



FAKULTA APLIKOVANÝCH VĚD
ZÁPADOČESKÉ UNIVERZITY
V PLZNI

KATEDRA INFORMATIKY
A VÝPOČETNÍ TECHNIKY



Implementace modulu pro import údajů RÚIAN

Martin Schön





**FAKULTA APLIKOVANÝCH VĚD
ZÁPADOČESKÉ UNIVERZITY
V PLZNI**

**KATEDRA INFORMATIKY
A VÝPOČETNÍ TECHNIKY**

Implementace modulu pro import údajů RÚIAN

Martin Schön

© Martin Schön, 2024.

Všechna práva vyhrazena. Žádná část tohoto dokumentu nesmí být reprodukována ani rozšiřována jakoukoli formou, elektronicky či mechanicky, fotokopírováním, nahráváním nebo jiným způsobem, nebo uložena v systému pro ukládání a vyhledávání informací bez písemného souhlasu držitelů autorských práv.

Citace v seznamu literatury:

SCHÖN, Martin. *Implementace modulu pro import údajů RÚIAN*. Plzeň, 2024. Bakalářská práce. Západočeská univerzita v Plzni, Fakulta aplikovaných věd, Katedra informatiky a výpočetní techniky. Vedoucí práce Ing. Martin Bíkl, Ing. Petr Příbyl, Ing. Martin Zíma Ph.D.

Podklad pro zadání BAKALÁŘSKÉ práce studenta

Jméno a příjmení: **Martin SCHÖN**
Osobní číslo: **A22B0144P**
Adresa: **Rpety 42, Rpety, 26801 Hořovice, Česká republika**

Téma práce: **Implementace modulu pro import údajů RÚIAN**
Téma práce anglicky:
Jazyk práce: **Čeština**

Související osoby: **Ing. Martin Zíma, Ph.D. (Konzultant z univerzity)**
Katedra informatiky a výpočetní techniky
Ing. Martin Bíkl (Konzultant mimo univerzitu)
Katedra informatiky a výpočetní techniky
Ing. Petr Příbyl (Vedoucí)
Katedra informatiky a výpočetní techniky

Zásady pro vypracování:

1. Prostudujte datové schéma registru RÚIAN a možnosti získávání dat prostřednictvím datových služeb.
2. Prozkoumejte možnosti konfigurace řešení s přihlédnutím na mapování datových struktur.
3. Navrhněte konfigurační soubor, který bude umožňovat nastavení úrovně přenášených územních objektů a nastavení cílové databáze a cílových struktur.
4. Vytvořte aplikaci, která bude pravidelně synchronizovat veřejnou databázi RÚIAN do databázových struktur podle konfiguračního souboru. Synchronizace bude probíhat buď jako kompletní sada dat nebo přírůstkově. Jako úložiště využijte databázi Oracle, Microsoft SQL Server a PostgreSQL v posledních verzích.
5. Vytvořenou aplikaci ověřte na 3 konfiguračních souborech, zhodnoťte využitelnost daného řešení pro další databázové enginy a otestujte rychlost daného řešení pro kompletní i přírůstkovou sadu dat.

Seznam doporučené literatury:

Dodá vedoucí bakalářské práce.

Podpis studenta:

Datum:

Podpis vedoucího práce:

Datum:

Prohlášení

Prohlašuji, že jsem tuto bakalářskou práci vypracoval samostatně a výhradně s použitím citovaných pramenů, literatury a dalších odborných zdrojů. Tato práce nebyla využita k získání jiného nebo stejného akademického titulu.

Beru na vědomí, že se na moji práci vztahují práva a povinnosti vyplývající ze zákona č. 121/2000 Sb., autorského zákona v platném znění, zejména skutečnost, že Západočeská univerzita v Plzni má právo na uzavření licenční smlouvy o užití této práce jako školního díla podle § 60 odst. 1 autorského zákona.

V Plzni dne 31. prosince 2024

.....

Martin Schön

V textu jsou použity názvy produktů, technologií, služeb, aplikací, společností apod., které mohou být ochrannými známkami nebo registrovanými ochrannými známkami příslušných vlastníků.

Abstrakt

Text abstraktu v jazyce práce, tj. zde česky.

Abstract

The abstract text in a secondary language, here in English.

Klíčová slova

Poděkování

Text poděkování.

Obsah

1	Úvod	5
2	RÚIAN	7
2.1	Výměnný formát RÚIAN	7
2.2	Datové struktury	8
2.3	Formát dat	9
3	Databázové systémy	11
3.1	Microsoft SQL Server	12
3.2	PostgreSQL	12
3.3	Oracle Database	13
3.4	Komunikace s využitím SQL	13
3.5	JDBC	14
4	Návrh aplikace	15
4.1	Stahování dat	16
4.2	Zpracování dat	16
4.3	Komunikace s databází	17
4.4	Formát konfiguračního souboru	17
4.4.1	XML	17
4.4.2	JSON	17
4.4.3	YAML	18
4.4.4	INI	18
4.4.5	Volba formátu	18
4.5	Obsah konfiguračního souboru	19
5	Použité technologie a nástroje	21
5.1	REST API	21
5.2	Spring Framework	21
5.3	Maven	22
5.4	Plánovač	22

5.5	Docker	22
5.6	Grafický klient pro správu databáze	23
5.7	Log4j	23
6	Příprava databázových systémů	25
6.1	PostgreSQL	25
6.2	Microsoft SQL Server	25
6.3	Oracle Database	26
6.4	Odlišnost ve skriptech pro tvorbu databází	26
7	Implementace aplikace	29
7.1	Architektura aplikace	29
7.2	Inicializace aplikace	29
7.3	Mapování dat	30
7.3.1	Dto objekty	30
7.3.2	Repositáře	30
7.3.3	Service třídy	30
7.3.4	Připojení k databázi	31
7.3.5	Načtení konfigurace	31
7.4	Quartz Scheduler	33
7.4.1	Jobs	33
7.4.2	Triggers	34
7.4.3	Zajištění správného pořadí	34
7.5	Získávání dat z API	36
7.5.1	VdpDownload	36
7.5.2	VdpClient	36
7.6	Zpracování dat	37
7.6.1	Parsing dat	37
7.6.2	Atributy Objektů	40
7.6.3	Ukládání dat do databáze	42
7.7	Po zpracování dat	43
8	Testování	45
8.1	Rychlostní porovnání databázových systémů	45
8.2	Testování Regionů	46
8.3	Testování Přírůstků	49
9	Budoucí rozšíření a úpravy	53
10	Závěr	55

A	Uživatelská příručka	57
B	Struktura přiloženého zip souboru	65
	Bibliografie	67
	Seznam obrázků	71
	Seznam tabulek	73

Úvod

Problematika správy územní identifikace a prostorových dat a jejich synchronizace mezi různými systémy nabývá na významu s rostoucí digitalizací státní správy a soukromých sektorů. Jedním z klíčových zdrojů těchto dat v České republice je Registr územní identifikace, adres a nemovitostí (RÚIAN), který poskytuje rozsáhlé a aktuální informace o územních objektech, adresách a dalších klíčových entitách. Efektivní využití dat z RÚIAN vyžaduje nejen jejich přístup prostřednictvím datových služeb, ale také robustní řešení pro mapování, konfiguraci a synchronizaci datových struktur.

Cílem této bakalářské práce je analyzovat datové schéma registru RÚIAN a možnosti získávání dat prostřednictvím nabízených datových služeb. Dále bude provedena analýza a návrh konfiguračního řešení, které umožní nastavit úroveň přenášených územních objektů a cílové databázové struktury. V rámci práce bude navržena a implementována aplikace, která umožní pravidelnou synchronizaci dat z veřejné databáze RÚIAN do cílových databázových struktur s podporou databází Oracle, Microsoft SQL Server a PostgreSQL. Aplikace bude schopna provádět synchronizaci kompletních datových sad i přírůstkových změn podle zadané konfigurace.

RÚIAN

RÚIAN je zkratkou pro Registr územní identifikace, adres a nemovitostí. Jedná se o státní informační systém v České republice, který obsahuje informace o adresách, budovách, parcelách a dalších objektech. Systém je spravován Českým úřadem zeměměřickým a katastrálním (ČÚZK). Data jsou využívána v mnoha oblastech, například v urbanistickém plánování, geodézii nebo při správě nemovitostí. Jednotlivé prvky jsou zobrazovány na mapách státního mapového díla a digitální mapě veřejné správy.

Data z RÚIAN jsou veřejně dostupná a lze je získat z webových služeb na adrese <https://vdp.cuzk.gov.cz/vdp/ruian>. Lze stahovat data ve formátu XML, která obsahují základní nebo úplné informace o územních prvcích. Mezi tyto prvky patří Stát, VÚSC (Vyšší územní samosprávný celek), ORP (Obec s rozšířenou působností), Obec, Část obce, Ulice, Adresa atd. Data lze vyhledávat, ověřovat a stahovat dle jednotlivých územních prvků, které jsou uloženy v databázi RÚIAN.

2.1 Výměnný formát RÚIAN

Výměnný formát RÚIAN (VFR) je jednou ze služeb, které poskytuje ČÚZK. Tento formát slouží k přenosu dat mezi různými informačními systémy.

Je možné stahovat data podle zadaných formátů: **Standardní**, **Historický** a **Speciální**. Dále je možné si vybrat mezi přírůstkovými daty a úplnou kopií. Přírůstky je možné vyhledávat podle data – od zvoleného dne až do současnosti. Úplná kopie obsahuje všechna data a je možné ji také časově vymežit. Tato data se aktualizují jednou měsíčně.

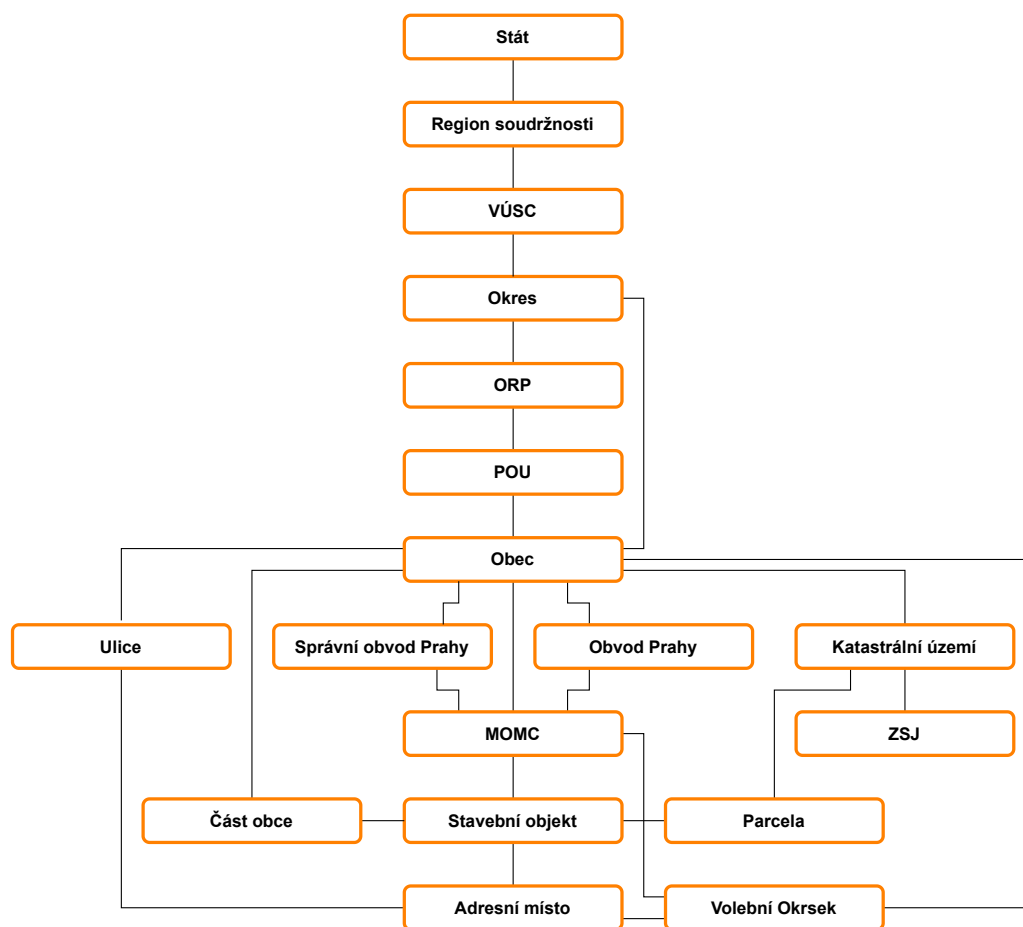
Každý formát navíc nabízí další parametry, které lze nastavit. Data z VFR jsou ve formátu XML. Každý XML element obsahuje atributy, které nesou informace o dané entitě (tabulce).

- **Standardní** – obsahuje úplná nebo přírůstková data.
 - Časový rozsah: přírůstky od data / úplná kopie,
 - Územní prvky: Stát až ZSJ / Obec a podřazené,
 - Datová sada: základní / kompletní,
 - Výběr údajů: základní údaje / generované hranice, originální hranice, vlajky a znaky,
 - Územní omezení: ČR / kraj (VÚSC) / ORP / obec,
- **Historický** – obsahuje historická data.
 - Časový rozsah: přírůstky od data / úplná kopie,
 - Územní prvky: Stát až ZJS / Obec a podřazené,
 - Územní omezení: ČR / kraj (VÚSC) / ORP / obec,
- **Speciální** – obsahuje speciální datové sady.
 - Časový rozsah: přírůstky od data / úplná kopie,
 - Výběr údajů: číselníky / vazby / vazby a číselníky,
 - Kategorie: všechny / geodetické body / nerostné bohatství,

2.2 Datové struktury

Data z RÚIAN jsou rozdělena do několika datových struktur neboli entit. Jak je vidět na obrázku 2.1 [1], každá entita obsahuje specifické informace. Mezi tyto informace patří například název státu, kód státu, geografické souřadnice, datum vzniku a další.

Stát představuje nejvyšší úroveň hierarchie, pod kterou spadají další entity závislé na ní. Příkladem je entita **VÚSC**, která obsahuje informace o vyšších územních samosprávných celcích. Jednotlivé entity na sebe navzájem odkazují pomocí cizích klíčů.



Obrázek 2.1: Schéma prvků RÚIAN

2.3 Formát dat

Data z RÚIAN jsou ve formátu XML. Tento formát je strukturovaný a umožňuje přenášet data mezi různými systémy. Bude proto třeba implementovat parser, který převede data z XML do dále zpracovávaného formátu pro databáze.

Databázové systémy

Databázový systém je softwarový nástroj, který slouží k efektivnímu ukládání, organizaci a vyhledávání dat. Díky své strukturované povaze umožňuje správu velkých objemů dat a poskytuje funkcionality pro zajištění konzistence, bezpečnosti a rychlého přístupu k uloženým informacím.

Databázový systém lze obecně rozdělit do dvou hlavních kategorií: **relační** a **objektové databázový systém**. Každý z těchto přístupů má své specifické vlastnosti a je vhodný pro odlišné typy aplikací.

Relační databázový systém

- Data jsou uložena v tabulkách.
- Každý řádek tabulky obsahuje jeden záznam.
- Každý sloupec tabulky obsahuje jeden atribut.
- Vztahy mezi tabulkami jsou definovány klíči.
- Využití jazyka SQL.

Relační databázový systém jsou vhodné pro strukturovaná data, která mají pevnou strukturu. Jedná se o nejčastěji používaný typ databázových systémů.

Objektové databázový systém

- Data jsou uložena jako objekty.
- Každý objekt obsahuje atributy a metody.
- Vztahy mezi objekty jsou definovány referencemi.
- Využití objektově orientovaného jazyka.

Objektové databázové systémy jsou vhodné pro nestrukturovaná data, která mají složitou strukturu. Jedná se o novější typ databázových systémů, který je vhodný pro moderní aplikace.

Vzhledem k zadání, kde je přímo specifikováno, které databázové systémy budou použity, není třeba vybírat mezi těmito dvěma systémy. Všechny tři databáze, které

budou popsány, jsou relační databázové systémy. Tyto systémy však mají mezi sebou rozdíly v použití, funkcích a možnostech.

Databáze bude potřebovat některé dodatečné funkce, které jsou nezbytné pro práci s daty. Mezi tyto funkce patří:

- Zpracování geometrických dat.
- Podpora JSON.
- Podpora vhodných datových typů (čas a datum, čísla, text).

3.1 Microsoft SQL Server

Microsoft SQL Server je relační databázový systém, který vyvinula společnost Microsoft a který se stal jedním z předních nástrojů pro ukládání, správu a analýzu dat. SQL Server je robustní a výkonný systém, který nabízí širokou škálu funkcí a možností přizpůsobení pro různé typy aplikací. Díky své dlouhodobé podpoře a integraci s dalšími produkty společnosti Microsoft, jako je Azure nebo Power BI, je SQL Server oblíbenou volbou pro velké a střední podniky.

Edice SQL Serveru zahrnují Standard, Enterprise a Express, které se liší funkcemi a cenou. Standard Edition je nejčastěji používanou edicí, která obsahuje všechny základní funkce a je vhodná pro většinu aplikací.

SQL Server byl původně navržen výhradně pro Windows, ale od verze SQL Server 2017 je dostupný také pro operační systém Linux. Tato multiplatformní podpora zvyšuje jeho použitelnost v různých IT prostředích. [2]

Pro tuto práci bude využita edice SQL Server 2017 Standard, která je dostupná pro Windows a Linux. Tato edice podporuje vhodné datové typy pro zpracování geometrických dat.

3.2 PostgreSQL

PostgreSQL je open source relační databázový systém, který je známý svou spolehlivostí, výkonem a rozšiřitelností. Původně byl vyvinut jako alternativa k proprietárním řešením, jako je SQL Server, a dnes se řadí mezi nejpokročilejší relační databázové systémy na trhu. Díky své otevřené povaze a aktivní komunitě uživatelů a vývojářů se stal oblíbenou volbou nejen mezi malými firmami, ale také ve středních a velkých organizacích.

PostgreSQL je dostupný zdarma a podporuje všechny hlavní operační systémy, včetně Windows, Linux a macOS. Tato multiplatformní dostupnost umožňuje snadnou integraci PostgreSQL do různých vývojových prostředí. [3]

PostgreSQL také podporuje jak formát JSON, tak všechny potřebné datové typy kromě geometrických dat. Pro ukládání geometrických dat je třeba použít *PostGIS*, což je rozšíření pro PostgreSQL, které přidává podporu pro prostorové a geometrické datové typy.

3.3 Oracle Database

Oracle Database, vyvinutá společností Oracle Corporation, patří mezi přední relační databázové systémy na světě. Tento systém je známý svou robustností, vysokým výkonem a schopností zvládat kritické podnikové aplikace a rozsáhlé datové sady. Oracle Database je navržena tak, aby poskytovala spolehlivé a efektivní řešení pro ukládání, správu a analýzu dat ve velkých i středních organizacích.

Edice Oracle Database zahrnují Standard Edition, Enterprise Edition a Express Edition, které se liší funkcemi a cenou. Enterprise Edition je nejkomplexnější edicí, která obsahuje všechny pokročilé funkce a je vhodná pro velké podniky s náročnými požadavky.

Oracle Database je kompatibilní s většinou hlavních operačních systémů, včetně Windows, Linux a Unix. Díky tomu může být nasazena v různorodých IT prostředích podle požadavků organizace. [4]

Pro tuto práci bude využita edice Oracle Express Edition, která je dostupná zdarma a je určena pro vývoj a testování. Tato edice podporuje všechny potřebné datové typy a také JSON. Ovšem není zde možnost uložení geometrických dat. Ta jsou obsažena v `SDO_GEOMETRY`, což je rozšíření pro Oracle Database.

3.4 Komunikace s využitím SQL

Všechny tři databázové systémy podporují dotazovací jazyka SQL (Structured Query Language). SQL je standardizovaný jazyk pro práci s relačními databázovými systémy, které umožňuje vytváření, čtení, aktualizaci a mazání dat. Pro komunikaci s databází je tedy třeba vytvořit SQL dotazy, které budou provádět operace nad daty. O tuto komunikaci se stará aplikační vrstva, která zajišťuje připojení k databázi a následné zpracování a odeslání SQL dotazů.

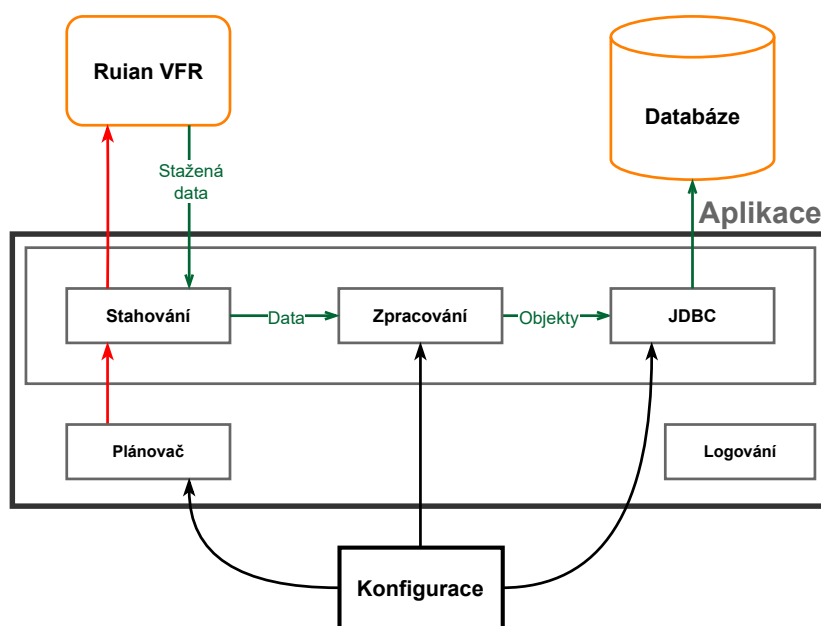
Vytvářet dotazy, pro ukládání různých dat bude velmi náročné a neefektivní. Ve výsledné aplikaci bude použito ORM (Object-Relational Mapping), které umožňuje práci s daty pomocí objektů. ORM je technika, která mapuje objekty v programovacím jazyce na tabulky v databázi a naopak. Konkrétně bude použito *Hibernate*, což je open-source ORM framework pro jazyk Java.

3.5 JDBC

Jedná se o Java API pro přístup k relačním databázovým systémům. JDBC (Java Database Connectivity) poskytuje standardní rozhraní pro připojení k databázím. Bude použito pro připojení k databázi a provádění SQL dotazů. Využívá různé ovladače pro různé databázové systémy, takže je možné připojit se k jakékoli databázi, která podporuje JDBC.

Návrh aplikace

Aplikace je prostředníkem mezi RÚIAN VFR a cílovou databází. Hlavním úkolem je stažení dat z RÚIAN VFR, jejich přečtení, zpracování a následné uložení do cílové databáze. Aplikace bude psána v jazyce Java. Bude třeba zajistit funkce pro zpracování a uložení dat. Dále bude třeba zajistit, aby se aplikace věděla, co bude dělat a jak se chovat. Pro toto bude třeba vytvořit konfigurační soubor, který bude obsahovat nastavení aplikace. Na obrázku 4.1 je znázorněna architektura aplikace. Zelená šipky popisují směr toku dat. Červené šipky popisují směr toku ovládacích signálů.



Obrázek 4.1: Architektura aplikace

4.1 Stahování dat

Stahování dat bude zajištěno pomocí knihovny *Apache HttpClient*, která je součástí balíčku *Apache HttpComponents*. Tato knihovna umožňuje snadné a efektivní stahování dat z webových stránek a API. V rámci této aplikace bude použita k stahování dat z RÚIAN VFR. Data budou stažena jako ZIP soubor, který bude následně rozbalen a zpracován. Soubory budou uloženy do dočasného adresáře, který bude po zpracování smazán z důvodu úspory volného místa.

4.2 Zpracování dat

Zpracování dat se bude skládat z několika částí: přečtení a mapování dat do objektů, které budou následně uloženy do databáze. Pro čtení dat bude třeba parser XML souborů, který přečte data a převede je do objektů. Vzhledem k velikosti dat bude třeba zajistit efektivní zpracování, aby nedocházelo k přetížení paměti a CPU. Pro čtení je možné využít knihovnu *Jackson* nebo *StAX*.

1. **Jackson** – knihovna pro zpracování JSON a XML dat. Umožňuje snadné mapování objektů na JSON a XML a naopak. Je velmi rychlá a efektivní, ale může být složitější na použití.
2. **StAX** – knihovna pro zpracování XML dat. Umožňuje čtení a zápis XML dat pomocí událostí. Je velmi rychlá a efektivní, ale může být složitější na použití.

Vybrán byl parser **StAX**, z důvodů snadného použití a rychlosti. Také nepotřebuje mít načtený celý soubor do paměti, ale čte jako proud.

Následně bude třeba vytvořit objekty pro mapování dat do objektů. Tyto objekty budou mít stejnou strukturu jako data z RÚIAN VFR. Pro mapování bude využita knihovna *JPA* (Java Persistence API), která umožňuje snadné mapování objektů na databázové tabulky a naopak. Pro tuto technologii je třeba vytvořit databázové entity, které budou mít stejnou strukturu jako v databázi. Dále bude třeba vytvořit repositáře pro práci s databází (CRUD operace). A nakonec bude třeba vytvořit služby, které budou sloužit pro práci s repositáři a pro zpracování dat.

Před uložením do databáze bude třeba provést validaci dat, aby nedocházelo k chybám při ukládání. Je třeba zajistit, aby se zabránilo ukládání dat, která jsou neplatná nebo nekompletní. Možné chyby při validaci dat:

- Chybějící primární klíče,
- Nevalidní cizí klíče,

4.3 Komunikace s databází

Pro ukládání dat je třeba se nejprve připojit k databázi. Pro připojení k databázi bude použita knihovna *JDBC* (Java Database Connectivity), která umožňuje připojení k různým databázím pomocí standardního API. Pro jednotlivé databáze budou použity různé JDBC ovladače, které umožňují připojení k vybraným databázím.

Pro připojení k databázi bude třeba vytvořit konfigurační soubor, který bude obsahovat informace o připojení k databázi.

4.4 Formát konfiguračního souboru

Je třeba zajistit, aby aplikace byla schopna načíst konfigurační soubor a podle něj se správně nakonfigurovat. Je tedy nutné zvolit vhodný formát tohoto souboru. Existuje několik možností, z nichž je třeba vybrat tu nejvhodnější:

- **XML (XSD)** – Extensible Markup Language,
- **JSON** – JavaScript Object Notation,
- **YAML** – YAML Ain't Markup Language,
- **INI** – Initialization File,

4.4.1 XML

XML je značkovací jazyk, na rozdíl od JSON a YAML. Jedná se o textový formát, který však není tak snadno čitelný pro člověka jako JSON nebo YAML. Podporuje víceúrovňovou strukturu a komentáře. Dále umožňuje použití schémat, což dovo-luje definovat strukturu a typy dat. XML je však zbytečně složité a náročné na čtení i psaní. Navíc je čtení velmi pomalé a náročné na výkon. [5]

4.4.2 JSON

JSON je formát pro výměnu dat, který je snadno čitelný jak pro člověka, tak pro počítač. Jedná se o textový formát, který se často používá k přenášení dat mezi serverem a klientem. Podporuje víceúrovňovou strukturu, ale nepodporuje komentáře. Je však velmi rozšířený a podporovaný většinou programovacích jazyků. Čtení i zápis dat je jednoduchý a rychlý. [5]

4.4.3 YAML

YAML je formát pro serializaci dat, který je čitelností a strukturovaností podobný JSON. Oproti JSON podporuje komentáře. Pro člověka však může být YAML složitější a náročnější na psaní. Rychlost čtení je srovnatelná s JSON, ale zápis je pomalejší. [5]

4.4.4 INI

INI je formát konfiguračních souborů, který je snadno čitelný pro člověka. Podporuje komentáře i víceúrovňovou strukturu. Je však často nejednotný a obtížně se s ním pracuje. Používá se spíše pro jednodušší konfigurační soubory a není vhodný pro složitější aplikace.

4.4.5 Volba formátu

Vzhledem k tomu, že aplikace bude obsahovat mnoho funkcionalit a bude vyžadovat nastavení řady parametrů, je třeba zvolit vhodný formát. Prozatímní verze aplikace si žádá formát, který bude čitelný, snadno rozšiřitelný a dostatečně flexibilní. Proto je nejvhodnější volbou formát **JSON**. JSON je snadno čitelný a zápis je přehledný, podporuje víceúrovňovou strukturu a je široce rozšířený. Přestože nepodporuje komentáře, lze jeho strukturu snadno pochopit a rozšířit o další parametry. Zároveň je JSON podporován většinou programovacích jazyků, což zajišťuje snadnou integraci.

4.5 Obsah konfiguračního souboru

Konfigurační soubor bude obsahovat nastavení aplikace. Jedná se o nastavení databáze, aplikace, plánovače a parametry pro stahování dat. Konfigurace se bude skládat z následujících částí:

1. **Databáze** – Nastavení připojení k databázi:

- Typ databáze – MSSQL, PostgreSQL, Oracle,
- Connection string – připojovací řetězec,
- Uživatelské jméno – přihlašovací jméno do databáze,
- Heslo – heslo do databáze,
- Název databáze – cílová databáze pro ukládání dat,

2. **Nastavení plánovače** – Nastavení plánovače, který bude stahovat data:

- Interval – interval stahování dat,
- Přeskočení – možnost přeskočit naplánované stahování,

3. **Parametry pro stahování dat** – Parametry pro stahování dat z webové služby:

- Seznam krajů – výčet krajů, pro které se budou data stahovat,
- Stahovat geometrii – ANO/NE (např. kvůli absenci geometrických typů v Oracle XE),
- Velikost commitů – počet záznamů na jeden commit do databáze,
- Konkrétní tabulky, sloupce a způsob práce s vybranými daty,

Toto nastavení bude uloženo v konfiguračním souboru, který bude načten při startu aplikace. Bude tedy třeba vytvořit třídu, která zajistí načtení tohoto souboru a zpracování jeho obsahu.

Použité technologie a nástroje

Některé technologie pro tuto práci již byly zmíněny (Java, JDBC, Hibernate atd.) v sekci 4.3. Pro účely této práce budou potřeba ještě další technologie, které umožní tvorbu aplikace.

5.1 REST API

REST API (Representational State Transfer Application Programming Interface) je architektonický styl pro návrh webových služeb, který umožňuje snadnou a efektivní komunikaci mezi klientem a serverem. Tento přístup se stal jedním z nejrozšířenějších způsobů integrace aplikací díky své jednoduchosti, flexibilitě a nezávislosti na platformě. REST API využívá standardní metody protokolu HTTP, jako jsou GET, POST, PUT a DELETE, k provádění různých operací s daty.

REST API poskytuje ideální prostředí pro implementaci přenosu dat díky své flexibilitě a schopnosti pracovat s různými datovými zdroji. V této práci bude kladen důraz na robustnost a spolehlivost API, což zahrnuje zpracování chybových stavů, zabezpečení komunikace (například prostřednictvím HTTPS) a optimalizaci výkonu. [6]

5.2 Spring Framework

Spring Framework je open-source framework pro vývoj aplikací v jazyce Java. V tomto frameworku bude vytvořena aplikace, která bude vše spojovat dohromady. Spring Framework obsahuje mnoho modulů, které usnadňují vývoj aplikací, jako například Spring Boot, Spring Data, Spring Security atd. Pro účely této práce bude využit modul Spring Boot, který umožňuje rychlé vytvoření aplikace s minimální konfigurací. [7]

5.3 Maven

Maven je nástroj pro správu projektů v jazyce Java, který usnadňuje správu závislostí, kompilaci a nasazení aplikací. Jedná se o jednoduchý a efektivní nástroj, který umožňuje správu projektů pomocí XML konfiguračního souboru. Konkrétně se jedná o soubor `pom.xml`, který obsahuje informace o projektu, jako jsou závislosti, pluginy a další nastavení. [8]

5.4 Plánovač

V popisu, co bude obsahovat konfigurační soubor, bylo zmíněno, že bude třeba zvolit plánovač. Plánovač je nástroj, který umožňuje spouštění úloh v pravidelných intervalech. Výhodou plánovače je, že umožňuje automatické spouštění úloh bez nutnosti manuálního zásahu.

Existuje několik možných plánovačů, které lze použít:

- **Cron** – Unix,
- **Task Scheduler** – Windows,
- **Quartz** – Java,
- **Apache Airflow** – Python,

Všechny tyto plánovače umožňují spouštění úloh v pravidelných intervalech. Pro účely této práce byl zvolen plánovač Quartz, který je napsán v jazyce Java a je podporován Spring Frameworkem.

Intervaly mohou být nastaveny v cron notaci, což umožňuje velkou flexibilitu při plánování úloh (například každý den ve 3:00, každý týden v pondělí v 8:00, každý měsíc první den ve 12:00 atd.).

5.5 Docker

Docker je open-source platforma pro vývoj, nasazení a provoz aplikací. Umožňuje vytváření kontejnerů, které obsahují všechny potřebné závislosti pro běh aplikace.

Pro každou databázi je třeba vytvořit instanci, která bude obsahovat potřebné tabulky, data a klienta pro komunikaci a práci s databází. Toto je úloha jako stvořená pro Docker, který umožňuje vytvoření kontejneru s databází, který bude obsahovat veškeré potřebné závislosti pro běh aplikace. [9]

5.6 Grafický klient pro správu databáze

Klient pro komunikaci s databází je nástroj, který umožňuje připojení k databázi a provádění dotazů. Je třeba vybrat klienta, který bude podporovat náhled do všech databází použitých v této práci.

Možnosti jsou:

- **DBeaver** – open-source nástroj pro správu databází, který podporuje mnoho různých databází.
- **SQL Developer** – nástroj pro správu databází Oracle.
- **Adminer** – open-source nástroj pro správu databází, který může zároveň běžet v Dockeru.

Všechny tyto nástroje umožňují připojení k databázi a provádění dotazů. Pro účely této práce byl zvolen DBeaver, který podporuje širokou škálu databází a edice Community Edition je open-source. Zároveň umožňuje export dat z databáze do různých formátů (CSV, Excel, JSON atd.). [10]

5.7 Log4j

Log4j je open-source knihovna pro logování v jazyce Java. Umožňuje snadné a efektivní logování událostí v aplikaci. Je možné nastavit různé úrovně logování (DEBUG, INFO, WARN, ERROR, FATAL) a také různé výstupy (soubor, konzole, databáze atd.). Tato technologie bude používána v průběhu vývoje aplikace pro logování událostí. [11]

Příprava databázových systémů

Každá databáze měla své problémy, ale všechny se nakonec podařilo vyřešit. V následujících částech bude popsáno, jak byly jednotlivé databáze implementovány a jaké problémy se objevily.

6.1 PostgreSQL

PostgreSQL se ukázala jako nejjednodušší databáze pro implementaci. Databáze běží na lokálním serveru v Dockeru. Pro vytvoření byl stažen oficiální image z Docker Hubu společně s knihovnou PostGIS, která je potřebná pro práci s geodaty.

```
docker pull postgres:latest
```

Po stažení image byl následně vytvořen a spuštěn kontejner pomocí docker-compose:

```
docker-compose up -d
```

Databáze je inicializována skriptem 'init.sql', který vytvoří potřebné tabulky a indexy. Tento skript byl napsán specificky pro databázi PostgreSQL a pomocí něj je vytvořena databáze se všemi potřebnými tabulkami.

6.2 Microsoft SQL Server

Microsoft SQL Server byla druhá databáze, která byla implementována. Při stahování oficiálního image z Docker Hubu se objevily problémy. Po úspěšném stažení image 'mssql/server:2017' byl kontejner spuštěn stejně jako u PostgreSQL pomocí docker-compose.

Menším rozdílem je způsob spouštění 'init' skriptu, který se nespouští při inicializaci databáze přes 'entrypoint', ale až následně pomocí příkazu 'sqlcmd'.

6.3 Oracle Database

Oracle byla poslední databáze, která byla implementována. Při stahování oficiálního image z Docker Hubu došlo k problémům. Když se však image podařilo stáhnout a kontejner spustit, databáze nefungovala správně. Nezbyvalo nic jiného než stáhnout Oracle XE a nainstalovat databázi na lokální stroj. Byla stažena verze 21c Express Edition, která je zdarma pro osobní použití. Tato verze byla stažena z oficiálních stránek Oracle <https://www.oracle.com/database/technologies/appdev/xe/quickstart.html>. Konkrétně byla stažena verze pro Windows 64-bit. Po nainstalování byla v databázi vytvořeno schéma **XEPDB1**, které je výchozím schématem pro databázi XE. V tomto schématu byl vytvořen uživatel *ruian_user* s heslem *12345*, který nadále bude přistupovat do databáze. V GUI DBeaver by následně použit inicializační skript *init.sql*, který vytvoří potřebné tabulky.

6.4 Odlišnost ve skriptech pro tvorbu databází

Každá databáze měla vlastní skript pro inicializaci databáze. Hlavní rozdíl byl v syntaxi SQL příkazů. Konkrétně se jednalo o rozdíly v datových typech, které byly pro každou databázi odlišné viz Tabulka 6.1.

Tabulka 6.1: Datové typy v různých databázích

	Oracle	PostgreSQL	MSSQL
Integer	NUMBER	INTEGER	INT
Long	NUMBER(19)	BIGINT	BIGINT
DateTime	DATE	TIMESTAMP	DATETIME
String	VARCHAR2(length)	VARCHAR(length)	NVARCHAR(length)
JSON	JSON	JSONB	NVARCHAR(MAX)
Boolean	NUMBER(1)	BIT	BIT
Geometry	SDO_GEOMETRY	GEOMETRY	GEOMETRY
Binary	BLOB	BYTEA	VARBINARY

Dalším rozdílem byly také cizí klíče (foreign key), které byly v každé databázi definovány odlišně. Zatím co v PostgreSQL6.1 a MSSQL6.2 se cizí klíče definují pomocí příkazu **REFERENCES**, a **FOREIGN KEY** přímo u sloupce, v Oracle6.3 se musí nejprve definovat sloupec a poté se cizí klíč definuje pomocí příkazu **CONSTRAINT** na konci tabulky nebo po určitém sloupci. **CONSTRAINT** navíc vytváří jméno vazby pro cizí klíč. Dále pokračuje podobně jako v PostgreSQL a MSSQL.

Zdrojový kód 6.1: PostgreSQL definice cizího klíče

```
1 <column_name> INTEGER REFERENCES <table>(<column_name>)
2 <column_name> BIGINT REFERENCES <table>(<column_name>)
```

Zdrojový kód 6.2: MSSQL definice cizího klíče

```
1 <column_name> INT FOREIGN KEY REFERENCES
2 <table>(<column_name>)
3 <column_name> BIGINT FOREIGN KEY REFERENCES
4 <table>(<column_name>)
```

Zdrojový kód 6.3: Oracle definice cizího klíče

```
1 <column_name> <column_type>,
2 CONSTRAINT <constraint_name> FOREIGN KEY (<column_name>)
3 REFERENCES <table>(<column_name>)
```

Implementace aplikace _

Před začátkem popisu implementace aplikace je vhodné upřesnit, jaký je její účel. Cílem aplikace je stahování dat z API RÚIAN a jejich následné zpracování. Výsledkem bude projekce těchto dat do databáze, která bude sloužit jako zdroj pro další zpracování. Aplikace je napsána v programovacím jazyce Java a využívá framework Spring Boot.

Projekce databáze

Projekce databáze představuje způsob zobrazení dat z jedné databáze do jiné. Podle konfigurace se nastavuje úroveň projekce, na jejímž základě se data zobrazují. Nastavení projekce je dále popsáno v sekci 7.3.5.

7.1 Architektura aplikace

Architektura je rozdělena do několika modulů, z nichž každý má svůj specifický úkol:

- **main** – hlavní modul aplikace, který obsahuje třídu Main, jež spouští celou aplikaci.
- **config** – stará se o konfiguraci aplikace a její inicializaci.
- **download** – zajišťuje stahování dat z API RÚIAN a následné zpracování těchto dat.
- **scheduler** – spravuje spouštění úloh a konfiguraci plánovače.

Všechny závislosti a knihovny jsou spravovány pomocí nástroje Maven.

7.2 Inicializace aplikace

Na začátku běhu aplikace se provede její inicializace. Důvodem, proč byl zvolen framework Spring Boot, je schopnost aplikace automaticky načítat všechny komponenty a konfigurace. Konkrétně se jedná o moduly označené anotacemi @Component, @Service, @Configuration, @Entity, @Repository a @Bean.

7.3 Mapování dat

Pro mapování dat mezi databází a aplikací se používá **Hibernate ORM**. **Hibernate** je objektově-relační mapovací framework, který umožňuje práci s databází pomocí objektů. Pro mapování byly použity **Dto objekty**, které reprezentují jednotlivé tabulky v databázi. Následně je zde **Repository** interface, které dědí z **JpaRepository**. A nakonec je zde **Service** třída, která obsahuje logiku pro práci s databází a komunikuje s **Repository**.

7.3.1 Dto objekty

Dto objekty nebo také **Data Transfer Object** jsou objekty, které slouží k přenášení dat mezi vrstvami aplikace. Pro jejich označení se používá anotace **@Entity**, která určuje, že se jedná o entitu mapovanou na tabulku v databázi. Dále mají anotaci **@Table**, která určuje název tabulky v databázi, a anotaci **@Id**, která určuje primární klíč tabulky. Pro kontrolu je zde anotace **ToString**, která slouží pro převod objektu na řetězec a anotace **@Data** z knihovny **@Lombok**, která generuje gettery a settery pro jednotlivé atributy. Pro každé atributy, které jsou ve výsledné databázi jako **JSON** je pro sloupec použita anotace **@JdbcTypeCode(SQLType.JSON)**.

7.3.2 Repositáře

Repositáře jsou rozhraní, která dědí z **JpaRepository** a poskytují metody pro práci s databází. Mezi tyto metody patří například **findAll**, **findById**, **save** a **delete**. **Repositáře** jsou označeny anotací **@Repository**, která určuje, že se jedná o komponentu pro práci s databází.

7.3.3 Service třídy

A nakonec jsou zde **Service** třídy, které obsahují logiku pro práci s databází a komunikaci s repositáři. Tyto třídy jsou označeny anotací **@Service**, která určuje, že se jedná o komponentu pro práci s databází. Zde se provádí veškerá logika před uložením do databáze:

- Kontrola, zda Dto má primární klíč.
- Kontrola, zda jsou cizí klíče validní a existují v databázi.
- Doplnění dat do objektu, podle již existujících dat v databázi.
- Úprava dat podle konfigurace.

Po provedení těchto kontrol se objekt uloží do databáze pomocí repositáře.

7.3.4 Připojení k databázi

O připojení k databázi se stará třída `DatabaseSource`, která zajišťuje navázání spojení s databází. Z konfiguračního souboru se načtou potřebné informace pro připojení:

- `type` – typ databáze (např. `postgresql`, `mssql`, `oracle`),
- `url` – adresa databáze (např. `localhost:5432` pro PostgreSQL),
- `dbname` – název databáze (např. `ruian`),
- `username` – uživatelské jméno pro připojení k databázi,
- `password` – heslo pro připojení k databázi.

Na základě těchto informací se vytvoří připojení k databázi, tzv. **DataSource**. Tento zdroj bude dále upravován v jiném modulu. Základem `DataSource` je nastavení základních parametrů připojení. Vytváří se výsledný connection string, který se upravuje podle typu databáze. K URL se připojí název databáze, uživatelské jméno, heslo a v případě MSSQL také certifikát pro zabezpečené připojení.

V dalším modulu se pro `DataSource` nastavují další parametry potřebné pro přístup k databázi a přenos dat. Dále se nastaví dialekt pro správnou syntaxi SQL příkazů podle použité databáze.

7.3.5 Načtení konfigurace

`DatabaseSource` se stará pouze o připojení k databázi, třída `AppConfig` načítá zbytek potřebné konfigurace. V sekci 4.5 jak bude konfigurace rozdělena do jednotlivých částí. Zde se budou popisovat jednotlivé atributy, jejich význam, a tedy i finální struktura konfiguračního souboru.

1. **Konfigurace úkolů pro Quartz Scheduler** Načítá se čas ve formátu cron pro spouštění přírůstkových dat a nastavení, zda přeskočit inicializaci hlavních územních prvků a krajů.

2. **Seznam krajů s příslušnými kódy** Každý řádek obsahuje kraj a jeho kód. Pokud je v konfiguraci nastaveno přeskočení inicializace krajů nebo je seznam prázdný, tento krok se přeskočí.
3. **Dodatečná nastavení** Například volba pro ignorování geometrických dat (některé databáze nepodporují geometrické typy) a nastavení velikosti jednotlivých commitů, které slouží pro optimalizaci výkonu.
4. **Nastavení zpracování jednotlivých tabulek** Základním parametrem je způsob zpracování tabulek: `all` nebo `selected`.
 - **all** – všechny tabulky budou zpracovány bez ohledu na konfiguraci,
 - **selected** – budou zpracovány pouze tabulky výslovně uvedené v konfiguraci.
5. **Konfigurace jednotlivých tabulek** Pokud je nastaven režim `selected`, zpracovávají se pouze specifikované tabulky. Každá tabulka může mít vlastní nastavení:

Zdrojový kód 7.1: Konfigurace jednotlivých tabulek

```
1      "<table_name>": {  
2          "howToProcess": "all | exclude |  
    include",  
3          "columns": ["<column_name>", ..., "<  
    column_name>"]  
4      }  
5
```

Možnosti zpracování:

- **all** – všechny sloupce budou zpracovány,
 - **exclude** – vybrané sloupce budou ignorovány,
 - **include** – vybrané sloupce budou zpracovány.
6. **Sloupce** u jednotlivé tabulky jsou definovány jako pole řetězců s názvy sloupců ve formátu lowercase bez mezer a speciálních znaků. Sloupce jsou odděleny čárkami. Sloupce nemusí být uvedeny, pokud je nastaveno `all` u konkrétní tabulky.

Možné chyby při načítání konfigurace:

- Pokud je nastaveno `exclude` nebo `include`, ale nejsou uvedeny žádné sloupce → neplatná konfigurace,

- Sloupce nebo tabulky, které neexistují v databázi → budou přeskočeny,
- Seznam krajů je povinný atribut, i pokud je prázdný. V takovém případě budou zpracovány pouze základní územní prvky.

7.4 Quartz Scheduler

Pro synchronizaci dat a spouštění úloh v určitých intervalech je použit Quartz Scheduler. Ten je přímo integrován do frameworku Spring Boot jako součást jedné z jeho knihoven. Samotný scheduler se spouští při startu aplikace.

7.4.1 Jobs

Joby jsou úlohy, které se spouštějí v určitých intervalech nebo na základě určité události. V našem případě se jedná o 3 úlohy:

- `InitStatAzZsj`
- `InitRegion`
- `Additions`

`InitStatAzZsj`

Tato úloha se spouští při startu aplikace a je zodpovědná za inicializaci základních územních prvků a krajů. Pokud je v konfiguraci nastaveno přeskočení inicializace těchto prvků, úloha se neprovede. Úloha nejprve stáhne data o základních územních prvcích a krajích z API RÚIAN za poslední měsíc. Následně se soubor zpracuje a data se uloží do databáze.

`InitRegion`

Tato úloha se spouští jako druhá, bezprostředně po dokončení úlohy `InitStatAzZsj`. Pokud je v konfiguraci nastaveno přeskočení inicializace regionů, tato úloha se přeskočí. Úloha ze zadané konfigurace načte všechny kraje, které byly uvedeny. Každý kraj má přiřazen seznam obcí, které se postupně stahují z API RÚIAN, zpracují a uloží do databáze.

`Additions`

Tato úloha se spouští na základě nastavení v konfiguračním souboru nebo po dokončení úlohy `InitRegion`. Každý den je na API RÚIAN vytvořen nový soubor s přírůstkovými daty za poslední den. Úloha tento soubor stáhne, zpracuje a uloží do databáze. Poté čeká na další spuštění podle časového nastavení v konfiguraci.

7.4.2 Triggers

Triggers jsou spouštěče, které určují, kdy se má daný job spustit. Jak bylo zmíněno výše, úlohy `InitStatAzZsj` a `InitRegion` se spouštějí jednorázově při startu aplikace a navazují na sebe. Úloha `Additions` se naproti tomu spouští pravidelně podle nastavení v konfiguračním souboru. Konkrétně je čas spuštění definován v cron formátu.

Cron formát

Cron formát je způsob, jak určit čas spuštění úloh. Původně byl použit v systémech Unix a stal se široce rozšířeným standardem. Skládá se z následujících částí:

- sekundy (0–59),
- minuta (0–59),
- hodina (0–23),
- den v měsíci (1–31),
- měsíc (1–12 nebo zkratky názvů měsíců),
- den v týdnu (0–6 nebo zkratky názvů dní).

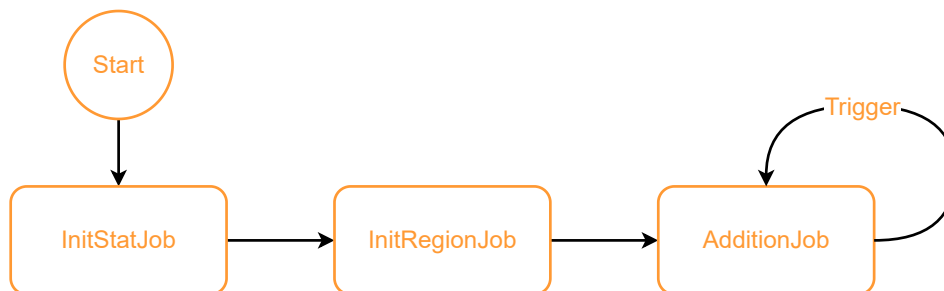
Hvězdička (*) znamená, že se úloha spouští v každém intervalu dané části. Otazník (?) znamená, že se úloha nespouští v žádném intervalu dané části. Zkratky měsíců a dnů se uvádějí v angličtině, typicky ve formátu tří písmen (např. `JAN`, `MON`). **Příklady:**

- `0 0 2 * * ?` – úloha se spustí každý den ve 2:00.
- `0 0 0 * * MON` – úloha se spustí každé pondělí o půlnoci.
- `30 14 3 JAN * ?` – úloha se spustí 3. ledna ve 14:30.

7.4.3 Zajištění správného pořadí

Je důležité zajistit, aby se úlohy spouštěly v správném pořadí. Kdyby se spustila úloha `Additions` před úlohou `InitRegion`, mohlo by dojít k problémům s neexistujícími daty. Tomu se předejde vytvořením závislostí mezi joby. Pořadí spuštění jobů je znázorněno na obrázku 7.4.3.

Obrázek 7.1: Pořadí jobů



Při startu aplikace se spustí job `InitStatAzZsj`, až úspěšně skončí, spustí se job `InitRegion`. Pokud je v konfiguraci nastaveno, že se job přeskočí tak job stále provede a dokončí jen bez zpracování dat. Stejně to platí pro job `InitRegionJob`. Problémem bylo, že job `InitStatAzZsj` a `Additions` se spouštěly vždy při startu aplikace, protože měly trigger pro start. Ovšem job `Additions` se má spustit až po úspěšném dokončení jobu `InitRegion`. Tento problém byl vyřešen pomocí **Semaforu**. Konkrétně byl použit Semaphore z balíčku `java.util.concurrent`. Vzhledem k tomu, že Quartz Scheduler má pro každý job vlastní vlákno, je možné použít Semaphore pro zajištění, že se job `Additions` může spustit až po úspěšném dokončení jobu `InitRegion`.

Semafor je inicializován na hodnotu 0, což znamená, že job `Additions` se nemůže spustit, dokud není semafor uvolněn. Ten je uvolněn následně na konci jobu `InitRegion`. Jakmile se job `Additions` spustí poprvé, je pak opakovaně spouštěn podle nastavení cron v konfiguračním souboru.

Semafor

Semafor je synchronizační primitivum, které umožňuje řídit přístup k sdíleným prostředkům v multithreadovém prostředí. Semafor může být binární (0 nebo 1) nebo počítací (libovolné celé číslo). Následně má dvě hlavní operace:

- `V()` – uvolnění semaforu, zvýšení hodnoty semaforu o 1.
- `P()` – zablokování semaforu, snížení hodnoty semaforu o 1, pokud je hodnota semaforu 0, vlákno se zablokuje a čeká na uvolnění semaforu.

Tyto dvě operace jsou atomické, což znamená, že jsou prováděny jako jedna operace a nelze je přerušit. V Javě je semafor implementován pomocí třídy `Semaphore` z balíčku `java.util.concurrent`. `P` a `V` operace jsou implementovány jako metody `acquire()` a `release()`. Více o semaforech a jeho použití je popsáno v [12].

7.5 Získávání dat z API

Stahování dat z API RÚIAN je realizováno pomocí tříd `VdpClient` a `VdpDownload`, které jsou součástí modulu `download`.

7.5.1 VdpDownload

Třída `VdpDownload` se stará o nastavení HTTP klienta s důvěrou ke všem certifikátům. Využívá se zde knihovna `Apache HttpClient`, která je součástí `Spring Boot`. V rámci metody `init()` se konfiguruje `SSLContext`, časové limity a případně HTTP proxy. Výsledný klient je uložen do proměnné `client`, která se následně používá pro volání metod `tryGet()` a `trySaveFilter()`. Tyto metody slouží ke stažení dat ze serveru VDP nebo k inicializaci filtru pomocí HTTP požadavku. Před ukončením aplikace je klient uzavřen pomocí metody označené anotací `@PreDestroy`.

7.5.2 VdpClient

`VdpClient` je hlavní komponenta pro komunikaci se službou Veřejného dálkového přístupu (VDP). Využívá celkem tři URL pro přístup k seznamům souborů a další tři URL pro stahování jednotlivých datových listů. Adresy jsou většinou nastaveny v konstantách této třídy, s výjimkou jedné, která si URL generuje dynamicky podle data předchozího dne. Pro samotné stahování dat používá instanci třídy `VdpDownload`. Třída obsahuje především tři metody, které jsou volány z dříve popsaných jobů:

- `zpracovatStatAzZsj()`
- `getListLinksObce()`
- `getAdditions()`

Každá z těchto metod nejprve stáhne textový soubor, který obsahuje seznam dostupných datových souborů. Tento seznam se následně parsuje a získané odkazy se využijí ke stažení souborů. Metody `zpracovatStatAzZsj()` a `getAdditions()` stahují pouze jeden konkrétní soubor, zatímco `getListLinksObce()` stahuje všechny soubory dostupné v seznamu.

Z tohoto důvodu byla doplněna metoda `downloadFilesFromLinks()`, která umožňuje stáhnout všechny soubory uvedené v daném seznamu. Získaná data jsou obvykle ve formátu ZIP, jsou automaticky rozbalena a předána jako `InputStream` dalším komponentám prostřednictvím funkčního rozhraní `Consumer`.

Třída rovněž implementuje opakování požadavků v případě selhání, logování a čištění dočasných souborů.

7.6 Zpracování dat

Jak bylo zmíněno výše, data jsou zpracovávána pomocí třídy `VdpClient` a pomocí funkčního rozhraní `Consumer` se předávají dalším komponentám. A to právě komponentě `VdpParser`, která se stará o parsování XML souboru v předaném `InputStream` dat a následné ukládání do databáze.

Na začátku se inicializuje `XMLStreamReader`, který se stará o parsování XML souboru. Reader je vytvořen pomocí třídy `XMLInputFactory`, která je součástí knihovny `StAX`.

Původně byl použit `DocumentBuilder`, ale ten byl nahrazen `StAX` parserem. Celé soubory si nejprve načítal do paměti a poté je parsoval, po prvcích (Node). To bylo v případě malých souborů (příklad při zpracování malých obcí). Problém nastal při zpracování velkých souborů (např. obec Praha), Soubor dosahoval velikosti přes 1 GB a bylo potřeba jej zpracovat po částech. Je ale velmi náročné dělit XML bez rozbití struktury souboru. Proto se později přešlo na `StAX` parser, který zpracovává XML po prvcích a nepotřebuje celou strukturu XML. `StAX` parser je mnohem efektivnější a umožňuje zpracovávat velké soubory bez nutnosti načítat je celé do paměti. Zároveň čte pouze události jako je začátek a konec elementu, což je pro zpracování dat dostačující. Stačí tedy nadefinovat jen název potřebných elementů a tyto elementy parsovat. U `DocumentBuilderu` bylo třeba číst všechny Listy a poté je zpracovávat.

7.6.1 Parsing dat

Jak už bylo řečeno výše, data jsou parsována ve třídě `VdpParser` s použitím `XMLStreamReader`. Zároveň ve stejném modulu se nachází i třída `VdpParserConsts`, která obsahuje konstanty pro názvy jednotlivých elementů. Při čtení XML se nejprve čte hlavička souboru. Ta obsahuje nepotřebné informace, které se ignorují. Dále se čtou jednotlivá data počínající elementem `vf:Data`. Nejprve je potřeba ale rozeznat jaký element viz 7.6.1. `XMLStreamReader` rozeznává několik eventů a u každého eventu se provádí jiná akce a získává jiná informace.

Tabulka 7.1: Eventy `XMLStreamReader`

Event	Hodnota Event	Důležitá informace	Příklad
<code>START_ELEMENT</code>	1	Název Elementu	<code><vf:Data></code>
<code>END_ELEMENT</code>	2	Název Elementu	<code></vf:Data></code>
<code>CHARACTERS</code>	4	Hodnoty	Data

Příkladem, pokud se vyskytne event `START_ELEMENT` tak z něj můžu získat název elementu a jeho atributy. Pokud se vyskytne event `CHARACTERS`, tak z něj získám text

uvnitř elementu, ale pokud se pokusím získat jméno elementu vyhodí to výjimku.

Cílem je tedy číst data v cyklu, dokud nenarazím na konec elementu `vf:Data`. V průběhu čtení se nachází další důležité elementy, které určují, co se zrovna čte za objekt. Jeden soubor může rozeznávat až 19 různých objektů, které se parsují do různých DTO objektů. Každý jeden objekt podléhá jinému seznamu objektů. Je tedy potřeba rozeznat, kdy začíná list a kdy jeden objekt 7.6.1.

Tabulka 7.2: Seznam objektů a jejich názvy

Název	List Počátek/konec	Objekt Počátek/Konec
<i>Stát</i>	<vf:Staty>	<vf:Stat>
<i>Region soudržnosti</i>	<vf:RegionySourdznosti>	<vf:RegionSourdznosti>
<i>VÚSC</i>	<vf:Vusc>	<vf:Vusc>
<i>Okres</i>	<vf:Okresy>	<vf:Okres>
<i>ORP</i>	<vf:Orp>	<vf:Orp>
<i>POU</i>	<vf:Pou>	<vf:Pou>
<i>Obce</i>	<vf:Obce>	<vf:Obec>
<i>Správní obvod</i>	<vf:SpravniObvody>	<vf:SpravniObvod>
<i>MOP</i>	<vf:Mop>	<vf:Mop>
<i>MOMC</i>	<vf:Momc>	<vf:Momc>
<i>Část obce</i>	<vf:CastiObci>	<vf:CastiObce>
<i>Katastrální území</i>	<vf:KatastralniUzemi>	<vf:KatastralniUzemi>
<i>Parcela</i>	<vf:Parcely>	<vf:Parcely>
<i>Ulice</i>	<vf:Ulice>	<vf:Ulice>
<i>Stavební objekt</i>	<vf:StavebniObjekty>	<vf:StavebniObjekt>
<i>Adresní místo</i>	<vf:AdresniMista>	<vf:AdresniMisto>
<i>ZSJ</i>	<vf:Zsj>	<vf:Zsj>
<i>VO</i>	<vf:VO>	<vf:VO>
<i>Zaniklý prvek</i>	<vf:ZaniklePrvky>	<vf:ZaniklyPrvek>

Zdrojový kód 7.2: Příklad struktury XML

```

1      <vf:Data>
2      <vf:Staty>
3          <sti:Stat>
4              <sti:Kod>...</sti:Kod>
5          </sti:Stat>
6      </vf:Staty>
7 </vf:Data>

```

Každý element je rozdělen na dvě části. Klasifikátor a název elementu. Klasifikátor je určen pro rozlišení jednotlivých elementů do jaké úrovně patří. Podle názvu elementu se určuje jaký objekt se parsuje. V příkladu 7.2 je uveden vf klasifikátor (*výměnný formát*), který je nejvyšší úrovní. Každý seznam objektů využívá určitý klasifikátor, přičemž každý objekt má přiřazený vlastní, specifický klasifikátor. Každý objekt dále obsahuje atributy, které odpovídají tomuto klasifikátoru. V případě, že je atribut cizím klíčem, používá se klasifikátor objektu, na který odkazuje.

Příkladem může být dvojice vf:Staty a sti:Stat.

Může se také vyskytnout atribut, který představuje další tabulku nebo seznam tabulek. V takovém případě se pro tyto tabulky používá klasifikátor com. Příklady klasifikátorů jsou uvedeny v příkladu 7.3.

Zdrojový kód 7.3: Příklad Klasifikátorů

```

1 <vf:Data>
2     <vf:Okresy>
3         <vf:Okres>
4             <oki:Kod>100</oki:Kod>
5             <oki:Nazev>...</oki:Nazev>
6         </vf:Okres>
7     </vf:Okresy>
8     <vf:Obce>
9         <vf:Obec>
10            <obi:Kod>...</obi:Kod>
11            <obi:Nazev>...</obi:Nazev>
12            <obi:Okres>
13                <oki:Kod>100<oki:Kod>
14            </obi:Okres>
15            <obi:MluvnickeCharakteristiky>
16                <com:Pad2>...</com:Pad2>
17                <com:Pad3>...</com:Pad3>
18            </obi:MluvnickeCharakteristiky>
19        </vf:Obec>
20    </vf:Obce>
21 </vf:Data>

```

7.6.2 Atributy Objektů

V tabulce 6.1 byly zmíněné všechny datové typy a jejich rozdíly mezi databázemi. Podle dokumentace RÚIAN VFR [1] každý atribut objektu má také svůj datový typ. Konkrétně se jedná o tyto datové typy:

- **String** – textový řetězec, který může obsahovat písmena, číslice a speciální znaky.
- **Integer** – celé číslo bez desetinné části.
- **Long** – reálné číslo s desetinnou částí.
- **Boolean** – logická hodnota, která může nabývat hodnoty true nebo false.
- **Binární data** – binární data, která budou primárně obsahovat obrázky.
- **DateTime** – datum a čas ve formátu YYYY-MM-DDTHH:MM:SS (příklad: 2007-12-03T10:15:30).
- **Kolekce** – seznam objektů nebo hodnot, které jsou uloženy v jednom atributu.

Datové typy jako String, Integer, Boolean a Long jsou standardní datové typy. DateTime je datový typ, který se používá pro ukládání data a času. Co se týče binárních dat, tak ty budou uloženy jako byte[], ale aplikace tyto data zatím nezpracovává. Ovšem zmíněné Kolekce jsou označeny všechny atributy, které obsahují více hodnot nebo cizí klíče. Kolekce budou do databáze ukládány jako **JSON**.

JSON Objekty

Proč byly zvoleny JSON objekty? Tabulky dle specifikace nebyly vhodné pro rozdělení na další tabulky a vytváření cizích klíčů. Některé kolekce totiž mohou být 1:1 nebo N:M. Proto byl zvolen JSON formát, který je vhodný pro ukládání takovýchto dat. Některé kolekce jsou uloženy jako JSON Objekt a některé jako JSON pole (Pole JSON Objektů). Pro každou JSON kolekci byla vytvořena specifická JSON mapovací metoda. Každá tato metoda funguje podobně jako XML parser. Hledá elementy, které jsou uloženy v kolekci a převede je na JSON objekt nebo pole. Všechny atributy jednotlivých JSON objektů jsou uloženy v konstantách v modulu jsonObjects.

Kolekce JSON Objekt

Tyto kolekce vždy obsahují pouze jeden JSON Objekt. Každá kolekce má svou vlastní metodu pro parsování. Jedná se o kolekce:

- **Mluvnické Charakteristiky** – metoda `readMCh()`,
- **Čísla domovní** – metoda `readCislaDomovni()`,
- **Nespravné Údaje** – metoda `readNespravneUdaje()`,

Kolekce JSON Pole

Tyto kolekce mohou obsahovat jednu nebo více Objektů. Pro tyto kolekce se používá sada dvou metod. Jedna pro kolekci a druhá pro jednotlivý objekt. Jedná se o kolekce:

- **Bonitované díly / Bonitovaný díl** – metoda `readBonitovaneDily()` a `readBonitovanyDil()`,
- **Způsoby ochrany / Způsob ochrany** – metoda `readZpusobyOchrany()` a `readZpusobOchrany()`,
- **DetailníTEA** – metoda `readDetailniTeas()` a `readDetailniTea()`,

V přírůstkových datech se nachází i další kolekce, a to **Nezjištěné údaje**. Ta je kompletně ignorována a neukládá se do databáze.

Kolekce s cizími klíči

Cizí klíče jsou uloženy jako kolekce, které obsahují pouze jeden cizí klíč. Z těchto kolekcí se tedy získává pouze Integer nebo Long. Vzhledem k tomu, že se jedná o cizí klíče, které odkazují na jiné objekty, byly vytvořeny metody pro parsování těchto cizích klíčů. Metoda `readFK()` a `readFKLong` dostane jako parametr název elementu a vrátí Integer nebo Long. Rozdělení na Integer a Long je z důvodu, že je zde objekt **Parcela**, která má primární klíč jako Long.

Geometrie

Geometrie je uložena speciální kolekce obsahující až 3 možné geometrické údaje. V základní datové sadě se nachází pouze Definiční bod specifikující střed daného objektu na mapě. V rozšířených se pak dodatečně nachází i Generalizované Hranice a Originální Hranice. Je zde i případ Definiční čáry, která je uložena jako `LineString` nebo `MultiLineString`. Tato geometrie se vyskytuje pouze u objektů **Ulice**. O problematiku parsování geometrie se stará třída `GeometryParser` v modulu `geometry`, která parsuje jednotlivé geometrické objekty. V současné práci se řeší pouze Definiční bod, který je uložen jako `Point` nebo `MultiPoint`. Generalizované Hranice a Originální Hranice jsou uloženy jako `Polygon` nebo `MultiPolygon`. Stejně jako u JSON objektů i geometrie se čte podle eventu a atributů v elementu. Názvy geometrických objektů jsou uloženy jako konstanty v modulu `geometryParserConsts` s příslušným jménem k objektu. Výstupem všech metod pro parsování geometrie je `Geometry` objekt z knihovny JTS.

7.6.3 Ukládání dat do databáze

Po úspěšném parsování jedné sady objektů (např. **Obce**) nám vznikne list Dto objektů dané sady, které se následně ukládají do databáze. Tento list se následně předává do metody `prepareAndSave()`, která se nachází v příslušné service třídě (např. `ObecService`) 7.3.3. Tato metoda se stará o přípravu a uložení dat do databáze. Metoda není použita v případě, pokud je v konfiguraci nastaveno `howToProcess` na `selected` a není v konfiguraci uvedena tabulka daného objektu.

Jak metoda `prepareAndSave()` funguje? Pro každý objekt se provede několik kontrol a úprav.

- **Bezpečnostní kontrola** – slouží k zabránění chyby nebo neplatným datům v databázi.
 - **Úprava dat** – slouží k úpravě dat podle konfigurace a podle již existujících dat v databázi.
1. (*Bezpečnostní kontrola*) Pokud objekt neobsahuje primární klíč, kontrola se zastaví a po zpracování všech objektů bude odstraněn ze seznamu.
 2. (*Úprava dat*) Když se objekt, již nachází v databázi:
 - a) Pokud je v konfiguraci u tabulky nastaveno `howToProcess` na `all`: Pokud je atribut u nového objektu `null` a v databázi je `notnull`, tak se použije hodnota z databáze.
 - b) Pokud je v konfiguraci u tabulky nastaveno `howToProcess` na `include`: Pokud je atribut u nového objektu `notnull` a v konfiguraci tento atribut není uveden, tak se použije hodnota z databáze, jinak se použije hodnota z nového objektu.
 - c) Pokud je v konfiguraci u tabulky nastaveno `howToProcess` na `exclude`: Pokud je atribut u nového objektu `notnull` a v konfiguraci tento atribut je uveden, tak se použije hodnota z databáze, jinak se použije hodnota z nového objektu.
 3. (*Úprava dat*) Když se objekt ještě nenachází v databázi:
 - a) Pokud je v konfiguraci u tabulky nastaveno `howToProcess` na `all`: Objekt se uloží do databáze tak jak je bez dodatečných úprav.
 - b) Pokud je v konfiguraci u tabulky nastaveno `howToProcess` na `include`: Objektu se dá hodnota `null` u atributů, které nejsou uvedeny v konfiguraci.

- c) Pokud je v konfiguraci u tabulky nastaveno `howToProcess` na `exclude`:
Objektu se dá hodnota `null` u atributů, které jsou uvedeny v konfiguraci.
- 4. (*Bezpečnostní kontrola*) Ověří, zda objekt obsahuje cizí klíč, ověří se jeho platnost a existence v databázi. Pokud cizí klíč neexistuje, kontrola se zastaví a objekt se přidá na list pro odstranění a po zpracování všech objektů bude odstraněn ze seznamu.

Úprava dat má dvě verze podle toho, zdali se jedná o nový nebo již existující objekt. Objekt z databáze se vybírá selekcí podle primárního klíče. Do paměti se načte objekt z databáze a podle toho se upraví nový objekt.

Během těchto kontrol se také loguje, kolik prvků z listu bylo zpracováno. Jedná se o milníky 25 %, 50 %, 75 % a 100 %.

Po úspěšném zpracování listu objektů se vypíše, kolik objektů bylo odstraněno a nebude uloženo do databáze.

Pak přijde na řadu ukládání dat do databáze. To probíhá tak, že si list objektů rozdělí na menší části podle velikosti `commitSize` z konfigurace nebo menší. Tento list se pak následně uloží do databáze pomocí metody `saveAll()`.

7.7 Po zpracování dat

Jakmile jsou data úspěšně přečtena a uložena do databáze, je potřeba provést další úkony. Ve třídě `VdpClient` u metody `unzipContent()` se po úspěšném rozbalení souboru a zpracování souboru v `try` bloku nachází `finally` blok který se postará o smazání dočasného souboru. Tento blok se vykoná i v případě že dojde k výjimce a soubor se neuloží do databáze.

Jakmile se soubor úspěšně vymaže z disku může začít nový job nebo zpracování dalšího souboru v případě že se jedná o job `InitRegionJob`. Dále po zpracování regionů se zpracují přírůstková data (job `AdditionJob`)

Testování

8.1 Rychlostní porovnání databázových systémů

Tento test byl proveden na stejném konfiguračním souboru pro všechny databázové systémy. Testovací konfigurace je nastavena na provedení stažení, zpracování a uložení Základní datové sady Stát až ZSJ. Tato datová sada byla vybrána z důvodu neměnnosti. Data v této datové sadě se mění jen velmi zřídka a je možné je stáhnout. Dále bylo nastavené, že se zpracují všechny tabulky uvedené v již zmíněné datové sadě. Geometrie bude zpracována a velikost jednotlivých commitů bude 2000. Vzhled konfiguračního souboru je uveden v příkladu 8.1 a bude konfigurován pro PostgreSQL. Testování probíhalo na stejném stroji, který měl instalované všechny databázové systémy.

Nad všemi databázovými systémy probíhaly testy 3×, aby se eliminovaly chyby způsobené jinými procesy na serveru. Každý test byl proveden na prázdné databázi. Měření začalo, když byla data připravena ke čtení z důvodu eliminace stahování dat z internetu. Konkrétně když se vypsala zpráva „Data processing started.“ a následně skončilo, když se vypsala zpráva „Data processing finished.“. Během tohoto testu byl také vyfiltrováno 1 ZSJ z důvodu neexistence cizího klíče v tabulce Katastrální území. Testování proběhlo s daty https://vdp.cuzk.gov.cz/vymenny_format/soucasna/20250331_ST_UZSZ.xml.zip Výsledky jsou uvedeny v tabulce 8.1 a mají formát HH:MM:SS.

Tabulka 8.1: Časy testů pro používané databáze

	PostgreSQL	MS SQL	Oracle
Test 1	0:35:59	0:22:46	1:04:19
Test 2	0:37:03	0:22:26	1:05:30
Test 3	0:35:53	0:21:26	1:04:57
Průměr	0:36:18	0:22:13	1:04:55

Jak je vidět v tabulce 8.1, ukázalo se, že MS SQL je nejrychlejší databázový systém pro zpracování datové sady. PostgreSQL je o přesně 14 minut pomalejší než MS

SQL a Oracle je o přesně 28 minut pomalejší než MS SQL. Je ale možné že Oracle je pomalejší z důvodu že se jedná o Express verzi, která je omezena na 2 GB RAM a 1 CPU. Dále je možné že rychlost byla omezena prostředím. Zatím co PostgreSQL a MS SQL běžely v Docker kontejneru na default nastavení, Oracle běžel přímo na hostitelském systému.

Zdrojový kód 8.1: Konfigurační soubor pro testování

```
1  {
2    "database": {
3      "type": "postgresql",
4      "url": "jdbc:postgresql://localhost:5432",
5      "dbname": "ruian",
6      "username": "postgres",
7      "password": "123"
8    },
9    "quartz": {
10     "cron": "0 0 2 * * ?",
11     "skipInitialRunStat": false,
12     "skipInitialRunRegion": true
13   },
14   "vuscCodes": {},
15   "additionalOptions": {
16     "includeGeometry": true,
17     "commitSize": 2000
18   },
19   "dataToProcess": {
20     "howToProcess": "all"
21   }
22 }
```

8.2 Testování Regionů

Vzhledem k výsledku předchozího testu byl tento test proveden pouze na MS SQL. Je to čistě z důvodu času, který by byl potřeba na provedení testu na všech databázových systémech. Tento test se zaměřil na nahrávání regionů do databáze. Regiony se vybírají na základě kódu kraje. Zpracovat všechny regiony v ČR by trvalo příliš dlouho. Proto byl vybrán kraj Karlovarský, který má kód 43. Tento kraj se řadí mezi menší kraje v ČR a jeho zpracování by mělo být rychlé. Zde bude použito několik souborů, které budou zpracovány, protože regiony jsou rozděleny do částí podle obcí. Tento test se provedl na 5 souborech, které byly staženy.

- https://vdp.cuzk.gov.cz/vymenny_format/soucasna/20250331_OB_503916_UZSZ.xml.zip
- https://vdp.cuzk.gov.cz/vymenny_format/soucasna/20250331_OB_506664_UZSZ.xml.zip
- https://vdp.cuzk.gov.cz/vymenny_format/soucasna/20250331_OB_530085_UZSZ.xml.zip
- https://vdp.cuzk.gov.cz/vymenny_format/soucasna/20250331_OB_530123_UZSZ.xml.zip
- https://vdp.cuzk.gov.cz/vymenny_format/soucasna/20250331_OB_530131_UZSZ.xml.zip
-

Test bude proveden na prázdné databázi. To znamená, že je třeba upravit konfigurační soubor. Bude třeba upravit konfigurační soubor tak, aby ignoroval atributy cizích klíčů. Kdyby se zpracovávaly cizí klíče, tak by po testu byla databáze prázdná. Všechny tabulky, vybrané v konfiguračním souboru, budou zpracovány bez cizích klíčů. Vzorový konfigurační soubor 8.2 obsahuje pouze důležité změny.

Výsledek testu bude opět čas zpracování těchto 5 souborů. Měření jednotlivých souborů bude probíhat stejně jako v předchozím testu. Výsledky jsou uvedeny v tabulce 8.2 a mají formát HH:MM:SS.

Tabulka 8.2: Časy zpracování jednotlivých souborů pro test regionů

	File 1	File 2	File 3	File 4	File 5	Celkem
Test 1	0:00:45	0:00:30	0:00:19	0:00:17	0:00:14	0:02:05
Test 2	0:00:46	0:00:29	0:00:18	0:00:17	0:00:12	0:02:02
Test 3	0:00:44	0:00:30	0:00:18	0:00:18	0:00:13	0:02:03

Výsledky ukazují, že zpracování jednotlivých souborů je velmi rychlé. Problém je s počtem souborů, které je třeba zpracovat. Jak bylo zmíněno, zde se zpracovalo pouze 5 souborů. I přes to, že byl vybrán pouze jeden kraj, a to Karlovarský, který je nejmenší kraj v ČR (mimo Prahu) tak celkový počet souborů, které se budou stahovat při plném zpracování je 134 (ke dni 26.4.2025).

Zdrojový kód 8.2: Konfigurační soubor pro test regionů

```
1  {
2    "vuscCodes": {"43": "Karlovarsky kraj"},
3    "dataToProcess": {
4      "howToProcess": "selected",
5      "tables": {
6        "obec": {
7          "howToProcess": "exclude",
8          "columns": ["okres", "pou"]
9        },
10       "castObce": {
11         "howToProcess": "exclude",
12         "columns": ["obec"]
13       },
14       "katastralniUzemi": {
15         "howToProcess": "exclude",
16         "columns": ["obec"]
17       },
18       "zsj": {
19         "howToProcess": "exclude",
20         "columns": ["katastralniuzemi"]
21       },
22       "ulice": {
23         "howToProcess": "exclude",
24         "columns": ["obec"]
25       },
26       "parcela": {
27         "howToProcess": "exclude",
28         "columns": ["katastralniuzemi"]
29       },
30       "stavebniObjekt": {
31         "howToProcess": "exclude",
32         "columns": ["momc", "castobce", "
identifikacniparcela"]
33       },
34       "adresniMisto": {
35         "howToProcess": "exclude",
36         "columns": ["stavebniobjekt", "ulice"]
37       }
38     }
39   }
40 }
```

8.3 Testování Přírůstků

Přírůstky jsou poslední test, který byl proveden. Opět byla použit MS SQL databáze. Toto testování probíhalo pouze na jednom souboru. Jedná se o soubor `https://vdp.cuzk.gov.cz/vymenny_format/soucasna/20250425_ST_ZZSZ.xml.zip`. Přírůstkové soubory jsou menší než základní datové sady. Testování probíhalo podobně jako v předchozím testu. Databáze byla prázdná a byly zpracovány pouze tabulky, které byly vybrány v konfiguračním souboru. Opět byly ignorovány cizí klíče. Konfigurační soubor je nyní nastaven na pouze přidávání zmíněných sloupců. Vzorový konfigurační soubor 8.3 obsahuje pouze důležité změny. Testování probíhalo dne 26.4.2025 a proto to aplikace stahovala data z dne 25.4.2025.

Tabulka 8.3: Časy zpracování jednoho přírůstkového souboru

Test 1	Test 2	Test 3	Průměr
0:01:39	0:01:33	0:01:32	0:01:35

Jak je vidět v tabulce 8.3, přírůstkový soubor byl zpracován velmi rychle. Pro porovnání dat uloženými v databázi, a těmi staženými je ukázka v příkladu 8.4 a 8.5. Oba příklady obsahují stejný stavební objekt a jsou ve formátu XML. Podle konfigurace v příkladu 8.3 se zpracovaly pouze vybrané sloupce. Všechny ostatní sloupce byly ignorovány a mají tedy prázdné hodnoty.

Zdrojový kód 8.3: Konfigurační soubor pro test přírůstků

```
1  {
2    "dataToProcess": {
3      "howToProcess": "selected",
4      "tables": {
5        "obec": {
6          "howToProcess": "include",
7          "columns": ["nazev", "platiod", "platido", "
geometriedefbod", "datumvzniku"]
8        },
9        "castObce": {
10         "howToProcess": "include",
11         "columns": ["nazev", "platiod", "platido", "
geometriedefbod", "mluvnickecharakteristiky"]
12       },
13       "katastralniUzemi": {
14         "howToProcess": "include",
15         "columns": ["nazev", "platiod", "platido", "
idtransakce", "globalniidnavrhuzmeny"]
16       },
17       "zsj": {
18         "howToProcess": "include",
19         "columns": ["nazev", "vymera", "datumvzniku"]
20       },
21       "ulice": {
22         "howToProcess": "include",
23         "columns": ["nazev", "geometriedefbod", "
idtransakce"]
24       },
25       "parcela": {
26         "howToProcess": "include",
27         "columns": ["vymeraparcely", "dokonceni", "
bonitovanedily"]
28       },
29       "stavebniObjekt": {
30         "howToProcess": "include",
31         "columns": ["cislodomovni", "dokonceni", "
detailnitea"]
32       }
33     }
34   }
35 }
```

Zdrojový kód 8.4: Ukázka dat ve Zdroji

```

1 <vf:StavebniObjekt gml:id="S0.150902140">
2   <soi:Kod>150902140</soi:Kod>
3   <soi:CislaDomovni>
4     <com:CisloDomovni>718</com:CisloDomovni>
5   </soi:CislaDomovni>
6   <soi:IdentifikacniParcela>
7     <pai:Id>99416819010</pai:Id>
8   </soi:IdentifikacniParcela>
9   <soi:TypStavebnihoObjektuKod>1</
soi:TypStavebnihoObjektuKod>
10  <soi:ZpusobVyuzitiKod>7</soi:ZpusobVyuzitiKod>
11  <soi:CastObce>
12    <coi:Kod>13901</coi:Kod>
13  </soi:CastObce>
14  <soi:PlatiOd>2025-04-16T00:00:00</soi:PlatiOd>
15  <soi:GlobalniIdNavrhuZmeny>4219723</
soi:GlobalniIdNavrhuZmeny>
16  <soi:IdTransakce>6647043</soi:IdTransakce>
17  <soi:IsknBudovaId>62586367010</soi:IsknBudovaId>
18  <soi:Dokonceni>2025-04-14T00:00:00</soi:Dokonceni>
19  <soi:DruhKonstrukceKod>1</soi:DruhKonstrukceKod>
20  <soi:ObestavenyProstor>1193</soi:ObestavenyProstor>
21  <soi:PocetBytu>1</soi:PocetBytu>
22  <soi:PocetPodlazi>1</soi:PocetPodlazi>
23  <soi:PodlahovaPlocha>165</soi:PodlahovaPlocha>
24  <soi:PripojeniKanalizaceKod>1</soi:PripojeniKanalizaceKod
>
25  <soi:PripojeniPlynKod>3</soi:PripojeniPlynKod>
26  <soi:PripojeniVodovodKod>1</soi:PripojeniVodovodKod>
27  <soi:VybaveniVytahemKod>2</soi:VybaveniVytahemKod>
28  <soi:ZastavenaPlocha>223</soi:ZastavenaPlocha>
29  <soi:ZpusobVytapeniKod>16</soi:ZpusobVytapeniKod>
30  <soi:Geometrie>
31    <soi:DefinicniBod>
32      <gml:Point gml:id="DS0.150902140" srsName="
urn:ogc:def:crs:EPSG::5514" srsDimension="2">
33        <gml:pos>-605872.15 -1202629.48</gml:pos>
34      </gml:Point>
35    </soi:DefinicniBod>
36  </soi:Geometrie>
37 </vf:StavebniObjekt>

```

Zdrojový kód 8.5: Ukázka dat v Databázi

```
1 <StavebniObjekt>
2   <DATA_RECORD>
3     <Kod>150,902,140</Kod>
4     <Nespravny></Nespravny>
5     <CisloDomovni>{"CisloDomovni1":"718"}</CisloDomovni>
6     <IdentifikacniParcela></IdentifikacniParcela>
7     <TypStavebnihoObjektuKod></TypStavebnihoObjektuKod>
8     <CastObce></CastObce>
9     <Momc></Momc>
10    <PlatiOd></PlatiOd>
11    <PlatiDo></PlatiDo>
12    <IdTransakce></IdTransakce>
13    <GlobalniIdNavrhuZmeny></GlobalniIdNavrhuZmeny>
14    <IsknBudovaId></IsknBudovaId>
15    <Dokonceni>2025-04-14 00:00:00.000</Dokonceni>
16    <DruhKonstrukceKod></DruhKonstrukceKod>
17    <ObestavenyProstor></ObestavenyProstor>
18    <PocetBytu></PocetBytu>
19    <PocetPodlazi></PocetPodlazi>
20    <PodlahovaPlocha></PodlahovaPlocha>
21    <PripojzeniKanalizaceKod></PripojzeniKanalizaceKod>
22    <PripojzeniPlynKod></PripojzeniPlynKod>
23    <PripojzeniVodovodKod></PripojzeniVodovodKod>
24    <VybaveniVytahemKod></VybaveniVytahemKod>
25    <ZastavenaPlocha></ZastavenaPlocha>
26    <ZpusobVytapeniKod></ZpusobVytapeniKod>
27    <ZpusobyOchrany></ZpusobyOchrany>
28    <DetailniTEA></DetailniTEA>
29    <GeometrieDefBod></GeometrieDefBod>
30    <GeometrieOriHranice></GeometrieOriHranice>
31    <NespravneUdaje></NespravneUdaje>
32  </DATA_RECORD>
33 </StavebniObjekt>
```

Budoucí rozšíření a úpravy

Na této aplikaci je stále co přidávat, vylepšovat a upravit. Proto se v následujících sekcích pokusím shrnout, co by se dalo do budoucna vylepšit a přidat.

GUI pro aplikaci

Jediný uživatelský pohled na aplikaci je v současnosti zajištěn pomocí logu, který je generován v textovém formátu. Jedná se pouze o výpis informací o prováděných úkonech, které aplikace dělá. Bylo by dobré přidat uživatelské rozhraní, které by bylo schopné zobrazit informace o prováděných úkonech v reálném čase.

GUI pro nastavení konfigurace

Při psaní konfigurace uživatelem je možné že dojde k chybě. V současnosti je možné že uživatel udělá chybu a aplikace se se nespustí. Bylo by užitečné vytvořit nějakou další aplikaci nebo nějaké GUI, při spuštění aplikace, ve kterém bude možné nastavit konfiguraci aplikace. Tímto způsobem by bylo možné uživateli pomoci s nastavením aplikace a zabránit chybám, které by mohly nastat při špatném nastavení aplikace.

Optimalizace

Aplikace je v současnosti napsána tak, že se snaží, aby byla schopná zpracovat jakákoliv data z RÚIANu. Ovšem rychlostně to není úplně ideální. Hlavním kamenem úrazu je především kontrola a porovnání nových dat s těmi, které jsou už v databázi. Bylo by dobré nějak zefektivnit tento proces, aby se obecně zrychlil celý proces zpracování dat.

Refektoring kódu

Tenhle oddíl vychází trochu z předchozího. Aplikace je napsána tak aby byla schopná zpracovat jakákoliv data z RÚIANu. Bylo by proto dobré udělat refaktoring kódu, aby se zjednodušil a zefektivnil celý proces zpracování dat. Konkrétně by se jednalo o Service a Dto třídy. Tyto třídy jsou v současnosti napsány negenericky, ale by bylo dobré je napsat genericky. Využití abstraktních tříd a rozhraní by mělo celý kód velmi zjednodušit.

Zpracování zbylých dat

V současnosti aplikace zpracovává pouze základní datové sady. To znamená zpracování zbývajících geometrických dat, které jsou v RÚIANu k dispozici. Dále zakomponovat obrázky (binární data), které se u některých objektů vyskytují. Databáze je v současnosti připravena na zpracování těchto dat, ale aplikace nemá implementovanou logiku pro jejich zpracování.

Závěr

Hlavním cílem této bakalářské práce bylo vytvořit aplikaci, která by byla schopná stahovat data podle konfigurace z RÚIANu a ukládat je do databáze. V rámci práce jsem navrhl prvotní verzi aplikace, která splňuje tento cíl, implementoval jsem ji a otestoval jsem ji na vzorových datech z RÚIANu pro všechny zadané databázové systémy.

Vytvořená aplikace je první verze aplikace se jménem **Ruian_Puller**, která má za úkol stahovat data z RÚIANu a ukládat je do databáze. Podporované databáze jsou PostgreSQL, MySQL a Oracle. Aplikace je napsaná v programovacím jazyce Java a využívá framework Spring Boot. Je možné použití plánování úloh pomocí Quartz Scheduleru. Aplikace se nastavuje podle konfiguračního souboru, který je ve formátu JSON. Nastavuje se připojení k databázi, jaké prvky se mají stahovat, či ignorovat a jak často se stahování má provádět. Aplikace je navržena jako základ pro budoucí rozšíření (viz kapitola 9).

Uživatelská příručka ---

Co je to Ruian Puller?

Tato aplikace stahuje data z <https://vdp.cuzk.cz/vdp/ruian/vymennyformat> a ukládá je do SQL databáze. Podporované databáze:

- **Microsoft SQL**
- **PostgreSQL**
- **Oracle**

Data jsou stahována ve formátu XML a převáděna do DTO (Data Transfer Object), které jsou následně ukládány do databáze pomocí JPA (Java Persistence API). Aplikace je napsána v Javě 21 a používá Spring Boot pro snadné nastavení a konfiguraci. Všechny závislosti jsou spravovány pomocí Maven. Aplikace dokáže stahovat a ukládat standardní datové sady:

- Stát až ZSJ (základní)
- Všechny kraje a obce (základní)
- Přírůstková data (základní)

Před prvním spuštěním

Je třeba mít zprovozněné SQL databáze pro ukládání dat. Databáze může běžet podle vlastních nastavení, nebo lze použít variantu s Dockerem. V adresáři database se nachází podadresáře pro každou z podporovaných databází. Každý podadresář obsahuje `docker-compose.yml` a `init.sql`, které vytvoří databázi a tabulky pro ukládání dat. Oracle databáze neobsahuje `docker-compose.yml`, protože byla vytvořena mimo Docker.

Docker MS SQL a PostgreSQL

Pro spuštění v Dockeru je zde adresář `db`, kde jsou dva podadresáře: PostgreSQL a MS SQL. V každém je příslušný `docker-compose.yml` a `init.sql`.

Pro vytvoření instancí spusťte v příslušném adresáři:

```
docker-compose up -d
```

Pro případné vyčištění:

```
docker-compose down -v
```

Oracle

Oracle databáze je vytvořena mimo Docker. Je třeba mít nainstalovanou Oracle Database 19c nebo vyšší. Pro vytvoření instance použijte `init.sql`. Chování databáze je následně stejné jako u předchozích dvou.

Aplikace

Aplikace se nachází v adresáři `app`, kde je zdrojový kód a Maven projekt. Je napsána v Javě 21 a používá Spring Boot. Závislosti jsou spravovány pomocí Maven. Projekt je nastaven na JDK 21 a Maven 3.8.6 a je spustitelný v IntelliJ IDEA nebo Eclipse.

Konfigurace

Konfigurační soubor `config.json` obsahuje nastavení připojení k databázi, seznam tabulek a další volby. Nachází se v adresáři `src/main/resources`.

Zdrojový kód A.1: Příklad konfiguračního souboru

```
{
  "database": {
    "type": "<database_type>",
    "url": "<connection_string>",
    "dbname": "<database_name>",
    "username": "<username>",
    "password": "<password>"
  },
  "quartz": {
    "cron": "0 0 2 * * ?",
    "skipInitialRunStat": false,
    "skipInitialRunRegion": true
  },
  "vuscCodes": {
    "<kod>": "<kraj_nazev>"
  },
  "additionalOptions": {
    "includeGeometry": true,
    "commitSize": 1000
  },
  "dataToProcess": {
    "howToProcess": "<all/selected>",
    "tables": {
      "<tablename>": {
        "howToProcess": "<all,include,exclude>",
        "columns": ["<column_name>", "<column_name>", ...]
      }
    }
  }
}
```

Popis jednotlivých částí

- **database** – nastavení připojení k databázi:
 - type: typ databáze (postgresql, mssql, oracle)
 - url: connection string
 - dbname, username, password
- **quartz** – plánovač stahování dat:
 - cron: výraz pro plánovač
 - skipInitialRunStat – přeskočit inicializaci Stát až ZSJ (true/false)
 - skipInitialRunRegion – přeskočit inicializaci krajů (true/false)
- **vuscCodes** – kódy krajů a obcí:

Tabulka A.1: Seznam krajů a jejich kódů

Kód	Kraj
19	Hlavní město Praha
27	Jihočeský kraj
35	Jihomoravský kraj
43	Karlovarský kraj
51	Kraj Vysočina
60	Královéhradecký kraj
78	Liberecký kraj
86	Moravskoslezský kraj
94	Olomoucký kraj
108	Pardubický kraj
116	Plzeňský kraj
124	Středočeský kraj
132	Ústecký kraj
141	Zlínský kraj

- **additionalOptions** – další volby:
 - includeGeometry – zahrnout geometrii (true/false)
 - commitSize – velikost dávky pro commit (default 1000)
- **dataToProcess** – nastavení pro zpracování dat:
 - howToProcess – jak zpracovat data (all/selected)
 - tables – nastavení pro jednotlivé tabulky:
 - * howToProcess – jak zpracovat tabulku (all/include/exclude)
 - * columns – seznam sloupců pro zpracování

Tabulka A.2: Seznam tabulek a sloupců

Tabulka	Sloupce
stat	nazev, nespravny, platiod, platido, idtransakce, globalniidnavrhuzmeny, nutsiau, geometriedefbod, geometriegenhranice, geometrieorihranice, nespravneudaje, datumvzniku
regionSoudrznosti	nazev, nespravny, stat, platiod, platido, idtransakce, globalniidnavrhuzmeny, nutsiau, geometriedefbod, geometriegenhranice, geometrieorihranice, nespravneudaje, datumvzniku
vusc	nazev, nespravny, regionsoudrznosti, platiod, platido, idtransakce, globalniidnavrhuzmeny, nutsiau, geometriedefbod, geometriegenhranice, geometrieorihranice, nespravneudaje, datumvzniku
okres	nazev, nespravny, kraj, vusc, platiod, platido, idtransakce, globalniidnavrhuzmeny, nutsiau, geometriedefbod, geometriegenhranice, geometrieorihranice, nespravneudaje, datumvzniku
orp	nazev, nespravny, spravniobeckod, vusc, okres, platiod, platido, idtransakce, globalniidnavrhuzmeny, geometriedefbod, geometriegenhranice, geometrieorihranice, nespravneudaje, datumvzniku

Tabulka	Sloupce
pou	nazev, nespravny, spravniobeckod, orp, platiod, platido, idtransakce, globalniidnavrhuzmeny, geometriedefbod, geometriegenhranice, geometrieorihranice, nespravneudaje, datumvzniku
obec	nazev, nespravny, statuskod, okres, pou, platiod, platido, idtransakce, globalniidnavrhuzmeny, mluvnicecharakteristiky, vlajkatext, vlajkaobrazek, znaktex, znakovobrazek, clenenisismrozsahkod, clenenisismtypkod, nutslau, geometriedefbod, geometriegenhranice, geometrieorihranice, nespravneudaje, datumvzniku
castObce	nazev, nespravny, obec, platiod, platido, idtransakce, globalniidnavrhuzmeny, mluvnicecharakteristiky, geometriedefbod, nespravneudaje, datumvzniku
mop	nazev, nespravny, obec, platiod, platido, idtransakce, globalniidnavrhuzmeny, geometriedefbod, geometrieorihranice, nespravneudaje, datumvzniku
spravniObvod	nazev, nespravny, spravnimomckod, obec, platiod, platido, idtransakce, globalniidnavrhuzmeny, geometriedefbod, geometrieorihranice, nespravneudaje, datumvzniku
momc	nazev, nespravny, spravniobvod, mop, obec, spravniobvod, platiod, platido, idtransakce, globalniidnavrhuzmeny, vlajkatext, vlajkaobrazek, znaktex, znakovobrazek, mluvnicecharakteristiky, geometriedefbod, geometrieorihranice, nespravneudaje, datumvzniku
katastralniUzemi	nazev, nespravny, existujedigitalnimapa, obec, platiod, platido, idtransakce, globalniidnavrhuzmeny, rizeniid, mluvnicecharakteristiky, geometriedefbod, geometriegenhranice, nespravneudaje, datumvzniku

Tabulka	Sloupce
parcela	nespravny, kmenovecisko, pododdelenicisko, vymeraparcely, zpusobyvyuzitipozemku, druhcislovanikod, druhpozemkukod, katastralniuzemi, platiod, platido, idtransakce, rizeniid, bonitovanedily, zpusobyochranypoziemku, geometriedefbod, geometrieorihranice, nespravneudaje
ulice	nazev, nespravny, obec, platiod, platido, idtransakce, globalniidnavrhuzmeny, geometriedefbod, geometriedefcara, nespravneudaje
stavebniObjekt	cislodomovni, identifikacniparcely, typstavebnihoobjektukod, castobce, momc, platiod, platido, idtransakce, globalniidnavrhuzmeny, isknbudovaid, dokonceni, druhkonstrukcekod, obestavenyprostor, pocetbytu, pocetpodlazi, podlahovaplocha, pripojenikanalizacekod, pripojeniplynkod, pripojenivodovodkod, vybavenivytahemkod, zastavenaplocha, zpusobvytapanikod, zpusobyochrany, detailnitea, geometriedefbod, geometrieorihranice, nespravneudaje
adresniMisto	nespravny, cislodomovni, cisloorientacni, cisloorientacnipismo, psc, stavebniobjekt, ulice, vokod, platiod, platido, idtransakce, globalniidnavrhuzmeny, geometriedefbod, nespravneudaje
zsj	nazev, nespravny, katastralniuzemi, platiod, platido, idtransakce, globalniidnavrhuzmeny, mluvnicecharakteristiky, vymera, charakterzsjkod, geometriedefbod, geometrieigenhranice, geometrieorihranice, nespravneudaje, datumvzniku
vo	platiod, platido, idtransakce, globalniidnavrhuzmeny, geometriedefbod, geometrieigenhranice, geometrieorihranice, nespravneudaje, cislo, nespravny, obec, momc, poznamka
zaniklyPrvek	typprvkukod, idtransakce

Poznámky:

1. Quartz Scheduler obsahuje dvě boolean nastavení pro přeskočení kroků inicializace.
2. Všechny tabulky kromě vo byly testovány. Je tedy možné, že při parsování dojde k chybě.

Průběh aplikace

1. Aplikace načte nastavení z `config.json`.
2. Pokud je `skipInitialRunStat = false`, stáhne se Stát až ZSJ.
3. Pokud je `skipInitialRunRegion = false`, stáhnou se vybrané kraje.
4. Následně se stahují přírůstková data.
5. Aplikace čeká na další běh dle Quartz Scheduleru.

Aplikace nepřepisuje data, pouze přidává nové nebo aktualizuje existující záznamy. Pokud nastane chyba typu Foreign Key, problémový objekt se přeskočí a pokračuje se dál.

Struktura přiloženého — zip souboru

Bibliografie

1. CZBAP-127. VFR. 2023. Dostupné také z: [https://cuzk.gov.cz/ruian/Poskytovani-udaju-ISUI-RUIAN-VDP/Vymenny-format-RUIAN-\(VFR\)/DL058RR2-v5-0-Struktura-a-popis-VFR_final.aspx](https://cuzk.gov.cz/ruian/Poskytovani-udaju-ISUI-RUIAN-VDP/Vymenny-format-RUIAN-(VFR)/DL058RR2-v5-0-Struktura-a-popis-VFR_final.aspx).
2. MICROSOFT. Microsoft SQL Server. *Microsoft Documentation*. 2025. Dostupné také z: <https://docs.microsoft.com/en-us/sql/sql-server/>. Accessed: 2025-05-22.
3. GROUP, PostgreSQL Global Development. PostgreSQL. *PostgreSQL Documentation*. 2025. Dostupné také z: <https://www.postgresql.org/docs/>. 2025-02-20.
4. ORACLE. Oracle Database. *Oracle Documentation*. 2025. Dostupné také z: <https://docs.oracle.com/en/database/>.
5. COMMUNITY, Cisco. *XML vs JSON vs YAML*. 2022. Dostupné také z: <https://community.cisco.com/t5/devnet-general-knowledge-base/xml-vs-json-vs-yaml/ta-p/4729758>. Accessed: 2022-11-29.
6. FIELDING, Roy T. *Architectural Styles and the Design of Network-based Software Architectures*. 2000. Dostupné také z: https://www.ics.uci.edu/~fielding/pubs/dissertation/rest_arch_style.htm.
7. SPRING. *Spring Framework Documentation*. 2025. Dostupné také z: <https://spring.io/projects/spring-framework>.
8. FOUNDATION, Apache Software. *What is Maven?* 2025. Dostupné také z: <https://maven.apache.org/what-is-maven.html>. Accessed: 2025-04-19.
9. DOCKER. *What is Docker?* 2025. Dostupné také z: <https://docs.docker.com/>.
10. DBEAVER. *DBeaver Documentation*. 2023. Dostupné také z: <https://dbeaver.com/docs/dbeaver/>.
11. FOUNDATION, Apache Software. *Apache Log4j*. 2025. Dostupné také z: <https://logging.apache.org/log4j/2.12.x/>. Accessed: 2021-12-28.

12. PEŠIČKA, Ladislav. *06. Mutexy, monitory*. 2024. Dostupné také z: <https://portal.zcu.cz/CoursewarePortlets2/DownloadDokumentu?id=16897>.

Seznam použitých zkratk

API Application Programming Interface.

DTO Data Transfer Object.

INI Initialization File.

JDBC Java Database Connectivity.

JSON JavaScript Object Notation.

JPA Java Persistence API.

JTS Java Topology Suite.

MOMC Městský obvod / městská část u statutárně členěných měst.

Mop Obvod Praha.

ORM Object-Relational Mapping.

ORP Obce s rozšířenou působností.

POU Pracovní území obce.

RÚIAN Registr územní identifikace adres a nemovitostí.

SQL Structured Query Language.

VO Volební okrsek.

VFR Veřejná funkční registrace.

VDP Veřejná datová služba.

VÚSC Vymezené územní samosprávné celky.

XML Extensible Markup Language.

YAML YAML Ain't Markup Language.

ZSJ Základní sídelní jednotka.

Seznam obrázků

2.1	Schéma prvků RÚIAN	9
4.1	Architektura aplikace	15
7.1	Pořadí jobů	35

Seznam tabulek

6.1	Datové typy v různých databázích	26
7.1	Eventy XMLStreamReader	37
7.2	Seznam objektů a jejich názvy	38
8.1	Časy testů pro používané databáze	45
8.2	Časy zpracování jednotlivých souborů pro test regionů	47
8.3	Časy zpracování jednoho přírůstkového souboru	49
A.1	Seznam krajů a jejich kódů	60
A.2	Seznam tabulek a sloupců	61

1101001
10101100001110010 1100001
101011010101 10



11010011101101001
0110000110101
111000101011101