# Traditional Data Warehouse Architectures

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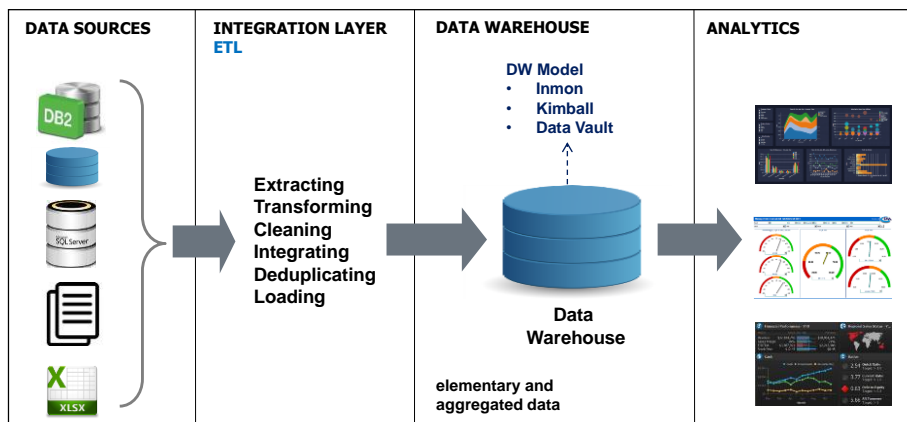**www.cs.put.poznan.pl/rwrembel**

# Outline

⊃ **Data Warehouse architectures**
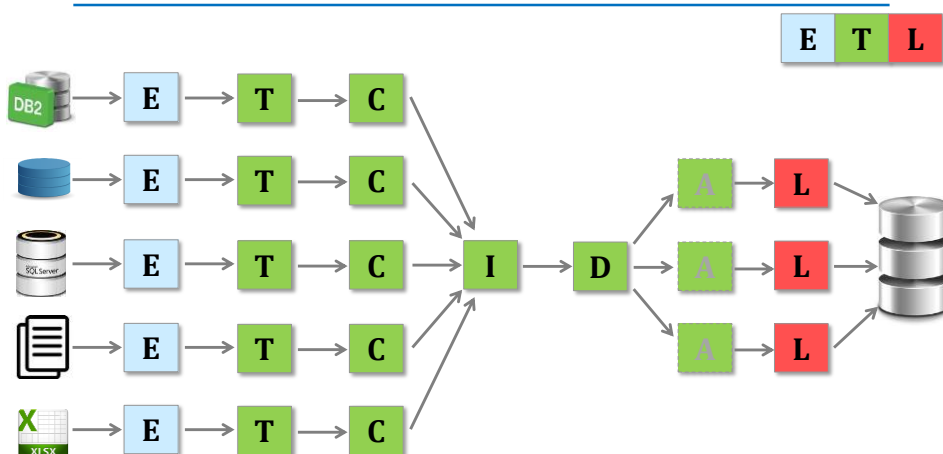⊃ **Data integration and loading: ETL vs. ELT**
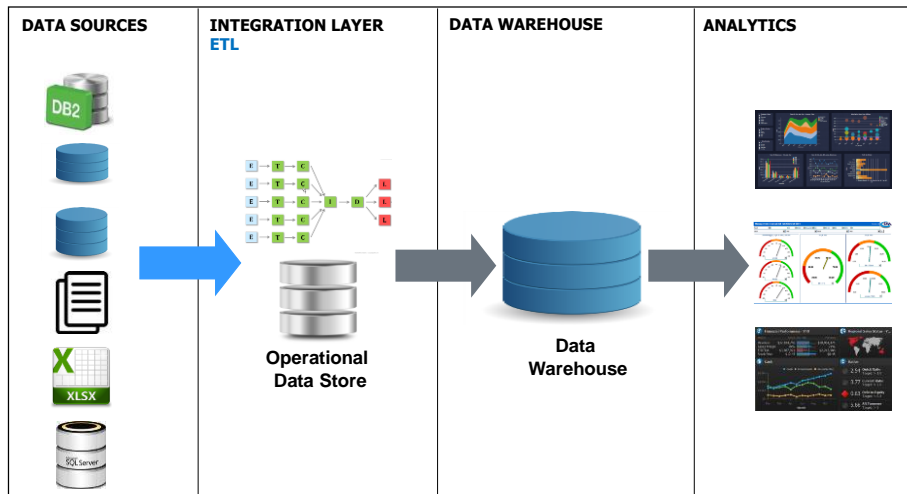
# DW Architecture 1 (basic)

| DATA SOURCES | INTEGRATION LAYER **ETL** | DATA WAREHOUSE | ANALYTICS |
|---|---|---|---|
| DB2 | **Extracting Transforming Cleaning Integrating Deduplicating Loading** | **DW Model** <br> • **Inmon** <br> • **Kimball** <br> • **Data Vault** <br><br> **Data Warehouse** <br><br> elementary and aggregated data | |

**typically OLTP data sources**

# ETL workflow

| E | T | L |
|---|---|---|

DB2 → E → T → C

→ E → T → C

SQLServer → E → T → C → I → D → A → L

→ E → T → C

XLSX → E → T → C

A → L

A → L

# DW Architecture 2

| DATA SOURCES | INTEGRATION LAYER **ETL** | DATA WAREHOUSE | ANALYTICS |
|---|---|---|---|



Operational Data Store
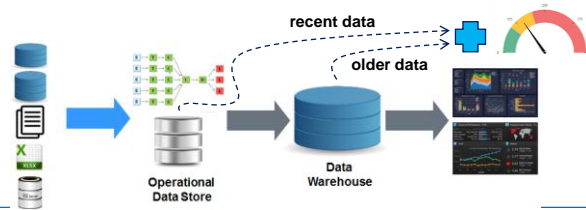
Data Warehouse

# Operational Data Store

- ⊃ = Staging Area = (Data) Stage
- ⊃ A repository for an ETL engine
- ⊃ To separate normal processing at DSs from data ingest
  - ▪ to separate transactional from batch processing
- ▪ Disk storage for processing large data volumes that will not fit in RAM
- ▪ To provide means for data provenance

# Operational Data Store

- ⊃ **To store intermediate results → to be shared (used) by multiple ETL tasks**
  - ▪ **re-using the same result datasets by multiple processes (optimization)**
  - ▪ **for recovery after crash of an ETL process**
    - • **re-executing a stopped process from a failed phase**
- ⊃ **Recent data can be accessed before a DW is refreshed**
- ⊃ **Implementation**
  - ▪ **database**
  - ▪ **(distributed) file system**

# DW Architecture 3

# DW Architecture (cd.)

| DATA SOURCES | ETL | DATA WAREHOUSE | ANALYTICS |
|---|---|---|---|
| DB2 | ETL / ODS | DW | |
| | | DM | |

single server or cluster

main memory appliance
(superserver)

single dedicated server
(typically for MOLAP)
or
the same server as for DW

# Allegro DW

| Production Oracle RAC | Logical Standby Oracle RAC | Load | DWH Staging Area Oracle RAC | ETL | DWH Production Oracle RAC |
|---|---|---|---|---|---|

Oracle Data Guard

Data Cleansing

OLAP  Datamart  Datamining

2 * IBM P590 Machine
(each 8 processors, 16 cores )
DS8300 Array
(56TB user space)

HP BL 460c
(2 Blades X 8 cores)
DS4800 Array
(11TB user space)

HP Oracle Database Machine
(64 Intel processor cores)
HP Oracle Exadata Storage Server
(168 TB of raw storage)
(14 GB/sec data bandwidth)
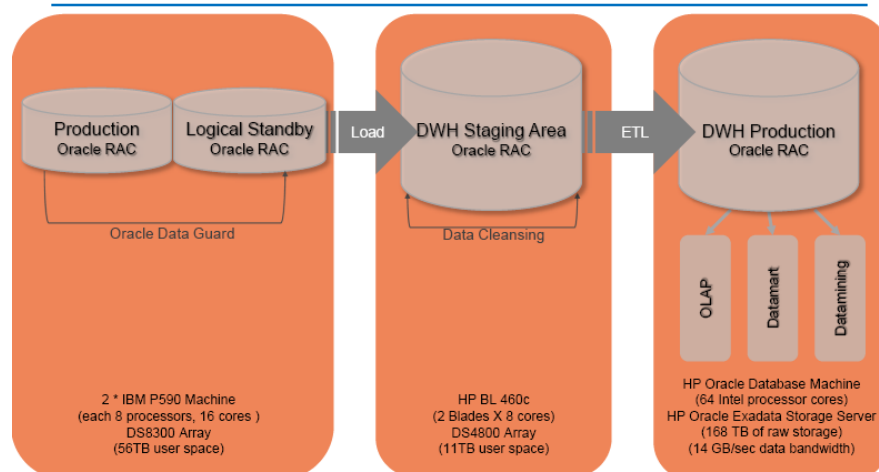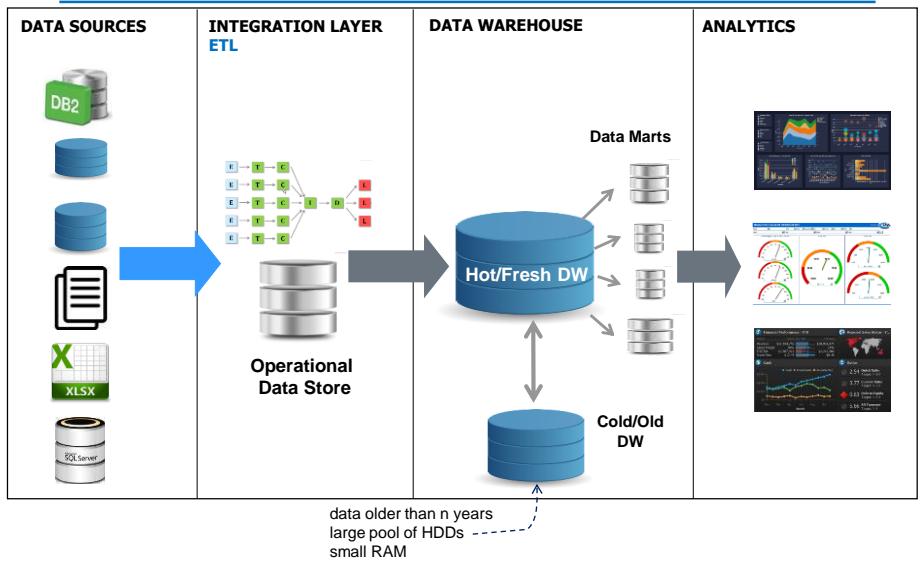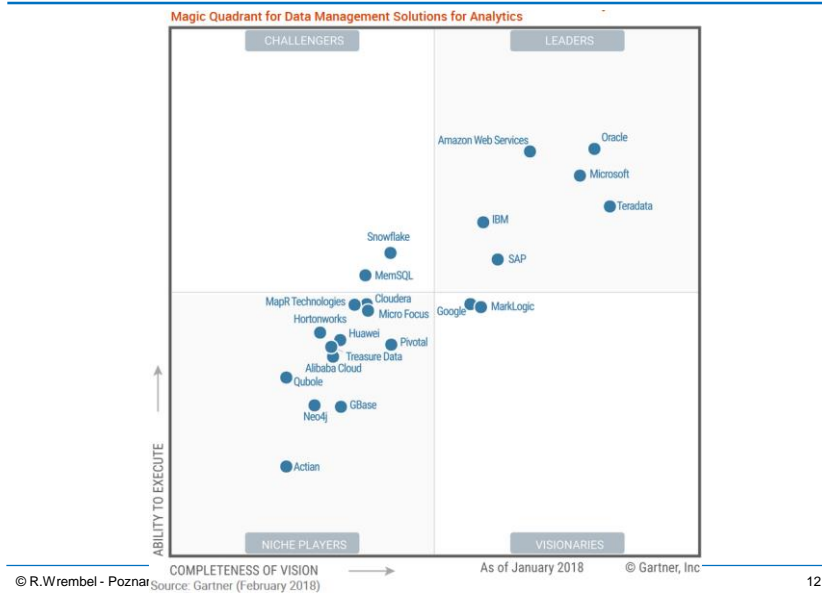
**C. Maar, R. Kudliński: Allegro on the way from XLS based controlling to a modern BI environment. National conference on Data Warehousing and Business Intelligence, Warsaw, 2008**
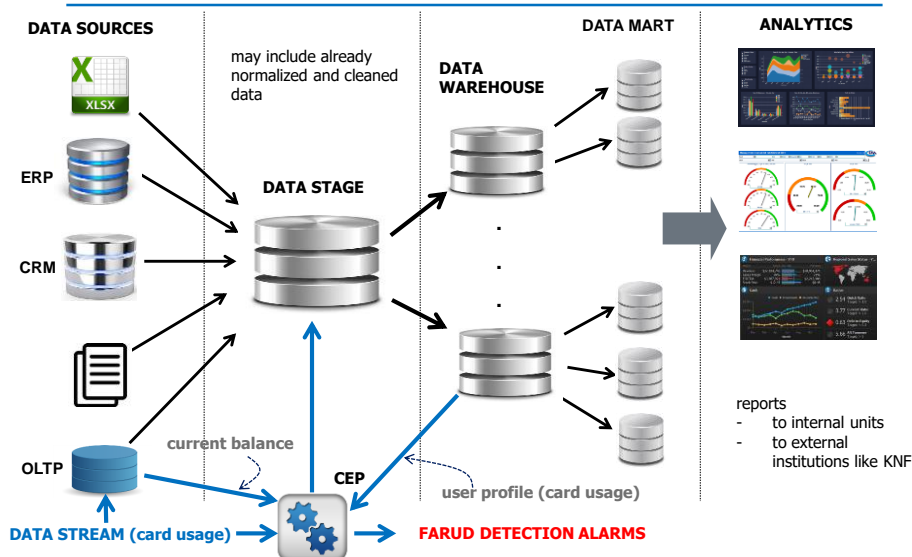
# DW Architecture 4



| DATA SOURCES | INTEGRATION LAYER ETL | DATA WAREHOUSE | ANALYTICS |
|---|---|---|---|

Operational Data Store

Data Marts

Hot/Fresh DW

Cold/Old DW

data older than n years
large pool of HDDs
small RAM

# Gartner Report: DW servers



Magic Quadrant for Data Management Solutions for Analytics

CHALLENGERS     LEADERS

Amazon Web Services    Oracle

Microsoft

IBM    Teradata

Snowflake

MemSQL    SAP

MapR Technologies   Cloudera    Google   MarkLogic
Hortonworks   Micro Focus
Huawei   Pivotal
Treasure Data
Alibaba Cloud
Qubole
   GBase
Neo4j

Actian

ABILITY TO EXECUTE

NICHE PLAYERS     VISIONARIES

COMPLETENESS OF VISION    As of January 2018    © Gartner, Inc

# Large DW Architectures

- ➲ **# data sources: 100 - 200**
- ➲ **Fact table: <span style="color:red">nn</span> * $10^9$ rows**
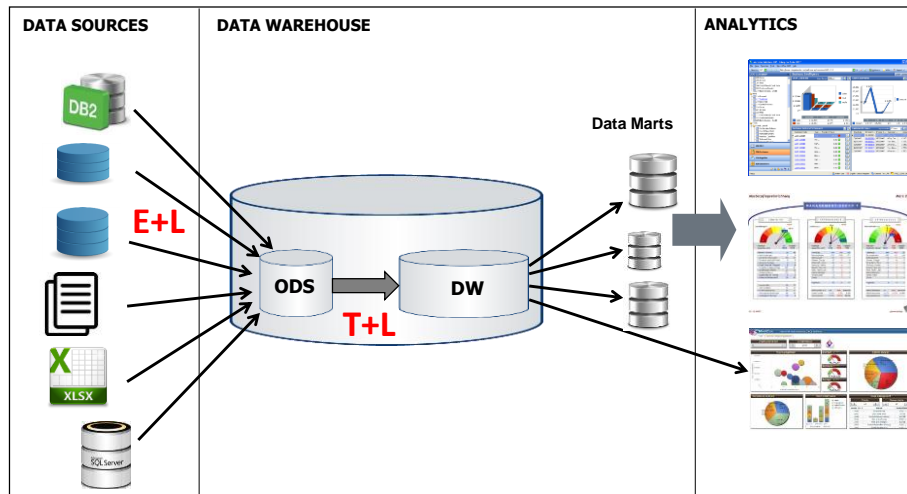- ➲ **Fact table: <span style="color:red">n</span> TB**
- ➲ **Multiple relational DWs in an organization**
  - ▪ **DW size: <span style="color:red">nn</span> TB**
- ➲ **Multiple data marts**
- ➲ **<span style="color:red">n</span> * $10^3$ to <span style="color:red">nnn</span> * $10^3$ ETL workflows**
- ➲ **DW composed of 100+ tables**
  - ▪ **on average 50+ attributes/table**

# DW in Bank

# DW Architecture 5: ELT (ELTL)

# ELT Architecture

➲ **Performance**
  ▪ **data stored in a DB ⇨ processing by means of: SQL, PL/SQL, SQL PL, Transact SQL**
  ▪ **data processed in a DB buffer cache ⇨ native DB environment**
  ▪ **advanced query optimization offered by DBMS**
  ▪ **single server for ELT and HD ⇨ heavier workload**

➲ **Functionality**
  ▪ **data provenance**
  ▪ **drill through**

➲ **Costs**
  ▪ **single DW server**
  ▪ **less software licences (OS, DBMS)**

# ETL vs. ELT (experiment 1)

➲ **Data sources**
- ▪ **topic: Internet auctions**
- ▪ **storage:**
  - **Oracle11g (Object-Relational model)**
  - **MySQL**
  - **PostgreSQL**
  - **XML**

➲ **Data warehouse: Oracle11g**

# ETL vs. ELT (experiment 1)

➲ **DW schema**

# ETL vs. ELT (experiment 1)

➲ **Transformations of data for:**
- **dimensions**
- **fact table**

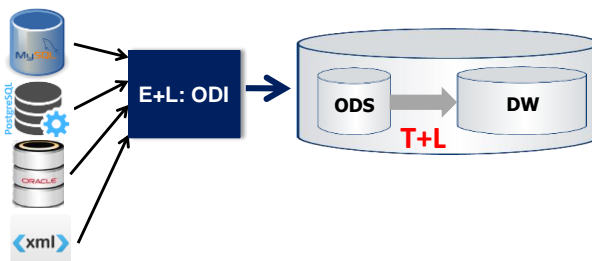➲ **ETL architecture ⇨ Oracle Data Integrator (ODI)**
- **ETL in a staging area on a separate server**

# ETL vs. ELT (experiment 1)

➲ **ELT architecture**
- **T+L in a staging area on the same server as a DW**
- **variant 1: E+L → ODI, T+L → implemented in ODI**
- **variant 2: E+L → ODI, T+L → implemented as materialized views (MVs)**
- **variant 3: E+L → ODI, T+L → implemented as stored packages (SPs)**
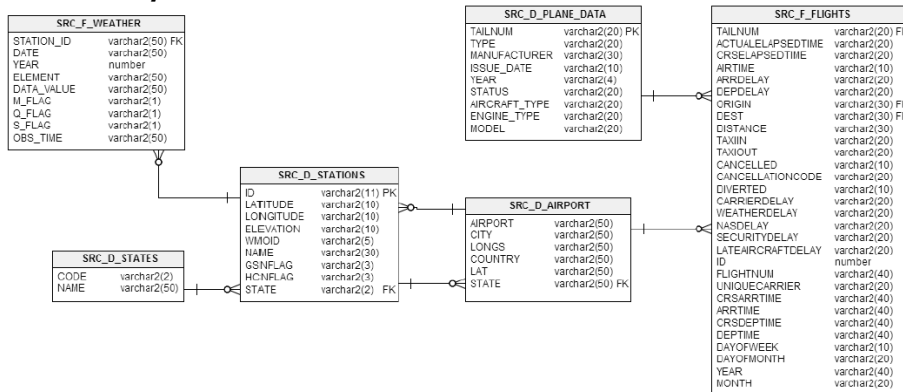- **variant 4: E+L → ODI, T+L → SPs + MVs**

# ETL vs. ELT (experiment 1)

# ETL vs. ELT (experiment 2)

➲ **Data source**
- ▪ **flight and weather data in the US, from 1986 until 2008**
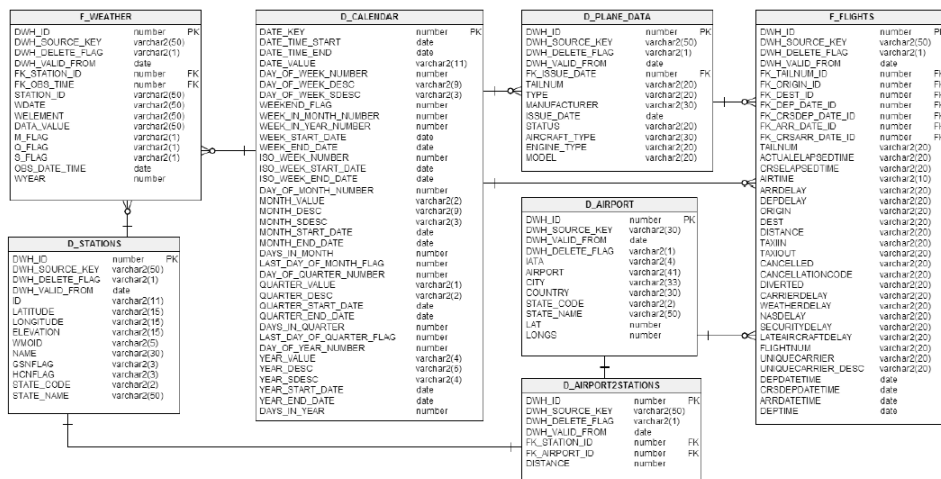- ▪ **6 tables in Oracle11g**

➲ **Data warehouse: Oracle11g**

➲ **ETL/ELT: Informatica**

# ETL vs. ELT (experiment 2)

➲ **DW schema (augmented with: calendar, airplane data, airport data)**

---

# ETL vs. ELT (experiment 2)

➲ **ETL ⇨ Informatica**
➲ **ELT ⇨ Informatica (E+L), DB views (T+L)**

# ETL vs. ELT (experiment 2)

➲ ETL ⇨ Informatica
➲ ELT ⇨ Informatica (E+L), DB views (**T+L**)

# Architecture for Data Science



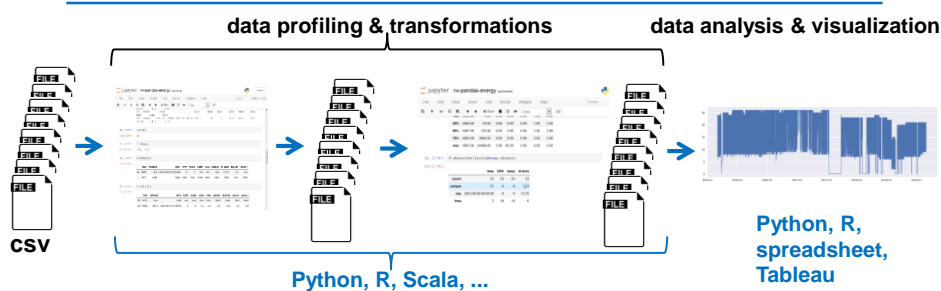data profiling & transformations      data analysis & visualization

CSV

Python, R, Scala, ...

Python, R, spreadsheet, Tableau

➲ **Data stored in files**
- **performance problem**
- **no backup & recovery**
- **no access control**

➲ **Data & code sharing is difficult**
- **re-usability problem**
- **low programming productivity**