



**VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ**

BRNO UNIVERSITY OF TECHNOLOGY

**FAKULTA INFORMAČNÍCH TECHNOLOGIÍ**

FACULTY OF INFORMATION TECHNOLOGY

**ÚSTAV INTELIGENTNÍCH SYSTÉMŮ**

DEPARTMENT OF INTELLIGENT SYSTEMS

**DATABÁZE PRO UKLÁDÁNÍ GENEALOGICKÝCH DAT**

DATABASE FOR GENEALOGICAL DATA

**BAKALÁŘSKÁ PRÁCE**

BACHELOR'S THESIS

**AUTOR PRÁCE**

AUTHOR

**MAREK KAFONĚK**

**VEDOUcí PRÁCE**

SUPERVISOR

**Ing. JAROSLAV ROZMAN, Ph.D**

BRNO 2021

## Zadání bakalářské práce



Student: **Kafoněk Marek**  
Program: Informační technologie  
Název: **Databáze pro ukládání genealogických dat**  
**Database for Genealogical Data**  
Kategorie: Databáze

### Zadání:

1. Nastudujte problematiku genealogie. Nastudujte relační databáze a grafovou databázi Neo4j.
2. Na základě nastudovaných znalostí navrhnete MySQL databázi pro ukládání genealogických dat (příbuzenské a další vztahy, vazba mezi původním genealogickým záznamem a záznamem v databázi). Dále navrhnete skripty pro přenos dat z Neo4j do MySQL databáze.
3. Navrženou databázi a potřebné skripty vytvořte a proveďte překopírování dat z Neo4j do MySQL.
4. Vytvořenou databázi otestujte.

### Literatura:

- Prostředníková Hana: Pokročilé zobrazování genealogických dat, bakalářská práce, Brno, FIT VUT v Brně, 2016.
- Valecký Dušan: Zobrazování genealogických dat, bakalářská práce, Brno, FIT VUT v Brně, 2017.

Pro udělení zápočtu za první semestr je požadováno:

- První dva body zadání

Podrobné závazné pokyny pro vypracování práce viz <https://www.fit.vut.cz/study/theses/>

Vedoucí práce: **Rozman Jaroslav, Ing., Ph.D.**

Vedoucí ústavu: Hanáček Petr, doc. Dr. Ing.

Datum zadání: 1. listopadu 2021

Datum odevzdání: 11. května 2022

Datum schválení: 3. listopadu 2021

## Abstrakt

Tato práce se zabývá možným převodem genealogických dat mezi grafovou a relační databází, kdy se zde navazuje na již existující data uložená v grafové databázi Neo4j. Možné převody budou sloužit k tomu, aby se s danými daty mohlo podle potřeby pracovat buď pomocí dotazovacího jazyku SQL v relační databázi, nebo pomocí jazyku Cypher v grafové databázi.

## Abstract

This work deals with the possible transfer of genealogical data between graph and relational databases, which builds on existing data stored in the graph database Neo4j. Possible transfers will be used to work with the data as needed using SQL in a relational database, or using Cypher in a graph database.

## Klíčová slova

Genealogy, MySQL, Relační databáze, Neo4j, Grafová databáze

## Keywords

Genealogy, MySQL, Relational database, Neo4j, Graph database

## Citace

KAFON K., Marek. *Databáze pro ukládání genealogických dat*. Brno, 2021. Bakalářská práce. Vysoké učení technické v Brně, Fakulta informačních technologií. Vedoucí práce Ing. Jaroslav Rozman, Ph.D.

# Databáze pro ukládání genealogických dat

## Prohlášení

Prohlašuji, že jsem tuto bakalářskou práci vypracoval samostatně pod vedením pana Ing. Jaroslava Rozmana Ph.D.. Uvedl jsem všechny literární prameny, publikace a další zdroje, ze kterých jsem čerpal.

.....  
Marek Kafoník  
18. května 2022

## Poděkování

Chtěl bych poděkovat mému vedoucímu panu Ing. Jaroslavovi Rozmanovi Ph.D. za ochotu a vstřícnost při řešení bakalářské práce

# Obsah

<b>1 Úvod</b>	<b>3</b>
<b>2 Studium problematiky</b>	<b>4</b>
2.1 Genealogie . . . . .	4
2.1.1 Vyhledávání genealogických dat . . . . .	4
2.1.2 Matriční zápisy . . . . .	5
2.1.3 Tvorba výstup . . . . .	6
2.2 Relační databáze . . . . .	7
2.2.1 Tvorba relační databáze . . . . .	7
2.2.2 Tvorba MySQL databáze . . . . .	10
2.3 Grafová databáze . . . . .	12
2.3.1 Neo4j . . . . .	12
<b>3 Návrh řešení</b>	<b>13</b>
3.1 Návrh MySQL databáze . . . . .	14
3.1.1 Matrika . . . . .	14
3.1.2 Matriční záznamy . . . . .	15
3.1.3 Osoba . . . . .	20
3.1.4 Evidence místa . . . . .	22
3.2 Návrh skript pro převod . . . . .	24
3.2.1 Převod z Neo4j do MySQL . . . . .	24
3.2.2 Převod z MySQL do Neo4j . . . . .	25
<b>4 Implementace a testování</b>	<b>27</b>
4.1 Implementace MySQL . . . . .	27
4.2 Implementace Neo4j . . . . .	30
4.3 Pomocný slovník . . . . .	33
4.4 Převod dat . . . . .	33
4.4.1 Neo4j do Mysql . . . . .	33
4.4.2 Mysql do Neo4j . . . . .	34
4.5 Testování . . . . .	35
4.5.1 Rychlost převodu . . . . .	36
4.5.2 Správnost převodu . . . . .	36
<b>5 Závěr</b>	<b>38</b>
<b>Literatura</b>	<b>39</b>

<b>A</b>	<b>Obsah přiloženého paměťového média</b>	<b>40</b>
<b>B</b>	<b>Návrh genealogické databáze</b>	<b>41</b>

# Kapitola 1

## Úvod

Při nedávné zmínce jednoho z vyukaujících na přednášce, ohledně vyhledávání genealogických dat, jsem začal přemýšlet, kdo vlastně byli mí předkové. Byli to venkované, bohatí hospodáři, měšťané nebo snad dokonce šlechta? Jsou mí předci spojeni pouze s životem na území naší republiky nebo mě mohou sahají i do ciziny?

Tyto otázky se prací na téma této bakalářské práce sice nedozvím, ale rozhodně mi pomohou, že si nastudovat informace ohledně práce s genealogickými daty a případně mít možnost využít výslednou implementaci pro mé vyhledávání.

Samotnou práci jsem si podle bodů zadání rozdělil na 4 části. První část se vnuje studiu problematiky tématu zadané práce, kde jsem popsal potřebné informace k nastudování tak, abych ve druhé části mohl provést návrh mého řešení. Za zmínku stojí, že už v návrhu je třeba zohlednit, že toto téma navazuje na již existující diplomovou práci Generování rodokmenů z matričních záznamů, provedenou na fakultě informatiky, od Ing. Lucie Tušimové [12]. Další částí je samotná implementace, kdy je pomocí předchozího návrhu vytvořena MySQL databáze, a k ní skripty potřebné k přesunutí dat z grafové databáze, která byla vytvořena ve zmínované návazné diplomové práci. Nakonec v poslední části je zmíněno testování výsledného implementovaného řešení.

## Kapitola 2

# Studium problematiky

K zadání mé práce bylo potřeba nastudovat problematiku témat související se zadáním. Jedná se celkem o tři problematiky a to genealogie, relační databáze MySQL a grafová databáze Neo4j. U Genealogie jsem, krom samotného významu tohoto oboru, studoval tvorbu možných genealogických výstupů (rodokmen, rozvod, vývod) a princip získávání genealogických dat. U MySQL databáze jsem se v novém studiu její funkcionality, tvorby ER diagramu a jejich převod do databáze a případně další práce s databází. V poslední době jsem se zaměřil na studium Neo4j, kde jsem, stejně jako u MySQL, musel nastudovat její funkcionality, princip ukládání dat, a pro mé téma důležité, export dat.

### 2.1 Genealogie

Jedná se o pomocnou vědu historickou, která zkoumá vztahy mezi lidmi. První zmínky o genealogii pocházejí již ze starověku, odkud se dochovaly například rodokmeny bohů starověkého egejského nebo pěstříbuzenské vztahy mezi faraony v Egyptě. Každopádně pro nás zajímavější jsou pozůstatky rodopis evropských královských rodin z raného středověku [1][7][8].

Genealogii dělíme na veřejnou a soukromou. Pro veřejnost je známější soukromá, které se také říká rodopis. Tato věda se omezuje na zájmy jednotlivců, kteří se zabývají historií své rodiny z osobních důvodů, poznání dějin, dědictví a dalších jimi podobných. V veřejné genealogii se zabývá rodinnými vztahy jedince, bez ohledu na jeho zájmy, a data z ní získaná se dále používají v rámci návazných veřejných oborů.

#### 2.1.1 Vyhledávání genealogických dat

K získávání genealogických dat slouží matriky, které se dají rozdělit na živé a neživé. Bádatelé, jak se nazývají lidé, kteří vyhledávají genealogické záznamy, mají přístup pouze k tzv. neživým, které jsou k dispozici v archívech. Živé matriky obsahují současně žijící lidi a nikdo by k nim nemohl mít přístup.

Původ slova matrika pochází z latinského slova matricula, což v překladu znamená seznamy duchovních. První zmínky o zápisu do matrik na území dnešní České republiky pochází z 16. století, ale jedná se pouze o zápisy evangelické církve. Katolická církev začala zapisovat údaje až v pozdějších dobách. První zápisy matrik byly pouze formou seznamu, který



umožňoval duchovnímu mít přehled o osobách v jeho farnosti. Do matrik kromě duchovních zapisovali i nevěřící lidé, kterých ale v té době nebylo tolik. V průběhu let, a s přibývajícím množstvím informací, docházelo k tomu, že se různé zápisy mohly lišit. Jelikož v té době nebylo přesně stanoveno jakou podobu má mít zápis do matriky, pochází odtud řada nejasností. A už duplicitní záznamy pro určitou osobu, zápis události do špatné matriky, nepřesné zápisy a podobně. Docházelo i k takovým případům, že farář odešel do jiné farnosti, vzal si s sebou matriku bývalé farnosti, kam poté začal zapisovat osoby z farnosti nové. Až od roku 1784 se matriky staly úředním dokumentem, který mohl vést ve správě především farnosti, protože byly zodpovědné jak za křty, svatby, tak i za pohřby. U zapisování do matrik se stanovilo jaké informace mají dané zápisy obsahovat, aby nedocházelo k chybám z předchozích let. I přes veškeré úsilí sjednotit zapisované informace stále docházelo k chybám při zápisu, protože farář i přesto zapisoval důležité informace, které dříve většina z nich zapisovala a ztížilo se tím bádání po rodinné historii.

V roce 1870 začaly vznikat civilní matriky, které byly vedeny u okresních soudů. Zpočátku byly určeny pouze pro nevěřící lid. Později se do těchto matrik začali zapisovat i věřící lidé, kteří uzavírali civilní sňatky.

Až od roku 1949 bylo zavedeno, že matriky připadnou státu a budou pod správou jednotlivých národních výborů. Část matrik, připadlých do roku 1875, byla předána do státních archivů a část zůstala na úřadech. I přes přesunu matrik do státní správy mají své vlastní zápisy farnosti, které i farářové. Na počátku 20. let 21. století se v archivu nachází matriky narozených do roku 1910 a matriky oddaných a zemřelých do roku 1930 [7].

### 2.1.2 Matriční zápisy

Zápisy se rozdělují na tři části, kterými jsou křestní zápisy, oddací zápisy a úmrtní zápisy. V křestních zápisech se evidují záznamy o narozeném dítěti. Přesněji jeho jméno, příjmení, datum a místo narození. Kromě informací o narozeném se uvádí i oba rodiče, případně jestli se jedná o nemanželské nebo adoptované dítě. Jak už název napovídá, součástí křestních zápisů jsou i informace o samotném křtu, konkrétně se zde uvádí datum a místo křtu, jméno křtitele a jména kmotrů. Ohledně křtu se zde evidují i informace o nouzovém křtu, který se konal předasně, a to v případě, že dítěti hrozila smrt. Dalšími možnými informacemi v těchto zápisech je náboženství narozeného, které se převážně dědí po svých rodičích, jméno porodní báby, která mohla provádět již zmíněný nouzový křest a další doplňující informace ohledně narozeného. Druhým typem zápisů jsou oddací zápisy. Uvádí se zde základní informace jako, datum a místo svatby, u snoubenců se evidují jejich jména, věk, bydliště, stav a rodiče. U oddacích zápisů dále nesmí chybět jméno oddávajícího a jméno svůdky. Dalšími možnými informacemi je náboženství prováděné svatby a ohláška, která má informovat o konání svatby. U oddacích záznamů se také uvádí informace o případném rozvodu manželů. Poslední, a z těchto zápisů nejnovější druh, jsou úmrtní zápisy. Evidují se zde informace o jméně, věku a příčině smrti zemřelého a také místo a datum úmrtí. Ohledně pohřbu jsou zde informace o zaopatření zemřelého, datumu a místu pohřbu. V některých případech se uvádí i svůdci pohřbu, a to převážně u pohřbu mrtvorozených, kteří se uvádějí jako již ve zmíněných křestních zápisech, tak i právě v úmrtních zápisech.

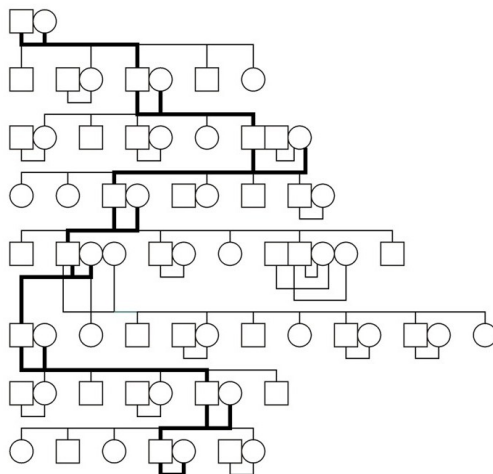
Kromě těchto zmíněných typů matričních zápisů existují i další doplňkové zápisy, které se využívají při zaznamenávání údajů o majetku osob. Mezi takové zápisy patří zápisy v gruntovních knihách, zápisy v urbářích, rubriky katastrů a soupisy poddaných.

### 2.1.3 Tvorba výstupů

Po bádání a získání potřebných dat je potřeba tato data prezentovat v nějaké srozumitelné podobě. Pro tento případ jsou k dispozici tři různé typy, které nyní stručně popíšu.

#### Rodokmen

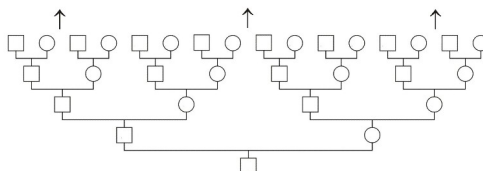
Pro širší veřejnost je nejznámějším druhem zápisu genealogických dat právě rodokmen [2.1]. Díky rodokmenu se dohledávají všichni předci v mužské linii. Pro jakoukoliv další práci s genealogickými daty slouží jako odrazový místeček a jeho pomocí lze provádět další pátrání a tvorbu ostatních druhů výstupů.



Obrázek 2.1: Schéma tvorby rodokmenu [11].

#### Vývod

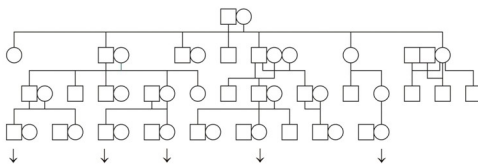
Jedná se o podobný druh jako rodokmen, rozdíl je v tom, že zde se zobrazují všichni předci výchozí osoby a to nejen mužští, ale i ženská linie. Vývod [2.2] je nejčastěji používaný tvar mezi badateli, kdy výsledkem jsou všechny linie rodinných předků. Výsledná podoba vývodu bývá zakreslena v podobě stromu.



Obrázek 2.2: Schéma tvorby vývodu [11].

## Rozrod

Nejsložitější variantou pro zobrazení genealogických dat je právě rozrod 2.3. Narozdíl od předchozích variant se zde postupuje opačně, tedy od nejstaršího dohledaného předka do současnosti. Dohledávají se všichni potomci stejného předka, potomci synů i nemanželské děti neprovdaných dcer. V případě tvorby rozrodu je často patrné postupné vymírání rodů nebo naopak jejich postupné rozrůstání.



Obrázek 2.3: Schéma tvorby rozrodu [11].

Existuje také možnost všechny zmíněvané varianty rozrodu kombinovat. Například se může pátrat pouze po mužské linii, a k ní přidat sourozence každé generace, nebo vyhledat předky ve všech liniích a k nim zahrnout všechny jakékoli příbuzné osoby.

## 2.2 Relační databáze

V rámci tohoto projektu jsem studoval relační databázi MySQL vyvíjenou společností Oracle, která ji poskytuje v rámci open source software.

Základním kamenem relační databáze je tabulka, která obsahuje sloupce sloužící jako názvy entit a jednotlivé řádky, které znázorňují samotné entity. Pro práci s databází se používá dotazovací jazyk SQL, pomocí kterého můžeme například zobrazovat, vkládat, upravovat nebo mazat jednotlivé záznamy v tabulkách.

### 2.2.1 Tvorba relační databáze

#### Konceptuální modelování

Po základní analýze požadavků pro tvorbu relační databáze se vytváří ER diagram, který uvádí vztah mezi entitami v klidovém režimu. Naopak v ER diagramu se neuvádí jednotlivé operace.

Jako entitu zde chápeme objekt reálného světa, který je odlišný od ostatních objektů. Jednotlivé entity jsou součástí entitních množin, které jsou následně použity ve výsledném ER diagramu. Název entitní množiny by měl jasně specifikovat jaké entity se budou do dané množiny ukládat. Vlastnosti entit se nazývají atributy, které se rozlišují na jednoduché (atomické) nebo složené (skládající se z více atributů). Každá entita musí obsahovat jeden atribut, který je potřeba k jednoznačné identifikaci. Zde se tento identifikátor nazývá jako primární klíč v modelu ER diagramu a označuje se «PK».

V ER diagramu se mezi entitními množinami používají binární, ternární nebo unární vztahy. Nejvíce používané jsou binární a unární. Ternární vztahy jsou speciální a jedná se především o vztahy mezi třemi entitními množinami, případně mezi více entitami. Jednotlivé vztahy by neměly mít stručný, ale výstižný název. Název vztahu nemusí být v případě, když vztah vyplývá z názvů entitních množin.

Po entitních množinách a vztazích entitních množin se při tvorbě ER diagramu používá kardinalita, která udává vždy minimální a maximální počet dat vstupující do vztahu s jinou entitou. Možnými kardinalitami jsou 1:1 (jedna entita je ve vztahu také pouze s jednou entitou), 1:N (jedna entita je ve vztahu s více entitami) a M:N (více entit je ve vztahu také s více entitami).

V rámci ER diagramu se mohou použít dvě rozšíření. Prvním rozšířením je slabá entitní množina («weak»), která je závislá na jiné entitní množině («identif»). Slabá entitní množina používá identifikátor dominantní silné entitní množiny a atribut, který identifikuje samotnou slabou entitní množinu. Tento atribut se nazývá diskriminátor a označuje se «D». Nejastěji je slabá entitní množina ve vztahu se silnou entitní množinou v poměru 1:N. Druhé rozšíření je generalizace/specializace, kde se jedná o vztah dvou entitních množin, kde jedna entitní množina rozšiřuje vlastnosti druhé množiny. V praxi to znamená, že pro entitní množiny je jedna generalizující entitní množina se společnými atributy a pak více specializací, které mají odlišné atributy a další atributy generalizující entitní množiny.

### **Transformace konceptuálního modelu na tabulky relační databáze**

Tato část se zaměřuje na návrh tabulek relační databáze z vytvořeného konceptuálního modelu, například z ER diagramu, kdy entitní množiny jsou jednotlivé tabulky. Vztahy mezi nimi jsou znázorněny vazbami mezi tabulkami pomocí primárních a cizích klíčů.

Při návrhu databáze dochází k jejímu chybnému návrhu hlavně v případě opakujících se informací (redundancí), nemožností zobrazovat určitou hodnotu, složitou kontrolu integritních omezení a návrhu zbytečných tabulek.

Na začátku návrhu je potřeba odstranění složených a vícehodnotových atributů na jednotlivé atributy. Následně je dobré se zabývat transformací silných entitních množin, které se dají převést na jednotlivé tabulky, kdy jednotlivé atributy určují sloupce tabulky, jednotlivé řádky jsou samotné entity a každý řádek má přiřazený jeden primární klíč, tak jako primární klíč u entity v entitní množině.

U vytváření vztahů mezi tabulkami se musí zohledňovat kardinalita z ER diagramu a podle ní přidávat tabulkám cizí klíče, které slouží právě pro vytvoření vazeb mezi tabulkami. U kardinality 1:1 je na návrháři, jaké tabulce přidat cizí klíč odkazující na primární klíč druhé tabulky a v případě existujícího dalšího atributu vztahové množiny, bude tento atribut přidán k tabulce s cizím klíčem. Narozdíl od předchozí kardinality, tak u kardinality 1:N, je dané, ke které tabulce bude vložený cizí klíč, případně další atributy vztahové množiny a to k tabulce s kardinální hodnotou N. U kardinality M:N je nutné vytvořit speciální vazební tabulku, které bude obsahovat jeden cizí klíč odkazující na primární klíč jedné tabulky a druhý cizí klíč odkazující na primární klíč druhé tabulky. Z těchto dvou cizích

klí vznikne složený primární klí vazební tabulky. Do vazební tabulky se přidávají atributy vazební množiny. V případě existence vztah vyššího stupně (tj. 3 a více), bude zapotřebí vždy vytvořit vazební tabulky spojující všechny tři samotné tabulky.

V případě existence slabé entitní množiny se bude postupovat jako u entitních množin s kardinalitou 1:N, kdy silná entitní množina se transformuje bez změny u slabé entitní množiny se transformují všechny její atributy do jedné tabulky společně s cizím klíčem, odkazujícím na primární klí tabulky pro silnou entitní množinu.

Při transformaci generalizace/specializace lze postupovat třemi způsoby. První je, že pro generalizaci se vytvoří jedna tabulka s jejími atributy a následně se vytvoří tabulky pro jednotlivé specializace, kdy každá specializace obsahuje svoje vlastní atributy a cizí klí odkazující na primární klí generalizační tabulky. Druhý způsob je vytvoření tabulek podle počet specializací, kdy každá tabulka bude obsahovat atributy generalizace a následně jednotlivé tabulky budou mít vlastní atributy jednotlivých specializací. Třetí a zároveň poslední možnost je, že budou vytvořeny dvě tabulky. Jedna bude obsahovat atributy generalizace a druhá bude obsahovat cizí klí odkazující na tabulku generalizace, potom atributy všech specializací a pomocný sloupec, který bude určovat, o kterou specializaci se jedná.

### Normalizace schématu databáze

Poslední část pro splnění použitelného návrhu databáze je normalizování. Normalizace se zaměřuje na úpravu chybného návrhu databáze, a to především na redundanci<sup>1</sup>, nemožnost reprezentovat určitou informaci a na složitou kontrolu integritních omezení. K tomuto se používají normální formy pro které platí, čím vyšší stupeň normální formy, tím kvalitnější databázový návrh. Normální formy dále definují, jaké vlastnosti musí být splněny v závislosti mezi atributy tabulky.

Při normalizaci databáze je zapotřebí zjistit funkční závislosti, mezi které patří triviální funkční závislost, plná funkční závislost a tranzitivní závislost. K dispozici máme 6 základních normálních forem v následující podobě, 1. normální formu (1NF), 2. normální formu (2NF), 3. normální formu (3NF), Boyce-Coddovu normální formu (BCNF), 4. normální formu (4NF) a 5. normální formu (5NF). Každá normální forma definuje určitou vlastnost, kterou musí tabulka relační databáze mít, aby tuto formu splnila. Kromě splnění určité vlastnosti jednotlivých normálních forem, musí být také splněny všechny předchozí normální formy. V případě, že tabulka nesplňuje kterou z normálních forem, je potřeba provést bezztrátovou dekompozici, při které rozdělíme tabulku na více tabulek bez přidání nových atributů nebo odebrání již existujících, zachováme-li závislosti jednotlivých atributů a zbavíme se případné redundance.

Pro splnění 1NF musí všechny jednoduché atributy obsahovat atomické hodnoty. Pro splnění 2NF je potřeba splnění vlastnosti, kdy každý neklíčový atribut je plně funkčně závislý na některém z kandidátních klíčů tabulky. 3NF je splněna za předpokladu, že neexistuje žádný neklíčový atribut, který je tranzitivně závislý na některém kandidátním klíči tabulky. Jelikož předchozí zmíněné normální formy se vztahují k závislosti neklíčových atributů, byla zavedena BCNF, která zajišťuje 3NF a odstraňuje závislosti mezi klíčovými atributy. 4NF a 5NF se moc nepoužívají, protože už po splnění BCNF je odstraněna re-

---

<sup>1</sup>Informace, které se opakují

dundance. Pro ur ení výsledné normální formy databáze se vybere nejmenší maximální normální forma[6, 4, 9, 10, 3].

## 2.2.2 Tvorba MySQL databáze

Po studiu návrhu databáze bych se cht l v této ásti v novat studii samotné tvorby MySQL databáze.

Pro práci s MySQL databází se používají dotazy v deklarativním programovacím jazyce SQL. Programovat databázi m žeme bu dotazováním p ímo v programovacím prost edí pro SQL nebo pomocí jiných jazyk a jejich p íslušných knihoven pro dotazování v SQL (nap . C, Python).

Prvním krokem p í tvorb databáze je její samotné vytvo ení, pro které se používá dotaz `CREATE DATABASE jmeno_databaze;`, který založí prázdnou databázi, s kterou je možné dále pracovat.

K tomu abychom ur ili s jakou databází chceme pracovat, je dobré vždy na za átku dotazovacího skriptu použít výb r databáze pomocí `USE jmeno_databaze;`. Výb r databáze se dá používat i jinými zp soby, a to nap íklad uvedením názvu databáze p ed její tabulkou.

Když je vytvo ená databáze a jsme k ní p ípojení, mohou se v ní vytvá et tabulky, a to pomocí dotazu `CREATE TABLE jmeno_tabulky;`. U vytvá ení tabulky je pot eba ješt definovat její sloupce, k nim datové typy a nastavit, jestli daný sloupec m že nebo nem že obsahovat hodnotu NULL. P í vytvá ení tabulky by se nem lo zapomenout nastavit primární klí tabulky a p ípadn p í adit cizí klí e, které budou sloužit k propojení s jinou tabulkou. Pro rychlejší vyhledávání je dobré v tabulce p í adit indexy, tím zaru íme, že si databázový systém zapamatuje, kde se p íbližn nachází indexovaný sloupec a nemusí procházet p í vyhledávání celou tabulku.

Pro p ípad, že je vytvo ená tabulka a je pot eba upravit nebo p ídat sloupec, slouží dotaz `ALTER TABLE jmeno_databaze prikaz;`. Jako prikaz je možné použít bu `ADD COLUMN` pro vložení nového sloupce daného datového typu na konec tabulky, nebo `DROP COLUMN` pro odstran ní daného sloupce.

K odstran ní tabulek z databáze se použije dotaz `DROP TABLE jmeno_tabulky;`. Pokud daná tabulka obsahuje sloupce ve form cizího klí e, který zaru uje propojení s jinou tabulkou v databázi, tak se tato tabulka neodstraní a dotaz skon í chybovou hláškou.

Když je tabulka vytvo ená, m že se s ní za ít pracovat, a to provád ním manipulaci s daty, pod které spadá nap íklad vkládání nových záznam , zm na existujících záznam , zobrazení záznam nebo mazání záznam .

Jednou ze zmín ných možností je vkládání nových záznam do tabulky pomocí p íkazu `INSERT INTO jmeno_tabulky (navez_sloupce) VALUES(hodnota_sloupce);`. Data se úsp šn vloží do tabulky pokud je p íazená hodnota všem sloupc m, které nemají povoleno vkládat hodnotu NULL a zároveň nemají nastavenou defaultní hodnotu nebo p í tvorb tabulky nebyl použit p íkaz `AUTO INCREMENT` pro automatické p ídlování hodnot pro daný sloupec.

Pokud je potřeba opravit hodnotu sloupce existujícího záznamu tabulky, poslouží k tomu dotaz `UPDATE nazev_tabulky SET jmeno_sloupce = nova_hodnota;`. Součástí tohoto dotazu bývá často i dotaz `WHERE`, kde po zadání vhodné podmínky se vybere určitý záznam z tabulky.

K zobrazení hodnot vybraných sloupců z dané tabulky, které splývají určité podmínky, slouží dotaz `SELECT nazev_sloupce FROM nazev_tabulky WHERE podminka;`. Při zobrazování záznamů pomocí dotazu `SELECT` je možné použít další doplňkové dotazy, jako je například `SELECT DISTINCT . . .`, který nezobrazí duplicitní záznamy. Je-li potřeba sloučit záznamy v tabulce podle daného sloupce, použije se dotaz `GROUP BY` a když je potřeba seřadit záznamy podle daného sloupce, použije se dotaz `ORDER BY`. Když je potřeba zobrazit záznamy z různých tabulek, musí se propojit tyto tabulky dotazem `JOIN`. K propojení tabulek je k dispozici více `JOIN` dotazů a to například `RIGHT JOIN`, `LEFT JOIN`, `OUTER JOIN`, `INNER JOIN` a `CROSS JOIN`. Pro případ, kdy je potřeba sloučit dva rozdílné `SELECT` dotazy, použije se dotaz `UNION`, kdy se vybraná data zobrazí za sebou.

Jako poslední dotazy pro manipulaci s daty bych zmínil dotaz pro mazání záznamů z tabulky, `DELETE FROM nazev_tabulky WHERE podminka;`. Pomocí kterého se odstraní záznam splývající danou podmínkou. Záznam tímto dotazem se ale odstraní pouze pokud neobsahuje žádná návazná data.

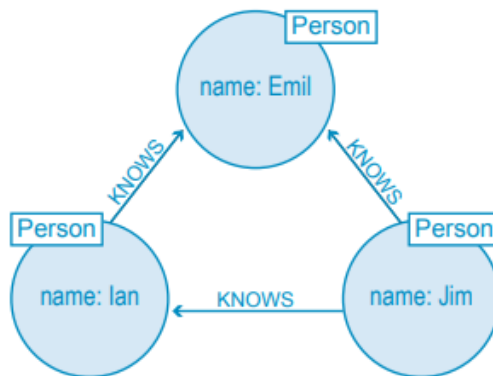
## 2.3 Grafová databáze

Jelikož zadání mé práce navazuje na už existující Diplomovou práci, kde výsledkem je nová vzniklá Neo4j databáze s genealogickými daty. Bylo potřeba nastudovat i princip fungování grafové databáze Neo4j, tak aby z ní mohl být proveden přechod do nově vzniklé MySQL databáze.

### 2.3.1 Neo4j

Neo4j je grafová databáze, která je k dispozici jako open-source. Jedná se o NoSQL nativní grafovou databázi, která je implementována v programovacím jazyce JAVA. Pro dotazování nad daty se používá primárně dotazovací jazyk Cypher, který je podobný jazyku SQL, akorát optimalizovaný pro grafy. S grafovou databází Neo4j je možné pracovat i v jiných programovacích jazycích, které mají potřebné knihovny. Kromě programovacího jazyku Java, to jsou například jazyky JavaScript a Python. Velkou výhodou u grafové databáze Neo4j je v tom, že její transakce mají ACID vlastnosti a nejenom díky této vlastnosti jí používají pro své podnikání velké společnosti.

V grafové databázi Neo4j se místo tabulek používají pro reprezentaci dat uzly a jejich vlastnosti [2, 4], které se mezi sebou spojují hrany a dohromady tvoří právě graf. Samotný graf této databáze je možné si představit jako když se nakreslí graf na papír. Na rozdíl od relační databáze a používání různých spojovacích dotazů tabulek, které je poměrně zdlouhavé a náročné, grafová databáze se soustředí obzvláště na vztahy a procházení mezi nimi [2, 5].



Obrázek 2.4: Ukázka grafu databáze Neo4j. Převzato z [2].



## Kapitola 3

### Návrh řešení

Po nastudování potřebných informací bylo potřeba v první řadě vytvořit návrh MySQL databáze pro ukládání genealogických dat, díky kterému bude možné provést následnou implementaci databáze. Návrh databáze vychází z řešení diplomové práce Lucie Tušimové z roku 2020 3.1, na kterou tato práce navazuje. Přesněji se jedná o návrh MySQL databáze, která bude odpovídat grafové databázi Neo4j z již zmínované diplomové práce. Společně s návrhem databáze bude potřeba navrhnout případné skripty pro přesun dat mezi grafovou databází a relační databází, ale i naopak.



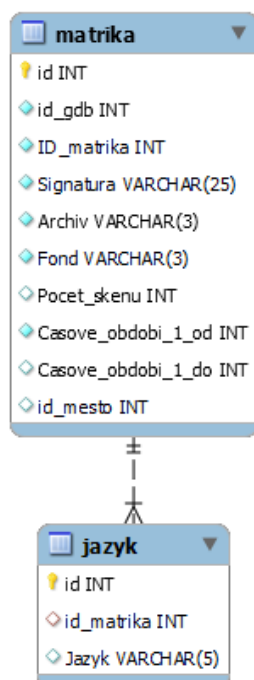
Obrázek 3.1: Příklad z provedení Neo4j databáze

## 3.1 Návrh MySQL databáze

Pro provedení návrhu databáze pomocí konceptuálního modelování, transformace modelu na tabulky relační databáze a jejich normalizaci jsem vytvořil následující databázi. Názvy sloupců tabulek z velké části navazují na názvy vlastností uzlů grafové databáze, zbylé sloupce jsou pojmenovány tak, aby jejich název odpovídal tomu, co se v nich zaznamenává. Pro přehlednost přesouvání dat se v každé tabulce nachází sloupec `id_gdb`, který obsahuje hodnotu jednoznačného identifikátoru přírodního uzlu. Samotný návrh databáze jsem pro přehlednost rozdělil na jednotlivé části, které zde představím a stručně popíši. Návrh celé databáze je součástí přílohy B.1.

### 3.1.1 Matrika

První část, kterou bych zde popsal, je hlavní tabulka `matrika` sloužící k evidenci matrik. V této tabulce se budou ukládat základní data jako označení matrice signaturou, v jakém archívu se matrice nachází, fond matrice, v případě zaznamenání, kolik má daná matrice naskenovaných stran, z jakého období pochází ve formě intervalu let od, do. Do tabulky se také přenesou data `ID_matrika`, která jsou evidována v přírodní grafové databázi. Jelikož jedna matrice může obsahovat více jazykových forem, je potřeba navrhnout další pomocnou tabulku, zde tabulku `jazyk`, kam se budou tyto jazykové formy evidovat a zároveň odkazovat k tabulce `matrika`. Kromě zmíněné pomocné tabulky, pro evidenci více jazyků k jedné matrice, je potřeba navrhnout další pomocné tabulky, které budou zároveň evidenci obcí, kde se matrice nachází. Tyto pomocné tabulky zároveň budou evidovány více obcí. Popis těchto tabulek se ale budu vnovat až v části 3.1.4, týkající se evidenci míst.



Obrázek 3.2: Návrh tabulek pro evidenci matrice.

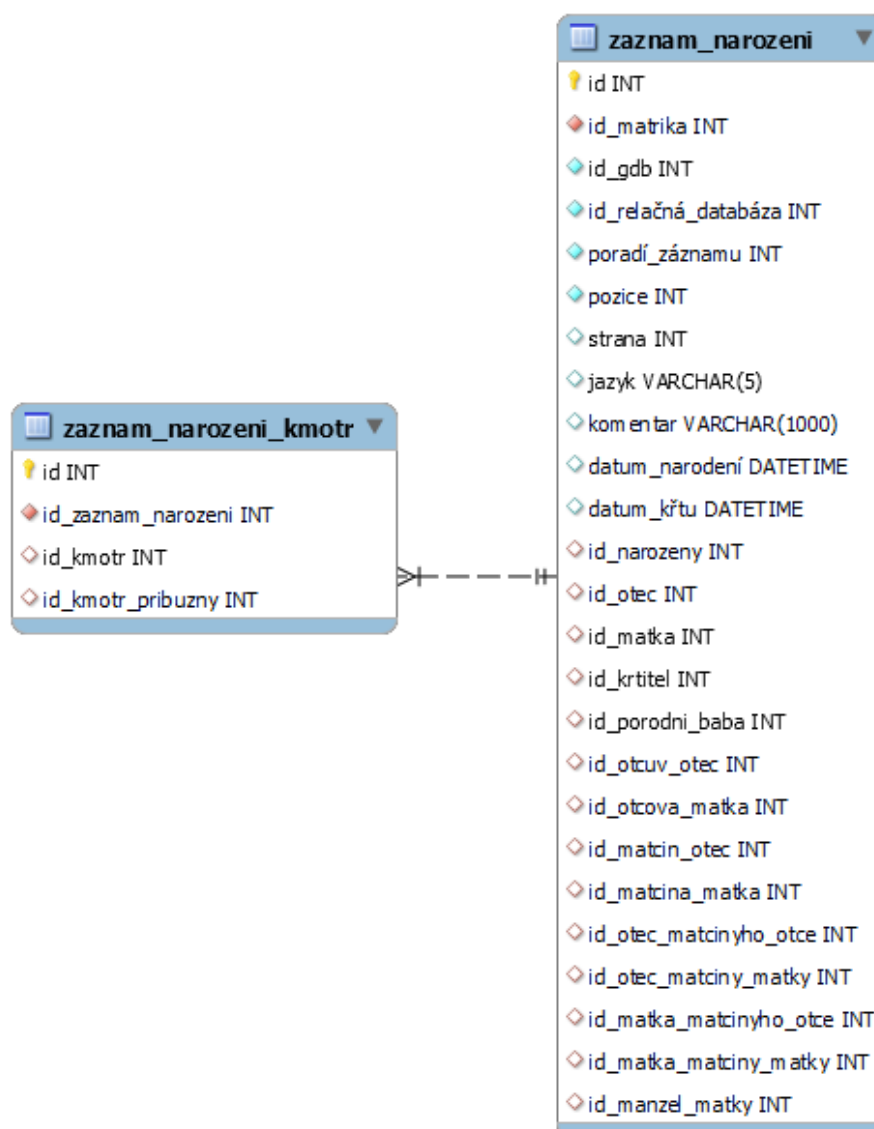
### 3.1.2 Matriční záznamy

Pro evidenci osob v matrikách slouží jednotlivé matriční záznamy. Tyto matriční záznamy se rozdělují do tří skupin pro evidenci narození, manželského sňatku a úmrtí. Do zde navržených tabulek se budou převádět data z uzlů `Záznam_o_k`te, `Záznam_o_svatbe` a `Záznam_o_úmrti` a jejich hrany propojení s uzly osob. Pro tyto úkony se u každého druhu matričního záznamu evidují základní informace týkající se pořadí záznamu, jeho pozice a strany kde se nachází. Dále se zde eviduje jazyk zápisu a případný komentář. V matričních záznamech je potřeba také navrhnout sloupec pro evidenci jednoznačného identifikátoru k povodní relační databázi, odkud se přesouvají data do grafové databáze Neo4j. Jelikož tabulky matričních záznamů navazují na tabulky matrik, je potřeba mít k dispozici sloupec, který bude formou cizího klíče odkazovat právě na návaznou tabulku matriky. Zbývá v tabulce sloupec, který bude také formou cizího klíče odkazovat na jednotlivé záznamy v tabulce osoba 3.1.3, která bude popsána ve zmiňované části.

#### Záznam narození

Tato tabulka 3.3, která slouží k evidenci nově narozených osob. Kromě již zmíněných sloupců pro základní evidenci matričních záznamů bude zde potřeba mít k dispozici sloupec pro evidenci data narození a to sloupec `datum_narodení`. Zbývající sloupce, jak už bylo zmíněno, budou formou cizího klíče odkazovat na záznamy v tabulce osob, čímž se vlastně nastaví veškeré vztahy mezi osobami a matričním záznamem. Prvním zmíněným sloupcem je sloupec `i_d_narozeny`, který k matričnímu záznamu popisuje narozenou osobu. Další dva sloupce se týkají rodičů, kdy pomocí sloupce `i_d_otec` je evidován otec narozené osoby a pomocí sloupce `i_d_matka` jeho matka. Po těchto poměrně jasných druhů osob k evidenci záznamu narození, je možné evidovat další návazné osoby a to k tetaře v podobě `i_d_krtitel`, porodní bábu v podobě `i_d_porodni_baba`. Dalšími sloupci pro evidenci návazných osob jsou sloupce určující prarodiče narozeného z otcovy strany a prarodiče i praprarodiče narozeného z matčiny strany. Posledním zde navrženým sloupcem je sloupec `i_d_manzel_matky` pro případnou evidenci nevlastního otce.

K záznamu narození se ještě eviduje kmotr a případně jeho příbuzný, ale protože může nastat situace, že k jednomu záznamu narození existuje více kmotrů a každý má svého příbuzného, je navržena pomocná tabulka `zaznam_narozeni_kmotr`.



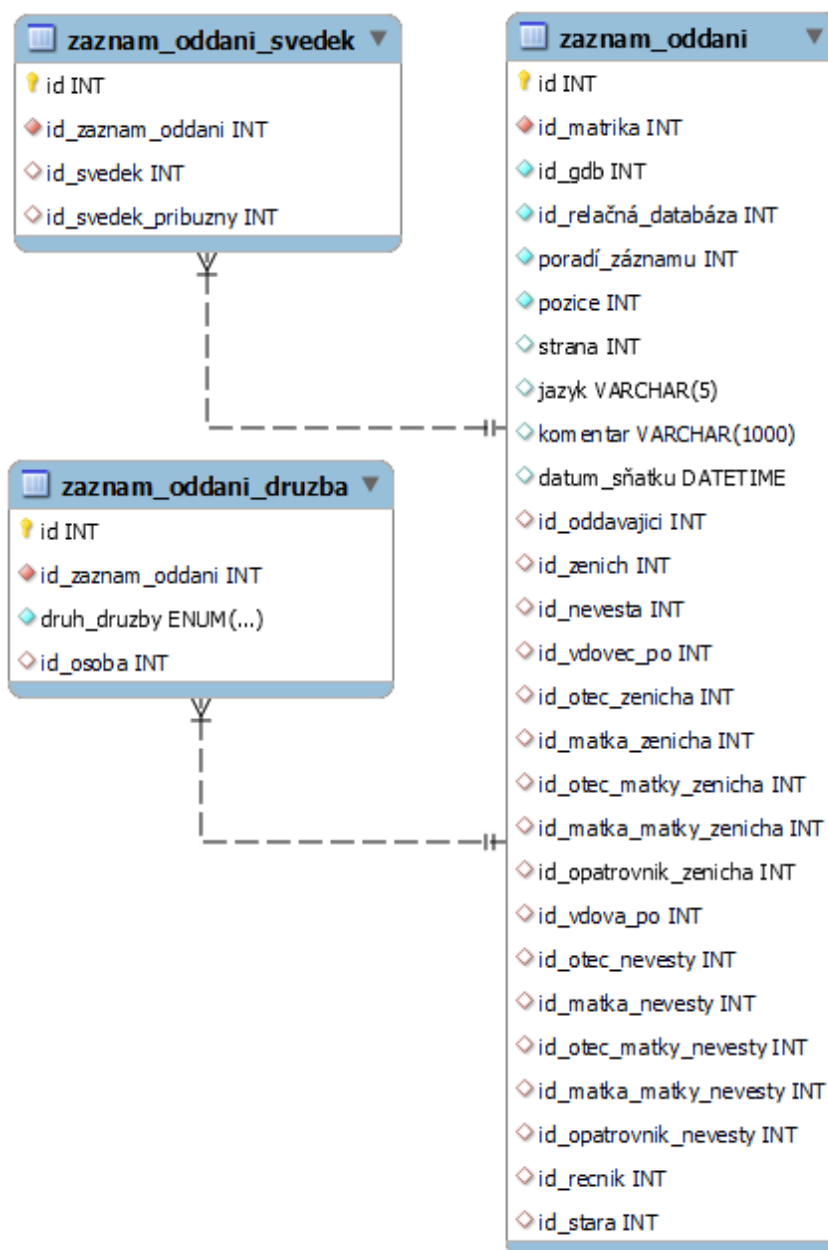
Obrázek 3.3: Návrh tabulek pro evidenci záznamů narození.

### Záznam oddání

V pořadí druhá tabulka pro ukládání matričních záznamů je tabulka evidující záznamy oddání 3.4. K zaznamenání data s datkem je navrhnutý sloupec datum\_s datku. Hlavními účastníky jsou zde v tšinou ženich s nevestou, kteří se zaznamenávají pomocí odkazů id\_zenich a id\_nevsta. Tak aby byli ženich s nevestou oddáni, je potřeba mít osobu, která by měla zaručit právoplatnost manželského s datku, pro tyto účely se zde eviduje oddávající ve sloupci id\_oddavajici. Jak ženich, tak i nevesta mohou být před s datkem ovdověni. Také na tyto případy je zde myšleno a pro evidenci zesnulé manželky ženicha slouží sloupec id\_vdovec\_po a pro zesnulého manžela nevesty sloupec id\_vdova\_po. Dalšími osobami, kteří se evidují u záznamu oddání, jsou rodinní příslušníci novomanželů, pro které jsou navrženy patřičné sloupce. Konkrétně se jedná o rodiče a prarodiče z matčiny strany jak

u ženicha, tak i u nevěsty. Jelikož nevěsta i ženich mohou mít píaze ní svého opatrovníka, tak i pro tyto píípady jsou navrhnuty vlastní sloupce k jejich evidenci. Pro opatrovníka ženicha sloupec i d\_opatrovník\_ženich a pro nevěstu i d\_opatrovník\_nevěsta. Posledními osobami evidovanými v tabulce zaznam\_oddani je e ník ve sloupci i d\_recník a stará ve sloupci i d\_stará.

K hlavní tabulce zaznam\_oddani jsou navrženy ještě další dvě pomocné tabulky, které zaručí, že u typu osob v těchto tabulkách může být k jednomu záznamu oddání více. Jako první bych zmínil tabulku pro evidenci svědků. Jedná se o tabulku zaznam\_oddani\_svedek kam je evidován každý svědek daného sátku, a to jak pro ženicha, tak i pro nevěstu ve sloupci i d\_svedek. Zároveň ke každému svědkovi může být evidován jeho pííbuzný, a to ve sloupci i d\_svedek\_příbuzny. Druhou pomocnou tabulkou je tabulka pro evidenci družeb a družek. Evidence záznamů v této tabulce je navrhnutá tak, že je zde sloupec i d\_osoba, sloužící jako odkaz na osobu, která je družbou nebo družkou daného matričního záznamu oddání. Pro rozeznání, jestli se jedná o družku nebo družbu slouží sloupec druh\_druzby s píednastavenými možnostmi výbírky s hodnotami 'druzba' a 'druzka'.



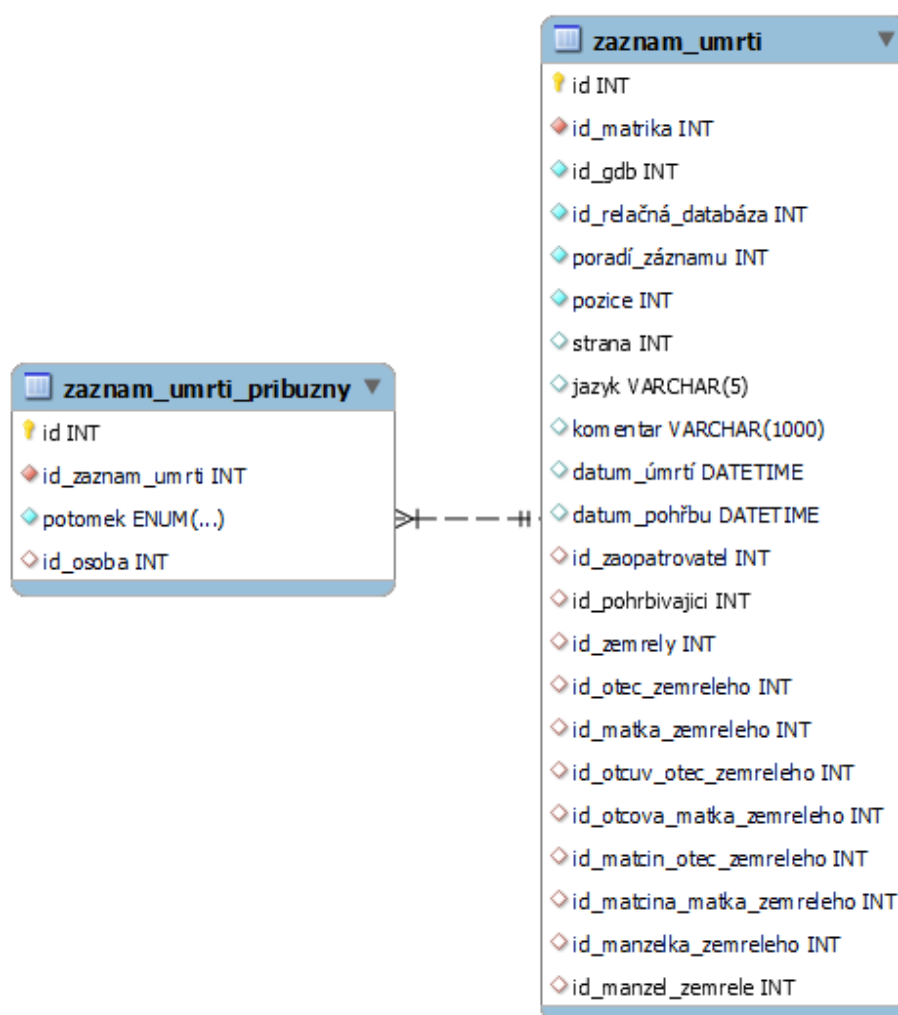
Obrázek 3.4: Návrh tabulek pro evidenci záznamů oddání.

### Záznam úmrtí

Poslední ze tří tabulek matričních záznamů je návrh tabulky `zaznam_umrti` 3.5, kam se budou evidovat, jak už název napovídá, jednotlivé záznamy ohledně úmrtí osob. Hlavní osoba, která je evidována matričním záznamem úmrtí je ta, která je pochována. Tato osoba se zde zaznamenává ve sloupci `id_zemrel` odkazem formou cizího klíče na tabulku osoby. Dalšími osobami evidovanými v popisovaném návrhu matričního záznamu je osoba zaopatřovatele ve sloupci `id_zaopatrovatel` a osoba pohřbívacího ve sloupci `id_pohrbivaji`. Jako u všech předchozích matričních záznamů, tak i zde se zaznamenává přesné datum

uskuteční dané události, zde konkrétní data úmrtí ve sloupci datum\_úmrtí. Další návrh evidence osob k matričnímu záznamu úmrtí je také podobný jako v předchozích matričních záznamech, a to u evidence rodinných příslušníků zesnulé osoby. Konkrétně se zde evidují přírodní rodiče a prarodiče zesnulé osoby. Zbývajících sloupci, navrženými v této tabulce, jsou id\_manželka\_zemreleho a id\_manzel\_zemrele k evidenci manželky nebo manžela zesnulé osoby.

U matričních záznamů úmrtí se také evidují přírodní potomci zesnulého nebo jeho příbuzní. Jelikož těchto osob může být k jednomu zesnulému více, je k této tabulce matričních záznamů navrhnutá pomocná tabulka evidující tyto osoby. Evidence je zde navrhnutá tak, že do sloupce id\_osoba se vloží odkaz na danou osobu a pomocí výběrového sloupce potomek se vybere z možností 'syn', 'dcera' a 'příbuzný', podle toho o jaký vztah k zesnulé osobě se jedná.



Obrázek 3.5: Návrh tabulek pro evidenci záznamů úmrtí.

### 3.1.3 Osoba

Po matrikách je pot eba mít samostatnou tabulku pro evidenci osob 3.6, která, jak už bylo zmín no, je propojená s tabulkami všech t í tabulek matri ních záznam . U každého záznamu osoby v této tabulce jsou evidovány následující data, pro které jsou vytvo eny podobn pojmenované sloupce. Datum narození, odhad narození od a odhad narození do ur ující interval, ve kterém se daná osoba mohla narodit, pokud její p esné datum narození není známé, datum k tu, datum viatiky, datum s atku, datum rozvodu, datum úmrtí a odhad úmrtí od a odhad úmrtí do, které ur ují v p ípad neznámého data úmrtí, interval ve kterém došlo k úmrtí dané osoby. M že se zdát, že sloupce pro evidenci data narození, s atku a úmrtí jsou zde zbyte né, práv proto že se tato data evidují p ímo v matri ních záznamech. Není tomu tak, tyto sloupce jsou zde navrženy zám rn , a to pro p ípad, že k dané osob není evidovaný n jaký ze t í matri ních záznam , ale máme k dispozici datum dané události.

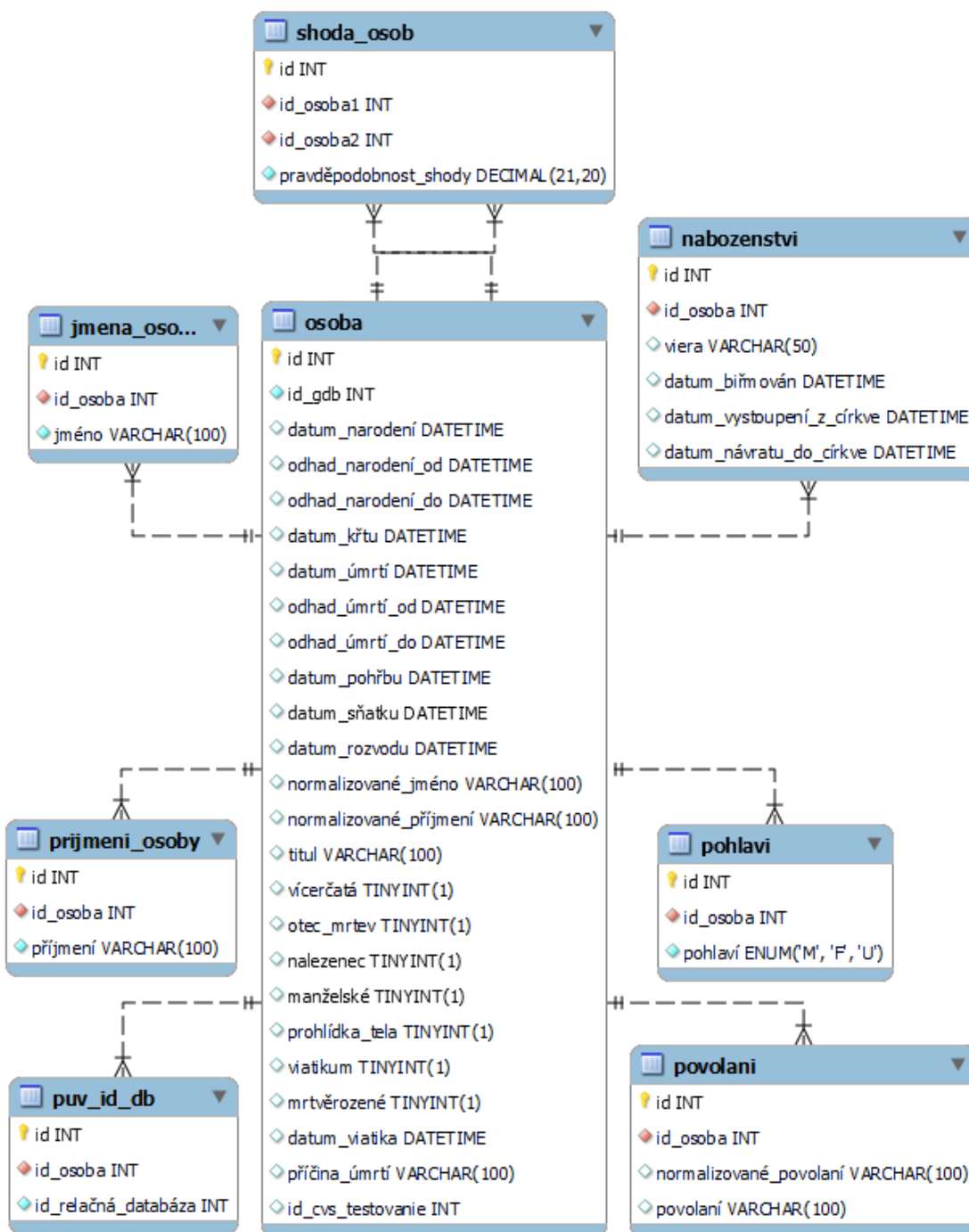
Jelikož jeden uzel pro osobu v grafové databázi m že zároveň obsahovat více jmen nebo p íjmení, jsou sou ástí tabulky sloupce normal i zované\_j méno pro zaznamenávání normalizovaného jména a normal i zované\_p íjmení pro zaznamenání normalizovaného p íjmení, které slouží práv pro jedno jméno a p íjmení, které jsou zároveň vyhodnoceny jako ty správné. Pro evidenci p ípadných titul je navrhnut sloupec titul , který bude obsahovat veškeré tituly dané osoby. V tabulce osoba jsou dále navrženy následující stavové sloupce. Sloupec ví cer atá ur ující, jestli daná osoba je sou ástí narozených vícer at, sloupec otec\_mrtev, který se eviduje v p ípad , má-li osoba zesnulého otce. nal ezenec je další ze skupiny stavových sloupc , tento je konkrétn navržen pro evidenci toho, jestli se nezná p vod dané osoby a byla nalezená. Stavový sloupec manžel ské je navržen pro evidenci, jestli se jedná nebo nejedná o nemanželské dít . Dalšími sloupci jsou prohl ídka\_tel a, ur ující, jestli u dané osoby byla provedena prohlídka t la, ví atí kum ur ující, jestli u dané osoby bylo viatikum provedeno a poslední stavový sloupec mrtv rozené ur ující, jesli se jedná o mrtvorozené dít .

V poslední ad je u osoby evidována p í ina úmrtí a jednozna ný identifikátor testovacího csv souboru související s testováním v p vodní grafové databázi Neo4j. K osobám je pot eba dále evidovat data ohledn pohlaví, náboženství, zmi ovaných více jmen nebo p íjmení, jednozna ný identifikátor z p vodní rela ní databáze, povolání a p ípadnou pravd podobnost shody s jinou osobou. K t mto zmín ným dat m je zapot ebí vytvo it další pomocné návazné tabulky, kdy všechny tyto pomocné tabulky jsou navrženy tak, že pro každou osobu m že být evidováno více záznam z jedné tabulky. K tomuto bude sloužit sloupec id\_osoba, která bude formou cizího klí e odkazovat na nad azenou tabulku osob.

#### **Jména, příjmení a původní ID osob**

Jak už jsem zmi oval v ásti pro popis nad azené tabulky osoby, tak pro jeden uzel osoby v grafové databázi m že existovat více záznam v p vodní rela ní databázi, tím pádem m že být sou ástí jednoho uzlu více jmen. Práv tato jména se zaznamenávají do tabulky jmena\_osoby 3.6. Na podobném principu jsou navrženy ještě další dv tabulky, a to tabulka pri jmeni \_osoby 3.6 pro evidenci p íjmení osoby a tabulka puv\_i d\_db 3.6 pro evidenci jednozna ných identifikátor osoby z p vodní rela ní databáze.





Obrázek 3.6: Návrh tabulky osoby a k ní návazných tabulek.

### Pohlaví osob

Další doplující tabulkou je tabulka pohlaví 3.6. Pro výběr pohlaví je zde navrhnutý výběrový sloupec s možnostmi výběru mužského, ženského nebo jiného pohlaví.

## Náboženství

Sou částí pomocných tabulek k tabulce osoby je tabulka náboženství 3.6, která je navržena pro evidenci názvu náboženství a případných dat s ním spojených. Mezi zmíněnými daty je datum bi mování, datum vystoupení z církve a datum návratu do církve.

## Povolání osob

Další pomocná tabulka, rozšiřující data k záznamům osob, je tabulka povolání 3.6, obsahující navržený sloupec povolání evidující originální název povolání ze kterého nemusí jít vždy poznat, o přesně které povolání se jedná, a to kvůli tomu, že je povolání zapsáno cizím jazykem nebo se jedná o nějaké historické povolání. Aby bylo povolání zapsáno i ve srozumitelné podobě, je zde navržený sloupec normalizované\_povolání pro evidenci normalizovaného názvu povolání.

## Shoda osob

Poslední ze skupiny pomocných tabulek k nadřazené tabulce osob je tabulka shoda\_osob 3.6. Na rozdíl od všech předchozích tabulek obsahuje dva sloupce odkazující na záznamy z tabulky osob, konkrétně se jedná o sloupce id\_osoba1 a id\_osoba2. Dalším sloupcem této tabulky je pravděpodobnost\_shody obsahující desetinné číslo s pravděpodobnostní shodou mezi odkázanými osobami.

### 3.1.4 Evidence místa

Poslední skupinou navržených tabulek jsou tabulky pro evidenci místa. Mezi skupinu těchto tabulek patří tabulky pro evidenci názvů ulic, čísel popisných, měst a další pomocná a vazební tabulka místa.

## Město

Jako první bych představil návrh tabulky mesto B.1, kde se kromě originálního názvu města ve sloupci nazov\_mesta eviduje i jeho normalizovaný název. Dalším navrženým sloupcem je id\_mesto, obsahující identifikátor záznamu města z primární grafové databáze.

K této tabulce bylo potřeba navrhnout pomocnou tabulku mesto\_gps B.1 pro evidenci GPS souřadnic. Kromě sloupce id\_mesto, který formou cizího klíče zaručuje propojení s nadřazenou tabulkou mesto, je zde sloupec gps\_souradnice evidující číselnou hodnotu dané souřadnice a sloupec poradí, zaručující pořadí zaevidované GPS souřadnice v souřadném systému GPS souřadnic.

## Ulice

Jako další bych představil návrh tabulky ulice B.1, kde se eviduje samotný název ulice a do sloupce id\_ulice identifikátor ulice z grafové databáze. V tabulce je dále navržený sloupec id\_mesto, který je formou cizího klíče a zajišťuje propojení s tabulkou mesto.

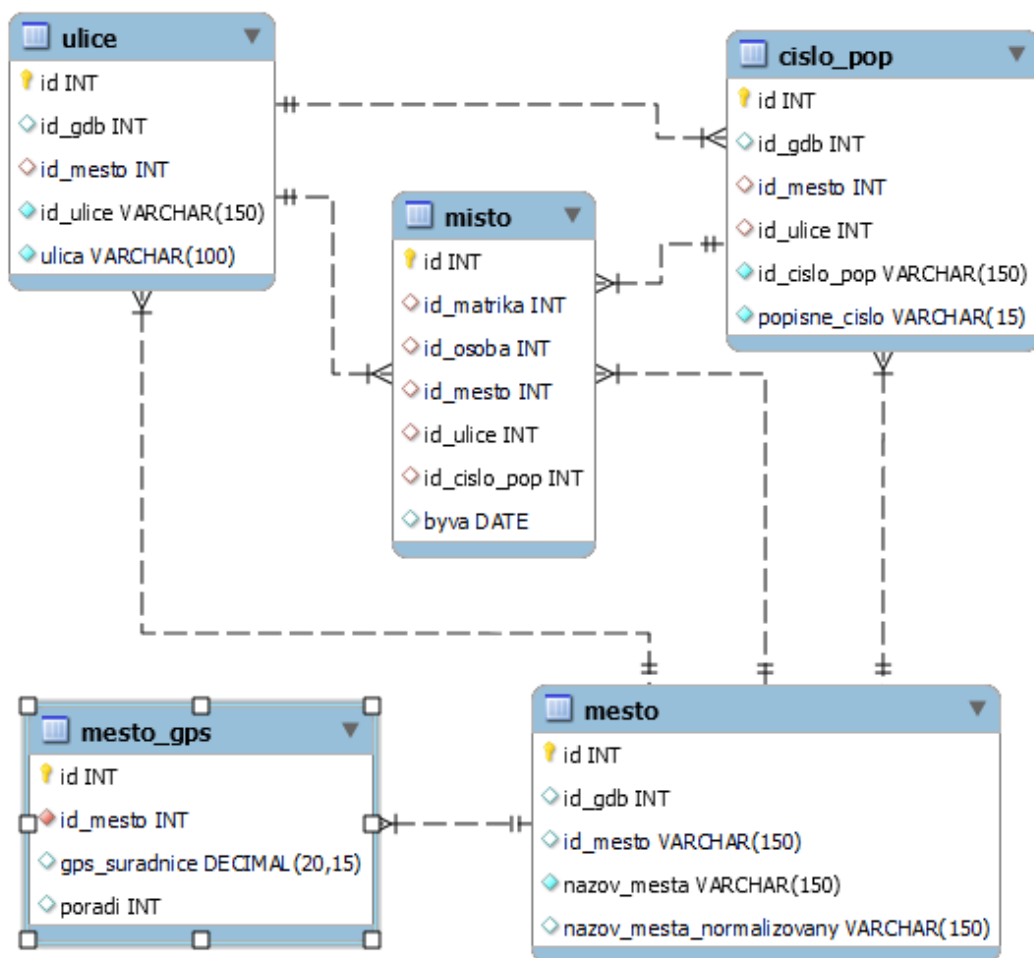
## Číslo popisné

Ve stylu všech předchozích tabulek této části je i návrh tabulky cislo\_pop B.1, kde se eviduje číslo popisné a do sloupce id\_cislo\_pop identifikátor čísla popisného z grafové databáze. Stejně jako v předchozí tabulce ulic, tak i zde je navržený sloupec id\_mesto,

kteřý je formou cizího klíče a zajistí propojení s tabulkou mesto. Podobně je navrhnutý sloupec id\_ulice, který zajistí propojení s tabulkou ulice.

## Místo

Poslední tabulkou v této části je vazební tabulka místo B.1, která slouží k propojení tabulek z této části míst 3.1.4 a tabulek matrik 3.1.1 a tabulek osob 3.1.3. Pro propojení s tabulkami mst slouží sloupec id\_mesto, s tabulkami ulic sloupec id\_ulice, s tabulkami čísel popisných sloupcem id\_cislo\_pop, s tabulkami osob sloupec id\_osoba a s tabulkami matrik sloupec id\_matrika. Posledním sloupcem je byva, kam se eviduje datum, které se zaznamenává při propojení tabulek míst s tabulkou osoby.



Obrázek 3.7: Návrh tabulky místa.

## 3.2 Návrh skript pro převod

Po vytvoření návrhu relační databáze, která bude sloužit ke kopii již existujících genealogických dat z grafové databáze, bude potřeba navrhnout samotný princip přenosu dat, a to jak z grafové databáze do relační, tak i obráceně. Právě možnému přenosu dat mezi těmito databázemi bych se chtěl v této části v novat a navrhnout možné řešení.

### 3.2.1 Převod z Neo4j do MySQL

Jako první bych představil popis návrhu převádní dat z grafové databáze Neo4j do relační MySQL databáze 3.8. Pro zajištění samotného převodu bude v první řadě potřeba se připojit k existující Neo4j databázi a připojit se k nově založené MySQL databázi i s vytvořenými tabulkami, které byly navrženy v předchozí části. Po úspěšném připojení k obou databázím půjde na řadu dotazování se a následný převod získaných dat. Dotazováním bude potřeba postupně procházet jednotlivé uzly každého typu a následně jejich data přenášet do konkrétních tabulek vytvořené relační databáze.

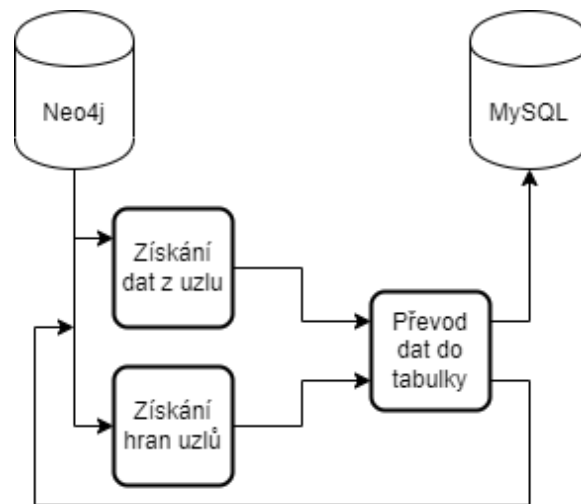
Začně se postupným dohledáváním dat z uzlů měst, ulic a čísel popisných. Tato data ze zmíněných uzlů se přenesou do patřičných tabulek v relační databázi. Data měst do tabulky mesto, ulice do tabulky ulice a číslo popisné do tabulky cislo\_pop.

Po dohledání a přenesutí všech míst se přejde na dotazování k uzlům matrik, odkud se budou přenášet data do tabulky matriky. Při každém zaznamenání matriky se zároveň zaeviduje i místo, přesněji se zde jedná pouze o název obce, případně více obcí. Tato evidence místa se provede způsobem, že se v případě neexistujícího názvu obce založí nový záznam s daným názvem obce v tabulce mesto. Následně se propojí obec s matrikou pomocí vazební tabulky m\_misto, kde se vyplní odkaz na záznam daného místa a dané matriky, zbylé sloupce odkazující na jiné tabulky zstanou v tomto případě prázdné. V případě, že obec, kde se daná matrika nachází, už existuje, nevytváří se další duplicitní záznam s tímto názvem obce, ale pouze se vytvoří záznam do vazební tabulky, která propojí dané záznamy obce s matrikou.

Po přenesení dat z uzlů týkajících se matrik se přejde k přesouvání dat ohledně jednotlivých osob. U osob se začne dotazováním k jejich uzlům a následný přesun dat do jejich tabulky osoba. Jelikož uzel osoby může obsahovat více jmen, příjmení nebo id z povodní relační databáze, provede se jejich přesun do navržených pomocných tabulek jmena\_osoby, prijmeni\_osoby nebo puv\_id\_db, které jsou návazné k tabulce osoba. Dalšími záznamy jsou data ohledně pohlaví, povolání a náboženství, budou přesouvány do pomocných tabulek pohlavi, povolani a nabozenstvi. Každý uzel osoby je hranou propojen s místem pobývání, které může být znázorněno městem, ulicí nebo číslem popisným. K propojení těchto záznamů slouží vazební tabulka m\_misto, kam se odkáže na místo pobývání osoby. V případě, že se jedná o ulici, vloží se odkaz do sloupce id\_ulice, pokud se jedná o číslo popisné, vloží se odkaz do sloupce id\_cislo\_pop, a pokud se jedná o název obce, vloží se odkaz do sloupce id\_mesto. Společně s jedním vyplněným sloupcem odkazující na nějaký druh adresy je potřeba do tabulky doplnit i odkaz na danou osobu do sloupce id\_osoba, čím se zajistí konečné propojení.

V této fázi přesunu, kdy je hotový přesun dat týkajících se matrik, dat všech osob, míst nalezení matrik a pobývání osob z grafové databáze, je potřeba převést data z matrikálních záznamů, které zaručují vztahy mezi jednotlivými osobami uložených v určitých matrikách. Jako první se provede přesun dat matrikálních záznamů narození. Jedná se o přeměnění matriky a osob k matrikálnímu záznamu, které jsou součástí daného záznamu a přenesení přesného data narození. Podobným způsobem se přesunou matrikální záznamy oddání a matrikální záznamy úmrtí, s tím rozdílem, že u oddání se eviduje pouze přesné datum oddání a u úmrtí přesné datum úmrtí.

Ve většině zmíněných případech se po dohledávání dotazy v grafové databázi, přenesou záznamy do relační databáze tak, že se vytvoří nový záznam. Jsou ale i případy, kdy záznam už existuje, jen je potřeba k danému záznamu přidat nebo upravit nějaká data. Pro tyto účely je potřeba si ukládat do proměnných určité identifikátory, pomocí kterých můžeme zaručit správně provedenou úpravu záznamů.



Obrázek 3.8: Návrh přesunu z Neo4j do MySQL databáze.

### 3.2.2 Převod z MySQL do Neo4j

Tato část není přímo součástí zadání mé bakalářské práce, ale po konzultaci s vedoucím mé práce jsme se domluvili, že ji tam přidám.

Princip tohoto převodu bude podobný jako u předchozího převodu 3.9. Stejně jako v předchozí části, tak i zde bude potřeba být v první fázi úsporně připojený k obou databázím, mezi kterými má dojít k přesunu dat. Na rozdíl od předchozí části se zde dohledávají data pomocí dotazů v MySQL databázi a ukládají pomocí dotazů Cypher do Neo4j databáze.

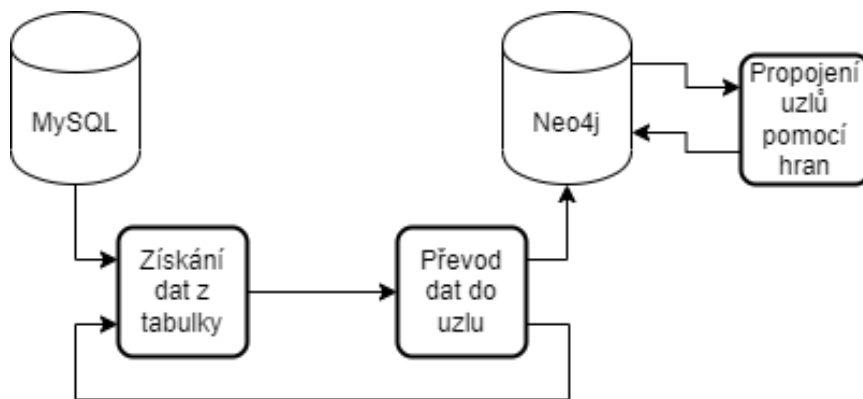
U tohoto převodu se začne převodem dat matrik, kdy se dotazem nad tabulkou matrik získají potřebná data, která se postupně přesunou do jednotlivých uzlů typu Matrika v grafové databázi.

Jakmile jsou provedené všechny záznamy matrik, začnou se postupně převádět data osob, a to z tabulky osoba do uzlu typu Osoba. Součástí tohoto převodu dat osob se provedou veškerá data z tabulky a zbylé vlastnosti se doplní až později.

Ještě než dojde k převodu samotných matričních záznamů, je potřeba převést data týkající se míst určení. Konkrétně se jedná o převod záznamů z tabulek mesto, ulice a číslo\_pop do uzlu typu Mesto, Ulice a Popisné\_ísto. S převodem samotných záznamů MySQL databáze bude potřeba ještě zaručit převod existujících vazeb mezi zmíněnými tabulkami místa určení a také převod vazeb mezi záznamy tabulek míst určení a záznamy v tabulce osoba.

Nyní se dostáváme k zásadnímu převodu matričních záznamů. Zde se pro všechny tři typy matričních záznamů provede první převod samotného záznamu dané tabulky do patřitého typu uzlu v grafové databázi. Poté dojde k převodu vazeb mezi matričním záznamem a osobami a matričním záznamem a matrikou, čím dostaneme výsledný vzhled schématu uzlu grafové databáze.

Jelikož ještě nejsou provedena data z pomocných tabulek k tabulce matrika a osoba. Proveďte se tento převod nyní



Obrázek 3.9: Návrh přenosu z MySQL do Neo4j databáze.

## Kapitola 4

# Implementace a testování

Po teoretické části návrhu bych zde chtěl představit i praktickou část samotné implementace. Jedná se o skripty sloužící pro práci jak s Neo4j, tak i s MySQL databází. Pro možnou práci s oběma databázemi zároveň jsou tyto skripty implementovány v programovacím jazyce Python. Přesněji se jedná o tyto různé skripty `dictonary.py` obsahující slovníky pomocí kterých dochází k propojování názvů vlastností a hran uzlů s názvy sloupců tabulek, `mysql Database.py` zahrnující obsluhu MySQL databáze, `graphDatabase.py` pro obsluhu grafové databáze Neo4j a `main.py`, který slouží pro propojení předchozích komunikací a tím zajišťující přesun dat mezi nimi.

Společně s těmito programy je potřeba mít ještě k dispozici soubor obsahující přihlašovací údaje k obou databázovým serverům. Tento soubor by měl obsahovat název `HOST-NAME`, číslo `PORTU`, uživatelské jméno, heslo a název databáze jak k MySQL, tak i k Neo4j serveru. V případě, že-li potřeba se připojit ke vzdálenému MySQL serveru jsou potřeba i přihlašovací údaje pro připojení se k SSH tunelu. V tomto případě musí externí soubor obsahovat i název `HOSTNAME`, `PORT`, uživatelské jméno a heslo pro SSH připojení. V následujících podčástech bych chtěl podrobněji představit popis implementace jednotlivých zmínovaných programů.

### 4.1 Implementace MySQL

Hlavní částí této práce bylo vytvořit samotnou MySQL databázi a především k tomuto slouží program `mysql Database.py`. V tomto programu je importována knihovna `pymysql`, která slouží pro možnost připojení se k MySQL serveru a implementaci MySQL dotazů v jazyce Python a knihovna `sshtunnel` pro zajištění SSH tunelu k případnému připojení ke vzdálenému MySQL serveru. Tento program je implementován pomocí vytvořené třídy `MysqlDB`.

Na začátku tohoto programu dojde ke čtení z externího souboru obsahující přihlašovací údaje. Tyto přihlašovací údaje se rozdělí do slovníku odkud je možné je následně používat a přihlásit se pomocí nich k databázovému serveru.

Jak už bylo zmíněno, tak pro případné připojení ke vzdálenému MySQL serveru je potřeba se nejdříve připojit k SSH tunelu, k tomuto se použije sekvence následujících příkazů :

```
tunnel = SSHTunnelForwarder((ssh_hostname, ssh_port),
                             ssh_username=ssh_username,
                             ssh_password=ssh_password,
                             remote_bind_address=(sql_hostname, sql_port)
                             )
tunnel.start()
```

Po úspěšném připojení k SSH tunelu nebo v případě připojování se k lokálnímu serveru a tím pádem nevyužití SSH tunelu dojde pomocí příkazu

```
pymysql.connect(hostname, port, user, password, database)
```

k připojení k databázi MySQL serveru a následně dojde k vytvoření kurzoru pomocí příkazu `cursor = connect.cursor()`, díky kterému jsou prováděny veškeré SQL dotazy.

Program následně obsahuje funkce, které pomocí MySQL dotazů mažou a vytvářejí jednotlivé tabulky. Zmínované mazání tabulek se provádí pokaždé před vytvářením tabulek, tak aby se v případě existence mohli tabulky korektně znovu založit. Protože se rovnou při tvorbě tabulek vytvářejí propojení tabulek pomocí cizích klíčů, musí se tabulky vytvářet v určitém pořadí, a to nejdříve tabulky `mesto`, `mesto_gps`, `ulice` a `ci_slopop`.

Po vytvoření tabulek sloužících k evidenci místa se vytváří tabulka `osoba`, která eviduje informace o osobách z genealogické databáze matrik. Jakmile je vytvořena tabulka `osoba`, přejde se na tvorbu tabulek, které jsou na tuto tabulku návazné. Jedná se o tabulku `povolani`, přiznávající povolání k dané osobě, tabulka `nabozenstvi`, evidující data týkající se náboženské víry, tabulka `pohlavi` určující pohlaví osoby. Jelikož jeden uzel osoby v databázi Neo4j může obsahovat více jmen, příjmení a id odkazující na původní záznam z relační databáze matrik, vytvoří se další tři tabulky pro evidenci právě těchto záznamů. Konkrétní tabulky `jmena_osoby`, `prijmeni_osoby` a `puv_id_db`. Poslední vytvořenou pomocnou tabulkou k tabulce `osoba`, je tabulka `shoda_osob`, kam se zaznamenává pravděpodobnostní shoda mezi osobami.

Jakmile jsou založené tabulky pro ukládání dat k osobám, začnou se zakládat tabulky týkající se samotných matrik. V první řadě se vytvoří samotná tabulka `matrika`, kam se budou ukládat data k jednotlivým matrikám. Následně se založí pomocná tabulka `jazyk` evidující jazyk zápisu k dané matrice. Právě teď se při zakládání tabulek nacházím v bodě, kdy máme založené tabulky pro evidenci míst, osob a jejich pomocných tabulek a matrik a její pomocné tabulky. To znamená, že je možné založit vazební tabulku mezi sloužícími k propojení tabulek míst s matrikou nebo osobou.

Jako poslední se vytvoří jedny z nejdůležitějších tabulek, a to tabulky pro jednotlivé záznamy matrik, které zahrnují nejen propojení osob s matričními záznamy, ale i nastavení vazeb mezi jednotlivými osobami. Jedná se o tabulky `matricnich_zaznam`, `zaznam_narozeni`, `zaznam_oddani` a `zaznam_umrti` a k nim pomocných tabulek `zaznam_narozeni_kmotr`, `zaznam_oddani_svedek`, `zaznam_oddani_druzba` a `zaznam_umrti_pribuzny`.



Po funkci zajišťující tvorbu tabulek se v tomto programu nachází funkce zajišťující volání MySQL dotazů. Jako první je funkce pro vytváření záznamů v MySQL tabulce, která má tři vstupní argumenty, a to příznak ignore, název tabulky, seznam s názvy sloupců a seznam s daty nově vytvořeného záznamu v tabulce. První argument této funkce v podobě pravdivostního příznaku slouží k tomu, jestli se bude volat dotaz INSERT INTO nebo INSERT IGNORE INTO, který zamezí vkládání duplicitních záznamů. Znamená to, že při volání funkce

```
insert_table_all (False, table_name, (column1, column2, column3),
                 (value1, value2, value3))
```

se vytvoří SQL dotaz

```
INSERT INTO table_name (column1, column2, column3)
VALUES(value1, value2, value3),
```

který v dané MySQL tabulce založí nový záznam.

Další funkce pro tvorbu MySQL dotazů slouží pro výpis dat z tabulky, přesněji se jedná o tvorbu SELECT dotazu. Zde se funkce volá s argumenty název tabulky, seznam s názvy sloupců a příznakem zec obsahující podmínku hledaných záznamů, hodnota hledané podmínky, název tabulky k možnému propojení, název sloupce, pomocí kterého dojde k propojení, jedné tabulky a druhé tabulky a typ propojení. Podmínka se zohledňuje pouze pokud je zadána, to stejné platí i pro případ propojení tabulek. Uvedl bych zde jednoduchý případ využití všech vstupních argumentů funkce.

```
select_table(table1, [table1.column1, table1.column2, table2.column3],
             table1.column2, value, table2, column2, column1)
```

Po zavolání funkce v této podobě se vygeneruje následující SQL dotaz.

```
SELECT table_name1.column1, table_name1.column2, table_name2.column3
FROM table_name1
JOIN table_name2 ON table_name2.column2 = table_name1.column1
WHERE table_name1.column2 = value,
```

Naopak tato funkce může sloužit i jen pro zobrazení hodnot určitých sloupců pro všechny záznamy jedné tabulky a v tomto případě by se volala následující funkce.

```
select_table(table_name, (column1, column2))
```

Následně se pomocí této funkce vygeneruje SQL dotaz.

```
SELECT column1, column2
FROM table_name
```

Jelikož v některých případech není možné poskládat potěbný SELECT dotaz. Jsou zde implementované další složitější dotazy, které se využívají zejména při testování a také naopak i jeden jednodušší dotaz, který slouží pouze pro výpis dané hodnoty, zde se konkrétně jedná o `only_select(column)`, kde právě argument `column` je daná hodnota. Výsledkem je vygenerovaný jednoduchý SQL dotaz `SELECT column`.

Kromě zakládání nových záznamů do tabulek a výpisu z nich je potřeba ještě mít k dispozici funkci pro úpravu záznamů v tabulkách. Pro tento účel slouží funkce `update_table`

obsahující podobné argumenty jako zmíněná funkce `insert_table`, s tím rozdílem, že argumenty pro podmínku jsou zde povinné a argumenty k propojení více tabulek zde nejsou k dispozici. Zmínil bych zde jednoduchou ukázkou volání této funkce

```
update_table(table, (column1, column2), (value1, value2), column3, value3)
```

Výsledkem této funkce bude vyvolání následujícího SQL dotazu.

```
UPDATE table_name1
SET column1 = value1, column2 = value2
WHERE column3 = value3,
```

Mezi poslední funkce této třídy patří funkce sloužící k provedení SQL dotazu a funkce pro odpojení kurzoru a ukončení spojení s MySQL databází.

Kromě posledních zmíněných funkcí se při úspěšném ukončení ostatních funkcí provede stručný zápis výsledku do textového souboru. Zároveň u všech těchto funkcí dochází k zachycování výjimek a v případě zachycení výjimky se do textového souboru zaznamená chybová hláška, uzavře se spojení s textovým souborem a ukončí se spojení pomocí poslední zmíněné funkce a následně se zobrazí chybová hláška i v terminálu, kde došlo ke spuštění.

## 4.2 Implementace Neo4j

Po seznámení se s programem obsluhující MySQL databázi, je potřeba ještě popsat funkci programu pro komunikaci s grafovou databází Neo4j. Jsou zde importovány následující knihovny. Knihovna `neo4j`, díky které se umožní práce s Neo4j databází, knihovna `datetime` pro práci s daty a je zde importovaný program `datetime.py` [4.3](#). Také tento program je implementovaný pomocí vytvořené třídy `GraphDatabase`.

Stejně jako u předchozího programu, tak i zde dojde na začátku k otevření souboru s přihlašovacími údaji, které využije následující funkce pro připojení ke grafové databázi Neo4j.

```
neo4j.GraphDatabase.driver(uri=uri,
                           auth=(username, password),
                           database=database)
```

Argument `uri` u zmíněné funkce připojení je URI adresa ke grafové databázi. Tato URI adresa je poskládána ještě před zavoláním funkce, a to z `HOSTNAME` a čísla portu z externího souboru s přihlašovacími údaji určené pro Neo4j. `Auth` slouží pro zadání přihlašovacích údajů a `database` pro výběr databáze, ke které se chceme připojit. Tento program dále obsahuje funkce pro manipulaci s daty grafové databáze Neo4j.

První ze zmíněných funkcí pro manipulaci dat je funkce sloužící pro zobrazení všech vlastností daného uzlu.

```
get_all_record(node_name)
```

Výsledkem této funkce bude vyvolání následujícího cypher dotazu.

```
MATCH (r:node_name) RETURN id(r), r
```

Zbylé funkce zobrazující data z grafové databáze jsou podobného principu, proto bych zde podrobněji zmínil pouze jednu a zbylé jen stručně popsal. Podrobněji bych popsal funkci `get_person_by_id`, která obsahuje 3 vstupní argumenty, a to argument `node` pro zadání názvu uzlu z grafové databáze a argumenty `id_person` a `r` sloužící jako vyhledávací podmínky dotazu. Výsledkem této funkce je výpis jednoznačného identifikátoru zadaného uzlu, který v určitém vztahu s daným uzlem osoby. Uvedu zde příklad volání této funkce.

```
get_person_by_id(node_name, value1, value2)
```

A následně výsledek vygenerovaného cypher dotazu.

```
MATCH (o:Osoba)-[r]->(z:node_name)
WHERE ID(o) = value1 and type(r) = 'value2'
RETURN ID(z)
```

Jak už jsem zmínil, další funkce pro výpis dat z grafové databáze si jsou hodně podobné. Mezi tyto funkce patří `get_record_by_id(node, id_record)`, kde vstupními argumenty je název uzlu a hledaná hodnota tohoto uzlu, ke kterému vede hrana z jiného uzlu `Osoba`. Tato funkce následně vrácí hodnoty všech vlastností uzlu `Osoba` a typ a datum hrany, kterou jsou tyto uzly spojeny.

K získání jednoznačného identifikátoru uzlu `Matrika` slouží další implementovaná funkce `get_registry_by_id(node, id_record)`, kde argument `node` určuje název uzlu který je propojený hranou typu `JE_V` s uzlem `Matrika` a druhý argument `id_record` značí hodnotu jednoznačného identifikátoru uzlu daného argumentem `node`.

Podobná funkce `get_lives_by_id(node, id_person)` se používá pro získání jednoznačného uzlu zadaným argumentem `node`, ke kterému vede hrana typu `BÝVÁ` z uzlu `Osoba` a právě druhý vstupní argument udává jednoznačný identifikátor uzlu `Osoba`.

Poslední funkce, pro výpis dat z grafové databáze, je funkce `get_persons_by_id()`, díky které se získají identifikátory obou uzlů, které jsou zde uzly `Osoba` a jsou propojeny hranou typu `POTENCIONALNÍ_ZHODA`, právě vlastnost tohoto uzlu `pravdopodobnost_shody` je totiž získaná hodnota z této funkce.

Jako poslední funkce spadající do skupiny získávání dat z grafové databáze je funkce `get_by_id(node1, node2, id_a)`. První dva argumenty této funkce značí první a druhý uzel, který jsou spolu propojeny hranou a argument `id_a` obsahuje hodnotu jednoznačného identifikátoru prvního uzlu. Výsledkem této funkce je hodnota jednoznačného identifikátoru druhého uzlu zadaného vstupním argumentem.

Od funkcí pomocí kterých se získávají data z grafové databáze Neo4j jsme se dostali k funkcím zapisujícím do grafové databáze. Nachází se zde funkce pro tvorbu uzlů se zápisem jejich vlastností a funkce pro tvorbu hran s příslušným zápisem jejich vlastností.

Začal bych u funkcí pro tvorbu uzlů, které vždy obsahují argument `properties` a `values`. Druhý argument obsahuje hodnoty ve formě matice, z tohoto důvodu se vždy na začátku každé funkce začne pomocí cyklu procházet touto maticí a poté je k dispozici pro každý krok cyklu pouze jedno pole hodnot. Následně se z prvního argumentu a získaného pole

hodnot vytvoří slovník, díky kterému se poté zapisují hodnoty k určitým vlastnostem podle klíče tohoto slovníku. Kromě všech vlastností uzlu, které jsou součástí původní databáze Neo4j, jsou součástí nové Neo4j databáze dvě další vlastnosti a to `id_gdb` určující jednoznačný identifikátor uzlu z původní Neo4j databáze a `id_mysql` obsahující jednoznačný identifikátor záznamu z tabulky MySQL databáze.

První z funkcí pro zakládání uzlu je `create_registry_nodes(properties, values)`, pomocí které se zakládají nové uzly Matrika se všemi vlastnostmi, tak jako v původní Neo4j databázi a to i se zadanými hodnotami. Podobná funkce je, která slouží pro zakládání uzlu Osoba a jeho vlastností je funkce `create_persons_nodes(properties, values)`. Dalšími funkcemi pro zakládání uzlu a jejich vlastností je funkce pro zakládání uzlu určující místo `create_place_nodes(node, properties, values)`, kde se argumentem `node` určuje, jestli se jedná o tvorbu uzlu Mesto, Ulice a Popisné sídlo. K zakládání jednotlivých uzlu matričních záznam slouží funkce `create_record_nodes(node, properties, values)`, kdy argument `node` určuje, o jaký druh matričního záznamu se jedná. Poslední funkce týkající se zápisu do uzlu je funkce `add_node_property(node, properties, values)`, která do uzlu zadaného pomocí argumentu `node` přidá do vlastnosti z druhé hodnoty seznamu `properties`, který je zadán jako argument funkce jeho hodnotu z argumentu `values`.

Druhou skupinou funkcí sloužících k zakládání záznamů v grafové databázi jsou funkce vytvářející hrany mezi jednotlivými uzly. Mezi tyto funkce patří například vytváření hran JE\_V pomocí `create_place_rel(node, properties, values)` a to mezi uzly míst Mesto, Ulice a Popisné sídlo, kdy výsledkem je jedna z následujících možností vytvořené hrany

```
Ulice-[r: JE_V]->Mesto,  
Popisné sídlo-[r: JE_V]->Ulice,  
Popisné sídlo-[r: JE_V]->Mesto.
```

K založení hrany BÝVÁ mezi uzlem Osoba a jedním z těchto typů uzlu určující místo slouží funkce `create_address_bond(rel_names, rel_values)`, zároveň se i zapíše vlastnost `date` této hrany, obsahující datum, kdy osoba pobývala na daném místě. K založení hran mezi uzly Záznam\_o\_ket, Záznam\_o\_svatbe nebo Záznam\_o\_úmrtí a uzlem Osoba je vytvořena funkce `create_record_bond(node, rel_names, rel_values, rel_property=None)`, která zároveň do vytvořené hrany vloží vlastnost `date`, určující datum konání dané události. Tyto zmíněné hrany se tvoří v cyklu tak, aby se každý jeden matriční záznam zároveň propojil se všemi osobami, které do něj patří. Před zahájením tohoto cyklu se ale ještě propojí matriční záznam hranou JE\_V k uzlu Matrika.

Jelikož kromě spojení mezi osobou a matričním záznamem je potřeba mít hrany i mezi jednotlivými uzly Osoba, bylo potřeba implementovat další funkci, a to první `create_persons_bond(rel_names, rel_values)`, kde se pouze vytvoří cyklus pro zobrazování jednotlivých polí hodnot, pomocí kterých se ještě s názvy sloupců vytvoří slovník a následně k samotnému vytváření hran mezi jednotlivými osobami se zavolá další funkce `create_persons_bond1(bond_dict, dictionaries, rel=None)`. Poslední zmíněná funkce se volá celkem pětkrát, kdy první argument `bond_dict` odkazuje na jednotlivé slovníky z programu `dictionaries.py` 4.3, druhý argument `dictionaries` obsahuje vždy slovník s hodnotami pro vkládání vytvořených v rodičovské funkci a poslední argument `rel` obsahuje případný název hrany mezi uzly. Poslední funkcí, která slouží k tvorbě záznamů v grafové databázi, konkrétně k tvorbě hran POTENCIONALNÍ\_ZHODA mezi uzly Osoba, je funkce

`create_person_probability(rel_names, rel_values)`, která zároveň vytvoří vlastnost `pravd_podobnost_shody` dané hrany.

Tak aby se mohly všechny Neo4j dotazy vygenerované v předchozích funkcích provést, je potřeba mít implementovanou následující funkci `run_graph_query()`. Tato funkce v podstatě pouze provede vygenerovaný dotaz a v případě chyby ukončí program zachycenou chybou.

### 4.3 Pomocný slovník

Tato část je implementovaná pro jednodušší práci při určování názvu hran pro jednotlivé sloupce MySQL databáze.

Slovník `table_col` udává název hrany mezi osobou a matřím záznamem v grafové databázi pro jednotlivé sloupce z tabulky `mesto` v relační databázi. Slovník `is_in_urban` je pro které spojení sloupců tabulek `mesto`, `ulice` a `ci_sloupop` z MySQL databáze se vytváří hrana `JE_V`, slovník `is_in_urban` je propojení sloupců z tabulek zmíněných u předchozího slovníku a tabulkou `osoba` hranou `BÝVÁ`. Zbývající slovníky slouží ke zjištění hran mezi sloupci tabulek `osoba` a to tak, že slovník `father` udává název hrany pro všechny jeho prvky `JE_OTEC`, slovník `mother` hranu s názvem `JE_MATKA`, slovník `married` název hrany `JSOU_MANŽELÉ`, slovník `related` název hrany `JE_PRÍBUZNÝ` a u posledního nezmíněného slovníku `other` je název hrany určen klíčem jednotlivých záznamů slovníku.

### 4.4 Převod dat

V předchozích částech implementace jsme si představili zvlášť komunikaci s MySQL databází a Neo4j databází. V této části bychom chtěli představit propojení zmíněných databází a zajištění převodu dat z Neo4j databáze do MySQL databáze a naopak. V první řadě se otevře pro zápis textový soubor určený k případnému logování běhu programu. Dále je potřeba se připojit jak k MySQL databázi, tak i k Neo4j databázi. Po úspěšném připojení se zavolá samotná funkce `main()`, díky které se zahájí samotný běh programu mezi databázemi. V případě, chceme-li přesouvat data z Neo4j databáze do MySQL, je potřeba zavolat funkci `main(neo4j_to_mysql=True)`, v opačném případě `main(neo4j_to_mysql=False)`.

#### 4.4.1 Neo4j do Mysql

Než dojde k samotnému převodu dat mezi databázemi, je potřeba založit tabulky v relační databázi MySQL, k tomu se využije funkce `create_all_table()`, která je součástí programu `mysql Database.py` 4.1, poté dojde k samotnému převodu dat a to v následujícím podkapitole.

Jako první se provede převod dat z uzlu `Mesto` do tabulky `mesto` a to z důvodu, že se jedná o data, které nenavazují na žádná jiná data. Na místě přímě navazují ulice a nakonec na ulice směřují popisná, to znamená, že nejdříve dojde k převodu dat z uzlu `Ulice` do tabulky `ulice` a následně propojení s tabulkou `mesto`. Jak už bylo zmíněno, další převod proběhne z uzlu `Popisná_ísla` do tabulky `ci_sloupop`, kdy se při převodu doplní i případné vazby s předchozími tabulkami určenými jaké konkrétní místo.

Další částí p evodu se týká dat matrik, kde se p evádí data z uzlů. Matrika do tabulky matrika. Krom vytvoření záznamů ve zmínované tabulce k evidenci matrik je nutné převést data ohledně jazyka a obcí výskytu do příslušných tabulek a propojit s hlavní matriční tabulkou.

Nyní, po p evodu dat matrik, dojde postupně k p evodu dat osob a všech tří druhů matričních záznamů a právě tyto jednotlivé p evody bych zde popsal podrobněji.

Jako první se provede p evod dat z uzlů Osoba do tabulky osoba. Princip je takový, že se nejdříve pomocí funkce z třídy graphDatabase 4.2 získají veškerá data ze všech uzlů Osoba. Z těchto získaných dat se pomocí vytvořených funkcí, zajišťujících potřebnou konverzi, za nů postupně, pomocí funkce z třídy MysqlDB 4.1, vytvářejí jednotlivé záznamy v tabulce osoba. Po každém vytvoření záznamu se do určené proměnné uloží jeho jednoznačný identifikátor, který bude sloužit k propojení záznamů z pomocných tabulek osoby s tabulkou osoba. Za nů se s p evodem dat do pomocných tabulek pohlaví, jména\_osoby, příjmení\_osoby, puv\_i\_d\_db, náboženství a povolání a to pomocí vytvořené funkce fill\_exp\_table, která zajistí správné vytvoření všech potřebných záznamů v zadané tabulce. Další nádu půjde volání funkce person\_links, která zaručí p enesení propojení hranou z grafové databáze mezi uzly Osoba a Mesto, Ulice nebo Popisné sídlo. Získaná data těchto propojení se zaznamenávají do vazební tabulky miasto, která zaručuje propojení tabulky osoba s tabulkami mesto, ulice a ci\_slo\_pop. Po vytvoření p evodu všech uzlů typu Osoba do tabulek relační databáze sloužící k evidenci osob, dojde k p enesení hran z grafové databáze týkající se pravděpodobnostní shody osob do pomocné tabulky shoda\_osob, návazně k hlavní tabulce osoba.

Dostáváme se k nejdležitější fázi p evodu, a to k p evodu samotných matričních záznamů a jejich příslušných hran. Na začátku p evodu u všech tří druhů matričních záznamů se získají hodnoty vlastností pro všechny uzly patřící k danému druhu matričního záznamu. Provede se konverze hodnot určitých vlastností, tak aby se tyto hodnoty mohli vkládat do tabulek relační databáze. Následně se pomocí funkcí z třídy graphDatabase 4.2 a MysqlDB 4.1 založí záznam pro daný matriční záznam i s propojením k nadřazené tabulce matrika. Jakmile je vytvořený záznam se základními hodnotami, za nů se přidávají jednotlivé hodnoty pro sloupce, které slouží k vytvoření vazeb mezi matričním záznamem a tabulkou osoba. Společně s vytvářením těchto vazeb se zaeviduje i datum uskutečnění dané události. Jelikož některých typů osob návazných na daný matriční záznam může být více, vytváří se do pomocných tabulek k danému matričnímu záznamu. Jako poslední v řadě se p evodou p íbuzné osoby, a to v pomocné tabulce zaznam\_narozeni\_kmotr p íbuzný pro kmotra a v pomocné tabulce zaznam\_oddani\_svedek p íbuzný pro svdka.

#### 4.4.2 Mysql do Neo4j

Nad rámec zadání byl implementovaný i p evod dat z navrhnuté MySQL databáze do původní Neo4j databáze. K provedení tohoto p evodu se znovu použijí obě dříve zmíněné třídy graphDatabase 4.2 i MysqlDB 4.1. Funkcionalita tohoto p evodu je poměrně jednoduchá a to, že se vždy z MySQL databáze pomocí určitého SELECT dotazu zobrazí potřebná data, která se pomocí Cypher dotazů p evodou do Neo4j databáze.

Na rozdíl od předchozího provedení, se zde začíná s provedením dat matrik, a to již zmíněným způsobem. Pomocí již implementované funkce se poskládá Cypher dotaz k vytvoření uzlu Matřice i se všemi vlastnostmi, ale zatím se vyplní pouze ty jejich hodnoty, které je možné získat také z už implementované funkce pro získávání dat z MySQL databáze, zde z tabulky matřice.

Jakmile jsou založeny všechny uzly matrik, přijde na tvorbu uzlu osoba a to podobným způsobem jako právě předšlé uzly matrik.

Nyní jsou vytvořeny všechny uzly typu Matřice a Osoba a může se začít s tvorbou uzlu k určení místa. Pro všechny typy uzlu patřící do této skupiny probíhá procedura, že se nejdříve založí všechny uzly i se zadanými vlastnostmi. Následně dojde k případnému propojení daných uzlu s uzly typu Osoba hranou BÝVÁ a v případě existence se k hraně přidá její vlastnost date. Poté u uzlu typu Mesto dojde k doplnění vlastnosti gps\_suradnice z pomocné MySQL tabulky mesto\_gps. U zbylých typech uzlu Ulice a Popisné síli k doplnění uzlu nedochází, ale dochází zde k vytváření hran JE\_V propojující buď uzly typu Ulice s Mesto nebo Popisné síli s Mesto a nebo Popisné síli s Ulice.

Po vytvoření všech předchozích základních typech uzlu, bude potřeba založit uzly jednotlivých matričních záznamů, následně jejich propojení pomocí hran s osobami a poté navazující propojení hran mezi osobami, které jsou spolu ve vztahu. U všech těchto druhů matričních záznamů se zmíněným provedením provádí stejným způsobem pomocí příslušných implementovaných funkcí jak ze třídy graphDatabase 4.2, tak i ze třídy MySQL DB 4.1. Následně u každého typu matričních záznamů dojde k přesunu dat z pomocných tabulek, konkrétně doplnění kmotrů a jejich příbuzných u typu uzlu Záznam\_o\_kte a Osoba, svdk, jejich příbuzných, družby a družek u typu uzlu Záznam\_o\_svatbe a Osoba a potomků a příbuzných u typu uzlu Záznam\_o\_úmrťi a Osoba. U všech těchto provedením z pomocných tabulek matričních záznamů se provede i doplnění hran mezi danými uzly typu Osoba.

V další fázi provedení proběhne doplnění chybějících vlastností a to doplnění vlastností Obce a Jazyky u uzlu Matřice a doplnění vlastností týkajících se povolání, náboženství, vlastností s názvem pohlaví, jméno, příjmení a id\_rela\_ná\_databáza u uzlu Osoba.

Nyní už je provedení skoro dokončené, zbývá jen převést hrany mezi jednotlivými osobami, konkrétně doplnění hran POTENCIONALNÍ\_ZHODA a doplnění její případné vlastnosti s názvem pravdopodobnost\_shody.

Tímto popsaným postupem je možné z vytvořené MySQL databáze, vytvořené pomocí předšlého provedení, vytvořit kopii provedení Neo4j databáze a je možné si tímto způsobem otestovat správnost nově vytvořené relační databáze. Každopádně tento provedení není primárně určený pro testování, k testování jsou vytvořené samostatné skripty, o kterých více v další části 4.5.

## 4.5 Testování

Pro výsledné naimplementované řešení jsem společně s vedoucím mé práce vymyslel následující testování. Otestovat rychlost provedení a to jak z Neo4j databáze do MySQL databáze, tak i naopak a poté tyto výsledky porovnat a zhodnotit. Hlavně ale bude potřeba otesto-

vat správnost samotného p evodu, a to nejen kontrola po tu p evedených záznam , ale i kontrola jejich hodnot. Jelikož hlavním úkolem mé práce byl p evod z Neo4j do MySQL, testoval jsem správnost pouze tohoto p evodu a také, jak už jsem zmi oval, i zp tný p e- vod se m že použít jako forma testování. Všechny testování jsem provád l s daty p vodní Neo4j databáze, kde se nachází 11374 uzl osmi r zných typ (Matri ka, Osoba, Mesto, Ul i ca, Popi sné\_ í slo, Záznam\_o\_k te, Záznam\_o\_svatbe a Záznam\_o\_úmrti) a 23696 hran, které navzájem tyto uzly propojují.

#### 4.5.1 Rychlost p evodu

Co se týká testování rychlosti p evodu, zde jsem si nasimuloval celkem 40 p evod . 10 z p vodní Neo4j do nov navržené MySQL databáze, a 10 p evod z navržené MySQL databáze do nové Neo4j databáze. Toto testování jsem provád l jak pro lokální databázi MySQL, tak i pro vzdálenou MySQL databázi na zp ístupn ěm školním serveru Radegast a databázi Neo4j jsem použil pouze lokáln , protože vzdálená databáze nebyla zprovozn na. U každého nasimulovaného p evodu jsem si zaznamenával jeho as a z t chto as jsem si následn vypo ítal pr m rnou dobu p evodu, zaznamenal nejdelší a nejkratší dobu trvání p evodu a výsledky zapsal do tabulky 4.1.

Typ p evodu	Po et p evod	Doba trvání p evodu [s]		
		Minimální	Maximální	Pr m rná
Neo4j -> Lokální MySQL	10	2110	2235	2165
Lokální MySQL -> Neo4j	10	772	799	783
Neo4j -> Vzdálený MySQL	10	3144	4083	3716
Vzdálený MySQL -> Neo4j	10	773	801	788

Tabulka 4.1: Rychlost p evodu dat mezi databázemi

Z nam ěných dat 4.1 jde jednozna n vid t, že p evod z MySQL databáze do Neo4j je asov m ěn ěr ě, a to jak p i p evodu s lokálním nebo vzdáleným databázovým serverem MySQL. Dokonce u tohoto p evodu nemá typ p ipojení k MySQL serveru žádný vliv. Protože se p i p evodu z MySQL databáze do Neo4j databáze používá pouze SQL dotaz SELECT, pro získání dat k p enesení a poté se zbytek d l ě v grafové databázi nám dokazuje, že Neo4j databáze zaru uje rychlejší manipulaci s daty, tedy alespo dotazy p i ukládání dat, než databáze rela ní. Naopak u p evodu z Neo4j do MySQL se krom vyhledávacích dotaz obou dotazovacích jazyk používají hlavn SQL dotazy INSERT a UPDATE, které z velké ásti zp sobují pomalejší p evod dat.

#### 4.5.2 Správnost p evodu

P edchozí testování slouží pouze pro informaci, jak p ibližn ě dlouho daný p evod trvá. Hlav- ním testováním je kontrola správnosti p evodu. Toto testování se provád ě pomocí vygene- rovaného csv souboru z p vodní Neo4j databáze a testovacího skriptu. Výsledky testování se zapíší do externího souboru ur ěného k zapisování pr b hu

Testování za íná tak, že se zkontroluje správnost p evedení všech uzl do jednotlivých záznam v tabulkách navržené rela ní databáze. Jedná se o kontrolu mezi následujícími typy uzl a k nim p íslušnou tabulkou, a to jak kontrolou po tu p evedených záznam ,



tak i správnost provedených hodnot záznamů. Typ uzlu Osoba s tabulkou osoba, Matřka s tabulkou matřka, Mesto s tabulkou mesto, Ulice s tabulkou ulice, Popisné číslo s tabulkou cislo\_pop, Záznam\_o\_křte s tabulkou zaznam\_narozeni, Záznam\_o\_svatbe s tabulkou zaznam\_oddani a Záznam\_o\_úmrti s tabulkou zaznam\_umrti.

Po kontrole správnosti provedu u hlavních tabulek se pějde na kontrolu pomocných tabulek povolani, pohlavi, nabozenstvi, jmena\_osoby, prijmeni\_osoby a puv\_id\_db s uzlem typu Osoba, tabulky mesto\_gps s uzlem Mesto a tabulek jazyk, mesto a mesto s uzlem Matřka. U těchto kontrol se postupně prochází testovací csv soubor a pomocí sloupce id\_gdb, který obsahuje jednoznačný identifikátor na provedení záznam v Neo4j databázi, se najde daný záznam z csv souboru a otestuje se správnost provedených hodnot.

Poslední fáze testování se týká kontroly správného provedu hran. Princip je zde takový, že se z testovacího csv souboru vyberou záznamy hran, kde se nachází informace o jaký typ hrany se jedná a odkud kam tato hrana směje. Následně se každá tato hrana projde a zjišťuje se, jestli se pěnesla správně a vznikla v MySQL databázi správná vazba.

Jak už jsem zmínil na konci pědchozí části 4.4.2, tak k otestování měže posloužit i zpětně vytvořený proved. Jelikož tento proved nebyl součástí zadání a byl dodán v průběhu práce, nepoužil jsem ho k testování a pro testování správnosti provedu jsem použil pouze testování pomocí csv souboru.

# Kapitola 5

## Závěr

Hlavním úkolem mé bakalářské práce bylo navrhnout MySQL databázi a skripty pro převod genealogických dat z grafové databáze Neo4j a následně tyto návrhy implementovat a otestovat.

Otestováním jsem zjistil, že by implementace v případě, že se k převodu použijí správná data z Neo4j databáze obsahující genealogická data, měla být funkční a měla splňovat potřeby pro každou tuto implementaci došlo. Po převodu máme k dispozici stejná data jako jsou v původní grafové databázi a je možné nad nimi pracovat pomocí SQL dotazovacích skriptů. V případě jaké úpravy nebo vytvoření nových záznamů v nově vytvořené MySQL databázi je možné použít zpětný převod z MySQL do Neo4j, který byl po domluvě s vedoucím mé práce vytvořen nad rámec zadání. Tímto zpětným převodem dostaneme novou Neo4j databázi i s případnými upravenými záznamy a využít veškeré výhody grafové databáze. Jednoduše řečeno, tato práce slouží k tomu aby se s genealogickými daty mohlo pracovat jak v relační, tak i grafové databázi a to protože některé vci jdou u jednoho typu databáze řešit jednodušeji a efektivněji než u druhého typu, nebo dokonce některé vci v jedné databázi nejdou a v druhé ano.

Věřím, že tato má práce bude v případě potřeby sloužit, tak jak má a bude tak zajišťovat na lepší kompatibilitu při vyhledávání genealogických dat mezi grafovou databází Neo4j a relační databází MySQL.

# Literatura

- [1] *Genealogie aneb rodopis*. 2015. [Online; accessed 28.12.2020]. Dostupné z: <http://www.genealogie.cz/genealogie/>.
- [2] BRYCE MERKL SASAKI, R. H. *Graph Databases for Beginners*. Neo4j, 2018. [https://go.neo4j.com/rs/710-RRR-335/images/Graph\\_Databases\\_for\\_Beginners.pdf](https://go.neo4j.com/rs/710-RRR-335/images/Graph_Databases_for_Beginners.pdf).
- [3] CORPORATION, O. *MySQL 8.0 Reference Manual*. Oracle Corporation, 2021.
- [4] HARRINGTON, J. L. *Relational Database Design and Implementation*. Todd Green, 2016. ISBN 978-0-12-804399-8.
- [5] IAN ROBINSON, E. E. *Graph Databases, 2nd Edition*. 2. vyd. Beijing: O'Reilly Media, Inc, 2015. ISBN 978-1-49-193089-2.
- [6] J., D. C. *An introduction to database systems. 8th Edition*. 8. vyd. Boston: Pearson/Addison Wesley, 2004. ISBN 978-0-32-119784-9.
- [7] LEDNICKÁ, B. *Sestavte si rodokmen pátráme po svých předcích*. Grada, 2012. ISBN 978-80-247-4069-0.
- [8] PETERKA, J. *Cesta k rodinným kořenům aneb Praktická příručka obanské genealogie*. Libri, 2006. ISBN 80-7277-307-0.
- [9] SILBERSCHATZ ABRAHAM, S. S. *Database system concepts, 6th Edition*. 6. vyd. New York: McGraw-Hill Higher Education, 2011. ISBN 978-0-07-352332-3.
- [10] SKŘIVAN, J. *Databáze a jazyk SQL*. 2000. [Online; accessed 20.12.2021]. Dostupné z: <https://www.interval.cz/clanky/databaze-a-jazyk-sql/>.
- [11] SPERÁT, P. I. *Rodokmeny - Rodokmeny a nakladatelství*. 2019. [Online; accessed 28.12.2020]. Dostupné z: <https://www.sperat.cz/rodokmeny/>.
- [12] TUŠIMOVÁ, L. *Generování rodokmenů z matričních záznamů*. Brno, CZ, 2020. Diplomová práce. Vysoké učení technické v Brně, Fakulta informačních technologií. Dostupné z: <https://www.fit.vut.cz/study/theses/22818/>.

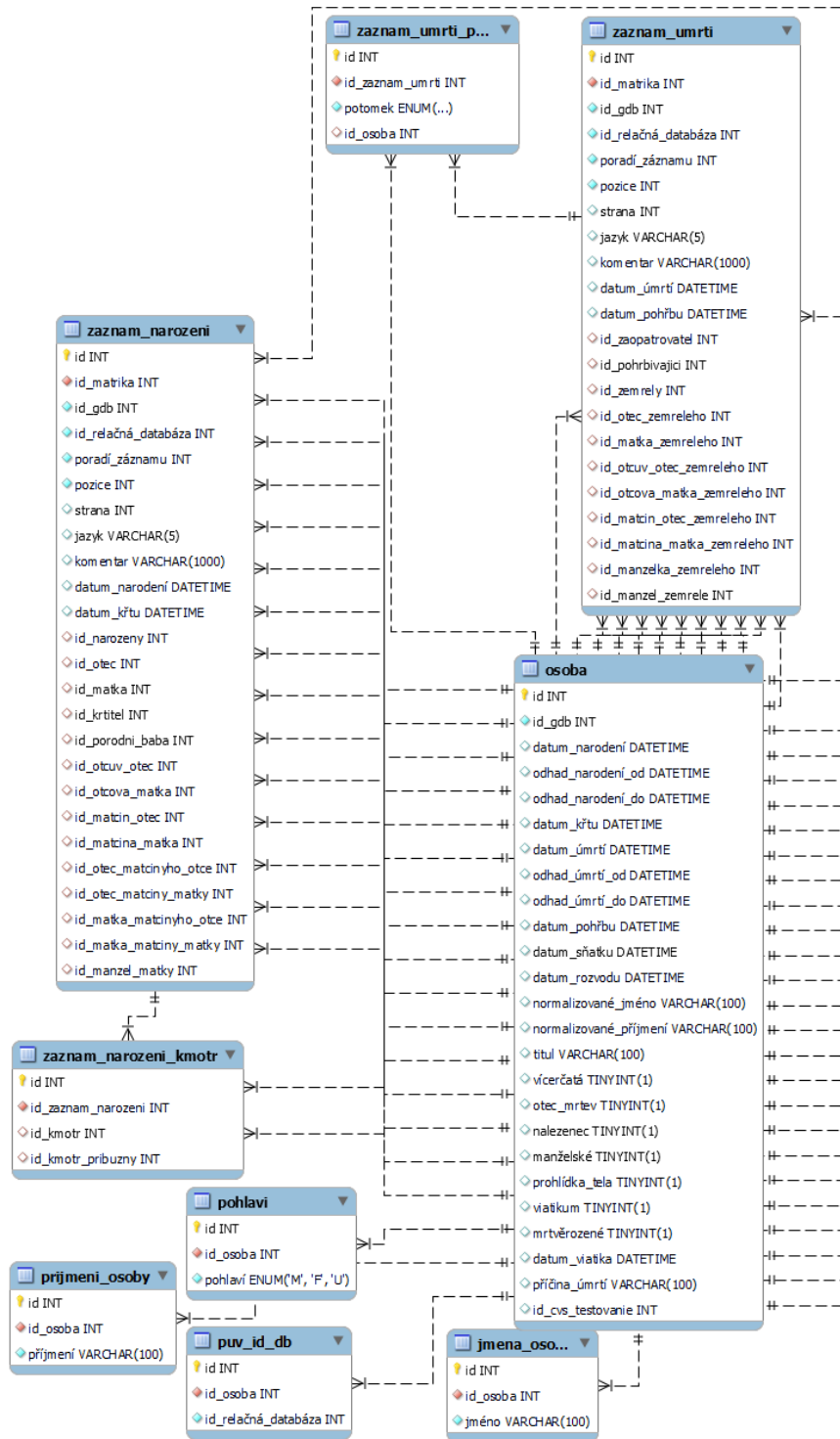
## Příloha A

# Obsah příloženého paměťového média

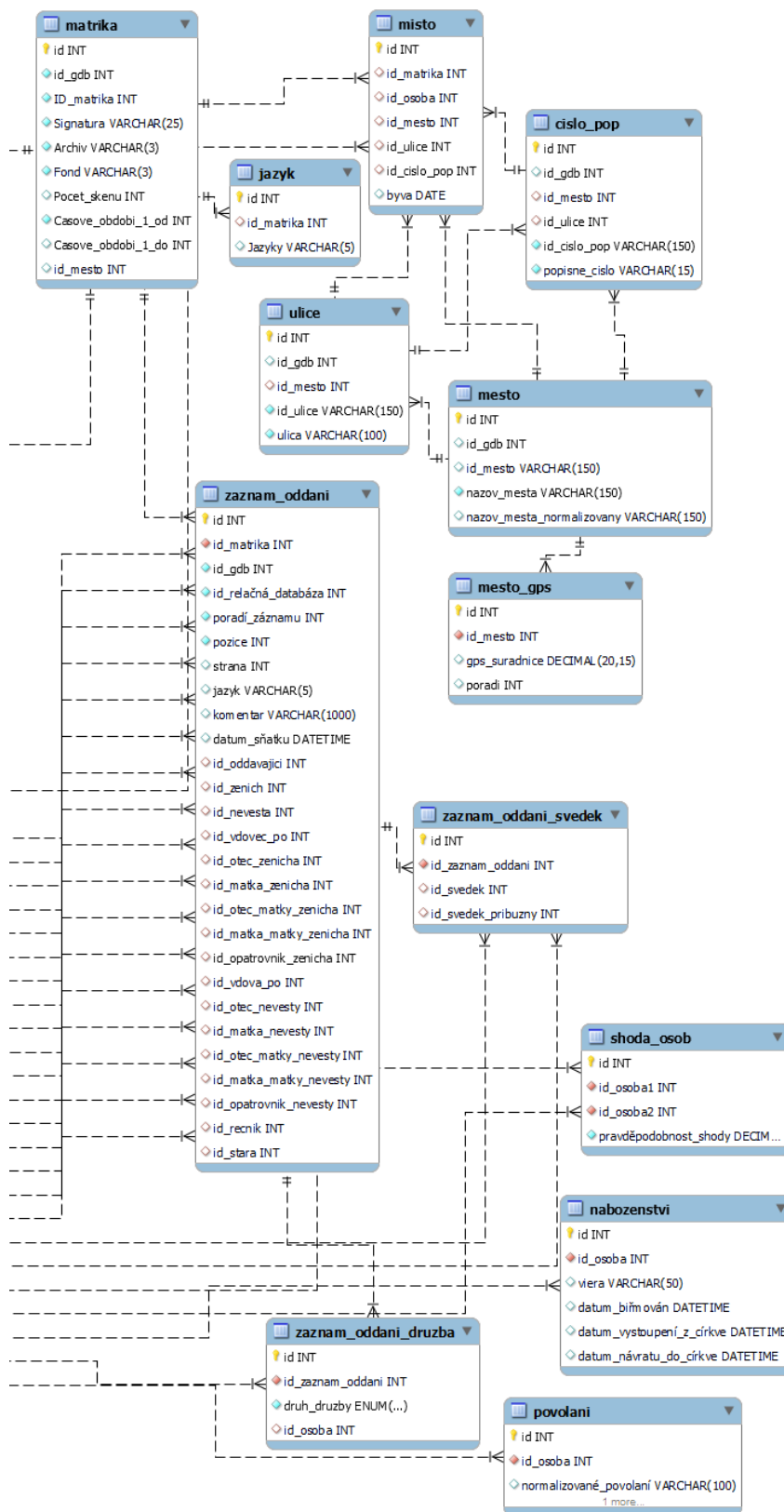
- Importní json soubor
- Testovací csv soubor
- Zdrojové kódy
- Návod k použití
- Zdrojové soubory L<sup>A</sup>T<sub>E</sub>X
- PDF technické zprávy
- Obrázek návrhu relační databáze

P íloha B

## Návrh genealogické databáze



Obrázek B.1: Návrh MySQL databáze 1/2.



Obrázek B.2: Návrh MySQL databaze 2/2.