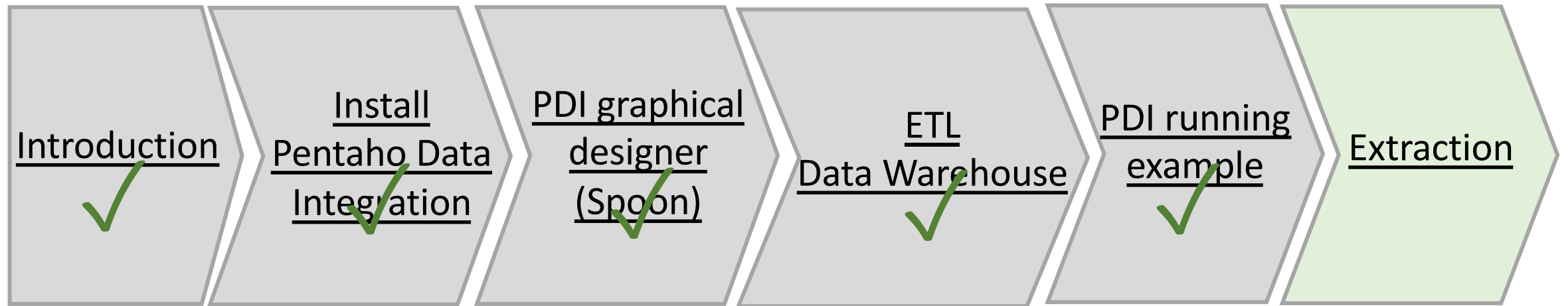


DATA ANALYTICS (Data Warehouse) Pentaho Data Integration

Luca Cinelli, PhD
luca.cinelli@unical.it

Outline



Extraction

Data sources



Example: Datasets

1. Sales Data
2. Customer Data
3. Product Data

Example: ETL problem scenario

1. Extraction

Customer data

CSV
TXT
Excel
Folder
Zip file

Product Data

XML
JSON

Sales Data

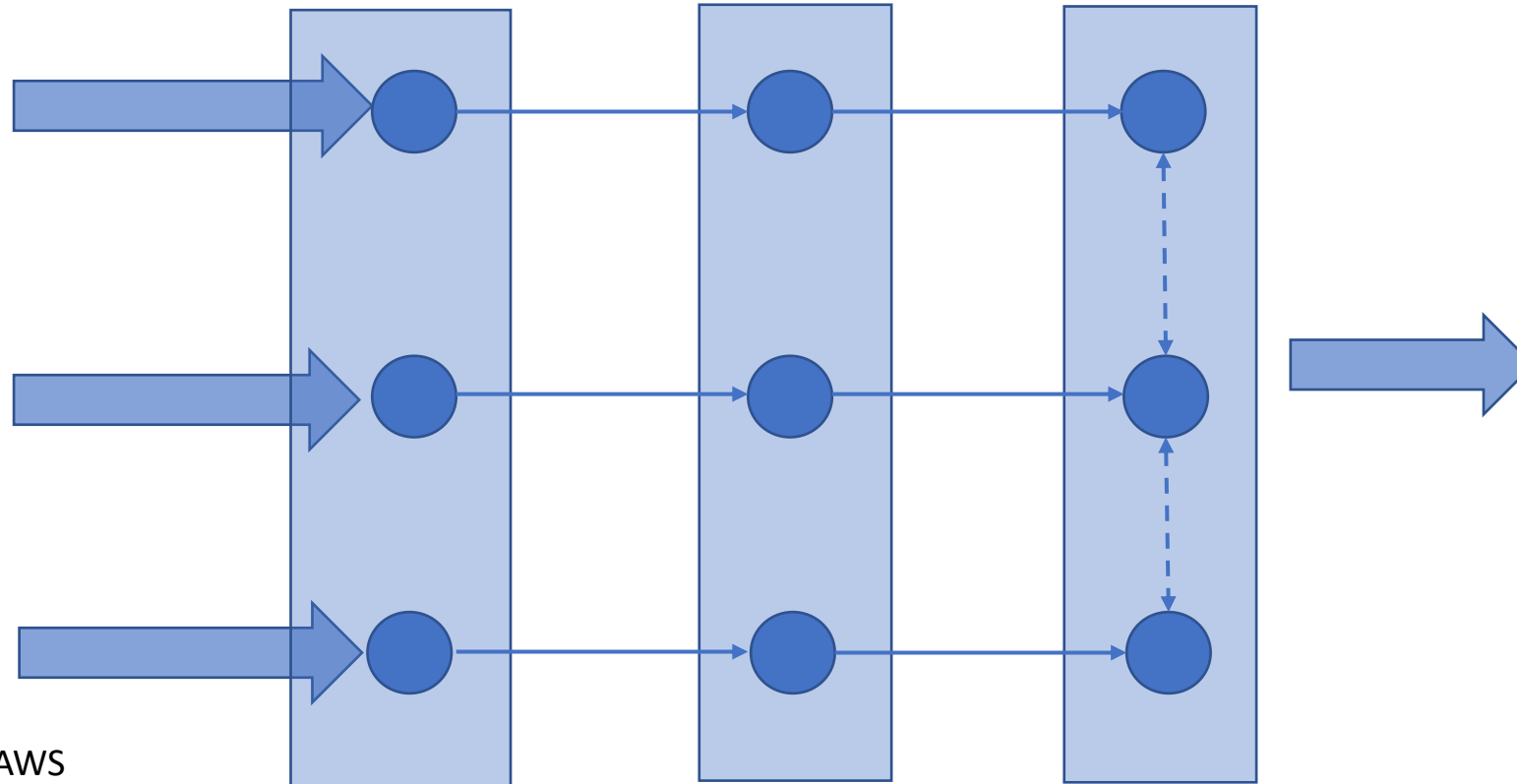
SQL based DB
Cloud based Storage AWS

3. Cleaning and Validation

2. Merging

4. Transformation, aggregation and Joining

5. Loading to create a DataMart



Extraction from tabular format

- Manually entering data into PDI (ManualInput.ktr)



Data grid

Step name: Manual Input

Meta Data

#	Name	Type	Format	Length	Precision	Currency	Decimal	Group	Null if	Set empty string?
1	Customer ID	String		10						N
2	Customer Name	String		50						N
3	Segment	String		20						N
4	Age	String		10						N
5	Country	String		25						N
6	City	String		50						N
7	State	String		50						N
8	Postal Code	String		10						N
9	Region	String		10						N

Help OK Preview Cancel

Extraction from tabular format

- Manually entering data into PDI
- Import data from text file (TxtInput.ktr, input: CustomerData_Central.txt)

Extraction from tabular format

- Manually entering data into PDI
- Import data from text file
- Import data from multiple CSV file (MultipleFiles.ktr, input: customer_data_multiple_files)

Extraction from tabular format

- Manually entering data into PDI
- Import data from text file
- Import data from multiple CSV file version 2
(MultipleFiles_withGetFileName.ktr, input:
customer_data_multiple_files)

Extraction from tabular format

- Manually entering data into PDI
- Import data from text file
- Import data from multiple CSV file
- Import data from excel file (ExcelInput.ktr, input: CustomerData_East.xlsx)

Extraction from tabular format

- Manually entering data into PDI
- Import data from text file
- Import data from multiple CSV file
- Import data from excel file
- Extract data from zip file (ZipInput.ktr, input: CustomerData_South.zip)

Extraction from no-tabular format

- Extract data from XML file (XMLInput.ktr, input: ProductDataAsXML.xml)

Extraction from no-tabular format

- Extract data from XML file
- Extract data from JSON file (JSONInput.ktr, input: ProductDataasJSON.js)