

Slovenská technická univerzita v Bratislave
Fakulta informatiky a informačných technológií
Ilkovičova 2, 842 16, Bratislava 4

Deep Search

Dokumentácia k tímovému projektu

(Inžinierske dielo)

Vedúci tímu: Ing. Nadežda Andrejčíkova, PhD.

Členovia tímu: Bc. Peter Berta, Bc. Matej Adamov, Bc. Michal Krempaský, Bc.
Oliver Macko, Bc. Bronislava Pečíková

Akademický rok : 2017/2018

Obsah

1. Úvod	3
2. Globálne ciele projektu	3
3. Celkový pohľad na systém	4
3.1. Modul analýzy vstupov	5
3.1.1. Výber databázy pre ukladanie neštruktúrovaných dát	7
3.1.2. Výber databázy pre ukladanie štruktúrovaných dát	7
3.2. Modul predspracovania textu	8
3.2.1. Požiadavky	8
3.2.2. Analýza prístupov cez jednotlivé služby	8
3.2.3. Návrh	12
3.2.4. Implementácia	12
3.3. Modul správa používateľov	15
3.4. Modul používateľské rozhranie	16
3.4.1. Prihlásenie	16
3.4.2. Registrácia	17
3.4.3. Hlavná obrazovka	17
3.4.4. Vyhľadávanie	18
3.4.5. Vyhľadávanie vzťahov	18
3.4.6. Spracovanie textu	19

1. Úvod

Informačná explózia so sebou prináša aj viacero problémov. Napriek tomu, že v dnešnej dobe si nemožno sťažovať na nedostatok informácií, máme často problém nájsť to, čo práve potrebujeme. Väčšina dokumentov je totiž v neštruktúrovanej podobe a získať z nich informácie typu: kto v danom období pôsobil v určitom regióne je prakticky nemožné. Fulltextové vyhľadávanie má jeden vážny nedostatok a to, že nezohľadňuje sémantiku daných kľúčových slov. V našom projekte sa zameriavame na spracovanie prirodzeného jazyka a stanovili sme si pomerne ambiciózne cieľ, ktorým je extrakcia štruktúrovaných dát z neštruktúrovaného textu so zachytením ich významu. Pričom zameriavať sa budeme najmä na životopisy a iné dokumenty, z ktorých budeme môcť extrahovať informácie o tom kde a kedy dané osoby študovali, prípadne pôsobili. Pokúsime sa tiež z textov získať vzťahy typu: kolega alebo spolužiak. Našou úlohou bude teda rozpoznávať a pokiaľ to bude možné aj jednoznačne identifikovať, entity typu osoba, korporácia, geografická jednotka, dátácia a zároveň identifikovať udalosť štúdium, prípadne pôsobenie a v rámci nich vzťahy medzi týmito entitami. Cieľom je tieto údaje uložiť v štruktúrovanej podobe tak, aby bolo možné v nich vyhľadávať a získať informáciu o tom kto a kedy v danom mieste študoval, s kým sa mohol poznať a vytvárať tak aj virtuálne komunity napr. pre určité zameranie.

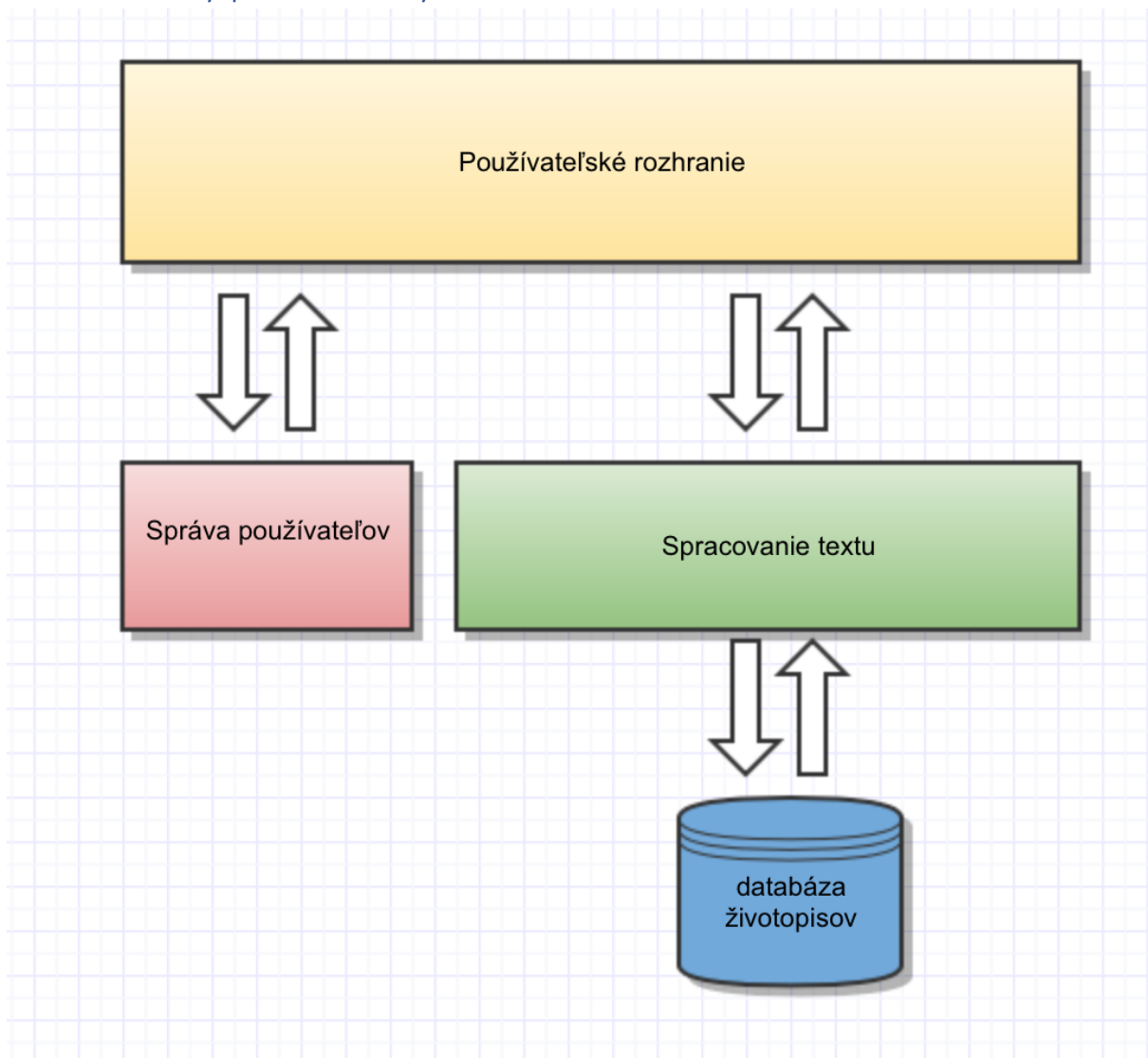
2. Globálne ciele projektu

Cieľom tohto semestra je vytvorenie databázy na základe neštruktúrovaných životopisov. Je potrebné identifikovať nasledovné polia:

- Narodenie – dátum a miesto
- Úmrtie – dátum a miesto
- Štúdium – dátum, miesto, profesori, spolužiaci

V ideálnom prípade by bolo vhodné taktiež vizualizovať súvislosti medzi osobnosťami pomocou grafu. Jednotlivé vrcholy budú predstavovať osobnosti a hrany budú predstavovať súvislosti medzi nimi.

3. Celkový pohľad na systém

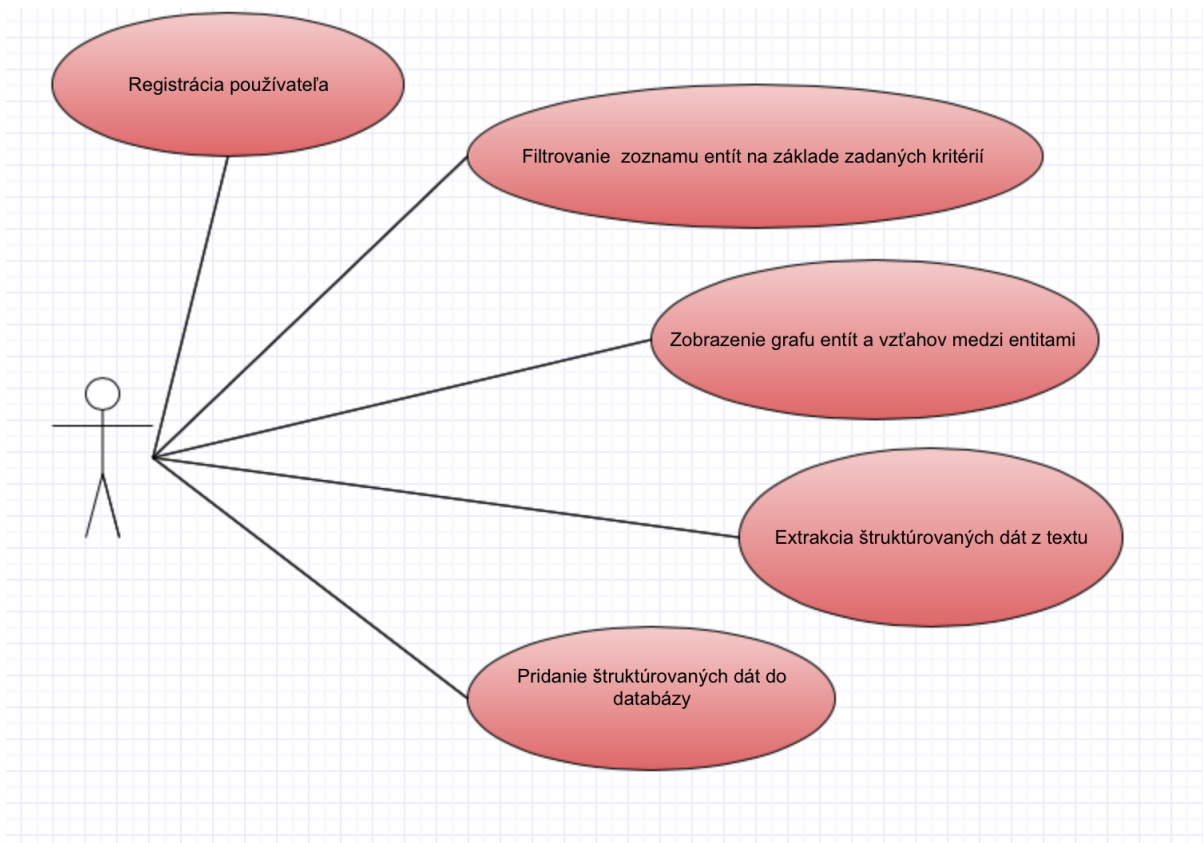


Obr. 1: Moduly systému

Systém pozostáva z troch modulov:

- používateľské rozhranie
- správa používateľov
- spracovanie textu

Používateľské rozhranie bude implementované prostredníctvom webovej aplikácie a bude ponúkať funkcionality ako sú vyhľadávanie v štruktúrovaných dátach, zobrazenie grafu vzťahov, pridanie životopisu... Správa používateľov bude riešiť prihlásenie, odhlásenie a registráciu používateľa. Modul spracovanie textu extrahuje zo životopisov údaje o štúdiu a zamestnaní a ukladá ich v štruktúrovanej podobe.



Obr. 2 diagram prípadov použitia

3.1. Modul analýzy vstupov

V rámci analýzy vstupných textov sme v druhom šprinte identifikovali kvalifikátory pre štúdium a zamestnanie. Manuálne sme spracovali okolo 100 životopisov, v ktorých sme zachytili rôzne vetné skladby v okolí identifikovaných kvalifikátorov. Z tejto množiny sme vytvorili univerzálnu štruktúru pre zachytenie výskytu kvalifikátorov spolu s pomenovanými entitami v ich okolí.

Pri zaznamenávaní pomenovaných entít sme použili dvojstupňovú kvalifikáciu pomenovaných entít pre český jazyk z článku "Czech Named Entity Corpus and SVM-based Recognizer"¹.

Skratky znázorňujú entity ktoré sa môžu nachádzať v okolí kvalifikátora, pričom čísla znázorňujú interval vzdialeností tejto entity od pozície kvalifikátora.

¹ KRAVALOVÁ, Jana; ŽABOKRTSKÝ, Zdeněk. Czech named entity corpus and SVM-based recognizer. In: Proceedings of the 2009 Named Entities Workshop: Shared Task on Transliteration. Association for Computational Linguistics, 2009. p. 194-201

Štúdium	Zamestnanie
studovat <IC -1 +2> <GU +2> <TY -1 +2> vystudovat <IC -1 +2> <GU +2> <TY -1 +2> absolvovat <IC +1> <GU +2> <TY -1> absolvování <IC +1> <GU> <TY> být žák <IC -1> <GU -2 +3> <TY -3 +4> studium <IC +2> <GU +1 +3> <TY -1 +1> nastoupit <IC +1> <GU +2> <TY -1> vzdělání <IC +1> <GU +3> <TY +2>	pracovat <IC +2> <GU +2> <TY -1> být pracovník <IC> <GU> <TY> <IF +1> působit <IC -1 +1> <GU +3> <TY -1 +1> <IF> člen <IC +1> <GU +2> <TY -1> <IF> zakladatel <IC +1> <GU -1 +1> <TY> <IF> spoluzakladatel <IC +1> <GU +2> <TY +3> <IF> zamestnanec <IC> <GU> <TY -1> <IF +1>

Tab. 1. Kvalifikátory

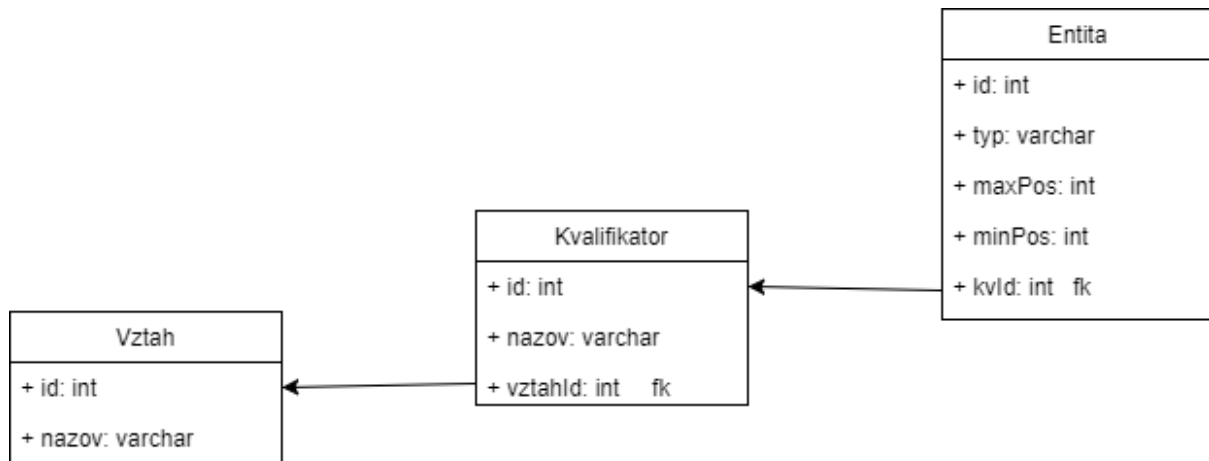
Túto štruktúru sme manuálne simulovali nad množinou nových životopisov a výsledky sme v nich farebne vizualizovali. Táto simulácia dosiahla 80% úspešnosť v rámci zachytenia kvalifikátora definovanou štruktúrou.

Hyliš, Petr,
 Sochař Petr Hyliš se narodil 6. 3. 1956 v Havlíčkově Brodě jako syn sochaře Karla Hyliše. V letech 1971-1975 vystudoval Střední průmyslovou školu kamenickou a sochařskou v Hořicích. Během středoškolského studia navštěvoval soukromé hodiny v ateliéru prof. Karla Lidického. Ve studiích pokračoval v letech 1975-1981 na Vysoké škole uměleckoprůmyslové v Praze u prof. Otto Eckerta. Vyučoval na Pedagogické fakultě Jihočeské univerzity v Českých Budějovicích. Je členem Unie výtvarných umělců České republiky, Sdružení výtvarných umělců Vysočiny, Sdružení sochařů Čech, Moravy a Slezska a tvůrčí skupiny EN FACE '91. Účastnil se mnoha kolektivních výstav u nás i v zahraničí. Velkou samostatnou výstavu uspořádal v r. 2001 v několika regionálních galeriích. Některá ze svých děl realizoval v architektuře ve veřejných prostorech. Žije a pracuje v Havlíčkově Brodě.

- kvalifikátor
- instituce
- geografické pomenovanie
- datácia
- problémová oblasť

Obr. 3: Simulácia

Vzhľadom na výsledky simulácie tejto štruktúry, sme navrhli model relačnej databázy, ktorý sme následne implementovali, a naplnili doposiaľ identifikovanými kvalifikátormi.



Obr. 4: Relačný model

3.1.1. Výber databázy pre ukladanie neštruktúrovaných dát

Pre ukladanie neštruktúrovaných dát sme zvažovali technológie *ElasticSearch*, *Postgre* a ukladanie v klasických textových súboroch.

3.1.2. Výber databázy pre ukladanie štruktúrovaných dát

Zvažovali sme *Neo4j* a *Cayley*, *Cache*, objektové databázy.

Štruktúrované dáta získane zo spracovania životopisov je potrebné v nejakej forme ukladať. K dispozícii sú grafové databázy, objektové alebo klasické relačné databázy. Výhodou grafových databáz je, že vedia dobre reprezentovať vzťahy medzi objektami, čo by bolo v našom prípade užitočné pre ukladanie vzťahov ako sú kolega, zamestnanec, spolužiak Z grafových databáz sme zvažovali najmä *Neo4j* a *Cayley*.

3.1.2.1. Neo4j

V dnešnej dobe asi jedna z najpopulárnejších *opensource* databáz vyvíjaná v jazyku JAVA. Vďaka dobrému *toolingu* pre analýzu a vizualizáciu grafov má dobré využitie v telekomunikačných službách, napríklad pri detekcii *fraudov*, sieťových a IT operáciách atď.

Výhody:

- dlho žijúci projekt s veľkou komunitou
- klient v mnohých jazykoch
- *end to end control* - napr. nie je závislá na externých úložiskách ako *Cayley*
- *RESTful api*

Nevýhody:

- nie je *cloud based* (nemusí sa vždy jednať o nevýhodu)
- nepodporuje *shardovanie*
- nie je zadarmo pre komerčné použitie

3.1.2.2. Cayley

Cayley je neoficiálny *opensource* produkt od spoločnosti *Google* vyvíjaný v jazyku *Go*, inšpirovaný grafovou databázou *Freebase*. Cieľom databázy *Cayley* je vytvorenie tzv. *toolboxu* pre prácu s prepojenými resp. grafovými dátami (sémantický web, sociálne siete, atď.).

Výhody:

- podporuje viac *query* jazykov ako napr. *Gizmo*, *GraphQL* alebo *MLQ*
- zabudovaný *query* editor a vizualizér
- *RESTful API*
- podpora viacerých *backend* úložísk: *kvl (Bolt, LevelDB)*, *NoSQL (MongoDB)*, *SQL (Postgres, CocroachDB, MySQL)*
- modulárny design (jednoduché rozširovanie s novými jazykmi a *backend* riešeniami)
- veľký dôraz na rýchlosť

Nevýhody:

- pomerne mladý projekt - ku dnešnému dňu má iba 1130 *commitov* a 59 prispievateľov
- nemá úplne kompletnú dokumentáciu
- *storage* rieši v externých databázach (môže byť aj výhoda)

3.2. Modul predspracovania textu

3.2.1. Požiadavky

Na to aby sme vedeli získavať z textu poznatky, ktoré by umožnili pochopenie kontextu vety, potrebujeme vytvoriť niekoľko nástrojov, ktoré umožnia zmysluplné dopytovanie nad textom.

Pre tento modul vznikli požiadavky na základné spracovanie neštruktúrovaného textu životopisov do viet a slov. Pre tieto účely potrebujeme niekoľko metód, ktoré by vedeli následne spracovávať text a vytvárať množinu poznatkov o slovách vo vetách.

- Tokenizácia - text potrebujeme vo forme tokenov - samostatných slov
- Lematizácia - potrebujeme určiť základný tvar slov aby sme mohli vyhľadávať kvalifikátory
- Morfológické značky a slovné druhy- potrebujeme zistiť pre každé slovo tieto značky a určiť slovný druh aby sme vedeli lepšie identifikovať entity

3.2.2. Analýza prístupov cez jednotlivé služby

Pre väčšinu jazykov existuje podpora prístupov k slovníkom obsahujúcim tieto vedomosti na úrovni webových služieb alebo voľne dostupných knižníc.

3.2.2.1. NLTK

Natural Language Tool Kit - je knižnica v *Pythone*, primárne robená pre účely spracovania anglického textu. Obsahuje v sebe mnoho slovníkov a funkcionalít, ktoré vedia generovať tagy, lemy, identifikovať slovné druhy alebo aj tokenizovať. Pre riešenie rozdeľovania slov vyhovuje z tohto slovníka najviac podpora tokenizácie textu pre rôzne jazyky.

3.2.2.2. Lemmagen

Je český voľne dostupný morfológický slovník a značkovač určený pre spracovávanie českého textu. Väčšina služieb tejto knižnice vykonáva morfológickú analýzu nad vloženým textom a vie vrátiť morfológické značky, tokeny a lematizované formy slov na základe natrénovaného lingvistického modelu.

Ponúka webové rozhrania s REST API ako aj knižnicu v C++ na lokálne použitie. Príklad použitia na webe je zobrazený na obrázku 1.



The screenshot shows a web interface for the Lemmagen tool. At the top, there is a text input field containing the sentence "Byl jsem venku". Below the input field is a blue button labeled "Process Input". Underneath is a green bar labeled "Output". Below the "Output" bar is a green button labeled "Save Output File". At the bottom, there is a table with three columns: "Token", "Lemma", and "Tag". The table contains three rows of data corresponding to the words in the input sentence.

Token	Lemma	Tag
Byl	být	VpYS---XR-AA---
jsem	být	VB-S---1P-AA---
venku	venku	Db-----

Obr. 5: Príklad použitia na webe

3.2.2.3. MorphoDita

Táto webová služba je dostupná aj cez rozhranie REST API, ktoré nám umožňuje na základe metód *tokenize*, *analyze* a *tag* vrátiť nad poslaným textom výstup vo forme JSON typu kde je každé slovo tokenizované s pridelenou *lemou* a morfológickou značkou.

```
model: "czech-morfflex-pdt-161115"
▼ acknowledgements:
  ▼ 0: "http://ufal.mff.cuni.cz/morphodita#morphodita_acknowledgements"
  ▼ 1: "http://ufal.mff.cuni.cz/morphodita/users-manual#czech-morfflex-pdt_acknowledgements"
▼ result:
  ▼ 0:
    ▼ 0:
      token: "Děti"
      ▼ analyses:
        ▼ 0:
          lemma: "dítě"
          tag: "POS=N|SubPOS=N|Gen=F|Num=P|Cas=1|Neg=A"
        ▼ 1:
          lemma: "dítě"
          tag: "POS=N|SubPOS=N|Gen=F|Num=P|Cas=4|Neg=A"
        ▼ 2:
          lemma: "dítě"
          tag: "POS=N|SubPOS=N|Gen=F|Num=P|Cas=5|Neg=A"
      space: " "
    ▼ 1:
      token: "pojedou"
      ▼ analyses:
        ▼ 0:
          lemma: "jet"
          tag: "POS=V|SubPOS=B|Num=P|Per=3|Ten=F|Neg=A|Voi=A"
      space: " "
```

Obr. 6: Příklad použitia metódy *analyze*

```

model: "czech-morfflex-pdt-161115"
▼ acknowledgements:
  ▼ 0: "http://ufal.mff.cuni.cz/morphodita#morphodita_acknowledgements"
  ▼ 1: "http://ufal.mff.cuni.cz/morphodita/users-manual#czech-morfflex-pdt_acknowledgements"
▼ result:
  ▼ 0:
    ▼ 0:
      token: "Děti"
      lemma: "dítě"
      tag: "NNFP1-----A----"
      space: " "
    ▼ 1:
      token: "pojedou"
      lemma: "jet-1_^(pohybovat_se,_ne_však_chůzí)"
      tag: "VB-P---3F-AA---"
      space: " "
    ▼ 2:
      token: "k"
      lemma: "k-1"
      tag: "RR--3-----"
      space: " "
    ▼ 3:
      token: "babičce"
      lemma: "babička"
      tag: "NNFS3-----A----"
    ▼ 4:
      token: "."
      lemma: "."
      tag: "Z:-----"
      space: " "

```

Obr. 7: Príklad použitia metódy tag

3.2.2.4. KonText

Je to rozhranie, ktoré umožňuje prístupovať ku slovníkom alebo korpusom českého jazyka. Každý z týchto korpusov reprezentuje vedeckú prácu výskumu akademikov českých univerzít. Tie následne poskytujú niekoľko operácií nad svojimi korpusmi. Nevýhoda tohto prístupu je, že nevieme nájsť knižnicu alebo prístup, ktorý by nám jednoducho umožnil lokálne sa dopytovať nad týmto rozhraním. Niektoré operácie dokonca vyžadujú prihlasovanie sa na webovej stránke. Tento koncept nateraz berieme ako nevýhodu a skúsime nájsť alternatívne prístupy.

3.2.2.5. Text.fiit

Je to stránka výskumu fakulty FIIT, ktorá sa zaoberá spracovaním slovenského jazyka. Na tejto stránke máme na výber postačujúce metódy tokenizácie, lematizácie, morfológických značiek a určovania slovných druhov. Tieto metódy sú prístupné vo forme webových služieb a tento prístup berieme ako vhodný pre spracovávanie slovenského jazyka.

3.2.3. Návrh

Z analýzy vyplýva možnosť použitia viacerých prístupov. Najideálnejšie je použitie morfodity na účely spracovania textu. Keďže ale odhadujeme rôznu postupnosť procesov spracovania textu, je teda vhodné navrhnúť volanie týchto služieb v samostatných pod moduloch. To nám umožní vyššiu flexibilitu pri definovaní procesu spracovania textu, kde môžeme implementovať rôzne adaptéri upravujúce vstupy a výstupy. Pre tento účel sa javí vhodné použitie architektonického štýlu dátovody a filtre, ktorý umožňuje rýchlu adaptáciu alebo aj konfiguráciu procesu spracovania textu.

3.2.4. Implementácia

Realizácia návrhu prebehla v duchu architektonického štýlu kde sme vytvorili spracovanie vo forme vstupných dát ako "slovník" a následne aplikujeme metódu, ktorá sa chová ako filter a spracuje text. Tento výstup následne použije pre vstup nasledujúceho filtra. Tento proces vieme konfigurovať.

Implementovali sme aj prístup k rozhraniu MorphoDity, kde je implementácia služieb volajúcich REST API, ktoré vrátia spracovaný JSON vo forme ako je vidieť v analýze.

Implementovali sme aj prototyp pre dopytovanie nad textom v databáze *elasticsearch*. V tomto prototypy sme použili plugin *Lemmagen*, ktorý umožňuje dopytovanie nad textom v lematizovanej forme. Tento prístup je vhodný na identifikáciu životopisov obsahujúcich kvalifikátory čo nám otvára prístup vytvorenia schémy uloženia životopisov v kombinácii s výstupom *MorphoDity*.

Nakoľko pre naše účely by nepostačoval automaticky vygenerovaný *mapping* vytvorili sme index a nastavili *mapping* a *analyzer* nasledovným spôsobom:

```

PUT /cv_cvicny
{
  "settings": {
    "index": {
      "analysis": {
        "filter": {
          "lemmagen_filter_en": {
            "type": "lemmagen",
            "lexicon": "cs"
          }
        },
        "analyzer": {
          "lemmagen_lowercase_en": {
            "type": "custom",
            "tokenizer": "uax_url_email",
            "filter": [
              "lemmagen_filter_en",
              "lowercase"
            ]
          }
        }
      }
    },
    "mappings": {
      "cv": {
        "properties": {
          "text": {
            "type": "text",
            "analyzer": "lemmagen_lowercase_en"
          }
        }
      }
    }
  }
}

```

Obr. 8: Mapping

Po vytvorení indexu a nastavenia *mappingu* a *analyzera* sme do indexu pridali dokumenty a následne sme sa mohli dopytovať dát. Príklad *query*, ktorá vyhľadáva životopisy obsahujúce kvalifikátor "študovať" je uvedený nižšie:

```

GET /cv/_search?pretty
{
  "query": {
    "bool": {
      "must": [
        {
          "match": {
            "text": "studovat"
          }
        },
        {
          "match_phrase": {
            "text": "studovat na"
          }
        }
      ]
    }
  },
  "highlight": {
    "fields": {
      "text": {"fragment_size" : 151}
    }
  },
  "_source": [
    "name"
  ],
  "size": 20
}

```

Obr. 9: Query

Príklad výstupu z uvedenej query:

```
    "_source": {
      "name": "Vokálek, Josef"
    },
    "highlight": {
      "text": [
        "Josef Vokálek se narodil 3. 1. 1887 v Čakovicích, zemřel 30. 6. 1969 v
        Sedlci. Malíř a profesor v Praze. <em>Studoval</em> <em>na</em> pražské akademii
        výtvarných",
        " umění u prof. Krattnera, v r. 1919 <em>studoval</em> <em>na</em> École
        des Beaux Arts v Paříži. Věnoval se krajinářství i figurální malbě, preferoval
        techniku pastelu."
      ]
    }
  }
```

Obr. 10: Výstup z query

3.3. Modul správa používateľov

Naša aplikácia nezahŕňa len vyhľadávanie nad štruktúrovanými dátami, ale aj samotné pridávanie nových textov na spracovanie naším systémom. Potrebujeme teda zabezpečiť, aby sa nám tam nedostali nesprávne a nepravdivé dáta. Toto dosiahneme tým, že zavedieme jednoduchú správu používateľov.

Definovali sme si 3 základné používateľské role:

- Admin

Admin je používateľ, ktorý spravuje systém. Prideluje práva registrovaným používateľom a validuje registráciu.

- Anonymný používateľ

Anonymný používateľ je používateľ systému, ktorý si chce vyhľadať nejakú konkrétnu informáciu v už spracovaných štruktúrovaných dátach. Takýto používateľ nemá práva na vkladanie vlastných dát do systému na spracovanie.

- Registrovaný používateľ

Registrovaný používateľ je taký používateľ, ktorý sa zaregistroval do systému, bola mu schválená registrácia adminom a následne sa prihlásil do systému. Takto prihlásený používateľ, ktorému boli nastavené práva adminom, dokáže do systému vložiť vlastné dáta, ktoré náš systém následne spracuje do štruktúrovanej podoby.

Registrácia používateľa prebieha pomocou jednoduchého formulára kde používateľ vyplní svoje meno, priezvisko, emailovú adresu a zvolí si používateľské meno a prihlasovacie heslo. Vyplnené dáta systém spracuje a vytvorí používateľovi záznam v databáze používateľov. Admin používateľovi registráciu buď schváli alebo zamietne. Pri prihlasovaní používateľ vyplní používateľské meno a heslo, systém skontroluje či existuje registrovaný používateľ s danými prihlasovacími údajmi a na základe výsledku vyhodnotí prihlásenie ako úspešné alebo neúspešné.

3.4. Modul používateľské rozhranie

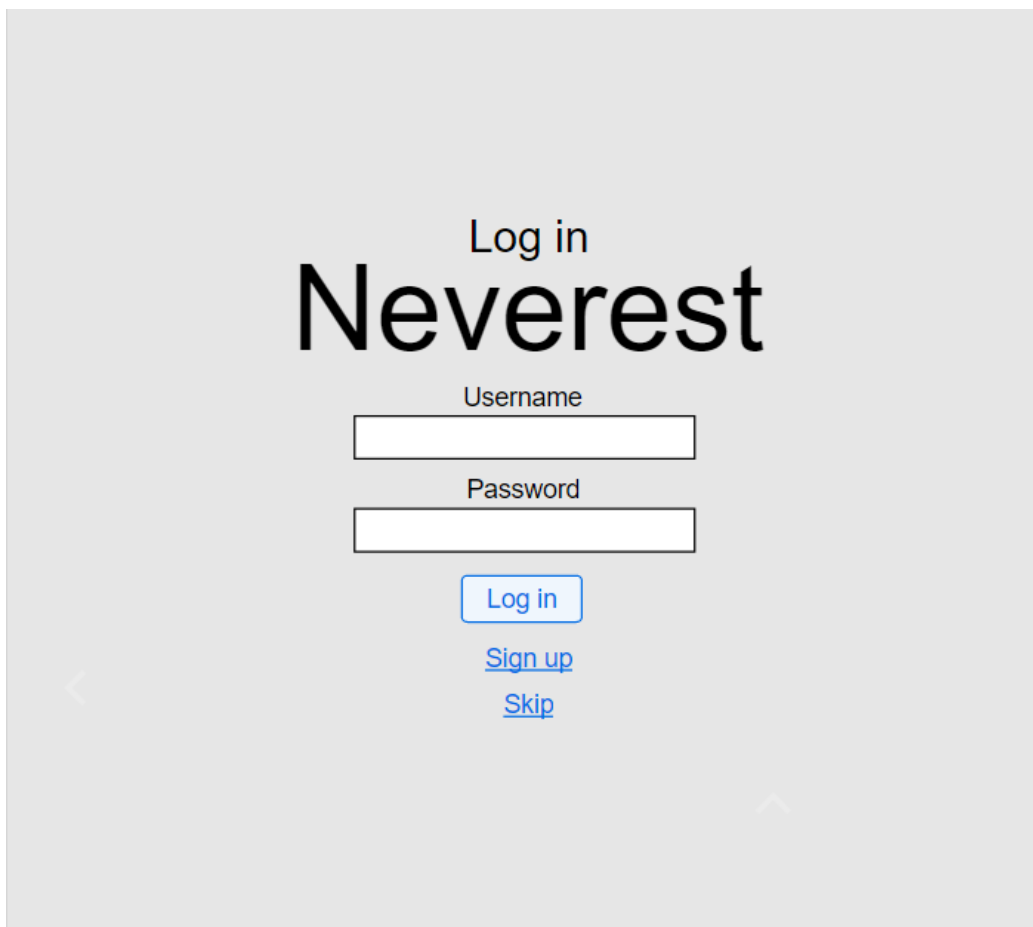
Systém bude sprístupnený používateľom v podobe webového rozhrania. Toto rozhranie bude ponúkať používateľom funkcionality medzi ktorú patri:

- vyhľadávanie informácií o zadaných entitách
- filtrovanie zoznamu entít na základe zadaných kritérií
- spracovanie a zobrazenie vlastného textu do štruktúrovanej podoby
- pridanie štruktúrovaných dát do databázy

Každá funkcionality je dostupná na základe typu účtu pod ktorým je používateľ prihlásený. Preto je potrebné vytvorenie obrazoviek na prihlasovanie a registráciu, ktoré budú spolu s modulom správy používateľov poskytovať možnosť prihlasovania a vedenia účtu.

3.4.1. Prihlásenie

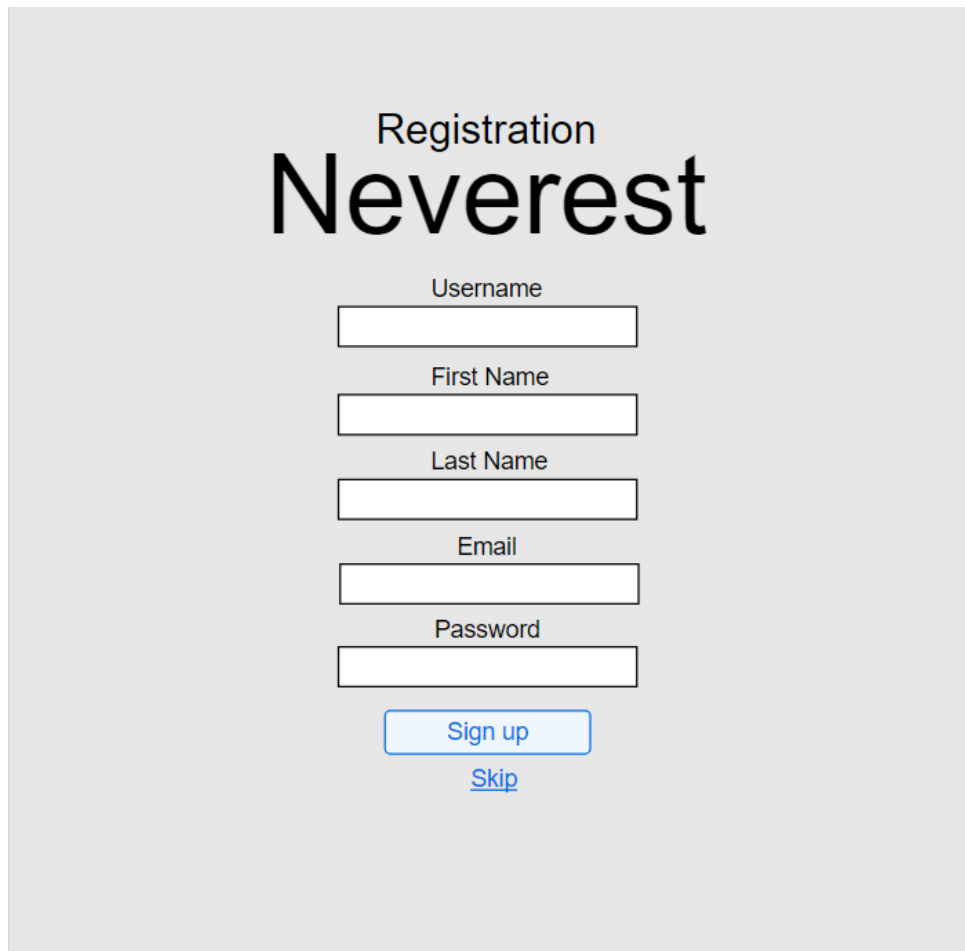
V obrazovke prihlasovania používateľ zadá svoje používateľské meno a heslo a pomocou tlačidla „Log in“ sa prihlási do systému. Nižšie pod týmto tlačidlom je dostupná registrácia nových používateľov, a taktiež preskočenie prihlásenia pre anonymných používateľov, ktorým je následne poskytnutá obmedzená funkcionality systému.



Obr. 11: Obrazovka prihlásenia

3.4.2. Registrácia

Registrácia používateľa sa vykonáva na základe zadaného používateľského mena, hesla a emailu. Používateľ do formulára vyplní tieto údaje, a odošle registráciu tlačidlom „Sign up“. V prípade, že si používateľ registráciu rozmyslí, má možnosť ju preskočiť a prísť k systému anonymne s obmedzenou funkcionalitou. Na tento úkon slúži odkaz „Skip“.

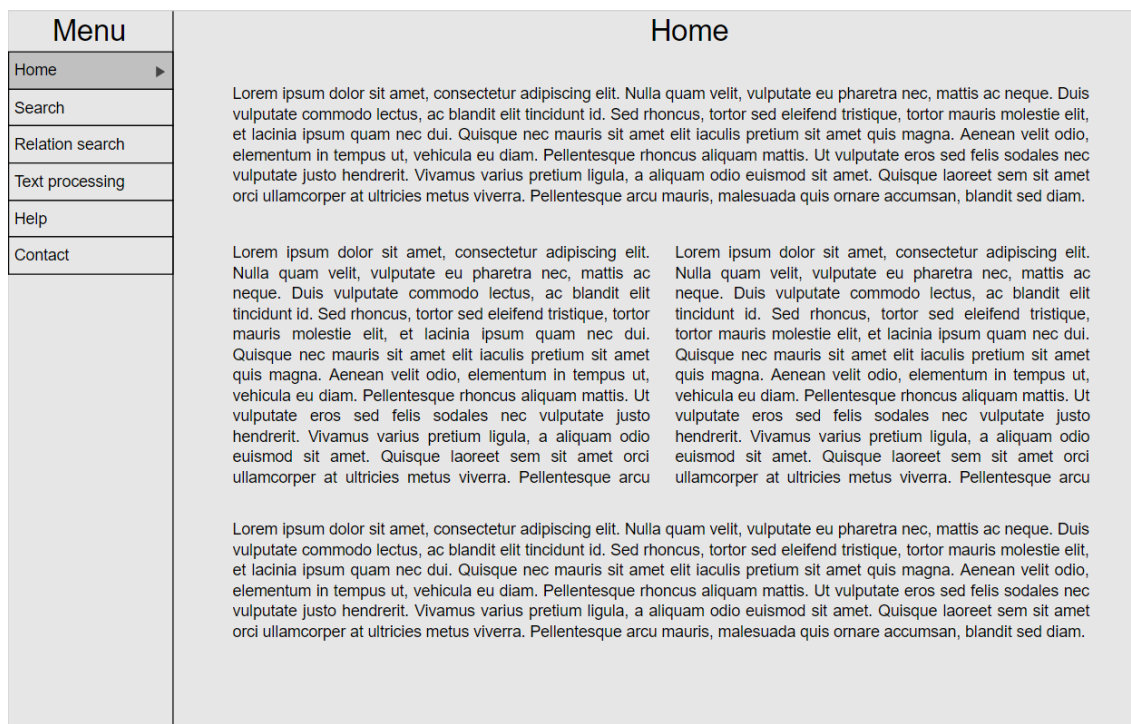


The image shows a registration form titled "Registration Neverest". The form is centered on a light gray background. It consists of the following elements from top to bottom: the title "Registration" in a smaller font and "Neverest" in a large, bold font; a text input field labeled "Username"; a text input field labeled "First Name"; a text input field labeled "Last Name"; a text input field labeled "Email"; a text input field labeled "Password"; a blue button with rounded corners labeled "Sign up"; and a blue text link labeled "Skip".

Obr. 12: Obrázok registrácie

3.4.3. Hlavná obrazovka

Hlavná obrazovka je vertikálne rozdelená na dva panely. V ľavom paneli sa nachádza menu na vyber požadovaných funkcií. Tie sa následne zobrazujú v pravom paneli.



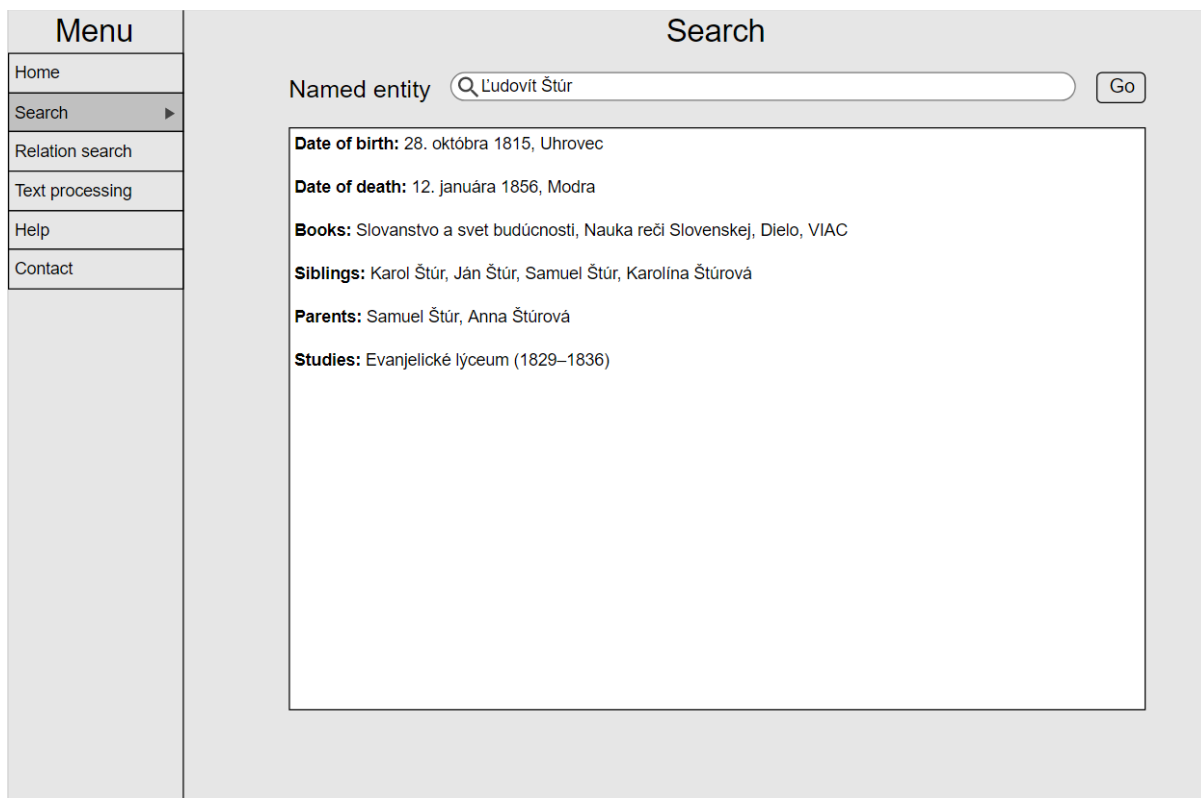
Obr. 13: Hlavná obrazovka

3.4.4. Vyhľadávanie

Funkcionalita vyhľadávania je v podobe jednoduchého vyhľadávacieho poľa v databáze na základe kľúčových slov zadaných používateľom. Používateľovi sa zobrazí zoznam výsledkov, v ktorom si následne môže zobrazíť detailne informácie o požadovaných entitách.

3.4.5. Vyhľadávanie vzťahov

Pri vyhľadávaní vzťahov si používateľ zadáva viacero kritérií vyhľadávania. Tieto kritéria sú v používateľskom rozhraní reprezentované filrami, ktoré je možné pridávať a odoberať podľa potreby. V týchto filtroch používateľ vie definovať časovú os, korporáciu, miesto, vzťah, osobu, podľa toho čo vyhľadáva. Výsledkom tohto vyhľadávania je kolekcia všetkých záznamov, ktoré zodpovedajú zadaným kritériám.



Obr. 14: Obrázok vyhľadávania

3.4.6. Spracovanie textu

Spracovanie textu je nástroj webového rozhrania, ktorý umožňuje používateľovi z neštruktúrovaného textu životopisu vyťažiť informácie o štúdiu. Výsledné informácie sú vo formáte XML. Na spracovanie textu slúži tlačidlo „Process“. Tlačidlo „Upload“ slúži na vloženie štruktúrovaných dát do databázy. Táto možnosť je dostupná iba používateľom s potrebnými právami prístupu.

Menu

- Home
- Search
- Relation search ▶
- Text processing
- Help
- Contact

Relation search

Named entity type ▼

Person ▼

Relation ▼

Study ▼ ✕

Date from ▼

4/22/2015 📅 ✕

Date to ▼

4/22/2016 📅 ✕

Corporation ▼

STU ✕

Number of people found: 4

▼ Name	▼ Surname	▼ Date	▼ Corporation	▼ Position
Matej	Adamov	2014-2017	FIIT STU	Student
Peter	Berta	2014-2017	FIIT STU	Student
Michal	Krempaský	2014-2017	FIIT STU	Student
Oliver	Macko	2014-2017	FIIT STU	Student
Broňa	Pečíková	2014-2017	FIIT STU	Student

Obr. 15: Vyhľadávanie podľa vzťahov

Menu	Text processing
Home	<p>Input - text:</p> <div style="border: 1px solid black; height: 100px; width: 100%;"></div> <p style="text-align: right;"><input type="button" value="Process"/></p>
Search	
Relation search	
Text processing ▶	
Help	
Contact	
	<p>Output - XML:</p> <div style="border: 1px solid black; height: 100px; width: 100%;"></div> <p style="text-align: right;"><input type="button" value="Upload"/></p>

Obr. 16: Obrazovka spracovania textu