

A solid orange vertical bar is positioned on the left side of the slide.

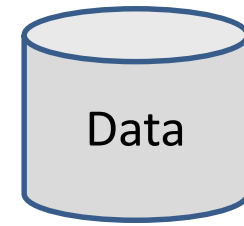
# ***Business Intelligence and Data Warehouse***

*Knut Hinkelmann*



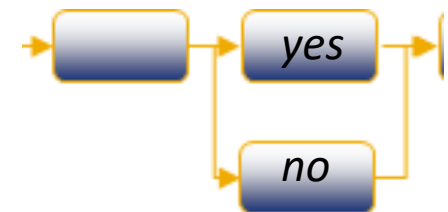
# Business Intelligence – Definition(s)

- *Sabherwal (2011)*: «We define BI as providing decision makers with valuable information and knowledge by leveraging a variety of sources of data as well as structured and unstructured information. [...] The key intellectual output of BI is **knowledge that enables decision making with information and data being the inputs.**»
- *Howson (2007)*: Business Intelligence allows people at all levels of an organisation to **access, interact with and analyse data to manage the business, improve performance, discover opportunities, and operate efficiently.**

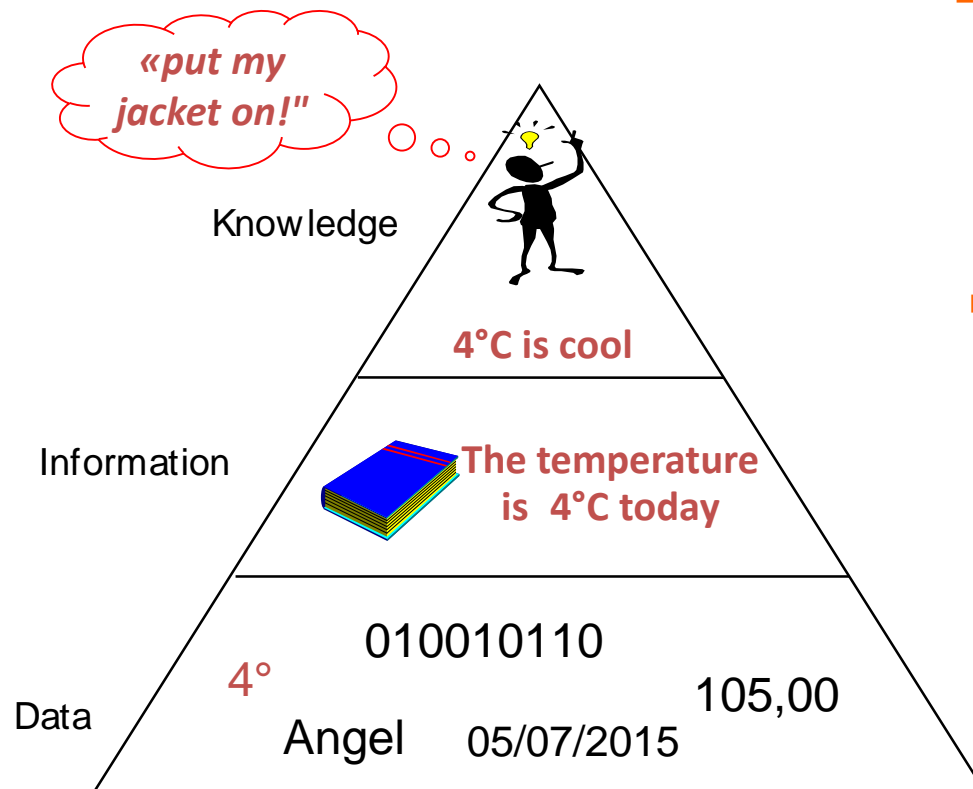


Alpha Corporation  
Sales in EUR

	'10	'11	ΔPY
Germany	84	87	+3
Austria	19	17	-2
France	28	27	-1
Rest	36	39	+3
Europe	167	170	+3



# Data, information and knowledge



- **Knowledge** enables decisions and actions
  - originates from messages (information), experience, insight
  - is embedded into the beliefs and opinions of its owner
- **Information** is an interpretation of data, often assembled in messages
  - influences the judgment and behaviour of the recipient and
  - that has a significance (relevance, purpose)
- **Data** is a set of facts and/or signals
  - Do not have meaning by itself
  - To understand data you need an interpretation

# BI overview

*Questions*

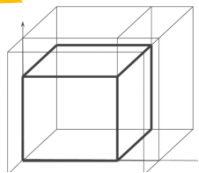
- strategic**
- What are our goals?
  - Are we reaching our goals?
  - If not, where is the problem?

- operative**
- Which customers are interested in the new product?
  - Which department has how many efforts that cannot be booked?

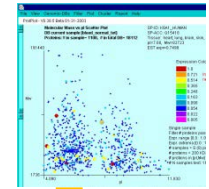
*Analyses*



*measure, aggregate, visualise*

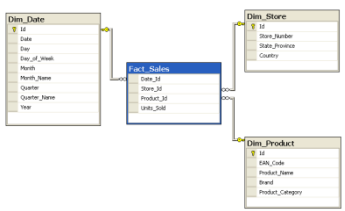


*Ad hoc queries, OLAP*



*find patterns (data mining)*

*dimensional modelling*

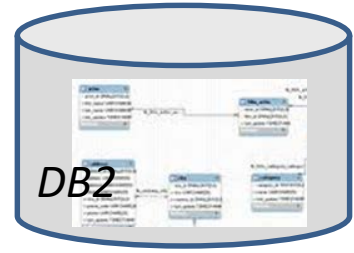
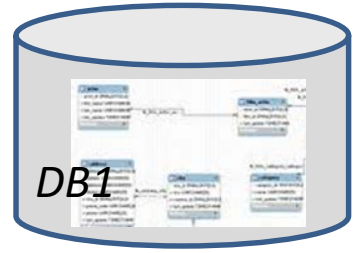


**ETL**

**ETL**

**IE**

*raw data*



# Analytic vs. transaction processing

- BI focuses on **analytic** processing instead of **transaction** processing
  - ◆ transaction processing supports *execution of core business processes*
    - *use knowledge*
  - ◆ analytic processing supports *insight and decision-making*
    - *create knowledge*



# Why introduce BI? – primary motivations

## ■ Drive company strategy

- ◆ being able to connect strategising/planning to measuring of impact (do not manage «blindly»)

## ■ Growth and competitiveness:

- ◆ anticipate market trends and adapt R&D accordingly
- ◆ better customer relationships through better-targeted offers
- ◆ better leverage of customer potential (cross-/up-selling)
- ◆ optimise business processes

## ■ Single point of truth

- ◆ no by-pass reporting, consistent data

## ■ Cost reduction

- ◆ faster access to information
- ◆ automation of reports, self-service BI
- ◆ no interference of analytics with operational systems

business  
drivers

technical  
drivers



# Decision making

- **Decision making** = *The action of selecting among alternatives to achieve a goal*
  - ◆ each alternative leads to a different future
  - ◆ what is needed is the ability to predict the futures
- **Options:**
  1. predict based on gut feeling
    - cheap in the first place
    - risk of low-quality decisions
  2. Experiment with real system (try out)
    - too risky
    - too time-consuming (only one set of conditions at a time)
  3. Predict based on the past:
    - Data collection is time-consuming
    - difficult to determine when to stop and make a decision
    - Too little or too much information



# Perspectives on BI – pain points

## MARKETING

*For targeted campaigns, we would urgently need a data basis that is harmonised with sales [...] ideally on an **integrated platform** where we can communicate with sales.*

## MANAGEMENT

*I told my people that I wanted to retrieve some numbers myself from my laptop. I then got **access to various (!) systems** [...] I finally gave up and now have an employee who does nothing but **create reports** for me [...]*

## SALES

*In most review meetings, we spend half the time discussing **which figures are the right ones** because everyone brings their own reporting. I have the impression that for any key figure it is possible to produce any value from the raw data.*

## ADVISORY BOARD

*Why weren't you able to **preview that trend**? All our competitors seem to have reacted long before we did!*



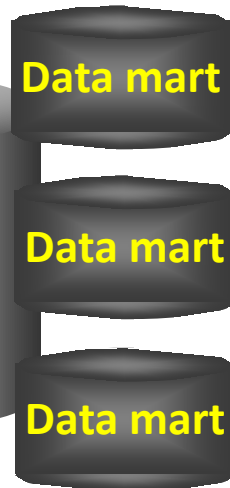
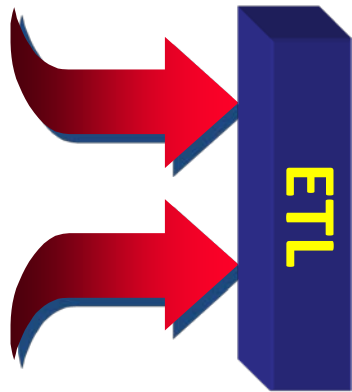


# Business Intelligence

## Data Sources



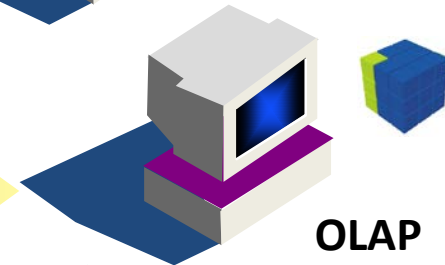
Operational data



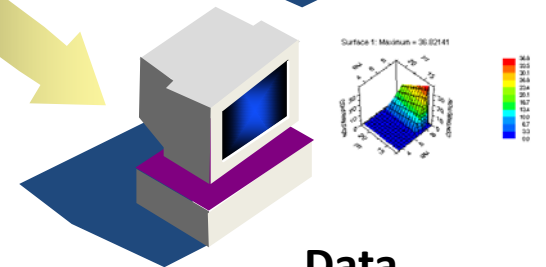
## Analysis and Use



Query & Reporting



OLAP



Data Mining



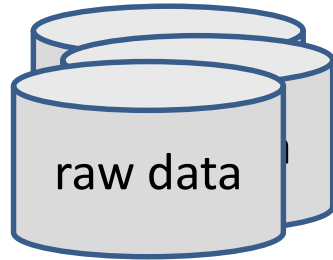
# Drivers for BI

# BI and fact-based decision making

- Fact-based decisions are based on information
- BI supports decision making by providing that information, usually in the following way:
  - ◆ the human decision maker (HDM) formulates the decision problem
  - ◆ the HDM identifies the questions that need to be answered in order to take an informed decision
  - ◆ the HDM consults a BI tool to get the answers, usually by querying or browsing (e.g. OLAP)
  - ◆ the HDM uses the answers to take an informed decision



# Data-driven vs. business-driven BI



*We have data.  
What can we do  
with it?*

ETL



Consolidate and  
integrate data



Analyze  
data

*We have questions.  
How do we get  
answers?*



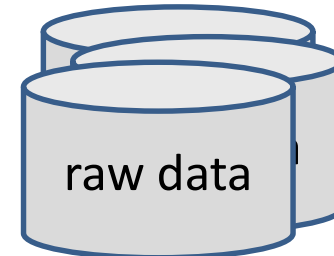
What data  
do we need to  
answer the  
questions?

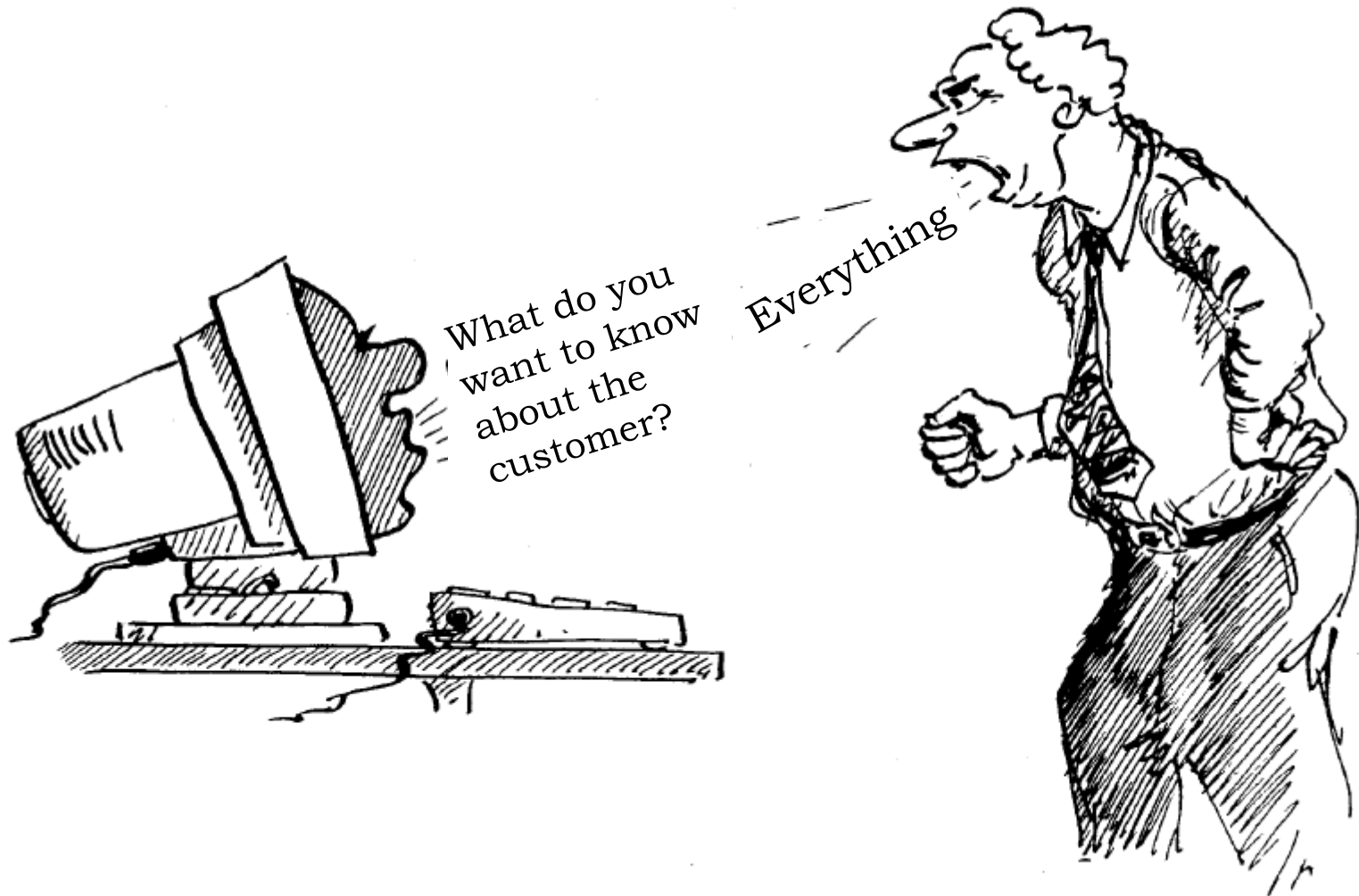


ETL



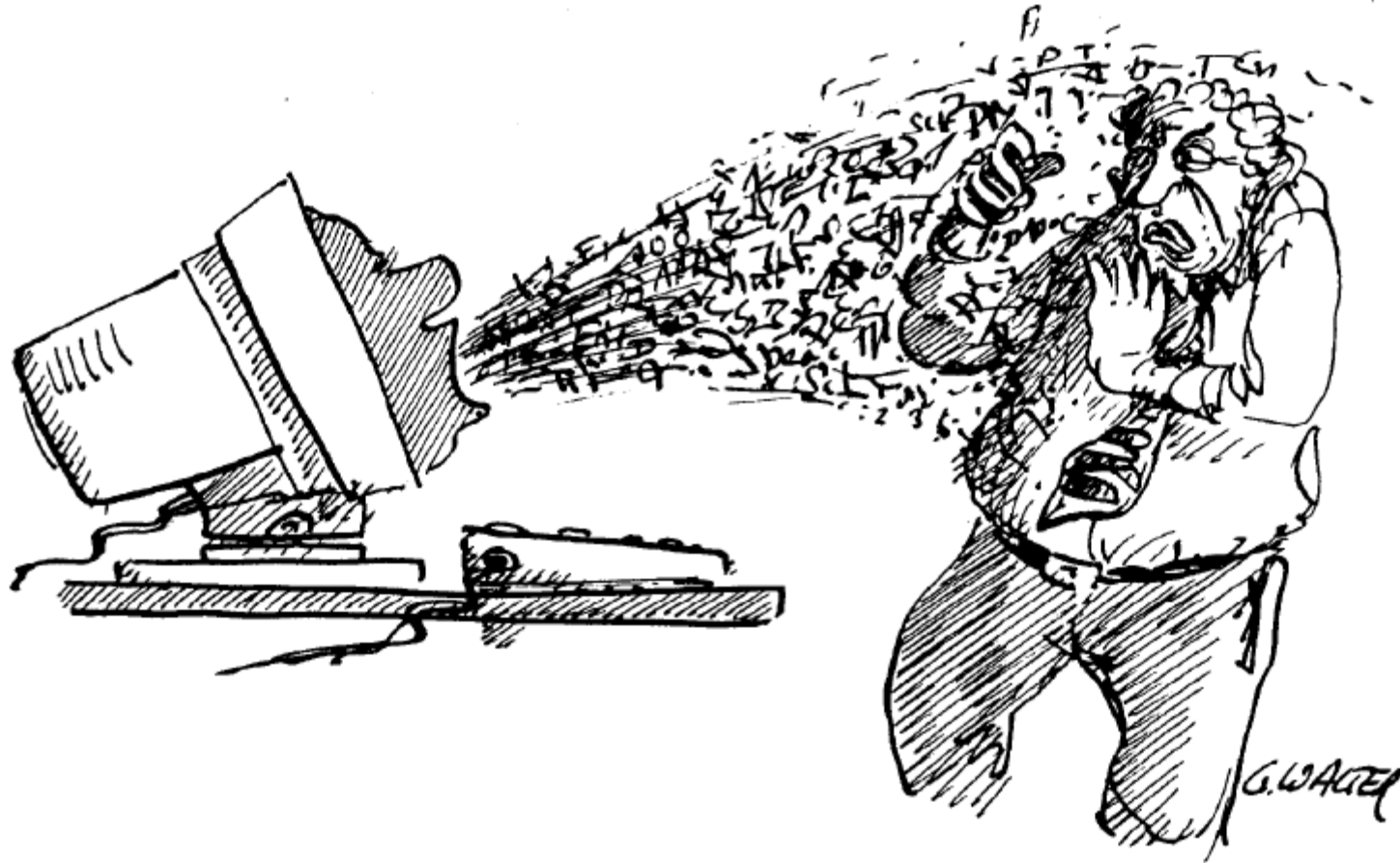
Collect and  
consolidate data





*adapted from slides by Dani Schneider*





*adapted from slides by Dani Schneider*



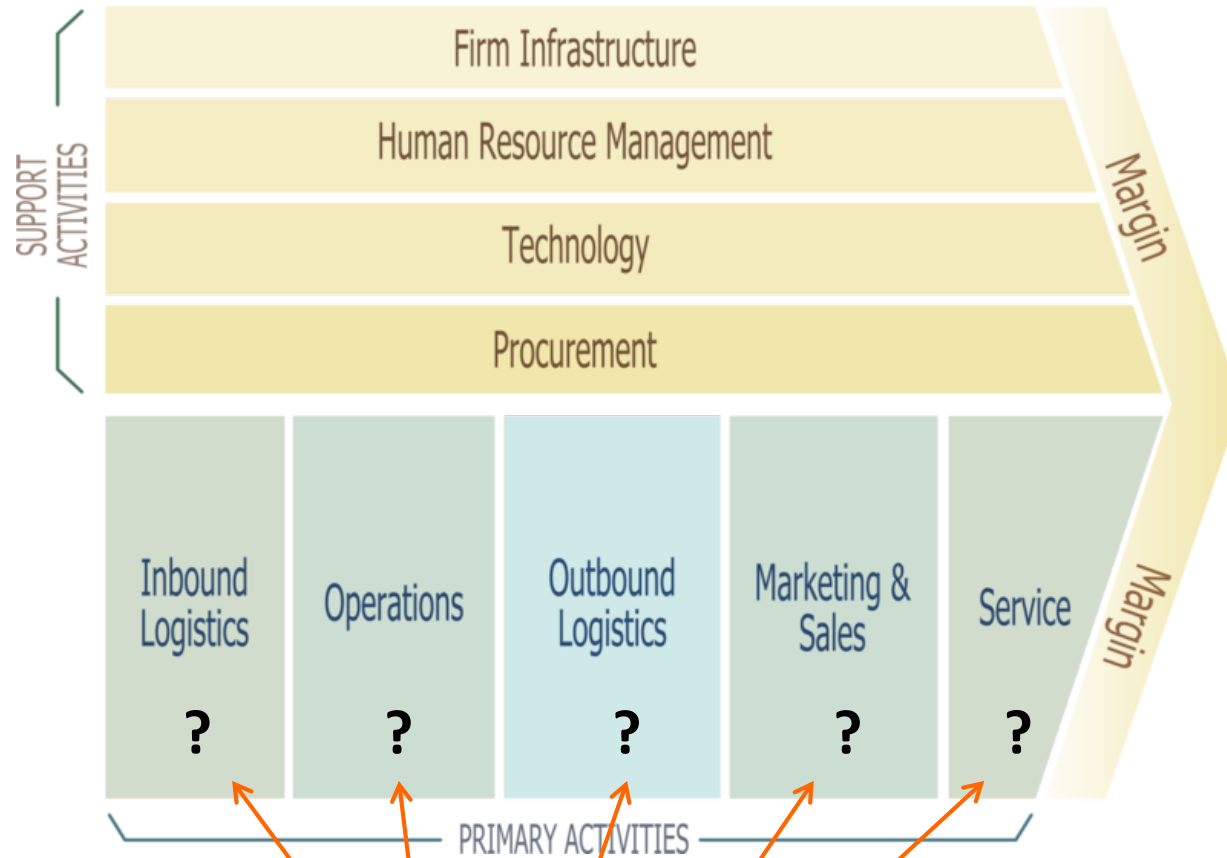
# Strategic decisions...

## ■ Business Performance Management:

- ◆ “how to perform better as a company?”
- ◆ BI helps to achieve that by enabling measurement of achievement of strategic goals via Key Performance Indicators (KPIs)
  1. Define **strategy**
  2. Define **goals**
    - e.g., identify key business processes to be improved, derive (concrete) strategic goals
    - for each goal, define KPIs and target values
  3. **Measure**
    - current values of KPIs (dashboard/cockpit)
    - analyse / compare current to targeted values
  4. **Decide...**
    - understand the (possible) deviation of KPI values from target!



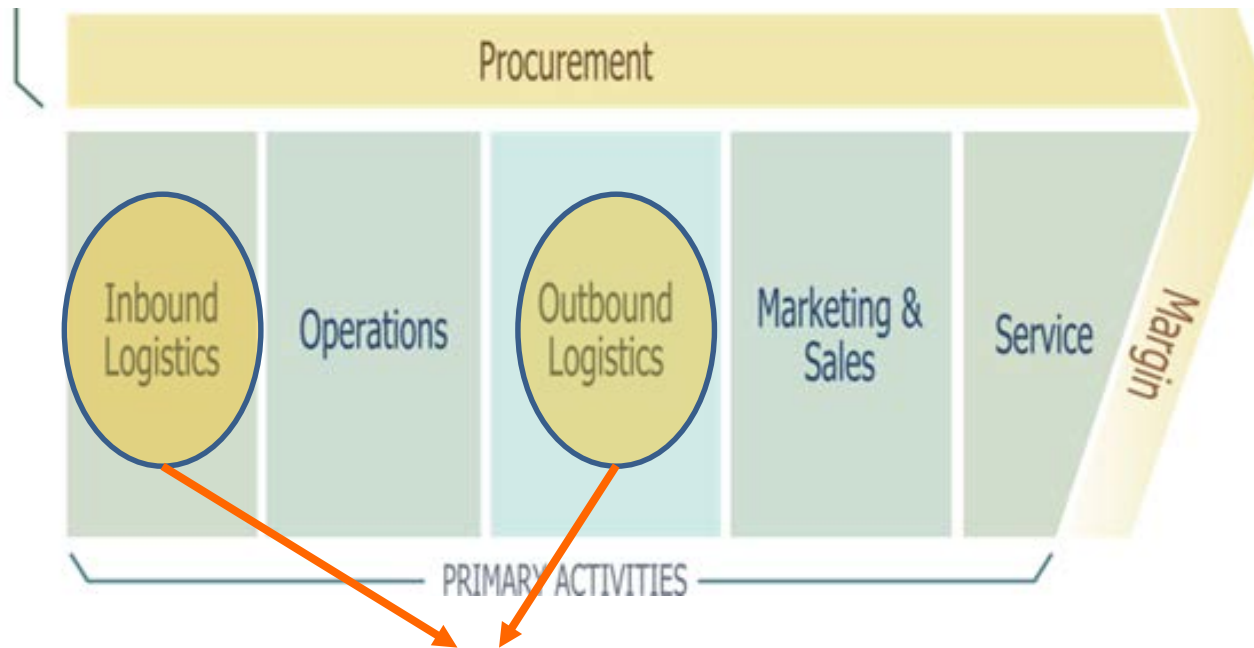
# Operative decisions: where BI creates value...



*decisions to be taken in corresponding business processes?*



# Operative decisions: where BI creates value...

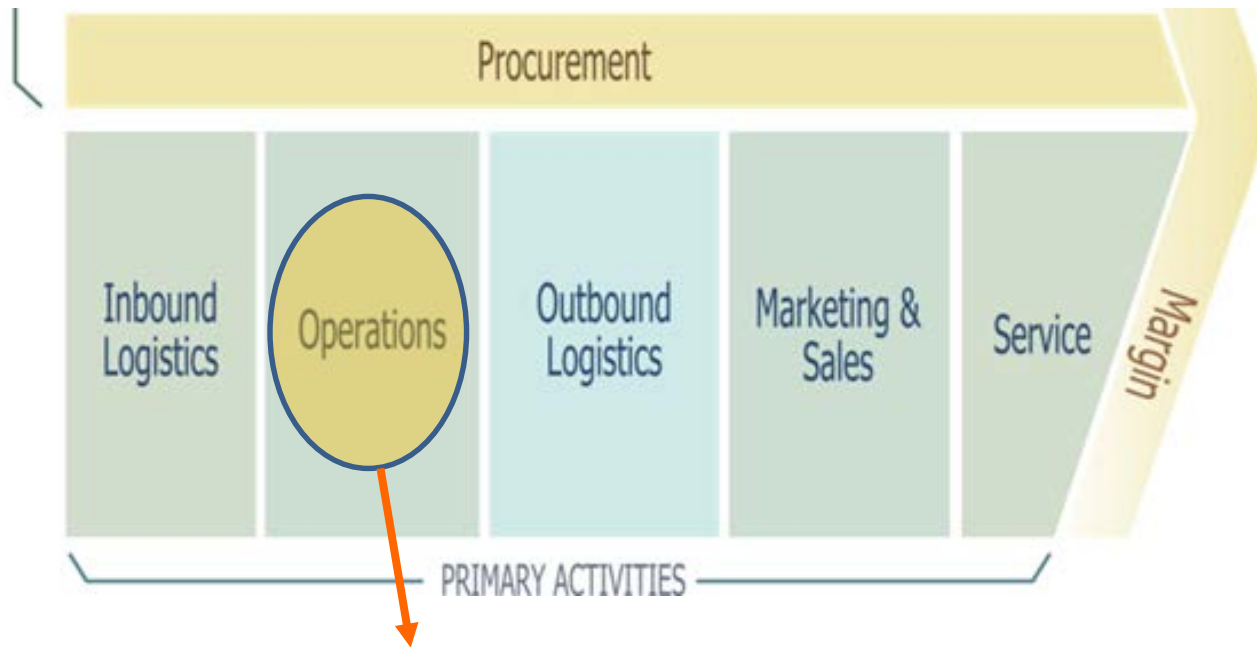


## **Logistics:**

*the process of planning, implementing and controlling the efficient, effective flow and storage of goods, services and related information from the point of origin to the point of consumption for the purpose of conforming to customer requirements*

- **how to best use resources (inbound)?**
  - which parts to order, in which quantity, at what time, from which supplier?
- **how to optimise processes (outbound)?**
  - which route/channel to use, how to schedule deliveries?

# Operative decisions: where BI creates value...

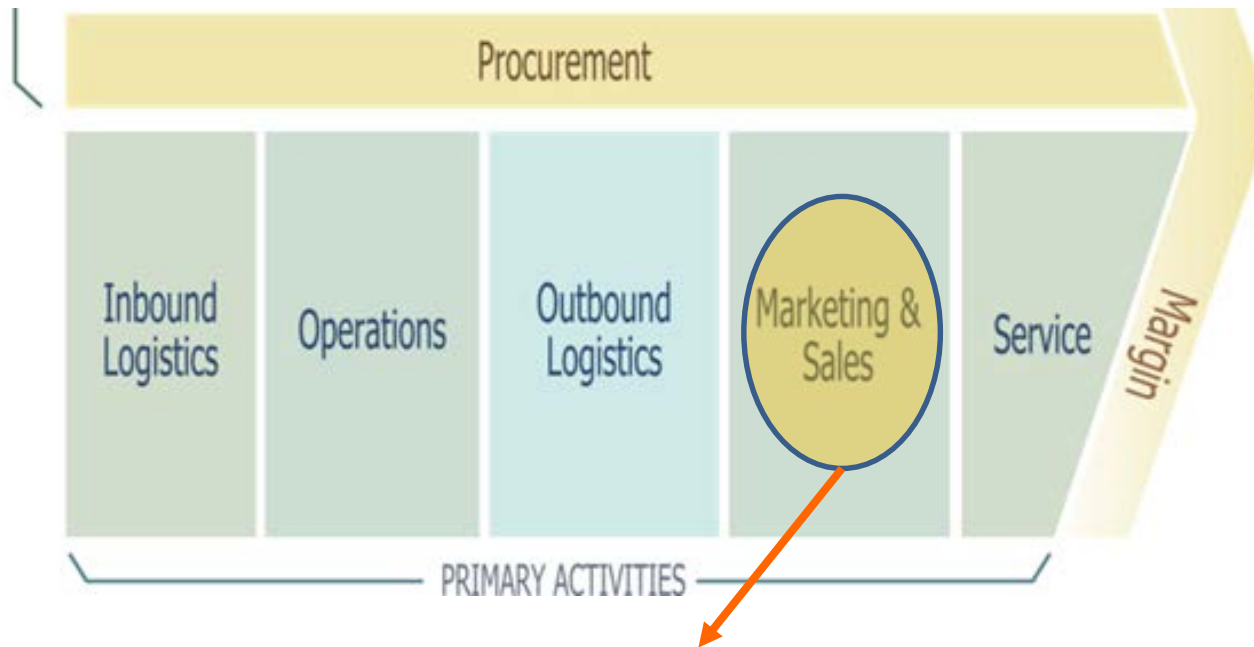


## **Operations:**

*activities associated with the functions of transforming inputs into the final product form, such as machining, packaging, assembly, equipment maintenance, testing, printing, and facility operations.*

- ***how to improve efficiency and effectiveness of processes?***
  - which resources to allocate, in which quantity, ...

# Operative decisions: where BI creates value...

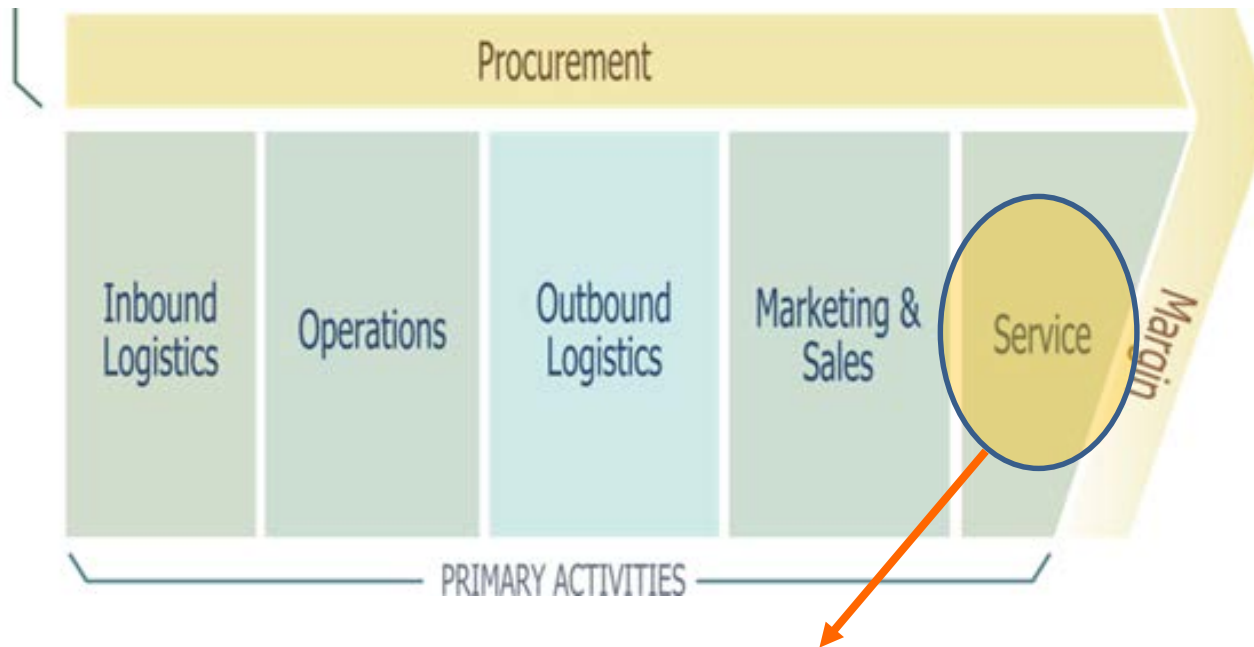


## **Marketing/Sales:**

*activities associated with the functions of providing the means by which buyers can purchase the product and inducing them to do so, such as advertising, promotion, quoting, pricing, channel and sales force management.*

- **how to understand and best address the market?**
  - which customers to approach with a campaign?
  - cross-selling: which offers to make?
  - where to place products in stores?
  - client profitability: which customers to treat with special care?
  - pricing decisions

# Operative decisions: where BI creates value...



## **Service:**

*activities associated with the functions of providing service to enhance or maintain the value of the product, such as installation, repair, training, parts supply, and product adjustment.*

- ***how to meet customer requirements and anticipate problems?***
  - which distribution channels to use for service delivery?
  - which quality problems to address first?
  - Attrition prediction: which customers to retain with special offers?

# Question types – summary

## ■ Types of questions identified:

- ◆ **query** for particular numbers or facts
  - *e.g. list of all policies that have been lost, list of all complaints, list of treatments that have been billed twice, list of high-value customers...*
- ◆ **compute a measure or KPI by aggregating numbers**
  - *e.g. cost, margin, turnover, profitability*
- ◆ **analyse KPIs / facts in different ways**
  - *e.g. sales/bookings by product/customer/sales rep/time*
  - *e.g. receipts/failures/stock by part/supplier*
  - *e.g. number of clicks/purchases by buyer/seller/page*
- ◆ **predict**
  - *e.g. predict fraudulent transactions/claims*
  - *e.g. predict if a customer will buy a product*
  - *e.g. detect types of customers or types of complaints*



# Where questions come from

- Generally speaking, companies need information to
  - ◆ monitor and improve performance
  - ◆ recognize and mitigate risks
  - ◆ recognize and seize opportunities
  
- All this can happen both on a strategic and an operative level

# Monitor and improve performance

- **Strategic level:** be able to measure if strategic goals are achieved
  - ◆ e.g. be able to measure the satisfaction of our customers over the last year
    - so that we can decide to change our customer service model
  
- **Operative level:** monitor performance within certain business processes, in small time intervals
  - ◆ e.g. find out that/why (individual) customers are not satisfied today
    - so that we can decide to call them and find a solution



# Recognise and mitigate risks

- **Strategic level:** be able to recognise general threats to our business
  - ◆ e.g. become aware that sales in certain product category are dropping dramatically (which is threatening our whole business)  
→ so that we can revise our product portfolio
  
- **Operative level:** be able to recognise risks related to individual processes, customers, suppliers, employees, ...
  - ◆ e.g. in telecommunications, be able to predict if a customer is going to cancel (or not renew) her contract  
→ so that we can decide to make a special offer to that customer





# Recognise and seize opportunities

- **Strategic level:** be able to recognise general opportunities for our business
  - ◆ e.g. become aware that (potential) customers are asking for a certain kind of product or product feature in social media  
→ so that we can decide to develop such a product
  
- **Operative level:** be able to recognise opportunities related to individual process instances, customers, suppliers, employees...
  - ◆ e.g. recognise that we can cross-sell a certain product to an existing customer  
→ so that we can decide to make the customer aware of that product



## Specific requirements

When analysing requirements for a certain company's future BI solution, usually at least the following need to be fixed:

- ◆ Which are the strategic goals → which are our **KPIs**?  
*e.g. revenue, delivery time, profitability, ...*
- ◆ By which criteria should KPI values be **grouped and/or filtered**?  
*e.g. by customer, by sales rep, by region, by date, ...*
- ◆ Which **drill-down paths** should be possible?  
*e.g. date: year → quarter → month → ...*  
*e.g. region: country → region → city*

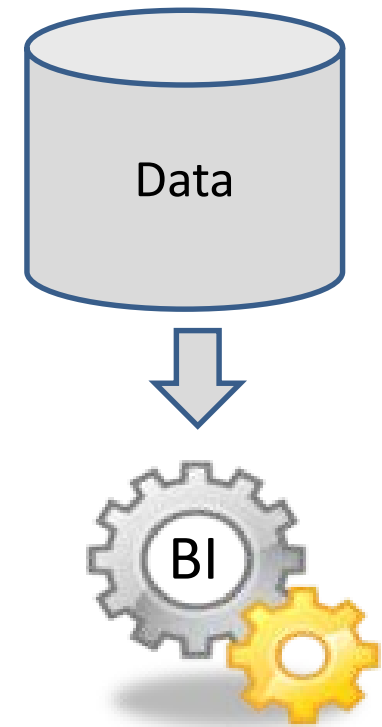


# Data

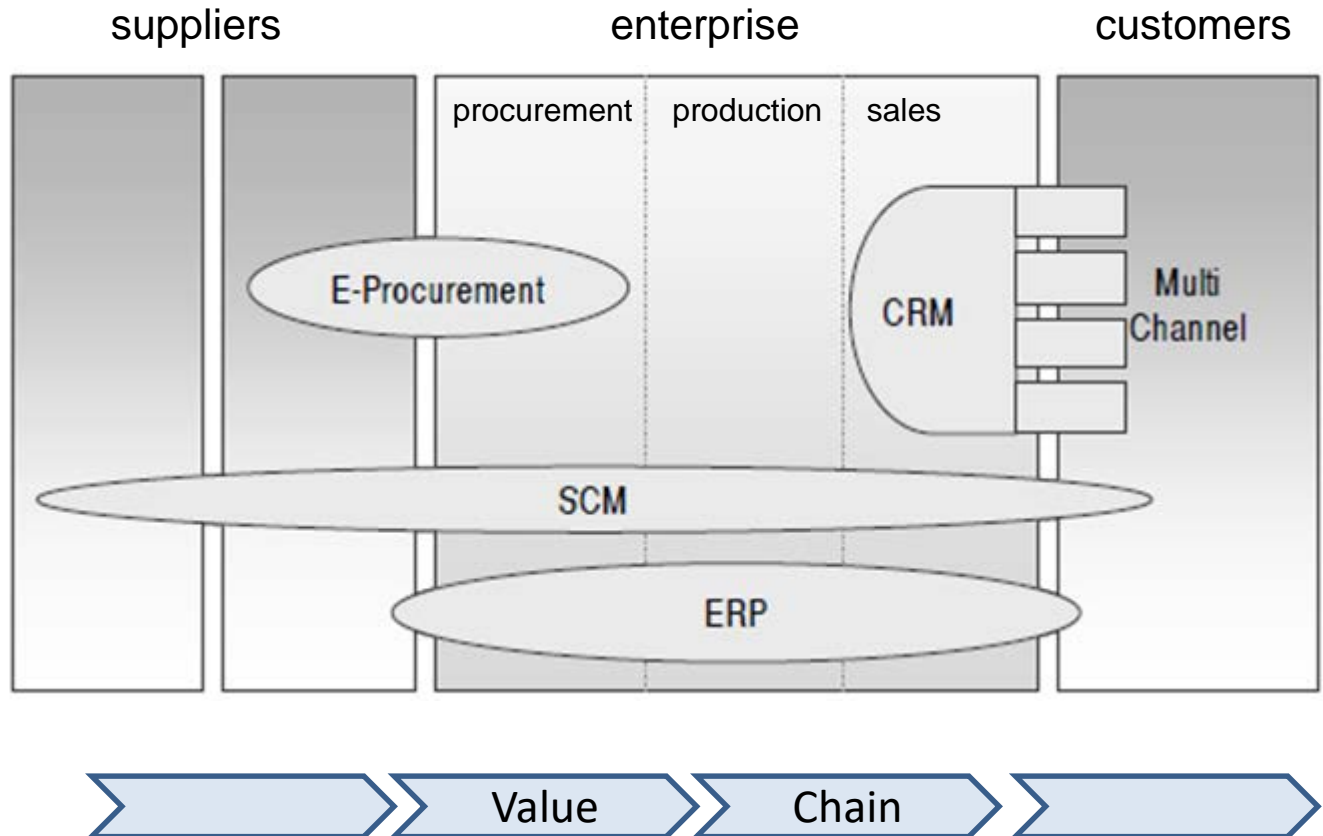


## Remember...

- ... transform **raw data** into meaningful and useful **information**...
- **Raw data** is the starting point!



# Where the data come from... (1)



CRM – Customer Relationship Management  
SCM – Supply Chain Management  
ERP – Enterprise Resource Planning

*adapted from Kemper et al. 2004*



## Where the data comes from (2)

### ■ Internal data sources:

- ◆ (Transactional) standard business applications: SCM, ERP, CRM, ...
- ◆ Legacy databases
- ◆ Web data: clickstreams from server logs, application logs
- ◆ textual documents (from DMS, CMS, intranet, email,...)



*structured-  
ness*

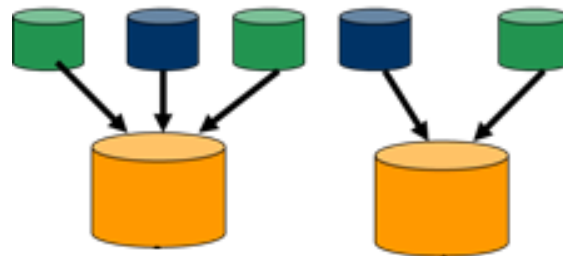
### ■ External data sources:

- ◆ Web and web 2.0

## BI tools – backend

### ■ Observations:

- ◆ many questions involve multiple (types of) data
- ◆ sometimes the data can be expected to originate from more than one source system
- ◆ for answering the questions, data from various sources needs to be connected
  - example: «Which is the best way to distribute product XYZ to customers?» → involves information about customers (e.g. profitability, behaviour) as well as about channels (e.g. cost of each channel)



# Planning Data vs. Operative Data (1)

- **operative data:** generated by and used in processing operational transactions (on-line transaction processing, OLTP)
  - ◆ many concurrent users access and modify the same data
  - ◆ focus on transactions
  - ◆ example: booking/reservation systems
- **planning data:** used for decision support
  - ◆ read-only data

*following Kemper et al. ch 2.1*





# Planning Data vs. Operative Data(2)

	<b>Operative data</b>	<b>Planning data</b>
<b>users</b>	clerk, IT professional	knowledge worker, decision maker
<b>Function/goal</b>	Support day to day operations (value adding business processes)	decision support
<b>DB design</b>	application-oriented	subject-oriented
<b>data</b>	current, up-to-date detailed information on business events, flat relational	historical, summarized, multidimensional integrated, consolidated
<b>usage</b>	Continuous, repetitive, concurrent	ad-hoc
<b>access</b>	read/write index/hash on primary key	lots of scans
<b>Queries</b>	Static, transactions embedded in application code	Ad-hoc, for changing information needs
<b># records accessed</b>	tens	millions
<b>#users</b>	thousands	hundreds
<b>DB size</b>	100MB-GB	100GB-TB
<b>metric</b>	transaction throughput	query throughput, response

*adapted from <http://www.slideshare.net/idnats/data-warehousing-and-data-mining-presentation-725476>*



# Data Warehouse – BI Backend

# Data warehouse

## ■ A data warehouse is

- ◆ “a **copy** of transaction data specifically structured for querying and reporting” (Kimball et al. 2008)

or

- ◆ “an environment [...] comprising a data store and [...] tools for data extraction, loading, storage, access, query and reporting [...] to **support decision-oriented management queries**” (Bashein/Markus, 2000)

## ■ 4 essential characteristics (Inmon 2005):

- ◆ **Subject oriented:** data are organized around sales, products, etc.
- ◆ **Integrated:** data are integrated to provide a comprehensive view
- ◆ **Time variant:** historical data are maintained
- ◆ **Nonvolatile:** data are not updated by users

# Integration, Time variance

- **Integration:** provide a «single version of the truth»
  - ◆ remove redundancy, inconsistency, semantic contradictions
  - ◆ The nasty bit... See ETL processes later
- **Time variance:** every record in the DWH has a specified moment or period when it was valid; data is collected over a long time period

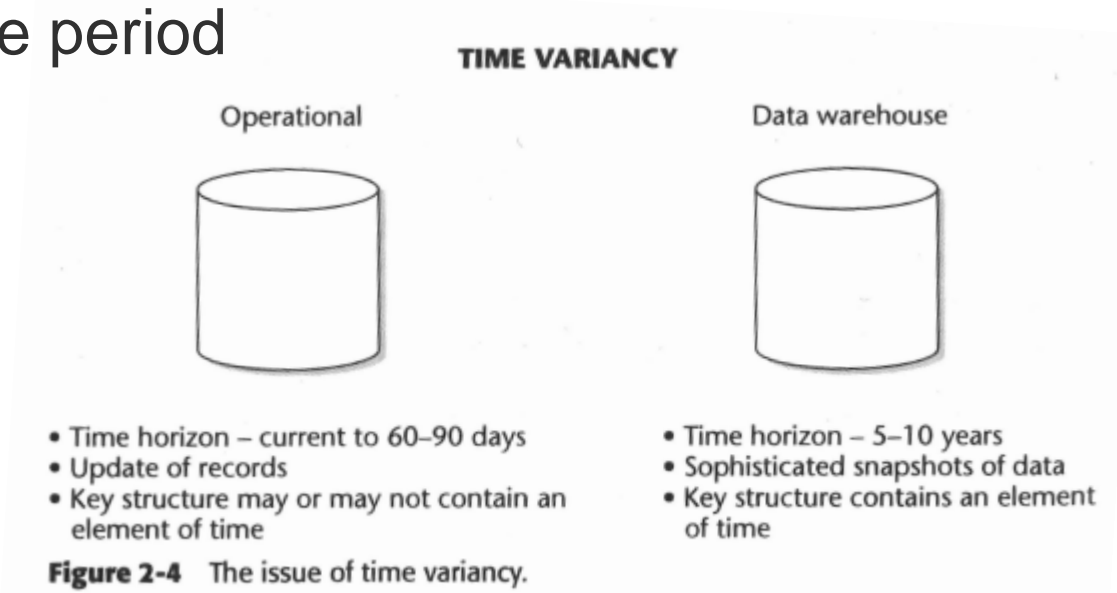


figure taken from B. Inmon: *Building the data warehouse.*

# Non-volatility

- **Non-volatility:** data is not updated by end users on a regular basis
  - ◆ bulk loading, «read-only» access

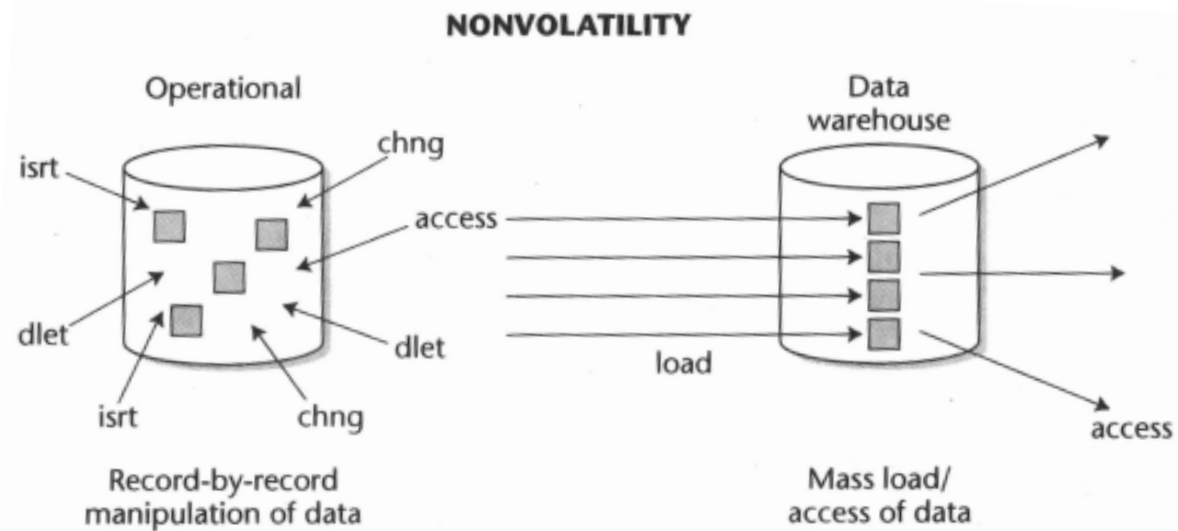


Figure 2-3 The issue of nonvolatility.

figure taken from B. Inmon: Building the data warehouse.



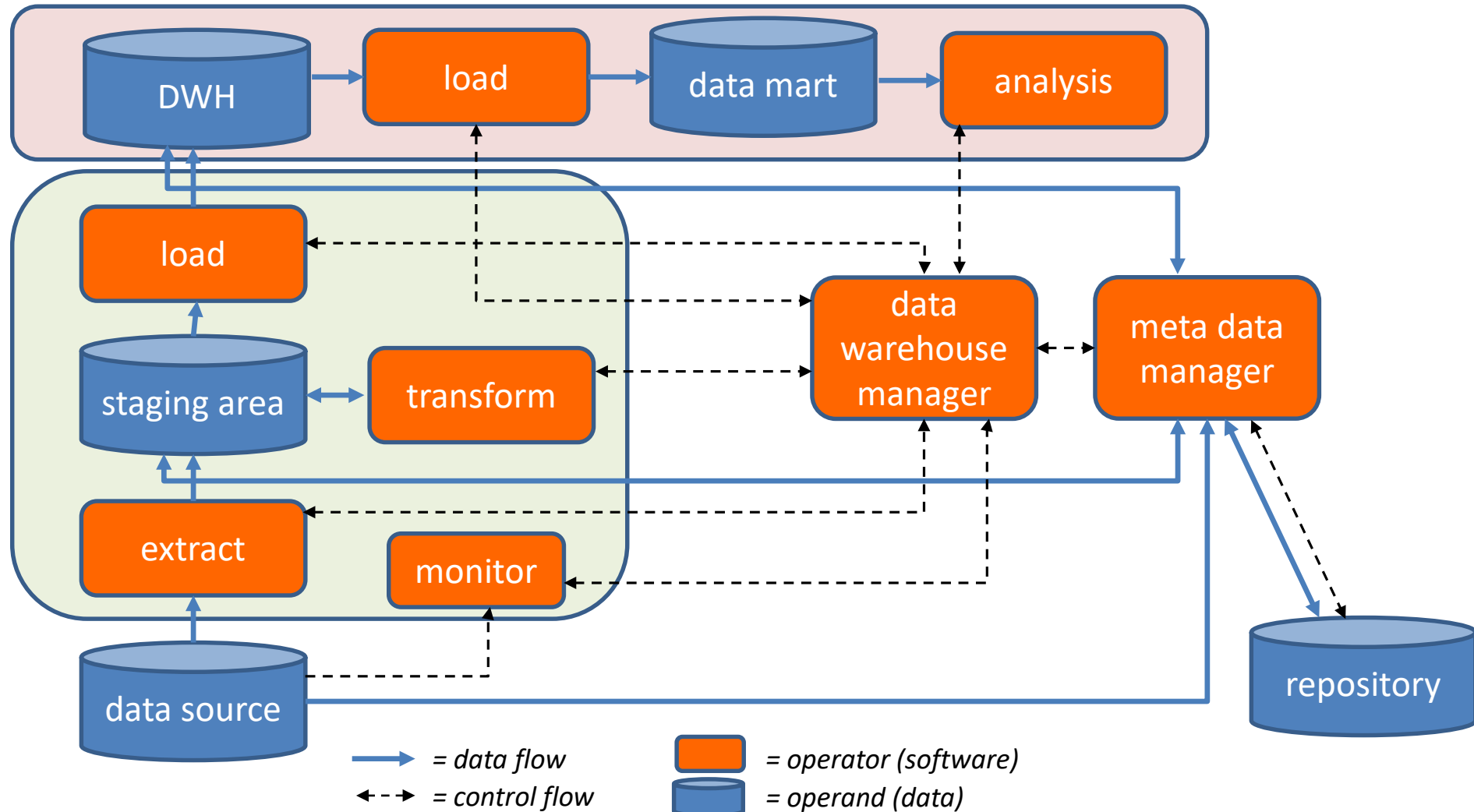
# Data warehousing

- *Data warehousing is the entire process of data **extraction, transformation, and loading** of data to the warehouse and the access of the data by end users and applications.*

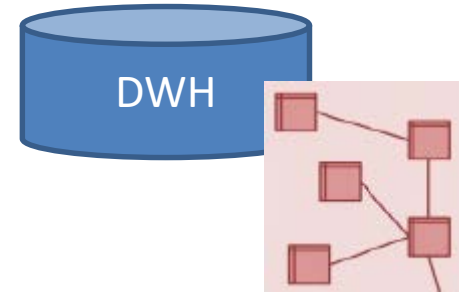
*following: [http://www.terry.uga.edu/~hwatson/dw\\_tutorial.ppt](http://www.terry.uga.edu/~hwatson/dw_tutorial.ppt)*



# Reference architecture – overview



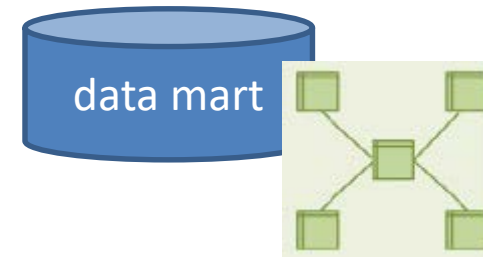
# Reference architecture – Data Warehouse (DWH)



- DWH = integrated data basis for all analyses
  - ◆ integration means both schema and data integration of various sources => «single point of truth» (no by-passes allowed!)
  - ◆ purpose: flexibility for re-using the data in multiple analyses, no focus on (and hence no pre-aggregations for) particular types of questions
  - ◆ not always present because of high cost: building data marts for specific analysis purposes directly from source data is often cheaper...



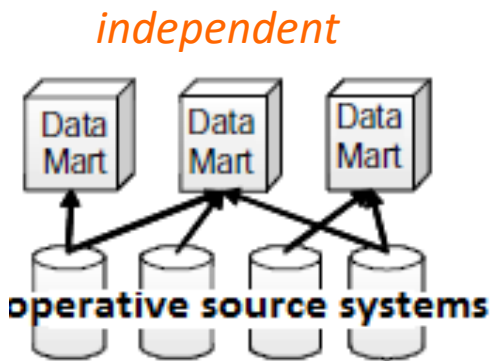
# Reference architecture – data mart



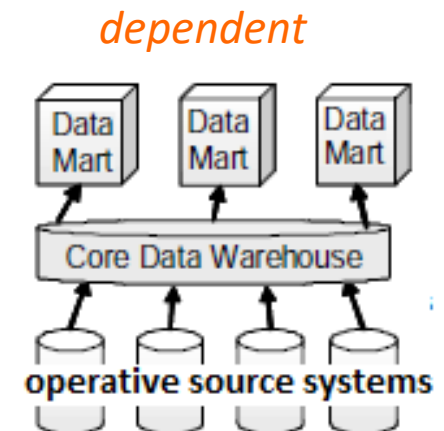
- Each data mart is designed for specific analyses
  - ◆ represents only a relevant subset of the DWH data,
  - ◆ transformed into target schema most suitable for the required analyses (very often a multidimensional model) and
  - ◆ sometimes aggregated
  - ◆ offers access to data as well as processing functionality (e.g. computation of sums, mean, variance,...) via query languages such as SQL or MDX

# Data Marts

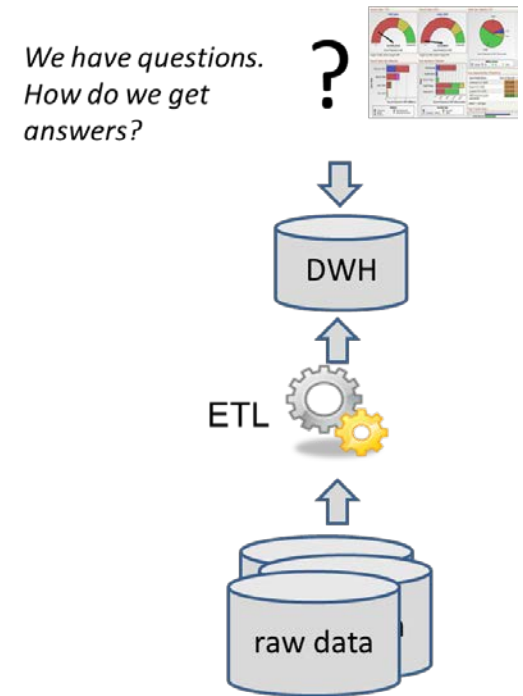
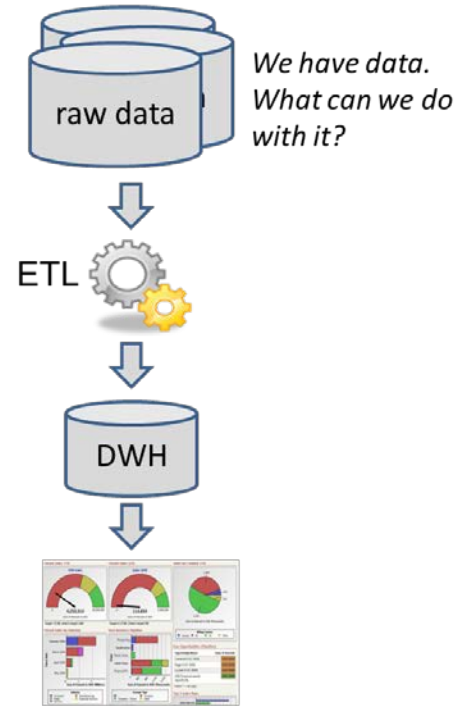
- A **data mart** stores data for a *limited* number of subject areas, such as marketing and sales data. It is used to support *specific* applications.
  - ◆ An **independent** data mart is created directly from source systems.
  - ◆ A **dependent** data mart is populated from a data warehouse.



following Kemper et al. fig. 2.4



# DWH architecture philosophies

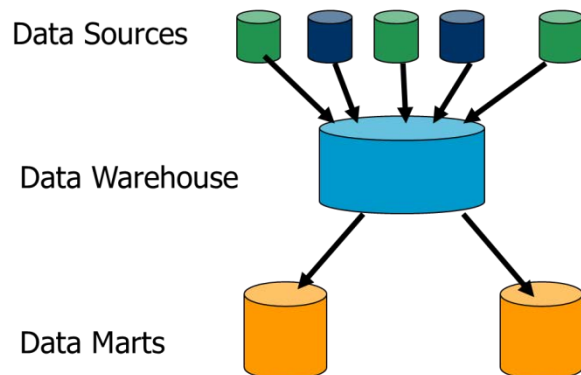


The requirements for the data warehouse cannot be known until it is partially populated and in use. [...] Therefore, data warehouses cannot be designed the same way as the classical requirements-driven system. On the other hand anticipating requirements is still important. Reality lies somewhere in between.  
*(from W.H. Inmon: «Building the Data Warehouse»)*

# DWH architectures

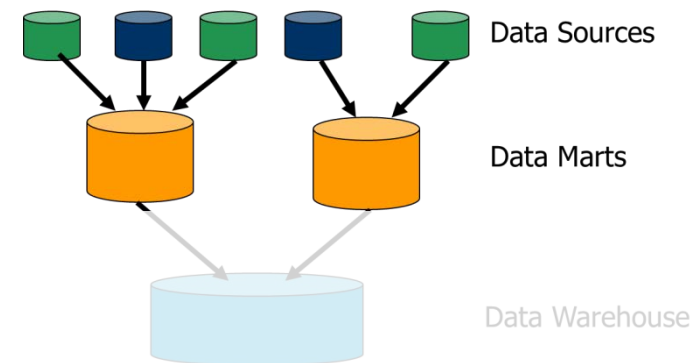
## Enterprise Datawarehouse

- Advocate: Bill Inmon
- also called «hub-and-spoke»
- **strategy:** aggregate all enterprise data into one core DWH, derive dependent data marts as subsets as needed

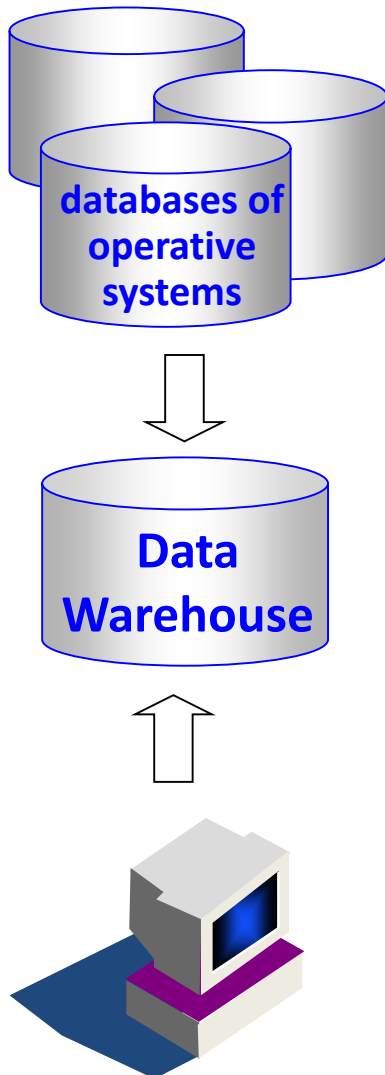


## Independent data marts

- Advocate: Ralph Kimball
- often done in reality («historic reasons»)
- **strategy:** build data marts from source systems, re-use dimension tables where possible. Possibly join marts into a centralised DWH later



# Data Warehouse



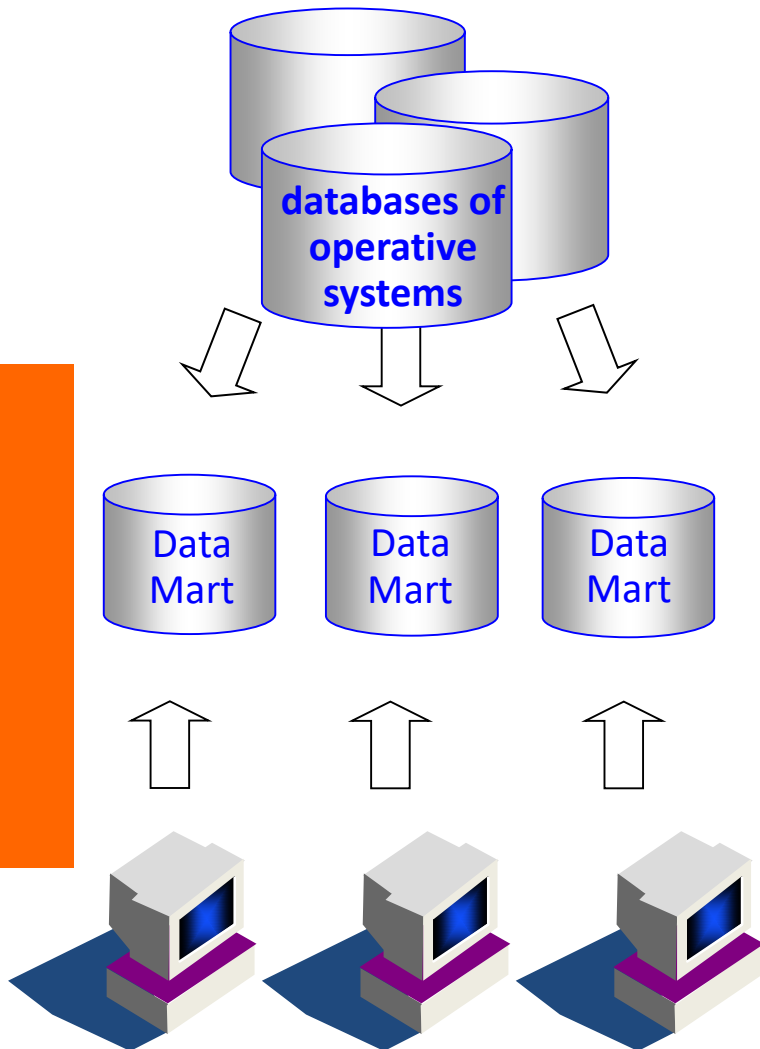
A Data Warehouse is a database which supports strategic decisions by providing ...

- ✓ high-volume and
- ✓ regular excerpts from
- ✓ operative databases
- ✓ by periods and
- ✓ often aggregated<sup>1</sup>
- ✓ also for ad hoc<sup>2</sup> analysis

1) combined, consolidated (e.g. als sum, average, indicators)

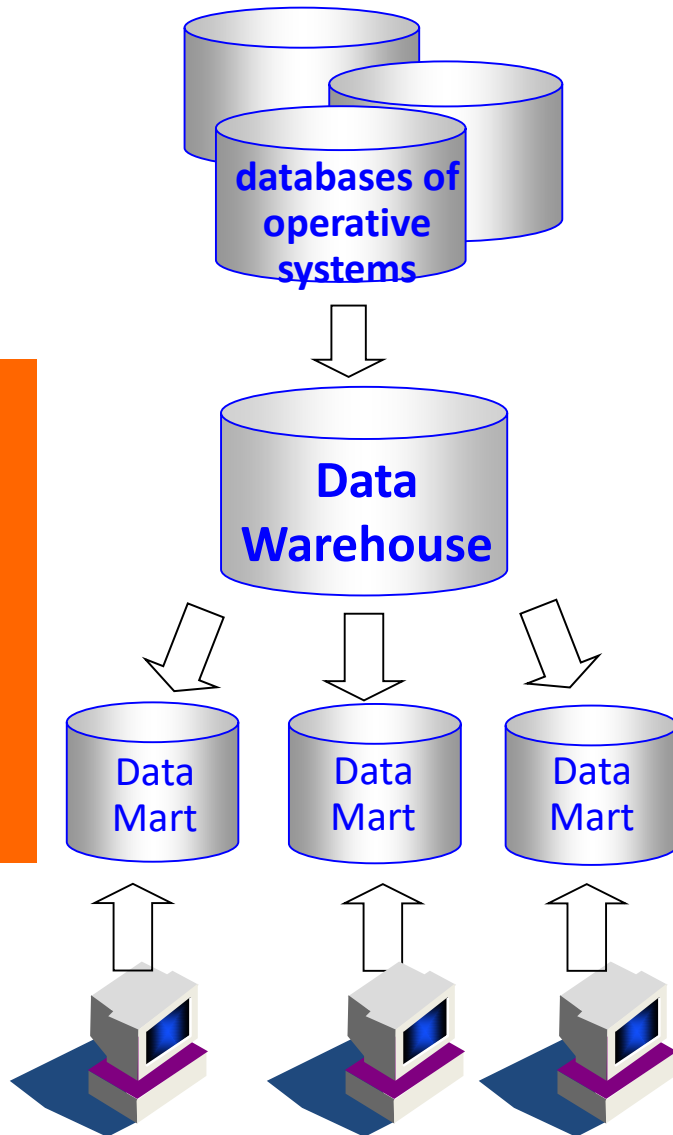
2) without preparation, in contrast to standardized analysis

# Modularisation: Enterprise Data Mart



- Distribute data to several Data Marts (no DWH)
- Advantages:
  - ◆ data model easier to understand
  - ◆ efficient access
- Problem: overall analyses
- Coordination:
  - ◆ Loading cycles
  - ◆ Data model:  
Attributes with equal meaning meaning should have the same identifier , key, datatype in all data marts

# Hierarchical Architecture



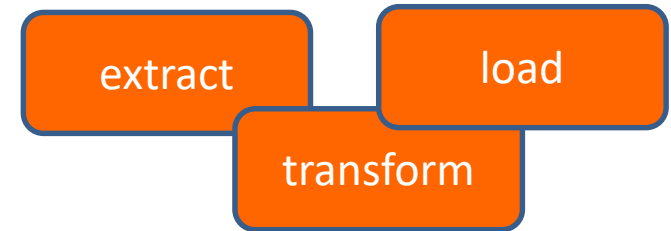
- Coordination of local Data Marts by an Enterprise Data Warehouse (EDW)
- Objective of the EDW:
  - ◆ Extraction, integration and distribution of the data
  - ◆ integrated data model
- Objective of the Data Marts:
  - ◆ Queries and analysis on parts of data
  - ◆ Adaptation to needs of an organisation unit or process
- Data distribution
  - ◆ time-based
  - ◆ event-driven
  - ◆ on demand

# Data Marts: Departmental vs. Enterprise?

- Question: should data marts be enterprise-wide or departmental?
  - ◆ Answer 1: Data marts should be organised around business processes (orders, invoices,...), not department boundaries!
  - ◆ Answer 2: ... but they don't necessarily have to be enterprise-wide (depends on the business process)!



# Reference architecture – ETL



## ■ Extract

- ◆ read source data into staging area as indicated by monitor
  - e.g. from file created by triggers or tables created by replication services
- ◆ control the selection of data that should be copied
- ◆ Extractions can happen...
  - ...periodically
  - ...on human request
  - ...event-based (e.g. when a certain number of changes has occurred)
  - ...upon each change

## ■ Transform: data cleaning and schema integration (more on this later)

## ■ Load: copy transformed data into DWH

# ETL – Extract, Transform, Load

# ETL process

- The process of
  - ◆ **extracting** relevant data from source systems (e.g. transaction-based ones)
  - ◆ **transforming** the data into the target format defined for the DWH or data mart
  - ◆ **loading** the data into the DWH
  
- The nasty, time-consuming and hence costly bit of data warehouse design
  - ◆ => do not underestimate the possible dirtiness of data!!!



# ETL process – transformation tasks

- **Transformation** = adapting data, data quality and schemas to the requirements of users
  - ◆ **Filtering:** remove syntactic and semantic defects of data
  - ◆ **Harmonisation:** map source schemas to the target schema of the DWH
    - syntactic harmonisation: schema integration + data integration
    - business harmonisation
  - ◆ **Aggregation:** aggregate data along dimension hierarchies (e.g. «customer», «customer segment», «all»)
  - ◆ **Enrichment:** pre-compute values of frequent interest and store as new attributes
    - on the basis of harmonised/aggregated data

# ETL- Filtering: Error Classes

	1. class: Automatic identification	2. class: (Semi-) automatic identification
Syntactic	<ul style="list-style-type: none"> <li>- Known formatting variants (abbreviations, date formatting etc.)</li> <li>- encoding problems</li> </ul>	<ul style="list-style-type: none"> <li>- Spelling variants/errors</li> </ul>
Semantic	<ul style="list-style-type: none"> <li>- Missing values (incompleteness)</li> <li>- redundancy (duplicates)</li> <li>- non-unique identifiers</li> <li>- missing referential integrity</li> </ul>	<ul style="list-style-type: none"> <li>- Incorrectness (e.g. outliers)</li> <li>- inconsistencies (violating business rules or contradictions)</li> <li>- dummy values</li> </ul>



# ETL- Filtering: Correction Measures

## ■ Correction measures

### ◆ 1st class:

- *incompleteness*: define rules to fill in missing values (e.g. replace sales values with ones from previous month or planned ones)
- *duplicate detection*: often there is a combination of values that unambiguously identifies a record => if these are the same, match!
- *formatting/encoding/non-unique id issues*: simple scripting

### ◆ 2nd class:

- *spelling variants/errors*: use string similarity, thesauri (extend as you go along)
- *general incorrectness*: hard to spot automatically, can define automatic sanity checks...
- *outliers*: statistic analyses
- *inconsistencies*: checks based on business rules

# ETL - Harmonization

- These are parts of tables that should be integrated in a DWH. What harmonisation tasks/problems do you see?

CustomerID	Name	City
11	Peter	Rom
15	Paul	Camerino
18	Mary	Olten
25	Joe	Bern

PurchaseID	CustomerID	Date	ProductID
1002	11	5 May 2015	SE4256
1003	18	5 May 2015	EA4516
1004	11	6 May 2015	EA4516
1005	25	6 May 2015	RG3452

ComplaintID	Complaint	Person
36536	Return	George
44363	Failure	Paul
46344	Failure	John

# ETL – Harmonisation: Schema integration

Problem	characteristics	Example: data source 1	Example: data source 2	Solution
Synonyms	Attributes with different names have identical meaning	Attribute «employee» contains employee name	Attribute «staff» contains employee name	Choose an attribute name
Homonyms	Same attribute name refers to attributes with different meaning	Attribute «partner» refers to name of customer	Attribute «partner» refers to name of supplier	Choose different attribute names



# ETL – Harmonisation: data integration (1)

Problem	characteristics	Example: data source 1	Example: data source 2	Solution
Deviating primary keys (synonyms)	Same entity has different id in different operational DBs	Customer «Smith» has id 376_ACC in accounting application	Customer «Smith» has id 7843_CC in call center application	Record linkage: identify identical entities via overlapping attribute values; use mapping table

- How to detect entity identity?

# ETL – Harmonisation: data integration (2)

- Mapping tables: allow to map updates in sources to DWH records

AD_SYS	...	customer	LOADTIME
AD-FX8257		Müller	31DEC2009:23:03:08
AD-FH2454		Meier	31DEC2009:23:03:08
AD-FX7059		Schulz	31DEC2009:23:03:08
AD-FT2567		Schmitz	31DEC2009:23:03:08
...	...	...	...

AC_SYS	customer	customerStatus
3857 ACC	Müller	A
3525 ACC	Meier	A
3635 ACC	Schulz	A
3566 ACC	Schmitz	B
...	...	...

CC_SYS	cust_grp	customer	LOADTIME
59235395	retail	Müller	31DEC2009:23:03:08
08485356	industry	Meier	31DEC2009:23:03:08
08555698	industry	Schulz	31DEC2009:23:03:08
85385386	retail	Schmitz	31DEC2009:23:03:08
...	...	...	...

AD=customer service  
CC = call center  
AC = accounting

Kunde ID	cust_id	...	AD_SYS	CC_SYS	AC_SYS	...	LOADTIME
0001	Müller		AD-FX8257	59235395	3857 ACC		31DEC2009:23:03:08
0002	Meier		AD-FH2454	08485356	3525 ACC		31DEC2009:23:03:08
0003	Schulz		AD-FX7059	08555698	3635 ACC		31DEC2009:23:03:08
0004	Schmitz		AD-FT2567	85385386	3566 ACC		31DEC2009:23:03:08
...	...	...	...	...	...	...	...

*adapted from Kemper et al.*



# ETL – business harmonisation

## ■ adjust figures/values

- ◆ consolidate figures from various databases based on their (business) meaning, e.g. apply rules to map location- or department-specific value deviations
- ◆ convert currencies and units (e.g. inch → cm)

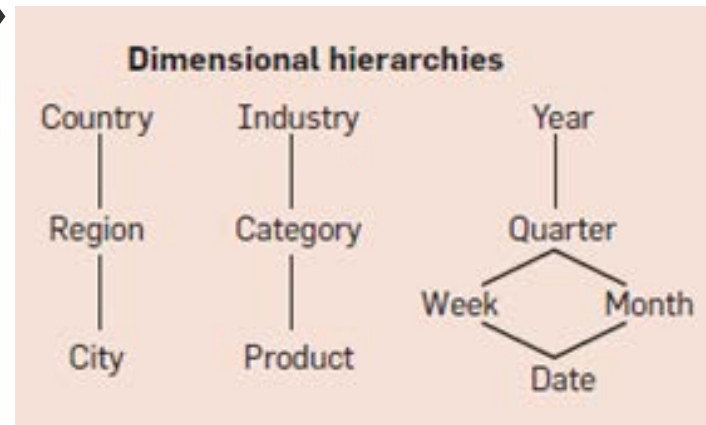
## ■ adjust granularity

- ◆ decide for a level of granularity (e.g. monthly or quarterly)
- ◆ harmonise according to period (source systems may have differing granularity, e.g. quarters vs. years)
- ◆ aggregate all values on that level (e.g. sum all records/receipts of one day together)



# ETL - Aggregation

- Aggregate data based on dimensional hierarchy
  - ◆ usually, aggregates are pre-computed for performance reasons
  - ◆ introduces «controlled redundancy»
  - ◆ aggregates become invalid when hierarchies and/or source data change...

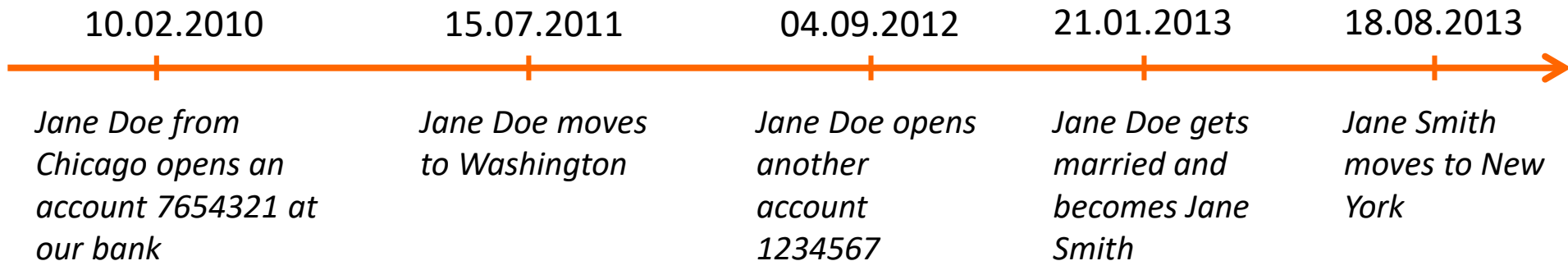


# ETL - Enrichment

- add new attributes that are functions of existing data; compute these functions and store the result
  - ◆ sums, averages or more complicated computations (e.g. profitability)
  - ◆ based on harmonised and/or aggregated data
  - ◆ same motivation as aggregation: performance
  - ◆ introduces another «controlled redundancy»

# Slowly Changing Dimensions

Example: customer dimension change




- who's the owner of the bank account 1234567?
  - ◆ as of today: Jane Smith from New York
  - ◆ as of 31.12.2012: Jane Doe from Washington
  - ◆ as of 31.12.2011: there is no such bank account



# Type I: no history

Cust_id	Cust_name	Cust_city	...
1	John Allan	Chicago	...
2	Chris Lee	Boston	...
3	Jane Doe	Chicago	...



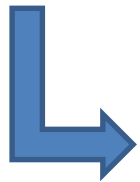
Cust_id	Cust_name	Cust_city	...
1	John Allan	Chicago	...
2	Chris Lee	Boston	...
3	Jane Doe	Washington	...

15.07.2011

old value is simple overwritten with new value

# Type II: full history

Cust_id	Cust_name	Cust_city	Valid from	Valid to
1	John Allan	Chicago	10.03.2008	
2	Chris Lee	Boston	02.06.2010	
3	Jane Doe	Chicago	10.02.2010	



18.08.2013

Cust_id	Cust_name	Cust_city	Valid from	Valid to
1	John Allan	Chicago	...	
2	Chris Lee	Boston	...	
3	Jane Doe	Chicago	10.02.2010	14.07.2011
3	Jane Doe	Washington	15.07.2011	20.01.2013
3	Jane Smith	Washington	21.01.2013	17.08.2013
3	Jane Smith	New York	18.08.2013	

- every intermediate state is documented, validity range of values is signalled via «valid from», «valid to» attributes
- «valid from» becomes part of primary key



## Type III: limited history

Cust_id	Previous Cust_name	Current Cust_name	Effective date cust_name	Previous Cust_city	Current Cust_city	Effective date cust_city
1		John Allan			Chicago	10.03.2008
2		Chris Lee			Boston	02.06.2010
3		Jane Doe			Chicago	10.02.2010

18.08.2013



Cust_id	Previous Cust_name	Current Cust_name	Effective date cust_name	Previous Cust_city	Current Cust_city	Effective date cust_city
1		John Allan			Chicago	10.03.2008
2		Chris Lee			Boston	02.06.2010
3	Jane Doe	Jane Smith	21.01.2013	Washington	New York	18.08.2013

- keeps the n previous values, each in a separate new column (in the example: n=1)
- effective date column(s) show(s) when the change occurred

