



Connecting the
Data-Driven Enterprise >



Participant Guide
Talend Studio Introduction

Version 6.5

Copyright 2018 Talend Inc. All rights reserved.

Information in this document is subject to change without notice. The software described in this document is furnished under a license agreement or nondisclosure agreement. The software may be used or copied only in accordance with the terms of those agreements. No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or any means electronic or mechanical, including photocopying and recording for any purpose other than the purchaser's personal use without the written permission of Talend Inc.

Talend Inc.
800 Bridge Parkway, Suite 200
Redwood City, CA 94065
United States
+1 (650) 539 3200

Welcome to Talend Training



Congratulations on choosing a Talend training course.

Working through the course

You will develop your skills by working through use cases and practice exercises using live software. Completing the exercises is critical to learning!

If you are following a self-paced, on-demand training (ODT) module, and you need an answer to proceed with a particular exercise, use the help suggestions on your image desktop. If you can't access your image, contact customercare@talend.com.

Exploring

You will be working in actual Talend software, not a simulation. We hope you have fun and get lots of practice using the software! However, if you work on tasks beyond the scope of the training, you could run out of time with the environment, or you could mess up data or Jobs needed for subsequent exercises. We suggest finishing the course first, and if you have remaining time, explore as you wish. Keep in mind that our technical support team can't assist with your exploring beyond the course materials.

For more information

Talend product documentation (help.talend.com)

Talend Community (community.talend.com)

Sharing

This course is provided for your personal use under an agreement with Talend. You may not take screenshots or redistribute the content or software.

Intentionally blank

S T E N E C O N T E X T

LESSON 1 Introduction to Talend Studio in Context	
Concepts	8
LESSON 2 Getting Started with Talend Studio	
Concepts	12
Overview	18
Creating a Project	19
Creating Your First Job	21
Running a Job	28
Review	30
LESSON 3 Using Talend Studio to Design a Data Integration Project	
Concepts	32
Overview	38
Importing External Resources into Your Project	39
Preparing Data Sources: Working with Structured Files	42
Preparing Data Sources: Reading Data from Databases	49
Joining Data From Multiple Data Sources Using tMap	55
Transforming Data	73
Review	81

Intentionally blank

L E S S O N

Introduction to Talend Studio in Context

In this lesson you have an overview of Talend Studio with basic data integration examples.

This chapter discusses:

Concepts	8
----------------	---

Concepts



OBJECTIVES



- Know the key benefits of using Talend Studio
- Define Data Integration
- Discuss common data processing tasks Talend Studio can help you with



Talend Data Fabric uses the same Eclipse environment for several applications: big data integration, data integration, application integration (ESB), mdm and cloud integration, and data preparation.

This course covers the basics of Talend Studio, introducing the most commonly used components with data integration examples.

WHY TALEND STUDIO?



Respond Faster to Business Needs



Talend Studio provides rich capabilities:

Design Faster: easy data integration with drag-and-drop user interface, no scripting.

Collaborate better: team support for developers, reliable continuous delivery of new functionalities.

Cleanse earlier: profiling and data matching studio to cleanse your data more accurately.

Manage more: centralize and manage deployment through a web-based administration console.

Scale easier: scale up with a high availability environment and advanced clustering.

KEY BENEFITS OF TALEND STUDIO



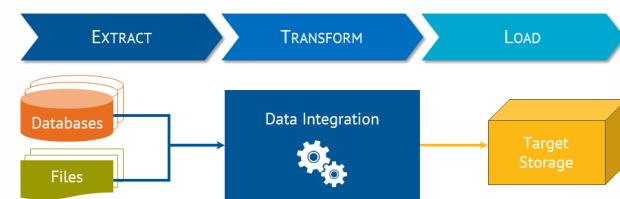
- Connectivity to source and target data systems
- Scalability and performance
- Transformation flexibility
- Data cleansing components
- Logging and exception handling
- Easy integration with web services
- Hundreds of components available

5

WHAT IS DATA INTEGRATION?



- Merge data coming from different systems
- Implementation is simplified by using an **ETL** tool



6

Data integration primarily supports the analytical processing of large, distributed data sets by aligning and combining different data sources through an **ETL** (Extract Transform Load) tool:

Extract: get the data from source systems as efficiently as possible.

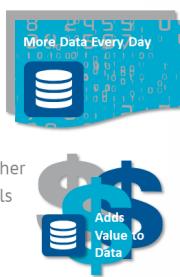
Transform: perform calculations/transformations on data. At this step, cleansing and standardization are also required most of the time.

Load: load the data to the target storage.

WHY DATA INTEGRATION IN TALEND STUDIO?



- Different formats
- Different database structures and models
- Different standards
- Different systems
- Removes mistakes and corrects data
- Adjusts data from multiple sources to be used together
- Structures and aggregates data to be used by BI tools



7

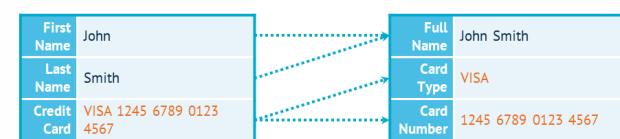
With Talend Studio you can get the most out of your data by managing and integrating a wide range of data formats, data structures, standards and systems. You can enhance and add value to your data by removing mistakes, correcting, organizing and aggregating data so that it can be processed by BI tools and many more.

COMMON DATA PROCESSING WITH TALEND STUDIO



- Splitting or merging fields
- Hundreds of components to assist you

Splitting or merging fields is easy using predefined functions and mapping facilities.

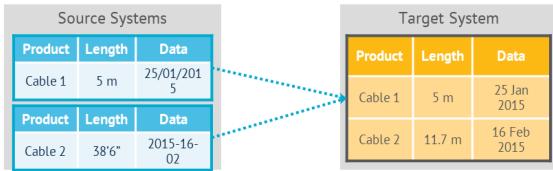


8

COMMON DATA PROCESSING WITH TALEND STUDIO



- Data standardization
 - Dates
 - Units of measurement



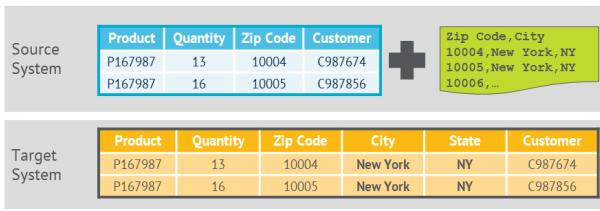
When merging data on cable sales registered by the company in different locations, you want to use the same standards for data representation and measurement units.

9

COMMON DATA PROCESSING WITH TALEND STUDIO



- Join data
 - Complete data with information from external sources



Join data from heterogeneous sources and load it to a target system.

10

COMMON DATA PROCESSING WITH TALEND STUDIO



- Aggregate data
 - Perform computations
 - Group and order elements as needed



Aggregate data to perform computations and statistics.

11

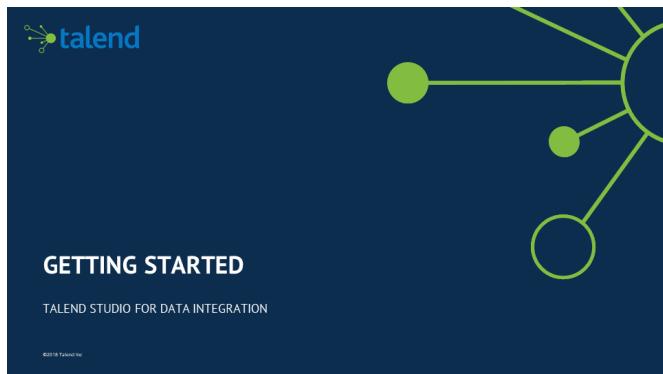
LESSON 2

Getting Started with Talend Studio

This chapter discusses:

Concepts	12
Overview	18
Creating a Project	19
Creating Your First Job	21
Running a Job	28
Review	30

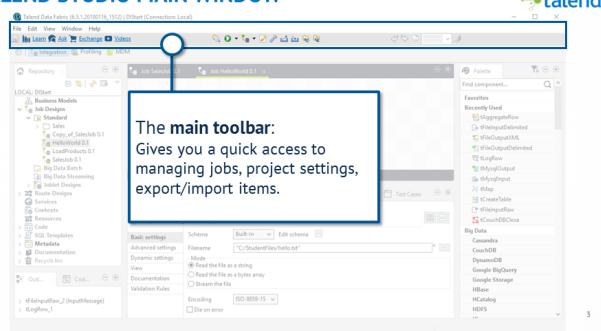
Concepts



OBJECTIVES

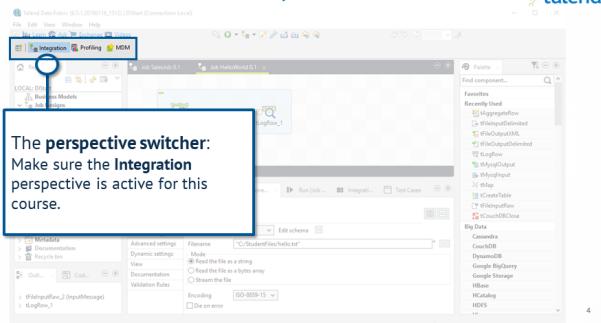
- Learn to use Talend Studio main window
- Create a Job
- Add components to a Job
- Connect and configure components
- Run a Job and view the results in the console

TALEND STUDIO MAIN WINDOW



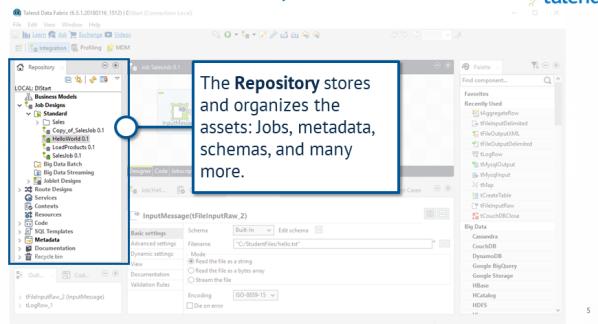
The **main toolbar** allows you to quickly find, run, or create Jobs, as well as importing or exporting items.

TALEND STUDIO MAIN WINDOW, CONTINUED



A **Perspective** defines the initial set and layout of views in the Talend Studio workbench. The selected perspective is always the active one. There are several perspectives in your training environment because the Talend Data Fabric can be used to create various types of projects, including data integration, Big Data, application integration, master data management, and data profiling. In this course, you use only the Integration perspective.

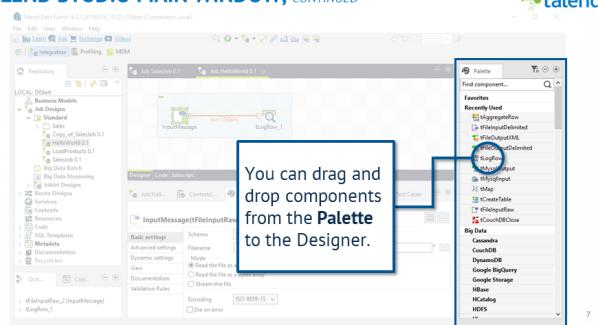
TALEND STUDIO MAIN WINDOW, CONTINUED



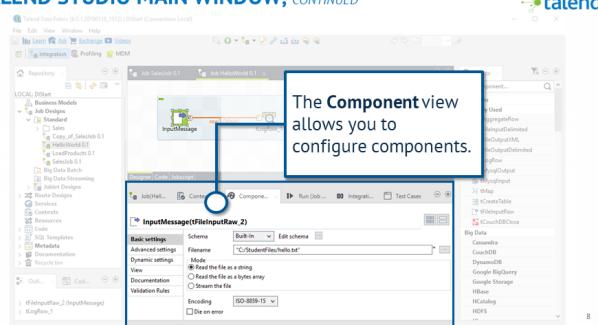
TALEND STUDIO MAIN WINDOW, CONTINUED



TALEND STUDIO MAIN WINDOW, CONTINUED



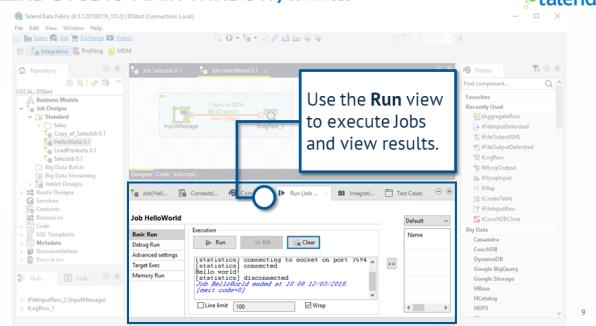
TALEND STUDIO MAIN WINDOW, CONTINUED



The **Palette** contains the technical components you use to design Jobs, grouped by families. Drag and drop components to the Designer to use them in your Job.

The **Component** view allows you to configure the properties of the component you selected in the Designer.

TALEND STUDIO MAIN WINDOW, CONTINUED



The **Run** view allows you to execute a Job and explore the results in the console.

WHAT IS A TALEND COMPONENT?

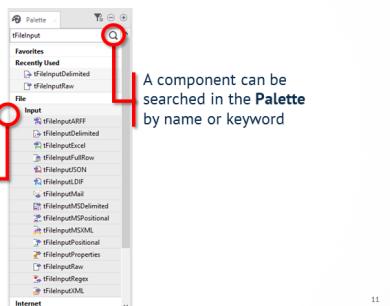
Components are predefined units that perform specific data integration operations and have a graphical representation.



Components are predefined units that perform specific data integration operations. They are implemented as Java code and have a graphical representation to allow easy configuration and simplified Job design.

COMPONENTS

Components are grouped by families and types.



A component can be searched in the **Palette** by name or keyword

Components are grouped by families (for example File, Databases) and types (for example Input, Process). Component main types are:

Input – Reads input from a source and returns output to the data flow

Process – Transforms data in the data flow

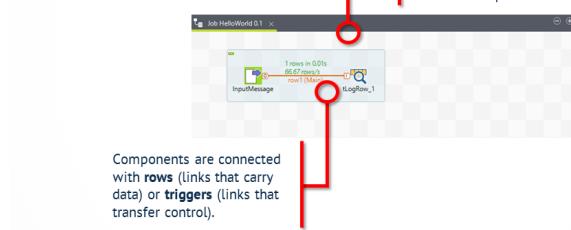
Output – Receives data from the data flow and outputs to a target

Virtual – Component constructed of two or more other components

WHAT IS A TALEND JOB?

Components are connected with **rows** (links that carry data) or **triggers** (links that transfer control).

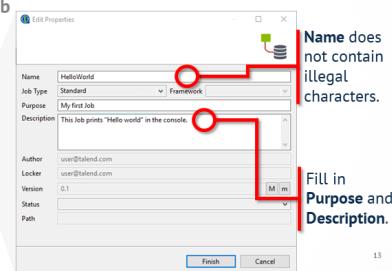
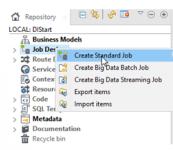
A **Job** consists of one or more connected components. You can execute it to perform a task.



Components are connected with **rows** (links that carry data) or **triggers** (links that transfer control).

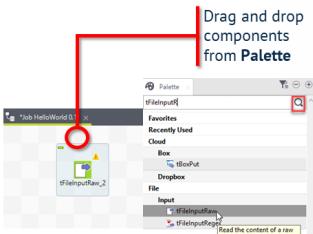
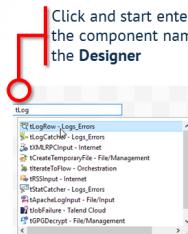
CREATE A JOB

- In the Repository, right-click Job Designs
- Select Create Standard Job



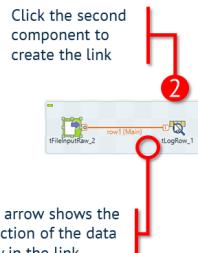
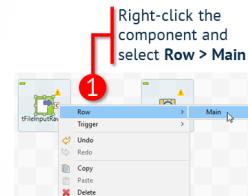
It is highly advised to fill in the **Purpose** and **Description** fields of a Job, for development and maintenance purpose.

ADD COMPONENTS TO A JOB



There are two methods to add components to your Job: directly from inside the Designer and using the Palette.

CONNECT COMPONENTS



The arrow shows the direction of the data flow in the link.

15

NAME A COMPONENT

Select the component you want to rename



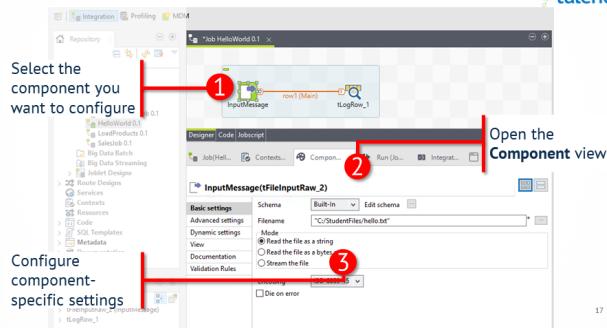
In the Component view, select the View tab



Enter a name in the Label format text box

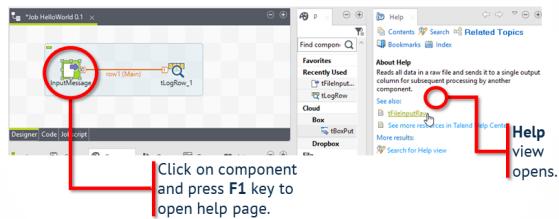
The **Label format** specifies the name of the component while the **Hint format** specifies the pop-up message displayed on mouse-over.

CONFIGURE A COMPONENT



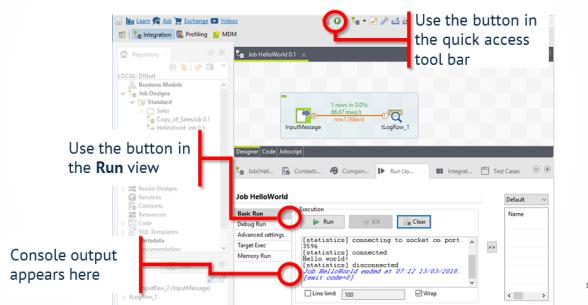
A double-click on a component will also open the **Component** view.

COMPONENT HELP



18

RUN A JOB



To execute a Job you can either use the Run icon in the Run view, or the one in the main toolbar. When you have several Job tabs open in the Designer, make sure to select the correct Job before using the Run icon.

EXERCISE OVERVIEW



Explore Talend Studio, design and run your first Job

- Learn to use Talend Studio
- Design and run the following Job:



20

LESSON SUMMARY



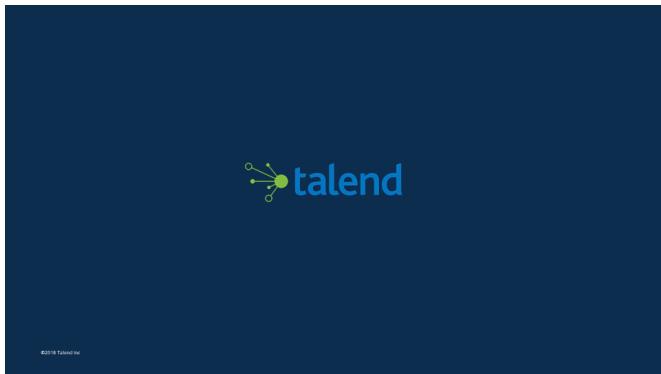
Learn to use Talend Studio main window

- Components, Job, The Designer, the Palette, the Repository etc.

Design a Job

- Create a Job
- Add components to a Job
- Connect and configure components
- Run a Job and visualize results in the console

21



Overview

Use Case

In this course, you are assigned a **Talend Data Fabric** training environment. The purpose of this lesson is to get you started with using Talend Studio. You create an empty project containing one simple job that displays "Hello world!" in the console.

Objectives

After completing this lesson, you will be able to:

- » Create a project
- » Create a Job
- » Add components to a Job
- » Link and configure components
- » Run a Job and visualize the results in the console

Next step

The next step is to create a project.

Creating a Project

Task outline

Before developing projects in Talend Studio, you need to start Talend Studio and create a project.

This exercise takes about 5 mins to complete.

Creating a project

1. START TALEND STUDIO

To start Talend Studio, double-click the Talend Studio icon on your desktop.

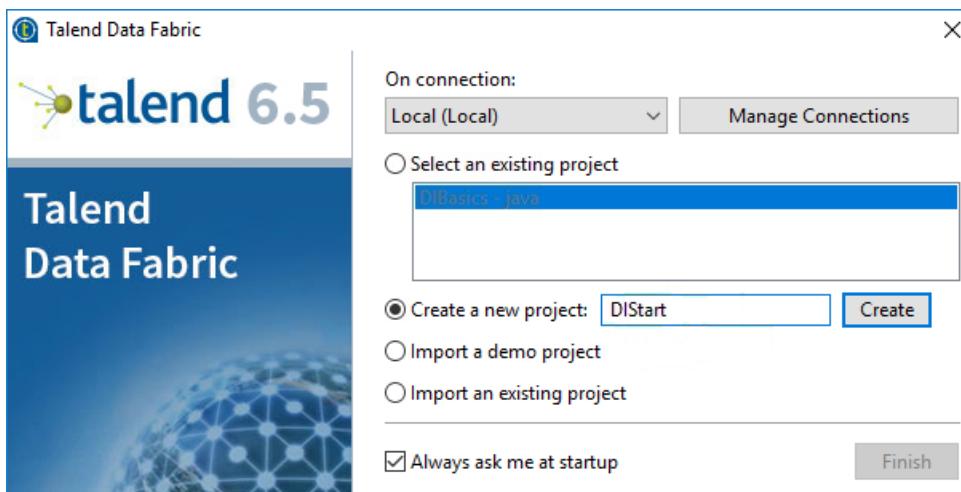


The **Talend Data Fabric** window opens allowing you to open an existing project or create a new one.

2. CREATE AN EMPTY PROJECT

To create an empty project:

- For **On connection**, select **Local (Local)**, which means that once created, the project is stored locally, on your machine
- Select the radio button **Create a new project** and name it *DIStart*. Click **Create** and wait until the project appears in the list of existing projects. The project you created is now selected by default.
- Click the **Finish** button.



NOTE:

- » If you are prompted to connect to the TalendForge community, select **Skip this Step**.
- » If the **Welcome** page appears, click **Start now!**

3. SWITCH TO THE INTEGRATION PERSPECTIVE

The Talend Data Fabric Studio allows you to create several types of projects: data integration, big data, application integration, master data management, and data profiling. Depending on the project type, you can use one or more perspectives to work with.

In this course you work on data integration use cases, therefore you need to select the **Integration** perspective, if it is not already selected. The selected Perspective is always the active one.



Next step

You have learned how to start Talend Studio and create a new project. You are ready to create your first Job.

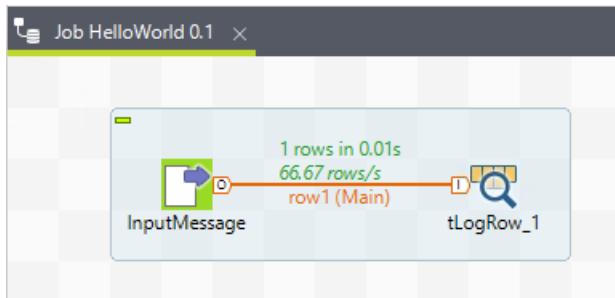
Creating Your First Job

Task outline

In this section, you create a Job that reads the string "Hello world!" from a file and sends it to the console.

This exercise takes about 15 mins to complete.

At the end, the Job looks like this:



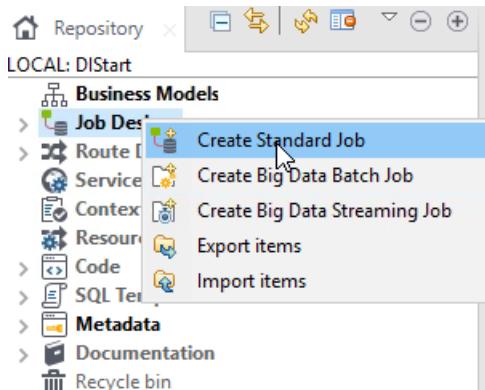
Creating a Job

Jobs are created and stored in the **Repository**, in **Job Designs**.

1. CREATE A STANDARD JOB

To create a Job:

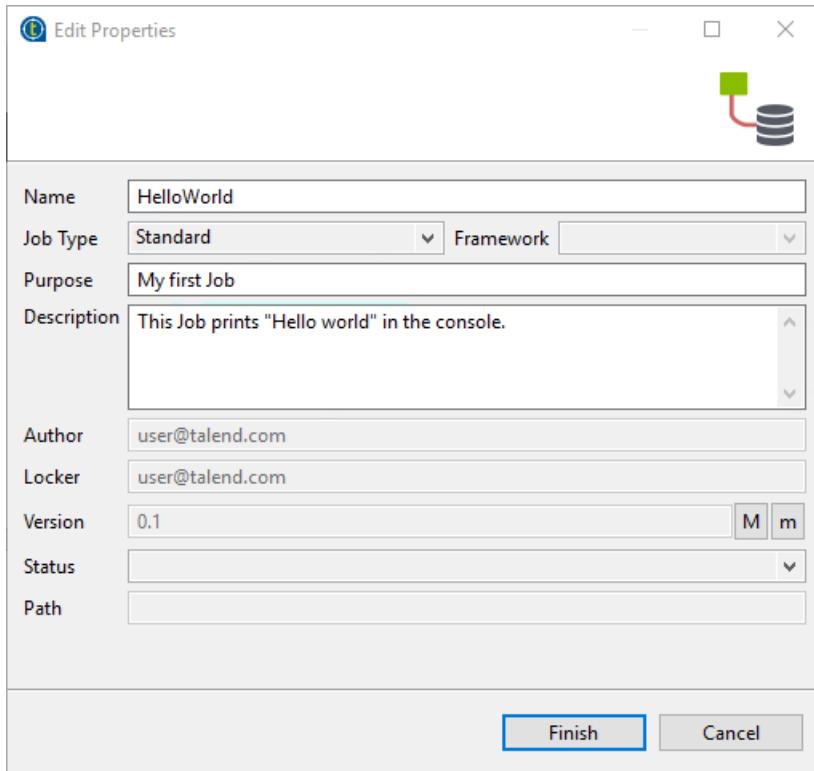
- Right-click **Repository** > **Job Designs**.
- Select **Create Standard Job**. The **New Job** dialog box opens.



2. DEFINE THE JOB PROPERTIES

Define the Job properties in the **New Job** dialog box:

- In the **Name** text box, enter *HelloWorld*.
- Fill-in the **Purpose** and **Description** text boxes.
- Leave the default values for the other parameters.
- Click **Finish**.



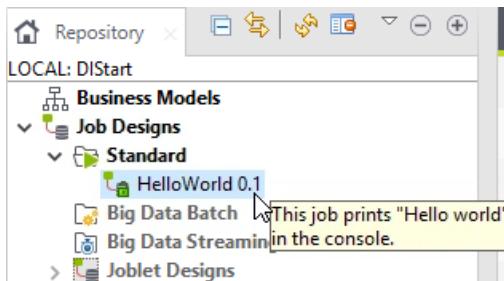
WARNING:

Avoid using common illegal filename characters. For example, *HelloWorld* is a valid name, but *Hello World* is incorrect due to the space between the words.

3. FURTHER EXPLANATIONS

Notice the changes in the **Designer**, **Palette**, and rest of the views.

New Jobs are automatically assigned version **0.1**, therefore, the Job you created is named "HelloWorld 0.1". A tab with the same name is open in the **Designer** and components are available in the **Palette**. You are ready to add components to your Job.



Adding components to your Job

The Job needs two components: a file containing the *Hello world!* string and a console to send the output to. To do this, you can use the **tFileInputRaw** and **tLogRow** components.

tFileInputRaw reads all data in a raw file and sends it to a single output column for subsequent processing by another component.

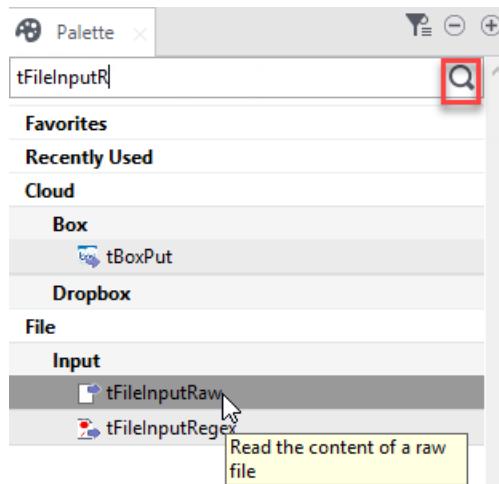
tLogRow displays data or results in the **Run** console to monitor data processed.

In the following, you explore two ways to add components to your Job: using the **Palette** and directly from inside the **Designer**.

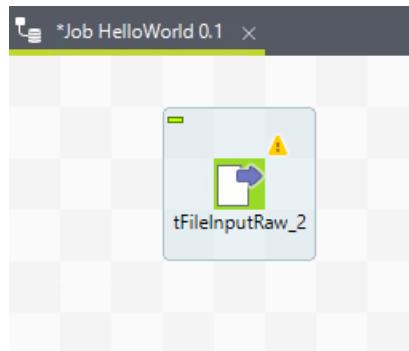
1. ADD THE tFileInputRaw COMPONENT USING THE PALETTE

To add a component using the **Palette**:

- a. In the Palette, look for the **Search** box.
- b. Start entering *tFileInputRaw* in the **Search** box. To list components matching your research, press the **Enter** key (or click the **Search** button). Notice that a short description is displayed next to each component name when you hover the mouse over it.
- c. Select **tFileInputRaw** on the list
- d. Drag the component and place it in the **Designer**.



When you are finished, the Job looks like this:



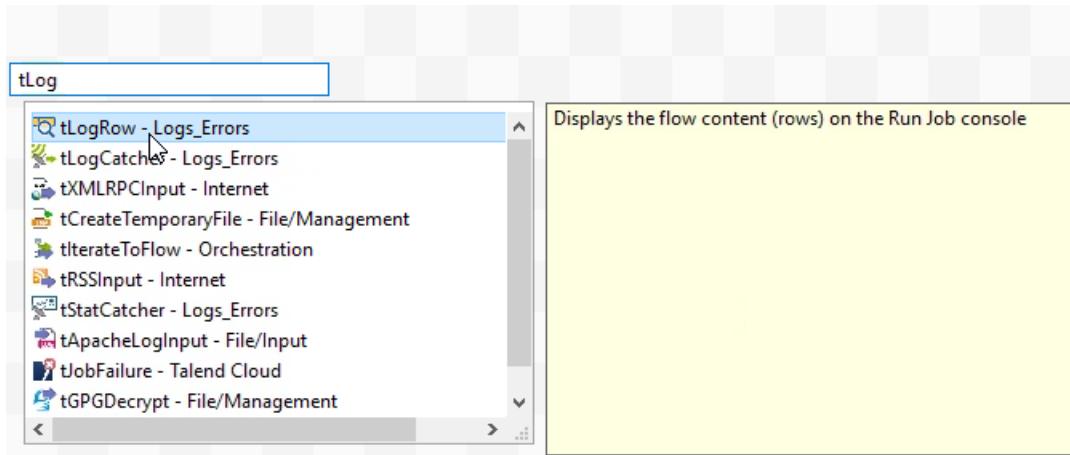
TIP:

To follow Talend Studio design best practices, place components in the **Designer** from the upper left to the lower right corner, as the data flow generally proceeds from left to right and from top to bottom. Therefore, place your input components on the upper left part of the **Designer**.

2. ADD THE tLogRow COMPONENT DIRECTLY FROM INSIDE THE DESIGNER

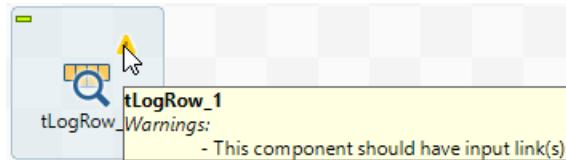
To add a component directly from inside the **Designer**:

- a. Click anywhere in the **Designer**.
- b. Start entering *tLogRow*. As you type, components matching your search are listed. Notice that a short description is displayed next to each component name when you select it with the mouse.
- c. Double-click **tLogRow** on the list. The component appears in the workspace.



3. INSPECT THE WARNINGS

Notice that the components are both displayed with a small warning icon. If you hover the mouse over the icon, a warning message says that the component should have inputs/outputs linked.



TIP:

Any time you need help on a component, just click on it and press the **F1** key.

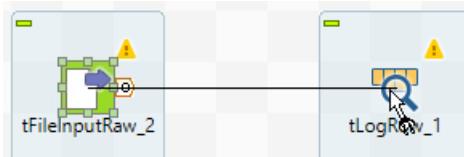
Connecting components

1. CONNECT THE COMPONENTS

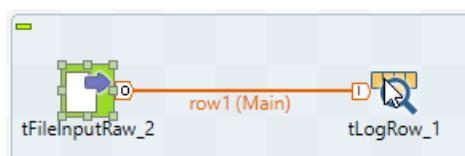
You can connect components with **rows**. These are links that allow data to flow between components.

To connect **tFileInputRaw** and **tLogRow** components:

- Right-click in the middle of the input component (**tFileInputRaw**). The component is highlighted in a frame.
- Select **Row > Main**. A connection link is visible between the input component and mouse pointer. You can now link the output component.



- Click in the middle of the output component (**tLogRow**) to create the link. The arrows on the link show the direction of the flow.



NOTE:

Notice that once you connect the components, they are part of a same subJob which is represented by a blue rectangle.

Configuring components

Remember that the **tFileInputRaw** component must read the "Hello world!" from a file, so you must first create a text file containing this string.

1. CREATE A TEXT FILE

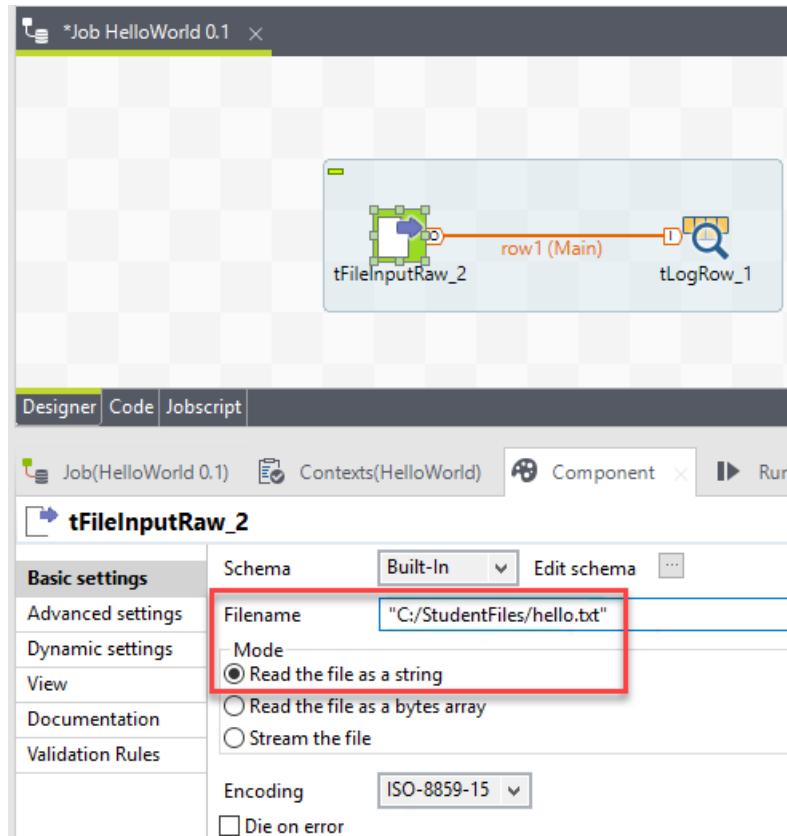
- a. Open Notepad++ (on the Windows taskbar, click the **Notepad++** icon, or go to **Windows Start Menu > NotePad++**).
- b. Create a new file (**File > New**).
- c. Enter *Hello world!*.
- d. Go to **File > Save as**. Name the file *hello.txt* and save it under C:/StudentFiles/
- e. Close Notepad++.

2. DEFINE COMPONENT BASIC SETTINGS

You can configure components using the **Component** view in the Designer. Depending on the component you use, the **Component** view displays different properties.

To configure a component, follow these steps:

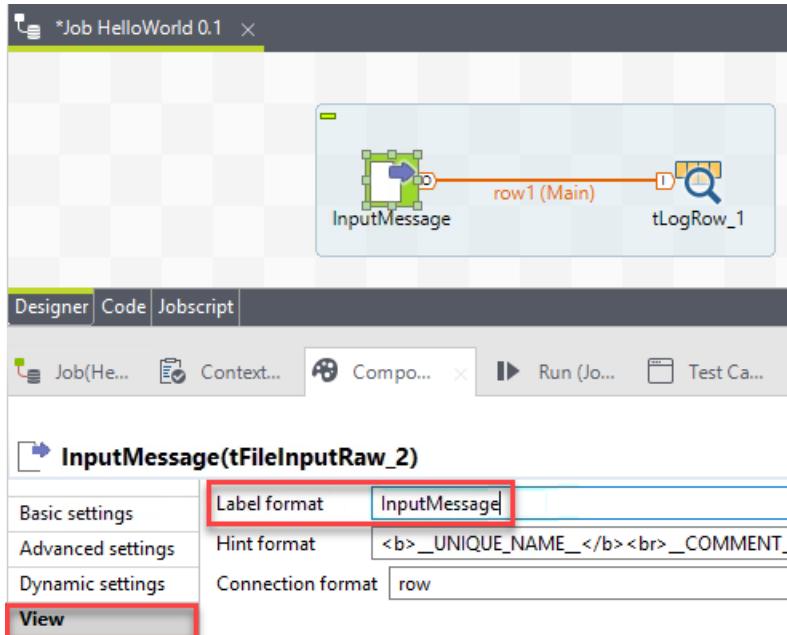
- a. Double-click in the middle of the **tFileInputRaw** component. The **Component** view opens.
- b. For **Filename**, click the [...] to select the file path *C:/StudentFiles/hello.txt*. Make sure the path you provide is the actual location where you saved your file.
- c. Select **Read the file as a string**



3. NAME THE COMPONENTS

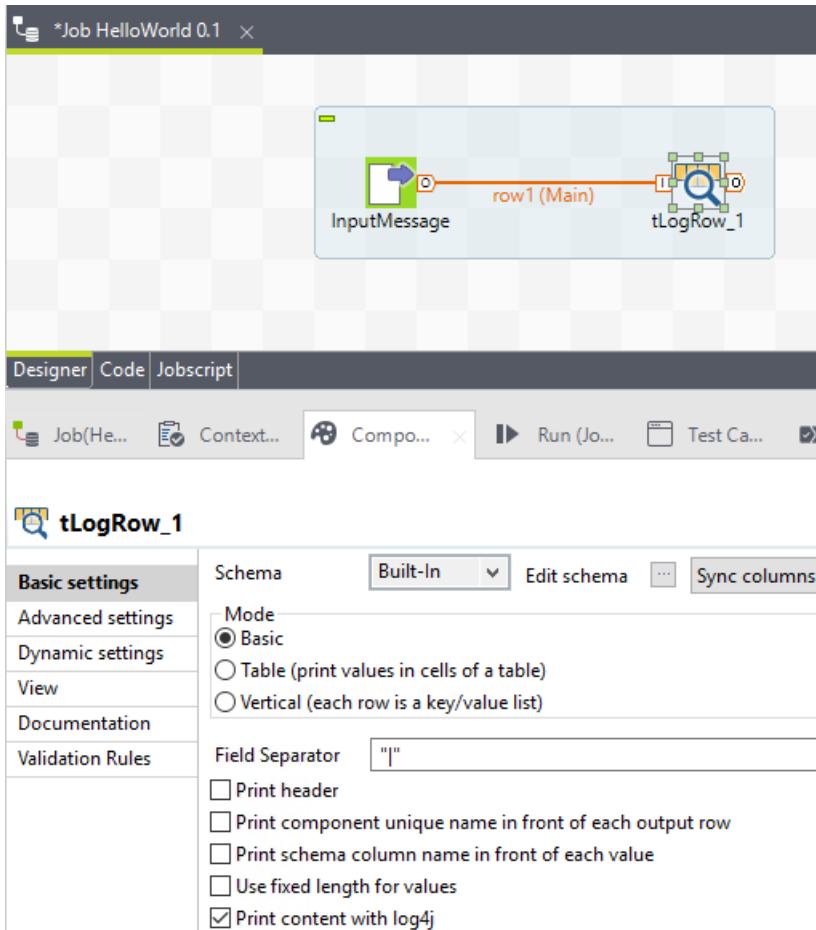
To name a component, in the **Component** view:

- a. Click the **View** tab.
- b. Replace the text in the **Label format** text box with *InputMessage*. This updates the component name in the **Designer**.



4. VERIFY THE tLogRow PROPERTIES

To check the settings for **tLogRow_1**, double-click the component and look at the settings in the **Component** view:



5. SAVE THE JOB

Notice that on the **Designer** tab, there is an asterisk (*) in front of your Job name (as you can see in the previous figure), which means the Job has unsaved changes. To save the Job, on the main toolbar, to the left, click the **Save** button ().

TIP:

You can also press **Ctrl + S** to save the Job.

Next step

You have created your first Job and now you are ready to run it.

Running a Job

Task outline

In this exercise, you learn how to run the Job you created and visualize the results in the console.

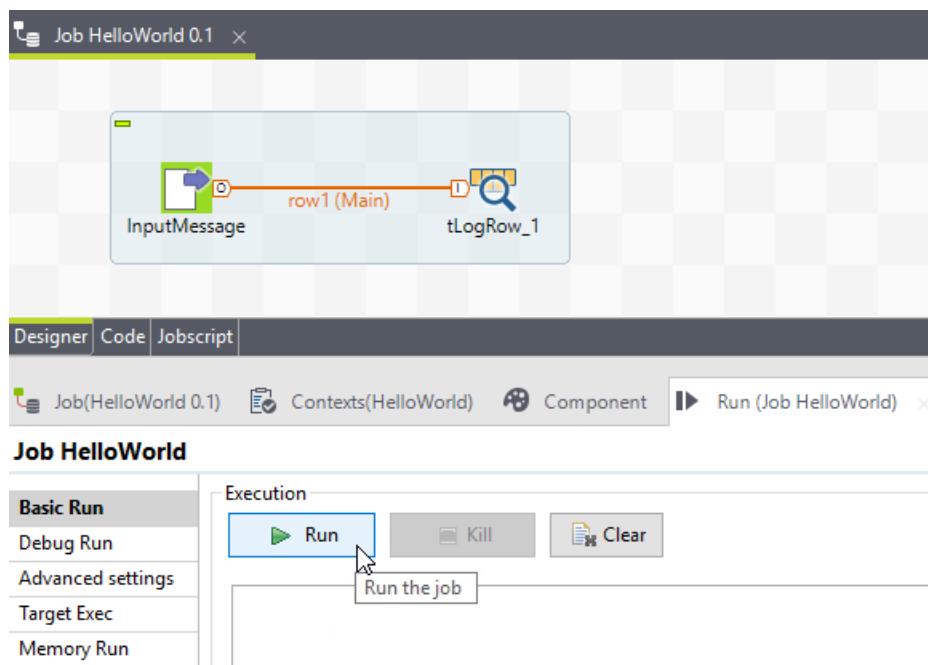
This exercise takes about 5 mins to complete.

Running a Job

1. OPEN THE RUN VIEW

To run the Job:

- In the **Designer**, open the **Run** view.
- Click **Run**.



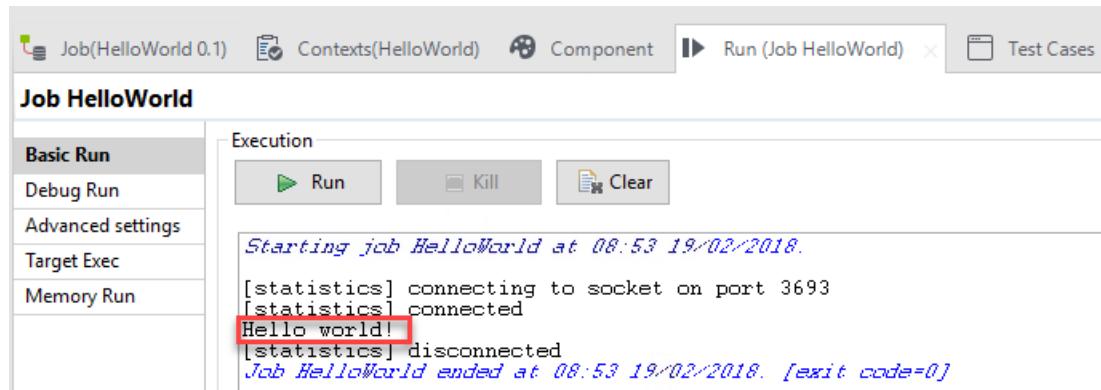
TIP:

You can also use the **Run** icon (▶) on the main toolbar to run the Job. When you have several Job tabs open in the **Designer**, make sure to select the correct Job before using the **Run** icon.

2. REVIEW THE JOB OUTPUT

Talend Studio first saves and then builds the Job. This may take a few moments. Then the Job is executed and the "Hello world!" message is displayed in the console in the **Run** view. The Job ends with an exit code = 0 which means that the Job

executed correctly.



Congratulations! You have built and run your first Talend Studio Job.

Next step

You have almost finished this section. Time for a quick review.

Review

In this lesson, you learned the basics required to create projects using Talend Studio.

You created your first Job, learned how to execute it and visualize the results in the console. You experienced different methods of adding and linking components, as well as how to configure them. You used the **Palette** to search for components, and you opened the component documentation help page.

You will have the opportunity to use other components and create more complex Jobs in the following lesson.



Using Talend Studio to Design a Data Integration Project

This chapter discusses:

Concepts	32
Overview	38
Importing External Resources into Your Project	39
Preparing Data Sources: Working with Structured Files	42
Preparing Data Sources: Reading Data from Databases	49
Joining Data From Multiple Data Sources Using tMap	55
Transforming Data	73
Review	81

Concepts



OBJECTIVES

- Import resources into your project
- Prepare the data sources: work with structured files
- Prepare the data sources: work with databases
- Join data coming from multiple sources using the **tMap** component
- Transforming data: aggregate, apply filters to data

USE CASE

Organize and transform sales data to generate report files

- Import resources into your project
- Work with structured files
- Work with databases
- Join data coming from multiple data sources
- Transform data: apply filters, expressions and aggregate functions



3

METADATA

- Stored in Repository > Metadata
- May be reused across jobs
- Metadata examples:
 - File schemas,
 - DB connections

Two screenshots of the Talend Studio interface. The left screenshot shows the 'Repository' view with a tree structure including 'LOCAL Disk', 'Code', 'SQL Templates', and 'Metadata'. The right screenshot is a detailed view of a database connection configuration for MySQL, showing fields like 'Db Type' (MySQL), 'Db Version' (MySQL 5), 'String of Connection' (jdbc:mysql://localhost:3306/training/noDatetimeStringSync=true), 'Login' (root), 'Password' (****), 'Server' (localhost), 'Port' (3306), and 'Database' (Training). There is also a 'Check' button and a 'Database Properties' section.

Metadata is stored in the **Repository > Metadata** folder and may be reused across Jobs. Indeed, some resources are good candidates for reuse, especially when they do not change from one project to the next, such as database connection settings or metadata (for example, a file schema that holds a file structure).

4

WORK WITH STRUCTURED FILES



Reading from a structured file

- Use **tFileInputDelimited**
- File name to read from
- File structure:
 - field/row separator and
 - number of header lines

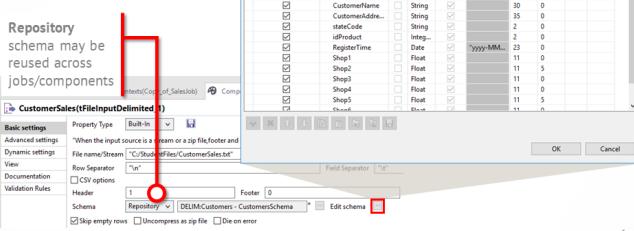


5

WORK WITH STRUCTURED FILES



- The **schema** defines how data is interpreted

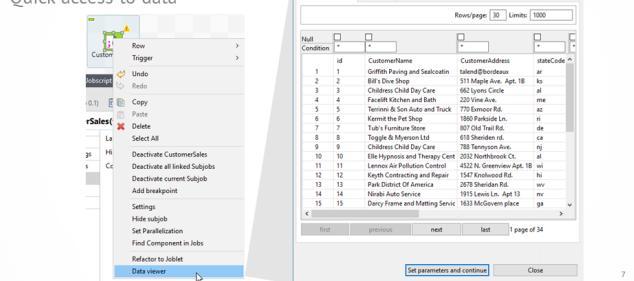


6

DATA VIEWER



Quick access to data



7

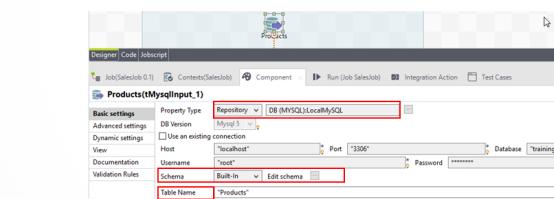
WORK WITH DATABASES



Read data from a database

- Use **tMySqlInput**
- Connect to the database using the Db connection metadata
- Edit schema and table name

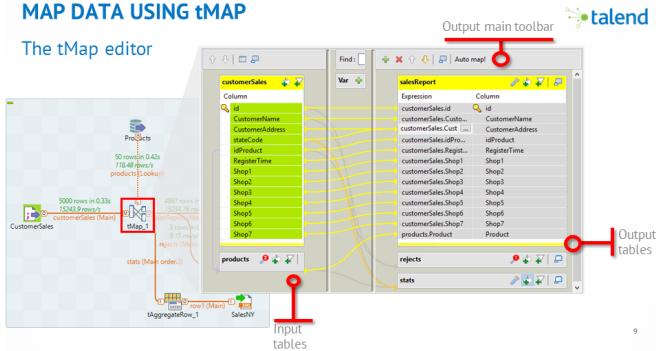
To read data from a database, you use a **tMySqlInput** component. You need to configure this component by specifying the database connection properties and database schema.



8

MAP DATA USING tMAP

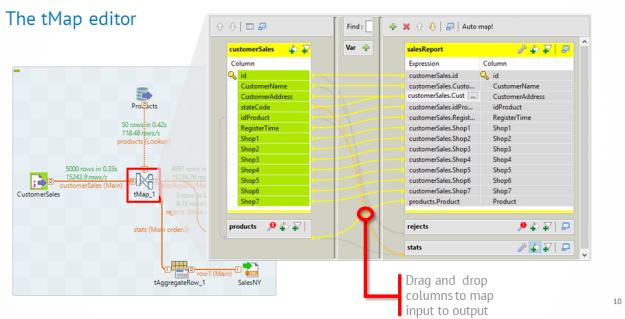
The tMap editor



The **tMap** component transforms and routes data from single or multiple sources to single or multiple destinations. Input and output are interpreted as tables.

MAP DATA USING tMAP

The tMap editor



The tMap editor allows you to map data easily from input to output by drag and drop. The arrows indicate the mapping of columns.

MAP DATA USING tMAP

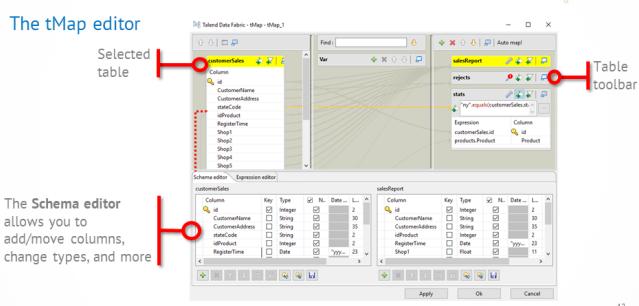
The tMap editor



The columns of the input table are an exact match of the rows entering the tMap component. In the output table, the **Expression** fields hold the matched data content, and the **Column** contains the data labels. Output rows are an exact match of the output table columns.

MAP DATA USING tMAP

The tMap editor



The Schema editor allows you to add/move columns, change types, and more

JOIN DATA SOURCES USING TMAP



- Use **tMap** to join and transform multiple inputs to multiple outputs
 - tMap is often part of a solution to various data challenges



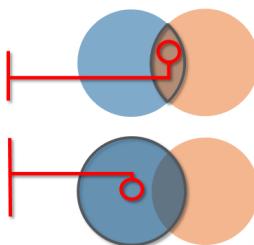
13

JOIN DATA SOURCES USING TMAP



Join Types

- A **join** combines rows from two or more tables
 - An **inner join** returns rows that only appear in both tables
 - A **left outer join** (or **left join**) returns all records from the left table with NULL column values if no match is found in the right table

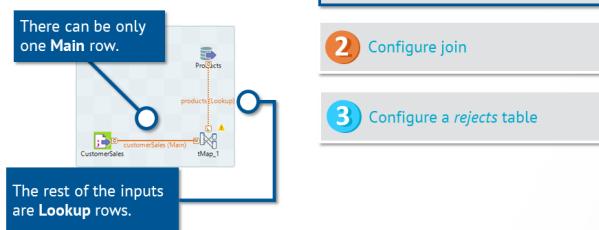


14

JOIN DATA USING TMAP

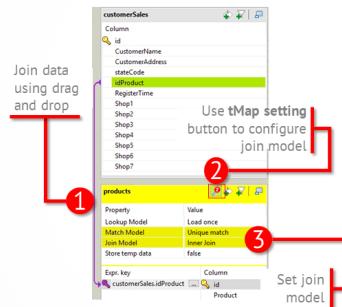


- tMap takes multiple data sources



15

JOIN DATA USING TMAP



16

JOIN DATA USING tMAP

Configure a *rejects* table to catch the columns that did not match the inner join condition.

Configure a *rejects* table to catch the columns that did not match the inner join condition.

TRANSFORM DATA USING tMAP

Create filters

Use the inline editor to create simple expressions.

Use the inline editor to create simple expressions.

TRANSFORM DATA USING tMAP

Create expressions

To easily create complex expressions, you use the **Expression builder**, which provides editing tools and a library of functions grouped by categories. In the editor, the **Variables** represent your Expression fields. Copy/paste variables to use them in expressions.

To easily create complex expressions, you use the **Expression builder**, which provides editing tools and a library of functions grouped by categories. In the editor, the **Variables** represent your Expression fields. Copy/paste variables to use them in expressions.

AGGREGATE DATA

The tAggregateRow component

Count the number of customers per sold product

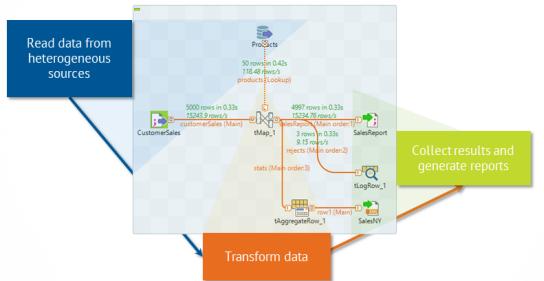
The **tAggregateRow** component receives a flow and aggregates it based on one or more columns. You can use it to create GROUP BY statements. In this example, you group data by product name and count the number of customer sales per product. You apply transformations on columns defined by **Input column position** and send the output to the columns defined by **Output column**.

EXERCISE OVERVIEW

Design a data integration project using Talend Studio



In this exercise, you learn how to design a data integration project using Talend Studio. Your role is to analyze sales data from a CRM system and output some reports that will help you prepare the next marketing campaign. You create a Job that joins data from the customer sales file and the products database and processes the results to extract statistics for the sales reports.



21

LESSON SUMMARY



tFileInputDelimited

- Reads a delimited file.

tMap

- Maps and transforms data using operators and predefined functions

tMysqlInput

- Reads from a database and extracts fields based on a query.

tAggregateRow

- Receives a flow and aggregates it based on one or more columns.

22



©2018 Talend Inc.

Overview

Use case

In this lesson, you learn how to design a data integration project using Talend Studio. You follow the classical phases of a data integration chain: extract data from heterogeneous sources, transform it and load the results in the target storage.

For this exercise, your role is to analyze sales data from a CRM system and output some reports that will help you prepare the next marketing campaign. You have been asked to focus on a specific list of products targeted by the marketing policy in your company.

The data you have at hand is:

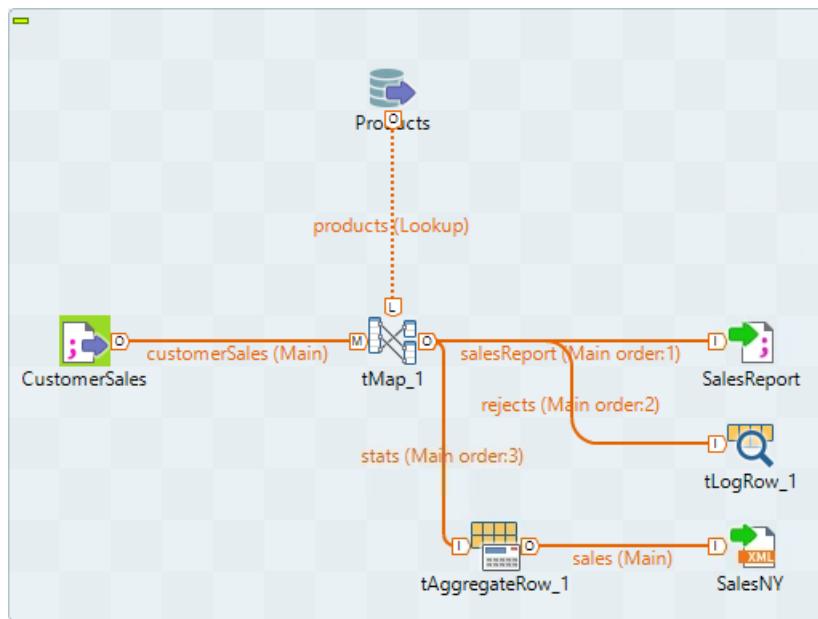
- » A list of targeted products that you can read from an existing Products database
- » Sales data from a CRM system stored in the CustomerSales.txt file

You must create a Job that joins data from the customer sales file and the products database and processes the results to extract statistics for the sales reports.

As a result, you want to:

- » Generate a global sales report
- » Generate a report containing statistics on the number of sold products for the city of New York
- » Display in the console products that were bought by customers but not targeted by the current marketing campaign

When you finish, the Job looks like this:



Objectives

After completing this lesson, you will be able to:

- » Import resources into your project
- » Work with structured files
- » Work with databases
- » Join data coming from multiple data sources
- » Process data flows

Next step

The next step is to prepare the data sources.

Importing External Resources into Your Project

Task outline

In this section, you learn how to import external resources into your project. Indeed, some resources are good candidates for reuse, especially when they do not change from one project to the next, such as database connection settings or metadata (for example, a file schema that holds a file structure).

For your convenience, all the metadata you need in this lesson has already been created. However, you need to import it into your project.

This exercise takes about 5 mins to complete.

Importing items into your project

Before you start this exercise, copy the content of the **C:/IntroTalendStudio** folder to **C:/StudentFiles** on your VM. The folder contains the files CustomerSales.txt, LoadProducts.zip and refProducts.csv. You are now ready to import items into your project.

The items to import are:

- » The **LocalMySQL 0.1** database connection settings for the connection to the local MySQL database instance
- » The **LoadProducts 0.1** Job, which allows you to create the Products table
- » The **Customers 0.1** and **Products 0.1** file schemas, which hold the structure of the Customers.csv file and schema of the Products database

To import resources from a different Studio project, you can use the Import items button on the main toolbar.



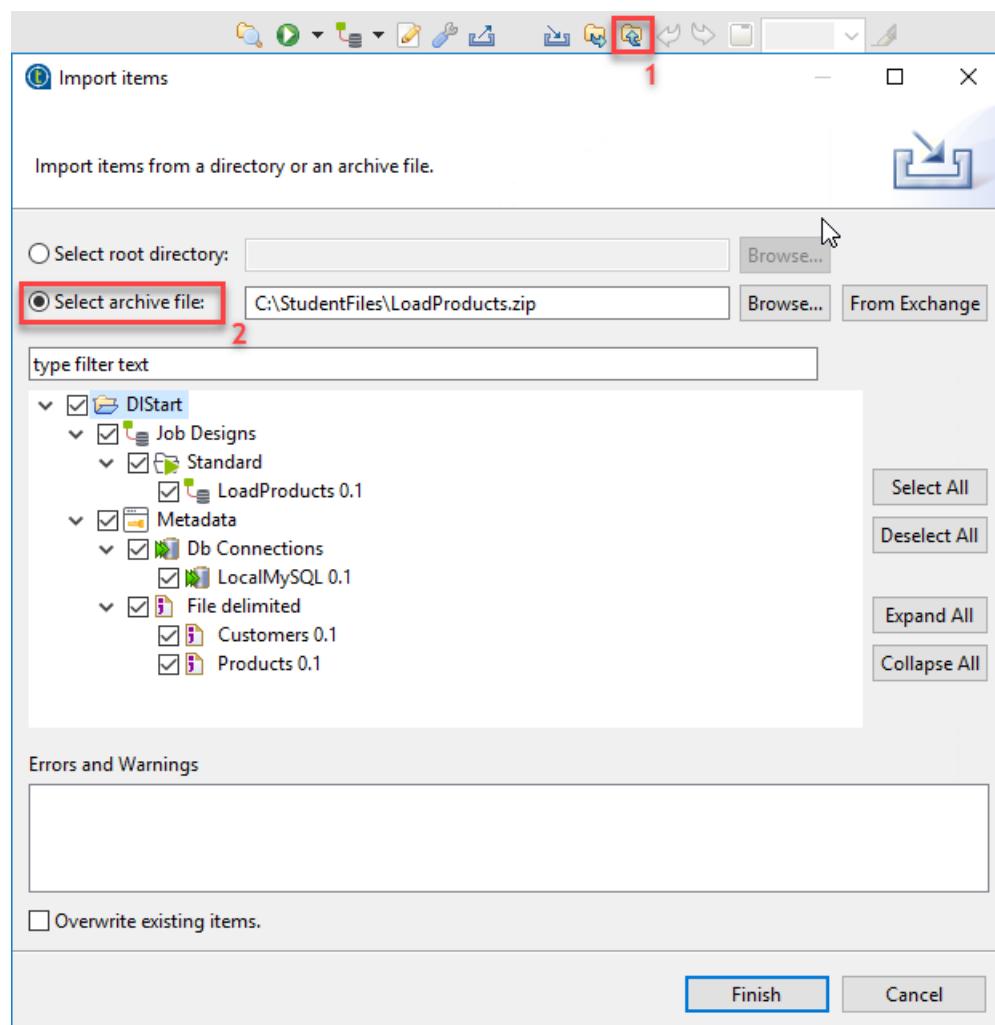
The metadata to import is located at: C:/StudentFiles/LoadProducts.zip.

1. IMPORT THE METADATA

To import the metadata:

- a. Click the **Import items** icon on the main Studio toolbar.
- b. In the dialog box that opens, select the check box **Select archive file**

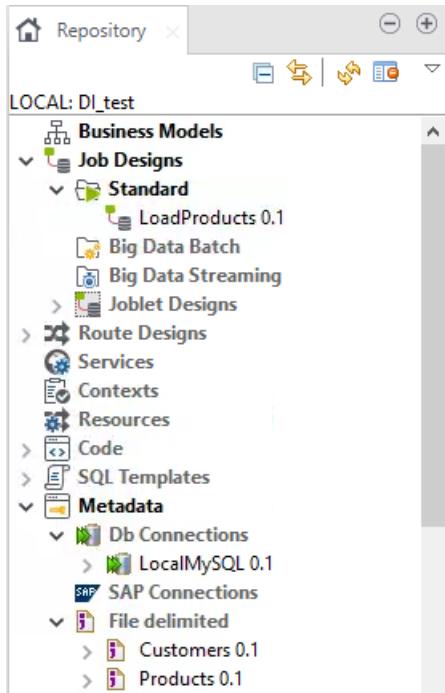
- c. Click **Browse** and select the zip file located at C:\StudentFiles\LoadProducts.zip.



- d. Select the elements to import: **LoadProducts 0.1**, **LocalMySQL 0.1**, **Customers 0.1** and **Products 0.1**.
e. Click **Finish**.

NOTE: When you import the metadata into your project, it is stored in the **Repository > Metadata** folder. Jobs are stored in the **Job Designs > Standard** folder.

When you finish this task, the **Repository** for your project looks like this:



Next step

You have imported the metadata you need and are ready to prepare the data sources.

Preparing Data Sources: Working with Structured Files

Task outline

In this lesson, you learn how to read and interpret data coming from a structured file (in this case, the CustomerSales.txt file from the CRM system).

This exercise takes about 15 minutes to complete.

Reading data from a structured file

To read data from a structured file (for example, CustomerSales.txt) you use a **tFileInputDelimited** component. You configure this component by specifying where the file is located and how the data inside is organized.

The **tFileInputDelimited** reads a delimited file row by row, splits them into fields, and sends the fields, as defined in the schema, to the next component.

Once the data is read from the file, you must give a meaning to it so that you can process it further. You can do this by assigning a schema to your component.

A schema holds information about the way data is interpreted, and is therefore considered metadata. A schema defines the columns/rows to be processed and passed on to the next component.

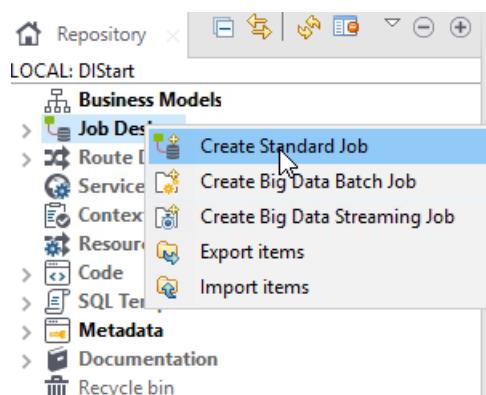
For your convenience, all the schemas you need for this lesson have been already created. Before continuing, make sure you have imported all necessary metadata into your project, as suggested in Import "Existing Resources into Your Project".

To get started, you must create a Job.

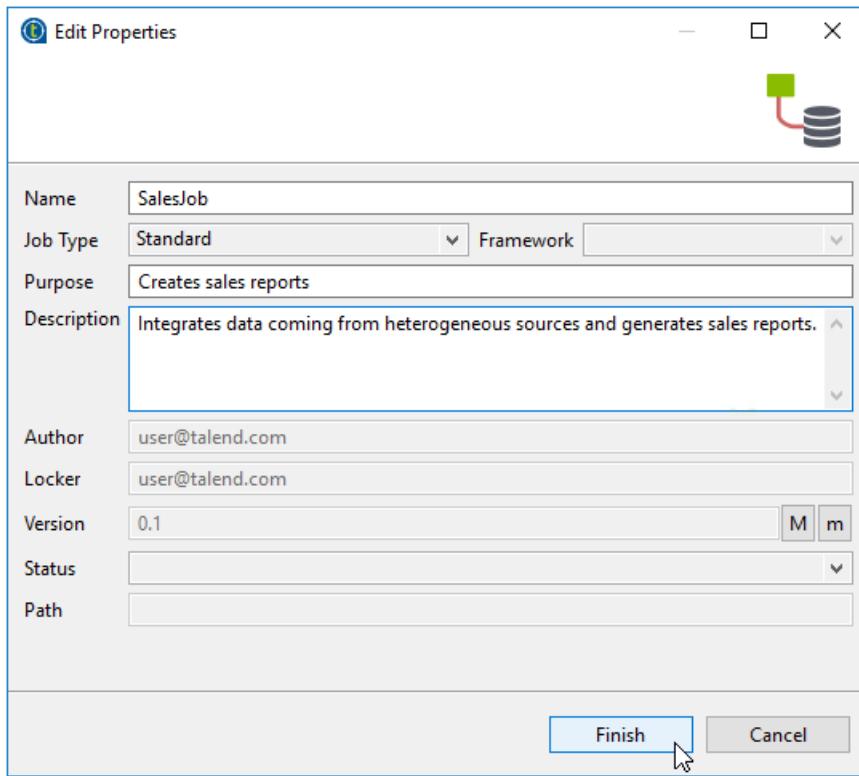
1. CREATE A JOB

To create a Job:

- Right-click **Repository > Job Designs**.
- Select **Create Standard Job**. The **New Job** dialog box opens.



- In the **Name** text box, enter *SalesJob*.
- Fill-in the **Purpose** and **Description** text boxes.
- Click **Finish**.



The Job opens in the Designer.

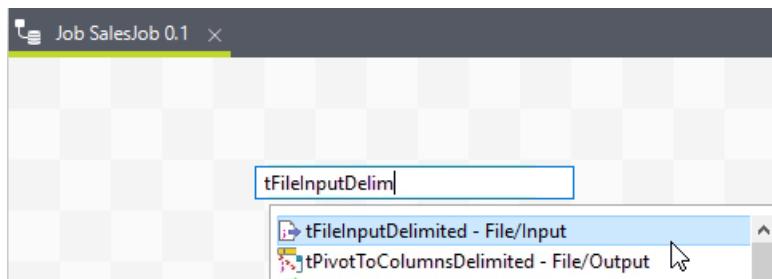
TIP:

When creating a Job, only the name is mandatory, but it is a good practice to complete the remaining fields to help document the Job.

2. ADD A tFileInputDelimited COMPONENT TO READ A STRUCTURED FILE

To add a **tFileInputDelimited** component:

- Click anywhere in the **Designer**
- Start entering *tFileInputDelimited*. As you type, components matching your search term appear.



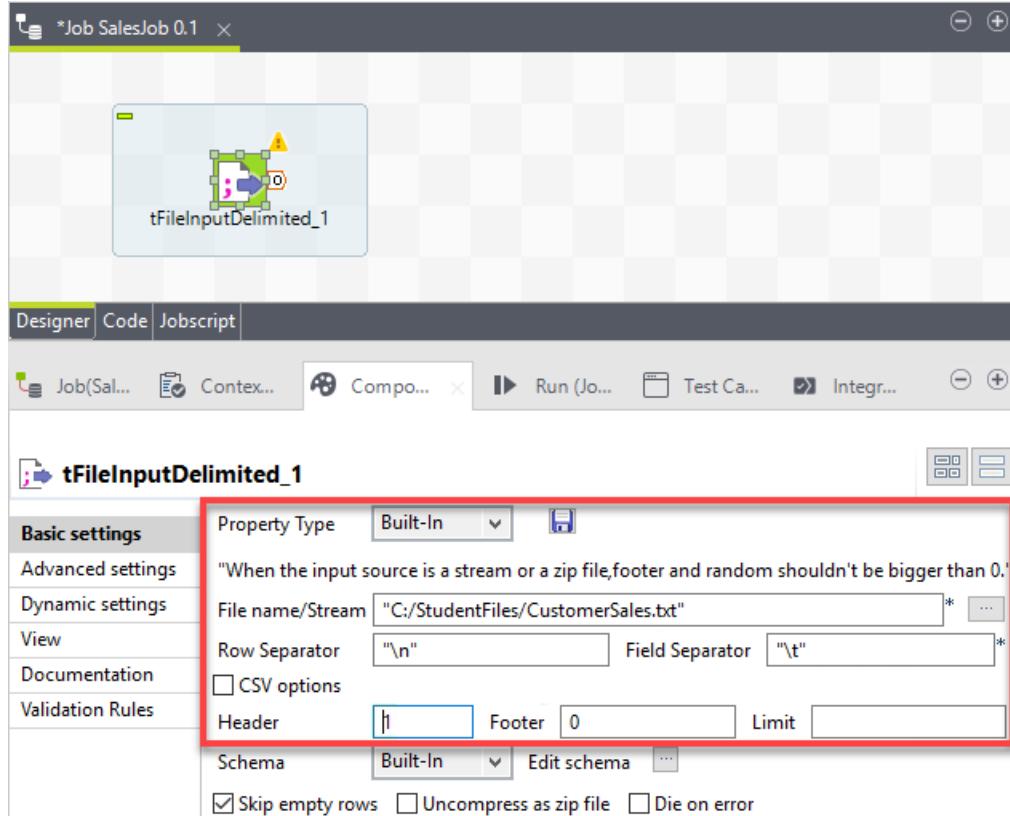
- Double-click **tFileInputDelimited** on the list. The component appears in the Designer.

3. SET BASIC FILE CONFIGURATION

To configure a **tFileInputDelimited** component, you must specify the name of the file to read from, and its structure.

- Double-click the **tFileInputDelimited_1** component to open the **Component** view.
- To set the **File name**, click [...] and select the *C:/StudentFiles/CustomerSales.txt* file path.

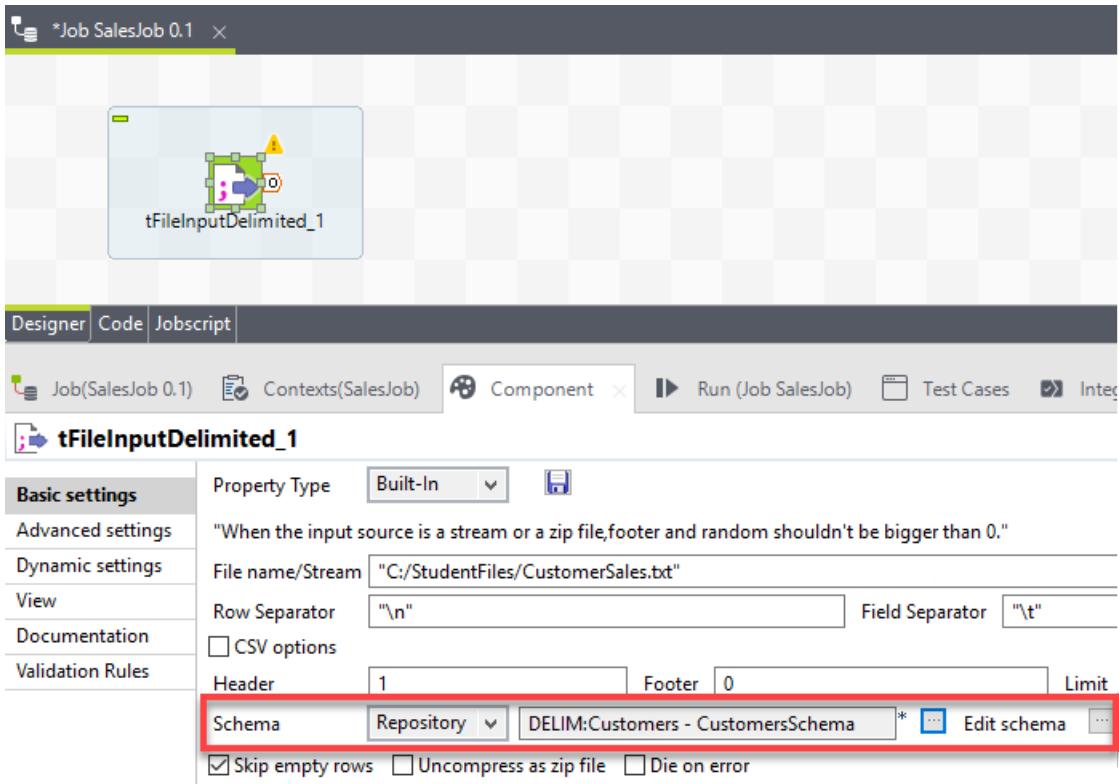
- c. Leave the default **Row Separator** as "\n".
- d. Set the tab character as **Field Separator**: "\t", to specify that the file contains tab-separated values.
- e. In the **Header** text box, enter 1. This means the first line in the file is ignored, which is what you want because this line does not contain actual data.



You are almost finished. You just need to configure the schema, which will let you process data further along.

4. CONFIGURE THE SCHEMA

- a. For **Schema**, select *Repository*.
- b. To its right, click [...] and select **Metadata > File delimited > Customers 0.1 > CustomersSchema**.



TIP:

In which case to use **Built-In** and **Repository**:

Built-In: Information is stored locally in the Job. Use Built-in for information that you need once or very rarely.

Repository: Information is stored in the repository. Use Repository for information that you need repeatedly in multiple components or Jobs.

5. VIEW THE SCHEMA

To visualize the schema:

- a. To the right of **Edit schema**, click [...]. A dialog box opens.
- b. Select **View schema** and click **OK**.

A dialog box showing the schema opens. Take a few moments to look at the schema. You can see that it is composed of columns with a given type and other related properties. Close the window.

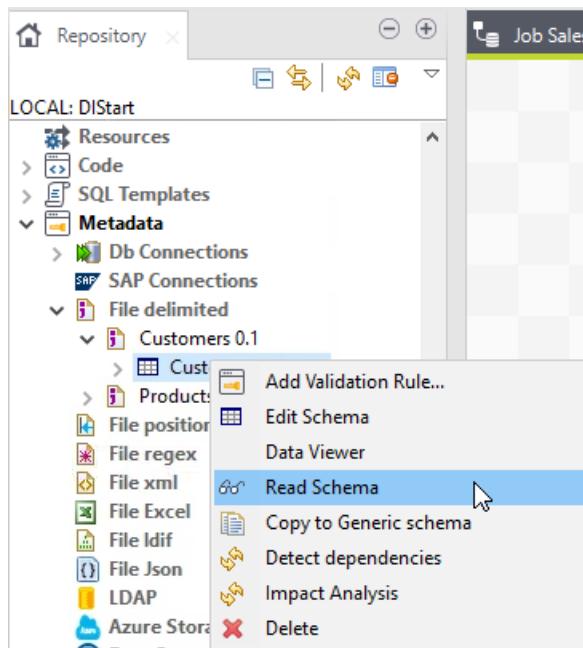
Schema of tFileInputDelimited_1

tFileInputDelimited_1

Used Column	Column	Key	Type	N..	Date Patter...	Length	Preci...	Def...	Com...
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> id	<input checked="" type="checkbox"/>	Integ...	<input checked="" type="checkbox"/>		2	0		
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> CustomerName	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		30	0		
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> CustomerAddress	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		35	0		
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> stateCode	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		2	0		
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> idProduct	<input type="checkbox"/>	Integ...	<input checked="" type="checkbox"/>		2	0		
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> RegisterTime	<input type="checkbox"/>	Date	<input checked="" type="checkbox"/>	"yyyy-MM...	23	0		
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> Shop1	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	0		
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> Shop2	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	5		
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> Shop3	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	0		
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> Shop4	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	0		
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> Shop5	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	5		
<input type="checkbox"/>	<input checked="" type="checkbox"/> Shop6	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	0		

OK **Cancel**

You can also view the schema right in the **Repository > Metadata**. Locate the schema you are interested in, right-click, and select **Read Schema**.

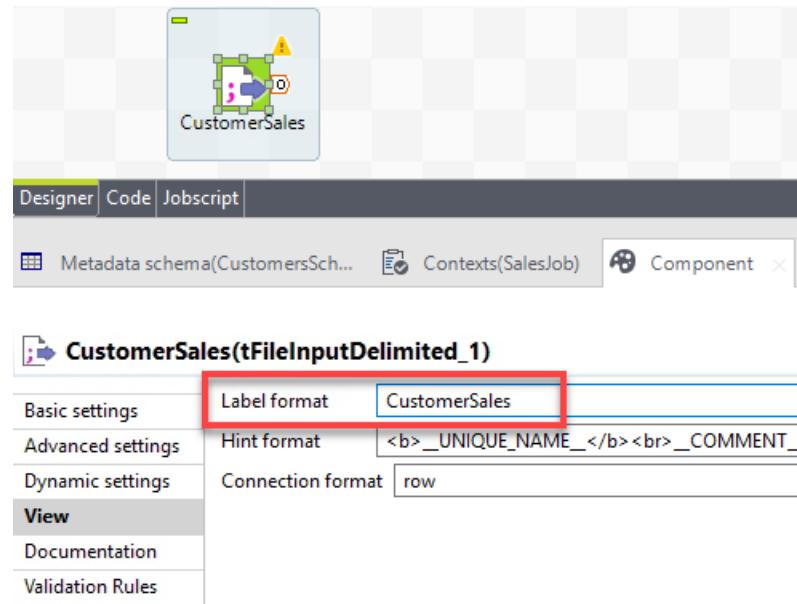


6. NAME THE COMPONENT

To name the component:

- Click the **View** tab in the **Component** view.
- Replace the text in the **Label format** text box with *CustomerSales*. This automatically updates the component

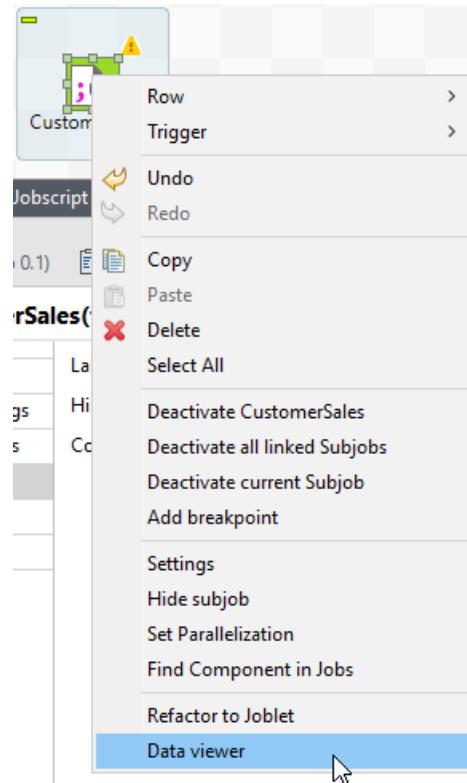
name in the design workspace.



7. VISUALIZE DATA USING THE DATA VIEWER

To visualize data using the Data viewer:

- Right-click the **CustomerSales** component
- Select **Data viewer**.



A window displays a preview of the data retrieved from the CustomerSales.txt file. The **File Content** tab displays the raw content of the file. Click **Close**.

The screenshot shows a 'Data Preview' window titled 'tFileInputDelimited_1'. The window has two tabs at the top: 'Result Data Preview' (selected) and 'File Content'. Below the tabs are buttons for 'Rows/page:' (set to 30) and 'Limits:' (set to 1000). The main area is a table with 15 rows of data. The columns are labeled: Null, Condition, id, CustomerName, CustomerAddress, and stateCode. The data includes various business names like 'Griffith Paving and Sealcoatin', 'Bill's Dive Shop', and 'Darcy Frame and Matting Servic'. The 'stateCode' column shows abbreviations like 'ar', 'ks', 'al', 'me', 'az', 'ri', 'de', 'ca', 'nj', 'al', 'wi', 'hi', 'wv', 'nv', and 'ga'. At the bottom of the table, there are navigation buttons: 'first', 'previous', 'next', 'last', and '1 page of 34'. Below the table, there are two buttons: 'Set parameters and continue' (highlighted with a blue border) and 'Close'.

Null	Condition	id	CustomerName	CustomerAddress	stateCode
1	*	1	Griffith Paving and Sealcoatin	talend@bordeaux	ar
2	*	2	Bill's Dive Shop	511 Maple Ave. Apt. 1B	ks
3	*	3	Childress Child Day Care	662 Lyons Circle	al
4	*	4	Facelift Kitchen and Bath	220 Vine Ave.	me
5	*	5	Terrinni & Son Auto and Truck	770 Exmoor Rd.	az
6	*	6	Kermit the Pet Shop	1860 Parkside Ln.	ri
7	*	7	Tub's Furniture Store	807 Old Trail Rd.	de
8	*	8	Toggle & Myerson Ltd	618 Sheridan rd.	ca
9	*	9	Childress Child Day Care	788 Tennyson Ave.	nj
10	*	10	Elle Hypnosis and Therapy Cent	2032 Northbrook Ct.	al
11	*	11	Lennox Air Pollution Control	4522 N. Greenview Apt. 1B	wi
12	*	12	Keyth Contracting and Repair	1547 Knolwood Rd.	hi
13	*	13	Park District Of America	2678 Sheridan Rd.	wv
14	*	14	Nirabi Auto Service	1915 Lewis Ln. Apt 13	nv
15	*	15	Darcy Frame and Matting Servic	1633 McGovern place	ga

Next step

You have set a **tFileInputDelimited** component to read data from a structured file, and now you are ready to read data from a database.

Preparing Data Sources: Reading Data from Databases

Task outline

In this lesson you learn how to read and interpret data from the Products table. You run the Job you imported to create the Products table and data, and then you can start using it.

This section takes about 15 minutes to complete.

Reading data from a database

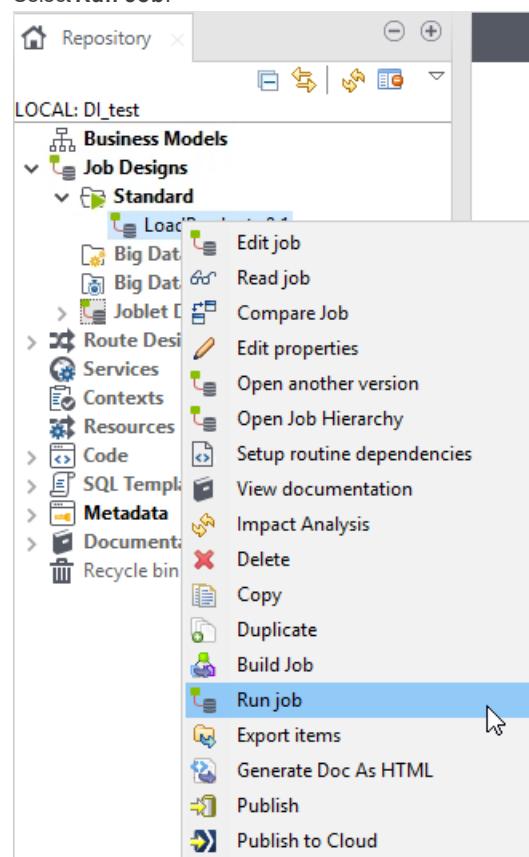
First you must create the **Products** table. The LoadProducts 0.1 Job you imported allows you to create the table and add some data to it. After you run this Job, the table is created inside the existing "training" database.

To read data from a database, you use a **tMysqlInput** component. You must configure this component by specifying the database connection properties and database schema. The **tMysqlInput** component executes a DB query with a strictly defined order which must correspond to the schema definition. Remember that you already imported this metadata into your project; there is no need to re-create it. You only need to assign it to the **tMysqlInput** component.

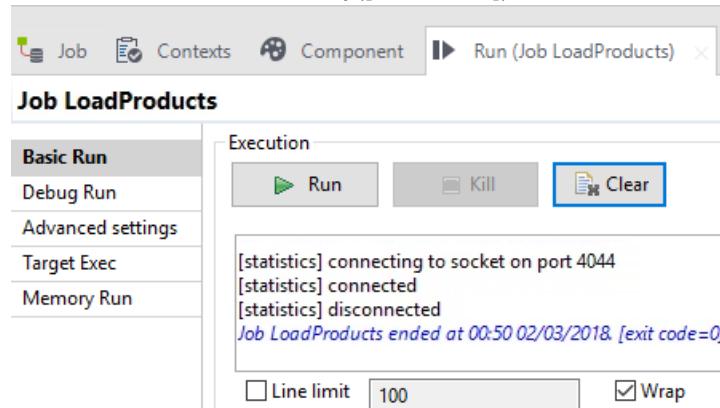
1. CREATE THE PRODUCTS TABLE AND DATA

To create the Products table, run the LoadProducts 0.1 Job:

- In **Repository > Job Designs > Standard**, right-click **LoadProducts 0.1**
- Select **Run Job**.



- c. Make sure the Job executes correctly ([exit code=0])

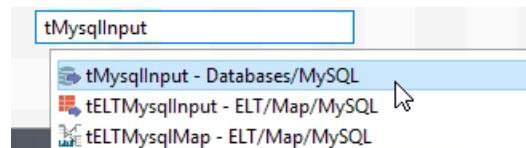


The Products table has been created.

2. ADD A tMysqlInput COMPONENT TO READ DATA FROM A DATABASE

To add a **tMysqlInput** component:

- In the **Designer**, click in the space to the right of the **CustomerSales** component.
- Start entering *tMysqlInput* to list matching components.

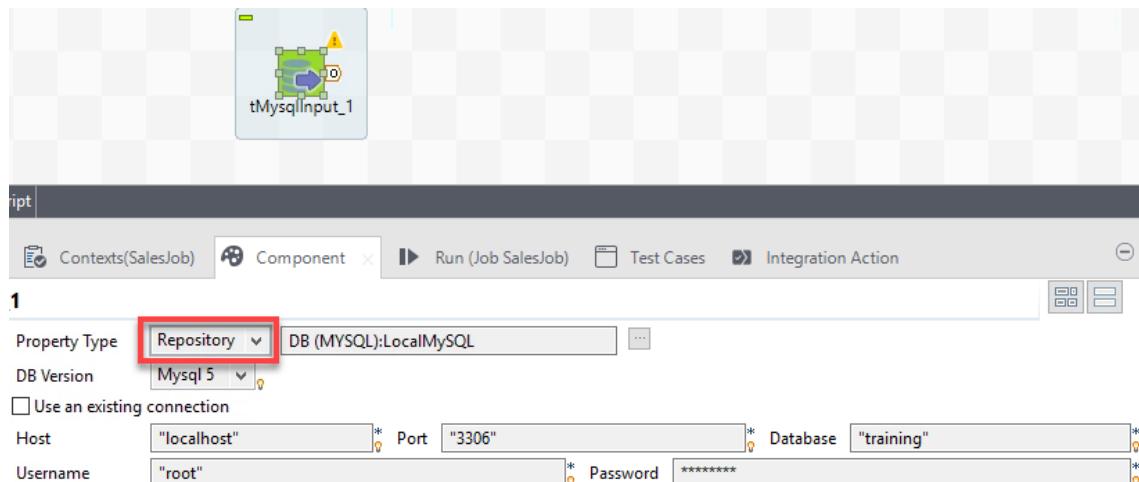


- On the list, double-click **tMysqlInput**. The component appears in the Designer.

3. CONFIGURE THE tMysqlInput COMPONENT

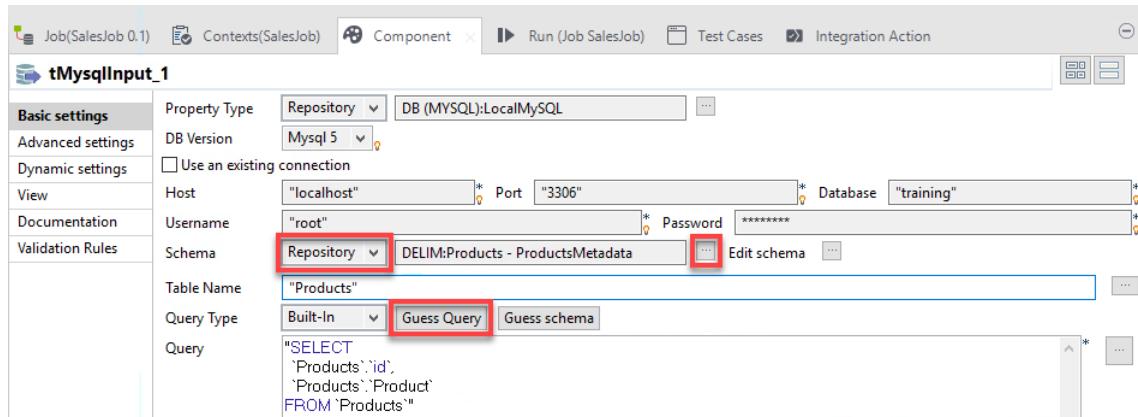
Configure the tMysqlInput component using the LocalMySQL database connection:

- Double-click in the middle of **tMysqlInput_1** to open the **Component** view.
- For **Property Type**, select *Repository*. This allows you to reuse existing metadata.
- Click [...] and select **Metadata > Db Connections > LocalMySQL 0.1**. Click **OK**. Notice that the database connection properties (host, database, username) are automatically filled in.
- For **DB Version**, select *Mysql 5*.



4. SET THE COMPONENT SCHEMA

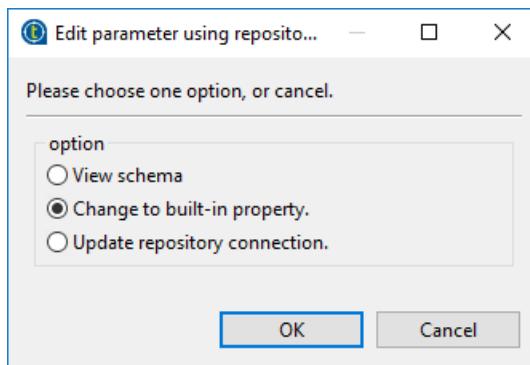
- For **Schema**, select *Repository*.
- To its right, click [...] and select **Metadata > File delimited > Products 0.1 > ProductsMetadata**
- In the **Table Name** text box, enter *Products*
- Click the **Guess Query** button. This generates a default query that retrieves all columns from the table, which is exactly what you want.



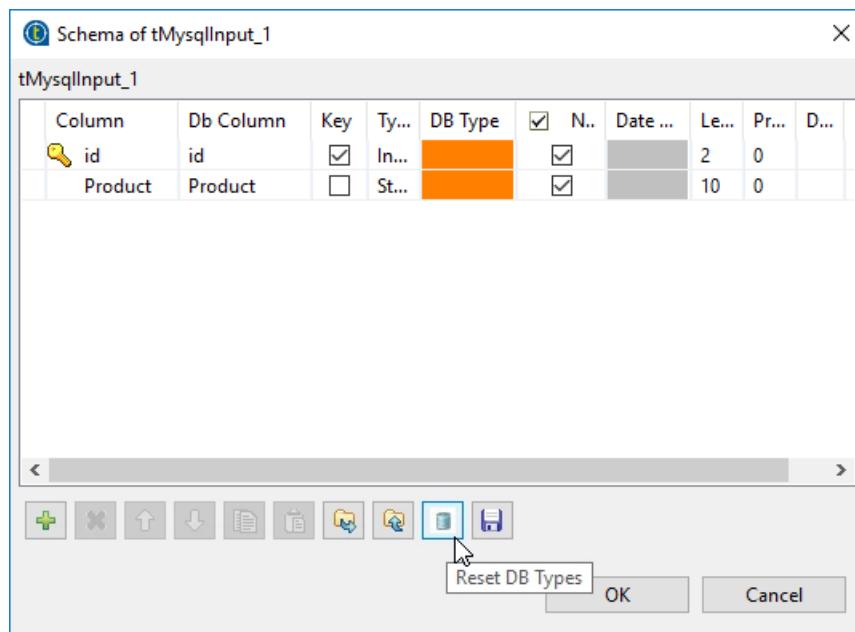
5. EDIT THE SCHEMA

Note that the **tMysqlInput_1** component has a warning that the schema database type (column **DB Type**) is not correct. This column in the schema defines the type of rows this component sends to the component connected to it. Therefore, you must specify a type for this column.

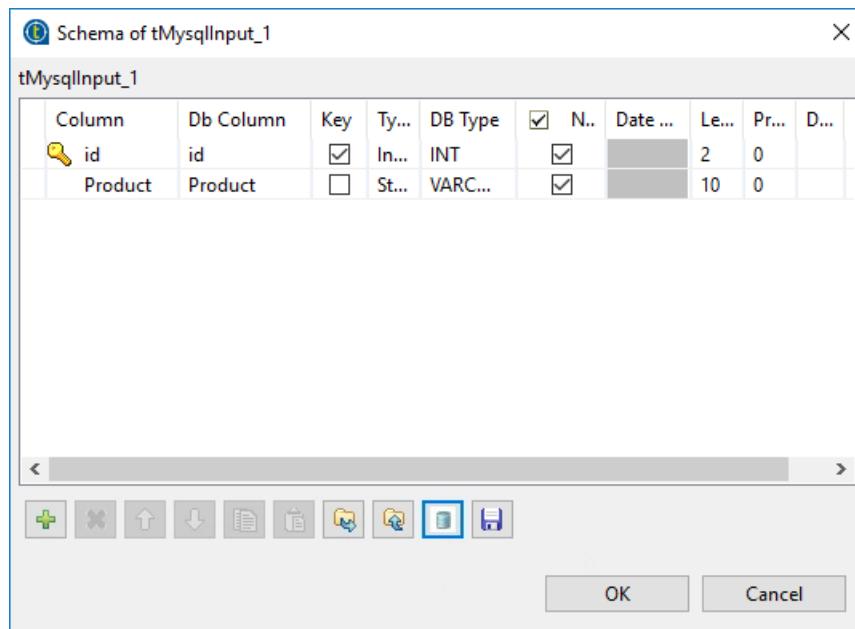
- To the right of **Edit schema**, click [...]. A dialog box opens with schema editing options.



- Select the radio button for **Change to built-in property**. This allows you to work on a local copy of the schema rather than the shared version stored in the repository. Click **OK**. A dialog box opens displaying the component schema. Notice that the **DB Type** column is empty and highlighted in orange.



- c. Click the **Reset DB Types** button as shown in the screenshot. The column types are automatically filled in. Click **OK**.



6. NAME THE COMPONENT

To name the component:

- a. In the **Component** view, select the **View** tab.
- b. Replace the text in the **Label format** text box with *Products*. This automatically updates the component name in the Designer.

Products(tMysqlInput_1)

- Basic settings
- Advanced settings
- Dynamic settings
- View**
- Documentation
- Validation Rules

Label format: Products

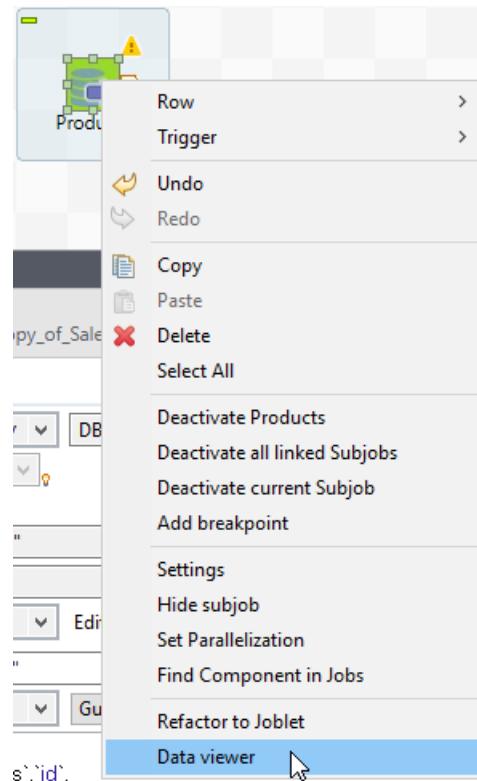
Hint format: _UNIQUE_NAME_
COMMENT

Connection format: row

1. VISUALIZE DATA USING THE DATA VIEWER

To visualize data using the **Data viewer**:

- Right-click in the middle of the **Products** component.
- Select **Data viewer**.



A window opens displaying a preview of the data retrieved from the Products table. Click **Close** when you are finished.

The screenshot shows a 'Data Preview' window titled 'tMysqlInput_1'. At the top, there are settings for 'Rows/page:' (30) and 'Limits:' (1000). Below this is a table with two columns: 'id' and 'Product'. The table contains 16 rows of data, each with an id from 1 to 16 and a corresponding product name. At the bottom of the table, there are navigation buttons: 'first', 'previous', 'next', and 'last', followed by the text '1 page of 2'. Below the table, there are two buttons: 'Set parameters and continue' (highlighted with a blue border) and 'Close'.

id	Product
1	Product 1
2	Product 2
3	Product 3
4	Product 4
5	Product 5
6	Product 6
7	Product 7
8	Product 8
9	Product 9
10	Product 10
11	Product 11
12	Product 12
13	Product 13
14	Product 14
15	Product 15
16	Product 16

Next step

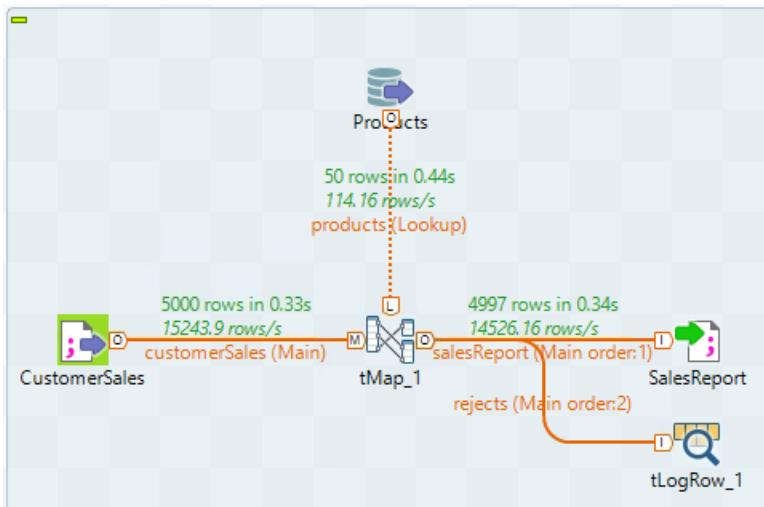
Earlier, you configured a component to read data from a structured file, and in this lesson, you learned to read data from a database using the **tMysqlInput** component. You are now ready to join data coming from multiple sources.

Joining Data From Multiple Data Sources Using tMap

Task outline

In this section, you learn to join data coming from multiple sources. Most enterprises have data in multiple locations and need to combine it, either in order to store it in a unified format or to process it consistently.

This exercise takes about 45 minutes to complete. At the end, the Job looks like this:



Mapping data using a tMap component

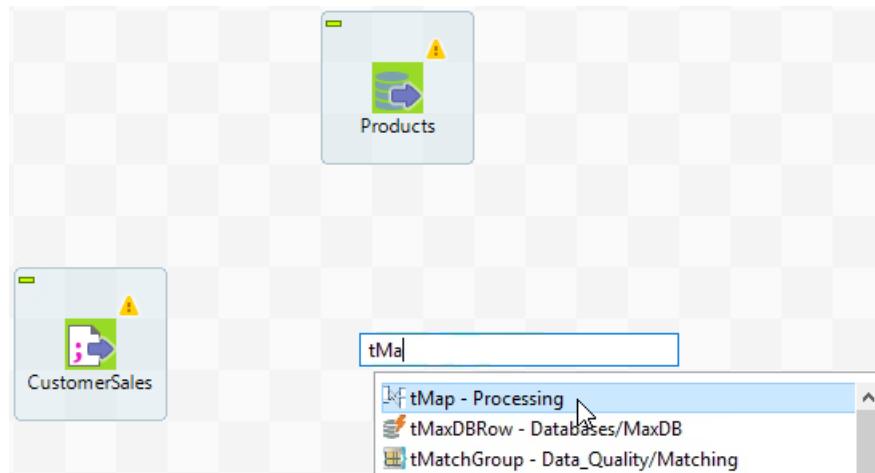
The **tMap** component maps data from input to output, with the capability to perform a variety of transformations on the data. It transforms and routes data from single or multiple sources to single or multiple destinations. You will find yourself using it as the centerpiece of almost all your Talend Data Integration Jobs.

In this exercise, you join data using **tMap** and output results to generate reports. You can output data to several structured file formats, for example **tFileOutputDelimited** and **tFileOutputXML**.

1. ADD A tMap COMPONENT

To add a tMap component:

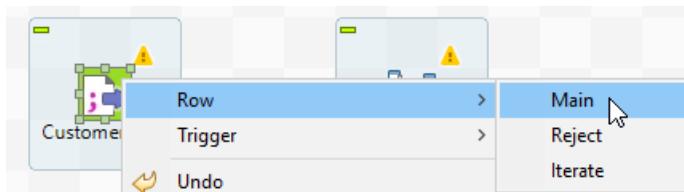
- Inside the **Designer**, click to the right of the **CustomerSales** component and below the **Products** component.
- Start entering *tMap* to bring up components matching your search term.



- On the list, double-click **tMap**. The component appears in the workspace.
2. CONNECT THE INPUT TO THE tMap COMPONENT

To connect the input to the **tMap** component:

- Right-click in the middle of **CustomerSales** and select **Row > Main**. A connection link is visible between the input component and mouse pointer. You can now link to the output component.



- Click the output component, **tMap**. The two components are connected.
- Right-click the **Products** component and select **Row > Main**.
- Click the output component (**tMap**) to link the components.

WARNING:

Notice that the second row you added automatically becomes a **Lookup** row. You can provide multiple input sources for a tMap component, but only one can be the **Main** row; the rest are **Lookup** rows.

Notice the warning displayed for the tMap component. The tMap component transforms data, so it needs an output destination.

3. NAME THE ROWS

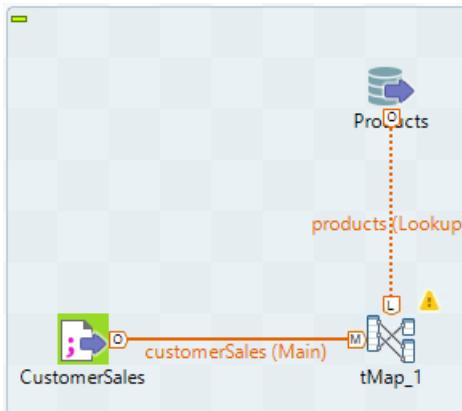
To name the rows between inputs and the **tMap** component:

- Click **row1** (the link between **CustomerSales** and **tMap_1**) to select the row. When selected, the row name is highlighted in a frame.
- Click once in the highlighted frame to activate editing and name the row *customerSales*.



- c. Repeat the same process for the **row2** (the link between **Products** and **tMap_1**) and name it *products*

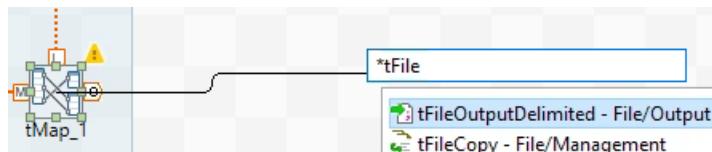
Your Job looks like this:



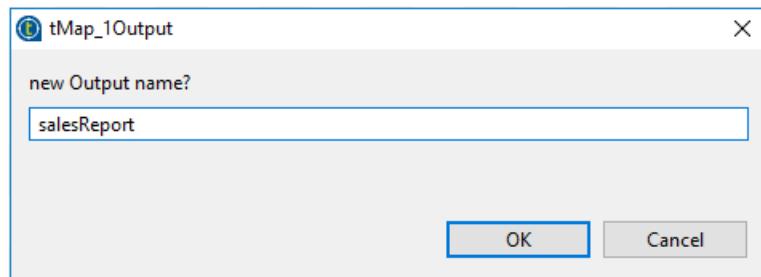
4. CREATE AN OUTPUT FILE

You want to create a file that is the sales report, containing the output you get from mapping the initial data. To write to a structured file, use a **tFileOutputDelimited** component. The **tFileOutputDelimited** component outputs the input data to a delimited file according to the defined schema.

- a. Right-click **tMap_1**. Holding the right-click, move the mouse pointer outside the component. When you release the mouse, a search box is displayed with a list of suggested components.



- b. Start entering **tFileOutputDelimited** in the search box.
 c. Double-click **tFileOutputDelimited** to add it to your Job. A dialog box opens, asking you to enter a name for the link between **tMap** and the output file
 d. Name the link *salesReport* and click **OK**.

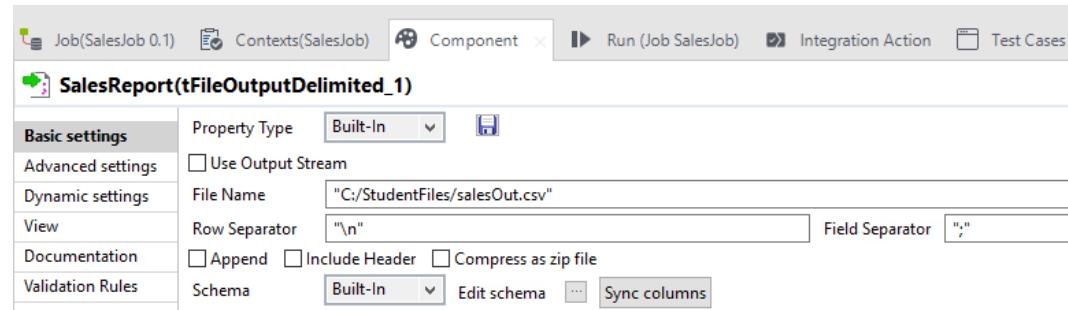


5. CONFIGURE THE OUTPUT FILE

To configure the output file:

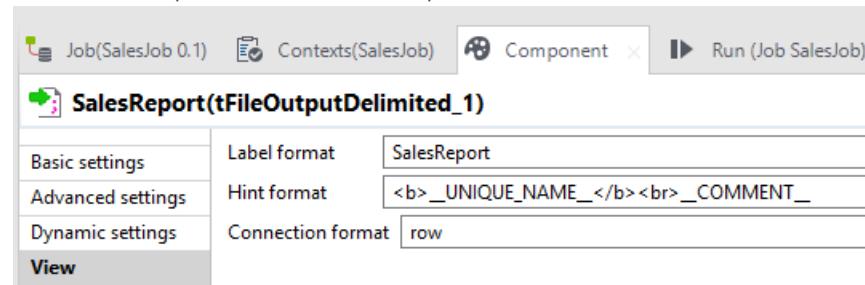
- a. Double-click the **tFileOutputDelimited_1** component to open the **Component** view.

- b. In the **File name** text box, enter "C:/StudentFiles/salesOut.csv"

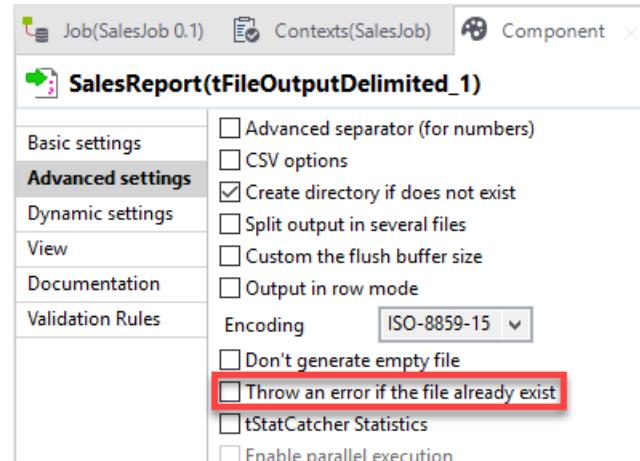


- c. Leave the other settings as configured.

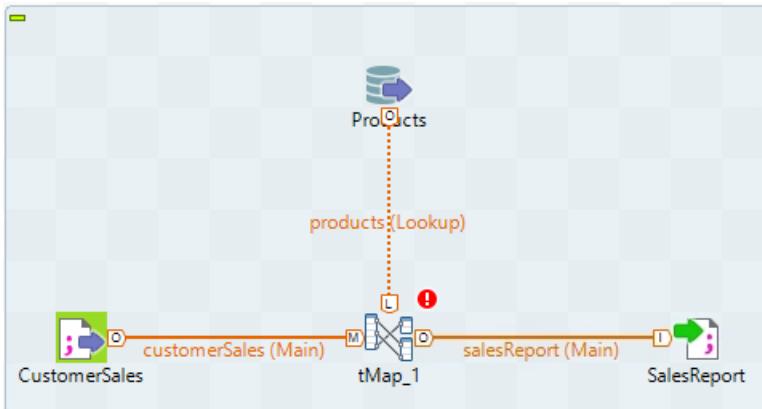
- d. To name the component, on the **View** tab, replace the text in the **Label format** text box with *SalesReport*



- e. By default, this type of component is configured to throw an error if the file already exists. Because you want to overwrite the file each time you run the Job, on the **Advanced settings** tab, deselect the check box **Throw an error if the file already exist**.



When you are finished, the Job looks like this:



Notice that the upper right corner of **tMap_1** shows an error icon. This is due to the fact that you did not define the output you expect from the mapping. You do this by configuring the tMap component.

Configuring the tMap output table

1. DEFINE THE OUTPUT TABLE

Data from the CustomerSales file and Products table is transmitted through the links customerSales and products, as rows, to the **tMap** component. The goal of this step is to define which incoming data to pass to the output.

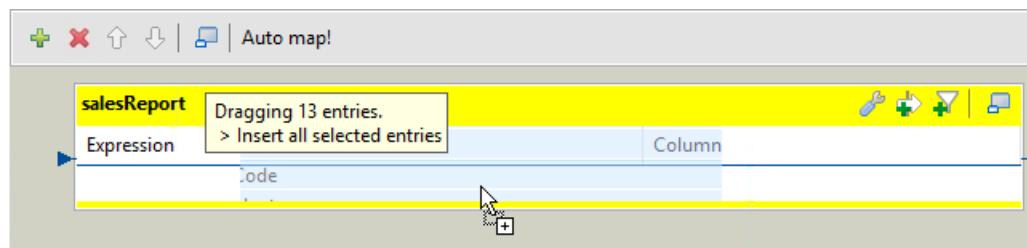
Define the output you expect from the mapping:

- Double-click **tMap_1** to open the mapping editor. Notice the **customerSales** and **products** tables on the left (the input flows). You can recognize their schemas because they match the schemas of the input components. Notice on the right that a table has been created with the same name as the output link **salesReport**. You now need to define the schema for this output table.
- Click on the **customerSales** toolbar at the top. The toolbar turns to yellow and the columns are all selected and highlighted in gray.

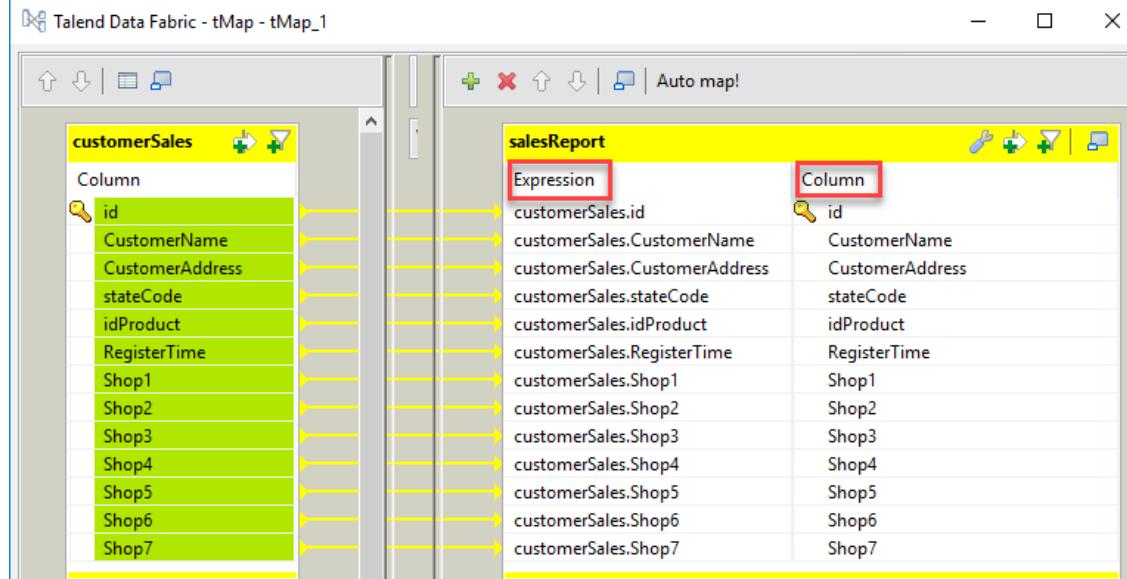
You can drag and drop them to the output table but before doing so, carefully read the instructions in the next step.

- Click anywhere in the grayed area to grab the columns and drag them to the **salesReport** table. A blue line with a pointed arrow shows the exact position of the mouse pointer. When the arrow is placed in an appropriate pos-

ition, release the mouse.



When you are finished, the tMap looks like this:



Notice the arrows indicating the mapping of columns. In the **salesReport** table, notice the **Expression** that holds the actual data content, and the **Column** which contains the data labels. At this point, the data and column names are an exact match to those in the **customerSales** table.

2. CREATE EXPRESSIONS

You want to merge the address with the state code and create a single field that contains the full address of the customer. Because some of the state codes in the file are in lowercase, you want to set them all to uppercase. The expression you want to build is:

```
customerSales.CustomerAddress + " " + StringHandling.UPCASE("customerSales.stateCode")
```

Talend Studio allows you to create expressions directly in the expression fields, expression editor and Expression builder.

NOTE:

To easily create complex expressions, you use the **Expression builder**, which provides editing tools and a library of functions grouped by categories.

To create an expression on an output column:

- In the salesReport table, click the output column **customerSales.CustomerAddress**. The [...] button appears on the right.

- b. Click [...] to open the Expression builder. On the left, the expression field already contains the column you selected.

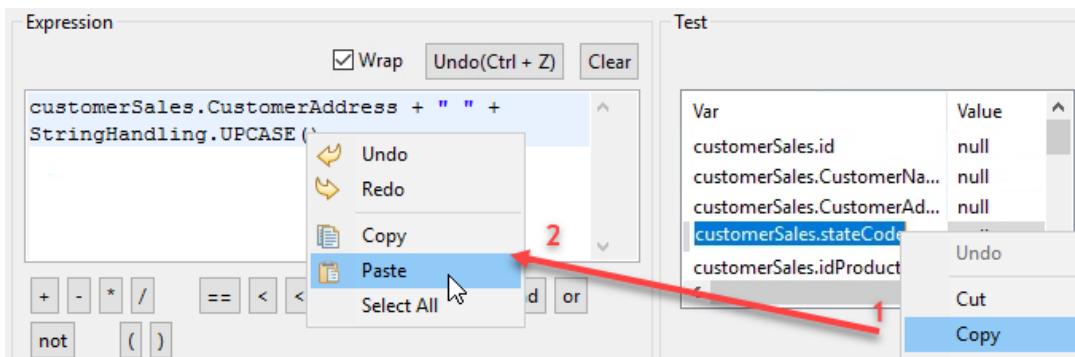
- c. Edit the field to change the expression to `customerSales.Address + " "`. You can now add a function that sets the state codes to uppercase.
d. In **Categories**, select **StringHandling**.
e. In **Functions**, scroll down and select **UPCASE**. Double-click **UPCASE** to add it to your expression. The expression looks like this:

- f. Delete the default "Hello" string and replace it with the variable `customerSales.stateCode`, which is of type String. When you are finished, the expression looks like this:

- g. Click **OK** to close the Expression builder.

TIP:

To avoid typing mistakes while using variables in expressions, you can directly select a variable from the list, right-click to copy and then paste in the expression editing box in the right location.



The **CustomerAddress** column in the **salesReport** table now contains both the address and the state code. The **stateCode** column is not useful anymore, you learn in the following how to delete it.

3. DELETE COLUMNS

To update the table schema (add, delete, move columns etc.), you can use the **Schema editor**, located at the bottom of the **tMap** editor. In the bottom right you find the **salesReport** schema.

To delete the **stateCode** column from the **salesReport** schema:

- Point the mouse to the left of the **stateCode** column. An arrow appears and the column is highlighted.

salesReport										
Column	Key	Type	N..	Date Pattern (Ctrl+Space)	Length	Precision	Default	Comment		
<input checked="" type="checkbox"/> id	<input checked="" type="checkbox"/>	Integer	<input checked="" type="checkbox"/>		2	0				
CustomerName	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		30	0				
CustomerAddress	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		35	0				
stateCode	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		2	0				
idProduct	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>		2	0				
RegisterTime	<input type="checkbox"/>	Date	<input checked="" type="checkbox"/>	"yyyy-MM-dd hh:mm:ss"	23	0				
Shop1	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	0				

- Click once to select the **stateCode** column. The column is highlighted in blue.
- Click the **Remove selected items** button

salesReport										
Column	Key	Type								
<input checked="" type="checkbox"/> id	<input checked="" type="checkbox"/>	Integer								
CustomerName	<input type="checkbox"/>	String								
CustomerAddress	<input type="checkbox"/>	String								
stateCode	<input type="checkbox"/>	String								
idProduct	<input type="checkbox"/>	Integer								
RegisterTime	<input type="checkbox"/>	Date								
Shop1	<input type="checkbox"/>	Float								

WARNING: Make sure you click the **Remove selected items** icon in the schema editor of **salesReport**, not to be confused with the similar **Remove selected output table** icon in the upper right part of the map editor. **Remove selected output table** deletes the entire output table!

You have configured the output from the **customerSales** table. You now want to join data coming from the Products database and update the output you expect in salesReport accordingly.

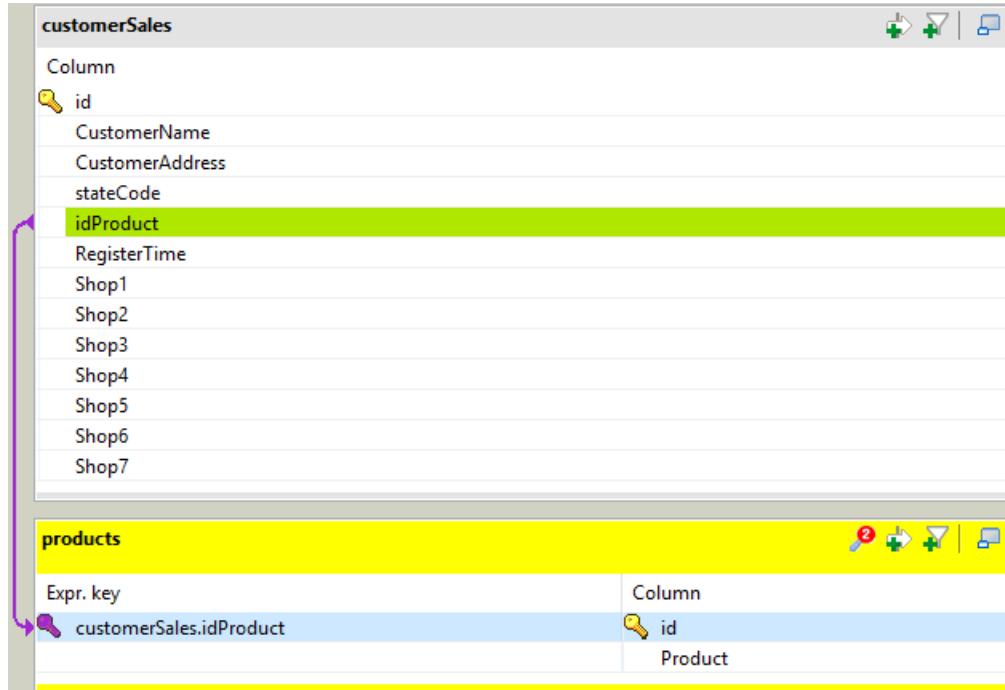
Configuring a join model in tMap

You can use **tMap** to create a join model that combines rows from the input tables, based on a related column between them. In this case, the related column is the product ID. As explained earlier, you are interested only in customers who bought products listed in the **Products** table, as they are the target of the marketing campaign. Therefore, you want to use an inner join that returns only records existing in both input tables.

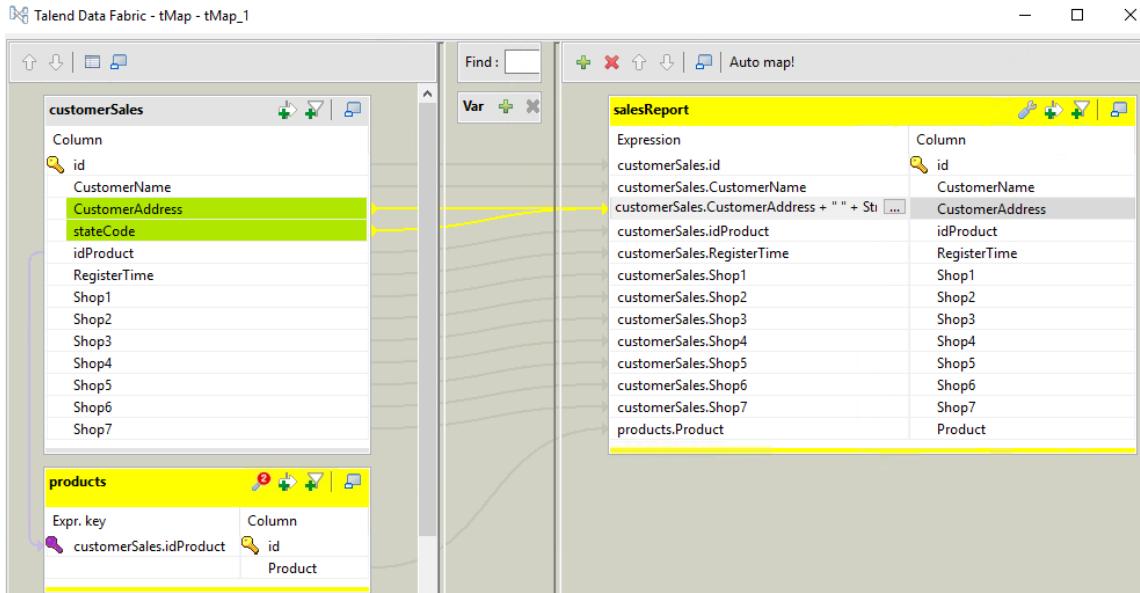
1. CREATE A JOIN

To create a join:

- Select the **idProduct** column in **customerSales**.
- Drag and drop it on the left side of the **id** column in the **products** table. This associates the **id** columns from the two tables.



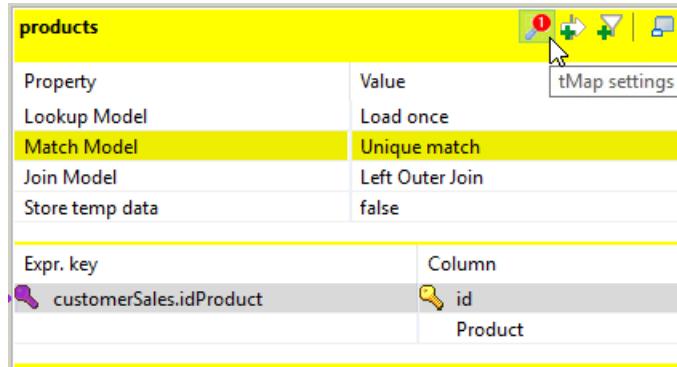
- Drag the **Product** column from the **products** table and drop it in the **salesReport** table.



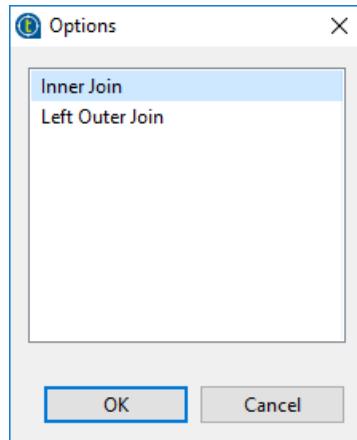
2. CONFIGURE THE JOIN MODEL

To create a join:

- On the **products** table toolbar, click the **tMap settings** icon. Notice that the default **Join Model** is **Left Outer Join**.



- Click **Left Outer Join**. Click the [...] button that appears in the right and in the dialog box that opens, select **Inner Join** and click **OK**.



- c. Make sure the **Join Model** is set to **Inner Join**.
- d. In the tMap editor window, click **OK** to save the changes.
- e. When asked if you want to propagate the changes, click **Yes**.

When you make changes to the schema of the output connection, you are prompted to propagate the changes (in other words, reflect those changes in the schema of the component on the other end of the connection). In the following you confirm that the schema was actually updated.

3. CHECK THE OUTPUT FILE SCHEMA

To check the schema:

- a. Double-click the **SalesReport** component to open the **Component** view.
- b. To the right of **Edit schema**, click [...].

A screenshot of the Spoon IDE's Component view for a 'SalesReport(tFileOutputDelimited_1)' component. The title bar shows 'Job(SalesJob 0.1)', 'Contexts(SalesJob)', 'Component', and 'Run (Job SalesJob)'. The main area has a tab bar with 'Basic settings' selected, followed by 'Advanced settings', 'Dynamic settings', 'View', 'Documentation', and 'Validation Rules'. Under 'Basic settings', there are fields for 'Property Type' (set to 'Built-In'), 'File Name' ('C:/StudentFiles/salesOut.csv'), 'Row Separator' ('\n'), and checkboxes for 'Append', 'Include Header', and 'Compress as zip file'. Below these is a 'Schema' dropdown set to 'Built-In' and an 'Edit schema' button, which is highlighted with a red rectangle. Other buttons in the row include '...', 'Sync columns', and a small icon.

- c. See if the schemas of the tMap output table and **SalesReport** are similar. If so, click **OK**.

Column	Key	Type	N.	Date ...	Le...	Pr...
id	<input checked="" type="checkbox"/>	Integer	<input checked="" type="checkbox"/>		2	0
CustomerName	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		30	0
CustomerAddress	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		35	0
idProduct	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>		2	0
RegisterTime	<input type="checkbox"/>	Date	<input checked="" type="checkbox"/>	"yyyy...	23	0
Shop1	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	0
Shop2	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	5
Shop3	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	0
Shop4	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	0
Shop5	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	5
Shop6	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	0

Column	Key	Type	N.	Date ...	Le...	Pr...
id	<input checked="" type="checkbox"/>	Integer	<input checked="" type="checkbox"/>		2	0
CustomerName	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		30	0
CustomerAddress	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>		35	0
idProduct	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>		2	0
RegisterTime	<input type="checkbox"/>	Date	<input checked="" type="checkbox"/>	"yyyy...	23	0
Shop1	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	0
Shop2	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	5
Shop3	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	0
Shop4	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	0
Shop5	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	5
Shop6	<input type="checkbox"/>	Float	<input checked="" type="checkbox"/>		11	0

NOTE:

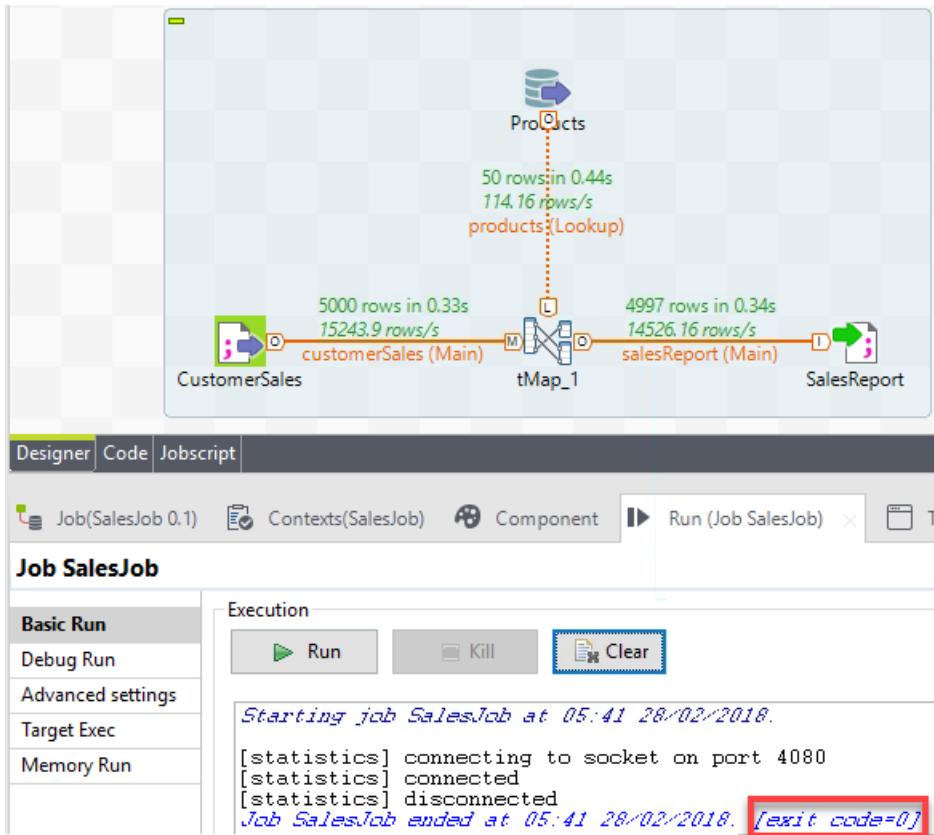
The **Sync columns** button in the **Component** view allows you to synchronize two schemas in one click whenever necessary.

Basic settings	Property Type	Built-In
Advanced settings	<input type="checkbox"/> Use Output Stream	
Dynamic settings	File Name: "C:/StudentFiles/salesOut.csv"	
View	Row Separator: "\n"	
Documentation	<input type="checkbox"/> Append <input type="checkbox"/> Include Header <input type="checkbox"/> Compress as zip file	
Validation Rules	Schema	<input type="button" value="Edit schema"/> <input type="button" value="Sync columns"/>

4. RUN THE JOB

To run the Job:

- In the **Run** view, click **Run**
- Make sure the Job executes correctly (**[exit code=0]**).



Notice that after execution, the number of processed rows transmitted from one component to another is displayed on the links. You can see that only 4997 of 5000 rows received at the tMap input were output in SalesReport. Some of the rows were rejected by the inner join.

5. VISUALIZE THE REPORT USING THE DATA VIEWER

To visualize data using the **Data viewer**:

- a. Right-click the **SalesReport** component.
- b. Select **Data viewer**. A window displays the results of the mapping.

The screenshot shows a Data Preview window titled "Data Preview: tFileOutputDelimited_1". The window has two tabs at the top: "Result Data Preview" (which is selected) and "File Content". Below the tabs, there are input fields for "Rows/page:" (set to 30) and "Limits:" (set to 1000). The main area contains a table with the following data:

Null	Condition	id	CustomerName	CustomerAddress
	*	1	Griffith Paving and Sealcoatin	talend@bordeaux AR
	*	2	Bill's Dive Shop	511 Maple Ave. Apt. 1B KS
	*	3	Childress Child Day Care	662 Lyons Circle AL
	*	4	Facelift Kitchen and Bath	220 Vine Ave. ME
	*	5	Terrinni & Son Auto and Truck	770 Exmoor Rd. AZ
	*	6	Kermit the Pet Shop	1860 Parkside Ln. RI
	*	7	Tub's Furniture Store	807 Old Trail Rd. DE
	*	8	Toggle & Myerson Ltd	618 Sheriden rd. CA
	*	9	Childress Child Day Care	788 Tennyson Ave. NJ
	*	10	Elle Hypnosis and Therapy Cent	2032 Northbrook Ct. AL
	*	11	Lennox Air Pollution Control	4522 N. Greenview Apt. 1B
	*	12	Keyth Contracting and Repair	1547 Knolwood Rd. HI
	*	13	Park District Of America	2678 Sheridan Rd. WV
	*	14	Nirabi Auto Service	1915 Lewis Ln. Apt 13 NV
	*	15	Darcy Frame and Matting Servic	1633 McGovern place GA

At the bottom of the window, there are buttons for "first", "previous", "next", "last", and "1 page of 34". Below the table, there are "Set parameters and continue" and "Close" buttons.

- c. Look at the data and click **Close**.

NOTE:

Alternatively, you can use a text editor like Notepad++ to visualize the resulting file located at C:/StudentFiles/salesOut.csv.

Capturing the join rejects

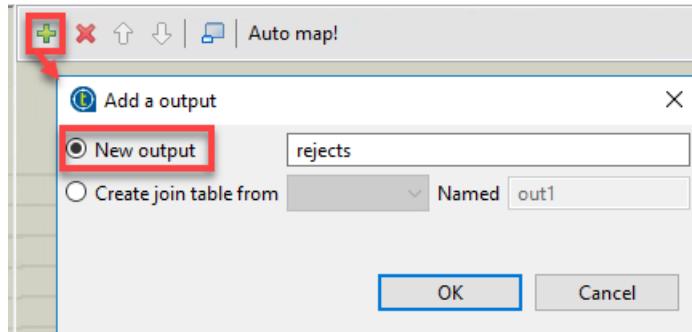
In the following section, you investigate which rows are rejected to make sure you did not exclude relevant data.

1. ADD AN OUTPUT TABLE TO CAPTURE JOIN REJECTS

To capture join rejects, create a second output table in the tMap configuration:

- a. Double-click **tMap_1** to open the mapping editor. Remember that the output tables are listed to the right of the editor window

- b. On the main toolbar on the right, click the **Add output table** button



- c. In the **New output** text box, enter the table name *rejects* and click **OK**. The rejects table is created.

2. CONFIGURE THE REJECTS TABLE

To configure the **rejects** table, you must define the table rows and set the **Catch lookup inner join rejects** property to **true**. You want the rejects table to contain the customer name and the product id columns. To define the table rows, you repeat the same process you used when creating the first output table.

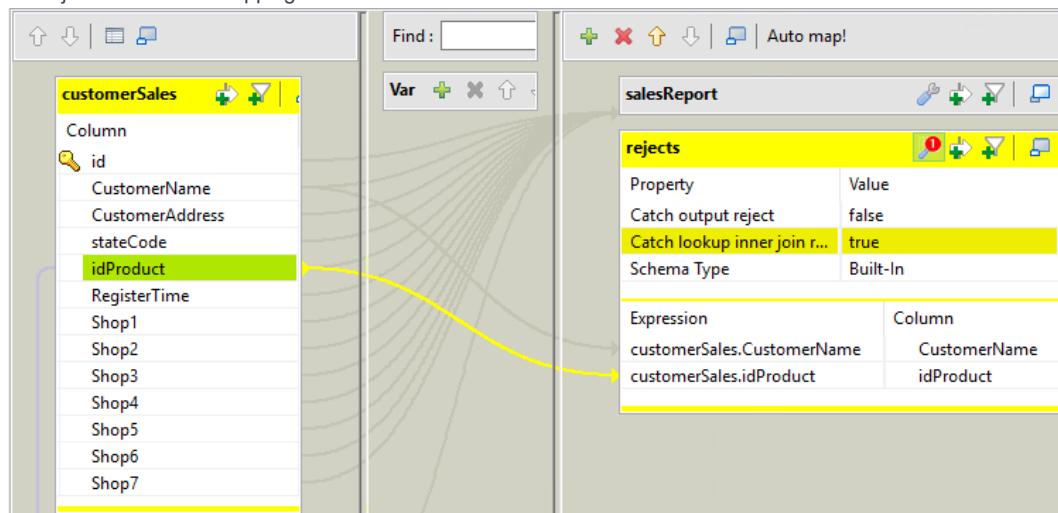
Follow these steps:

- In the **customerSales** table, select **CustomerName** and drag and drop it in the **rejects** table as soon as the blue arrow is positioned correctly.
- In the **customerSales** table, drag **idProduct** and drop it in the **rejects** table
- On the **rejects** table toolbar, click the **tMap settings** icon. Notice that the **Catch lookup inner join rejects** is set to **false**.

rejects	
Property	Value
Catch output reject	false
Catch lookup inner join reject	false
Schema Type	Built-In
Expression	Column
customerSales.CustomerName	CustomerName
customerSales.idProduct	idProduct

- To the right of **Catch lookup inner join reject**, click **false**. The [...] button appears to the right.

- e. Click on [...] and in the dialog that opens, select **true** and click **OK**. This allows you to catch the unmatched rows in the rejects table. Your mapping looks like this:



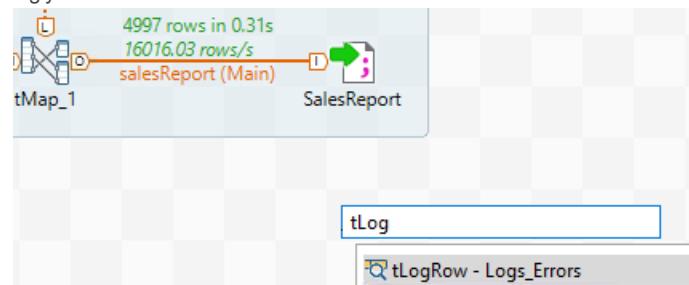
- f. In the tMap editor window, click **OK**.

You have configured the rejects table. You are ready to send the output to the console.

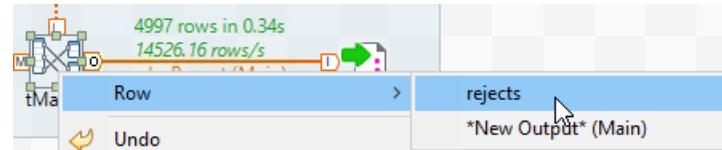
3. LOG REJECTED ROWS TO THE CONSOLE

To log failures to the console, use a **tLogRow** component:

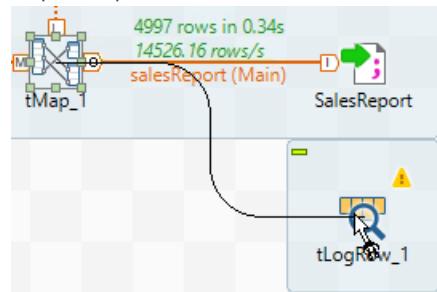
- a. In the **Designer**, under the **SalesReport** component, start entering **tLogRow** to bring up the components matching your search.



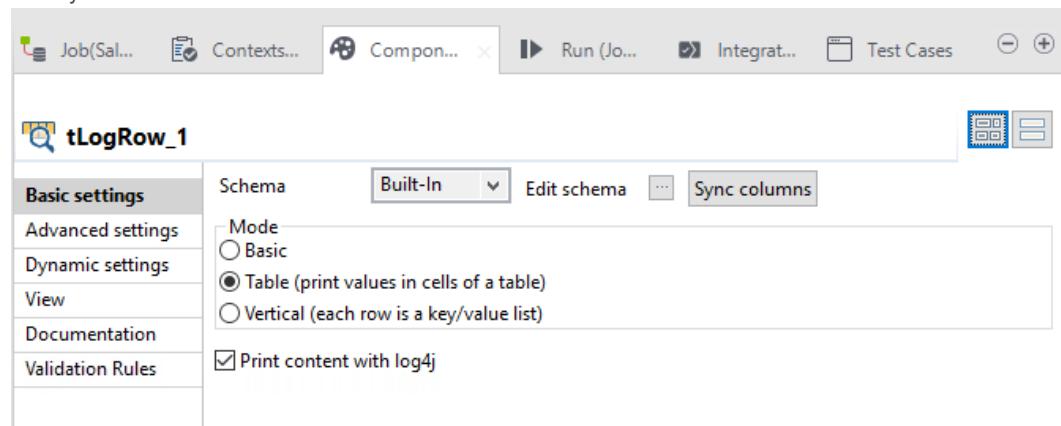
- b. Double-click **tLogRow** on the list. The component appears in the workspace.
 c. Right-click the **tMap_1** component and select **Row > rejects**.



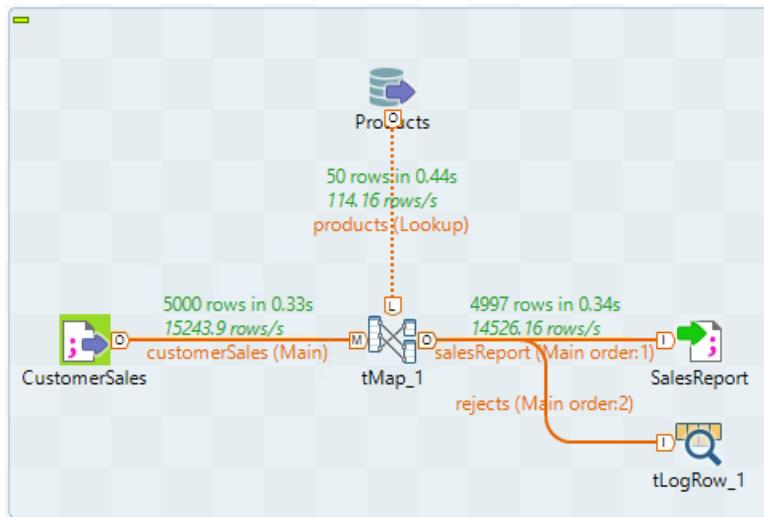
- d. A connection link is visible between tMap_1 and the mouse pointer. You can now drag the mouse pointer over to the output component.



- e. Click **tLogRow_1**. Notice that a **rejects** row is created between the two components.
f. Double-click **tLogRow_1** to open the **Component** view and select **Table**. This allows you to print results in a user friendly format.



Your Job looks like this:



You can now run the Job and see the results.

4. RUN THE JOB

To run the Job:

- On the **Run** tab, click **Run**
- Make sure the Job executed correctly (**[exit code=0]**)

```

Starting job SalesJob at 06:46 13/03/2018.
[statistics] connecting to socket on port 4075
[statistics] connected
+-----+
| tLogRow_1 |
+-----+
| CustomerName | idProduct |
|-----|
| Rythmics Ltd. | 58      |
| Jay's Cabinet Company | 58      |
| Pivot Point College | 58      |
+-----+
[statistics] disconnected
Job SalesJob ended at 06:46 13/03/2018. [exit code=0]

```

Line limit Wrap

Notice that three rows were sent to the log component and displayed in the console. For each row, the name of the customer and product ID are shown. As in the Products database, you have no product with identifier 58, so you can now understand why these rows were rejected.

Next step

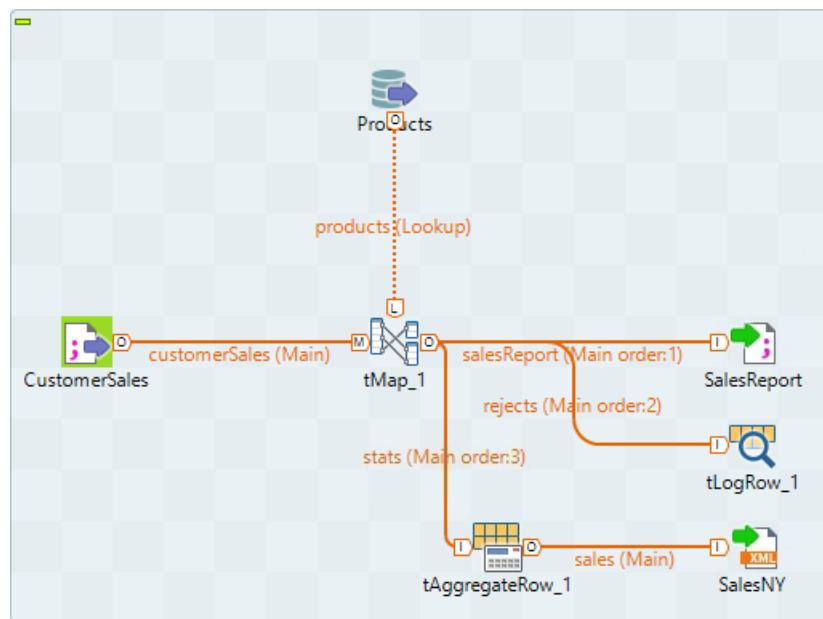
You have joined data using the tMap component and generated a report containing the join results. You have also investigated the rejects to make sure you do not omit relevant data. You are ready to learn more about transforming data.

Transforming Data

Task outline

In this exercise you learn to aggregate data, apply filters, and compute some statistics for the New York (NY) state in order to identify the best-selling products. At the end, you want to generate a report in a format that can be further processed by another system.

This section takes about 35 minutes to complete. At the end, the Job looks like this:



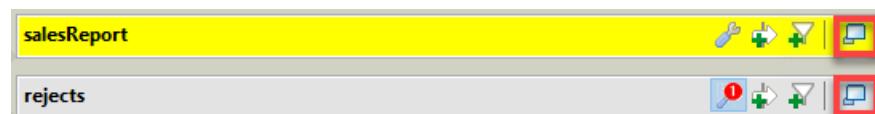
Aggregating and filtering data

You want to compute some statistics for the New York (NY) state in order to identify the best-selling products. To achieve this, you must create a new output table (*stats*) that contains the product name and customer ID coming from the input tables. However, you need to create a filter on the state so that only products sold in New York are considered. Then you use the **tAggregateRow** component to group data by product name and count the number of customer sales per product. You want to print the results in an XML file that can be transferred to another system for processing.

NOTE:

tAggregateRow: Receives a flow and aggregates it based on one or more columns. For each output line, are provided the aggregation key and the operations to compute (e.g. min, max, sum).

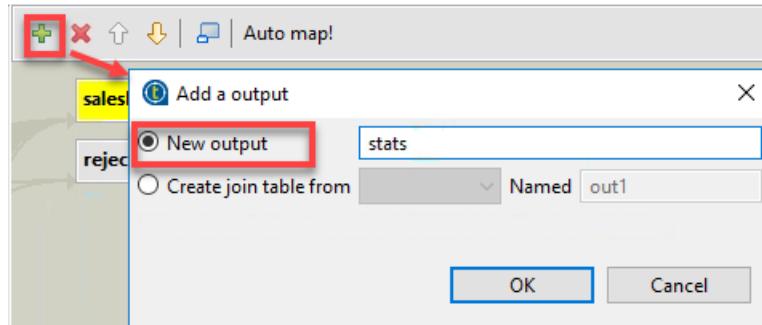
Before starting the exercise, double-click **tMap_1** to open the mapping editor. On the **salesReport** table toolbar, on the right, click the **Minimize** button to close the table editor. Repeat this step for the **rejects** table. Now you can focus on creating the new output table.



1. ADD THE STATS TABLE IN tMap

To create the stats table in the mapping editor:

- On the toolbar situated on the right, click the **Add output table** button.



- In the **New output** text box, enter the table name **stats** and click **OK**. The stats table is created.
- CONFIGURE THE STATS TABLE**
- To configure the **stats** table, you must define the table rows:
- In the **customerSales** table, select **id**. Drag it to the **stats** table and drop it as soon as the blue arrow is positioned correctly.
 - In the **products** table, select **Product** and drop it to the **stats** table.

stats	
Expression	Column
customerSales.id	id
products.Product	Product

You are ready to create a filter on the state column so that you only consider products sold in New York.

- APPLY A FILTER USING THE FILTER EXPRESSION IN tMap**

You want to create the filter expression:

```
"ny".equals(customerSales.stateCode)
```

To apply filters to data:

- On the stats table toolbar, click the **Enable/disable expression filter** button. The in-line editor opens.

stats	
Expression	
customerSales.id	id
products.Product	Product

- In the **customerSales** table, select and drag the **stateCode** column into the in-line editor.
- In the in-line editor, complete the expression `"ny".equals(customerSales.stateCode)`. Alternatively, you can click [...] on the right to open the Expression builder dialog box that you used earlier.

stats	
"ny".equals(customerSales.stateCode)	
Expression	Column
customerSales.id	id
products.Product	Product

- In the tMap editor window, click **OK**.

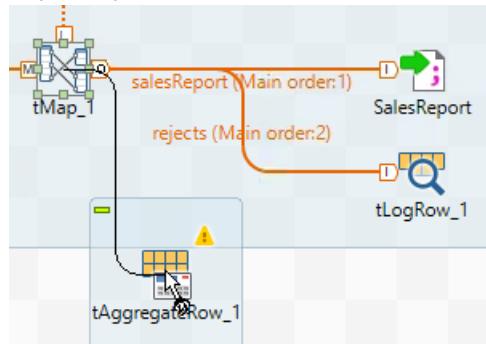
- ADD A tAggregateRow COMPONENT**

To add a **tAggregateRow** component:

- In the **Designer**, under the **tMap_1** component, start entering **tAggregateRow** to bring up components matching your search.
- Double-click **tAggregateRow** on the list. The component appears in the workspace.
- Right-click the **tMap_1** component and select **Row > stats**.



- A connection link is visible between **tMap** and the mouse pointer. You can now drag the mouse pointer over to the output component.

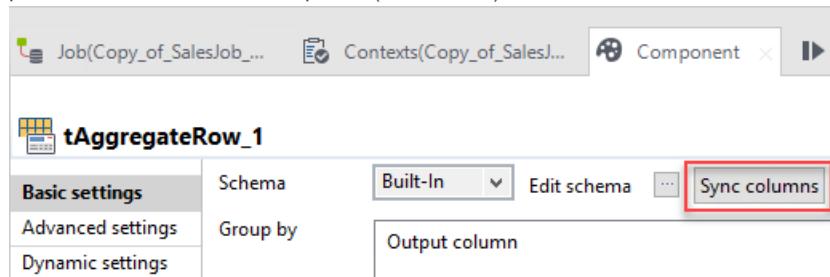


- Click **tAggregateRow_1**. Notice that a **stats** link is created between the two components.

5. CONFIGURE THE tAggregateRow SCHEMA

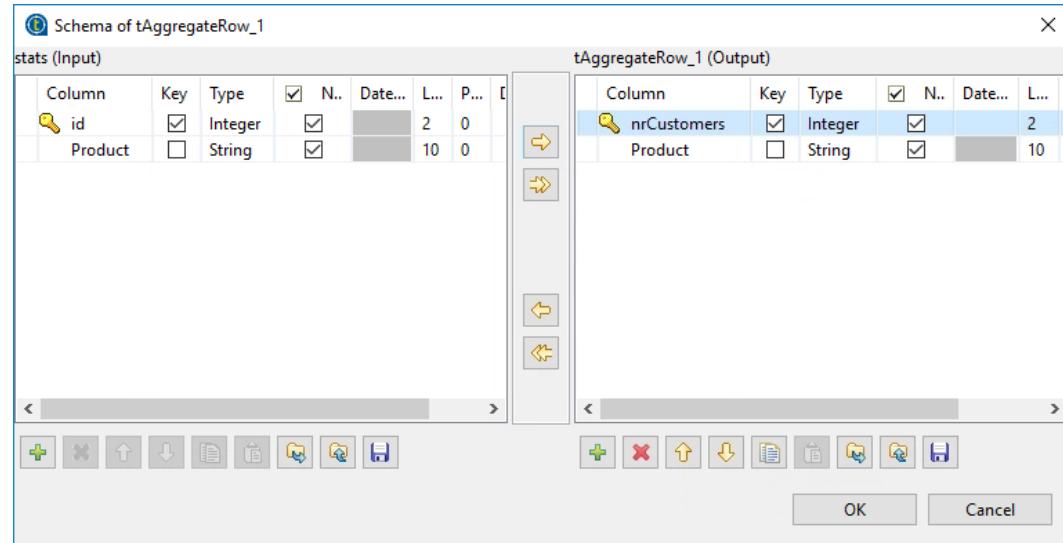
To configure the **tAggregateRow** schema:

- Double-click **tAggregateRow_1** to open the **Component** view.
- Click the **Sync columns** button. Remember that this button allows you to synchronize the schema of your component with the schema of the input link (the **stats** link).

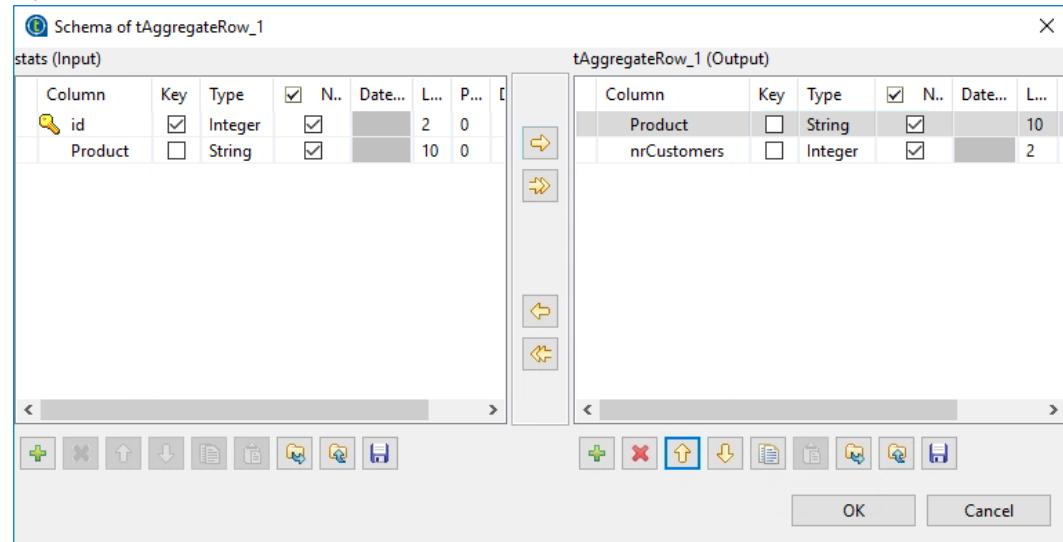


- To the right of **Edit schema**, click [...] and notice that the two schemas are exactly the same.

- d. In the Output schema **tAggregateRow_1**, click **id**. Rename the column **nrCustomers**. This gives you a more user friendly display in the file report you create.



- e. To the right of the **nrCustomers** column, deselect **Key**. You don't need this column to be a primary key.
f. In the Output schema **tAggregateRow_1**, select the **Product** column. In the toolbar at the bottom, click the **Move up selected items** button, to place the Product column on top. This gives you a more user friendly display in the file report.



g. Click **OK** to close the window.

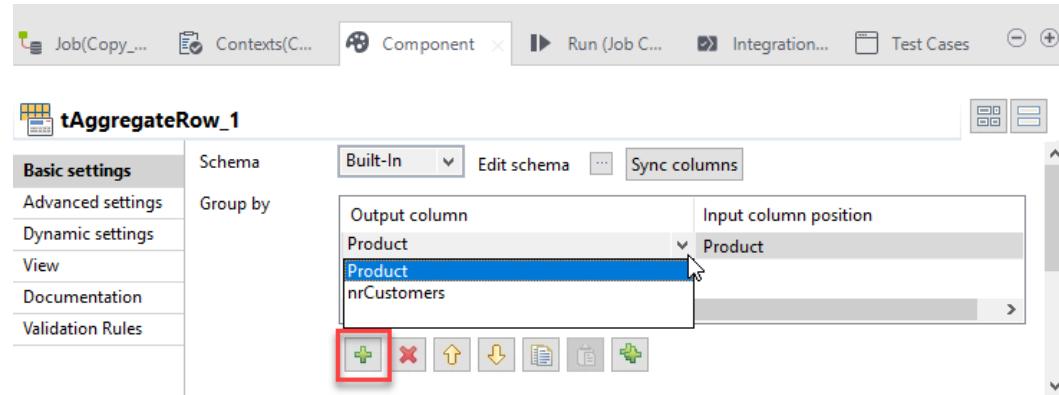
6. CREATE A GROUP BY CLAUSE WITH AGGREGATE FUNCTIONS

To create a GROUP BY clause with aggregate functions, you configure the Output column and the Input column position in the Group by and Operations tables.

Output column: allows you to select columns from the output schema

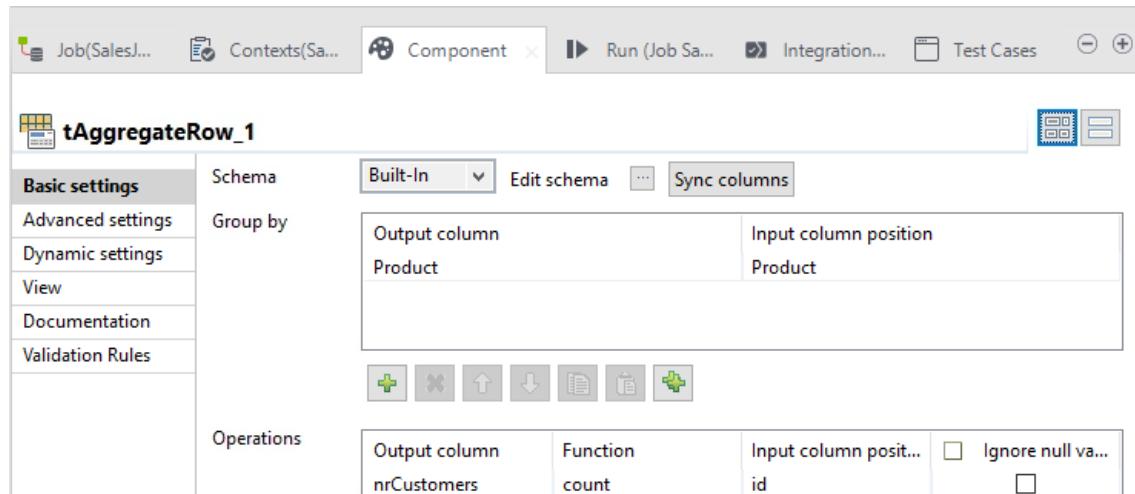
Input column position: allows you to select columns from the input schema

- a. Below the **Group by** section, click the **[+]** button. Notice that the Output column and the Input column position both display the **Product** column by default. If this is not the case, just click on the column to select Product, as shown in the screenshot.



- b. For **Operations**, click the **[+]** button. Make sure the **Output column** is set to *nrCustomers* and the **Input column position** to *id*

When you are finished, the configuration looks like this:



Generating the NY sales report

You want to print the results of the aggregation in an XML file that can be transferred to another system for processing. For this purpose you can use the **tFileOutputXML** component.

NOTE:

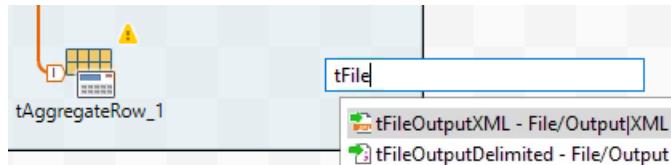
tFileOutputXML: Writes an XML file with separated data values according to a defined schema.

1. ADD AN XML FILE TO WRITE THE RESULTS

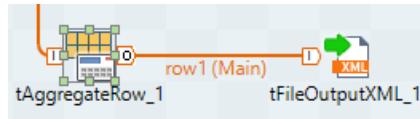
To write the results to an XML file, your report, you can use a **tFileOutputXML** component.

To add a **tFileOutputXML** component:

- In the **Designer**, to the right of the **tAggregateRow_1** component, start entering **tFileOutputXML** to bring up components matching your search.



- b. Double-click **tFileOutputXML** on the list. The component appears in the workspace.
- c. Right-click the **tAggregateRow_1** component and select **Row > Main**.
- d. A connection link is visible between tMap and the mouse pointer. You can now drag the mouse pointer over to the output component.
- e. Click **tFileOutputXML_1**. Notice that a link is created between the two components.

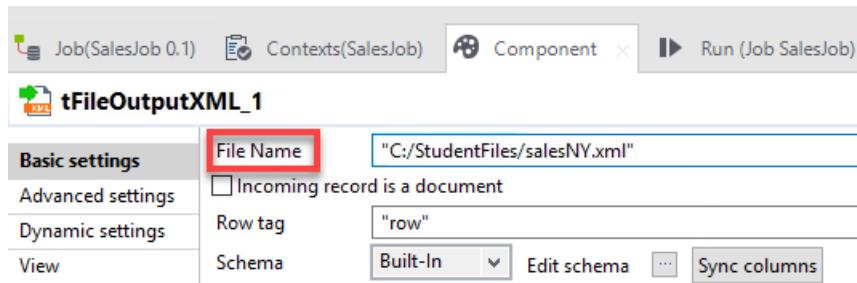


- f. Click on **row1** to edit and rename the link sales

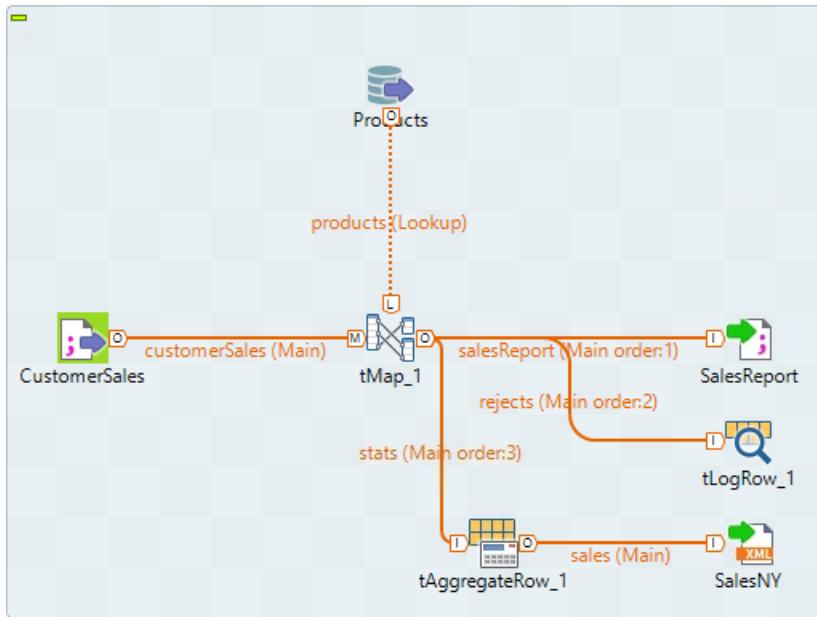
2. CONFIGURE THE OUTPUT COMPONENT

To configure the **tFileOutputXML_1** component:

- a. Double-click **tFileOutputXML_1** to open the **Component** view.
- b. In the **File Name** text box, enter "C:/StudentFiles/salesNY.xml"
- c. On the **View** tab, in the **Label format** text box, enter SalesNY to rename the file.



When you finish, your Job looks like this:



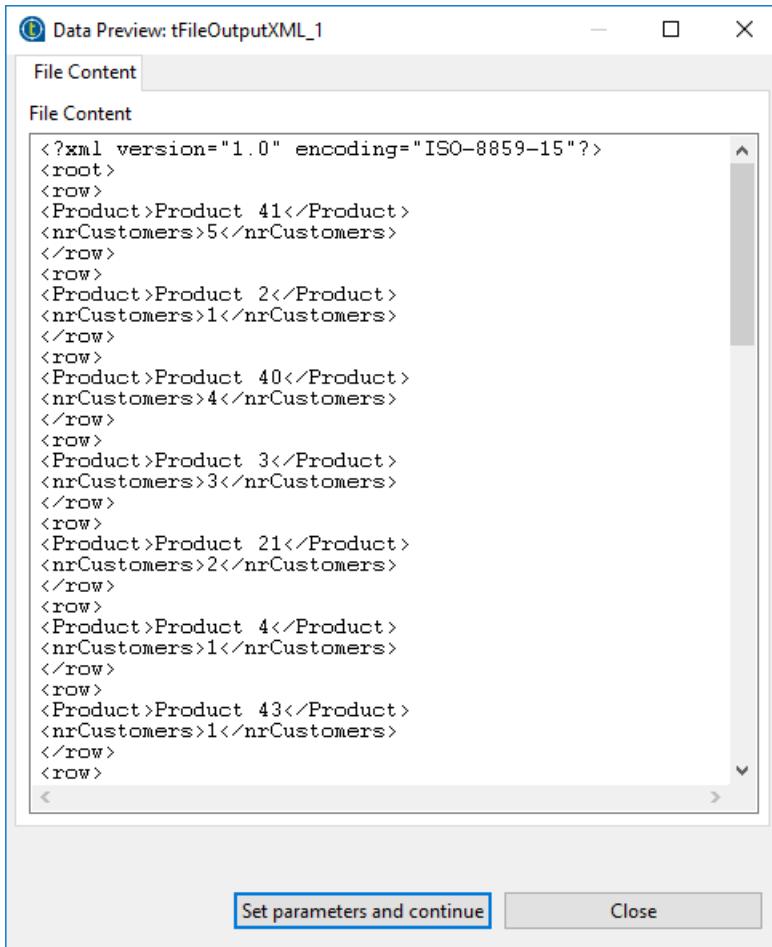
You can now run the job and generate the SalesNY.xml report.

3. RUN THE JOB

To run the Job:

- On the **Run** tab, click **Run**
- Make sure the Job executes correctly (**[exit code=0]**).

To check the results, right-click the **SalesNY** component and select **Data viewer**.



The screenshot shows a 'Data Preview' window titled 'tFileOutputXML_1'. The main area is labeled 'File Content' and displays the following XML code:

```
<?xml version="1.0" encoding="ISO-8859-15"?>
<root>
<row>
<Product>Product 41</Product>
<nrCustomers>5</nrCustomers>
</row>
<row>
<Product>Product 2</Product>
<nrCustomers>1</nrCustomers>
</row>
<row>
<Product>Product 40</Product>
<nrCustomers>4</nrCustomers>
</row>
<row>
<Product>Product 3</Product>
<nrCustomers>3</nrCustomers>
</row>
<row>
<Product>Product 21</Product>
<nrCustomers>2</nrCustomers>
</row>
<row>
<Product>Product 4</Product>
<nrCustomers>1</nrCustomers>
</row>
<row>
<Product>Product 43</Product>
<nrCustomers>1</nrCustomers>
</row>
<row>
```

At the bottom of the window are two buttons: 'Set parameters and continue' (highlighted with a blue border) and 'Close'.

Congratulations! You have designed and run your first data integration Job using Talend Studio.

Next step

You have almost finished this section. Time for a quick review.

Review

In this lesson, you have learned how to design and run your first data integration project using Talend Studio.

You followed the classical phases of a data integration chain. You extracted data from structured files using **tFileInputDelimited** and from databases using **tMysqlInput**. You explored the functionnalities offered by the **tMap** component to join and process data. You applied different transformations on the data such as filters and common aggregation techniques using the **tAggregateRow** component. You learned how to construct expressions using the inline editor and the Expression builder. Then you loaded the result in the appropriate target storage.

Intentionally blank