

## 2-12 Napredne teme baza podataka

Sve situacije do sada su obuhvaćale bazu podataka ograničene veličine. Tj. nisu razmatrani slučajevi gdje baza podataka postane dovoljno velika da ne može stati na prostor dostupan u jednom računar. U ovim slučajevima potrebno je podijeliti bazu podataka na više računara. Zbog ovih potreba postoje *data* centri, koji skladište podatke na hiljadama računara.

Npr. WhatsApp razmjenjuje 65 milijardi poruka dnevno, Facebook generiše 4000 TB podataka dnevno, služeći 2.3 milijarde korisnika<sup>9</sup>. Procjenjuje se da Facebook skladišti 300.000 TB podataka ili 300 PetaB podataka. Svrha skladištenja ove količine podataka je bolji uvid u načine na koji korisnici koriste neki sistem, kako bi se od toga izvukli korisni i profitabilni podaci.

Baze podataka koje skladište podatke na više računara, ne nužno na istoj fizičkoj lokaciji, zovu se distribuirane baze podataka. One su jedna vrsta distribuiranih sistema koje su posebno polje izučavanja u računarskim naukama i informatici. Osim zahtjeva za skladištenje, potrebno je zadovoljiti zahtjeve za pretragu svih tih podataka.

Distribuiranje sistema na više računara donosi niz izazova, kao što su potencijalni nestanci struje, nefunkcionisanje mreže kojom su računari povezani, problemi u radi mrežnih uređaja, problemi prestanka rada diskova, ili drugih komponenti računara. Npr. jedan primjer crkavanja diskova je 0.14% na godišnjoj bazi<sup>10</sup> stoga ako centar ima oko 150.000 diskova to znači da će se pokvariti oko 210 diskova na godišnjem nivou, što znači da će trebati mijenjati disk skoro svaki dan.

Ovo su primjeri različitih potreba i različitih problema. Zbog dinamičnosti razvoja sistema, postojanje strukture podataka kakva postoji u SQL bazama, donosi niz potencijalnih problema radi organizacije modela podataka. Spašavanje svakog klika ili pokreta miša korisnika, kao i svih podataka desetina ili stotina miliona korisnika, donosi jedinstvene izazove za svaki sistem. Za ove i ostale svrhe, napravljene su druge vrste podataka kako bi se mogli spasiti podaci sa dinamičkom strukturom.

Poseban problem je pretraživanje i analiziranje svih ovih podataka. Npr. analiziranje nekog skupa podataka obično se odvija nad istom kolonom. Kao npr. dobavljanje prosječne ocjene svih studenata. U tom slučaju nisu potrebne ostale kolone, ali SQL baze podataka spašavaju redove na disk u jednom nizu, tako da bi se dobavili podaci jedne kolone iz svih redova, baza podataka i disk, moraju preći sve dostupne redove. Da bi se ovaj problem riješio napravljene su baze podataka koje spašavaju podatke sekvencijalno po koloni, tako da su svi redovi jedne kolone spašeni u nizu na disku, te je pretraživanje podataka brže i habaju disk manje.

To je jedan primjer NoSQL baze podataka, mada se pojam NoSQL odnosi za sve baze podataka koje nemaju fiksnu strukturu. Dati primjer analize podataka se zove *Analitičko Procesiranje* ili

---

<sup>9</sup> <https://medium.com/@srank2000/how-facebook-handles-the-4-petabyte-of-data-generated-per-day-ab86877956f4#:~:text=Hive%20is%20Facebook's%20data%20warehouse,map%2Dreduce%20jobs%20per%20day>.

<sup>10</sup> <https://www.backblaze.com/blog/backblaze-hard-drive-stats-q2-2020/>

OLAP (*Online Analytical Processing*). Nije ni strano da se podaci dupliraju u dvije različite baze podataka, koje su optimizovane za različite svrhe. SQL baze podataka, podržavaju tzv. *Transakcijsko Procesiranje* ili OLTP (*Online Transactional Processing*), tj. garantovanje procesiranja transakcija koje zadovoljavaju ACID pravila, tj. garanciju konsistentnosti podataka. Međutim zadovoljavanje ovih uslova, onemogućava bazu podataka da brzo procesira podatke, pogotovo ne brzo onoliko kada je procesiranje potrebno milionama korisnika.

Još jedan naziv za OLAP tj. podnaziv je i *Data Mining* ili kopanje informacija, jer postoje osobe zadužene za pronalazak novih znanja u moru podataka. Tako raznorazne kompanije mogu pronaći navike kupca kao i korisnicima omiljene stvari. Npr. je li vam se ikad desilo da potražite na Googl-u ili Ebay-u neki proizvod a da bi naknadno vidjeli reklame za isti proizvod na potpuno drugoj stranici.

Kako bi se svi ovi podaci spasili u distribuirane baze podataka, postoje razni pristupi rješavanju ovog problema. Svi ovi pristupi pokušavaju zadovoljiti CAP teoremu. CAP teorema govori da bilo koja distribuirana baza podataka može zadovoljiti dva od naredna tri uslova: Konsistentnost (Consistency), Dostupnost (Availability) i Toleranciju na pracionisanje (Partition tolerance). Drugim riječima, distribuirana baza podataka ne može zadovoljiti absolutne zahtjeve, stoga će pokušati zadovoljiti neke od njih. Konsistentnost ovih sistema ne mora biti na prvom mjestu ali će se konsistentnost eventualno pokušati zadovoljiti.

Jedan od načina korištenja distribuiranih baza podataka je *batch-processing* pristup. Ovim pristupom, vrši se procesiranje određene količine podataka, ili svih, tako što se pokrene program u zakazano vrijeme i obradi te podatke. Drugi, noviji, način je koristeći programa koji konstantno procesiraju beskonačni niz podataka, kako dolaze budu obrađeni i zove se *stream-processing*.

Osim ovih tipova baza podataka, postoje ostale specijalizovane baze podataka. Jedan primjer je baza podataka za pretraživanje teksta i sličnosti teksta. Npr. potreba za pretragom artikala na oglasniku čak ni kada naziv ne podudara striktno, jer je možda neko napravio grešku u kucanju. Primjeri te vrste podataka su Apache Lucene i Elasticsearch. Još jedan primjer specijalizovane baze podataka su baze za Geografske Informacione Sisteme (GIS) gdje je potrebno brzo pretraživanje 2D ili 3D podataka tipa udaljenost od grada, ili raspon između GPS koordinata. Iako baze podataka poput MongoDB ili PostgreSQL imaju ove sposobnosti, za određene potrebe postoje specijalizovane baze podataka.

---

Postoji još mnogo tema za izučavanje baza podataka kao i sistema koji ih koriste poput: 2-phase-commit, replikacija, particionisanje, nivoi izolacije, skladišta podataka, key-value store, linearizibilnost, distribuirane transakcije itd, međutim ovdje se fokusiralo na osnovne funkcije SQL i objektnih baza podataka.

Za potrebe velikog broja jednostvnih sistema, ova dokumentacija može biti dovoljna podloga, a za sve ostale detalje tu su dokumentacije svake pojedinačne baze podataka koje detaljno opisuju sve dostupne funkcionalnosti.