

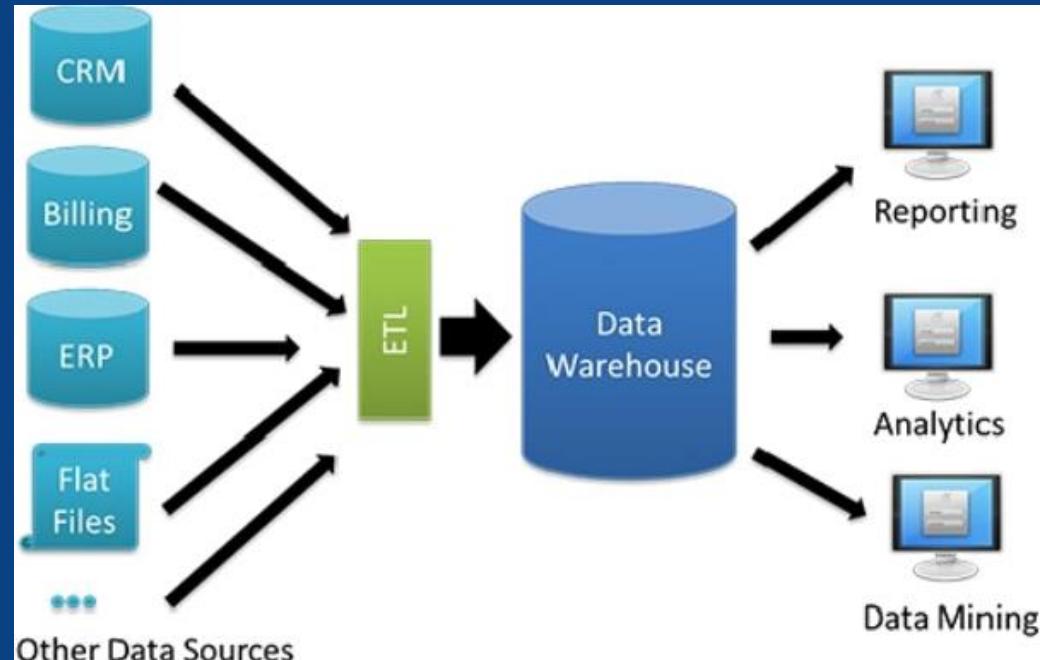


THE UNIVERSITY OF
MELBOURNE

Data Warehousing

Database Systems & Information Modelling
INFO90002

Week 9 – DW
Dr Tanya Linden
Dr Renata Borovica-Gajic
David Eccles





This Lecture Discusses

The differences between **transactional** (operational) and **informational** (dimensional) databases

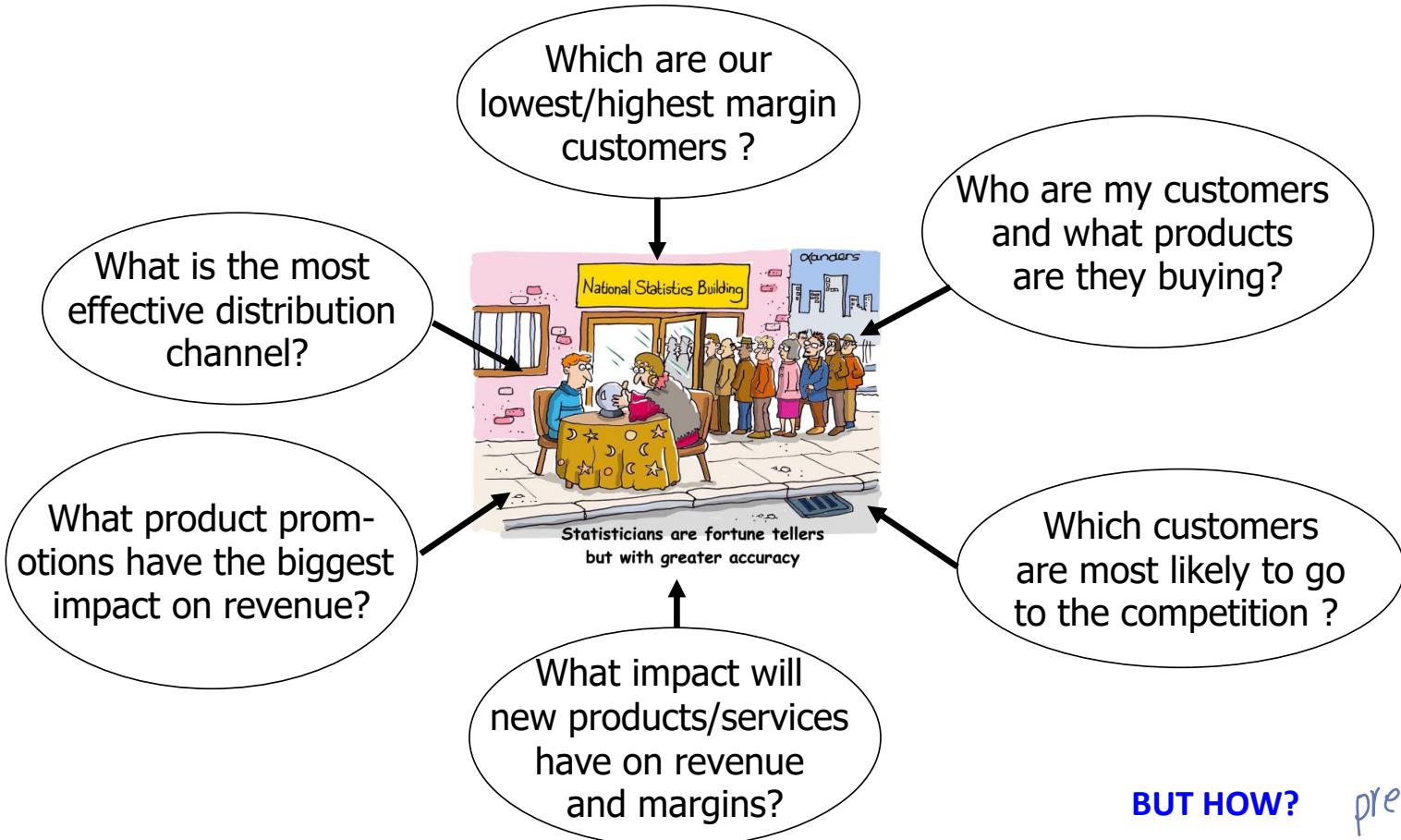
Characteristics of a Data Warehouse

Overall architecture of a Data Warehouse

Star Schemas

for management
↳ where data is group
catergorised in dataware
house

Motivations: A manager wants to know....



Relational Databases for Operational Processing

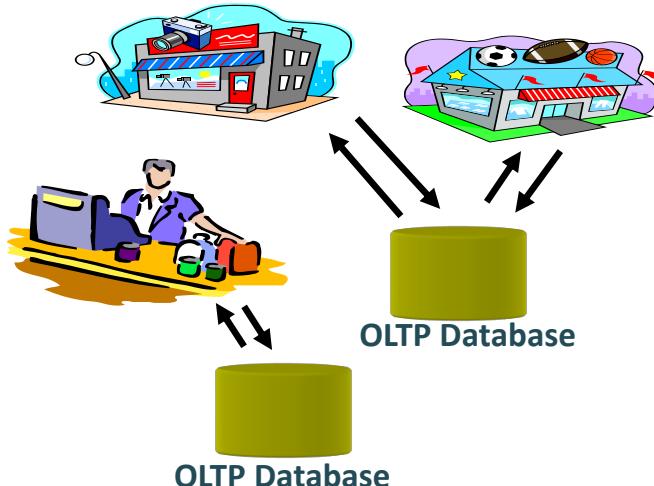
Used to run day to day business operations

Automation of routine business processes

- Accounting
- Inventory
- Purchasing
- Sales

Created huge efficiencies

We use relational database for operational processing



OLTP Databases

① OLTP = “OnLine Transaction Processing”

Transaction processing supports daily (routine, repetitive) operations

- Mundane but crucial
- Become even more important with the growth of the internet

②

Definition:

1. Collection of read/write operations

OLTP is a collection of read and write operations.

2. Processed as one unit

3. Reliably and efficiently processed

4. No data loss due to interference and failures (operating system, program, disk, ...)

OLTP Data Characteristics

3

Characteristics of data:

- Transaction oriented

It's transaction oriented, so it is data manipulated language.

- DML

4

Inserts/Updates/Deletes

- May be inconsistent and incomplete

- Data may not be in its final form

- Different entry standards and formats
mm/dd/yyyy vs dd/mm/yyyy
- Missing / inaccurate / incomplete
Online customers accidentally or intentionally enter inaccurate details

4.2

- Volatile – continually changing (ex: person change the surname)

- Data maybe subject to change

4.3

- Current

- Data related to the operation of the business TODAY!

companies based on their needs might use different types of database.

Databases are great, BUT ...

Too many of them

- Everybody wanted one, or two, or more
- Production, Marketing, Sales, Accounting ...

Everybody got what was best for them

- IBM, Oracle, Access, Microsoft

Eventually this re-created the problem databases were meant to solve

- ① Duplicated data → the same person could be in different database, when person update the data, the data should be propagated across all different database.
(ex: change surname)
?
- ② Inaccessible data
- ③ Inconsistent data



But data is useful for analysis and decision making

some company are still use multiple system, because it's more interoperable now.
And each company need single point of truth, a way to aggregate



all data from entire organisations. from CRM, CSM from every system they're using.

What can be done about it? SPOT!

Need an integrated way of getting the ENTIRE organisational data

Need an Informational Database, rather than a Transactional Database

- * A single database that allows *all* of the organisations' data to be stored in a form that can be used to support **organisational decision processes**

So, information database is a centralised repository for decision making.

A centralised repository for decision making

- Populated from operational databases and external data sources
- Integrated and transformed data
- Optimised for reporting

And it's optimised for reporting, this mean we don't need every individual transaction. We are happy to take transactions about

this product for on daily or weekly base. We don't need every minutes what happened with this product.

CRM: customer relationship management

SCM: supply chain management

transaction database that integrity is also important. Database need to be normalized.
And consistency is also important.

OLTP Versus OLAP Systems

OLTP = “OnLine Transaction Processing” (Transactional, Operational)

- Maintaining data integrity and effectiveness while dealing with numerous transactions simultaneously

OLAP = Online analytical processing (Informational, Dimensional) analytical database are different, it's read only. There is not writing data, because it's historical data, we shouldn't change it. So, the query are much complex. And quite often it's combination of selection and aggregation because we want to know total across certain period of time. to know like is the product popular or not,

- Purpose-built, optimised for handling data analysis queries
- Deals with fewer, but complex queries
- SELECT and aggregate queries

OLTP and OLAP complement each other.



So one data cannot exist without

another, because analytical database comes from transaction database.



Data Warehouse: An Informational Database

Data Warehouse:

- ① • A single repository of historic organisational data
 - ② • Integrates data from multiple sources
 - Extracts data from source systems, transforms, loads into the warehouse
 - ③ • Makes data available to managers/users
 - ④ • Supports analysis and decision-making
 - ⑤ • Read only
 - ⑥ • Involve a large data store (often several Terabytes, Petabytes of data)
- data warehouse is an analytical database, it's a single repository of historic organisational data.
the data might come from multiple sources, and grouped by:
If it's a global appear, different level of granularity will be available.

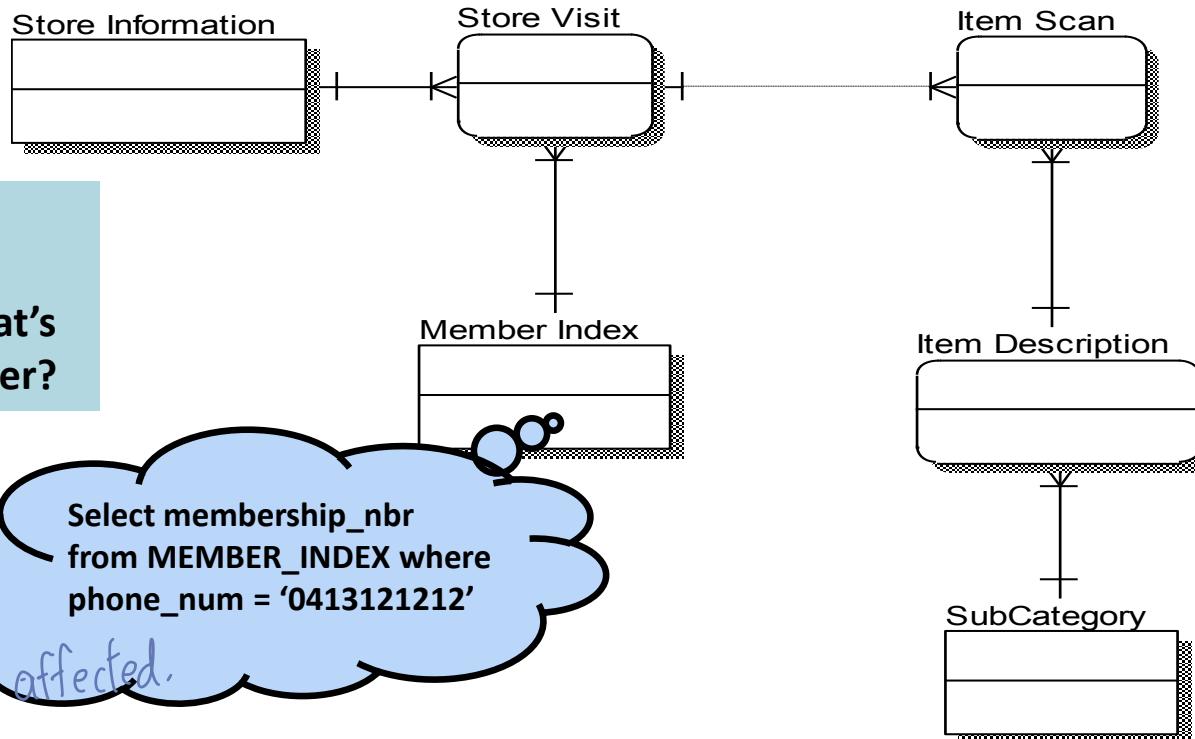
Difference between Transactional and Informational Systems

Characteristic	Transactional	Informational
Primary Purpose	Run the day-to-day business	Support decision making
Type of Data	Current data – representing the state of the business	Historical data – snapshots and predictions
Primary Users	Supports thousands/millions of users Customers, clerks and other employees	Supports hundreds of users Managers, analysts
Scope of Usage	Narrow, planned, fixed interfaces Fast and effective query processing and ensuring data integrity in multi-access environments	Broad, ad hoc, complex interfaces Emphasis on the response time to executing complex queries on large amounts of historical data aggregated from many rows
Design Goal	Performance and availability Data is stored in 3NF third normalization	Flexible use and data accessibility Data is denormalized to improve query performance
Volume	Many constant updates and queries on a few tables or rows	Periodic batch updates, complex querying on multiple or all rows <i>read only</i>

Transactional (Operational) Questions



Customer Service:
Help! I forgot my
membership card! What's
my membership number?

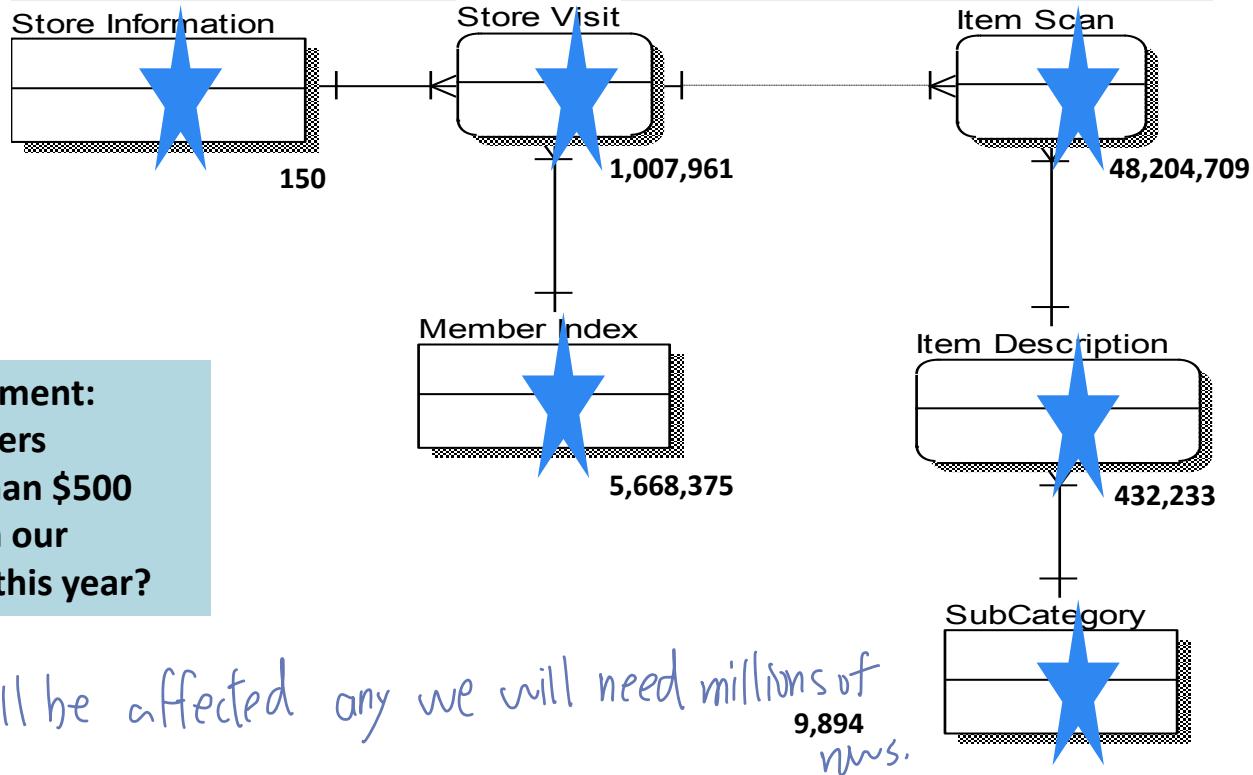


It's one query from one
table, and one table is affected.

How many tables affected? 1

How many rows have to be accessed? 1 (with an index in place)

Analytical Questions



 **Campaign Management:**
How many customers
purchased more than \$500
worth of alcohol in our
Melbourne stores this year?

many tables will be affected any we will need millions of
rows.

How many tables affected? 6
How many rows? millions

DW Supports Analytical Queries

A manager may be interested in **numerical aggregations**

- How **many**?
- What is the **average**?
- What is the **total cost**?

A manager may be interested in understanding **dimensions**

- Sales **by state by customer type**
- Sales **by product by store by quarter**

★ And there always have time dimensions.

Data Warehouse will help answer these questions

when we look into dimensions and numerical aggregations

They definitely have validated and integrated data. They move anything that is not complete.

↳ They deal in some ways with incomplete data.

Characteristics of a DW

1 Subject oriented

- Data warehouses are organised around particular subjects (sales, customers, products)

2 Validated, Integrated data

- Data from different systems converted to a common format: allows comparison and consolidation of data from different sources
- Data from various sources validated before storing it in a data warehouse

Problems

3

Incomplete Errors

- Missing Fields
- Records or Fields that, by design, are not recorded, e.g. the type of people that buy Big Issue from Big Issue Vendors when a sale is made

4

incorrect Errors

- Wrong data entered into source system
 - E.g. manual entering of data will always have a percentage of incorrect data → human error

Characteristics of a DW

5

Time variant

- Historical data
- Trend analysis crucial for decision support: requires historical data
- Data consists of a series of “snapshots” which are time stamped

6

Non-volatile

- Users have Read access only – all updating done automatically by ETL* process and periodically by a DBA

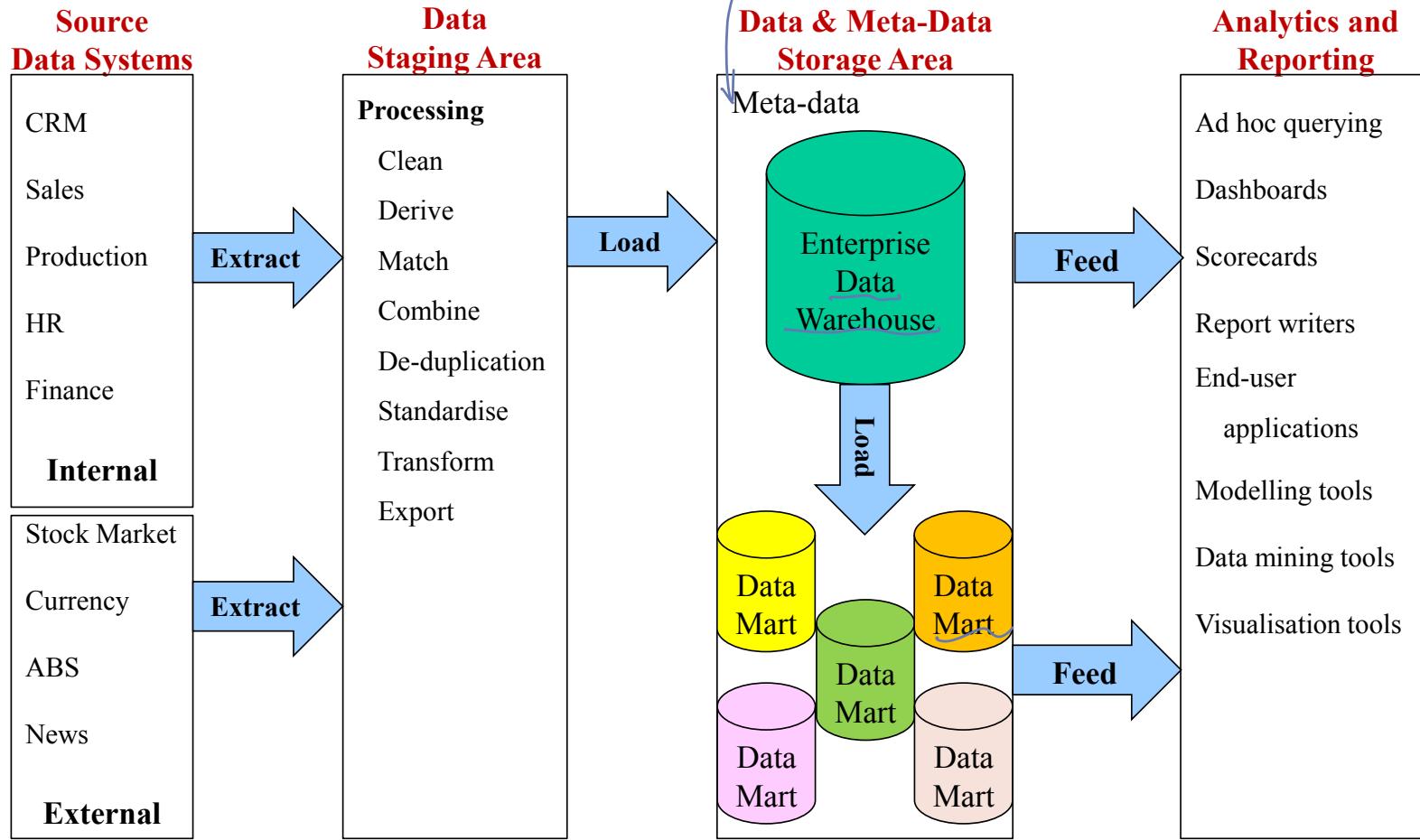
The whole idea of data warehouse is to look at the data patterns across certain periods of time.

↳ they're looking into a group. No one keeps tracking every transaction individually, transactions are grouped (like snapshots with time stamp)

Updates are done by software, not here.

*ETL stands for “extract, transform, load” - the three processes that, in combination, move data from one or more databases, or other sources to a unified repository, typically a data warehouse

A DW Architecture



Data marts and data mining

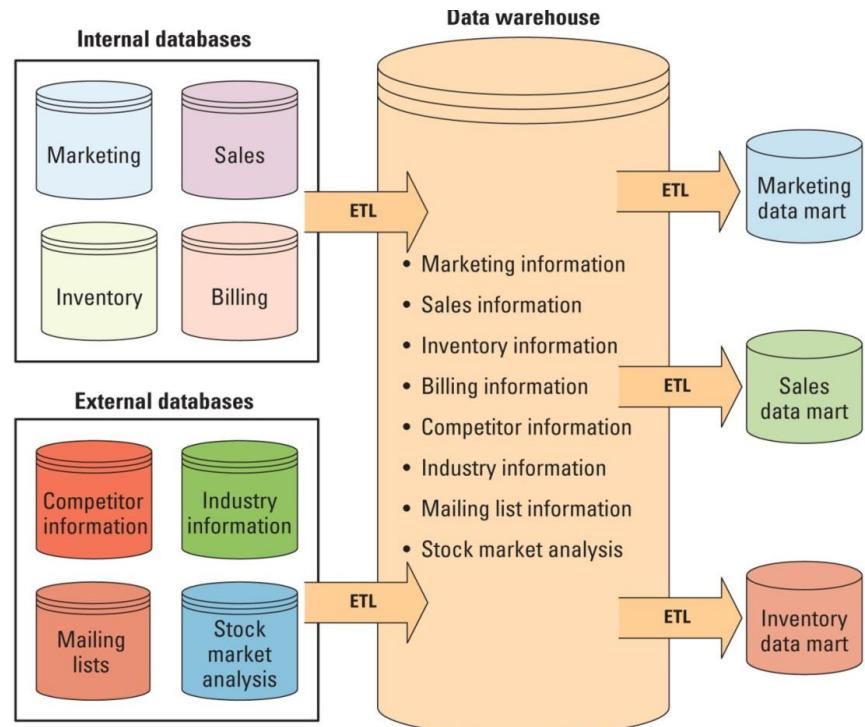
Data mart

- contains a subset of data warehouse information

when you work on data mart, the actual process when algorithms are applied to uncover things, it's called data mining.

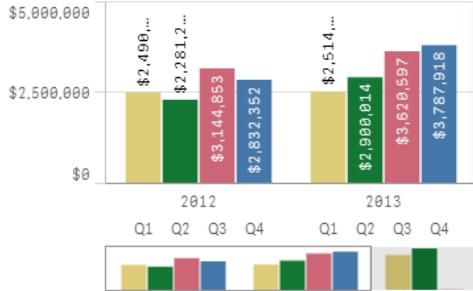
Data-mining

- A process in which algorithms are applied to information to uncover patterns and relationships which otherwise are difficult to find



Business Intelligence Dashboard

Total Sales = \$31,314.1K

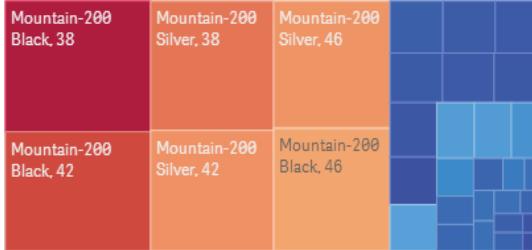


Profit Margin

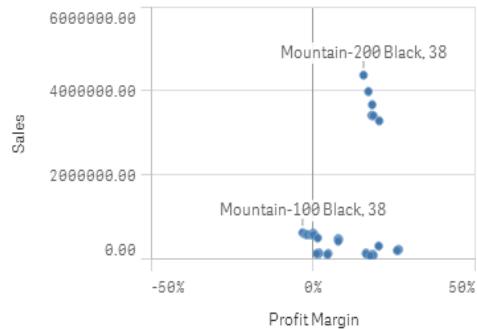


Sales

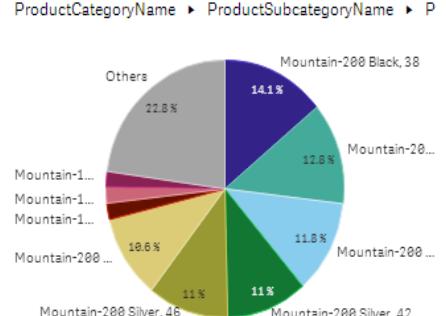
* red = most ordered



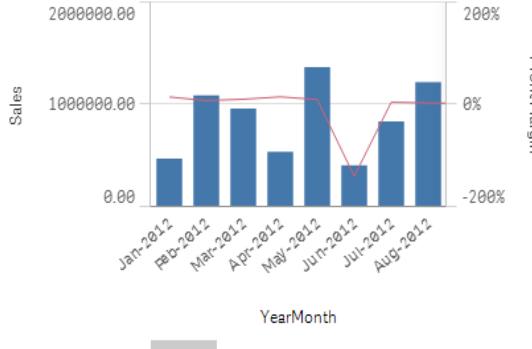
Sales vs Profit Margin

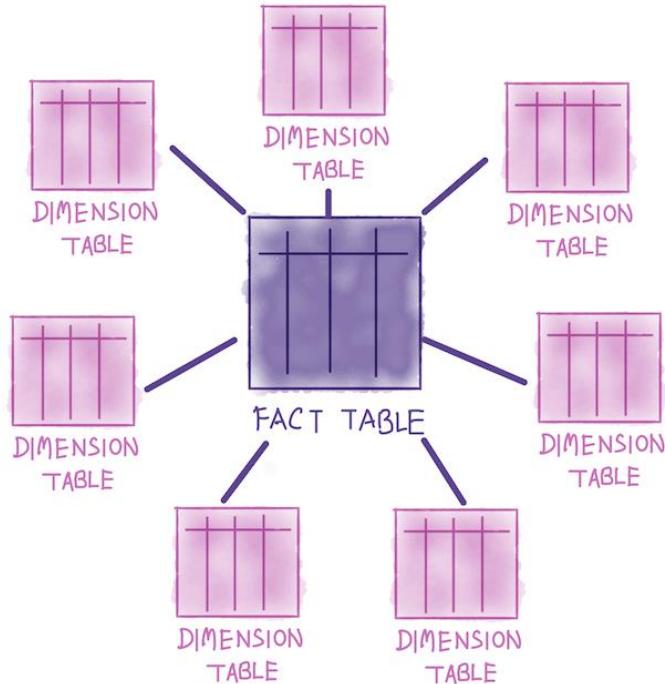


% of Total Sales



Sales and Profit Margin by Year-Month





Dimensional Modelling

Introduction to Dimensional Modelling

Popularised by Ralph Kimball in the 1990s

Based on the *multi-dimensional* model of data and designed for retrieval-only databases

Very simple, intuitive, and easily-understood structure

Also known as star schema design

A dimensional model consists of:

- Fact table
- Several dimensional tables
- (Sometimes) hierarchies in the dimensions



the central is the fact table.
Fact is what we need to analyse to know.

Dimension is what contributes to our facts.

Essentially a simple and restricted type of ER model

Business Analyst World

Fact Dimension
How much **revenue** did the **product G** generate in the last three months, broken down by

Dimension
month for the south eastern sales **region**, by individual **stores**, broken down by

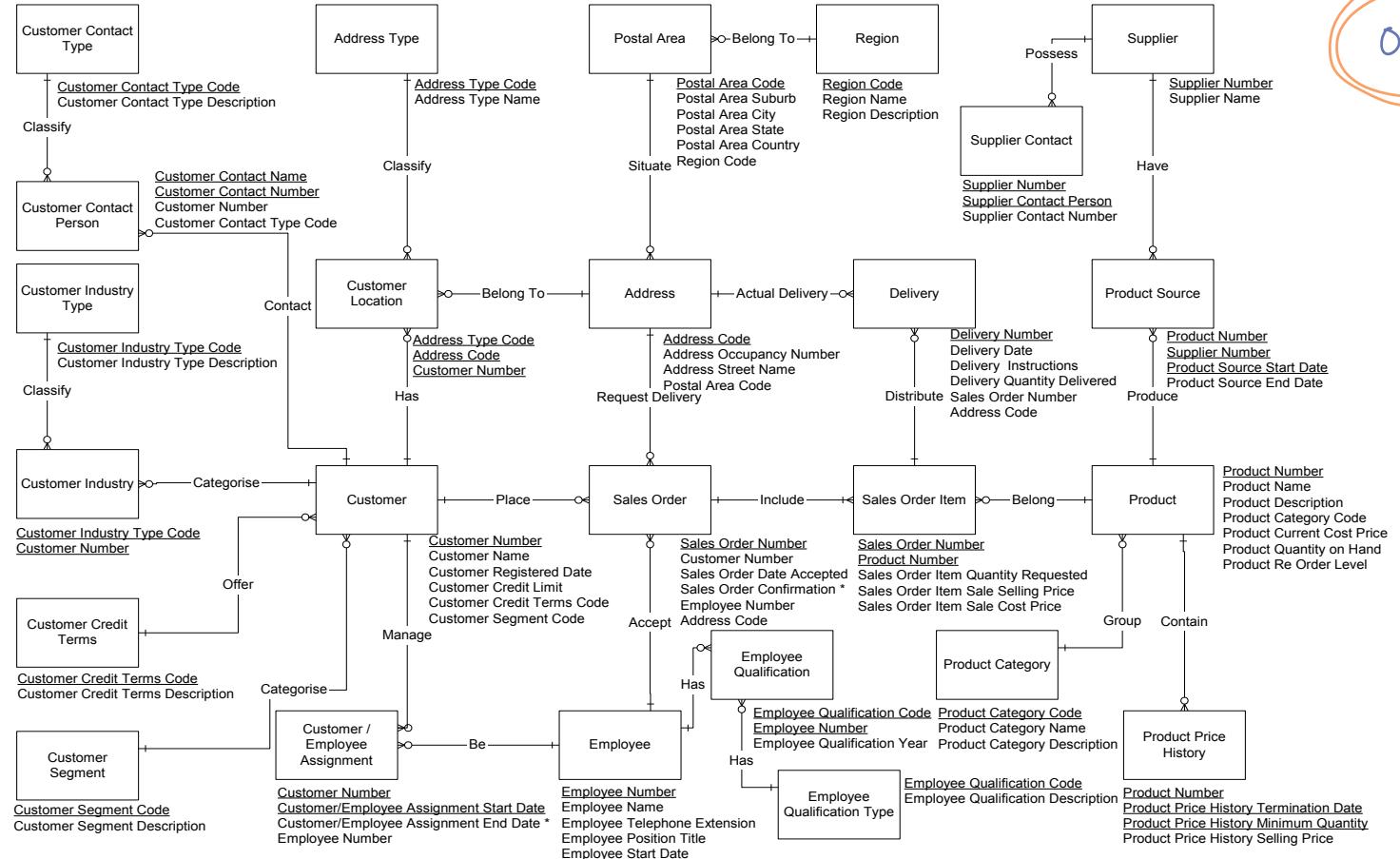
Dimension
promotions, compared to estimates and to the previous version of the product

- Analysis starts usually with a single indication of something strange, then goes deep into the data, left to a new dimension, right to another, up to the summary, back down and left and right again, until the problem is identified...

- Dimensional Analysis: To support business analysts view
 - Revenue per product per customer per location?



Example ER model



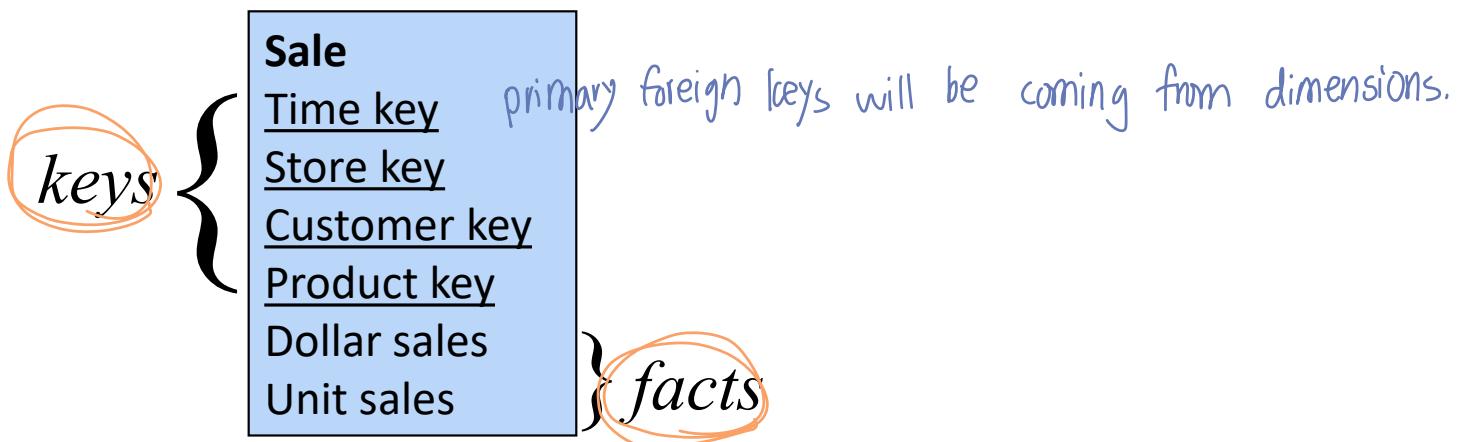
Level 1

OLTP

Fact Table

A fact table contains the actual business measures (additive, aggregates), called ***facts***

The fact table also contains *foreign keys* pointing to ***dimensions***



Fact Table - example

Actual data might look like this

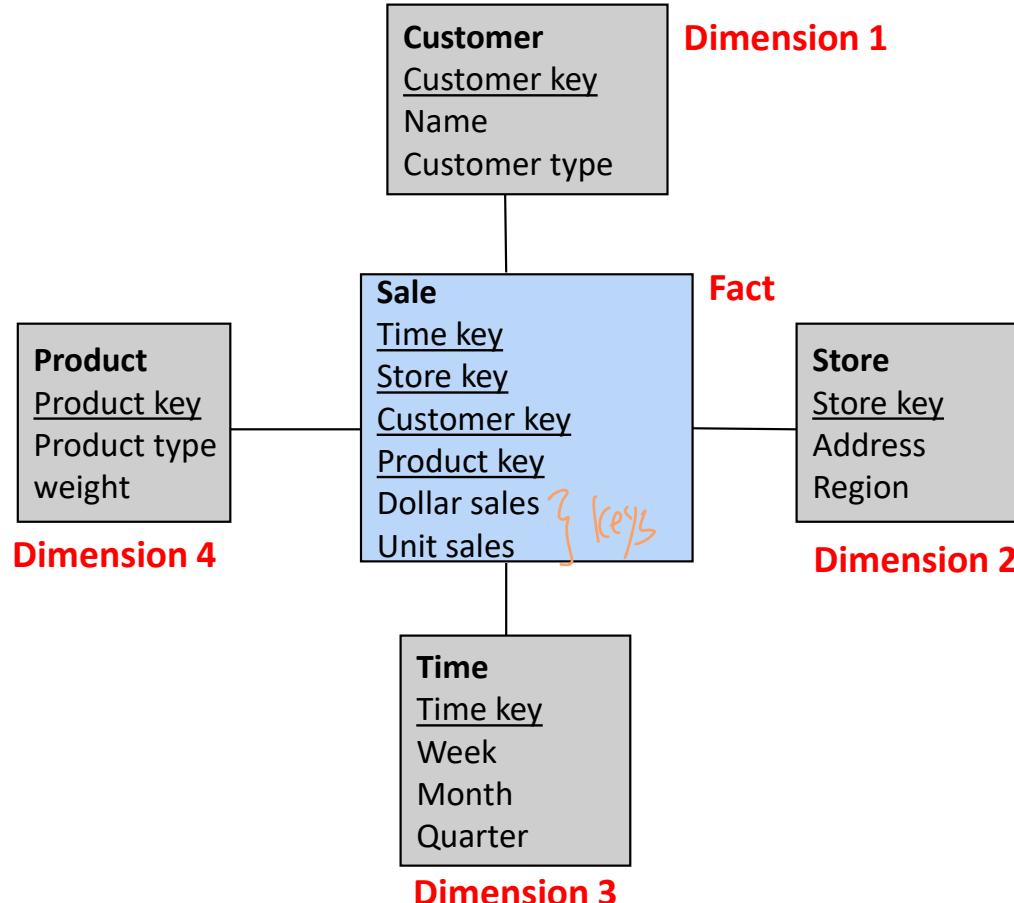
Granularity, or level of detail, is a key issue

- Finest level of detail for a fact table, determined by the finest level of each dimension

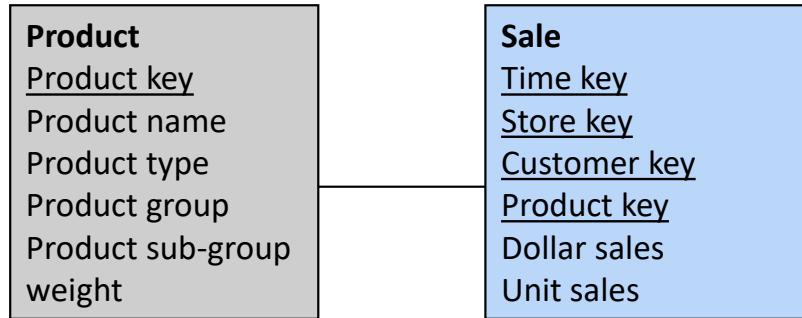
there are time stamp.

<i>Time-id</i>	<i>Store-id</i>	<i>Cust-id</i>	<i>Prod-id</i>	<i>Dollar sales</i>	<i>Unit Sales</i>
T100	S303	C101	P98	\$120,000	5,000
T101	S303	C256	P98	\$240000	10,000
T102	S387	C101	P10	\$456,000	27,899
T100	S234	C400	P56	\$100,200	5,600

Star schema – dimensional model



Dimension Hierarchies



- | | |
|----------------------|---------------|
| Product name | e.g. Hammer |
| - Product type | e.g. Tool |
| - Product group | e.g. Hardware |

Dimension Table - example

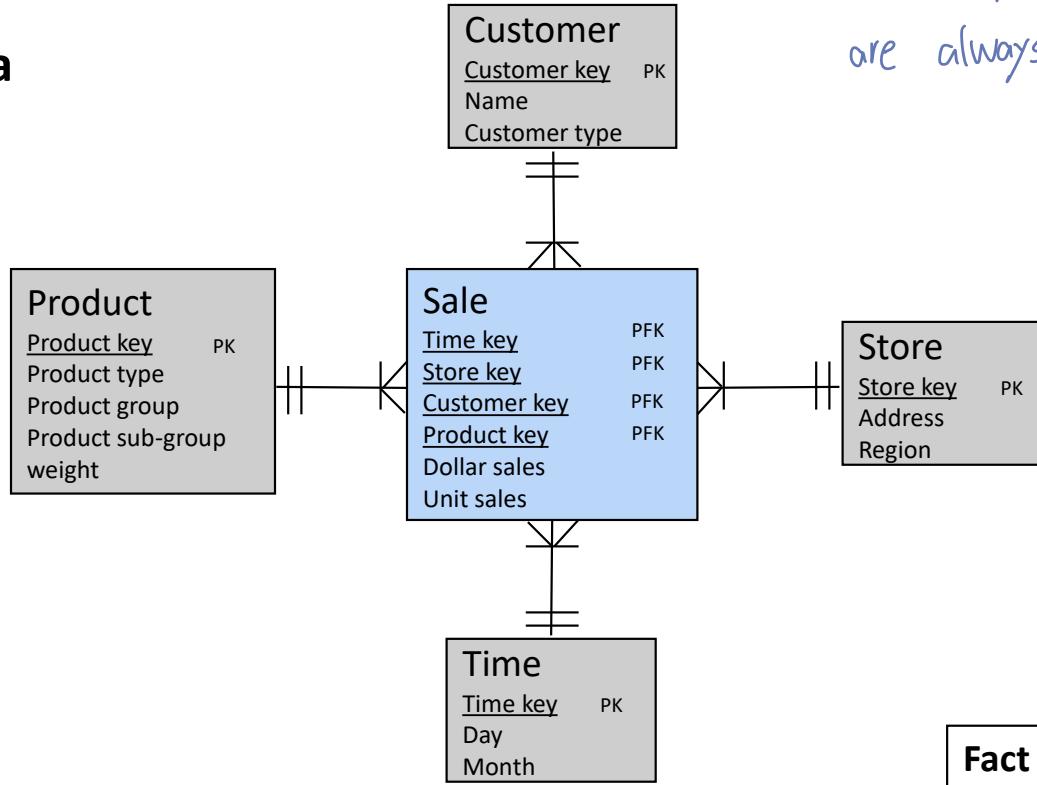
Actual data might look like this

Hierarchy evident in data

<i>Prod-id</i>	<i>Prod-Name</i>	<i>Prod-Group</i>	<i>Prod-Subgroup</i>	<i>Weight</i>
P10	Hammer	Hardware	Tool	5kg
P56	10cm Nails	Hardware	Nails	1kg
P98	Plastic Pipe	Plumbing	Pipe	1kg

Dimensional model as an ER model

Star schema



relationships between dimensions and facts are always one to many.

Fact table is an intersection table

Designing a Dimensional Model



Steps:

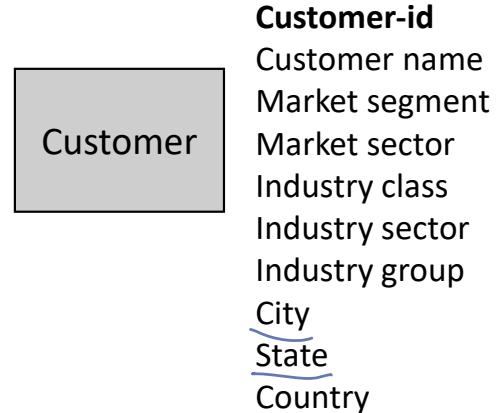
1. Choose a Business Process what the manager could need.
2. Choose the measured facts (usually numeric, additive quantities)
3. Choose the granularity of the fact table (usually numeric quantities or dollars)
4. Choose the dimensions and whether your dimensions table will have sub dimensions.
decide the granularity in the fact table
5. Complete the dimension tables

(Kimball, 1996)

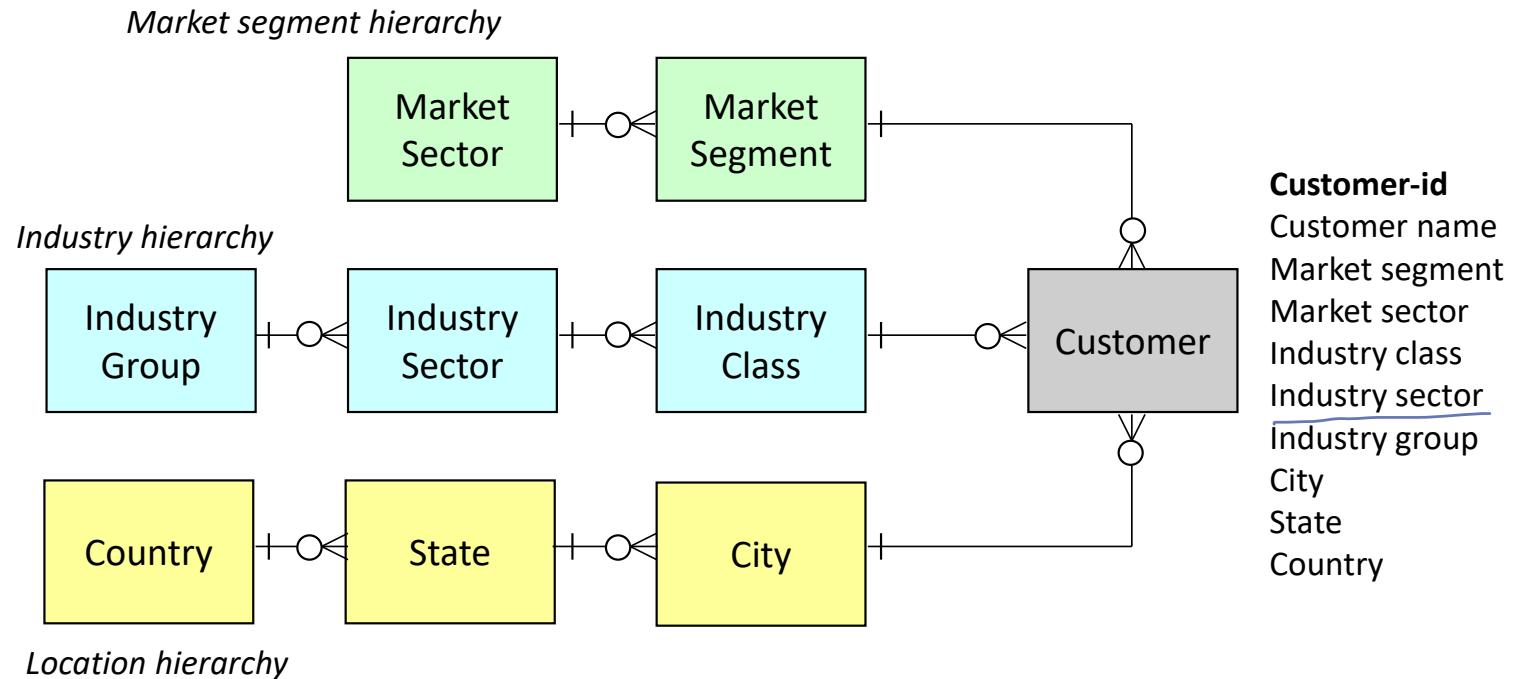


Embedded Hierarchies in Dimensional Tables

sub-dimension

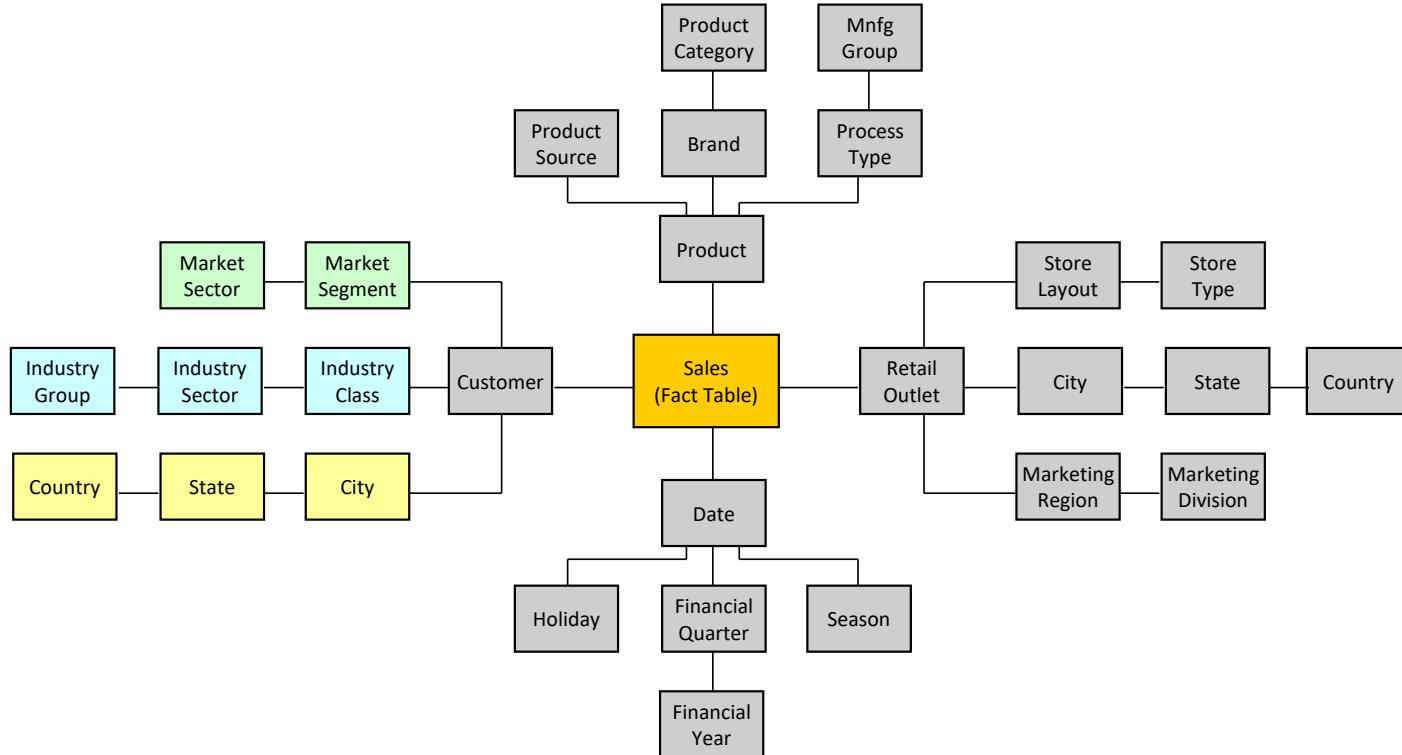


Embedded Hierarchies in Dimensional Tables



Snowflake Schema: hierarchy in dimensions

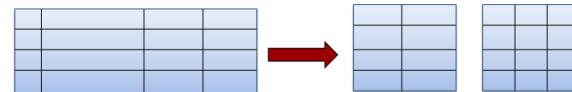
end up in real life.



Design Outcomes: Normalised or Denormalised?

Normalisation

- Eliminates redundancy
- Storage efficiency
- Referential Integrity



Denormalisation

- Fewer tables (fewer joins)
- Fast querying
- Design is tuned for end-user analysis

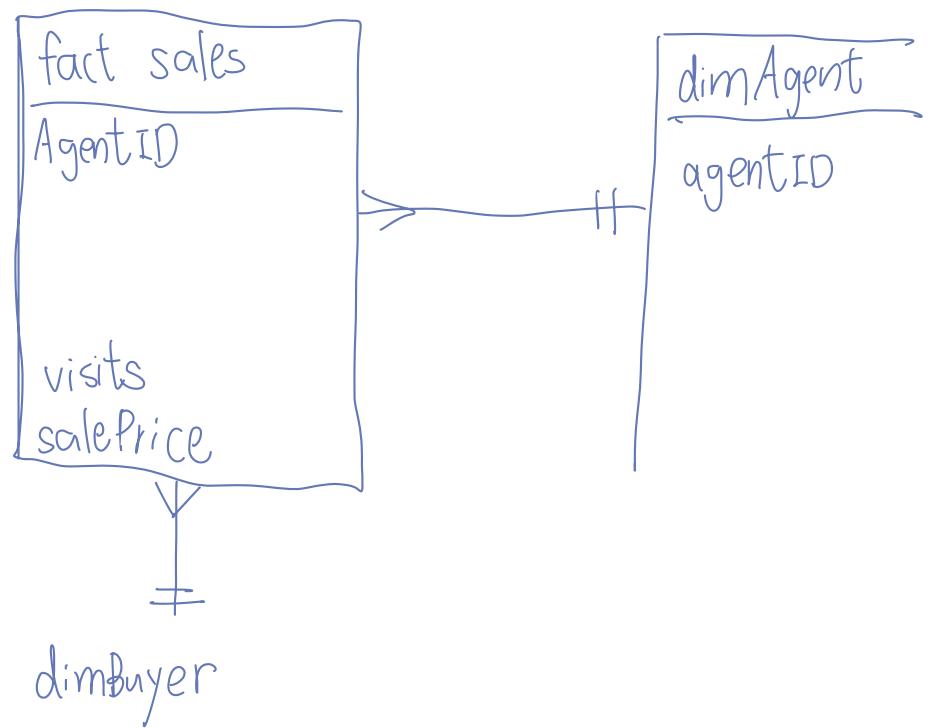




Exercise

We are making a data warehouse for a real estate agency. The company wants to track information about the **selling** of their properties. This warehouse keeps information about the **agents** (license#, first name, last name, phone #), **buyers** that come in (buyer id, first name, last name, phone #), and **property** (property#, property address, price). The information managers want to be able to find is **the number of times a property is viewed, sales price**. The information needs to be broken down **by rental agent, by buyer, by property and for different time** (day, week, month, quarter and year).

Draw a star schema to support the design of this data warehouse.





What's examinable

- Differences between transactional and informational databases
- Modelling a star schema
- Identifying the best grain level
- Defining facts and dimension tables



THE UNIVERSITY OF
MELBOURNE

Thank you