

Univerzita Jana Evangelisty Purkyně
v Ústí nad Labem
Přírodovědecká fakulta



Využití open-source a komerčních nástrojů pro
vizualizaci a analýzu dat na datové platformě
Portabo

BAKALÁŘSKÁ PRÁCE

Vypracoval: Ladislav Tahal

Vedoucí práce: Ing. Roman Vaibar, Ph.D., MBA

Studijní program: Aplikovaná informatika

Studijní obor:

ÚSTÍ NAD LABEM 2025

Namísto žlutých stránek vložte digitálně podepsané zadání kvalifikační práce poskytnuté vedoucím katedry.

Zadání musí zaujímat právě dvě strany.

Zadání je nutno vložit jako PDF pomocí některého nástroje, který umožňuje editaci dokumentů (se zachováním elektronického podpisu).

V Linuxe lze například použít příkaz `pdftk`.

Prohlášení

Prohlašuji, že jsem tuto bakalářskou práci vypracoval samostatně a použil jen pramenů, které cituji a uvádím v přiloženém seznamu literatury.

Byl jsem seznámen s tím, že se na moji práci vztahují práva a povinnosti vyplývající ze zákona č. 121/2000 Sb., ve znění zákona č. 81/2005 Sb., autorský zákon, zejména se skutečností, že Univerzita Jana Evangelisty Purkyně v Ústí nad Labem má právo na uzavření licenční smlouvy o užití této práce jako školního díla podle § 60 odst. 1 autorského zákona, a s tím, že pokud dojde k užití této práce mnou nebo bude poskytnuta licence o užití jinému subjektu, je Univerzita Jana Evangelisty Purkyně v Ústí nad Labem oprávněna ode mne požadovat přiměřený příspěvek na úhradu nákladu, které na vytvoření díla vynaložila, a to podle okolností až do jejich skutečné výše.

V Ústí nad Labem dne 24. října 2025

Podpis:

Děkuji vedoucímu bakalářské práce Ing. Romanu Vaibarovi, Ph.D., MBA za jeho odborné vedení a praktické podněty, které výrazně přispěly k analýze dat, hledání efektivních řešení a k celkovému úspěšnému dokončení bakalářské práce.

Abstrakt:

Bakalářská práce se zabývá srovnáním open-source a komerčních nástrojů pro vizualizaci a analýzu dat, datové sklady a OLAP technologie s cílem zhodnotit jejich vhodnost pro využití v datové platformě Portabo. V praktické části byla realizována implementace několika variant datových řešení, zahrnujících kombinace databázových systémů MSSQL, MariaDB, ClickHouse a PostgreSQL s nástroji Power BI, Superset a Cube.js. Bylo provedeno testování výkonu, analýza technických požadavků, nároků na uživatelské dovednosti a srovnání ekonomických aspektů provozu. Výsledkem práce je komplexní přehled výhod a nevýhod jednotlivých přístupů, měření jejich efektivity při zpracování velkých objemů dat a doporučení optimálního řešení pro organizace usilující o efektivní datovou analytiku bez nutnosti vysoké IT expertizy.

Klíčová slova: Business Intelligence, datové sklady, olap, vizualizace dat, Portabo

UTILIZATION OF OPEN-SOURCE AND COMMERCIAL TOOLS FOR DATA VISUALIZATION AND ANALYSIS ON THE PORTABO DATA PLATFORM

Abstract:

This bachelor's thesis focuses on the comparison of open-source and commercial tools for data visualization and analysis, data warehouses, and OLAP technologies, with the aim of evaluating their applicability within the Portabo data platform. In the practical part, several data architecture variants were implemented, combining database systems such as MSSQL, MariaDB, ClickHouse, and PostgreSQL with visualization tools including Power BI, Superset, and Cube.js. The analysis covers performance testing, technical requirements, user skill demands, and economic aspects of operation. The outcome of the thesis is a comprehensive overview of the advantages and limitations of both open-source and commercial solutions, performance evaluation on large data volumes, and recommendations for organizations seeking efficient data analytics without requiring extensive IT expertise.

Keywords: Business Intelligence, data warehouses, OLAP, data visualization, Portabo

Obsah

| | |
|---|-----------|
| Úvod | 13 |
| 1. Přehled současného stavu problematiky | 15 |
| 1.1. Rešerše v oblasti datových skladů, OLAP technologií a nástrojů pro vizualizaci dat | 15 |
| 1.2. Přehled současných open-source a komerčních ekosystémů | 16 |
| 1.3. Analýza trendů v implementaci | 20 |
| 2. Teoretická část | 23 |
| 2.1. Úvod do problematiky datových skladů a OLAP | 23 |
| 2.2. Nástroje Business Intelligence | 29 |
| 2.3. Přehled a charakteristika využitých technologií | 29 |
| 3. Praktická část | 31 |
| 3.1. Návrh metodiky porovnání | 31 |
| 3.2. ETL procesy | 34 |
| 3.3. Příprava dat a návrh datového skladu | 36 |
| 3.4. OLAP a vizualizace | 37 |
| 3.5. Provedení srovnávací analýzy | 39 |
| 3.6. Zhodnocení výsledků | 40 |
| 4. Sazba ukázek kódu | 41 |
| 5. Citace | 43 |
| 5.1. Označování citací | 44 |
| 5.2. Bibliografický záznam | 45 |
| 5.3. Často kladené otázky | 49 |
| 6. Zhodnocení | 53 |
| 7. Závěr | 55 |
| A. Externí přílohy | 59 |
| B. Další přílohy | 61 |

Úvod

V současné době organizace generují a shromažďují obrovské množství dat, která se stávají klíčovým zdrojem pro strategické i operativní rozhodování. Společnosti napříč odvětvími se proto zaměřují na to, jak tato data efektivně zpracovávat, ukládat, analyzovat a vizualizovat tak, aby přinášela skutečnou informační hodnotu i uživatelům bez hlubokého technického zázemí. S rozvojem datových technologií vzniká široké spektrum nástrojů, které umožňují tvorbu reportů, analytických modelů a vizualizací nad velkými objemy dat. Tyto nástroje mohou být jak komerční, nabízející komplexní podporu a integrované služby, tak open-source, které poskytují flexibilitu, otevřenost a nižší náklady na implementaci.

Zároveň však s tímto rozvojem vyvstává otázka, jaké řešení je pro konkrétní organizaci nejvhodnější – z pohledu technického, uživatelského i ekonomického. Volba správného nástroje či architektury datové platformy zásadně ovlivňuje nejen efektivitu zpracování dat, ale také dostupnost a interpretaci výsledků pro koncové uživatele.

Tato bakalářská práce se zabývá využitím open-source a komerčních nástrojů pro vizualizaci a analýzu dat na datové platformě Portabo. Cílem práce je provést srovnávací analýzu vybraných řešení z hlediska technických parametrů, uživatelských požadavků a ekonomických aspektů jejich provozu. Praktická část je zaměřena na implementaci několika variant datové architektury – zahrnujících kombinace databázových systémů MSSQL, MariaDB, ClickHouse a PostgreSQL s nástroji Power BI, Superset a Cube.js. Součástí analýzy je také měření výkonu při práci s rozsáhlým datovým souborem, testování zátěže a vyhodnocení celkových nákladů na provoz (TCO).

První kapitola práce shrnuje současný stav problematiky a základní pojmy z oblasti datových skladů, OLAP a vizualizačních nástrojů. Následující teoretická část rozebírá principy a architekturu zvolených technologií. Praktická část popisuje návrh metodiky, implementaci jednotlivých řešení a provedení testování. Závěrečná kapitola obsahuje shrnutí zjištěných výsledků, jejich interpretaci a doporučení vhodného řešení pro organizace využívající datovou platformu Portabo.

1. Přehled současného stavu problematiky

1.1. Rešerše v oblasti datových skladů, OLAP technologií a nástrojů pro vizualizaci dat

Tato kapitola přináší přehled současného stavu v oblasti zpracování, ukládání a vizualizace dat. Cílem je zmapovat, jak jsou dnes řešeny datové platformy určené pro analytické účely, jaké nástroje se používají v praxi a jaké trendy určují směr vývoje v oblasti datové analytiky.

Datové sklady a vývoj analytických databází

Datové sklady a technologie OLAP patří k dlouhodobě stabilním základům datové analytiky. Již několik desetiletí tvoří klíčovou infrastrukturu pro zpracování a konsolidaci dat z různých zdrojů, přičemž jejich architektura se neustále vyvíjí s ohledem na rostoucí objem dat a potřebu vyšší výpočetní efektivity.

Původně byly analytické systémy budovány nad relačními databázemi, jako jsou *Microsoft SQL Server*, *Oracle Database* nebo *IBM DB2*. Tyto systémy dominovaly zejména podnikovému prostředí a poskytovaly základní podporu pro tvorbu datových skladů. Postupně se však začala prosazovat specializovaná řešení určená přímo pro analytické zpracování velkých objemů dat.

Moderní analytické databáze, jako jsou *ClickHouse*, *Snowflake*, *Amazon Redshift* či *Google BigQuery*, využívají kolumnární uložení dat, paralelní zpracování a škálovatelnou cloudovou architekturu [1, 2]. Výhodou těchto systémů je možnost provádět rozsáhlé analytické dotazy v reálném čase a minimalizovat potřebu předběžných agregací. Vedle komerčních řešení se rozvíjí i open-source alternativy, které kombinují vysoký výkon s nízkými náklady na provoz – například *ClickHouse* či *Apache Druid*.

OLAP technologie a analytické přístupy

OLAP technologie zůstávají jedním z hlavních způsobů, jak efektivně analyzovat a agregovat data pro potřeby rozhodování. Tradiční přístupy založené na předpočítaných datových kostkách jsou

dnes doplňovány flexibilnějšími modely, které umožňují ad-hoc analýzy nad velkými datovými sadami.

V současnosti lze pozorovat trend integrace OLAP funkcionality přímo do datových skladů. Například platformy jako *Snowflake*, *ClickHouse* nebo *PostgreSQL* s rozšířením *TimescaleDB* umožňují analytické dotazování bez nutnosti samostatné OLAP vrstvy [1]. Výsledkem je zjednodušená architektura a nižší náklady na údržbu.

OLAP se zároveň posouvá směrem k real-time zpracování, kde je možné kombinovat streamovaná a historická data. Tento přístup umožňuje okamžitou reakci na události, což je zásadní například v oblasti IoT, průmyslového monitoringu nebo dopravní analýzy.

Nástroje pro vizualizaci a business intelligence

Vizualizační a BI nástroje představují uživatelskou vrstvu datové analytiky. Umožňují nejen tvorbu reportů a dashboardů, ale i interaktivní průzkum dat bez nutnosti znalosti programování.

Na trhu existuje široké spektrum nástrojů, od komerčních platforem po open-source řešení. Mezi nejpoužívanější komerční systémy patří *Microsoft Power BI*, *Tableau* a *Qlik Sense*. Tyto produkty nabízejí komplexní ekosystém pro tvorbu vizualizací, datových modelů a propojení s různými datovými zdroji.

Z open-source nástrojů se výrazně prosadily *Apache Superset*, *Grafana* a *Metabase*, které poskytují flexibilní možnosti přizpůsobení, automatizace a integrace. Zajímavým směrem vývoje je i koncept tzv. *sémantické vrstvy* (*semantic layer*), reprezentovaný například frameworkem *Cube.js*, který umožňuje firmám zpřístupnit datový sklad i uživatelům bez hluboké znalosti SQL. [3, 4, 5].

1.2. Přehled současných open-source a komerčních ekosystémů

V současné době existuje řada komplexních ekosystémů určených pro správu, zpracování a analýzu dat. Tyto ekosystémy zpravidla zahrnují nástroje pro integraci dat (ETL/ELT), datové sklady, OLAP vrstvy a vizualizační rozhraní. Cílem této části je poskytnout rešeršní přehled nejrozšířenějších komerčních a open-source řešení, která se využívají při budování datových platforem v podnicích i ve veřejné sféře.

Microsoft Data Platform (MS SQL, SSIS, SSAS, Power BI)

Ekosystém společnosti Microsoft představuje jedno z nejkomplexnějších a nejrozšířenějších řešení pro podnikové zpracování a analýzu dat. Jeho základ tvoří *Microsoft SQL Server*, který funguje jako relační databázový systém s podporou datového skladu. Pro procesy integrace a transformace

dat je využíván modul *SQL Server Integration Services* (SSIS), jenž umožňuje efektivní ETL/ELT zpracování z různých zdrojů [6].

Analytická nadstavba *SQL Server Analysis Services* (SSAS) poskytuje OLAP a tabulární modely, které umožňují definovat hierarchie, metriky a předpočítané agregace. Vizuální vrstvu pak tvoří *Power BI*, nástroj pro tvorbu reportů, interaktivních dashboardů a samoobslužnou analytiku [7].

Silnou stránkou tohoto ekosystému je jeho vzájemná provázanost – od správy dat přes analytické modely až po prezentaci. Díky hluboké integraci s prostředím *Microsoft Azure* lze tento systém provozovat v hybridní nebo plně cloudové architektuře (např. *Azure Synapse Analytics*, *Fabric*). Nevýhodou může být proprietární charakter a závislost na licenční politice Microsoftu.

Google Cloud Data Analytics (BigQuery, Dataflow, Looker Studio)

Společnost Google nabízí ucelený cloudový ekosystém zaměřený na masivně paralelní zpracování dat a škálovatelnou analytiku. Jeho jádro tvoří *Google BigQuery*, kolumnární databázový systém optimalizovaný pro analytické dotazy nad rozsáhlými datovými sadami [8].

Datové toky je možné zpracovávat pomocí nástrojů *Dataflow* (pro stream a batch processing) a *Dataprep* (pro datové čištění). Vizualizační vrstvu představuje *Looker Studio* (dříve *Google Data Studio*), které umožňuje vytvářet interaktivní dashboardy a reporty propojené přímo s BigQuery či dalšími zdroji.

Ekosystém Google se vyznačuje vysokou úrovní automatizace, integrací s nástroji pro strojové učení (např. *Vertex AI*) a dostupností v rámci pay-as-you-go modelu. Na druhé straně je určen převážně pro cloudové prostředí, což může být limitující pro organizace preferující on-premise infrastrukturu.

Amazon Web Services (AWS Analytics Stack)

Dalším významným hráčem je společnost Amazon, jejíž cloudová platforma *Amazon Web Services* (AWS) nabízí rozsáhlý ekosystém pro datové inženýrství a analytiku. Klíčovou roli zde zastává *Amazon Redshift* – cloudový datový sklad postavený na kolumnární architektuře [9].

Doplňkové služby jako *AWS Glue* umožňují automatizované ETL procesy, zatímco *Amazon QuickSight* poskytuje prostředí pro interaktivní vizualizace a reporting. AWS zároveň podporuje integraci s open-source systémy, například *Apache Spark* (prostřednictvím *EMR*) nebo *Presto*.

Výhodou AWS je modularita – jednotlivé komponenty lze kombinovat podle potřeb projektu. Slabinou může být větší komplexita nastavení a vyšší provozní náklady při rozsáhlém využívání více služeb současně.

Open-source datové ekosystémy (modulární přístup)

Na rozdíl od komerčních platforem, které nabízejí uzavřená a plně integrovaná řešení, jsou open-source datové ekosystémy založeny na principu modularity. Jednotlivé komponenty – databázové systémy, nástroje pro modelování dat, semantické vrstvy a vizualizační rozhraní – lze často volně kombinovat a integrovat podle potřeb konkrétního projektu. Tento přístup umožňuje vysokou míru flexibility a přizpůsobení architektury bez závislosti na jednom dodavateli.

Typickým příkladem je kombinace relační databáze (např. *MariaDB*, *PostgreSQL*, či *ClickHouse*) používané jako datový sklad, s nadřazenou semantickou vrstvou v podobě *Cube.js* nebo *dbt Semantic Layer*. Tyto nástroje definují metriky, dimenze a vztahy mezi tabulkami a umožňují jednotnou logiku pro přístup k datům napříč celým analytickým systémem [3, 10].

Nad touto vrstvou je pak možné postavit vizualizační a analytické prostředí, které interpretuje modelovaná data do uživatelsky přívětivé podoby. V praxi se často používají nástroje jako *Apache Superset*, *Metabase* nebo *Grafana*, které se dokážou k semantické vrstvě připojit přímo pomocí SQL nebo prostřednictvím API [11, 12].

Tento modulární přístup umožňuje organizacím skládat vlastní analytickou platformu „na míru“ – například kombinaci *MariaDB* pro uložení dat, *Cube.js* pro definici logiky a metrik a *Superset* pro vizualizaci výsledků.

Výhodou těchto řešení je transparentnost, možnost úprav zdrojového kódu a nezávislost na licenčních poplatcích. Slabinou může být potřeba vyšší technické znalosti při konfiguraci a údržbě systému, zejména při integraci více komponent od různých vývojářů. Přesto jsou open-source datové ekosystémy v praxi stále častěji využívány – zejména v akademickém prostředí, start-upech a menších organizacích, které preferují flexibilitu a kontrolu nad celým datovým procesem.

Hybridní a vícevrstvá řešení

Moderní datové prostředí se často neomezuje na jediný ekosystém. V praxi je běžné kombinovat open-source a komerční komponenty – například provozovat *MariaDB*, *PostgreSQL* či *ClickHouse* jako datový sklad a nad nimi využívat vizualizační nástroje typu *Power BI* nebo *Tableau*.

Tento přístup umožňuje organizacím optimalizovat náklady a současně využít výhod obou světů – otevřenost a flexibilitu open-source řešení v kombinaci s uživatelským komfortem a technickou podporou komerčních platforem.

V posledních letech lze v oblasti datové analytiky pozorovat trend tzv. *composable data stack*, tedy modulární architektury, která umožňuje skládat jednotlivé komponenty platformy podle konkrétních potřeb projektu [2]. Namísto používání jednoho monolitického systému organizace kombinují specializované nástroje – například *dbt* pro transformaci dat, *Cube.js* jako sémantickou vrstvu a *Superset* nebo *Power BI* pro vizualizaci.

Tento přístup přináší vyšší míru škálovatelnosti a nezávislosti na jednom dodavateli, což umožňuje reagovat na technologické změny a rozšiřovat platformu o nové nástroje.

Stejně jako u čistě open-source ekosystémů však i hybridní architektury narážejí na problém vyšší technické složitosti. Integrace komponent z odlišných prostředí vyžaduje nejen znalost různých technologií, ale i pochopení jejich odlišných licenčních modelů, způsobů autentizace a řízení přístupu. Zatímco open-source řešení bývají náročná na konfiguraci a údržbu, hybridní přístup přidává další vrstvu komplexity v podobě nutnosti zajištění kompatibility mezi komerčními a otevřenými systémy.

Pro úspěšnou implementaci hybridního modelu je proto klíčová standardizace rozhraní, pečlivé řízení verzí a využití integračních nástrojů, které umožňují monitorovat a spravovat datové toky napříč prostředím. Pokud jsou tyto výzvy zvládnuty, může hybridní architektura představovat efektivní kompromis mezi nákladovou efektivitou open-source řešení a robustní podporou komerčních platforem.

1.3. Analýza trendů v implementaci

V posledních letech dochází v oblasti implementace datových platforem a systémů pro analýzu dat k významným změnám. Tyto změny jsou důsledkem rostoucích požadavků na rychlost zpracování, automatizaci a dostupnost analytických nástrojů i mimo oblast IT. Současný vývoj ukazuje posun od uzavřených, monolitických řešení směrem k otevřeným, modulárním a škálovatelným architekturám, které umožňují flexibilnější integraci technologií a snadnější rozšiřování podle potřeb organizace [13, 14].

Automatizace a ELT přístup

Jedním z hlavních trendů posledních let je přechod od tradičního přístupu *ETL* (Extract–Transform–Load) k modernějšímu konceptu *ELT* (Extract–Load–Transform). Transformace dat se tak přesouvá z úrovně samostatných procesů do samotného datového skladu, který disponuje dostatečným výkonem pro jejich zpracování. Tento přístup zjednodušuje datové toky, zvyšuje jejich transparentnost a umožňuje verzování datových modelů. Rozšířené je využití nástrojů jako *dbt*, *Apache Airflow* nebo *Fivetran*, které podporují automatizované a reprodukovatelné datové pipeline.

Cloudová infrastruktura a škálovatelnost

Cloudové technologie se staly standardem pro moderní datové platformy díky své flexibilitě, dostupnosti a ekonomickému modelu *pay-as-you-go*. Platformy jako *Snowflake*, *Google BigQuery*, *Amazon Redshift* či *Azure Synapse* umožňují dynamicky přizpůsobovat výkon podle aktuálních potřeb. Open-source komunita reaguje podobně – vznikají cloudové varianty nástrojů, například *ClickHouse Cloud* nebo *Timescale Cloud*, které kombinují otevřenost s jednoduchostí správy.

Roste i využívání tzv. *hybridních cloudových modelů*, které kombinují lokální a cloudové komponenty. Tento přístup umožňuje citlivá data uchovávat v on-premise prostředí, zatímco výpočetně náročné úlohy lze provádět v cloudu.

Datová governance a kvalita dat

S rostoucím objemem dat získává na významu jejich správa, kvalita a dohledatelnost. Moderní přístupy kladou důraz na principy *data governance*, které zahrnují řízení přístupových práv, dokumentaci datových toků (*data lineage*) a sledování kvality dat. K rozšířeným nástrojům patří například *Apache Atlas*, *Amundsen* či *DataHub*, které umožňují centralizovanou správu metadat. V komerční sféře pak obdobné funkce zajišťují systémy jako *Microsoft Purview* nebo *Collibra*. Implementace těchto principů je klíčová nejen z hlediska efektivity, ale i souladu s legislativními požadavky, například GDPR nebo ISO 27001.

LLM

V posledních letech lze zároveň pozorovat výrazný rozvoj velkých jazykových modelů (LLM, *Large Language Models*) a jejich integraci do analytických platform prostřednictvím aplikačních rozhraní (API). Tento trend se nachází ve fázi intenzivního vývoje, avšak již nyní naznačuje značný potenciál pro zjednodušení interakce s datovými systémy.

Například společnost Google nabízí volně dostupné rozhraní *Gemini API*, které umožňuje vývoj aplikací schopných převádět přirozený jazyk uživatele do analytických dotazů. V kombinaci s knihovnamy, jako je *PyAudio* (pro zpracování hlasového vstupu) či *Pyadomd* (pro komunikaci s Microsoft SQL Server Analysis Services), lze vytvářet aplikace, které interpretují hlasové požadavky uživatele a transformují je do odpovídajících dotazů nad datovým skladem.

Díky těmto technologiím se stává reálným scénář, v němž uživatel zadá systémům požadavek v přirozeném jazyce, například: „Zjisti, kolik tun papíru bylo vyrobeno za poslední týden,“ a následně obdrží okamžitou odpověď bez nutnosti manuální práce s dashboardy či reporty. Po dosažení dostatečné spolehlivosti a bezpečnosti mohou podobná řešení zásadně ovlivnit způsob, jakým uživatelé pracují s datovými platformami a nástroji Business Intelligence.

2. Teoretická část

Tato kapitola se zaměřuje na teoretické vymezení klíčových pojmů a principů, které tvoří základ pro praktickou část práce. Cílem je popsat architekturu moderních datových platforem, jejich vrstvy a související technologie využívané pro ukládání, zpracování a analýzu dat. Pozornost je věnována především konceptům *Datového jezera*, *Datového skladu* a *OLAP* technologií, které představují jádro analytických systémů. Součástí kapitoly je také přehled nástrojů Business Intelligence a teoretické vymezení *sémantické vrstvy*, která propojuje datový sklad s prezentační částí platformy a umožňuje jednotnou interpretaci dat napříč organizací.

2.1. Úvod do problematiky datových skladů a OLAP

Datové jezero

Datové jezero (*Data Lake*) představuje logický koncept centralizovaného úložiště určeného pro uchovávání obrovského množství *surových dat* (*raw data*) v jejich nativním formátu (např. JSON, XML, binární soubory).

Ačkoliv je Datové jezero často asociováno s cloudovými službami (např. AWS S3, Azure Data Lake Storage), tento koncept lze efektivně implementovat i v lokálním (on-premise) prostředí.

V rámci řešeného projektu je Datové jezero implementováno na relační databázi (viz konfigurační soubor „docker-compose3.yml“). Přestože relační databáze vyžaduje definici schématu tabulky, je princip *schema-on-read* zachován. To je realizováno tak, že veškerá variabilní data jsou uložena v jediném sloupci (např. `payload`), který je definován jako řetězcový (textový) typ (TEXT nebo VARCHAR). Tím se fyzicky vytvoří pevné schéma pro tabulku, avšak **logický datový typ jednotlivých vnitřních prvků** dat (např. JSON) se určuje staticky / dynamicky **až v rámci transformačního (ETL) procesu** při jejich parsování a vkládání do Datového skladu. Tato volba zachovává maximální flexibilitu, ačkoliv pokročilejší implementace by mohla využít nativní datové typy pro nestrukturovaná data, jako je JSONB v PostgreSQL.

Klíčové role a cíle Datového jezera:

1. **Konsolidace a Vstupní Zóna (*Landing Zone*):** Primární funkcí je konsolidovat data z mnoha heterogenních zdrojů na jedinou platformu, sloužící jako **vstupní zóna** pro surová data před jejich dalším zpracováním.

2. **Oddělení Integrovaní Logiky:** Umožňuje **oddělit logiku připojení a sběru dat** od vlastního Datového skladu. Tím se zajišťuje, že náročné transformační procesy (ETL/ELT) v DWH nejsou bezprostředně zatíženy problémy s konektivitou, dostupností zdrojů nebo dynamickou strukturou dat.
3. **Audit a Rodokmen Dat (*Data Lineage*):** Uchovávání dat v jejich **původním (surovém) formátu** umožňuje kdykoliv ověřit, z jakých zdrojových dat byla odvozena data v Datovém skladu. To je nezbytné pro **auditní účely** a pro **reprodukcí analytických výsledků** při změnách transformačních pravidel.

Datový sklad

Datový sklad (*Data Warehouse*, DWH) je centrální, časově závislé úložiště historických i aktuálních dat. Jeho primárním účelem je podpora rozhodovacích procesů. Na rozdíl od provozních databází (OLTP) je struktura DWH optimalizována pro rychlé čtení, agregace a dotazování na velkých objemech dat.

Datové tržiště (*Data Mart*) je v korporátním prostředí definováno jako **podložka datového skladu** zaměřená na data a metriky potřebné pro specifickou obchodní oblast (např. výroba, kvalita, prodej). Jedná se o menší, tématicky specializovanou entitu, která usnadňuje reportování a analýzu pro konkrétní skupinu uživatelů.

Konceptuální návrh DWH se opírá o dvě hlavní, avšak protichůdné, metodiky:

1. **Metodika Billa Inmona (*Top-Down Approach*):** Inmonova metodika je označována jako **Top-Down (shora dolů)**, protože začíná návrhem centrálního, podnikového datového skladu (*Enterprise Data Warehouse*, EDW).
 - **Schema:** EDW je modelováno ve vysoce **normalizované** formě (typicky 3. normální forma, 3NF), což zajišťuje nízkou datovou redundanci a maximální integritu.
 - **Tok dat a tržiště:** Data jsou nejprve ETL procesem načtena do detailního, normalizovaného EDW (centrální zdroj pravdy). **Datová tržiště** se vytváří až **sekundárně** z dat v EDW a jsou denormalizovaná, aby sloužila pro rychlé reportování.
2. **Metodika Ralpha Kimballa (*Bottom-Up Approach*):** Kimballova metodika je označována jako **Bottom-Up (zdola nahoru)**, protože se zaměřuje na rychlé dodání řešení pro specifické obchodní procesy.
 - **Schema:** Využívá **dimenzionální modelování** (schéma *Hvězda* nebo *Sněhová vločka*), které je záměrně **denormalizované**. To zjednodušuje dotazování a maximalizuje výkon pro OLAP úlohy.
 - **Tok dat a tržiště:** Data jsou transformována a ukládána **přímo** do dimenzionálních modelů, které **představují Datová tržiště**. Podnikový datový sklad je pak **logickou unií** (sjednocením) těchto jednotlivých tržišť.

V praxi se však často objevují hybridní systémy kombinující prvky obou přístupů. Takovýto hybridní přístup jsem zvolil i já.

Datové modely

Základními modely používanými pro návrh datového skladu v rámci dimenzionálního modelování (Kimball) jsou:

- **Schéma Hvězda (*Star Schema*):** Jedná se o nejjednodušší a nejčastěji používaný dimenzionální model. Skládá se z centrální tabulky **Faktů** (*Fact Table*), která je obklopena několika tabulkami **Dimenzí** (*Dimension Tables*). Všechny dimenze jsou přímo napojeny na tabulku faktů, čímž vzniká struktura připomínající hvězdu. Vysoká redundance je vyvážena extrémní rychlostí dotazování.
- **Schéma Sněhová vločka (*Snowflake Schema*):** Jedná se o rozšíření schématu Hvězda, kde některé dimenze jsou **normalizovány** do několika souvisejících tabulek. Tím se snižuje redundance dat, ale na úkor zvýšení složitosti dotazování (je potřeba více JOIN operací), což může mírně zhoršit výkon.

Datové vrstvy v implementaci DWH

DWH z pravidla obsahuje:

- **Staging Area:** Slouží jako dočasné úložiště, kam jsou data přenesena pomocí ETL (Extract, transform and load). V této vrstvě se provádí **harmonizace a validace dat** před aplikací dimenzionálního modelu. Příkladem je tabulka `Stg.CameraCamea` ve Vaší implementaci (viz `bilina_kamery_lake_to_staging.py`).
- **DWH Vrstva (Dimenze a Fakta):** Hlavní vrstva organizovaná dle Kimballova modelu (tedy Datový Mart pro tuto analytickou doménu). Skládá se z:
 1. **Dimenze (*Dimensions*):** Uchovávají kontext, atributy a popisné detaily (např. `DimCity`, `DimSensor`).
 2. **Fakta (*Facts*):** Uchovávají měřitelné hodnoty a cizí klíče k dimenzím (např. `FactCameraDetection`, viz `bilina_kamery_staging_to_fact.py`).
- **Sémantická vrstva (*Semantic Layer*):** Tato vrstva nahrazuje tradiční prezentační vrstvu a slouží jako **jednotná definice obchodní logiky** a metrik nad daty z DWH. Je reprezentována nástrojem **Cube.js** (viz `docker-compose3.yml`), který zajišťuje předagregaci, cachování a standardizované výpočty. Koncovým BI nástrojům (jako je Apache Superset) pak zpřístupňuje data prostřednictvím **SQL API** nebo **GraphQL**, čímž se zjednodušuje dotazování a zvyšuje výkon.

Datové vrstvy v implementaci DWH

Základní vrstvy DWH, které odrážejí tok dat ve Vašich ETL skriptech, jsou:

- **Vrstva získávání/Staging (*Staging Area*):** Slouží jako dočasné úložiště, kam jsou data přenesena po jejich transformaci a čištění z Datového jezera. V této vrstvě se provádí **harmonizace a validace dat** před aplikací dimenzionálního modelu. Příkladem je tabulka `Stg.CameraCamea` ve Vaší implementaci (viz `bilina_kamery_lake_to_staging.py`).
- **DWH Vrstva (Dimenze a Fakta):** Hlavní vrstva organizovaná dle Kimballova modelu. Skládá se z:
 1. **Dimenze (*Dimensions*):** Uchovávají kontext, atributy a popisné detaily (např. `DimCity`, `DimSensor`).
 2. **Fakta (*Facts*):** Uchovávají měřitelné hodnoty a cizí klíče k dimenzím (např. `FactCameraDetection`).
- **Sémantická vrstva (*Semantic Layer*):** Tato vrstva nahrazuje tradiční prezentační vrstvu a slouží jako **jednotná definice obchodní logiky** a metrik nad daty z DWH. Je reprezentována nástrojem **Cube.js** (viz `docker-compose3.yml`), který zajišťuje předagregaci, cachování a standardizované výpočty. Koncovým BI nástrojům (jako je Apache Superset) pak zpřístupňuje data prostřednictvím **SQL API** nebo **GraphQL**, čímž se zjednodušuje dotazování a zvyšuje výkon.

Použité technologie pro DWH

Mezi nejpoužívanější technologie pro datové sklady patří relační databáze jako **PostgreSQL** s rozšířením **TimescaleDB** (využívané pro ukládání časových řad, jak je patrné z `pg-warehouse`) nebo **Microsoft SQL Server** (využívaný ve Vašich ETL skriptech). Pro extrémně rychlé analytické dotazy se používají sloupcově orientované databáze, jako je **ClickHouse** (viditelné v `docker-compose2.yml`).

OLAP technologie

OLAP (*Online Analytical Processing*) představuje přístup ke zpracování dat, který umožňuje provádět rychlé analytické dotazy nad velkými objemy informací. Základní myšlenkou OLAP je možnost pohledu na data z více dimenzí – typicky podle času, lokality, zařízení, typu senzoru nebo jiných analytických kritérií. Tento princip umožňuje uživatelům provádět tzv. *slice and dice* operace, tj. analyzovat data z různých perspektiv, agregovat je nebo filtrovat v reálném čase.

Rozlišují se tři tradiční typy OLAP řešení:

- **MOLAP (Multidimensional OLAP)** – pracuje s vlastní vícerozměrnou strukturou dat uloženou mimo relační databázi. Tento přístup poskytuje velmi rychlé odezvy při výpočtech a agregacích díky předpočítaným datovým strukturám (tzv. *cubes*), avšak vyžaduje rozsáhlejší přípravu a větší prostorové nároky.
- **ROLAP (Relational OLAP)** – využívá klasickou relační databázi, nad kterou jsou definovány pohledy a agregační dotazy. Tento přístup je flexibilnější, lépe škálovatelný a umožňuje

přímou práci s aktuálními daty bez nutnosti jejich předpočítávání. Moderní nástroje, jako například **Microsoft SSAS Tabular** nebo **Apache Druid**, reprezentují evoluční formu ROLAP – tzv. *in-memory tabulární modely*, které kombinují relační přístup se sloupcovým uložením dat a jazykem DAX pro definici metrik.

- **HOLAP (Hybrid OLAP)** – kombinuje výhody obou předchozích přístupů, tedy rychlé předpočítané agregace z MOLAP a flexibilitu dotazování nad detailními daty z ROLAP. Tento přístup je často využíván v moderních BI řešeních, kde se pro nejčastěji používané metriky udržují souhrnné cache.

Mezi nejpoužívanější nástroje pro implementaci OLAP vrstev patří například **Microsoft SQL Server Analysis Services (SSAS Tabular)** jako zástupce komerčního řešení, a z open-source prostředí pak zejména **Apache Kylin**, **Mondrian** (součást Pentaho) nebo **Cube.js**, který poskytuje moderní API pro analytické dotazy. Některé databázové systémy, jako například **ClickHouse**, navíc integrují funkce OLAP přímo do svého jádra, což umožňuje provádět agregace v reálném čase bez nutnosti vytvářet samostatnou analytickou vrstvu.

V současnosti je trendem přechod od klasických multidimenzionálních modelů k tabulárním řešením, která nabízejí vyšší flexibilitu, lepší integraci s nástroji pro vizualizaci dat a nižší nároky na údržbu datového modelu.

Sémantická vrstva

Sémantická vrstva (*Semantic Layer*) představuje logickou nadstavbu nad datovým skladem, jejímž hlavním cílem je sjednotit přístup k datům napříč celou organizací. Slouží jako prostředník mezi technickou strukturou databáze a uživatelskými nástroji Business Intelligence. Namísto přímé práce s databázovými tabulkami umožňuje uživatelům přistupovat k datům prostřednictvím předdefinovaných metrik, dimenzí a vztahů, které odpovídají obchodní logice organizace.

Tímto způsobem sémantická vrstva abstrahuje technické detaily a umožňuje vytváření analytických dotazů bez nutnosti znalosti SQL či fyzického modelu dat. V praxi tento koncept implementují moderní frameworky jako *Cube.js*, *dbt Semantic Layer* nebo *Looker Semantic Model*, které poskytují rozhraní pro standardizovaný přístup k datům prostřednictvím API nebo SQL proxy [3, 10, 15].

K hlavním přínosům sémantické vrstvy patří:

- **Konzistence metrik a výpočtů:** Veškeré reporty a dashboardy využívají jednotně definované ukazatele, čímž se eliminuje riziko rozdílné interpretace dat.
- **Zjednodušení analytické práce:** Uživatelé mohou tvořit vizualizace nebo reporty bez nutnosti znalosti komplexní struktury DWH.
- **Zvýšení výkonu a bezpečnosti:** Sémantická vrstva často zajišťuje cachování a řízení přístupů na úrovni obchodních objektů, což zlepšuje odezvu systému a kontrolu nad daty.

Z pohledu architektury datových platforem tak sémantická vrstva představuje klíčový prvek mezi datovým skladem a vizualizační vrstvou, který propojuje technologickou přesnost s obchodní

srozumitelností.

2.2. Nástroje Business Intelligence

Nástroje pro vizualizaci a analýzu dat

Vizualizace dat představuje klíčovou část procesu Business Intelligence (BI). Umožňuje uživatelům rychle pochopit trendy, vztahy a odchylky v datech.

K nejrozšířenějším nástrojům patří:

- **Open-source řešení:** Grafana, Apache Superset, Metabase, Redash – nabízejí integraci s většinou databázových systémů a jsou vhodné pro interaktivní dashboardy.
- **Komerční řešení:** Microsoft Power BI, Tableau, Qlik Sense – poskytují rozšířené funkce pro datovou transformaci, prediktivní analýzy a sdílení výstupů v rámci organizace.

Důležitými kritérii při výběru vizualizačního nástroje jsou:

1. Uživatelská přívětivost a intuitivní rozhraní,
2. Podpora přímého připojení na datový sklad nebo OLAP vrstvu,
3. Možnost automatizace a plánování aktualizací dat,
4. Licenční model a celkové náklady na provoz (TCO).

2.3. Přehled a charakteristika využitých technologií

3. Praktická část

Praktická část této práce se zaměřuje na návrh, implementaci a ověření metodiky pro srovnání open-source a komerčních nástrojů využitelných pro analýzu a vizualizaci dat na datové platformě Portabo. Cílem bylo vytvořit plně funkční datový sklad s podporou OLAP analýz, který umožňuje zpracovávat reálná data poskytovaná datovým centrem Ústeckého kraje a připravit prostředí pro následné výkonové a funkční porovnání vizualizačních nástrojů.

V rámci řešení byly navrženy a implementovány datové toky, které respektují jednotnou architekturu systému složeného z vrstev *datového jezera*, *datového skladu* (*staging vrstva*, *dimenze*, *fakta*), *OLAP vrstvy* a *reportingu*. Tato struktura byla zachována napříč všemi testovanými technologiemi, přičemž pro jednotlivé implementace byl použit stejný princip ETL zpracování, stejný způsob transformace a obdobné rozvržení datového modelu založeného na dimenzionálním schématu.

Data využitá pro testování pocházejí převážně z datového toku města Bílina, konkrétně z MQTT témat obsahujících telemetrické údaje o kamerových detekcích vozidel. V případě implementace nad PostgreSQL/TimescaleDB byla navíc použita širší množina dat zahrnující také Wi-Fi senzory, elektrické měřiče a další zdroje. Tyto zprávy byly zpracovány pomocí implementovaných ETL procesů a uloženy do datového skladu vytvořeného nad několika databázovými technologiemi (PostgreSQL/TimescaleDB, MariaDB, Microsoft SQL Server a ClickHouse). Každá z těchto technologií představuje odlišný přístup – od komerčního řešení po výkonné open-source systémy určené pro analytické zpracování dat.

3.1. Návrh metodiky porovnání

Cílem metodiky porovnání je zajistit objektivní a opakovatelné vyhodnocení open-source a komerčních technologií využitelných pro analýzu a vizualizaci dat v prostředí platformy Portabo. Porovnání je založeno na principech jednotného datového toku, jednotného datového modelu a shodného způsobu zpracování dat. Tím je zaručeno, že rozdíly ve výsledcích budou dány samotnými vlastnostmi testovaných technologií, nikoli odlišnostmi v implementaci.

Každá z porovnávaných variant byla hodnocena z následujících hledisek:

- **Výkonové parametry** – rychlost načítání a zpracování dat, doba odezvy analytických dotazů, efektivita indexace a agregací;
- **Funkční vlastnosti** – podpora OLAP analýz, kompatibilita s vizualizační vrstvou a možnosti rozšiřitelnosti;

- **Uživatelská přívětivost** – náročnost správy, dostupnost nástrojů a srozumitelnost konfigurace;
- **Ekonomické aspekty** – licenční politika, náklady na provoz, údržbu a škálování.

Metodika tedy určuje rámec, v němž budou jednotlivé technologie nasazeny, naplněny shodnými daty a následně testovány.

Načtení a analýza zdrojových dat

V úvodní fázi práce byla vybrána konkrétní datová doména, nad kterou bude systém postaven. V tomto případě se jedná o kamerový systém města Bílina. Cílem bylo ověřit, zda je nejdříve možné tato data efektivně zpracovávat a následně ověřit možnosti jejich dynamického zpracování v reálném čase.

Poskytnutý MariaDB SQL dump obsahoval následující sloupce: `id` (`bigint(11)`), `time` (`timestamp`), `topic` (`tinytext`) a `payload` (`text`). Tyto údaje představovaly surová telemetrická data, původně pocházející z MQTT komunikace, která byla exportována do databázového dumpu. Tento dump sloužil jako vstupní datová sada pro analýzu a testování datového toku.

[illegible]

Obrázek 3.1.: Struktura tabulky `mqttentries` obsahující zdrojová telemetrická data

Z obrázku je patrná struktura vstupních dat. Ze sloupce `topic` lze částečně odvodit, o jaký typ senzoru se jedná. Na první pohled by se mohlo zdát, že na základě názvů témat lze vytvořit logické vazby mezi jednotlivými senzory. Podrobnější analýza však ukázala, že názvy nejsou konzistentní. Například kamera z Bíliny může mít téma `/Bilina/kamery/comea/BI-TP-02`, zatímco vodoměr je označen jednoduše `/vodomery/decin`. V databázi se rovněž vyskytují senzory, které vykonávají stejnou funkci, avšak používají odlišné pojmenování. Některé dokonce uvádějí na konci názvu i měřenou veličinu, přičemž jejich JSON obsahuje pouze jedinou hodnotu.

Z těchto důvodů bylo vyhodnoceno, že vytvoření skriptu, který by dynamicky parsoval JSON struktury na základě názvů témat, není vzhledem k nekonzistenci dat reálně možné. K dosažení jednotného přístupu by bylo nezbytné sjednotit pojmenování témat na úrovni zdroje.

Následně byla provedena analýza samotné JSON struktury. Byl zvolen přístup, který seskupoval témata do logických celků podle podobnosti jejich JSON schémat ve sloupci `payload`. K tomu byl vytvořen skript `analyze_json.py`, který využíval Jaccardovu podobnost pro porovnání struktury jednotlivých zpráv.

Analýza s prahovou hodnotou 100 % shody odhalila větší množství identických struktur, avšak některé senzory se lišily pouze v několika attributech — pravděpodobně kvůli odlišným verzím zařízení. Proto byla hodnota podobnosti postupně snižována až na 50 %, čímž se podařilo identifikovat širší spektrum vzájemně příbuzných JSON struktur.

Na základě výsledků byly pro další zpracování vybrány skupiny témat označené žlutě a oranžově, které vykazovaly největší logickou podobnost.

| GroupID | NumTopics | Topics | Keys |
|---------|-----------|--|--|
| 0 | 7 | /Bilina/comea/kamero/ve/test-tom, /Bilina/comea/test-tom, /Bilina/comea/test-tom/testik, /Bilina/comea/test-tom/testik | data:object:value |
| 1 | 31 | /Bilina/kamery/comea/BI-MO-11, /Bilina/kamery/comea/BI-MO-11B, /Bilina/kamery/comea/BI-MO-11B | detectiontype:object:value, ilpc:object:value, lp:object:value |
| 2 | 1 | /cez/odecty | batchpos:object:value, batchsize:object:value, ear:object:value |
| 3 | 1 | /DCUK/dochazka/telegram/ | action:object:value, location:object_latitude:object:value |
| 4 | 4 | /Decin/comea/kamero/ve/testy, /mve/kresin/kresinovati, /mve/libochovice/test&A-k, /mve/test-lib | subjson:object_data:object:value |
| 5 | 2 | /detektory/video/brna/brnaIF, /detektory/video/brna/cyklotrasa | format:object:value, msg:object_countdata:object:value |
| 6 | 1 | /detektory/video/brna/cyklotrasa-hailo | format:object:value, msg:object_countdata:object:value |
| 7 | 1254 | /Energ/DCUK/INEPRO-Pro380-Mod/inepro1-2/ApparentPower, /Energ/DCUK/INEPRO-Pro380-M | value:value |
| 8 | 2 | /Energ/DCUK/SML133/SML133-01, /Energ/DCUK/SML133/SML133-01/act | subjson:object_3i:object:value, subjson:object_3i:object:value |
| 9 | 1 | /Energ/DCUK/SML133/SML133-01/elm | subjson:object_delta:object:value, subjson:object_delta:object:value |
| 10 | 2 | /Energ/DCUK/SML133/SML133-01/har, /Energ/DCUK/SML133/SML133-01/osc | subjson:object_delta:object:value, subjson:object_delta:object:value |
| 11 | 1 | /Energ/DCUK/SML133/SML133-01/io | subjson:object_delta:object:value, subjson:object_delta:object:value |
| 12 | 1 | /Energ/DCUK/SML133/SML133-01/set | subjson:object_baudrate:object:value, subjson:object_baudrate:object:value |
| 13 | 2 | iot-2/type/mt/id/MVELibochovice/evt/topics_update/fmt/json, iot-2/type/mt/id/MVE_Libochovice | d:object_topics:object_array:compression:object:value |
| 14 | 1 | /klimakomora/memmertC110-01 | value:scalar |
| 15 | 1 | /mve/kresin | d:object_hladina jez:object_array:value, ts:object_array:value |
| 16 | 2 | /mve/kresin/pokus, /tmp/pokus | subjson:object_data_zony:object:value |
| 17 | 3 | /mve/MVELibochovice, /mve/libochovice, /mve/patek | gen1vykcin:object_array:value, gen2vykcin:object_array:value |
| 18 | 3 | /mve/MVELenesice/NS:_DI_hl_horni, /mve/MVELenesice/TG2:_I_prutok, /mve/MVELenesice/TG2:_I_prutok | alias:object:value, archived:object:value, compou:object:value |
| 19 | 1 | /povodnova-cidla | data:object_subjson:object_battery:object:value, data:object_subjson:object_battery:object:value |
| 20 | 1 | /RVtraffics_status-1 | cas:object:value, presenter:object:value, soubor:object_array:value |
| 21 | 12 | /senzory/wifi/co2/11, /senzory/wifi/co2/12, /senzory/wifi/co2/13, /senzory/wifi/co2/14, /senzory/wifi/co2/14 | end_device_ids:object_device_id:object:value, end_device_ids:object_device_id:object:value |
| 22 | 288 | /ttndata/eui-24e124136c489089-02, /ttndata/eui-24e124136c489130-03, /ttndata/eui-24e124136c489130-03 | correlation_ids:object_array:value, end_device_id:object_value |
| 23 | 7 | /ttndata/eui-24e124713d321639-vzdalenost-05, /ttndata/eui-24e124785c462751-07, /ttndata/eui-24e124785c462751-07 | correlation_ids:object_array:value, end_device_id:object_value |
| 24 | 1 | /ttndata/eui-d4d4dae09f5cffff-prach-03 | correlation_ids:object_array:value, end_device_id:object_value |
| 25 | 14 | /ttndata/voda/eui-24e124713d167798, /ttndata/voda/eui-24e124713d321242-vzdalenost-06, /ttnapp_id:object:value, battery:object:value, dev_id:object:value | app_id:object:value, battery:object:value, dev_id:object:value |
| 26 | 169 | /udp1881/15102001, /udp1881/15102002, /udp1881/15102004, /udp1881/15102005, /udp1881/15102005 | _msgid:object:value, _session:object_id:object:value |
| 27 | 1 | /vodomery/decin | misto:object:value, stav:object:value |

Obrázek 3.2.: Výstup skriptu `analyze_json.py` – logické skupiny JSON struktur s 50 % Jaccardovou podobností

3.2. ETL procesy

V rámci praktické části byly navrženy a implementovány samostatné ETL skripty pro jednotlivé databázové technologie – SQL Server, MariaDB, TimescaleDB a kombinaci MariaDB s ClickHouse. Každý z těchto skriptů zajišťuje přenos dat z tzv. *landing zóny* (surová MQTT data) do *staging vrstvy*, kde jsou data očištěna a standardizována. Následně jsou zpracovaná data přesunuta do faktové tabulky `FactCameraDetection`, která tvoří jádro analytické vrstvy systému.

Všechny implementace sdílejí jednotný princip inkrementálního načítání dat. Každý běh ETL skriptu začíná načtením posledního zpracovaného identifikátoru z tabulky `ETL_IncrementalControl` (sloupec `LastLoadedID`). Tímto způsobem se při každém dalším běhu zpracovávají pouze nově přijaté zprávy. Do produkčního nasazení je vhodné zvážit také možnost periodického plného nahrání dat, které by zajistilo synchronizaci všech historických záznamů. Průběh každého běhu je zapisován do tabulky `ETL_RunLog`, kde jsou evidovány údaje jako název úlohy, čas spuštění, stav (`RUNNING`, `SUCCESS` či `FAILED`) a zpracovávaný MQTT topic. Tento přístup umožňuje transparentní správu, kontrolu historie načítání i snadnou detekci případných chyb.

Data jsou zpracovávána v dávkách (tzv. *batchích*), jejichž velikost je řízena parametrem `BATCH_SIZE`. Každý skript používá databázové kurzory a příkazy `fetchmany()`, aby se minimalizovalo zatížení paměti. V rámci každé dávky jsou načteny nové záznamy z tabulky `mqttentries`, která obsahuje surové MQTT zprávy. Klíčovým polem v těchto datech je `payload`, obsahující JSON strukturu s detaily detekce. U většiny implementací (SQL Server, MariaDB a MariaDB+ClickHouse) je z názvu MQTT topicu navíc extrahováno město pomocí funkce `parse_city_from_topic`. Implementace

pro TimescaleDB tuto logiku nevyužívá, jelikož její hlavní cíl spočívá v univerzálním zpracování a automatické adaptaci struktury dat bez vazby na konkrétní topologii MQTT témat. !! TOTO DODĚLAT !! Upravíme dle úpravy scriptu. Pravděpodobně udělám dva. Jeden na bílinu a druhý na zbytek.

Následně je JSON dekodován a rozparsován na jednotlivé atributy – typ detekce, registrační značka, senzor, třída vozidla, rychlost a další hodnoty. Při zpracování dat jsou prováděny převody datových formátů (například převod časových značek nebo číselných hodnot z textu). Skripty zároveň ošetřují nekorektní nebo chybějící hodnoty – NULL, prázdné řetězce či nevalidní JSON. Tyto hodnoty jsou nahrazovány výchozími konstantami, případně je záznam přeskočen. Každý skript je obalen konstrukcí `try/except`, která zajišťuje zachycení výjimek a zapsání chybového stavu do logu běhu. Všechny transakce jsou explicitně potvrzovány příkazem `commit()`, což zabraňuje částečnému načtení dat v případě přerušení.

SQL Server (`bilina_kamery_lake_to_staging.py`)

ETL proces pro SQL Server využívá knihovnu `pyodbc` a ODBC ovladač. Po připojení k databázím se načte hodnota `LastLoadedID`, poté se založí záznam o běhu v tabulce `ETL_RunLog`, přičemž identifikátor běhu (`RunID`) je okamžitě vrácen pomocí klauzule `OUTPUT inserted.RunID`. Následně jsou po dávkách (10 000 záznamů) načítána data z tabulky `mqttentries`. Každý záznam je rozparsován z JSON formátu, doplněn o město extrahované z MQTT topicu a vložen do tabulky `[Stg].[CameraCamea]`. Po každé úspěšné dávce je aktualizován `LastLoadedID` v tabulce kontrolních hodnot. Tento skript je doplněn o ošetření výjimek a zápis chyb do logu běhu.

MariaDB (`maria_bilina_kamery_lake_to_staging.py`)

Verze pro MariaDB je implementována pomocí knihovny `pymysql`. Logika odpovídá variantě pro SQL Server, avšak využívá syntaxi MySQL. Každý běh je logován do tabulky `ETL_RunLog` a potvrzován metodou `commit()`. Skript načítá data po dávkách o velikosti 10 000 řádků, zpracovává JSON payload, provádí datové konverze a výsledky ukládá do tabulky `Stg_CameraCamea`. Záznamy jsou validovány a doplněny o odvozené informace, například město detekce.

TimescaleDB (`pg_timescale_lake_to_staging.py`)

ETL proces pro TimescaleDB využívá pokročilou automatickou správu schématu. Pomocí funkcí `flatten_json()` a `infer_pg_type()` se JSON payload rozbaluje do jednotlivých atributů a dynamicky se vytvářejí chybějící sloupce podle zjištěných datových typů (`BOOLEAN`, `BIGINT`, `DOUBLE PRECISION`, `JSONB`, `TIMESTAMPTZ` apod.). Díky funkcím `ensure_table_and_columns()` a `ensure_new_columns()` není nutné znát strukturu předem – systém se adaptuje na příchozí data. Vkládání probíhá dávkově pomocí `execute_batch()`, čímž je dosaženo vysoké efektivity. Tento přístup se ukázal jako ideální v prostředí, kde se struktura MQTT zpráv často mění. Na

rozdíl od ostatních implementací zde není aplikována logika pro extrakci města z názvu topicu, protože tento proces se soustředí primárně na obecné zpracování různorodých JSON struktur bez závislosti na konkrétním MQTT modelu.

MariaDB + ClickHouse

(`maria_click_bilina_kamery_lake_to_staging.py`)

Nejvýkonnější varianta ETL procesu kombinuje připojení na MariaDB (landing) a ClickHouse (staging). Použitý je server-side kurzor (SSCursor), který umožňuje streamové načítání dat bez zatížení operační paměti. Data jsou přenášena po dávkách (5 000 řádků) a zapisována do tabulky `Stg_CameraCamea` v ClickHouse. Každý běh je logován přímo v tabulkách `ETL_RunLog` a `ETL_IncrementalControl` uložených v ClickHouse. Z důvodu výkonu je používán binární protokol ClickHouse, který výrazně zrychluje přenos i zápis dat. Skript zároveň zahrnuje ošetření výjimek a bezpečné ukončení v případě chyby.

| StgID | LandingID | LoadDttm | OriginalTime | City | DetectionType | Utc | LP | Sensor | VehClass | ILPC | Velocity |
|-------|-------------|-------------------------|-------------------------|--------|---------------|-------------------------|------------|-----------|----------|------|----------|
| 201 | 127,537,681 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:23.141 | A660631FA0 | BI-TP-O2 | 0 CZ | | -1 |
| 202 | 127,537,682 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:23.675 | 4C9AD46758 | BI-TP-I2 | 0 CZ | | -1 |
| 203 | 127,537,683 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:24.106 | E7293B72AB | BI-MO-O1B | 0 CZ | | -1 |
| 204 | 127,537,684 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:24.519 | 2E3848ABD5 | BI-MO-I1B | 0 CZ | | -1 |
| 205 | 127,537,685 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:24.457 | 44DB566FE6 | BI-MO-I1 | 0 CZ | | -1 |
| 206 | 127,537,686 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:25.759 | 24BA7F0631 | BI-MO-I1 | 0 CZ | | -1 |
| 207 | 127,537,687 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:25.696 | 93421DECAC | BI-TP-I2B | 0 CZ | | -1 |
| 208 | 127,537,688 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:25.904 | C038CE1C23 | BI-MO-O1B | 0 CZ | | -1 |
| 209 | 127,537,689 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:26.068 | A660631FA0 | BI-TP-O2B | 0 CZ | | -1 |
| 210 | 127,537,690 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:25.892 | A5302FB387 | BI-MO-O1 | 0 PL | | -1 |
| 211 | 127,537,691 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:27.060 | 4C9AD46758 | BI-TP-I2B | 0 CZ | | -1 |
| 212 | 127,537,692 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:28.139 | E149B5F383 | BI-TP-I2 | 0 CZ | | -1 |
| 213 | 127,537,693 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:28.487 | D3E45C666A | BI-MO-I1 | 0 CZ | | -1 |
| 214 | 127,537,694 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:28.487 | 44DB566FE6 | BI-MO-I1B | 0 CZ | | -1 |
| 215 | 127,537,695 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:28.511 | DBC725AE85 | BI-TP-I1 | 0 CZ | | -1 |
| 216 | 127,537,696 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:29.686 | C62A6E353C | BI-MO-O1B | 0 PL | | -1 |
| 217 | 127,537,697 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:29.789 | 2DCE6849A2 | BI-MO-I1 | 0 CZ | | -1 |
| 218 | 127,537,698 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:29.851 | 24BA7F0631 | BI-MO-I1B | 0 CZ | | -1 |
| 219 | 127,537,699 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:30.557 | 242ED7377F | BI-TP-I2 | 0 CZ | | -1 |
| 220 | 127,537,700 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:31.214 | DBC725AE85 | BI-TP-I1B | 0 CZ | | -1 |
| 221 | 127,537,701 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:32.207 | 65026ACFEB | BI-MO-I1 | 0 CZ | | -1 |
| 222 | 127,537,702 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:32.082 | E149B5F383 | BI-TP-I2B | 0 CZ | | -1 |
| 223 | 127,537,703 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:32.418 | 68BA3FB107 | BI-TP-I2 | 0 CZ | | -1 |
| 224 | 127,537,704 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:32.579 | D3E45C666A | BI-MO-I1B | 0 CZ | | -1 |
| 225 | 127,537,705 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:33.447 | 6E4A519807 | BI-MO-I1 | 0 CZ | | -1 |
| 226 | 127,537,706 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:34.005 | 2DCE6849A2 | BI-MO-I1B | 0 CZ | | -1 |
| 227 | 127,537,707 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:34.216 | 11B864B44C | BI-TP-I2 | 0 CZ | | -1 |
| 228 | 127,537,708 | 2025-04-04 13:42:11.000 | 2024-06-26 02:00:15.000 | Bilina | detector | 2024-06-25 07:29:33.530 | 242ED7377F | BI-TP-I2B | 0 CZ | | -1 |

Obrázek 3.3.: Staging tabulka v datovém skladu

3.3. Příprava dat a návrh datového skladu

Zdrojová data byla získávána ze systému Portabo, kde jsou ukládána jako MQTT zprávy v JSON formátu. Data mi byli poskytnuty zašifrované, abychom je mohli v rámci bakalářské práce využít. Tyto zprávy obsahují informace o detekcích vozidel – například registrační značku, typ detekce, rychlost, město, čas a senzor. Data z vodoměrů, elektroměrů a dalších senzorů. V rámci ETL procesu byla provedena tato fáze:

- **Data lake:** surová data načtená přímo z MariaDB SQLDump,
- **DW Staging:** očištěná, formátovaná a rozparsovaná data,

- **DW Faktová tabulka:** data transformovaná do hvězdicového modelu s vazbami na dimenze.

Každá implementace využívá dávkové zpracování (`BATCH_SIZE`) a kontrolu posledního zpracovaného záznamu (`LastLoadedID`), aby bylo možné ETL proces spouštět opakovaně a inkrementálně.

3.4. OLAP a vizualizace

V rámci praktické části byly nad vytvořenými datovými sklady vystavěny dvě samostatné analytické vrstvy – každá odpovídající jedné z testovaných technologií. Cílem bylo vytvořit funkční OLAP modely umožňující tvorbu interaktivních vizualizací, a současně porovnat rozdíly mezi komerčním řešením na platformě Microsoft SQL Server a open-source přístupem postaveným na nástroji Cube.js.

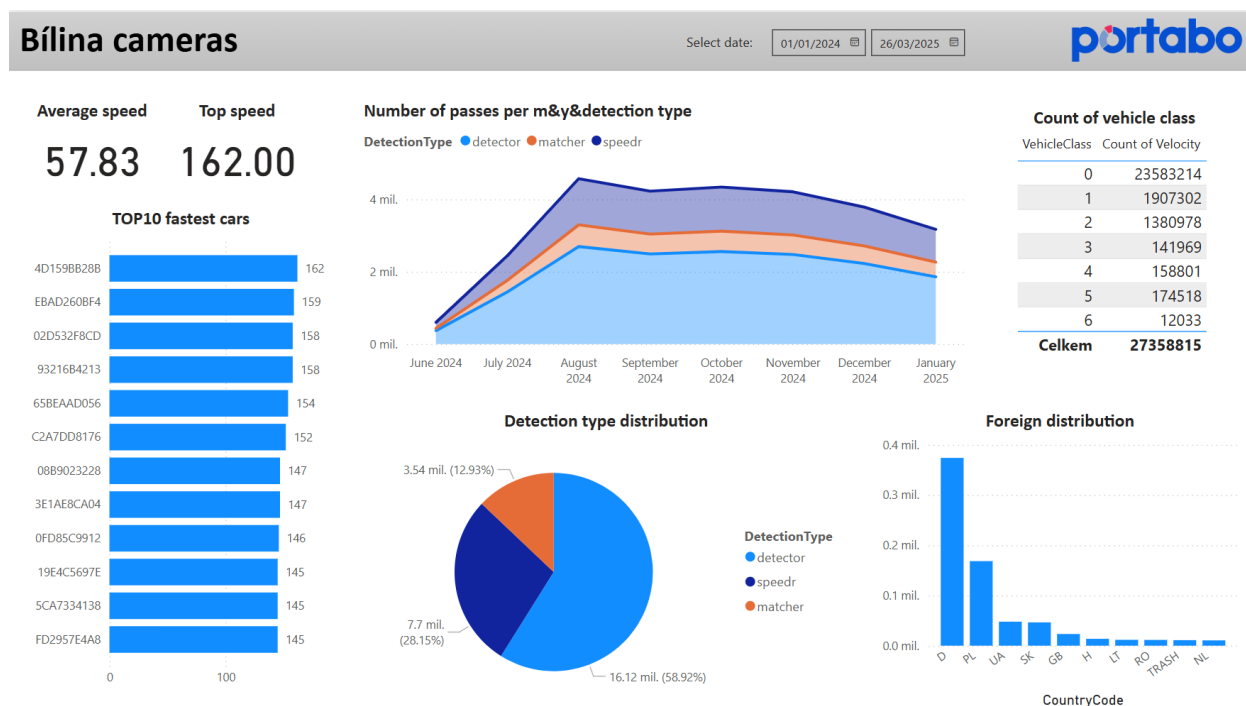
Komerční varianta – SSAS Tabular

V komerční větvi datové platformy byl vytvořen datový model v prostředí **SQL Server Analysis Services (SSAS) Tabular**. Model byl vystavěn nad tabulkami `FactCameraDetection` a dimenzemi `DimCity`, `DimSensor`, `DimVehicleClass`, `DimDetectionType`, `DimLP` a `DimTime`. Propojení tabulek bylo definováno prostřednictvím relačních vazeb podle cizích klíčů, přičemž datový model respektuje hvězdicovou topologii (*star schema*).

Po načtení dat z datového skladu byla v prostředí SSAS vytvořena sada vypočítaných metrik v jazyce DAX. Jednalo se například o:

- `Detection Count` – celkový počet detekcí,
- `Average Velocity` – průměrná rychlost vozidel (vylučující nulové hodnoty),
- `Unique Plates` – počet unikátních registračních značek,
- `Detections Over Time` – časová agregace počtu detekcí podle dne a hodiny.

Součástí modelu byla také implementace základních hierarchií (například *Rok* → *Měsíc* → *Den*) a vybraných filtračních atributů pro přehlednější práci v Power BI. Po ověření funkčnosti byl model nasazen na instanci SSAS Tabular a zpřístupněn uživatelům prostřednictvím **Microsoft Power BI** a **Excel PivotTables**. V Power BI byly vytvořeny interaktivní reporty zobrazující přehled detekcí podle města, typu vozidla, senzoru a času. Některé metriky byly navíc doplněny o logiku pro kontrolu konzistence dat a filtrování extrémních hodnot. Výhodou tohoto přístupu byla plná integrace s ekosystémem Microsoft, vysoký výkon in-memory výpočtů a možnost implementace bezpečnostních pravidel na úrovni datového modelu.



Obrázek 3.4.: Ukázka výsledného PBI reportu

Open-source varianta – Cube.js a Apache Superset

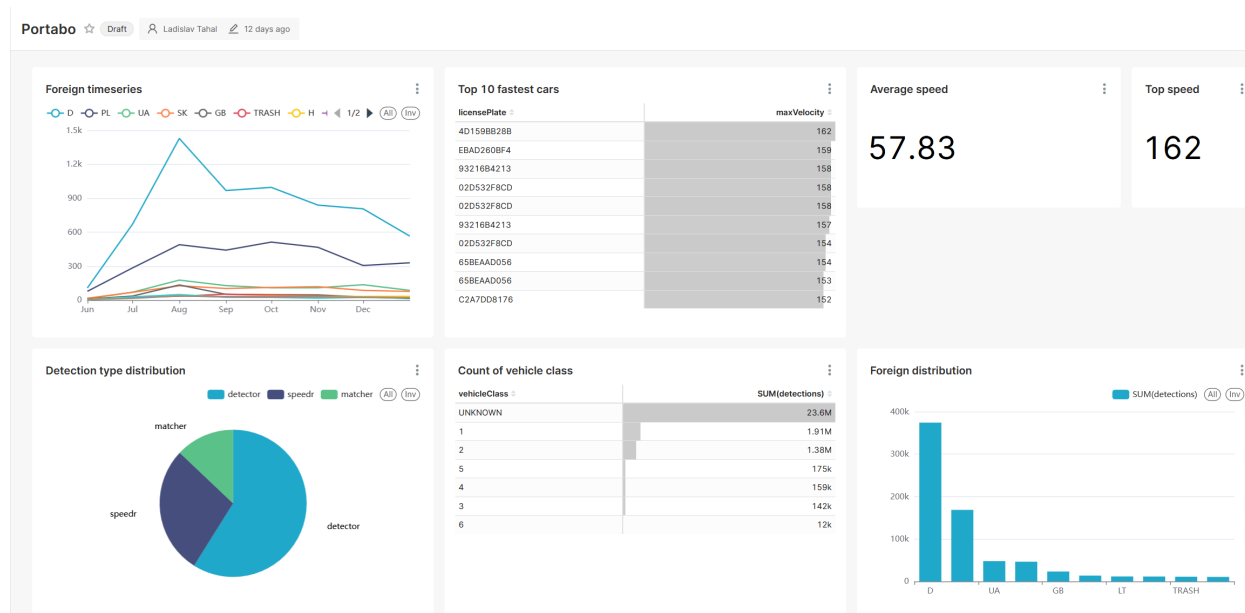
V open-source větvi datové platformy byla analytická vrstva implementována pomocí nástroje **Cube.js**, který byl nasazen nad databázemi **ClickHouse**, **MariaDB** a **PostgreSQL (TimescaleDB)**. Cube.js v tomto řešení plnil roli mezivrstvy typu OLAP, která sjednocovala přístup k datům z různých databázových systémů a poskytovala standardizované rozhraní pro vizualizaci.

Definice datového modelu byla realizována formou JavaScriptových souborů, přičemž hlavní model byl popsán v souboru `factcameradetection.js`. V tomto souboru byla definována analytická kostka `FactCameraDetection`, která obsahovala popis faktové tabulky a vazby na příslušné dimenze (`DimCity`, `DimSensor`, `DimVehicleClass`, `DimDetectionType`, `DimLP` a `DimTime`). Součástí definice byly také klíčové metriky odpovídající komerční variantě: celkový počet detekcí, průměrná rychlost, minimální a maximální rychlost, počet unikátních SPZ a počet detekcí s nenulovou rychlostí. Model využíval vazby typu `JOIN` na cizí klíče a odpovídal hvězdicovému schématu použitému v datovém skladu.

Po sestavení modelu byl **Cube.js** spuštěn jako samostatná `Node.js` služba, která poskytovala rozhraní typu **PostgreSQL API**. Toto rozhraní umožňuje klientským aplikacím komunikovat s Cube.js stejným způsobem, jako by šlo o běžnou SQL databázi – tedy pomocí SQL dotazů. Díky tomu bylo možné na Cube.js napojit nástroj **Apache Superset**, který byl použit jako hlavní vizualizační platforma open-source řešení.

V rámci Superset byly vytvořeny interaktivní dashboardy zobrazující klíčové metriky: počty detekcí v čase, rozložení typů vozidel, přehled rychlostí. Data byla načítána prostřednictvím SQL dotazů směřovaných na PostgreSQL endpoint poskytovaný Cube.js, čímž byla zajištěna kompatibilita bez

nutnosti dalších úprav. Pro zvýšení výkonu byly v Cube.js aktivovány funkce *pre-aggregations* a *query caching*, které umožňovaly ukládat výsledky často opakovaných dotazů.



Obrázek 3.5.: Ukázka výsledného Superset reportu

Porovnání přístupů

Oba přístupy poskytují plnohodnotnou OLAP analytiku, avšak liší se způsobem implementace i správy. **SSAS Tabular** umožňuje centrální řízení, pokročilé DAX výpočty a přímou integraci s Power BI, což zjednodušuje nasazení ve firemním prostředí. Naopak **Cube.js** nabízí flexibilitu, otevřenost a možnost úprav datového modelu přímo v kódu, což je výhodné zejména v agilním vývoji nebo při potřebě integrovat analytiku do webové aplikace.

Obě řešení tak umožnila plnohodnotnou vizualizaci a analýzu dat z detekčních systémů, přičemž volba konkrétní technologie závisí především na prostředí, infrastruktuře a požadavcích na rozšiřitelnost systému.

3.5. Provedení srovnávací analýzy

V této fázi byly implementované datové sklady naplněny daty a připojeny k vizualizační vrstvě. Pomocí nástroje Cube.js byla připravena sada standardizovaných dotazů, které umožní srovnání výkonu jednotlivých databází při:

- agregacích nad časem (měsíční a hodinové sumáře),
- filtrování podle města, typu detekce nebo třídy vozidla,
- výpočtu metrik (průměrná rychlost, počet detekcí, počet unikátních SPZ).

Pro samotné měření a vizualizaci výsledků bude použit software poskytnutý vedoucím práce. Tento nástroj umožní jednotným způsobem zaznamenávat časy dotazů, vytížení databází a odezvu vizualizační vrstvy.

3.6. Zhodnocení výsledků

Na základě provedené analýzy budou jednotlivé technologie porovnány z hlediska:

- **Technických parametrů** – výkon při dotazech, rychlost načítání dat, náročnost správy,
- **Uživatelské přívětivosti** – jednoduchost integrace s nástrojem Cube.js a možnost tvorby OLAP analýz,
- **Ekonomických aspektů** – licenční podmínky, náklady na provoz a škálování.

Závěrečná část shrne výhody a nevýhody jednotlivých řešení a poskytne doporučení vhodné platformy pro využití v rámci Portabo.

4. Sazba ukázek kódu

Sazba ukázek kódu je klíčová pro bakalářské práce věnované implementaci nějaké aplikace či knihovny.

Základní zásady pro sazbu ukázek kódu:

- 1.

5. Citace

Citace tvoří jeden ze základních pilířů závěrečné práce. Platí zde základní pravidlo: pokud použijete jakoukoliv zdroj informací, pak je nutné tento zdroj citovat, tj. uvést příslušný zdroj.

Zdrojem je ve většině případů text, ale může to být i obrázek, audiovizuální materiál či ve speciálních případech i ústní sdělení. V případě informatických prací je častým zdrojem u zdrojový kód.

Informaci ze zdroje můžete použít dvěma různými způsoby:

- přímo převzít (u textů je to známé Ctrl-C, Ctrl-V)
- použít jako základ vlastního intelektuálního vytvoření (textu, grafiky, programu, apod.), tj. použijete jen informaci, ale její formu změníte.

První druh tzv. přímé citace by měli mít v informatických závěrečných pracích jen velmi omezený rozsah (méně než stránka), neboť jejich přínos pro hodnocení práce je diskutabilní. Přesto jsou však případy, kdy jsou vhodné:

1. matematické definice a tvrzení (věty, axiomy)
2. definice termínů z neinformatických oborů (např. společenských věd)
3. citace norem resp. standardů

Citace mají tři základní cíle:

1. určují, co je váš vlastní intelektuální přínos a co jste pouze převzali
2. pomáhají určovat primárního autora (resp. autory)
3. definují kontext vaší práce resp. mohou usnadňovat nalezení dalších souvisejících informací

Jakákoliv vědecká práce nevzniká na zelené louce a tak jsou citace její nezbytnou součástí. V rámci práce běžně navazujeme na existující výzkumy, projekty, technologie apod. Stejně tak se můžeme odkazovat na autority či s nimi polemizovat.

První dva cíle také úzce souvisejí s plagiátorstvím. Pokud v práci použijete myšlenku či údaj bez citování je vaše práce plagiátem. I když se to v obecném mínění vztahuje jen na přímé kopírování, není tomu tak. Přímé kopírování se jen snadněji vyhledává a prokazuje. Je také obtížnější jej kvantifikovat.

V případě přímého kopírování, jež není označeno jako přímá citace, postačuje i relativně malý rozsah (například věta, jeden obrázek, jedna procedura), aby byla práce označena jako plagiát.

Plagiát není možno obhájit a v případě většího rozsahu hrozí i vyloučení ze studia. Za přímé kopírování se považují i případy, kde je změna jen formální (změna slovosledu, náhrada synonym, zkrácení, vložení textové výplně, změna barvy či afinní transformace obrázku).

V případě převzetí myšlenky jde o zjevné plagiování, pokud je tato myšlenka důležitou částí práce (podílí se na splnění cílů).

To, že není plagiátorství odhaleno před obhajobou práce, není důkazem, že se nejedná o plagiát. Pokud je plagiátorství zjištěno později, může vám být odebrán titul i zpětně (a jak jste si jistě všimli, plagiátorství je běžně využíváno v politickém boji).

V každém případě si uvědomte, že plagiátorství je druh krádeže a že ani vy nechcete, aby někdo vaše myšlenky nebo dokonce váš text vydával za vlastní.

5.1. Označování citací

Označování citace má dvě části. Za prvé je nutno označit, jaká část práce je citací (rozsah) a jaký je původní zdroj.

Zdroj je vždy určen odkazem na bibliografický záznam, které jsou v případě bakalářské práce uvedeny v kapitole *Použité zdroje* na konci práce. Odkaz může mít různý tvar, ale preferovaný styl je uvedení čísla záznamu v hranatých závorkách. V případě použití našeho latexovského stylu stačí použít příkaz `\cite{id-zaznamu}`. V případě potřeby lze zdroj zpřesnit uvedením např. stránky či kapitoly, jež se uvádí za číslem záznamu (po čárce a ještě před uzavírající hranatou závorkou). V \LaTeX u lze využít nepovinný parametr příkazu `cite`.

Označení rozsahu se poněkud liší u přímých a nepřímých citací.

U přímých citací je označení rozsahu kritické. V případě citací, které jsou kratší než odstavec je nutné text vyznačit kurzívou a zahrnout do uvozovek. Odkaz musí následovat hned za označným textem.

U citací v rozsahu odstavce či více odstavců se využívá zvětšení okrajů na levé i pravé straně odstavců (viditelné na první pohled). V \LaTeX u lze použít prostředí `quote` nebo `quotation`. Text by měl být navíc v uvozovkách. Kurzíva je možná, ale u rozsáhlejších citací není příliš vhodná. Odkaz se umísťuje na konec posledního převzatého odstavce.

Speciální případ je možný v případě matematických definic a vět. Pokud jsou v rámci kapitoly převzaty jen z jediného zdroje, lze na počátku kapitoly uvést hromadný odkaz například v podobě věty: Všechny definice a věty uvedené v této kapitole jsou převzaty z [X].

V případě převzatých obrázků se odkaz umísťuje na konec popisku. Aby však bylo zřejmé, že se jedná o přímé převzetí (kurzívu ani uvozovky nelze použít) je nutné explicitně vyjádřit, že obrázek byl převzat ze zdroje bez podstatných změn například: (převzato z [X]) nebo (překresleno z [X]).

U nepřímých citací je vyznačování rozsahu volnější. Nejjednodušší je uvádění holé citace na konci vět (před tečkou) nebo konci odstavce (za poslední tečkou). V mnoha případech je ale možné citace uvádět explicitněji a stylistiky je provázet s okolním textem.

příklady:

- zajímavá alternativa je popsána v [x]
- údaj je převzat z [x]
- použití návrhového vzoru poprvé popsal N.N v [x]
- volně přeloženo z [x]
- řešení bylo navrženo uživatelem N v [x] (vhodné např. pro stackoverflow a podobné zdroje)

Explicitnější vyjádření je nutno použít i v případě, že rozsah citace přesahuje odstavec.

- následující příklad je převzat z [x]
- výčet vychází z [x] je však doplněn o ...

Teoreticky lze podobné řešení využít i u celých sekcí či kapitol (kapitola je zpracována na základě [x]). V tomto případě je však nutné předpokládat, že v dané kapitole není žádná autorská myšlenka, a že autor se nesnažil najít alternativní pohledy či zdroje (a hodnotit tak, lze pouze autorovu schopnost výběru informací či stylistiky).

Výjimečně lze uvádět i několik citací se shodným či překrývajícím se rozsahem např. *následující specifikace je převzata z [x] a [y]*. To je však tolerovatelné jen v případě, v kdy by oddělení oddělení zdrojů bylo obtížné nebo nepřehledné a spojení nepřináší problémy s intelektuálním vlastnictvím (mají stejného autora či copyright). Zcela nepoužitelné jsou v případě většího rozsahu citace (např. na úrovni sekcí či kapitol)!

V případě obrázků je vhodné uvést explicitnější specifikaci, jak byl originální obrázek pozměněn resp. rozšířen.

Příklad:

- (převzato z [x] a doplněno)
- (převzato z [x], přeloženo)
- (upraveno z [x] pro novou verzi technologie ...)
- (inspirováno diagramem [x])
- (viz také [x] pro data X)

5.2. Bibliografický záznam

Bibliografický záznam je datová struktura, jenž má dvě základní funkce:

1. jednoznačné identifikování zdroje
2. určení primární odpovědnosti (typicky je to autor resp. autoři, u webových zdrojů to však často bývá korporace).

Pro každý typ zdrojového dokumentu (zdroje) existuje množina klíčových atributů, které by měly být specifikovány (ne zcela vhodně označované jako povinné) a další, které hrají jen pomocnou roli.

V praxi však může nastat situace, kdy není zřejmé, jaký typ dokumentu pro daný zdroj zvolit resp. nelze zjistit hodnoty klíčových atributů. V tomto případě je nutné improvizovat a snažit se, aby záznam plnil v maximální míře obě funkce.

Struktura bibliografického záznamu je v zásadě dána těmito dimenzemi:

médium – základní dělení je na tištěné dokumenty a online dokumenty (dokumenty na elektronických nosičích tvoří jakási přechod mezi oběma typy dokumentů)

samostatnost – zdroj může být samostatný nebo součást rozsáhlejšího zdroje

periodičnost – periodický dokument vychází po jednotlivých částech, přičemž počet částí není předem znám (např. časopis).

Tištěné samostatné dokumenty neperiodické

Typickým příkladem samostatného tištěného dokumentu je kniha či monografie.

Základním zdrojem informací pro bibliografický záznam u knih je tzv. tiráž, tj. soupis vydavatelských údajů uvedený na konci knihy či na stránce za titulem. Využít lze i další zdroje (např. katalogy knihoven či knižní e-shopy, bibliografické záznamy v jiných dokumentech), ale v tomto případě je nutné provádět kontrolu, neboť tyto sekundární zdroje často obsahují chyby.

klíčové atributy:

ISBN : ISBN je celosvětový jedinečný identifikátor neperiodických tištěných dokumentů. Pokud ho kniha má, pak je dokument jednoznačně identifikován (a další identifikace už hraje jen sekundární roli). Pomlčky v ISBN nejsou součástí identifikátoru a lze je vynechávat (i když občas se jedno ISBN přiděluje více svazkům). Navíc existují ve dvou podobách ISBN-10 s deseti číslicemi a ISBN-13 s třinácti. Pokud jsou k dispozici oba je vhodnější uvádět ISBN-13 (i když ISBN-10 lze snadno mapovat na ISBN-13).

název : název knihy je povinný údaj a měl by být vždy vyplněn. Použit by měl být vždy originální název bez úprav. Jedinou přípustnou úpravou je změna velkých písmen (verzálék) na malá, které by mělo odpovídat pravidlům příslušného jazyka.

podnázev : některé knihy mají i podnázev. Někdy je těžké rozeznat, co je název a podnázev. Zde platí pravidlo, že název by neměl obsahovat dvojtečku, tečku, středník apod. Od podnázvu je potřeba odlišit název edice. Podnázev je nepovinný (doporučuji uvádět pokud obsahuje klíčové informace).

autoři : v bibliografickém záznamu by měli být uvedeni všichni primární autoři (tj. není potřeba uvádět překladatele, ilustrátory, apod.)

vydání : označení konkrétního vydání. Je důležité především tehdy, když není známo ISBN a existuje více odlišných vydání (s různým obsahem)

nakladatel : uvádí se jméno nakladatelství, a to především z důvodů odpovědnosti

místo vydání : uvádí se jméno města, popřípadě stát, především tehdy pokud není jednoznačné (např. Cambridge) a to ve stručné podobě (např. stačí *United Kingdom*). Podobně stručný by měl být název nakladatelství (tj. bez označení typu společnosti, apod., rodičovské společnosti, apod.) V dnešní době globalizace je tento údaj v mnoha případech nevýznamný (tj. ho lze vynechat, především tehdy pokud je nakladatelství neznámé).

rok vydání : rok vydání přesněji identifikuje dokument. Pokud ho nelze zjistit, lze jej nahradit rokem copyrightu (v tomto případě je uvozen znakem c např. c2022)

edice : kniha může být vydána v rámci edice. Edici doporučuji neuvádět, výjimkou jsou edice, které jsou všeobecně známe.

URL : uvádí se pouze v případě, že je kniha dostupná online a to oficiálně a bez poplatků. Uvedení URL v tomto případě usnadňuje její získání (v tomto případě je ale často lepší citovat ji jako elektronickou knihu).

příklad:

Následující biblografický záznam byl získán z katalogu systému knihovny UJEP (volba Citace Pro v dolní části výpis záznamu).

RASCHKA, Sebastian a Vahid MIRJALILI. *Python machine learning: machine learning and deep learning with Python, scikit-learn, and TensorFlow*. Second edition. Birmingham: Packt, 2017. Expert insight. ISBN 978-1-78712-593-3.

Tento záznam splňuje základní požadavky, neboť obsahuje údaje týkající se odpovědnosti i jednoznačnou identifikaci dokumentu (a to jak ISBN tak přesným určením vydání). Zahrnutí podnázvu je vhodné, neboť obsahuje dodatečné informace (jména frameworků). Nakladatelství je uvedeno ve stručné podobě (tj. *Packt*) je uvedeno ve stručné podobě. Nadbytečné je jen uvedení edice (*Expert insight*).

Online samostatné dokumenty neperiodické

Typickým příkladem je online PDF dokument (včetně elektronické knihy). Dalším příkladem je webové sídlo (*web site*) tj. typicky hierarchický systém více stránek (nikoliv tedy jedna konkrétní web stran).

medium : u online zdrojů se jako médium uvádí slovo *online*.

URL : klíčový údaj pro online zdroje. Některé systémy (např. Wikipedia) poskytují tj. fixní URL, které odkazují na konkrétní verzi dokumentu, resp. stránek. I když jsou tato URL obecně delší, je nutné jim dát přednost, neboť zaručují jedinečnost.

název : název nelze vynechat i když ne vždy je jasné, co je hlavním názvem. V tomto případě je možné využít obsahu elementu title v hlavičce HTML (pokud je zdroj v HTML) nebo jiná metadata (například jak je zdroj pojmenován v odkazu).

autoři : autor nebývá u mnoha online dokumentů dohledatelný (a v tomto případě je nutné ho vynechat). Rozhodně však věnujte čas zjištění autorství (může být uvedeno i mimo dokument).

odpovědná korporace : typicky je to držitel intelektuálních práv (copyrightu). Důležitý je především v případě, že není znám autor, ale uvádějte ho ve všech případech, kdy je dohledatelný. Většina bibliografických stylů tento atribut nepodporuje resp. ho běžně nezobrazuje. Proto je vhodné pro tento účel využívat atribut *nakladatel* (i když to není totéž).

verze/čas poslední aktualizace : nahrazuje rok vydání. V případě, že není použit fixní odkaz, je klíčovým zdrojem informací, jaká z verzí dokumentu byla použita jako zdroj. Online dokumenty se mění často, a tak je vhodné uvádět, co nejpřesnější specifikaci (číslo verze, čas poslední aktualizace). Jen v případě, že dokument není verzován a nelze zjistit přesnější čas poslední modifikace, lze využít vrocení (stejně jako u knih může být odhadnuto z copyrightu).

datum použití : je to povinný údaj i když důležitý je jen v případě, kdy nelze určit přesnější verzi. Měl by být v ISO formátu tj. ve tvaru RRRR-MM-DD. Toto datum běžně generuje editor bibliografických citací podle data vytvoření záznamu. V každém případě by mělo ležet v časovém intervalu od poslední modifikace zdroje (je-li uvedeno) do data odevzdání závěrečné práce.

příklad:

Dílčí tištěné dokumenty

U dílčích tištěných dokumentů je typické, že kromě identifikace dílčí části obsahují i identifikaci dokumentu jako celku.

Klasickým příkladem jsou články ve sborníku nebo vědeckém časopise. Kapitoly v knize (monografii) se citují, jen případě, že každou z nich vytvořil jiný autor (či kolektiv autorů)

V zásadě platí tato pravidla:

- uvádí se jen autoři dílčí části, nikoliv například editoři sborníku nebo časopisu
- uvádí se pozice části v celém dokumentu nejlépe pomocí rozsahu stránek
- pokud má dílčí část vlastní jednoznačný identifikátor (například DOI), není potřeba uvádět identifikátor knihy nebo periodika.

Dílčí online dokumenty

Tento typ citací se používá pro webové stránky, jež jsou součástí webového sídla například pro konkrétní stránky s dokumentací nebo dokumenty uložené na GitHubu. Pro jiné elektronické dokumenty, pokud nejsou výslovně součástí webového sídla (např. elektronického sborníku) je vhodnější použít záznam samostatného dokumentu (viz výše).

Název stránky je doplněn jménem webového sídla (to je typicky uvedeno v záhlaví každé stránky resp. na hlavní stránce webového sídla. Autoři se vztahují ke stránce zatímco korporátní odpovědnost je typicky vztažena k celému sídlu (pokud jsou známy autoři i korporátní odpovědnost je vhodné uvést oba údaje, vždy však musí být uveden alespoň jeden z těchto údajů). Všechny ostatní atributy se vztahují

příklady:

What's New In Python 3.9: Summary – Release highlights. *Python 3.9.0 documentation [online]*. Python Software Foundation, October 14, 2020 [cit. 2020-10-15]. Dostupné z: <https://docs.python.org/3/whatsnew/>

Záznam obsahuje název dílčí části a také název celého webového sídla (v kurzívě). Odpovědná organizace je uvedena na místě nakladatele (autor není uveden a tak je tato informace klíčová). Verze je určena datem poslední modifikace (je uvedeno přímo ve tvaru použitém na stránce). Datum citování je povinné, ale v tomto případě nenese žádnou přidanou informaci (jen to, že citace byla vytvořena jen den po poslední modifikaci. Poslední součástí je URL.

Python nonlocal statement. *Stack Overflow [online]*. Stack Exchange, 2022-03 [cit. 2022-07-27]. Dostupné z: <https://stackoverflow.com/questions/1261875/python-nonlocal-statement>

Struktura záznamu je stejná. Čas poslední modifikace byl určen z informace, že poslední modifikace proběhla před čtyřmi měsíci (je uvedena v ISO formátu, ale odpovídající by byl i údaj například ve tvaru *březen 2022* nebo *March 2022*).

5.3. Často kladené otázky

Co není potřeba citovat?

Obecně platí, že citovat není potřeba znalosti, které jste získali v průběhu studia a to jak při výuce tak i z učebních materiálů (opor, skript). Citovat není potřeba ani zdroj formálních údajů (např. významu zkratk), pokud je lze snadno získat (například na Wikipedii).

To jest není nutné uvádět citaci při uvedení zkratky HTTP (zkratka je všeobecně známá a běžně využívána v mnoha kurzech). Podobně není nutné odkazovat pojmy jako Internet, počítačová síť, programovací jazyk, procesor, apod.

Běžně se také necitují (původní) myšlenky vedoucího práce, pokud si vedoucí práce nevyžádá jinak. Pokud vám zprostředkuje nepůvodní myšlenku, měl by vám pomoci najít originální zdroj (který uvedete v citaci).

Citování není možné v případě, kdy není znám původní zdroj, resp. je v podobě, kterou není možné citovat (lidová říčení, apod.) Pravděpodobnost výskytu takových textů v informatické bakalářské práci je však velmi nízká.

Jak citovat informace z (podnikových) školení

Pokud se jedná o evidentní výtvar školitele, můžete odkazovat příslušný výukový materiál (i když je neveřejný). Pokud je informace nepůvodní, pak je vhodné citovat primární resp. alespoň dostatečně autoritativní zdroj.

Jak citovat ústní sdělení?

Ústní sdělení je potřeba citovat jen tehdy, když je od autoritativní osoby v oblasti její odbornosti. Pokud například píšete práci o nasazení databáze, pak je autoritativní osobou například správce databázového systému (který vám sdělí například zkušenosti s nasazením).

Jak je uvedeno výše, ve většině případů se necitují ústní sdělení učitelů, školitelů, vedoucího práce a dalších sekundárních zdrojů.

Pokud citujete ústní sdělení je vhodné s tím danou osobu seznámit či získat alespoň neformální souhlas, neboť sdělené informace nemusí být veřejné.

Navzdory důležitosti ústních sdělení v některých typech prakticky zaměřených prací, není citace ústních sdělení standardizována. Jednoduchý návod nabízí například blog na citace.com [XXX].

Ústí sdělení je však neověřitelné a nelze ho jednoznačně identifikovat. Proto je lepší pokud se sdělení děje například e-mailem. Citace e-mailové komunikace i dalších netradičních zdrojů shrnuje dokument [XXX].

Je možno citovat Wikipedii?

Citování Wikipedie se obecně nedoporučuje, neboť se jedná o terciární zdroj (encyklopedia vytvořená na základě druhotných informací) a její kvalita je značně kolísavá.

Na druhou stranu Wikipedia (především v anglické verzi) často obsahuje i hodnotný a jinak jen obtížně dostupný materiál, a tak nelze citování z Wikipedie striktně zakázat.

Základní doporučení pro citování z Wikipedie:

- citujte jen tehdy, pokud nemáte k dispozici primární zdroje (ty jsou často odkazovány přímo z Wikipedie)
- citujte jen kvalitní články (které nejsou označeny jako problematické), které se v oblasti informatiky a matematiky objevují spíše na anglické Wikipedii
- citace z Wikipedie by měly tvořit jen malou část zdrojů (typicky méně než 10

Z Wikipedie rozhodně necitujte články věnované běžně známým technologiím a poznatkům, které jsou běžnou součástí kurzů.

6. Zhodnocení

7. Závěr

Závěr je klíčovou kapitolou, která může nejvíce ovlivnit vaši obhajobu. Základní částí závěru je přehledné shrnutí výstupů práce tj. co jste udělali pro dosažení cílů práce. Je nutné se vyhnout hodnocení, zda tím byli splněny cíle práce, či nikoliv (to je úkol posudků a především komise).

Seznam použitých zdrojů

1. BRAGIN, Tanya. *The Unbundling of the Cloud Data Warehouse*. ClickHouse Inc., 2023-11. Dostupné také z: <https://clickhouse.com/blog/the-unbundling-of-the-cloud-data-warehouse>. Accessed 2025-10-23.
2. PRIEBE, Torsten; NEUMAIER, Sebastian; MARKUS, Stefan. Von Data Warehouse bis Data Mesh: Ein Wegweiser durch den Dschungel analytischer Datenarchitekturen. *arXiv preprint arXiv:2212.03612*. 2022. Dostupné také z: <https://arxiv.org/abs/2212.03612>. Accessed 2025-10-23.
3. CUBE.DEV TEAM. *Semantic Layer and AI: The Future of Data Querying with Natural Language*. Cube.dev, 2024-08. Dostupné také z: <https://cube.dev/blog/semantic-layer-and-ai-the-future-of-data-querying-with-natural-language>. Accessed 2025-10-23.
4. SHARMA, Ayush. *Introduction to Cube.js – The Semantic Layer for Building Data Apps*. 2023-02. Dostupné také z: https://heyayush.com/post/2023-02-22_cube-js. Accessed 2025-10-23.
5. GEBRIM, Josue Luzardo. *Cube: Creating a Semantic Data Layer!* Medium, 2023-03. Dostupné také z: <https://medium.com/codex/cube-creating-a-semantic-data-layer-a947fd0c11f6>. Accessed 2025-10-23.
6. MICROSOFT CORPORATION. *Microsoft SQL Server Overview*. 2024. Dostupné také z: <https://learn.microsoft.com/en-us/sql/sql-server/>. Accessed 2025-10-23.
7. MICROSOFT CORPORATION. *Power BI Documentation*. 2024. Dostupné také z: <https://learn.microsoft.com/en-us/power-bi/>. Accessed 2025-10-23.
8. GOOGLE CLOUD. *BigQuery Documentation*. 2024. Dostupné také z: <https://cloud.google.com/bigquery/docs/>. Accessed 2025-10-23.
9. AMAZON WEB SERVICES. *Amazon Redshift Overview*. 2024. Dostupné také z: <https://aws.amazon.com/redshift/>. Accessed 2025-10-23.
10. DBT LABS. *dbt Semantic Layer Documentation*. 2024. Dostupné také z: <https://docs.getdbt.com/docs/use-dbt-semantic-layer/dbt-sl>. Accessed 2025-10-23.
11. APACHE SOFTWARE FOUNDATION. *Apache Superset Documentation*. 2024. Dostupné také z: <https://superset.apache.org/docs/>. Accessed 2025-10-23.
12. METABASE INC. *Metabase Open Source Business Intelligence*. 2024. Dostupné také z: <https://www.metabase.com/>. Accessed 2025-10-23.

13. SERVICES, Amazon Web; TALEND. *Data Warehouse Modernization: From Monolith to Modular*. Amazon Web Services, 2019. Dostupné také z: https://pages.awscloud.com/rs/112-TZM-766/images/GLOBAL_PTNR_STPM_DWM_Talend_eBook_Aug-2019.pdf. Accessed 2025-10-23.
14. DATABRICKS. *The Modern Data Stack: How the Evolution of Data Architecture Led to the Data Intelligence Platform*. Databricks, 2024. Dostupné také z: <https://www.databricks.com/blog/modern-data-stack-how-evolution-data-architecture-led-data-intelligence-platform>. Accessed 2025-10-23.
15. CLOUD, Google. *Looker Semantic Model — Unifying Business Logic for Data Experiences*. Google LLC, 2024. Dostupné také z: <https://cloud.google.com/looker/docs/semantic-model-overview>. Accessed 2025-10-23.

A. Externí přílohy

Externí přílohy této bakalářské práce jsou umístěny na adrese:

https://github.com/Jiri-Fiser/thesis_ki_ujep.

Na úložišti GitHub mohou být uloženy tyto externí přílohy:

- **zdrojové kódy**
- **doplňkové texty** (například jak instalovat aplikaci, manuály aplikace)
- **schémata** (především, pokud se nevejdou na stranu A4 a jejich vytištění je tak problematické)
- **screenshoty** (v textu práce lze použít jen omezený počet snímků obrazovky, které navíc nemusí být při černobílém tisku příliš přehledné)
- **videa** (například ovládání aplikace)

V každém případě by to však měli být pouze materiály, které jste vytvořili sami. Materiály jiných autorů uvádějte v seznamu použité literatury (včetně případných odkazů na jejich originální umístění).

V této kapitole stačí uvést pouze základní strukturu úložiště (co se kde nalézá a jakou má funkci) například v podobě tabulky.

| | |
|------------------------------|--|
| ki-thesis.pdf | text práce v PDF |
| ki-thesis.tex | zdrojový kód práce v \LaTeX |
| kitheses.cls | definice třídy dokumentů (rozšířená třída scrbook) |
| thesis.bib | bibliografická databáze (exportována z citace.com) |
| LOGO_PRF_CZ_RGB_standard.jpg | logo fakulty s českým textem |
| LOGO_PRF_EM_RGB_standard.jpg | logo fakulty s anglickým textem |

Všechny tyto soubory jsou potřeba pro překlad dokumentu (logo stačí jedno v příslušné jazykové verzi).

B. Další přílohy

Výjimečně může práce obsahovat i další tištěné přílohy. Obecně však dávejte přednost elektronickým přílohám umístěným na GitHubu (tato kapitola tak bude úplně chybět).