

Žilinská univerzita v Žiline

Záverečná správa

Aktivita 1.4

Názov aktivity: Vybudovanie dátového úložiska



OBSAH

OBSAH	ii
1 ÚVOD	1
1.1 Popis aktivity	1
1.2 Sumárne zhrnutie výsledkov	3
Top level architektúra systému	3
Jednoduchý popis komponentov systému	4
asynchrónna komunikácia	5
Zabezpečenie QOS	5
Rozhrania medzi jednotlivými komponentami systému	7
Main Database ⇔ Main WF Environment.....	7
Main WF Environment ⇔Local DB Provider	8
Main WF Environment ⇔VM Provider	9
Podrobný popis jednotlivých komponentov	11
Hlavné úložisko (Main DB)	11
Main WF Environment.....	12
VM Provider.....	12
Local DB Provider.....	13
Definovanie vlastných WF.....	14
Main workflow	14
Hlavná funkcionálna (MFcn - Main Functionality)	15
Operácia.....	15
Processing task	15
2 Server a úložisko	16
3 Dátové úložisko	23

4	Rozsah spracovania dát z pohľadu hierarchie dátových tokov	29
4.1	Import surových dát	29
4.2	Import spracovaných produktov	30
4.3	Spracovanie a uloženie meta údajov do databázy	30
4.4	Zabezpečenie dát v dátovom úložisku - súborová úroveň.....	38
4.5	Zabezpečenie dát v dátovom úložisku - Webový portál dátového úložiska.....	38
4.6	Prístup do dátového úložiska cez VPN.....	39
4.7	Žurnálový prehľad aktivity užívateľov	39
4.8	Vyhľadávanie spracovaných a uložených údajov	39
4.9	Vyhľadávací formulár.....	40
4.10	Stiahnutie súborov ako *.zip	42
4.11	Stiahnutie súborov pomocou *.cmd	42
4.12	Spracovanie *.las pomocou lastools	42
4.13	Nahratie shape file pre lasclip	42
4.14	Žurnálový formulár.....	43
4.15	Formulár pre import surových dát	43
4.16	Formulár pre import spracovaných produktov	43
4.17	Automatické mazanie starých súborov	44
4.18	Rozlišovanie oprávnení, podľa autentifikačného servera.....	44
5	Realizácia sofistikovaného prepojenia na externé zdroje dát.....	45
5.2	PostgreSQL.....	49
6	Realizácia sofistikovaných prístupov pre akceleráciu dátových operácií	52
6.1	Riešenie pre dátové úložisko	52
6.2	Architektúra riešenia.....	53
7	Zoznam merateľných ukazovateľov	57



1 ÚVOD

1.1 POPIS AKTIVITY

Cieľom aktivity bolo vytvorenie spoločného úložiska dát, ktoré bolo vytvárané pre potreby dosiahnutia špecifických cieľov projektu a slúži ako základná údajová báza pre proces získavania znalostí. Pre dosiahnutie definovaného cieľa aktivity boli stanovené nasledovné čiastkové aktivity:

- *Vytvorenie úložiska transakčných dát :*

Cieľom bolo vytvorenie primárneho úložiska transakčných dát. Úložisko okrem klasických dát ukladaných konvenčným spôsobom v relačnej databáze pracuje aj s dátami, ktoré nie je možné priamo uložiť v relačnej resp. objektovej databáze. Z dôvodu účelného využitia maximálneho množstva údajov bolo nutné vytvoriť rozšírenia databáz tak, aby bolo možné efektívne využívať všetky dáta vznikajúce ako výstupy v jednotlivých aktivitách projektu.

Pri riešení aktivity sme vychádzali z predpokladu silne distribuovaného charakteru dát, ako aj značnej veľkosti výsledných databáz – rádovo TB až PB. Tomuto bolo prispôsobené celkové riešenie transakčného úložiska.

- *Vytvorenie sémantického úložiska :*

Pre účely uloženia komplexných štruktúr identifikovaných v rámci projektu sa ukázalo ako vhodné vytvorenie spoločnej ontológie pre popis zásadných údajov a procesov vznikajúcich v rámci riešených aktivít. Vytvorená ontológia umožnila vytvorenie sémantického úložiska, ktoré môže spoločne s transakčnou bázou vytvorenou v predchádzajúcej podaktivite slúžiť ako základný zdroj informácií pre konzumentov informácií (fyzických alebo virtuálnych).

- *Vytvorenie analytickej databázy :*

Vytvorenie základnej analytickej bázy bolo realizované ako vhodné z pohľadu konsolidácie veľkého množstva dát do informácií z definovaného úložiska. V rámci



aktivity bolo potrebné vytvoriť základnú analytickú bázu ako predpoklad pre ďalšie, kontinuálne spracovávanie dát vznikajúcich v kompetenčnom centre.

Obstarané technické a programové vybavenie slúži na vytvorenie dátového úložiska pre veľké databázy (predpoklad uloženia peta bytov dát). Úložisko funguje ako základný zdroj údajov pre transakčné a analytické spracovanie v rámci aktivít kompetenčného centra.

Na základe analýzy potrieb na technické a programové vybavenie za účelom vytvorenie veľkého úložiska dát budú vyšpecifikované požadované parametre technického a programového vybavenia ktoré bolo realizované špecializovanou spoločnosťou (dodávateľom) v súlade so zákonom č. 25/2006 Z.z o verejnom obstarávaní. Riešitelia organizačne dohliadali na dodanie vybavenia a jeho príslušenstva v požadovaných termínoch a požadovanej kvalite.

Všetky realizované aktivity spolu navzájom súvisia, výsledkom ich realizácie je vybudovanie kompetenčného centra.

Všetky aktivity sú v súlade so špecifickým cieľom z ktorého vychádzajú, ako aj so strategickým cieľom predkladaného projektu.

Na začiatku projektu boli definované nasledovné riziká:

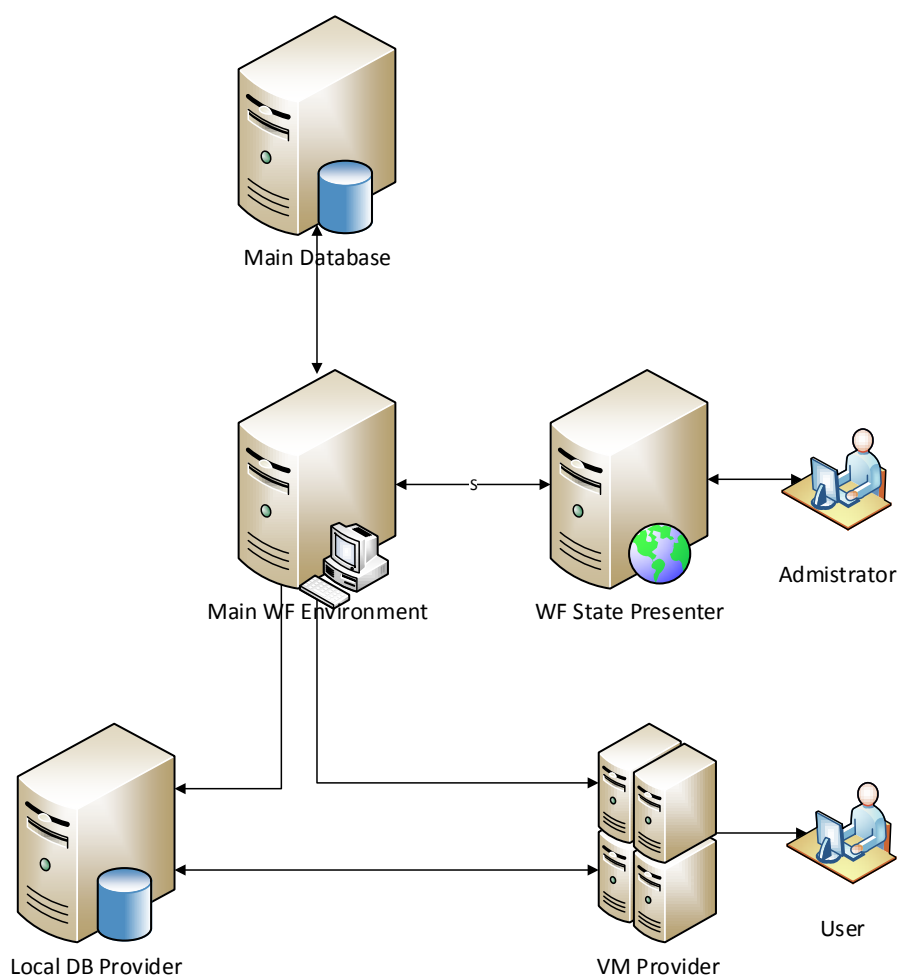
- zlyhanie ľudského faktora (choroby, preťaženosť, nával povinností, cesty do zahraničia, neovládanie zodpovednosti, a iné)
- projektové (týkajúce sa všetkých častí riadenia projektu, t.j. času, nákladov, rozsahu a kvality ukazovateľov výstupu a dopadu, zosúladenie výstupov s ostatnými časťami projektu, odborná úroveň členov projektového tímu a pod.)
- vonkajšie (zmeny na úrovni organizačnej štruktúry univerzity, vplyv rektora a iné)
- obchodné – týkajúce sa vhodnosti využitia navrhovaného prístrojového vybavenia pre pokrytie požiadaviek výskumu (externí dodávatelia, zmena cien, drahé doplnky, poddimenzovanie kapacít a schopností využívania prístrojov a pod.)
- technické – riziká spojené s uplatnením zvolených metodických postupov a zvoleného riešenia v prostredí laboratórneho i exteriérneho pracoviska (zmena platforiem, postupov, štandardov, zákonov a iné)
- etické (dodržiavanie pracovnej disciplíny, etických a morálnych zásad a iné)



Výskyt predpokladaných rizík sme eliminovali dostatočným naplánovaním personálnych kapacít na riešenej uvedenej problematiky. Zároveň sme naplánovali dostatočnú časovú rezervu tak, aby prípadné časové sklzy v riešení projektu nenarušili celkový harmonogram projektu.

1.2 SUMÁRNE ZHRNUTIE VÝSLEDKOV

TOP LEVEL ARCHITEKTÚRA SYSTÉMU



JEDNODUCHÝ POPIS KOMPONENTOV SYSTÉMU

HLAVNÉ ÚLOŽISKO (MAIN DB)

Hlavné databázové úložisko, v ktorom sú uložené všetky existujúce dáta. Toto úložisko je vytvorené s dôrazom na veľkosť, dĺžku uloženia a bezpečnosť uložených dát. V tomto úložisku nie je potrebné podporovať štandardné funkcie databázového úložiska (zložité výbery, filtrovanie, ...) postačujúce sú len funkcie na čítanie bloku dát.

Súčasťou tohto úložiska musí byť aj dočasné úložisko slúžiace na ukladanie odovzdávaných výsledkov.

MAIN WF ENVIRONMENT

Prostredie pre spúšťanie hlavných workflow. Toto prostredie je reprezentované jedným alebo viacerými počítačmi s nainštalovaným AppFabric (možnosť využitia týchto počítačov ako farmy). V tomto prostredí sa spúšťajú dodané workflowy, ktoré obsahujú funkcionality vykonávanú nad definovanými dátami. Tieto workflowy sú spúšťané s definovanou podmnožinou dát zo základnej databázy, nad ktorými sa má daná funkcionality vykonať. Tieto dáta sú určené len na čítanie. Výstupné dáta z vykonaného WF sa uložia do dočasného úložiska, ktoré bude súčasťou hlavnej databázy.

LOCAL DATA PROVIDER

Poskytovateľ úložiska pre lokálne dáta. Tento poskytovateľ musí podporovať dátové úložiská, ktoré sú potrebné pre vykonanie jednotlivých krokov postupu definovaného v hlavnom WF. V tomto úložisku sa ukladajú aktuálne dáta, ktoré môžu byť ľubovoľne spracovávané na rozdiel od dát v hlavnom úložisku.

Tento poskytovateľ je potrebný hlavne pre zvýšenie bezpečnosti hlavného úložiska a pre lepšiu schopnosť práce s dátami.

VIRTUAL MACHINES PROVIDER (VM PROVIDER)

Poskytovateľ, ktorý spravuje virtuálne stroje a poskytuje ich jednotlivým WF pre spracovanie časovo a výpočtovo náročných krokov. Tieto stroje sú vytvárané na základe predpripravených šablón a s možnosťou definovania výpočtových kapacít vytváraného stroja (veľkosť pamäti,

procesor, ...). Poskytovateľ takisto umožňuje definovať obmedzenia vytvárania a súbežného behu strojov podľa niektorých šablón (napr. kvôli licenčným podmienkam).

WF STATE MANAGEMENT

Základné rozhranie na prácu s bežiacimi WF s možnosťou kontroly stavu, spúšťania nových WF s možnosťou definovania ich parametrov, a podobne. Toto rozhranie slúži hlavne na monitorovanie a kontrolu, nemá prioritne slúžiť na spúšťanie, pretože to by sa malo diať automatizovane.

Takisto toto rozhranie slúži na správu existujúcich uložených riešení (predchádzajúce úspešné behy WF s rôznymi parametrami).

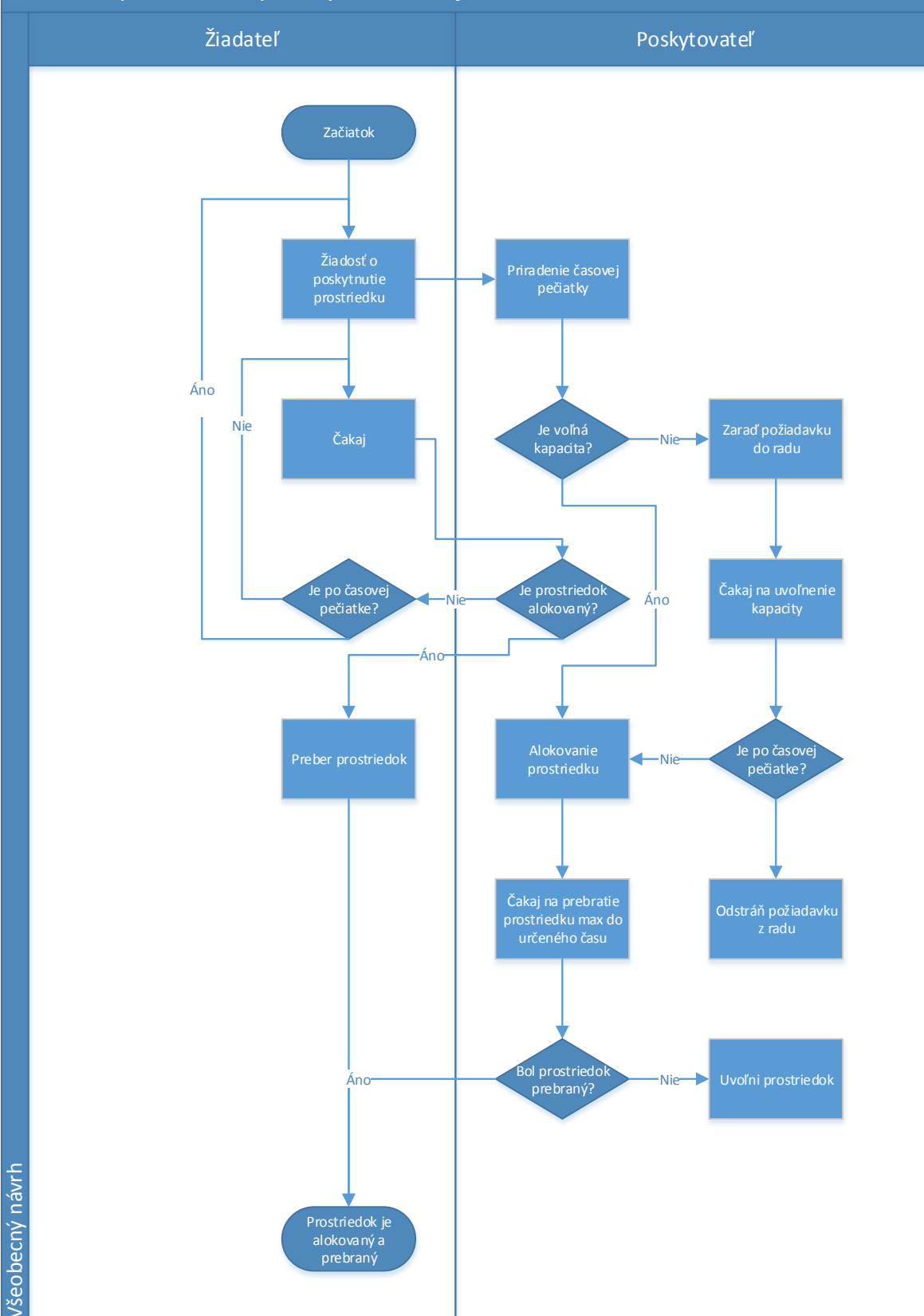
ASYNCHRÓNNA KOMUNIKÁCIA

ZABEZPEČENIE QOS

Pri výpadkoch WF, alebo ľubovoľného z poskytovateľov môže nastať to, že na jednej, alebo druhej strane niekto nepočúva. Vtedy treba riešiť to, ako zabezpečiť jednak poskytnutie prostriedkov, ak sú ešte žiadané, alebo uvoľnenie prostriedkov, ak už nie sú potrebné. Toto vieme riešiť časovou značkou, ktorá určuje životnosť daného alokovaného prostriedku.

Tým pádom je prostriedok držaný len po určitý čas a ak do vtedy nie je prebraný, tak je uvoľnený. Týmto sa zamedzí zbytočnému držaniu alokovaných prostriedkov ak bol napríklad žiadateľ zrušený, alebo je uspaný a podobne. Takisto sa týmto zabezpečí to, že žiadateľ dostane prostriedky v konečnom čase (požiadavky sú zaraďované do radu postupne a sú im pridelené časové značky na strane poskytovateľa).

Riešenie problémov pri asynchrónnej komunikácii



Všeobecný návrh

ROZHRANIA MEDZI JEDNOTLIVÝMI KOMPONENTAMI SYSTÉMU

MAIN DATABASE ⇔ MAIN WF ENVIRONMENT

VYŽIADANIE DÁT Z DOČASNÉHO ÚLOŽISKA

Vstupné parametre:

Parametre definujúce rozsah vstupných dát a identifikátor workflow. Identifikátor workflow musí byť rovnaký ako identifikátor použitý na uloženie tohto workflow.

Výstup:

Na základe vstupných parametrov bude pripravený endpoint, z ktorého bude možné postupne (sekvenčne) čítať pripravené dáta. Tento endpoint bude živý a pripravený na čítanie minimálne po dobu definovaného časového intervalu (možnosť zadať v nastaveniach) od poslednej vykonanej operácie. Počas doby existencie endpointu nesmú byť tieto dočasné dáta vymazané. V prípade, že neexistujú dáta pre zadaný vstup, oznámi sa chyba o nedostupnosti dát.

ČÍTANIE DÁT Z DLHODOBÉHO ÚLOŽISKA

Vstupné parametre:

Definovanie filtra vstupných dát pomocou definovania rozsahov existujúcich stĺpcov

Výstup:

Na základe vstupných parametrov bude pripravený endpoint, z ktorého bude možné postupne (sekvenčne) čítať pripravené dáta. Tento endpoint bude živý a pripravený na čítanie minimálne po dobu definovaného časového intervalu (možnosť zadať v nastaveniach) od poslednej vykonanej operácie.

OZNÁMENIE O UKONČENÍ ČÍTANIA DÁT Z DATABÁZY

Vstup:

Jednoznačný identifikátor endpointu, z ktorého bolo čítané



ULOŽENIE DÁT DO DOČASNÉHO ÚLOŽISKA

Vstup:

Jednoznačný identifikátor WF pomocou ktorého boli vygenerované tieto dáta a definovaný rozsah vstupu, nad ktorým bol tento WF spustený. Taktiež bude vstupom endpoint, ktorý slúži na čítanie ukladaných dát do dočasného úložiska. Tento endpoint bude živý a pripravený na čítanie minimálne po dobu definovaného časového intervalu (možnosť zadať v nastaveniach) od poslednej vykonanej operácie.

Ďalším vstupom bude endpoint, na ktorý bude oznámené ukončenie čítania a ukladania dát.

ULOŽENIE PARAMETROV VYKONANÉHO WF

Vstup:

Jednoznačný identifikátor ukončeného WF spolu so všetkými špecifickými krokmi vykonanými počas behu. Definovanie parametrov dát, nad ktorými prebehol daný WF.

MAIN WF ENVIRONMENT ⇔ LOCAL DB PROVIDER

ALOKOVANIE MIESTA

Popis

Príprava miesta pre WF, ktorý je identifikovaný jednoznačným identifikátorom. Toto miesto musí byť dostupné minimálne počas stanoveného časového úseku (možnosť definovania v nastaveniach) od poslednej operácie.

Vstup

Jednoznačný identifikátor asociovaného WF

VYŽIADANIE KONKRÉTNÉHO TYPU DÁTOVÉHO ÚLOŽISKA

Popis

Výstupom tejto metódy je alokované miesto v dátovom úložisku podľa požadovaných parametrov. Zároveň musí byť vytvorený používateľ na čítanie (prípadne aj administráciu) a zápis. Taktiež musí byť vygenerovaný spôsob pripojenia sa na dané úložisko.

Vstup

Jednoznačný identifikátor WF, ktorý bol použitý na alokáciu, typ dátového úložiska, parametre dátového úložiska.

Výstup

Spôsob pripojenia na toto úložisko (connection string, IP adresa, ...). Login pre používateľa s právom zapisovať a používateľa s právom čítať.

UVOĽNENIE ALOKOVANÉHO MIESTA

Popis

Na konci celého behu WF je uvoľnené lokálne dátové úložisko spolu so všetkými asociovanými úložiskami (vyžiadanými pomocou rovnakého jednoznačného identifikátora).

Vstup

Jednoznačný identifikátor asociovaného WF.

MAIN WF ENVIRONMENT ⇔ VM PROVIDER

VYŽIADANIE INŠTANCIE VIRTUÁLNEHO STROJA

Popis

Pre asociovaný WF vytvorí poskytovateľ virtuálny stroj podľa vstupných parametrov z existujúcej šablóny a spustí ho. Po naštartovaní stroja oznámi žiadateľovi IP adresu tohto



stroja. V prípade, že je použitá šablóna obmedzená počtom paralelne bežiacich kópií, žiadateľ musí byť zaradený do zoznamu žiadateľov. Po vytvorení a spustení stroja je tento udržiavaný minimálne po definovaný časový interval (možnosť definovania v nastaveniach) až po potvrdenie prijatia tohto stroja. Po potvrdení prijatia musí byť tento stroj spustený počas definovaných hodín až do uvoľnenia.

Vstup

Jednoznačný identifikátor asociovaného WF, parametre virtuálneho stroja a použitá šablóna.

Výstup

Jednoznačný identifikátor virtuálneho stroja

POTVRDENIE PREBRATIA VM

Popis

Potvrdenie prebratia virtuálneho stroja a začiatok jeho používania.

Vstup

Jednoznačný identifikátor VM

Výstup

IP adresa daného stroja

UVOĽNENIE VIRTUÁLNEHO STROJA

Popis

Uvoľnenie konkrétnej inštancie virtuálneho stroja identifikovaného pomocou jednoznačného identifikátora.

Vstup

Jednoznačný identifikátor virtuálneho stroja

UVOĽNENIE VŠETKÝCH VM ALOKOVANÝCH V RÁMCI WF

Popis

Ukončenie činnosti všetkých VM alokovaných v rámci behu WF. Toto uvoľňovanie nastáva len v prípade nepredvídaných okolností a chýb.

Vstup

Jednoznačný identifikátor WF asociovaného s alokovanými VM.

Výstup

PODROBNÝ POPIS JEDNOTLIVÝCH KOMPONENTOV

HLAVNÉ ÚLOŽISKO (MAIN DB)

Hlavné databázové úložisko, v ktorom sú uložené všetky existujúce dáta. Toto úložisko je vytvorené s dôrazom na veľkosť, dĺžku uloženia a bezpečnosť uložených dát. V tomto úložisku nie je potrebné podporovať štandardné funkcie databázového úložiska (zložité výbery, filtrovanie, ...) postačujúce sú len funkcie na čítanie bloku dát.

Súčasťou tohto úložiska musí byť aj dočasné úložisko slúžiace na ukladanie odovzdávaných výsledkov.

UKLADANIE DOČASNÝCH VÝSLEDKOV

V tomto úložisku sa ukladajú dáta, ktoré vznikli ako výsledok behu WF nad určitou podmnožinou vstupných dát. Tieto dáta slúžia ako výstup pre používateľa. Tieto dáta sa uchovávajú v tomto úložisku po určitú dobu v závislosti na tom, ako často sú vyžadované.

UKLADANIE BEHU WF

Pre ukladanie behu WF je potrebné uložiť jednoznačný identifikátor vytvoreného WF a celú postupnosť krokov a vstupných parametrov pre daný dotaz. Dotaz je definovaný filtrom vstupných dát z hlavnej databázy.

Uložené predchádzajúce behy je potrebné manažovať z dôvodu možnej zmeny prislúchajúceho WF. Používateľ musí mať možnosť uchovať predchádzajúce behy v prípade, že zmeny WF nespôsobili znemožnenie predchádzajúceho behu, vymazať predchádzajúce behy prípadne vymazať len niektoré.

MAIN WF ENVIRONMENT

Prostredie, ktoré slúži na spúšťanie a správu WF. Tieto WF slúžia na spracovanie dát z hlavnej databázy a výstup z týchto WF sa tiež ukladá do dočasnej hlavnej databázy.

REGISTROVANIE WF

Každý vytvorený WF musí byť do prostredia registrovaný cez používateľské prostredie, pretože je potrebné asociovať registrovaný WF s jednoznačným identifikátorom. Tento identifikátor je jedinečný pre každú verziu WF (z dôvodu potreby ukladania existujúcich behov pre potrebu ich opakovania). Takisto v prípade zmeny už existujúceho WF je potrebné zmenu zaregistrovať v prostredí.

Pre registrovaný WF sa tiež vytvorí endpoint, ktorý slúži na jeho spúšťanie v prípade potreby.

VM PROVIDER

Poskytovateľ virtuálnych strojov slúžiacich na spúšťanie jednotlivých krokov postupu. V rámci jedného bežiacieho WF je možné mať alokovaných viacero virtuálnych strojov v jednom čase. Tieto stroje sú vytvárané na základe šablón. Žiadosť o alokovanie stroja sa dejú z hlavného workflow, nie z používateľskej časti, čím je zabezpečené bezpečné uvoľňovanie týchto strojov.

ŠABLÓNY (VM TEMPLATES)

Poskytovateľ musí podporovať viacero rôznych šablón pre vytvárané virtuálne stroje. Šablóna predstavuje určitý uložený stav bežiacieho systému, spolu s nainštalovanými aplikáciami.

Príklad: Windows 7 64 s nainštalovaným Supervizorom a Geomediou.

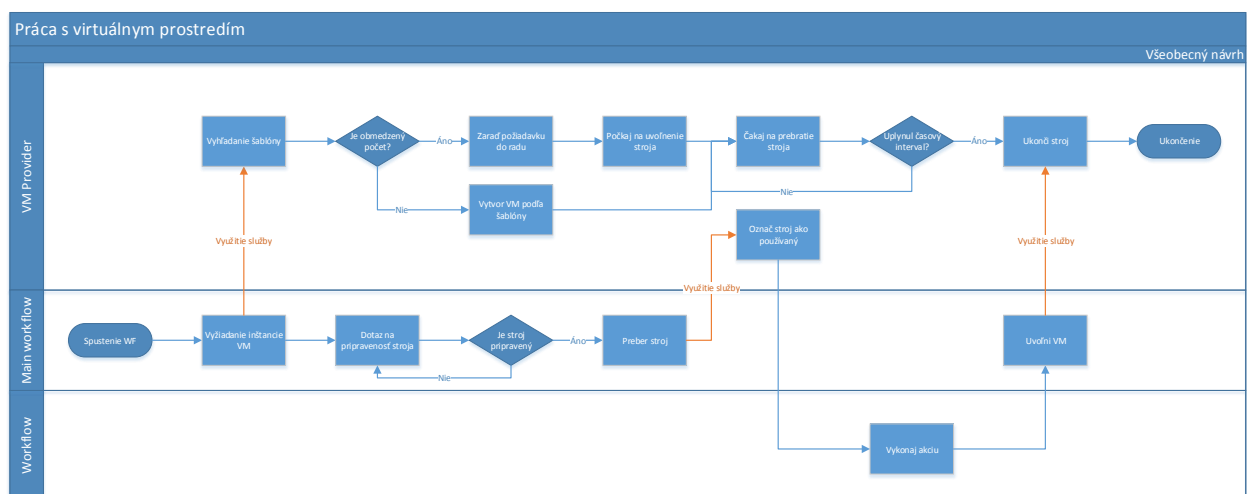
V rámci šablóny musí byť možné aplikovať obmedzenie počtu paralelne bežiacich virtuálnych strojov vytvorených na základe danej šablóny.

Príklad: Šablóna z predchádzajúceho príkladu obsahuje nainštalovanú Geomediú, čo znamená, že počet paralelných inštancií musí byť obmedzený na počet kúpených licencií.

Môže existovať špeciálna šablóna, ktorá je naviazaná priamo na konkrétny fyzický stroj. Takáto šablóna sa využíva v prípade nutnosti špecifického zariadenia, napríklad HW kľúča. Šablóna je identifikovaná na základe prideleného jednoznačného identifikátora, ktorý je pevne naviazaný na šablónu v rámci inštalácie. V rámci jedného poskytovateľa sa nesmú vyskytovať rôzne šablóny s rovnakým identifikátorom.

Poskytovateľ musí podporovať pridávanie a uberanie šablón, definovanie vlastností týchto šablón (obmedzenie počtu používateľov, minimálne požiadavky na VM, ...).

ZJEDNODUŠENÝ ŽIVOTNÝ CYKLUS VIRTUÁLNEHO STROJA



LOCAL DB PROVIDER

Poskytovateľ úložiska pre lokálne dáta. Tento poskytovateľ musí podporovať dátové úložiská, ktoré sú potrebné pre vykonanie jednotlivých krokov postupu definovaného v hlavnom WF.

V tomto úložisku sa ukladajú aktuálne dáta, ktoré môžu byť ľubovoľne spracovávané na rozdiel od dát v hlavnom úložisku.

Tento poskytovateľ je potrebný hlavne pre zvýšenú bezpečnosť hlavného úložiska a pre lepšiu schopnosť práce s dátami.

Ako základné úložisko údajov sa bude používať sekvenčná databáza bez nutnosti licencií, napríklad SQL Express. V prípade, že použitý proces bude potrebovať iný typ databázy, tak bude nutné vytvoriť ďalší modul to daného poskytovateľa.

Alokácie sa vykonávajú v základnom WF, ktorý sa postará aj o oznámenie na uvoľnenie prostriedkov, ale tieto prostriedky sa uvoľňujú až na konci. Všetky alokované prostriedky ostávajú alokované kvôli tomu, aby bolo možné kedykoľvek sa vrátiť na predchádzajúce kroky.

DEFINOVANIE VLASTNÝCH WF

Je potrebné v dodávanom systéme vytvoriť možnosť jednoduchého definovania vlastných WF nad spracovávanými dátami. Preto je dodávaný systém tvorený z viacerých modulov, ktoré spolupracujú.

MAIN WORKFLOW

Jeden z modulov je základný workflow, ktorý slúži na zabalenie celkovej funkcionality. Táto funkcionality môže byť tvorená z viacerých existujúcich spracovateľských úloh (processing task), ktoré sú vytvorené ako balík pozostávajúci z reprezentácie v rámci WF, z definície vstupných a výstupných dát a zo samotnej funkcionality, ktorá je vykonávaná vo VM.

Celková funkcionality pozostáva z „byrokracie“ ktorú treba vykonať v rámci WF, ako je alokovanie lokálnej databázy a jej finálne uvoľnenie, logovanie prípadných chýb, korektné ukončenie a uvoľnenie všetkých alokovaných prostriedkov a hlavne vykonanie hlavnej funkcionality ⇔ definovanej logiky. Tá pozostáva z viacerých operácií, ktoré sú definované používateľom.

HLAVNÁ FUNKCIONALITA (MFCN - MAIN FUNCTIONALITY)

Funkčnosť, ktorá je definovaná používateľom. Na začiatku budeme určite túto funkčnosť definovať u nás, ale v neskoršom čase je možné presunúť vytváranie týchto funkčností do rúk používateľom. Toto sa bude vytvárať tak, že používateľ vytvorí workflow, kde sa sebou zoradí operácie s definovanými parametrami podľa požadovanej funkčnosti. Potom tento workflow zaregistruje do systému, aby mohol byť vykonaný.

Takže používateľ len zoraduje existujúce kroky za sebou a len im definuje vstupné parametre.

OPERÁCIA

Operácia zabaľuje vykonanie jednotlivých processing taskov a všetkého, čo je ku tomu potrebné, ako napríklad vytvorenie LDB, skopírovanie dát, vyžiadanie VM, uvoľnenie VM a samotné spustenie processing tasku.

PROCESSING TASK

Toto je gro celého výpočtu. Processing task v sebe zahŕňa reprezentáciu v rámci WF, definíciu metadát a samotnú logiku, ktorú je možné spustiť na virtuálnej mašine.

Processing task bude potrebné registrovať do systému, aby ho bolo možné použiť. To zodpovedá zaregistrovaniu do systému, kde sa vytvárajú WF a taktiež zaregistrovanie na miesto, z kadiaľ sa spúšťa logika na VM.

Processing tasky budeme zo začiatku vytvárať my ale neskôr je možné porozmýšľať nad tým, že sa uvoľní SDK na vytváranie týchto processing taskov. Vytvorenie PT bude vlastne implementácia nejakého interface pre execution na VM, ďalej implementácia nejakého interface pre WF.

2 SERVER A ÚLOŽISKO

V rámci projektu bolo zaobstarané nasledovné hardvérové vybavenie.

2.1.1 DÁTOVÝ SERVER SO 4X BLADE SERVER

BladeCenter S CID	1
Switch Moduls Rloor	1
Procelar Bladr Filer	5
BladeCenter S Codo	1
BladeCenter S Disk Storage Module Filler	2
BladeCenter S Battery Backup Until Filler	2
From Borel ASM	1
BladeCenter S Label (SATA Optical)	1
IBM BGb STP + SW Optical Transerver	R
1.8m, IBA/230N, C 13 to CEE 7-WI Europe Line Cord	4
2.0m, 13A/125-10A/250V, C13 to IEC 320-C14 Rack Fower Cable	4
Ciscu Catalist Switch Module 3012 for IBM BladeCenter	2
Cloud 20-port 4/8 Gb SAN Switch Moduls for IBM BladeCenter	2
System Base (DATA Optical)	1
IBM BladeCenter S C14 950W/1450W Auto – sending Power Supply (Standard)	2
System Documentation and Software UK English	1
IBM BludeCenter 5 C14 SSOW/1450W Auto –seting Power Supplies 3 - 4	2
BladeCenter S Packaging WW	1
IBM Ultra Slim Enhanced SATA Multi Barner	1
Back mloment- Folory Integrale	1

BladeCenter KVM/Advanced Management Module	1
IBM 3000 VA LCD 3U Rack LIPS (2 30V)	1
C19 4.3 meter Line Cord – Europe	1
3 Year Onsite a Mega 24x74 Hour Response	1
HS22 CTO	4
2GB (1x2GB, 1Rx8, 1.35V), PC 3L-1D600 CL SI CC DOR3 1333MHz. VUP RDMM	4
Customer provided and installed	1
Unknown or not required	1
Blade Cover	1
Packaging – IU Blade WW	1
Labels for HS22 Blade Base	1
Dummy DIMM for Improved airflow	8
CPU Heat Sink Filler	1
Integrated SAS Mirroring – 2 Identical HDDs required	1
2/4 Port External Expansion Card (CFFN) for IBM BladeCenter	1
IBM 146GB 15K 8GB hpe SAS 2.5" SFF G2HS HDD	2
Clogic 8Gb Fibre Channel Expansion Card (CIOV) for IBM BladeCenter	1
Intel Xeon Processor E5670 2.93 GHz, 12MB Cache 1066MHz 80 W	1
System Documentation and Software – UK English	1
Blade Base 2 for Intel Xeon 5600 series	1
3 Year Onsite Regal 24x7 4 Hour Response	4
Úložisko dát (datastorage)	
IBM System Storage DS3512 Express Dual Controller Storage System	5
2TB 3.5 inch 7.2 KNL SAS HDD	60
8Gb FC 4 Port Daughter Card	10
8Gb FC SW STP Transporting Card	10

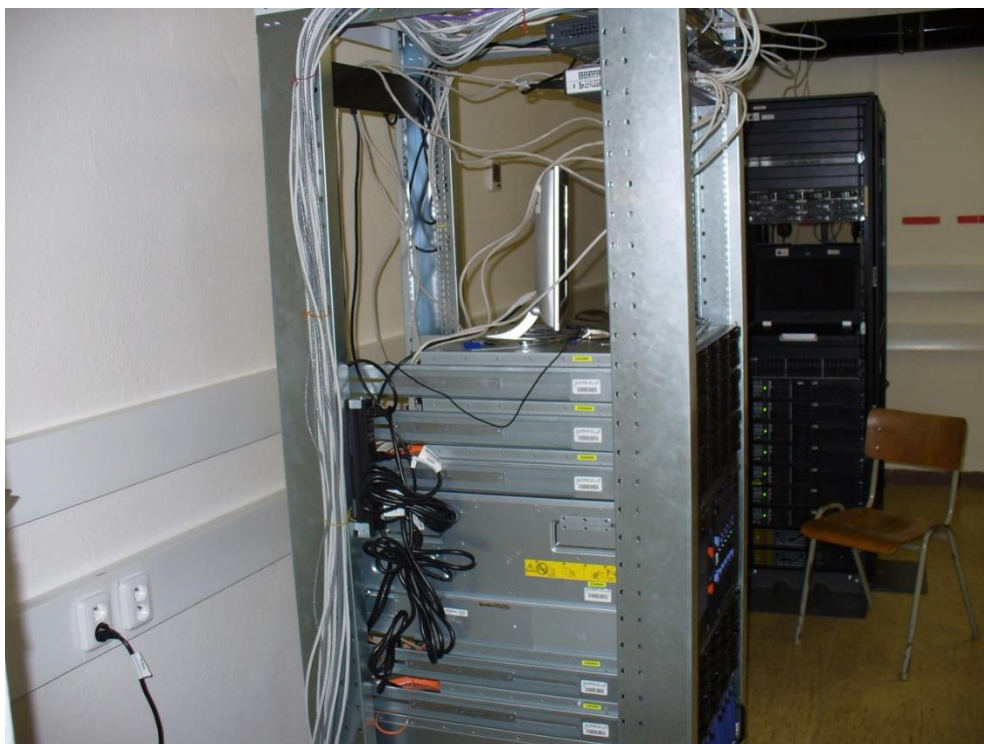
Sm Fiber Optic Cable ICAC	20
3 Year Onsite Repair 24x7 Same Business Day	5
Záložná údajová jednotka	
IBM System Storage DS3512 Express Dual Controller Storage System	1
2TB 3.5 km 7.2 KNL SAS HDD	12
BGb FC 4 Port Daughter Card	2
8Gb FC SW STP Transporting Card	2
Sm Fiber Optic Cable ICAC	4
3 Year Onsite Repair 24x7 Same Business Day	1
Server IBM x3650	
X3650 MB, Xe 4C E5606 80 W, 2 13 GHz/1066 MHz/BiB, 1x4GB, OxBy HS 2 Sin SAS/SATA, SR M1015, 4GOW p/S, Rack	1
4GB (1x4GB, 1Rx4, 1,35V), PC 34-10600 CLOSPCCDORF3 1333MHzUP ROOIMM	1
IBM L46GN 2.5n SFF Slim H5 15K 6Gips SA5HDD	2
IBM 460W Redundant Power Supply Unit	1
IBM Until Slim dolanced SATA DVD-ROM	1
3 Year Orsalte Repair 34x7 4 Hour Response	1



Obrázok 1 Pohľad na serverovňu



Obrázok 2 Informačná tabuľa



Obrázok 3 HW vybavenie



Obrázok 4 HW vybavenie



Obrázok 5 HW vybavenie - pohľad z blízka



Obrázok 6 HW vybavenie v serverovni



Obrázok 7 Bočný pohľad



Obrázok 8 Bočný pohľad

3 DÁTOVÉ ÚLOŽISKO

Úvodnou úlohou v rámci projektu bol výber vhodného dátového úložiska. Diskové (dátové úložisko) je základné označenie pre archiváciu dát v elektromagnetickej alebo inej forme pre použitie počítačom alebo iným zariadením. Rozdielne typy dátových úložísk hrajú rozdielne úlohy vo výpočtovom prostredí. Novou formou dátového úložiska sú vzdialené dátové úložiska, ako napríklad cloud computing, ktoré prinášajú revolúciu v spôsobe prístupu k dátam.

Okrem bežných vlastností diskových úložísk, ako veľkosť, rýchlosť, sú veľmi dôležité ďalšie vlastnosti, ako zálohovanie dát v úložisku, replikácia medzi rozdielnymi geolokalitami. V dnešnej dobe sa pod dátovým úložiskom teda môže rozumieť kombinácia nástrojov a hardvéru, ktoré umožňujú všetky tieto funkcionality (vlastnosti).

Dôležitým faktorom pri výbere správnej technológie pre ukladanie dát je aj štruktúra a následné použitie ukladaných dát.

V tomto projekte sú ukladané 2 základné typy dát:

- Surové nalietané LIDAR dáta
- Spracované produkty (las, geotiff, jpg, png, atď..)

Keďže sa jedná o geografické dáta, tak aj priestor v dátovom úložisku, ktorý tieto dáta budú zaberáť je markantný a je dôležité na to myslieť pri samotnom návrhu dátového úložiska. Odhadovaný objem jedného náletu povrchu zeme s plochou 5 km², je cca 100 GB. Čo pri ploche Slovenska znamená cca 10 000 takýchto náletov pre pokrytie celého územia. Takýto počet náletov znamená, že pre uloženie surových dát z jediného náletu územia SR, je potrebné uložiť cca 1 PB uložených surových nespracovaných dát.

Pre prácu s týmito dátami je potrebné následne zabezpečiť uloženie týchto dát a samozrejme umožniť ich spracovanie pre vytvorenie spracovaných produktov. Nakoľko bol spracúvaný veľký objem dát, je potrebné takto uložené dáta v prvom rade jednoducho identifikovať. Na to boli použité metadáta vo forme XML, ktoré budú identifikovať všetky dôležité atribúty ukladaných geografických dát.



- Metadáta.

3.1.1 SUROVÉ DÁTA

Surové dáta sú získavané technológiou **Lidar**. Lidar je technológia diaľkového prieskumu, ktorá meria vzdialenosť osvetľovaním cieľa pomocou lasera a následnou analýzou odrazeného svetla. Lidar je veľmi populárna technológia pre vytváranie máp s vysokým rozlíšením.

LAS je formát súborov, ktorý sa používa na vymieňanie troj-dimenzionálnych mrakov bodov medzi užívateľmi. Primárne bol vyvinutý na výmenu mračna bodov z Lidar-u. Tento formát podporuje výmenu akéhokoľvek troj-dimenzionálneho x, y, z súradnicového systému. Tento binárny formát je alternatívou pre iné proprietárne formáty alebo základný ASCII systém na výmenu dát.

Tieto dáta sa skladajú z veľkého počtu súborov (v rôznych adresároch), ktoré je potrebné uložiť. Z tohto dôvodu navrhujeme uloženie takto nalietaných dát do 1 adresára s príslušným meta súborom popisujúcim tieto nalietané dáta.

3.1.2 SPRACOVANÉ DÁTA

Spracované dáta sa získavajú analýzou surových Lidar dát. Pod pojmom spracované dáta si pre potreby tohto projektu môžeme predstaviť ako skupinu niekoľkých súborov s rozdielnou príponou. Z tohto dôvodu bolo navrhnuté ukladanie takto vytvorených dát tak, že všetky súbory majú rovnaký názov, ale rozdielnu príponu a jeden xml súbor, ktorá tieto dáta popisuje.

3.1.3 XML SÚBOR

Pre jednoznačnú identifikáciu uložených dát v dátovom úložisku je potrebné zabezpečiť meta údaje o týchto uložených dátach, ktoré budú tieto dáta jednoznačne identifikovať. XML súbor obsahuje viacero informácií než je reálne potrebné pre potreby tohto projektu. Z tohto dôvodu boli zvolené nasledujúce meta údaje ako kľúčové:



Tabuľka 3.1 Dôležité atribúty XML

Názov atribútu	Popis atribútu
Meta údaje o meta údajoch	
ID	Identifikátor meta údajového záznamu
Nadradené ID	Identifikátor nadradeného metaúdajového záznamu
Jazyk	Jazyk meta údajov
Identifikácia	
Názov	Charakteristický, často jedinečný názov, pod ktorým je zdroj známy
Abstrakt	Stručne popísané zhrnutie obsahu zdroja
Typ	Typ zdroja, ktorý je popisovaný meta údajmi
Identifikátor	Hodnota, ktorá zdroj jednoznačne definuje
Dátum	Referenčný dátum vytvorenia, aktualizácie či revízie údajov
Účel	Dôvod vytvorenia údajovej sady (napr. citácia zákona, vyhlášky a pod.)
INSPIRE	Zatriedenie zdroja k témam INSPIRE
Podmienky prístupu a použitia	Podmienky prístupu a použitia údajových sád
Obmedzenie verejného prístupu	Obmedzenie verejného prístupu k údajovým sadám
Téma	Tematická kategória predstavuje najširšiu klasifikáciu využívanú pri zoskupovaní a tematickom vyhľadávaní údajových sád
Geometria	Typ priestorovej prezentácie - metóda použitá k priestorovej reprezentácii geografickej informácie
Projekcia	Súradnicový referenčný systém
Jazyk	Jazyk údajovej sady alebo série
Znaková sada	Znaková sada zdroja
Aktualizácia	Informácie o početnosti aktualizácie
Mierka	Priestorové rozlíšenie - mierka mapy, z ktorej boli údaje odvodené

Vzdialenosť	Priestorové rozlíšenie - vzdialenosť - presnosť rozlíšenia údajovej sady v metroch
Rozsah	
Priestorový rozsah	
Zemepisná dĺžka - západ	
Zemepisná dĺžka - východ	
Zemepisná šírka - juh	
Zemepisná šírka - sever	
Časový rozsah	
Od	
Do	
Distribúcia	
Formát	Distribučný formát, v akom sú údaje uložené alebo distribuované
Verzia	Verzia formátu
Odkaz	Položka definuje odkaz na zdroj a/alebo odkaz na ďalšie informácie o zdroji
Kvalita	
Pôvod	Vyjadrenie histórie spracovania a/alebo celkovej kvality súboru priestorových údajov
Súlad	Prvok uvádza, či je údajová sada v súlade so špecifikáciami INSPIRE

Špecifikácia	Špecifikácia INSPIRE - oproti ktorej je súlad údajovej sady hodnotený
Kontakt - zdroj	Povinná osoba - organizácia zodpovedná za vytvorenie, správu, údržbu a sprístupnenie zdrojových evidencií
Organizácia	Názov organizácie
Dodací bod	Ulica a číslo
Mesto	Mesto alebo sídlo
PSČ	Poštové smerovacie číslo
Štát	Štát
telefón	Telefónne číslo
email	E-mailová adresa
www	WWW stránka organizácie
Rola	Rola organizácie vo vzťahu ku zdroju
Kontakt – meta údaje	Popis organizácie zodpovednej za vytvorenie a aktualizáciu meta údajov.
Organizácia	Názov organizácie
Dodací bod	Ulica a číslo
Mesto	Mesto alebo sídlo
PSČ	Poštové smerovacie číslo
Štát	Štát
telefón	Telefónne číslo
email	E-mailová adresa
www	WWW stránka organizácie
Rola	Rola organizácie vo vzťahu ku zdroju

Tieto zvolené atribúty sú následne v dátovom úložisku ukladané a párované ako jednoznačný identifikátor dát k surovým dátam alebo k spracovaným dátam.

Spracované a xml dáta sa musia ukladať ako výsledné súbory s príslušnými metadátami, avšak surové dáta dávajú priestor pre analýzu týchto dát.



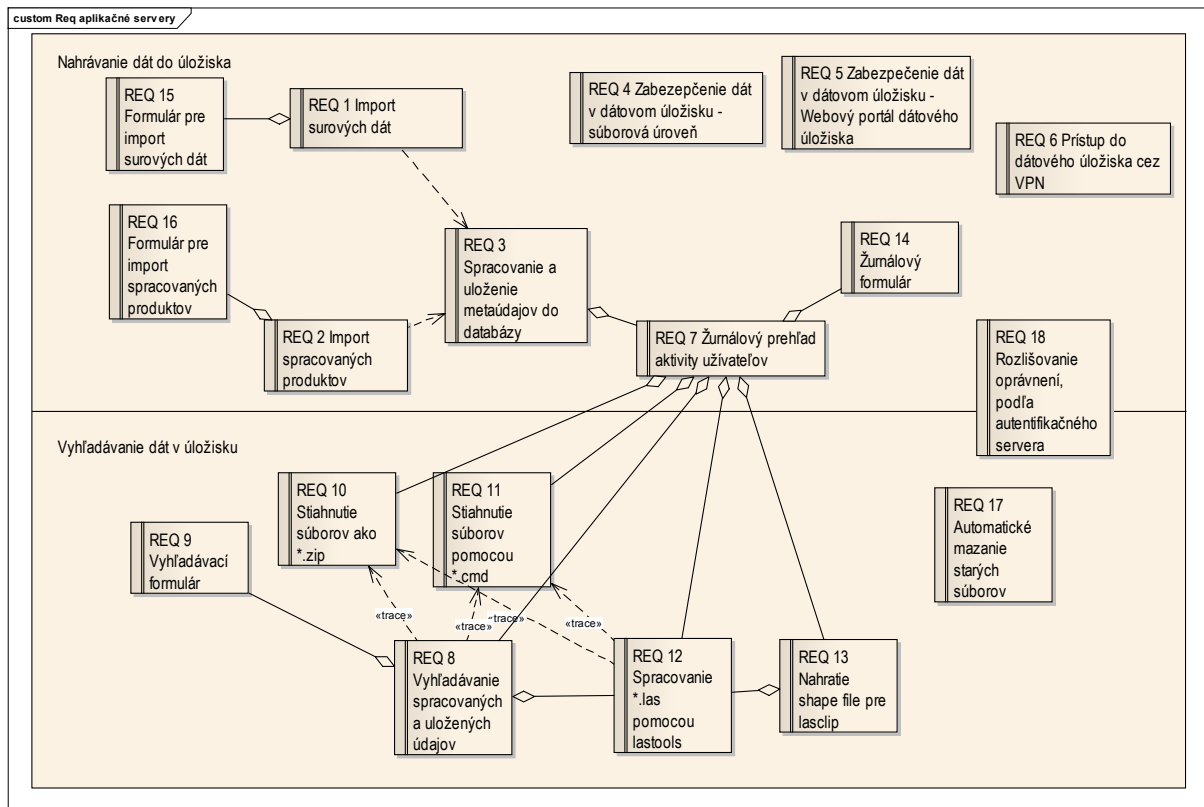
Z tohto dôvodu boli do testovania pre výber vhodného úložiska dát zvolené tieto riešenia:

1. Hadoop
2. Postgres + GIS = PostGIS
3. GlusterFS

4 ROZSAH SPRACOVANIA DÁT Z POHLĎADU HIERARCHIE DÁTOVÝCH TOKOV

Aplikačné servery pracujúce s dátovým úložiskom

Na obrázku sú zobrazené všetky požiadavky na aplikačné servery, ktoré pracujú s dátami v dátovom úložisku.



Obr. 4.1 Požiadavky na aplikačné servery

4.1 IMPORT SUROVÝCH DÁT

Ako vyplynulo z analýzy, surové dáta sú dáta, ktoré sú nespracované a predstavujú veľký počet súborov v rôznych adresároch. Tieto surové dáta popisuje 1 XML súbor, preto je potrebné do dátového úložiska tieto dáta ukladať iným spôsobom ako sú ukladané spracované produkty. Preto sme navrhli konvenciu pre nazývanie súborov/adresárov nasledovne:

- Názov adresára so surovými dátami, napr. test,

- Názov XML súboru popisujúceho dané surové dáta test.xml.

Takto zabezpečíme to, že dané dáta budú oddelené, nakoľko je dôležité, aby nevznikli surové dáta s rovnakým názvom je nutné, aby každý názov bol jednoznačný a unikátny.

4.2 IMPORT SPRACOVANÝCH PRODUKTOV

Z analýzy vyplynulo, že spracované produkty predstavujú jeden alebo niekoľko súborov, ktoré majú rovnaký názov, ale rozdielnu príponu. Z tohto dôvodu a najmä kvôli jednotnosti nazývania súborov bolo navrhnuté, aby sa súbory volali nasledovne:

- Názov súboru napr.:
 - o test1.tiff
 - o test1.jpg
 - o test1.las
- Názov XML súboru popisujúceho daný spracovaný produkt (všetky súbory):
 - o test1.xml

Takto zabezpečíme, že všetky súbory, ktoré patria k jednému produktu budú mať unikátny názov a tak isto je nežiadúce, aby existovali produkty s 2 rovnakými názvami.

4.3 SPRACOVANIE A ULOŽENIE META ÚDAJOV DO DATABÁZY

Meta údaje o každom type súborov nahrávaných do dátového úložiska sa nachádzajú v súbore XML, ktorý má svoju predpísanú štruktúru. Z toho dôvodu je dôležité načítať dané XML a podľa príslušných a relevantných atribútov ich vložiť do databázy. Pri ukladaní týchto údajov je dôležité uložiť okrem daných údajov obsiahnutých v XML aj cestu v destination_files, kde sú súbory fyzicky uložené, aby sme pri vyhľadávaní mohli dané dáta ďalej spracovať alebo stiahnuť.

Atribúty, ktoré je potrebné ukladať do databázy z XML aj s ich popisom sa nachádzajú v tabuľke nižšie.

V tabuľke sú zobrazené presné cesty v štruktúre XML ako sa dostať k hodnotám týchto atribútov.

Tabuľka 4.1 Cesty k hodnotám atribútov v XML

Názov atribútu	lokalizácia v xml súbore
Meta údaje o meta údajoch	
ID	gmd:fileIdentifier
Nadradené ID	gmd:parentIdentifier
Jazyk	gmd:language
Identifikácia	
Názov	gmd:identificationInfo, gmd:citation, gmd:CI_Citation, gmd:title
Abstrakt	gmd:identificationInfo, gmd:abstract
Typ	gmd:hierarchyLevel, gmd:MD_ScopeCode
Identifikátor	gmd:identificationInfo, gmd:citation, gmd:CI_Citation, gmd:identifier, gmd:RS_Identifier, gmd:code
Dátum	gmd:identificationInfo, gmd:citation, gmd:CI_Citation, gmd:date, gmd:CI_Date, gmd:date
Účel	gmd:identificationInfo, gmd:purpose
INSPIRE	gmd:identificationInfo, gmd:descriptiveKeywords, gmd:MD_Keywords, gmd:keyword
Podmienky prístupu a použitia	gmd:identificationInfo, gmd:resourceConstraints, gmd:MD_Constraints, gmd:useLimitation
Obmedzenie verejného prístupu	gmd:identificationInfo, gmd:resourceConstraints, gmd:MD_LegalConstraints, gmd:otherConstraints
Téma	gmd:identificationInfo, gmd:topicCategory
Geometria	gmd:identificationInfo, gmd:spatialRepresentationType, gmd:MD_SpatialRepresentationTypeCode
Projekcia	gmd:referenceSystemInfo, gmd:MD_ReferenceSystem, gmd:referenceSystemIdentifier, gmd:RS_Identifier, gmd:code

Jazyk	gmd:identificationInfo, gmd:language
Znaková sada	gmd:identificationInfo, gmd:characterSet
Aktualizácia	gmd:identificationInfo, gmd:resourceMaintenance, gmd:MD_MaintenanceInformation, gmd:maintenanceAndUpdateFuency
Mierka	gmd:identificationInfo, gmd:spatialResolution, gmd:MD_Resolution, gmd:equivalentScale, gmd:MD_RepresentativeFraction, gmd:denominator
Vzdialenosť	gmd:identificationInfo, gmd:spatialResolution, gmd:MD_Resolution, gmd:distance
Rozsah	
Priestorový rozsah	
Zemepisná dĺžka - západ	gmd:identificationInfo, gmd:extent, gmd:EX_Extent, gmd:geographicElement, gmd:EX_GeographicBoundingBox, gmd:westBoundLongitude
Zemepisná dĺžka - východ	gmd:identificationInfo, gmd:extent, gmd:EX_Extent, gmd:geographicElement, gmd:EX_GeographicBoundingBox, gmd:eastBoundLongitude
Zemepisná šírka - juh	gmd:identificationInfo, gmd:extent, gmd:EX_Extent, gmd:geographicElement, gmd:EX_GeographicBoundingBox, gmd:southBoundLatitude
Zemepisná šírka - sever	gmd:identificationInfo, gmd:extent, gmd:EX_Extent, gmd:geographicElement, gmd:EX_GeographicBoundingBox, gmd:northBoundLatitude
Časový rozsah	
Od	gmd:identificationInfo, gmd:extent, gmd:EX_Extent, gmd:temporalElement, gmd:EX_TemporalExtent, gmd:extent, gml:TimePeriod, gml:beginPosition
Do	gmd:identificationInfo, gmd:extent, gmd:EX_Extent,

	gmd:temporalElement, gmd:EX_TemporalExtent, gmd:extent, gml:TimePeriod, gml:endPosition
Distribúcia	
Formát	gmd:distributionInfo, gmd:MD_Distribution, gmd:distributionFormat, gmd:MD_Format, gmd:name
Verzia	gmd:distributionInfo, gmd:MD_Distribution, gmd:distributionFormat, gmd:MD_Format, gmd:version
Odkaz	gmd:distributionInfo, gmd:MD_Distribution, gmd:transferOptions, gmd:MD_DigitalTransferOptions, gmd:onLine, gmd:CI_OnlineResource, gmd:linkage
Kvalita	
Pôvod	gmd:dataQualityInfo, gmd:DQ_DataQuality, gmd:lineage, gmd:LI_Lineage, gmd:statement
Súlاد	gmd:dataQualityInfo, gmd:DQ_DataQuality, gmd:report, gmd:DQ_DomainConsistency, gmd:result, gmd:DQ_ConformanceResult, gmd:pass
Špecifikácia	gmd:dataQualityInfo, gmd:DQ_DataQuality, gmd:report, gmd:DQ_DomainConsistency, gmd:result, gmd:DQ_ConformanceResult, gmd:specification, gmd:CI_Citation, gmd:title
Kontakt zdroj	
Organizácia	gmd:pointOfContact, gmd:CI_ResponsibleParty, gmd:organisationName
Dodací bod	gmd:pointOfContact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:address, gmd:CI_Address, gmd:deliveryPoint
Mesto	gmd:pointOfContact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:address, gmd:CI_Address, gmd:city
PSČ	gmd:pointOfContact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:address, gmd:CI_Address, gmd:postalCode

Štát	gmd:pointOfContact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:address, gmd:CI_Address, gmd:country
telefon	gmd:pointOfContact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:phone, gmd:CI_Telephone, gmd:voice
email	gmd:pointOfContact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:address, gmd:CI_Address, gmd:electronicMailAddress
www	gmd:pointOfContact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:onlineResource, gmd:CI_OnlineResource, gmd:linkage
Rola	gmd:contact, gmd:CI_ResponsibleParty, gmd:role
Kontakt – meta údaje	
Organizácia	gmd:contact, gmd:CI_ResponsibleParty, gmd:organisationName
Dodací bod	gmd:contact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:address, gmd:CI_Address, gmd:deliveryPoint
Mesto	gmd:contact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:address, gmd:CI_Address, gmd:city
PSČ	gmd:contact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:address, gmd:CI_Address, gmd:postalCode
Štát	gmd:contact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:address, gmd:CI_Address, gmd:country
telefon	gmd:contact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:phone, gmd:CI_Telephone, gmd:voice
email	gmd:contact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:address, gmd:CI_Address, gmd:electronicMailAddress
www	gmd:contact, gmd:CI_ResponsibleParty, gmd:contactInfo, gmd:CI_Contact, gmd:onlineResource, gmd:CI_OnlineResource, gmd:linkage
Rola	gmd:contact, gmd:CI_ResponsibleParty, gmd:role

Dátový typ jednotlivých atribútov v XML súbore sa nachádza v tabuľke ďalej.

Tabuľka 4.2 Dátové typy atribútov v XML

Názov atribútu	Typ	Xsd
Meta údaje o meta údajoch		
ID	gco:CharacterString	string
Nadradené ID	gco:CharacterString	string
Jazyk	gmd:LanguageCode	string
Identifikácia		
Názov	gco:CharacterString	string
Abstrakt	gco:CharacterString	string
Typ	gmd:MD_ScopeCode	string
Identifikátor	gco:CharacterString	string
Datum	gco:Date	date
Účel	gco:CharacterString	string
INSPIRE	gco:CharacterString	string
Podmienky prístupu a použitia	gco:CharacterString	string
Obmedzenie verejného prístupu	gco:CharacterString	string
Téma	gmd:MD_TopicCategoryCode	string
Geometria	gmd:MD_SpatialRepresentationTypeCode	string
Projekcia	gmd:MD_SpatialRepresentationTypeCode	string
Jazyk	gmd:LanguageCode	string
Znaková sada	gmd:MD_CharacterSetCode	string
Aktualizácia	gmd:MD_MaintenanceFrequencyCode	string
Mierka	gco:Integer	Integer
Vzdialenosť	gco:Distance	double
Rozsah		

Priestorový rozsah		
Zemepisná dĺžka - západ	gco:Decimal	decimal
Zemepisná dĺžka - východ	gco:Decimal	decimal
Zemepisná šírka - juh	gco:Decimal	decimal
Zemepisná šírka - sever	gco:Decimal	decimal
Časový rozsah		
Od	gml:beginPosition	date
Do	gml:endPosition	date
Distribúcia		
Formát	gco:CharacterString	string
Verzia	gco:CharacterString	string
Odkaz	string	
Kvalita		
Pôvod	gco:CharacterString	string
Súlad	gco:Boolean	boolean
Špecifikácia	gco:CharacterString	string
Kontakt - zdroj		
Organizácia	gco:CharacterString	string
Dodací bod	gco:CharacterString	string
Mesto	gco:CharacterString	string
PSČ	gco:CharacterString	string
Štát	gco:CharacterString	string
telefon	gco:CharacterString	string
email	gco:CharacterString	string
www	gco:CharacterString	string
Rola	gmd:CI_RoleCode	string
Kontakt - metaúdaje		

Organizácia	gco:CharacterString	string
Dodací bod	gco:CharacterString	string
Mesto	gco:CharacterString	string
PSČ	gco:CharacterString	string
Štát	gco:CharacterString	string
telefon	gco:CharacterString	string
email	gco:CharacterString	string
www	gco:CharacterString	string
Rola	gmd:CI_RoleCode	string

Po načítaní všetkých týchto atribútov je dôležité podľa načítaných hodnôt dáta presunúť na správne „miesto“ v dátovom úložisku. Pre jednoduchšiu orientáciu administrátorov, ktorí budú priamo pristupovať do dátového úložiska je dôležité, aby táto forma bola ľudske pochopiteľná.

Adresárová štruktúra, podľa ktorej sa musia nahraté dáta v dátovom úložisku ukladať a následne z daného umiestnenia sprístupňovať je zobrazená v nasledujúcej tabuľke.

Tabuľka 4.3 Adresárová štruktúra pre surové dáta

Názov produktu	
Surove_data_2013	
	boot15020301 – názov adresára
	boot13020301.xml – názov XML

Tabuľka 4.4 Adresárová štruktúra pre spracované produkty

Názov produktu	Identifikátor lokality	Formát	
Ortofotomapy_2013_0.25			
	Zilina		
		GEOTIFF	
			Zilina_1_8.tiff

			Zilina_1_8.xml
		JPG	
			Zilina_1_8.jpeg
			Zilina_1_8.jgw
			Zilina_1_8.xml
	Kosice		
		GEOTIFF	
			Kosice_2_8.tiff
			Kosice_2_8.xml

4.4 ZABEZPEČENIE DÁT V DÁTOVOM ÚLOŽISKU - SÚBOROVÁ ÚROVEŇ

Zabezpečenie súborov na súborovej úrovni je zabezpečené prostredníctvom nastavení vlastníka a skupiny pre ukladané adresáre/súbory a následne aj začepčením súborov, pred prístupom užívateľa bez dostatočných systémových oprávnení.

4.5 ZABEZPEČENIE DÁT V DÁTOVOM ÚLOŽISKU - WEBOVÝ PORTÁL DÁTOVÉHO ÚLOŽISKA

Zabezpečenie dát v dátovom úložisku z pohľadu webového portálu je definované tým, že webový server pristupuje k adresárom nazdieľaným cez sambu. Takto je zabezpečené, že akýkoľvek užívateľ nahrá dáta cez webovú aplikáciu, dáta budú mať vždy rovnaké oprávnenia, ako keby boli nahraté užívateľom priamo pripojeným na sambu. Do webového portálu sa prihlásia len užívatelia, ktorí sú overení voči centrálnemu manažmentu, ktorí po prihlásení a overení či užívateľ môže nahrávať alebo prehľadávať dané dáta, mu zobrazí alebo nezobrazí odkaz na webovú aplikáciu dátového úložiska. Priamy prístup do webových aplikácii dátového úložiska je zakázaný.

4.6 PRÍSTUP DO DÁTOVÉHO ÚLOŽISKA CEZ VPN

Do dátového úložiska je možné sa pripojiť len prostredníctvom VPN. Bez správnych nastavení VPN sa nikto do daného úložiska nepripojí, nakoľko to je prístupné len cez VPN privátnu sieť.

4.7 ŽURNÁLOVÝ PREHEAD AKTIVITY UŽÍVATEĽOV

Z dôvodu zabezpečenia maximálnej kontroly toho, čo užívatelia v systéme vykonávajú, alebo čo samotný systém robí, je dôležité robiť žurnálovacie záznamy toho, čo sa vo webovej aplikácii deje. Preto boli navrhnuté algoritmy, ktoré zabezpečia túto funkcionálnosť a budú zaznamenávať tieto stavy:

- Užívateľ, ktorý operáciu vykonal.
- Samotná vykonaná akcia.
- Dátum a čas, kedy táto operácia bola vykonaná.
- Ak bol stiahnutý súbor, zaznamenať veľkosť stiahnutých dát.

4.8 VYHLÁDÁVANIE SPRACOVANÝCH A ULOŽENÝCH ÚDAJOV

Dátové úložisko musí zabezpečiť vyhľadávanie uložených a spracovaných dát. Dáta je možné vyhľadávať podľa parametrov vo vyhľadávacom formulári. Po vyhľadaní musí dátové úložisko zobrazovať spolu s vyhľadaným súborom aj nasledujúce informácie:

- Názov súboru (bez relatívnej cesty v systéme).
- Veľkosť súboru.
- Dátum vytvorenia z metadát.
- Odkaz na XML súbor s metadátami.
- Možnosť stiahnuť označené a vyhľadané dáta (samotný súbor, viac súborov ako zip alebo cmd), ak užívateľ neprekročil limit na stiahnutie dát.
- Možnosť ďalej spracovať *.las súbory.

4.9 VYHLADÁVACÍ FORMULÁR

Vyhľadávanie dát uložených v dátovom úložisku je rozdelené na základné a rozšírené vyhľadávanie.

Základné vyhľadávanie poskytuje vyhľadávanie na základe týchto kritérií:

- Názov produktu.
- Identifikátor mapového listu – možnosť zadať pri vyhľadávaní viacero mapových listov oddelených čiarkou alebo výberom z mapy.
- Dátum vytvorenia – možnosť zadať rozsah dátumov vybraním z kalendára.
- Formát súboru a verzia.
- Priestorové rozlíšenie – možnosť zadať vyhľadávacie kritérium > < >= <= = <> od – do.

Rozšírené vyhľadávanie poskytuje sadu zvyšných kritérií podľa rozsahu XML metadát:

- ESPG.
- Typ priestorovej reprezentácie.
- Tematická kategória.
- Téma INSPIRE.

V databázovom návrhu boli zohľadnené niektoré polia ako číselníky a úložisko musí zabezpečovať identifikáciu podľa dátového typu v databáze.

Polia ukladané ako číselníky budú:

- Názov produktu
- Formát súboru
- Verzia formátu súboru
- EPSG
- Typ priestorovej reprezentácie
- Tematická kategória
- Téma INSPIRE

Pri číselníkoch boli už pred vyplnené údaje, ktoré sa majú zobrazovať vo webovej aplikácii namiesto anglických názvov v XML. Číselník téma musí mať takto namapované produkty podľa tabuľky nižšie.



Tabuľka 4.5 Mapovanie meta údajov k téme v číselníkoch

Téma	Téma_metaúdaje
poľnohospodárstvo	farming
letecké snímky, základné mapy	imageryBaseMapsEarthCover
obrana, vojsko	intelligenceMilitary
vodné hospodárstvo	inlandWaters
lokalizácia, navigácia	location
more	oceans
kataster, územné plánovanie	planningCadastre
spoločnosť	society
občianska vybavenosť	structure
doprava	transportation
siete	utilitiesCommunication
biota	biota
správne rozdelenie	boundaries
ovzdušie, meteorológia	climatologyMeteorologyAtmosphere
ekonomika	economy
výškopis	elevation
životné prostredie	environment
geológia, geofyzika	geoscientificInformation
zdravie	health

V tabuľke nižšie sú zobrazené namapované hodnoty v číselníku projekcia.

Tabuľka 4.6 Mapovanie meta údajov k projekcii v číselníkoch

Projekcia	Projekcia_metaúdaje
WGS 84	4326
WGS 84 / UTM zone 34N	32634
WGS 84 / UTM zone 33N	32633



Vyhľadávací formulár musí tak isto obsahovať mapový komponent. Tento komponent poskytuje možnosť vybrať mapové listy z mapy a následne musí vyhľadávací formulár zabezpečiť prebratie týchto údajov do vyhľadávania.

4.10 STIAHNUTIE SÚBOROV AKO *.ZIP

Po vyhľadaní a nájdení produktov, ktoré vyhovujú zadaným kritériám dátové úložisko umožňuje označiť viac ako 1 produkt a takto označené súbory umožniť stiahnuť ako 1 *.zip súbor, na ktorého stiahnutie však užívateľ musí mať oprávnenie a nesmie mať prekročený limit na sťahovanie súborov.

4.11 STIAHNUTIE SÚBOROV POMOCOU *.CMD

Tak ako má možnosť po vyhľadaní užívateľ stiahnuť nájdené súbory pomocou zip, má možnosť stiahnuť tieto súbory aj priamo z dátového úložiska pomocou batch súboru. Na to, aby to však užívateľ mohol spraviť, musí mať prístupové údaje na sambu dátového úložiska a musí vyplniť cieľový adresár v jeho počítači, kam budú dané dáta nakopírované.

4.12 SPRACOVANIE *.LAS POMOCOU LASTOOLS

Ak sa medzi nájdenými výsledkami budú nachádzať aj *.las súbory, aplikácia umožňuje spracovanie týchto las súborov. Aplikácia umožňuje automatické načítanie dostupných lastools a umožňuje zadať jednotlivé parametre pre analýzu las súborov pomocou týchto lastools.

4.13 NAHRATIE SHAPE FILE PRE LASCLIP

Ak užívateľ vyberie pri spracovaní las súborov možnosť lasclip, má možnosť do systému nahráť shape súbory, ktoré definujú oblasť, v ktorej sa má analýza daného las súboru vykonať.



4.14 ŽURNÁLOVÝ FORMULÁR

Žurnálovací formulár umožňuje vyhľadávanie v žurnálovacej tabuľke podľa nasledujúcich kritérií:

- Meno užívateľa
- Akcia
- Dátum od
- Dátum do

Žurnálovací formulár pri zobrazení výsledkov používa stránkovanie, tak aby sa na stránke nezobrazovalo príliš veľa výsledkov.

4.15 FORMULÁR PRE IMPORT SUROVÝCH DÁT

Nakoľko surové dáta sú skupinou súborov a adresárov, nie je možné cez webový prehliadač takéto súbory do dátového úložiska nahráť. Surové dáta zaberajú najväčší objem diskového priestoru, preto je potrebné, aby bola možnosť surové dáta nahrávať do dátového úložiska pomocou priameho samba prístupu, do vopred špecifikovaného užívateľského adresára. Tento formulár umožňuje vygenerovanie *.cmd súboru, ktorý pripojí vzdialený samba adresár, pod vopred zvoleným písmenom.

4.16 FORMULÁR PRE IMPORT SPRACOVANÝCH PRODUKTOV

Formulár pre import spracovaných produktov disponuje rozhraním, ktoré poskytuje nasledujúcu funkcionálnosť:

- Nahratie súborov pomocou tlačidla vyber súbory.
- Nahratie súborov pomocou drag & drop.
- Prerušenie prebiehajúceho nahrávania.
- Nahratie pridaných súborov jedným tlačidlom.
- Po nahratí umožniť automaticky ďalšie nahrávanie.
- Zobrazovať informácie o prebiehajúcom nahrávaní.
- Zobrazovať informáciu o aktuálne nahratom objeme súborov z celkového objemu.

- Nahrávať súboru pomocou tzv. „chunk“ súborov.
- Zabezpečiť, aby práve nahrávaný súbor mal názov *.part, ktorý bude identifikovať súbor aktuálne nahrávaný.

4.17 AUTOMATICKÉ MAZANIE STARÝCH SÚBOROV

Dátové úložisko disponuje mechanizmom, ktorý zabezpečí, aby sa súbory (*.zip, *.cmd, *.shp a ostatné dočasné súbory, poprípade nespracované nahraté súbory) po uplynutí vopred stanoveného času mazali z dátového úložiska. Tak isto tento mechanizmus zabezpečuje žurnálovací záznam o zmazaní týchto súborov.

4.18 ROZLIŠOVANIE OPRÁVNENÍ, PODĽA AUTENTIFIKAČNÉHO SERVERA

Webové aplikácie dátového úložiska zobrazujú a poskytujú funkcionality len podľa toho, aké oprávnenia má daný užívateľ podľa centrálného autentifikačného servera.

Webové aplikácie majú tieto oprávnenia:

- DataSeeker – užívateľ, ktorý môže v dátovom úložisku dáta len hľadať ale nemôže sťahovať ani spracovať *.las,
- DataReader – užívateľ, ktorý môže v dátovom úložisku dáta vyhľadávať, spracovať a sťahovať, avšak len podľa nastavenia koľko dát za aký čas,
- RawDataImporter – užívateľ, ktorý môže nahrávať len surové dáta a spracovať ich,
- ProductDataImporter – užívateľ, ktorý môže nahrávať spracované dáta (produkty) a spracovať ich (uložiť) do dátového úložiska,

5 REALIZÁCIA SOFISTIKOVANÉHO PREPOJENIA NA EXTERNÉ ZDROJE DÁT

Pre výber optimálneho uloženia údajov v úložisku boli testované vybrané technológie:

1. Hadoop
2. PostGIS
3. GlusterFS

Na základe vyhodnotenia výsledkov bolo ako vhodné riešenie vybrané GlusterFS

5.1.1 HLAVNÉ VÝHODY GLUSTERFS

Inovativnosť – eliminuje metadáta a môže dramaticky zlepšiť výkon, čo nám pomáha lepšie unifikovať dáta a objekty.

Elasticita – prispôsobený na rast a redukciu veľkosti dát.

Lineárna škálovateľnosť – má schopnosť narásť na petabajty a oveľa viac.

Jednoduchosť – je jednoduchý na správu a nezávislý od jadra aj keď beží v užívateľskom prostredí.

5.1.2 ČO ROBÍ GLUSTERFS ODLIŠNÝ OD INÝCH DISTRIBUOVANÝCH SYSTÉMOV?

Predajný – absencia meta dátového servera poskytuje oveľa rýchlejší súborový systém.

Dostupný – beží na bežne dostupnom hardvéri.

Flexibilný – GlusterFS je iba softvérový súborový systém. Skutočné dáta sú ukladané v natívnych súborových systémoch ako napr. ext4, xfs a pod..

Open source – momentálne ho spravuje Red Hat Inc..

5.1.3 KONCEPT UKLADANIA DÁT

Brick – je v základe hocijaký adresár, ktorý je nazdieľaný naprieč dôveryhodným diskovým poolom.

Dôveryhodný diskový pool – je kolekcia týchto zdieľaných súborov/adresárov, ktoré sú postavené na navrhnutom protokole.

Blokové úložisko – existujú zariadenia, pomocou ktorých sú dáta presúvané naprieč systémami vo forme blokov.

Klaster – aj klaster aj dôveryhodný diskový pool znamenajú to isté pri kolaborácii úložných serverov postavených na definovanom protokole.

Distribuovaný súborový systém – je súborový systém, kde sú dáta rozmiestnené medzi niekoľko uzlov a užívatelia môžu k nim pristupovať bez toho aby vedeli kde sa dáta reálne nachádzajú. Užívatelia si neuvedomujú vzdialený prístup k dátam.

FUSE – je nahrateľný modul jadra, ktorý umožňuje užívateľom vytvoriť súborový systém na jadrom operačného systému, bez použitia akéhokoľvek kódu jadra.

Glusterd – je manažment démon, ktorý je chrbticou súborového systému, ktorý bež stále bez ohľadu na to či sú servery v aktívnom stave.

POSIX – je to rodina štandardov definovaných IEEE ako riešenie pre kompatibilitu medzi variantami Unix-u vo forme API.

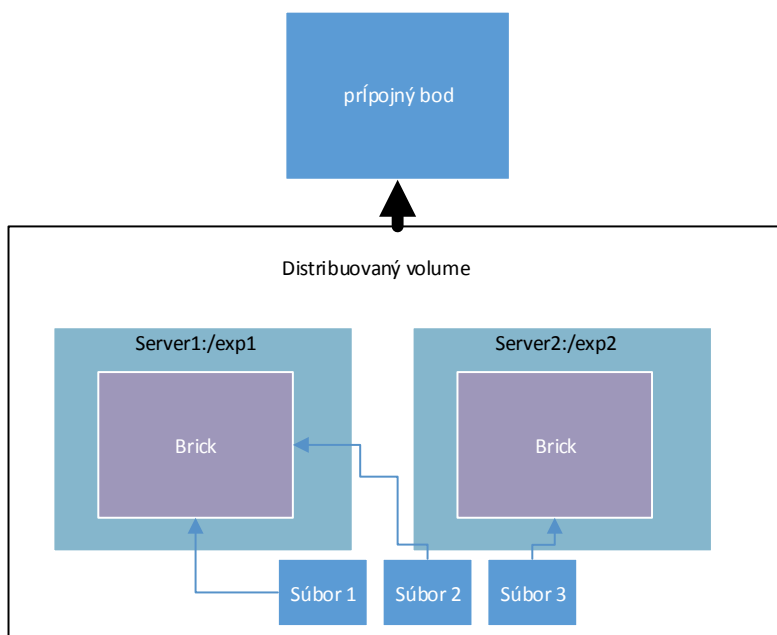
RAID – je technológia, ktorá poskytuje zvýšenú spoľahlivosť diskového úložiska pomocou redundancie.

Subvolume – brick spracovaný minimálne jedným prekladačom.

Prekladač – je základným kusom kódu, ktorý zabezpečuje základné akcie vyvolané užívateľom z prípojného bodu. Pripája jednu alebo viac subvolume.

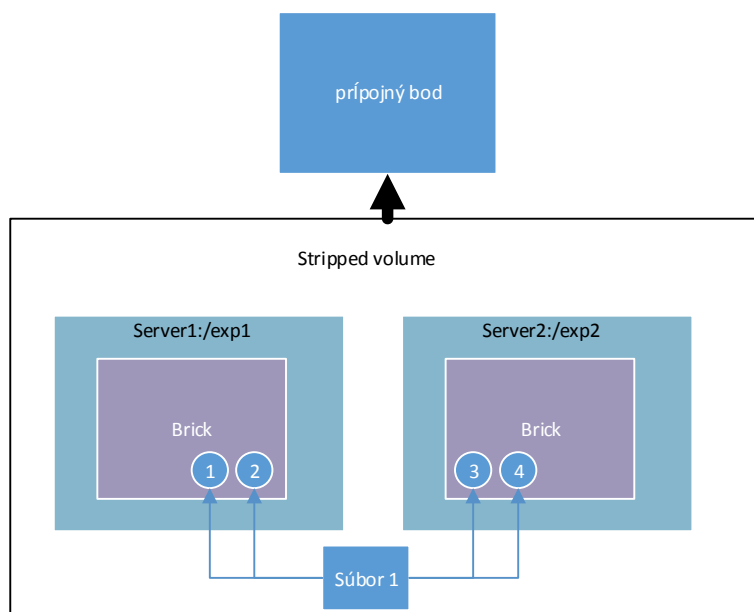
Volume – je logická kolekcia brick. Všetky operácie sú postavené na rozdielnych typoch volume vytvorených užívateľom.

Na obrázku je bloková schéma pre distribuovaný volume.



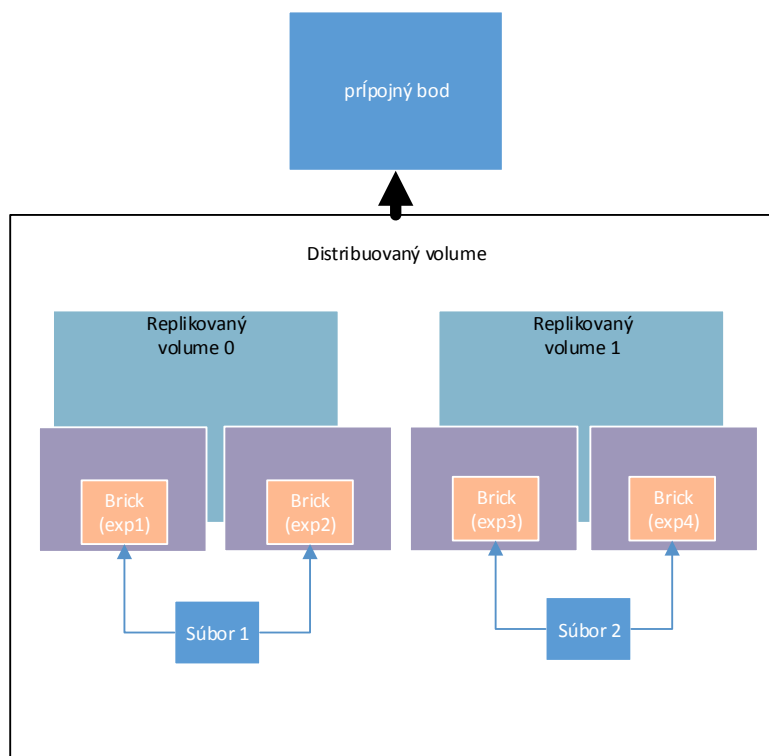
Obr. 5.1 Distribučovaný volume GlusterFS

Na obrázku je znázornený striped volume.



Obr. 5.2 GlusterFS striped volume

Na obrázku je znázornený distribuovaný a replikovaný volume.



Obr. 5.3 GlusterFS distribuovaný replikovaný volume

Nakoľko je GlusterFS súborový systém, ktorý je možné do akéhokoľvek OS pripojiť pomocou jednotného prípojného bodu, poskytuje nám jednoduchú prácu a najmä prístup z pohľadu akéhokoľvek programu, ktorý k nemu bude pristupovať. Nakoľko tento prípojňý bod navonok vystupuje ako hocikáky klasický adresár v operačnom systéme. Vďaka podporovaným funkcionalitám je možné pomocou GlusterFS riešiť aj replikáciu dát v rámci jednej lokality ale tak isto bez rozdielu aj replikáciu medzi rozdielnymi lokalitami. Práca s týmto distribuovaným systémom je oveľa jednoduchšia ako pri Hadoop a najmä je menej náchylný na výpadok, než Hadoop, nakoľko ak pri Hadoop-e spadol namenode alebo sa poškodil, došlo k strate dát.

Do takto pripraveného diskového úložiska je možné ukladať dáta a cez kontrolovaný prístup limitovať užívateľov, ktorí k nim môžu pristupovať. Čiže pomocou GlusterFS sme schopní ukladať surové, spracované dáta a aj fyzické súbory metadát. Nakoľko ale samotne uložené meta dáta súbory nepopisujú uložené dáta, je potrebné ich pri použití GlusterFS v aplikácii importovať do databázy, nad ktorou sa budú robiť volania a tie budú vracať cestu k súboru uloženému v GlusterFS. Ako databázový systém bolo vyhodnotené použiť PostgreSQL.

5.2 POSTGRESQL

PostgreSQL je silný, open source objektovo orientovaný databázový systém. Má za sebou viac ako 15 rokov aktívneho vývoja a overenú architektúru, ktorá získala silné renomé pre svoju spoľahlivosť, integritu dát a korektnosť. Beží na väčšine operačných systémov ako napríklad Linux, Unix (AIX, BSD, HP-UX, SGI IRIX, Mac OS X, Solaris, Tru64) a Windows. Je plne ACID kompatibilný a má plnú podporu pre cudzie kľúče, vnorenia, pohľady, trigre a uložené procedúry (v niekoľkých jazykoch). Obsahujú väčšinu SQL:2008 dátových typov ako napr. integer, numeric, boolean, char, varchar, date, interval a timestamp. Tak isto podporuje ukladanie veľkých binárnych objektov ako obrázky, zvuk a video. Má natívne programovacie rozhranie pre C/C++, Java, .Net, Perl, Python, Ruby, Tcl, ODBC a veľa ďalších.

Zároveň je to databáza podnikovej triedy, PostgreSQL prináša sofistikované vlastnosti ako Multi-Version Concurrency Control (MVCC), obnova z bodu času, tablespaces, asynchrónna replikácia, vnorené transakcie, online/hot zálohy, sofistikované dotazy plánovača/optimizéra a logovanie typu write ahead pre toleranciu výpadkov. Podporuje medzinárodné znakové sady a formátovanie. Je vysoko škálovateľný aj v smere úplného zväčšovania objemu dát, ktorý je schopný spravovať, ale tak isto aj v počte naraz pripojených užívateľov používajúcich databázu. Existujú nasadené aktívne PostgreSQL inštalácie, ktoré manažujú v najvyššom vytížení 4 TB dát.

PostgreSQL sa hrdí dodržovaním štandardov. Jeho SQL implementácia silne vyhovuje štandardu ANSI-SQL:2008. Má plnú podporu pre subqueries (zahŕňajúce subselecty z FROM volaní), potvrdenie čítania a serializácia úrovni transakcií. Zatiaľ čo PostgreSQL má plne relačný systémový katalóg, ktorý sám o sebe podporuje niekoľko schém v databáze, jeho katalóg taktiež prístupný cez Informačnú Schému ako je to definované v SQL štandarde. Vlastnosti dátovej integrity obsahujú (zlučujú) primárne kľúče, cudzie kľúče s reštrikciou a kaskádovaním aktualizácií/mazaní, overovaní obmedzení, unikátne obmedzenia a nenulové obmedzenia.



Tak isto môže hosťovať rozšírenia a pokročilé vlastnosti. Podľa konvencie sú automaticky inkrementované bunky pomocou sekvencií a LIMIT/OFFSET umožňujúci vrátenie čiastočných sád výsledkov. PostgreSQL podporuje zlúčené, unikátne, čiastočné a funkčné indexy, ktoré môžu byť použité v B-tree, R-tree, hash alebo GiST metódach ukladania.

GiST (Generalized Search Tree) – indexovanie je pokročilý systém, ktorý spája dohromady široké pole rozdielnych zoraďovacích a vyhľadávacích algoritmov zahŕňajúcich B-tree, B+-tree, R-tree, čiastočné sum stromy, zoradené B+-tree a veľa ďalších. Tak isto poskytuje rozhranie, ktoré umožňuje vytvorenie vlastných dátových typov a zároveň aj rozšíriteľné metódy volaní, pomocou ktorých ich môžeme vyhľadávať. Preto GiST ponúka flexibilitu pre špecifikáciu toho čo ukladáte, ako to ukladáte a ponúka možnosť definovať nové spôsoby ako to hľadať.

GiST slúži ako nadácia pre veľa verejných projektov, ktoré používajú PostgreSQL, ako napríklad OpenFTS a PostGIS. OpenFTS (Open Source Full Text Search engine) poskytuje online indexovanie dát a ich relevantné hodnotenie pre databázové vyhľadávanie. PostGIS je projekt, ktorý pridáva podporu pre geografické objekty v PostgreSQL, umožňujúci ich využitie ako priestorovú databázu pre geografické informačné systémy (GIS), veľmi podobne ako ESRI SDE alebo Oracle Spatial rozšírenie.

Iné pokročilé vlastnosti zahŕňajú dedenie tabuliek, systém pravidiel a databázové udalosti. Dedenie tabuliek pridáva objektovo orientovaný sklon pri vytvorení tabuľky, čo umožňuje databázovým vývojárom vytvárať nové tabuľky z iných tabuliek.

Systém pravidiel, nazvaný aj query rewrite system, umožňuje databázovým dizajnérom vytvoriť pravidlá, ktoré identifikujú špecifické operácie pre zvolenú tabuľku alebo pohľad a ich dynamickú transformáciu na alternatívne operácie, keď sú spracované.

Systém udalostí je medzi procesný komunikačný systém, v ktorom správa a udalosti môžu byť vymieňané medzi klientmi použitím LISTEN alebo NOTIFY príkazov. Umožňuje aj jednoduchú peer to peer komunikáciu, ale taktiež aj pokročilú koordináciu medzi databázovými udalosťami. Nakoľko môžu byť notifikácie volané pomocou trigrov a uložených procedúr, PostgreSQL klienti môžu monitorovať databázové udalosti ako aktualizácie tabuliek, vložení alebo zmazaní, ak sa stanú.

PostgreSQL spúšťa uložené procedúry vo viac ako tucte programovacích jazykoch zahŕňajúcich Java, Perl, Python, Tcl, C/C++ a jeho vlastnom PL/pgSQL, ktorý je podobný PL/SQL od Oracle. V jeho základnej knižnici funkcií sú stovky funkcií, ktoré sú od základných matematických a reťazcových operácií až po kryptografické a Oracle kompatibilné. Trigre a uložené procedúry môžu byť napísané v jazyku C a nahraté do databázy ako knižnica, ktoré umožňuje vysokú rozšíriteľnosť vlastností. Podobne obsahuje PostgreSQL rozhranie, ktoré umožňuje vývojárom definovať a vytvoriť ich vlastné dátové typy spolu s podporujúcimi funkciami a operáciami, ktoré definujú ich správanie. Výsledkom je vytvorenie serverov s pokročilými dátovými typmi, ktoré sú od geometrických a priestorových primitív až po adresy a ISBN/ISSN (International Standard Book Number/International Standard Serial Number) dátové typy, ktoré môžu byť dodané do systému.

Tak ako je podporovaných veľa procedurálnych jazykov, tak je aj veľa rozhraní pre knižnice v jazykoch umožňujúcich rôzne kompilované a interpretované rozhrania s PostgreSQL. Existujú rozhrania pre Java (JDBC), ODBC, Perl, Python, Ruby, C, C++, PHP, Lisp, Scheme, QT, ktoré sú len časťou podporovaných rozhraní.

Najlepšou vecou je, že zdrojový kód PostgreSQL je dostupný pod liberal open source licenciou: PostgreSQL License. Licencia nám dáva voľnosť použitia, modifikovania a distribuovania PostgreSQL v akejkoľvek forme, otvorené alebo uzavretý kód. Všetky modifikácie, vylepšenia alebo zmeny, ktoré urobíte sú Vaše. PostgreSQL nie je len silný databázový systém schopný bežať v podnikovej sfére, ale je zároveň aj vývojovou platformou, na ktorej je možné vytvárať softvér „na mieru“ alebo komerčný softvér, ktorý potrebuje podporu pre RDBMS.

PostgreSQL je podľa vyššej spomenutej analýzy vhodný pre použitie na ukladanie meta dát extrahovaných z metadáta XML súborov. Použitie PostgreSQL nám dáva priestor dotazovať takto extrahované údaje z rôznych platforiem. Veľkosť tabuliek a databázy, ktorá môže byť použitá v PostgreSQL je taktiež dostačujúca. Rýchlosť odoziev je niekoľko násobne nižšia ako pri Hadoop, čo je samozrejme spojené aj z dôvodu rozdielneho primárneho použitia týchto dvoch systémov.



6 REALIZÁCIA SOFISTIKOVANÝCH PRÍSTUPOV PRE AKCELERÁCIU DÁTOVÝCH OPERÁCIÍ

VYBRANÉ TECHNOLOGIE POUŽITÉ PRE DEPLOY DÁTOVÉHO ÚLOŽISKA

Z analýzy ukladaných súborov vyplynulo, že potrebujeme ukladať tri rozdielne typy súborov (bez ohľadu na príponu, či binárny typ súboru):

- Surové dáta (veľa súborov v rôznych adresároch) – veľkosť cca 100 GB.
- Spracované dáta (súbory s rozdielnymi príponami, ale rovnaký názov) – stovky MB.
- XML metadáta (súbor s rovnakým názvom ako surové dáta, či spracované dáta, príponou XML, ktorý obsahuje meta údaje popisujúce surové alebo spracované dáta).

Na základe tohto členenia nám vyplývajú tieto požiadavky, ktoré musí dátové úložisko spĺňať:

- Uloženie surových, spracovaných dát a XML súboru bezo zmeny do dátového úložiska.
- Rozparsovanie XML súboru, jeho import do databázy s jednoznačnou identifikáciou, ktoré súbory popisuje.
-

6.1 RIEŠENIE PRE DÁTOVÉ ÚLOŽISKO

Pre dátové úložisko bola použitá technológia GlusterFS, nakoľko ako je spomenuté vyššie pre prácu s ním nie je potrebné volať špeciálne príkazy a v systéme vystupuje ako klasický disk.

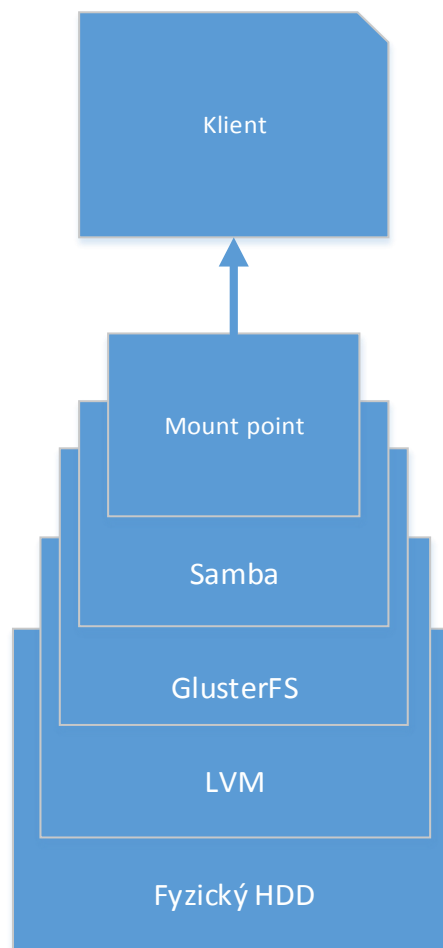
GlusterFS spĺňa tieto požiadavky:

- Replikácia dát
- Tolerancia voči výpadkom
- Georeplikácia dát
- Lineárna škálovateľnosť

Nakoľko GlusterFS je softvérový súborový systém a beží nad klasickým súborovým systémom operačného systému, je potrebné myslieť na bezpečnosť dát a predídeniu ich strate, už pri návrhu tejto vrstvy celého systému. Preto pre konfiguráciu zariadení, ktoré budú priamo pracovať s GlusterFS odporúčame použiť LVM, ktoré je bližšie spomenuté a opísané v kapitole 6.3.

6.2 ARCHITEKTÚRA RIEŠENIA

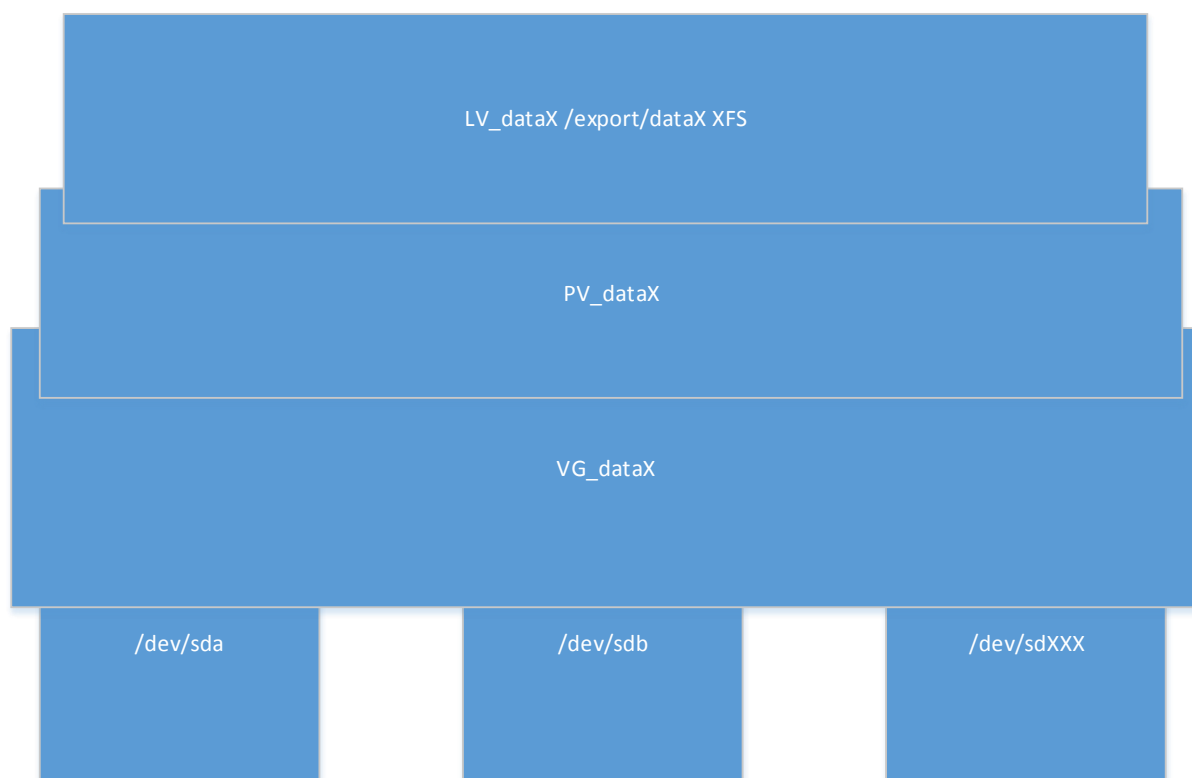
Na obrázku 6.2 je znázornená architektúra diskového úložiska. Na obrázku je vidno jednotlivé úrovne, v ktorých sú jednotlivé časti navrhovaných komponentov, pre samotné uloženie súborov do dátového úložiska.



Obr. 6.1 Architektúra dátového úložiska

Pod klientom sa rozumie, či už operačný systém, alebo samba klient, Windows klient alebo akýkoľvek iný typ klienta, ktorý bude pristupovať k dátovému úložisku. Z návrhu vyplýva aj riešenie samotnej bezpečnosti uložených dát, nakoľko pre klientov je povolený len 1 kontrolovaný prístup do dátového úložiska, a to cez samba politiky, ktoré špecifikujú povolený rozsah IP adries, vlastníka a skupinu vytvorených adresárov a súborov, tak isto aj oprávnenia na adresáre a súbory.

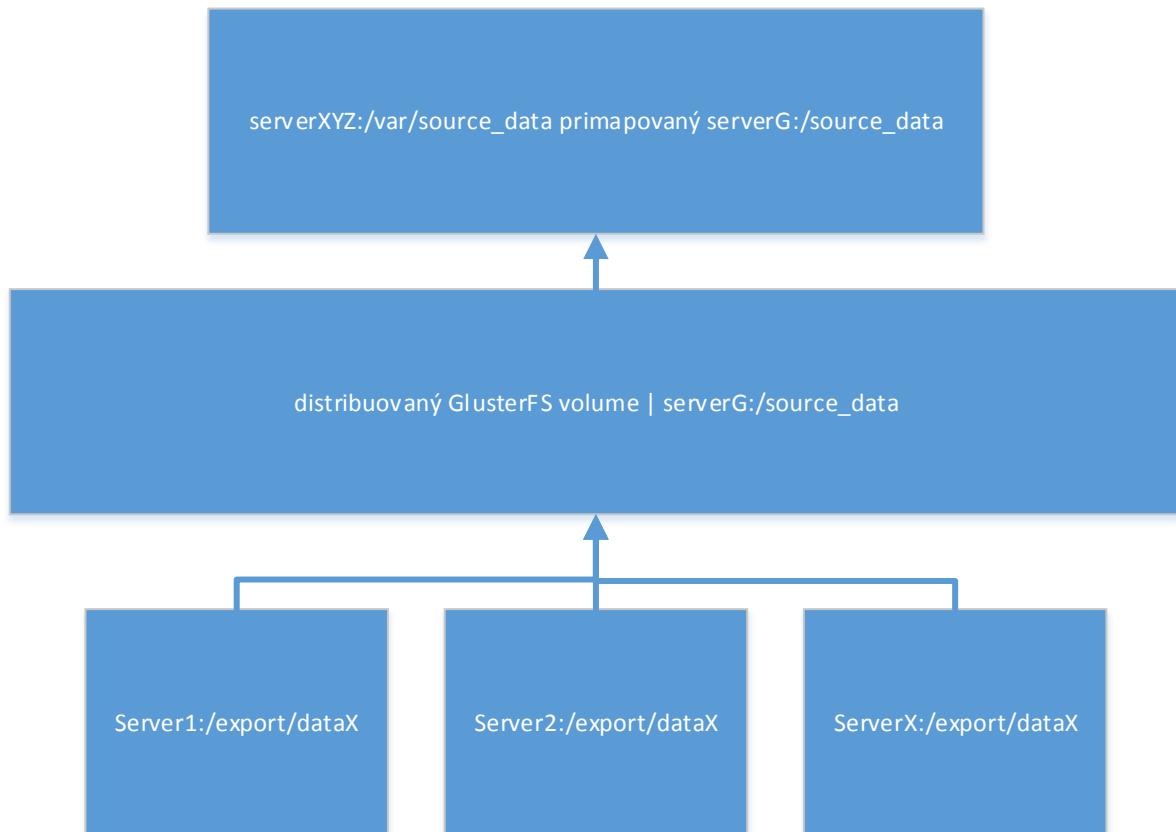
Na obrázku je znázornené plánované zapojenie fyzických diskov do LVM a ich primapovanie do systému.



Obr. 6.2 LVM konfigurácia diskov v dátovom úložisku

Ako z obrázku vyplýva nad fyzickými diskami (či už lokálne pripojenými k serveru, alebo pripojenými pomocou FC, či inej technológie) vytvoríme podľa LVM pravidiel `VG_dataX`, do ktorého bude možné ľubovoľne pridávať ďalšie disky alebo ich odoberať a tým je zabezpečená bez výpadková údržba systému. Nad touto vytvorenou volume group vytvoríme physical volume a v ňom vytvoríme logical volume, ktorý bude primapovaný do systému do adresára `/export/dataX` a s nastaveným súborovým systémom XFS.

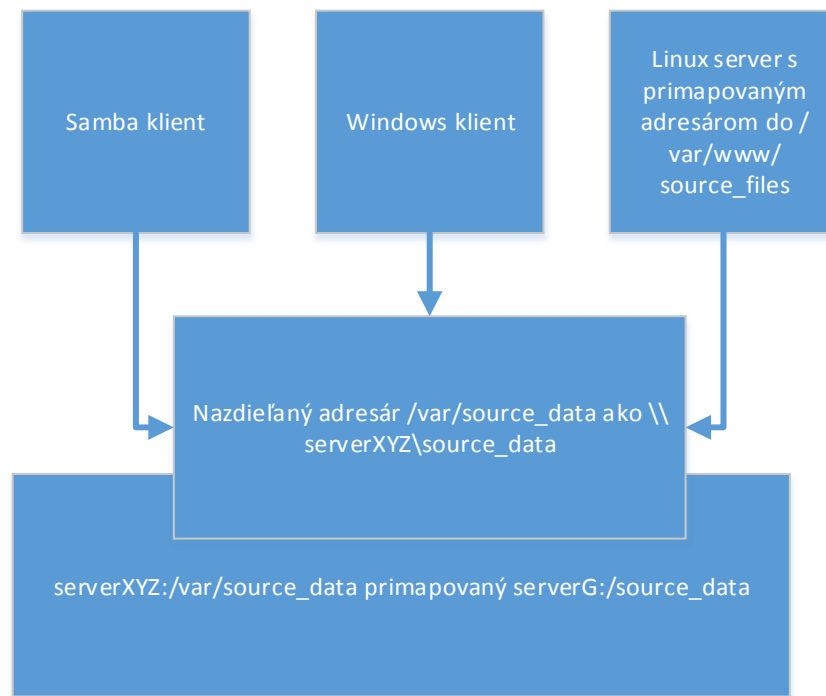
Na obrázku je znázornené zapojenie vytvorených LVM diskov na jednotlivých serveroch do samotného GlusterFS, ktorý bude medzi nimi vytvárať potrebné replikácie dát, kvôli maximálnemu zabezpečeniu voči strate.



Obr. 6.3 Pripojenie LVM diskov do GlusterFS

Z obrázku vyplýva, že LVM nastavené disky použijeme na vytvorenie GlusterFS distribuovaného volume. Tento volume bude obsahovať adresáre z 3 serverov, kde 1 server bude použitý ako master a zvyšné budú použité na replikáciu dát. Z master servera bude následne vyexportovaný takto vytvorený disk, ktorý bude pripojený na každý server do adresára. Tento adresár bude vstupným bodom pre každý server, čiže ak tam nakopíruje systém alebo užívateľ dáta spätnou cestou sa tieto dáta dostanú na fyzické disky a nastane ich replikácia. Servery, na ktorých sú zobrazované replikované dáta, dané dáta neuvidia do momentu, pokiaľ nebudú úspešne synchronizované s hlavným master serverom.

Na obrázku 6.5 je zobrazené vyexportovanie adresára pre priamy prístup užívateľov do dátového úložiska, alebo pre prístup podporných serverov do daného úložiska.



Obr. 6.4 Architektúra pripojenia Samba klientov

Z obrázku vyplýva, že akýkoľvek klient, ktorý bude chcieť priamo pristupovať do dátového úložiska mimo samotných serverov, ktorý budú komunikovať pomocou gluster klienta, budú pristupovať cez 1 vstupný bod. Tým bodom bude samba server, ktorý bude tieto dáta zdieľať. Tým je zabezpečené to, že žiaden iný užívateľ mimo užívateľov špecifikovaných v nastaveniach samby, nezíska prístup k daným dátam.

7 ZOZNAM MERATEĽNÝCH UKAZOVATEĽOV

<i>Študenti doktorandského štúdia vlastnej organizácie a partnerov v projekte, ktorí využívajú poskytnutú podporu - muži</i>						
Názov partnera	Merná jednotka	Východisková hodnota	Rok	Plánovaná hodnota	Rok	Podiel v %
Hlavný partner Žilinská univerzita v Žiline	počet	0	2011	1 splnené	2014	100,00
Spolu	počet	0	2011	1 splnené	2014	100,00

<i>Výskumníci do 35 rokov vlastnej organizácie a partnerov, ktorí využívajú poskytnutú podporu - muži</i>						
Názov partnera	Merná jednotka	Východisková hodnota	Rok	Plánovaná hodnota	Rok	Podiel v %
Hlavný partner Žilinská univerzita v Žiline	počet	0	2011	1 splnené	2014	100,00
Spolu	počet	0	2011	1 splnené	2014	100,00

<i>Výskumníci nad 35 rokov vlastnej organizácie a partnerov, ktorí využívajú poskytnutú podporu - muži</i>						
Názov partnera	Merná jednotka	Východisková hodnota	Rok	Plánovaná hodnota	Rok	Podiel v %
Hlavný partner Žilinská univerzita v Žiline	počet	0	2011	1 splnené	2014	100,00
Spolu	počet	0	2011	1 splnené	2014	100,00

<i>Počet publikácií v nekarentovaných časopisoch</i>						
Názov partnera	Merná jednotka	Východisková hodnota	Rok	Plánovaná hodnota	Rok	Podiel v %
Hlavný partner Žilinská univerzita v Žiline	počet	0	2011	1 splnené	2014	50,00
Partner. č. 3 YMS, a.s.	počet	0	2011	1 splnené	2014	50,00
Spolu	počet	0	2011	2	2014	100,00

<i>Počet prác publikovaných v nerecenzovaných vedeckých periodikách a zborníkoch</i>						
Názov partnera	Merná jednotka	Východisková hodnota	Rok	Plánovaná hodnota	Rok	Podiel v %
Hlavný partner Žilinská univerzita v Žiline	počet	0	2011	1 splnené	2014	50,00
Partner. č. 3 YMS, a.s.	počet	0	2011	1 splnené	2014	50,00
Spolu	počet	0	2011	2	2014	100,00

<i>Objem finančných prostriedkov poskytnutých na projekty venované problematike životného prostredia</i>						
Názov partnera	Merná jednotka	Východisková hodnota	Rok	Plánovaná hodnota	Rok	Podiel v %
Hlavný partner Žilinská univerzita v Žiline	Eur	0	2011	408 300,00	2014	97,84
Partner. č. 3 YMS, a.s.	Eur	0	2011	9 000,00	2014	2,16
Spolu	Eur	0	2011	417 300,00	2014	100,00



- [1] Ihring, Hronček, Holubec, 2014: Možnosti využitia diferenciálnej geometrie pre analýzu digitálnych geografických dát. Aerožurnál II/2014, 1338-8215.
- [2] Holubec, Bobál, 2013: Full-Waveform Lidar Data. INAIR 2013 (1).
- [3] Hoger, Holubec, Otčenášová, Szabová, 2014: Overenie možností mapovania koridorov vonkajších elektrických vedení leteckým lidarom v podmienkach SR. GIS Ostrava 2014, 1213-239X.