

Data Warehousing and Decision Support

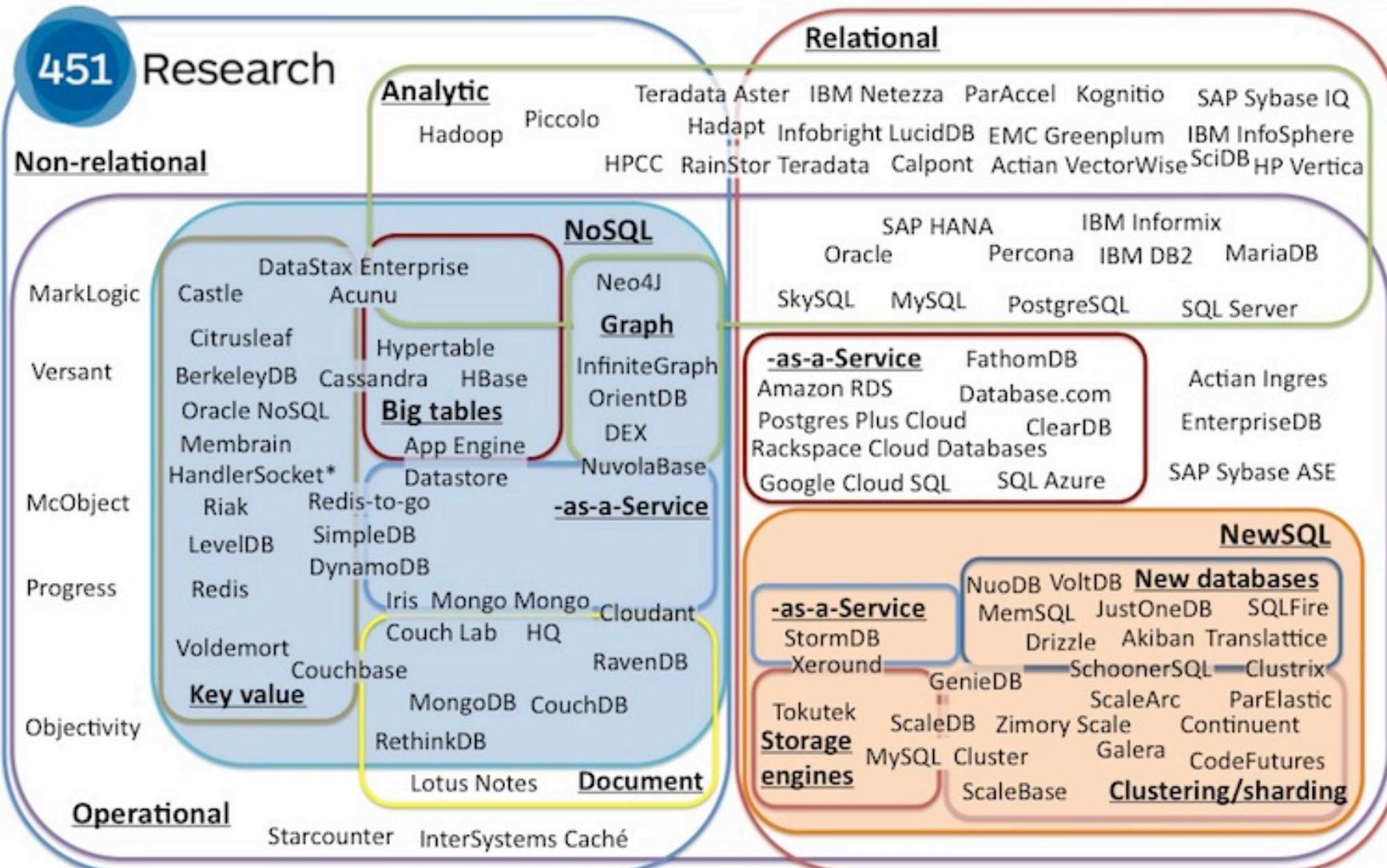
[R&G] Chapter 23, Part A

Machine Learning is not a business

- ❖ Databases/data warehousing is!
- ❖ ~\$50000/TB/year for DWH.
- ❖ Stable for decades.
- ❖ You can't sell a machine learning solution unless you package it with a data management system.
- ❖ But then, there's plenty of money to be made and nobody competes on price!

The evolving database landscape

451 Research

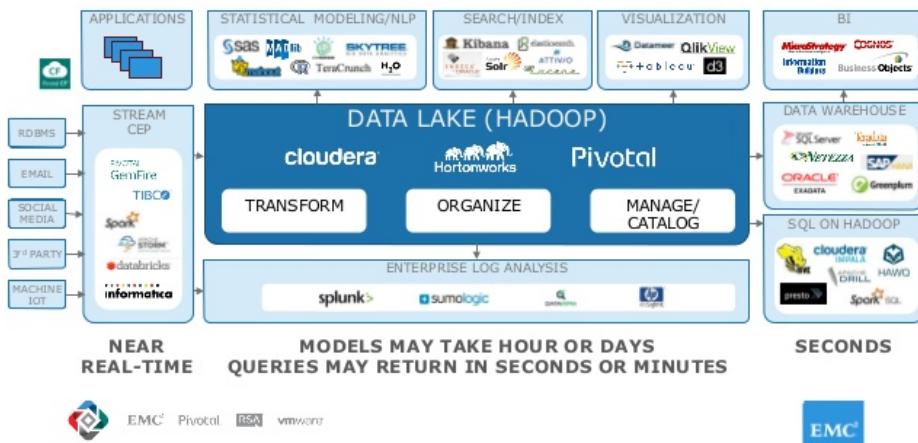


© 2012 by The 451 Group. All rights reserved.

Data warehousing is big business

- ❖ So lots of hype is generated constantly to keep the topic hot.
- ❖ Buzzwords
 - OLAP, Data warehouses, Data Marts, Data Lakes, Business Intelligence,...

THE BIG DATA LANDSCAPE



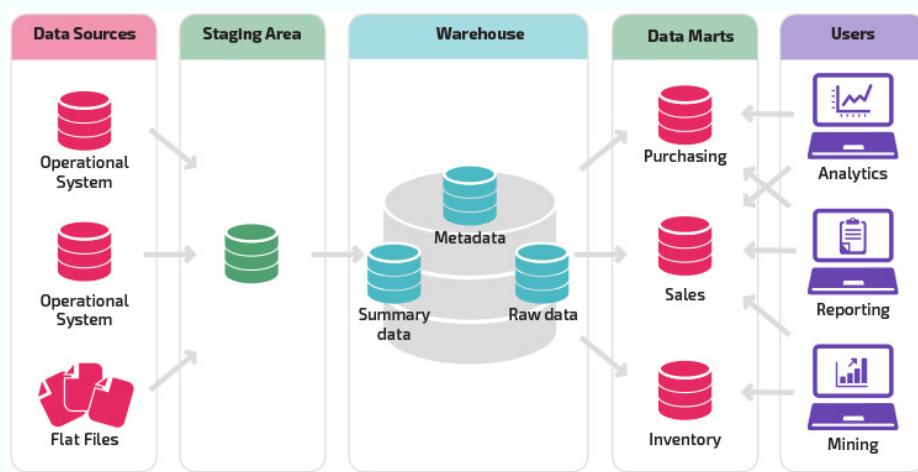
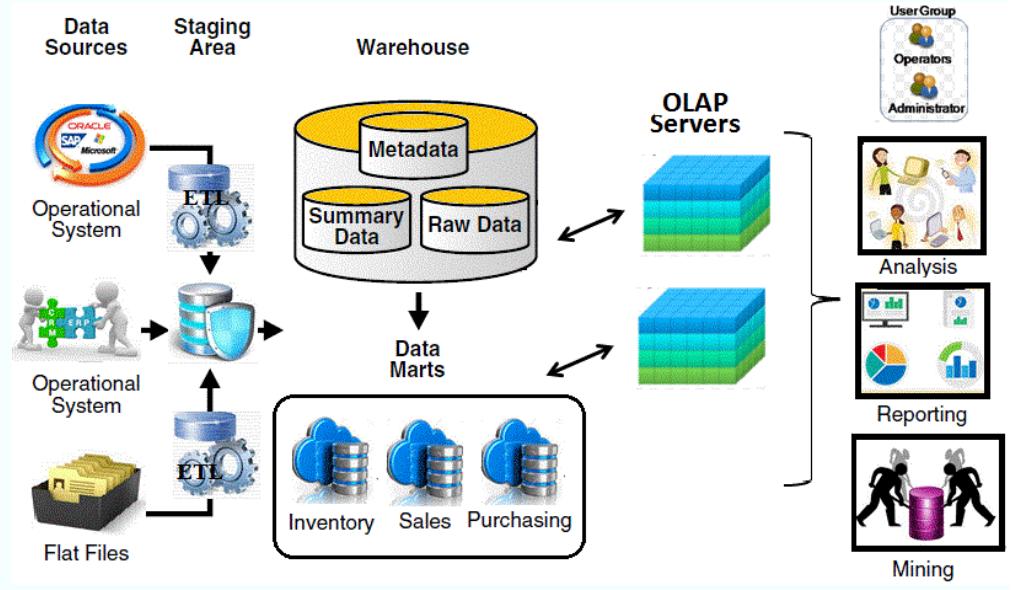
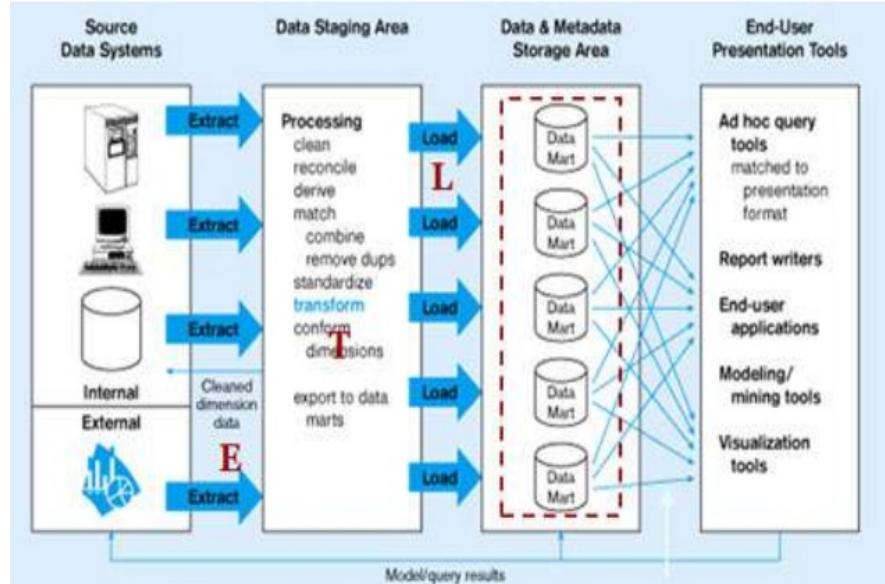
Azure Data Lake

Microsoft



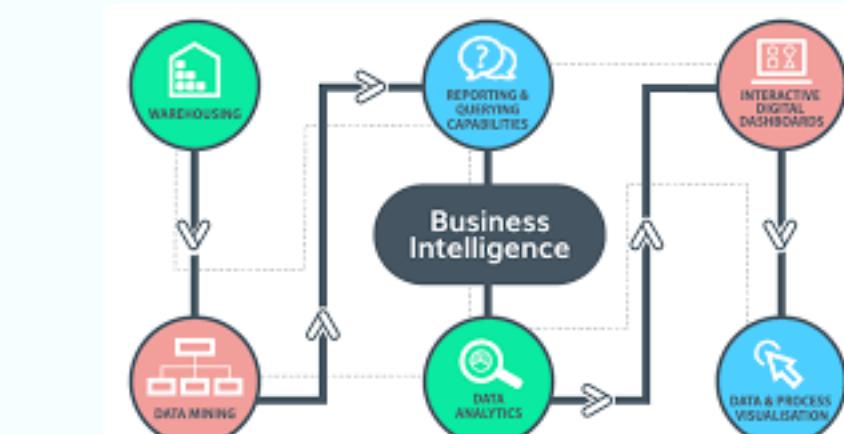
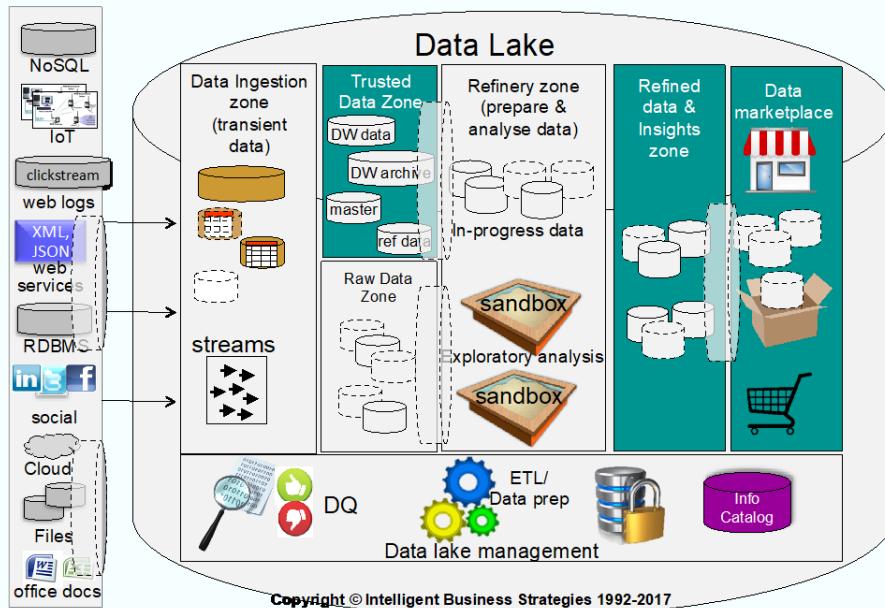
Most Important Use Group & Use-Cases	Time-to-Market Questions & Solutions	Cost Implementation & Ownership	Users (# & Types)	Data Growth Volume & Variety
Data Lake	Predictive & Advanced Analytics	Weeks - Months	\$\$\$\$	Bar chart showing high volume and variety
Data Warehouse	Multi-Purpose Enabler of Operational & Performance Analytics	Hours - Days	\$\$\$\$	Bar chart showing moderate volume and variety
Data Mart	Line of Business Specific Reporting & Analytics	Minutes - Hours	\$\$\$\$	Bar chart showing low volume and variety

Independent or stand alone Data Mart

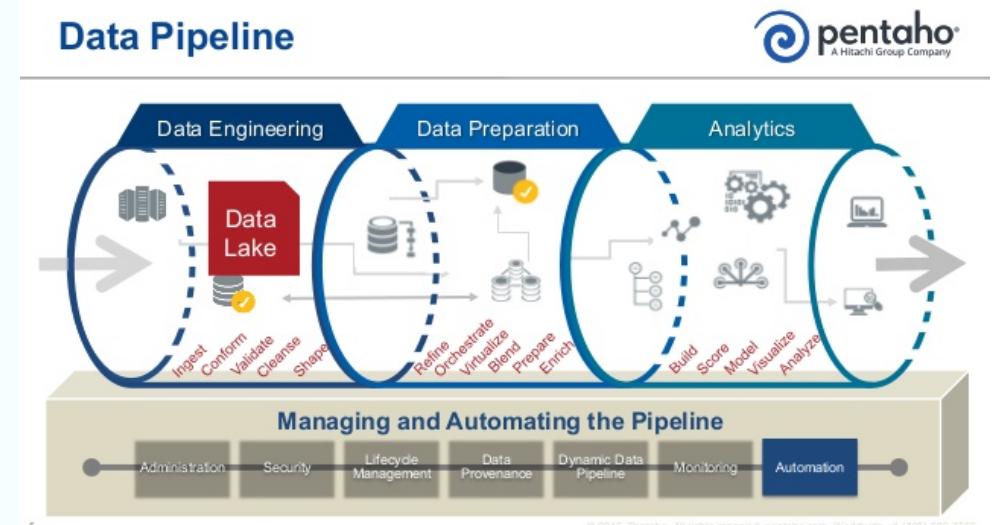


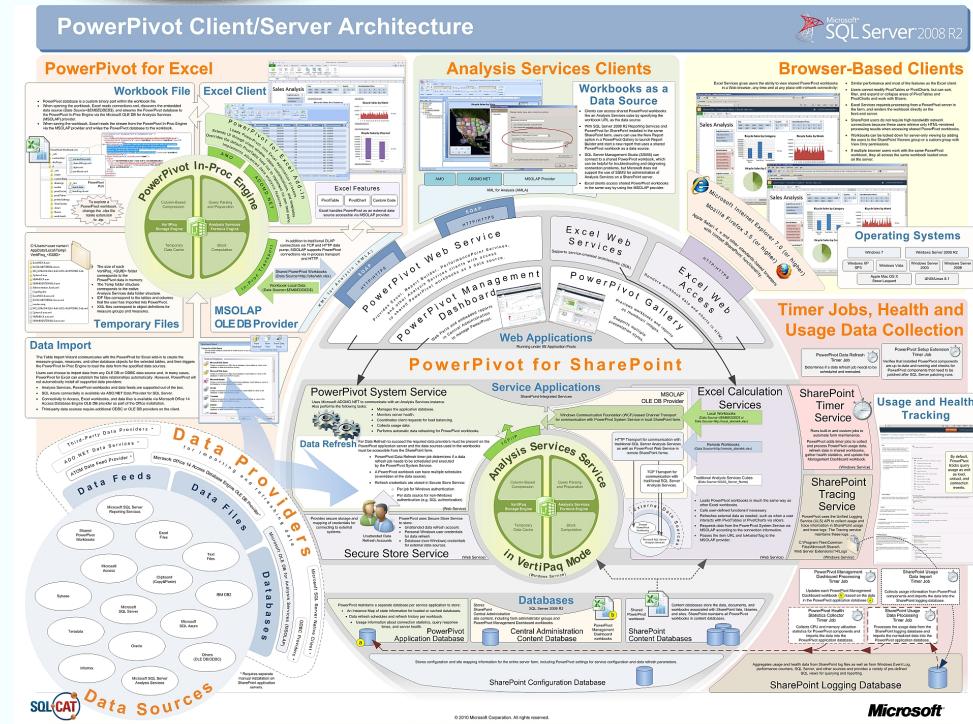
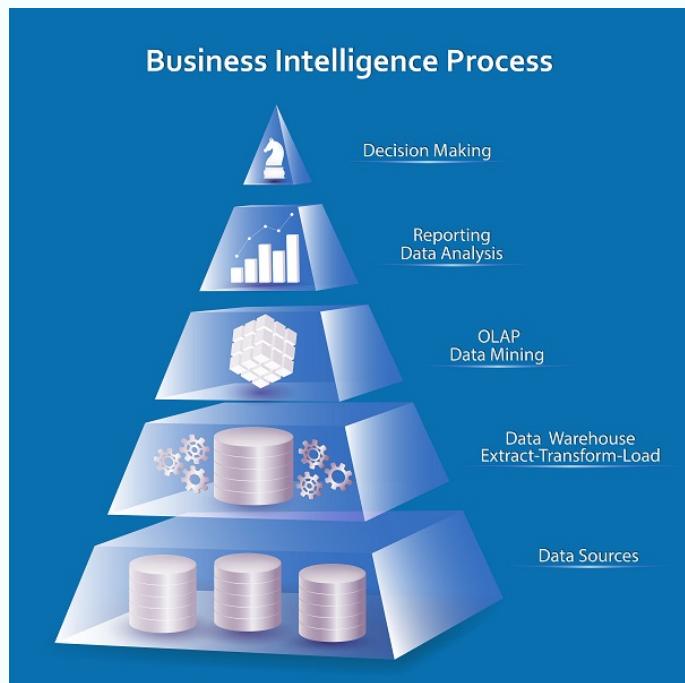
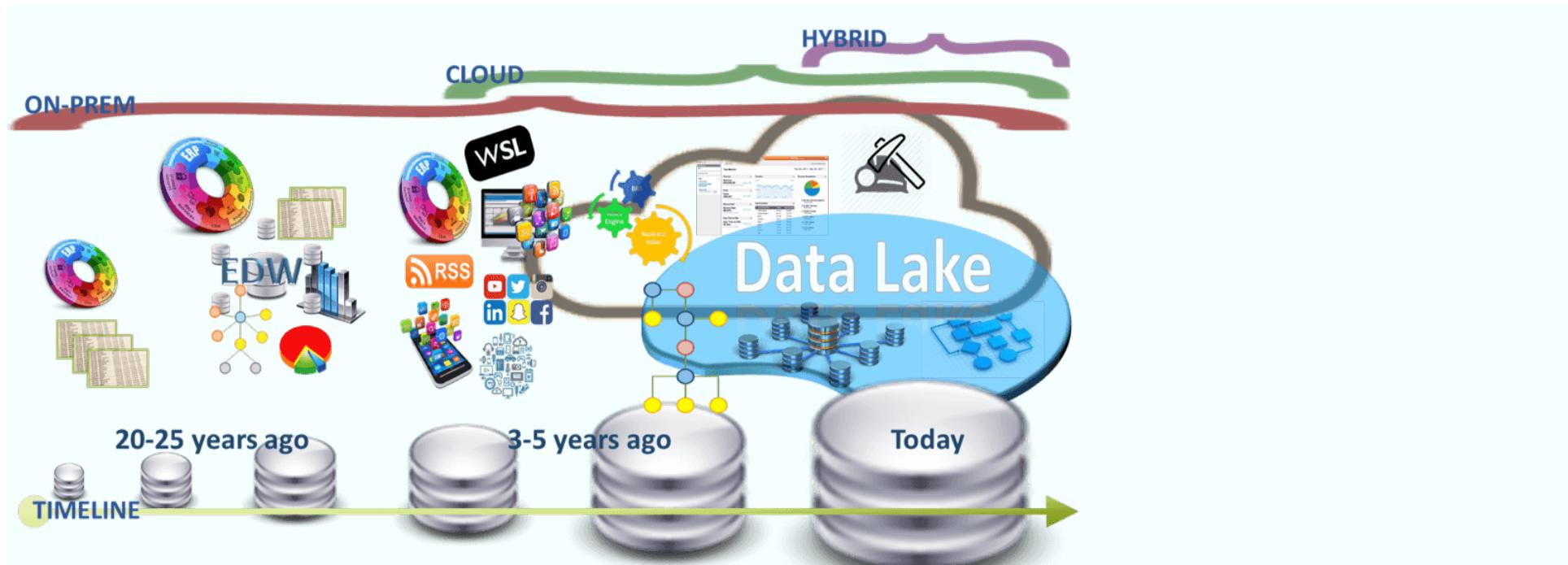


Data Virtualisation Can Help Introduce Agility Into A Data Lake While Reducing Data Copying



Data Pipeline





Introduction

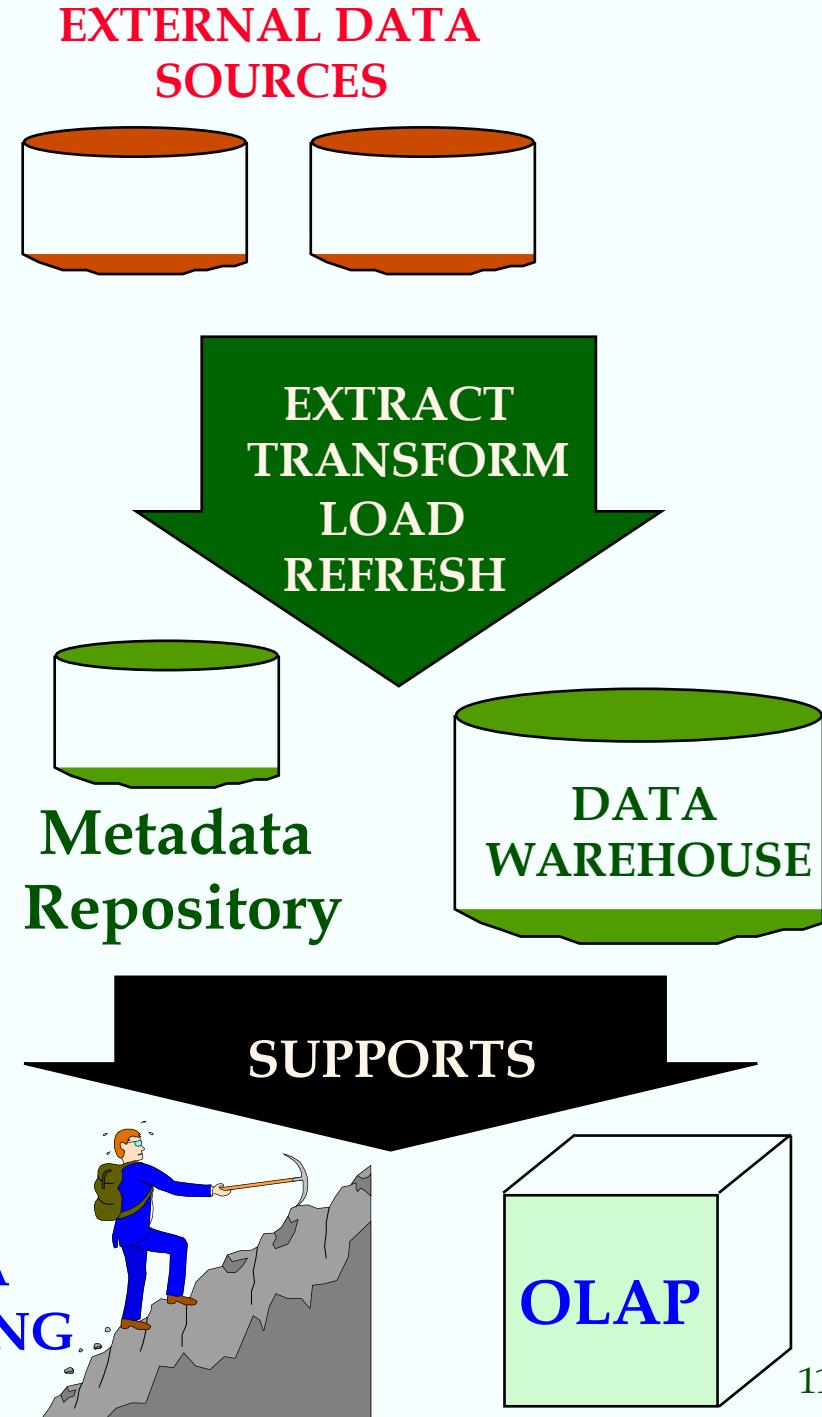
- ❖ Increasingly, organizations are analyzing current and historical data to identify useful patterns and support business strategies.
- ❖ Emphasis is on complex, interactive, exploratory analysis of very large datasets created by integrating data from across all parts of an enterprise; data is fairly static.
 - Contrast such **On-Line Analytic Processing (OLAP)** with traditional **On-line Transaction Processing (OLTP)**: mostly long queries, instead of short update Xacts.

Three Complementary Trends

- ❖ **Data Warehousing:** Consolidate data from many sources in one large repository.
 - Loading, periodic synchronization of replicas.
 - Semantic integration.
- ❖ **OLAP:**
 - Complex SQL queries and views.
 - Queries based on spreadsheet-style operations and “multidimensional” view of data.
 - Interactive and “online” queries.
- ❖ **Data Mining:** Exploratory search for interesting trends and anomalies.

Data Warehousing

- ❖ Integrated data spanning long time periods, often augmented with summary information.
- ❖ Several terabytes common.
- ❖ Interactive response times expected for complex queries; ad-hoc updates uncommon.



Warehousing Issues

- ❖ **Semantic Integration:** When getting data from multiple sources, must eliminate mismatches, e.g., different currencies, schemas.
- ❖ **Heterogeneous Sources:** Must access data from a variety of source formats and repositories.
 - Replication capabilities can be exploited here.
- ❖ **Load, Refresh, Purge:** Must load data, periodically refresh it, and purge too-old data.
- ❖ **Metadata Management:** Must keep track of source, loading time, and other information for all data in the warehouse.

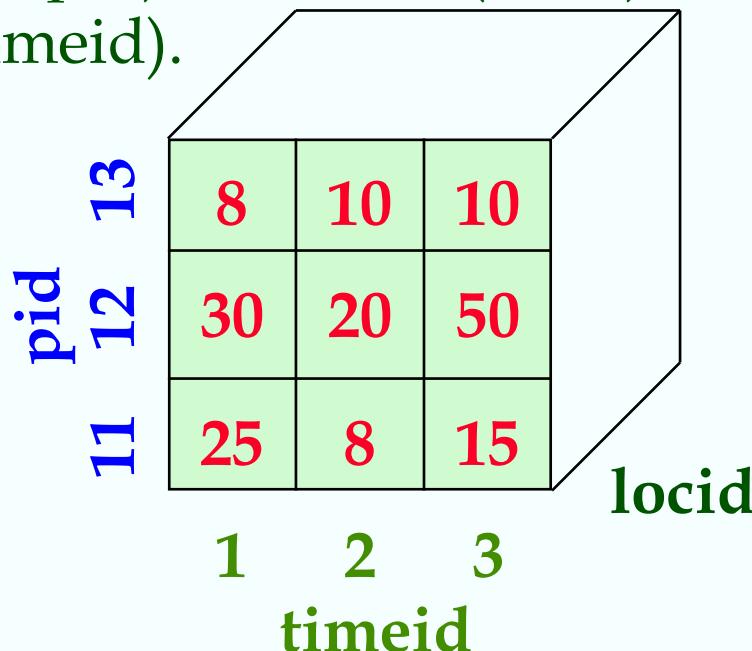
Data Quality; Data Cleaning

- ❖ Needs metrics; needs to be addressed proactively
 - Redesign processes for data gathering.
 - Feedback loops
 - Assign \$ cost to low quality data.
 - Assign people responsible for data quality
- ❖ Key problems:
 - Erroneous and missing data
 - Truncation and censoring
 - Entity resolution/deduplication/record linkage

Multidimensional Data Model

- ❖ Collection of numeric measures, which depend on a set of dimensions.
 - E.g., measure **Sales**, dimensions **Product** (key: pid), **Location** (locid), and **Time** (timeid).

Slice locid=1
is shown:



pid	timeid	locid	sales
11	1	1	25
11	2	1	8
11	3	1	15
12	1	1	30
12	2	1	20
12	3	1	50
13	1	1	8
13	2	1	10
13	3	1	10
11	1	2	35

MOLAP vs ROLAP

- ❖ Multidimensional data can be stored physically in a (disk-resident, persistent) array; called **MOLAP** systems. Alternatively, can store as a relation; called **ROLAP** systems.
- ❖ The main relation, which relates dimensions to a measure, is called the **fact table**. Each dimension can have additional attributes and an associated **dimension table**.
 - E.g., **Products(pid, pname, category, price)**
 - Fact tables are *much* larger than dimensional tables.

Dimension Hierarchies

- ❖ For each dimension, the set of values can be organized in a hierarchy:

PRODUCT

category
|
pname

TIME



LOCATION

country
|
state
|
city

OLAP Queries

- ❖ Influenced by SQL and by spreadsheets.
- ❖ A common operation is to aggregate a measure over one or more dimensions.
 - Find total sales.
 - Find total sales for each city, or for each state.
 - Find top five products ranked by total sales.
- ❖ Roll-up: Aggregating at different levels of a dimension hierarchy.
 - E.g., Given total sales by city, we can roll-up to get sales by state.

OLAP Queries

- ❖ Drill-down: The inverse of roll-up.
 - E.g., Given total sales by state, can drill-down to get total sales by city.
 - E.g., Can also drill-down on different dimension to get total sales by product for each state.
- ❖ Pivoting: Aggregation on selected dimensions.
 - E.g., Pivoting on Location and Time yields this cross-tabulation:
- ❖ Slicing and Dicing: Equality and range selections on one or more dimensions.

	WI	CA	Total
1995	63	81	144
1996	38	107	145
1997	75	35	110
Total	176	223	339

Comparison with SQL Queries

- ❖ The cross-tabulation obtained by pivoting can also be computed using a collection of SQLqueries:

```
SELECT SUM(S.sales)
FROM Sales S, Times T, Locations L
WHERE S.timeid=T.timeid AND S.timeid=L.timeid
GROUP BY T.year, L.state
```

```
SELECT SUM(S.sales)
FROM Sales S, Times T
WHERE S.timeid=T.timeid
GROUP BY T.year
```

```
SELECT SUM(S.sales)
FROM Sales S, Location L
WHERE S.timeid=L.timeid
GROUP BY L.state
```

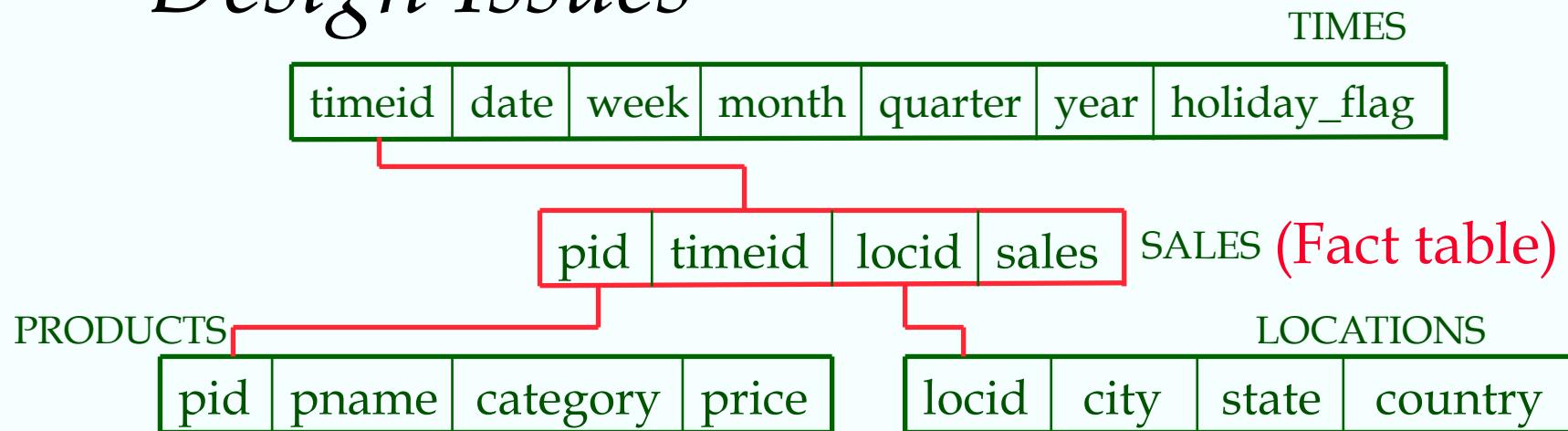
The CUBE Operator

- ❖ Generalizing the previous example, if there are k dimensions, we have 2^k possible SQL GROUP BY queries that can be generated through pivoting on a subset of dimensions.
- ❖ **CUBE pid, locid, timeid BY SUM Sales**
 - Equivalent to rolling up Sales on all eight subsets of the set {pid, locid, timeid}; each roll-up corresponds to an SQL query of the form:

Lots of work on optimizing the CUBE operator!

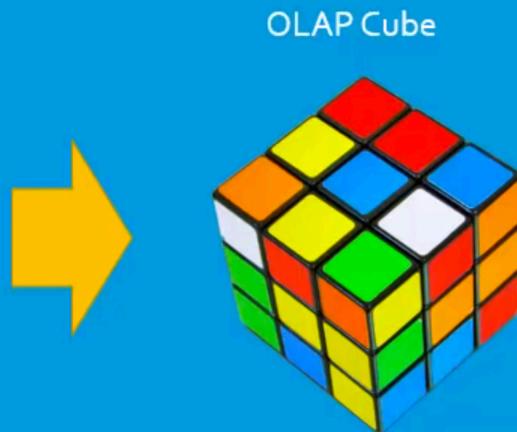
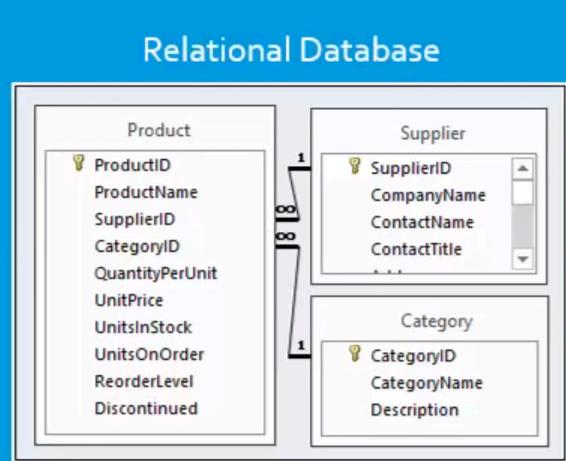
```
SELECT SUM(S.sales)
FROM   Sales S
GROUP BY grouping-list
```

Design Issues



- ❖ Fact table in BCNF; dimension tables un-normalized.
 - Dimension tables are small; updates/inserts/deletes are rare. So, anomalies less important than query performance.
- ❖ This kind of schema is very common in OLAP applications, and is called a **star schema**; computing the join of all these relations is called a **star join**.

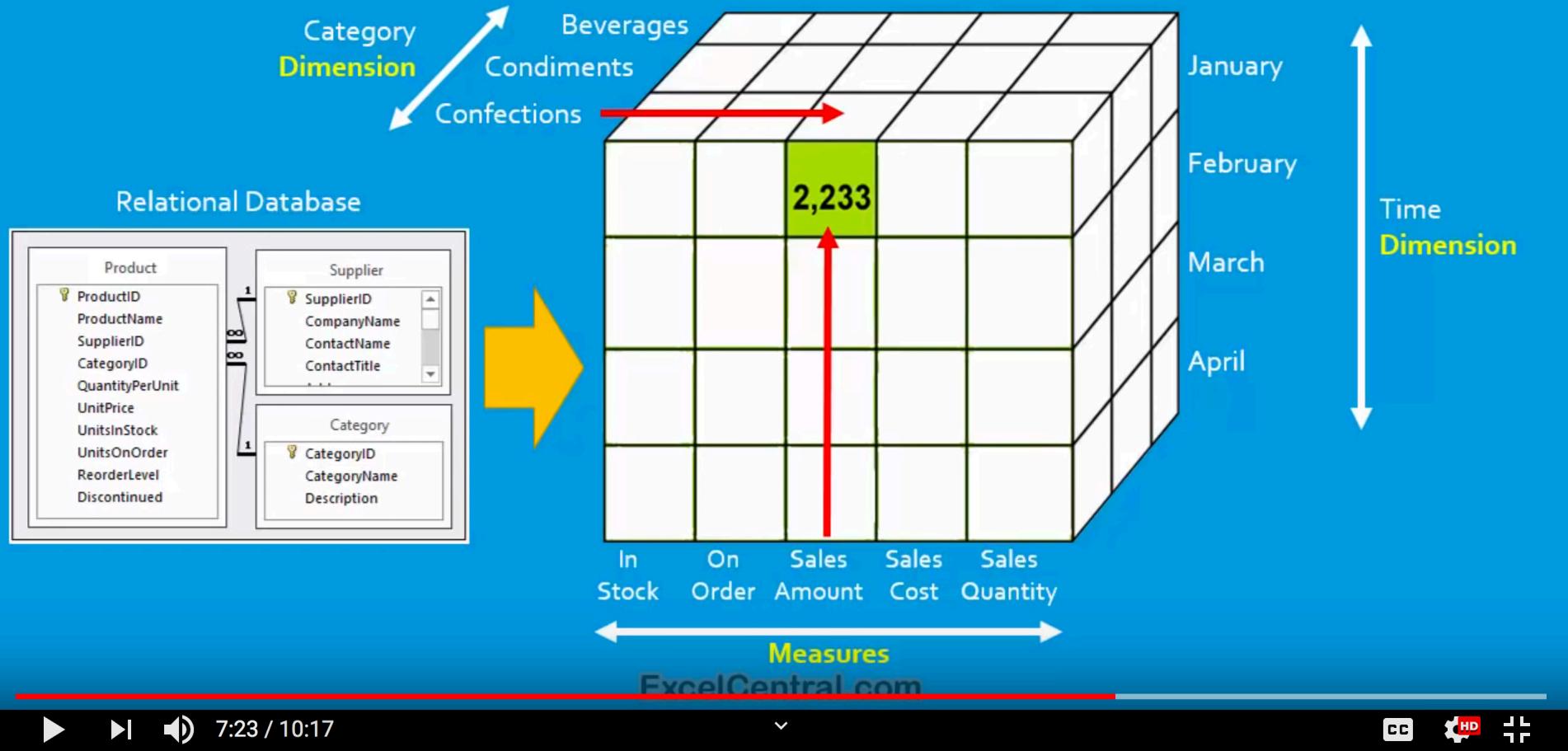
MDX (multi-dimensional expressions)



OLAP Pivot Table

	A	B	C
1			
2	Row Labels	Sum of UnitsInStock	Sum of UnitsOnOrder
3	Beverages	559	60
4	Condiments	507	170
5	Confections	386	180
6	Dairy Products	393	140
7	Grains/Cereals	308	90
8	Meat/Poultry	165	0
9	Produce	100	20
10	Seafood	701	120
11	Grand Total	3119	780

How an OLAP cube works



Introduction to Pivot Tables, Charts, and Dashboards in Excel (Part 1)

Pivot Table Areas Diagram

Revenue	Column	Q1	Q2	Q3	Q4
Row Labels	Year	2014			
Andrew Cencin	Year	300	1,590	13,920	1,740
Anne Larsen	Year	483	414	1,598	
Jan Kotas	Year	35	552	1,536	4,200
Laura Giussan	Year	1,950	450	1,136	250

Drag fields between areas below:

- FILTERS
- COLUMNS
- ROWS
- VALUES

Sheet8 Data Pivot Table Diagram

READY 6:43 / 14:47 100%

Introduction to Pivot Tables, Charts, and Dashboards in Excel (Part 1)

Sales Data for December 2014.xlsx - Excel

youtube.com befindet sich jetzt im Vollbildmodus. Vollbild beenden (esc)

FILE HOME INSERT PAGE LAYOUT FORMULAS DATA REVIEW VIEWS

From Access From Web From Text From Other Sources Existing Connections Refresh All Edit Links Filter Advanced Sort Text to Columns Flash Fill Remove Duplicates Data Validation Consolidate Relationships What-If Analysis Group Ungroup Subtotal Outline

A3 Order ID

Order Details for December 2014

Order ID	Order Date	Customer ID	Customer Name	Address	City	State	ZIP/Postal Code	Country/Region	Salesperson	Region	Shipped Date	Shipper Name	Shipper ID	Comments	
1368	12/27/14	27	Company AA	789 27th Street	Las Vegas	NV	99999	USA	Mariya Sergienko	West	12/29/14	Shipping Company B	Kar	Customer B	
1369	12/27/14	27	Company AA	789 27th Street								Shipping Company B	Kar	Customer B	
1370	12/04/14	4	Company D	123 4th Street								Shipping Company A	Chi	Customer A	
1371	12/04/14	4	Company D	123 4th Street								Shipping Company A	Chi	Customer A	
1372	12/04/14	4	Company D	123 4th Street								Shipping Company A	Chi	Customer A	
1373	12/12/14	12	Company L	123 12th Street								Shipping Company B	Joh	Customer B	
1374	12/12/14	12	Company L	123 12th Street								Shipping Company B	Joh	Customer B	
1375	12/08/14	8	Company H	123 8th Street								Shipping Company C	Eliz	Customer C	
1376	12/04/14	4	Company D	123 4th Street								Shipping Company C	Chi	Customer C	
1377	12/29/14	29	Company CC	789 29th Street								Shipping Company B	Soc	Customer B	
1378	12/03/14	3	Company C	123 3rd Street								Shipping Company B	Tho	Customer B	
1379	12/06/14	6	Company F	123 6th Street	Milwaukee	WI	99999	USA	Michael Kellper	North	12/08/14	Shipping Company B	Fra	Customer F	
1380	12/28/14	28	Company BB	789 28th Street	Memphis	TN	99999	USA	Anne Larsen	South	12/30/14	Shipping Company C	Am	Customer B	
1381	12/08/14	8	Company H	123 8th Street	Portland	OR	99999	USA	Nancy Freehafer	North	12/10/14	Shipping Company C	Eliz	Customer C	
1382	12/10/14	10	Company J	123 10th Street	Chicago	IL	99999	USA	Laura Giussani	East	12/12/14	Shipping Company B	Rol	Customer E	
1383	12/07/14	7	Company G	123 7th Street	Boise	ID	99999	USA	Nancy Freehafer	North			Mir	Customer F	
1384	12/10/14	10	Company J	123 10th Street	Chicago	IL	99999	USA	Laura Giussani	East	12/12/14	Shipping Company A	Rol	Customer E	
1385	12/10/14	10	Company J	123 10th Street	Chicago	IL	99999	USA	Laura Giussani	East	12/12/14	Shipping Company A	Rol	Customer E	
1386	12/10/14	10	Company J	123 10th Street	Chicago	IL	99999	USA	Laura Giussani	East	12/12/14	Shipping Company A	Rol	Customer E	
1387	12/11/14	11	Company K	123 11th Street	Miami	FL	99999	USA	Anne Larsen	South			Shipping Company C	Pet	Customer C

Data Sales by Rep Pivot Table Diagram

READY COUNT: 26

2:37 / 14:47

CC HD

Tabular Format:
One row of headers that describe the data.

Introduction to Pivot Tables, Charts, and Dashboards in Excel (Part 1)

Sales Data for December 2014.xlsx - Excel

Jon Acampora

PivotTable Recommended PivotTables Tables Illustrations Apps Charts Tours Reports Sparklines Filters Links

C4 : X ✓ fx 27

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
3	Order ID	Order Date	Customer ID	Customer Name	Address	City	State	ZIP/Postal Code	Country/Region	Salesperson	Region	Shipped Date	Shipper Name	Ship	
4	1368	12/27/14	27	Company AA	789 27th Street	Las Vegas	NV	99999	USA	Mariya Sergienko	West	12/29/14	Shipping Company B	Kar	
5	1369	12/27/14	27	Company AA	789 27th Street	Las Vegas	NV	99999	USA	Mariya Sergienko	West	12/29/14	Shipping Company B	Kar	
6	1370	12/04/14	40					99999	USA	Andrew Cencini	East	12/06/14	Shipping Company A	Chr	
7	1371	12/04/14	40					99999	USA	Andrew Cencini	East	12/06/14	Shipping Company A	Chr	
8	1372	12/04/14	40					99999	USA	Andrew Cencini	East	12/06/14	Shipping Company A	Chr	
9	1373	12/12/14	12					99999	USA	Mariya Sergienko	West	12/14/14	Shipping Company B	Joh	
10	1374	12/12/14	12					99999	USA	Mariya Sergienko	West	12/14/14	Shipping Company B	Joh	
11	1375	12/08/14	8					99999	USA	Nancy Freehafer	North	12/10/14	Shipping Company C	Eliz	
12	1376	12/04/14	40					99999	USA	Andrew Cencini	East	12/06/14	Shipping Company C	Chr	
13	1377	12/29/14	29					99999	USA	Jan Kotas	West	12/31/14	Shipping Company B	Soc	
14	1378	12/03/14	3					99999	USA	Mariya Sergienko	West	12/05/14	Shipping Company B	Tho	
15	1379	12/06/14	6					99999	USA	Michael Neipper	North	12/08/14	Shipping Company B	Fra	
16	1380	12/28/14	28					99999	USA	Anne Larsen	South	12/30/14	Shipping Company C	Am	
17	1381	12/08/14	8					99999	USA	Nancy Freehafer	North	12/10/14	Shipping Company C	Eliz	
18	1382	12/10/14	10					99999	USA	Laura Giussani	East	12/12/14	Shipping Company B	Rol	
19	1383	12/07/14	7					99999	USA	Nancy Freehafer	North			Min	
20	1384	12/10/14	10					99999	USA	Laura Giussani	East	12/12/14	Shipping Company A	Rol	
21	1385	12/10/14	10					99999	USA	Laura Giussani	East	12/12/14	Shipping Company A	Rol	
22	1386	12/10/14	10	Company J	123 10th Street	Chicago	IL	99999	USA	Laura Giussani	East	12/12/14	Shipping Company A	Rol	
23	1387	12/11/14	11	Company K	123 11th Street	Miami	FL	99999	USA	Anne Larsen	South		Shipping Company C	Pet	
24	1388	12/11/14	11	Company K	123 11th Street	Miami	FL	99999	USA	Anne Larsen	South		Shipping Company C	Pet	
25	1389	12/01/14	1	Company A	123 1st Street	Seattle	WA	99999	USA	Nancy Freehafer	North			Ani	

Create PivotTable ? X

Choose the data that you want to analyze

Select a table or range
Table/Range: Data!\$A\$3:\$Z\$68

Use an external data source
Choose Connection...
Connection name:

Choose where you want the PivotTable report to be placed

New Worksheet
 Existing Worksheet
Location:

Choose whether you want to analyze multiple tables

Add this data to the Data Model

OK Cancel

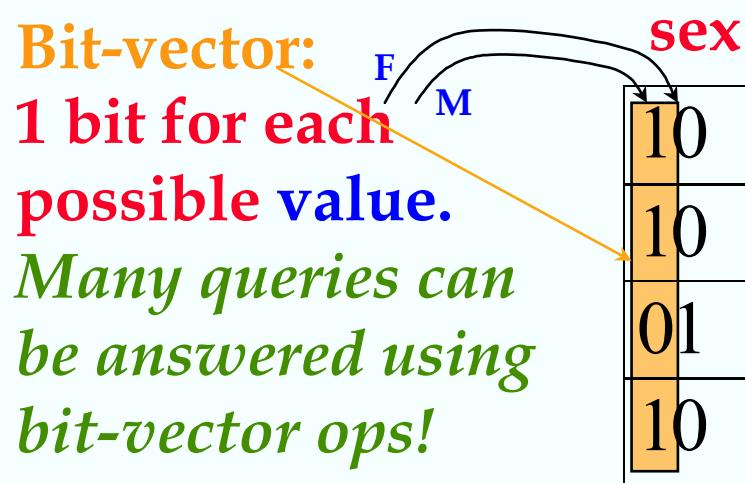
Data Sales by Rep Pivot Table Diagram +

ENTER

3:59 / 14:47

Implementation Issues

- ❖ New indexing techniques: Bitmap indexes, Join indexes, array representations, compression, precomputation of aggregations, etc.
- ❖ E.g., Bitmap index:

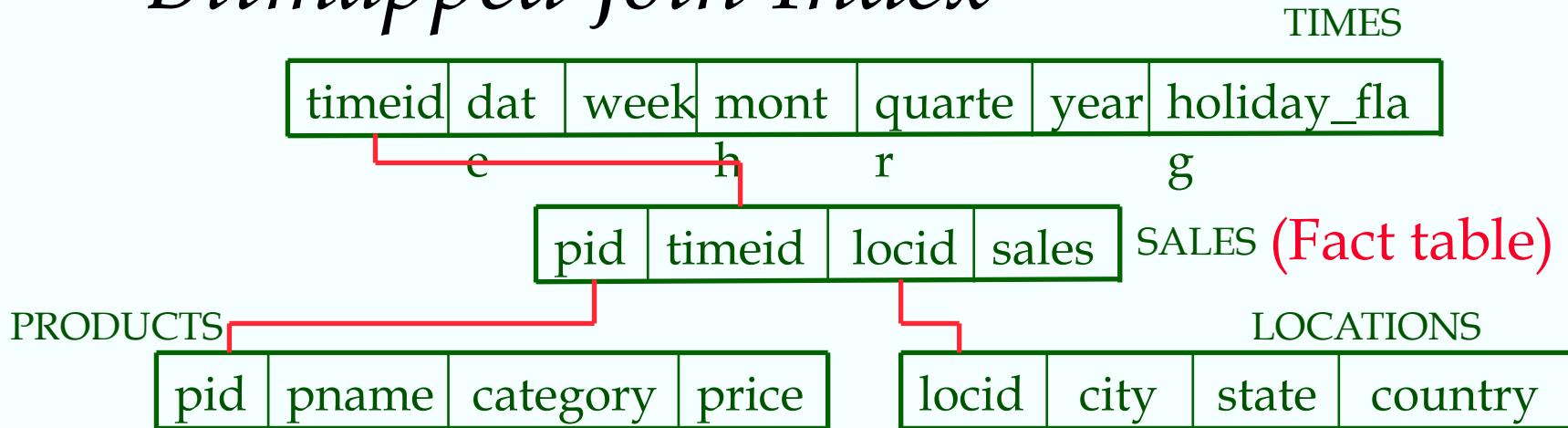


custid	name	sex	rating	rating
112	Joe	M	3	00100
115	Ram	M	5	00001
119	Sue	F	5	00001
112	Woo	M	4	00010

Join Indexes

- ❖ Consider the join of Sales, Products, Times, and Locations, possibly with additional selection conditions (e.g., $\text{country} = \text{"USA"}$).
 - A **join index** can be constructed to speed up such joins. The index contains $[s,p,t,l]$ if there are tuples (with sid) s in Sales, p in Products, t in Times and l in Locations that satisfy the join (and selection) conditions.
- ❖ **Problem:** Number of join indexes can grow rapidly.
 - A variation addresses this problem: For each column with an additional selection (e.g., country), build an index with $[c,s]$ in it if a dimension table tuple with value c in the selection column joins with a Sales tuple with sid s ; if indexes are bitmaps, called **bitmapped join index**.

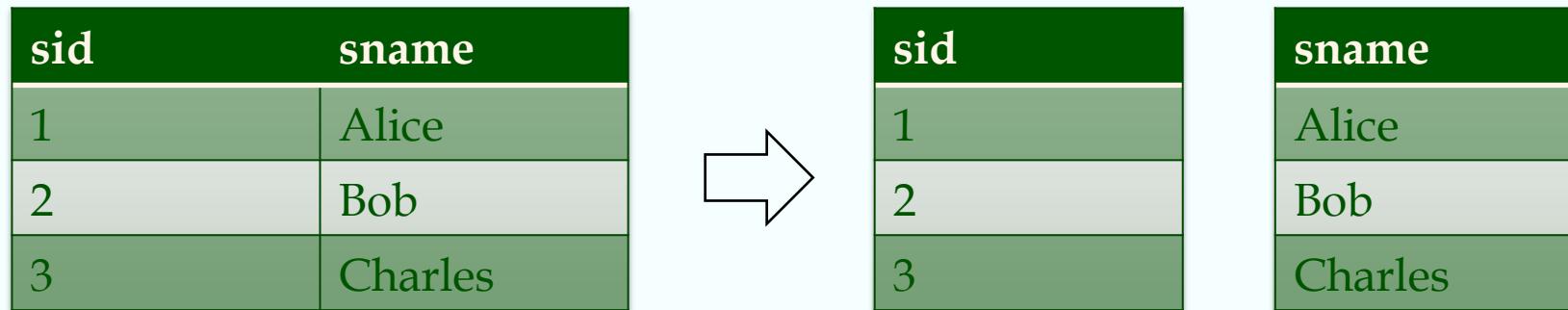
Bitmapped Join Index



- ❖ Consider a query with conditions $\text{price}=10$ and $\text{country}=\text{"USA"}$. Suppose tuple (with sid) s in Sales joins with a tuple p with $\text{price}=10$ and a tuple l with $\text{country}=\text{"USA"}$. There are two join indexes; one containing $[10,s]$ and the other $[\text{USA},s]$.
- ❖ Intersecting these indexes tells us which tuples in Sales are in the join and satisfy the given selection.

Column Stores

- ❖ ROLAP approach in which relations are stored column by column, rather than row by row.



- ❖ Order and duplicates matter! Like (lossless) vertical decomposition with implicit rowid as key.
- ❖ Optimizations: column compression usually works better than table compression.
- ❖ Simple idea, but very fashionable right now. (HP Vertica, SAP HANA, ...)

Row Stores vs Column Stores

- ❖ RDBMS are traditionally implemented as row stores. Can be implemented as column stores.
- ❖ Optimize for OLTP => row store
 - Optimizes data locality for updating. A tuple is in one place only
- ❖ Optimize for analytics => column store
 - Assuming each query uses only some of the columns => less data to scan