

POZNAN UNIVERSITY OF TECHNOLOGY

# **Modeling Data Warehouse**

Robert Wrembel Poznan University of Technology Institute of Computing Science Robert.Wrembel@cs.put.poznan.pl www.cs.put.poznan.pl/rwrembel







## Models

### Conceptual

- mutlidimensional data model (MD)
- Logical (implemenation)
  - relational →ROLAP
  - multidimensional → MOLAP
  - hybrid → HOLAP
- Slowly Changing Dimensions



© Robert Wrembel (Poznan University of Technology, Poland)



# **Categories of Data**

### ⇒ Facts

- data to be analyzed in a given context
  - sales, phone calls
  - characterized quantitatively by measures
  - number of intems sold, phone call duration time

### Dimensions

- define the context of an analysis
  - chocolate sales (product) in Auchan (shop) in months (time)
- typically composed of levels that form hierachies





### ROLAP

### **C ROLAP**

- star schema
- snowflake schema
- fact constellation schema
- starflake schema

© Robert Wrembel (Poznan University of Technology, Poland)



© Robert Wrembel (Poznan University of Technology, Poland)



# Structure of dimension



© Robert Wrembel (Poznan University of Technology, Poland)



# **ROLAP - snowflake schema**

Countries #country_id name  Regions #region_id name  country_id	<ul> <li>Dimension</li> <li>⇒ levels (level tables)</li> <li>⇒ hierarchies</li> </ul>
Cities Cities Cities Feity_id name region_id Shops Shop_id name address City_id	
Groups #group_id group_id group_id group_id group_id #products #product id product_id product	Days #time_key day_no day_name  month_id #onth_id duarter_id duarter_id  guarter_id year   







© Robert Wrembel (Poznan University of Technology, Poland)



# Which schema to use?

### Data redundancy

- e.g., dimension TIME: 1sec granularity and time span of 10 years ⇒ 300 000 000 rows
  - telecom, physics

Query performance

© Robert Wrembel (Poznan University of Technology, Poland)



### Which schema to use?

### Hints

- attributes from different hier. levels frequently used together in roll-up ⇒ store them in one table
- high level attributes rarely used ⇒ store them in separate high level tables





### Stores facts that typically represent events





## **Dimension TIME**

### Exists in all DW schemas

### Typical granularity - day

- Time\_key
  - artificial ID: values 1, 2, ..., n
  - date-time type
  - timestamp type
  - numerical type: 11032008 (11-03-2008)

#### Date\_id (PK) Date Day\_name Day\_no\_week Day\_no\_month Day\_no\_year Last\_day\_month Last\_day\_year Week\_no\_year Month\_name Month\_no Quarter Year Fiscal\_day\_no\_week Fiscal\_day\_no\_month Fiscal\_day\_no\_year Fiscal\_...

Holiday Holiday\_type Shop\_open\_holiday Weekend\_day

© Robert Wrembel (Poznan University of Technology, Poland)



# **Dimension TIME**

# Registering time with granulairty > days

### timestamp in a fact table





# **Dimension roles**

# The same dimension may play different roles (give contexts) for a fact table

e.g., the TIME dimension





### **Dimension roles**

### ⇒ Fact dimension: a dimension in a fact table

Dim: TIME	EUR_exchange
Time_key (PK)	Time_key (FK)
	Currency_symbol
	exch_rate

### Measure being a dimension

 trip length → discretization of values: short, medium, long



# **Artificial IDs**

### By an ETL process

### ⇒ ID types

- numerical
  - efficiency in processing
  - chronology represented by values
  - no semantics
- alphanumerical
  - less efficient in processing
  - may have semantics
  - longer than numerical
  - typically concatenation of natural key and timestamp

© Robert Wrembel (Poznan University of Technology, Poland)



© Robert Wrembel (Poznan University of Technology, Poland)



# Aggregability

#### Aggregability (summarizability)

to be able to aggregate measures on a higher level (City) in a dim. hierarchy based on aggregates from a lower level (Shop)

#### Criteria of correct aggregability

#### disjointnes of level instances

- 1:M relationship between an upper and a lover level in a dim. hierarchy
- a lower level instance is related to only one upper level instance
- completeness
  - all level instances belong to a dimension instance
  - every lower level instance is related to an upper level instance
- right aggregate function to compute
  - using an adequate aggregate function to a given type of a measure

© Robert Wrembel (Poznan University of Technology, Poland)

ABURNINA POZNUJANA KORA

### **Types of measures**

- Correct aggregation for all dimensions
  - e.g., nb of items sold ⇒ possible aggregation in dim. Time, Customer, Product, Shop, Supplier, ...
- Semiadditive (stock, level) ⇒ correct aggregation for some dimensions
  - e.g., nb of items in a store ⇒ possible aggregation in dim. (Store → City → Region)
  - aggregating in Timt results in uninterpretable results
- Solution Nonadditive (value-per-unit) ⇒ aggregated values are non-interpretable in any dimension
  - e.g., net price, exchange rate, opening tick for SUM



### **Types of measures**

#### Distributive

- partial results of aggreg. function F on n subsets of set S can be aggregated into a result that is identical to applying F on the whole S
- upper level aggregate can be computed from lower level aggregates (sales by cities aggregate to sales in a country)
- count, min, max, sum

#### Algebraic

- computed using distributive functions
- avg, stdev

#### Holistic

- can be computed only on all elementary data
- median

© Robert Wrembel (Poznan University of Technology, Poland)

23



Malinowski E., Zimanyi E.: Advanced Data Warehouse Design. From Conventional to Spatial and Temporal Applications. Springer Verlag, 2008



# **Balanced hierarchy**

- Includes only one aggregation path
- **C** 1:M both-side obligatory relationship between an upper level and a lower level



© Robert Wrembel (Poznan University of Technology, Poland)



- Includes only one aggregation path
- **C** 1:M relationship between an upper level and a lower level, obligatory only from an upper level





# **Unbalanced recursive hierarchy**

### Special case



© Robert Wrembel (Poznan University of Technology, Poland)





# **Non-strict hierarchy**

# Dimension schema includes at least one M:N relationships between an upper and a lower level



#### Problem: aggregating twice the same product

© Robert Wrembel (Poznan University of Technology, Poland)

Non-strict hierarchy

### Solution

- adding a new cagetory e.g., smartphone
- distributing a measure value between categories
   e.g., telephone: 50%, digital camera: 50%
- transforming a non-strict into a strict hierarchy → assigning a product to one (main) category



# **Non-strict hierarchy**

- **⊃** Non-strict  $\rightarrow$  strict transformation
  - we must know the distribution of a measure value between categories
  - e.g., how much an employee earns in a given project





#### Problem

- count employees by organization unit
  - Proj1 is run by OU1, Proj2 is run by OU2
  - Emp1 works in Proj1 and Proj2 → will be counted for OU1 and OU2





# **Alternative hierarchy**

- **Composed of hierarchies that share at least one level**
- Aggregating through any path from facts to a shared level yields the same results at the shared level



© Robert Wrembel (Poznan University of Technology, Poland)

Pa

### **Parallel hierarchies**

- Control Parallel independent hierarchies ⇒ a standard case
  - hierarchies don't share levels
  - each hierarchy is an independent context of an analysis, e.g., Prodcts, Customers, Time



© Robert Wrembel (Poznan University of Technology, Poland)



# **Parallel hierarchies**

### Parallel dependent hierarchies

- hierarchies share levels
- each hierarchy is an independent context of an analysis



given city

35

- Sales $\rightarrow$ Shops $\rightarrow$ Cities: sales by shops in a given city
- both analyses yield different results

© Robert Wrembel (Poznan University of Technology, Poland)









© Robert Wrembel (Poznan University of Technology, Poland)





# **Drill-through**

Accessing a central DW to get more detailed data





# HOLAP

- THE MERICAN AND THE MERICAN AN
  - Control Detailed data, aggregated data (fine grained) → ROLAP
  - Copic-wise data, aggregated data (fine grained to coarse) → MOLAP
  - Central HD + data marts → HOLAP

© Robert Wrembel (Poznan University of Technology, Poland)



# **DW modeling: some remarks**

- ⊃ Identify facts → key types of transactions
  - commerce: sales transactions
  - banks: financial operations on accounts
  - stock exchange: sell/buy quotations
  - insurance: buying a policy, damage payment
  - telecom: phone calls
- Identify key dimensions
- Design a fact table
- Design dimension table(s)



# **DW modeling: some remarks**

### Fact data: how detailed?

- storing every single product purchase record
- storing the value of a whole basket
- derived attributes
  - net, vat, gross ⇔ net, vat, and dynamically computed gross
- storing only necessary attributes
  - dim table Customer with 8\*10<sup>6</sup> records
  - each customer makes daily 2 phone calls
  - one year time span in a fact table
  - decreasing the length of each fact row by 10B  $\rightarrow$  size decrease by 54GB

© Robert Wrembel (Poznan University of Technology, Poland)



© Robert Wrembel (Poznan University of Technology, Poland)



# Source-driven

### Semi-automatic discovery of a DW schema

- Song I.-Y., Khare R., Dai B.: SAMSTAR: a semi-automated lexical method for generating star schemas from an entity-relationship diagram. DOLAP, 2007
- Romero O., Abello A.: Automating Multidimensional Design from Ontologies. DOLAP, 2007
- Phipps D., Davis K.: Automating Data Warehouse Conceptual Schema Design and Evaluation. DMDW, 2002
- Jensen M, Holmgren R., Pedersen T.B.: Discovering Multidimensional Structure in Relational Data. DAWAK, 2004

© Robert Wrembel (Poznan University of Technology, Poland)

HUNNER POZNER

# SAMSTAR

### ⇒ SAMSTAR

- output: ER DW schema (generated semi-automatically)
- input: ERD of a source system
- Steps
  - find fiacts and direct dimension levels (automat.)
  - find transitive levels (automat.)
  - select facts (human)
  - select dimensions for selected facts (automat.)



### SAMSTAR

С

D

49

### **C** Rules for selecting facts and dim. levels

- entities at cardinality of M ⇒ fact candidates
- entities at cardinality of 1 ⇒ level candidates



#### **CTV** (Connection Topology Value)

- sum of weighted direct and indirect levels (transitive)
- the weight of a direct level > the weight of an indirect level
- CTV is computed for each entity in a source schema

$$CTV(e) = 1 \cdot w_{direct} + w_{indirect} \cdot \sum_{i=1}^{n} CTV(i)$$

© Robert Wrembel (Poznan University of Technology, Poland)





### SAMSTAR

### ℑ Entities for which CTV > ϑ represent facts

$$artheta = rac{\sum_{i=1}^{n} CTV(i)}{n} + k \cdot StDev$$

- k coefficient whose value is given by a designer; the greater its value the less candidates for fact entities (in practice 1.5 - 1.75)
- **c** Other entites are candidates for dimension
- WordNet for discovering synonyms of entity names

© Robert Wrembel (Poznan University of Technology, Poland)

FULL NUKA POSICIAL AND ADDRESS OF THE POSICIAL ADDRESS OF THE

# SAMSTAR

- Dim. candidates:
  - entities representing direct levels (at the 1 side of a relationship, where M connects to a fact entity)
  - entities representing indirect levels
- Designing dimensions
  - using WordNet for synonyms
  - using Dimensional Design Pattern (DDP) for defining a dimension structure
  - DDP ightarrow a collection of over 100 patterns of typical dimensions
  - adding the TIME dimension from DDP
- S Manual acceptance or modification of the result schema



### **Real case**

#### Retail chain (real data)

- a total number of products in stock, including one-timesales and seasonal ones: 15000
- a total number of regular products (always in stock): at least 1400
- maximum rate of a value change, e.g., a price: once per day
- minimum rate of a value change, e.g., a product dimension instance: once per 6 months
- average rate of a value change: once every 2 weeks
- a number of individual items sold per year: at least 15 300 000 000
- maximum length of a record describing a product: 82B
- maximum length of a record describing a sold item: 64B

© Robert Wrembel (Poznan University of Technology, Poland)

**Real case** Products Sales EAN clientID article number itemEAN article\_description categories categories quantity price price discount total value sales\_unit tax classification sales date



# **Real case**



### **C** An example instance of dimension Product

© Robert Wrembel (Poznan University of Technology, Poland)

# **SCD: Introduction**

- Dimension instances change in time
- The change could be for
  - correcting an error
    - e.g., misspelling of a customer's name
  - real change in the business world
     e.g., changing product price
- Such changes are called Slowly Changing Dimensions, as they are less frequent than changes of facts



### Real case: data change

### Case 1: value change

- price change
- it can change every day as the result of:
  - promotion campaign
  - entering a particular holiday period (e.g., prices of chocolates for Christmas and Easter)
  - reacting to price changes of competitors
  - reaching a 'best before date'
  - ...

© Robert Wrembel (Poznan University of Technology, Poland)

AUTONIA POZNARA

### Real case: data change

- Case 2: reference to units of sales
  - a given product, say yoghurt A, is sold in different package units
    - in January yoghurt A was sold in 2-packs
    - in February  $\rightarrow$  in 4-packs
    - in March  $\rightarrow$  in 6-packs
    - in April  $\rightarrow$  in 8-packs
  - in each of these sales periods, a database registers the number of units sold (i.e., n-packs), and not individual yoghurt jars



## Real case: data change

### Case 3: dimension instance change

- the instance of dimension Product changes its hierarchical structure as the result of:
- reclassifying products to other categories
  - milk  $\rightarrow$  diary  $\rightarrow$  food changed to milk  $\rightarrow$  liquids  $\rightarrow$  food
  - electronics  $\rightarrow$  7% VAT changed to electronics  $\rightarrow$  22% VAT
- removing some instances of levels 2, 3, or 4 (typically one instance level is removed)

© Robert Wrembel (Poznan University of Technology, Poland)

THE MULTINE POLICY OF THE

### Real case: data change

- Other cases: dimension instance change
  - administrative changes to a country
    - voivodships in Poland
    - East Germany + West Germany
    - Yugoslavia split
    - Czech-Slovakia split
- New products or services offered



### **SCD: Introduction**

- The value of an attribute of a dimension table changes in time ⇒ need to record the history of changes
- R. Kimball proposed seven (SCD1-SCD7, 3 basic and 4 extended) techniques to record the history of dimension attributes' value changes
  - R. Kimball, M. Ross: The Data Warehouse Toolkit, 3rd Edition, Wiley, 2013
  - Slowly Changing Dimensions: http://www.kimballgroup.com/2013/02/designtip-152-slowly-changing-dimension-types-0-4-5-6-7/

© Robert Wrembel (Poznan University of Technology, Poland)



# Type 1: Overwrite the old value

- The old attribute value is overwritten by the new one
- No history is kept
- Typically, this technique is used for correcting the erroneous data e.g., correcting spelling errors



## Type 1: Overwrite the old value

 Consider that the home city of employee 'Ahmed' is changed from 'Poznan' to 'Brussels'

EmpID	EmpName	EmpCity
111	John	Poznan
222	Doe	Wroclaw
333	Ahmed	Poznan

Employee table before change

EmpID	EmpName	EmpCity
111	John	Poznan
222	Doe	Wroclaw
333	Ahmed	Brussels

Employee table after change

© Robert Wrembel (Poznan University of Technology, Poland)

THE WHITE POINT

# Type 1: Overwrite the old value

- No history of the overwritten field is maintained which may be unacceptable in some systems e.g., financial systems
- Aggregates on the overwritten field must be recomputed
- All distributed copies of the overwritten dimension must be updated simultaneously



- For every change in the attribute value of a dimension member, a new dimension record is added
- Each record version is identified by an artificial primary key called a surrogate key
- A pair of timestamps is used to identify the validity period of the record version

© Robert Wrembel (Poznan University of Technology, Poland)



# Type 2: Add new dimension record

 Consider employee 'Ahmed' changed his city from Poznan to Brussels in November 2015

EmpKey	EmpID	EmpName	EmpCity	FromDate	ToDate
1	111	John	Poznan	01/01/2015	31/12/9999
2	222	Doe	Wroclaw	01/01/2015	31/12/9999
3	333	Ahmed	Poznan	01/01/2015	31/10/2015
4	333	Ahmed	Brussels	01/11/2015	31/12/9999

Employee table using Type 2 Changes

- <u>EmpKey</u> is a surrogate key
- <u>FromDate</u> and <u>ToDate</u> together represent the validity period of the record version
- <u>ToDate</u> value of the current record is set to a date in distant future, i.e. 31/12/9999

© Robert Wrembel (Poznan University of Technology, Poland)



# Type 2: Add new dimension record

- It keeps the full track of history, but
- If the changes are frequent the table size may grow very fast
- Storage and performance may become a concern
- Complicate ETL processes

© Robert Wrembel (Poznan University of Technology, Poland)

AUTO AUTO AUTO

# Type 3: Add a new field

- There are two columns for a particular attribute of interest to represent it's historic (prior) and current values
- Recall our example of address change

EmpID	EmpName	CurrentCity	PreviousCity	EffectiveDate
111	John	Poznan	Poznan	01/02/2015
222	Doe	Wroclaw	Poznan	01/04/2015
333	Ahmed	Brussels	Poznan	01/11/2015

Employee dimension using type 3 changes

© Robert Wrembel (Poznan University of Technology, Poland)



- Type 3 technique does not provide the complete history
  - if the address changes again, the information that the employee has lived in 'Poznan' will be lost
- Upon every change, the values in the current and history columns are overwritten which provokes all the caveats in the case of Type 1 changes

© Robert Wrembel (Poznan University of Technology, Poland)

HUNNEA PORTAN

### **SCDs: Summary**

SCD Type	Dimension Table Action	Impact on Fact Analysis
Туре 0	No change	Facts associated with the attribute's original value
Type 1	Overwrite attribute value	Facts associated with the attribute's current value
Type 2	Add new dimension row for each member with new value	Facts associate with the attribute value in effect when facts occurred
Туре 3	Add new column to preserve attribute's current and historic value	Facts associate with both current and historic attribute value
Type 4	Add mini-dimension	Facts associate with rapidly changing attributes when facts occurred

© Robert Wrembel (Poznan University of Technology, Poland)



# **SCDs: Summary**

SCD Type	Dimension Table Action	Impact on Fact Analysis
Type 5	Add type 4 mini-dimension along with overwritten type 1 mini- dimension key in the base dimension	Facts associate with the attribute's value at the time facts occurred, plus rapidly changing attributes' current value
Type 6	Add type 1 attribute to type 2 base dimension and overwrite all prior values on change	Facts associate with the attribute's value at the time facts occurred, plus the current value
Type 7	Add type 2 dimension row with new values plus the view limited to the current rows and/or attribute values	Facts associate with the attribute's value at the time facts occurred, plus the current value

© Robert Wrembel (Poznan University of Technology, Poland)