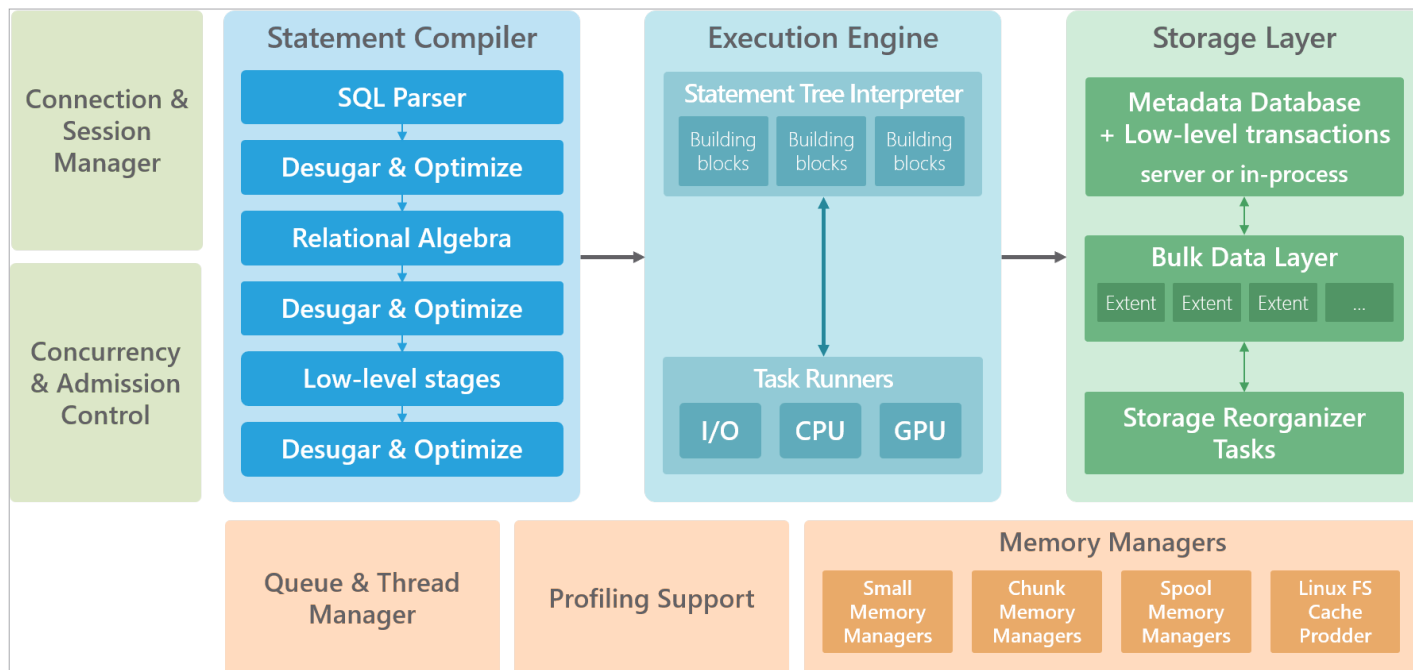


# Brza analiza velikih sku

**Modernizacijom poslovanja**, u skladu s trenutačnim tehnološkim zahtjevima svjetskog tržišta, brzina obrade velike količine digitalnih podataka ulazi u fokus svake tvrtke



Interna arhitektura SQream DB baze

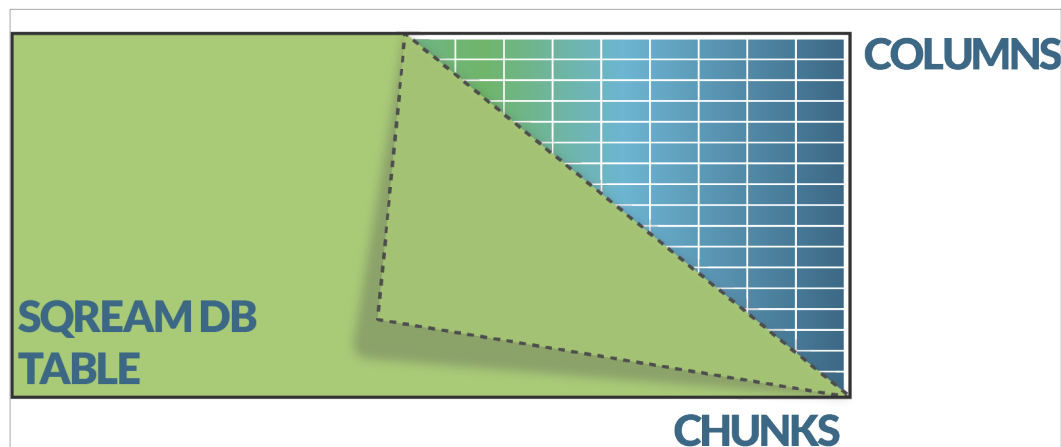
Vedran Podubski, konzultant za poslovna rješenja, Megatrend poslovna rješenja

**P**odaci se obrađuju u svakoj industriji i oni čine dio svakodnevnog poslovanja, neovisno o primarnoj djelatnosti tvrtke. Uporabom i obradom podataka iz raznih izvora – iz skladišta podataka, iz raznih oblika samostalnih datoteka (Excel, Open Office), IoT podataka (Internet Of Things), računovodstvenih softvera s vlastitim bazama podataka, itd. – količina informacija koje treba analizirati eksponencijalno se povećava, a sama brzina njihove obrade pada. "Tradicionalne" baze podataka postaju izvor sve većih troškova, zbog spore obrade i nužne dodatne optimizacije informacija te zbog stalnih ulaganja u hardver i ljudske potencijale.

Sva ta dodatna ulaganja, kako vremena tako i materijalnih resursa, izvode se u svrhu "popravka" performansi baza po-

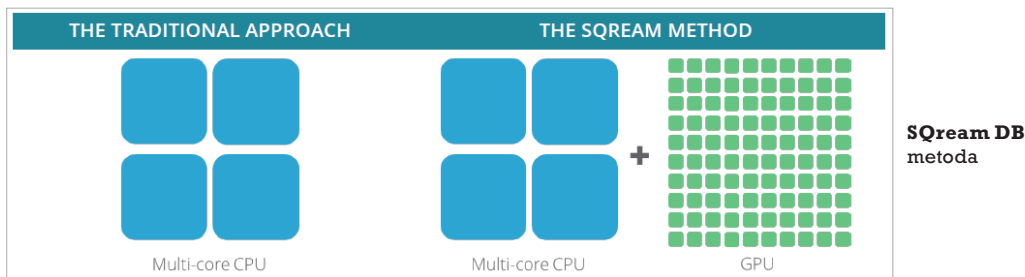
data. Tradicionalne RDBMS baze (Oracle, MS SQL Server, itd.), koje se koriste kao skladišta za analizu podataka, obra-

duju informacije u procesoru, u realnom vremenu, čitajući i pišući podatke na diskove serijski. Zbog takvog načina funkcioniranja – kad rade s većom količinom podataka i kad se suoče s većom kompleksnošću



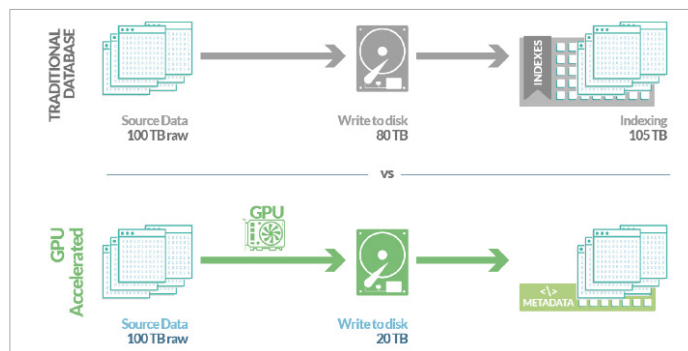
Particioniranje baze

# pova podataka



upita koji se izvršavaju - tradicionalnim bazama podataka padaju performanse.

Kako se to ne bi dogodilo, stalno treba ulagati u brže diskove, nove procesore i procesorske jezgre, u više radne memorije, itd. - kako bi se kompenzirala vremenska komponenta obrade. Uz to, svaka nadogradnja hardvera i seljenje podataka na nove servere, dovodi do zastoja u radu samog skladišta podataka.



Način zapisivanja podataka

## SQREAM DB

Izvođenje analitičkih operacija nad skupom digitalnih informacija unutar tvrtke - bilo u svrhu analize povijesnih podataka poslovanja, nadzora poslovanja ili planiranja budućnosti, bilo u automatiziranim procesima ili od tima analitičara - ne bi smjelo biti točki zagušenja. Stoga se mnoge tvrtke okreću suvremenim rješenjima, točnije, bazama podataka napravljenima baš u svrhu ubrzanja analitičkih operacija u velikim skupovima podataka. Za potrebe obrade i skladištenja velike količine podataka, točnije za veliki promet podataka na tjednoj i mjesečnoj razini, možemo preporučiti vrlo brzu bazu podataka - SQream DB.

Riječ je o suvremenoj RDBMS bazi, primarno zamišljenoj za skladištenje velike količine podataka (Big Data), bržoj jer za obradu podataka, umjesto običnih, koristi grafičke procesore (GPU). Prva inačica te

baze predstavljena je još 2014. godine u Silicijskoj dolini. Cilj je bio napraviti rješenje koje će ubrzati analizu velikih skupova podataka uz pomoć višejezgrenih procesora Nvidijinih grafičkih kartica, paralelnim izvršavanjem upita nad bazom. Rješenje SQream DB napravljeno je iznova - točnije, u njegovu dizajnu nije korišten nijedan postojeći sustav kao temelj za razvoj, primjerice Hadoop ili Postgres. Izvršavanje upita na

grafičkim procesorima tehnologija je koja je vrlo slična sustavima korištenima za rudarenje podataka kod kriptovaluta te omogućuje masivnu paralelnu obradu podataka na svakoj jezgri procesora grafičke kartice.

To postiže koristeći višu frekvenciju grafičke memorije, koja je mnogo brža od uobičajene RAM memorije na matičnoj ploči. U standardnim, često korištenim skladištima podataka sve komponente unutar sustava

usko su povezane i zajednički koriste hardverske resurse, što kod velikog protoka podataka i velikog broja korisnika otežava skaliranje i stvara probleme s performansama. SQream DB taj problem rješava inteligentnom internom arhitekturom koristeći odvojeni kompajler, izvršni dio i spremnik podataka, kako bi bolje optimizirao protok podataka i njihovu obradu.

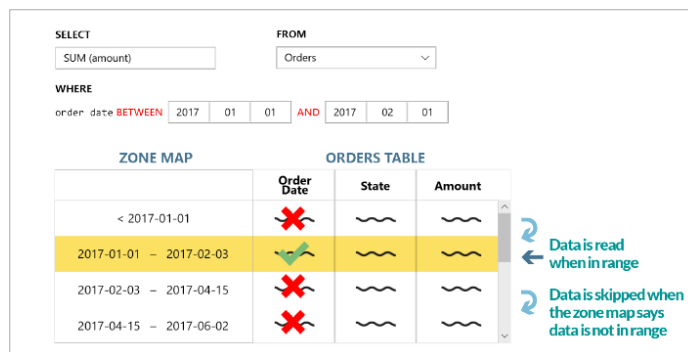
## DVA NAČINA PARTICIONIRANJA

Sljedeći korak ubrzanja performansi unutar SQream DB baze postiže se partitioniranjem podataka. To se izvodi na dva načina. Prvo je hiperpartitioniranje i namijenjeno je većoj kompresiji podataka i ubrzanju njihova protoka te se izvršava potpuno automatski. Taj dio partitioniranja je vertikalni, točnije kolumnarni, i omogućuje selektivni pristup određenim podskupovima kolumni u bazi, čime se smanjuje potreba za čestim pisanjem/čitanjem podataka s diskova.

Ta vrsta partitioniranja savršena je za paraleliziranu obradu podataka, primjerice, preko grafičkog procesora. Drugi dio partitioniranja je horizontalan. Izvodi se podjelom podataka na komade manjih opsega (engl. chunks and extents). Horizontalna podjela podataka na manje podskupove omogućuje bolju iskorištenost hardvera i relativno male količine GRAM-a (RAM na grafičkoj kartici), što se postiže spajanjem podataka (engl. spooling) i inteligentnim korištenjem predmemorije (cache).

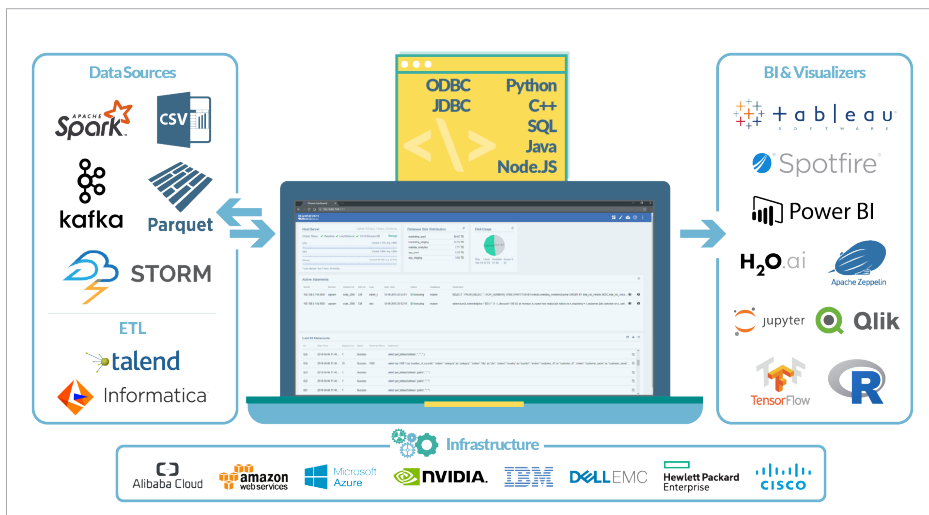
## METAPODACI I NAČINI ZAPISIVANJA

Standardne data warehouse baze podataka koriste isključivo procesorske jezgre i RAM memoriju za obradu, upisivanje i dohvaćanje podataka. Kod SQream DB baze taj proces je proširen i na inteligentno korištenje kombinacije dostu-



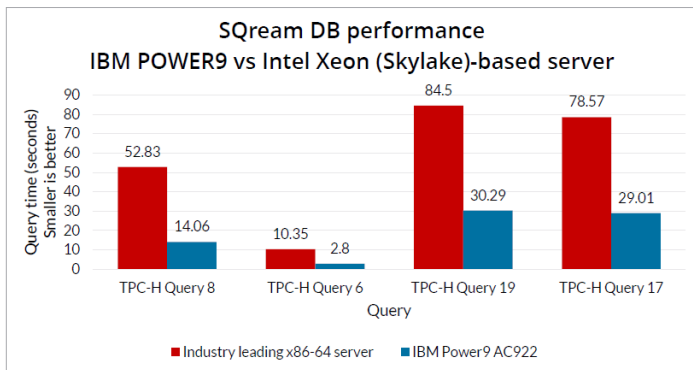
Metadata zonske mape

EKOSUSTAV



SQream DB ekosustav


**SQream DB** može raditi na većini standardnog serverskog x86 - 64 hardvera s Nvidijinim grafičkim karticama, pa čak i na komercijalnim laptopima opremljenim takvim hardverom, no za najbolji radni učinak preporučuju se 2x grafičke kartice Nvidia Tesla (K80, P40, P100 itd.), a za još veće ubrzanje IBM-ovi procesori POWER9 na kojima se performanse podižu i do 3,7 puta, prema testiranjima. Na nezavisnim testiranjima performansi SQream DB baze, u sustavu mobilnog operatera, pri "probavljanju" 1,6 TB podataka tjedno, performanse u usporedbi s konkurentskom bazom podataka (Greenplum), pokazuju od 5 do 18 puta veću brzinu, uključujući unos i kompresiju podataka i brzinu izvršavanja upita. Baza je dostupna u obliku softvera koji se može instalirati na standardnu x86 - 64 ili IBM-ovu POWER9 arhitekturu s Nvidijinim grafičkim karticama, kao servis u *cloudu* (Amazon P2 / P3 with NVIDIA Tesla, Azure NCv3 with Tesla V100) i IBM-ove Bluemix bare-metal sustave. Za dodatne informacije i stručne savjete oko SQream DB baza i/ili IBM-ova POWER9 sustava slobodno se obratite našim stručnjacima na poslovna.rjesenja@megatrend.com. ◀

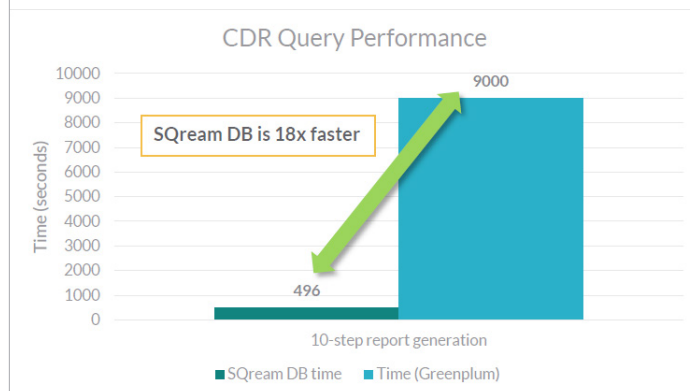
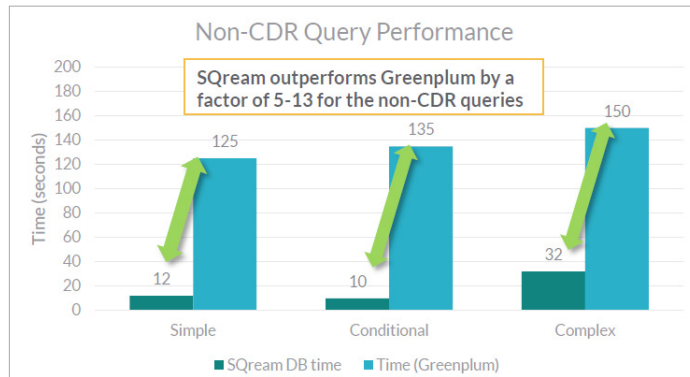


IBM POWER9 vs. Intel

pnih resursa procesora, RAM memorije i grafičkih procesora. Primjerice, interni sustav u bazi automatski koristi centralni procesor (CPU) ako bi kopiranje podatka u grafički procesor (GPU) uzelo previše vremena i usporilo upit/obradu. Tako se postiže mnogo brža obrada podataka. Još jedan revolucionarni pristup spremanju podataka kod SQream DB-a je i inteligentno korištenje metapodataka generiranih obradom preko grafičkih procesora.

Metapodaci sadrže opisne podatke o opsegu (engl. range) i vrijednosti svakog skupa podataka (engl. chunks), te su spremljeni zasebno od stvarnih podataka, čime se omogućuje inteligentno preskakanje nepotrebnih opsega podataka prilikom izvršenih upita. Tako

se stvaraju tzv. zonske mape, što za rezultat ima smanjenje uporabe hardverskih resursa. Uz to, SQream DB je potpuno ANSI - 92 SQL kompatibilna i lako se implementira u sve informatičke ekosustave jer ima ugrađenu podršku za sve tipične ODBC i JDBC konektore, uključujući i Python, C#, .NET, C++, Java i druge. Izvorna podrška za SQL jezik omogućuje korištenje bilo kojeg ETL alata i ostalih aplikacija nad bazom, čime se smanjuje vrijeme implementacije na minimum. Prema mjerenjima iz prakse, količina "probavljenih" podataka (naravno, u ovisnosti od hardvera) može biti i do 3,5 TB na sat, iz raznih izvora, te se može implementirati kao sloj između Apache Kafke i Apache Sparka te služi kao sloj za analitiku između to dvoje. 



SQream DB vs. Greenplum