti je reálnym číslom v intervale  $(0, 1 >$ , pričom hodnotenie 1 priradzuje maximálnu dôležitosť. Proces výpočtu výsledného skóre pozostáva z dvoch hlavných operácií.

Prvá operácia predstavuje spojenie čiastkových výsledkov. Pre každý predikát je sumované čiastkové skóre a celkový počet výskytu vo vyšetrovaných zdrojoch. Nech  $c_p$  je počet výskytu predikátu  $p$  a  $S_{p_i}$  je čiastkové skóre predikátu  $p$  na pozícii  $i$  potom sumu celkového čiastkového skóre predikátu  $p$  definujeme ako sumu všetkých čiastkových skór. Vzťah je znázornený na obrázku 2.8

$$S_{T_p} = \sum_{i=1}^{c_p} (S_{p_i} \times \frac{1}{\max(c_p, 1)})$$

Obr. 2.8: Definícia celkového čiastkového skóre predikátu  $p$  triedy  $T$ 

Druhá operácia načítava výsledky globálneho skóre všetkých predikátov danej triedy a pripája ich k aktuálnym celkovým čiastkovým výsledkom predikátov. Celkové výsledne skóre pre predikáty je stanovené spojením celkového čiastkového skóre a globálneho skóre. Vzťah stanovenia výsledného skóre je znázornený na obrázku 2.9.

$$R_{T_p} = \begin{cases} \frac{S_{T_p} + G_{T_p}}{2} & , S_{T_p} > 0 \wedge G_{T_p} > 0 \\ \max(S_{T_p}, G_{T_p}) & , \text{inak} \end{cases}$$

Obr. 2.9: Definícia výsledného skóre predikátu  $p$  triedy  $T$ 

Tento vzťah definuje všeobecný spôsob stanovenia výsledného skóre. Na základe aktuálneho návrhu však neuvažujeme iné než literálové predikáty v globalnej štatistike a opačne pri čiastkových výsledkoch. Preto sa bude v praxi jednať o priradenie jednotlivých výsledkov.

Navrhovaný model poskytuje výslednú bázu znalosti, okrem toho poskytuje vedľajšie výsledky. Medzi tie najpodstatnejšie patrí zoznam identických predikátov a globálna štatistika literálových predikátov. Okrem toho dokáže poskytnúť aj jednotlivé štatistiky zvolených zdrojov. Na základe týchto skutočností predpokládame, že model dokáže pokryť potrebné ciele tejto práce pre stanovenie miery dôležitosti predikátov. Zároveň poskytne možné uplatnenie vedľajších výsledkov aj v rámci iných prác. Pre potreby čo možno najľahšieho použitia modelu bude nutná vhodná implementácia nástroja, ktorý aplikuje navrhnutý model. Problematiku implementácie nástroja navrhnutého modelu budeme riešiť v nasledujúcej kapitole.

## 2.4 Implementácia nástroja

V tejto kapitole sa budeme venovať celému procesu implementácie nástroja, ktorý aplikuje navrhnutý model. V úvode si popíšeme zvolený programovací jazyk a spôsob akým zabalíme nástroj v rámci použiteľného balíčka. Následne navrhne štruktúru nástroja a uvedieme niektoré zaujímavé príklady z jeho implementácie.

### 2.4.1 Ruby

Ruby je skriptovací interpretovaný programovací jazyk. Vďaka tomu, že Ruby má celkom jednoduchú syntax, stále je však dosť výkonný na to, aby sa vedel vyrovnat' známejším programovacím jazykom. Napríklad Perl a Python. Ruby je plne objektovo orientovaný. Autorom Ruby je len jeden človek, Yukihiro Matsumoto (známi ako „Matz“). Nevyhovovali mu ostatné programovacie jazyky, a preto sa rozhodol vziať z viacerých jazykov ich užitočné stránky a vytvoriť svoj vlastný programovací jazyk. Tak v roku 1995 vzniká prvá verzia Ruby. Štandardná implementácia Ruby („Matz“ interpret) je napísaná v jazyku C, pričom má podobu klasického interpretra. Existuje niekoľko zaujímavých alternatív. Patrí medzi ne napríklad natívna implementácia v jazyku Java. JRuby umožňuje beh kódu Ruby v JVM prostredí. To samozrejme poskytuje všetky výhody JVM. [37]

Spôsob akým umožňuje jazyk Ruby vytvoriť ucelenú časť kódu v rámci určitého balíčka je vytvorenie tzv. gemu (angl. gem). Gemy sú knižnice pripraveného Ruby kódu, ktorý rieši určitú funkcionálnosť. Pre správu gemov existuje ruby balíčkovací manažér *RubyGems*. Tento manažér poskytuje štandardný spôsob distribúcie Ruby programov, jednoduchú inštaláciu a správu gemov. RubyGems používa sémantické verzovanie. Verzia sa uvádza ako reťazec bodkami oddelenými troch nezáporných čísel. Jednotlivé čísla sa definujú ako "major"."minor"."fix". [38] Každý gem obsahuje tieto základné časti:

- Zdrojový kód. Obsahuje výkonú časť funkcionality gemu, vytvorenú vývojárom.
- Gemspec súbor. Súbor ktorý obsahuje informácie o gеме ako sú napríklad verzia, platforma, popis gemu. Obsahuje taktiež definície pre načítanie všetkých potrebných súborov z časti zdrojového kódu. Predatavuje súbor na základe, ktorého je riadený proces vytvárania gemu (angl. build gem).
- Súbor README.md. Obsahuje základný návod pre použitie daného gemu.
- Súbor LICENSE.md. Obsahuje informácie o licenciách.

[38]

### 2.4.2 Architektúra nástroja

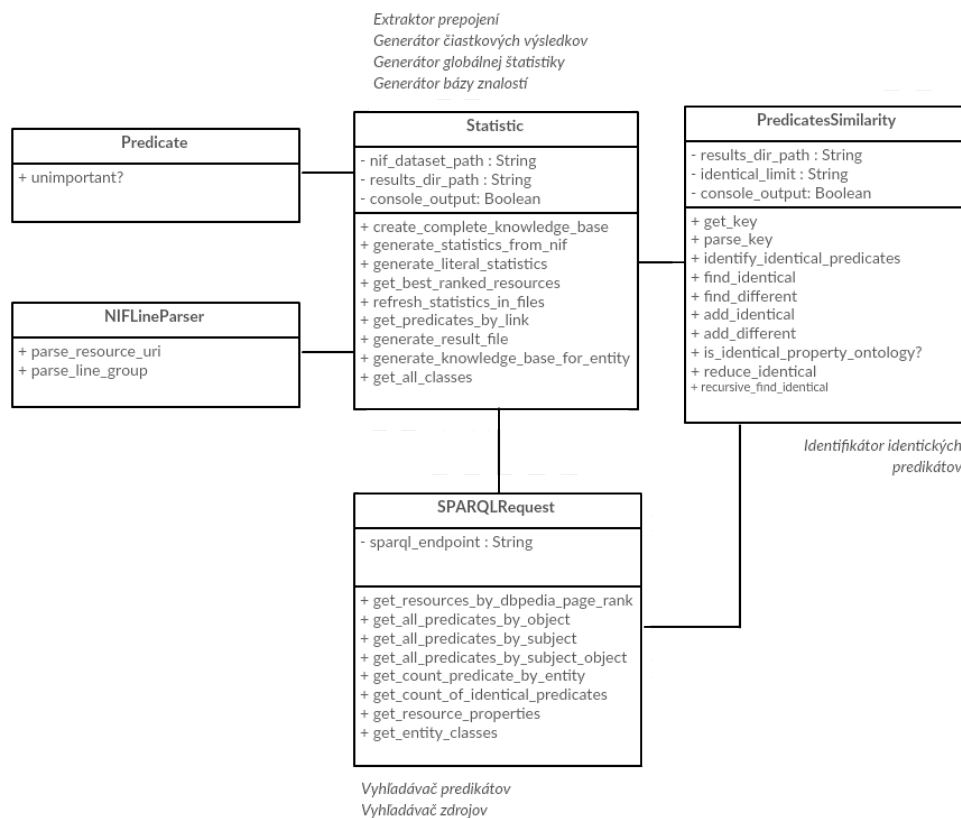
V zmysle definovaného modelu, v tejto časti predstavíme architektúru vytvoreného nástroja. Hlavnými časťami nástroja sú triedy. Rozlišujeme pomocné triedy a aplikačné triedy, ktoré zodpovedajú za hlavnú funkcionálnosť nástroja. Pomocné triedy majú na starosti hlavne prácu s ukladaným dočasným súborom, parsovaním určitých štruktúr a pod. Diagram aplikačných tried je znázornený na obrázku 2.10. Nástroj obsahuje 3 hlavné aplikačné triedy, ktoré v

sebe majú implementované všetky komponenty z navrhnutého modelu. Sú to triedy:

**Statistic** Trieda v sebe obsahuje funkčnosti spojené s generovaním všetkých častí štatistik ako aj výslednej bázy znalostí.

**PredicateSimilarity** Trieda ma na starosti funkcionality pre komponentu identifikátor identických predikátov. Vyhodnocuje identickosť a stará sa o uchovanie získaných informácií o identickosti.

**SPARQLRequest** Predstavuje triedu, ktorá má na starosti komunikáciu s SPARQL endpointom. Ma definované metódy pre získanie potrebných výsledkov. Jednotlivé metódy aplikujú potrebné dotazy.



Obr. 2.10: Základný diagram tried nástroja

### 2.4.3 Charakteristika tried

Trieda *Statistic*, predstavuje najpodstatnejšiu triedu. Zabezpečuje zároveň najnáročnejšie operácie v rámci nástroja. Jedným z problémov, ktorý sa v

rámci jej implementácie objavil bol spôsob použitia veľkého datasetu. Aj napriek tomu že v rámci jazyka existujú gemy, ktoré sú určené pre prácu s RDF datasetmi, nedokázali sme získať taký, ktorý by neriešil použitie datasetu spôsobom, kedy dochádza k načítaniu celého súboru do pamäte. Preto bolo nutné použiť vlastné riešenie pre vyhľadávanie tripletov v datsete. Nutnosťou bolo spracovanie pomocou streamu. Pri takomto použití načítavania sme získali postupne jednotlivé riadky súboru. V rámci načítavania riadkov sme následne identifikovali jednotlivé tvrdenia (angl. statement) ako postupnosť sedmíc riadkov. Tieto skupiny riadkov boli pre jednoduchosť parsované pomocou Ruby metódy *scan*, pre stringové hodnoty, s využitím nami definovanými regulárnymi výrazmi. Na základe týchto potrieb vznikla trieda *NIFLineParser*. Jej hlavnou úlohou je získanie potrebných údajov z jednotlivých skupín riadkov. Týmto sa vyriešil spôsob extrakcie informácií z NIF DBpedia datasetu. Po extrakcii informácií trieda pomocou metód triedy *SPARQLRequest* získa potrebné predikáty. Pre predikáty vypočíta čiastkové skóre. Všetky výsledky uloží v cieľovom adresari.

Trieda *PredicatesSimilarity* ma jedinú úlohu. Rozhodnúť o miere identickosti predikátov. Pre potreby rozhodnutia o príslušnosti k skupine na základe miery identickosti definovanej v 2.7 bolo nutné stanoviť prahovú hodnotu. V prvých testoch sme použili vyššiu hodnotu 0.9. Následne sme na základe výsledkov usúdili, že je táto hodnota príliš striktná a nezískali sme takmer žiadne identické predikáty. Znížili sme danú konštantu na úroveň 0.8. Pri takto nastavenej úrovni sa zvýšilo množstvo očakávaných identických predikátov, pričom zároveň sa počas testov neobjavili žiadne nevhodne priradené. Preto pokladáme hodnotu 0.8 za prijateľnú pre naše potreby. Ďalšou dôležitou konštantou, ktorú bolo nutné zaviesť je maximálny počet predikátov k overeniu. Keďže proces identifikácie spočíva v overení každej kombinácií dvojíc, je nárast výpočetného času výrazný. Stanovili sme hornú hranicu predikátov na 250, odpovedá to 31 125 kombináciám dvojíc. Tento limit je aplikovaný pri každom volaní metódy *identify\_identical\_predicates*

V rámci triedy *Predicate* je definovaná jediná metóda, určená pre vylúčenie predikátov, ktoré sú nepodstatné. Táto príslušnosť bola definovaná manuálne. Nepredstavuje teda žiaden automatický proces. Predikáty sú vylúčené na základe myšlienky, že nepodávajú informácie k odlišeniu jednotlivých tried medzi sebou. V rámci predikátov jednotlivých zdrojov, existujú také, ktoré môžu mať vysokú mieru informácie, nepridávajú však žiadnu pridanú hodnotu pre informovanie používateľa. Medzi takéto predikáty patria napríklad identifikátory Wikipedia článkov, z ktorých boli extrahované. Takéto predikáty potom obsahuje každý zdroj. Vplyvom na spôsob stanovenia výsledného skóre dôležitosti, sa vždicky tieto predikáty umiestňovali na prvých miestach. Preto sa zaviedla množina tzv. nedôležitých predikátov. Medzi nedôležité predikáty patria:

- <http://xmlns.com/foaf/0.1/primaryTopic>
- <http://dbpedia.org/ontology/wikiPageRedirects>



- <http://dbpedia.org/ontology/wikiPageDisambiguates>
- <http://dbpedia.org/ontology/wikiPageRevisionID>
- <http://dbpedia.org/ontology/wikiPageID>
- <http://www.w3.org/2002/07/owl#sameAs>
- <http://www.w3.org/2000/01/rdf-schema#seeAlso>
- <http://www.w3.org/2002/07/owl#differentFrom>
- <http://dbpedia.org/ontology/wikiPageExternalLink>
- <http://xmlns.com/foaf/0.1/depiction>

Na druhej strane sme zároveň identifikovali predikáty, u ktorých nemá zmysel uvažovať ich dôležitosť. Preto sú aj tieto predikáty vylúčené a nebudú použité v rámci sumarizácie. Patria sem hlavne základne popisné predikáty. Zvolili sme množinu týchto, ako takzvaných všeobecne dôležitých predikátov:

- <http://dbpedia.org/ontology/thumbnail>
- <http://xmlns.com/foaf/0.1/name>
- <http://www.w3.org/2000/01/rdf-schema#label>
- <http://dbpedia.org/property/name>
- <http://dbpedia.org/property/commonName>
- <http://dbpedia.org/property/title>
- <http://www.w3.org/2000/01/rdf-schema#comment>
- <http://dbpedia.org/ontology/abstract>

Ako sme už uviedli trieda *SPARQLRequest* obstaráva komunikáciu s endpointom. Má preto zároveň preddefinovanú štruktúru dotazov pre potrebné výsledky. Pre vyjasnenie spôsobu získavania jednotlivých informácií z DBpedia SPARQL endpointu, si predstavíme podstatné dotazy.

**Získanie zdrojov na základe DBpedia PageRank.** V rámci tohto dotazu je potrebné definovať parameter *entity* a *limit* pričom entity definuje pre akú triedu budú získané zdroje. Napríklad je to <http://dbpedia.org/ontology/Person>. Limit je defaultne nastavený na hodnotu 10. Dotaz je znázornený na obrázku 2.11.

## 2. METÓDA PRE SUMARIZÁCIU A HODNOTENIE VÝZNAMNOSTI INFORMÁCIÍ

---

```
PREFIX rdf:<http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX dbo:<http://dbpedia.org/ontology/>
PREFIX vrank:<http://purl.org/voc/vrank#>

SELECT ?entity ?rank
FROM <http://dbpedia.org>
FROM <http://people.aifb.kit.edu/ath/#DBpedia_PageRank>
WHERE {
    ?entity rdf:type #{entity_type}.
    ?entity vrank:hasRank/vrank:rankValue ?rank.
}
ORDER BY DESC(?rank) LIMIT #{limit}
```

Obr. 2.11: SPARQL dotaz pre získanie zdrojov na základe DBpedia PageRank

```
SELECT DISTINCT ?property
WHERE {
    #{subject} ?property #{object}.
}
```

Obr. 2.12: SPARQL dotaz pre získanie predikátov zdroja na základe extrahovaného prepojenia

### Získanie predikátov zdroja na základe extrahovaného prepojenia.

Vstupom pre tento dotaz je vyšetřovaný zdroj ako *subject* a extrahované prepojenie ako *object*. Výsledok je použitý ako tzv. striktné predikáty vyšetřovaného zdroja. Dotaz je znázornený na obrázku 2.12.

### Získanie predikátov triedy zdroja na základe extrahovaného prepojenia.

Získanie týchto predikátov je logicky rozdelené do dvoch zvláštných dotazov v prvom sa získavú kompletne všetky predikátky z tripltov, v ktorých je *object* extrahované prepojenie. V druhom dotaze sa pre každý predikát vyhledá počet všetkých použití predikátu v triede ako *entityClass*. Uvažuje sa vstupné aj výstupné prepojenie. Dotaz je znázornený na obrázku 2.15 respektíve 2.14.

```
SELECT DISTINCT ?property
WHERE {
    ?subject ?property #{object}.
}
```

Obr. 2.13: SPARQL dotaz pre získanie všetkých predikátov na základe extrahovaného prepojenia

### Získanie predikátov pre globálnu štatistiku triedy.

Vstupom pre dotaz je zdroj ako *subject*. V prípade že sú požadované iba predikaty v rámci

```

SELECT DISTINCT COUNT(?subject) as ?count
WHERE {
  ?subject a #{entity_class} .
  {?subject #{predicate} ?a .} UNION {?b #{predicate} ?subject .}
}
ORDER BY DESC(?count)

```

Obr. 2.14: SPARQL dotaz pre získanie počtu použitia predikátu triedy.

tripletovej s literalmi, doplní sa do dotazu filtračné kritérium `only_literal` ako `FILTER(isLiteral(?object))`. Dotaz je znázornený na obrázku 2.15.

```

SELECT DISTINCT ?property
WHERE {
  #{subject} ?property ?object .
  #{only_literal}
}

```

Obr. 2.15: SPARQL dotaz pre získanie predikátov globálnej štatistiky triedy

## 2.5 Výstup nástroja

Výstupom nástroja sú 4 typy súborov v zmysle navrhnutého modelu. Všetky súbory sú vo formáte JSON. Voči modelu je pozmenený spôsob ukladaných hodnôt. Pre jednotlivé predikáty v rámci všetkých typov okrem výslednej bázy znalosti, sú priradené počty použitia. V prípade aplikácie výpočtu by sa uchovávali len relatívne hodnoty, ktoré pri zmenách a pokusoch nebude možné späťne prepočítať na pôvodne počty. Preto sa všetky typy skôr prepočítavajú až pri finálnom generovaní bázy znalosti. Ďalej si ukážeme štruktúru jednotlivých výsledných súborov.

**Súbor s globálnou štatistikou.** Súbor obsahuje JSON objekt so štruktúrou kľúč *trieda ontológie* a hodnota objekt s štruktúrou kľúč *predikát* a hodnota *počet všetkých výskytov*. Štruktúra je znázornená na obrázku 2.16.

**Súbor s identickými predikátmi.** Súbor obsahuje pole reťazcov. V rámci jedného reťazca sú uchované predikáty, ktoré sú si navzájom identické. Každý predikát je ohraničený znakmi «ä »". Štruktúra je znázornená na obrázku 2.17.

**Súbor s čiastkovým výsledkom zdroja.** Súbor obsahuje na prvej úrovni 2 základne kľúče *resource\_uri* s hodnotou URI zdroja, ktorému patria

## 2. METÓDA PRE SUMARIZÁCIU A HODNOTENIE VÝZNAMNOSTI INFORMÁCIÍ

```
{
  "City": {
    "http://dbpedia.org/ontology/PopulatedPlace/areaMetro": 415,
    ...
  }
  ...
}
```

Obr. 2.16: Štruktúra súboru s globálnou štatistikou

```
[
  "<http://dbpedia.org/ontology/raceHorse><http://dbpedia.org/property/horses>",
  ...
]
```

Obr. 2.17: Štruktúra súboru s identickými predikátmi

výsledky a *nif\_data*, ktorý obsahuje pole s objektmi pre každé nájdené prepojenie. Objekt prepojenia obsahuje kľúče *link* s hodnotou nájdeného prepojenia v NIF datasete, *anchor* s hodnotou reťazca, pre ktorý je dané prepojenie mapované, *indexes* obsahuje pole s počiatočným a koncovým indexom daného reťazca, *section* s hodnotou z akého paragrafu v abstrakte je získane prepojenie, *weight* s hodnotou čiastkového skóre (čiastkové skóre je zhodné pre všetky predikáty obsiahnuté v rámci daného objektu prepojenia, výpočet je založený na definícii 2.5). Ďalej obsahuje *properties* a *strict\_properties*, pričom prvému sú priradené hodnoty predikátov nájdené pomocou daného prepojenia pre celú triedu a druhému sú priradené predikáty, nájdené pomocou prepojenia priamo k vyšetrovanému zdroju *resource\_uri*. Obe priradenia sú zároveň zanorené pod triedu, do ktorej patrí daný zdroj. Štruktúra je znázornená na obrázku 2.18.

**Súbor s bázov znalosti.** Súbor obsahuje objekt s kľúčmi pre jednotlivé triedy.

Každý triede je priradené pole objektov. Pole predstavuje zoradený zoznam na základe miery dôležitosti každého záznamu. Objekty obsahujú výsledné skóre dôležitosti *score* a pole *predicates* s predikátmi, ktoré sú priradené na danú pozíciu. Štruktúra je znázornená na obrázku 2.19.

Nástroj bol vo finálnej verzii zabalený ako gem na základe štandardu RubyGems. Zároveň bol pre potreby neskoršieho použitia publikovaný pomocou RubyGems manažéra. Gem je dostupný na adrese [https://rubygems.org/gems/browser\\_web\\_data\\_entity\\_sumarization](https://rubygems.org/gems/browser_web_data_entity_sumarization). Spôsob ako je možné gem získať je jeho inštalácia. Spôsob inštalácie a zavedenia gemu pre použitie je znázornený na obrázku 2.20. Po vykonaní *require* je celá funkcionálnosť gemu sprístupnená.

```
{
  "resource_uri": "http://dbpedia.org/resource/Abraham_Lincoln",
  "nif_data": [
    {
      "link": "http://dbpedia.org/resource/Whig_Party_(United_States)",
      "anchor": "Whig Party leader",
      "indexes": ["588", "605"],
      "section": "paragraph_435_1442",
      "properties": {
        "Person": {
          "http://dbpedia.org/ontology/occupation": 329669.0,
          "http://dbpedia.org/ontology/otherParty": 3584.0,
          "http://dbpedia.org/ontology/politicalPartyInLegislature": 4.0,
          "http://dbpedia.org/property/otherparty": 3620.0,
          "http://dbpedia.org/property/party": 78771.0
        }
      },
      "weight": 0.8715,
      "strict_properties": {
        "Person": {
          "http://dbpedia.org/ontology/party": 80846.0
        }
      }
    },
    ...
  ]
}
```

Obr. 2.18: Štruktúra súboru s čiastkovým výsledkom zdroja

```
{
  "Game": [
    {
      "score": 0.8912,
      "predicates": [
        "http://dbpedia.org/ontology/designer",
        "http://dbpedia.org/property/designer"
      ]
    },
    ...
  ],
  ...
}
```

Obr. 2.19: Štruktúra súboru bázy znalostí

```
V CMD:
gem install browser_web_data_entity_sumarization

V Ruby projekte – gemfile:
gem 'browser_web_data_entity_sumarization'

V Ruby skripte:
require 'browser_web_data_entity_sumarization'
```

Obr. 2.20: Spôsob inštalácie nástroja

## 2.6 Použitie nástroja

V tejto kapitole si popíšeme, ako sme generovali výslednú bázu znalosti. Poskytneme takto stručný návod použitia nástroja.

## 2. METÓDA PRE SUMARIZÁCIU A HODNOTENIE VÝZNAMNOSTI INFORMÁCIÍ

---

Pre potreby vygenerovania bázy znalosti v zmysle sumarizácie sme analyzovali dostupné NIF datasety. Vychádzajúc zo štruktúry tak, ako je uvedená na obrázkoch 2.1, 2.2 a 2.3 sme za kľúčový dataset zvolili *nif-text-links*. Zvolený typ datasetu sme získali z publikácie DBpedia 2016-04. Dataset je dostupný na adrese <http://downloads.dbpedia.org/2016-04/ext/nif-abstracts/en/>. Samotný dataset ma približnú veľkosť 45 GB. V ďalšom kroku sme potrebovali vybrať, pre ktoré triedy ontológie budeme generovať výsledky. Ako sme už uviedli uvažujeme triedy z ontológie DBpedia. Hierarchia ako aj zoznam týchto tried je verejne dostupný na <http://mappings.dbpedia.org/server/ontology/classes/>. V úvodných testoch sme spracovali len triedy <http://dbpedia.org/ontology/City> a <http://dbpedia.org/ontology/Person>. V závere, po vyladení nedostatkov, sme sputili generovanie pre všetky triedy tejto ontológie. V čase generovania obsahovala 685 tried. Reálne však neboli pre každú definovanú triedu nájdené zdroje na základe ich DBpedia PageRanku. Vo výsledku sme spracovali sumarizáciu pre 401 tried. Dôležitým výberom bol počet zástupcov každej triedy. Zvolili sme si konštantný počet 100 zdrojov pre každú triedu. Ako sme už uviedli, pre niektoré triedy neboli získaní zástupcovia. Zároveň u niektorých tried bol počet získaných zástupcov nižší než bol stanovený požadovaný počet. Vo výsledku sme teda reálne spracovali 32 060 zdrojov. V prieme sa teda reálne pre každú triedu použilo približne 80 zdrojov.

Spôsob akým sa spúšťa generovania bázy znalosti, je volanie metódy *create\_complete\_knowledge\_base* inštancie triedy *Statistic*. Metóda na vstupe vyžaduje hash s parametrami:

- *entity\_types* - hodnotou je pole tried ontológie, pre ktoré budú následne získaní zástupci.
- *entity\_count* - hodnotou je počet, koľko zástupcov ma byť pre každú triedu získany.
- *demand\_reload* - označenie, či sa požaduje pregenerovanie čiastkových výsledkov pre zdroje, ku ktorým už aktuálne sú uložené.
- *identify\_identical\_predicates* - označenie, či sa požaduje v závere identifikácia identických predikátov.

Nami vygenerovaná báza znalosti mala aplikovane všetky konštrukcie a výpočty, ktoré sme si v tejto kapitole definovali. Pokladáme ju preto za kompletnú a vhodnú pre použitie v rámci webovej aplikácie za účelom zlepšenia prehľadania informácií. Preto sa v nasledujúcej kapitole budeme venovať použitiu tejto bázy znalosti.

---

## Webová aplikácia

V rámci tejto kapitoly si popíšeme proces návrhu implementácie a zavedenia webovej aplikácie, ktorá dokáže inerpretovať vygenerovanú bázu znalosti. V prvej časti sa pozrieme na spôsob intergacie bázy znalosti v rámci webovej aplikácie. Ďalej si popíšeme návrh a implementáciu.

### 3.1 Integrácia výsledkov sumarizácie entít

Báza znalosti je reprezentovaná vo formate JSON, preto nepredstavuje integrácia žiaden problém. Pre zjednotenie načítavania bázy znalosti je obsiahnutá v rámci vytvoreného nástroja. Tento nástroj je použitý na strane severa webovej aplikácie. Server využíva niektoré jeho pomocné metódy a získava spomínanú bázu znalosti. Nástroj ma preto implementovanú pomocnú triedu *CacheHelper*, pomocou ktorej okrem iného načítava bázu znalosti. O tento proces sa stará metóda *load\_knowledge*, ktorá na základe požadovaného typu výsledného súboru ho načíta z priečinku */knowledge*. Okrem už spomínanej bázy znalosti je v rámci tohto priečinku uložená hierarchia tried ontológie na základe, ktorej bola báza znalosti vygenerovaná a súbor so zoznamom všeobecne dôležitých predikátov. Všetky tieto výsledky sú práve potrebné pre aplikovanie výsledkov vo webovej aplikácii.

### 3.2 Návrh prehliadača

Pri našom návrhu prehliadača musíme vychádzať zo stanoveného cieľa práce a síce, sprehľadnenie informácií o zdrojoch DBpedia. Preto budeme potrebovať vytvoriť užívateľský jednoduchú aplikáciu, ktorá však musí zároveň poskytnúť dostatočnú účinnosť prehliadania informácií. Z hľadiska grafického návrhu budeme potrebovať užívateľský intuitívne prostredie. Aby bolo prostredie dostatočne intuitívne, potrebujeme si stanoviť aký je hlavný účel našej aplikácie. Hlavným účelom, bude samozrejme vyhľadávanie. Uvažujme teda používateľa

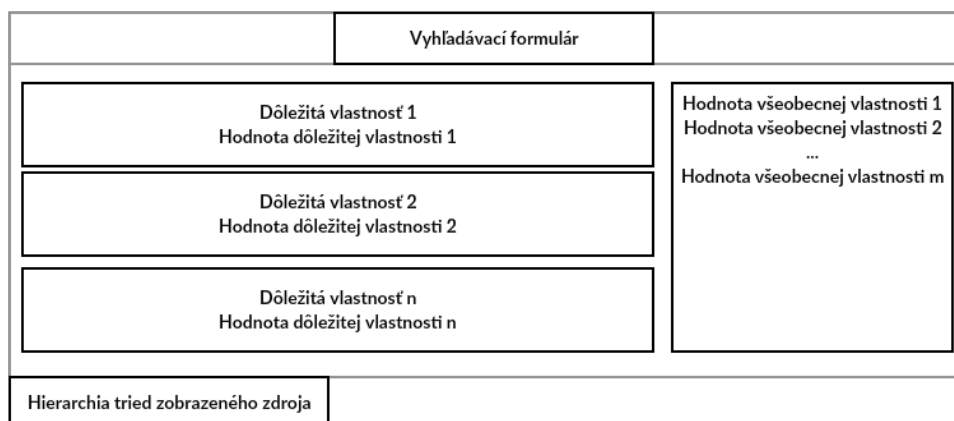
našej aplikácie, ktorý je posadený pred našu aplikáciu. Komponenta, ktorú bude najviac očakávať je vyhľadávacie pole (formulár), kde uvedie aký zdroj ma v záujme vyhľadať. V tomto štádiu nastáva určitý problém, ktorý sme si uvedomili. A síce ako umožniť používateľovi zadanie iba kľúčového výrazu, bez nutnosti poznania URI jednotlivých zdrojov. Bude jednoznačne veľkým prínosom ak mu poskytneme určitý našepkávač (angl. autocomplete) zoznam s návrhmi zdrojov. Na tento problém budeme myslieť v ďalšej časti implementácie. Máme teda základnú komponentu celej aplikácie, určenej k interakcii s používateľom. Mala by sa umiestniť v priestore, kde bude ľahko nájditeľná a zároveň používateľom očakávaná. Uvažujme teda základnú navigačnú lištu, ktorej dominantnou komponentou bude spomínaný vyhľadávací formulár. Pre ďalšie potreby uvažujme, že tam ďalej bude prepínanie základných stránok aplikácie.

Potom, ako používateľ odošle žiadosť, pre získanie informácií zvoleného zdroja, mali by sa zobrazíť v hlavnej časti stránky. Samotná voľba by mala byť, čo možno najrýchlejšia. Bez nutnosti stlačiť tlačidlo a pod. Povedzme teda, že sa tak stane po výbere navrhovaného zdroja, ktorý zaručene existuje na DBpedii. Jednotlivé informácie budú predstavovať určité dvojice. Ide o dvojice typu tvrdenie a význam. Vždy to bude názov vlastnosti a jej hodnota. Tieto výsledné informácie budú zobrazené pod navigačnou lištou.

Je ešte potrebné uvažovať, či existuje nejaký rozdiel v jednotlivých možných informáciách, ktoré sa budú zobrazovať. Jedným takým rozdielom môže napríklad byť účel informácie. Väčšina zdrojov obsahuje špecifické informácie, ktoré sú vlastné danému typu triedy, prípadne tried, do ktorých patrí. Avšak takmer u každého zdroja sú taktiež tzv. spoločné typy informácií. Myslíme tým typ informácií, ktoré sa dajú nájsť u každého zdroja. Ak by sme si dokázali zvoliť malé množstvo takýchto jednoznačných informácií mohli by sme návrh zobrazenia informácií upraviť na rozdelenie dvoch slúpcov. V pravom budú zobrazené len všeobecné hodnoty. Medzi všeobecné hodnoty, bude zaradený veľmi malý počet informácií. Budú to informácie, ktoré vždy účinne popisujú dané zdroje. Medzi takéto informácie patrí názov, alebo titulok. Ďalšou informáciou je náhľadový obrázok a popis. O týchto informáciách nemá zmysel diskutovať či sú podstatné pre zobrazenie.

Za účelom demonštrácie zmeny dôležitosti informácií, pri zmene konkrétnej zaradenia zdroja k typu triedy, bude v spodnej časti obrazovky zobrazená hierarchia tried. Hierarchia bude odpovedať aktuálne zobrazenému zdroju. Pričom bude možné, pomocou hierarchie meniť spôsob interpretácie informácií o zdroju podľa zvolenej triedy. Automaticky bude vždy aplikované zobrazenie pre najkrokretnejšiu triedu hierarchie. Schématický návrh vizuálnej podoby prehliadača je zobrazený na obrázku 3.1.





Obr. 3.1: Návrh dizajnu webovej aplikácie

### 3.3 Implementácia prehliadača

Implementácia prehliadača spočíva z dvoch základných častí. Ide o serverovú a klientsku časť. Pre vytvorenie servera sme použili jazyk Ruby a framework *Sinatra* a klient je založený na JavaScript technológií *React*. Server bude obsluhovať požiadavky klienta. Ako prvé budeme potrebovať vyriešiť základne zobrazenie prehliadača. Ruby poskytuje jedeno z rozhraní určených pre jednoduchú konfiguráciu webového servera, rozhranie *Rack*. Základnou konfiguráciou pre *Rack* je konfigurácia súboru *config.ru* v koreňovom priečinku aplikácie. Ide o vstupný súbor, ktorý mapuje URL požiadavky z webového prehliadača na server. Okrem iného je v tomto súbore zvyčajne riešene načítanie všetkých potrebných tried a konfigurácií.

#### 3.3.1 Architektúra servera

Pre lepšie pochopenie aké funkčnosti sú implementované na strane serveru, si v tejto časti uvedieme popis základných tried. Diagram tried servera je zobrazený na obrázku 3.2.

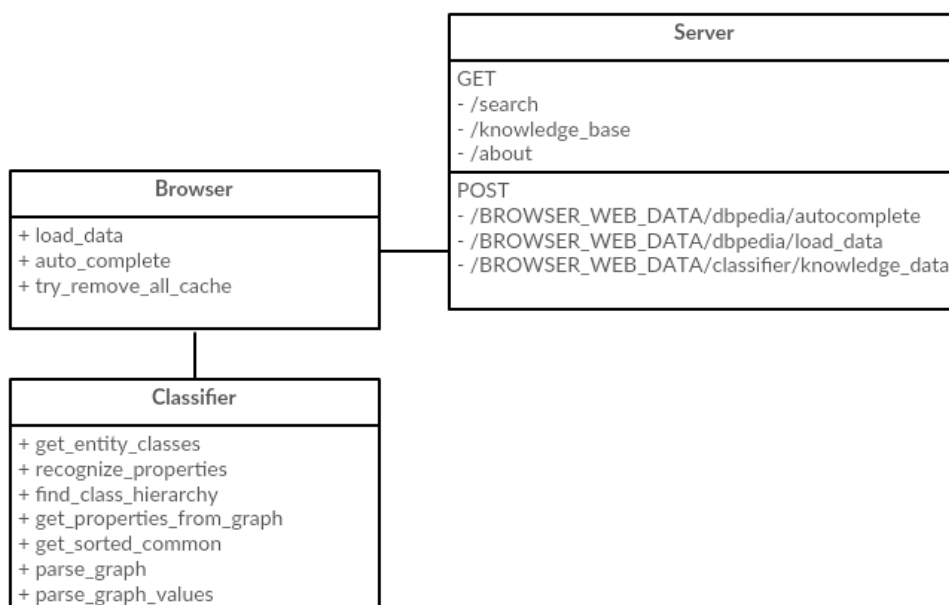
#### 3.3.2 Charakteristika tried servera

Hlavnou triedou servera je trieda *Server* ktorá definuje samotný server našej aplikácie. V rámci neho sú definované tri hlavné požiadavky typu POST. Vstupom pre všetky požiadavky je JSON so základnou štruktúrou zobrazenou na obrázku 3.3. Na prefixe adresy `/BROWSER_WEB_DATA/` sú dostupné tieto volania:

**dbpedia/autocomplete** Vstupom požiadavku je JSON základnej štruktúry s kľúčom *find\_value* a hodnotou reťazca, pre ktorý majú byť nájdené ná-

### 3. WEBOVÁ APLIKÁCIA

---



Obr. 3.2: Základný diagram tried servera

```
{
  "data": {
    "body": {
      KEY: VALUE
    }
  }
}
```

Obr. 3.3: Základná štruktúra JSON požiadavku na server

vrhy. Návratovou hodnotou je JSON základnej štruktúry s kľúčom *items* a hodnotou je pole objektov. Tieto objekty obsahujú kľúče a hodnoty *label* s reťazcom pre označenie zdroja a *uri* s reťazcom URI zdroja.

**dbpedia/load\_data** Vstupom požiadavku je JSON základnej štruktúry s kľúčom *dbpedia\_resource* a hodnotou reťazca, s URI zdroja, pre ktorý majú byť zobrazené dôležité informácie. Návratou hodnotou je JSON základnej štruktúry s kľúčmi a hodnotami *entity* s reťazcom najkonkrétnejšej triedy ontologie, do ktorej daný zdroj patrí, *entities* s zoradeným polom reťazcov hierarchie tried ontologie daného zdroja, *properties* s zoradeným polom výsledkov pre zobrazenie a *common* s zoradeným polom všeobecne dôležitých informácií pre zobrazenie.

**classifier/knowledge\_data** Vstupom požiadavku je prázdny JSON základnej štruktúry. Návratovou hodnotou je báza znalosti vo formáte JSON,

tak ako bola popísaná na obrázku 2.19.

V rámci triedy *Browser* sú definované 2 hlavné metódy ako implementácia obstarávania požiadavkov *dbpedia/autocomplete* a *dbpedia/load\_data* na servery. Trieda zároveň poskytuje použitie metód inštancie triedy *Classifier* pomocou atribútu *classifier*. Trieda zároveň vytvára dočasné súbory (angl. cache) s výsledkami jednotlivých zdrojov, pre zníženie záťaže a zrýchlenie prehľadávania.

Trieda *Classifier* rieši všetky potrebné operácie potrebné k získaniu tripletov požadovaného zdroja, aplikácií výberu dôležitých informácií a konštrukcie výsledku v podobe vhodnej pre zobrazenie na strane klienta. V rámci potrebných dotazov na stranu DBpedia endpointu využíva existujúcu triedu nástroja *SPARQLRequest*. V rámci tejto triedy nástroja boli pre potreby prehladača definované ďalšie dotazy:

**Získanie všetkých tried zdroja.** Vstupom pre dotaz je zdroj ako *resource*. Vráti sa všetky triedy aké ma zdroj priradené.

```
SELECT DISTINCT ?entity_class
WHERE {
  #{resource} a ?entity_class .
  ?entity_class a owl:Class .
}
```

Obr. 3.4: SPARQL dotaz pre získanie tried zdroja

**Získanie všetkých vlastností zdroja.** Vstupom pre dotaz je zdroj ako *resource*. Vráti sa všetky predikáty ako *predicate* s označením *predicate\_label* a hodnoty ako *value* s označením *value\_label*. Dotaz zároveň umožňuje filtrovanie na základe jazyka, parameter *lang*, pričom základnou hodnotou je *en*.

```
PREFIX dbo: <http://dbpedia.org/ontology/>
PREFIX dbp: <http://dbpedia.org/property/>
SELECT DISTINCT ?predicate, ?predicate_label, ?value, ?value_label
WHERE {
  { #{resource} ?predicate ?value . } UNION { ?value ?predicate #{
    resource} . }

  OPTIONAL{
    ?value rdfs:label ?value_label .
    FILTER (lang(?value_label) = '#{lang}')
  }

  ?predicate rdfs:label ?predicate_label .
  FILTER (lang(?predicate_label) = '#{lang}')
}
```

Obr. 3.5: SPARQL dotaz pre získanie vlastností zdroja

V rámci návratovej hodnoty pre získanie vlastnosti zdroja sú jednotlivé výsledky mapované na všetky triedy zdroja, ku ktorým patrí. Týmto spôsobom sa pri jednom požiadavku na stranu klienta dostanú všetky údaje potrebné k prepínaniu pod akou triedou sú zobrazené výsledky zdroja. Hlavným dôvodom tohto riešenia je fakt, že v procese získavania vlastnosti už metóda pozná všetky vlastnosti. Preto je veľmi jednoduché ich namapovať na všetky triedy. Týmto sa zároveň zrýchľuje prezentácia výsledkov na strane klienta.

#### 3.3.3 Architektúra klienta

Klient webovej aplikácie je postavený predovšetkým na technológií *JavaScriptu* s využitím *REACTu*. Kaskadové štýly sú vytvorené pomocou frameworku *Sass*. Pre potreby použitia ďalších balíčkov je použitý balíčkovacieho nástroj *npm*. JavaScriptové súbory sú písané s použitím *JSX* harmony. Použité sú aj niektoré konštrukcie vyplývajúce z definície *ECMAScript6*. Keďže je klientský projekt vzhľadom na použité technológie zložitejší je aj jeho štruktúra rozdelená na vývojovú a produkčnú časť. V rámci vývojovej časti sú predovšetkým jednotlivé *REACT* komponenty a kaskadové štýly. Pri tvorbe štýlov je použitý prístup návrhu *Mobile First* pre zabezpečenie podpory zobrazenia webovej aplikácie na rôznych zariadeniach. Medzi ďalšie použité frameworky patrí taktiež *Bootstrap*. JavaScriptové súbory sú pomocou nástroja *babel* kompilované na čistý interpretovateľný JavaScript. Výsledne súbory sú minifikované. Pre prácu s úlohami kompilácie a minifikácie je použitý nástroj *Grunt*. Zobrazenie základných komponent klienta je na obrázku 3.6.

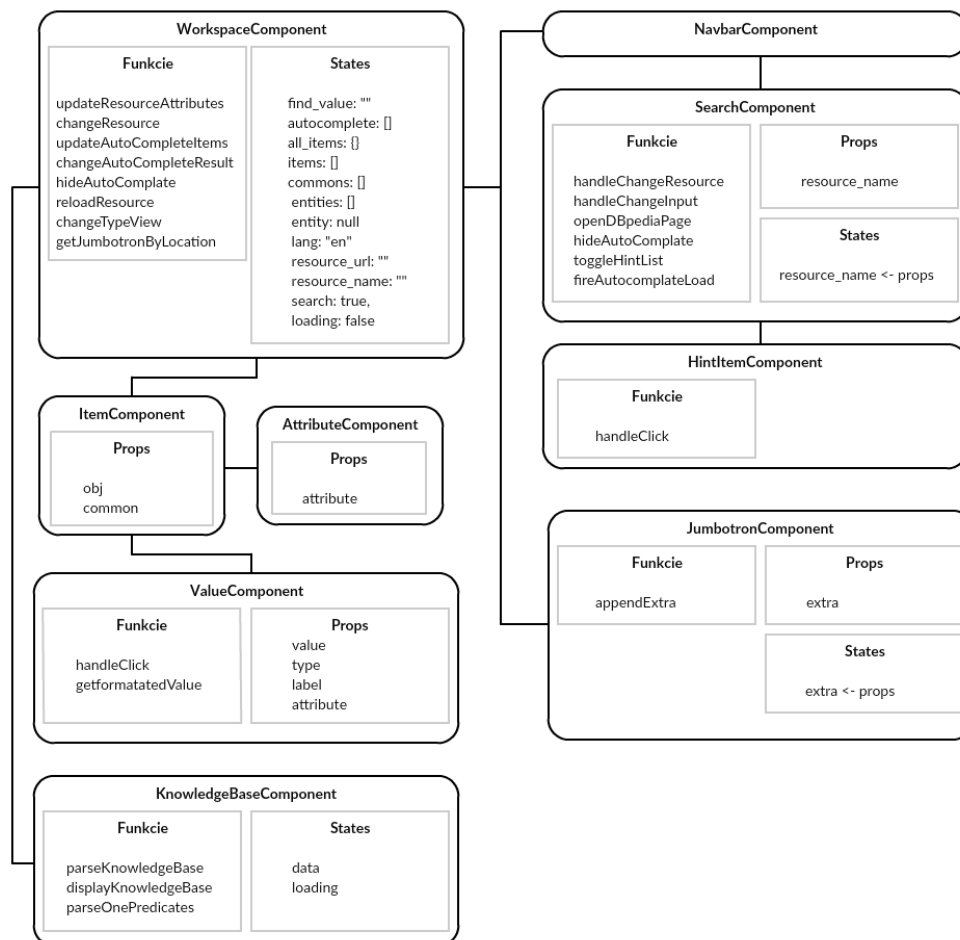
#### 3.3.4 Charakteristika komponent klienta

Každá z komponent uvedených na obrázku 3.6 zabezpečuje vykreslenie určite časti prehliadača.

**WorkspaceComponent** Predstavuje celú oblasť prehliadača. Zabezpečuje pripájanie a vykreslenie podradených komponent na základe aktuálnych potrieb. V prípade, že nie sú zobrazené informácie žiadneho zdroja, je zobrazená informačná komponenta *JumbotronComponent*. V spodnej časti je zobrazená hierarchia tried zdroja. Zároveň je umožnené prepínanie medzi danými triedami.

**NavbarComponent** Definuje len spôsob vykreslenia navigačnej lišty. V navigačnej lište sú prepnutia stránok a vyhľadávacie pole *SearchComponent*.

**SearchComponent** Komponenta umožňuje hlavnú interakciu s používateľom. Obsahuje vstupné pole pomocou ktorého vyhľadáva zdroj na *DBpedia*. Poskytuje zároveň našepkávač pre zadaný výraz. Našepkávač je vytvorený zo získaných možností na základe callback volania požiadavkom *dbpedia/autocomplete* na server.



Obr. 3.6: Základná štruktúra komponent klienta

**HintItemComponent** Predstavuje jednu položku z našepkávača. Zobrazuje textové označenie, ktoré je priradené určitej URI zdroja. Ak si používateľ zvolí zdroj na základe tejto komponenty, premietne sa URI zdroja do komponenty *SearchComponent* a automaticky sa spustí vyhľadávanie.

**KnowledgeBaseComponent** Zabezpečuje vykreslenie bázy znalosti. Sú zobrazené príslušne triedy ontológie a v rámci nich predikáty spolu s ich hodnotením dôležitosti.

**ItemComponent** Predstavuje jednotlivé záznamy vlastnosti určené pre vykreslenie. V rámci komponenty sa rozlišuje, či sa jedná o všeobecne dôležitú komponentu tzv. *common*. Ak je to *common* bude vykreslená v pravej časti stránky, pričom nebude vykreslený názov vlastnosti, ale len samotná hodnota. V prípade, že je to bežná vlastnosť bude vykreslená

### 3. WEBOVÁ APLIKÁCIA

---

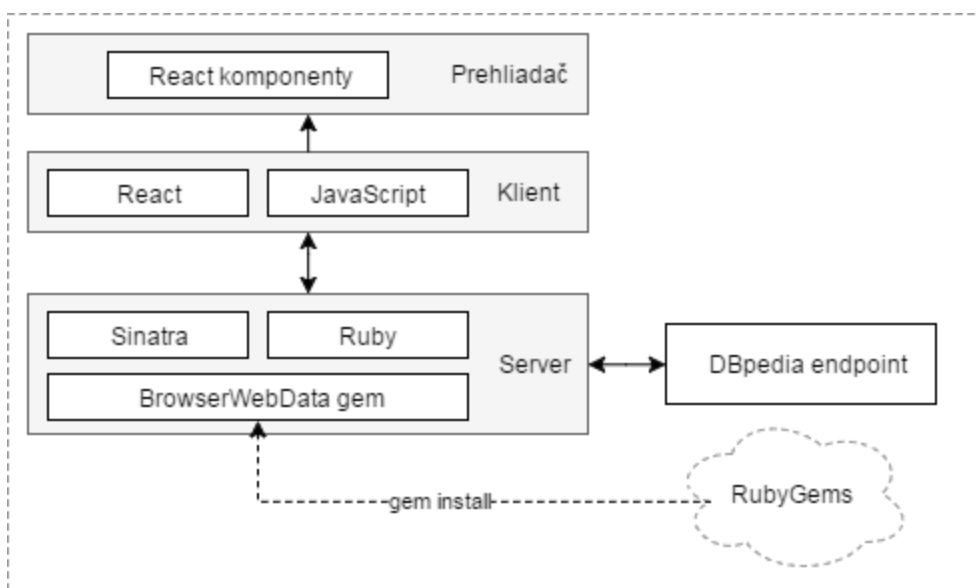
v hlavnej časti. Komponenta vykresľuje označenie vlastnosti pomocou komponenty *AttributeComponent* a hodnoty pomocou *ValueComponent*.

**AttributeComponent** Zabezpečuje spôsob vykreslenie označenia vlastnosti.

**ValueComponent** Zabezpečuje spôsob vykreslenie hodnôt vlastnosti. Pre jednu vlastnosť môže byť priradené viacero hodnôt, v takom prípade sa vždy zobrazí len prvých 5 hodnôt. Všetky ďalšie sú skryté. Ich zobrazenie je možné pomocou tlačidla. Vykreslenie hodnôt zároveň aplikuje formátovanie na základe prideleného typu (napr. číslo, dátum a ďalšie).

## 3.4 Architektúra prehliadača

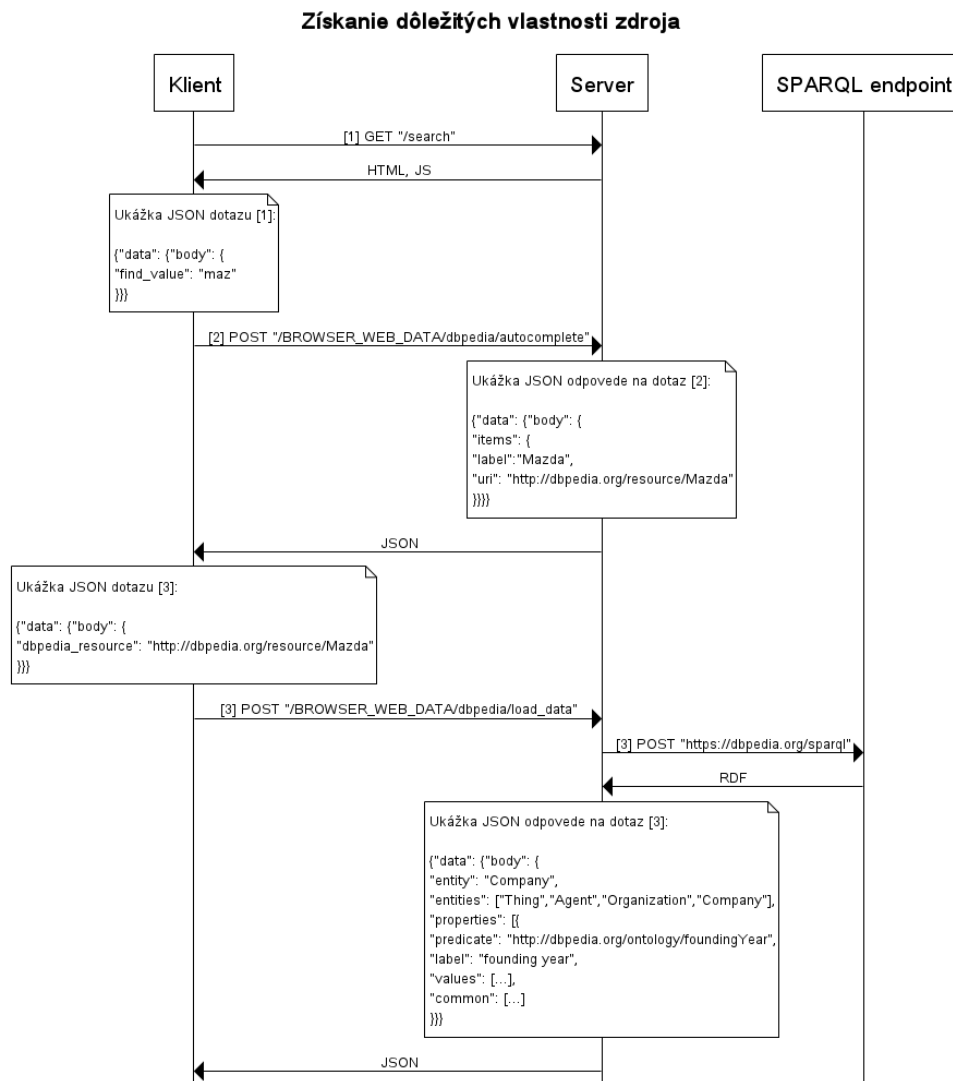
Prehliadač obsahuje 3 základne stránky. Úvodnou a hlavnou stránkou je */search*, ktorá umožňuje prehľadanie informácií zdrojov. Ďalšou stránkou je */knowledge\_base*. Na tejto stránke je prehľadnou formou zobrazená používaná báza znalosti. Poslednou stránkou je stránka */about* so základnými informáciami o prehliadači ako aj celom projekte. Základná architektúra celej webovej aplikácie je znázornená na obrázku 3.7.



Obr. 3.7: Základná architektúra webového prehliadača

Základná komunikácia medzi klientom a serverom za účelom získania dôležitých informácií je zobrazená na obrázku 3.8. Webová aplikácia bola, okrem lokálneho testovania, nasadená pomocou služby *Heroku*. Tento proces bol riadený pomocou technológie *git*. Repozitár s aplikáciou, ktorý bol vytvorený v

službe Heroku sa pri operácii *push* nasadil ako celá webová aplikácia. Tým sa zároveň overilo jej reálne použitie v praxi.



Obr. 3.8: Základná komunikácia klienta so serverom





## Experimenty a vyhodnotenie

V tejto kapitole sa pozrieme na získané výsledky. Hlavným výsledkom je vygenerovaná báza znalosti. Obsahuje zoradené predikáty na základe ich skóre. Na základe získaných hodnotení predikátov vykonáme porovnanie voči iným existujúcim riešeniam. Budeme hlavne porovnávať pozície jednotlivých predikátov.

Pozičné porovnanie si vyžiadalo zvoliť prístup akým budú jednotlivé umiestnenia našich výsledkov porovnané voči iným riešeniam. Porovnanie bude vychádzať z prechodu zoradeným zoznamom predikátov iných riešení. Pre každý predikát bude nájdená pozícia v rámci našej bázy znalosti. Pri niektorých riešeniach bude nutné premapovanie predikátov. Znamená to, aký predikát zdrojov DBpedia odpovedá svojím popisom predikátu daného riešenia. Pre stanovenie skóre identity umiestnenia sme si definovali vzťah, v ktorom je znižovaná presnosť v miere relatívneho vzdialenia sa nášho predikátu. Nech je  $c_o$  počet overovaných predikátov cudzieho riešenia,  $i_{p_a}$  index umiestnenia predikátu  $p$  cudzieho riešenia a  $i_{p_b}$  index umiestnenia predikátu  $p$  nášho riešenia potom skóre miery identity umiestnenia nášho predikátu v báze znalosti definujeme pomocou vzťahu znázornenom na obrázku 4.1.

$$S_p = 1 - \frac{|i_{p_a} - i_{p_b}|}{c_o}$$

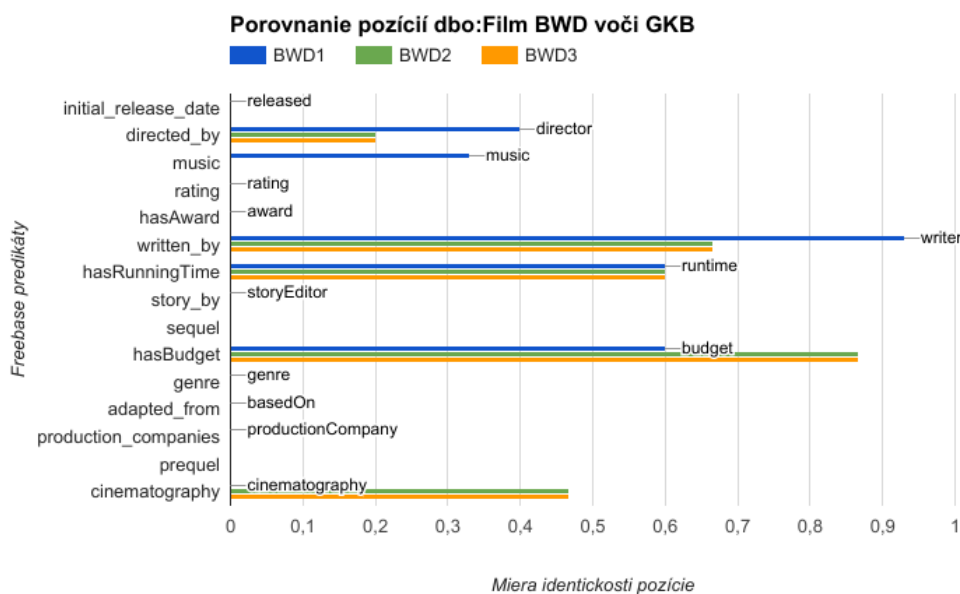
Obr. 4.1: Definícia miery identity pozície predikátu  $p$ .

Všetky porovnania sú realizované jednotným spôsobom. Porovnania sú reprezentované pomocou grafov. Prvým typom grafu je graf porovnania miery identity pozície predikátov. Druhým je graf porovnania miery identity hodnotenia dôležitosti predikátov. Za kompletne porovnanie je možné považovať priemernú hodnotu z identity pozícií predikátov. Budeme ho preto označovať pod názvom *kompletné hodnotenie*, budeme ho uvádzať ako percentá. V rámci experimentu sme zároveň porovnávali 3 druhy vygenerovaných báz znalosti. Prvou je kompletná báza znalosti so všetkým procesmi. Druhou

#### 4. EXPERIMENTY A VYHODNOTENIE

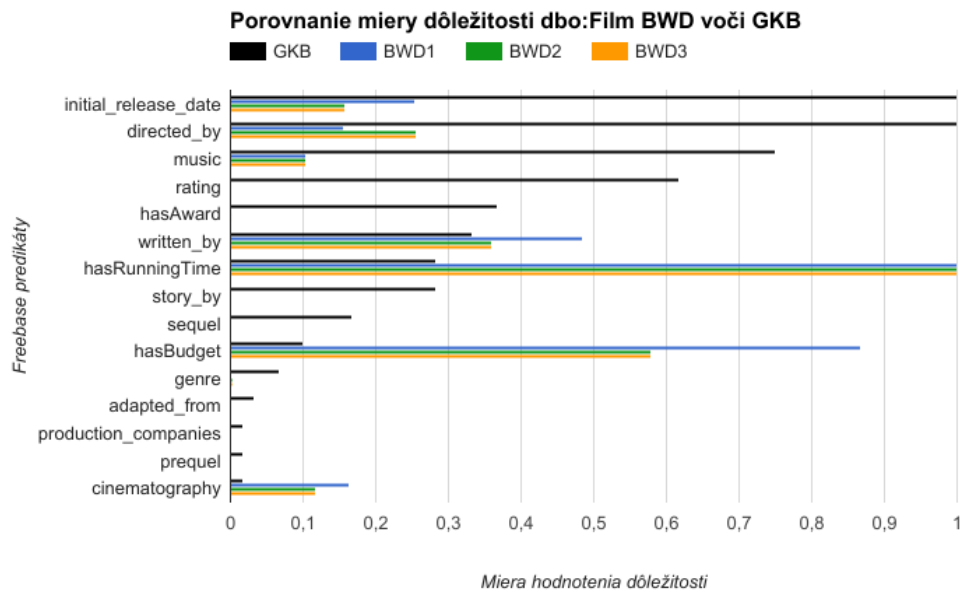
je báza znalosti vygenerovaná na základe iba striktných predikátov a identifikácií identických predikátov. Poslednou je pri použití striktných predikátov a triedných predikátov (podobne ako v hlavnej báze znalosti), ale bez identifikácie identických predikátov. Porovnávalo sa prvých 15 predikátov. Výsledky našej práce sú označené ako BrowserWebData (ďalej ako BWD).

Prvým overením výsledkov je porovnanie voči projektu založenom na summarizácii entít pomocou vedomostnej hry (popísaná v kapitole 1.5.2). Riešenie používalo predikáty z datasetu Freebase. Prehľadná tabuľka výsledkov danej práce je na obrázku 1.3. Tie niesú zhodné s predikátmi použitými v DBpedii. Vytvorili sme si preto mapovanie ktoré je zobrazené v tabuľke B.1 prílohy B. Porovnávanými výsledkami sú hodnotenia triedy `http://dbpedia.org/ontology/Film`. Vo svojej práci overovali výsledky voči *Google knowledge base*. Preto sme pomocou ich práce realizovali porovnanie voči ich výsledkom pod označením UBES a Google knowledge base pod označením GKB. Porovnanie pozícií BWD voči GKB je na obrázku 4.2, porovnanie miery dôležitosti BWD voči GKB je na obrázku 4.3. Kompletným hodnotením je 22%, je rovnakým pre všetky 3 druhy uvažovaných báz.

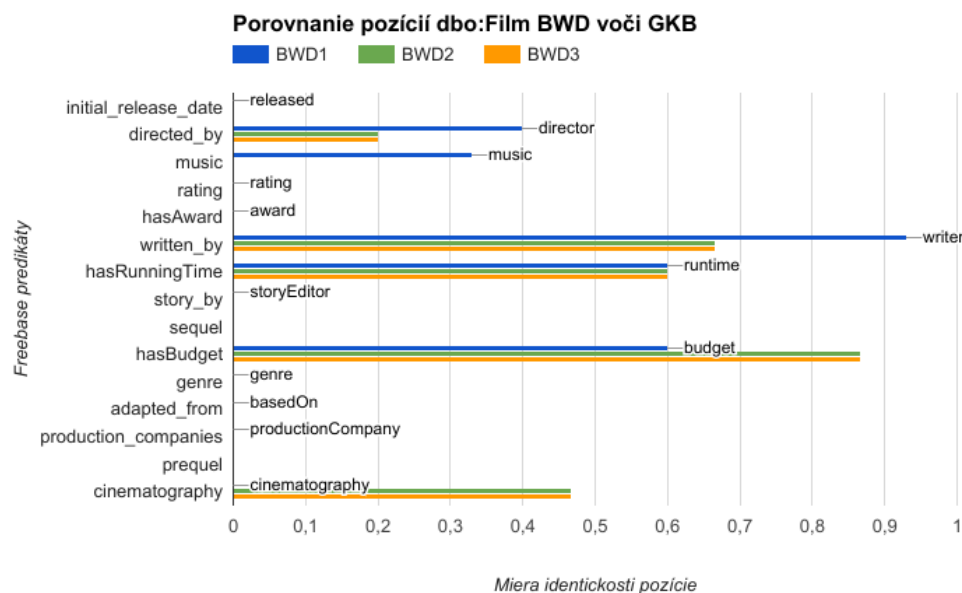


Obr. 4.2: Graf porovnania pozícií BWD voči GKB pre triedu Film

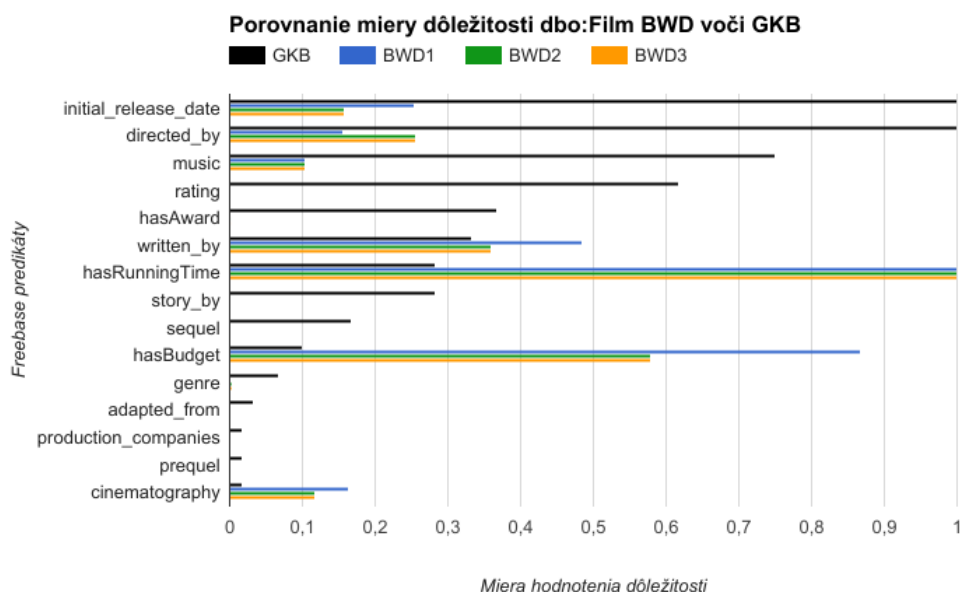
Porovnanie pozícií BWD voči UBES je na obrázku 4.4, porovnanie miery dôležitosti BWD voči UBES je na obrázku 4.5. Kompletným hodnotením je 37%, pre kompletnú bázu znalostí a 19% pre ďalšie 2 z experimentálnych báz.



Obr. 4.3: Graf porovnania miery dôležitosti BWD voči GKB pre triedu Film



Obr. 4.4: Graf porovnania pozícií BWD voči UBES pre triedu Film



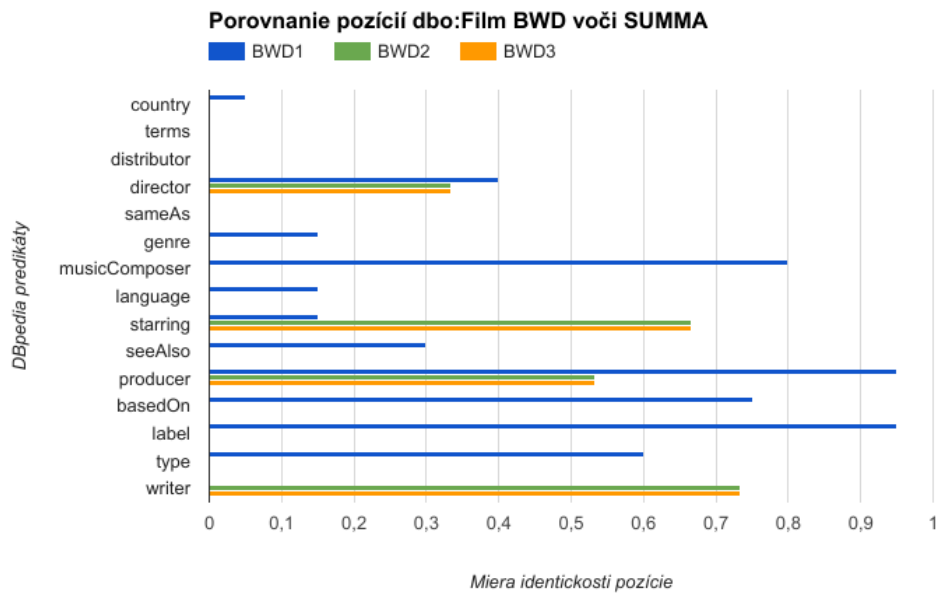
Obr. 4.5: Graf porovnania miery dôležitosti BWD voči UBES pre triedu Film

Poslednou prácou, ktorej výsledky porovnáme s našimi je API rozhranie pre sumarizáciu (popísaná v kapitole 1.5.3). Tieto výsledky sú ďalej označené ako SUMMA. Pomocou ich API sme postupne získali hodnotenia zvolených zdrojov. Porovnávali sme viacero tried. Zvolili sme si rovnakú triedu ako v prvom porovnaní <http://dbpedia.org/ontology/Film>. Ďalej sme si zvolili 2 triedy <http://dbpedia.org/ontology/Person> a jej podtriedu <http://dbpedia.org/ontology/President>. Porovnanie podtriedy je za účelom demonštrácie zmeny dôležitosti predikátov pri viac špecifickom zaradení zdroja.

Porovnanie pozícií triedy <http://dbpedia.org/ontology/Film> BWD voči SUMMA je na obrázku 4.6, porovnanie miery dôležitosti BWD voči SUMMA je na obrázku 4.7. Kompletným hodnotením je 35%, pre kompletnú bázu znalostí a 15% pre ďalšie 2 z experimentálnych báz.

Porovnanie pozícií triedy <http://dbpedia.org/ontology/Person> BWD voči SUMMA je na obrázku 4.8, porovnanie miery dôležitosti BWD voči SUMMA je na obrázku 4.9. Kompletným hodnotením je 12%, pre kompletnú bázu znalostí, 14% pre striktných predikátov a 11% u bázy bez použitia identifikácie identických predikátov.

Porovnanie pozícií triedy <http://dbpedia.org/ontology/President> BWD voči SUMMA je na obrázku 4.10, porovnanie miery dôležitosti BWD voči SUMMA je na obrázku 4.11. Kompletným hodnotením je 12%, pre kompletnú bázu znalostí a striktné predikáty. Pre bázu bez použitia identifikácie identických

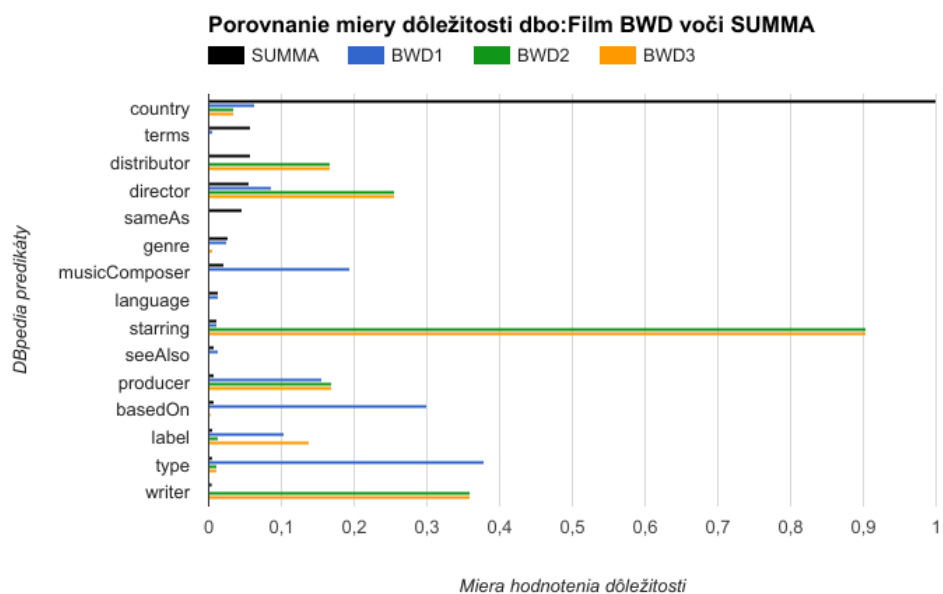


Obr. 4.6: Graf porovnania pozícií BWD voči SUMMA pre triedu Film

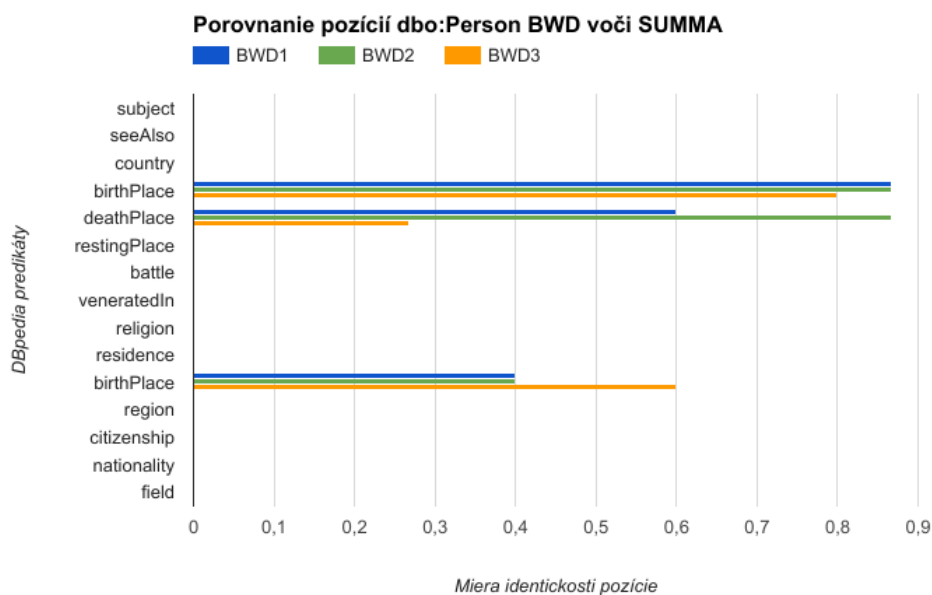
kých predikátov 8%.

Zhrnutím celého porovnania je skutočnosť, že pri použití kompletnej bázy znalosti získavame identickejšie pozície voči iným riešeniam. Zároveň platí, že bez použitia identifikácie identických predikátov sa znižuje presnosť pozícií. Aj napriek nízkym kompletným hodnoteniam pokladáme naše riešenia za dostatočne úspešne. Zároveň napríklad u výsledkov SUMMA sú vo výsledkoch aj duplicity, ktoré negatívne vplyvajú na porovnanie pozícií.

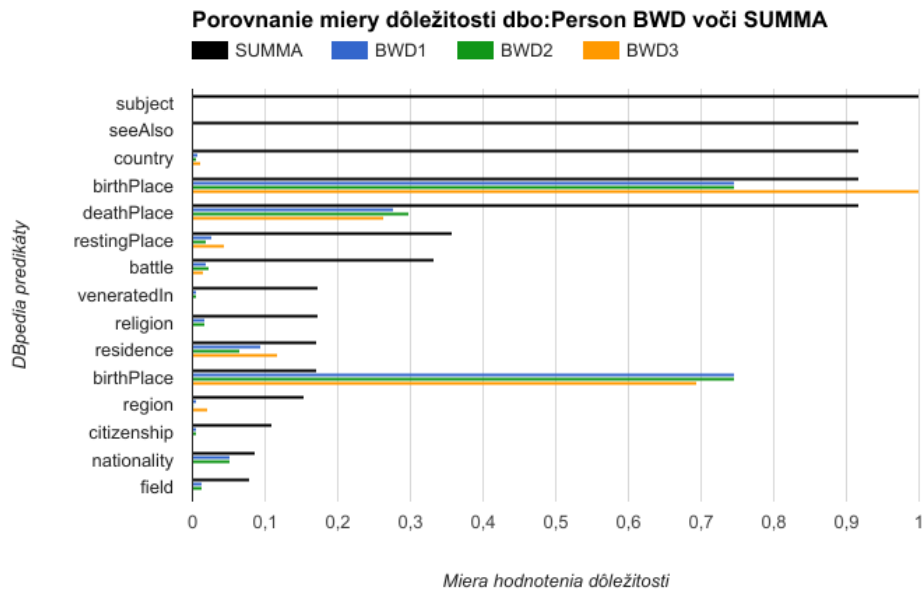
#### 4. EXPERIMENTY A VYHODNOTENIE



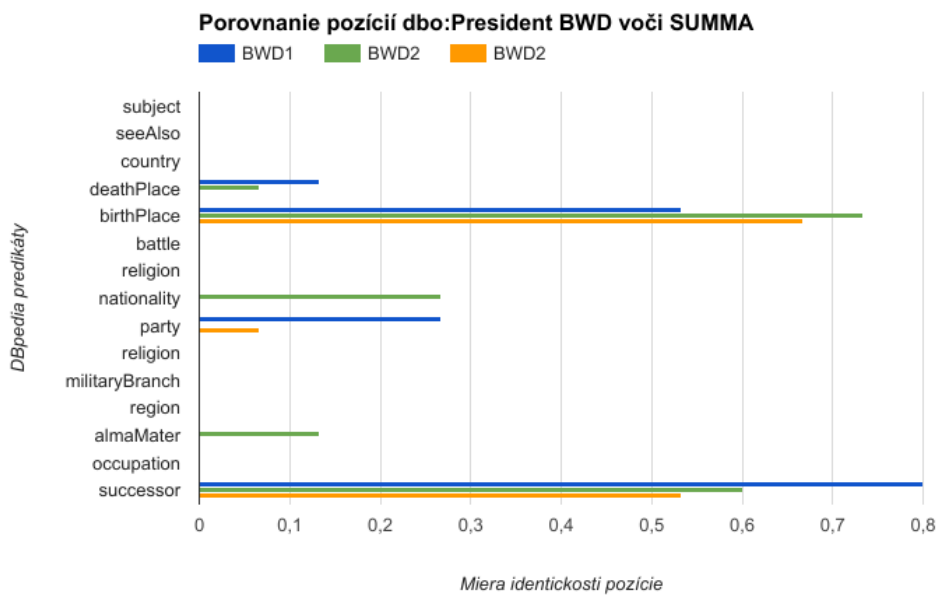
Obr. 4.7: Graf porovnania miery dôležitosti BWD voči SUMMA pre triedu Film



Obr. 4.8: Graf porovnania pozícií BWD voči SUMMA pre triedu Person



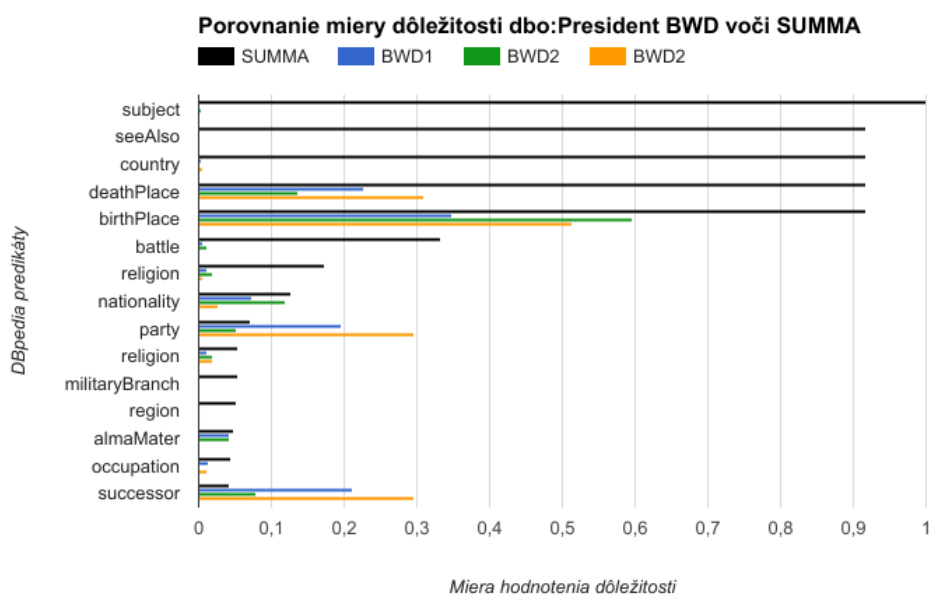
Obr. 4.9: Graf porovnania miery dôležitosti BWD voči SUMMA pre triedu Person



Obr. 4.10: Graf porovnania pozícií BWD voči SUMMA pre triedu President

#### 4. EXPERIMENTY A VYHODNOTENIE

---



Obr. 4.11: Graf porovnania miery dôležitosti BWD voči SUMMA pre triedu President



---

## Záver

V závere práce by sme radi zhodnotili získané výsledky práce. Pozreli sme sa na problematiku sumarizácie entít, identifikovali podobné riešenia. V časti návrhu sme podrobnejšie rozobrali a navrhli vlastný model pre sumarizáciu entít na základe NIF DBpedia abstract datasetu. Model bol aplikovaný pri vytvorení nástroja. Nástroj sme publikovali pomocou verejného manažéra balíčkov *RubyGems*. Je dostupný na adrese [https://rubygems.org/gems/browser\\_web\\_data\\_entity\\_sumarization](https://rubygems.org/gems/browser_web_data_entity_sumarization). Spôsob ako použiť nástroj sme v krátkosti uviedli v tejto práci. Podrobnejší návod je súčasťou publikovaného kódu nástroja spolu s našimi výsledkami na adrese [https://github.com/MarekFiltes/entity\\_sumarization](https://github.com/MarekFiltes/entity_sumarization). Navrhli a implementovali sme webovú aplikáciu založenú na *Ruby* frameworku *Sinatra*. Klient webovej aplikácie je založený na technológii *React*. Webová aplikácia aplikuje získane výsledky z bázy znalosti. Poskytli sme zároveň vedľajšie výsledky práce. Najpodstatnejším je zoznam identických predikátov. Pomocou neho je možné znížiť mieru šumových duplicitných predikátov.

Medzi zlepšenia práce patrí rozšírenie použitých tried ontológií o triedy z iných ontológií, napríklad triedy *yago*. Zlepšenie procesu identifikácie identických predikátov. Často sa aktuálne nezaradia ako identické, predikáty, ktorých hodnoty sú rozdielným spôsobom priradené, ale majú rovnaký význam. Možným rozšírením práce do budúca je poskytnutie Web API pre získanie zoznamu dôležitých predikátov a ich hodnôt.

Všetky ciele, ktoré sme si v úvode stanovili pokladáme za splnené. Naše riešenie poskytuje model a nástroj určený pre zlepšenie prehľadnosti informácií zdrojov na DBpedií. Rieši problém identifikácie, ktoré z informácií sú dôležitejšie než iné. Pri prehliadaní informácií na základe našej bázy znalosti sa znižuje miera náročnosti identifikácie informácií používateľom. Zlepšuje sa použitie informácií zdrojov DBpedia pre bežného používateľa, ktorý nepotrebuje všetky informácie. Výsledná báza znalostí má široký záber. Obsahuje predikáty mapované na všetky triedy DBpedia ontológie.



---

## Literatúra

- [1] Mani, I. *Automatic Summarization*. Natural Language Processing, John Benjamins Publishing Company, 2001, ISBN 9789027299109. Available from: <https://books.google.cz/books?id=CSHxU0fb5bwC>
- [2] Page, L.; Brin, S.; Motwani, R.; et al. The PageRank Citation Ranking: Bringing Order to the Web. Technical Report 1999-66, Stanford InfoLab, November 1999, previous number = SIDL-WP-1999-0120. Available from: <http://ilpubs.stanford.edu:8090/422/>
- [3] Thalhammer, A.; Stadtmüller, S. *SUMMA: A Common API for Linked Data Entity Summaries*. Cham: Springer International Publishing, 2015, ISBN 978-3-319-19890-3, pp. 430–446, doi:10.1007/978-3-319-19890-3\_28. Available from: [http://dx.doi.org/10.1007/978-3-319-19890-3\\_28](http://dx.doi.org/10.1007/978-3-319-19890-3_28)
- [4] LLC, M. NIF Abstract Datasets. 2016, [cit. 2017-05-01]. Available from: <http://wiki.dbpedia.org/nif-abstract-datasets>
- [5] Kahn, B. K.; Strong, D. M.; Wang, R. Y. Information Quality Benchmarks: Product and Service Performance. *Commun. ACM*, volume 45, no. 4, Apr. 2002: pp. 184–192, ISSN 0001-0782, doi: 10.1145/505248.506007. Available from: <http://doi.acm.org/10.1145/505248.506007>
- [6] Thalhammer, A.; Knuth, M.; Sack, H. *Evaluating Entity Summarization Using a Game-Based Ground Truth*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, ISBN 978-3-642-35173-0, pp. 350–361, doi:10.1007/978-3-642-35173-0\_24. Available from: [http://dx.doi.org/10.1007/978-3-642-35173-0\\_24](http://dx.doi.org/10.1007/978-3-642-35173-0_24)
- [7] DBpedia version 2016-04. 2016, [cit. 2017-05-02]. Available from: <http://wiki.dbpedia.org/dbpedia-version-2016-04/OnlineAccess>

- [8] Cover, T. M.; Thomas, J. A. *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*. Wiley-Interscience, 2006, ISBN 0471241954.
- [9] Šlapák, O. Data, informace, znalosti. *E-LOGOS ELECTRONIC JOURNAL FOR PHILOSOPHY*, 2003.
- [10] Hoffmannová, J. *Stylistika a: současná situace stylistiky*. Trizonia, 1997, ISBN 9788085573671. Available from: <https://books.google.cz/books?id=2qj1AAAAMAAJ>
- [11] Stalnaker, R. *Context and Content: Essays on Intentionality in Speech and Thought*. Oxford cognitive science series, Oxford University Press, 1999, ISBN 9780198237082. Available from: <https://books.google.cz/books?id=0hhMRF2dz0IC>
- [12] Jaromír, J. *Psychologické základy verbální komunikace*. 2015, ISBN 9788024798509.
- [13] Ding, L.; Pan, R.; Finin, T.; et al. *Finding and Ranking Knowledge on the Semantic Web*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2005, ISBN 978-3-540-32082-1, pp. 156–170, doi:10.1007/11574620\_14. Available from: [http://dx.doi.org/10.1007/11574620\\_14](http://dx.doi.org/10.1007/11574620_14)
- [14] Lachica, R.; Karabeg, D.; Rudan, S. Quality, relevance and importance in information retrieval with fuzzy semantic networks. *Proc. TMRA*, 2008.
- [15] Kagolovsky, Y.; Mohr, J. R. A new approach to the concept of "relevance in information retrieval (IR)". *Studies in health technology and informatics*, no. 1, 2001: pp. 348–352.
- [16] Grešková, M. Empirický výskum interakcie človeka s agentom. *Nová paradigma spracovania a využívania informácií*, 2007: p. 24.
- [17] Radev, D. R.; Hovy, E.; McKeown, K. Introduction to the Special Issue on Summarization. *Comput. Linguist.*, volume 28, no. 4, Dec. 2002: pp. 399–408, ISSN 0891-2017, doi:10.1162/089120102762671927. Available from: <http://dx.doi.org/10.1162/089120102762671927>
- [18] Steinberger, J.; Ježek, K. Text summarization and singular value decomposition. In *International Conference on Advances in Information Systems*, Springer, 2004, pp. 245–254.
- [19] *Sémantika programovacích jazyků*. Skripta UK, 1997, ISBN 80-7184-327-X.

- 
- [20] Berners-Lee, T.; Hendler, J.; Lassila, O. The Semantic Web. *Scientific American*, volume 284, no. 5, 2001: pp. 34–43, ISSN 0036-8733. Available from: <http://www.scientificamerican.com/article.cfm?id=the-semantic-web>
- [21] RDF Primer. 2014, [cit. 2017-04-14]. Available from: <https://www.w3.org/TR/2004/REC-rdf-primer-20040210/>
- [22] SPARQL Query Language for RDF. 2008, [cit. 2017-04-16]. Available from: <https://www.w3.org/TR/rdf-sparql-query/>
- [23] Svátek, V. Ontologie a WWW. In *Sborník konference Datakon*, 2002, pp. 27–55.
- [24] Berners-Lee, T. Linked Data. 2006, [cit. 2017-04-17]. Available from: <http://www.w3.org/DesignIssues/LinkedData.html>
- [25] Bizer, C.; Heath, T.; Berners-Lee, T. Linked data - the story so far. *Int. J. Semantic Web Inf. Syst.*, volume 5, no. 3, 2009: p. 1–22.
- [26] Wikipedie. 2014, [cit. 2017-04-16]. Available from: <https://cs.wikipedia.org/wiki/Wikipedie>
- [27] Auer, S.; Bizer, C.; Kobilarov, G.; et al. Dbpedia: A nucleus for a web of open data. *The semantic web*, 2007: pp. 722–735.
- [28] DBpedia. 2017, [cit. 2017-05-02]. Available from: <https://en.wikipedia.org/wiki/DBpedia>
- [29] Cheng, G.; Tran, T.; Qu, Y. *RELIN: Relatedness and Informativeness-Based Centrality for Entity Summarization*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, ISBN 978-3-642-25073-6, pp. 114–129, doi:10.1007/978-3-642-25073-6\_8. Available from: [http://dx.doi.org/10.1007/978-3-642-25073-6\\_8](http://dx.doi.org/10.1007/978-3-642-25073-6_8)
- [30] Brümmer, M.; Dojchinovski, M.; Hellmann, S. DBpedia Abstracts: A Large-Scale, Open, Multilingual NLP Training Corpus. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, edited by N. C. C. Chair; K. Choukri; T. Declerck; M. Grobelnik; B. Maegaard; J. Mariani; A. Moreno; J. Odijk; S. Piperidis, Paris, France: European Language Resources Association (ELRA), May 2016. Available from: [https://svn.aksw.org/papers/2016/LREC\\_DBpedia\\_Abstracts/public.pdf](https://svn.aksw.org/papers/2016/LREC_DBpedia_Abstracts/public.pdf)
- [31] Hellmann, S.; Lehmann, J.; Auer, S. NIF: An ontology-based and linked-data-aware NLP Interchange Format. *Working Draft*, 2012.

- [32] Hellmann, S.; Lehmann, J.; Auer, S.; et al. Integrating NLP using linked data. In *International Semantic Web Conference*, Springer, 2013, pp. 98–113.
- [33] Accessing the DBpedia Data Set over the Web. 2016, [cit. 2017-05-01]. Available from: <http://wiki.dbpedia.org/OnlineAccess>
- [34] Auer, S.; Bizer, C.; Kobilarov, G.; et al. *DBpedia: A Nucleus for a Web of Open Data*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2007, ISBN 978-3-540-76298-0, pp. 722–735, doi:10.1007/978-3-540-76298-0\_52. Available from: [http://dx.doi.org/10.1007/978-3-540-76298-0\\_52](http://dx.doi.org/10.1007/978-3-540-76298-0_52)
- [35] Thalhammer, A.; Rettinger, A. PageRank on Wikipedia: Towards General Importance Scores for Entities. In *The Semantic Web: ESWC 2016 Satellite Events, Heraklion, Crete, Greece, May 29 – June 2, 2016, Revised Selected Papers*, Cham: Springer International Publishing, Oct. 2016, ISBN 978-3-319-47602-5, pp. 227–240, doi:10.1007/978-3-319-47602-5\_41. Available from: [http://dx.doi.org/10.1007/978-3-319-47602-5\\_41](http://dx.doi.org/10.1007/978-3-319-47602-5_41)
- [36] DBpedia 2016-04 Statistics. 2016, [cit. 2017-05-03]. Available from: <http://wiki.dbpedia.org/dbpedia-2016-04-statistics>
- [37] About Ruby. 2017, [cit. 2017-05-05]. Available from: <https://www.ruby-lang.org/en/about>
- [38] What is gem. 2017, [cit. 2017-05-05]. Available from: <http://guides.rubygems.org/what-is-a-gem>

## Obsah priloženého CD

	readme.txt.....	stručný popis obsahu CD
	exe.....	adresár so spustiteľnou formou implementácie
	src	
	impl .....	zdrojové kódy implementácie
	thesis.....	zdrojová forma práce vo formáte L <sup>A</sup> T <sub>E</sub> X
	text .....	text práce
	thesis.pdf .....	text práce vo formáte PDF





## **Doplňujúce informácie testovania**

Tabuľka B.1: Tabuľka mapovania Freebase predikatov na DBpedia predikáty

<b>Predikáty Freebase</b>	<b>Predikáty DBpedia</b>
<a href="http://rdf.freebase.com/ns/film.film.initial_release_date">http://rdf.freebase.com/ns/film.film.initial_release_date</a>	<a href="http://dbpedia.org/ontology/releaseDate">http://dbpedia.org/ontology/releaseDate</a>
<a href="http://rdf.freebase.com/ns/film.film.directed_by">http://rdf.freebase.com/ns/film.film.directed_by</a>	<a href="http://dbpedia.org/ontology/director">http://dbpedia.org/ontology/director</a>
<a href="http://rdf.freebase.com/ns/film.film.music">http://rdf.freebase.com/ns/film.film.music</a>	<a href="http://dbpedia.org/property/music">http://dbpedia.org/property/music</a>
<a href="http://rdf.freebase.com/ns/film.film.rating">http://rdf.freebase.com/ns/film.film.rating</a>	<a href="http://dbpedia.org/ontology/rating">http://dbpedia.org/ontology/rating</a>
<a href="http://test.com/#hasAward">http://test.com/#hasAward</a>	<a href="http://dbpedia.org/ontology/award">http://dbpedia.org/ontology/award</a>
<a href="http://rdf.freebase.com/ns/film.film.written_by">http://rdf.freebase.com/ns/film.film.written_by</a>	<a href="http://dbpedia.org/ontology/writer">http://dbpedia.org/ontology/writer</a>
<a href="http://test.com/#hasRunningTime">http://test.com/#hasRunningTime</a>	<a href="http://dbpedia.org/ontology/runtime">http://dbpedia.org/ontology/runtime</a>
<a href="http://rdf.freebase.com/ns/film.film.story_by">http://rdf.freebase.com/ns/film.film.story_by</a>	
<a href="http://rdf.freebase.com/ns/film.film.sequel">http://rdf.freebase.com/ns/film.film.sequel</a>	
<a href="http://test.com/#hasBudget">http://test.com/#hasBudget</a>	<a href="http://dbpedia.org/ontology/budget">http://dbpedia.org/ontology/budget</a>
<a href="http://rdf.freebase.com/ns/film.film.genre">http://rdf.freebase.com/ns/film.film.genre</a>	<a href="http://dbpedia.org/ontology/genre">http://dbpedia.org/ontology/genre</a>
<a href="http://rdf.freebase.com/ns/media_common.adaptation.adapted_from">http://rdf.freebase.com/ns/media_common.adaptation.adapted_from</a>	<a href="http://dbpedia.org/ontology/basedOn">http://dbpedia.org/ontology/basedOn</a>
<a href="http://rdf.freebase.com/ns/film.film.production_companies">http://rdf.freebase.com/ns/film.film.production_companies</a>	<a href="http://dbpedia.org/ontology/productionCompany">http://dbpedia.org/ontology/productionCompany</a>
<a href="http://rdf.freebase.com/ns/film.film.prequel">http://rdf.freebase.com/ns/film.film.prequel</a>	
<a href="http://rdf.freebase.com/ns/film.film.cinematography">http://rdf.freebase.com/ns/film.film.cinematography</a>	<a href="http://dbpedia.org/ontology/cinematography">http://dbpedia.org/ontology/cinematography</a>
<a href="http://rdf.freebase.com/ns/film.film.subjects">http://rdf.freebase.com/ns/film.film.subjects</a>	<a href="http://purl.org/dc/elements/1.1/subject">http://purl.org/dc/elements/1.1/subject</a>
<a href="http://rdf.freebase.com/ns/film.film.film_series">http://rdf.freebase.com/ns/film.film.film_series</a>	<a href="http://dbpedia.org/ontology/series">http://dbpedia.org/ontology/series</a>
<a href="http://rdf.freebase.com/ns/film.film.film_festivals">http://rdf.freebase.com/ns/film.film.film_festivals</a>	
<a href="http://rdf.freebase.com/ns/film.film.film_casting_director">http://rdf.freebase.com/ns/film.film.film_casting_director</a>	
<a href="http://rdf.freebase.com/ns/film.film.featured_film_locations">http://rdf.freebase.com/ns/film.film.featured_film_locations</a>	


## Ukážky webovej aplikácie

Search Knowledge base About

**Birth Date**  
1946-08-19

**Birth Place**  
[United States](#)  
[Arkansas](#)  
[Hope, Arkansas](#)

**Term Start**  
1977-01-03  
1979-01-09  
1983-01-11  
1993-01-20



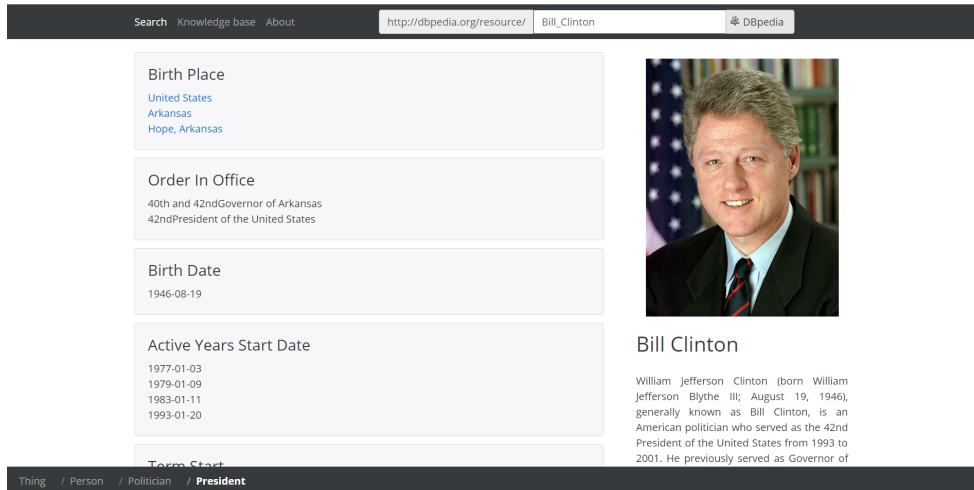
**Bill Clinton**

William Jefferson Clinton (born William Jefferson Blythe III; August 19, 1946), generally known as Bill Clinton, is an American politician who served as the 42nd President of the United States from 1993 to 2001. He previously served as Governor of

Thing / **Person** / Politician / President

Obr. C.1: Ukážka zobrazenia zdroja triedy Person

## C. UKÁŽKY WEBOVEJ APLIKÁCIE

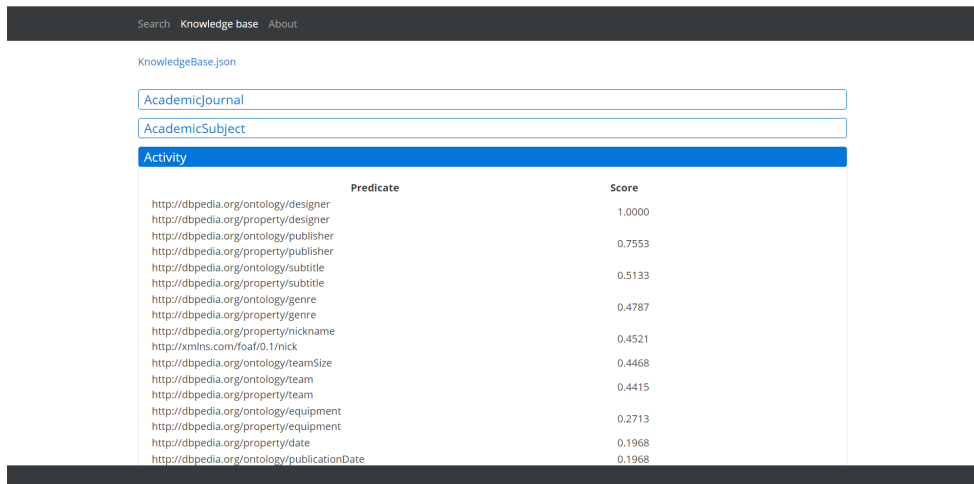


The screenshot shows the DBpedia profile for Bill Clinton. The page includes a search bar at the top with the URL [http://dbpedia.org/resource/Bill\\_Clinton](http://dbpedia.org/resource/Bill_Clinton). The profile is organized into several sections:

- Birth Place:** United States, Arkansas, Hope, Arkansas
- Order In Office:** 40th and 42nd Governor of Arkansas, 42nd President of the United States
- Birth Date:** 1946-08-19
- Active Years Start Date:** 1977-01-03, 1979-01-09, 1983-01-11, 1993-01-20

A portrait of Bill Clinton is displayed on the right side of the page. Below the portrait, the name "Bill Clinton" is followed by a short biographical text: "William Jefferson Clinton (born William Jefferson Blythe III; August 19, 1946), generally known as Bill Clinton, is an American politician who served as the 42nd President of the United States from 1993 to 2001. He previously served as Governor of..."

Obr. C.2: Ukážka zobrazenia zdroja triedy President



The screenshot shows the KnowledgeBase.json file for the class "Activity". The file lists various predicates and their corresponding scores:

Predicate	Score
<a href="http://dbpedia.org/ontology/designer">http://dbpedia.org/ontology/designer</a>	1.0000
<a href="http://dbpedia.org/property/designer">http://dbpedia.org/property/designer</a>	
<a href="http://dbpedia.org/ontology/publisher">http://dbpedia.org/ontology/publisher</a>	0.7553
<a href="http://dbpedia.org/property/publisher">http://dbpedia.org/property/publisher</a>	
<a href="http://dbpedia.org/ontology/subtitle">http://dbpedia.org/ontology/subtitle</a>	0.5133
<a href="http://dbpedia.org/property/subtitle">http://dbpedia.org/property/subtitle</a>	
<a href="http://dbpedia.org/ontology/genre">http://dbpedia.org/ontology/genre</a>	0.4787
<a href="http://dbpedia.org/property/genre">http://dbpedia.org/property/genre</a>	
<a href="http://dbpedia.org/property/mickname">http://dbpedia.org/property/mickname</a>	0.4521
<a href="http://xmlns.com/foaf/0.1/nick">http://xmlns.com/foaf/0.1/nick</a>	0.4468
<a href="http://dbpedia.org/ontology/teamSize">http://dbpedia.org/ontology/teamSize</a>	0.4415
<a href="http://dbpedia.org/ontology/team">http://dbpedia.org/ontology/team</a>	
<a href="http://dbpedia.org/property/team">http://dbpedia.org/property/team</a>	
<a href="http://dbpedia.org/ontology/equipment">http://dbpedia.org/ontology/equipment</a>	0.2713
<a href="http://dbpedia.org/property/equipment">http://dbpedia.org/property/equipment</a>	
<a href="http://dbpedia.org/property/date">http://dbpedia.org/property/date</a>	0.1968
<a href="http://dbpedia.org/ontology/publicationDate">http://dbpedia.org/ontology/publicationDate</a>	0.1968

Obr. C.3: Ukážka zobrazenia bázy znalosti

## **Ukážky bázy znalostí**

## D. UKÁŽKY BÁZY ZNALOSTÍ

---

Tabuľka D.1: Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Plant

Predikát	Skóre
<a href="http://dbpedia.org/property/r1Surface">http://dbpedia.org/property/r1Surface</a>	1.0000
<a href="http://dbpedia.org/property/regnum">http://dbpedia.org/property/regnum</a>	0.7658
<a href="http://dbpedia.org/ontology/synonym">http://dbpedia.org/ontology/synonym</a>	0.7611
<a href="http://dbpedia.org/property/synonyms">http://dbpedia.org/property/synonyms</a>	0.7606
<a href="http://dbpedia.org/property/r2Surface">http://dbpedia.org/property/r2Surface</a>	0.7148
<a href="http://dbpedia.org/ontology/genus">http://dbpedia.org/ontology/genus</a> <a href="http://dbpedia.org/property/genus">http://dbpedia.org/property/genus</a>	0.5003
<a href="http://dbpedia.org/ontology/binomialAuthority">http://dbpedia.org/ontology/binomialAuthority</a> <a href="http://dbpedia.org/property/binomialAuthority">http://dbpedia.org/property/binomialAuthority</a>	0.2135
<a href="http://dbpedia.org/property/statusSystem">http://dbpedia.org/property/statusSystem</a>	0.1673
<a href="http://dbpedia.org/ontology/conservationStatusSystem">http://dbpedia.org/ontology/conservationStatusSystem</a>	0.1673
<a href="http://dbpedia.org/ontology/conservationStatus">http://dbpedia.org/ontology/conservationStatus</a>	0.1661
<a href="http://dbpedia.org/ontology/status">http://dbpedia.org/ontology/status</a> <a href="http://dbpedia.org/property/status">http://dbpedia.org/property/status</a>	0.166
<a href="http://dbpedia.org/ontology/plant">http://dbpedia.org/ontology/plant</a>	0.1069
<a href="http://dbpedia.org/property/plant">http://dbpedia.org/property/plant</a>	0.0764
<a href="http://dbpedia.org/ontology/species">http://dbpedia.org/ontology/species</a> <a href="http://dbpedia.org/property/species">http://dbpedia.org/property/species</a>	0.0589
<a href="http://dbpedia.org/property/unrankedClassis">http://dbpedia.org/property/unrankedClassis</a>	0.0525
<a href="http://dbpedia.org/property/unrankedDivisio">http://dbpedia.org/property/unrankedDivisio</a>	0.0493
<a href="http://dbpedia.org/ontology/ingredient">http://dbpedia.org/ontology/ingredient</a>	0.0474
<a href="http://dbpedia.org/property/subgenus">http://dbpedia.org/property/subgenus</a>	0.0408
<a href="http://dbpedia.org/ontology/associatedBand">http://dbpedia.org/ontology/associatedBand</a> <a href="http://dbpedia.org/ontology/associatedMusicalArtist">http://dbpedia.org/ontology/associatedMusicalArtist</a>	0.0315
<a href="http://dbpedia.org/property/unrankedOrdo">http://dbpedia.org/property/unrankedOrdo</a>	0.0289

Tabuľka D.2: Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Animal

Predikát	Skóre
<a href="http://dbpedia.org/property/regnum">http://dbpedia.org/property/regnum</a>	1.0000
<a href="http://dbpedia.org/ontology/kingdom">http://dbpedia.org/ontology/kingdom</a>	0.8903
<a href="http://dbpedia.org/ontology/genus">http://dbpedia.org/ontology/genus</a> <a href="http://dbpedia.org/property/genus">http://dbpedia.org/property/genus</a>	0.8017
<a href="http://dbpedia.org/ontology/order">http://dbpedia.org/ontology/order</a>	0.6659
<a href="http://dbpedia.org/ontology/family">http://dbpedia.org/ontology/family</a> <a href="http://dbpedia.org/property/familia">http://dbpedia.org/property/familia</a>	0.5094
<a href="http://dbpedia.org/ontology/synonym">http://dbpedia.org/ontology/synonym</a>	0.4756
<a href="http://dbpedia.org/property/synonyms">http://dbpedia.org/property/synonyms</a>	0.4754
<a href="http://dbpedia.org/ontology/conservationStatus">http://dbpedia.org/ontology/conservationStatus</a>	0.4651
<a href="http://dbpedia.org/ontology/status">http://dbpedia.org/ontology/status</a> <a href="http://dbpedia.org/property/status">http://dbpedia.org/property/status</a>	0.4636
<a href="http://dbpedia.org/property/statusSystem">http://dbpedia.org/property/statusSystem</a>	0.4549
<a href="http://dbpedia.org/ontology/conservationStatusSystem">http://dbpedia.org/ontology/conservationStatusSystem</a>	0.4549
<a href="http://dbpedia.org/ontology/phylum">http://dbpedia.org/ontology/phylum</a> <a href="http://dbpedia.org/property/phylum">http://dbpedia.org/property/phylum</a>	0.3532
<a href="http://dbpedia.org/property/classis">http://dbpedia.org/property/classis</a>	0.3187
<a href="http://dbpedia.org/property/ordo">http://dbpedia.org/property/ordo</a>	0.214
<a href="http://dbpedia.org/ontology/binomialAuthority">http://dbpedia.org/ontology/binomialAuthority</a> <a href="http://dbpedia.org/property/binomialAuthority">http://dbpedia.org/property/binomialAuthority</a>	0.1694
<a href="http://dbpedia.org/ontology/class">http://dbpedia.org/ontology/class</a> <a href="http://dbpedia.org/property/class">http://dbpedia.org/property/class</a>	0.1100
<a href="http://dbpedia.org/property/subordo">http://dbpedia.org/property/subordo</a>	0.0771
<a href="http://dbpedia.org/property/superfamilia">http://dbpedia.org/property/superfamilia</a>	0.0671
<a href="http://dbpedia.org/ontology/species">http://dbpedia.org/ontology/species</a> <a href="http://dbpedia.org/property/species">http://dbpedia.org/property/species</a>	0.0192
<a href="http://dbpedia.org/property/filename">http://dbpedia.org/property/filename</a>	0.0143

Tabuľka D.3: Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Person

Predikát	Skóre
<a href="http://dbpedia.org/ontology/birthDate">http://dbpedia.org/ontology/birthDate</a> <a href="http://dbpedia.org/property/birthDate">http://dbpedia.org/property/birthDate</a>	0.9341
<a href="http://dbpedia.org/ontology/birthPlace">http://dbpedia.org/ontology/birthPlace</a> <a href="http://dbpedia.org/property/birthPlace">http://dbpedia.org/property/birthPlace</a>	0.7464
<a href="http://dbpedia.org/ontology/cinematography">http://dbpedia.org/ontology/cinematography</a> <a href="http://dbpedia.org/property/cinematography">http://dbpedia.org/property/cinematography</a>	0.7182
<a href="http://dbpedia.org/ontology/deathDate">http://dbpedia.org/ontology/deathDate</a> <a href="http://dbpedia.org/property/deathDate">http://dbpedia.org/property/deathDate</a>	0.5573
<a href="http://dbpedia.org/ontology/activeYearsStartYear">http://dbpedia.org/ontology/activeYearsStartYear</a>	0.4234
<a href="http://dbpedia.org/property/placeOfBirth">http://dbpedia.org/property/placeOfBirth</a>	0.3589
<a href="http://dbpedia.org/ontology/birthYear">http://dbpedia.org/ontology/birthYear</a>	0.3474
<a href="http://dbpedia.org/ontology/birthName">http://dbpedia.org/ontology/birthName</a> <a href="http://dbpedia.org/property/birthName">http://dbpedia.org/property/birthName</a>	0.3271
<a href="http://dbpedia.org/property/voices">http://dbpedia.org/property/voices</a>	0.3253
<a href="http://dbpedia.org/ontology/voice">http://dbpedia.org/ontology/voice</a> <a href="http://dbpedia.org/property/voice">http://dbpedia.org/property/voice</a>	0.3051
<a href="http://dbpedia.org/ontology/deathPlace">http://dbpedia.org/ontology/deathPlace</a> <a href="http://dbpedia.org/property/deathPlace">http://dbpedia.org/property/deathPlace</a>	0.2765
<a href="http://dbpedia.org/ontology/starring">http://dbpedia.org/ontology/starring</a> <a href="http://dbpedia.org/property/starring">http://dbpedia.org/property/starring</a>	0.2518
<a href="http://dbpedia.org/property/yearsActive">http://dbpedia.org/property/yearsActive</a>	0.2384
<a href="http://dbpedia.org/ontology/activeYearsEndYear">http://dbpedia.org/ontology/activeYearsEndYear</a>	0.2236
<a href="http://dbpedia.org/ontology/deathYear">http://dbpedia.org/ontology/deathYear</a>	0.1984
<a href="http://dbpedia.org/property/termStart">http://dbpedia.org/property/termStart</a>	0.1759
<a href="http://dbpedia.org/property/termEnd">http://dbpedia.org/property/termEnd</a>	0.1711
<a href="http://xmlns.com/foaf/0.1/givenName">http://xmlns.com/foaf/0.1/givenName</a>	0.159
<a href="http://dbpedia.org/ontology/occupation">http://dbpedia.org/ontology/occupation</a> <a href="http://dbpedia.org/property/occupation">http://dbpedia.org/property/occupation</a>	0.1454
<a href="http://dbpedia.org/property/reign">http://dbpedia.org/property/reign</a>	0.1321



Tabuľka D.4: Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:President

Predikát	Skóre
<a href="http://dbpedia.org/ontology/orderInOffice">http://dbpedia.org/ontology/orderInOffice</a>	1.0000
<a href="http://dbpedia.org/ontology/birthDate">http://dbpedia.org/ontology/birthDate</a> <a href="http://dbpedia.org/property/birthDate">http://dbpedia.org/property/birthDate</a>	0.8971
<a href="http://dbpedia.org/ontology/activeYearsStartDate">http://dbpedia.org/ontology/activeYearsStartDate</a> <a href="http://dbpedia.org/property/termStart">http://dbpedia.org/property/termStart</a>	0.8636 0.8627
<a href="http://dbpedia.org/ontology/activeYearsEndDate">http://dbpedia.org/ontology/activeYearsEndDate</a> <a href="http://dbpedia.org/property/termEnd">http://dbpedia.org/property/termEnd</a>	0.8284 0.8275
<a href="http://dbpedia.org/property/birthname">http://dbpedia.org/property/birthname</a>	0.5152
<a href="http://dbpedia.org/ontology/deathDate">http://dbpedia.org/ontology/deathDate</a> <a href="http://dbpedia.org/property/deathDate">http://dbpedia.org/property/deathDate</a>	0.4706
<a href="http://dbpedia.org/property/shortDescription">http://dbpedia.org/property/shortDescription</a> <a href="http://purl.org/dc/elements/1.1/description">http://purl.org/dc/elements/1.1/description</a>	0.3817
<a href="http://xmlns.com/foaf/0.1/surname">http://xmlns.com/foaf/0.1/surname</a>	0.3568
<a href="http://xmlns.com/foaf/0.1/givenName">http://xmlns.com/foaf/0.1/givenName</a>	0.3568
<a href="http://dbpedia.org/ontology/birthPlace">http://dbpedia.org/ontology/birthPlace</a> <a href="http://dbpedia.org/property/birthPlace">http://dbpedia.org/property/birthPlace</a>	0.349
<a href="http://dbpedia.org/ontology/birthYear">http://dbpedia.org/ontology/birthYear</a> <a href="http://dbpedia.org/property/dateOfBirth">http://dbpedia.org/property/dateOfBirth</a>	0.3309 0.3309
<a href="http://dbpedia.org/property/order">http://dbpedia.org/property/order</a>	0.3073
<a href="http://dbpedia.org/ontology/office">http://dbpedia.org/ontology/office</a> <a href="http://dbpedia.org/property/office">http://dbpedia.org/property/office</a>	0.2377
<a href="http://dbpedia.org/ontology/deathPlace">http://dbpedia.org/ontology/deathPlace</a> <a href="http://dbpedia.org/property/deathPlace">http://dbpedia.org/property/deathPlace</a>	0.2262
<a href="http://dbpedia.org/ontology/successor">http://dbpedia.org/ontology/successor</a> <a href="http://dbpedia.org/property/successor">http://dbpedia.org/property/successor</a>	0.2109
<a href="http://dbpedia.org/ontology/president">http://dbpedia.org/ontology/president</a> <a href="http://dbpedia.org/property/president">http://dbpedia.org/property/president</a>	0.2015
<a href="http://dbpedia.org/ontology/party">http://dbpedia.org/ontology/party</a> <a href="http://dbpedia.org/property/party">http://dbpedia.org/property/party</a>	0.1965

## D. UKÁŽKY BÁZY ZNALOSTÍ

---

Tabuľka D.5: Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Monarch

Predikát	Skóre
<a href="http://dbpedia.org/ontology/title">http://dbpedia.org/ontology/title</a>	1.0000
<a href="http://dbpedia.org/ontology/activeYearsStartYear">http://dbpedia.org/ontology/activeYearsStartYear</a>	0.939
<a href="http://dbpedia.org/property/reign">http://dbpedia.org/property/reign</a>	0.9317
<a href="http://dbpedia.org/ontology/activeYearsEndYear">http://dbpedia.org/ontology/activeYearsEndYear</a>	0.8255
<a href="http://dbpedia.org/ontology/successor">http://dbpedia.org/ontology/successor</a> <a href="http://dbpedia.org/property/successor">http://dbpedia.org/property/successor</a>	0.5307
<a href="http://dbpedia.org/ontology/predecessor">http://dbpedia.org/ontology/predecessor</a> <a href="http://dbpedia.org/property/predecessor">http://dbpedia.org/property/predecessor</a>	0.4649
<a href="http://dbpedia.org/property/shortDescription">http://dbpedia.org/property/shortDescription</a> <a href="http://purl.org/dc/elements/1.1/description">http://purl.org/dc/elements/1.1/description</a>	0.2365
<a href="http://dbpedia.org/property/fullName">http://dbpedia.org/property/fullName</a>	0.2346
<a href="http://dbpedia.org/ontology/deathYear">http://dbpedia.org/ontology/deathYear</a>	0.1664
<a href="http://dbpedia.org/property/dateOfDeath">http://dbpedia.org/property/dateOfDeath</a>	0.0974
<a href="http://dbpedia.org/ontology/birthYear">http://dbpedia.org/ontology/birthYear</a>	0.0945
<a href="http://dbpedia.org/property/dateOfBirth">http://dbpedia.org/property/dateOfBirth</a>	0.094
<a href="http://xmlns.com/foaf/0.1/surname">http://xmlns.com/foaf/0.1/surname</a>	0.0909
<a href="http://xmlns.com/foaf/0.1/givenName">http://xmlns.com/foaf/0.1/givenName</a>	0.0909
<a href="http://dbpedia.org/ontology/deathDate">http://dbpedia.org/ontology/deathDate</a> <a href="http://dbpedia.org/property/deathDate">http://dbpedia.org/property/deathDate</a>	0.0821
<a href="http://dbpedia.org/ontology/birthDate">http://dbpedia.org/ontology/birthDate</a> <a href="http://dbpedia.org/property/birthDate">http://dbpedia.org/property/birthDate</a>	0.0472
<a href="http://dbpedia.org/ontology/alias">http://dbpedia.org/ontology/alias</a> <a href="http://dbpedia.org/property/alias">http://dbpedia.org/property/alias</a>	0.0339
<a href="http://dbpedia.org/property/father">http://dbpedia.org/property/father</a>	0.0326
<a href="http://dbpedia.org/ontology/deathPlace">http://dbpedia.org/ontology/deathPlace</a> <a href="http://dbpedia.org/property/deathPlace">http://dbpedia.org/property/deathPlace</a>	0.0315
<a href="http://dbpedia.org/ontology/placeOfBurial">http://dbpedia.org/ontology/placeOfBurial</a> <a href="http://dbpedia.org/property/placeOfBurial">http://dbpedia.org/property/placeOfBurial</a>	0.0272

Tabuľka D.6: Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Artist

Predikát	Skóre
<a href="http://dbpedia.org/ontology/birthDate">http://dbpedia.org/ontology/birthDate</a> <a href="http://dbpedia.org/property/birthDate">http://dbpedia.org/property/birthDate</a>	0.9832
<a href="http://dbpedia.org/ontology/background">http://dbpedia.org/ontology/background</a> <a href="http://dbpedia.org/property/background">http://dbpedia.org/property/background</a>	0.8549
<a href="http://dbpedia.org/ontology/activeYearsStartYear">http://dbpedia.org/ontology/activeYearsStartYear</a> <a href="http://dbpedia.org/property/yearsActive">http://dbpedia.org/property/yearsActive</a>	0.7691 0.732
<a href="http://xmlns.com/foaf/0.1/givenName">http://xmlns.com/foaf/0.1/givenName</a>	0.7132
<a href="http://dbpedia.org/ontology/birthPlace">http://dbpedia.org/ontology/birthPlace</a> <a href="http://dbpedia.org/property/birthPlace">http://dbpedia.org/property/birthPlace</a>	0.6269
<a href="http://dbpedia.org/ontology/birthName">http://dbpedia.org/ontology/birthName</a> <a href="http://dbpedia.org/property/birthName">http://dbpedia.org/property/birthName</a>	0.4192
<a href="http://dbpedia.org/ontology/associatedBand">http://dbpedia.org/ontology/associatedBand</a> <a href="http://dbpedia.org/ontology/associatedMusicalArtist">http://dbpedia.org/ontology/associatedMusicalArtist</a>	0.3485
<a href="http://dbpedia.org/ontology/hometown">http://dbpedia.org/ontology/hometown</a> <a href="http://dbpedia.org/property/hometown">http://dbpedia.org/property/hometown</a>	0.3214
<a href="http://dbpedia.org/ontology/deathDate">http://dbpedia.org/ontology/deathDate</a> <a href="http://dbpedia.org/property/deathDate">http://dbpedia.org/property/deathDate</a>	0.3028
<a href="http://dbpedia.org/ontology/alias">http://dbpedia.org/ontology/alias</a> <a href="http://dbpedia.org/property/alias">http://dbpedia.org/property/alias</a>	0.301
<a href="http://dbpedia.org/property/placeOfBirth">http://dbpedia.org/property/placeOfBirth</a>	0.2822
<a href="http://dbpedia.org/ontology/origin">http://dbpedia.org/ontology/origin</a> <a href="http://dbpedia.org/property/origin">http://dbpedia.org/property/origin</a>	0.2564
<a href="http://dbpedia.org/property/associatedActs">http://dbpedia.org/property/associatedActs</a>	0.2353
<a href="http://dbpedia.org/ontology/genre">http://dbpedia.org/ontology/genre</a> <a href="http://dbpedia.org/property/genre">http://dbpedia.org/property/genre</a>	0.2327
<a href="http://dbpedia.org/ontology/activeYearsEndYear">http://dbpedia.org/ontology/activeYearsEndYear</a>	0.2248
<a href="http://dbpedia.org/ontology/recordLabel">http://dbpedia.org/ontology/recordLabel</a> <a href="http://dbpedia.org/property/label">http://dbpedia.org/property/label</a>	0.2206
<a href="http://dbpedia.org/property/shortDescription">http://dbpedia.org/property/shortDescription</a> <a href="http://purl.org/dc/elements/1.1/description">http://purl.org/dc/elements/1.1/description</a>	0.1835
<a href="http://dbpedia.org/ontology/deathPlace">http://dbpedia.org/ontology/deathPlace</a> <a href="http://dbpedia.org/property/deathPlace">http://dbpedia.org/property/deathPlace</a>	0.1819
<a href="http://dbpedia.org/ontology/birthYear">http://dbpedia.org/ontology/birthYear</a>	0.1734

## D. UKÁŽKY BÁZY ZNALOSTÍ

Tabuľka D.7: Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Film

Predikát	Skóre
<a href="http://dbpedia.org/ontology/runtime">http://dbpedia.org/ontology/runtime</a>	1.0000
<a href="http://dbpedia.org/ontology/Work/runtime">http://dbpedia.org/ontology/Work/runtime</a>	1.0000
<a href="http://dbpedia.org/property/runtime">http://dbpedia.org/property/runtime</a>	0.9997
<a href="http://dbpedia.org/ontology/starring">http://dbpedia.org/ontology/starring</a>	0.6816
<a href="http://dbpedia.org/property/starring">http://dbpedia.org/property/starring</a>	
<a href="http://dbpedia.org/ontology/gross">http://dbpedia.org/ontology/gross</a>	0.6582
<a href="http://dbpedia.org/property/gross">http://dbpedia.org/property/gross</a>	0.6582
<a href="http://dbpedia.org/ontology/budget">http://dbpedia.org/ontology/budget</a>	0.5786
<a href="http://dbpedia.org/property/budget">http://dbpedia.org/property/budget</a>	
<a href="http://dbpedia.org/ontology/writer">http://dbpedia.org/ontology/writer</a>	0.3796
<a href="http://dbpedia.org/property/writer">http://dbpedia.org/property/writer</a>	
<a href="http://dbpedia.org/ontology/producer">http://dbpedia.org/ontology/producer</a>	0.3009
<a href="http://dbpedia.org/property/producer">http://dbpedia.org/property/producer</a>	
<a href="http://dbpedia.org/ontology/imdbId">http://dbpedia.org/ontology/imdbId</a>	0.279
<a href="http://dbpedia.org/property/released">http://dbpedia.org/property/released</a>	0.2539
<a href="http://dbpedia.org/property/id">http://dbpedia.org/property/id</a>	0.2452
<a href="http://dbpedia.org/ontology/distributor">http://dbpedia.org/ontology/distributor</a>	0.194
<a href="http://dbpedia.org/property/distributor">http://dbpedia.org/property/distributor</a>	
<a href="http://dbpedia.org/ontology/releaseDate">http://dbpedia.org/ontology/releaseDate</a>	0.158
<a href="http://dbpedia.org/ontology/director">http://dbpedia.org/ontology/director</a>	0.1553
<a href="http://dbpedia.org/property/director">http://dbpedia.org/property/director</a>	
<a href="http://dbpedia.org/property/screenplay">http://dbpedia.org/property/screenplay</a>	0.1236
<a href="http://dbpedia.org/ontology/musicComposer">http://dbpedia.org/ontology/musicComposer</a>	0.1033
<a href="http://dbpedia.org/property/music">http://dbpedia.org/property/music</a>	
<a href="http://dbpedia.org/ontology/Engine/length">http://dbpedia.org/ontology/Engine/length</a>	0.0861
<a href="http://dbpedia.org/ontology/Infrastructure/length">http://dbpedia.org/ontology/Infrastructure/length</a>	
<a href="http://dbpedia.org/ontology/MeanOfTransportation/length">http://dbpedia.org/ontology/MeanOfTransportation/length</a>	
<a href="http://dbpedia.org/ontology/Weapon/length">http://dbpedia.org/ontology/Weapon/length</a>	
<a href="http://dbpedia.org/ontology/length">http://dbpedia.org/ontology/length</a>	
<a href="http://dbpedia.org/property/length">http://dbpedia.org/property/length</a>	
<a href="http://dbpedia.org/ontology/language">http://dbpedia.org/ontology/language</a>	0.0856
<a href="http://dbpedia.org/property/language">http://dbpedia.org/property/language</a>	
<a href="http://dbpedia.org/ontology/editing">http://dbpedia.org/ontology/editing</a>	0.0707
<a href="http://dbpedia.org/property/editing">http://dbpedia.org/property/editing</a>	

Tabuľka D.8: Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Country

Predikát	Skóre
<a href="http://dbpedia.org/property/conventionalLongName">http://dbpedia.org/property/conventionalLongName</a>	1.0000
<a href="http://dbpedia.org/ontology/country">http://dbpedia.org/ontology/country</a> <a href="http://dbpedia.org/property/country">http://dbpedia.org/property/country</a>	1.0000
<a href="http://dbpedia.org/ontology/dissolutionYear">http://dbpedia.org/ontology/dissolutionYear</a>	0.8837
<a href="http://dbpedia.org/property/yearEnd">http://dbpedia.org/property/yearEnd</a>	0.8811
<a href="http://dbpedia.org/ontology/foundingYear">http://dbpedia.org/ontology/foundingYear</a>	0.8709
<a href="http://dbpedia.org/property/yearStart">http://dbpedia.org/property/yearStart</a>	0.868
<a href="http://dbpedia.org/property/nativeName">http://dbpedia.org/property/nativeName</a>	0.6563
<a href="http://dbpedia.org/ontology/birthPlace">http://dbpedia.org/ontology/birthPlace</a> <a href="http://dbpedia.org/property/birthPlace">http://dbpedia.org/property/birthPlace</a>	0.5696
<a href="http://dbpedia.org/ontology/foundingDate">http://dbpedia.org/ontology/foundingDate</a>	0.3981
<a href="http://www.georss.org/georss/point">http://www.georss.org/georss/point</a>	0.3722
<a href="http://www.w3.org/2003/01/geo/wgs84_pos#geometry">http://www.w3.org/2003/01/geo/wgs84_pos#geometry</a>	0.3722
<a href="http://www.w3.org/2003/01/geo/wgs84_pos#long">http://www.w3.org/2003/01/geo/wgs84_pos#long</a>	0.3722
<a href="http://www.w3.org/2003/01/geo/wgs84_pos#lat">http://www.w3.org/2003/01/geo/wgs84_pos#lat</a>	0.3722
<a href="http://dbpedia.org/property/placeOfBirth">http://dbpedia.org/property/placeOfBirth</a>	0.3348
<a href="http://dbpedia.org/ontology/dissolutionDate">http://dbpedia.org/ontology/dissolutionDate</a>	0.3273
<a href="http://dbpedia.org/property/nationalAnthem">http://dbpedia.org/property/nationalAnthem</a>	0.1638
<a href="http://dbpedia.org/ontology/capital">http://dbpedia.org/ontology/capital</a> <a href="http://dbpedia.org/property/capital">http://dbpedia.org/property/capital</a>	0.152
<a href="http://dbpedia.org/ontology/leaderTitle">http://dbpedia.org/ontology/leaderTitle</a>	0.1468
<a href="http://dbpedia.org/ontology/location">http://dbpedia.org/ontology/location</a> <a href="http://dbpedia.org/property/location">http://dbpedia.org/property/location</a>	0.1455
<a href="http://dbpedia.org/ontology/nationality">http://dbpedia.org/ontology/nationality</a> <a href="http://dbpedia.org/property/nationality">http://dbpedia.org/property/nationality</a>	0.1449

## D. UKÁŽKY BÁZY ZNALOSTÍ

---

Tabuľka D.9: Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Company

Predikát	Skóre
<a href="http://dbpedia.org/ontology/foundingYear">http://dbpedia.org/ontology/foundingYear</a>	1.0000
<a href="http://dbpedia.org/ontology/recordLabel">http://dbpedia.org/ontology/recordLabel</a> <a href="http://dbpedia.org/property/label">http://dbpedia.org/property/label</a>	0.8116
<a href="http://dbpedia.org/ontology/numberOfEmployees">http://dbpedia.org/ontology/numberOfEmployees</a> <a href="http://dbpedia.org/property/numEmployees">http://dbpedia.org/property/numEmployees</a>	0.5023
<a href="http://dbpedia.org/property/website">http://dbpedia.org/property/website</a>	0.5011
<a href="http://dbpedia.org/property/foundation">http://dbpedia.org/property/foundation</a>	0.357
<a href="http://dbpedia.org/ontology/revenue">http://dbpedia.org/ontology/revenue</a> <a href="http://dbpedia.org/property/revenue">http://dbpedia.org/property/revenue</a>	0.3422
<a href="http://dbpedia.org/ontology/industry">http://dbpedia.org/ontology/industry</a> <a href="http://dbpedia.org/property/industry">http://dbpedia.org/property/industry</a>	0.2239
<a href="http://dbpedia.org/ontology/netIncome">http://dbpedia.org/ontology/netIncome</a> <a href="http://dbpedia.org/property/netIncome">http://dbpedia.org/property/netIncome</a>	0.2226
<a href="http://dbpedia.org/property/products">http://dbpedia.org/property/products</a>	0.217
<a href="http://dbpedia.org/ontology/product">http://dbpedia.org/ontology/product</a> <a href="http://dbpedia.org/property/product">http://dbpedia.org/property/product</a>	0.19
<a href="http://dbpedia.org/property/keyPeople">http://dbpedia.org/property/keyPeople</a>	0.1873
<a href="http://dbpedia.org/ontology/distributor">http://dbpedia.org/ontology/distributor</a> <a href="http://dbpedia.org/property/distributor">http://dbpedia.org/property/distributor</a>	0.1688
<a href="http://dbpedia.org/ontology/location">http://dbpedia.org/ontology/location</a> <a href="http://dbpedia.org/property/location">http://dbpedia.org/property/location</a>	0.1655
<a href="http://dbpedia.org/ontology/operatingIncome">http://dbpedia.org/ontology/operatingIncome</a> <a href="http://dbpedia.org/property/operatingIncome">http://dbpedia.org/property/operatingIncome</a>	0.1639
<a href="http://dbpedia.org/ontology/type">http://dbpedia.org/ontology/type</a> <a href="http://dbpedia.org/property/type">http://dbpedia.org/property/type</a>	0.1628
<a href="http://dbpedia.org/ontology/assets">http://dbpedia.org/ontology/assets</a> <a href="http://dbpedia.org/property/assets">http://dbpedia.org/property/assets</a>	0.1467
<a href="http://dbpedia.org/property/founded">http://dbpedia.org/property/founded</a>	0.1441
<a href="http://dbpedia.org/ontology/extinctionYear">http://dbpedia.org/ontology/extinctionYear</a>	0.1394
<a href="http://dbpedia.org/property/defunct">http://dbpedia.org/property/defunct</a>	0.1262
<a href="http://dbpedia.org/ontology/fate">http://dbpedia.org/ontology/fate</a> <a href="http://dbpedia.org/property/fate">http://dbpedia.org/property/fate</a>	0.125

Tabuľka D.10: Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Sport

Predikát	Skóre
<a href="http://dbpedia.org/ontology/sport">http://dbpedia.org/ontology/sport</a> <a href="http://dbpedia.org/property/sport">http://dbpedia.org/property/sport</a>	1.0000
<a href="http://dbpedia.org/property/nickname">http://dbpedia.org/property/nickname</a> <a href="http://xmlns.com/foaf/0.1/nick">http://xmlns.com/foaf/0.1/nick</a>	0.9884
<a href="http://dbpedia.org/ontology/teamSize">http://dbpedia.org/ontology/teamSize</a> <a href="http://dbpedia.org/property/sports">http://dbpedia.org/property/sports</a>	0.9767 0.5461
<a href="http://dbpedia.org/ontology/team">http://dbpedia.org/ontology/team</a> <a href="http://dbpedia.org/property/team">http://dbpedia.org/property/team</a>	0.4826
<a href="http://dbpedia.org/ontology/equipment">http://dbpedia.org/ontology/equipment</a> <a href="http://dbpedia.org/property/equipment">http://dbpedia.org/property/equipment</a>	0.3772
<a href="http://dbpedia.org/ontology/sportGoverningBody">http://dbpedia.org/ontology/sportGoverningBody</a> <a href="http://dbpedia.org/property/union">http://dbpedia.org/property/union</a>	0.2093
<a href="http://dbpedia.org/ontology/category">http://dbpedia.org/ontology/category</a> <a href="http://dbpedia.org/property/category">http://dbpedia.org/property/category</a>	0.1835
<a href="http://dbpedia.org/property/ball">http://dbpedia.org/property/ball</a>	0.1512
<a href="http://dbpedia.org/ontology/occupation">http://dbpedia.org/ontology/occupation</a> <a href="http://dbpedia.org/property/occupation">http://dbpedia.org/property/occupation</a>	0.0378
<a href="http://dbpedia.org/ontology/event">http://dbpedia.org/ontology/event</a> <a href="http://dbpedia.org/property/event">http://dbpedia.org/property/event</a>	0.0154
<a href="http://dbpedia.org/ontology/birthDate">http://dbpedia.org/ontology/birthDate</a> <a href="http://dbpedia.org/property/birthDate">http://dbpedia.org/property/birthDate</a>	0.0116
<a href="http://xmlns.com/foaf/0.1/surname">http://xmlns.com/foaf/0.1/surname</a>	0.0116
<a href="http://xmlns.com/foaf/0.1/givenName">http://xmlns.com/foaf/0.1/givenName</a>	0.0116
<a href="http://dbpedia.org/property/dateOfBirth">http://dbpedia.org/property/dateOfBirth</a>	0.0116
<a href="http://dbpedia.org/property/fileName">http://dbpedia.org/property/fileName</a>	0.0116
<a href="http://dbpedia.org/property/description">http://dbpedia.org/property/description</a>	0.0116
<a href="http://dbpedia.org/property/mostChamps">http://dbpedia.org/property/mostChamps</a>	0.0116
<a href="http://dbpedia.org/property/founded">http://dbpedia.org/property/founded</a>	0.0116
<a href="http://dbpedia.org/property/teams">http://dbpedia.org/property/teams</a>	0.0116

Tabuľka D.11: Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Work

Predikát	Skóre
<a href="http://dbpedia.org/ontology/runtime">http://dbpedia.org/ontology/runtime</a>	1.0000
<a href="http://dbpedia.org/ontology/Work/runtime">http://dbpedia.org/ontology/Work/runtime</a>	1.0000
<a href="http://dbpedia.org/property/runtime">http://dbpedia.org/property/runtime</a>	0.805
<a href="http://dbpedia.org/ontology/releaseDate">http://dbpedia.org/ontology/releaseDate</a>	0.8039
<a href="http://dbpedia.org/ontology/genre">http://dbpedia.org/ontology/genre</a> <a href="http://dbpedia.org/property/genre">http://dbpedia.org/property/genre</a>	0.569
<a href="http://dbpedia.org/property/firstAired">http://dbpedia.org/property/firstAired</a>	0.4665
<a href="http://dbpedia.org/property/showName">http://dbpedia.org/property/showName</a>	0.4356
<a href="http://dbpedia.org/property/numEpisodes">http://dbpedia.org/property/numEpisodes</a>	0.433
<a href="http://dbpedia.org/ontology/numberOfEpisodes">http://dbpedia.org/ontology/numberOfEpisodes</a>	0.4318
<a href="http://dbpedia.org/ontology/producer">http://dbpedia.org/ontology/producer</a> <a href="http://dbpedia.org/property/producer">http://dbpedia.org/property/producer</a>	0.3835
<a href="http://dbpedia.org/property/gross">http://dbpedia.org/property/gross</a>	0.3519
<a href="http://dbpedia.org/ontology/gross">http://dbpedia.org/ontology/gross</a>	0.3515
<a href="http://dbpedia.org/ontology/numberOfSeasons">http://dbpedia.org/ontology/numberOfSeasons</a>	0.3508
<a href="http://dbpedia.org/ontology/budget">http://dbpedia.org/ontology/budget</a> <a href="http://dbpedia.org/property/budget">http://dbpedia.org/property/budget</a>	0.3307
<a href="http://dbpedia.org/property/lastAired">http://dbpedia.org/property/lastAired</a>	0.3033
<a href="http://dbpedia.org/ontology/completionDate">http://dbpedia.org/ontology/completionDate</a>	0.3024
<a href="http://dbpedia.org/property/numSeasons">http://dbpedia.org/property/numSeasons</a>	0.3022
<a href="http://dbpedia.org/ontology/starring">http://dbpedia.org/ontology/starring</a> <a href="http://dbpedia.org/property/starring">http://dbpedia.org/property/starring</a>	0.2959
<a href="http://dbpedia.org/ontology/recordLabel">http://dbpedia.org/ontology/recordLabel</a> <a href="http://dbpedia.org/property/label">http://dbpedia.org/property/label</a>	0.2612
<a href="http://dbpedia.org/ontology/writer">http://dbpedia.org/ontology/writer</a> <a href="http://dbpedia.org/property/writer">http://dbpedia.org/property/writer</a>	0.231



Tabuľka D.12: Ukážka prvých 20 záznamov bázy znalosti pre triedu dbo:Eukaryote

Predikát	Skóre
<a href="http://dbpedia.org/property/regnum">http://dbpedia.org/property/regnum</a>	1.0000
<a href="http://dbpedia.org/ontology/kingdom">http://dbpedia.org/ontology/kingdom</a>	0.7810
<a href="http://dbpedia.org/ontology/genus">http://dbpedia.org/ontology/genus</a>	0.7000
<a href="http://dbpedia.org/property/genus">http://dbpedia.org/property/genus</a>	
<a href="http://dbpedia.org/ontology/synonym">http://dbpedia.org/ontology/synonym</a>	0.547
<a href="http://dbpedia.org/property/synonyms">http://dbpedia.org/property/synonyms</a>	0.5463
<a href="http://dbpedia.org/property/classis">http://dbpedia.org/property/classis</a>	0.5177
<a href="http://dbpedia.org/ontology/conservationStatus">http://dbpedia.org/ontology/conservationStatus</a>	0.3544
<a href="http://dbpedia.org/ontology/status">http://dbpedia.org/ontology/status</a>	0.3529
<a href="http://dbpedia.org/property/status">http://dbpedia.org/property/status</a>	
<a href="http://dbpedia.org/property/statusSystem">http://dbpedia.org/property/statusSystem</a>	0.3452
<a href="http://dbpedia.org/ontology/conservationStatusSystem">http://dbpedia.org/ontology/conservationStatusSystem</a>	0.3452
<a href="http://dbpedia.org/ontology/class">http://dbpedia.org/ontology/class</a>	0.3168
<a href="http://dbpedia.org/property/class">http://dbpedia.org/property/class</a>	
<a href="http://dbpedia.org/ontology/division">http://dbpedia.org/ontology/division</a>	0.2591
<a href="http://dbpedia.org/ontology/binomialAuthority">http://dbpedia.org/ontology/binomialAuthority</a>	0.2365
<a href="http://dbpedia.org/property/binomialAuthority">http://dbpedia.org/property/binomialAuthority</a>	
<a href="http://dbpedia.org/ontology/phylum">http://dbpedia.org/ontology/phylum</a>	0.1605
<a href="http://dbpedia.org/property/phylum">http://dbpedia.org/property/phylum</a>	
<a href="http://dbpedia.org/ontology/order">http://dbpedia.org/ontology/order</a>	0.1131
<a href="http://dbpedia.org/property/ordo">http://dbpedia.org/property/ordo</a>	0.0801
<a href="http://dbpedia.org/ontology/species">http://dbpedia.org/ontology/species</a>	0.0525
<a href="http://dbpedia.org/property/species">http://dbpedia.org/property/species</a>	
<a href="http://dbpedia.org/property/divisio">http://dbpedia.org/property/divisio</a>	0.0296
<a href="http://dbpedia.org/property/subgenus">http://dbpedia.org/property/subgenus</a>	0.0221
<a href="http://dbpedia.org/property/unrankedDivisio">http://dbpedia.org/property/unrankedDivisio</a>	0.0195