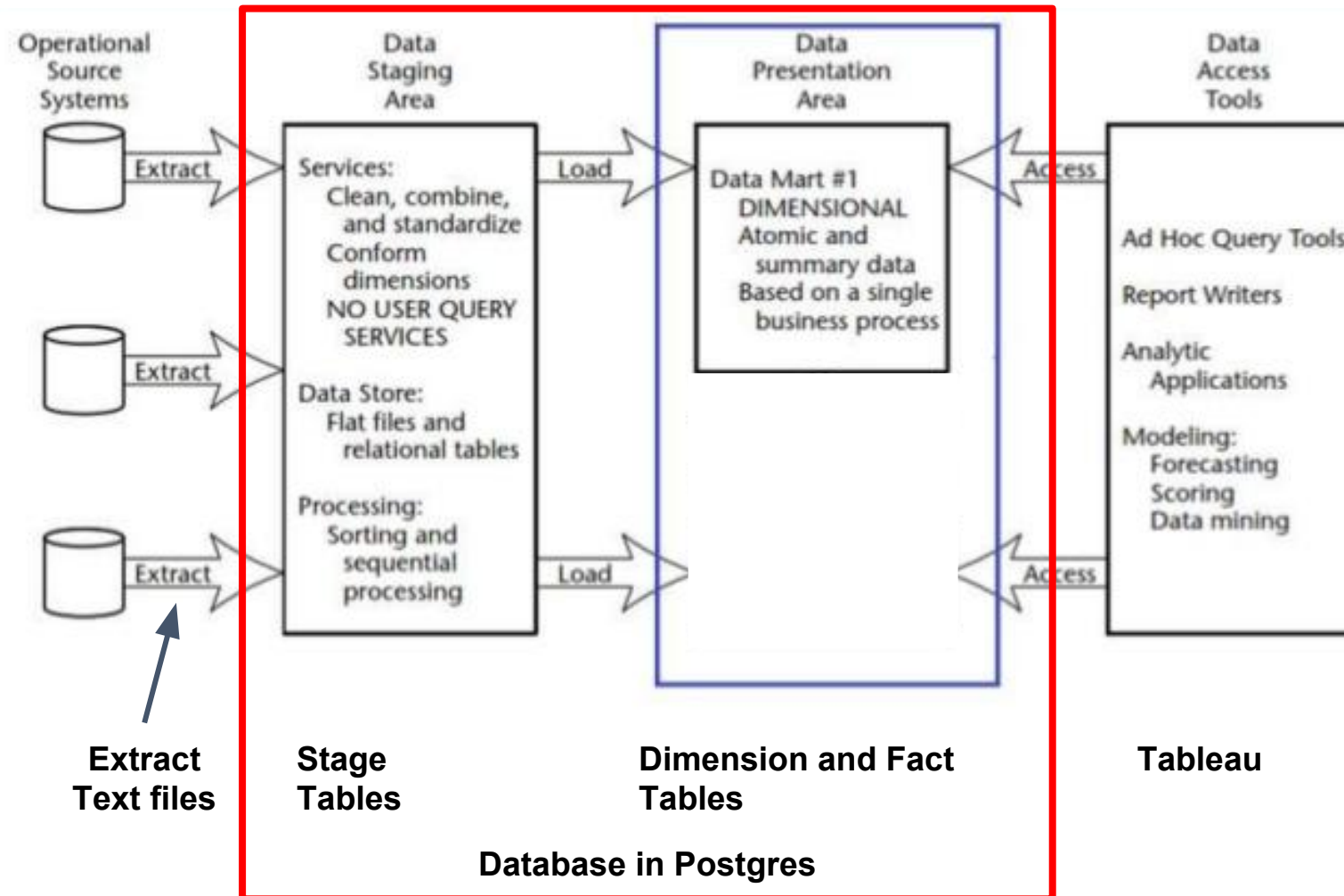# Building Sales Data mart Using **Pentaho** **Part 2**

By Naheed Anjum Arafat

# Sales Data mart: High Level View

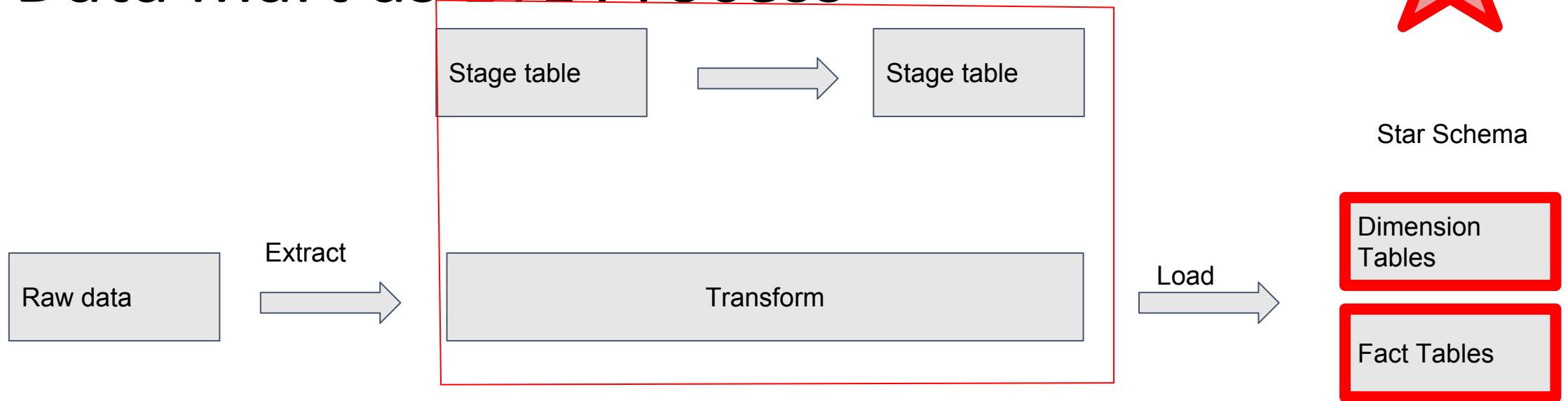# Data mart Architecture



Operational Source Systems — Extract → Data Staging Area

Services:
Clean, combine, and standardize
Conform dimensions
NO USER QUERY SERVICES

Data Store:
Flat files and relational tables

Processing:
Sorting and sequential processing

Load → Data Presentation Area

Data Mart #1
DIMENSIONAL
Atomic and summary data
Based on a single business process

Access ← Data Access Tools

Ad Hoc Query Tools

Report Writers

Analytic Applications

Modeling:
Forecasting
Scoring
Data mining

**Extract Text files** → **Stage Tables** → **Dimension and Fact Tables** → **Tableau**

**Database in Postgres**

# Data mart as ETL Process

```
Raw data  --Extract-->   [Stage table]  -->  [Stage table]
                         [        Transform        ]  --Load-->   Star Schema
                                                                  Dimension Tables
                                                                  Fact Tables
```
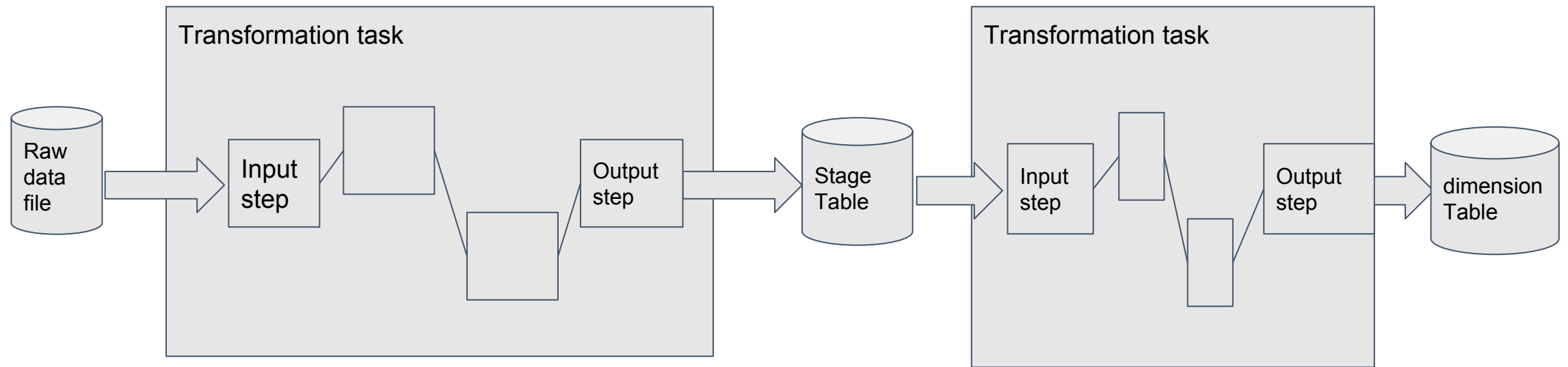
- Extract from source

- **Start fresh**
- Cleanse data
- Align to data type
- Field renaming
- Transformations
- Aggregates
- Splitting into tables

- **Persistent storage**
- Keys and measures
- Performance oriented
- Enforce Grain
- Enforce Unique keys

# ETL Process (for dimensions) in Spoon



Raw data file → Transformation task [Input step → → Output step] → Stage Table → Transformation task [Input step → → Output step] → dimension Table

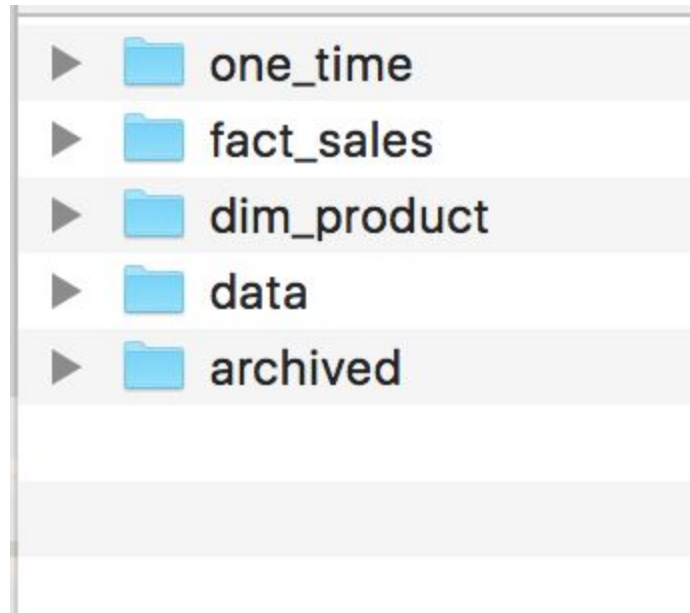- Cleanse data
- Align to data type
- Field renaming
- Transformations
- Aggregates

- Enforce Slowly changing Dimension functionality

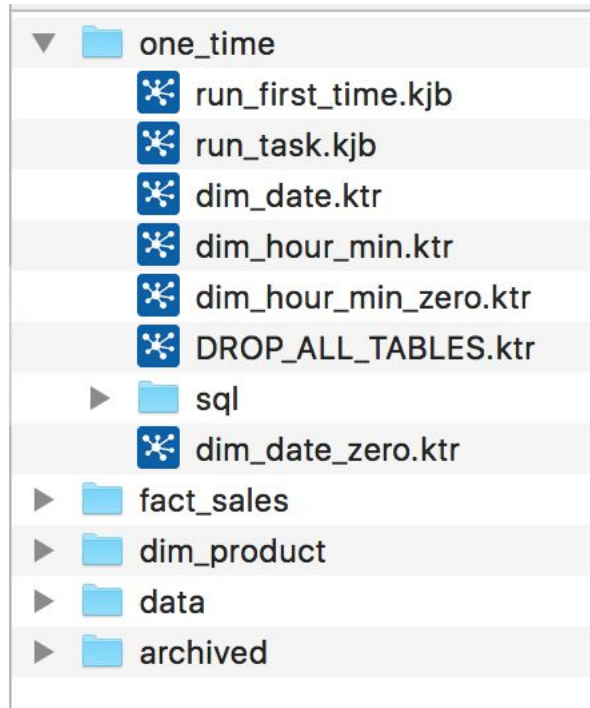# Example of a transformation

# Setup

# Contents of the lecture Pack

| | |
|---|---|
| ▶ 📁 one_time | |
| ▶ 📁 fact_sales | |
| ▶ 📁 dim_product | |
| ▶ 📁 data | |
| ▶ 📁 archived | |

- Transformation/jobs to be run once

- Task 2 folder

- Task 1 folder

- Datasets (to be used in class)

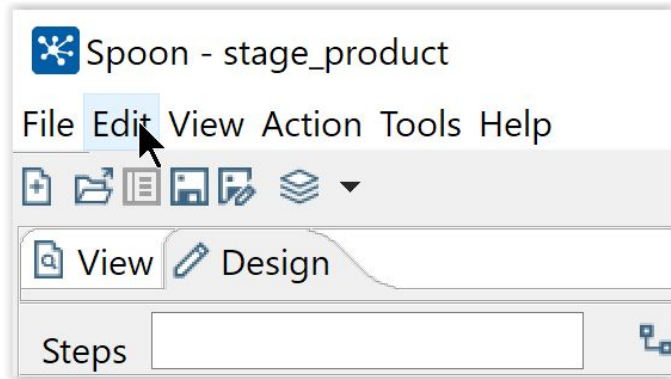- Archive directory

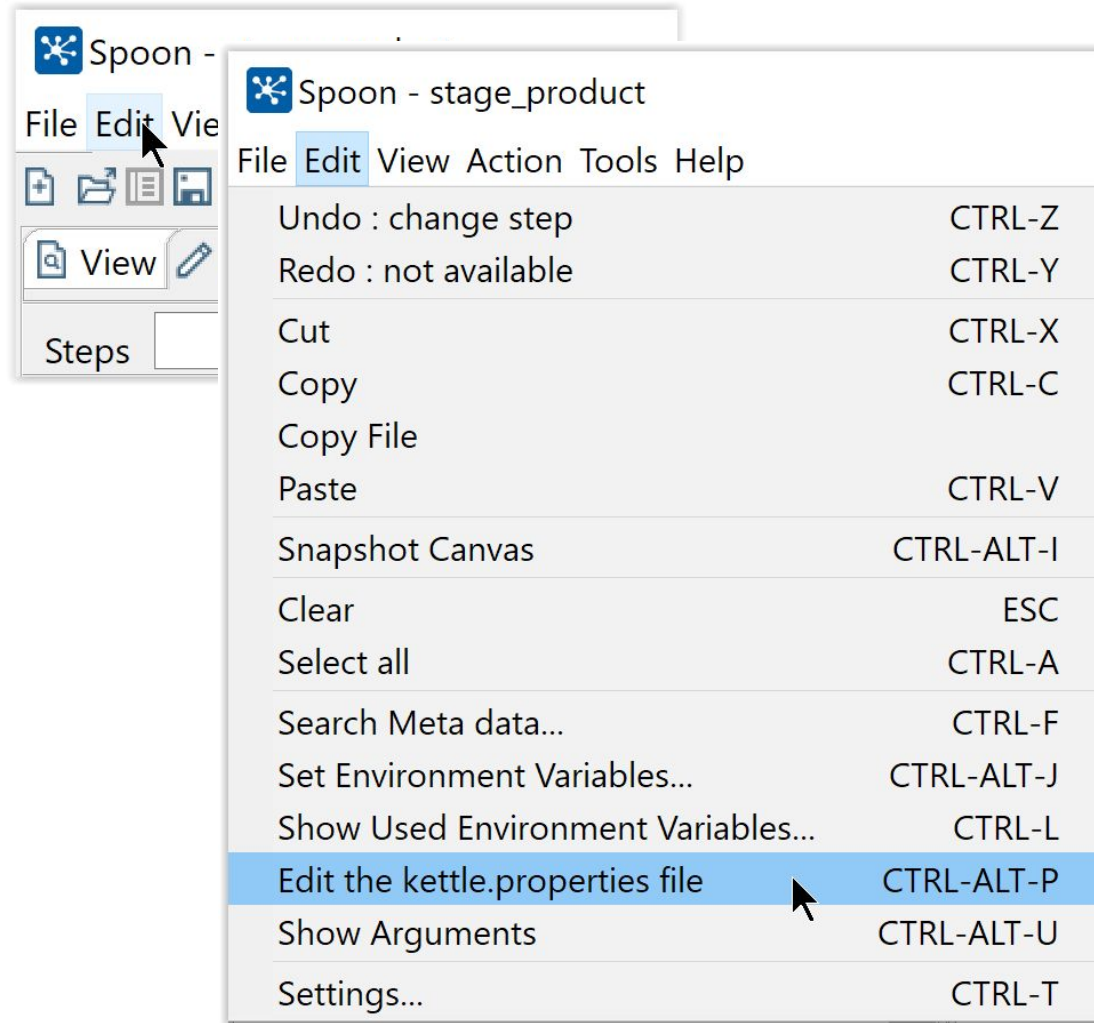# Contents of the lecture Pack

# Preparation

- Download the Lecture pack from IVLE
- Start up postgres server from manager-osx (Bitnami M/W/LAPP stack)
- Open PDI



Pentaho Data Integration

**Version 7.0**
General Availability Release - 7.0.0.0-25
Build Date: November 5, 2016 03:35:36

Copyright (C) 2007 - 2017 Pentaho Corporation. All rights reserved.

Licensed under the Apache License, Version 2.0 (the "License");
you may not use this application and all files except in compliance
with the License.  You may obtain a copy of the License at

http://www.apache.org/licenses/LICENSE-2.0

Unless required by applicable law or agreed to in writing, software
distributed under the License is distributed on an "AS IS" BASIS,
WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
See the License for the specific language governing permissions and
limitations under the License.

**pentaho**
A Hitachi Group Company

# Global Variables

# Global Variables

# Global Variables

```
data_folder
project_folder
```

| 68 | PENTAHO_METAST... | |
|----|-------------------|---|
| 69 | data_folder | C:\Users\~~herma\OneDrive\Documents\Projects\NUS-pentaho-tutorial~~\retail_sales\data |
| 70 | project_folder | C:\Users\~~herma\OneDrive\Documents\Projects\NUS-pentaho-tutorial~~\project_test |

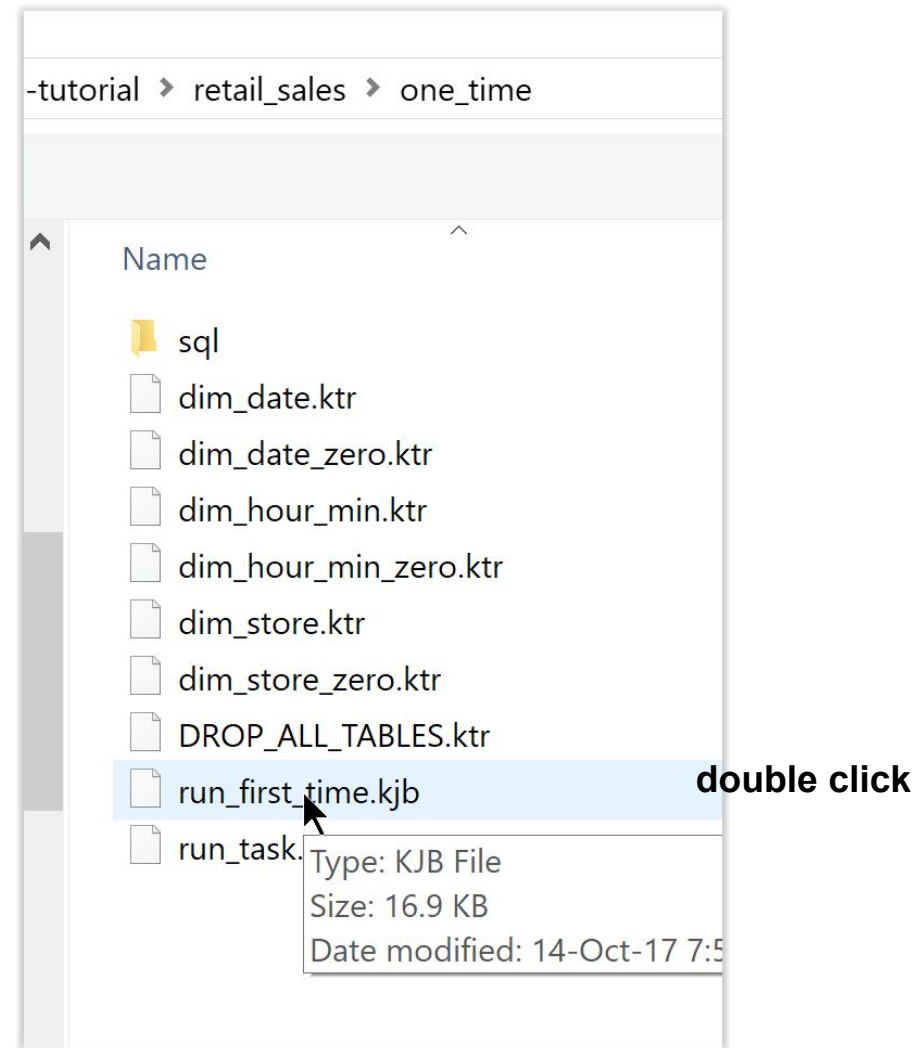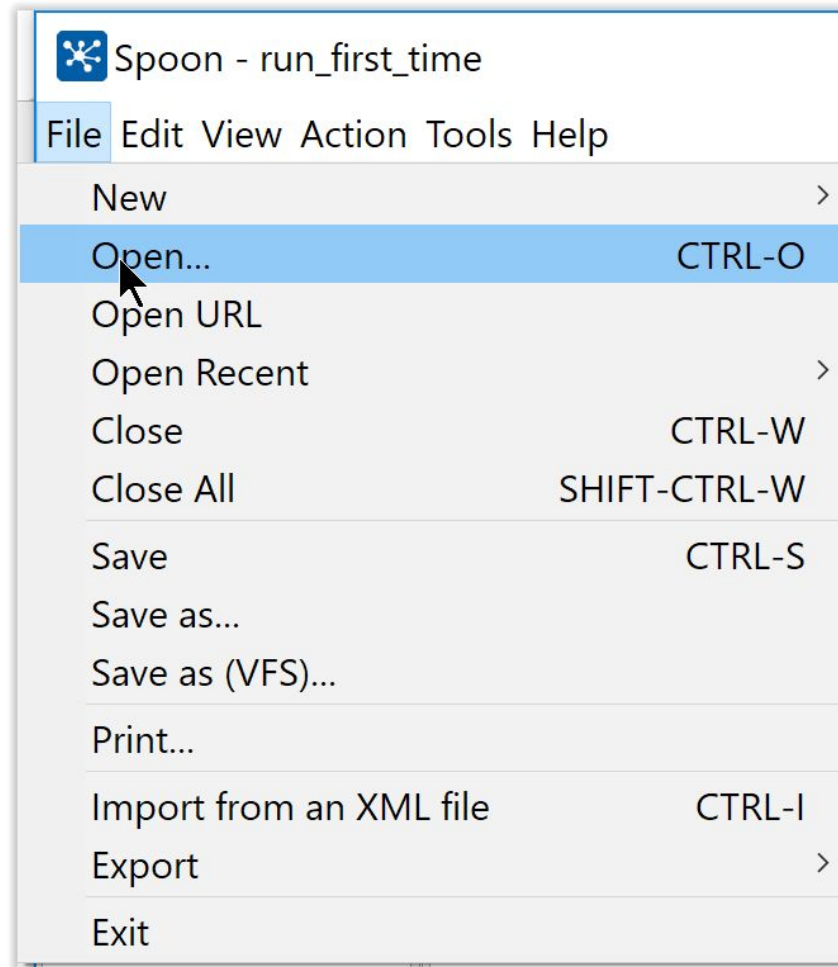**Use variable** `${data_folder}` **or** `${project_folder}`

**Any dialog field with**                    **icon**

# How to use Global Variables



Step name: extract_product

elds | Additional output fields

File or directory: ${data_folder}/extract_product.txt

# Preparation

- Make sure you have set up data_folder and project_folder variables
- Open "run_first_time.kjb" Job in Spoon

# Preparation

Spoon - run_first_time

File Edit View Action Tools Help

| New | > |
| **Open...** | **CTRL-O** |
| Open URL | |
| Open Recent | > |
| Close | CTRL-W |
| Close All | SHIFT-CTRL-W |
| Save | CTRL-S |
| Save as... | |
| Save as (VFS)... | |
| Print... | |
| Import from an XML file | CTRL-I |
| Export | > |
| Exit | |

-tutorial > retail_sales > one_time

Name

📁 sql

📄 dim_date.ktr

📄 dim_date_zero.ktr

📄 dim_hour_min.ktr

📄 dim_hour_min_zero.ktr

📄 dim_store.ktr

📄 dim_store_zero.ktr

📄 DROP_ALL_TABLES.ktr

📄 run_first_time.kjb **double click**

📄 run_task.
Type: KJB File
Size: 16.9 KB
Date modified: 14-Oct-17 7:5

# Connect to Postgres

# Set up Database Connection

# Set up Database Connection

# Set up Database Connection

# Edit/Test Database Connection



**Set the values of the following configuration fields:**

**Host Name:**
**Database Name:**
**Port Number:**
**User Name:**
**Password:**

**Set it according to your configuration of the Postgres installation**

Test the connection

# Check the contents of Database



Navigate to Tables
See that the database is where you want the data mart tables to be

Click on Explore

# Share Database Connection



Sharing a connection enables
all transformations and jobs to use the same database.

If the Database Connection is shared, it will be in **BOLD**

# How to set Connection to "Shared"



The database connection is not in BOLD
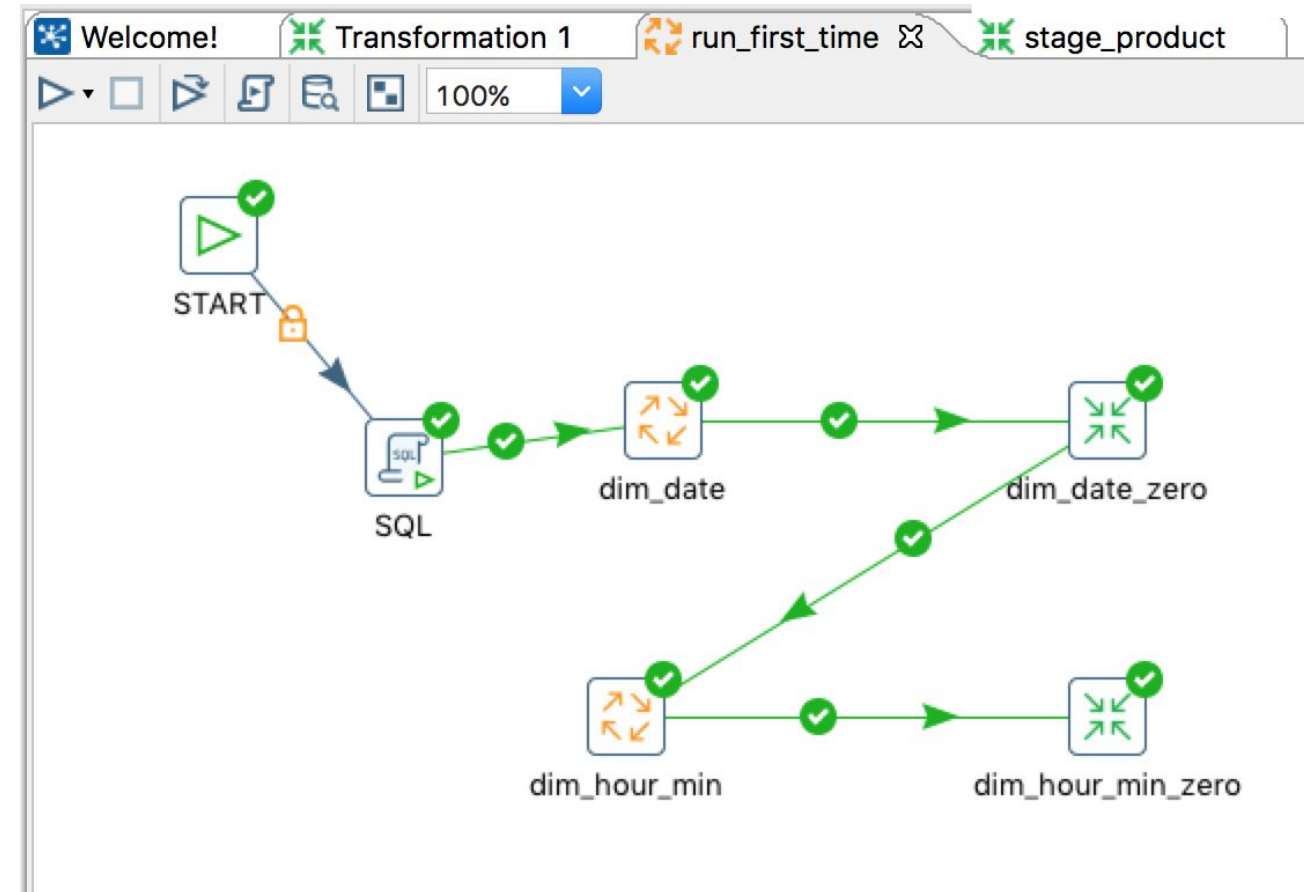so it is not shared and will not be
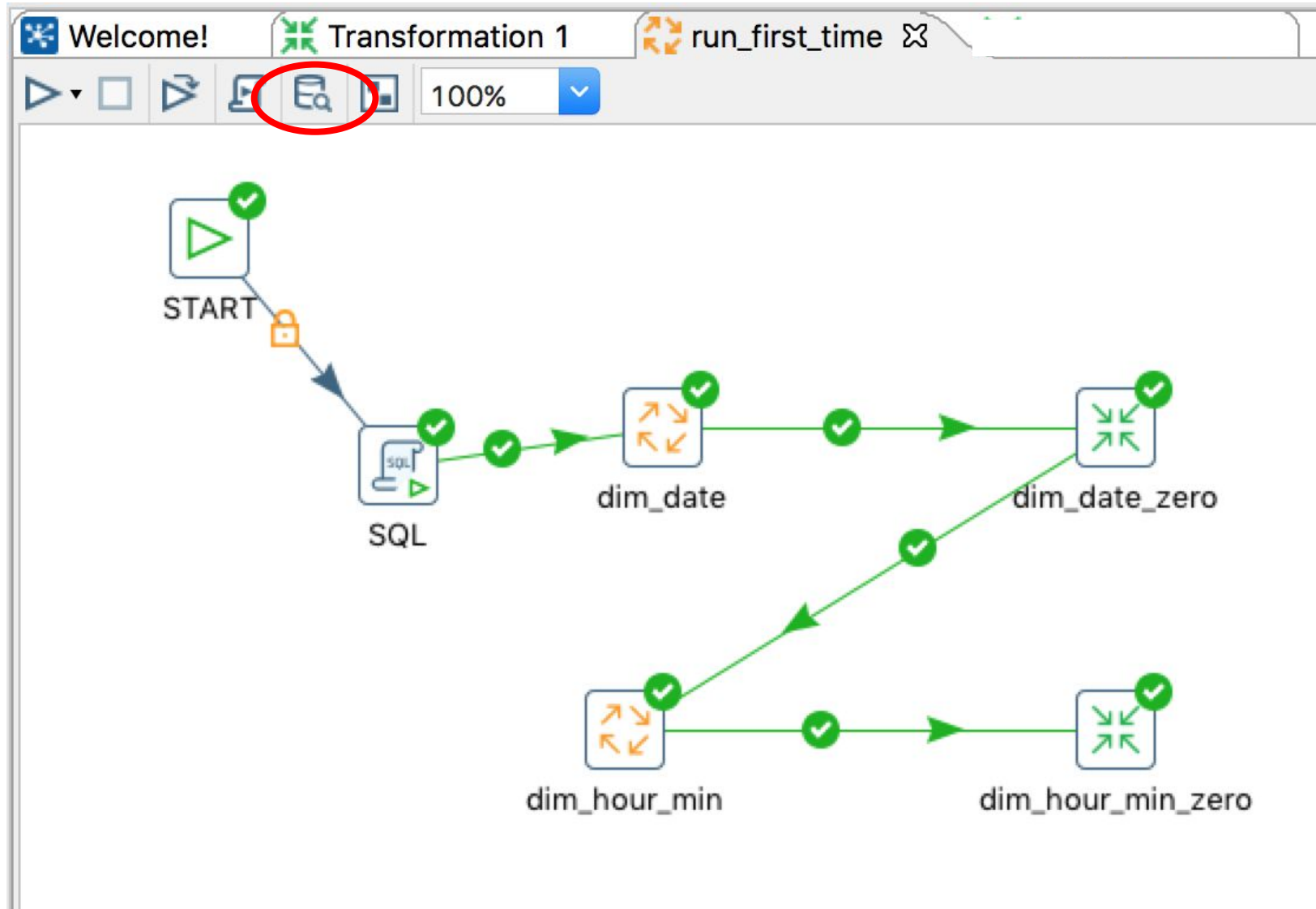accessed by other transformations and jobs

# Run First Time Job

- Creates 2 dimensions for you
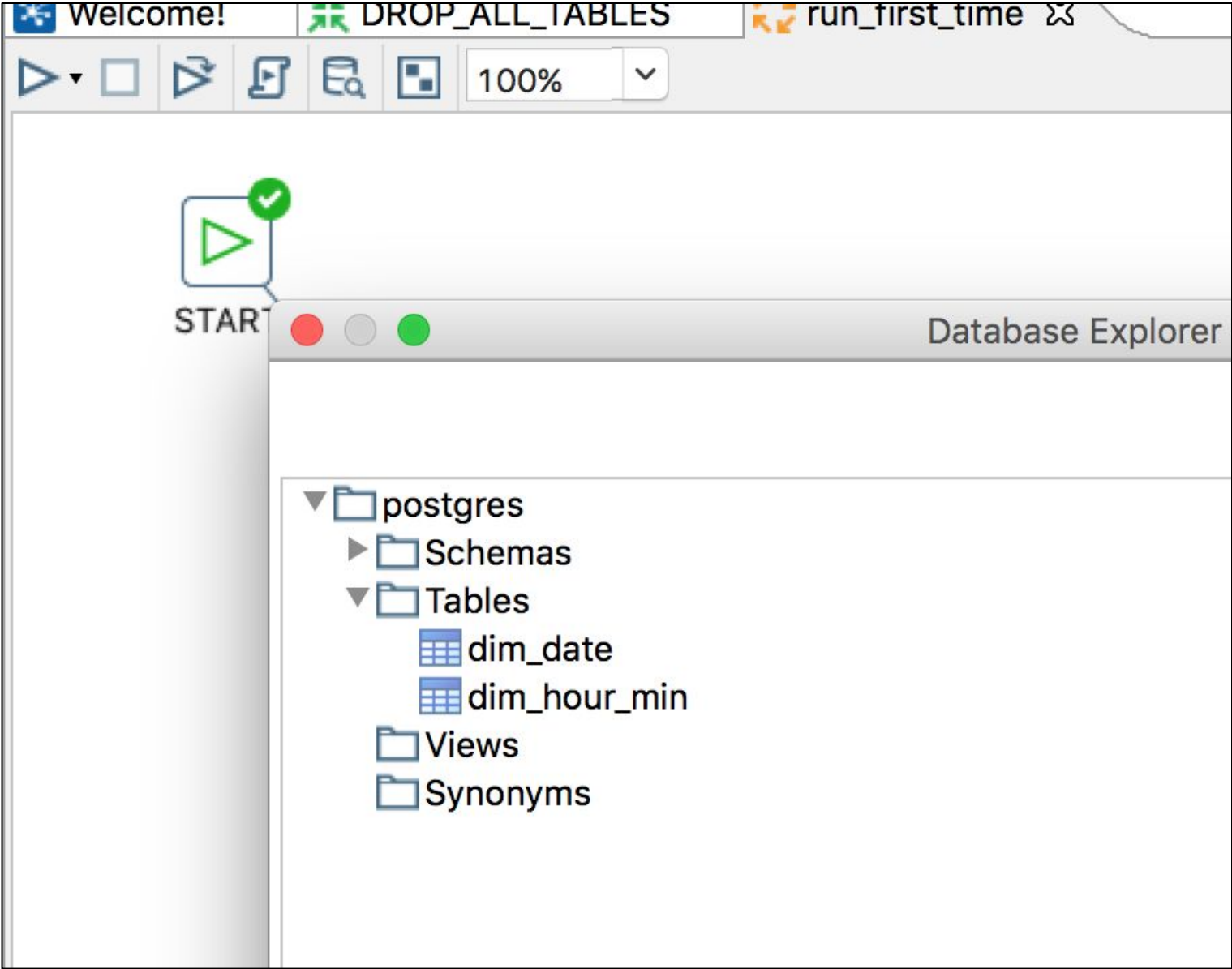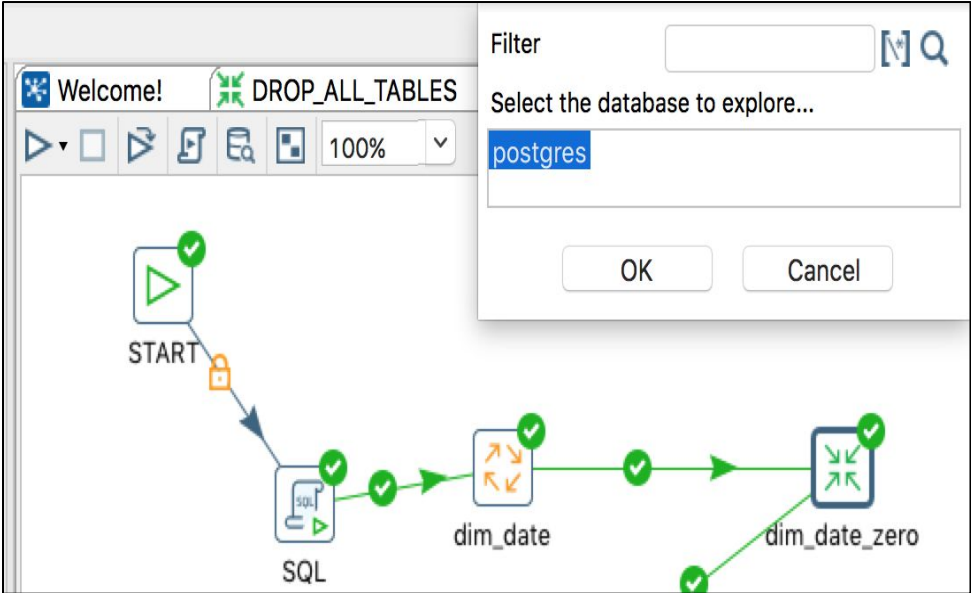  - date dimension
  - hour-minute dimension

# Run the Job



Check the database using the Explore option.

You should see 2 tables.

# Check Database

# Todays tasks