POZNAN UNIVERSITY OF TECHNOLOGY

# Modeling Data Warehouse
## Part 1

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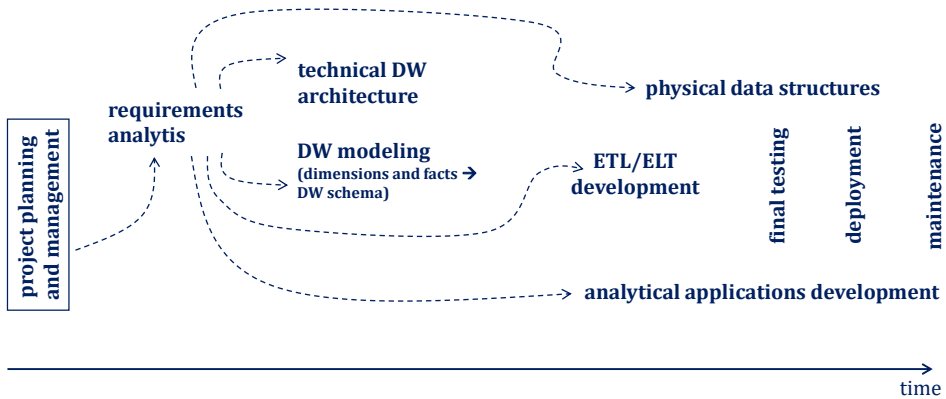www.cs.put.poznan.pl/rwrembel

---

# Outline

- ➲ **Data warehouse modeling pipeline**
- ➲ **Conceptual mulitdimensional data model**
- ➲ **Logical multidimensional data models**
- ➲ **Relational DW schemas**

# DW development pipeline

# DW Data models

- ⊃ **Conceptual**
  - ▪ **mutlidimensional data model (MD)**
- ⊃ **Logical (implemenation)**
  - ▪ **relational →ROLAP**
  - ▪ **multidimensional → MOLAP**
  - ▪ **hybrid → HOLAP**
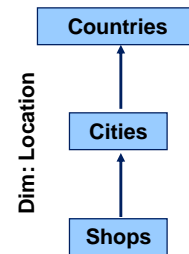
# Conceptual DW data model

⊃ **Only two 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**

Dim: Location

| Countries |
|---|

↑

| Cities |
|---|

↑

| Shops |
|---|

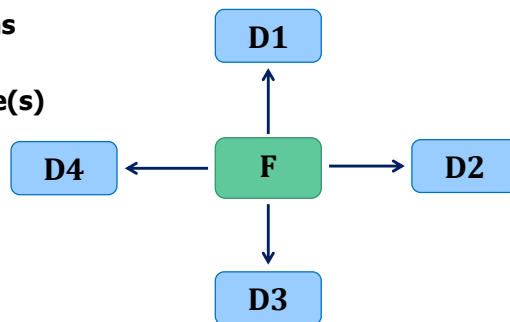# Logical DW data model

⊃ **ROLAP**

- **star schema**
- **snowflake schema**
- **star-flake schema**
- **fact constellation schema**

# 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)**

```
            D1
            ↑
  D4  ←     F     →  D2
            ↓
            D3
```

# 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 $8*10^6$ records
  - each customer makes daily 2 phone calls
  - one year time span of a fact table
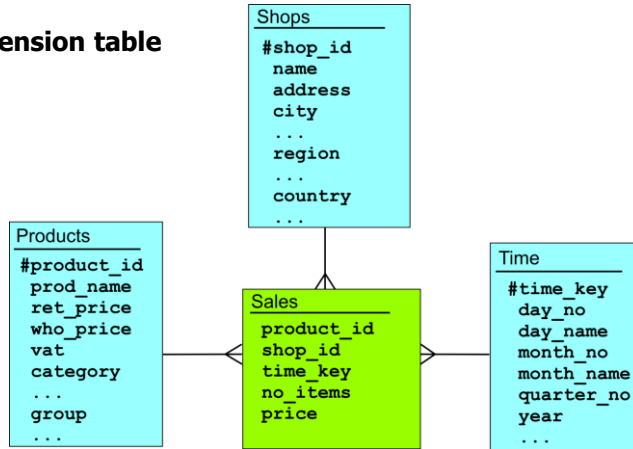  - decreasing the length of each fact row by 10B → size decrease by 54GB

# ROLAP - star schema

- ⊃ **One fact table**
  - ▪ **measures**
- ⊃ **At least one dimension table**

**Shops**
```
#shop_id
 name
 address
 city
 ...
 region
 ...
 country
 ...
```

**Products**
```
#product_id
 prod_name
 ret_price
 who_price
 vat
 category
 ...
 group
 ...
```

**Sales**
```
 product_id
 shop_id
 time_key
 no_items
 price
```

**Time**
```
#time_key
 day_no
 day_name
 month_no
 month_name
 quarter_no
 year
 ...
```

# Dimension table

| | Dimension table | Fact table |
|---|---|---|
| artificial ID | ID (PK) | Dim1 (FK) |
| code used by business | BCode | Dim2 (FK) |
| natural ID (may come from a source system) | ATR1 (NK) | ..... |
| | ..... | Dimn (FK) |
| | ATRn (NK) | M1 |
| descriptors | Atr1 | M2 |
| | Atr2 | ..... |
| | ..... | Mn |
| | Atrn | |

measures

# ROLAP - snowflake schema

Countries
#country_id
name
...

Regions
#region_id
name
...
country_id

Cities
#city_id
name
...
region_id

Shops
#shop_id
name
address
...
city_id

➲ **Dimension**
  ➲ **levels (level tables)**
  ➲ **hierarchies**

Groups
#group_id
group
...

Categories
#categ_id
category
...
group_id

Products
#product_id
prod_name
ret_price
who_price
vat
categ_id

Sales
product_id
shop_id
time_key
no_items
price

Days
#time_key
day_no
day_name
...
month_id

Months
#month_id
month_no
month_name
...
quarter_id

Quarters
#quarter_id
...
year

Years
#year
...

# ROLAP - star-flake schema

Shops
#shop_id
name
address
city
...
region
...
country
...

Groups
#group_id
group
...

Categories
#categ_id
category
...
group_id

Products
#product_id
prod_name
ret_price
who_price
vat
categ_id

Sales
product_id
shop_id
time_key
no_items
price

Days
#time_key
day_no
day_name
...
month_id

Months
#month_id
month_no
month_name
...
quarter_id

Quarters-Years
#quarter_id
...
year
...

# Fact constellation schema

```
                    ┌─────────────┐
                    │ Shops       │        ┌──────────────┐
                    │ #shop_id    │        │ Marketing    │
                    │  name       │        │  campaign_id │
                    │  address    │        │  product_id  │
                    │  city       │        │  shop_id     │
                    │  ...        │        │  time_key    │
                    │  region     │        │  cost        │
                    │  ...        │        │  ...         │
                    │  country    │        └──────────────┘
                    │  ...        │
                    └─────────────┘
┌──────────┐ ┌─────────────┐ ┌─────────────┐ ┌─────────────┐ ┌─────────────┐ ┌─────────────┐ ┌─────────────┐
│ Groups   │ │ Categories  │ │ Products    │ │ Sales       │ │ Days        │ │ Months      │ │ Quarters-   │
│ #group_id│ │ #categ_id   │ │ #product_id │ │  product_id │ │ #time_key   │ │ #month_id   │ │  Years      │
│  group   │ │  category   │ │  prod_name  │ │  shop_id    │ │  day_no     │ │  month_no   │ │ #quarter_id │
│  ...     │ │  ...        │ │  ret_price  │ │  time_key   │ │  day_name   │ │  month_name │ │  ...        │
│          │ │  group_id   │ │  who_price  │ │  no_items   │ │  ...        │ │  ...        │ │  year       │
│          │ │             │ │  vat        │ │  price      │ │  month_id   │ │  quarter_id │ │  ...        │
│          │ │             │ │  categ_id   │ │             │ │             │ │             │ │             │
└──────────┘ └─────────────┘ └─────────────┘ └─────────────┘ └─────────────┘ └─────────────┘ └─────────────┘
```

# Fact table

**Sales**

| Time_key (FK) |
| --- |
| **Product_id (FK)** |
| **Shop_id (FK)** |
| **Manager_id (FK)** |
| **Promo_id (FK)** |
| **Payment_type_id (FK)** |

⎱ dimensions

| **Transaction_time** |
| --- |

measures ⎰

| **Nb_items_sold** |
| --- |
| **Net_price** |
| **VAT** |

# Factless fact

⊃ **No explicit measure attribute**

⊃ **Stores facts that typically represent events**

**Accident**

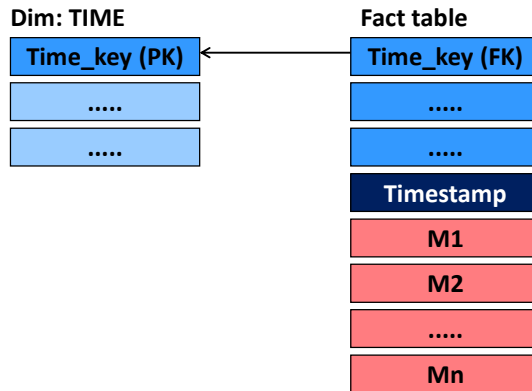| Accident_type_id (FK) |
|---|
| Insurance_NO (FK) |
| Time_key (FK) |

---

# 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** |
| **...** |

# Dimension TIME

➲ **Registering time with granulairty > days**
  ▪ **timestamp in a fact table**

| Dim: TIME | Fact table |
|---|---|
| **Time_key (PK)** ← | **Time_key (FK)** |
| ..... | ..... |
| ..... | ..... |
|  | **Timestamp** |
|  | **M1** |
|  | **M2** |
|  | ..... |
|  | **Mn** |

# Dimension roles

➲ **The same dimension may play different roles (give contexts) for a fact table**
  ▪ **e.g., the TIME dimension**

| Dim: TIME | Fact table |
|---|---|
| **Time_key (PK)** | **Order_time (FK)** |
| ..... | **Delivery_time (FK)** |
| ..... | **Payment_time (FK)** |
|  | **Return_time (FK)** |
|  | **M1** |
|  | **M2** |
|  | ..... |
|  | **Mn** |

# Dimension roles

⮑ **Fact dimension:** dimension in 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

⮑ **Created by an ETL process**
⮑ **ID types**
- **numerical**
  - **efficiency in processing**
  - **chronology represented by values**
  - **typically no semantics**
- alphanumerical
  - **less efficient in processing**
  - **may have semantics**
  - **longer than numerical**
  - **may be constructed as concatenation of natural key and timestamp**

# Which schema to use?

⊃ **Advantages of star schema**
   - **less tables → less joins**
   - **simpler DW loading procedure**
⊃ **Disadvantages of star schema**
   - **bigger tables → more I/O to read**
   - **bigger indexes → more I/O to read**
   - **dimension hierarchies may not be visible**

# Which schema to use?

⊃ **Star schema**
   - **denormalized dimension Time → data redundancy**
   - **1 sec granularity: 24*60*60 times the same date is stored**
   - **1 sec granularity and time span of 10 years ⇨ 300 000 000 rows**
   - **DATE datatype occupies 7B**
     - **lost space: $7B * 300*10^6$**

# Which schema to use?

⊃ **Advantages of snowflake schema**
- **normalized tables → less storage**
- **clearly visible dimension hierarchies**

⊃ **Disadvantages of snowflake schema**
- **more tables → more joins**
- **more complex DW loading procedure**

# Which schema to use?

⊃ **Possible solutions**
- **attributes from different hierarchy levels frequently used together in roll-up → store them in the same level table**
- **high level attributes rarely used → store them in separate high level tables**

⊃ **Star-flake schema**

# Designing DW schema



source driven

user driven

# Designing DW schema

➲ **Mappings between DS and DW**
- **DS1.table1 → DW.Dim_table1**
- **DS1.table2 → DW.Dim_table1**
- **...**
- **DS1.table1.attribute1 → DW.Dim_table1.attributeA**
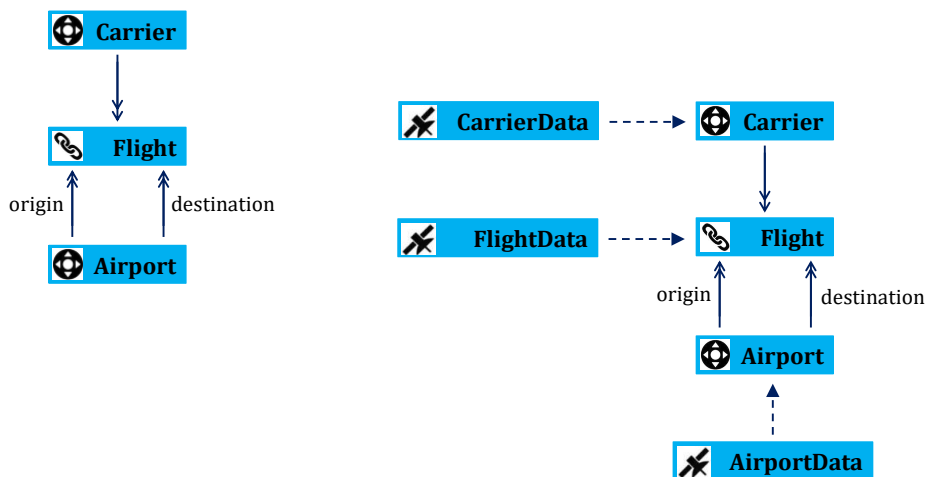- **DS1.table3.attribute5 → DW.F_table1.attributeX**
- **...**

➲ **Model entities**

- **Hub** ⊕
  - **to store business keys of business objects with some other data**
- **Link** 🔗
  - **connects two or more hubs**
  - **may be a candidate for a fact**
  - **may contain its own features (attributes)**
- **Satellite** 🛰
  - **store attributes that belong to a business key (in a hub) or relationship (in a link)**
  - **attached to only one hub or link**

> D. Linstedt, M. Olschimke: Building a Scalable Data Warehouse with Data Vault 2.0. Morgan Kaufman, 2016

# Data Vault Modeling

➲ **Hub content**
- **business key** (value), may be composed of multiple attributes, like a composite primary key
  - e.g., a natural identifier
- **hash key** (used to reference the business object in links and satellites; used in joins; plays the role of a primary key)
  - e.g., an artificial identifier
- **record source**
- **business key load date** (date + time)
- last seen in a source (optional attribute)

⊕ **Carrier**          ⊕ **Airport**

---

# Data Vault Modeling

➲ **Hub examples**
- **Carrier**          ⊕ **Carrier**
  - **CarrierID (e.g., IATA code)**
  - **CarrierHashKey**
  - **LoadDate**
  - **RecordSource**
- **Airport**          ⊕ **Airport**
  - **AirportCode (e.g., IATA code)**
  - **AirportHashKey**
  - **LoadDate**
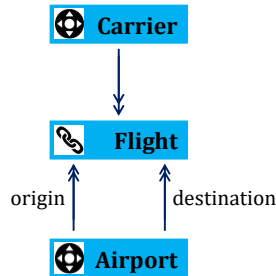  - **RecordSource**

# Data Vault Modeling

⊃ **Link content**
  - **HashKey (counterpart of PK)**
  - **hash keys of connected hubs (counterparts of FKs)**
  - **LoadDate**
  - **RecordSource**
  - **LastSeen**
⊃ **Link example**
  - **Flight**
    - **HashKey**
    - **LoadDate**
    - **RecordSource**
    - **CarrierHashKey**
    - **OriginAirportHashKey**
    - **DestinationAirportHashKey**
    - **LastSeen**

⊕ **Carrier**

🔗 **Flight**

origin       destination

⊕ **Airport**

---

# Data Vault Modeling

⊃ **Satellite**
  - **stores every change to raw data → stores data that evolve in time (like SCDs)**
  - **recommended: to distribute data among various satellites**
    - **split raw data first by source system and**
    - **second by rate of change**

✴ **CarrierData**　　　　✴ **FlightData**

# Data Vault Modeling

⊃ **Satellite content**
- **parent hash key (counterpart of FK; part of PK)**
- **original business features of a satellite**
- **load date (part of PK)**
- **load end date (when the record becomes invalid)**
- **hash diff (hash value of business features; helps in identifying if a record value changed by comparing it with hash of source values)**
- **extract date (optional; when a record was ingested from a source)**

# Data Vault Modeling

⊃ **Satellite example**
- **FlightData**
  - **FlightHashKey (FK, PK)**
  - **DepartureTime**
  - **ArrivalTime**
  - **Distance**
  - **LoadDate (PK)**
  - **LoadEndDate**
  - **RecordSource**
  - **HashDiff**

FlightData ----> Flight