

Západočeská univerzita v Plzni  
Fakulta aplikovaných věd  
Katedra informatiky a výpočetní techniky

## **Diplomová práce**

# **Tvorba rozsáhlých úložišť patentových dat**

Místo této strany bude  
zadání práce.

# Prohlášení

Prohlašuji, že jsem diplomovou práci vypracoval samostatně a výhradně s použitím citovaných pramenů.

V diplomové práci jsou použity názvy programových produktů, firem apod., které mohou být ochrannými známkami nebo registrovanými ochrannými známkami příslušných vlastníků.

V Plzni dne 3. května 2022

Bc. Vojtěch Danišík

# Poděkování

Děkuji panu Doc. Ing. Daliboru Fialovi, Ph.D. za ochotu při vedení diplomové práce a rady s jejím vypracováním.

## **Abstract**

Creation of large-scale patent data repositories. The aim of the diploma thesis is to get acquainted with the available sources of patent data and to create extensive local repositories of patent data enabling their effective searching and mining. The first part of the thesis thoroughly describes the types of patents, existing data sources and file formats in which patents are stored. Subsequently, the applicable technologies for searching and mining are described. The second part of the thesis is devoted to the selection of usable data and the implementation of selected technologies. Several queries and scenarios have been created to test efficient mining. The results of the testing are part of this work.

## **Abstrakt**

Cílem diplomové práce je seznámit se s dostupnými zdroji dat o patentech a vytvořit rozsáhlá lokální úložiště patentových dat umožňující jejich efektivní prohledávání a vytěžování. První část práce důkladně popisuje typy patentů, existující zdroje dat a formáty souborů, ve kterých se patenty ukládají. Následně jsou popsány použitelné technologie pro prohledávání a vytěžování. Druhá část práce se věnuje výběru použitelných dat a implementaci vybraných technologií. Pro otestování efektivního vytěžování bylo vytvořeno několik dotazů a scénářů. Výsledky testování jsou součástí této práce.

# Obsah

<b>1</b>	<b>Úvod</b>	<b>1</b>
<b>2</b>	<b>Patent</b>	<b>2</b>
2.1	Typy patentů . . . . .	2
2.1.1	Užitný patent . . . . .	2
2.1.2	Návrhový patent . . . . .	3
2.1.3	Patent rostlin . . . . .	4
2.2	Příhláška vs Publikace . . . . .	4
2.3	Obory . . . . .	4
<b>3</b>	<b>Databáze</b>	<b>6</b>
3.1	Systém řízení báze dat . . . . .	6
3.2	Komponenty databáze . . . . .	7
3.3	Typy databází . . . . .	7
3.3.1	Relační databáze . . . . .	7
3.3.2	Objektově-orientovaná databáze . . . . .	11
3.3.3	NoSQL databáze . . . . .	13
3.3.4	Databáze Klíč-Hodnota . . . . .	13
3.3.5	Grafová databáze . . . . .	14
3.3.6	Databáze dokumentů . . . . .	16
3.4	Existující řešení . . . . .	18
3.4.1	MySQL . . . . .	18
3.4.2	PostgreSQL . . . . .	18
3.4.3	LevelDB . . . . .	19
3.4.4	MongoDB . . . . .	20
3.4.5	Neo4j . . . . .	20
3.5	Jazyky . . . . .	21
3.5.1	Structured Query Language . . . . .	23
3.5.2	MongoDB Query Language . . . . .	23
3.5.3	Cypher Query Language . . . . .	25
<b>4</b>	<b>Návrh úložiště</b>	<b>26</b>
4.1	Výběr patentů . . . . .	26
4.1.1	Zdroje dat . . . . .	27
4.1.2	Atributy . . . . .	30
4.1.3	Závěr průzkumu . . . . .	34

4.2	Výběr databáze . . . . .	36
4.2.1	Výběr typu databáze . . . . .	36
4.2.2	Výběr z existujících řešení . . . . .	38
<b>5</b>	<b>Implementace úložiště</b>	<b>39</b>
5.1	Adresářová struktura . . . . .	39
5.1.1	Patenty . . . . .	39
5.1.2	Docker . . . . .	39
5.2	Implementace databáze . . . . .	39
5.2.1	MySQL . . . . .	39
5.2.2	MongoDB . . . . .	39
5.3	Výsledné úložiště . . . . .	39
5.3.1	Technologické požadavky . . . . .	39
5.3.2	Docker . . . . .	39
5.3.3	Inicializace MySQL . . . . .	39
5.3.4	Inicializace MongoDB . . . . .	39
<b>6</b>	<b>Rozšiřitelnost úložiště</b>	<b>40</b>
6.1	Přidávání nových patentů . . . . .	40
6.2	Zjišťování autorů pro české patenty . . . . .	40
6.3	Automatické stahování dat z ověřených zdrojů . . . . .	41
<b>7</b>	<b>Ověření efektivního vytěžování</b>	<b>43</b>
7.1	Mongo + ElasticSearch . . . . .	43
7.2	MySQL . . . . .	43
7.2.1	Scénář č.1 . . . . .	43
7.2.2	Scénář č.2 . . . . .	44
7.2.3	Scénář č.3 . . . . .	45
7.2.4	Scénář č.4 . . . . .	46
7.2.5	Scénář č.5 . . . . .	47
7.2.6	Scénář č.6 . . . . .	47
7.2.7	Scénář č.7 . . . . .	48
7.2.8	Scénář č.8 . . . . .	49
7.2.9	Scénář č.9 . . . . .	49
<b>8</b>	<b>Závěr</b>	<b>51</b>
	<b>Zkratky</b>	<b>52</b>
	<b>Literatura</b>	<b>53</b>
<b>A</b>	<b>Uživatelská dokumentace</b>	<b>57</b>





# 1 Úvod

První odstavec by byl o patentu. Jednoduše by se popsalo proč vlastně je potřeba patent, proč ho chce zadávající, lehce popsat co to je atd.

Druhý odstavec by byl o databázi, lehký popis jako k čemu to je, jak je to např. rozšířený atp.

Cílem této práce je se seznámit s dostupnými zdroji dat o patentech a vytvořit rozsáhlá lokální úložiště patentových dat umožňující jejich efektivní vytěžování. Zdroje dat musí poskytovat své databáze patentů (žádosti i publikace) zdarma a patenty musí obsahovat předem stanovené povinné atributy, aby je bylo možné použít. Získaná data budou následně uložena do specifického typu databáze, která bude umožňovat co nejefektivnější vytěžování uložených dat. To znamená rychlé vyhledávání správných výsledků v relativně krátkém čase pro miliony (až desítky milionů) záznamů. Import dat bude řešen pomocí jednoduché aplikace, která bude procházet všechna data a filtrovat ty patenty, kterým chybí některé povinné údaje (i přes to, že struktura obsahuje elementy, ve kterých se údaj má nacházet), nebo jsou nevalidní.

## 2 Patent

todo

Test [22, 30]

### 2.1 Typy patentů

Patenty jsou klasifikovány do tří hlavních typů: **Užitný patent**, **Návrhový patent** a **Patent rostlin**. V následujících kapitolách jsou tyto typy podrobně popsány a srovnány s jejich právně jednoduššími protějšky.

#### 2.1.1 Užitný patent

Užitný patent, nebo také patent na vynález, poskytuje právní ochranu lidem, kteří vynalézají nový a užitečný proces, výrobní předmět, stroj, složení hmoty nebo užitečné vylepšení. Užitný patent přetrvává 20 let od data podání, pokud jsou placeny udržovací poplatky. Užitné patenty jsou nejběžnějším typem patentů, a lze ho dále rozdělit na 2 podtypy: **Trvalý** a **Dočasný**. [24, 30].

#### Struktura přihlášky užitného patentu

Struktura užitného patentu obsahuje sedm různých sekcí (v případě USPTO), které patent musí obsahovat (pro dočasný patent jsou povinné pouze dvě sekce) [39]:

- **Pozadí** - Sekce **Pozadí** popisuje současný stav techniky týkající se navrhovaného vynálezu a neměl by se o tomto vynálezu zmiňovat. **Pozadí** musí být velmi stručné a nesmí být zveřejněno více informací, než je požadováno.
- **Stručný přehled** - Sekce **Stručný přehled** shrňuje oddíl **Podrobný popis** do několika odstavců.
- **Stručný popis nákresů** - Sekce **Stručný popis nákresů** obsahuje široký přehled nákresů v oddílu **Nákresy**. Lze zde zmínit a identifikovat části vynálezu v každém z obrázků nebo například uvádět různé pohledy nákresů (perspektivní, pohled zleva a zprava, ...).
- **Podrobný popis** - Sekce **Podrobný popis** popisuje různá alternativní provedení a vlastnosti vynálezu a měla by obsahovat referenční

číslice, které odkazují na obrázky v oddílu **Nákresy**. Tato je povinná i pro dočasnou přihlášku.

- **Nároky** - Sekce **Nároky** popisuje rozsah ochrany poskytované patentem nebo ochranu požadovanou v patentové přihlášce.
- **Abstrakt** - Sekce **Abstrakt** slouží jako krátké shrnutí o čem patent je. Abstrakt nesmí být delší jak 150 slov.
- **Nákresy** - Sekce **Nákresy** obsahují nákresy popisující vlastnosti vynálezu v různých náhledech (perspektivní, ...). Tato je povinná i pro dočasnou přihlášku.

### **Užitný patent vs Užitný vzor**

Patentová ochrana není jedinou formou průmyslově právní ochrany vynálezu. Pro vynálezy s nižší vynálezeckou úrovní je možné zvolit **Užitný vzor**, který je jednodušší, rychlejší a méně nákladnou alternativou k užitému patentu.

#### **2.1.2 Návrhový patent**

Návrhový patent, jinak zvaný jako Průmyslový vzor, je patent vydaný pro originální, nové a ornamentální vzory pro vyráběné výrobky. Vzhledem k tomu, že se návrhový patent projevuje vzhledem, atk se předmět patentové přihlášky může týkat konfigurace nebo tvaru předmětu, povrchové výzdoby aplikované na předmět nebo kombinace konfigurace a povrchové výzdoby. Návrh povrchové výzdoby je neoddělitelný od předmětu, na který je aplikován, a nemůže existovat sám. Lze získat návrhový i užitný patent na jeden produkt v případě, že z technologického a designového hlediska jsou neoddělitelné. Návrhový patent má platnost (stejně jako u ostatních patentů) od 5 až do 25 let (každá země může mít jinak nastavenou platnost) [21].

#### **Struktura přihlášky návrhového patentu**

Struktura přihlášky návrhového patentu by měla obsahovat tyto prvky (v případě USPTO):

- Preambule s uvedením jména přihlašovatele, názvu návrhového patentu a stručného popisu použití předmětu, ve kterém je tento návrh ztělesněn.
- Křížový odkaz na související užitkovou žádost.

- Prohlášení týkající se federálně sponzorovaného výzkumu nebo vývoje.
- Popis obrázků a výkresů.
- Popis funkce.
- Jeden nárok na návrh.
- Výkresy a fotografie.
- Vykonaný slib nebo přísaha.

### Návrhový patent vs Ochranná známka

Ochranná známka je označení pro schopné grafické znázornění, tvořené zejména slovy, písmeny, číslicemi, barvou, kresbou nebo tvarem výrobku či jeho obalu, určené k rozlišení výrobků nebo služeb[26], zatímco návrhový patent se vztahuje na okrasný vzhled výrobku jako takového, který není spojen s identifikací zdroje zboží [31].

#### 2.1.3 Patent rostlin

Patent na rostlinu je vydáván na novou nebo odlišnou odrůdu rostliny, která byla vynalezena nebo objevena a asexuálně rozmnožována. Tyto patenty se udělují na 20 let od data podání přihlášky a neplatí se žádné udržovací poplatky [23]. Účelem těchto patentů je zamezení duplikování rostlin nebo jejich nabízení k prodeji (jako celek nebo po částech).

#### Struktura přihlášky patentu rostlin

Struktura patentu pro rostliny (v případě USPTO) je až jednu výjimku stejná jako pro užitný patent (viz kapitola č. 2.1.1). Výjimka se týká podrobného popisu, který musí obsahovat kompletní botanický popis dané rostliny a rozdíly, které tuto rostlinu odlišují od ostatních, již existujících rostlin.

## 2.2 Přihláška vs Publikace

todo

## 2.3 Obory

Každý patent je po podání přihlášky klasifikován do jednoho oboru podle klasifikačního systému IPC. **International Patent Classification** (IPC) je

hierarchický systém, který používá jazykově nezávislé symboly pro klasifikaci patentů a užitných vzorů podle oblasti technologie, ke které se vztahují [9]. Tento systém je používán ve více než 100 zemí.

IPC schéma definuje celkem osm hlavních oblastí technologie, do kterých lze patent přiřadit. Každá oblast je symbolizována jedním písmenem (viz tabulka č. 2.1).

<b>Kód</b>	<b>Popisek</b>
A	Lidské potřeby - jídlo, léky, oblečení, ...
B	Operace a Doprava - tisk, auta, koleje, nanotechnologie, ...
C	Chemie a Hutnictví - sklo, cement, železo, ...
D	Textílie a Papír - výroba papíru, provazy, tkalcovství
E	Pevné konstrukce - stavby, dveře, zámky, dolování, ...
F	Strojírenství, Osvětlení, Vytápění, Zbraně, Odstřelování
G	Fyzika - měření, testování, optika, nukleární fyzika, ...
H	Elektřina - základní elektrické elementy (kabely, rezistory, ...), techniky elektrické komunikace, ...

Tabulka 2.1: Základní IPC klasifikace oborů [12].

Klasifikační symboly definované IPC jsou tvořeny písmenem označující hlavní oblast technologie (neboli sekce), následovaným dvěma číslicemi označující třídu, poté písmenem označujícím podtřídu. Po podtřídě následuje proměnná čítající jednu až čtyři číslice, která označuje hlavní skupinu. Následuje dopředné lomítko a číslo čítající dvě až šest číslic označující podskupinu [9]. Příklad klasifikace lze vidět v tabulce č. 2.2.

<b>Divize</b>	<b>Počet divizí</b>	<b>Symbol</b>	<b>Titulek</b>
Sekce	8	G	Fyzika
Třída	120	G01	Měření; Testování
Podtřída	628	G01G	Vážení
Skupina	+ - 6 900	G01G 21	Podrobnosti o vážení aparátu
Podskupina	+ - 62 100	G01G 21/02	Podrobnosti o vážení ložisek s ostřím nože

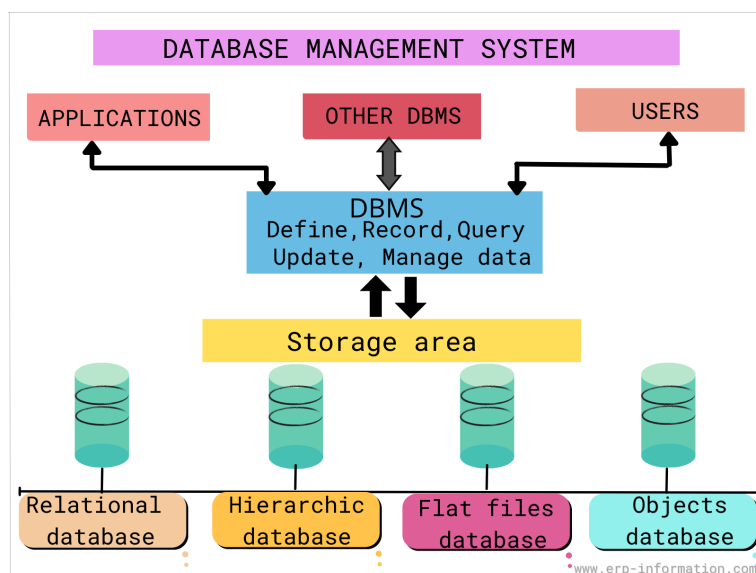
Tabulka 2.2: Příklad IPC klasifikace [28].

## 3 Databáze

Termín databáze označuje organizovanou kolekci strukturovaných informací nebo dat, která jsou typicky ukládána elektronicky v počítačovém systému. Data / informace lze nazvat jako fakta vztahující se k libovolnému uvažovanému objektu. Typický příklad objektu je člověk, jehož fakta jsou: jméno, věk, výška, váha a mnoho dalších [34].

### 3.1 Systém řízení báze dat

Pro správu dat v databázi a její řízení je potřeba komplexní software, který se nazývá **Systém Řízení Báze Dat** (SŘBD, anglicky DBMS). SŘBD slouží jako interface mezi samotnou databází a koncovým uživatelem (může být i program), umožňující jak vytěžování a aktualizaci dat, tak i možnosti nastavení záloh a jiných administrativních operací [1]. V dnešním světě existuje několik různých DBMS (například Relační DBMS, Objektově-orientované DBMS).



Obrázek 3.1: Systém řízení báze dat [10].

## 3.2 Komponenty databáze

Všechny databáze sestávají z pěti základních komponent, nehledě na použitý typ databáze [33, 34]:

- **Hardware** - Fyzické stroje (počítače, servery, pevné disky, ...) na kterých běží databázový software.
- **Software** - Databázový software poskytuje uživateli / programu kontrolu nad databází. Zahrnuje to samotný databázový software, operační systém, software pro správu sdílení dat mezi uživateli a programy pro přístup k datům v databázi.
- **Data** - Nezpracované a neorganizované fakty, které je potřeba zpracovat. Administrátor databáze organizuje tyto data a dává jim význam. Data se obecně skládají hlavně z faktů, observací, percepce, čísel, znaků a mnoho dalších.
- **Jazyk** - Typický příklad použití jazyku je přístup k datům, přidávání nových dat, úpravu již existujících dat z databáze. Uživatel / program napíše specifické příkazy v jazyku pro přístup k datům (Database Access Language) a tyto příkazy následně pošle databázi ke zpracování. Více viz kapitola č. 3.5.
- **Procedury** - Procedura obsahuje předpřipravený seznam příkazů, které se následně vykonávají po zavolání dané procedury.

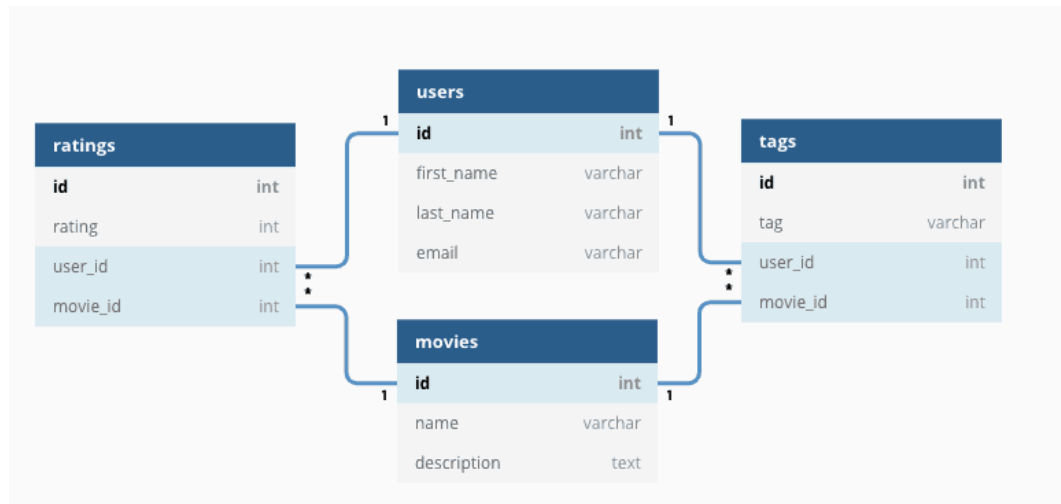
## 3.3 Typy databází

V dnešním světě existuje mnoho různých typů databází. Výběr nejlepšího typu databáze pro konkrétní organizaci závisí na tom, jak organizace zamýšlí data používat. V této kapitole je vypsáno pouze pár typů, protože vznikají stále nové, méně známé typy databází, které jsou tvořeny pro specifické požadavky (například finanční, vědecké) [1, 13].

### 3.3.1 Relační databáze

Název relační databáze pochází ze způsobu, jakým jsou data uložena, a to ve více souvisejících tabulkách. Data v tabulkách jsou uložena v řádcích a sloupcích. Relační databáze jsou velice spolehlivé a podporují všechny čtyři žádoucí vlastnosti databázových transakcí ACID. Pro co nejeфекtivnější využití tohoto typu databáze je potřeba ukládat pouze dobře strukturovaná

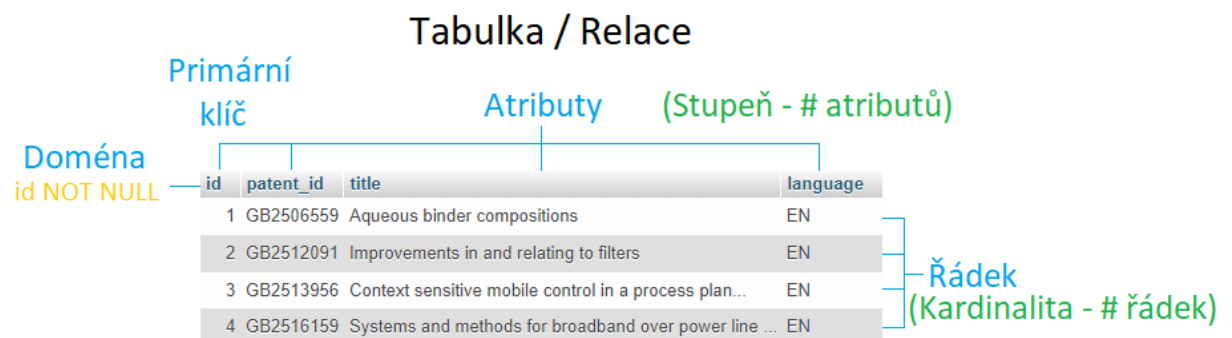
data, pro částečně strukturovaná či nestrukturovaná data je vhodné použít například grafové nebo dokumentově založené databáze. Typické relační databáze jsou například: Microsoft SQL Server, Oracle Database, MySQL. Ukázkou relační databáze lze vidět na obrázku č. 3.2.



Obrázek 3.2: Ukázka relační databáze.

## Datový model

Relační datový model obsahuje několik fundamentálních konceptů. Koncepty lze vidět na obrázku č. 3.3.



Obrázek 3.3: Koncepty relačního datového modelu [8].

První z konceptů se nazývá **Relace**, což je dvou-dimenzionální tabulka, která se používá pro ukládání kolekce datových elementů. Tabulka je tvořena řádky a sloupce, kde řádky reprezentují záznamy a sloupce reprezentují



atributy.

**Řádka** je další koncept relačního modelu, která pouze reprezentuje jeden záznam v tabulce.

Další z konceptů je **Atribut**, který reprezentuje sloupec v tabulce, neboli vlastnosti jednotlivých řádků (například jméno, příjmení, věk, ...).

Koncept **Doména atributů** slouží k definici vlastností pro každou hodnotu daného atributu. Pomocí domény lze určit, zda hodnoty daného atributu mohou být prázdné, budou dlouhé maximálně 50 znaků nebo například určit datový typ atributu (textová hodnota, číslo, ...).

Další z konceptů je **Stupeň**, který pouze určuje počet atributů v dané relaci.

**Kardinalita** určuje počet řádků / záznamů existujících v dané relaci.

Koncept **Relační schéma** popisuje návrh a strukturu relace. Obsahuje názvy tabulek, jejich atributy a typy atributů. Relační schéma pro naši tabulku lze vidět na obrázku č. 3.4.

```
PATENTS(id INT NOT NULL,  
        patent_id VARCHAR(50),  
        title VARCHAR(300),  
        language VARCHAR(2)  
        )
```

Obrázek 3.4: Relační schéma pro tabulku na obrázku č. 3.3.

**Relační instance** reprezentuje kolekci záznamů, které jsou uloženy v tabulce v určitém čase.

Poslední koncept **Relační klíč** je atribut / seznam atributů, které lze využít jako unikátní identifikátor jedné entity v tabulce, případně k určení vazby mezi dvěma relacemi. Existuje šest typů relačních klíčů - kandidátní, super, složený, primární, cizí, sekundární / alternativní [27].

## Výhody

Výhody relační databáze jsou [25]:

- **Jednoduchost modelu** - Při porovnávání ostatních typů databází s relačním, relační databáze je o mnoho jednodušší. Díky tomu, že zde neprobíhá žádné zpracování dat, tak není potřeba využívat žádné složité dotazy.
- **Snadné použití** - Uživatelé mohou jednoduše přistupovat a získávat všechny potřebné informace v rámci sekund bez ohledu na složitost

databáze.

- **Přesnost** - Relační databáze jsou dobře uspořádané a velice striktně definované. I za pomoci primárních a cizích klíčů se v databázi udržuje unikátnost hodnot, takže se zde nevyskytují žádné duplikáty.
- **Integrita dat** - Integrita dat zajišťuje konzistentnost všech tabulek v databázi, díky čemuž lze dosáhnout vlastností jako přesnost a snadné použití.
- **Normalizace** - Normalizace je metoda, pomocí které lze rozdělit jednu informaci do několika bloků za účelem snížení velikosti.
- **Spolupráce** - Více uživatelů může přistupovat k datům ve stejný čas i v případě, že část dat je upravována.
- **Bezpečnost** - Bezpečnost je zajištěna autorizací uživatelů, kdy pouze uživatelé s právy a přístupovými údaji mohou přistupovat k datům.

## Nevýhody

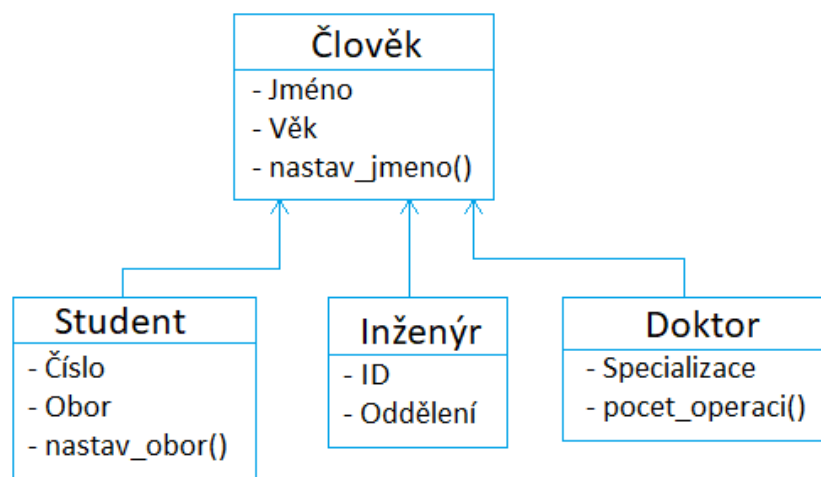
Nevýhody relační databáze jsou [25]:

- **Problém s údržbou** - Údržba relační databáze se stává postupem času náročnější vzhledem ke zvýšenému počtu uložených dat.
- **Cena** - Systém relační databáze je drahý k pořízení i pro správu. Samotná prvotní cena systému je relativně drahá pro menší byznys, ale zhoršuje se při zohlednění najímání profesionálních techniků, které musí mít dobré znalosti ohledně používaného systému.
- **Fyzické úložiště** - Relační databáze jsou složeny z řádků a sloupců, které potřebují hodně fyzické paměti, protože každá provedená operace závisí na samostatném úložišti.
- **Malá škálovatelnost** - Při používání relační databáze na více serverech se její struktura mění a stává se obtížně zvládnutelnou, zejména při velkém objemu dat. Jak se databáze zvětšuje nebo více distribuuje s větším počtem serverů, tak se zvětšuje latence a problémy s dostupností, které ovlivňují celkový výkon.
- **Složitost struktury** - Relační databáze dokáží ukládat data pouze v tabulkové formě, která neumožňuje vyobrazit složitější vazby mezi objekty. Toto může být velký problém u dost aplikací, u kterých data nelze reprezentovat pouze jednou tabulkou díky jejich aplikační logice.

- **Snížení výkonu postupem času** - S větším množstvím uložených dat a tabulek se zvětšuje i složitost systému, díky čemuž bude systém reagovat pomaleji na dotazy, případně může i spadnout v případě více dotazů od více uživatelů.

### 3.3.2 Objektově-orientovaná databáze

Objektově-orientovaná databáze je založena na objektově-orientovaném programování, kdy data a všechny jejich atributy a metody jsou svázány dohromady jako objekt. Stejně jako relační databáze, i objektově-orientované databáze odpovídají standardům ACID. Typické příklady jsou například: ObjectStore, ConceptBase. Ukázku objektově-orientované databáze lze vidět na obrázku č. 3.5.



Obrázek 3.5: Ukázka objektově-orientované databáze.

### Datový model

V objektově-orientovaném modelu jsou data a jejich vztahy mezi sebou uloženy v jediné struktuře, která se jinak nazývá objekt. Všechny objekty mají mezi sebou vícenásobné vztahy. Jednoduše řečeno, objektově-orientovaný datový model je spojením relačního databázového modelu a objektově-orientovaného programování [8].

V datovém modelu existují tyto komponenty:

- **Objekt** - Objekt je abstrakcí jakékoliv entity z reálného světa, jinak zvaná jako instance jedné třídy. Objekt zapouzdřuje data a funkční kód do celku, který poskytuje pouze datovou abstrakci, zatímco schovává

implementační detaily od uživatele. Příklad objektů z obrázku č. 3.5: *Student*, *Doktor* a *Inženýr* jsou instancí celku *Člověk*.

- **Atribut** - Atribut popisuje vlastnosti objektu. Například *Student* obsahuje atributy *Číslo* a *Obor*.
- **Metoda** - Metoda reprezentuje chování objektu. Například *Student* obsahuje metodu s názvem *nastav\_obor*, pomocí které můžeme získat studovaný obor daného studenta.
- **Třída** - Třída je vlastně kolekce podobných objektů, které sdílejí strukturu (neboli atributy) a chování (neboli metody).
- **Dědičnost** - Vytvořenému objektu se říká instance třídy, která zdědí kopie / instance všech atributů a metod dané třídy. *Student*, *Doktor* i *Inženýr* dědí od celku *Člověk* atributy *Jméno*, *Věk* a metodu *nastav\_jmeno*.

## Výhody

Výhody objektově-orientované databáze jsou [37]:

- **Výkonnost** - Mnohonásobně výkonnější než relační databáze.
- **Rozšiřitelnost** - lze vytvářet nové datové typy z již existujících. Jako příklad lze uvést vytvořením super třídy, která bude obsahovat všechny společné atributy a metody. Tímto lze snížit redundanci v systému.
- **Podpora velkého množství datových typů** - Oproti ostatním typům databáze, objektově-orientovaná databáze podporuje ukládání různých typů dat, jako například obrázky, zvuky, video a mnoho dalších.
- **Podpora vývoje schématu** - Těsné propojení mezi daty a aplikacemi činí vývoj schématu více proveditelnějším.
- **Podpora pro dlouhotrvající transakce** - Objektově-orientovaná databáze využívá jiný protokol pro zpracovávání dlouhotrvajících transakcí než relační databáze.

## Nevýhody

Nevýhody objektově-orientované databáze jsou [37]:

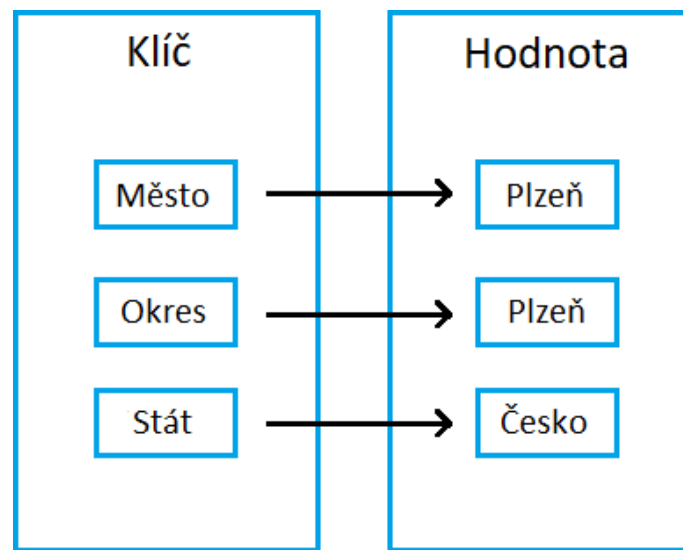
- **Neexistující univerzální datový model** - V dnešní době stále neexistuje univerzální datový model, navíc většině modelů chybí teoretický základ.
- **Nedostačující standardy** - Pro objektově-orientované databáze neexistují žádný univerzální datový model, stejně jako standardní dotazovací jazyk.
- **Složitost** - Funkcionality jako například dlouhotrvající transakce, zpráva verzí nebo evoluce schémat činí výsledný systém mnohonásobně složitější, což vede k vyšší ceně a složitějšímu používání.
- **Zabezpečení** - V databázi neexistuje adekvátní zabezpečovací systém, který by mohl přiřazovat přístupová práva na objekty nebo třídy.

### 3.3.3 NoSQL databáze

NoSQL je široká kategorie databází, které nepoužívají SQL jako svůj primární jazyk pro přístup k datům. Tyto typy databází jsou také někdy označovány jako nerelační databáze. V NoSQL databázích se pracuje s nestrukturovanými a polostrukturovanými sadami distribuovaných dat. Jednou z výhod je, že vývojáři mohou provádět změny databáze za běhu, aniž by to ovlivnilo aplikace, které databázi používají.

### 3.3.4 Databáze Klíč-Hodnota

Databáze klíč-hodnota poskytuje nejjednodušší možný NoSQL datový model. Data jsou uložena jako pár klíč - hodnota ve slovníku / mapě, kdy klíč je indexem. Hodnota může být například celé číslo, řetězec, struktura JSON nebo pole. Z vlastností databáze vyplývá, že zde není potřeba žádný dotazovací jazyk pro získávání výsledků. Typické příklady jsou: Redis, Riak, LevelDB. Ukázkou databáze klíč-hodnota lze vidět na obrázku č. 3.6.



Obrázek 3.6: Ukázka databáze klíč-hodnota.

### Výhody

Výhody této databáze jsou [4]:

- **Škálovatelnost** - Databázi lze škálovat jak vertikálně, tak i horizontálně.
- **Redundance** - Zabudovaná redundance zapříčiňuje větší spolehlivost databáze.
- **Rychlost** - Reakční čas je velice rychlý díky jednoduchosti struktury a jednoduchých příkazů (vložit, smazat, získat).

### Nevýhody

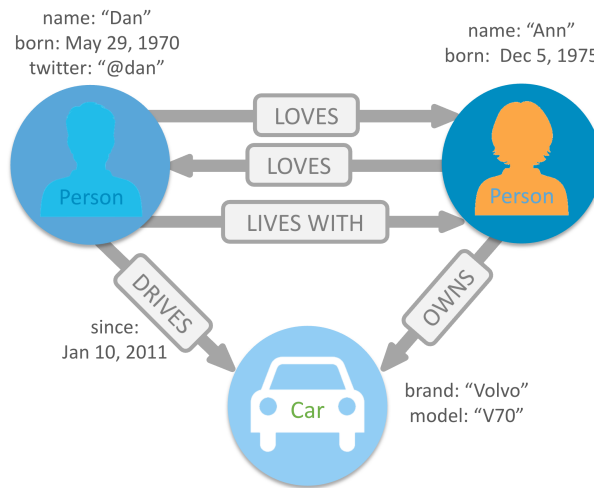
Nevýhody této databáze jsou [4]:

- **Optimalizace dat** - Optimalizace je provedena pouze pro data, kde je pouze jeden klíč a jedna hodnota. V případě ukládání složitějších struktur je potřeba parser.
- **Složitě dotazy** - Nelze používat složité dotazy, pomocí kterých lze vyhledávat specifické hodnoty.

### 3.3.5 Grafová databáze

Grafová databáze je typem NoSQL databáze, která je založená na teorii grafů. Data jsou reprezentována jako uzly, hrany zase reprezentují vztahy

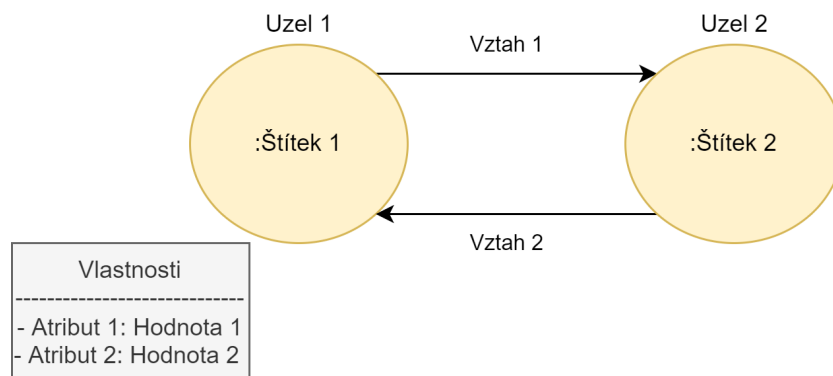
mezi daty. Graf lze procházet podél určitých typů hran nebo přes celý graf. Procházení spojení nebo relací je velmi rychlé, protože vztahy mezi uzly se nepočítají v době dotazu, ale jsou v databázi trvalé. Typické příklady jsou: Neo4j, OrientDB, Microsoft Azure CosmosDB. Ukázkou grafové databáze lze vidět na obrázku č. 3.7.



Obrázek 3.7: Ukáзка grafové databáze [18].

## Datový model

Datový model grafové databáze se skládá ze čtyř komponent (viz obrázek č. 3.8) [7]:



Obrázek 3.8: Komponenty grafu.

- **Uzel** - Uzel je jeden ze dvou fundamentálních komponent který vytváří graf. Uzly slouží k reprezentaci entit nebo jiných doménových komponent.

- **Vztah** - Vztah propojuje dva uzly a dovoluje nám vyhledávat související uzly. Uzel, ze kterého vztah začíná, se jmenuje zdrojový, zatímco uzel, ve kterém vztah končí, se nazývá cílový (šipka ukazuje směr vztahu). Vztahy musí mít vždy jen jeden zdrojový a jeden cílový uzel, proto při mazání uzlů se mažou i všechny jeho závislosti (vstupující a vystupující vztahy).
- **Štítek** - Štítek slouží k zařazování uzlů do skupin. Všechny uzly, které jsou označeny stejným štítkem, patří do jedné skupiny. Uzel může obsahovat libovolné množství štítků (0 až nekonečno). Při vyhledávání může databáze pracovat nejen s celým grafem, ale i s množinou uzlů patřící do jedné skupiny.
- **Vlastnosti** - Vlastnost je množina dvojic klíč - hodnota, které lze ukládat s každým uzlem a vztahem. Jsou podporovány skoro všechny datové typy.

## Výhody

Výhody grafové databáze jsou [3]:

- Struktury jsou flexibilní a přizpůsobivé.
- Reprezentace vztahů mezi entitami je zřetelné.
- Dotazy poskytují výsledky v reálném čase. Rychlost závisí na počtu relací.

## Nevýhody

Nevýhody grafové databáze jsou [3]:

- Neexistuje žádný standardizovaný jazyk. Jazyk závisí na použité platformě.
- Grafy jsou nevhodné pro transakční systémy.
- Je těžké najít podporu, protože uživatelská základna je velice malá.

### 3.3.6 Databáze dokumentů

Databáze dokumentů jsou typem NoSQL databáze a jsou navrženy pro ukládání, načítání a správu informací orientovaných na dokumenty. Typické příklady jsou: MongoDB, Amazon DocumentDB, Elasticsearch. Ukázku dokumentové databáze lze vidět na obrázku č. 3.9.



Document 1	Document 2	Document 3
<pre>{   "id": "1",   "name": "John Smith",   "isActive": true,   "dob": "1964-30-08" }</pre>	<pre>{   "id": "2",   "fullName": "Sarah Jones",   "isActive": false,   "dob": "2002-02-18" }</pre>	<pre>{   "id": "3",   "fullName": {     "first": "Adam",     "last": "Stark"   },   "isActive": true,   "dob": "2015-04-19" }</pre>

Obrázek 3.9: Ukázka dokumentově orientované databáze [32].

## Datový model

Základním prvkem dokumentové databáze je **Dokument**. Definice dokumentů se liší podle konkrétní implementace databáze, ale jedno mají společné: dokumenty kódují zapouzdřená data či informace do nějakého standardního formátu nebo kódování. Mezi typy kódování patří například JSON, XML. Dokument nemusí dodržovat pevně definovanou strukturu, takže v databázi mohou existovat dva dokumenty stejného formátu s rozdílnou strukturou dat [38].

## Výhody

Výhody dokumentové databáze jsou [2]:

- **Bez schématu** - Neexistují zde žádná omezení ve formátu a struktuře ukládání dat, proto lze bez problémů uchovávat data i ve stále se měnícím systému s obrovským množstvím dat.
- **Údržba** - Po vytvoření dokumentu je vyžadována minimální údržba.
- **Nezávislost dokumentů** - Dokumenty na sobě jsou nezávislé kvůli absenci cizích klíčů.
- **Otevřené formáty** - K popisu dokumentů lze použít například formát XML, JSON a mnoho dalších.
- **Věstavené verzování** - Díky verzování lze snižovat konflikty.

## Nevýhody

Nevýhody dokumentové databáze jsou [2]:

- **Kontrola konzistence** - Je zde omezená kontrola konzistence, takže se v databázi můžou vyskytovat duplikáty.
- **Neexistence atomicity** - V případě úpravy, která ovlivňuje dvě kolekce, je potřeba spustit dva samostatné dotazy (jeden pro každou kolekci).

## 3.4 Existující řešení

Pro vybrané typy databáze existují mnoho databázových řešení, které lze zmínit. V této kapitole se budeme zabývat především těmi nejznámějšími pro daný typ databáze, a které jsou zdarma ke stažení a používání. Pro každý typ databáze bylo vybráno vždy jedno z nejznámějších řešení.

### 3.4.1 MySQL

MySQL je multiplatformní databáze uplatňující relační databázový model. Komunikace s databází (získávání dat, vytváření objektů, ...) probíhá pomocí jazyka SQL, který je rozšířen o nové funkce. Nejnovější verze MySQL je open-source, což znamená, že kdokoli může používat a libovolně upravovat MySQL systém, aniž by musel cokoli platit. V případě změny zdrojových kódů je potřeba nastudovat podmínky užívání definované licencí GPL [17].

Od samých počátků bylo MySQL optimalizováno především na rychlost i za cenu některých zjednodušení (například způsob zálohování dat). Díky tomuto lze provozovat jednoduché servery na počítači společně s jinými aplikacemi, případně jiné databáze. Server lze nakonfigurovat tím způsobem, že může využívat veškerou paměť, procesorový čas i vstupně výstupní kapacity.

MySQL server může být využit dvěma způsoby:

- **Klient / server** - Vícevláknový SQL server, který podporuje různé back-endy, několik různých klientských programů a knihoven a mnoho dalšího.
- **Věstavěná knihovna** - Vícevláknová věstavěná knihovna, kterou lze propojit do své aplikace a získat tím menší, rychlejší a snadněji spravovatelný samostatný produkt.

### 3.4.2 PostgreSQL

PostgreSQL je open-source objektově-relační databázový systém, který vznikl spojením relačního a objektově-orientovaného databázového systému. Post-

greSQL je velice silný nástroj, který používá rozšíření jazyka SQL společně s mnoha funkcemi, které bezpečně skladují a škálují většinu nejsložitějších datových úloh.

PostgreSQL přichází s mnoha funkcemi, které jsou zaměřené na pomoc vývojářům při vytváření aplikací, správu a bezpečnost dat a odolnost proti chybám v systému. Jako další výhody lze zmínit velkou rozšiřitelnost (lze tvořit vlastní datové typy a funkce), psaní kódu v jiných programovacích jazycích, podpora ACID a možnost provozovat server na všech hlavních operačních systémech [20].

PostgreSQL obsahuje mnoho funkcí, které může uživatel využít. Zde je výpis pouze část z nich:

- **Datové typy** - Primitivní (číslo, text, ...), strukturované (datum, pole, ...), dokumenty (JSON, XML, ...), geometrie (bod, kruh, ...) a vlastní datové typy.
- **Celistvost dat** - Unikátní hodnoty, primární a cizí klíče, zámky.
- **Výkonnost** - Indexování pomocí stromů, výrazů. Základní a vnořené transakce.
- **Zabezpečení** - Vícefaktorová autentikace s certifikáty, sloupcové a řádkové zabezpečení.
- **Textové vyhledávání** - Full-textové vyhledávání.

### 3.4.3 LevelDB

LevelDB je open-source databáze typu Klíč-hodnota, která se používá hlavně u malých přenosných aplikací a nepotřebují žádné API (rozhraní). Databáze byla vytvořena dvěma programátory Googlu, kteří byli inspirováni již existující databází Bigtable (databáze typu klíč-hodnota, která je součástí platformy Google Cloud), ale chtěli vytvořit jednoduchou, lehce přenosnou databázi, kterou lze distribuovat zároveň s aplikací, která ji využívá.

Algoritmus pro ukládání databáze funguje tak, že dočasně ukládá data v *MemTable* (mezipaměť pro zpětný zápis řádků, ve které lze hledat pomocí klíče), ze které se data postupně přesouvají do *SSTable* (Sorted String Table), což je tabulka seřazených řetězců, které nelze měnit. Neměnná data jsou ukládána na disk, který může být sdílen s více clustery [29].

Výhody LevelDB jsou:

- **Jednoduché operace** - LevelDB má tři základní jednoduché operace *Get* (vrací hodnotu podle klíče), *Put* (vkládá dvojici klíč-hodnota) a *Delete* (mazání dvojice klíč-hodnota).
- **Bytové pole** - Klíče a hodnoty lze ukládat i do bytového pole, což může být užitečné v případě, kdy nechceme ukládat hodnoty jako řetězce.
- **Atomické operace** - LevelDB podporuje atomické operace, to znamená, že lze použít více operací najednou v jednom nepřerušeném volání.

### 3.4.4 MongoDB

MongoDB je dokumentově-orientovaná databáze, která se používá zejména tam, kde je potřeba uchovávat velké množství dat. Firma MongoDB, Inc poskytuje oficiální ovladače ke všem populárním programovacím jazykům jako je například C#, Java, C++ a mnoho dalších. Existují i neoficiální ovladače vytvořené komunitou, které pokrývají ještě více programovacích jazyků.

MongoDB je vlastně ve skutečnosti server, který umožňuje vytvářet a udržovat několik databází najednou. Každá databáze může mít své vlastní kolekce, které sdružují dokumenty. MongoDB podporuje vnořená data, což umožňuje vytvářet složité vztahy mezi dokumenty a ukládat je do stejného dokumentu, což činí práci a načítání dat extrémně efektivní ve srovnání s SQL.[14]

Výhody MongoDB jsou [5]:

- Dokumenty lze mapovat na objekty v kódu aplikace, takže se s nimi dá jednodušeji pracovat.
- Indexování, shlukování v reálném čase a ad-hoc dotazy poskytují velice výkonné způsoby přístupu k datům a jejich analýzy.
- MongoDB nabízí vysokou dostupnost, horizontální škálování a geografickou distribuci a to díky tomu, že je ve svém jádře distribuovanou databází.

### 3.4.5 Neo4j

Neo4j je open-source nativní grafová databáze, která efektivně implementuje vlastnosti grafového modelu až na úroveň úložiště, které je velmi výkonné.

Pomocí Neo4j lze naimplementovat každý graf, který dokážeme nakreslit na tabuli, za pomoci ukazatelů. Stejně jako pro MongoDB, i zde existují ovladače pro populární programovací jazyky, jako například Java, .NET a mnoho dalších.

Neo4j je velice populární právě z důvodu konstantních časových přechodů ve velkých grafech jak pro prohledávání do šířky, tak i do hloubky, díky efektivní reprezentaci a škálování uzlů a vztahů mezi nimi. Databáze navíc umožňuje vytvářet flexibilní schéma vlastností grafu, které se může v průběhu času přizpůsobovat, díky čemuž lze přidávat nové vztahy pro vytváření zkratk mezi uzly pro zrychlení práce s daty. Neo4j také poskytuje úplné databázové charakteristiky, které zahrnují i ACID vlastnosti, podpory clusterů a převzetí služeb při selhání za běhu [35].

## 3.5 Jazyky

Databázové jazyky, jinak známé jako dotazovací jazyky, jsou klasifikací programovacích jazyků, které se používají k definování a přístupu k databázím. Pomocí těchto jazyků dokáže uživatel získávat nebo spravovat data v databázích. V dnešní době se jazyky (například SQL) mohou skládat ze čtyř podjazyků, kdy každý slouží k jinému účelu v rámci vykonávání příkazů [11, 36]:

- **Data definition language (DDL)** - DDL umí vytvářet jednotlivé komponenty databázového schématu (tabulky, soubory, indexy, ...), které tvoří strukturu reprezentující organizaci dat v databázi. Dostupné příkazy pro jazyk DDL:
  - **CREATE** - Vytvoření nového objektu (tabulka, index, ...).
  - **ALTER** - Změna struktury objektu.
  - **DROP** - Smazání objektu.
  - **RENAME** - Změna názvu objektu.
  - **TRUNCATE** - Smazání podobjektů v objektu (například záznamy v tabulce).
- **Data manipulation language (DML)** - DML slouží pro manipulaci s daty, které se nachází v již existující databázi. Dostupné příkazy pro jazyk DML:
  - **SELECT** - Získání záznamů (dat) z tabulky.
  - **INSERT** - Vložení nového záznamu (dat) do tabulky.

- **UPDATE** - Úprava existujícího záznamu v tabulce.
- **DELETE** - Smazání záznamu z tabulky.
- **Data control language (DCL)** - Pomocí DCL lze kontrolovat přístupy a práva k datům, které jsou uloženy v databázi. Uživatel může nastavit práva k jednotlivým DML příkazům nad tabulkami / procedurami (například uživatel bude mít přístup pouze k příkazu SELECT nad tabulkou "TABULKA"). Dostupné příkazy pro jazyk DCL:
  - **GRANT** - Přidání práv uživateli nad danou tabulkou / procedurou.
  - **REVOKE** - Odebrání práv uživateli nad danou tabulkou / procedurou.
- **Transaction control language (TCL)** - TCL spravuje transakce v databázi. Transakce obsahuje jeden či více DML příkazů nad tabulkami, které se vykonávají po sobě. Všechny příkazy musí být úspěšně provedeny, aby bylo možné transakci označit za úspěšnou. Ukázka jedné transakce viz obrázek č. 3.10. Dostupné příkazy pro jazyk TCL:
  - **COMMIT** - Potvrzení transakce, změny provedené v transakci jsou permanentní a nejdou vzít zpět.
  - **ROLLBACK** - Vezme zpět veškerou práci v aktuální transakci. Lze se vrátit na začátek transakce nebo k SAVEPOINTu.
  - **SAVEPOINT** - Nastavení bodu v transakci, ke kterému se lze v budoucnu vrátit pomocí ROLLBACK.

```
SQL> SAVEPOINT SP1;
Savepoint created.
SQL> DELETE FROM CUSTOMERS WHERE ID=1;
1 row deleted.
SQL> SAVEPOINT SP2;
Savepoint created.
SQL> DELETE FROM CUSTOMERS WHERE ID=2;
1 row deleted.
SQL> SAVEPOINT SP3;
Savepoint created.
SQL> DELETE FROM CUSTOMERS WHERE ID=3;
1 row deleted.
SQL> ROLLBACK TO SP2;
Rollback complete.
```

Obrázek 3.10: Ukázka jedné transakce (bez commitu)

Níže v kapitolách jsou popsány příklady dnešních jazyků.

### 3.5.1 Structured Query Language

Structured Query Language (SQL) je jazyk pro komunikaci s databázema, v dnešní době standard pro relační databázové systémy. Pomocí SQL příkazů lze například vytvářet nové objekty v databázi, upravovat existující data v tabulkách nebo vytvářet různá integritní omezení a triggery. Většina existujících databázových systémů používá upravený SQL jazyk, který navíc obsahuje dodatečná rozšíření pro splnění požadavků v jejich systémech.

#### Syntax

Syntaxe SQL se skládá z unikátního seznamu pravidel a směrnic. Při psaní příkazů zde nehraje roli citlivost písma (příkazy `select` a `SELECT` jsou záměnné). Dotazy lze psát na jednu nebo více řádek, které musí / mohou být zakončené středníkem (záleží na pravidlech používaného systému). Na obrázku č. 3.11 lze vidět příklad dotazu, který získá jména a příjmení uživatelů z tabulky `'user'` s datem narození po roce 2000.

```
SELECT firstname, lastname FROM user WHERE birthdate_year > 2000;
```

Obrázek 3.11: Příklad SQL dotazu.

Dotazy lze zanořovat do sebe, kdy výsledek jednoho dotazu jde použít jako podmínka pro druhý dotaz, viz obrázek č. 3.12.

```
SELECT * FROM user WHERE user.id = (SELECT id_user FROM meal_orders WHERE id = 1);
```

Obrázek 3.12: Příklad zanořeného SQL dotazu.

### 3.5.2 MongoDB Query Language

MongoDB Query Language (MQL) je jazyk pro získávání dat z MongoDB dokumentových databází. Dotazy zde poskytují jednoduchost v procesu načítání dat z databáze, stejně jako tomu je u SQL. Při provádění dotazů lze také použít kritéria nebo podmínky, kterými lze načíst konkrétní data z databáze. Jazyk také podporuje CRUD operace. Výsledky můžeme třídit, seskupovat, filtrovat a spočítat jejich četnost za pomoci agregační pipeline (zřetěženého zpracování). MQL podporuje transakce více dokumentů [15].

## Syntax

Syntaxe MQL je intuitivní a jednoduchá na používání i pro velice složité dotazy, protože ta samá syntaxe se používá i pro uložené dokumenty v databázi. Příklad syntaxe pro vytváření, čtení, úpravu a mazání dokumentů (CRUD) lze vidět na obrázku č. 3.13.

### (a) Create

```
db.users.insertOne(  ← collection
  {
    name: "sue",      ← field: value
    age: 26,          ← field: value
    status: "pending" ← field: value } document
  }
)
```

### (b) Read

```
db.users.find(           ← collection
  { age: { $gt: 18 } },   ← query criteria
  { name: 1, address: 1 } ← projection
).limit(5)               ← cursor modifier
```

### (c) Update

```
db.users.updateMany(      ← collection
  { age: { $lt: 18 } },    ← update filter
  { $set: { status: "reject" } } ← update action
)
```

### (d) Delete

```
db.users.deleteMany(      ← collection
  { status: "reject" }     ← delete filter
)
```

Obrázek 3.13: Příklady CRUD operací v MongoDB [16].



### 3.5.3 Cypher Query Language

Cypher je dotazovací jazyk pro grafovou databázi Neo4j a umožňuje získávat data z grafů. Tento jazyk byl inspirován hlavně SQL - uživatel se zaměřuje pouze na to, jaká data chce získat, ne jak je má získat. Cypher je unikátní v tom, že poskytuje vizuální způsob, jak sladit vzory a vztahy [6].

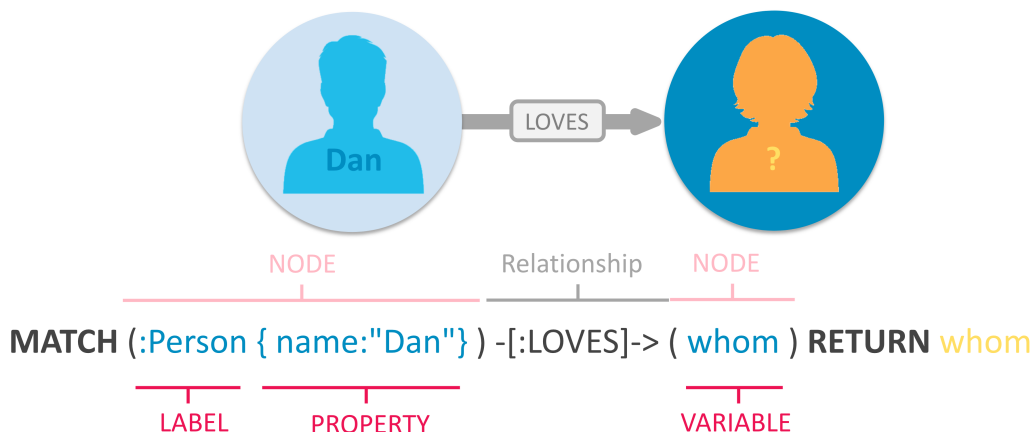
#### Syntax

Cypher využívá ASCII-art typ syntaxe, což je umění, které pracuje s počítačovým textem jako s výtvarným médiem (například obrázky se skládají ze znaků kódu ASCII). Syntaxi lze vidět na obrázku č. 3.14. Pro jednotlivé uzly se používají kruhové závorky, pro vztah se používá šipka s hranatými závorkama obsahující vztah prvního uzlu s druhým uzlem.

```
(nodes) - [ :ARE_CONNECTED_TO ] -> (otherNodes)
```

Obrázek 3.14: Syntaxe jazyka Cypher.

Na obrázku č. 3.15 lze vidět jednoduchý dotaz, který hledá výsledný uzel pro vstupní uzel, kterým je člověk se jménem 'Dan', a vztahu 'LOVES' mezi uzly.



Obrázek 3.15: Dotaz v jazyce Cypher [6].

## 4 Návrh úložiště

Hlavní motivací pro vytvoření této práce je vytěžování patentů pro účely zjišťování existence například různých technologických vynálezů či algoritmů. Pomocí těchto informací lze zjistit, zda například má smysl vymýšlet nový algoritmus pro určitý problém a neexistuje k němu jiné, lepší řešení, případně vymyslet modifikaci, která zajistí lepší výsledky.

Dále je potřeba definovat, co vlastně znamená pojem efektivní vytěžování. Vytěžování lze označit za efektivní, pokud budou splněny tyto podmínky:

- **Rychlost** - Vyhledávání musí probíhat v rámci jednotek až desítek sekund (případně jednotky minut, záleží na celkovém počtu patentů a na hardwarové konfiguraci serveru).
- **Stabilita** - Server musí být stabilní a nesmí padat při práci s velkým množstvím dat, zvláště při vyhledávání s použitím složitých dotazů (například hledání přes více tabulek).

V následujících kapitolách bude popsán postup výběru zdrojů patentů a patentových dat. Následně budou vybrány typy databáze, které budou vhodné pro uložení vybraných patentových dat a následné zvolení existujících řešení.

### 4.1 Výběr patentů

Při výběru patentů byly stanoveny čtyři podmínky, které museli být splněny:

- **Dostupnost** - Patenty musí být dostupné z online stránek / databází bez poplatků.
- **Datum** - Patentová přihláška nebo publikace patentu musí být podána alespoň v roce 2000, všechny ostatní patenty budou vyfiltrovány.
- **Atributy** - Všechny patenty musí obsahovat povinné atributy (viz kapitola č. 4.1.2).
- Alternativy jako Užitný a Průmyslový vzor (Návrhový patent) nebudou brány v potaz.

### 4.1.1 Zdroje dat

V dnešním světě existuje několik desítek až stovek patentových zdrojů dat, od webových vyhledávačů v databázi až po plný export databáze s patenty. Velké organizace, jako například EPO, WIPO, USPTO, udržují jedny z největších patentových databází (desítky až stovky milionů patentů), ve kterých lze vyhledávat velké množství informací zdarma za použití webových vyhledávačů na dané stránce organizace. Lze zde najít všechny typy patentů (příhlášky, publikace), národní patenty i patenty registrované například u EPO. V případě exportu databází, USPTO poskytuje plný export svých databází veřejnosti pro libovolné používání, zcela zdarma. Využití těchto zdrojů dat by bylo určitě skvělé, ale tyto zdroje byly nedávno použity a rozebrány v jiné diplomové práci, proto je vhodné se spíše zaměřit na národní zdroje dat patentů.

Národní databáze patentů dané země obsahuje všechny národní patenty, některé dokonce i patenty z jiných zemí registrovaných u EPO.

Při průzkumu bylo zkoumáno celkem 51 národních zdrojů dat (patentových úřadů). V tabulkách č. 4.1 a 4.2 lze vidět název země, název patentového úřadu v dané zemi, zkratku patentového úřadu (pokud nějakou má) a jestli patentový úřad poskytoval data nebo ne. U každého patentového úřadu byl procházen její oficiální web a zkoumán na dostupnost patentových dat. Většina úřadů má na svých stránkách vyhledávač pro procházení vlastní databáze patentů, ale jen zlomek z nich poskytoval použitelná data zadarmo. Tyto data byla většinou schována pod neodpovídajícím názvem článku / příspěvku, a některé dokonce poskytovaly odkazy ke stažení dat na svých stránkách pouze v národním jazyce (neexistující článek s daty v anglické verzi webu).

<b>Země</b>	<b>Patentový úřad</b>	<b>Zkratka</b>	<b>Data</b>
Anglie	Intellectual Property Office	IPO	ANO
Arménie	Intellectual Property Office	-	NE
Austrálie	IP Australia	-	ANO
Bělorusko	National Center of Intellectual Property	NCIP	NE
Bulharsko	Patent Office of Republic of Bulgaria	-	NE
Česko	Industrial Property Office of the Czech Republic	-	ANO
Čína	China National Intellectual Property Administration	CNIPA	NE
Dánsko	Danish Patent and Trademark Office	-	NE
Egypt	Egyptian Patent Office	-	NE
Estonsko	The Estonian Patent Office	-	NE
Filipíny	Intellectual Property Office of the Philippines	IPOPHL	NE
Finsko	Finnish Patent and Registration Office	PRH	NE
Francie	National Institute of Industrial Property	INPI	ANO
Hong Kong	Intellectual Property Department	-	NE
Chorvatsko	State Intellectual Property Office of the Republic of Croatia	SIPO	NE
Indie	Office of the Controller General of Patents, Designs and Trade Marks	-	NE
Indonésie	Directorate General of Intellectual Property	DGIP	NE
Irsko	Intellectual Property Office of Ireland	IPOI	NE
Island	Icelandic Intellectual Property Office	ISIPO	NE
Israel	The Israel Patent Office	ILPO	ANO
Itálie	Directorate General for the Protection of Industrial Property	-	ANO*
Japonsko	Japan Patent Office	JPO	ANO
Jižní Korea	Korean Intellectual Property Office	KIPO	ANO
Kanada	Canadian Intellectual Property Office	CIPO	ANO

Tabulka 4.1: Národní patentové úřady a jejich zkratky, část první.

<b>Země</b>	<b>Patentový úřad</b>	<b>Zkratka</b>	<b>Data</b>
Kuba	Cuban Industrial Property Office	OCPI	NE
Litva	State Patent Bureau of the Republic of Lithuania	-	ANO
Lotyšsko	Patent Office of the Republic of Latvia	-	NE
Maďarsko	Hungarian Intellectual Property Office	HIPO	NE
Malajsie	Intellectual Property Corporation of Malaysia	MyIPO	NE
Mexiko	Instituto Mexicano De La Propiedad Industrial	IMPI	ANO
Moldova	State Agency on Intellectual Property	AGEPI	NE
Německo	German Patent and Trade Mark Office	DPMA	ANO
Nizozemsko	Netherlands Patent Office	-	NE
Norsko	Norwegian Industrial Property Office	NIPO	NE
Nový Zéland	Intellectual Property Office of New Zealand	IPONZ	ANO
Peru	National Institute for the Defense of Competition and Protection of Intellectual Property	INDECOPI	ANO
Polsko	Urząd Patentowy Rzeczypospolitej Polskiej	UPRP	ANO
Portugalsko	Portuguese Institute of Industrial Property	-	ANO
Rakousko	Austrian Patent Office	-	NE
Rumunsko	State Office for Inventions and Trademarks	OSIM	NE
Rusko	Federal Service for Intellectual Property	Rospatent	ANO
Řecko	Hellenic Industrial Property Organization	HIPO	NE
Singapur	Intellectual Property Office of Singapore	IPOS	NE
Slovensko	Industrial Property Office of the Slovak Republic	-	NE
Slovinsko	Slovenian Intellectual Property Office	SIPO	NE
Srbsko	Intellectual Property Office of the Republic of Serbia	-	NE
Španělsko	Spanish Patent and Trademark Office	OEPM	ANO
Švédsko	Swedish Intellectual Property Office	PRV	ANO
Švýcarsko	Swiss Federal Institute of Intellectual Property	-	NE
Turecko	Turkish Patent and Trademark Office	Turkpatent	NE
Ukrajina	Ukrainian Intellectual Property Institute	Ukrpatent	NE

Tabulka 4.2: Národní patentové úřady a jejich zkratky, část druhá.

Z celkových 51 patentových zdrojů nám pouze 19 zdrojů poskytuje data. V případě Itálie nám data neposkytuje přímo patentový úřad, ale výzkumný úřad PATIRIS, který poskytuje patentová data z univerzit a veřejných výzkumných ústavů v Itálii.

Bohužel ne všechny patentové úřady poskytují svá data zdarma. Celkem tři úřady - Austrálie, Německo, Nový Zéland chtěli za svá data zaplatit.

Ještě je potřeba zmínit Japonsko, které svá data poskytuje, ale je potřeba vyplnit formulář, ve kterém bylo potřeba naskenovat oficiální dokument potvrzující adresu školy. Z tohoto důvodu jsme bylo Japonsko jako zdroj dat zavrhnuto. Z původních 19 zdrojů dat poskytující data zůstalo nakonec jen 14 zdrojů dat, které poskytují svá data zadarmo.

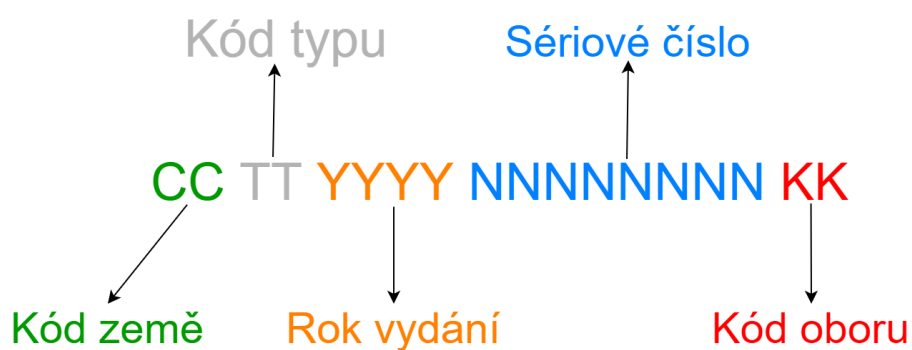
### 4.1.2 Atributy

Při zadávání práce byly definovány podmínky pro výběr platných zdrojů dat, a jednou z nich bylo i povinnost mít důležité atributy ve struktuře dat patentu. Celkem byly definovány čtyři povinné atributy společně s osmi nepovinnými. Povinné a nepovinné atributy jsou podrobně popsány níže.

#### Povinné atributy

Se zadavatelem bylo domluveno, že validní zdroj patentových dat musí poskytovat patenty obsahující tyto atributy:

- **Titulek** - Titulek patentu, který říká o čem daný patent.
- **Rok přihlášky / publikace** - Patent musí obsahovat alespoň rok přihlášky / publikace. Publikace / přihláška musí být minimálně z roku 2000, patenty před rokem 2000 budou zamítnuty.
- **Autor** - U patentu bude nutné vědět jeho autor (jméno autora, případně název instituce).
- **ID patentu** - Patent musí mít nějaké kódové označení / identifikátor, podle kterého ho lze vyhledávat. Identifikátor se bude držet formátu viz obrázek č. 4.1, ale nemusí obsahovat všechny položky, protože každá země může některé položky zanedbávat, případně měnit počet znaků v položce, viz tabulka č. 4.3.



Obrázek 4.1: Základní formát pro identifikátor patentu [19].

<b>Země</b>	<b>CC</b>	<b>TT</b>	<b>YYYY</b>	<b>NNNNNNNN</b>	<b>KK</b>
Austrálie	AU		4 znaky	6 znaků	ANO
Kanada	CA			7 znaků	ANO
Čína	CN	1 znak		8 znaků	ANO
EPO	EP			7 znaků	ANO
Německo	DE	2 znaky	4 znaky	6 znaků	ANO
Francie	FR			7 znaků	ANO
Velká Británie	GB			7 znaků	ANO
Nizozemsko	NL			7 znaků	ANO
Japonsko	JP		4 znaky	6 znaků	ANO
Korea	KR	2 znaky	4 znaky	7 znaků	ANO
Rusko	RU		4 znaky	6 znaků	ANO
USA	US		4 znaky	7 znaků	ANO
WIPO	PCT		4 znaky	6 znaků	ANO

Tabulka 4.3: Aktuálně používané formáty pro patenty z různých zemí [19].

V tabulce č. 4.4 lze vidět patenty, poskytované národními patentovými úřady zdarma, obsahují povinné atributy.

Země	Titulek patentu	Rok přihlášky / publikace	Autor	ID patentu
Kanada	x	x	x	x
Česko	x	x	-	x
Litva	x	x	x	x
Portugalsko	x	x	x	x
Španělsko	x	x	x	x
Švédsko	-	x	-	x
Izrael	x	x	x	x
Itálie	x	x	x	x
Mexiko	x	x	x	x
Polsko	x	x	-	-
Anglie	x	x	x	x
Rusko	x	x	x	x
Peru	x	x	x	x
Francie	x	x	x	x

Tabulka 4.4: Povinné atributy nacházející se v dostupných patentech.

### Nepovinné atributy

Při průzkumu byly zjišťovány i nepovinné atributy, které nemají vliv na výběr zdrojů dat, ale je dobré vědět co který patent z daného patentového zdroje poskytuje za atributy. Nepovinné atributy jsou:

- **Abstrakt** - Stručný výtah patentu, který popisuje o čem daný patent je.
- **Klíčová slova** - Klíčová slova nebo fráze spojené s patentem. Mohou sloužit při vyhledávání patentů se stejným zaměřením.
- **Reference** - Reference na podobné typy patentů nebo na související patenty (například odkaz na základní verzi algoritmu).
- **Žadatel** - Žadatel a autor může být tatáž osoba, ale v některých případech je žadatelem někdo jiný (například autor je zaměstnanec firmy, žadatelem je samotná firma).
- **Adresa autora / instituce** - Adresa autora nebo instituce.
- **Rodina patentů** - Rodina patentů je kolekce patentových žádostí, které se zaměřují na stejný nebo alespoň podobný technický obsah.



- **Obor** - Obor, který daný patent pokrývá.
- **Full-text** - Zdroje dat poskytují veškerá data o patentu (nejenom to co je v přihláškách / publikacích, například různé poznámky, obrázky).

V tabulkách č. 4.5, č. 4.6 lze vidět patenty, poskytované národními patentovými úřady zdarma, obsahující nepovinné atributy.

Země	Abstrakt	Klíčová slova	Reference	Žadatel
Kanada	-	-	-	x
Česko	x	-	x	-
Litva	x	-	-	x
Portugalsko	x	-	-	x
Španělsko	x	-	-	x
Švédsko	x	-	-	-
Izrael	-	-	-	x
Itálie	-	-	-	x
Mexiko	x	-	-	x
Polsko	x	x	-	-
Anglie	-	-	-	-
Rusko	-	-	-	-
Peru	-	-	-	-
Francie	x	-	x	x

Tabulka 4.5: Nepovinné atributy nacházející se v dostupných patentech, část první.

<b>Země</b>	<b>Adresa</b>	<b>Rodina patentů</b>	<b>Obor</b>	<b>Full-text</b>
Kanada	x	x	x	-
Česko	-	x	x	-
Litva	-	x	x	-
Portugalsko	-	-	x	-
Španělsko	x	-	x	x
Švédsko	-	x	x	x
Izrael	x	x	-	-
Itálie	-	-	-	-
Mexiko	-	-	x	-
Polsko	-	-	-	-
Anglie	-	-	x	-
Rusko	-	-	-	-
Peru	-	-	x	-
Francie	-	x	x	-

Tabulka 4.6: Nepovinné atributy nacházející se v dostupných patentech, část druhá.

### 4.1.3 Závěr průzkumu

Průzkum národních zdrojů zahrnoval celkem 51 národních patentujících institucí, ze kterých pouze 10 poskytovalo svá data zdarma a splňovala všechny podmínky. V tabulce č. 4.7 lze vidět souhrn výsledků.

<b>Popis</b>	<b>Počet</b>	<b>Poměr</b>
Nedostupné	33	64,70 %
Nepoužitelné	4	7,85 %
Za peníze	4	7,85 %
Použitelné	10	19,60 %

Tabulka 4.7: Souhrn průzkumu národních patentujících institucí.

Pro všechny použitelné národní zdroje bylo kromě výskytu atributů dále sledováno: formát uložených dat, počet patentů po roce 2000 (včetně roku 2000) obsahující všechny povinné atributy, počet patentů před rokem 2000 a počet duplikátů. Duplikátem se myslí jiná verze daného patentu, protože v průběhu let se mohl měnit obsah patentu a v databázi nechceme ukládat žádné starší verze jednoho patentu (výsledky vyhledávání v databázi nebudou validní). V tabulce č. 4.8 lze vidět všechny validní národní zdroje.

<b>Země</b>	<b>Formát dat</b>	<b>Počet patentů (rok &gt;= 2000)</b>	<b>Počet patentů (rok &lt; 2000)</b>	<b>Počet duplikátů</b>
Anglie	XML	88 032	141	19
Francie	XML	746 899	192 630	1 140 084
Israel	XML	116 380	0	9 956
Itálie	SQL	17 622	3 728	0
Kanada	XML	936 464	130 279	499 682
Litva	XML	869	0	0
Peru	XLSX	1 805	21 352	0
Portugalsko	XML	69	0	0
Rusko	CSV	614 256	139 714	15 166 816
Španělsko	XML	381 713	37 596	39 033
<b>Souhrn</b>		<b>2 904 109</b>	<b>525 440</b>	<b>16 855 590</b>

Tabulka 4.8: Seznam všech validních národních zdrojů.

V tabulce č. 4.9 lze vidět výsledný počet patentů (po filtraci všech patentů neobsahující povinné atributy - ID, titulek, autor, datum).

<b>Země</b>	<b>Počet patentů před filtrací</b>	<b>Filtrace</b>	<b>Počet patentů po filtraci</b>
Anglie	88 032	1	88 031
Francie	746 899	474 850	272 049
Israel	116 380	291	116 089
Itálie	17 622	10 532	7 090
Kanada	936 464	119 696	816 768
Litva	869	0	869
Peru	1 805	0	1 805
Portugalsko	69	0	69
Rusko	614 256	223	614 033
Španělsko	381 713	308 029	73 684
<b>Souhrn</b>	<b>2 904 109</b>	<b>913 622</b>	<b>1 990 487</b>

Tabulka 4.9: Výsledný počet patentů po filtraci.

## 4.2 Výběr databáze

V kapitole č. 3 bylo podrobně popsáno co databáze je, jaké typy databází dnes existují (krátký výčet) i nejznámější existující řešení pro popsané typy databází. V této kapitole budou podrobně popsány rozdíly mezi jednotlivými řešeními a následně se vybere nejlepší typ databáze pro ukládání velkého objemu patentových dat. Následně, podle vybraného typu databáze, se vybere nejvhodnější existující řešení.

### 4.2.1 Výběr typu databáze

Abychom zajistili co nejefektivnější vytěžování, tak je potřeba vybrat co nejvhodnější typ databáze vzhledem k povaze úlohy. V budoucnu se očekává, že počet skladovaných patentů bude v řádech jednotek až desítek milionů (nelze vyvrátit i stovky milionů v případě, že se budou ukládat i patenty z jiných než národních zdrojů).

#### Relační databáze

Relační databáze není vhodným kandidátem pro ukládání nestrukturovaných patentových dat. V databázi sice existuje datový typ **BLOB**, který umožňuje ukládat binární soubory (v našem případě soubor s patentovými daty), ale nelze to pokládat za nejlepší řešení, když existují například dokumentové databáze. Lze zmínit i datový typ **TEXT** / **LONGTEXT**, který umožňuje ukládat velké množství textu a lze ho procházet pomocí full-text vyhledávání, ale výkonnostně a rychlostně se stejně nevyrovná NoSQL databázím.

Využití relační databáze by mělo smysl pouze v případě vytváření statistik (například počet v patentů v Kanadě za rok 2020). Tento přístup by ale vyžadoval extrahovat specifická data (například jen povinné atributy) ze souboru pomocí parseru a následné uložení hodnot do tabulek.

#### Objektově-orientovaná databáze

Objektově-orientovaná databáze, stejně jako relační databáze, není vhodným kandidátem pro ukládání nestrukturovaných dat. Lze argumentovat vytvořením objektů odpovídající struktuře dokumentu, ale při vložení dokumentu s jinou strukturou nastává problém s uložením atributů, které se nenachází v objektu. V některých případech může vysoká složitost systému zpomalovat vyhledávání. Velká výhoda objektově-orientované databáze spočívá v jednoduchém mapování objektů při práci s objektově-orientovaným programováním,

které ale v našem případě nemá využití. V případě vytěžování statistik je objektově-orientovaná databáze horší volbou než relační databáze.

### **Databáze klíč-hodnota**

Databáze klíč-hodnota je jednoduchá a velice rychlá databáze, která umožňuje ukládat i nestrukturovaná data a nepotřebuje k tomu velké množství paměti. Její nevýhoda je ale v ukládání složitých struktur, které soubory s patenty mají. Lze uložit celou strukturu patentu jako hodnotu, ale následné vyhledávání hodnot pomocí názvů parametrů je nemožné. Použití databáze klíč-hodnota v našem případě není moc vhodné.

### **Grafová databáze**

Grafová databáze není vhodným kandidátem pro ukládání patentů, protože se zaměřuje hlavně na vztahy mezi jednotlivými daty, což u patentů nelze a ani není potřeba sledovat.

### **Databáze dokumentů**

Databáze dokumentů, jak už název napovídá, je databáze pro efektivní ukládání dokumentů a jejich vytěžování. Umožňuje ukládat velké množství nestrukturovaných dat, její udržba je snadná a akceptuje dokumenty v několika datových formátech. Její velká nevýhoda je v kontrole konzistence, takže se v databázi mohou vyskytovat duplikáty. I přes tuto nevýhodu je dokumentová databáze vhodným kandidátem pro ukládání patentových dat.

### **Závěr**

Dokumentová databáze bude použita jako primární databáze, protože umožňuje ukládat nestrukturované dokumenty velice efektivně. Zároveň podporuje vkládání dokumentů ve více formátech, což v případě mnoha národních zdrojů, kdy každý zdroj ukládá dat v jiném formátu, je velice vhodná vlastnost. Její nevýhoda, kontrola konzistence, může být odstraněna pomocí jednoduché aplikace, která bude kontrolovat výskyt patentu v databázi podle jeho identifikátoru (ID). Dokumentová databáze podporuje i full-text vyhledávání pro efektivnější a rychlejší vyhledávání.

Relační databáze bude použita jako sekundární databáze pro vytěžování, zaměřená hlavně na tvorbu statistik. Relační databáze je vhodnou volbou pro získávání statistik, protože vyhledávání je velice rychlé a snadné. SQL dotazy pro statistiky se můžou uložit do pohledů, které zjednoduší uživateli práci se zadáváním dotazů. Lze zmínit i nevýhody jako třeba problém s

údržbou nebo potřeba velkého množství paměti. Tyto nevýhody ale nehrají velkou roli v případě jednotek až desítek milionů záznamů.

#### 4.2.2 Výběr z existujících řešení

V dnešní době existuje mnoho dokumentových i relačních databází, placených i zdarma poskytovaných. Placené verze oproti verzím zdarma mají výhodu v lepší podpoře ze strany vývojářů, obsahují více užitečných funkcí a mají lepší zabezpečení. Pro naše účely bohatě postačí verze zdarma.

Jako existující řešení databáze dokumentů byla vybrána komunitní verze databáze MongoDB. Komunitní verze je zdarma a server lze provozovat jak lokálně, tak i na cloudu, kde MongoDB poskytuje zdarma uložení o velikosti 512 MB. Pro lepší a spolehlivější vyhledávání v datech bude MongoDB spojena s vyhledávačem Elasticsearch. Elasticsearch je full-textový open-source vyhledávač, který nabízí vysokou dostupnost, rychlost a škálovatelnost. MongoDB sice obsahuje vlastní full-textový vyhledávač, který ale není tak výkonný jako Elasticsearch.

Jako existující řešení relační databáze byla vybrána komunitní verze databáze MySQL. MySQL je skvělá databáze, která se používá hlavně pro čtení dat. Zároveň je to jedna z nejpoužívanějších relačních databází, což znamená, že je pro ní k dispozici více nástrojů třetích stran.

# 5 Implementace úložiště

todo

## 5.1 Adresářová struktura

todo

### 5.1.1 Patenty

### 5.1.2 Docker

## 5.2 Implementace databáze

todo

### 5.2.1 MySQL

vypsat počet hodnot v tabulkách

### 5.2.2 MongoDB

## 5.3 Výsledné úložiště

todo

### 5.3.1 Technologické požadavky

### 5.3.2 Docker

Images, scripty

Import dat (skripty, data soubory, ...)

### 5.3.3 Inicializace MySQL

### 5.3.4 Inicializace MongoDB

## 6 Rozšiřitelnost úložiště

Zadání diplomové práce sice splněno bylo, ale v blízké budoucnosti mohou být požadavky na modul změněny. Jako příklad lze uvést podporu přidávání nových patentů do databází, zjištění autorů pro české patenty, automatické stahování dat z již ověřených patentových zdrojů. V této kapitole jsou popsány 3 možné návrhy na rozšíření modulu ohledně importu dat do již existujících databází.

### 6.1 Přidávání nových patentů

Cílem tohoto rozšíření by bylo automatické přidávání patentů z datových souborů jak do MySQL databáze, tak i do Mongo.

Rozšíření by se dalo realizovat jako aplikace ve vyšším programovacím jazyku (např. Java, C), kdy vstupem do aplikace by byl soubor v datovém formátu JSON/XML/CSV a jiné. Vstupní soubor by se následně:

- převedl na JSON řetězec (v případě že soubor není ve formátu JSON) a vložil do Mongo databáze.
- rozparsoval a extrahovali by se všechny atributy, které se ukládají v MySQL databázi (viz kapitola č. 5.2.1).

Jelikož je dost časté, že každý národní zdroj dat používá odlišnou strukturu patentu, tak bude potřeba aplikaci neustále upravovat (ať už v rámci přidávání nových zdrojů, nebo v případě změny struktury patentu u již podporovaných zdrojů).

Jako další velký problém lze zmínit extrakci atributů patentu ze souborů. Tím, že různé patentové soubory mají odlišnou strukturu, to znamená hloubku zanoření specifických elementů, jiné názvy elementů, tak bude obtížné naimplementovat řešení extrakce pro všechny soubory. Tento problém by se dal řešit tak, že se vytvoří soubory se slovníkama, které by obsahovaly názvy elementů pro daný atribut. Slovníky by se následně použily při extrakci.

### 6.2 Zjišťování autorů pro české patenty

Český národní patentový úřad poskytuje data o českých patentech, které ale neobsahují autora ani instituci. Pro zjištění autora nebo instituce, která



patent registrovala, je nutné použít oficiální vyhledávač. Cílem tohoto rozšíření by bylo vytvořit aplikaci ve vyšším programovacím jazyku, která se pro všechny české patenty bude snažit najít jejich autory za pomoci využití prohlédávačů webů (web crawler). Postupů řešení může být mnoho:

- Zjišťování autorů by se provedlo pro všechny existující české patenty v databázi. Z MySQL databáze se zjistí všechny identifikátory pro české patenty, které se následně použijí jako vstup pro web crawler.
- Zjišťování autorů by se provedlo pro patent/y uložené v souboru, kdy aplikace by pro všechny patenty v souboru zjistila autory a následně je dopsala do příslušného elementu patentu v daném souboru.
- Stejný postup jako předchozí s tím rozdílem, že po zjištění autora se patent rovnou přidá do MySQL i Mongo databáze.

## 6.3 Automatické stahování dat z ověřených zdrojů

Cílem tohoto rozšíření by bylo automatické stahování dat (případně i jejich parsování) z ověřených zdrojů. Ověřené zdroje by byly uloženy například v XML souboru, kdy každý zdroj by měl tyto položky:

- **Název země**
- **URL** - URL zdroje dat, na které lze stáhnout data.
- **XPath** - XPath výraz, pomocí kterého lze ze stránky vyfiltrovat a získat odkazy ke stažení dat  
(např. `/html/body//a[contains(@href,'example')]/@href`)
- **Poslední verze** - Název / číslo poslední stažené verze.

XML soubor by byl následně zpracován pomocí aplikace (např. Java, C#), která by následně pro každý zdroj dat provedla následující kroky:

1. Získání seznamu odkazů na zdroje dat.
2. Stažení všech zdrojů dat, jejichž verze je větší než aktuálně uložená verze v XML.
3. V tomto bodě se dá naimplementovat cokoliv - např. lze uložená data extrahovat ze ZIP souborů, importovat patenty do databází (viz kapitola č. 6.1), pouze notifikace o stažení několika nových souborů z daty a mnoho dalšího.

#### 4. Aktualizace verze v XML souboru.

Automatizace stahování dat by spočívala ve spouštění aplikace pro stahování dat v pravidelných intervalech (např. každé druhé úterý v 17:00). Jako příklad lze uvést použití pipeline na Jenkins serveru, který bude spouštět z lokálního uložště spustitelnou aplikaci v daný čas (pomocí CRON). Po vykonání celého procesu může Jenkins poslat email o stavu posledního spuštění (zda se spuštění povedlo, kolik souborů byl schopen stáhnout pro jaké země, ...). Samozřejmě bohatě postačí i použití plánovače v operačním systému.

## 7 Ověření efektivního vytěžování

K ověření efektivního vytěžování bylo připraveno několik scénářů jak pro SQL, tak i pro Mongo s využitím vyhledávače Elasticsearch. Z výsledků scénářů můžeme poté usoudit, jak moc efektivní vytěžování je vzhledem k technickým parametrům stroje a použitým technologiím.

### Testovací stroj - technické parametry

Testovací stroj, který byl použit pro kontrolu efektivity vytěžování, má tyto technické parametry:

- **Procesor** - Intel(R) Core(TM) i5-8250U CPU @ 1.60GHz (8 CPUs).
- **RAM** - 8192 MB
- **Disk** - SSD 256 GB
- **Grafická karta 1** - Intel(R) UHD Graphics 620
- **Grafická karta 2** - NVIDIA GeForce 940MX
- **Operační systém** - Windows 10 Home 64bit (10.0, build 19044).

### 7.1 Mongo + ElasticSearch

todo

### 7.2 MySQL

Pro MySQL bylo připraveno devět scénářů, které testují všechny vytvořené tabulky v databázi. Každý scénář obsahuje textový popis, SQL příkaz, rychlost vykonání příkazu a ukázkou výsledků. Každý scénář odpovídá jednomu pohledu v databázi.

#### 7.2.1 Scénář č.1

**Textový popis:** Deset nejčastěji patentujících institucí v Izraeli v roce 2015.  
**SQL:**

```
select count(*), count(*) * 100.0 / ((select count(*) from inventors left
  outer join patents on inventors.id_patent = patents.id where YEAR(
    patents.patent_date) = 2015 and patents.patent_id like '%IL%') * 1.0)
  as percentage, inventors.inventor from inventors left outer join
  patents on inventors.id_patent = patents.id where YEAR(patents.
    patent_date) = 2015 and patents.patent_id like '%IL%' group by
    inventors.inventor order by count(*) desc, percentage desc LIMIT 10;
```

Rychlost vykonání dotazu: +- 1,4 sekundy

Výsledek dotazu:

count(*)	percentage	inventor
39	1.00309	DOW AGROSCIENCES LLC
29	0.74588	NOVARTIS AG
29	0.74588	F. HOFFMANN-LA ROCHE AG
29	0.74588	GENENTECH, INC.
24	0.61728	RAYTHEON COMPANY
22	0.56584	QUALCOMM INCORPORATED
20	0.51440	BIOSENSE WEBSTER (ISRAEL) LTD.
17	0.43724	MICROSOFT CORPORATION
17	0.43724	SANOFI
16	0.41152	YERKES, CARLA N.

Obrázek 7.1: Ukázka výsledku dotazu pro scénář č.1

### 7.2.2 Scénář č.2

**Textový popis:** Tři nejméně patentované obory v Kanadě od roku 2010.

**SQL:**

```
select count(*), count(*) * 100.0 / ((select count(*) from classification
  left outer join patents on patents.id = classification.id_patent
  where YEAR(patents.patent_date) >= 2010 and patents.patent_id like '%
  CA%') * 1.0) as percentage, classification.section from
  classification left outer join patents on patents.id = classification
  .id_patent where YEAR(patents.patent_date) >= 2010 and patents.
  patent_id like '%CA%' group by classification.section order by count
  (*) asc, percentage asc LIMIT 5;
```

Rychlost vykonání dotazu: +- 6,2 sekund

Výsledek dotazu:

count(*)	percentage	section
7044	0.93164	D
41792	5.52742	E
60066	7.94434	F
71423	9.44642	H
101433	13.41556	G

Obrázek 7.2: Ukázka výsledku dotazu pro scénář č.2

### 7.2.3 Scénář č.3

**Textový popis:** Dvacet nejčastějších klasifikací patentů za rok 2008 ve Španělsku.

**SQL:**

```
select count(*), count(*) * 100.0 / ((select count(*) from classification
  left outer join patents on patents.id = classification.id_patent
 where YEAR(patents.patent_date) = 2008 and patents.patent_id LIKE '%
ES%') * 1.0) as percentage, classification.section, classification.
class, classification.subclass from classification left outer join
patents on patents.id = classification.id_patent where YEAR(patents.
patent_date) = 2008 and patents.patent_id LIKE '%ES%' group by
classification.section, classification.class, classification.subclass
order by count(*) desc, percentage desc LIMIT 20;
```

**Rychlost vykonání dotazu:** +- 1,3 sekundy

**Výsledek dotazu:**

count(*)	percentage	section	class	subclass
71	2.99831	B	65	D
56	2.36486	E	04	G
47	1.98480	A	61	K
44	1.85811	E	04	B
40	1.68919	G	01	N
33	1.39358	F	24	J
32	1.35135	E	06	B
30	1.26689	A	23	L
27	1.14020	C	02	F
27	1.14020	B	60	R
27	1.14020	E	04	H
26	1.09797	D	06	F
26	1.09797	A	47	L
25	1.05574	F	03	D
25	1.05574	A	47	C
24	1.01351	A	01	K
24	1.01351	B	01	D
22	0.92905	B	65	B
22	0.92905	E	04	C
21	0.88682	G	02	B

Obrázek 7.3: Ukázka výsledku dotazu pro scénář č.3

#### 7.2.4 Scénář č.4

**Textový popis:** Deset autorů s největším počtem patentů ze všech zemí.

**SQL:**

```
select count(*), count(*) * 100.0 / ((select count(*) from inventors) *
1.0) as percentage, inventors.inventor from inventors group by
inventors.inventor order by count(*) desc, percentage desc LIMIT 10;
```

**Rychlost vykonání dotazu:** +- 8 sekund

**Výsledek dotazu:**

count(*)	percentage	inventor
26987	0.65261	Квасенков Олег Иванович (RU)
5123	0.12389	Щепочкина Юлия Алексеевна (RU)
1932	0.04672	Кочетов Олег Савельевич (RU)
1675	0.04051	QUALCOMM INCORPORATED
1430	0.03458	Квасенков Олег Иванович
1158	0.02800	ASTRAZENECA AB
1113	0.02692	NOVARTIS AG
1035	0.02503	F. HOFFMANN-LA ROCHE AG
948	0.02293	BASF SE
904	0.02186	Consiglio Nazionale delle Ricerche - CNR

Obrázek 7.4: Ukázka výsledku dotazu pro scénář č.4

### 7.2.5 Scénář č.5

**Textový popis:** Pět nejméně používaných jazyků pro patenty za rok 2003.  
**SQL:**

```
select count(*), count(*) * 100.0 / ((select count(*) from patents where
patents.language not like '%-%') * 1.0) as percentage, patents.
language from patents where patents.language not like '%-%' group by
patents.language order by count(*) asc, percentage asc LIMIT 5;
```

**Rychlost vykonání dotazu:** +- 5,3 sekund

**Výsledek dotazu:**

count(*)	percentage	language
69	0.00348	PT
869	0.04381	LT
75489	3.80605	ES
299732	15.11205	FR
614033	30.95865	RU

Obrázek 7.5: Ukázka výsledku dotazu pro scénář č.5

### 7.2.6 Scénář č.6

**Textový popis:** Deset Institucí / autorů s patenty pokrývající největší množství oborů ve Španělsku.

**SQL:**

```
select count(distinct classification.section), count(*) * 100.0 / ((
select count(*) from inventors left outer join classification on
classification.id_patent = inventors.id_patent left outer join
patents on patents.id = inventors.id_patent where section is not null
```

```

and patents.country like '%ES%') * 1.0) as percentage, inventors.
inventor from inventors left outer join classification on
classification.id_patent = inventors.id_patent left outer join
patents on patents.id = inventors.id_patent where section is not null
and patents.country like '%ES%' group by inventors.inventor order by
count(distinct classification.section) desc, percentage desc LIMIT
10;

```

Rychlost vykonání dotazu: +- 3 sekundy

Výsledek dotazu:

count(distinct classification.section)	percentage	inventor
8	0.08400	TRENCH ROCA LLUIS
8	0.07000	ALET VIDAL JOSEP
7	0.11201	CORMA CANOS AVELINO
7	0.09100	PORRAS VILA FCO. JAVIER
6	0.09450	PORRAS VILA F. JAVIER
6	0.03500	GUTIERREZ MIGUELEZ ANGEL
6	0.03150	LLOVERAS MACIA JOAQUIM
5	0.13651	LLORENTE GONZALEZ JOSE IGNACIO
5	0.13301	MONTERDE AZNAR FERNANDO
5	0.12601	OH JANG-KEUN

Obrázek 7.6: Ukázka výsledku dotazu pro scénář č.6

## 7.2.7 Scénář č.7

Textový popis: Pět zemí s nejvíce patenty od roku 2018.

SQL:

```

select count(*), count(*) * 100.0 / ((select count(*) from patents where
YEAR(patents.patent_date) >= 2018) * 1.0) as percentage, patents.
country from patents where YEAR(patents.patent_date) >= 2018 group by
patents.country order by count(*) desc, percentage desc LIMIT 5;

```

Rychlost vykonání dotazu: +- 2,4 sekundy

Výsledek dotazu:

count(*)	percentage	country
122299	39.44658	RU
112202	36.18987	CA
56954	18.37007	FR
7295	2.35294	IL
5769	1.86075	UK

Obrázek 7.7: Ukázka výsledku dotazu pro scénář č.7



### 7.2.8 Scénář č.8

**Textový popis:** Nejvíce používaný typ patentu ve Francii.

**SQL:**

```
select count(*), count(*) * 100.0 / ((select count(*) from patents where
    patents.patent_id like '%FR%' and patents.kind not like '%-%') * 1.0)
    as percentage, patents.kind from patents where patents.patent_id
    like '%FR%' and patents.kind not like '%-%' group by patents.kind
    order by count(*) desc, percentage desc;
```

**Rychlost vykonání dotazu:** +- 5,5 sekund

**Výsledek dotazu:**

count(*)	percentage	kind
264824	97.65293	A

Obrázek 7.8: Ukázka výsledku dotazu pro scénář č.8

### 7.2.9 Scénář č.9

**Textový popis:** Patnáct nejčastěji patentujících institucí / autorů v Anglii v textilním oboru za rok 2013.

**SQL:**

```
select count(*), count(*) * 100.0 / ((select count(*) from inventors left
    outer join patents on patents.id = inventors.id_patent left outer
    join classification on classification.id_patent = patents.id where
    classification.section like '%D%' and patents.patent_id like '%GB%'
    and YEAR(patents.patent_date) = 2013) * 1.0) as percentage, inventors
    .inventor from inventors left outer join patents on patents.id =
    inventors.id_patent left outer join classification on classification.
    id_patent = patents.id where classification.section like '%D%' and
    patents.patent_id like '%GB%' and YEAR(patents.patent_date) = 2013
    group by inventors.inventor order by count(*) desc, percentage desc
    LIMIT 15;
```

**Rychlost vykonání dotazu:** +- 1,3 sekundy

**Výsledek dotazu:**

count(*)	percentage	inventor
8	6.20155	Gould Nigel
4	3.10078	Philips Andrew
4	3.10078	Lee Sangik
3	2.32558	Gordon Gregory Charles
3	2.32558	Trokhan Paul Dennis
3	2.32558	Weisman Paul Thomas
3	2.32558	Dreher Andreas Josef
3	2.32558	Sivik Mark Robert
3	2.32558	Park Bio
3	2.32558	Kim Jeongyun
3	2.32558	Hamad-Ebrahimpour Alyssandrea Hope
3	2.32558	Kim Seonghwan
2	1.55039	Fontaine Gregory
2	1.55039	Lee Yongju
2	1.55039	Watson David John

Obrázek 7.9: Ukázka výsledku dotazu pro scénář č.9

## 8 Závěr

---

todo

# Zkratky

**ACID** Atomicity, Consistency, Isolation, Durability 7, 11, 19, 21

**API** Application Programming Interface 19

**ASCII** American Standard Code for Information Interchange 25

**CRUD** Create, Read, Update, Delete 23, 24

**CSV** Comma-separated values 40

**DBMS** Database Management Systems 6

**EPO** European Patent Office 27

**GPL** GNU General Public License 18

**IPC** International Patent Classification 4, 5

**JSON** JavaScript Object Notation 13, 17, 19, 40

**SQL** MongoDB Query Language 23, 24

**SQL** Structured Query Language 13, 18, 19, 21, 23, 25, 37

**URL** Uniform Resource Locator 41

**USPTO** United States Patent and Trademark Office 2–4, 27

**WIPO** World Intellectual Property Organization 27

**XML** Extensible Markup Language 17, 19, 40–42

# Literatura

- [1] *What Is a Database?* [online]. Oracle. [cit. 21.04.2022]. Dostupné z: <https://www.oracle.com/database/what-is-database/>.
- [2] *What is a Document Database?* [online]. phoenixNAP, 2021. [cit. 27.04.2022]. Dostupné z: <https://phoenixnap.com/kb/document-database>.
- [3] *What Is a Graph Database?* [online]. phoenixNAP, 2021. [cit. 27.04.2022]. Dostupné z: <https://phoenixnap.com/kb/graph-database>.
- [4] *What is a MongoDB Query?* [online]. GeeksforGeeks, 2021. [cit. 27.04.2022]. Dostupné z: <https://www.geeksforgeeks.org/key-value-data-model-in-nosql/>.
- [5] *What Is MongoDB?* [online]. MongoDB. [cit. 28.04.2022]. Dostupné z: <https://www.mongodb.com/what-is-mongodb>.
- [6] *Cypher Query Language* [online]. Neo4j. [cit. 26.04.2022]. Dostupné z: <https://neo4j.com/developer/cypher/>.
- [7] *Graph Modeling Guidelines* [online]. Neo4j. [cit. 27.04.2022]. Dostupné z: <https://neo4j.com/developer/guide-data-modeling/>.
- [8] *Basic Object Oriented Data Model* [online]. GeeksforGeeks, 2021. [cit. 27.04.2022]. Dostupné z: <https://www.geeksforgeeks.org/basic-object-oriented-data-model/>.
- [9] *International Patent Classification (IPC)* [online]. Espacenet, 2016. [cit. 02.05.2022]. Dostupné z: [https://is.espacenet.com/help?locale=en\\_IS&method=handleHelpTopic&topic=ipc](https://is.espacenet.com/help?locale=en_IS&method=handleHelpTopic&topic=ipc).
- [10] *What is Database Management System (DBMS)? – FAQ, Types, and Details* [online]. erp-information. [cit. 27.04.2022]. Dostupné z: <https://www.erp-information.com/database-management-system.html>.
- [11] *Types of Database Languages and Their Uses (Plus Examples)* [online]. indeed, 2021. [cit. 21.04.2022]. Dostupné z: <https://www.indeed.com/career-advice/career-development/database-languages>.
- [12] *IPC Publication* [online]. WIPO, 2022. [cit. 01.05.2022]. Dostupné z: <https://ipcpub.wipo.int/?menulang=en>.

- [13] *The Types of Databases (with Examples)* [online]. Matillion, 2018. [cit. 22.04.2022]. Dostupné z: <https://www.matillion.com/resources/blog/the-types-of-databases-with-examples>.
- [14] *What is MongoDB – Working and Features* [online]. GeeksforGeeks, 2021. [cit. 28.04.2022]. Dostupné z: <https://www.geeksforgeeks.org/what-is-mongodb-working-and-features/>.
- [15] *What is a MongoDB Query?* [online]. GeeksforGeeks, 2021. [cit. 26.04.2022]. Dostupné z: <https://www.geeksforgeeks.org/what-is-a-mongodb-query/>.
- [16] *MongoDB Query Language* [online]. Devopedia, 2021. [cit. 27.04.2022]. Dostupné z: <https://devopedia.org/mongodb-query-language>.
- [17] *What is MySQL?* [online]. MySQL. [cit. 26.04.2022]. Dostupné z: <https://dev.mysql.com/doc/refman/8.0/en/what-is-mysql.html>.
- [18] *What is a Graph Database?* [online]. Neo4j. [cit. 27.04.2022]. Dostupné z: <https://neo4j.com/developer/graph-database/>.
- [19] *PATENT NUMBER, PUBLICATION NUMBER* [online]. IamIP, 2019. [cit. 01.05.2022]. Dostupné z: <http://patentwiki.iamip.com/publication-number>.
- [20] *About PostgreSQL* [online]. PostgreSQL. [cit. 28.04.2022]. Dostupné z: <https://www.postgresql.org/about/>.
- [21] *Design Patent Application Guide* [online]. USPTO, 2017. [cit. 02.05.2022]. Dostupné z: <https://www.uspto.gov/patents/basics/types-patent-applications/design-patent-application-guide#def>.
- [22] *General information concerning patents* [online]. USPTO, 2018. [cit. 02.05.2022]. Dostupné z: <https://www.uspto.gov/patents/basics/general-information-patents>.
- [23] *General Information About 35 U.S.C. 161 Plant Patents* [online]. USPTO, 2017. [cit. 02.05.2022]. Dostupné z: <https://www.uspto.gov/patents/basics/types-patent-applications/general-information-about-35-usc-161#heading-1>.
- [24] *Nonprovisional (Utility) Patent Application Filing Guide* [online]. USPTO, 2018. [cit. 02.05.2022]. Dostupné z: <https://www.uspto.gov/patents/basics/types-patent-applications/nonprovisional-utility-patent#heading-1>.

- [25] AKHTAR, Z. *Relational Database Benefits and Limitations (Advantages & Disadvantages)* [online]. DatabaseTown, 2021. [cit. 27.04.2022]. Dostupné z: <https://databasetown.com/relational-database-benefits-and-limitations/>.
- [26] FRANKLOVÁ, M. *Patenty a jejich klasifikace* [online]. ČVUT. [cit. 03.05.2022]. Dostupné z: [http://pspev.cvut.cz/PSPEV\\_CD/V11/main.html?ID=0](http://pspev.cvut.cz/PSPEV_CD/V11/main.html?ID=0).
- [27] GULATI, V. *Relational Model in DBMS* [online]. Scaler, 2022. [cit. 26.04.2022]. Dostupné z: <https://www.scaler.com/topics/dbms/relational-model-in-dbms/>.
- [28] JOHNSTONE, e. a. N. Renewable Energy Policies and Technological Innovation: Evidence Based on Patent Counts. *Environmental and Resource Economics*. November 2009, 3, 1, s. 38. doi: 10.1007/s10640-009-9309-1. Dostupné z: [https://www.researchgate.net/publication/225430825\\_Renewable\\_Energy\\_Policies\\_and\\_Technological\\_Innovation\\_Evidence\\_Based\\_on\\_Patent\\_Counts](https://www.researchgate.net/publication/225430825_Renewable_Energy_Policies_and_Technological_Innovation_Evidence_Based_on_Patent_Counts).
- [29] KARKI, S. *What Is LevelDB* [online]. C# Corner, 2021. [cit. 28.04.2022]. Dostupné z: <https://www.c-sharpcorner.com/article/what-is-leveldb2/>.
- [30] KENTON, W. *Patent* [online]. Investopedia, 2021. [cit. 02.05.2022]. Dostupné z: <https://www.investopedia.com/terms/p/patent.asp>.
- [31] LEBER, e. a. C. *The Difference Between Trademarks and Design Patents: What you need to know* [online]. Alt Legal, 2021. [cit. 03.05.2022]. Dostupné z: <https://www.altlegal.com/blog/the-difference-between-trademarks-and-design-patents-what-you-need-to-know/>.
- [32] LOBEL, L. *Relational Databases vs. NoSQL Document Databases* [online]. WordPress, 2015. [cit. 27.04.2022]. Dostupné z: <https://lennilobel.wordpress.com/2015/06/01/relational-databases-vs-nosql-document-databases/>.
- [33] LUTKEVICH, B. *database (DB)* [online]. Tech Target, 2021. [cit. 21.04.2022]. Dostupné z: <https://www.techtarget.com/searchdatamanagement/definition/database>.
- [34] PETERSON, R. *What is a Database? Definition, Meaning, Types with Example* [online]. Guru99, 2022. [cit. 21.04.2022]. Dostupné z: <https://www.guru99.com/introduction-to-database-sql.html>.

- [35] SEKHON, S. *What is Neo4j?* [online]. DEV Community, 2020. [cit. 28.04.2022]. Dostupné z: <https://dev.to/sukhbirsekhon/what-is-neo4j-8jc>.
- [36] SINGH, C. *DBMS languages* [online]. BeginnersBook, 2015. [cit. 21.04.2022]. Dostupné z: <https://beginnersbook.com/2015/04/dbms-languages/>.
- [37] THAKUR, D. *What is Object Oriented Database (OODB)? Advantages and Disadvantages of OODBMS.* [online]. ComputerNotes. [cit. 27.04.2022]. Dostupné z: <https://ecomputernotes.com/database-system/adv-database/object-oriented-database-oodb>.
- [38] WILLIAMS, A. *NoSQL Document-Oriented Databases: A Detailed Overview* [online]. RavenDB, 2021. [cit. 27.04.2022]. Dostupné z: <https://ravendb.net/articles/nosql-document-oriented-databases-detailed-overview>.
- [39] YANG, J. *Parts of a Utility Patent Application (Chapter 11)* [online]. OC Patent Lawyer, 2018. [cit. 03.05.2022]. Dostupné z: <https://ocpatentlawyer.com/lesson/basics-utility-patent-application-sections/>.



# A Uživatelská dokumentace

## B Vzhled modulů