

UNIVERZITET U BEOGRADU
MATEMATIČKI FAKULTET



Luka B. Đorović

ANALIZA SLUČAJEVA UPOTREBE
RELACIONIH I KOLONSKI ORIJENTISANIH
NERELACIONIH BAZA PODATAKA

master rad

Beograd, 2024.

Mentor:

dr Saša MALKOV, vanredni profesor
Univerzitet u Beogradu, Matematički fakultet

Članovi komisije:

dr Nenad MITIĆ, redovni profesor
Univerzitet u Beogradu, Matematički fakultet

dr Ivana TANASIJEVIĆ, docent
Univerzitet u Beogradu, Matematički fakultet

Datum odbrane: 15. januar 2016.

Ovaj rad posvećujem...

Naslov master rada: Analiza slučajeva upotrebe relacionih i kolonski orijentisanih nerelacionih baza podataka

Rezime:

Ključne reči: analiza, geometrija, algebra, logika, računarstvo, astronomija

Sadržaj

1	Uvod	1
2	Modeli podataka	3
2.1	Relacioni model	3
2.2	Kolonski-orijentisani model	5
2.3	Glavne razlike između relacionog i kolonski-orijentisanog modela	12
3	Slučajevi upotrebe	17
3.1	Onlajn transakciono procesiranje (OLTP)	17
3.2	Onlajn analitičko procesiranje (OLAP)	18
3.3	Distribuirano okruženje	18
4	Merenje performansi po modelima	21
4.1	Opis platforme za testiranje	21
4.2	Merenje performansi u OLTP okruženju	22
4.3	Merenje performansi u OLAP okruženju	30
4.4	Merenje performansi u distribuiranom okruženju	38
5	Analiza rezultata	41
5.1	Analiza rezultata kod testiranja u OLTP okruženju	41
5.2	Analiza rezultata kod testiranja u OLAP okruženju	42
5.3	Analiza rezultata kod testiranja u distribuiranom okruženju	44
6	Zaključak	47
	Bibliografija	48

Glava 1

Uvod

Podaci su najstabilniji deo svakog sistema. Oni su reprezentacija činjenica i instrukcija u formalizovanom stanju spremnom za dalju interakciju, interpretaciju ili obradu od strane korisnika ili mašine. Iako kroz svoju istoriju računarstvo važi za oblast koja uvodi nove tehnologije i alate neverovatnom brzinom, to nije slučaj za svaku njenu granu. Postoje oblasti koje se kroz istoriju nisu menjale, ili su se slabo menjale i proširivale. Primera za to ima puno i oni su uglavnom usko vezani za funkcionalne principe koji se prožimaju kroz računarske mreže, kompilatore, operativne sisteme, sisteme za upravljanje podacima itd.

Kada je reč o istoriji sistema za upravljanje podacima, mogu se izdvojiti tri faze: period pre relacionih sistema, vreme neprikosnovene vladavine relacionih sistema i nastanak alternativa relacionim sistemima pod grupnim nazivom *NoSQL*.

Do nastanka relacionih sistema za upravljanje podacima, rukovanje podacima izvodilo se kroz pisanje i čitanje iz datoteka operativnog sistema. Rukovanje većim količinama podataka nije bilo standardizovano ni na koji način, već su se konvencije uvodile na nivou organizacija. Apolo sletanje na Mesec realizovano je koristeći ovakav vid rada sa podacima, što ovaj poduhvat čini utoliko neverovatnim [7].

S obzirom da je ovaj vid rada sa podacima imao mnogobrojne mane, među kojima je jedna od glavnih bila komplikovan pristup podacima, javile su se potrebe za unapređenjem. Najuspešniji je bio Edgar F. Codd ¹ koji je 1970. godine objavio rad pod imenom „*A Relational Model of Data for Large Shared Data Banks*” kao rezultat sopstvenih istraživanja i teorija o organizaciji podataka. Kao dokaz da je njegov model moguće implementirati pokrenut je System R ², čiji je rezultat bio

¹Edgar Frank „Ted” Codd (19 Avgust 1923 – 18 April 2003) Američki računarski naučnik

²System R - sistem za rad sa podacima napravljen kao deo istraživačkog projekta IBM-a

i pojava SQL-a (*Structured Query Language*) kao standardizovanog jezika za rad sa podacima. Nakon toga pojavili su se Oracle i IBM sa svojim komercijalnim proizvodima za upravljanje relacionih baza podataka. Naredni period obeležio je rad sa podacima koristeći relacioni model.

Ubrzana digitalizacija, povećana dostupnost interneta, donela je sa sobom potrebu za obradom veće količine podataka. Sve ovo je pokazalo pojedine slabosti dosadašnjih sistema zasnovanih na relacionim modelima, koji nisu mogli u svim segmentima da odgovore na zahteve modernog doba. Ovi problemi poznati pod imenom: problemi velikih podataka (engl. *BigData problems*), doveli su do pojave niza novih modela i principa za čuvanje podataka, kojima je dodeljen grupni naziv: nerelacione baze podataka (engl. *NoSQL*). Sistematizovanje ogromne količine fizičkog prostora na disku na kojem se podaci mogu čuvati i kasnije koristiti, kao i fleksibilnost strukture podataka sa kojima se radi, glavni su problemi tog vremena na koje su se fokusirale tehnologije nastale u *NoSQL* pokretu. Decenije vladavine relacionih sistema za čuvanje podataka ostavile su dubok trag u praksama rada sa podacima, i sa razlogom predstavljaju standard i dan danas, te je eventualno usvajanje tehnologija nastalih u ovoj fazi i dalje česta dilema mnogih stručnjaka.

Kao važna grupa nerelacionih baza podataka izdvajaju se kolonski-orijentisane baze podataka. One su uvele tada nekonvencionalne koncepte čuvanja podataka po kolonama. To podrazumeva sekvencijalno skladištenje vrednosti jedne kolone na disku, sa referencom na red kojem pripada. To sa sobom donosi razne mogućnosti za optimizaciju ali i nove pristupe modelovanja i organizacije podataka. Ovakav način skladištenja ispitivan je još davnih sedamdesetih godina XX veka, međutim u ranim godinama XXI veka došlo je do obnove interesovanja u akademskim ali i industrijskim krugovima.

Nijedan od navedenih koncepata nije univerzalno rešenje, zato je bitno postojanje sadržaja koji se bave analizom slučajeva upotrebe tih tehnologija. Pored teorijske analize koja se može pronaći u relevantnim javnim dokumentacijama korisno je imati i konkretne implementacije testova čiji se rezultati mogu iskoristiti kako bi se povukle paralele u skladu sa potrebama realnih sistema.

Cilj ovog rada je analiza i upoređivanje slučajeva upotrebe relacionih i kolonski orijentisanih baza podataka. Rad će se sastojati iz teorijskog opisa navedenih tehnologija kao i opisa konkretnih predstavnika baza podataka koji će biti korišćeni. Na osnovu teorijskih izvora i istraživanja biće analizirani različiti slučajevi upotrebe.

Glava 2

Modeli podataka

2.1 Relacioni model

Opšte karakteristike

Relacioni model je najpopularniji model za rad sa podacima. On podatke kao i veze između njih predstavlja kroz skup relacija koje predstavljaju skupove torki. Da bi jedan skup torki ili vrsta bila validna relacija u relacionom modelu, on mora ispunjavati sledeće uslove [7]:

- Presek kolone i vrste jedinstveno određuje vrednosnu ćeliju.
- Sve vrednosne ćelije jedne kolone pripadaju nekom zajedničkom skupu.
- Svaka kolona ima jedinstveno ime.
- Ne postoje dve identične vrste jedne relacije.

Iako ovakva formalizacija relacije jeste intuitivna (usled istorijskog uticaja koji je relacioni model ostavio na ideju organizacije podataka) ona je neophodna za definisanje složenijih pojmova.

Koncept ključa relacionog modela

Skup kolona relacije za koji važi da dva reda te relacije nemaju identične vrednosti za svaku kolonu iz tog skupa naziva se *natključ* relacije. Svaki minimalan natključ je *ključ kandidat*. Svaka relacija može imati više ključeva kandidata, a jedan on njih se bira za *primarni ključ* koji mora imati definisanu vrednost

za svaku njegovu kolonu. *Strani ključ* je kolona ili skup kolona čije vrednosti predstavljaju referencu na određeni red neke druge relacije. On uzima vrednost primarnog ključa torke na koju pokazuje.

Primarni i strani ključ igraju veliku ulogu u očuvanju integriteta baze podataka o čemu će biti reči u nastavku.

Integritet relacionog modela

Integritet relacionog modela predstavlja uslove koje podaci treba da zadovolje kako bi stanje u bazi ostalo konzistentno [13]. On se drugačije naziva i „unutrašnja konzistentost” s obzirom da predstavlja aspekte koji mogu da se provere bez konsultovanja domena (npr. ne može se proveriti da li je ime studenta u tabeli ispravno bez konsultovanja domena, ali može se garantovati da će neophodni podaci biti prisutni, uzimati vrednosti iz predviđenog skupa vrednosti i sl.). Provera integriteta se izvršava implicitno ili eksplicitno prilikom svakog ažuriranja baze podataka. Postoji više vrsta integriteta u relacionom modelu: *integritet entiteta*, *integritet domena*, *integritet nepostojeće vrednosti* i *referencijalni integritet*.

Integritet entiteta kaže da svaka torka mora imati definisan primarni ključ bez nedostajućih vrednosti.

Integritet domena predstavlja uslov da za svaku kolonu postoji unapred poznati skup vrednosti koje ona može uzimati.

Integritet nepostojeće vrednosti dodeljuje se kolonama koje ne smeju uzimati nedostajuću vrednost.

Referencijalni integritet nalaže da svaki strani ključ mora imati vrednost primarnog ključa relacije na koju pokazuje.

PostgreSQL

PostgreSQL je objektno-relacioni sistem za upravljanje bazama podataka koji je nastao, a kasnije i bio razvijan na Berkliju, Univerzitet Kalifornija. PostgreSQL je otvorenog koda sa velikom SQL podrškom kao i modernim funkcionalnostima poput: okidača, izmenjivih pogleda, transakcionog integriteta i mnogih drugih [6]. Postgres nudi širok spektar proširenja od strane korisnika poput dodavanja novih tipova podataka, funkcija, operatora, agregatnih funkcija itd.

Kao takav, PostgreSQL je pogodan sistem za čuvanje najkompleksnijih podataka i veza između njih. Mogućnost kreiranja procedura na samoj bazi u integrisanoj

SQL sintaksi, daje široke mogućnosti optimizacije aplikacija.

PostgreSQL koristi server-klijent model funkcionisanja. Sastoji se iz serverskog i klijentskog dela procesa. Serverski deo rukuje fajlovima baze podataka, prihvata konekcije, izvršava konkretne operacije nad bazom. Klijentski deo predstavlja aplikaciju kojom korisnik može da komunicira i rukuje podacima na serverskom delu. Klijent i server komunkiraju preko TCP/IP protokola. Serverski deo može raditi sa više konekcija istovremeno tako što svaka klijentska konekcija radi kao zaseban proces [12].

PostgreSQL iza sebe ima razvijenu društvenu zajednicu, pa samim tim ima dosta izvora i dokumentacije koje mogu olakšati učenje ovog sistema.

2.2 Kolonski-orijentisani model

Opšte karakteristike

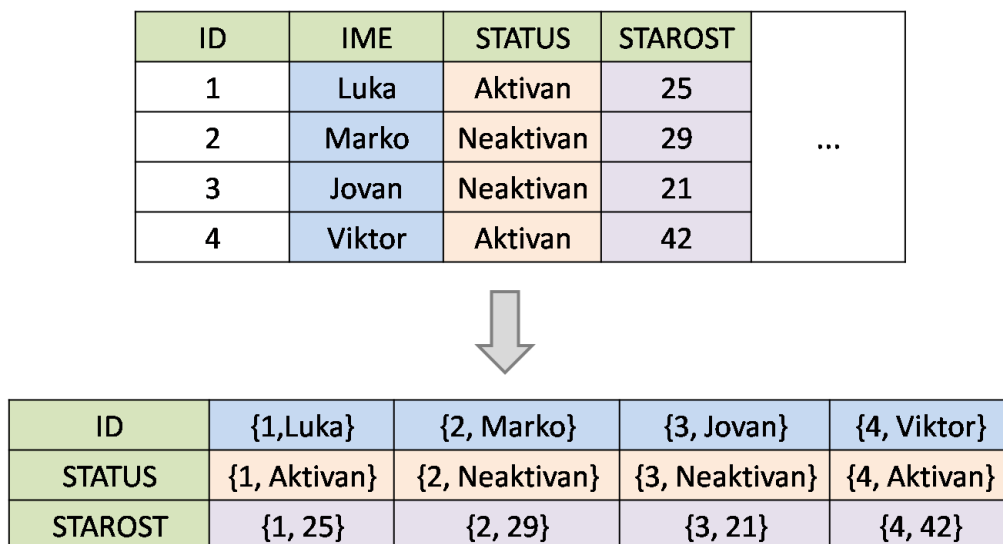
Susret sa problemom *Big Data* doveo je do potrebe za tabelama koje imaju ogroman broj kolona, i ogroman broj redova u okviru tih tabela. Novonastali zahtevi ukazali su na problem kod postojećih relacionih modela. Svaki upit nad tabelom podrazumevao je dohvaćanje svih kolona jednog reda, gde bi se filtriranje nepotrebnih kolona izvršavalo nakon što su se sve kolone učitale u memoriju. Ovo je bila samo jedna od motivacija za implemetanciju sistema zasnovanih na kolonski orijentisanom modelu koji je dizajniran tako da ovakav problem izbegne i uz to donese i druga poboljšanja o kojima će biti reči u nastavku.

Kolonski orijentisan model, podatke na disku skladišti po kolonama, a ne po redovima kao što je to slučaj kod relacionih modela, slika 2.1. Sve vrednosti kolone svih redova skladište se jedna do druge, a na konkretnu vrednosnu ćeliju referiše se pomoću ključa konkretnog reda kao i kolone čiju vrednost želimo da pročitamo. Ovakav dizajn doveo je do toga da za dohvaćanje određenog skupa kolona nema potrebe da čitamo sve vrednosti tog reda, već je dovoljno da znamo konkretan ključ tog reda kao i imena kolona čije vrednosti želimo da pročitamo i tako izbegnemo višak operacija čitanja sa diska.

Ovakav vid skladištenja podataka sa sobom nosi veliki potencijal za primenu raznih algoritama za kompresiju podataka. Kompresija nad sličnim podacima koji se nalaze na uzastopnim adresama u memoriji, omogućava izbegavanje čuvanja složenih meta informacija u okviru struktura koje se koriste za tu kompresiju, što

ovaj model čini posebno pogodnim za njihovu primenu.

Kolonski orijentisan model kao i većina ostalih nerelacionih modela, nudi fleksibilnost strukture podataka koja se ogleda u neograničenom broju kolona, što daje dosta prostora za eksperimentisanje sa dizajnom baze podataka. Primer toga biće prikazan u okviru analize OLAP slučaja upotrebe.



Slika 2.1: Kolonski orijentisan format

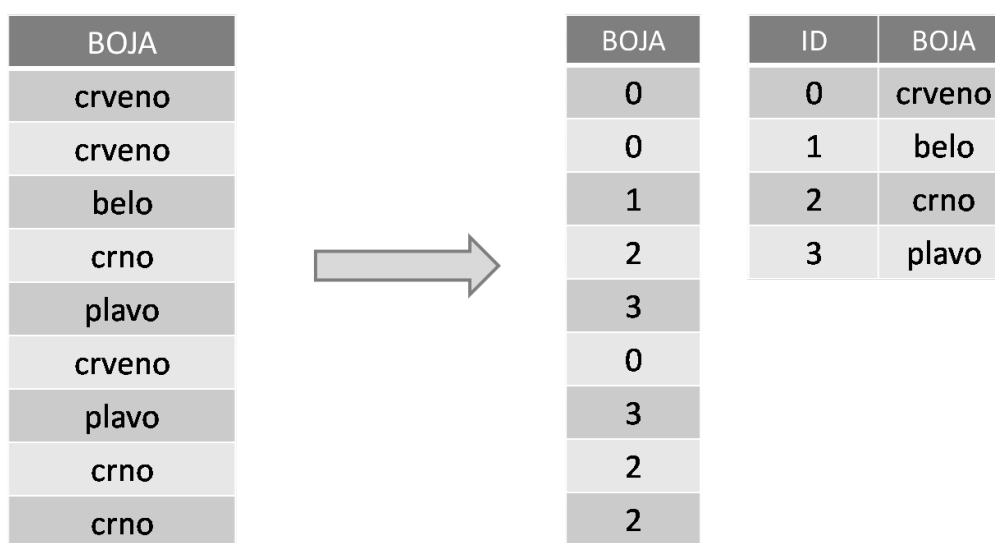
Popularni algoritmi kompresije kolonski orijentisanog modela

Neke od najpoznatijih algoritama kompresije koje kolonski orijentisan model koristi i koji će biti opisani u nastavku jesu: enkodiranje zasnovano na rečniku, enkodiranje po broju ponavljanja i delta enkoding.¹

¹Važno je napomenuti da relacioni modeli imaju svoje mehanizme kompresije podataka koji u ovom radu nisu analizirani.

Enkodiranje zasnovano na rečniku (engl. *Dictionary based encoding*), slika 2.2., funkcioniše tako što se napravi mapa vrednosti koja sadrži svaku vrednost kolone koja je prisutna među podacima. Kao vrednost kolone tada se ne upisuje konkretna vrednost, već ključ iz rečnika koji je mapiran na tu vrednost. Veličina ključa je srazmerna veličini mape, te je ovaj vid kompresije najpogodniji za kolone koje imaju mali broj vrednosti koje se ponavljaju.

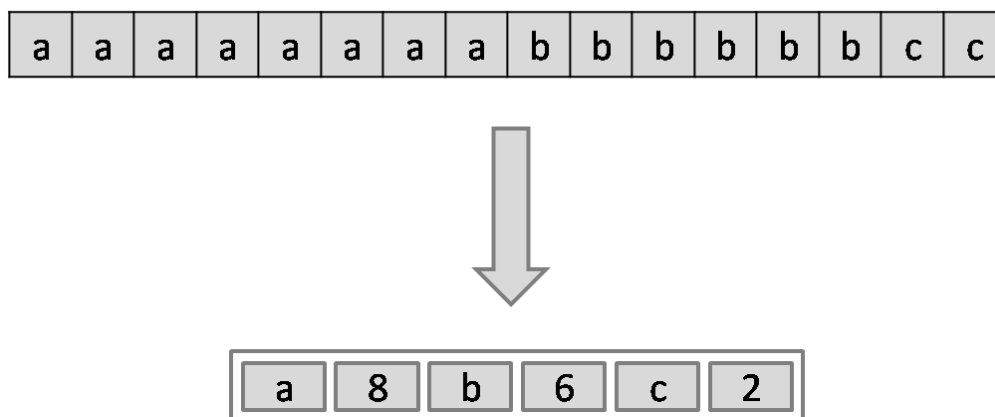
ENKODIRANJE ZASNOVANO NA REČNIKU



Slika 2.2: Enkodiranje zasnovano na rečniku

Enkodiranje po broju ponavljanja (engl. *Run Length Encoding*), slika 2.3., funkcioniše tako što se uz svaku vrednost koja se ponavlja čuva i broj ponavljanja te vrednosti. Na taj način se izbegava pojava duplikata na uzastopnim adresama u memoriji. Ovaj vid kompresije najpogodniji je za kolone koje su sortirane i imaju ponavljajuće vrednosti.

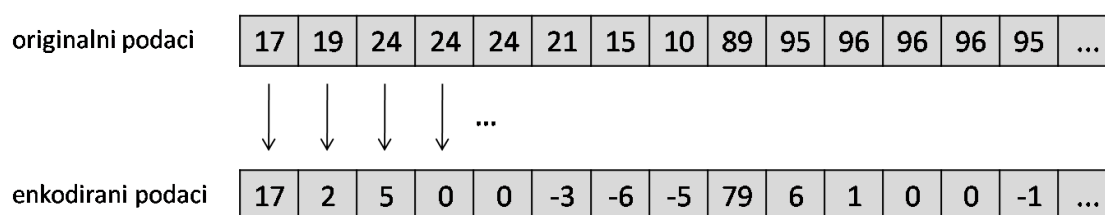
ENKODIRANJE PO BROJU PONAVLJANJA



Slika 2.3: Enkodiranje po broju ponavljanja

Delta enkodiranje algoritam kompresije, slika 2.4. funkcioniše tako što se u kolonama ne čuvaju same vrednosti već razlike između uzastopnih vrednosti. Očigledan primer primene ove kompresije je datumska kolona. U tom slučaju je dovoljno da izaberemo neki referentni datum i da za ostale vrednosti čuvamo razliku u odnosu na njega.

DELTA ENKODIRANJE



Slika 2.4: Delta enkodiranje

HBase

HBase je kolonski orijentisana nerelaciona baza podataka nastala 2007. godine kao prototip *BigTable* baze koja je modelovana u okviru Google-ovog članka iz 2006.godine[9].

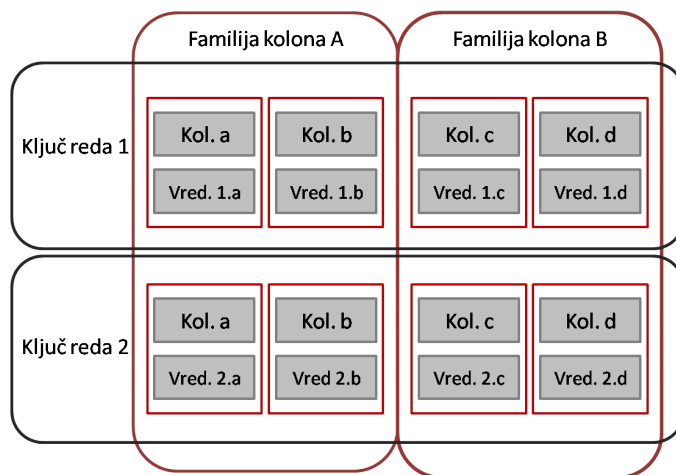
Model podataka koji HBase koristi podrazumeva da se svaka tabela sastoji iz familije kolona, a da svaka familija kolona sadrži određeni skup kolona. S obzirom da će se kolone koje pripadaju jednoj familiji sladištiti blizu na disku, cilj je da atributi, odnosno kolone koje su po prirodi slične, pripadaju istoj familiji kolona, kako bi se nad njima mogli primeniti algoritmi kompresije. HBase nudi fleksibilnost strukture podataka, što znači da, da bismo neki podatak skladištili ne moramo unapred da definišemo skup kolona koji pripada nekoj tabeli, ali moramo definisati skup familija kolona te tabele.

Redovi u HBase-u, odnosno njihovi ključevi sortirani su rastuće, tako da se za pretragu po ključu koristi binarna pretraga. Ovo svojstvo daje na važnosti dizajniranju ključa jer je poželjno da svako čitanje podataka ide po ključu ili njegovom prefiksu [11].

Vrednosnu ćeliju u HBase tabeli određuje ključ reda, ime kolone te vrednosne ćelije, kao i familija kojoj kolona pripada, slika 2.5.

HBase nije ACID baza podataka, ali nudi sledeće garancije [10]:

- Atomičnost pri radu sa jednim redom tabele koja se ogleda u tome što će svaka izmena reda u potpunosti uspeti, ili u potpunosti propasti.
- Svako čitanje reda iz tabele vratiće stanje reda koje je bilo aktuelno najranije u trenutku kada je čitanje započeto.



Slika 2.5: HBase model

HBase podatke može čuvati na lokalnom fajl sistemu ili na *Hadoop* distribuiranom fajl sistemu(HDFS). HDFS je distribuiran fajl sistem koji ima visok prag

tolerancije na greške (*fault tolerance*). HDFS klaster sastoji se iz master čvorova (*Namenode*) i čvorova sa podacima (*DataNode*). Master čvor radi sa metainformacijama samog fajl sistema, kao što su izmena direktorijuma, brisanje direktorijuma i sl. Fajlovi na HDFS su interno podeljeni na blokove, gde se svaki od blokova može replikovati i skladištiti na bilo kom od čvorova sa podacima. O mapiranju blokova fajlova na čvorove sa podacima brine master čvor. Čvorovi sa podacima odgovorni su za izvršavanje operacija čitanja i pisanja fajlova na zahtev klijenata. Pored toga čvorovi sa podacima mogu kreirati, brisati blokove a i replikovati ih, u skladu sa instrukcijama master čvora [2].

Arhitektura HBase klastera sastoji se iz tri glavne komponente: master server, region server i *zookeeper*.

Master server je komponenta HBase klastera koja se bavi metainformacijama tabela HBase-a. On region serverima dodeljuje regione. Pored toga, master server bavi se i balansiranjem opterećenja klastera (*load balancing*).

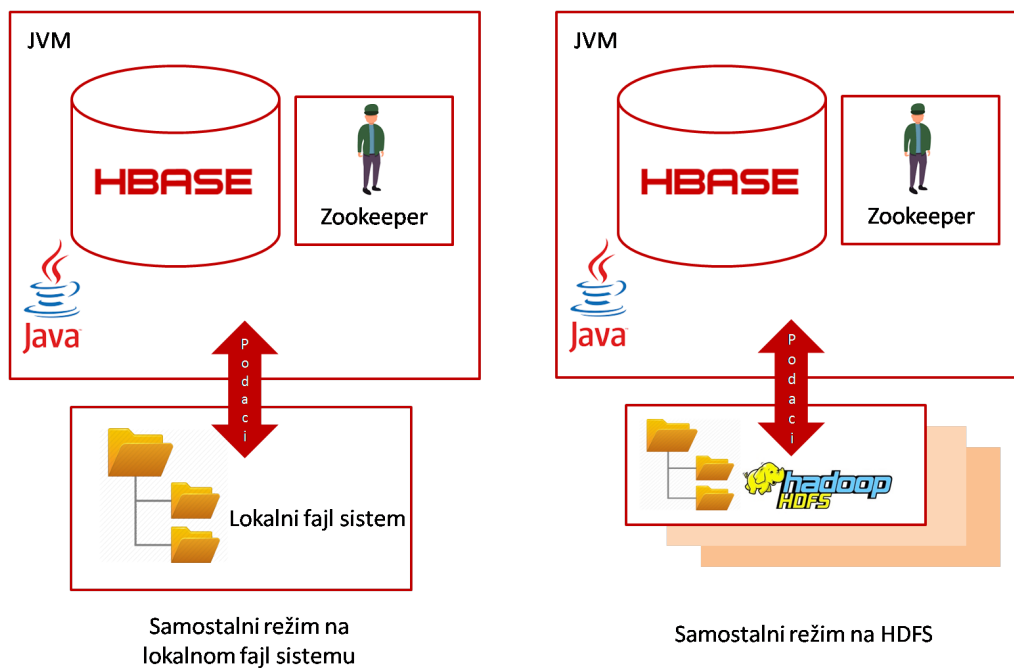
Region server zadužen je za rad sa dodeljenim regionima od strane master servera. Region je skup redova HBase-a u nekom rasponu ključeva. Kako bi region server mogao da piše i čita podatke regiona, on održava WAL (engl *Write Ahead Log*) fajl, MemStore fajl kao i HFile fajlove. Svaka izmena, pre nego što je sačuvana na disku, upisuje se u vidu loga (*commit log*) u WAL, tako da u slučaju pada sistema pre nego što su se izmene sačuvale, logovi koji se nalaze WAL fajlu mogu rekonstruisati stanje pre pada sistema [8]. MemStore je fajl u koji se upisuju podaci pre nego što se sačuvaju na disku. Predstavlja vid bafera koji kada se popuni reflektuje podatke na disk, odnosno na HFile. HFile fajl sadrži konkretne podatke u predviđenom formatu. Oni se mogu skladištiti na lokalnom fajl sistemu ili na HDFS, u zavisnosti od konfiguracije.

Zookeeper predstavlja most u komunikaciji komponenti HBase klastera [1]. Zookeeper servis sinronizuje sve master i region servere, ima evidenciju o tome koji je master server aktivan a koji pasivan, koji region serveri više nije dostupni itd.

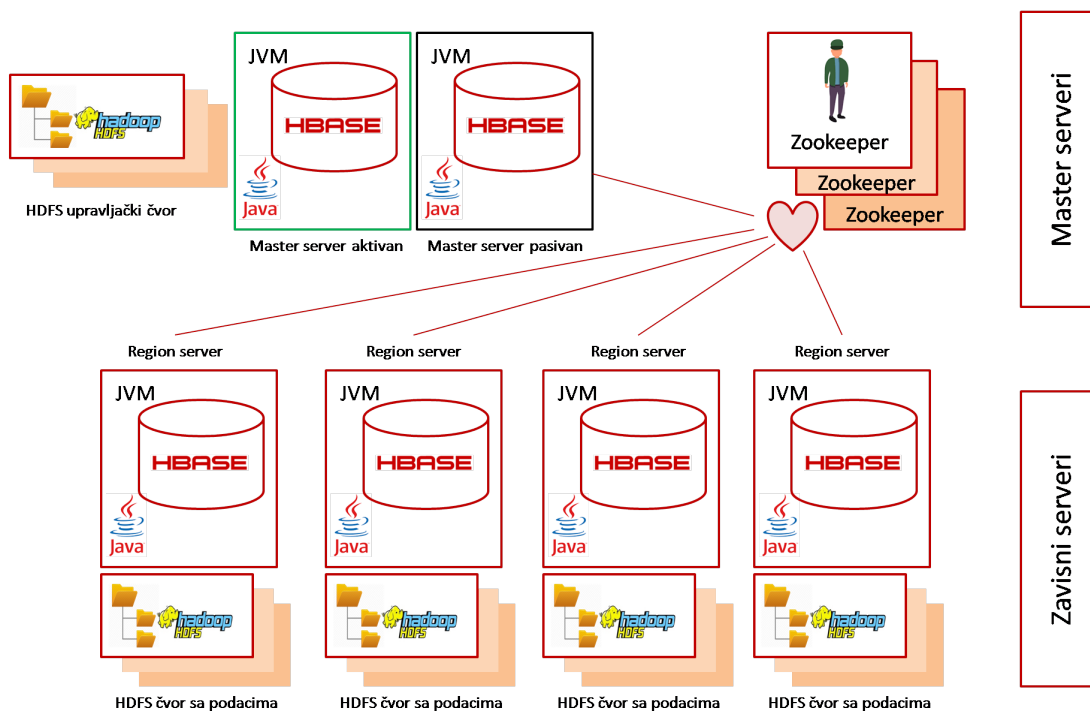
HBase dozvoljava dva režima: samostalan i distribuirani režim.

U samostalnom režimu HFile fajlovi se mogu čuvati na lokalnom fajl sistemu i korišćenje HDFS-a je opciono, slika 2.6.

Sa druge strane u distribuiranom režimu neophodno je korišćenje HDFS-a za skladištenje, slika 2.7.



Slika 2.6: HBase - samostalan režim



Slika 2.7: HBase -distribuirani režim

2.3 Glavne razlike između relacionog i kolonski-orijentisanog modela

Normalizacija i denormalizacija

U relacionim modelima često se radi na izbegavanju *redudantnosti* u podacima. Redudantni podaci zauzimaju višak prostora na disku i otežavaju kasnije održavanje sistema. Kako bi se izbegla redudantnost, postoje postupci koji nam pomažu da organizujemo podatke tako da redudantnost umanjimo. Proces izmene logičkog modela baze podataka u cilju oslobađanja od redundantosti podataka naziva se normalizacija. U zavisnosti od toga koja pravila zadovoljava određena relacija, dodeljuje joj se odgovarajuća normalna forma. Neke od normalnih formi relacionog modela su: 1. normalna forma, 2. normalna forma, 3. normalna forma, normalna forma elementarnog ključa, Bojs-Kodova normalna forma, 4. normalna forma, normalna forma esencijalnih torki, normalna forma bez redundansi, normalna forma superključeva, 5. normalna forma, normalna forma domena i ključa.

Normalizovani modeli obično raspolažu velikim brojem stranih ključeva što dovodi do povećanja broja tabela kojima se pristupa u upitima, a time i povećanja broja operacija čitanja sa diska, što može uticati na performanse.

Denormalizacija je strategija koja se koristi kod modela kod kojih je neophodno ubrzati operacije čitanja podataka, odnosno umanjiti broj tabela kojima je neophodno pristupiti kako bi se neki skup podataka pročitao iz baze podataka.

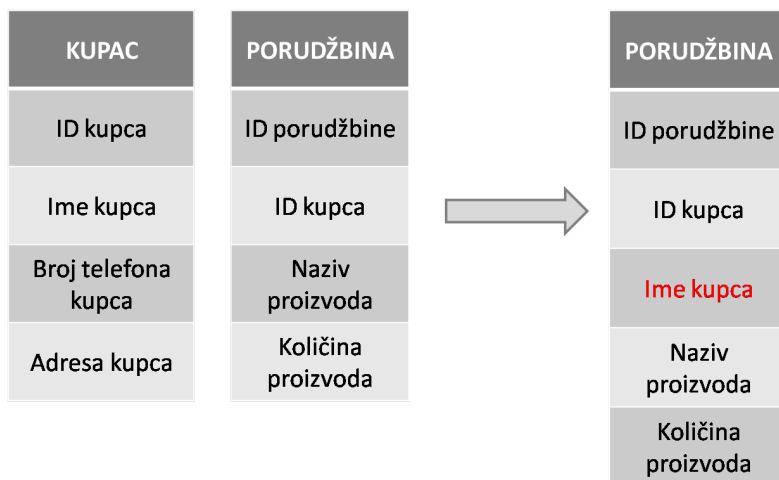
Neke od tehnika denormalizacije su: spajanje kolone, horizontalna podela tabele, vertikalna podela tabele i uvođenje izvedene kolone.

Spajanje kolone, slika 2.8. je dodavanje kolone kojoj bi se često pristupalo preko stranog ključa.

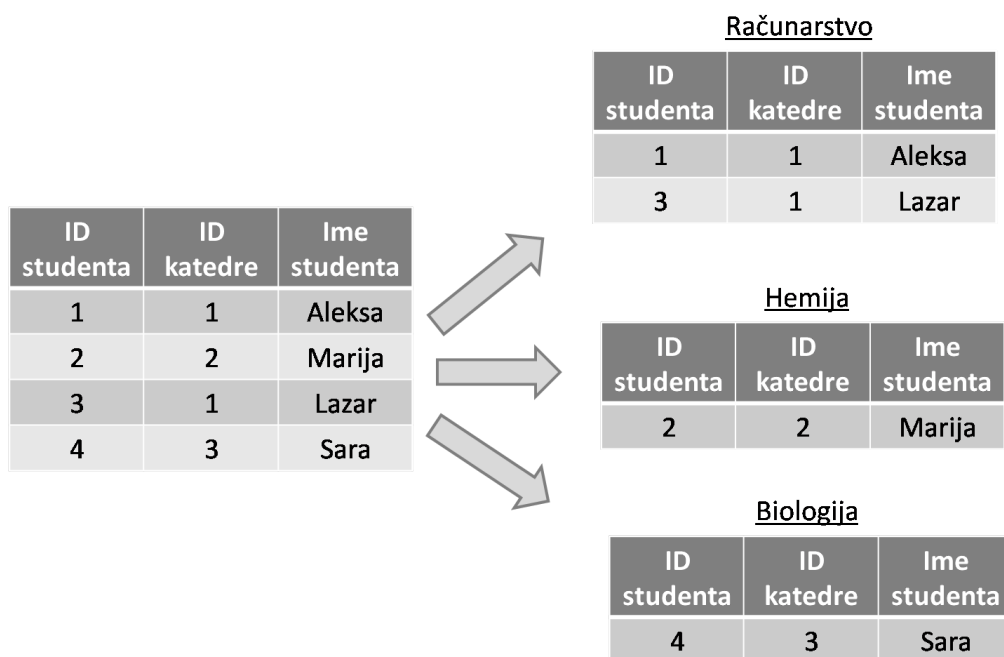
Horizontalno deljenje tabele, slika 2.9. podrazumeva da se na osnovu prirode podataka jedna tabela podeli na više tabela tako da se čitanje svede samo na čitanje grupe redova.

Vertikalno deljenje tabele, slika 2.10. podrazumeva da se tabela podeli grupisanjem kolona koje se često čitaju zajedno.

Uvođenje izvedene kolone, slika 2.11. predstavlja dodavanje kolone koja čuva rezultat neke agregatne funkcije. Time se izbegava da se pri svakom čitanju ta agregatna funkcija izvršava, već se pri ažuriranju stanja ta vrednost ažurira, da bi se rezultat kasnije mogao samo pročitati.



Slika 2.8: Spajanje kolone

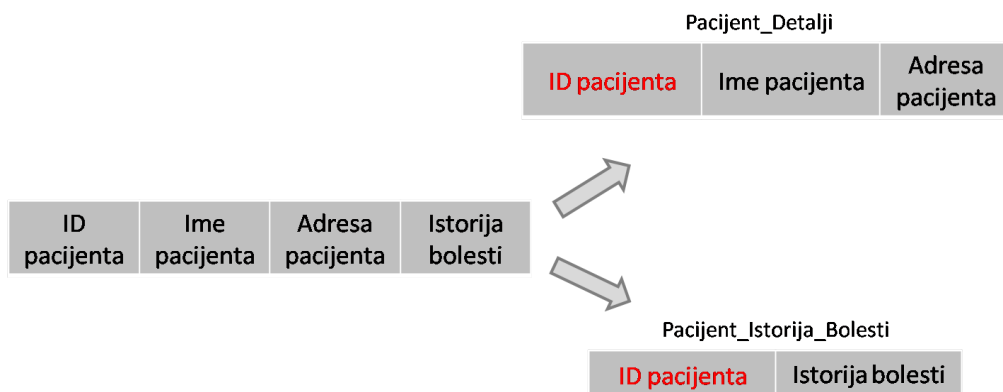


Slika 2.9: Horizontalna podela tabele

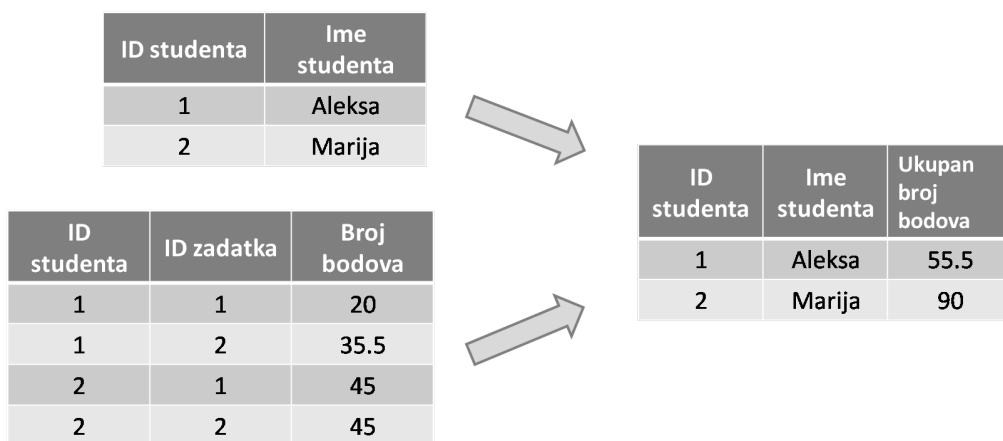
ACID i BASE

Transakcija je logička jedinica posla pri radu sa podacima [13]. Kod relacionih modela jednu transakciju karakteriše: atomičnost, konzistentost, izolovanost i trajnost (ACID).

Atomičnost transakcije se može objasniti pravilom: Jedna transakcija se izvršava u celini ili se ne izvršava nijedan njen deo, odnosno dejstvo transakcije je



Slika 2.10: Vertikalna podela tabele



Slika 2.11: Uvođenje izvedene kolone

nedeljivo.

Konzistentost transakcije znači da dejstvo transakcije ne može ostaviti stanje koje narušava integritet baze podataka.

Izolovanost transakcije čini da transakcije ne mogu uticati međusobno jedna na drugu, odnosno, kada se jedna transakcija pokrene, pa sve dok se ne završi, za nju, izmena neke druge transakcije neće biti vidljiva.

Trajnost transakcije garantuje da će kompletirana transakcija u slučaju prekida rada sistema pre nego što su izmene upisane na disk, biti upamćena i izvršena nakon restarta sistema. Svaka izmena se upisuje u log fajl pre nego što je upisana na disk, kako bi se operacije mogle poništiti u slučaju poništavanja transakcije.

Ova svojstva obično karakterišu transakcije u relacionom modelu, kada konzistentost baze ima veći prioritet od brzine i dostupnosti. Kod kolonski orijentisanih baza podataka, posebno u distribuiranom režimu, dostupnost i brzina često imaju

veći prioritet od stalne konzistentosti. Kod njih se koristi alternativni pristup karakterizacije transakcije - BASE. BASE transakcije zadovoljavaju sledeća svojstva: suštinska raspoloživost, postojanje mekog stanja i konvergentna konzistentnost.

Suštinska raspoloživost omogućava maksimalnu dostupnost čitanja i pisanja, ali bez garancije konzistentosti.

Meko stanje daje mogućnost da se stanje baze podataka menja čak i kada nijedna transakcija nije u toku, i to u periodu postizanja konzistentnosti.

Konvergentna konzistentnost obezbeđuje da će baza podataka u slučaju da nema novih upisa, postati konzistentna kroz neki vremenski period.

CAP teorema

Kako bi se postigla skalabilnost, umanjilo kašnjenje odgovora (engl *latency*) i povećala dostupnost sistema u slučaju nepredviđenih okolnosti, podaci često bivaju distribuirani na više povezanih jedinica. Takve baze podataka nazivamo distribuiranim. Kvalitet distribuiranih baza podataka ogleda se kroz tri svojstva: konzistentost, raspoloživost i tolerancija razdvojenosti.

Konzistentost u ovom kontekstu razlikuje se od konzistentosti transakcije. U kontekstu distribuiranih sistema, konzistentnost je svojstvo baze podataka, koje garantuje da će odgovor koji se šalje sa bilo kog čvora klastera (ukoliko odgovora ima) biti ažuran.

Raspoloživost je svojstvo koje obezbeđuje da će baza uvek vratiti nekakav odgovor.

Tolerancija razdvojenosti znači da će u slučaju delimičnih otkaza unutar klastera, sistem i dalje funkcionisati.

CAP teorema koju je formulisao Eric Brewer ², navodi da distribuirana baza podataka ne može istovremeno ispunjavati sva tri svojstva, slika 2.12.

Obzirom da je konzistentost i dostupnost u praksi skoro nemoguće dostići, distribuirane baze podataka organizuju se u skladu sa tim da li se veći prioritet daje dostupnosti ili stalnoj konzistentosti. Priroda sistema koji koriste relacioni model obično je takva da daju prioritet konzistentnosti, pa se od relacionog modela očekuje da u distribuiranom okruženju osim konzistentnosti nudi i toleranciju razdvojenosti, za razliku od kolonski orijentisane nerelacione baze koja konzistentnost ne garantuje (garantuje konvergentnu konzistenciju), ali uz toleranciju razdvojenosti nudi stalnu dostupnost.

²Eric Allen Brewer profesor računarских nauka na Univerzitetu Kalifornija, Berkli



Slika 2.12: CAP teorema

Glava 3

Slučajevi upotrebe

Kako bi se postigao dovoljan dokaz koncepta (engl. *proof of concept*) kreirano je više scenarija koji na različit način pristupaju radu sa podacima. Scenariji, odnosno slučajevi upotrebe koji su testirani u ovom radu jesu:

1. Onlajn transakciono procesiranje (OLTP)
2. Onlajn analitičko procesiranje (OLAP)
3. Distribuirano okruženje

3.1 Onlajn transakciono procesiranje (OLTP)

Onlajn transakciono procesiranje predstavlja procesiranje podataka koje se sastoji iz velikog broja transakcija. Svaka pojedinačna transakcija obično deluje na manjem skupu podataka i obično uključuje njihovu izmenu, a kada se radi čitanje podataka, ono je obično po ključu [5]. Kod ovakvog vida procesiranja, gde su izmene podataka česte, čuvanje integriteta podataka prilikom tih upisa je prioritet, pa je OLTPs često sastavni deo finansijskih sistema, sistema za rezervacije i sl. Konkretan scenario koji će u radu simulirati OLTP okruženje je prebacivanje sredstava sa jednog računa na drugi. Svako prebacivanje novca predstavljaće jednu poslovnu transakciju koja se sastoji iz više transakcija baze podataka. Definicija poslovne transakcije koja se koristi za testiranje delom je preuzeta iz TPC-C specifikacije [3].

Jedna poslovna transakcija sastoji se iz tri transakcije baze podataka: kreiranje transfera sredstava (*createFXTransaction*), izvršavanje uplate (*executePayment*) i provera statusa transfera (*checkTransactionStatus*).

CreateFXTransaction sastoji se iz provere sredstava na računu korisnika i unosa novog transfera u tabelu.

ExecutePayment radi ažuriranje računa korisnika i menja status transfera

CheckTransaction dohvata status transfera.

3.2 Onlajn analitičko procesiranje (OLAP)

Onlajn analitičko procesiranje sačinjeno je od velikog broja čitanja podataka i manjeg broja izmena podataka (koje su obično masivne - *bulk*). Upiti koji se koriste obično imaju parametre, visok nivo kompleksnosti i visok procenat podataka kojima pristupaju. OLAP je obično sastavni deo sistema koji nude podršku za donošenje poslovnih odluka, generisanje složenih izveštaja kao i bilo koji vid korišćenja velike količine podataka koristeći složene upite i operacije. Scenario koji se simulira u ovom radu u svrhu testiranja, jeste primer funkcionisanja trgovinskog lanca sa skupom svojih mušterija, proizvođa, dobavljača i narudžbina. Korišćen model delom je preuzet iz TPC-H specifikacije [4]. Test će se sastojati iz paralelnog izvršavanja složenog upita koji iz tabele sa stavkama narudžbina izvlači vrednosti agregatnih funkcija određenih kolona grupisanih po statusu. Parametri upita su dan slanja narudžbine i status po kojima filtriramo rezultat.

3.3 Distribuirano okruženje

Distribuirani sistemi su sistemi kod kojih se podaci nalaze na više povezanih jedinica, odnosno čvorova kako bi se postigla skalabilnost, umanjilo kašnjenje odgovora (engl *latency*) i povećala dostupnost sistema u slučaju problema. Dva načina distribuiranja podataka su: **replikacija** i **particionisanje**.

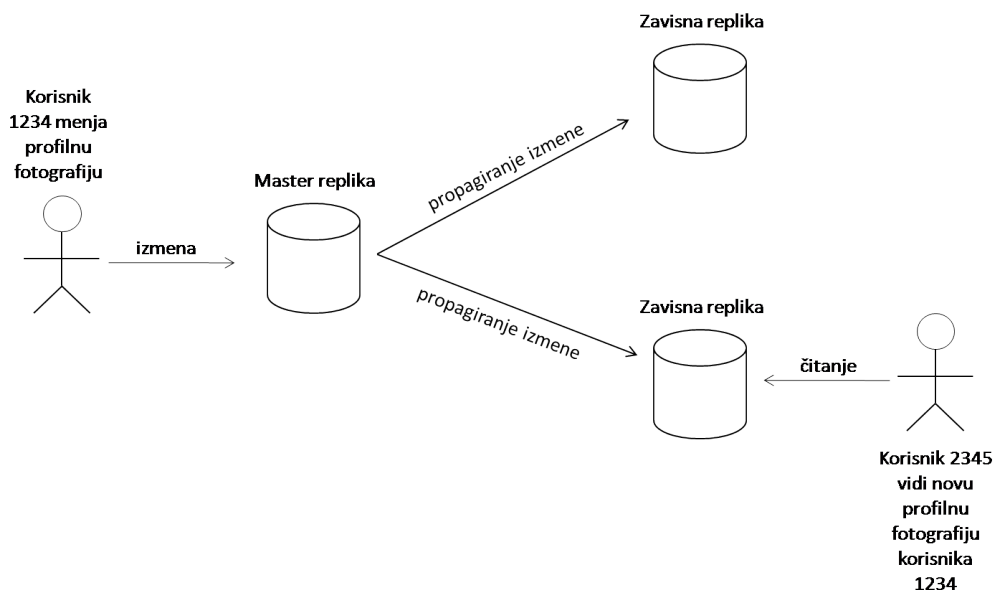
Particionisanje distribuira podatke tako što se veliki skupovi podataka podele na manje celine i čuvaju na različitim čvorovima.

Replikacija podrazumeva da se više kopija podataka čuva na različitim čvorovima. Ukoliko neki od čvorova postane nedostupan, klijentski zahtevi mogu biti obrađeni na nekom od preostalih čvorova. Kako bi svaka replika bila ažurna neophodno je da se svaka izmena podataka propagira do svakog čvora. Najpoznatiji model po kojem se to realizuje je *master-slave*, slika 3.1. (postoje i alternative: *multi-master* model i *masterless* model, ali usled njihove složenosti oni se koriste samo u nekim specifičnim slučajevima). Po tom modelu jedna od replika izabrana

je da bude master replika, i svaka izmena podataka mora ići kroz nju, a onda ona izmenu propagira do ostalih zavisnih (*slave*) replika. Čitanje podataka se sa druge strane, može raditi sa bilo kog čvora.

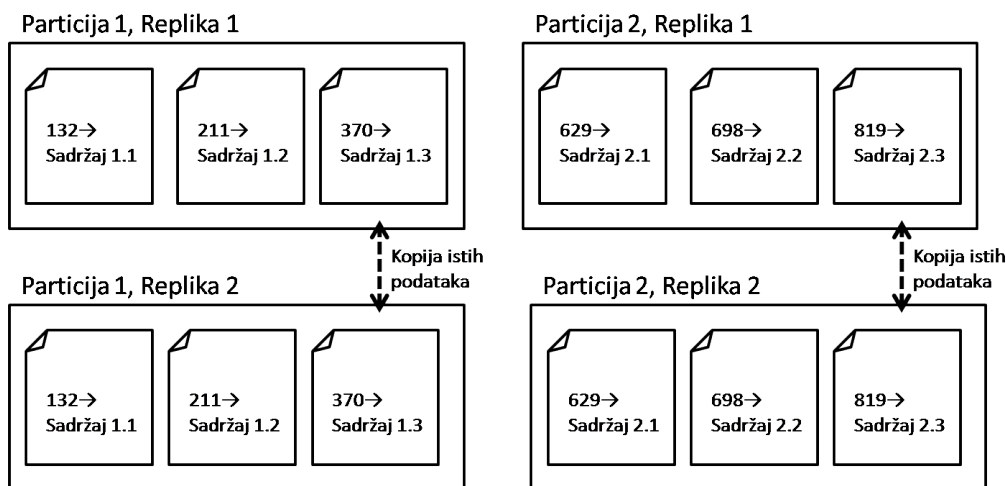
Replikacija može biti sinhrona i asinhrona. Kod sinhronne replikacije kada izmena dođe do master čvora, on pre nego što klijentu vrati odgovor da je izmena uspešno sačuvana, sačeka da mu svaki zavisni čvor potvrdi da je izmena uspešno upisana kod njega. Sinhrona replikacija je karakteristična za sisteme koji garantuju da će klijent pri čitanju podatka sa bilo koje replike imati ažurno stanje, međutim, cena toga je umanjena dostupnost, s obzirom da ukoliko neki od zavisnih čvorova postane nedostupan, on ne može master čvoru potvrditi da je izmena upisana na njega, pa samim tim ni master ne može klijentu potvrditi da je izmena uspešno upisana. Kod asinhronne replikacije, kada izmena podataka dođe do mastera, i on tu izmenu propagira ka ostalim zavisnim čvorovima, master ne čeka nikakav odgovor od zavisnih čvorova, već ukoliko je izmena uspešno upisana na njega on klijentu vraća potvrdu o sačuvanoj izmeni. Asinhrona replikacija obično nudi bolje performanse i veću dostupnost servisa, po cenu toga da čitanje sa neke od replika može vraćati zastarele podatke.

Kako relacije baze podataka daju prioritet konzistentosti (CP iz CAP teoreme) najčešće se kod njih koristi sinhrona replikacija, dok je kod modela kojima je dostupnost (AP iz CAP teoreme) prioritet, asinhrona replikacija model po kojem se podaci distribuiraju.



Slika 3.1: Master-slave replikacija

Replikacija i particionisanje obično idu zajedno, tako što se podaci prvo podele na particije, a svaka particija replicuje na više čvorova [7], slika 3.1.



Slika 3.2: Replikacija i particionisanje

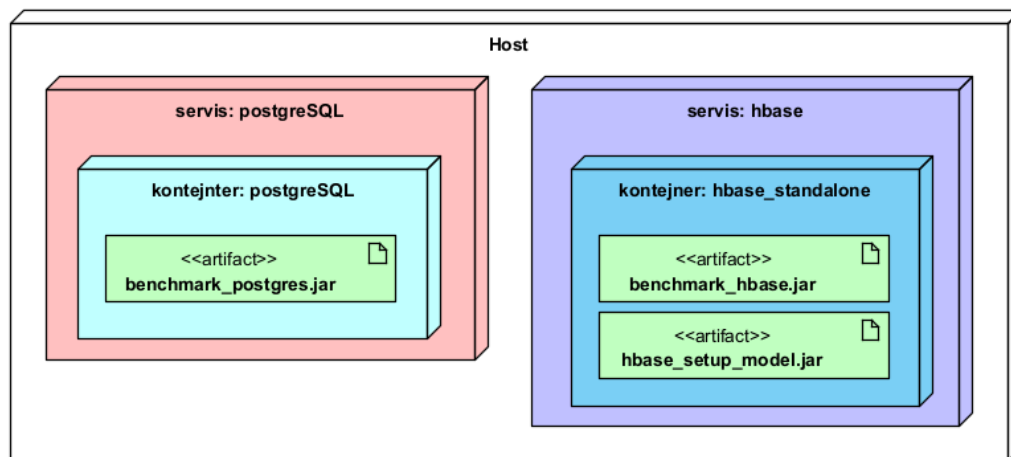
U svrhu demonstracije i merenja performansi HBase-a u distribuiranom okruženju, biće korišćeni testovi iz testiranja OLTP i OLAP scenarija. Usled složenosti implementacije distribuiranog relacionog sistema, za distribuirano okruženje primer je dat samo za HBase, a upoređivanje je dato na osnovu teorijskih istraživanja.

Glava 4

Merenje performansi po modelima

4.1 Opis platforme za testiranje

Za predstavnike baza podataka izabrani su HBase (kao predstavnik kolonski orijentisane nerelacione baze podataka) i PostgreSQL (kao predstavnik relacione baze podataka) ¹. Kao okruženje za izvršavanje testova korišćen je *docker*. Serveri oba predstavnika biće pokrenuti kao nezavisni kontejneri.



Slika 4.1: Platforma testiranja

¹Napomena: Kako su za predstavnike izabrani PostgreSQL i HBase, rezultati dobijeni u nastavku su rezultati poređenja konkretnih predstavnika, i nisu opšti za sve relacione i kolonski orijentisane nerelacione baze podataka.

4.2 Merenje performansi u OLTP okruženju

Model baze podataka za testiranje u OLTP okruženju koji se koristi u radu, sastoji se iz četiri tabele:

- **fxrates**: Sadrži informacije o kursu valutnih parova. Sadrži 16 redova.
- **fxuser**: Podaci o korisnicima. Sadrži 30 000 redova.
- **fxaccount**: Podaci o računima korisnika. Sadrži 120 000 redova.
- **fxtransaction**: Sadrži informacije o transferima, odnosno transakcijama sredstava.

setup-postgres-model.sql

```
1
2 create table postgresdb.fxrates (
3     currency_from varchar(50) not null,
4     currency_to varchar(50) not null,
5     rate real not null,
6     primary key(currency_from,currency_to)
7 );
8
9 create table postgresdb.fxuser (
10     id integer primary key,
11     username varchar (50) unique not null,
12     password varchar (50) not null,
13     start_balancecurrency varchar (10) not null,
14     start_balance real not null,
15     firstname varchar(100) not null,
16     lastname varchar(100) not null,
17     street varchar(100) not null,
18     city varchar(100) not null,
19     state varchar(100) not null,
20     zip varchar(50) not null,
21     phone varchar(30) not null,
22     mobile varchar(30) not null,
23     email varchar(50) unique not null,
24     created timestamp not null
25 );
26
27 create table postgresdb.fxaccount (
```

```
28         id integer primary key,
29         fxuser integer not null references postgresdb.fxuser(id),
30         currency_code varchar(10) not null,
31         balance real not null,
32         created timestamp not null,
33         unique (fxuser, currency_code)
34 );
35
36 create table postgresdb.fxtransaction (
37     id integer primary key,
38     fxaccount_from integer references postgresdb.fxaccount(id),
39     fxaccount_to integer references postgresdb.fxaccount(id),
40     amount numeric(15,2) not null,
41     status varchar(50) not null,
42     entry_date timestamp not null
43 );
```

setup-hbase-model.sh

```
1 create fxrates, 'data';
2 create fxuser, 'data';
3 create fxaccount, 'data';
4 create fxtransaction, 'data';
```

Kao što je ranije spomenuto, poslovna transakcija koja se testira jeste prebacivanje sredstava sa jednog računa na drugi. Ona se sastoji iz tri dela: kreiranje transfera sredstava (*createFxTransaction*), izvršavanja uplate(*executePayment*) i provera statusa transfera (*checkTransactionStatus*).

CreateFXTransaction prvo pročitava stanje naloga sa kojeg treba preneti sredstva, nakon toga čita kurs odgovarajućeg valutnog para i ukoliko ima dovoljno sredstava na računu, u tabelu sa transferima upisuje transfer u statusu *NEW*.

CreateFXTransaction

```
1
2 select fa.balance
3 from fxaccount fa
4 where fa.id = ?;
5
6 select fr.rate
7 from fxrates fr
8 where fr.currency_to = ? and fr.currency_from = ?;
9
10 insert into fxtransaction
11 (id,fxaccount_from, fxaccount_to, amount, status, entry_date)
12 values(?,?,?,?,?,?,?)
```

ExecutePayment prvo pročita stanja sa naloga koji učestvuju u transferu, menja im balans i nakon toga transferu menja status.

ExecutePayment

```
1
2 select balance
3 from fxaccount
4 where id = ?
5
6 update fxaccount
7 set balance = ?
8 where id = ?
9
10 select balance
11 from fxaccount
12 where id = ?
13
14 update fxaccount
15 set balance = ?
16 where id = ?
17
18 update fxtransaction
19 set status = ?
20 where id = ?
```

CheckTransactionStatus za odgovarajuću transakciju čita status po primarnom ključu.

CheckTransactionStatus

```
1 select status
2 from fxtransaction
3 where id = ?
```

Test ima podršku za paralelno izvršavanje poslovnih transakcija. Parametri testa su **broj klijenata** (*numOfClients*) i **broj poslovnih transakcija koje treba izvršiti** (*totalTransactions*). Oba parametra postavljaju se prilikom pokretanja testa, kroz standardni ulaz. Politika dodeljivanja poslovnih transakcija svakom od klijenata je da se ukupan broj poslovnih transakcija ravnomerno podeli svim klijentima, a eventualni ostatak pri podeli dodeljuje se nekom od njih.

Implementacija politike podele posla klijentima

```
1 List<Integer> transClientList = new ArrayList<>();
2 int transToAssign = totalTransactions;
3 int transPerClient = transToAssign / numOfClients;
4
5 for(int i = 0; i<numOfClients;i++){
6     transClientList.add(transPerClient);
7     transToAssign--transPerClient;
8 }
9 if (transToAssign > 0) {
10     int transForLast = transPerClientList.get(numOfClients - 1);
11     transClientList.set(numOfClients - 1, transForLast + transToAssign);
12 }
13
14 assert numOfClients==transClientList.size();
15 Thread[] threads = new Thread[numOfClients];
16 for(int i = 0;i<numOfClients;i++){
17     threads[i] = new Thread(
18         new BenchmarkSingleClientExecutor(
19             i*transClientList.get(i),transClientList.get(i)
20         )
21     );
22 }
```

Nakon što se klijentu dodeli skup poslovnih transakcija koje treba da obradi, on krene da izvršava sve tri faze svake poslovne transakcije koja mu je dodeljena. Svaki transfer koji učestvuje u poslovnoj transakciji generisan je na osnovu rednog broja poslovne transakcije koju klijent procesira. To je garancija da se ne može desiti da se isti transfer obrađuje više puta, kao i to da dva klijenta ne mogu

obrađivati isti transfer.

BenchmarkSingleClientExecutor.java

```
1
2 public class BenchmarkSingleClientExecutor implements Runnable {
3
4     private final CountDownLatch endSignal;
5     private final BenchmarkOLTPUtility oltpUtil;
6     private final int numOfT;
7     private final int startFrom;
8
9     private final Object connection;
10
11     @Override
12     public void run() {
13
14         try {
15             for (int i = this.start; i < this.start + this.numOfT; i++) {
16                 FXTransaction fxTransaction = DataGenerator.seedTransacton(i);
17                 ExecutePaymentInfo executePaymentInfo =
18                     oltpUtil.createFXTransaction(connection);
19                 oltpUtil.executePayment(,executePaymentInfo);
20                 oltpUtil.checkTransactionStatus(fxT);
21             }
22             endSignal.countDown();
23         } catch (Exception e) {
24             throw new IllegalStateException(e);
25         }
26     }
27 }
```

Priprema okruženja za testiranje podrazumeva kompilaciju java testova, pokretanje docker kontejnera i prebacivanje kompiliranih testova na odgovarajuće kontejnere. Dodatan korak za HBase jeste da se na HBase kontejneru kreira struktura baze podataka. Skripta u nastavku sadrži sve neophodne komande za pokretanje okruženja.

prepare-env.sh

```
1 #!/bin/bash
2 echo 'PREPARING ENVIRONMENT...';
3
4 export JAVA_HOME="$JAVA_8";
5 mvn -f ./hbase_setup_model clean compile assembly:single;
6 mvn -f ./benchmark_hbase clean compile assembly:single;
7 export JAVA_HOME="$JAVA_17";
8 mvn -f ./benchmark_postgres clean compile assembly:single;
9
10 docker-compose -f docker-compose.yml up --build -d;
11 docker exec -it hbase-master-1 sh -c "java -jar hbase_setup_model.jar";
```

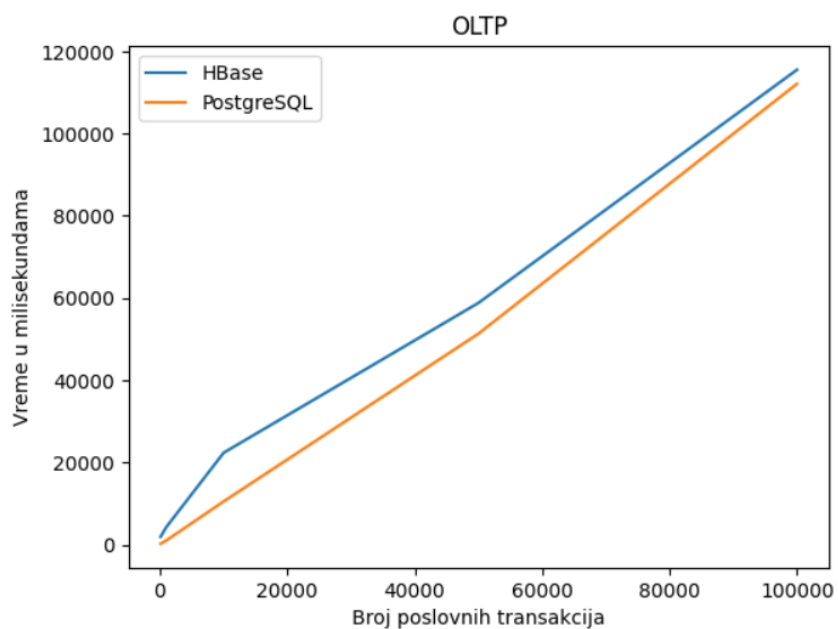
docker-compose.yml

```
1 services:
2   postgres:
3     container_name: postgres
4     ports:
5       - "5433:5432"
6     volumes:
7       - ./setup_model.sql:/docker-entrypoint-initdb.d/create_script.sql
8       - ./benchmark_postgres.jar:/benchmark_postgres.jar
9     environment:
10       - POSTGRES_PASSWORD=postgres
11       - POSTGRES_USER=postgres
12       - POSTGRES_DB=postgresdb
13     build:
14       context: .
15       dockerfile: ./Dockerfile_postgres
16
17   hbase:
18     image: bde2020/hbase-standalone:1.0.0-hbase1.2.6
19     container_name: hbase
20     volumes:
21       - hbase_data:/hbase-data
22       - hbase_zookeeper_data:/zookeeper-data
23       - ./hbase_setup_model.jar:/hbase_setup_model.jar
24       - ./benchmark_hbase.jar:/benchmark_hbase.jar
25     ports:
26       - 16000:16000
27       - 16010:16010
28       - 16020:16020
```

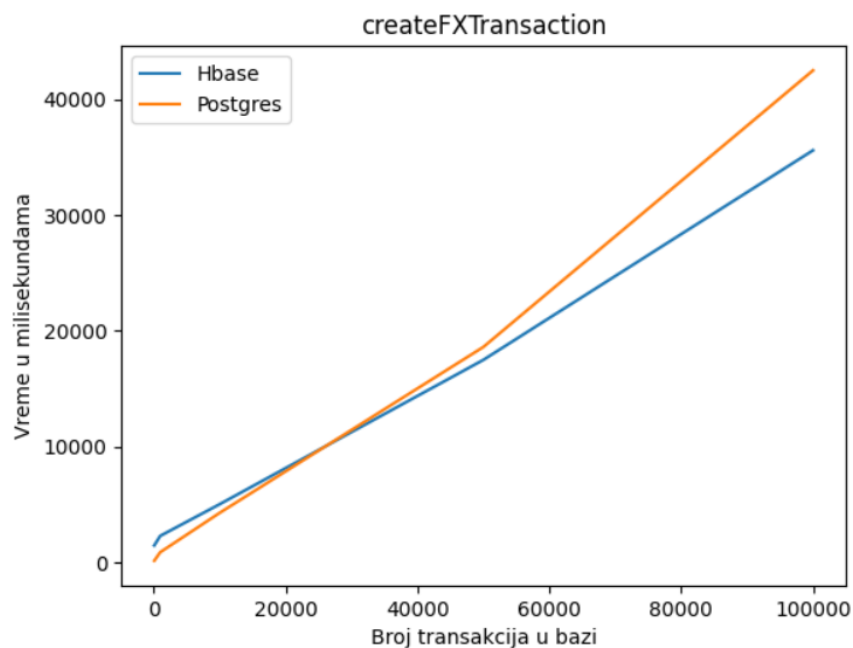


```
29      - 16030:16030
30      - 2888:2888
31      - 3888:3888
32      - 2181:2181
33
34 volumes:
35   hbase_data:
36   hbase_zookeeper_data:
```

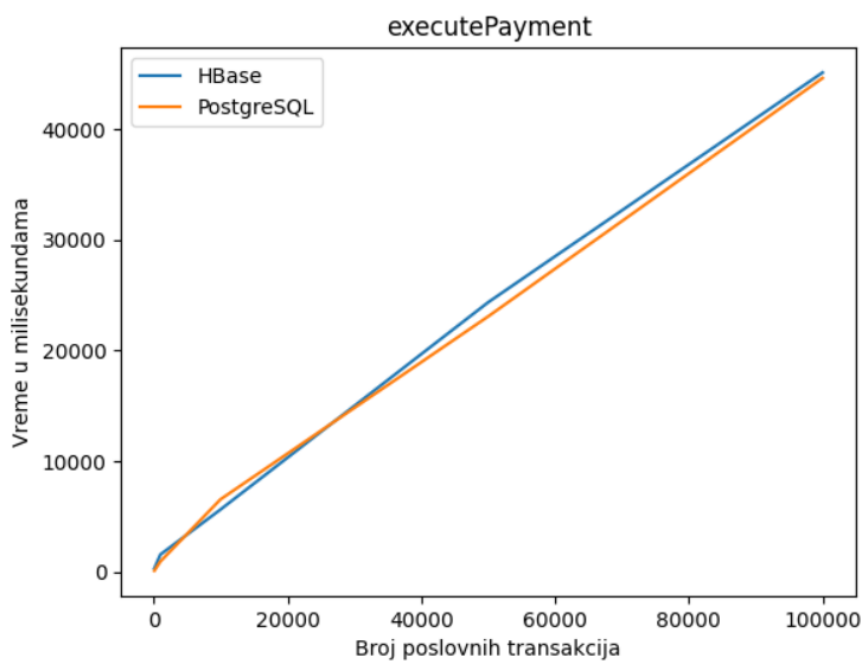
Rezultati merenja prikazani u nastavku nastali kao rezultat pokretanja testova sa 100, 1 000, 10 000, 50 000, 100 000 poslovnih transakcija i pet klijenata koji te poslovne transakcije paralelno obrađuju.



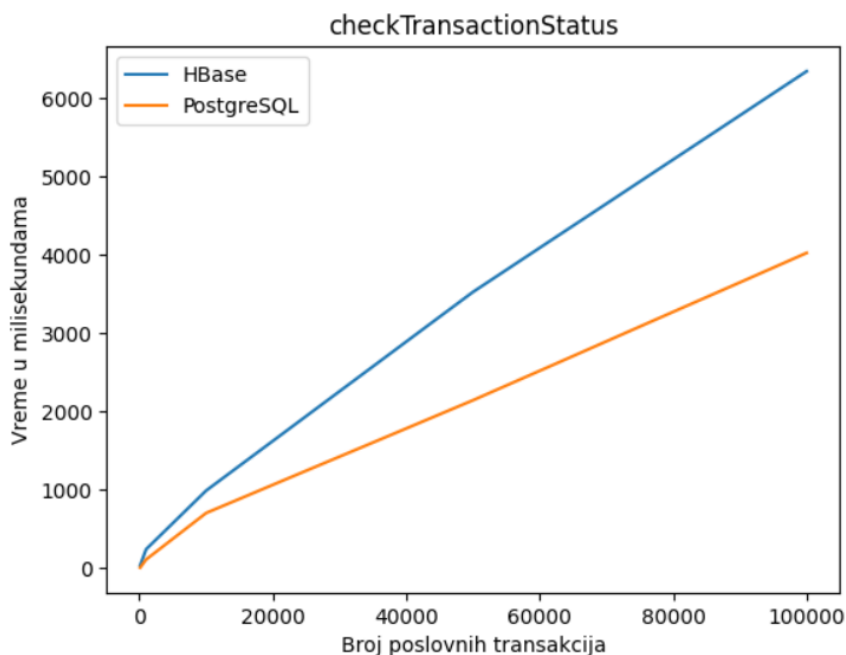
Slika 4.2: Rezultati merenja u OLTP okruženju



Slika 4.3: Rezultati merenja createFXTransaction dela poslovne transakcije



Slika 4.4: Rezultati merenja executePayment dela poslovne transakcije



Slika 4.5: Rezultati merenja checkTransactionStatus dela poslovne transakcije

4.3 Merenje performansi u OLAP okruženju

Model baze podataka koji se koristi za testiranje u OLAP okruženju sastoji se iz sledećih tabela:

- **product:** Sadrži informacije o proizvodima. Sadrži 3 000 000 redova.
- **supplier:** Podaci o dovaljačima. Sadrži 1 000 000 redova.
- **productsupplier:** Vezna tabela između dobavljača i proizvoda. Sadrži 5 000 000 redova.
- **customer:** Informacije o mušterijama. Sadrži 1 500 000 redova.
- **order:** Informacije o narudžbinama. Sadrži 1 500 000 redova.
- **orderitem:** Informacije o pojedinim stavkama narudžbine. Sadrži 6 000 000 redova.

setup-postgres-model.sql

```
1 create table product (
2     id integer not null primary key,
```

```
3     name varchar(50) not null,
4     brand varchar(50) not null,
5     type varchar(50) not null,
6     size integer not null,
7     container varchar(50) not null,
8     price varchar(50) not null,
9     comment varchar(50)
10 );
11
12 create table supplier (
13     id integer primary key,
14     name varchar(50) not null,
15     address varchar(200) not null,
16     phone varchar(50) not null
17 );
18
19 create table productsupplier (
20     id integer not null primary key,
21     product integer not null,
22     supplier integer not null,
23     available integer not null,
24     supply_cost real not null,
25     comment varchar(200),
26     constraint fk_product
27     foreign key(product) references postgresdb.product(id),
28     constraint fk_supplier
29     foreign key(supplier) references postgresdb.supplier(id)
30
31 );
32
33 create table customer(
34     id integer primary key,
35     name varchar(50) not null,
36     address varchar(200) not null,
37     phone varchar(50) not null,
38     comment varchar(200)
39 );
40
41 create table order(
42     id integer primary key,
43     customer integer not null,
44     status varchar(20) not null,
```

```

45     total_price real not null,
46     entry_date date not null,
47     priority varchar(20) not null,
48     comment varchar(200),
49     constraint fk_customer
50     foreign key(customer) references postgresdb.customer(id)
51 );
52
53 create table order_item(
54     order_id integer not null,
55     product integer not null,
56     supplier integer not null,
57     order_no integer not null,
58     quantity integer not null,
59     base_price real not null,
60     discount real not null,
61     tax real not null,
62     status varchar(20) not null,
63     ship_date date not null,
64     commit_date date not null,
65     comment varchar(200),
66     primary key(order_id,product,supplier),
67     constraint fk_order
68     foreign key(order_id) references postgresdb.order(id),
69     constraint fk_product
70     foreign key(product) references postgresdb.product(id),
71     constraint fk_supplier
72     foreign key(supplier) references postgresdb.supplier(id)
73 );

```

hbase-setup-model

```

1
2 create product, 'data';
3 create supplier, 'data';
4 create productsupplier, 'data';
5 create customer, 'data';
6 create order, 'data';
7 create orderitem, 'data';

```

OLAP test će obuhvatiti upunjavanje tabela iz csv fajlova (*bulk load*), kao i iz izvršavanja upita koji će klijenti paralelno izvršavati više puta. Upit će biti parametrizovan i uzimaće dva parametra. Prvi parametar je dan narudžbine, a

drugi je njen status. Na osnovu tih parametara dohvataju se agregirane vrednosti kolona iz tabele *orderitem*, grupisane po statusu.

bulkLoad

```
1 copy product from "./product.csv";
2 copy supplier from "./supplier.csv";
3 copy productsupplier from "./productsupplier.csv";
4 copy customer from "./customer.csv";
5 copy order from "./order.csv";
6 copy order_item from "./order_item.csv";
```

executeOLAPQuery

```
1
2 select
3     oi.status status,
4     sum(oi.quantity) as sum_qty,
5     sum(oi.base_price) as sum_base_price,
6     sum(oi.base_price*(1-oi.discount)) as sum_disc_price,
7     sum(oi.base_price*(1-oi.discount)*(1+oi.tax)) as sum_charge,
8     avg(oi.quantity) as avg_qty,
9     avg(oi.base_price) as avg_price,
10    avg(oi.discount) as avg_disc,
11    count(*) as count_order
12 from
13     postgresdb.order_item oi
14 where
15     oi.ship_date = to_date(?, 'dd.mm.yyyy') and
16     oi.status = ?
17 group by oi.status;
```

Parametri testa su **broj klijenata** (*numOfClients*) koji će paralelno izvršavati olap upit i **ukupan broj izvršavanja upita** (*totalIterations*). Oba parametra postavljaju se prilikom pokretanja testa, kroz standardni ulaz. Politika dodeljivanja broja izvršavanja upita svakom od klijenata ista je kao i dodeljivanja broja poslovnih transakcija klijentima, kod testiranja OLTP okruženja.

Implementacija politike podele posla klijentima

```
1 List<Integer> iterClientList = new ArrayList<>();
2 int itersToAssign = totalIterations;
3 int itersPerClient = itersToAssign / numOfClients;
4
```

```

5  for(int i = 0; i<numOfClients;i++){
6      iterClientList.add(itersPerClient);
7      itersToAssign-=itersPerClient;
8  }
9  if (itersToAssign > 0) {
10     int itersForLast = iterClientList.get(numOfClients - 1);
11     iterClientList.set(numOfClients - 1, itersForLast + itersToAssign);
12 }
13
14 assert numOfClients==iterClientList.size();
15 Thread[] threads = new Thread[numOfClients];
16 for(int i = 0;i<numOfClients;i++){
17     threads[i] = new Thread(
18         new BenchmarkSingleClientExecutor(
19             i*iterClientList.get(i),iterClientList.get(i)
20         )
21     );
22 }

```

Svaki klijent nakon što mu je dodeljen broj iteracija, kreće da izvršava upit onoliko puta koliko mu je iteracija dodeljeno. Svako izvršavanje OLAP upita ima iste parametre.

BenchmarkSingleClientExecutor.java

```

1  public class BenchmarkSingleClientExecutor implements Runnable {
2
3  private final CountDownLatch endSignal;
4  private final BenchmarkOLAPUtility olapUtil;
5  private final int numOfIters;
6  private final int startFrom;
7
8  private final Object connection;
9  @Override
10 public void run() {
11
12     try {
13         for (int i = this.start; i < this.start + this.numOfIters; i++) {
14             olapUtil.executeOLAPQuery(connection);
15         }
16         endSignal.countDown();
17     } catch (Throwable e) {
18         throw new IllegalStateException(e);

```

```
19     }
20   }
21 }
```

Priprema okruženja za testiranje podrazumeva generisanje csv fajlova koji kasnije treba da budu učitani u tabele koristeći podršku za *bulk load*, zatim kompilaciju java testova, pokretanje docker kontejnera i prebacivanje kompiliranih testova kao i izgenerisanih csv fajlova na odgovarajuće kontejnere. Dodatan korak za HBase jeste da se na HBase kontejneru kreira struktura baze podataka. Skripta u nastavku sadrži sve neophodne komande za pokretanje okruženja.

prepareEnv.sh

```
1  #!/bin/bash
2  echo 'PREPARING ENVIRONMENT...';
3  echo 'PREPARING HBASE BENCHMARK JARS...';
4  export JAVA_HOME="$JAVA_8";
5  mvn -f olap_benchmark_hbase clean compile assembly:single;
6  mvn -f hbase_setup_olap_model clean compile assembly:single;
7  mvn -f hbase_bulk_load_setup clean compile assembly:single;
8
9  echo 'PREPARING HBASE BULK LOAD RESOURCES..';
10 java -jar ./hbase_bulk_load_setup.jar;
11
12
13 echo 'PREPARING POSTGRES BENCHMARK JARS...';
14 export JAVA_HOME="$JAVA_17";
15 mvn -f olap_benchmark_postgres clean compile assembly:single;
16 mvn -f postgres_bulk_load_setup clean compile assembly:single;
17
18 echo 'PREPARING POSTGRES BULK LOAD RESOURCES..';
19 java -jar postgres_bulk_load_setup.jar;
20
21 docker-compose -f docker-compose.yml up --build -d;
22 winpty docker exec -it hbase sh -c "java -jar setup_olap_model.jar";
```

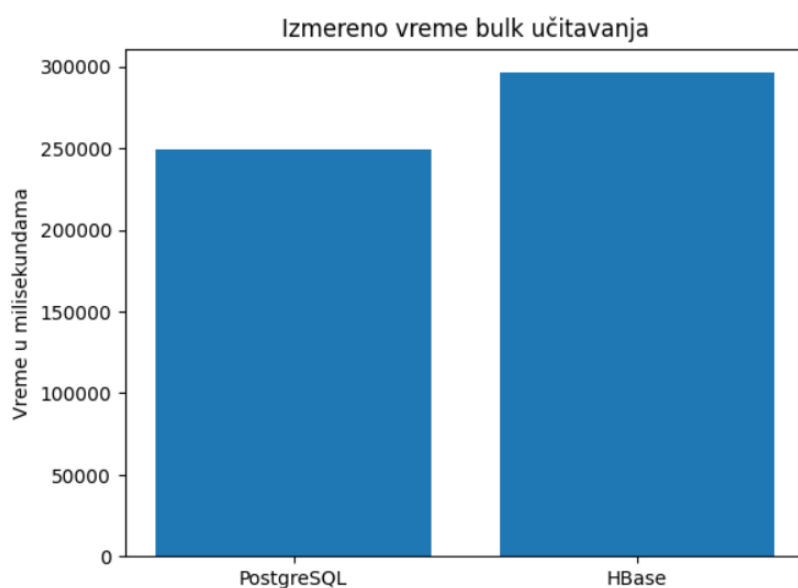
docker-compose.yml

```
1
2 services:
3   postgres:
4     container_name: postgres
5     ports:
```

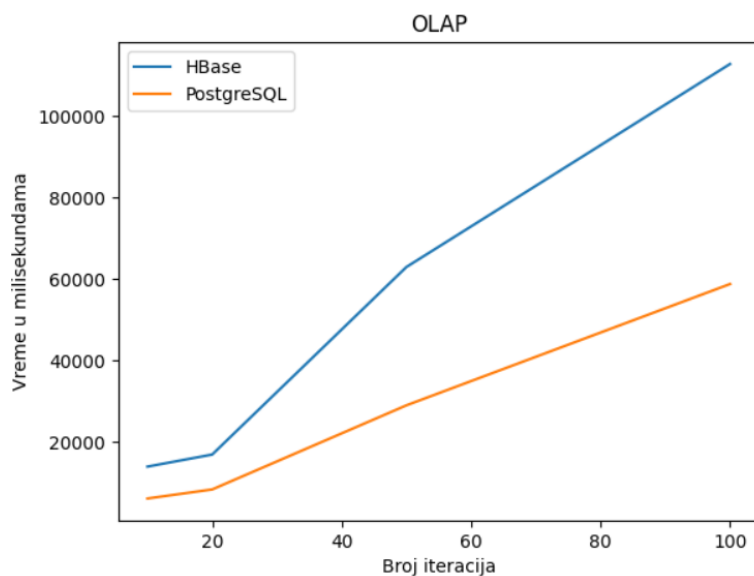


```
6     - "5433:5432"
7 volumes:
8     - ./setup_model.sql:/docker-entrypoint-initdb.d/create_script.sql
9     - ./benchmark_postgres.jar:/benchmark_postgres.jar
10 environment:
11     - POSTGRES_PASSWORD=postgres
12     - POSTGRES_USER=postgres
13     - POSTGRES_DB=postgresdb
14 build:
15     context: .
16     dockerfile: ./Dockerfile_postgres
17
18 hbase:
19     image: bde2020/hbase-standalone:1.0.0-hbase1.2.6
20     container_name: hbase
21     volumes:
22     - ./productsupplierHB.csv:/productsupplier.csv
23     - ./productHB.csv:/product.csv
24     - ./supplierHB.csv:/supplier.csv
25     - ./customerHB.csv:/customer.csv
26     - ./orderHB.csv:/order.csv
27     - ./orderitemHB.csv:/orderitem.csv
28     - ./orderitemStatsHB.csv:/orderitemstats.csv
29     - hbase_data:/hbase-data
30     - hbase_zookeeper_data:/zookeeper-data
31     - ./hbase_setup_model.jar:/hbase_setup_model.jar
32     - ./benchmark_hbase.jar:/benchmark_hbase.jar
33     ports:
34     - 16000:16000
35     - 16010:16010
36     - 16020:16020
37     - 16030:16030
38     - 2888:2888
39     - 3888:3888
40     - 2181:2181
41
42 volumes:
43     hbase_data:
44     hbase_zookeeper_data:
```

Rezultati merenja prikazani u nastavku nastali kao rezultat pokretanja testova sa 10, 20, 50, 100 iteracija izvršavanja OLAP upita i pet klijenata koji te upite paralelno izvršavaju. Na slici se može primetiti da u legendi uz HBase stoji reč naivni (*naive*). To se odnosi na model koji se koristi, vse reči o razlogu zašto se ovako napravljen model za HBase u ovom okruženju može smatrati naivnim, biće u sledećem poglavlju.



Slika 4.6: Bulk load podataka



Slika 4.7: OLAP

4.4 Merenje performansi u distribuiranom okruženju

Kao što je napomenuto u 3.4 za merenje performansi biće iskorišćen OLTP test iz poglavlja 3.1. Jedina razlika će biti priprema okruženja za testiranje. U ovom slučaju koristićemo lokalnu mašinu sa ranije instaliranim HBase-om i Hadoop-om.

Priprema okruženja podrazumeva najpre pokretanje hadoop servisa sa postavljenim faktorom replikacije blokova, kao i konfigurisanom čvorom podataka:

hdfs-site.xml

```
1 <configuration>
2   <property>
3     <name>dfs.replication</name>
4     <value>2</value>
5   </property>
6   <property>
7     <name>dfs.datanode.data.dir</name>
8     <value>/home/lukadj/hadoop/data2</value>
9   </property>
10 </configuration>
```

start-hadoop.sh

```
1 #!/bin/bash
2 $HADOOP_HOME/sbin/start-dfs.sh
3 $HADOOP_HOME/sbin/start-yarn.sh
```

Nakon toga pokrećemo HBase klaster sa sledećom konfiguracijom:

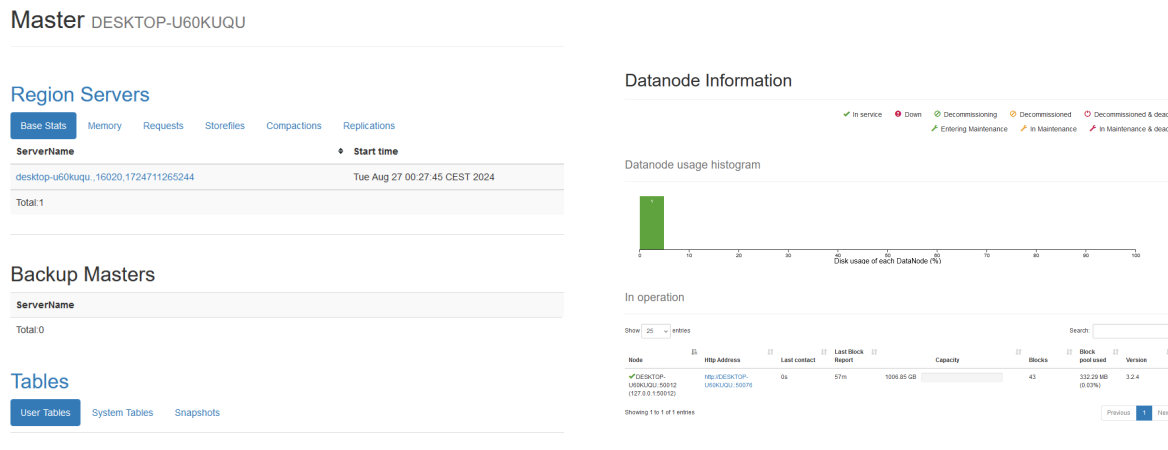
hbase-site.xml

```
1 <property>
2   <name>hbase.cluster.distributed</name>
3   <value>true</value>
4 </property>
5 <property>
6   <name>hbase.rootdir</name>
7   <value>hdfs://localhost:9000/hbase</value>
8 </property>
9 <property>
10  <name>hbase.zookeeper.property.clientPort</name>
11  <value>10231</value>
12 </property>
```

GLAVA 4. MERENJE PERFORMANSI PO MODELIMA

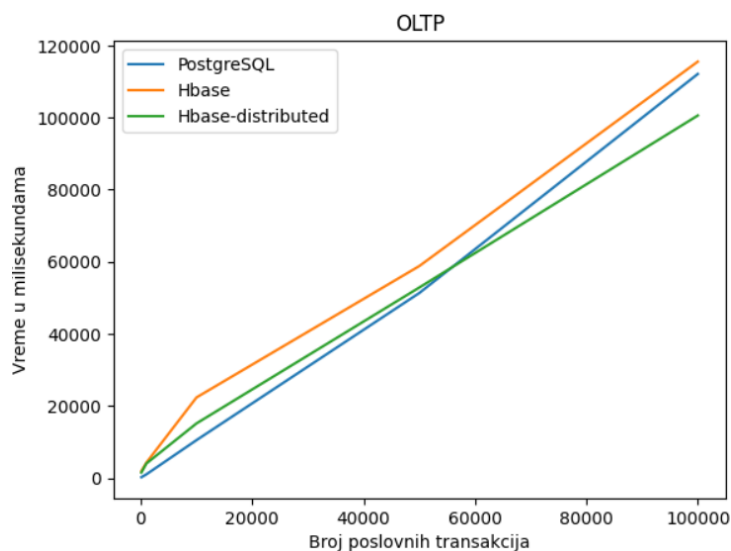
start-hbase.sh

```
1 #!/bin/bash
2 $HBASE_HOME/bin/start-hbase.sh
```



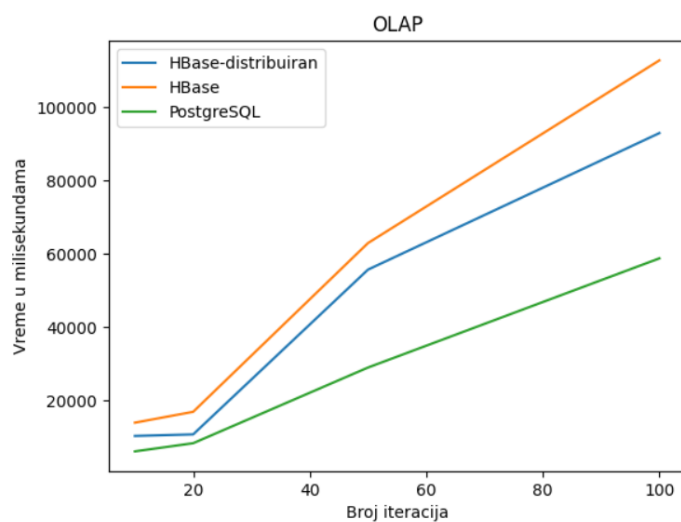
Slika 4.8: Stranice za praćenje hbase master servera (levo) i HDFS čvorova podataka(desno)

Poređenje rezultata merenja distribuiranog HBase klastera sa samostalim HBase-om i PostgreSQL-om prikazano u nastavku nastalo je kao rezultat pokretanja testova sa 100, 1 000, 10 000, 50 000, 100 000 iteracija izvršavanja OLTP.



Slika 4.9: Upoređivanje rezultata dobijenih u distribuiranom okruženju i rezultata iz testiranja OLTP okruženja

Poređenje rezultata merenja distribuiranog HBase klastera sa samostalim HBase-om i PostgreSQL-om prikazano u nastavku nastalo je kao rezultat pokretanja testova sa 10, 20, 50, 100 iteracija izvršavanja OLAP upita



Slika 4.10: Upoređivanje rezultata dobijenih u distribuiranom okruženju i rezultata iz testiranja OLAP okruženja

Glava 5

Analiza rezultata

Analiza rezultata obuhvata interpretaciju rezultata koji su dobijeni merenjima performansi, kao i navođenje eventualnih unapređenja i primedbi koje treba imati u vidu kada se radi sa datim tehnologijama.

5.1 Analiza rezultata kod testiranja u OLTP okruženju

Rezultati merenja ukazuju da kako raste broj poslovnih transakcija koje treba obraditi u testu, inicijalna prednost PostgreSQL-a u odnosu na HBase opada. Konkretno po fazama poslovne transakcije, izdvaja se `createFXTransaction` transakcija baze podataka, kod koje na postavljenih 50000 i 100000 poslovnih transakcija, HBase ima prednost. Sa druge strane, na primeru `checkTransactionStatus`-a vidi se da pretraga po primarnom ključu kod PostgreSQL-a radi brže nego što je to slučaj kod HBase-a. I pored toga ukupni rezultati ukazuju na tendenciju da bi sa porastom broja poslovnih transakcija (na milion ili deset miliona) PostgreSQL imao veća usporenja, relativno u odnosu na HBase, međutim za testiranje takvog okruženja neophodno je koristiti računar koji je sposoban da sprovede tako zahtevan test.

Važno je naglasiti da HBase ne garantuje konzistentost nad svim podacima kakvu nudi PostgreSQL. HBase nudi jedan vid konzistentosti, i to konzistentnost u radu sa jednom redom tabele (više o tome u 2.2). To dovodi do zaključka da ukoliko je neophodno implementirati transakciju koja uključuje rad sa podacima više tabela ili više redova jedne tabele, ukoliko koristimo HBase, moramo na apli-

kativnom sloju voditi računa o očuvanju eventualne konzistentnosti. PostgreSQL sa druge strane kao ACID baza podataka, sama garantuje konzistentost, te ne zahteva dodatan napor kao u slučaju HBase-a.

5.2 Analiza rezultata kod testiranja u OLAP okruženju

Rezultati merenja u okviru OLAP rezultata jasno ističu prednost PostgreSQL-a u performansama. Prednost HBase-a u odnosu na PostgreSQL jeste fleksibilnost sheme koja u ovakvom slučaju može doći do izražaja. Vid jednostavne denormalizacije koja se u slučaju HBase-a može primeniti može doneti dramatična poboljšanja u performansama čitanja HBase-a čak i u odnosu na Postgres.

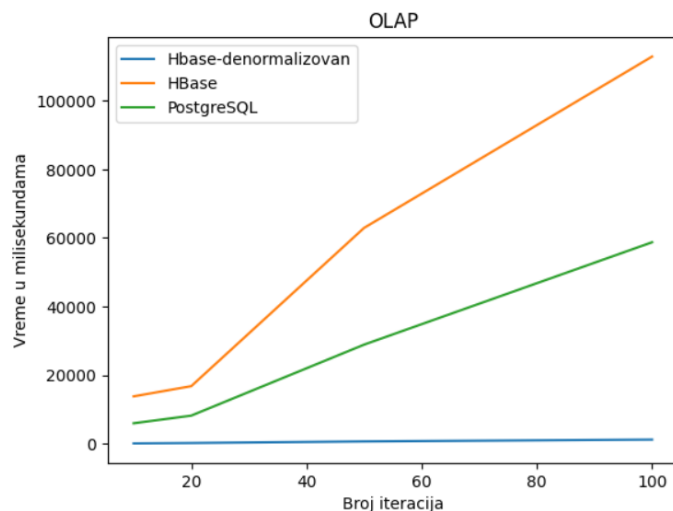
Kreiranje tabele (u ovom slučaju `orderitemstats`) koja za vrednost ključa ima status, a za kolone ima vrednosti svake potrebne agregirane funkcije po danima, uticalo bi na to da se čitanje izveštaja svodi na čitanje po ključu sa izdvajanjem potrebnih kolona¹. Primer dodate tabele može se videti na slici 5.1, a uticaj na performanse na slici 5.2.

status	sum_qty01072021	sum_base01072021	sum_disc01072021	sum_chg01072021	...	sum_qty12092024	sum_base12092024	avg_qty12092024	...
Status_01	1818.00	234 921.00	32 921.00	32 921.00		219.00	89 321.00	32.50	
Status_02	1829.00	242 021.00	34 331.00	32 987.00		230.00	75 012.00	65.20	
...									

Slika 5.1: Tabela `orderitemstats`

Sa druge strane PostgreSQL takođe podržava mehanizme koji u ovom slučaju mogu poboljšati performanse. Konkretno za ovaj test, možemo kreirati materijalizovani pogled `mv_order_item_bydate` sa rezultatima ranije izračunatih statističkih funkcija. Tada bi se OLAP upit iz testa pretvorio u čitanje podataka iz materijalizovanog pogleda umesto izračunavanja agregiranih funkcija. Takva optimizacija ima veliki uticaj na performanse, slika 5.3.

¹Razlog zašto se ovakva modifikacija ne može primeniti u slučaju PostgreSQL-a jeste ograničen broj kolona, kao i to što skup kolona mora biti unapred definisan pri kreiranju modela. Kod HBase-a takvi limiti ne postoje.



Slika 5.2: Uporedna analiza performansi postgressa i hbase-a

mvorder_item_bydate

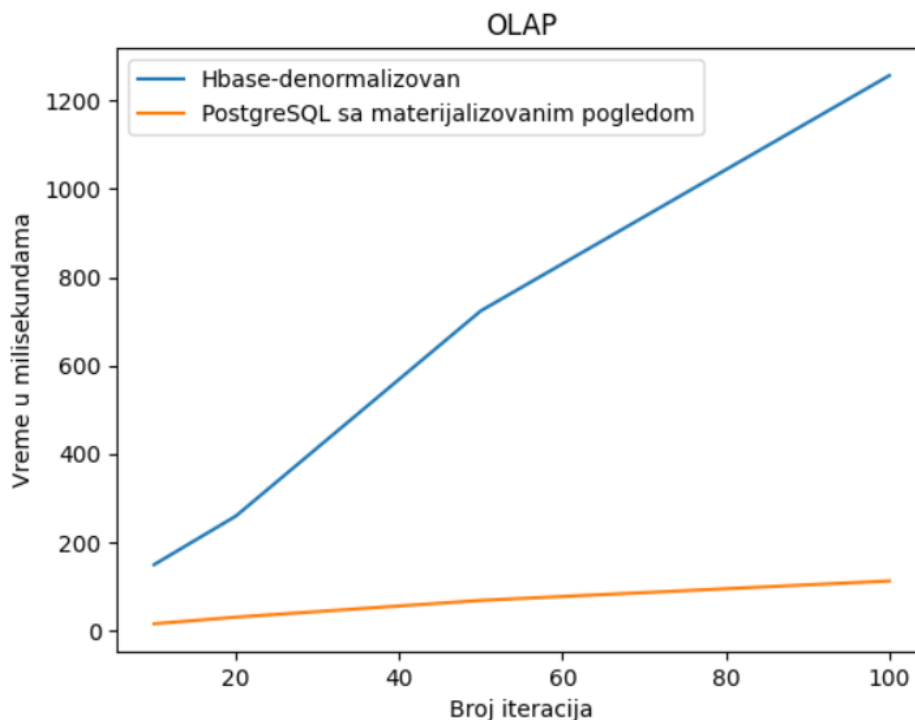
```

1 create materialized view postgresdb.mvorder_item_bydate
2   as (select
3       oi.ship_date,
4       sum( oi.quantity ) as sum_qty,
5       sum(OI.BASE_PRICE) as sum_base_price,
6       sum(OI.BASE_PRICE*(1-OI.discount)) as sum_disc_price,
7       sum(OI.BASE_PRICE*(1-OI.discount)*(1+OI.TAX)) as sum_charge,
8       avg(OI.QUANTITY) as avg_qty,
9       avg(OI.BASE_PRICE) as avg_price,
10      avg(OI.DISCOUNT) as avg_disc,
11      count(*) as count_order
12      from postgresdb.order_item oi
13      group by oi.status, oi.ship_date)
14      with data;
```

executeOLAPQuery

```

1 select *
2     from postgresdb.mv_order_item_bydate
3     where ship_date = TO_DATE(?, 'DD.MM.YYYY')
4     and status = ?
```

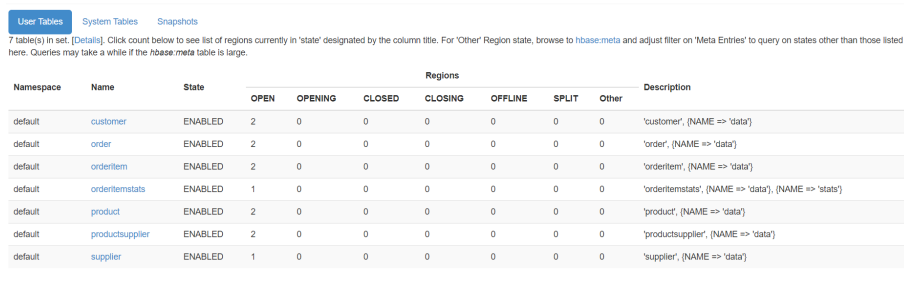



Slika 5.3: Upporedna analiza performansi koristeći optimizaciju kod HBase-a i PostgreSQL-a

5.3 Analiza rezultata kod testiranja u distribuiranom okruženju

Pokrenuti HBase klaster sa jednim HDFS čvorom podataka pokazuje nešto bolje rezultate nego samostalni HBase koji podatke skladišti na lokalnom fajl sistemu. Svakako, različito okruženje u kojem su testovi pokrenuti može uticati na rezultate, pa ih treba uzeti sa rezervom. O okviru analize veća pažnja biće posvećena načinu na koji su organizovane komponente u okviru distribuiranog HBase klastera. Na slici 5.4. prikazana je organizacija tabela po regionima u okviru klastera. Broj regiona srazmeran je broju redova (a ne količini podataka) koje tabela sadrži.

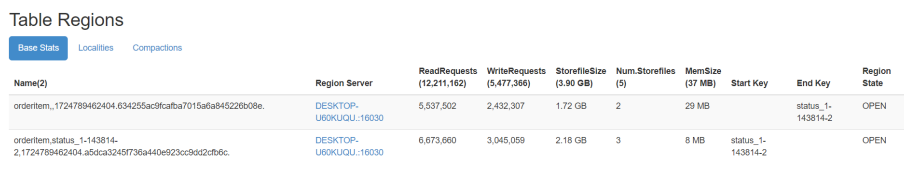
GLAVA 5. ANALIZA REZULTATA



Namespace	Name	State	Regions							Description
			OPEN	OPENING	CLOSED	CLOSING	OFFLINE	SPLIT	Other	
default	customer	ENABLED	2	0	0	0	0	0	0	'customer', (NAME => 'data')
default	order	ENABLED	2	0	0	0	0	0	0	'order', (NAME => 'data')
default	orderitem	ENABLED	2	0	0	0	0	0	0	'orderitem', (NAME => 'data')
default	orderitemstats	ENABLED	1	0	0	0	0	0	0	'orderitemstats', (NAME => 'data'), (NAME => 'stats')
default	product	ENABLED	2	0	0	0	0	0	0	'product', (NAME => 'data')
default	productsupplier	ENABLED	2	0	0	0	0	0	0	'productsupplier', (NAME => 'data')
default	supplier	ENABLED	1	0	0	0	0	0	0	'supplier', (NAME => 'data')

Slika 5.4: Organizcija tabela po regionima koje obuhvataju

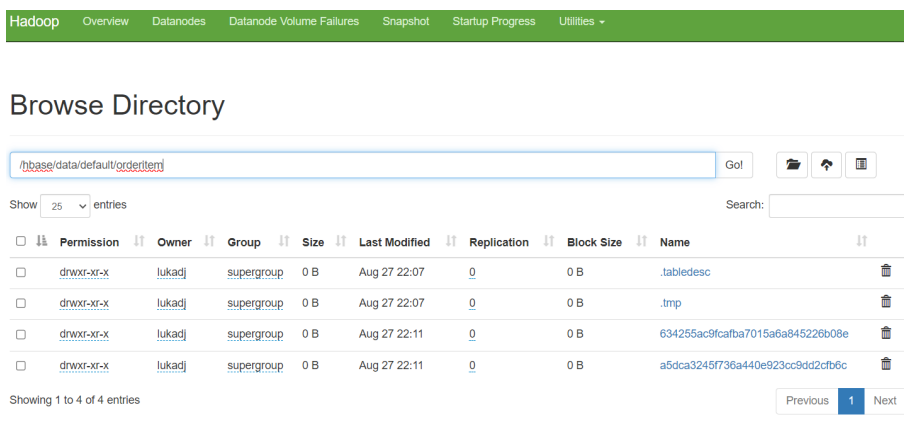
Listu i opis regiona možemo videti na primeru regiona tabele *orderitem* na slici 5.5.



Name(2)		Region Server	ReadRequests (12,211,162)	WriteRequests (5,477,366)	StorefileSize (3.90 GB)	Num.Storefiles (5)	MemSize (37 MB)	Start Key	End Key	Region State
orderitem,,1724789462404.634255ac9fca7015a6a845226b08e.		DESKTOP- U69KUQU-16030	5,537,502	2,432,307	1.72 GB	2	29 MB		status_1- 143814-2	OPEN
orderitem.status_1-143814- 2,1724789462404.a5dca3245f736a440e923cc9dd2cfb6c.		DESKTOP- U69KUQU-16030	6,673,660	3,045,059	2.18 GB	3	8 MB	status_1- 143814-2		OPEN

Slika 5.5: Regioni tabele orderitem

Na HDFS čvoru podaci su organizovani tako da svaka tabela ima svoj direktorijum. U okviru direktorijuma tabele nalaze se njeni regioni. Svaki od tih regiona podeljen je na blokove fiksne veličine koji bivaju replikovani unutar klastera.



Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
drwxr-xr-x	lukadj	supergroup	0 B	Aug 27 22:07	0	0 B	.tabledesc
drwxr-xr-x	lukadj	supergroup	0 B	Aug 27 22:07	0	0 B	.tmp
drwxr-xr-x	lukadj	supergroup	0 B	Aug 27 22:11	0	0 B	634255ac9fca7015a6a845226b08e
drwxr-xr-x	lukadj	supergroup	0 B	Aug 27 22:11	0	0 B	a5dca3245f736a440e923cc9dd2cfb6c

Slika 5.6: Orderitem direktorijum na HDFS čvoru sadrži regione u vidu poddirektorijuma

Hadoop

Overview

Datanodes

Datanode Volume Failures

Snapshot

Startup Progress

Utilities

Browse Directory

/hbase/data/default/orderitem/a5dca3245f736a440e923cc9dd2cfb6c/data

Go!

Show

25

entries

Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
<input type="checkbox"/>	-rW-r--r--	lukadj	supergroup	210.55 MB	Aug 27 22:13	2	128 MB	a3465fb1770d45e0ae25434951f9c78a	
<input type="checkbox"/>	-rW-r--r--	lukadj	supergroup	730.96 MB	Aug 27 22:13	2	128 MB	c9448e6d30f5432b94359edbb8ab69f4	
<input type="checkbox"/>	-rW-r--r--	lukadj	supergroup	1.26 GB	Aug 27 22:12	2	128 MB	fa3da9f4e682439bbd40f786eb8c6686	

Showing 1 to 3 of 3 entries

Previous

1

Next

Slika 5.7: Region je interno podeljen na blokove

46

Glava 6

Zaključak

Bibliografija

- [1] Apache ZooKeeper. on-line at: <https://zookeeper.apache.org/>.
- [2] Hdfs architecture. on-line at: <https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-hdfs/HdfsDesign.html#Introduction>.
- [3] Tpc-c. on-line at: <https://www.tpc.org/tpcc/>.
- [4] Tpc-h. on-line at: <https://www.tpc.org/tpch/>.
- [5] What is oltp? on-line at: <https://www.oracle.com/database/what-is-oltp/>.
- [6] A brief history of postgresql. 2024. on-line at: <https://www.postgresql.org/docs/current/history.html>.
- [7] Designing data-intensive applications. the big ideas behind reliable, scalable and maintainable systems. 2024.
- [8] John D Book. Martin Fowler. on-line at: <https://martinfowler.com/articles/patterns-of-distributed-systems/write-ahead-log.html>.
- [9] Sanjay Ghemawat Fay Chang, Jeffrey Dean. Bigtable: A Distributed Storage System for Structured Data. *SIAM Journal on Computing*, 16:486–502, 2006. on-line at: <https://static.googleusercontent.com/media/research.google.com/fr//archive/bigtable-osdi06.pdf>.
- [10] Free Software Foundation. ApacheHBase, 2013. on-line at: <https://hbase.apache.org/acid- semantics.html>.
- [11] Amandeep Khurana. Introduction to hbase Schema Design. on-line at: http://0b4af6cdc2f0c5998459-c0245c5c937c5dedcca3f1764ecc9b2f.r43.cf2.rackcdn.com/9353-login1210_khurana.pdf.

BIBLIOGRAFIJA

- [12] Regina O. Obe and Leo S. Hsu. PostgreSQL: Up and Running.
- [13] Gordana Pavlović-Lažetić. Uvod u relacione baze podataka. 1999.

Biografija autora

Vuk Stefanović Karadžić (*Tršić, 26. oktobar/6. novembar 1787. — Beč, 7. februar 1864.*) bio je srpski filolog, reformator srpskog jezika, sakupljač narodnih umotvorina i pisac prvog rečnika srpskog jezika. Vuk je najznačajnija ličnost srpske književnosti prve polovine XIX veka. Stekao je i nekoliko počasnih doktorata. Učestvovao je u Prvom srpskom ustanku kao pisar i činovnik u Negotinskoj krajini, a nakon sloma ustanka preselio se u Beč, 1813. godine. Tu je upoznao Jerneja Kopitara, cenzora slovenskih knjiga, na čiji je podsticaj krenuo u prikupljanje srpskih narodnih pesama, reformu ćirilice i borbu za uvođenje narodnog jezika u srpsku književnost. Vukovim reformama u srpski jezik je uveden fonetski pravopis, a srpski jezik je potisnuo slavenosrpski jezik koji je u to vreme bio jezik obrazovanih ljudi. Tako se kao najvažnije godine Vukove reforme ističu 1818., 1836., 1839., 1847. i 1852.