

SKLADIŠTE PODATAKA - UVOD

Definicije

Skladište podataka (data warehouse): Domenski orijentisan, integriran, vremensko promenljiv i neuništiv skup podataka namenjen podršci odlučivanju kod upravljanja nekim sistemom.

Podskladište podataka (data mart): Poseban izdvojeni deo skladišra podataka namenjen potrebama nekog dela sistema.

Web skladište podataka (data webhouse): Distribuirano skladište podataka implementirano preko web-a (za koje ne postoji centralizovano čuvanje podataka).

U sva tri slučaja, reč je o strateškom IS, za razliku od operacionog IS.

Pogodnosti

Visok stepen povraćaja ulaganja

Povećanje konkurentnosti

Povećanje produktivnosti odlučivanja

Povećanje kvaliteta odlučivanja

Poređenje

Operacioni (OLTP) IS

Trenutni podaci stanja i prometa

Detaljne podatke

Dinamički podaci

Ponavljajuća predefinisana obrada

Visok nivo transakcione aktivnosti

Predvidiv način korišćenja

Transakciono orijentisan

Aplikativno orijentisan

Podrška svakodnevnom odlučivanju

Opslužuje veliki broj korisnika

Strateški (OLAP) IS

Istorijski podaci stanja i prometa

Detaljni i srednje i visoko svodni podaci

Ugaljnom statički podaci

Od hoc obrada po zahtevu

Nizak / srednji nivo transakcione aktivnosti

Nepredvidiv način korišćenja

Analitički orijentisan

Domenski orijentisan

Podrška strateškom odlučivanju

Opslužuje manji broj korisnika

Problemi koji se javljaju u vezi skladišta podataka

Podcenjivanje resursa potrebnih za punjenje podacima

Skriveni problemi unutar izvornih IS

Neobuhvatanje neophodnih podataka unutar izvornih IS

Semantika i homogenizacija podataka

Visoki zahtevi za resursima

Vlasništvo/pristup podacima

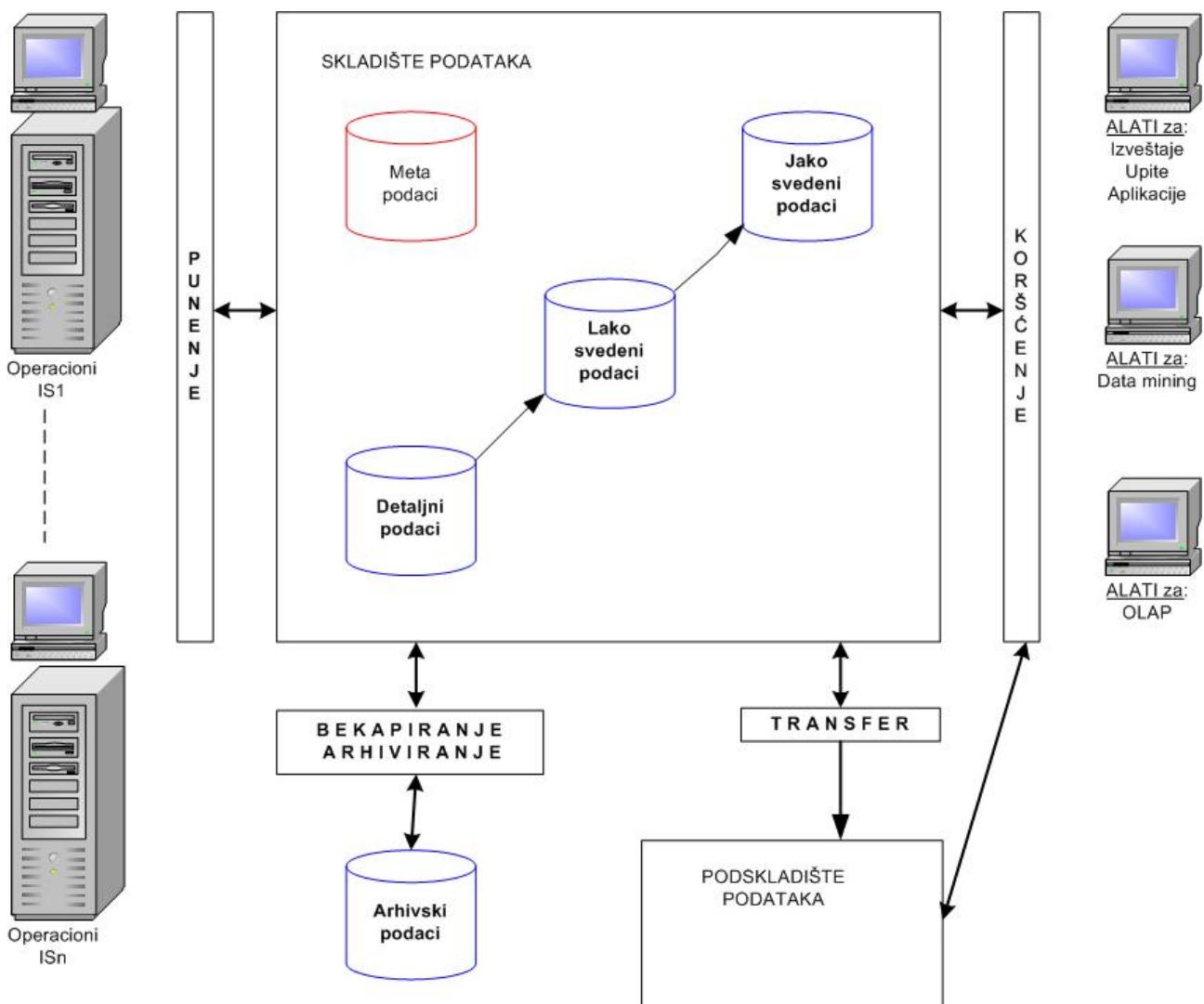
Obimno naknadno održavanje

Dugoročnost projekta (≥ 3 godine)

Kompleksnost integracije sistema

SKLADIŠTE PODATAKA – USTROJSTVO I RAD

Struktura skladišta podataka



DATA WAREHOUSE

DATA MART

- skladište podataka

- podskladište podataka

(globalna namena)

(specifična namena)

Implementacija skladišta podataka

Uobičajeno je da se vremenom iz jednog skladišta podataka formiraju podskladišta podataka, pri čemu isti podaci i dalje ostaju u skladištu podataka.

U praksi se dešava obrnuto - prvo se realizuju pojedina kritična podskladišta podataka koja se direktno pune iz operacionih IS, a zatim se naknadno formira skladište podataka. Od tog trenutka, tok podataka iz operacionih IS je ka skladištu podataka, a podskladišta podataka se pune transferom iz skladišta podataka.

Implementacija skladišta podataka je preko baze podataka (uglavnom relacione) sa visokim stepenom redundanse, a mehanizam suođenja je zasnovan na okidačima.

Punjjenje skladišta podataka

U principu postoje dva načina punjenja podacima iz operacionih IS:

- ✓ **Totalno punjenje:** U određenim vremenskim trenucima, skladište se isprazni a zatim ponovo napuni podacima iz operacionih IS.
- ✓ **Inkrementalno punjenje:** Prilikom punjenja, u skladište se prenose samo izmene nastale u operacionim IS nakon prethodnog punjenja.

Postoje dve varijante inkrementalnog punjenja:

- ✓ **Paketno inkrementalno punjenje:** Vrši se u određenim vremenskim trenucima. Zahteva izmene u operacionom IS (bazi podataka) koje će implementirati mehanizam prepoznavanja nastalih izmena.
- ✓ **Neprekidno inkrementalno punjenje:** Vrši se neprekidno. Nakon svake promene u operacionim IS mehanizmom okidača vrši se prenos podataka ka skladištu podataka.

Mogući problemi kod punjenja skladišta podataka, naročito izraženi kod heterogenih operacionih IS, su:

- ✓ **Netačnost podataka iz operacionih IS:** Pri punjenju je neophodno filtriranje podataka, odnosno odbacivanje netačnih podataka.
- ✓ **Neusaglašenost podataka po tipu/preciznosti:** Pri punjenju je neophodno usaglašavanje po oba osnova.

SKLADIŠTE PODATAKA – DIMENZIONI MODELI

Dimenziono modeliranje

Polazi se od toga da se atributi entiteta (tabela) mogu podeliti na dve grupe:

- ✓ **atributi mere**: atributi čija vrednost odražava meru (veličinu) nečega; bitna osobina je mogućnost svođenja;
- ✓ **atributi dimenzije**: atributi čija vrednost služi kao osnov klasifikacije i svođenja; mogu biti bez i sa varijabilnom granulacijom (nivoima svođenja).

Dva specijalna atributa dimenzije sa varijabilnom granulacijom: [Prostorni](#) i [Vremenski](#)

| Vreme | Prostor |
|---------------|---------|
| Dan | Opština |
| Nedelja | Grad |
| Mesec | Oblast |
| Kvartal | Država |
| Godina | Regija |
| Dan-u-nedelji | |
| Dan-u-mesecu | |
| Deo-dana | |
| Sat-u-danu | |

Prostorni dimenzioni atribut ima jednu šemu hijerarhije. Najniži nivo rezolucije je GPS geografska pozicija iz koje se dalje izvode viši nivoi rezolucije. U tom pogledu očekuje se uvođenje novog standardnog tipa podatka (POSITION?) i odgovarajućih standardnih funkcija izvođenja.

Vremenski dimenzioni atribut ima 5 paralelnih šema hijerarhije. Najniži nivo rezolucije je vreme (sat-minut-sekunda...) i za to postoji standardni tip podatka TIMESTAMP. Postoji i čitav niz standardnih funkcija koje vrše izvođenje po raznim šemama hijerarhija (TIME, DATE, DAY_OF_WEEK, MONTH itd.).

Kod dimenzionog modeliranja može istovremeno da postoji više vremenskih dimenzionih atributa (na primer, Dan-u-nedelji i Sat), pod uslovom da su one nezavisne, odnosno da nizu izvodive jedna iz druge (kombinacija Dan i Dan-u-mesecu nema smisla).

Dve vrste tabela u dimenzionom modelu:

- ✓ **tabela fakata** (FT, jedna): sastoji se iz složenog primarnog ključa u koji ulaze ID-ovi svih dimenzija i jednog ili više atributa mere;
- ✓ **tabele dimenzija** (DTi, više): sastoje se iz prostog primarnog ključa koji odgovara jednoj komponenti primarnog ključa tabele fakata i jednog ili više atributa dimenzije.

Unutar jednog skladišta ili podskladišta podataka može se nalaziti više dimenzionalnih modela.

Varijante dimenzionog modeliranja

Šema "zvezda" (star): jedan dimenzioni model

Sadrži jednu FT i za svaku dimenziju tačno jednu DT. DT mogu da sadrže denormalizovane podatke (neključne funkcije zavisnosti koje nisu u 3. NF). Maksimalna dubina referisanja je 1.

FT --> DT_i

Šema "pahuljica" (snowflake): jedan dimenzioni model

Sadrži jednu FT i niz DT koje ne mogu da sadrže denormalizovane podatke, nego umesto toga sadrže reference na dodatne informacione tabele IT. Pri tome IT mogu da sadrže denormalizovane podatke. Maksimalna dubina referisanja je 2, pri čemu pojedine dimenzione tabele ne moraju imati reference na IT.

FT --> DT_i --> IT_i
DT_j

Šema "zvezda-pahuljica" (starflake): jedan dimenzioni model

Uz prethodni uslov ima i ograničenje da ni jedna IT ne sme da sadrži denormalizovane podatke, što dovodi do proizvoljne dubine referisanja

FT --> DT_i --> IT_ia --> IT_ib ...
DT_j

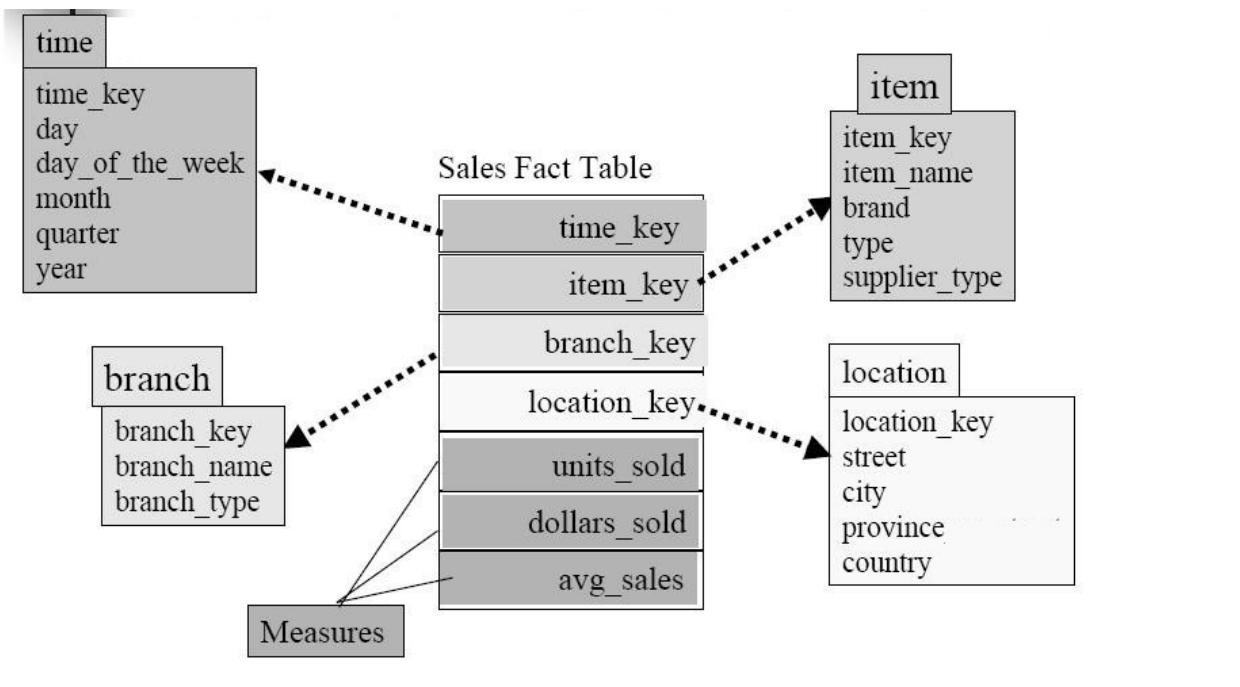
Šema sazvežđe (constellation): bar dva dimenziona modela

Situacija kada postoje bar dva dimenziona modela a njihove FT referišu bar jednu DT. Pri tome dimenzioni modeli mogu biti bilo kog od prethodna tri tipa.

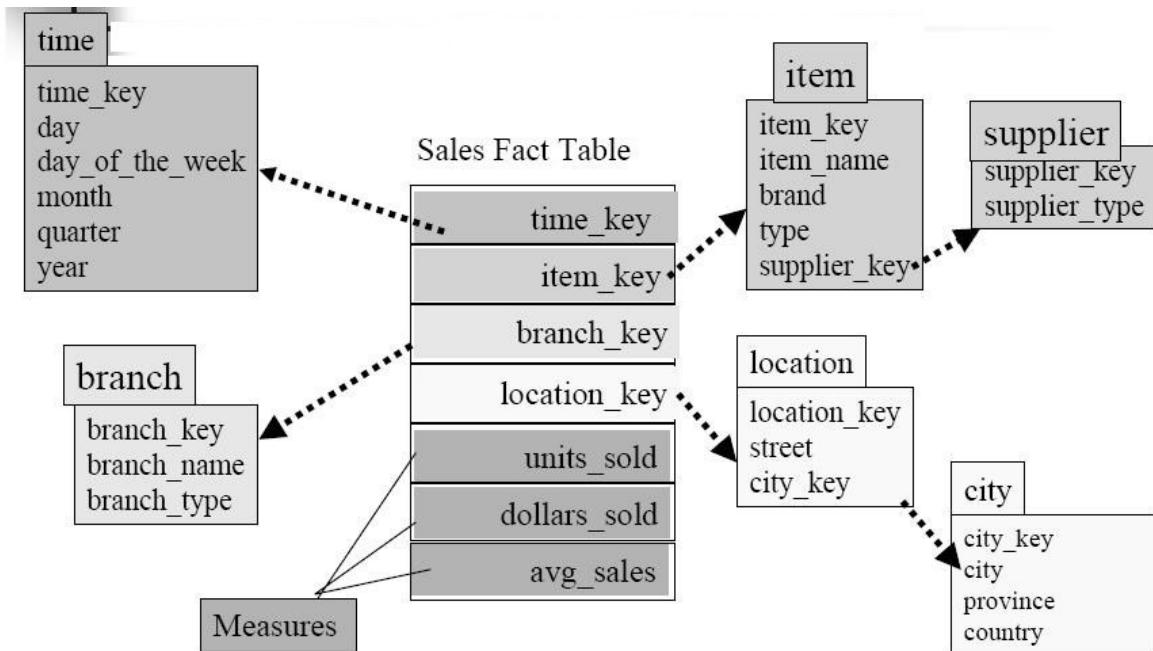
Prednost šeme "Zvezda" je smanjenje broja operacija spajanja, a manu povećanje memoriskog prostora.

Prednost ostalih šema je smanjenje memoriskog prostora, a manu povećanje broja operacija spajanja.

Primer "Zvezde": Prodaja artikala



Primer "Pahuljice": Prodaja artikala

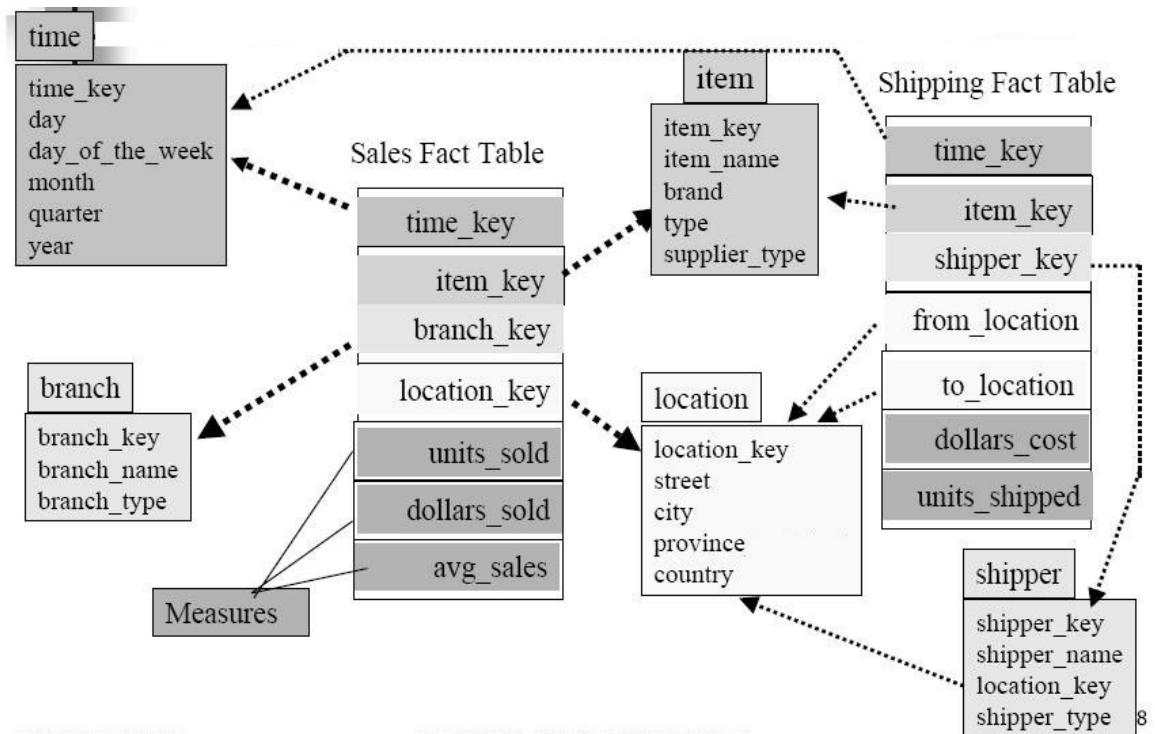


Napomene:

Kod oba modela dimenzije "branch" i "location" nisu potpuno nezavisne.

Poslednji model bi u varijanti "Starflake" imao dodatne informacione tabele province i country i lanac referenci location->city->province->country. Potpuna normalizacija po dimenziji time se izuzetno zbog kompaktnosti podataka ne primenjuje.

Primer "Sazvežđa":



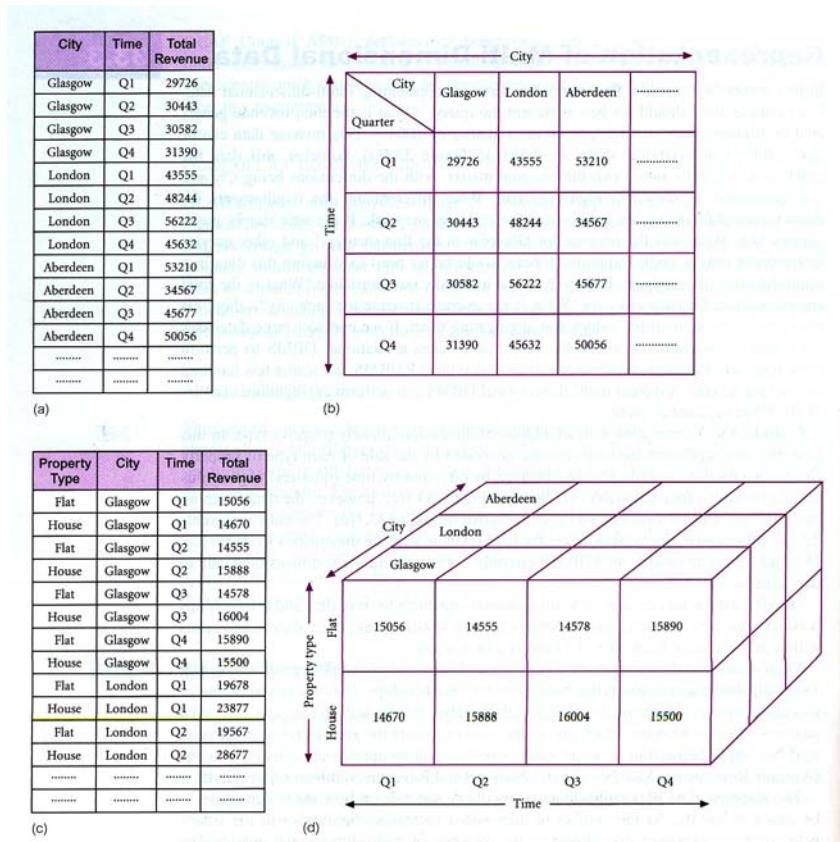
Napomene:

Prodaja (Sales) i otprema (Shipping) su dve odvojene tabele fakata zato što evidentiraju nezavisne i odvojene poslovne događaje.

U dimenzionom modelu otpreme dimenzija lokacija se javlja u dva svojstva (od i do) pa je referenca ka tabeli dimenzije lokacija.

SKLADIŠTE PODATAKA – DIMENZIONA ANALIZA (OLAP)

Osnovu OLAP dimenzione analize predstavlja OLAP kocka (CUBE) koja može biti sa 1 ili n dimezija (direktna vizualizacija je moguća do n=3).

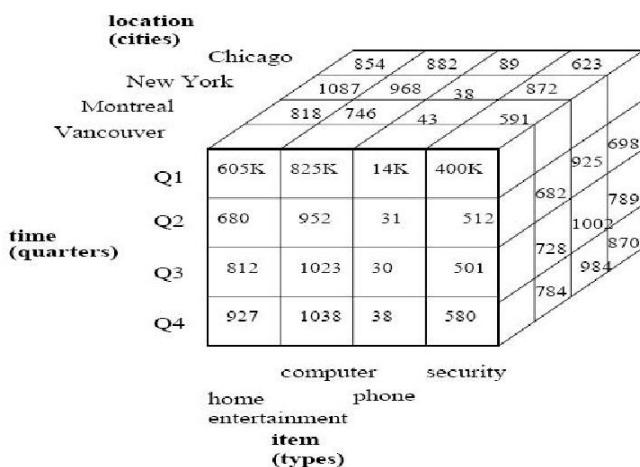


Primer: dvodimenziona (b) i trodimenziona (d) OLAP kocka za promet nekretnina. (a) i (c) su odgovarajuće tabelarne predstave.

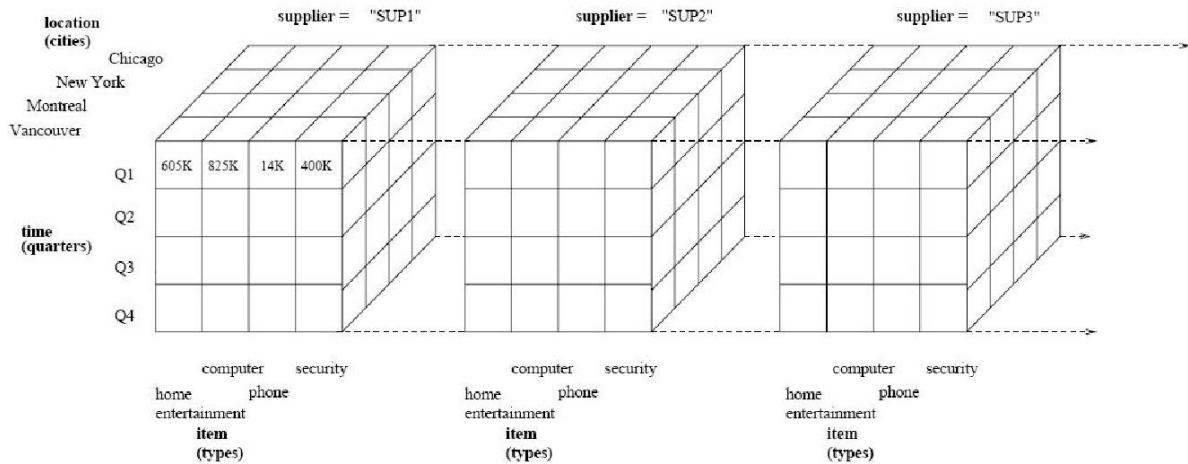
Primer: Prodaja artikala

| t i m e | location = "Vancouver" | | | | location = "Montreal" | | | | location = "New York" | | | | location = "Chicago" | | | |
|------------------|------------------------|-------|-------|------|-----------------------|-------|-------|------|-----------------------|-------|-------|-------|----------------------|-------|-------|------|
| | item | | | | item | | | | item | | | | item | | | |
| | home | comp. | phone | sec. | home | comp. | phone | sec. | home | comp. | phone | sec. | home | comp. | phone | sec. |
| Q1 | 605K | 825K | 14K | 400K | 818K | 746K | 43K | 591K | 1087K | 968K | 38K | 872K | 854K | 882K | 89K | 623K |
| Q2 | 680K | 952K | 31K | 512K | 894K | 769K | 52K | 682K | 1130K | 1024K | 41K | 925K | 943K | 890K | 64K | 698K |
| Q3 | 812K | 1023K | 30K | 501K | 940K | 795K | 58K | 728K | 1034K | 1048K | 45K | 1002K | 1032K | 924K | 59K | 789K |
| Q4 | 927K | 1038K | 38K | 580K | 978K | 864K | 59K | 784K | 1142K | 1091K | 54K | 984K | 1129K | 992K | 63K | 870K |

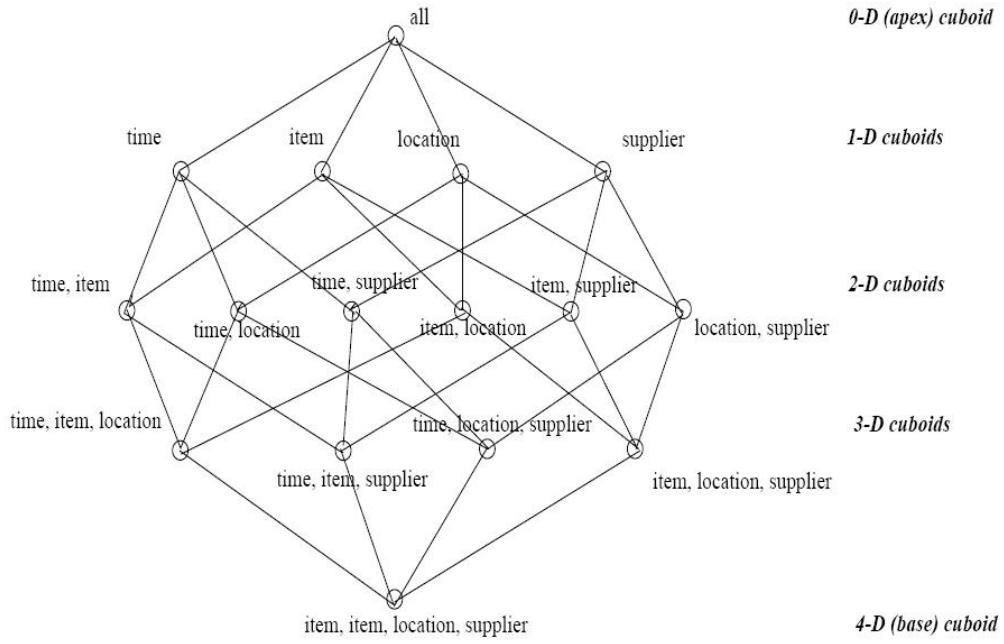
Tabelarni podaci za slučaj tri dimenzije (lokacija, vrsta artikla, vreme) i jedne mere (Iznos).



Trodimenziona OLAP kocka za prethodne tabelarne podatke.



Primer: Vizualizacija četvorodimenzione OLAP kocke preko trodimenzionih OLAP kocki od kojih svaka odgovara jednoj vrednosti četvrte dimenzije (isporučilac).



"Latisa kuboida" - sve moguće kocke za sve dimenzije (prethodni primer).

Operacije transformacije nad OLAP kockom (implementira ih OLAP Browser):

| | |
|------------|--|
| SLICE | Izdvajanje svodnih podataka za dati uslov po jednoj dimenziji. Rezultat je OLAP podkocka. |
| DICE | Izdvajanje svodnih podataka po datim uslovima dve ili više dimenzija. Rezultat je OLAP podkocka. |
| PIVOT | Operacija vizualizacije koja obrće dimenzione ose radi alternativnog prikaza podataka. |
| ROLL UP | Svođenje OLAP kocke, bilo penjanjem po hijearhiji dimenzije bilo izostavljanjem jedne dimenzije. |
| DRILL DOWN | Detaljizacija OLAP kocke, bilo spuštanjem po hijerahiji dimenzije bilo uvođenjem jedne nove dimenzije. Obrnuto od ROLL UP. |

U OLAP kocki se pri analizi uz izabrane dimenzije javlja samo jedna izabrana mera.

Ilustracija OLAP operacija:

