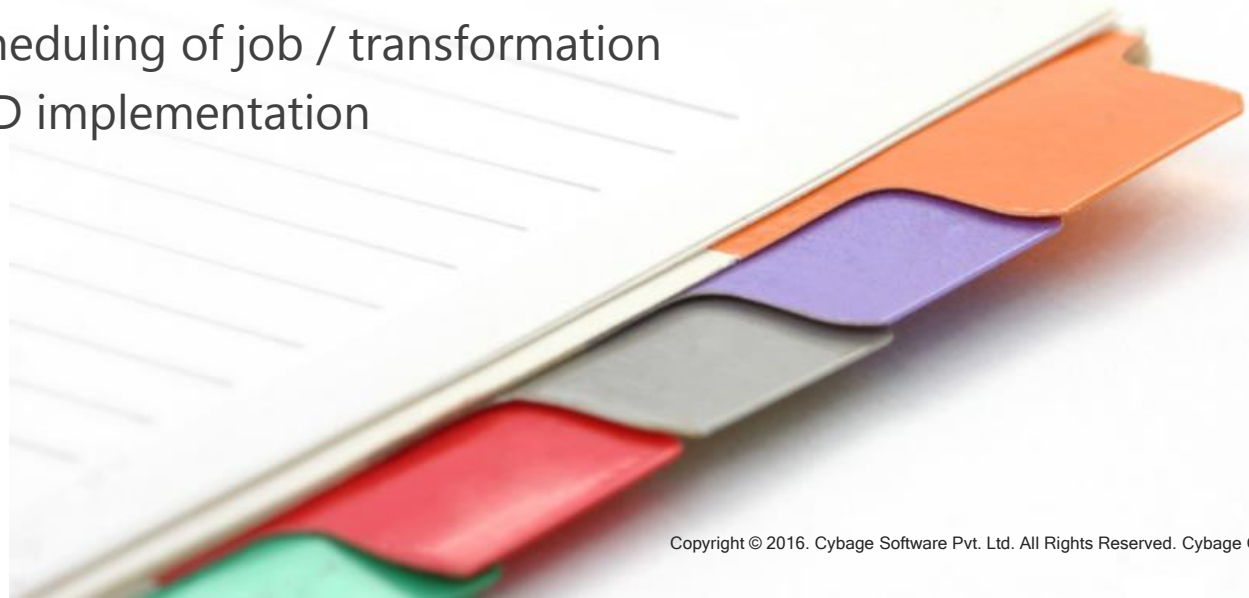# Pentaho Kettle Training

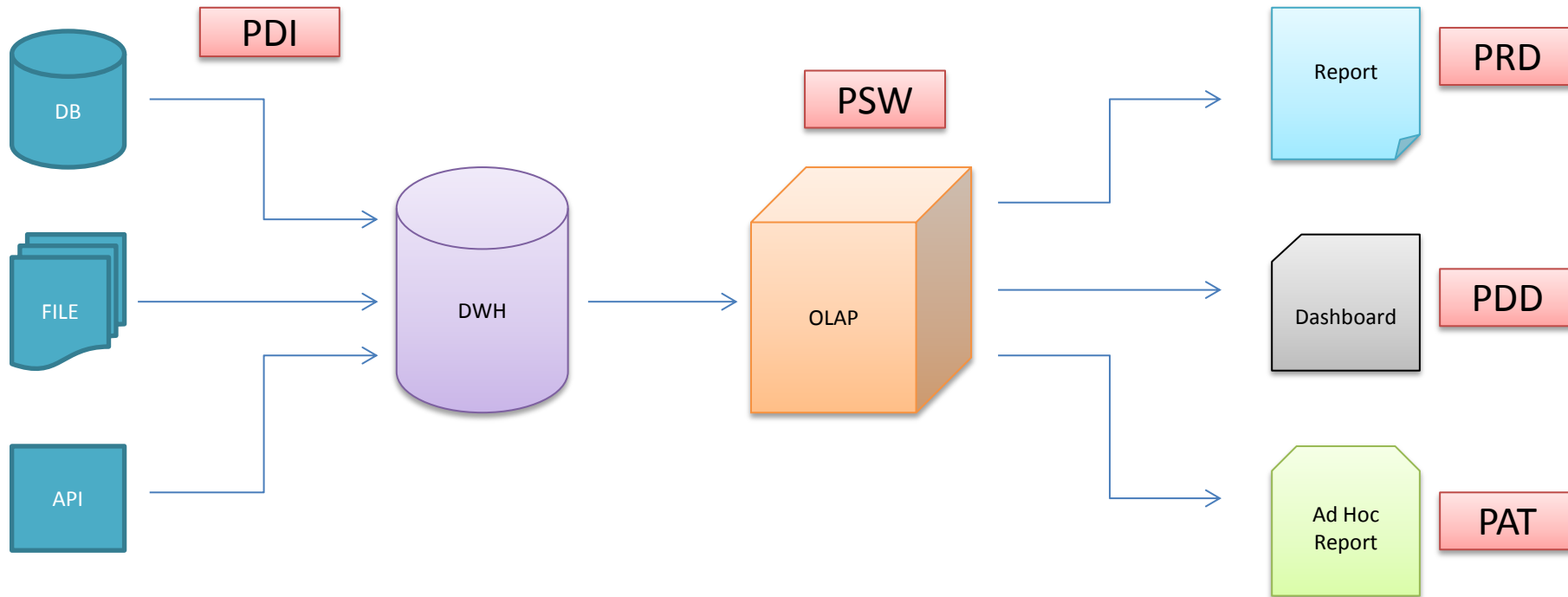Presented by   : Prashant Tikone

# Agenda

- Introduction to Pentaho BA suite
- Various component of PDI-Kettle
- Dealing with repository
- Transformations
- Understanding PDI Steps in Transformation
- Jobs
- Variables and Parameters
- Scheduling of job / transformation
- SCD implementation

# Introduction to Pentaho BA suite

- Pentaho Data Integration

- Pentaho Report Designer

- Pentaho Schema Workbench

- Pentaho BA Server

- Pentaho Dashboard Designer

- Pentaho Analyzer Tool

# ETL Architecture and Pentaho
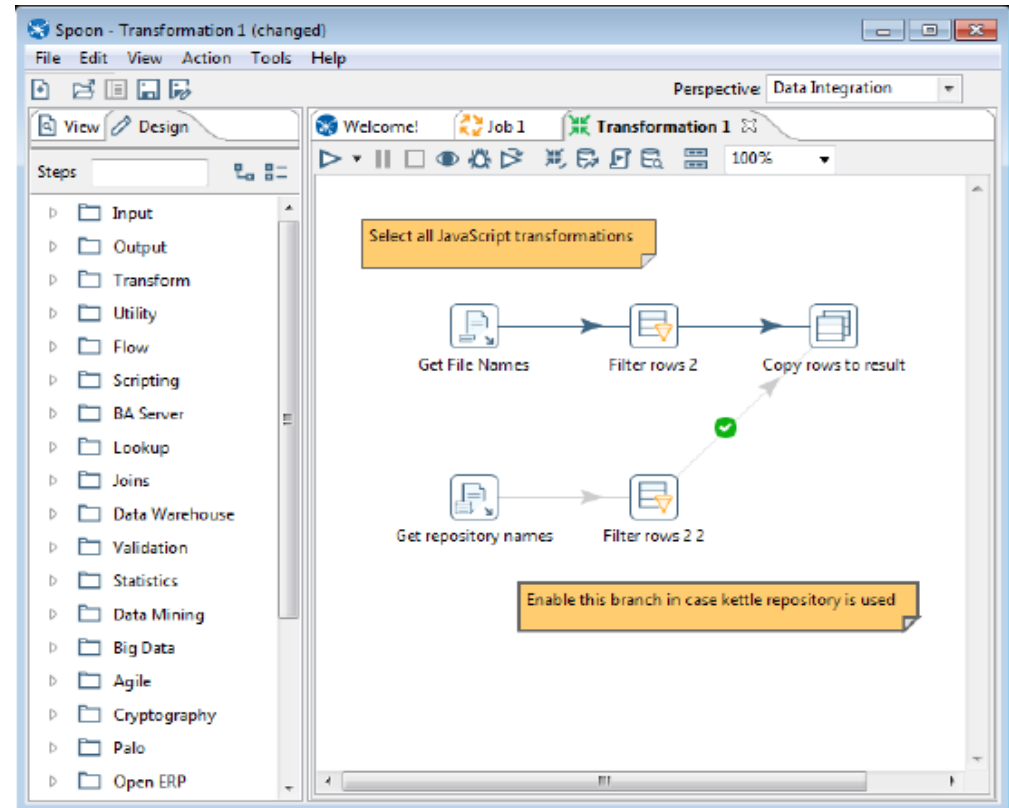
# Various components of PDI

- Data Integration Server

- Spoon

- Pan

- Kitchen

- Carte

# Data Integration Server

- Centrally Store for transformations and jobs

- Pentaho Repository

- Processing engine

- Security and authentication,

- Scheduling.

# Spoon

- User Interface to design the Pentaho Jobs

- Drag and Drop interface

- Uses library of more than 300 pre-built transformations

- Build Workflows using a series of data integration processing entries.

# Pan, Kitchen & Carte

- Pan - Execute PDI transformations, which represent independent data processing tasks

- Kitchen – Execute PDI jobs, which contain transformations and other job entries as part of a larger business process.

- Carte - Set up cluster of PDI servers.  Helps to execute data transformations within a cluster of Carte cluster nodes.

# Hardware Requirements

## Pentaho Data Integration Server

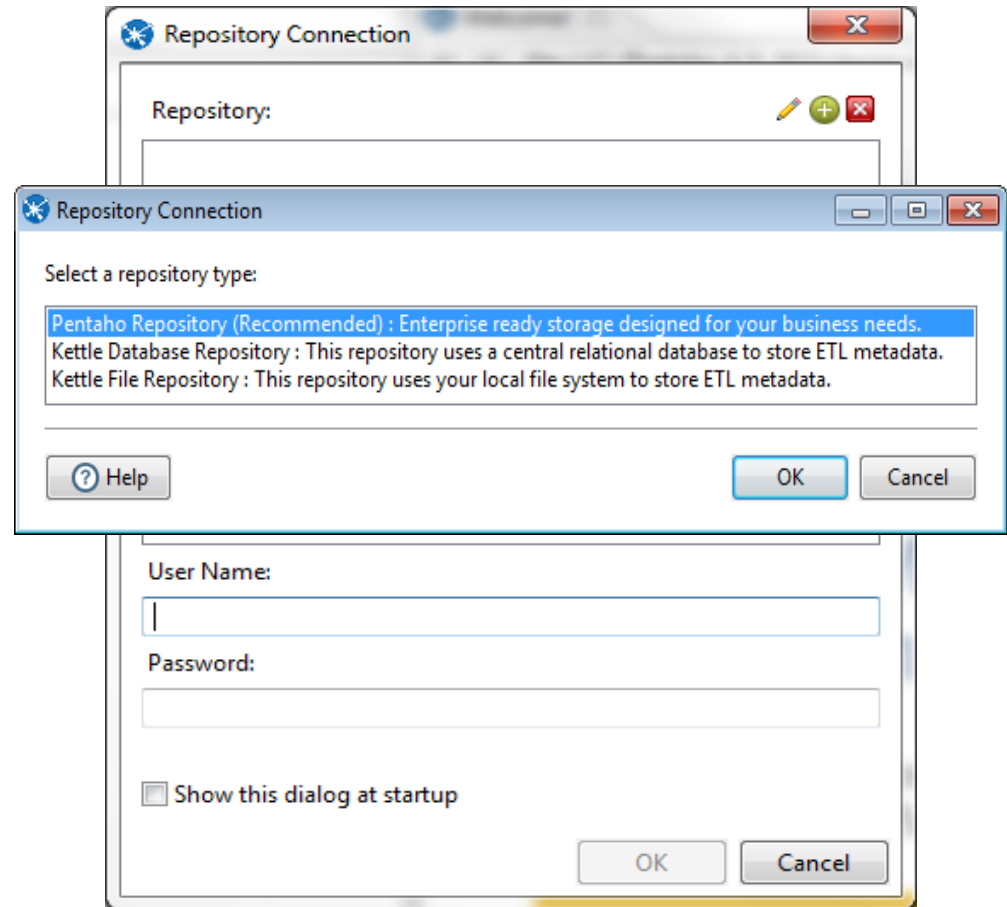| Hardware—64 bit | Operating System—64 bit |
|---|---|
| •Processor: Apple Macintosh Pro Quad-Core or Macintosh Mini Quad-Core<br>•Intel EM64T or AMD64 Dual-Core<br>RAM: 8 GB with 4 GB dedicated to Pentaho servers<br>Disk Space: 20 GB free after installation | •Microsoft Windows 2008 Server R2 & 2012 Server<br>•CentOS 6 & 7<br>•Red Hat Enterprise 6 & 7<br>•Ubuntu Server 12.04 LTS & 14.04 LTS<br>•OSX 10.10 & 10.11<br>•Suse Linux SLES 11 (SP3+) |

## Pentaho Data Integration – Spoon

| Hardware—64 bit | Operating System—64 bit |
|---|---|
| •Processors: Apple Macintosh Dual-Core<br>•Intel EM64T or AMD64 Dual-Core<br>RAM: 2 GB RAM for most of the design tools, PDI requires 2 GB dedicated Disk Space: 2 GB free after installation<br>Minimum Screen Size: 1280 x 960 | •Microsoft Windows 7 & 10<br>•Ubuntu Desktop 12.04 LTS & 14.04 LTS<br>•OSX 10.10 & 10.11<br>•iOS 8.x |

# Repository

What is repository?

- Meta data Storage
- Provides Revision history
- Track changes
- compare revisions

- Types:
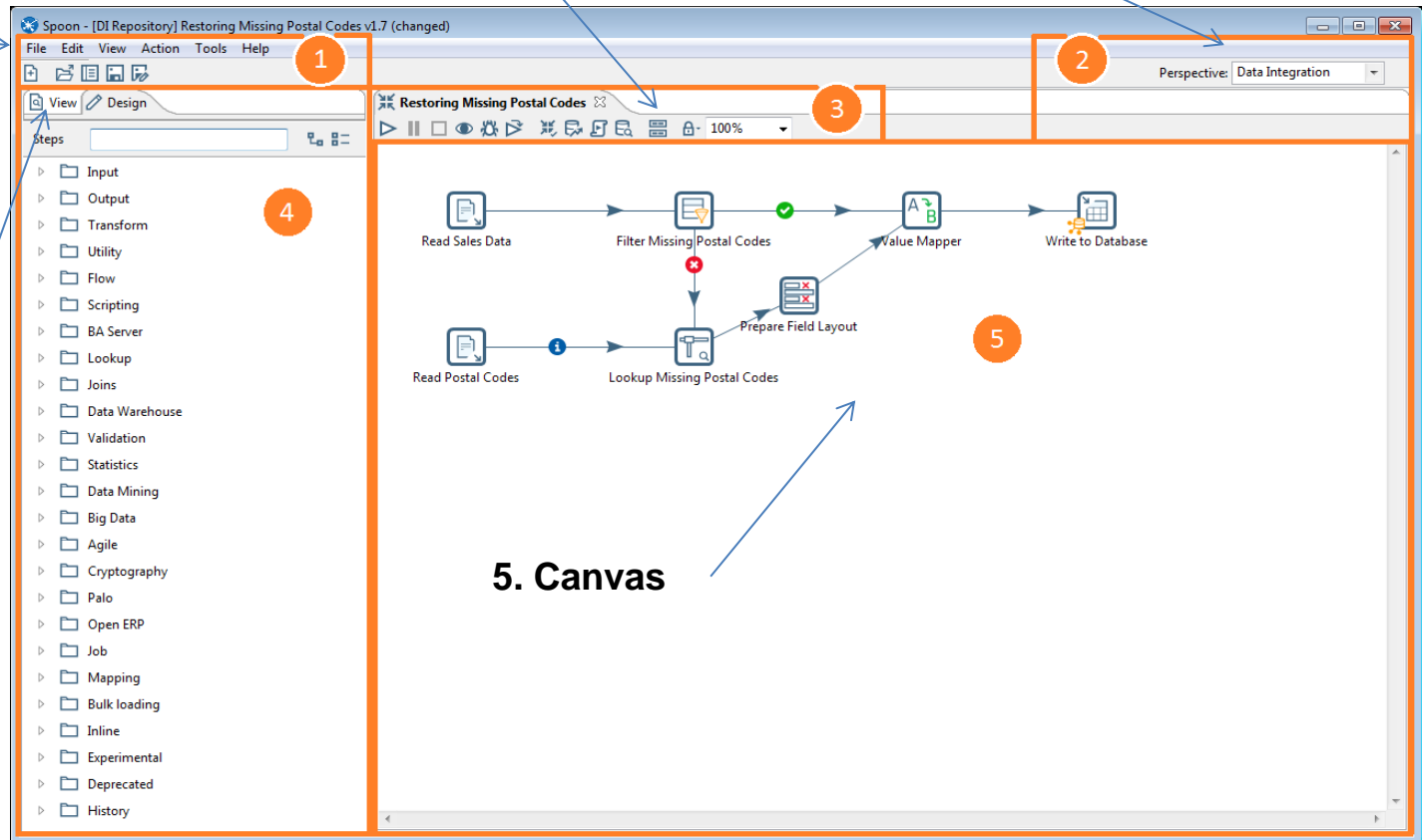  1. File
  2. Database
  3. Enterprise Repository

# Spoon

**3. Sub- Toolbar**

**2. Perspective**
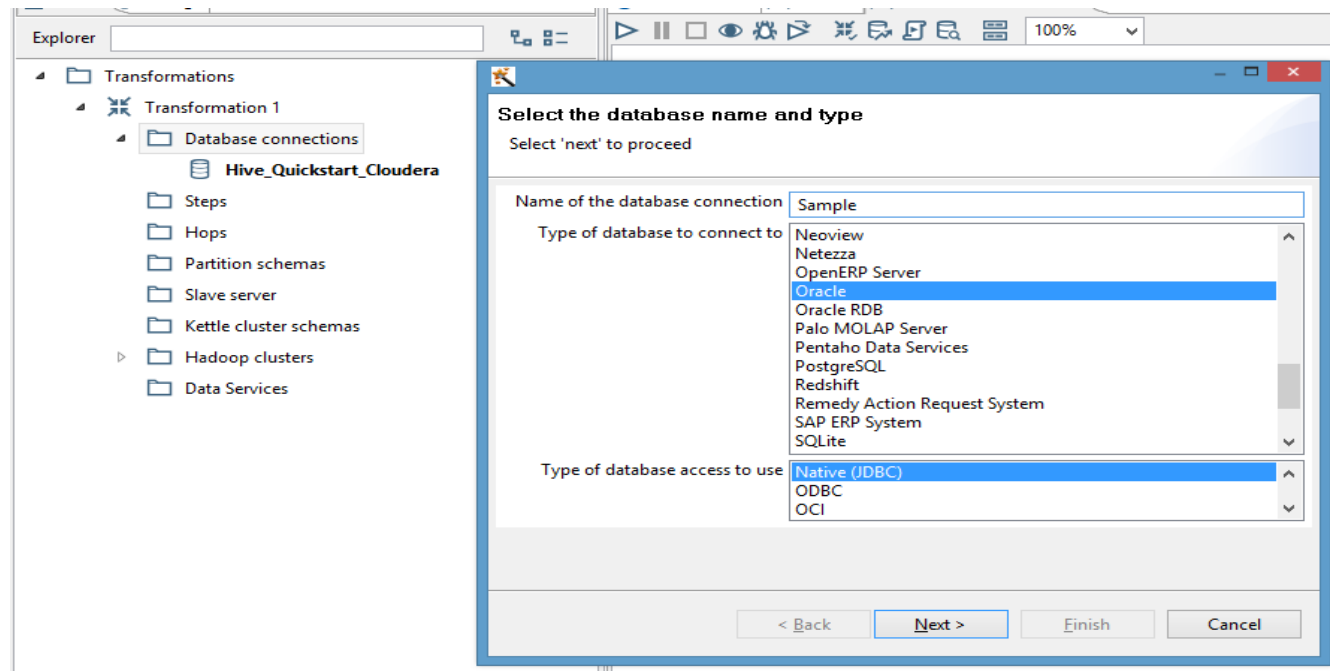
**1. Toolbar**

**4. Design & View Tabs**
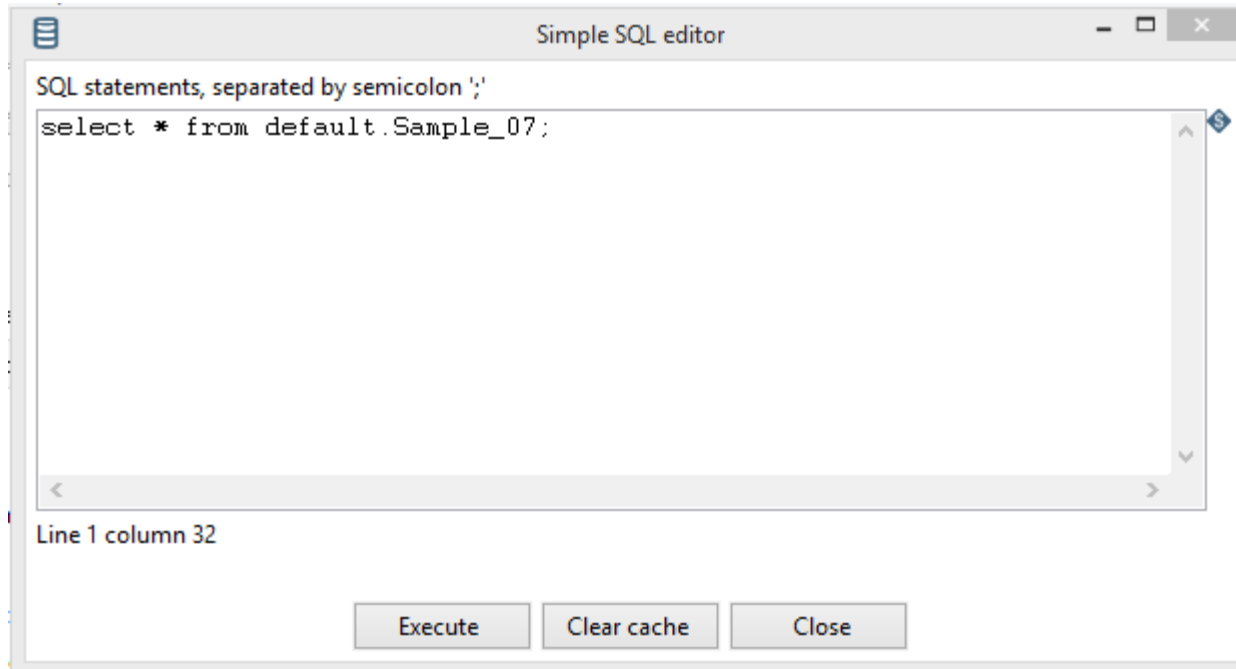
**5. Canvas**

# Pentaho Workflow

- Input steps

- Output steps

- Transformation steps

- Flow controls available in PDI

- Lookup data at various sources
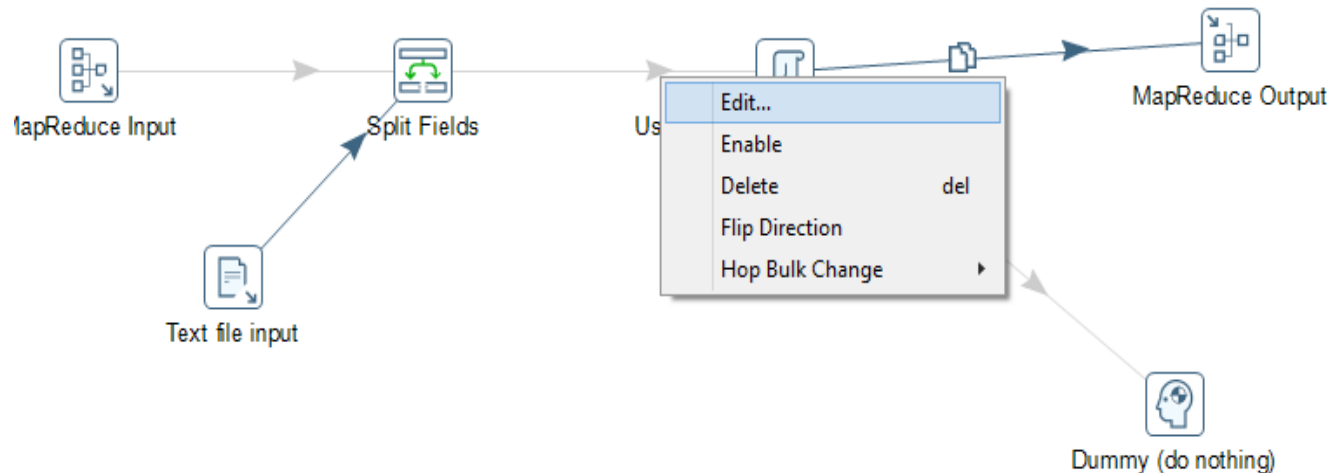
- Data Validation

# Database Connection

# SQL Editor

# Hops

- Hop connects one transformation step or job with another.
- Direction of the data flow is indicated with an arrow.
- A hop can be enabled or disabled.

# Variables

**Definition**

- Set Variable step in a transformation
- Using kettle.properties file in the directory
- Syntax UNIX -  ${VARIABLE}  Windows -  %%VARIABLE%%

**Types**

**Environment variables**

- set an environment variable
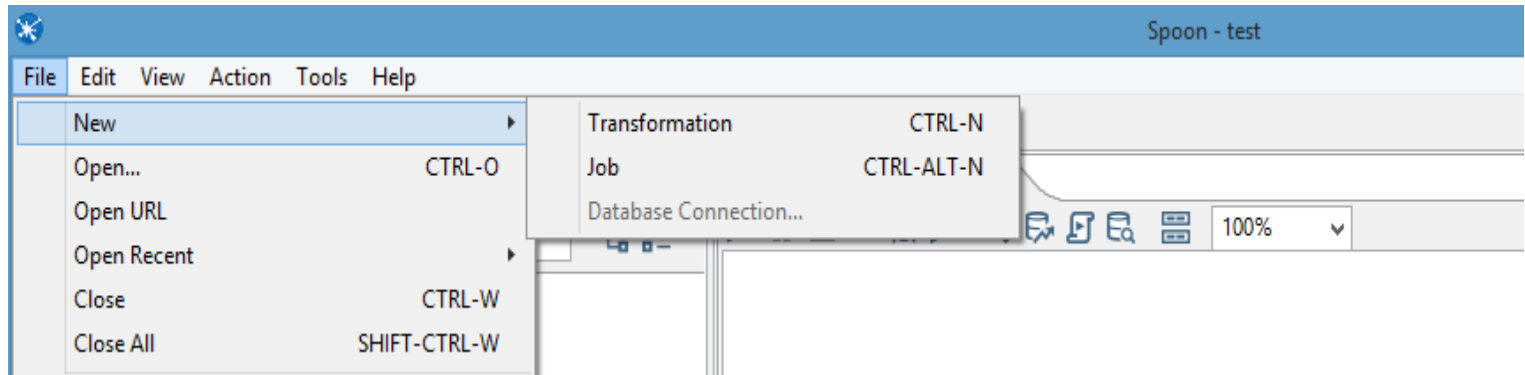- Java Virtual Machine (JVM) with the -D option

**Kettle variables**

- local to the job
- "Set Variable" step in a transformation

**Internal variables**

- Internal.Kettle.Version
- Internal.Job.Name
- Internal.Job.Filename.Name

# Transformation

# Transformations
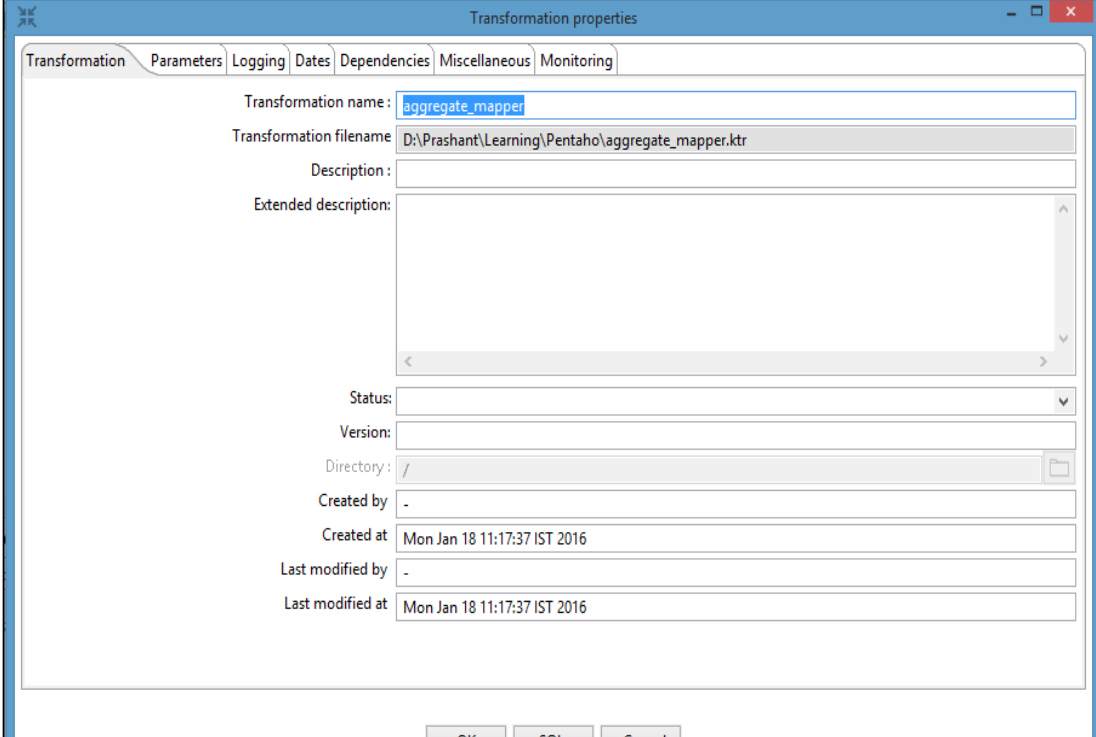
Transformation Tab

Parameters
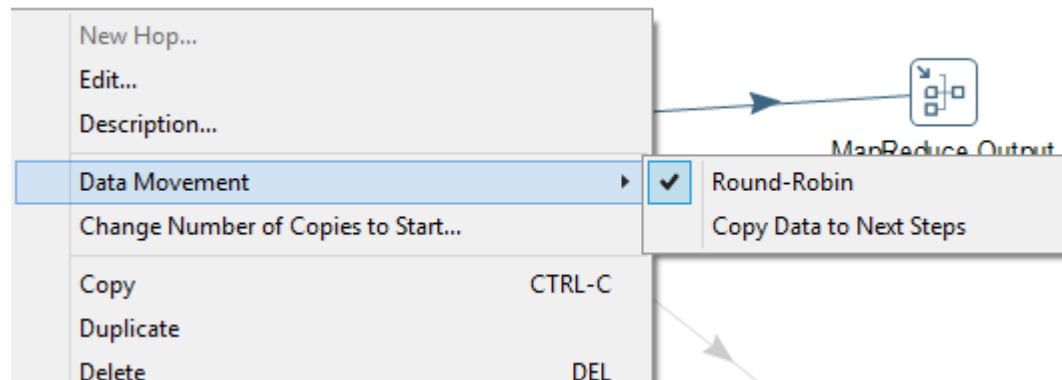
Logging

Dates

Dependencies

Miscellaneous

Monitoring

# Transformation Steps

**Change number of copies to start**
• Launch same step several times to minimize the latency.
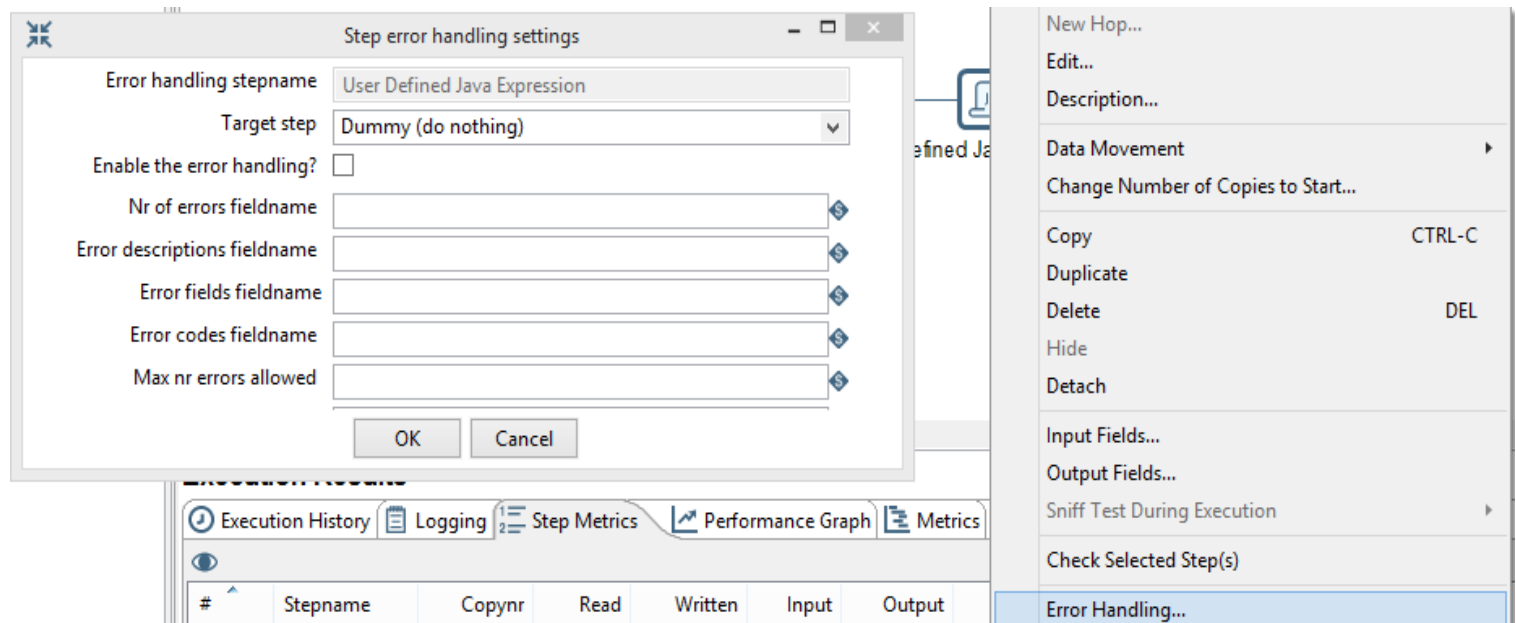
**Distribute or Copy the data –**
• By Default Round-Robin
• Copy data option will copy the data to all target steps

# Transformation Steps

**Step error handling settings –**

Allows you to configure a step so that instead of halting a transformation when an error occurs, the rows that caused an error are passed to a different step.

## Demos

- Transformations
- Joins
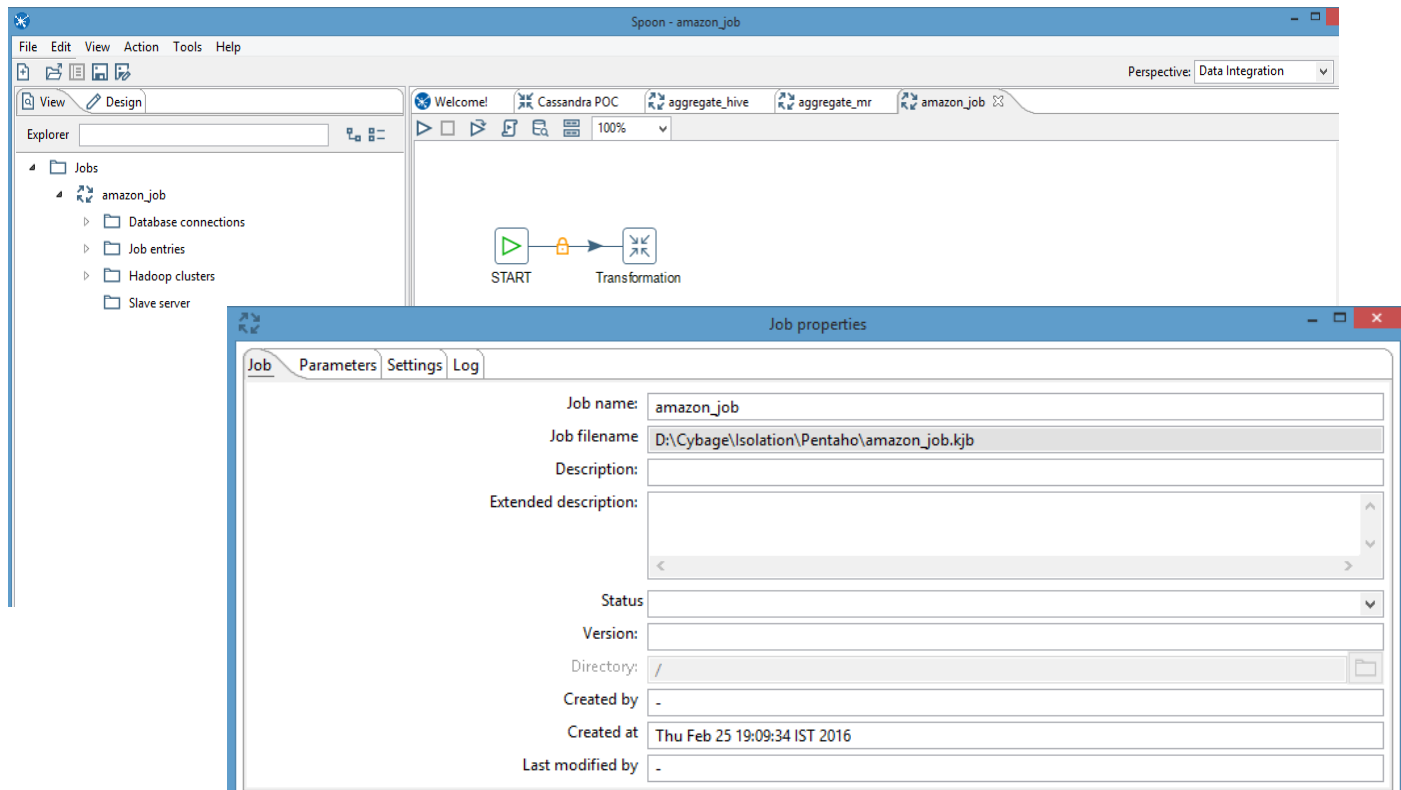- Lookup

# Jobs

## Job Settings

# Jobs

## Job Design Process



Adding Steps or Job Entries → Transformation Step Options → Job Entry Options → Adding Hops → Running Job

# Jobs

## Adding Transformation Steps

# Jobs

## Transformation Options

# Run Transformations / Jobs

- **Local Execution**
- **Execute remotely**
- **Execute clustered**

# Logging



- Nothing: Don't show any output
- Error: Only show errors
- Minimal: Only use minimal logging
- Basic: This is the default basic logging level
- Detailed: Give detailed logging output
- Debug: For debugging purposes, very detailed output.
- Row level: Logging at a row level, this can generate a lot of data.

## Demos

- Jobs

# Building Dimensional Model

The four key decisions made during the design of a dimensional model

- Identify the Source Data for business process.

- Define the grain of data .

- Identify the dimensions.

- Identify the facts.

# Operational vs Reporting Databases

- Relational databases are typically either:

**Operational**

**Reporting**

# Features of an Operational Database

- Operational databases:

  – are designed to maximize accuracy and minimize redundancy

  – are optimized for writing/updating data rather than reading data

  – often result in monolithic designs with multiple joins

  – Large queries can perform slowly.

# Identify Issues with Operational Databases

- "Show all customer types that bought from a product line."

- The query must check data in seven tables before returning a result set.

# Reporting Databases (Star Schema Design)



**Sales Rep**

1..1
0..n

**Customer** 1..1 **Order Fact** 1..1 **Product**
0..n 0..n

0..n
1..1

**Date**

- Transactional data is stored in a fact table
- Reference data is stored in separate dimension tables

- **same information, but five tables instead of nine**

# Create a Star Schema

- Collapse the relationships to form dimensions (perspectives).

# Examine Operational Data

- Data is normalized

**Product Line Table**

| PL# | PL_Desc |
|-----|--------------|
| a | Classic Tents |
| b | Moose Boots |

**2 rows**

**Product Type Table**

| PL# | PT# | PT_Desc |
|-----|-----|-------------|
| a | 1 | Pup Tents |
| a | 2 | Family Tents |
| b | 11 | Child Boots |
| b | 12 | Adult Boots |

**4 rows**

**Product Table**

| PT# | Prod# | Prod_Desc |
|-----|-------|-----------|
| 1 | 101 | Green |
| 1 | 102 | Black |
| 2 | 201 | Yellow |
| 2 | 203 | Brown |
| 11 | 1101 | Blue |
| 12 | 1102 | Blue |

**6 rows**

**Before collapsing into a star schema dimension**

# Examine Reporting Data

- Data is de-normalized

## Product Dimension Table

| PL# | PL_Desc | PT# | PT_Desc | Prod# | Prod_Desc |
|-----|---------|-----|---------|-------|-----------|
| A | Classic Tents | 1 | Pup Tents | 101 | Green |
| A | Classic Tents | 1 | Pup Tents | 102 | Black |
| A | Classic Tents | 2 | Family Tents | 201 | Yellow |
| A | Classic Tents | 2 | Family Tents | 203 | Brown |
| B | Moose Boots | 11 | Child Boots | 1101 | Blue |
| B | Moose Boots | 12 | Adult Boots | 1102 | Blue |

**6 rows**

**After collapsing into a star schema dimension**

# Fact Tables

- Fact tables contain the (usually additive) numbers by which a company measures itself:

  – Standard Selling Price - not additive

  – Sale Amount - additive

**Dimension Tables**

**Fact Table**

**Product**

**Measures** → Sales Revenue
Quantity
..........................
**Foreign Keys** → Product Key
Customer Key
Time Key

**Customer**

**Time**

# Dimension Tables

- Dimension tables provide descriptive information.
- Dimension tables may be "conformed" so that they are applicable to multiple fact tables.

# Dimension Types

- What is SCD?

- SCD Type 1

- SCD Type 2

- SCD Type 3

# Dimension Types

Slowly Changing Dimension – Type 1: Overwrite

Before:

| Supplier_Key | Supplier_Code | Supplier_Name | Supplier_State |
|---|---|---|---|
| 123 | ABC | Acme Supply Co | CA |

After:

| Supplier_Key | Supplier_Code | Supplier_Name | Supplier_State |
|---|---|---|---|
| 123 | ABC | Acme Supply Co | IL |

# Dimension Types

Slowly Changing Dimension –

Type 2: Add new row

Before**:**

| Supplier_Key | Supplier_Code | Supplier_Name | Supplier_State |
|---|---|---|---|
| 123 | ABC | Acme Supply Co | CA |

After**:**

| Supplier_Key | Supplier_Code | Supplier_Name | Supplier_State | Start_Date | End_Date |
|---|---|---|---|---|---|
| 123 | ABC | Acme Supply Co | CA | 01-Jan-2000 | 21-Dec-2004 |
| 124 | ABC | Acme Supply Co | IL | 22-Dec-2004 | |

# Dimension Types

Slowly Changing Dimension –

•Type 3: Add new attribute

Before**:**

| Supplier_Key | Supplier_Code | Supplier_Name | Supplier_State |
| --- | --- | --- | --- |
| 123 | ABC | Acme Supply Co | CA |

After**:**

| Supplier_Key | Supplier_Code | Supplier_Name | Original_Supplier_State | Effective_Date | Current_Supplier_State |
| --- | --- | --- | --- | --- | --- |
| 123 | ABC | Acme Supply Co | CA | 22-Dec-2004 | IL |

## Fact Types

**Factless fact tables**

Most Fact Tables are used to capture numerical results, but it is possible that the event merely records a set of dimensional entities coming together at a moment in time.

Such Fact table will have foreign keys from all related dimension tables without having any particular fact entry.

Example, an event of a student attending a class on a given day may not have a recorded numeric fact

## Fact Types

**Aggregate fact tables**

- *Aggregate fact tables* are simple numeric rollups of atomic fact table data.

- Achieve improved query performance.

- Materialized views can serve as aggregate facts

- BI tools can choose appropriate (aggregated or atomic) aggregate level at query time.

# Demos

Datawarehouse
- Building Dimension
- Building Fact Tables

# Important Links

- Download Link: https://sourceforge.net/projects/pentaho/

- Documentation:
  http://wiki.pentaho.com/display/EAI/Latest+Pentaho+Data+Integration+%28aka+Kettle%29+Documentation

# Any Questions?

Thank you!